

| Network o | f | Networks | on | Impact                   | Evaluation                 | ł |
|-----------|---|----------|----|--------------------------|----------------------------|---|
|           |   |          |    | <ul> <li>Call</li> </ul> | 117.147.147.Lat" L.d" I da | × |

DAC Evaluation Network Evaluation Cooperation Group International Organization for Cooperation in Evaluation UN Evaluation Group

### IMPACT EVALUATIONS AND DEVELOPMENT

## **NONIE Guidance on Impact Evaluation**

Frans Leeuw

Jos Vaessen

Draft Version for Discussion at the Cairo conference March-April, 2009

### Acknowledgments

We would like to thank the steering committee of this project (Andrew Warner, David Todd, Zenda Ofir and Henri Jorritsma) for their suggestions. We also would like to thank Antonie de Kemp for exchanging ideas on design questions and Patricia Rogers for reviewing an earlier version. This document draws on previous impact evaluation guidance documents produced by NONIE members. We would like to thank the authors for their efforts, laying the groundwork for the present Guidance.

Frans Leeuw Jos Vaessen

## Table of contents

| Introduction   | 5        |
|--|----------|
| PART I - METHODOLOGICAL AND CONCEPTUAL ISSUES IN IMPACT EVALUATION                           | 8        |
| 1. Identify the (type and scope of the) intervention   | 8        |
| 1.1. The Impact Evaluation landscape and the scope of IE                                     | 8        |
| 1.2. Impact of what?   | 9        |
| 1.3. How to identify interventions   | 12       |
| 1.4. Impact on what?   | 13       |
| 1.5. Key Message   | 16       |
| 2. Agree on the objectives of the intervention that are valued                               | 17       |
| 2.1. Stakeholder values in IE  | 17       |
| 2.2. Intended versus unintended effects  | 18       |
| 2.3. Short- term versus long-term effects  | 18       |
| 2.4. The sustainability of effects   | 19       |
| 2.5. Key Message   | 19       |
| 3. Carefully articulate the theories linking interventions to outcomes                       | _20      |
| 3.1. Seeing interventions as theories: the black box and the contribution problem            | _20      |
| 3.2. Articulating intervention theories on impact  | 20       |
| 3.3. Testing intervention theories on impact   | 24       |
| 3.4. Key message   | _25      |
| 4. Address the attribution problem   | _26      |
| 4.1. The attribution problem   | 26       |
| 4.2. Methodological approaches addressing the attribution problem                            | -28      |
|  | _30      |
| 4.2.2. Propensity score matching   |          |
| 4.2.3. Flopensity score matching   | יכ       |
| 4.2.4. Double difference (difference in difference)  | 22       |
| 4.2.5. Bouble difference (difference)  | <br>     |
| 4.2.0. Regression analysis and double difference   | 24<br>24 |
| 4.2.8 Regression Discontinuity Analysis  | 25<br>25 |
| A 3 Applicability of quantitative methods for addressing the attribution problem             | رر<br>۶6 |
| 4.4. Other approaches  | تر<br>8د |
| 4.4.1. Participatory approaches  | <br>38   |
| 4.4.2. Useful methods for data collection and analysis which are often part of IE designs    | <br>40   |
| 4.5. Key message   | 41       |
| 5. Build on existing knowledge relevant to the impact of interventions                       | '<br>    |
| 5.1. Review and synthesis approaches as methods for analyzing existing evidence on impact    | ·<br>42  |
| 5.2. Key message   | <br>44   |
| 6. Use a mixed methods approach: the logic of the comparative advantages of methods          | 46       |
| 6.1. Different methodologies have comparative advantages in addressing particular concerns a | nd       |
| needs in impact evaluation   | _46      |
| 6.2. Advantages of combining different methods and sources of evidence                       | 47       |
| 6.3. Key message   | 50       |
| PART II – MANAGING IMPACT EVALUATIONS  | 51       |
| 7. Determine if an IE is feasible and worth the cost   | 51       |
| 7.1. Evaluability  | 51       |
| 7.2. Key message   | 52       |
| 8. Start early – getting the data  | 53       |
| 8.1. Timing of data collection   | 53       |
| 8.2. Data availability   | 53       |
| 8.3. Quality of the data   | 55       |
| 8.4. Dealing with data constraints   | 56       |
| 8.5. Key message   | 57       |
| 9. Front-end planning is important   | 59       |

| 9.1. Front-end planning | 59 |
|-------------------------|----|
| 9.2. Key message        | 63 |
| References              | 64 |
| Appendices              | 73 |

## Introduction

The Network of Networks for Impact Evaluation (NONIE) was established in 2006 to foster more and better impact evaluations by its membership. NONIE uses the DAC definition, defining impacts as (OECD-DAC, 2002: 24) "the positive and negative, primary and secondary long-term effects produced by a development intervention, directly or indirectly, intended or unintended. These effects can be economic, sociocultural, institutional, environmental, technological or of other types".

The impact evaluations pursued by NONIE are expected to reinforce and complement the broader evaluation work by NONIE members. The DAC definition refers to the *"effects produced by"*, stressing the attribution aspect. This implies an approach to impact evaluation which is about attributing impacts rather than assessing what happened. In most contexts, adequate empirical knowledge about the effects produced by an intervention requires at least an accurate measurement of what *would have occurred* in the absence of the intervention and a comparison with *what has occurred* with the intervention implemented.

The purpose of NONIE is to promote more and better impact evaluations among its members. Issues relating to evaluations in general are more effectively dealt with in the parent networks, and are thus not the primary focus of NONIE. NONIE will focus on sharing of methods and learning-by-doing to promote the practice of IE (impact evaluation). The current Guidance document was developed for supporting those purposes.

Impact evaluation in development assistance has received considerable attention over the past few years. The major reason is that many outside of development agencies believe that achievement of results has been poor, or at best not convincingly established. Many development interventions appear to leave no trace of sustained positive change after they have been terminated and it is hard to determine the extent to which interventions are making a difference. However, the development world is not 'alone' in attaching increasing importance to impact evaluations. In fields like crime and justice, education and social welfare, IEs have over the last decade become more and more important<sup>1</sup>. Evidence-based (sometimes 'evidence-informed') policies are high on the (political) agenda and some even refer to the 'Evidence Movement' (Rieper et al, 2009). This includes the development of knowledge repositories where results of impact evaluations are summarized. In some fields like criminology and in some professional associations like the Campbell Collaboration, methodological standards and scales are used to grade IEs<sup>2</sup>, although not without discussion (Leeuw, 2005; Worral, 2002; 2007).

<sup>&</sup>lt;sup>1</sup> The history of impact evaluations in some countries goes back many decades (Oakley, 2000). <sup>2</sup> The Maryland Scientific Methods Scale (MSMS) is for example used in parts of criminology and in several countries. See Leeuw (2005). RCTs (randomized controlled trials) are believed to be the top design (level 5).

Important reasons for doing impact evaluations are the following:

- IEs provide evidence on 'what works and what doesn't' (under what circumstances) and how large is the impact. As OED (2005) puts it: measuring outcomes and impacts of an activity and distinguishing these from the influence of other, external factors is one of the rationales behind IE;
- Measuring impacts and relating the changes in dependent variables to developmental policies and programs is not something that can be done 'from an armchair'. IE is *the* instrument for these tasks;
- Impact evaluation can gather evidence on the sustainability of effects of interventions;
- IEs produce information that is relevant from an accountability perspective; IEs disclose knowledge about the (societal) effects of programs which can be linked to the (financial) resources used to reach these effects;
- Individual and organizational learning can be stimulated by doing impact evaluations. This is true for organizations in developing countries, but also for donor organizations. Informing decision makers on whether to expand, modify or eliminate projects, programs and policies is linked to this point, as is the OED (2005) argument that IEs enable sponsors, partners and recipients to compare the effectiveness of alternative interventions.

We believe that the ultimate reason for promoting impact evaluations is to learn about 'what works and what doesn't and why' and thus contribute to the effectiveness of (future) development interventions. In addition to this fundamental motive, impact evaluations have a key role to play in the international drive for better evidence on results and development effectiveness. They are particularly well suited to answering important questions about whether development interventions made a difference (and how cost-effective they were). Well-designed impact evaluations also shed light on why an intervention did or did not work, which can vary across time and space.

Decision makers need better evidence on impact and its causes to ensure that resources are allocated where they can have most impact and to maintain future public funding for international development. The pressures for this are already strong and will increase as resources are scaled up for international development. Without such evidence there is a risk of the case for aid and future funding sources being undermined.

Using the word 'effects' and 'effectiveness' implies that the changes in the 'dependent variable[s]' that are measured within the context of an impact evaluation (IE), are *caused* by the intervention under study. The concept of 'goal achievement' is used when causality is *not* necessarily present. Goals can also be achieved *independent* from the intervention. Changes in financial or economic situations, in the world of health and agriculture or in other social conditions can

help realize goal achievement, even in a situation where the 'believed-to-be-effective' intervention under review, is not working.

The question whether or not IE should always attempt to measure all possible impacts is not easy to answer. IE involves finding the appropriate balance between the desire to understand and measure the full range of effects in the most rigorous manner possible, and the practical need to delimit and prioritize on the basis of interests of stakeholders as well as resource constraints.

#### Key issues addressed in the Guidance Document

The Guidance is structured around nine key issues in impact evaluation:

- 1. Identify the (type and scope of the) intervention
- 2. Agree on the objectives of the intervention that are valued
- 3. Carefully articulate the theories linking interventions to outcomes
- 4. Address the attribution problem
- 5. Build on existing knowledge relevant to the impact of interventions

6. Use a mixed methods approach: the logic of the comparative advantages of methods

- 7. Determine if an IE is feasible and worth the cost
- 8. Start early getting the data
- 9. Front-end planning is important

The discussion of these nine issues constitutes the structure of this Guidance document. The first part, comprising the first six issues deals with methodological and conceptual issues in IE and constitutes the core of the Guidance document. In addition, a shorter second part focuses on managing IE and addresses aspects of evaluability, benefits and costs of IE and planning.

There is no universally accepted definition of rigorous impact evaluation. There are some who equate rigorous IE with particular methods and designs. In our view rigorous impact evaluation is more than a methodological design. Rigorous impact evaluation requires addressing the issues described above in an appropriate manner, especially the core methodological and conceptual issues described in part I.

## PART I - METHODOLOGICAL AND CONCEPTUAL ISSUES IN IMPACT EVALUATION

### 1. Identify the (type and scope of the) intervention

#### 1.1. The Impact Evaluation landscape and the scope of IE

In international development, impact evaluation is principally concerned with final results of interventions (programs, projects, policy measures, reforms) on the welfare of communities, households and individuals (citizens, taxpayers and voters). Impact is often associated with progress at the level of the Millennium Development Goals, which primarily comprise indicators of welfare of (these) households and individuals. The renewed attention for results- and evidence-based thinking and ensuing interest in impact evaluation provides a new momentum for applying rigorous methods and techniques in assessing the impact of interventions.

There is today more than ever a 'continuum' of interventions. At one end of the continuum are relatively simple projects characterized by single 'strand' initiatives with explicit objectives, carried out within a relatively short timeframe, where interventions can be isolated, manipulated and measured. An impact evaluation in the agricultural sector for example, will seek to attribute changes in crop yield to an intervention such as a new technology or agricultural practice. In a similar guise, in the health sector, a reduction in malaria will be analyzed in relation to the introduction of bed nets. For these types of interventions, experimental and quasi-experimental designs may be appropriate for assessing causal relationships, along with attention to the other tasks of impact evaluation. At the other end of the continuum are comprehensive programs with an extensive range and scope (increasingly at country, regional or global level), with a variety of activities that cut across sectors, themes and geographic areas, and emergent specific activities. Many of these interventions address aspects that are assumed to be critical for effective development yet difficult to define and measure, such as human security, good governance, political will and capacity, sustainability, and effective institutional systems.

Some evidence of this continuum is provided in Appendix 1 in which two examples of IEs are presented, implemented at different (institutional) levels and based on divergent methodologies with different time frames (see also Figure 1.1.).

The endorsement in 2000 of the Millennium Development Goals by all heads of state, together with other defining events and occurrences has propelled new action which challenges development evaluation to enter new arenas. There is a shift away from fragmented, top-down and asymmetrical approaches.

Increasingly, ideals such as 'harmonization', 'partnership', 'participation', 'ownership' and 'empowerment' are being emphasized by stakeholders.

However, this trend in policy is not yet reflected in evaluative practices, including IE. Especially, institutional policies such as anticorruption policies but also regional and global policy networks and public private partnerships with their different forms and structures<sup>3</sup> appear to be less often part or goal of Impact Evaluations, when compared to (top-down driven) small programs for specific groups of beneficiaries. Ravallion (2008: 6) is of the opinion that there is "a 'myopia bias' in our knowledge, favoring *development projects that yield quick results*"<sup>4</sup>. In the promotion of more rigorous IE, development agencies, national governments, civil society organizations and other stakeholders in development should be aware of this bias in focus, keeping in mind the full range of policy interventions that (eventually) affect the welfare of developing societies.

Evaluating the impact of policies with their own settings and levels requires appropriate methodological responses. These can be usefully discussed under the banner of two key issues: the impact of what and the impact on what? These two issues point at a key challenge in impact evaluation: the scope of the impact evaluation.

#### 1.2. Impact of what?

What is the independent variable (intervention) we are looking at? In recent years, we have seen a broadening in the range of policy interventions which should/could be subject to impact evaluation.

One of the trends in development is that donors are moving up the *aid chain*. Whereas in the past donors were very much involved in 'micro-managing' their own projects and (sometimes) bypassing government systems, nowadays a sizeable chunk of aid is allocated to national support for recipient governments. Conditionality to some extent has shifted from micro-earmarking (e.g. donor money destined for an irrigation project in district *x*) to meso-earmarking (e.g. support for the agricultural sector) or macro-earmarking (e.g. support for the government budget to be allocated according to country priorities).

Besides a continued interest in the impact of individual projects, donors, governments and nongovernmental institutions are increasingly interested in the impact of comprehensive programs, sector strategies or country strategies, often comprising multiple instruments, stakeholders, sites of intervention and target groups.

<sup>&</sup>lt;sup>3</sup> An interesting overview of public private partnerships and their evaluation is given by Utce Ltd and Japan Pfi Association (2003).

<sup>&</sup>lt;sup>4</sup> "We probably also under-invest in evaluative research on types of interventions that tend to have diffused, wide-spread benefits" (Ravallion, 2008: 6). See also Jones et al. (2008) who have identified geographical and sectoral biases in IE.

There is a growing demand for assessing the impact of new instruments and modalities such as:

- international treaties governing the actions of multiple stakeholders (e.g. the Paris Declaration, the Kyoto Protocol);
- new aid modalities such as sector budget support (SBS) or general budget support (GBS);
- Instruments such as institutional capacity-building, institutional reform, partnership development and stakeholder dialogues at national or regional levels.

In most countries donor organizations are (still) the main promoters of IE. The shift of the unit of analysis to the macro and (government) institutional level requires impact evaluators to pay more attention to complicated and more complex interventions at national, sector or program level. Multi-site, multi-governance and multiple (simultaneous) causal strands are important elements of this (see Rogers, 2008).

At the same time, the need for more rigorous IE at the 'project level' remains urgent. The majority of aid money is (still) micro-earmarked money for particular projects managed by donors in collaboration with national institutions. Furthermore, the ongoing efforts in capacity-building on national M&E systems (see Kusek and Rist, 2005) and the promotion of country-led evaluation efforts stress the need for further guidance on IE at 'single' intervention level.

Earlier we referred to a continuum of interventions. At one end of the continuum are relatively simple projects characterized by single 'strand' initiatives with explicit objectives, carried out within a relatively short timeframe where interventions can be relatively easy isolated, manipulated and measured. Examples of these kinds of interventions include building new roads, repairing roads, reducing the price of fertilizer for farmers, providing clean drinking water at lower cost, etc. It is important to be precise in what the interventions are and what they focus on. In the case of new roads or the rehabilitation of existing ones, the goal often is a reduction in journey time and therefore reduction of societal transaction costs.

At the other end of the continuum are comprehensive programs with an extensive range and scope (increasingly at country, regional or global level), with a variety of activities that cut across sectors, themes and geographic areas, and emergent specific activities. Rogers (2008) has outlined several aspects of what constitutes complicated interventions (see Tables 1.1. and 1.2.): alternative and multiple causal strands, recursive causality and emergent outcomes.

Table 1.1. Aspects of complication in interventions

| Aspect of complication  | Simple intervention                         | Complicated intervention   |
|---|---|--|
| I. Governance and location  | Single organization                         | Multiple agencies, often<br>interdisciplinary and<br>cross-jurisdictional                              |
| <ol> <li>Simultaneous causal strands</li> <li>Alternative causal strands</li> </ol> | Single causal strand<br>Universal mechanism | Multiple simultaneous causal strands<br>Different causal mechanisms<br>operating in different contexts |

#### Source: Rogers (2008)

#### Table 1.2. Aspects of complexity in interventions

| Aspect of complexity  | Simple intervention                             | Complex intervention   |
|---|---|--|
| <ol> <li>Recursive causality and<br/>disproportionate effect</li> </ol> | Linear, constant dose–<br>response relationship | Recursive, with feedback loops, including<br>reinforcing loops; disproportionate<br>effects at critical limits |
| 2. Emergent outcomes  | Pre-identified outcomes                         | Emergent outcomes  |

#### Source: Rogers (2008)

Rogers (2008: 40) recently argued that "the greatest challenge [for the evaluator] comes when interventions have both complicated aspects (multi-level and multisite) and complex aspects (emergent outcomes)". These aspects often converge in interventions in the context of public private partnerships or new aid modalities which have become more important in the development world. Demands for accountability and learning about results at country, agency, sector or program and strategy levels are also increasing, which has made the need for appropriate methodological frameworks to assess their impact more pressing.

Pawson (2005) has distinguished five principles on complex programs that can be helpful when designing impact evaluations of aid:

- Locate key program components. Evaluation should begin with a comprehensive scoping study mapping out the potential conjectures and influences that appear to shape the program under investigation. One can envisage stage-one mapping as the hypothesis generator. It should alert the evaluator to the array of decisions that constitute a program, as well as providing some initial deliberation upon of their intended and wayward outcomes.
- Prioritize among program components. The general rule here is to concentrate on: i) those components of the program (intervention) theory which seem likely to have the most significant bearing on overall outcomes, and ii) those segments of program theory about which least is known.

- 3. Evaluate program components by subsets. This principle is about when and where to locate evaluation effort in relation to a program. The evaluation should take *on sub-sets* of program theory. Evaluation should occur in ongoing portfolios rather than one-off projects. Suites of evaluations and reviews should track program theories as and wherever they unfold.
- 4. Identify bottlenecks in the program network. "Theories-of-change" analysis perceives programs as implementation chains and asks, 'what are the flows and blockages as we put a program into action?' The basic strategy is to investigate how the implementation details sustain or hinder program outputs. The main analytic effort is directed at configurations made up of selected segments of the implementation chains across a limited range of program locations.
- 5. Provide feedback on the conceptual framework. What the theory-based approach initiates is a process of 'thinking though' the tortuous pathways along which a successful program has to travel. What would be described are the main series of decision points through which an initiative has proceeded and the findings would be put to use in alerting stakeholders to the caveats and considerations that should inform those decisions. The most durable and practical recommendations that evaluators can offer come from research that begins with theory and ends with a refined theory.

If interventions are complicated in that they have multiple active components, it is helpful to state these separately and treat the intervention as a package of components. Depending on the context, the impact of intervention components can be analyzed separately and/or as part of a package<sup>5</sup>. The separate analysis of intervention components implies interventions being *unpacked* in such a way that the most important social and behavioral mechanisms believed to make the 'package' work, are spelled out.

Although complex interventions are becoming more important and therefore should be subject to impact evaluation, this should not imply a reduction of interest in evaluating the impact of *relatively simple, single strand interventions*. The sheer number of these interventions makes doing robust IEs of great importance.

#### 1.3. How to identify interventions

To a large extent interventions can be identified and categorized on the basis of the main theme addressed. Examples of thematic areas of interventions are: roads and railroads, protected area management, alternative livelihoods and research on innovative practices.

<sup>&</sup>lt;sup>5</sup> For example Elbers et al. (2008) directly assess the impact of a set of policy variables (i.e. the equivalent of a multi-stranded program), by means of a regression-based evaluation approach (see section 4), on outcome variables.

A second way to identify interventions is to find out which *generic policy instruments and their combinations* constitute the intervention: economic incentives (e.g. tax reductions, subsidies), regulations (e.g. laws, restrictions), or information (e.g. education, technical assistance). As argued by authors such as Pawson (2006), Salamon (1981) and Vedung (1998), using this relative simple classification helps to identify the interventions. "Rather than focusing on individual programs, as is now done, or even collections of programs grouped according to major 'purpose' as is frequently proposed, the suggestion here is that we should concentrate on the generic tools of government that come to be used, in varying combinations in particular public programs" (Salamon, 1981: 256). Acknowledging the central role of policy instruments enables evaluators to take into account lessons from the application of particular (combinations of) policy interventions elsewhere (see Bemelmans-Videc and Rist, 1998).

#### 1.4. Impact on what?

This topic concerns the 'dependent variable problem'. Interventions often affect multiple institutions, groups and individuals. What level of impact should we be interested in?

The causality chain linking policy interventions to ultimate policy goals (e.g. poverty alleviation) can be relatively direct and straightforward (e.g. the impact of vaccination programs on mortality levels) but also complex and diffuse. Impact evaluations of for example sector strategies or general budget support potentially encompass multiple causal pathways resulting in long-term direct and indirect impacts. Some of the causal pathways linking interventions to impacts might be 'fairly' straightforward<sup>6</sup> (e.g. from training programs in alternative income generating activities to employment and to income levels), whereas other pathways are more complex and diffuse in terms of going through more intermediate changes, and being contingent upon more external variables (e.g. from stakeholder dialogue to changes in policy priorities to changes in policy implementation to changes in human welfare).

Given this diversity we think it is useful for purposes of 'scoping' to distinguish between two principal levels of impact: *impact at the institutional level* and *impact at the beneficiary level*<sup>7</sup>. It broadens impact evaluation beyond either simply measuring whether objectives have been achieved or assessing direct effects on intended beneficiaries. It includes the full range of impacts at all levels of the results chain, including ripple effects on families, households and communities, on institutional, technical or social systems, and on the environment. In terms of a simple logic model, there can be multiple intermediate (short and medium term) outcomes over time that eventually lead to impact – some or all of which may be included in an evaluation of impact at a specific moment in time.

<sup>&</sup>lt;sup>6</sup> Though not necessarily easy to measure.

<sup>&</sup>lt;sup>7</sup> Please note that the two levels should not be regarded as a dichotomy. In fact, a particular intervention might induce 'a cascade' of processes of change at different institutional levels (e.g. national, provincial government, cooperatives) before finally affecting the welfare of individuals.

Interventions that can be labeled as *institutional* primarily aim at changing secondorder conditions (i.e. the capacities, willingness, and organizational structures enabling institutions to design, manage and implement better policies for communities, households and individuals). Examples are policy dialogues, policy networks, training programs, institutional reforms, and strategic support to institutional actors (i.e. governmental, civil society institutions, private corporations, hybrids) and public private partnerships.

Other types of interventions directly aim at/affect communities, households, individuals, including voters and taxpayers. Examples are fiscal reforms, trade liberalization measures, technical assistance programs, cash transfer programs, construction of schools, etc.

#### BOX 1.1. 'Unpacking' the aid chain

The importance of distinguishing between different levels of impact is also discussed by Bourguignon and Sundberg (2007) who 'unpack' the aid effectiveness box by differentiating between three essential links between aid and final policy outcomes:

- Policies to outcomes; how do policies, programs and projects affect investment, production, growth, social welfare and poverty levels? (beneficiary level impact);
- Policy makers to policies; how does the policymaking process at national and local levels lead to 'good policies'? This is about governance (institutional capacities, checks and balances mechanisms, etc.) and is likely to be affected by donor policies and aid. (institutional level impact);
- External donors and international financial institutions to policy makers; how do external institutions influence the policymaking process through financial resources, dialogue, technical assistance, conditionalities, etc. (institutional level impact).

The above links can be perceived as channels through which aid eventually affects beneficiary level impact. At the same time, the processes triggered by aid generate lasting impacts at institutional levels.

Source: Bourguignon and Sundberg (2007)

Figure 1.1. graphically presents different levels of intervention and levels of impact. The differentiation between impact at institutional level and impact at beneficiary level<sup>8</sup> can be useful in the discussion on scope and method choice in impact evaluation.

<sup>&</sup>lt;sup>8</sup> A *third and fourth level of impact*, more difficult to pinpoint, respectively refer to the replicatory impact and the wider systemic effects of interventions. Both replicatory and systemic effects can result from processes of change at institutional or beneficiary levels. With respect to the first, evaluations that cover replicatory effects are quite scarce. This is in dire contrast with the manifest presence of replication (and the related concept of scaling-up) as explicit objectives in many policy interventions. For further discussion on replication see for example GEF (2007). These dimensions can be addressed in a theory-based impact evaluation framework (see section 3).



Figure 1.1. Levels of intervention, programs and policies and types of impact

Having illustrated this differentiation, it is important to note that for many in the development community, impact assessment is essentially about impact at beneficiary level. The *main* concern is how (sets of) policy interventions directly or indirectly affect the welfare of beneficiaries and to what extent changes in welfare can be attributed to these interventions. In line with this interpretation of impact evaluation<sup>9</sup>, throughout this Guidance document we will focus on impact assessment at the beneficiary level (see the red oval in Figure 1.1.), addressing key methodological concerns, and methodological approaches as well as the choice of methodological approach in a particular evaluation context. Where necessary other levels and settings of impact will be addressed (see the *blue oval* in Figure 1.1.). The implication is that with respect to the impact evaluation of , for example, New Aid Modalities (e.g. general budget support, sector budget support), this will only be discussed as far as interventions financed through these modalities (aim to) affect the lives of households and individuals<sup>10</sup>. We do not address the

<sup>&</sup>lt;sup>9</sup> This is the interpretation that has received most attention in methodological guidelines of international organizations working on impact evaluation such as for example the World Bank or the Asian Development Bank.

<sup>&</sup>lt;sup>10</sup> In this context one can distinguish between the effect of aid modalities on 'the way business is being done' (additionality of funding, direction of funding, public sector performance, coherence of policy changes, quality of intervention design, etc.; see for example Lawson et al. 2005), i.e.

question of how to do impact evaluations of New Aid modalities as such (see Lister and Carter, 2006; Elbers et al., 2008).

#### 1.5. Key Message

Identify the (scope and type of the) intervention. Interventions range from single ('strand') initiatives with explicit objectives to complex institutional policies. Look closely at the nature of the intervention, for example on the basis of the main theme addressed or by the generic policy instruments used. If interventions are complex in that they have multiple active components, state these separately and treat the intervention as a package of components that should be unpacked. Although complex interventions, sometimes of an institutional nature, are becoming more important and therefore should be subject to impact evaluation, this should not imply a reduction of interest in evaluating the impact of relatively simple, single strand interventions. The sheer number of these interventions makes doing robust IEs of great importance. In addition, one should be clear about the level of impact to be evaluated. Although most policy makers and stakeholders are primarily interested in beneficiary level impact (e.g. impact on poverty), specific policy interventions are primarily geared at inducing sustainable changes at (government) institutional level ('second order'-effects) with only indirect effects at beneficiary level.

what we call institutional level impact, and subsequently the impact of interventions funded (*in part*) by General Budget Support, Sector Budget Support or Debt Relief funds at beneficiary level. In the latter case, we are talking about impact evaluation as it is understood in most of the literature.

# 2. Agree on the objectives of the intervention that are valued

IE requires finding a balance between taking into account the values of stakeholders and paying appropriate attention to the empirical complexity of processes of change induced by an intervention. Some of this complexity has been unpacked in the discussion on the topic of scope of the IE, where we distinguished between levels of impact that neatly capture the often complex and diffuse causal pathways from intervention to different outcomes and impact: institutional level, beneficiary level and replicatory impact. It is recommended to try as much as possible to translate objectives into measurable indicators, but at the same time not losing track of important aspects that are difficult to measure.

After addressing the issue of stakeholder value we briefly discuss three dimensions that are particularly important and at the same time challenging to capture in terms of measurable indicators: intended versus unintended effects, short-term versus long-term effects, and the sustainability of effects.

#### 2.1. Stakeholder values in IE

Impact evaluation needs to assess the value of the results derived from an intervention. This is not only an empirical question but inherently a question about values – which impacts are judged as significant (whether positive or negative), what types of processes are valued in themselves (either positive or negative), and what and whose values are used to judge the distribution of costs and benefits of interventions.

First of all, stakeholder values are reflected in the objectives of an intervention as stated in the official documents produced by an intervention. However, interventions evolve and objectives might change. In addition, stakeholders groups, besides funding and implementing agencies, might harbor expectations not adequately covered by official documents. Impact evaluations need to answer questions related to 'for whom' the impacts have been intended, and how context influences impacts of interest. One of the tasks of an impact evaluation is therefore to be clear about who decides what the right aims are and to ensure that the legitimate different perspectives of different stakeholders are given adequate weight. Where there are multiple aims, there needs to be agreement about the standards of performance required in the weighting of these – for example, can an intervention be considered a success overall if it fails to meet some of the targets but does well in terms of the main intended outcome?

Depending on the evaluation context, there are different ways available to evaluators to address stakeholder values:

informal consultation with representatives from different stakeholder groups;

- using values inquiry<sup>11</sup> (Henry, 2002) as a basis for more systematic stakeholder consultation;
- using a participatory evaluation approach to include stakeholder values in the evaluation (see for example Cousins and Whitmore, 1998).

#### 2.2. Intended versus unintended effects

In development programs and projects intended effects are often translated into measurable indicators as early as the design phase. IE should go beyond assessing the expected effects given an intervention's logical framework and objectives. Interventions often change over time with consequences for how they affect institutional and people's realities. Moreover, effects are mostly context-specific, where different contexts trigger particular processes of change. Finally, in most cases the full scope of an intervention's effects is not known in advance. A well-articulated intervention theory can help to anticipate some of the unintended effects (see below, section 3).

Classic impact evaluations assume that there are no impacts for non-participants, but this is unlikely to be true for most development interventions. Spillover effects or replicatory effects (see section 1) can stem from market responses (given that participants and non-participants trade in the same markets), the (non-market) behavior of participants/non-participants or the behavior of intervening agents (governmental/NGO). For example, aid projects often target local areas, assuming that the local government will not respond; yet if one village gets the project, the local government may well cut its spending on that village, and move to the control village (Ravallion, 2008).

#### 2.3. Short- term versus long-term effects

In some types of interventions, impacts emerge quickly. In others impact may take much longer, and change over time. The timing of the evaluation is therefore important. Development interventions are usually assumed to contribute to longterm development (with the exception of humanitarian disaster and emergency situations). However, focusing on short-term or intermediate outcomes often provides for more useful and immediate information for policy- and decisionmaking. Intermediate outcomes may be misleading, often differing markedly from those achieved in the longer term. Many of the impacts of interest from development interventions will only be evident in the longer-term, such as environmental changes, or changes in social impacts on subsequent generations. Searching for evidence of such impacts too early might mistakenly conclude that they have failed.

In this context, the *exposure time* of an intervention to be able to make an impact is an important point. A typical agricultural innovation project that tries to change

<sup>&</sup>lt;sup>11</sup> "Values inquiry refers to a variety of methods that can be applied to the systematic assessment of the value positions surrounding the existence, activities, and outcomes of a social policy and program" (Mark et al., 1999: 183).

farmers' behavior with incentives (training, technical assistance, credit) is faced with time lags in both the adoption effect (farmers typically are risk averse and face resource constraints and start adopting innovations on an experimental scale) as well as the diffusion effect (other farmers want to see evidence of results before they copy). In such gradual non-linear processes of change with cascading effects, the timing of the ex post measurement (of land use) is crucial. Ex post measurements just after project closure could either underestimate (full adoption/diffusion of interesting practices has not taken place yet) or overestimate impact (as farmers will stop investing in those land use practices that are not attractive enough to be maintained without project incentives).

#### 2.4. The sustainability of effects

Focusing on short or intermediate term outcomes may underestimate the importance of designs that are able to measure effects (positive or negative) in the long term. An example is that an effective strategy to reduce child malnutrition in a certain population may quite quickly produce impressive results, yet fail soon after in the absence of systems, resources and capacities to maintain the work, or follow-up work, after termination of the intervention.

Few impact evaluations will probably provide direct evidence of long-term impacts, and in any case results are needed before these impacts become evident to inform decisions on continuation, next phases and scaling-up. Impact evaluations therefore need to identify short-term impacts and, where possible, indicate whether longer-term impacts are likely to occur.

In order to detect negative impacts in the long term, early warning indicators can be important to include. A well-articulated intervention theory (see section 3), that also addresses the time horizons over which different types of outcomes and impacts could reasonably be expected to occur, can help to identify impacts which can and should be explored in an evaluation. The sustainability of positive impacts is also likely to be only evident in the longer term. Impact evaluations therefore can focus on other impacts that will be observable in the short term, such as the institutionalization of practices and the development of organizational capacity, that are likely to contribute to the sustainability of impacts for participants and communities in the longer term<sup>12</sup>.

#### 2.5. Key Message

Agree on the objectives of the intervention that are valued. Select objectives that are important. Do not be afraid of selecting one objective; focus and clarity is a virtue not a vice. As much as possible try to translate objectives into measurable indicators, but at the same time do not lose track of important aspects that are difficult to measure. In addition, keep in mind the dimensions of exposure time and the sustainability of changes.

<sup>&</sup>lt;sup>12</sup> For a discussion on different dimensions of sustainability in development intervention, see Mog (2004).

# 3. Carefully articulate the theories linking interventions to outcomes

#### 3.1. Seeing interventions as theories: the black box and the contribution problem

When evaluators talk about the black box 'problem', they are usually referring to the practice of viewing interventions primarily in terms of effects, with little attention paid to how and why those effects are produced. The common thread underlying the various versions of theory-based evaluation is the argument that 'interventions are theories incarnate' and evaluation constitutes a test of intervention theory or theories. Interventions are embodiments of theories in at least two ways. First, they comprise an expectation that the introduction of a program or policy intervention will help ameliorate a recurring social problem. Second, they involve an assumption or set of assumptions about how and why program activities and resources will bring about changes for the better. The underlying theory of a program is often not surfaced and remains hidden; typically in the minds of policy architects and staff. Policies, be it relatively small-scale direct interventions like information campaigns, training programs or subsidization, meso-level interventions like public private partnerships and social funds, or macro-level interventions such as 'General Budget Support' rest on social, behavioral and institutional assumptions indicating why 'this' policy intervention will work, which at first view are difficult to uncover.

By seeing interventions as theories, and by using insights from theory-based evaluations, it is possible to open up the *black box*. Development policies and interventions, in one way or another have to do with changing behavior/ intentions/knowledge of households, individuals and organizations (grass roots, private, and public sector). Crucial for understanding what can change behavior is information on *behavioral and social mechanisms*. An important insight from theory-based evaluations is that policy interventions are (often) believed to address and trigger certain social and behavioral responses among people and organizations while in reality this may not necessarily be the case.

#### 3.2. Articulating intervention theories on impact

Program theory (or intervention theory) can be identified (articulated) and expressed in many ways – a graphic display of boxes and arrows, a table, a narrative description and so on. The methodology for constructing intervention theory, as well as the level of detail and complexity, also varies significantly (e.g. Connell et al., 1995; Leeuw, 2003; Lipsey, 1993; McClintock, 1990; Rogers et al., 2000; Trochim, 1989; Wholey, 1987).

Too often the role of methodology is neglected, assuming that 'intervention theories' are like *manna* falling out of the sky. That is not the case. Often the underlying theory has to be digged up. Moreover, much of what passes as theory-

based evaluation today is simply a form of "analytic evaluation [which] involves no theory in anything like a proper use of that term" (Scriven, 1998: 59).

The intervention theory provides an overall framework for making sense of potential processes of change induced by an intervention. Several pieces of evidence can be used for articulating the intervention theory, for example:

- an intervention's existing logical framework provides a useful starting point for mapping causal assumptions linked to objectives; other written documents produced within the framework of an intervention are also useful in this respect;
- insights provided by as well as expectations harbored by policy makers and staff (and other stakeholders) on how they think the intervention will affect/is affecting/has affected target groups;
- (written) evidence on past experiences of similar interventions (including those implemented by other organizations);
- research literature on mechanisms and processes of change in certain institutional contexts, for particular social problems, in specific sectors, etc.

#### Box 3.1. Social funds and government capacity: competing theories

Proponents of social funds argue they will develop government capacity in several ways. Principle amongst these are that the social fund will develop superior means of resource allocation and monitoring which will be transferred to government either directly through collaborative work or indirectly through copying the procedures shown to be successful by the social fund. But critics argue that social funds bypass normal government channels and so undermine government capacity, an effect reinforced by drawing away government's best people by paying a project premium. Hence these are rather different theories of how social funds affect government capacity. Carvalho and White (2004) refer to both sets of assumptions in terms of 'theory' and 'anti-theory'. Their study found that well-functioning, decentralized social funds, such as ZAMSIF in Zambia, worked through, rather than parallel to, existing structures, and that the social fund procedures were indeed adopted more generally by district staff. But at national level there was generally little evidence of either positive or negative effects on capacity – with some exceptions such as the promotion of poverty mapping in some countries.

Source: Carvalho and White (2004)

An example of what an impact theory might look like is the following. Consider the case of a small business development project that provides training to young managers who have started up a business. The goal is to help making small businesses financially sustainable and indirectly to generate more employment in the region where the project is being implemented. Closer scrutiny reveals that the project might have a positive influence on the viability of small businesses in two ways. First by training young people in basic management and accounting skills the project intends to have a positive effect on financial viability and ultimately on the growth and sustainability of the business. Second, by supporting the writing of a business plan, the project aims to increase the number of successful applications for credit with the local bank, which previously excluded the project's target group due to the small loan sizes (high transaction costs) and high risks involved. Following this second causal strand, efficient and effective spending of the loan is also expected to contribute to the strength of the business. Outputs are measured in terms of the number of people trained by the project and the number of loans received from the bank (see Figure 3.1.).



Figure 3.1. Basic intervention theory of a fictitious small business support project

Any further empirical analysis of the impact of the project requires insight into the different factors besides the project itself that affect small business development and employment generation. Even in this rather simple example the number of external variables that affect the impact variables either directly or by moderating the causal relations specified in Figure 3.1. are manifold. We will restrict ourselves to some examples on the output-impact part of the chain:

- short-term demands on the labor efforts of business owners in other livelihood activities may lead to suboptimal strategic choices jeopardizing the sustainability of the business;
- inefficient or ineffective use of loans due to short term demands for cash for other expenditures might jeopardize repayment and financial viability of the business;
- deteriorating market conditions (in input or output markets) may jeopardize the future of the business;
- the availability and quality of infrastructure or skilled labor at any point may become constraining factors on business development prospects;
- the efforts of other institutions promoting small business development or any particular aspect of it might positively (or negatively) affect businesses;
- etc.

Methods for reconstructing the underlying assumptions of project/program/policy theories are the following (see Leeuw, 2003):

- a policy-scientific method, which focuses on interviews, documents and argumentation analysis;
- a strategic assessment method, which focuses on group dynamics and dialogue, and;
- An elicitation method, which focuses on cognitive and organizational psychology.

Central in all three approaches is the search for mechanisms that are believed to be 'at work' when a policy is implemented.

## Box 3.2. Social and behavioral mechanisms as heuristics for understanding processes of change and impact

Hedström (2005: 25) has defined the concept of "[social] mechanisms" as "a constellation of entities and activities that are organized such that they regularly bring about a particular type of outcome". Mechanisms form the 'nuts and bolts' (Elster, 1989) or the 'engines' (Leeuw, 2003) of interventions (policies and programs), making them work, given certain contexts (Pawson and Tilley, 1997). Hedström and Swedberg (1998: 296-98), building on the work of Coleman (1990), discuss three types of mechanisms: situational mechanisms, action formation mechanisms and transformational mechanisms.

Examples of situational mechanisms are:

- self-fulfilling and self-denying prophecies;

- crowding-out (e.g. by striving to force people who are already largely compliant with laws and regulations into full compliance, the opposite is realized, because due to the extra focus on laws and regulation the internal motivation of people to comply is reduced).

Action-formation mechanisms are the heuristics that people develop to deal with their bounded rationality such as:

- framing and the endowment effect — "the fact that people often demand much more to give up an object than they would be willing to pay to acquire it", but also the tendency for people to have a stronger preference for more immediate payoffs relative to later payoffs, the closer to the present both payoffs are;

- types of learning (social learning, vicarious learning);

- 'game-theoretical' mechanisms such as using the 'grim strategy' ( to repeatedly refuse to cooperate with another party as a punishment for the other party's failure to cooperate previously), and the shadow of the future /shadow of the past-mechanism;

- mechanisms like the 'fight-or-flight-response' to stress and the 'tend-and befriend-mechanism' are other examples.

*Transformational mechanisms* illuminate how processes and results of interacting individuals and groups are 'transformed' into collective outcomes. Examples are the following:

- Cascading is a process by which people influence one another, so much so that participants ignore their private knowledge and rely instead on the publicly stated judgments of others. The bandwagon phenomenon (the tendency to do (or believe) things because many other people do (or believe) the same) is related to this as are group think, the common knowledge effect and herd behavior;

- 'Tipping points' "where a small additional effort can have a disproportionately large effect, can be created through virtuous circles, or be a result of achieving certain critical levels" (Rogers, 2008: 35).

#### The relevance of mechanisms for impact evaluations

Development policies and interventions, in one way or another have to do with changing behavior/ intentions/knowledge of households, individuals and organizations (grass roots, private, and public sector). Crucial for understanding what can change behavior is information about these mechanisms. *The mechanisms underlying processes of change might not be necessarily those that are assumed to be at work by policy makers, programs designers and staff.* Creating awareness on the basis of (public) information campaigns is not always leading to behavioral change. Subsidies and other financial incentives run the risk of causing unintended side effects such as benefit snatching but also create the 'Mitnahme-effect' (people already tended to behave in a way the incentive wanted them to behave before there was an incentive). Mentoring drop outs in education might cause 'learned helplessness' and therefore will increase drop out rates. Many other examples are available in the literature. The relevance of knowing which social and behavioral mechanisms are believed to do the work increases the more complicated and complex interventions are.

A focus on mechanisms helps evaluators and managers to open up and test the theory underlying an intervention. Spending time and money on programs based on 'pet theories' of policy makers or implementation agents that are not corroborated by relevant research, should probably not be high on the agenda. If a policy intervention is based on mechanisms that are known not to work (in a given context or more in general), then that is a signal that the intervention probably will not be very effective. This can be found out on the basis of desk research as a first test of the relevance and 'validity' of an intervention theory, i.e. by confronting the theory with existing knowledge about mechanisms. That knowledge stems from synthesis and review studies (see section 5). Further empirical IE can generate more contextualized and precise tests of the intervention theory.

#### 3.3. Testing intervention theories on impact

After articulating the assumptions on how an intervention is expected to affect outcomes and impacts, the question arises to what extent these assumptions are valid. In practice, evaluators have at their disposal a wide range of methods and techniques to test the intervention theory. We can broadly distinguish between two broad approaches. The first is that the theory constitutes the basis for constructing a 'causal story' about how and to what extent the intervention has produced results. Usually different methods and sources of evidence are used to further refine the theory in an iterative manner until a credible and reliable causal story has been generated. The second way is to use the theory as an explicit benchmark for testing (some of) the assumptions in a formal manner. Besides providing a benchmark, the theory provides the template for method choice, variable selection and other data collection and analysis issues. This approach is typically applied in statistical analysis, but is not in any way restricted to this type of method. In short, theory-based methodological designs can be situated anywhere in between 'telling the causal story' to 'formally testing causal assumptions'.

The systematic development and corroboration of the causal story can be achieved through *causal contribution analysis* (Mayne, 2001) which aims to demonstrate whether or not the evaluated intervention is one of the causes of observed change. Contribution analysis relies upon chains of logical arguments that are verified through a careful analysis. Rigor in causal contribution analysis involves systematically identifying and investigating alternative explanations for observed impacts. This includes being able to rule out implementation failure as an explanation of lack of results, and developing testable hypotheses and predictions to identify the conditions under which interventions contribute to specific impacts.

The causal story is inferred from the following evidence:

1. There is a reasoned theory of change for the intervention: it makes sense, it is plausible, and is agreed by key players.

- 2. The activities of the intervention were implemented.
- 3. The theory of change—or key elements thereof— is verified by evidence: the chain of expected results occurred.
- 4. Other influencing factors have been assessed and either shown not to have made a significant contribution or their relative role in contributing to the desired result has been recognized.

The analysis is best done iteratively, building up over time a more robust assessment of causal contribution. The overall aim is to reduce the uncertainty about the contribution the intervention is making to the observed results through an increased understanding of why the observed results have occurred (or not) and the roles played by the intervention and other factors. At the level of impact this is the most challenging, and a 'contribution story' has to be developed for each major strategy that is part of an intervention, at different levels of analysis. They would be linked, as each would treat the other strategies as influencing factors.

One of the key limitations in the foregoing analysis is to pinpoint the exact causal effect from intervention to impact. Despite the potential strength of the causal argumentation on the links between the intervention and impact, and despite the possible availability of data on indicators, as well as data on contributing factors (etc.), there remains uncertainty about the *magnitude* of the impact as well as *the extent* to which the changes in impact variables are really due to the intervention or due to other influential variables. This is called the attribution problem and is discussed in section 4.

#### 3.4. Key message

Carefully articulate the theories linking interventions to outcomes. What are the causal pathways linking intervention outputs to processes of change and impact? Be critical if an 'intervention theory' appears to assert or assume changes without much explanation. The focus should be on dissecting the causal (social, behavioral and institutional) mechanisms that make interventions 'work'.

## 4. Address the attribution problem

#### 4.1. The attribution problem

Multiple factors can affect the livelihoods of individuals or the capacities of institutions. For policy makers as well as for stakeholders it is important to know what the added value is of the policy intervention apart from these other factors. The attribution problem is often referred to as the central problem in impact evaluation. The central question is to what extent can changes in outcomes of interest be *attributed* to a particular intervention? Attribution refers both to isolating and measuring accurately the particular contribution of an intervention and ensuring that causality runs from the intervention to the outcome.

The changes in welfare for a particular group of people can be observed by doing "before and after" studies, but these only rarely accurately measure impact. Baseline data (before the intervention) and end-line data (after the intervention) give facts about the development over time and describe "the factual" for the treatment group (not the counterfactual). But changes observed by comparing before-after (or pre-post) data are rarely caused by the intervention alone since other interventions and processes influence developments, both in time and space. There are some exceptions in which before versus after will suffice to determine impact. For example, supplying village water pumps reduces time spent fetching water. If nothing else of importance happened during the period under study, attribution is so clear that there is no need to resort to anything other than before versus after to determine this impact.

In general the observed changes are only partly caused by the intervention of interest. Other interventions inside or outside the core area will often interact and strengthen/reduce the effects of the intervention of interest for the evaluation. In addition other unplanned events or general change processes will often influence development, be it for example nature catastrophes, urbanization, the growing economies of China and India, business cycles, war or long term climate change. For example, in evaluating the impact of microfinance on poverty, we have to control for the influences of changing market conditions, infrastructure developments or climate shocks such as droughts, and so on.

A discussion that often comes up in IE is the issue of *attribution of what*? This issue is complementary to the independent variable question discussed in chapter one. How the impact of the intervention is measured may be stated in several ways:

- 1. what is the impact of an additional dollar of funding to program  $x^{13}$ ?
- 2. what is the impact of country y's contribution to a particular intervention?
- 3. what is the impact of intervention z?

<sup>&</sup>lt;sup>13</sup> Economists employ several useful techniques for estimating the marginal impact of an extra dollar invested in a particular policy intervention. See for example Appendix 1, second example. We consider these methods to be complementary to IE and beyond the scope of this Guidance.

In this guidance we will focus on the third level of attribution: what is the impact of a particular policy intervention (from very simple to complex) independent of the specific monetary and non-monetary contributions of the (institutional) actors involved? The issue of attributing impact to a particular intervention can be a quite complicated issue in itself (especially when talking about complex interventions such as sector strategies or programs). Additional levels of attribution such as tracing impact back from interventions to specific (financial) contributions of different donors are either meaningless or too complicated to achieve in a pragmatic and cost-effective manner.

Proper analysis of the attribution problem is to compare the situation 'with' an intervention to what would have happened in the absence of an intervention, the 'without' situation (the *counterfactual*). Such comparison of the situation "with *and without*" the intervention is challenging since it is not possible to observe how the situation would have been without the intervention, and has to be constructed by the evaluator. The counterfactual is illustrated in Figure 4.1. The value of a target variable (point a) after an intervention should not be regarded as the intervention's impact, nor is it simply the difference between the before and after situation (a-b, measured on the vertical axis). The net impact (at a given point in time) is the difference between the target variable's value after the intervention and the value the variable would have had in case the intervention would not have taken place (a-c).

#### Figure 4.1. Graphical display of the net impact of an intervention



value target variable

The starting point for an evaluation is a good account of the factual – what happened in terms of the outputs/outcomes targeted by the intervention? A good account of the factual requires articulating the intervention theory (or theories) connecting the different causal assumptions from intervention outputs to outcomes and impacts as discussed earlier in section 3. As to the counterfactual, in this guidance we will discuss several options for measuring the counterfactual. Evaluations can either be experimental as when the evaluator purposely collects data and designs evaluations in advance or guasi-experimental as when data are collected to mimic an experimental situation. Multiple regression analysis is the all-purpose technique that can be used in virtually all settings; when the experiment is organized in such a way that no controls are needed, a simple comparison of means can be used instead of a regression since it will give the same answer. Experimental and guasi-experimental approaches will be discussed in section 4.2. We briefly introduce the general principles and the most common approaches. The idea of (quasi-) experimental counterfactual analysis is that the situation of a participant group (receiving benefits from/affected by an intervention) is compared over time with the situation of an equivalent comparison group that is not affected by the intervention. Several designs exist of combinations of ex ante and ex post measurements of participant and control group (see section 4.2.). Randomization of intervention participation is considered to be the best way to create equivalent groups. Random assignment to the participant and control group guarantees that the two groups will have similar average characteristics both for observables and non-observables except for the intervention. As a second-best alternative several matching techniques (e.g. propensity score matching) can be used to create control groups that are as similar to participant groups as possible (see below).

#### 4.2. Methodological approaches addressing the attribution problem

Three related problems that quantitative impact evaluation techniques attempt to address are the following:

- the establishment of a *counterfactual*: What would have happened in the absence of the intervention(s);
- the elimination of *selection effects*, leading to differences between intervention group (or treatment group) and control group;
- a solution for the problem of *unobservables*: the omission of one or more unobserved variables, leading to biased estimates.

Selection effects occur for example when those in the intervention group are more motivated or less motivated than those in the control group. It is particularly a problem when the variable in question, in this case motivation, is not easily observable. As long as selection is based on *observable* characteristics, and these are measured in the evaluation, they may be included, and thus controlled for, in the regression analysis. However, not all relevant characteristics are observed or measured. This problem of *selection on unobservables* is one of the main problems in impact evaluation.

In the following sections we will discuss different techniques of quantitative impact evaluation, thereby mainly focusing our discussion on the selection bias issue. (Quasi-) experimental design-based approaches such as the randomized controlled trial (RCT) or the pipeline approach, in trying to deal systematically with selection effects can be compromised by two sets of problems: contamination and unintended behavioral responses. We briefly discuss the two.

#### Contamination

Contamination (or contagion, treatment diffusion) refers to the problem of groups of people that are not supposed to be exposed to certain project benefits are in fact benefiting from these in one or more ways. Contamination comes from two possible sources. The first is contamination from the intervention itself as a result of spill-over effects. Interventions are most often planned and implemented within a delimited space (a village, district, nation, region or institution). The influence zone of an intervention may, however, be larger than the core area where the intervention takes place or is intended to generate results (geographical spill-over effects). To avoid contamination, control and comparison groups must be located outside the influence zone. Second, the selected comparison group may be subject to similar interventions implemented by different agencies, or even somewhat dissimilar interventions but which affect the same outcomes. The counterfactual is thus a different type of intervention rather than no intervention. This problem is often overlooked. A good intervention theory as a basis for designing a good measurement instrument that records the different potential problems of contamination is a good way to address this problem.

#### Unintended behavioral responses

Several unintended behavioral responses not caused by the intervention or 'normal' conditions might disrupt the validity of comparisons between groups and hence the ability to attribute changes to project incentives. Important possible effects are the following (see Shadish et al., 2002, Rossi et al., 2004)):

- Expected behavior or compliance behavior: participants react in accordance with intervention staff expectations for reasons such as compliance with the established contract, or due to certain expectations about future benefits from the organization (not necessarily the project).
- Compensatory equalization: discontent among staff or recipients with inequality between incentives might result in compensation of groups that receive less than other groups.
- Compensatory rivalry: differentiation of incentives to groups of people might result in social competition between those receiving (many) intervention benefits and those that receive less or no benefits.
- Hawthorne effect: the fact of being part of an experiment rather than the intervention as such causing people to change their behavior.
- Placebo effect: the behavioral effect is not the result of the incentives provided by the intervention but people's perception of the incentives and the subsequent anticipatory behavior.
- Other effects (see Shadish et al., 2002).

These problems are relevant in most experimental and quasi-experimental design approaches that are based on ex ante participant and control/comparison group

designs<sup>14</sup>. They are less relevant in regression-based approaches that use statistical matching procedures or do not rely on the participant-control group comparison for counterfactual analysis<sup>15</sup>.

#### 4.2.1. Randomized Controlled Trial<sup>16</sup>

The safest way to avoid selection effects is a *randomized selection* of intervention group and control *before* the experiment starts. When the experimental group and the control group are selected randomly from the same eligible population, both groups will have similar average characteristics except that the experiment group received the intervention. That is why a simple comparison of average outcomes in the two groups solves the attribution problem and yields accurate estimates of the impact of the intervention: by design, the only difference between the two groups was the intervention.

To determine if the intervention had a statistically significant impact, one simply performs a test of equality between the mean outcomes in the experiment and control group. Statistical analysis will tell you if the impact is significant and how large it is. Of course, with larger samples, the statistical inferences will be increasingly precise; but if the impact of an intervention really is large, it can be detected and measured even with a relatively small sample.

A proper Randomized Controlled Trial (RCT) nicely solves many attribution issues, but has to be managed carefully to avoid contamination. Risks of a RCT are a) different rates of attrition in the two groups, for instance caused by a high dropout in one of the two groups, b) spillover effects (contamination) resulting in the control group receiving some of the treatment, and c) unintended behavioral responses.

#### 4.2.2. Pipeline approach

One of the problems for the evaluation of development projects or programs is that evaluators rarely get involved early enough to design a good evaluation (although this is changing). Often, households or individuals are selected for a specific project, while not everybody participates (directly) in the project. A reason may be, for instance, a gradual implementation of the project. Large projects (such as in housing or construction of schools) normally have a phased implementation.

In such a case, it may be possible to exploit this phasing of the project by comparing the outcomes of households or communities who actually participate

<sup>&</sup>lt;sup>14</sup> It is difficult to identify general guidelines for avoiding these problems. Evaluators have to be aware of the possibility of these effects affecting the validity of the design. For other problems as well as solutions see Shadish et al. (2002).

<sup>&</sup>lt;sup>15</sup> For further discussion on the approaches discussed below see Appendices 3 to 6.

<sup>&</sup>lt;sup>16</sup> We like to thank Antonie de Kemp of IOB for insightful suggestions.

(the experiment group) with households or communities who are selected, but do not participate yet (the comparison group). A specific project (school building) may start for instance in a number of villages and be implemented later on in other villages. This creates the possibility to evaluate the effect of school building on enrolment. One has to be certain, of course, that the second selection – the actual inclusion in the project – does not introduce a selection bias. If for instance, at the start of the project a choice is made to start construction in a number of specific villages, the (relevant) characteristics of these villages must be similar to other villages that are eligible for new schools. Self-selection (due to villages that are eager to participate) or other selection criteria (starting in remote areas or in urban areas) may introduce a selection bias.

#### 4.2.3. Propensity score matching

When no comparison group has been created at the start of the project or program, a comparison group may be created ex post through a matching procedure: for every member of the treatment group one or more members in a control group are selected on the basis of similar observed (and relevant) characteristics. Suppose we have two groups, a relatively small intervention group, consisting for instance of 100 pupils who will receive a specific reading program. If we want to analyze the effects of this program, we must compare the results of the pupils in the program with other pupils who were not included in the program. We cannot select just any control group, because the intervention group may be selected on the basis of specific characteristics (pupils with relatively good results or relatively bad results, pupils from rural areas, from private schools or public schools, boys, girls, orphans, etc.). Therefore, we need to select a group with similar characteristics. One way of doing this would be to find for every boy, aged 10 years, from a small rural school with a high pupil teacher ratio in a poor district, another boy with the same observed characteristics. This would be an enormously time consuming procedure, especially when you have to do this for a hundred pupils.

An alternative way to create such a control group is the method of *propensity score matching*. This technique involves forming pairs, not by matching every characteristic exactly, but by selecting groups that have similar *probabilities* of being included in the sample as the treatment group. The technique uses all *available* information in order to construct a control group (see box 4.1.)<sup>17</sup>. In 1983, Rosenbaum and Rubin (1983) showed that this method made it possible to create a control group ex post with characteristics that are similar to the kind of intervention and control groups that would have been created had they been selected randomly before the beginning of the project.

Box 4.1. Using propensity scores to select a matched comparison group - the Viet Nam Rural Roads Project

The survey sample included 100 project communes and 100 non-project communes in the

<sup>&</sup>lt;sup>17</sup> For an explanation, see Wooldridge (2002), chapter 18.

same districts. Using the same districts simplified survey logistics and reduced costs, but communes were still far enough apart to avoid "contamination" (control areas being affected by the project). A logit model of the probability of participating in the project was used to calculate the propensity score for each project and non-project commune. Comparison *communes* were then selected with *propensity scores* similar to the project communes. The evaluation was also able to draw on commune-level data collected for administrative purposes that cover infrastructure, employment, education, health care, agriculture and community organization. These data will be used for contextual analysis and to construct commune-level indicators of welfare and to test program impacts over time. The administrative data will also be used to model the process of project selection and to assess whether there are any selection biases.

Source: Van De Walle and Cratty, 2005 (literal citation: Bamberger, 2006)

It must be noted that the technique only deals with selection bias on observables and does not solve potential endogeneity bias (see Appendix 4) that results from the omission of unobserved variables. Nevertheless propensity score matching may be combined with the technique of double differencing in order to correct for the influence of time invariant unobservables (see below). Moreover, the technique may require a large sample for the selection of the comparison group and this may be an issue when the researcher cannot rely on existing secondary data for this procedure.

#### 4.2.4. Judgmental matching<sup>18</sup>

A less precise method for selecting control groups uses descriptive information from (e.g.) survey data to construct comparison groups.

*Matching areas on observables.* The researcher, in consultation with clients and other knowledgeable persons, identifies characteristics on which comparison and project areas should be matched (e.g., access to services, type or quality of house construction, economic level, central location or remoteness, types of agricultural production). The researcher then combines information from maps (and sometimes Geographic Information System (GIS) data and aerial photographs), observation, secondary data (censuses, household surveys, school records, etc) and key informants to select comparison areas with the best match of characteristics. When operating under real-world constraints it will often be necessary to rely on easily observable or identifiable characteristics such as types of housing and infrastructure. While this may expedite matters, it is important to keep in mind the potential for unobservable differences, to address these as far as possible through qualitative research, and to attach the appropriate caveats to the results.

Matching individuals or households on observables. Similar procedures are used to match individuals and households. Sample selection can sometimes draw on previous or ongoing household surveys, but in many cases researchers must

<sup>&</sup>lt;sup>18</sup> This sub section is largely a literal citation from Bamberger (2006).

develop their own ways to select the sample. Sometimes the selection is based on observable physical characteristics (type of housing, distance from water and other services, type of crops or area cultivated) while in other cases selection is based on characteristics that require screening interviews, such as economic status, labor market activity, or school attendance. In these latter cases the interviewer is given quotas of subjects with different characteristics to be located and interviewed (quota sampling).

4.2.5. Double difference (difference in difference)

Differences between intervention group and control group may be unobserved and therefore problematic. Nevertherless, even though such differences cannot be measured, the technique of double difference (or difference in difference) deals with these differences as long as they are time invariant. The technique measures differences between the two groups, before and after the intervention (hence the name: double difference):

|                        | Intervention Group             | Control Group                  | Difference across groups       |
|------------------------|--------------------------------|--------------------------------|--------------------------------|
| Baseline               | Ι <sub>ο</sub>                 | Co                             | I <sub>o</sub> -C <sub>o</sub> |
| Follow-up              | l <sub>1</sub>                 | C <sub>1</sub>                 | I <sub>1</sub> -C <sub>1</sub> |
| Difference across time | l <sub>1</sub> -l <sub>o</sub> | C <sub>1</sub> -C <sub>0</sub> | Double-difference:             |
|                        |                                |                                | $(I_1 - C_1) - (I_0 - C_0) =$  |
|                        |                                |                                | $(I_1 - I_0) - (C_1 - C_0)$    |

Table 4.1. Double difference and other designs

Source: adapated from: Maluccio and Flores (2005).

Suppose there are two groups, an intervention group I and a control group C. One measures for instance enrolment rates before (0) and after (1) the intervention. According to this method, the effect is

 $(I_1 - I_0) - (C_1 - C_0)$  or  $(I_1 - C_1) - (I_0 - C_0)$ 

For example, if enrolment rates at t=0 would be 80% (for the intervention group) and 75% for the control group and at t=1 these rates would be respectively 90% and 75%, then the effect of the intervention would be: (90% - 80%) - (75% - 70%) = 5%.

The techniques of propensity score matching (see above) and double difference may be combined. Propensity score matching increases the likelihood that the treatment group and control group have similar characteristics, but cannot guarantee that all relevant characteristics are included in the selection procedure. The double difference technique can completely eliminate the effects of an unobserved selection bias but this technique may work better when differences between intervention group and control group are eliminated as much as possible. The approach eliminates *initial* differences between the two groups (for instance differences in enrolment rates) and therefore gives an unbiased estimate of the effects of the intervention, as long as these differences are time invariant. When an unobserved variable is time variant (changes over time), the measured effect will still be biased.

#### 4.2.6. Regression analysis and double difference

In some programs, the interventions are all or nothing (a household or individual is subjected to the intervention or not); in others they vary continuously over a range as when programs vary the type of benefit offered to target groups. Take for example a cash transfer program or a microfinance facility where the amount transferred or lent may depend on the income of the participant. Improved drinking water facilities are another example. These facilities differ in capacity and are implemented in different circumstances with beneficiaries living at different distances to these facilities.

In addition to the need to deal with both discrete and continuous interventions, we also need to control for other factors that affect the outcome other than the magnitude of the intervention. The standard methodology for such an approach is a regression analysis. One of the reasons for the popularity of regression-based approaches is their flexibility: they may deal with the heterogeneity of treatment, multiple interventions, heterogeneity of characteristics of participants, interactions between interventions and interactions between interventions and specific characteristics, as long as the treatment (or intervention) and the characteristics of the subjects in the sample are observed (can be measured). With a regression approach, it may be possible to estimate the contribution of a specific interventions. The analysis may include an explicit control group.

We must go beyond a standard regression-based approach when there are unobserved selection effects or endogeneity (see next section). A way to deal with unobserved selection effects is the application of the 'difference in difference' approach in a regression model (see Appendix 4). In such a model we do not analyze the (cross section) effects between groups, but the changes (within groups) over time. Instead of taking the specific values of a variable in a specific year, we analyze the *changes* in these variables over time. In such an analysis, unobserved time invariant variables drop from the equation. The approach is similar to a *fixed effects regression* model that uses deviations from individual means in order to deal with (unobserved) selection effects.

Again, the quality of this method as a solution depends on the validity of the assumption that unobservables are time invariant. Moreover, the quality of the method also depends on the quality of the underlying data. The method of first differencing is more vulnerable than some other methods to the presence of measurement error in the data,

4.2.7. Instrumental variables

An important problem when analyzing the impact of an intervention is the problem of *endogeneity*. The most common example of endogeneity is when a third variable causes two other variables to correlate without there being any causality. For example, Doctors are observed to be frequently in the presence of people with fevers, but Doctors do not cause the fevers, it's the third variable (the illness) that causes the two other variables to correlate (people with fevers and the presence of Doctors). In econometric language, when there is endogeneity an explanatory variable will be correlated with the error term in a mathematical model (see Appendix 4). When an explanatory variable is endogenous, it is not possible to give an unbiased estimate of the causal effect of this variable.

Selection effects also give rise to bias. Consider the following example. Various studies in the field of education find that repeaters produce lower test and examination results than pupils who have not repeated class levels. A preliminary and false conclusion would be that repetition does not have a positive effect on student performance and that it is simply a waste of resources. But such a conclusion neglects the endogeneity of repetition: intelligent children with well-educated parents are more likely to perform well and therefore do not repeat. Less intelligent children, on the other hand, will probably not achieve good results and are therefore more likely to repeat. So, both groups of pupils (i.e. repeaters and non-repeaters) have different characteristics, which at first view makes it impossible to draw conclusions based on a comparison between them.

The technique of instrumental variables is used to address the endogeneity problem. An instrumental variable (or instrument) is a third variable that is used to get an unbiased estimate of the effect of the original endogenous variable (see Appendix 4). A good instrument correlates with the original endogenous variable in the equation, but not with the error term. Suppose a researcher is interested in the effect of a training. Actual participation may be endogenous, because (for instance) the most motivated employees may subscribe to the training. Therefore, one cannot compare employees who had the training with employees who had not without incurring bias. The effect of the training may be determined if a subset were assigned to the training by accident or through some process unrelated to their own personal motivation. In this case, the instrumental variables procedure essentially only uses data from that subset to estimate the impact of training.

#### 4.2.8. Regression Discontinuity Analysis

The basic idea of regression discontinuity analysis is simple. Suppose program participation depends on income. On the left side of the *cut off point* people (or households) have an income that is just low enough to be eligible for participation; on the right side of the cut off point people are no longer allowed to participate, even though their income is just slightly higher. There may be more criteria that define the threshold and these criteria may be either explicit or implicit. Regression discontinuity analysis compares the treatment group with the

control group at the cut off point. At that point, it is unlikely that there are unobserved differences between the two groups.



Figure 4.2. Regression discontinuity analysis

Suppose we want to analyze the effect of a specific program to improve the learning achievements of pupils. This program focuses on the poorest households: the program includes only households with an income below a certain maximum. We know that learning achievements are correlated with income and therefore we cannot compare households participating in the program with households who do not participate. Second, other factors may induce an endogeneity bias (such as differences in the educational background of parents or the distance to the school). Nevertheless, at the cut off point, there is no reason to assume that there are systematic differences between the two groups of households (apart from small differences in income). Estimating the impact can now be done for example by comparing the mean difference between the regression line of learning achievements in function of income *before* the intervention with the regression line *after* (see Figure 4.2.).

A major disadvantage of a regression discontinuity design is that the method assesses the marginal impact of the program only around the cut-off point for eligibility. Moreover, it must be possible to construct a specific threshold and individuals should not be able to manipulate the selection process (ADB, 2006: 14). Many researchers prefer regression discontinuity analysis above propensity score matching, because the technique generates a higher likelihood that estimates will not be biased by unobserved variables<sup>19</sup>.

#### 4.3. Applicability of quantitative methods for addressing the attribution problem

<sup>&</sup>lt;sup>19</sup> With instrumental variables one may try to get rid of an expected bias, but the technique cannot guarantee that endogeneity problems will be solved completely (the instrumental variable may also be endogenous). Moreover, with weak instruments the precision of the estimate may be low.
There are some limitations to the applicability of these techniques. We briefly highlight some of the more important ones (for a more comprehensive discussion see for example Bamberger and White, 2007). First, in general counterfactual estimation is not applicable in full coverage interventions such as price policies or regulation on land use, which affect everybody (although to different degrees). In this case there are still possibilities to use statistical 'counterfactual-like' analyses such as those that focus on the variability in exposure/participation in relation to changes in an outcome variable (see for example Rossi et al., 2004). Second, there are several pragmatic constraints to apply this type of analysis especially with respect to randomization and other design-based techniques. For example there might be ethical objections to randomization or lack of data representing the baseline situation of intervention target groups.

An important critique on the applicability of the abovementioned methods refers to the nature of the intervention and the complexity of the context in which the intervention is embedded. The methodological difficulties of evaluating complex and complicated interventions to some extent can be 'neutralized' by deconstructing complex interventions into their 'active ingredients' (see for example Vaessen and Todd, 2008)<sup>20</sup>. Consider the example of School Reform in Kenya described by Duflo and Kremer (2005). School Reform constitutes a set of different simultaneous interventions at different levels ranging from revisions in and decentralization of the budget allocation process to addressing links between teacher pay and performance to vouchers and school choice. While the total package of interventions constituting School Reform represents an impressive landscape of causal pathways of change at different levels, directly and indirectly affecting individual school, teacher and student welfare in different ways, it can be unpacked into different (workable) components such as teacher incentives and their effects on student performance indicators or school vouchers and their effects on student performance.

Some final remarks on attribution are in order. Given the centrality of the attribution issue in impact evaluation we concur with many of our colleagues that there is scope for more quantitative impact evaluation, as these techniques offer a comparative advantage of formally addressing the counterfactual. However, at the same time it is admitted that given the limitations discussed above, the application of experimental and quasi-experimental design-based approaches will necessarily be limited to only a part of the total amount of interventions in development.

<sup>&</sup>lt;sup>20</sup> Alternatively, IE in the case of complex interventions or complex processes of change can rely on several statistical modeling approaches to capture the complexity of a phenomenon. For example, an extension of reduced form regression-based approaches to IE referred to earlier are structural equation models which can be used to model some of the more complex causal relationships that underlie interventions, using for example an intervention theory as a basis.

The combination between the theory-based evaluation approach and quantitative impact evaluation provides a powerful methodological basis for rigorous impact evaluation for several reasons:

- the intervention theory will help indicating which of the intervention components are amenable to quantitative counterfactual analysis through for example quasi-experimental evaluation and how this part of the analysis relates to other elements of the theory;

the intervention theory approach will help identifying key determinants of impact variables to be taken into account in quantitative impact evaluation;
the intervention theory can help strengthen the interpretation of findings generated by quantitative impact evaluation techniques.

This symbiosis between theory-based evaluation and quantitative impact evaluation has been acknowledged by a growing number of authors in both the general impact evaluation literature (e.g. Cook, 2000; Shadish et al., 2002: Rossi et al., 2004; Morgan and Winship, 2007) as well as in the literature on development impact evaluation (e.g. Bamberger et al., 2004, Bourguignon and Sundberg, 2007; Ravallion, 2008). When this combination is not feasible in practice, alternative methods embedded in a theory-based evaluation framework should be applied.

### 4.4. Other approaches

In this section we introduce a range of methodological approaches which can be used in specific phases of an impact evaluation or address particular aspects of the impact evaluation<sup>21</sup>.

### 4.4.1. Participatory approaches

Nowadays, participatory methods have become 'mainstream' tools in development in almost every area of policy intervention. The roots of participation in development lie in the rural sector, where Chambers (1995) and others developed the now widely used principles of participatory rural appraisal (PRA). Participatory evaluation approaches (see for example, Cousins and Whitmore, 1998) are built on the principle that stakeholders should be involved in some or all stages of the evaluation. In the case of impact evaluation this includes aspects such as the determination of objectives, indicators to be taken into account, as well as stakeholder participation in data collection and analysis.

Methodologies commonly included under this umbrella include: Appreciative Inquiry (AI), Citizen Report Cards (CRCs), Community Score Cards (CSCs), Beneficiary Assessment (BA), Participatory Impact Monitoring (PIM, see Box 4.2.), the Participatory Learning and Action (PLA) family including Rapid Rural Appraisal (RRA), Participatory Rural Appraisal (PRA), and Participatory Poverty Assessment (PPA), Policy and Social Impact Analysis (PSIA), Social Assessment (SA), Systematic Client Consultation (SSC), Self-esteem, associative strength,

<sup>&</sup>lt;sup>21</sup> See Appendices 7 and 8 for brief discussions on additional approaches applicable to IE problems in multi-level settings.

resourcefulness, action planning and responsibility (SARAR), and Objectives-Oriented Project Planning (ZOPP) (see for example Mikkelsen, 2005; Salmen and Kane, 2006).

These methods rely on different degrees of participation ranging from consultation to collaboration to joint decision-making. In general, the higher the degree of participation, the more costly and difficult it is to set up the impact evaluation. In addition, a high degree of participation might be difficult to realize in case of large-scale comprehensive interventions such as sector programs<sup>22</sup>.

More importantly, a difference we suggest should be made between:

- stakeholder participation as a process and;
- stakeholder perceptions and views as sources of evidence<sup>23</sup>.

Some of the advantages of an IE involving both of these aspects of participation are the following:

- By engaging a range of stakeholders, a more comprehensive and/or appropriate set of *valued* impacts are likely to be identified;
- Involving stakeholders in providing evidence or gathering evidence can result in more ownership and a better level of understanding among stakeholders;
- More in-depth attention to stakeholder views and opinions can help evaluators to better understand processes of change and the ways in which interventions affect people.

Disadvantages of participatory approaches to IE are:

- Limitations to the validity of information based on stakeholder perceptions (only); this problem is related to the general issue of shortcomings in individual and group perceptional data;
- Strategic responses, manipulation or advocacy by stakeholders can influence the validity of the data collection and analysis;
- limitations to the applicability of IE with a high degree of participation especially in large-scale, comprehensive, multi-site interventions (aspects of time and cost).

### Box 4.2. Participatory Impact Monitoring in the context of the Poverty Reduction Strategy process

Participatory Impact Monitoring (PIM) builds on the voiced perceptions and assessments of poor men and women and aims at strengthening these as relevant factors in decision-making at national and subnational level. In the context of PRS monitoring it will provide systematic and fast feedback on the implementation progress, early indications of outcomes, impact and on the unintended effects of policies and programs.

The purposes are the following:

- to increase the voice and the agency of poor people through participatory monitoring and

<sup>&</sup>lt;sup>22</sup> Although in such cases particular case studies of localized intervention activities within the sector program might be conducted in a participatory manner.

<sup>&</sup>lt;sup>23</sup> There are several methods for capturing these views (see for example Mikkelsen, 2005).

evaluation;

- so as to enhance the effectiveness of poverty oriented policies and programs in PRSP countries, and;
- to contribute to methodology development, strengthen the knowledge base and facilitate cross-country learning on the effective use of Participatory Monitoring on policy level, and in the context of PRS processes in particular.

Conceptually, the proposed PIM approach combines (1) the analysis of relevant policies and programs on national level, leading to an inventory of 'impact hypotheses', with (2) extensive consultations on district/local government level, and (3) joint analysis and consultations with poor communities on their perceptions of change, their attributions to causal factors and their contextualized assessments of how policies and programs effect their situation.

Source: Booth and Lucas (2002)

4.4.2. Useful methods for data collection and analysis which are often part of IE designs

In this section we distinguish a set of methods that are useful:

- for testing/refining particular parts (i.e. assumptions) of the impact theory but not specifically focused on impact assessment as such;
- for strengthening particular lines of argumentation with additional/ detailed knowledge, useful for triangulation with other sources of evidence;
- for deepening the understanding of the nature of particular relationships between intervention and processes of change.

The literature on (impact) evaluation methodology as any other field of methodology is riddled with labels representing different (and sometimes not so different) methodological approaches. In essence however, methodologies are built upon specific methods. Survey data collection and (descriptive) analysis, semi-structured interviews, focus-group interviews are but a few of the specific methods that are found throughout the landscape of methodological approaches to IE.

Evaluators, commissioners and other stakeholders in IE should have a basic knowledge about the more common research techniques<sup>24</sup>:

 Descriptive statistical techniques (e.g. of survey or registry data): the statistician Tukey (e.g. Tukey, 1977) argued for more attention to exploratory data analysis techniques as powerful and relatively simple ways to understand patterns in data. Examples are: univariate and bivariate statistical analysis of primary or secondary data using graphical analysis and simple statistical summaries (e.g. for univariate analysis:

<sup>&</sup>lt;sup>24</sup> Please note that the following methods rely on different types of data collection techniques. For example, quantitative descriptive analysis (preferably) relies on sample data based on random (simple, stratified, clustered) samples or on census data. In contrast, many qualitative methods rely on non-random sampling techniques such as purposive or snowball sampling, or do not rely on sampling at all as they might focus on a relatively small number of observations.

mean, standard deviation, median, interquartile range; e.g. for bivariate analysis: series of boxplots, scatterplots, odds ratios).

- Inferential statistical techniques (e.g. of survey or registry data): univariate analysis (e.g. confidence intervals around the mean; t-test of the mean), bivariate analysis (e.g. t-test for difference in means) and multivariate analysis (e.g. cluster analysis, multiple regression) can be rather useful in the estimation of impact effects or testing particular causal assumptions of the intervention theory. These techniques (including the first bullet point) are also used in the (quasi-)experimental and regression-based approaches described in section 3.2. For more information see for example Agresti and Finlay (1997) or Hair et al. (2005) or more specifically for development contexts see for example Casley and Lury (1987) or Mukherjee et al. (1998).
- 'Qualitative methods': including widely used methods such as semistructured interviews, open interviews, focus group interviews, discourse analysis, but also less conventional approaches such as mystery guests, unobtrusive measures (e.g. through observation)<sup>25</sup>, etc. For more information see for example Patton (2002) or more specifically for development contexts see for example Mikkelsen (2005) or Roche (1999).

### 4.5. Key message

Address the attribution problem. Although there is no single method that is best in all cases (a gold standard), some methods are indeed best in specific cases. When empirically addressing the attribution problem, experimental and quasiexperimental designs embedded in a theory-based evaluation framework have clear advantages over other designs. NONIE supports the use of randomized controlled trials where appropriate. If addressing the attribution problem can only be achieved by doing a contribution analysis, be clear about that and specify the limits and opportunities of this approach. For impact evaluations, quantitative methods are usually preferable and should be pursued when possible, and qualitative techniques should be used to evaluate the important issues for which quantification is not feasible or practical.

<sup>&</sup>lt;sup>25</sup> See for example Webb et al. (2000).

# 5. Build on existing knowledge relevant to the impact of interventions

## 5.1. Review and synthesis approaches as methods for analyzing existing evidence on impact

Review and synthesis approaches are commonly associated with systematic reviews and meta-analyses. Using these methods, comparable interventions evaluated across countries and regions can provide the empirical basis to identify 'robust' performance goals and to help assess the relative effectiveness of alternative intervention designs under different country contexts and settings. These methods can lead to increased emphasis on the rigor of impact evaluations so they can contribute to future knowledge-building as well as meet the information needs of stakeholders. These methods can also lead to a more selective approach to extensive impact evaluation, where existing knowledge is more systematically reviewed before undertaking a local impact evaluation.

The systematic review is a term which is used to indicate a number of methodologies that deal with synthesizing lessons from existing evidence. In general, one can define a systematic review as a synthesis of primary studies which contains an explicit statement of objectives and is conducted according to a transparent, systematic and replicable methodology (Greenhalgh et al., 2004). Typical features of a protocol underlying a systematic review are the following (Oliver et al., 2005):

- 1. Defining the review question(s)
- 2. Developing the protocol
- 3. Searching for relevant bibliographic sources
- 4. Defining and applying criteria for including and excluding documents
- 5. Defining and applying criteria for assessing the methodological quality of the documents
- 6. Extracting information<sup>26</sup>
- 7. Synthesizing the information into findings.

A meta-analysis is a quantitative aggregation of effect scores established in individual studies. The synthesis is often limited to a calculation of an overall effect score expressing the impact attributable to a specific intervention or a group of interventions. In order to arrive at such a calculation, meta-analysis involves a strict procedure to search for and select appropriate evidence on the impact of single interventions. The selection of evidence is based on an assessment of the methodology of the single intervention impact study. In this type of assessment usually a hierarchy of methods is applied in which randomly controlled trials rank highest and provide the most rigorous sources of evidence

<sup>&</sup>lt;sup>26</sup> This step may rely on statistical methods (meta-analysis) for analyzing and summarizing the results of included studies, if quantitative evidence at the level of single interventions studies is available and if interventions are considered similar enough.

for meta-analysis. Meta-analysis differs from multi-center clinical trials in the sense that in the former case the evaluator has no control over the single intervention evaluations as such. As a result, despite the fact that homogeneity of implementation of similar interventions is a precondition for successful metaanalysis, inevitably meta-analysis is confronted with higher levels of variability in individual project implementation, context and evaluation methodology than in the case of multi-center clinical trials.

Meta-analysis is most frequently applied in professional fields as medicine, education, and (to a lesser extent) criminal justice and social work (Clarke, 2006). Knowledge repositories like the Campbell Collaboration and Cochrane Society rely heavily on meta-analysis as a rigorous tool for knowledge management on what works. Both from within these professional fields as well as from other fields criticism has emerged. In part, this criticism reflects a resistance to the idea of a 'gold standard' underlying the practice of meta-analysis. The discussion has been useful in the sense that is has helped to define the boundaries of applicability of meta-analysis and the idea that, given the huge variability in parameters characterizing evaluations, there is no such thing as a gold standard (see Clarke, 2006).

Partly as a response to the limitations in applicability of meta-analysis as a synthesis tool, more comprehensive methodologies of systematic review have been developed. An example is a systematic review of health behavior amongst young people in the UK involving both quantitative and qualitative synthesis (see Oliver et al., 2005). The case shows that meta-analytic work on evidence stemming from what the authors call 'intervention studies' (evaluation studies on similar interventions) can be combined with qualitative systematic review of 'non-intervention studies', mainly research on relevant topics related to the problems addressed by the intervention. Regarding the latter, similar to the quantitative part, a systematic procedure for evidence search, assessment and selection is applied. The difference lies mostly in the synthesis part which in the latter case is a qualitative analysis of major findings. The two types of review can subsequently be used for triangulation purposes, reinforcing the overall synthesis findings.

Other examples of review and synthesis approaches are the narrative review and the realist synthesis. A narrative review is a descriptive account of intervention processes and/or results covering a series of interventions (see Box 5.1.). Often, the evaluator relies on a common analytical framework which serves as a basis for a template that is used for data extraction from the individual studies. In the end, the main findings are summarized in a narrative account and/or tables and matrices representing key aspects of the interventions. A realist synthesis is a theory-based approach that helps synthesizing findings across interventions. It focuses on the question which mechanisms are assumed to be at work in a given intervention, taking into account the context the intervention operates in. (see Appendix 9). Although interventions often appear different at face value, they not seldom rely on strikingly similar mechanisms. Recognition of this can broaden the range of applicable evidence from other studies.

Combinations of meta-approaches are also possible. In a recent study on the impact of public policy programs designed to reduce and/or prevent violence in the public arena, Van der Knaap et al. (2008) have shown the relevance of *combining* synthesis approaches (see Appendix 10).

### Box 5.1. Narrative review and synthesis study: Targeting and impact of community-based development initiatives

The study was performed by Mansuri and Rao (2004) who reviewed the evidence on communitybased development projects (CBD) funded by the World Bank. At the time it was estimated that an estimated US\$ 7 billion of World Bank projects are about CBD.

#### **Review questions:**

- 1. Does community participation improve the targeting of private benefits, like welfare or relief?
- 2. Are the public goods created by community participation projects better targeted to the poor?
- 3. Are they of higher quality, or better managed, than similar public goods provided by the government?
- 4. Does participation lead to the empowerment of marginalized groups—does it lessen exclusion, increase the capacity for collective action, or reduce the possibility that locally powerful elites will capture project benefits?
- 5. Do the characteristics of external agents—donors, governments, nongovernmental organizations (NGOs), and project facilitators—affect the quality of participation or project success or failure?
- 6. And finally, can community participation projects be sustainably scaled up?

In order to obtain relevant and reliable evidence on community-based development projects the reviewers decided to restrict the review process to peer-reviewed publications, or studies conducted by independent researchers. This provided an exogenous rule that improves the quality and reduces the level of potential bias while casting a wide-enough net to let in research from a variety of disciplinary perspectives on different types of CBD projects. The following sources of evidence were included: impact evaluations, which use statistical or econometric techniques to assess the causal impact of specific project outcomes; ethnographic or case studies, which use anthropological methods such as participant observation, in-depth interviews, and focus group discussions.

#### Some conclusions:

- Projects that rely on community participation have not been particularly effective at targeting the poor; there is some evidence that CBD/CDD projects create effective community infrastructure, but not a single study establishes a causal relationship between any outcome and participatory elements of a CBD project.
- A naïve application of complex contextual concepts like "participation", "social capital" and "empowerment" is endemic among project implementers and contributes to poor design and implementation.

Source: Mansuri and Rao (2004)

#### 5.2. Key message

Build on existing knowledge relevant to the impact of interventions. Review and synthesis methods can play a pivotal role in IE. Although interventions often appear

different at face value, they not seldom rely on strikingly similar mechanisms. Recognition of this can broaden the range of applicable evidence. As there are several meta-approaches available, it is worthwile to try to combine (some of) them. Review and synthesis work can provide a useful basis for empirical impact analysis of a specific intervention and in some cases may even take away the need for further indepth IE.

# 6. Use a mixed methods approach: the logic of the comparative advantages of methods

## 6.1. Different methodologies have comparative advantages in addressing particular concerns and needs in impact evaluation

The work by Campbell and others on validity and threats to validity within experiments and other types of evaluations have left deep marks on the way researchers and evaluators have addressed methodological challenges in impact evaluation (see Campbell, 1957; Campbell and Stanley, 1963; Cook and Campbell, 1979; Shadish et al, 2002). Validity can be broadly defined as the "truth of, or correctness of, or degree of support for an inference" (Shadish et al., 2002: 513).

Campbell distinguished between four types of validity which can be explained in a concise manner by looking at the questions underlying the four types:

- internal validity: how do we establish that there is a causal relationship between intervention outputs and processes of change leading to outcomes and impacts?
- construct validity: how do we make sure that the variables that we are measuring adequately represent the underlying realities of development interventions linked to processes of change?
- external validity: how do we (and to what extent can we) generalize about findings to other settings (interventions, regions, target groups, etc.)?
- statistical conclusion validity: How do we make sure that our conclusion about the existence of a relationship between intervention and impact variable is in fact true? How can we be sure about the magnitude of change<sup>27</sup>?

Applying the logic of comparative advantages makes it possible for evaluators to compare methods on the basis of their relative merits on addressing particular aspects of validity. This provides a useful basis for methodological design choice; given the evaluation's priorities, methods that better address particular aspects of validity of interest are selected in favor of others. In addition, the logic of comparative advantages can support decisions on combining methods in order to be able to simultaneously address multiple aspects of validity.

We will illustrate this logic using the example of Randomized Controlled Trials. Internal validity usually receives (and justifiably so) a lot of attention in IE as it lies at the heart of the attribution problem; is there a causal link between intervention outputs and outcomes and impacts? Arguably RCTs (see section 4.2.) are viewed by many as the best method for addressing the attribution problem *from the point of view of internal validity*. Random allocation of project benefits ensures that there are no systematic (observable and unobservable) differences between

<sup>&</sup>lt;sup>27</sup> This dimension is only addressed by quantitative impact evaluation techniques.

those that receive benefits and those that do not. However, this does not make it necessarily the best method *overall*. For example, RCTs control for differences between groups within the particular setting that is covered by the study. Other settings have other characteristics that are not controlled, hence there are limitations of *external validity* here. To resolve this issue Duflo and Kremer (2005) propose to undertake series of RCTs on the same type of instrument in different settings. However, as argued by Ravallion "the feasibility of doing a sufficient number of trials -sufficient to span the relevant domain of variation found in reality for a given program, as well as across the range of policy options- is far from clear. The scale of the randomized trials needed to test even one large national program could well be prohibitive" (Ravallion, 2008: 19).

Another limitation of RCTs (and in this case also valid for other approaches discussed in section 4.2.) lies in the realm of construct validity. Does the limited set of indicators adequately represent the impact of a policy on a complex phenomenon such as poverty? In-depth qualitative methods can more adequately capture the complexity and diversity of aspects that define (and determine) poverty than the singular or limited set of impact indicators taken into account in RCTs. Consequently, the latter have a *comparative advantage* in addressing *construct validity* concerns<sup>28</sup>. However, a downside of most qualitative approaches is that the focus is local and findings are very context-specific with limited external validity. *External validity* can be adequately addressed by, for example, quantitative quasi- and non-experimental approaches that are based on large samples covering substantial diversity in context and people.

Theory-based evaluation provides the *basis* for combining different methodological approaches that have comparative advantages in addressing validity concerns. In addition, the intervention theory as such, as a structure for making explicit causal assumptions, generalizing findings and in-depth analysis of specific assumptions, can help to strengthen internal, external and construct validity claims.

To conclude:

- there is no single best method in IE, one that can address the different aspects of validity always better than others;
- methods have particular advantages in dealing with particular validity concerns; this provides a strong rationale for combining methods.

### 6.2. Advantages of combining different methods and sources of evidence

In principle, each impact evaluation in some way is supported by different methods and sources of evidence. For example, even the quite technical

<sup>&</sup>lt;sup>28</sup> We do not want to repeat the full set of arguments in favor and against RCTs here. Worthwhile mentioning is the argument, often used against RCTs, that they are only applicable to a narrow range of interventions. However, this does not preclude that there is more scope for doing RCTs within this 'narrow' range of interventions. For further discussion on the strengths and limitations of RCTs in development see for example Duflo and Kremer (2005) or Bamberger and White (2007).

quantitative approaches described in section 4.2. include other modes of inquiry such as the research review to identify key variables that should be controlled for in for example a quasi-experimental setting. Nevertheless, there is a growing literature on the explicit use of *multiple methods* to strengthen the quality of the analysis<sup>29</sup>. At the same time the discordance between the practice and 'theory' of mixed method research (Bryman, 2006) suggests that mixed method research is often more an art than a science.

Triangulation is a key concept that embodies much of the rationale behind doing mixed method research and represents a set of principles to fortify the design, analysis and interpretation of findings in  $IE^{30}$ . Triangulation is about looking at things from multiple points of view, a method "to overcome the problems that stem from studies relying upon a single theory, a single method, a single set of data [...] and from a single investigator" (Mikkelsen, 2005: 96). As can be deducted from the definition there are different types of triangulation. Broadly, these are the following (Mikkelsen, 2005):

- data triangulation (to study a problem using different types of data, different points in time, different units of analysis);
- investigator triangulation (multiple researchers looking at the same problem);
- discipline triangulation (researchers trained in different disciplines looking at the same problem);
- theory triangulation (using multiple competing theories to explain and analyze a problem);
- methodological triangulation (using different methods, or the same method over time, to study a problem).

As can be observed from this list, particular methodologies already embody aspects of triangulation. Quantitative double difference IE (see section 4.2.) for example embodies aspects of methodological and data triangulation. Theorybased evaluation often involves theory triangulation (see section 3; see also Carvalho and White (2004) who refer to competing theories in their study on social funds). Moreover, it also allows for methodological and data triangulation by relying on different methods and sources of evidence to test particular causal assumptions. Discipline triangulation and theory triangulation both point at the need for more diversity in perspectives for understanding processes of change in IE. Strong pleas have recently been made for development evaluators to recognize and make full use of the wide spectrum of frameworks and methodologies that have emerged from different disciplines and provide evaluation with a rich arsenal of possibilities (Kanbur 2003; White, 2002;

<sup>&</sup>lt;sup>29</sup> The most commonly used term is mixed methods (see for example Tashakkori and Teddlie, 2003). In case of development research and evaluation see for example Bamberger (2000) and Kanbur (2003).

<sup>&</sup>lt;sup>30</sup> This is true for the broad interpretation of the concept of triangulation as used by for example Mikkelsen (2005). Other authors use the concept in a more restrictive way (e.g. Bamberger (2000) uses triangulation in the more narrow sense of validating findings by looking at different data sources).

Bamberger and White, 2007). For example, when doing IEs, evaluators can benefit from approaches developed in different disciplines and sub-disciplines. Among others, neo-institutionalist economists have shown ways to study the impact of institutions as 'rules of the game' (see North, 1990), and interventions such as policies can be considered as attempts to establish specific rules of the game with the expectation (through a 'theory of change') of generating certain impacts (Picciotto and Wiesner, 1997). In addition, the literature on behavioral and social mechanisms (see Appendix 9) provides a wealth on explanatory insights that help evaluators to better understand and frame processes of change triggered by interventions.

Good methodological practice in IE is to encourage applying these principles of triangulation as much as possible. Advantages of mixed method approaches to IE are the following:

- A mix of methods can be used to assess important outcomes or impacts of the intervention being studied. If the results from different methods converge, then inferences about the nature and magnitude of these impacts will be stronger. For example, triangulation of standardized indicators of children's educational attainments with results from an analysis of samples of children's academic work yields stronger confidence in the educational impacts observed than either method alone (especially if the methods employed have offsetting biases).
- A mix of methods can be used to assess different facets of complex outcomes or impacts, yielding a broader, richer portrait than can one method alone. For example, standardized indicators of health status could be mixed with onsite observations of practices related to dietary nutrition, water quality, environmental risks or other contributors to health, jointly yielding a richer understanding of the intervention's impacts on targeted health behaviors. In a more general sense, quantitative IE techniques work well for a limited set of pre-established variables (preferably determined and measured ex ante) but less well for capturing unintended less expected (indirect) effects of interventions. Qualitative methods or descriptive (secondary) data analysis can be helpful in better understanding the latter.
- One set of methods could be used to assess outcomes or impacts and another set to assess the quality and character of program implementation, including program integrity and the experiences during the implementation phase.
- Multiple methods can help to ensure that the sampling frame and the sample selection strategies cover the whole of the target intervention and comparison populations. Many sampling frames leave out important sectors of the population (usually the most vulnerable groups or people who have recently moved into the community), while respondent selection procedures often under-represent women, youth or the elderly or ethnic minorities. This is critical because important positive or negative impacts on the vulnerable groups (or other important sectors) are completely ignored if they do not even get included in the sample. This is particularly

important (and frequently ignored) where the evaluation uses secondary data sets, as the evaluator often does not have access to information on how the sample was selected.

Appendix 11 presents four interesting examples of IEs that are based on a mixed method perspective:

- Case 1: Combining qualitative and quantitative descriptive methods Ex post impact study of the Noakhali Rural Development Project in Bangladesh;
- Case 2: Combining qualitative and quantitative descriptive methods Mixed method impact evaluation of IFAD projects in Gambia, Ghana and Morocco
- Case 3: Combining qualitative and quantitative descriptive methods Impact evaluation: agricultural development projects in Guinea
- Case 4: A theory-based approach with qualitative methods GEF Impact Evaluation 2007

### 6.3. Key message

Use a mixed methods design. Bear in mind the logic of the comparative advantages of designs and methods. A mix of methods can be used to assess different facets of complex outcomes or impacts, yielding more breadth, depth and width in the portrait than can one method alone. One set of methods could be used to assess outcomes or impacts and another set to assess the quality and nature of intervention implementation, thus enhancing impact evaluation with information about program integrity and program experiences.

### **PART II – MANAGING IMPACT EVALUATIONS**

### 7. Determine if an IE is feasible and worth the cost

### 7.1. Evaluability

Managers and policymakers sometimes assume that impact evaluation is synonymous with any other kind of evaluation. They might request an 'impact evaluation' when the real need is for a quite different kind of evaluation (for example to provide feedback on the implementation process, or to assess the accessibility of program services to vulnerable groups). Ensuring clarity in the information needed and for what purposes is a prerequisite to defining the type of evaluation to be conducted.

Moreover, IE is not 'the' alternative but, draws on, and complements rather than replaces other types of monitoring and evaluation activities. It should therefore be seen as one of several in a cycle of potentially useful evaluations in the lifetime of an intervention. The rather traditional difference between ex ante and ex post impact evaluations remains important, where the ex ante impact assessment is, by nature, largely an activity in which 'predictions' are made of any effects and side effects a particular intervention might have, when being implemented. Ex post IE, or simply 'IE' as defined by the development community (and elsewhere) can test whether or not and to what extent these ex ante predictions have been correct. In fact, one of the potential uses of IE, not yet frequently applied in the field of development intervention, could be to strengthen the process of ex ante impact assessments.

When should an impact evaluation ideally be conducted?

- When there is an *articulated need* to obtain the information from an IE in order to know whether the intervention worked, to learn from it, to increase transparancy of the intervention and to know its 'value for money';
- When a 'readiness assessment' shows that political, technical, resource and other practical considerations are adequate and it is feasible to do an IE. More specifically, this would include the following conditions:
  - The evaluation has a clearly defined purpose and agreed upon intended use, appropriate to its timing and with support of influential stakeholders.
  - There is clarity about the evaluation design. The evaluation design has to be clearly described and well justified after due consideration of alternatives and constraints.
  - The evaluation design has a chance to be credibly executed given the nature and context of the intervention, the data and

information needs and the availability of adequate resources and expertise to conduct the evaluation.

Impact evaluations may not be appropriate when:

- Other valuable forms of evaluation will yield more useful information to support decisions to be made or other purposes;
- It moves too much resources and attention away from the need to develop and use a rich spectrum of evaluation approaches and capacities;
- Political, technical, practical or resource considerations are likely to prevent a credible, rigorous and useful evaluation;
- There are signs that the evaluation will not be used (or may be misused, for example for political reasons).

Not all interventions should be submitted to elaborate and costly IE exercises. Rather, those sectors, regions and intervention approaches about which less is known (including new, innovative ideas) should receive funding and support for impact evaluation. Ideally, organizations should pool their resources and expertise to select interventions of interest for rigorous and elaborate impact evaluation and consequently contribute jointly to the public good of knowledge on impact of (under-evaluated) interventions.

### 7.2. Key message

Determine if an impact evaluation is feasible and worth the cost. What are the benefits of the IE? In what ways does the IE contribute to accountability, learning and information about the value for money' about what works? What is the likely added value of an IE in relation to what is already known about a particular intervention? What are the costs of the IE? What are the costs of estimating or measuring what would have happened without the intervention? Is the likelihood of getting accurate information on impact high enough to justify the cost of the evaluation?

### 8. Start early - getting the data

Although issues of data and data collection like availability and quality often sound like 'mere' operational issues that only need to be discussed on a technical level, it should not be forgotten that these aspects are of crucial importance for any impact evaluation (and any evaluation in general). Data are needed to test whether or not there have been changes in the dependent variables or to represent the *counterfactual* estimate of what would have been the situation of the project population if the project had not taken place. The data issue is strongly linked to the type of method.

### 8.1. Timing of data collection

Ideally, impact evaluations should be based on data presenting the situation before an intervention took place and after the intervention has been implemented<sup>31</sup>. An important question is if the baseline period or end-line period is representative or normal. If the baseline or end-line year (or season) are not normal, then this affects the observed change over time. If for example the baseline year is influenced by unusually high/low agricultural production, or a natural disaster, then the observed change up to the end-line year can be strongly biased, when conclusions are drawn on the impact of an intervention during those years. In most cases it is the timing of the intervention, or the impact evaluation, which determines the timing of the baseline and end-line studies. This timing is not random, and it is necessary that evaluators investigate if the baseline/end-line data are representative for 'normal' years/periods, before drawing conclusions. If not, even rigorous evaluations may produce unreliable conclusions about impacts. An additional issue concerns short term versus long term effects, discussed earlier. Depending on the intervention and its context, at the time of ex post data collection some effects might not have occurred or not be visible yet, whereas others might whither over time. The evaluator should be aware of how this affects conclusions about impact.

### 8.2. Data availability

In practice, impact evaluation starts with an appraisal of existing data, the data which have been produced during the course of an intervention on inputs, processes, outputs (and outcomes). This inventory is useful for several reasons, to mention but a few:

- available data are useful for reconstructing the intervention theory that further guides primary and secondary data collection efforts;
- available data might affect the choice of methodological design or options for further data processing and analysis; for example, ex ante and ex post data sets of target groups might be complemented with other (existing)

<sup>&</sup>lt;sup>31</sup> In some cases talking about the 'end' of an intervention is not or less applicable, for example in case of institutional reforms, new legislation, fiscal policy, etc.

data sets to construct useful control groups; the amount and type of data available might influence the choice of whether or not to organize additional primary data collection efforts;

• available data from different sources allow for triangulation of findings.

In addition, evaluators can rely on a variety of data from other sources that can be used in the evaluation process. We briefly mention the following sources:

- national census data;
- general household surveys such as Living Standards Measurement Surveys (LSMS);
- specialized surveys such as Demographic and Health Surveys (DHS);
- administrative data collected by line ministries and other public agencies (e.g. on school enrolment, use of health facilities, market prices for agricultural produce);
- studies conducted by donor agencies, non-government organizations and universities;
- administrative data from agency, ministries or other organizations;
- mass media (newspapers, television documentaries, etc.); these can be useful, among other things, for understanding the local economic and political context of an intervention.

Appendix 12 describes an example of an impact evaluation implemented by IEG. In 1986 the Government of Ghana embarked on an ambitious program of educational reform, shortening the length of pre-University education from 17 to 12 years, reducing subsidies at the secondary and tertiary levels, increasing the school day and taking steps to eliminate unqualified teachers from schools. There was no clearly defined 'project' for this study, but the focus was World Bank support to the sub-sector through four large operations. These operations had supported a range of activities, from rehabilitating school buildings to assisting in the formation of community-based school management committees. The impact evaluation heavily relied on existing data sets such as the *Ghana Living Standards Survey* for impact analyses.

A useful stepwise approach for assessing data availability is the following:

- Make an inventory of the availability of data and assess the quality of available data. Sometimes secondary data can be used to carry out the whole impact study. This is especially true when evaluating national or sector-wide interventions. More usually, secondary data can be used to buttress other data.
- 2. Analyze, from the perspective of the intervention theory, the necessity of additional data. The process of data gathering must be based on the evaluation design which is, in turn, (partly) based on the intervention theory. Data must be collected across the results chain, not just on outcomes.
- 3. Assess the best way(s) to get the additional data.
- 4. A comparison group sample must be of adequate size, and subject to the same, or virtually the same, questionnaire or other data collecting instruments. Whilst some intervention-specific questions may not be

appropriate, similar questions of a more general nature can help test for contagion.

- 5. It is necessary to check if other interventions, unexpected events or other processes have influenced developments in the comparison group or the treatment group (i.e. check if the comparison group is influenced by other processes than in case of the treatment group).
- 6. Multiple instruments (e.g. household and facility level) are usually desirable, and must be coded in such a way that they can be linked.
- 7. *Baseline data* must cover the relevant welfare indicators, but preferably also the main determinants of the relevant welfare elements, so it will be easier to investigate later on if other processes than the intervention have influenced welfare developments over time. *End-line data* must be collected across the results chain, not just on intended outcomes.

When there is *no baseline*, the option of a field survey using recall on the variables of interest may be considered. Many commentators are critical of relying on recall. But all survey questions in the end are recall, so it is a question of degree. The evaluator needs to use his or her judgment (and knowledge about cognitive processes) as to what are credible data given a respondent's capacity to recall.

### 8.3. Quality of the data

The quality of data can make or break any impact evaluation. Mixed methods and triangulation are strategies to reduce the problem of data quality. Yet, they are insufficient in terms of the quality control that is needed to ensure that evaluation findings are not (heavily) biased due to data quality problems.

Several questions should be asked by the evaluator:

- What principles should we follow to improve the quality of data (collection)<sup>32</sup>? Some examples of subquestions:
  - How to address missing data (missing observations in a data set; missing variables)?
  - How to address measurement error? Does the value of a variable or the answer to a question represent the true value?
  - How to address specification error? Does the question asked or variable measured represent the concept that it was intended to cover?
- Does the quality of the data allow for (advanced) statistical analysis? New advances in and the more widespread use of quasi-experimental evaluation and multivariate data analysis are promising in the light of IE. Yet, often data quality is a constraining factor in terms of the quality of the findings (see for example Deaton, 2005).
- In case of secondary data, what do we know about the data collection process that might strengthen or weaken the validity of our findings<sup>33</sup>?

<sup>&</sup>lt;sup>32</sup> Please note that data quality is correlated with the type of data and the method(s) generating the data.

<sup>&</sup>lt;sup>33</sup> For example in the case of secondary data sets, what do we know about the quality of the data collection (e.g. sampling errors, training and supervision of interviewers), data processing (e.g.

De Leeuw et al. (2008) discuss data quality issues in survey data analysis. Much of their discussion on measurement error (errors resulting from respondent, interviewer, method and question-related sources or a combination of these; examples are recall problems or the sensitivity of certain topics) is equally relevant for semi-structured interviews and similar techniques in qualitative research. With respect to group processes in qualitative research Cooke (2001) discusses three of the most widely cited problems: risky shift, groupthink and coercive persuasion. A detailed discussion of these issues is beyond the scope of this guidance. However, they lead us to some important points:

- data based on the perceptions, views and opinions of people on the causes and effects of an intervention (e.g. target groups) do not necessarily adequately reflect the real causes of an intervention; data that is collected through observation, measurement or counting (e.g. assets, farm size, infrastructure, profits) in general is less prone to measurement error (but is not always easy to collect nor sufficient to cover all information needs);
- the quality of data is more often than not a constraining factor in the overall quality of the impact evaluation; it cannot be solved by sophisticated methods, it might be solved in part through triangulation between data sources.

### 8.4. Dealing with data constraints

According to Bamberger et al. (2004: 8): "Frequently, funds for the evaluation were not included in the original project budget and the evaluation must be conducted with a much smaller budget than would normally be allocated for this kind of study. As a result, it may not be possible to apply the desirable data collection instruments (tracer studies or sample surveys, for example), or to apply the methods for reconstructing baseline data or creating control groups." Data problems are often correlated with or compounded by time and budget constraints. The following scenario's can occur (see Table 8.1.).

### Table 8.1. Evaluation scenario's with time, data and budget constraints

dealing with missing values, weighting issues)? We cannot simply take for granted that a data set is free from error and bias. Lack of information on the process of generating the data base inevitably constraints any subsequent data analysis efforts.

| The Constraints Under<br>Which the Evaluation<br>Must be Conducted |        |      |  |
|--|--------|------|--|
| Time   | Budget | Data | Typical Scenarios  |
| x  |        |      | The evaluator is called in late in the project and is told that the<br>evaluation must be completed by a certain date so that it can be used<br>in a decision making process or contribute to a report. The budget<br>may be adequate but it may be difficult to collect or analyze survey<br>data within the time-frame.  |
|  | Х      | v    | The evaluation is only allocated a small budget, but there is not<br>necessarily excessive time pressure. However, it will be difficult to<br>collect sample survey data because of the limited budget.  |
|  |        | х    | Consequently no baseline survey has been conducted either on the<br>project population or on a control group. The evaluation does have an<br>adequate scope, either to analyze existing household survey data or<br>to collect additional data. In some cases the intended project impacts<br>may also concern changes in sensitive areas such as domestic<br>violence, community conflict, women's empowerment, community<br>leadership styles, or corruption on which it is difficult to collect<br>reliable data—even when time and budget are not constraints. |
| Х  | х      |      | The evaluator has to operate under time pressure and with a limited<br>budget. Secondary survey data may be available but there is little<br>time or resources to analyze it.  |
| Х  |        | Х    | The evaluator has little time and no access to baseline data or a control group. Funds are available to collect additional data but the survey design is constrained by the tight deadlines.   |
|  | Х      | Х    | The evaluator is called in late and has no access to baseline data or control groups. The budget is limited but time is not a constraint.  |
| х  | х      | х    | The evaluator is called in late, is given a limited budget, has no access<br>to baseline survey data and no control group has been identified.   |

Source: Bamberger et al. (2004)

The latter part of the article by Bamberger et al. (2004) describes scenarios for working within these constraints. For example, the implications for quasi-experimental designs are that evaluators have to rely on less robust designs such as ex post comparisons only (see Appendix 13).

### 8.5. Key message

Start early. Good baseline data are essential to understanding and estimating impact. Depending on the type of intervention, the collection of baseline data, as well as the setup of other aspects of the IE, requires an efficient relationship between the impact evaluators and the implementers of the intervention that is being evaluated. Policy makers and commissioners need to involve experts in impact evaluation as early as possible in the intervention design in order to be able to design high-quality impact evaluations. Ensuring high-quality data collection should be part and parcel of every IE. When working with secondary data, a lack of information on the quality of data collection can restrict data analysis options and validity of findings. Take notice of and deal effectively with the restrictions under which an impact evaluation has to be carried out (time, data and money).

### 9. Front-end planning is important

### 9.1. Front-end planning

Front-end planning refers to the initial planning and design phase of an IE. Ad hoc commissioned IEs usually do not have a long period of time for planning, thereby risking a sub optimally planned and executed impact evaluation process. As good IE relies on good data, preferably including baseline data, attention for proper front-end planning of IEs should be a priority issue. Ideally, front-end planning of IE should be closely articulated to the initial design and planning phase of the policy intervention. Indeed, this articulation is most clearly visible in an RCT, in which intervention and IE are inextricably linked. In this section we briefly recapitulate some of the priority issues in front-end planning of IE.

### Planning tools

Clear definition of scope (sections 1 and 2 in this Guidance) and sound methodological design (sections 3 to 6) cannot be captured in standardized frameworks. Decision trees on assessing data availability (see section 8.2.) and method choice (see Appendix 6) are useful though provide only partial answers to methodological design choice issues. Pragmatic considerations of time, budget and data (see section 8.4.; Bamberger et al., 2004) but also culture and politics play a role. Two tools that are particularly helpful in the planning phase of an IE are the *approach paper* and the *evaluation matrix*.

The *approach paper* outlines what the evaluation is about and how it will be implemented. This document can be widely circulated and gives stakeholders and others a chance to comment and improve upon the intended evaluation design from an early stage. It also helps to generate broad 'buy-in' or at worst to define the main grounds of potential disagreement between evaluators and practitioners. In addition, it is wise to use an *evaluation matrix* when planning and executing the work. Such a matrix ensures that key questions are identified, together with the ways to address them, sources of data, role of theory, etc. This can also play an important role in stakeholder consultation to ensure that important elements are not omitted.

### Staffing and resources

There is no such thing as a free lunch. This is also true for evaluation projects. Resources are important and the spending should be safeguarded up front. The longer the time horizon of a study, the more difficult this is. Resources are also important to realize the much-needed independence of an evaluator ('s team). A template for assessing the independence of evaluation organizations can be downloaded from <u>http://www.ecgnet.org/docs/ecg.doc</u>. It specifies a number of criteria and questions that can be asked.

Evaluation is not only a *financial resources' business* but even more a *people's business*. So is the *planning of an evaluation*. As evaluation projects usually no longer are *lonely hunter-activities*, *staffing* is crucial. So when starting the preparation of the study, a crucial point concerns addressing questions like:

- Who are the people that do the evaluation?
- Under which (contractual) conditions are they 'doing the job'?
- What is their expertise? And;
- Which roles will they be carrying out?

Topics that deserve attention are the following.

- The mix of disciplines and traditions that are brought together in the team;
- The competencies the team has 'in stock'. Competencies range from methodological expertise to negotiating with institutional actors and stakeholders, getting involved in 'hearing both sides' (evaluand and principal) and in the clearance of the report;
- The structure of the evaluation team. In order for the evaluation to be planned and carried out effectively, the roles of the project director, staff, and other evaluators must be made clear to all parties;
- The responsibilities of the team members;
- The more an evaluation is linked to a political 'hot spot', the more it is necessary that at least one member of the team has a 'political nose'. Not primarily to deal with administrators and (local) politicians, but to understand when an evaluation project becomes too much of what is known as a 'partnerial evaluation' (Pollit, 1999);
- Also, staff should be active in realizing an adequate documentation and evaluation trail.

A range of skills is needed in evaluation work. The quality and eventual utility of the impact evaluation can be greatly enhanced with coordination between team members and policymakers from the outset. It is therefore important to identify team members as early as possible, agree upon roles and responsibilities, and establish mechanisms for communication during key points of the evaluation.

The balance between independence and collaboration between evaluators and stakeholders

One of the questions on the agenda within the world of impact evaluations is what degree of institutional separation to put in place between the evaluation providers and the evaluation users. There is much to be gained from the objectivity provided by having the evaluation carried out independently of the institution responsible for the project being evaluated. Pollitt (1999) warned against 'partnerial evaluations', where positions of stakeholders, commissioners and evaluators blurred (too much)<sup>34</sup>. However, evaluations can often have

<sup>&</sup>lt;sup>34</sup> An example from Europe stresses this point. In some situations, educational evaluators of the Danish Evaluation Institute (EVA) discussed their reports with up to 20-plus stakeholders before the report is cleared and published (Leeuw, 2003).

multiple goals, including building evaluation capacity within government agencies and sensitizing program operators to the realities of their projects once these are carried out in the field. At a minimum, the evaluation users, who can range from government agencies in client countries to, bilateral and multilateral donors, international NGOs, and grass roots / civil society organizations, must remain sufficiently involved in the evaluation to ensure that the evaluation process is recognized as legitimate and that the results produced are relevant to their information needs. Otherwise, the evaluation results are less likely to be used to inform policy. The evaluation manager and his or her clients must achieve the right balance between involving the users of evaluations and maintaining the objectivity and legitimacy of the results (Baker, 2000).

#### Ethical issues

It is important to take the ethical objections and political sensitivities seriously. There can be ethical concerns with deliberately denying a program to those who need it and providing the program to some who do not; this applies to both experimental and non-experimental methods. For example, with too few resources to go around, randomization may be seen as a fair solution, possibly after conditioning on observables. However, the information available to the evaluator (for conditioning) is typically a partial subset of the information available 'on the ground' (including voters/taxpayers). The idea of 'intention-totreat' helps alleviate these concerns; one has a randomized assignment, but anyone is free to not participate. But even then, the 'randomized out' group may include people in great need. All these issues must be discussed openly, and weighed against the (potentially large) longer-term welfare gains from better information for public decision-making (Ravallion, 2008).

### Norms and standards

As said before, often impact evaluations are designed, implemented, analyzed, disseminated and used under budget, time and data constraints while facing diverse and often competing political interests. Given these constraints, the management of a real-world evaluation is much more complicated than textbook descriptions. Evaluations sometimes fail because the stakeholders were not involved, or the findings were not used because they did not address the priorities of the stakeholders. Others fail because of administrative or political difficulties in getting access to the required data, being able to meet with all of the individuals and groups that should be interviewed, or being able to ask all the questions that the evaluator feel are necessary. Many other evaluations fail because the sampling frame, often based on existing administrative data, omits important sectors of the target population - often without anyone being aware of this. In other cases the budget was insufficient, or was too unpredictable to permit an adequate evaluation to be conducted. Needless to say that evaluations also fail because of emphasizing stakeholders' participation too much (leading to 'partnerial evaluations' (Pollitt, 1999)) and because of insufficient methodological and theoretical expertise.

While many of these constraints are presented in the final evaluation report as being completely beyond the control of the evaluator, in fact their effects could very probably have been reduced by more effective management of the evaluation. For example, a more thorough scoping analysis could have revealed many of these problems and the client(s) could then have been made aware of the likely limitations on the methodological rigor of the findings. The client(s) and evaluator could then strategize to either seek ways to increase the budget or extend the time, or agree to limit the scope of the evaluation and what it promises to deliver. If clients understand that the current design will not hold up under the scrutiny of critics, they can find ways to help address some of the constraints.

For the sake of honest commitment to development, evaluators and evaluation units should ensure that impact evaluations are designed and executed in a manner that limits manipulation of processes or results towards any ideological or political agenda. They should also ensure there are realistic expectations of what can be achieved by a single evaluation within time and resource constraints, and that findings from the evaluation are presented in ways that are accessible to the intended users. This includes finding a balance between simple, clear messages and properly acknowledging the complexities and limitations of the findings.

International evaluation standards (such as the OECD-DAC / UNEG Norms and Standards and/or the standards and guidelines developed by national or regional evaluation associations) should be applied where appropriate (Picciotto, 2004).

Greater emphasis on impact evaluation for evidence-based policy-making can create greater risk of manipulation aimed at producing desirable results (positive or negative) (House, 2008). Impact evaluations require an honest search for the truth and thus place high demands on the integrity of those commissioning and conducting them. For the sake of honest commitment to development, evaluators and evaluation units need to ensure that impact evaluations are designed and executed in a manner that limits manipulation of processes or results towards any ideological or political agenda.

### Ownership and capacity-building

*Capacity-building* at the level of governmental or non-governmental agencies involved should be an explicit purpose in IE. In cases where sector-wide investment programs are financed by multi-donor co-financing schemes, the participating donors would make natural partners for a joint evaluation of that sector program<sup>35</sup>. Other factors in selecting other donors as partners in joint evaluation work may be relevant as well. Selecting those with similar development philosophies, organizational cultures, evaluation procedures and techniques, regional affiliations and proximity (etc.) may make working together

<sup>&</sup>lt;sup>35</sup> See OECD-DAC (2000). This paragraph is copied from these guidelines.

easier. Another issue may be limiting the total number of donors to a 'manageable' number. In cases where a larger group of donors is involved, a key group of development partners (including national actors) may assume management responsibilities where the role of others is more limited. Once appropriate donors are identified that have a likely stake in an evaluation topic, the next step is to contact them and see if they are interested in participating. In some cases, there may already be an appropriate donor consortium or group where the issue of a joint evaluation can be raised and expressions of interest can be easily solicited. The DAC Working Party on Aid Evaluation, the UN Evaluation Group (UNEG) and the Evaluation Cooperation Group (ECG) have a long tradition of cooperation, shared vision on evaluation principles, and personal relationships built over the years and have fostered numerous joint evaluations.

The interaction between the international development evaluation community, the countries/regions themselves and the academic evaluation communities should also be stimulated as it is likely to affect the pace and quality of capacity-building in IE. Capacity-building will also strengthen (country and regional) *ownership* of impact evaluation. Providing a space for consultation and agreement on IE priorities among the different stakeholders of an intervention will also help to enhance utilization and ownership.

### 9.2. Key message

Front-end planning is important. It can help to manage the study, its reception and its use. When managing the evaluation keep a clear eye on items such as costs, staffing, ethical issues and level of independence (of the evaluator(s' team) versus level of collaboration with stakeholders). Pay attention to country and regional ownership of impact evaluation and capacity-building and promote it. Providing a space for consultation and agreement on IE priorities among the different stakeholders of an intervention will help to enhance utilization and ownership.

### References

ADB (2006) Impact evaluation – methodological and operational issues, Economics and Research Department, Asian Development Bank, Manila.

Agresti, A. and B. Finlay (1997) Statistical Methods for the Social Sciences, Prentice Hall, New Jersey.

Baker, J.L. (2000), Evaluating the impact of development projects on poverty, The World Bank, Washington D.C.

Bamberger, M. (2000) "Opportunities and challenges for integrating quantitative and qualitative research", in: M. Bamberger (ed.) *Integrating quantitative and qualitative research in development projects*, World Bank, Washington D.C.

Bamberger, M. (2006) Conducting Quality Impact Evaluations under Budget, Time and Data Constraints, World Bank, Washington D.C.

Bamberger, M., J. Rugh, M. Church and L. Fort (2004) "Shoestring Evaluation: Designing Impact Evaluations under Budget, Time and Data Constraints", *American Journal of Evaluation* 25(1), 5-37.

Bamberger, M. J. Rugh and L. Mabry (2006) *RealWorld Evaluation Working Under Budget, Time, Data, and Political Constraints*, Sage Publications, Thousand Oaks.

Bamberger, M. and H. White (2007) "Using strong evaluation designs in developing countries: Experience and challenges", *Journal of Multidisciplinary Evaluation* 4(8), 58-73.

Bemelmans-Videc, M.L. and R.C. Rist (eds.) (1998) *Carrots, Sticks and Sermons:* Policy Instruments and their Evaluation, Transaction Publishers, New Brunswick.

Booth, D. and H. Lucas, H. (2002) "Good Practice in the Development of PRSP Indicators", *Working Paper* 172, Overseas Development Institute, London.

Bourguignon, F. and M. Sundberg (2007), "Aid effectiveness, opening the black box", American Economic Review 97(2), 316-321.

Bryman, A. (2006) "Integrating quantitative and qualitative research: how is it done?" *Qualitative Research* 6(1), 97-113.

Bunge, M. (2004) "How Does It Work? The Search for Explanatory Mechanisms", *Philosophy of the Social Sciences* 34(2), 182-210.

Campbell, D.T. (1957) "Factors relevant to the validity of experiments in social settings", *Psychological Bulletin* 54, 297-312.

Campbell, D.T. and J.C. Stanley (1963) "Experimental and quasi-experimental designs for research on teaching", in: N. L. Gage (ed.) *Handbook of research on teaching*, Rand McNally, Chicago.

Carvalho, S., and H. White (2004) "Theory-based evaluation: The case of social funds", American Journal of Evaluation 25(2), 141–160.

Casley, D.J. and D.A. Lury (1987) *Data* Collection in Developing Countries, Oxford University Press, New York.

CGD (2006) When will we ever learn? Improving lives through impact evaluation, Report of the Evaluation Gap Working Group, Center for Global Development, Washington, DC.

Chambers, R., (1995) "Paradigm Shifts and the Practice of Participatory Research and Development", in: S. Wright and N. Nelson (eds.) *Power and Participatory Development: Theory and Practice*, Intermediate Technology Publications, London.

Clarke, A. (2006) "Evidence-Based Evaluation in Different Professional Domains: Similarities, Differences and Challenges", in: I.F. Shaw, J.C. Greene and M.M. Mark (eds.) The SAGE Handbook of Evaluation, Sage Publications, London.

Coleman, J.S. (1990) Foundations of Social Theory, Belknap Press, Cambridge.

Cook, T.D. (2000) "The false choice between theory-based evaluation and experimentation", in: P.J. Rogers, T.A. Hacsi, A. Petrosino and T.A. Huebner (eds.) (2000) *Program theory in evaluation: challenges and opportunities*, New Directions for Evaluation, 87, Jossey-Bass, San Francisco.

Cook, T.D. and D.T. Campbell (1979) Quasi-Experimentation: Design and Analysis for Field Settings, Rand McNally, Chicago.

Cooke, B. (2001) "The Social Psychological Limits of Participation?", in: B. Cooke and U. Kothari (eds.) *Participation: The New Tyranny*?, Zed Books, London.

Connell, J.P., A.C. Kubisch, L.B. Schorr and C.H. Weiss (eds.) (1995) *New approaches to evaluating community initiatives*, 1, The Aspen Institute, Washington D.C.

Cousins, J.B. and E. Whitmore (1998) "Framing Participatory Evaluation", in: E. Whitmore (ed.) Understanding and Practicing Participatory Evaluation, New Directions for Evaluation 80, Jossey-Bass, San Francisco.

Deaton, A. (2005) "Some remarks on randomization, econometrics and data", in: G.K. Pitman, O.N. Feinstein and G.K. Ingram (eds.) *Evaluating development effectiveness*, Transaction publishers, New Brunswick. Dehejia, R. (1999) "Evaluation in multi-site programs", Working paper, Columbia University and NBER,

http://emlab.berkeley.edu/symposia/nsf99/papers/dehejia.pdf (last consulted January 12, 2009).

De Leeuw, E.D., J.J. Hox and D.A. Dillman (eds.) (2008) *International Handbook of Survey Methodology*, Lawrence Erlbaum Associates, London.

Duflo, E. and M. Kremer (2005) "Use of randomization in the evaluation of development effectiveness", in: G.K. Pitman, O.N. Feinstein and G.K. Ingram (eds.) *Evaluating development effectiveness*, Transaction publishers, New Brunswick.

Elbers, C., J.W. Gunning and K. De Hoop (2008) "Assessing sector-wide programs with statistical impact evaluation: a methodological proposal", *World Development* 37(2), 513-520.

Elster, J. (1989) Nuts and Bolts for the Social Sciences, Cambridge University Press, Cambridge.

Elster, J. (2007) Explaining Social Behavior - More Nuts and Bolts for the Social Sciences, Cambridge University Press, Cambridge.

Farnsworth, W. (2007) The Legal Analyst - A Toolkit for Thinking about the Law, University of Chicago Press, Chicago.

GEF (2007) "Evaluation of the Catalytic Role of the GEF", Approach Paper, GEF Evaluation Office, Washington D.C.

Gittinger, J.P. (1982) Economic analysis of agricultural projects, Johns Hopkins University Press, Baltimore.

Greenhalgh, T., G. Robert, F. Macfarlane, P. Bate and O. Kyriakidou (2004) "Diffusion of Innovations in Service Organizations: Systematic Review and Recommendations", *The Milbank Quarterly* 82(1), 581-629.

Hair, J.F., B. Black, B. Babin, R.E. Anderson and R.L. Tatham (2005) *Multivariate Data Analysis*, Prentice Hall, New Jersey.

Hansen, H.F. and Rieper, O. (2009) "Institutionalization of second-order evidence producing organizations", in: O. Rieper, F.L. Leeuw and T. Ling (eds.) *The Evidence Book: concepts, generation and use of evidence*, Transaction Publishers, New Brunswick.

Hedström, P. (2005) Dissecting the Social: On the Principles of Analytical Sociology, Cambridge University Press, Cambridge.

Hedström, P. and R. Swedberg (1998) Social Mechanisms: An Analytical Approach to Social Theory, Cambridge University Press, Cambridge.

Henry, G.T. (2002) "Choosing Criteria to Judge Program Success – A Values Inquiry", *Evaluation* 8(2), 182-204.

House, E. (2008) "Blowback: Consequences of Evaluation for Evaluation", *American Journal of Evaluation* 29(4), 416-426.

Jones, N., C. Walsh, H. Jones and C. Tincati (2008) Improving impact evaluation coordination and uptake - A scoping study commissioned by the DFID Evaluation Department on behalf of NONIE, Overseas Development Institute, London.

Kanbur, R. (ed.) (2003) Q-Squared: Combining Qualitative and Quantitative Methods in Poverty Appraisal, Permanent Black, Delhi.

Kellogg Foundation (1991) Information on Cluster Evaluation, Kellogg Foundation, Battle Creek.

Kraemer, H.C. (2000) "Pitfalls of Multisite Randomized Clinical Trials of Efficacy and Effectiveness." Schizophrenia Bulletin 26, 533-541.

Kruisbergen, E.W. 2005. Voorlichting: doen of laten? Theorie van afschrikwekkende voorlichtingscampagnes toegepast op de casus van bolletjesslikkers, *Beleidswetenschap* 19(3), 38 -51.

Kusek, J and R.C. Rist (2005) Ten Steps to a Results-Based Monitoring and Evaluation System: A Handbook for Development Practitioners, World Bank, Washington D.C.

Lawson, A., D. Booth, M. Msuya, S. Wangwe and T. Williamson (2005) *Does general budget support work? Evidence from Tanzania*, Overseas Development Institute, London.

Leeuw, F.L. (2003) "Reconstructing Program Theories: Methods Available and Problems to be Solved", American Journal of Evaluation 24(1), 5-20.

Leeuw, F.L. (2005) "Trends and Developments in Program Evaluation in General and Criminal Justice Programs in Particular", *European Journal on Criminal Policy and Research* 11, 18-35.

Levinsohn, J., S. Berry, and J. Friedman (1999) "Impacts of the Indonesian Economic Crisis: Price Changes and the Poor", *Working Paper* 7194, National Bureau of Economic Research, Cambridge. Lipsey, M.W. (1993) "Theory as Method: Small Theories of Treatments," in: L.B. Sechrest and A.G. Scott (eds.), *Understanding Causes and Generalizing about Them*, New Directions for Program Evaluation 57, Jossey-Bass, San Franscisco.

Lister, S. and R. Carter (2006) *Evaluation of general budget support: synthesis report*, joint evaluation of general budget support 1994 –2004, Department for International Development, University of Birmingham.

Maluccio, J. A. and R. Flores (2005) "Impact evaluation of conditional cash transfer program: The Nicaraguan Red de Protección Social", International Food Policy Research Institute, Washington D.C.

Mansuri, G & V.Rajo, Community-Based and -Driven Development: A Critical Review, *The World Bank Research Observer* 19(1), 1-39.

Mark, M.M., G.T. Henry and G. Julnes (1999) "Toward an Integrative Framework for Evaluation Practice", American Journal of Evaluation 20, 177-198.

Mayne, J. (2001) "Addressing Attribution through Contribution Analysis: Using Performance Measures Sensibly", *Canadian Journal of Program Evaluation* 16(1), 1-24.

Mayntz, R. (2004) "Mechanisms in the Analysis of Social Macro-phenomena", *Philosophy of the Social Sciences* 34(2), 237-259.

McClintock, C. (1990) "Administrators as applied theorists", in: L. Bickman (ed.) Advances in program theory, New directions for evaluation 47, Jossey-Bass, San Francisco.

Mikkelsen, B. (2005) Methods for development work and research, Sage Publications, Thousand Oaks.

Mog, J.M. (2004) "Struggling with sustainability: A comparative framework for evaluating sustainable development programs", *World Development* 32(12), 2139–2160.

Morgan, S.L. and C. Winship (2007) Counterfactuals and causal inference – methods and principles for social research, Cambridge University Press, Cambridge.

Mukherjee, C., H. White and M. Wuyts (1998) Econometrics and Data Analysis for Developing Countries, Routledge, London.

North, D.C. (1990) Institutions, Institutional Change and Economic Performance, Cambridge University Press, New York.

Oakley, A. (2000) Experiments in knowing: Gender and method in the social sciences, Polity Press, Cambridge.

OECD-DAC (2000) Effective Practices in Conducting a Multi-donor Evaluation, OECD-DAC, Paris.

OECD-DAC (2002) Glossary of Key Terms in Evaluation and Results Based Management, OECD-DAC, Paris.

OED (2005) OED and impact evaluation: A discussion note, Operations Evaluation Department, World Bank, Washington D.C.

Oliver, S., A. Harden, R. Rees, J. Shepherd, G. Brunton, J. Garcia and A. Oakley (2005) "An Emerging Framework for Including Different Types of Evidence in Systematic Reviews for Public Policy", *Evaluation* 11(4), 428-446.

Patton, M.Q. (2002), *Qualitative research & evaluation methods*, Sage Publications, Thousand Oaks.

Pawson, R. (2002) "Evidence-based Policy: The Promise of 'Realist Synthesis", *Evaluation* 8(3), 340-358.

Pawson, R. (2005) "Simple Principles for The Evaluation of Complex Programmes", in: A. Killoran, M. Kelly, C. Swann, L. Taylor, L. Milward and S. Ellis (eds.) *Evidence Based Public Health*, Oxford University Press, Oxford.

Pawson, R. (2006) Evidence-based policy: A realist perspective, Sage Publications, London.

Pawson, R. and N. Tilley (1997) *Realistic Evaluation*, Sage Publications, Thousand Oaks.

Perloff, R. (2003) "A potpourri of cursory thoughts on evaluation", Industrial-Organizational Psychologist 40(3), 52-54.

Picciotto, R. (2004) "The value of evaluation standards: a comparative assessment" Paper presented at the European Evaluation Society's 6<sup>th</sup> biennial Conference on Democracy and Evaluation, Berlin.

Picciotto, R. and E. Wiesner (eds.) (1997) *Evaluation and development: the institutional dimension*, World Bank Series on Evaluation and Development, Transaction Publishers, New Brunswick.

Pollitt, C. (1999) "Stunted by stakeholders? Limits to collaborative evaluation", *Public Policy and Administration* 14 (2), 77-90.

Ravallion, M. (2008) "Evaluation in the practice of development", Policy Research Working Paper, 4547, World Bank, Washington D.C.

Rieper, O., F.L. Leeuw and T. Ling (eds.) (2009) The Evidence Book: concepts, generation and use of evidence, Transaction Publishers, New Brunswick.

Rist, R. and N. Stame (eds.) (2006) From Studies to Streams – Managing Evaluative Systems, Transaction Publishers, New Brunswick.

Robillard, A.S., F. Bourguignon and S. Robinson (2001) "Crisis and Income Distribution: A Micro-Macro Model for Indonesia", International Food Policy Research Institute, Washington D.C.

Roche, C. (1999) Impact Assessment for Development Agencies: Learning to Value Change, Oxfam, Oxford.

Rogers, P. J. (2008) "Using programme theory for complex and complicated programs", *Evaluation* 14(1), 29-48.

Rogers, P.J., Hacsi, T.A., Petrosino, A., and Huebner, T.A., (Eds.) (2000) Program theory in evaluation: Challenges and opportunities, New directions for evaluation 87, Jossey-Bass, San Francisco.

Rosenbaum, P.R. and D.B. Rubin (1983) "The central role of the propensity score in observational studies for causal effects", *Biometrika* 70, 41-55.

Rossi, P.H., Lipsey, M.W., and Freeman, H. E. (2004) *Evaluation:* A systematic *approach*, Sage Publications, Thousand Oaks.

Salamon, L. (1981) "Rethinking public management: Third party government and the changing forms of government action", *Public Policy* 29(3), 255–275.

Salmen, L. and E. Kane (2006) Bridging Diversity: Participatory Learning for Responsive Development, World Bank, Washington D.C.

Scriven, M. (1976) "Maximizing the Power of Causal Investigations: The Modus Operandi Method", in: G. V. Glass (ed.) *Evaluation Studies Review Annual*, Vol. 1, Sage Publications, Beverly Hills.

Scriven, M. (1998) "Minimalist theory: The least theory that practice requires", American Journal of Evaluation 19(1), 57-70.

Scriven, M. (2008) "Summative Evaluation of RCT Methodology & An Alternative Approach to Causal Research", *Journal of Multidisciplinary Evaluation* 5(9), 11-24.

Shadish, W. R., T.D. Cook and D.T. Campbell (2002) Experimental and quasiexperimental designs for generalized causal inference, Houghton Mifflin Company, Boston.

Sherman, L.W., D.C. Gottfredson, D.L. MacKenzie, J. Eck, P. Reuter and S.D. Bushway (1998) "Preventing Crime: what works, what doesn't, what's promising", *National Institute of Justice Research Brief*, July 1998, Washington D.C.

Snijders, T. and R. Bosker (1999) Multilevel Analysis: An introduction to basic and advanced multilevel modeling, Sage Publications, London.

Straw, R.B. and J.M. Herrell (2002) "A Framework for Understanding and improving Multisite Evaluations", in: J.M. Herrell and R.B. Straw (eds.), Conducting Multiple Site Evaluations in Real-World Settings, New Directions for Evaluation 94, Jossey-Bass, San Francisco.

Swedberg, R. (2005) Principles of Economic Sociology, Princeton University Press, Princeton.

Tashakkori, A., and C. Teddlie (eds.) (2003) Handbook of mixed methods in social and behavioral research, Sage Publications, Thousand Oaks.

Trochim, W.M.K. (1989) "An introduction to concept mapping for planning and evaluation", *Evaluation and Program Planning* 12, 1-16.

Tukey, J.W. (1977) Exploratory Data Analysis, Addison-Wasley, Reading.

Turpin, R.S. and J.M. Sinacore (eds.) (1991) *Multisite Evaluations*. New Directions for Evaluation 50, Jossey-Bass, San Francisco.

Utce Ltd & Japan Pfi Association (2003) Impact Evaluation Study on Public-Private Partnerships: the case of Angat Water Supply Optimization Project and the Metropolitan Waterworks and Sewerage System, Republic of the Philippines.

Vaessen, J. and J. De Groot (2004) "Evaluating Training Projects on Low External Input Agriculture: Lessons from Guatemala", *Agricultural Research & Extension Network Papers* 139, Overseas Development Institute, London.

Vaessen, J. and D. Todd (2008) "Methodological challenges of evaluating the impact of the Global Environment Facility's biodiversity program", *Evaluation and Program Planning* 31(3), 231-240.

Van der Knaap, L.M., F.L. Leeuw, S. Bogaerts and L.T.J. Nijssen (2008) "Combining Campbell standards and the realist evaluation approach – the best of two worlds?", American Journal of Evaluation 29(1), 48-57.

Van De Walle, D. and D. Cratty (2005) "Do Donors Get What They Paid For? Micro Evidence on the Fungibility of Development Project Aid", *World Bank Policy Research Working Paper* 3542, World Bank, Washington D.C.

Vedung, E. (1998) "Policy instruments: Typologies and theories", In: M. L. Bemelmans- Videc, and R. C. Rist (Eds.), *Carrots, sticks and sermons: Policy instruments and their evaluation*, Transaction Publishers, New Brunswick.

Webb, E.J., D.T. Campbell, R.D. Schwartz and L. Sechrest (2000) *Unobtrusive measures,* Sage Publications, Thousand Oaks.

Weiss, C.H. (1998) Evaluation – Methods for Studying Programs and Policies, Prentice Hall, New Jersey.

Welsh, B. and D.P. Farrington (eds.) (2006) Preventing crime: What works for children, offenders, victims and places, Springer, Berlin.

White, H. (2002) "Combining quantitative and qualitative approaches in poverty analysis", World Development 30(3), 511-522.

White, H. (2006) Impact Evaluation Experience of the Independent Evaluation Group of the World Bank, World Bank, Washington D.C.

White, H. and G. Dijkstra (2003), Programme aid and development, beyond conditionality, Routledge, London.

Wholey, J.S. (1987) "Evaluability Assessment: Developing Program Theory", in: L. Bickman (ed.) Using Program Theory in Evaluation, New Directions for Program Evaluation 33, Jossey-Bass, San Francisco.

Wooldridge, J.M. (2002), Econometric analysis of cross section and panel data, The MIT Press, Cambridge.

World Bank (2003) A User's Guide to Poverty and Social Impact Analysis, Poverty Reduction Group and Social Development Department, World Bank, Washington D.C.

Worrall, J. (2002) "What evidence in Evidence-based medicine?", Philosophy of science, 69, 316-330.

Worrall, J. (2007) "Why there's no cause to randomize", The British Journal for the Philosophy of Science 58(3), 451-488.

Worthen, B.R. and C.C. Schmitz (1997) "Conceptual Challenges Confronting Cluster Evaluation." *Evaluation* 3(3), 300-319.

Yang, H., J. Shen, H. Cao and C. Warfield (2004) "Multilevel Evaluation Alignment: An Explication of a Four-Step Model", *American Journal of Evaluation* 25(4), 493-507.
## Appendices

#### Appendix 1. Diversity in impact evaluation, some examples

## Example 1. Evaluating the impact of an EU-funded training project on Low External Input Agriculture (LEIA) in Guatemala

Within the Framework of an EU-funded integrated Rural Development Project, financial support was provided to a training project aimed the promotion of Low External Input Agriculture (LEIA) as a viable agricultural livelihood approach for small farmers in the highlands of Western Guatemala.

The impact evaluation design of this project was based on a quasi-experimental design and complemented by qualitative methods of data collection (Vaessen and De Groot, 2004). An intervention theory was reconstructed on the basis of field observations and relevant literature to make explicit the different causal assumptions of the project, facilitating further data collection and analysis. The quasi-experimental design included data collection on the ex ante and ex post situation of participants complemented with ex post data collection involving a control group (based on judgmental matching using descriptive statistical techniques). Without complex matching procedures and with limited statistical power, the strength of the quasi-experiment relied heavily on additional qualitative information. This shift in emphasis should not give the impression of a lack of rigor. Problems such as the influence of selection bias were explicitly addressed, even if not done in a formal statistical way.

The evaluation's findings include the following. Farmers' adoption behavior after the termination of the project can be characterized as selective and partial. Given the particular circumstances of small farmers (e.g. risk aversion, high opportunity costs of labor) it is not realistic to assume that a training project will bring about a complete transformation from a conventional farming system to a LEIA farming system (as assumed in the objectives). In line with the literature, the most popular practices (in this case for example organic fertilizers and medicinal plants) were those that offer a clear short term return while not requiring significant investments in terms of labor or capital. Finally, the evaluation argued that an ideological faith in the absolute supremacy of LEIA practices is not in the best interest of the farmer. Projects promoting LEIA should focus on the complementary effects of LEIA practices and conventional farming techniques, encouraging each farmer to choose the best balance fitted to his/her needs.

## Example 2. Assessing the impact of Swedish program aid

White and Dijkstra (2003) analyzed the impact of Swedish program aid. Their analysis accepted from the start that it is impossible to separate the impact of Swedish money from that of other donors' money. Therefore, the analysis focuses on all program aid with nine (country) case studies which trace how program aid has affected macro-economic aggregates (like imports and government spending) and (through these indicators) economic growth. The authors discern two channels for influencing policy: money and policy dialogue. The main evaluation questions are: 1. How has the policy dialogue affected the pattern and pace of reform (and what has been the contribution of program aid to this process)?

2. What has been the impact of the program aid funds (on imports, government expenditure, investment etc)?

3. What has been the impact of reform programs?

Their analytical model treats donor funds and the policy dialogue as inputs, specific economic, social and political indicators as outputs and the main program objectives (like economic growth, democracy, human rights and gender equality) as outcomes and poverty reduction as the overall goal.

The analysis focuses on marginal impact and uses a combination of quantitative and qualitative approaches (interviews, questionnaires and e-mail enquiries). The analysis of the impact of aid is largely quantitative, while the analysis of the impact of the policy dialogue is mainly qualitative.

An accounting approach is used to identify aid impact on expenditure levels and patterns using a number of ad hoc techniques, such as analyzing behavior during surges and before versus after breaks in key series, searching the data for other explanations of the patterns observed.

Moreover the authors analyze the impact of aid on stabilization through:

- a) the effect on imports;
- b) its impact on the markets for domestic currency and foreign exchange;
- c) the reduction of inflationary financing of the government deficit.

In terms of the impact of program aid on reform, they conclude that domestic political considerations are a key factor in determining reform: most countries have initiated reform without the help from donors, and have carried out some measure of reform not required by them, while ignoring others that gave been required.

#### Appendix 2. The General Elimination Methodology as a basis for causal analysis

This line of work consists of two steps. The first is to develop the counterfactual by using one or more well-corroborated theories from which the 'what would have happened-situation' can be *deduced*. The second step is to apply '*General Elimination Methodology*' to try to falsify this '(theoretical) prediction'. If the (theoretical) counterfactual survives the elimination / falsification efforts, the more one can rely on it.

An example is the following. Suppose one wants to study the impact of an anticorruption intervention implemented in five poor countries. There are believed to be two crucial ingredients. The first is that high-level politicians, civil servants, and respected members of the civil society publicly sign pledges indicating that they themselves are not involved in any fraud or corruption and have never been. Of the signing of the pledges press releases are issued. The intervention theory is that signing a pledge in public has serious reputation costs for those that sign the pledge, but that later turn out to be involved in corruption activities. The other part of this 'theory' is that officials do not want a negative reputation and therefore only sign the pledges if and when it is true what they sign. Implicitly, a naming and shaming-mechanism is believed to be at work, leading to a situation where all 'public pledges' can be assumed to be free of corrupt behavior. The second crucial ingredient is the diffusion mechanism, i.e. that information and communication about the signing of the pledges diffuses to society at large and prompts, as a signal, other people to also act in a non-corrupt way. Assume that data on court cases dealing with corruption are collected over a period of 5 year, as are data on public prosecution. The findings indicate that corruption has gone down, which makes the question crucial if and to what extent the pledge-policy caused 'the' difference.

A theoretically informed counterfactual now can be developed as follows. The first question is what is known about the symbolic and behavioral impact of signing pledges in public? Secondly, what is known in the social and behavioral sciences about the effects of and the condition under which naming and shaming makes an impact? And third, what is known from communication and diffusion studies regarding the diffusion of publicly uttered statements in general and pledges in particular. What is the current knowledge fund on the conditions under which naming and shaming, diffusion and behavioral adaptation (= i.e. refrain from corrupt behavior) can and will take place and how that knowledge fund is related to the specific situation in the (five) poor countries. If the answer is that, in the respective countries, these conditions are not fulfilled, then the assumption that the reduction in numbers of prosecution and court cases on corruption can be attributed to the 'pledge policy', is weak. This implies that the counterfactual probably is that not much would have been different, had there been *no* pledge policy. If the answer is that these conditions indeed are fulfilled, then a more 'positive' answer is possible.

So far for step 1. Step 2 is applying the 'General Elimination Methodology'. Although General Elimination Methodology epistemologically goes back to Popper's falsification principle, it was Scriven (1976; 2008) to add to the methodology of (impact) evaluations this approach. Using GEM involves the evaluator in setting up a 'competition' between the hypothesis that pledges matter and possible other hypotheses on what causes the changes in the dependent variable(s). If GEM is applied, then the hypothesis that the intervention is the (major) cause of the changes measured, will be tested extra. The stronger the tests, the larger the probability that, when surviving, the intervention hypothesis is indeed (probably) *the* or one of the 'real' causes of the changes in the dependent variable(s).

## Appendix 3. Overview of quantitative techniques of impact evaluation

|   |   | ANALYSIS OF INTERVENTION(S) |                                     |
|---|---|-----------------------------|-------------------------------------|
|   |   | EXPLICIT COUNTERFACTUAL     | ANALYSIS OF                         |
|   |   | (WITH / WITHOUT)            | MULTIPLE INTERVENTIONS AND          |
|   |   |                             | INFLUENCES                          |
| S | 0 |                             |                                     |
| E | В |                             |                                     |
| L | S | Propensity Score            | Regression analysis                 |
| E | E |                             |                                     |
| С | R |                             |                                     |
| Т | V |                             |                                     |
| 1 | E |                             |                                     |
| 0 | D |                             |                                     |
| Ν |   |                             |                                     |
|   | U | Randomized Controlled Trial |                                     |
| E | Ν | Pipeline Approach           |                                     |
| F | 0 |                             |                                     |
| F | В | Double difference           | Difference in difference regression |
| E | S | (Difference in difference)  | Fixed effects regression            |
| C | E |                             |                                     |
| Т | R |                             | Instrumental variables              |
| S | V | Regression Discontinuity    |                                     |
|   | E |                             |                                     |
|   | D |                             |                                     |

#### Appendix 4. Technical aspects of quantitative impact evaluation techniques

#### Endogeneity

The selection on unobservables is an important cause of *endogeneity*, a correlation of one of the explanatory variables with the error term in a mathematical model. This correlation occurs when an omitted variable has an effect at the same time on the dependent variable and an explanatory variable<sup>36</sup>.



When a third variable is not included in the model, the effect of the variable becomes part of the error term and contributes to the 'unexplained variance'. As long as this variable does not have an effect at the same time on one of the explanatory variables in the model, this does not lead to biased estimates. However, when this third variable has an effect on one of the explanatory variables, this explanatory variable will 'pick up' part of the error and therefore will be correlated with the error. In that case, omission of the third variable leads to a biased estimate.

Suppose we have the relation:

 $Y_i = a + bP_i + cX_i + e_i$ ,

where  $Y_i$  is the effect,  $P_i$  is the programme or intervention,  $X_i$  is an unobserved variable, and  $e_i$  is the error term. Ignoring X we try to estimate the equation:

 $Y_i = a + bP_i + e_i$ , while in effect we have:

 $Y_i = a + bP_i + (e_i + e_x),$ 

where  $e_i$  is a random error term and  $e_x$  is the effect of the unobserved variable. P and  $e_x$  are correlated and therefore P is *endogenous*. Ignoring this correlation results in a biased estimate of b. When the source of the selection bias (X) is known, inclusion of this variable (or these variables) leads to an unbiased estimate of the effect:

 $Y_i = a + bP_i + cX_i + e_i$ ,

<sup>&</sup>lt;sup>36</sup> In traditional usage, a variable is endogenous if it is determined within the context of a model. In econometrics, it is used to describe any situation where an explanatory variable is correlated with the disturbance term. Endogeneity arises as a result of omitted variables, measurement error or in situations where one of the explanatory variables is determined along with the dependent variable (Wooldridge 2002: 50).

An example is the effect of class size on learning achievements. The school choice of motivated (and probably well-educated) parents is probably correlated with class size, as these parents tend to send their children to schools with low pupil teacher ratios. The neglect of the *endogeneity* of class size may lead to biased estimates (with an overestimation of the real effect of class size). When the selection effects are observable, a regression-based approach may be used to get an unbiased estimate of the effects.

Figure 1 gives the relation between class size and learning achievements for two groups of schools: the left side of the figure shows private schools in urban areas with pupils with relatively rich and well educated parents; the left side shows public schools with pupils from poor remote rural areas. A neglect of the differences between the two schools leads to a biased estimate as shown by the black line. Including these effects in the equation leads to the smaller effect of the dotted lines.





#### Double difference and regression analysis

The technique of 'double differencing' can also be applied in a regression analysis. Suppose that the anticipated effect (Y) is a function of participation in the project (P) and of a vector of background characteristics. In a regression equation we may estimate the effect as:

$$Y_i = a + bP_i + cX_i + e_i$$

Where e is the error term and a, b en c the parameters to be estimated. When we analyse changes over time, we get (taking the *first differences* of the variables in the model):

$$(Y_{i,1} - Y_{i,0}) = a + b(P_{i,1} - P_{i,0}) + c(X_{i,1} - X_{i,0}) + e_i$$

When the (unobserved) variables X are time invariant,  $(X_{i,1} - X_{i,0})=0$ , and these variables drop from the equation. Suppose, for instance that a variable X denotes the 'year of birth'. For every individual the year of birth in year 1 = year of birth in year and therefore  $(X_{i,1} - X_{i,0})=0$ . So, if we expect that the year of birth is correlated with the probability of being included in the programme and with the anticipated effect of the programme, but we have no data on the year of birth, we may get an unbiased estimate through taking the first differences of the original variables. This technique helps to get rid of the problem of 'unobservables'<sup>37</sup>.

#### Instrumental variables

The use of instrumental variables is another technique to get rid of the endogeneity problem. A good instrument correlates with the (endogenous) intervention, but not with the error term. This instrument is used to get an unbiased estimate of the effect of the endogenous variable.

In practice researchers often use the method of two stage least squares: in the first stage an exogenous variable (Z) is used to give an estimate of the endogenous intervention-variable (P):

 $P_{i}' = a + dZ_{i} + e_{i}$ 

In the second stage this new variable is used to get an unbiased estimate of the effect of the intervention:

 $Y_i = a + bP'_i + cX_i + e_i$ .

#### The computation of propensity scores

The method of *propensity score matching* involves forming pairs by matching on the *probability* that subjects have been part of the treatment group. The method uses all *available* information in order to construct a control group. A standard way to do this is using a *probit* or *logit* regression model. In a logit specification, we get:

 $\ln (p_i / (1-p_i)) = a + bX_i + cY_i + dZ_i + e_i,$ 

where p<sub>i</sub> is the probability of being included in the intervention group and X, Y and Z denote specific *observed* characteristics. In this model, the probability is a function of the observed characteristics. Rosenbaum and Rubin proved that when subjects in the control group have the same probability of being included in the treatment group as subjects who actually belong to the treatment group, treatment group and control group will have similar characteristics.

<sup>&</sup>lt;sup>37</sup> The approach is similar to a fixed effects regression model, using deviations from individual means.

## Appendix 5. Evaluations using quantitative impact evaluation approaches (White, 2006)<sup>38</sup>

## AGRICULTURE AND RURAL DEVELOPMENT

#### Case Study: Pakistan

The projects: Irrigation in Pakistan suffers from the "twin menaces" of salinity and waterlogging. These problems have been tackled through Salinity Control and Reclamation Projects (SCARPs), financed in part by the Bank. Whilst technically successful, SCARP tubewells imposed an unsustainable burden on the government's budget. The project was to address this problem in areas with plentiful groundwater by closing public tubewells and subsidizing farmers to construct their own wells.

*Methodology:* IEG commissioned a survey in 1994 to create a panel from two earlier surveys undertaken in 1989 and 1990. The survey covered 391 farmers in project areas and 100 from comparison areas. Single and double differences of group means are reported.

*Findings:* The success of the project was that the public tubewells were closed without the public protests that had been expected. Coverage of private tubewells grew rapidly. However, private tubewells grew even more rapidly in the control area. This growth may be a case of contagion, though a demonstration effect. But it seems more likely that other factors (e.g. availability of cheaper tubewell technology) were behind the rapid diffusion of private water exploitation. Hence the project did not have any impact on agricultural productivity or incomes. It did however have a positive rate of return by virtue of the savings in government revenue.

## Case study: Philippines

The project: the Second Rural Credit Projects (SRCP) operated between 1969 and 1974 with a US\$12.5 million loan from the World Bank. SRCP was the continuation of a pilot credit project started in 1965 and completed in 1969. As its successful predecessor, SRCP aimed at providing credit to small and medium rice and sugar farmers for the purchase of farm machinery, power tillers, and irrigation equipment. Credits were to be channeled through 250 rural banks scattered around the country. An average financial contribution to the project of 10% was required from both rural banks and farmers. The SRCP was followed by a third loan of US\$22.0 million from 1975-77, and by a fourth loan of US\$36.5 million that was still in operation at the time of the evaluation (1983).

*Methodology:* the study uses data of a survey of 738 borrowers (nearly 20% of total project beneficiaries) from seven provinces of the country. Data were collected through household questionnaires on land, production, employment and measures of standard of living. In addition, 47 banks were surveyed in order to measure the impact on their profitability, liquidity, and solvency. The study uses before-after comparisons of means and ratios to assess the project impact on farmers. National level data are often used to validate the effects observed. Regarding the rural banks, the study compares measures of financial performance before and after the project taking advantage of the fact that the banks surveyed joined the project at different stages.

*Findings:* the mechanization of farming did not produce an expansion of holding sizes (though the effect of a contemporaneous land reform should be taken into account). Mechanization did not change cropping patterns, and most farmers were concentrating on a single crop at the time of the interviews. No change in cropping intensity was observed, but production and productivity were found to be higher at the end of the project. The project increased the demand for both family and hired labor. Farmers reported an increase in incomes and savings, and in several other welfare indicators, as a result of the project. Regarding the project impact on rural banks, the study observes an increase in the net income of the sample banks from 1969 to 1975, and a decline

 $<sup>^{38}</sup>$  For further examples see White (2006).

thereafter. Banks' liquidity and solvency position was negatively affected by poor collection and loan arrears.

#### HEALTH, NUTRITION AND POPULATION

#### **Case Study: India**

The project: The Tamil Nadu Integrated Nutrition Project (TINP) operated between 1980 and 1989, with a credit of US\$32 million from IDA. The overall objective of the project was to improve the nutritional and health status of pre-school children, pregnant women and nursing mothers. The intervention consisted of a package of services including: nutrition education, primary health care, supplementary feeding, administration of vitamin A, and periodic de-worming. The project was the first to employ Growth Monitoring and Promotion (GMP) on a large scale. The evaluation is concerned with the impact of the project on nutritional status of children.

Methodology: The study uses three cross-sectional rounds of data collected by the TINP Monitoring Office. Child and household characteristics of children participating in the program were collected in 1982, 1986, and 1990, each round consisting of between 1000 and 1500 observations. The study uses before-after comparisons of means, regression analysis, and charts to provide evidence of the following: frequency of project participation, improvement in nutritional status of participating children over time, differential participation and differential project impact across social groups. Data on the change in nutritional status in project areas are compared to secondary data on the nutritional status of children outside the project areas. With some assumptions, the use of secondary data, make the findings plausible.

*Findings*: The study concludes that the implementation of Growth Monitoring and Promotion programs on a large scale is feasible, and that this had a positive impact on nutritional status of children of Tamil Nadu. More specifically, these are the findings of the study:

- Program participation: Among children participating in GMP, all service delivery indicators (age at enrolment, regular attendance of sessions, administration of vitamin A, and deworming), show a substantial increase between 1982 and 1986, though subsequently declined to around their initial levels. Levels of service delivery, however, are generally high.
- Nutritional status: mean weight and malnutrition rates of children aged between 6 and 36 months and participating in GMP have improved over time. Data on non-project areas in Tamil Nadu, and all-India data, that show a smaller improvement over the same time period. Regression analysis of nutritional status on a set of explanatory variables, including the participation in a cotemporaneous nutrition project (the National Meal Program) shows that the latter had no additional benefit on nutritional outcomes. Positive associations are also found between nutritional status and intensive participation in the program, and complete immunization.
- Targeting: using tabulations and regression analysis, it is shown that initially girls have benefited more from the program, but that at the end of the program boys have benefited more. Children from scheduled caste are shown to have benefited more than other groups. Nutritional status was observed to be improving at all income levels, the highest income category benefiting slightly more than the lowest.



Decision tree for IE design using quantitative IE techniques (continued)

- 1. If the evaluation is being designed before the intervention (ex-ante), is randomization possible? If the treatment group is chosen at random then a random sample drawn from the sample population is a valid control group, and will remain so provided they are outside the influence zone and contamination is avoided. This approach does not mean that targeting specific analytical units is not possible. The random allocation may be to a subgroup of the total population, e.g. from the poorest districts.
- 2. If randomisation is not possible, are all selection determinants observed? If they are, then there are a number of regression-based approaches which can remove the selection bias.
- 3. If the selection determinants are unobserved and if they are thought to be time invariant then using panel data will remove their influence, so a baseline is essential (or some means of substituting for a baseline).
- 4. If the study is done afterwards (ex post) so it is not possible to get information for exactly the same units (a panel of persons, households, etcetera) and selection is determined by unobservables, then some means of observing the supposed unobservables should be sought. If that is not possible, then a pipeline approach can be used if there are as yet untreated beneficiaries. For example, ADB's impact study of microfinance in the Philippines matched treatment areas with areas which were in the program but had not yet received the intervention.
- 5. If none of the above mentioned procedures is possible, then the problem of selection bias cannot be addressed. The impact evaluation will have to rely heavily on the intervention theory and triangulation to build an argument by plausible association.

### Appendix 7. Hierarchical modeling and other statistical approaches

This group of approaches covers a quite diverse set of quite advanced modeling and statistical approaches. Detailed discussion of these technical features is beyond the scope of this Guidance. The common element that binds these approaches is the purpose modeling and estimating direct and indirect effects of interventions at various levels of aggregation (from micro to macro). At the risk of substantial oversimplification we briefly mention a few of the approaches. In hierarchical modeling, evaluators and researchers look at the interrelationships between different levels of a program. The goal is "to measure the true and often intertwined effects of the program. In a typical hierarchical linear model analysis, for example, the emphasis is on how to model the effect of variables at one level on the relations occurring at another level. Such analyses often attempt to decompose the total effect of the program into the effect across various program levels and that between program sites within a level (Dehejia, 1999)" (Yang et al., 2004: 494).

Also part of this branch of approaches is a range of statistical approaches such as nested models, models with latent variables, multi-level regression approaches and others (see for example, Snijders and Bosker (1999). Other examples are typical economist tools such as partial equilibrium analyses, general computable equilibrium models are often used to assess the impact of (e.g.) macroeconomic policies on markets and subsequently on household welfare (see Box 1).

## Box 1. Impact of the Indonesian Financial Crisis on the Poor: Partial Equilibrium Modeling and CGE Modeling with Micro-simulation

General equilibrium models permit the analyst to examine explicitly the indirect and second-round consequences of policy changes. These indirect consequences are often larger than the direct, immediate impact, and may have different distributional implications. General equilibrium models and partial equilibrium models may thus lead to significantly different conclusions. A comparison of conclusions reached by two sets of researchers, examining the same event using different methods, reveals the differences between the models. Levinsohn et al. (1999) and Robillard et al. (2001) both look at the impact of the Indonesian financial crisis on the poor—the former using partial equilibrium methods, the latter using a CGE model with microsimulation. The Levinsohn study used consumption data for nearly 60,000 households from the 1993 SUSENAS survey, together with detailed information on price changes over the 1997–98 crisis period, to compute household-specific cost-of-living changes. It finds that the poorest urban households were hit hardest by the shock, experiencing a 10–30 percent increase in the cost of living (depending on the method used to calculate the change). Rural households and wealthy urban households actually saw the cost of living fall.

These results suggest that the poor are just as integrated into the economy as other classes, but have fewer opportunities to smooth consumption during a crisis. However, the methods used have at least three serious drawbacks. First, the consumption parameters are fixed, that is, no substitution is permitted between more expensive and less expensive consumption items. Second, the results are exclusively *nominal*, in that the welfare changes are due entirely to changes in the price of consumption, and do not account for any concomitant change in income. Third, this analysis cannot control for other exogenous events, such as the El Niño drought and resulting widespread forest fires.

Robillard, Bourguignon, and Robinson use a CGE model, connected to a micro-simulation model. The results are obtained in two steps. First, the CGE is run to derive a set of parameters for prices, wages, and labor demand. These results are fed into a micro-simulation model to estimate the effects on each of 10,000 households in the 1996 SUSENAS survey. In the microsimulation model, workers are divided into groups according to sex, residence, and skill. Individuals earn factor income from wage labor and enterprise profits, and households accrue profits and income to factors in proportion to their endowments. Labor supply is endogenous. The micro-simulation model is constrained to conform to the aggregate levels provided by the CGE model. The Robillard team finds that poverty did increase during the crisis, although not as severely as the previous results suggest. Also, the increase in poverty was due in equal parts to the crisis and to the drought. Comparing their micro-simulation results to those produced by the CGE alone, the authors find that the representative household model is likely to *underestimate* the impact of shocks on poverty. On the other hand, ignoring both substitution and income effects, as Levinsohn, Berry, and Friedman do, is likely to lead to

overestimating the increase in poverty, since it does not permit the household to reallocate resources in response to the shock.

Source: literal citation from World Bank (2003)

#### Appendix 8. Multi-site evaluation approaches

Multi-site evaluation approaches involve primary data collection processes and analyses at multiple sites or interventions. They usually focus on programs encompassing multiple interventions implemented in different sites (Turpin and Sinacore, 1991; Straw and Herrell, 2002). Although often referred to as a family of methodologies, in what follows, and in line with the literature, we will use a somewhat more narrow definition of multi-site evaluations alongside several specific methodologies to address the issue of aggregation and cross-site evaluation of multiple interventions.

Straw and Herrell (2002) use the term multi-site evaluation both as an overarching concept, i.e. including cluster evaluation and multi-center clinical trials, as well as a particular type of multi-level evaluation distinguishable from cluster evaluation and multi-center clinical trials. Here we use the latter definition, the term multi-site evaluation referring to a particular (though rather flexible) methodological framework applicable to the evaluation of comprehensive multilevel programs addressing health, economic, environmental or social issues.

The *multi-center clinical trial* is a methodology in which empirical data collection in a selection of homogenous intervention sites is systematically organized and coordinated. Basically it consists of a series of randomized controlled trials. The latter are experimental evaluations in which treatment is randomly assigned to a target group while a similar group not receiving the treatment is used as a control group. Consequently, changes in impact variables between the two groups can be traced back to the treatment as all other variables are assumed to be similar at group level. In the multi-center clinical trial sample size is increased and multiple sites are included in the experiment in order to strengthen the external validity of the findings. Control over all aspects of the evaluation is very tight in order to keep as many variables constant over the different sites. Applications are mostly found in the health sector (see Kraemer, 2000).

*Multi-site evaluation* distinguishes itself from cluster evaluation in the sense that its primary purpose is summative, the assessment of merit and worth after the completion of the intervention. In addition, multi-site evaluations are less participatory in nature vis-à-vis intervention staff. In contrast to settings in which multi-center clinical trials are applied, multi-site evaluations address large-scale programs which because of their (complex) underlying strategies, implementation issues or other reasons are not amenable to controlled experimental impact evaluation designs. Possible variations in implementation among interventions sites, and variations in terms of available data require a different more flexible approach to data collection and analysis than in the case of the multi-center clinical trials. A common framework of questions and indicators is established to counter this variability, enabling data analysis across interventions in function of establishing generalizable findings (Straw and Herrell, 2002).

*Cluster evaluation* is a methodology that is especially useful for evaluating large-scale interventions addressing complex societal themes such as education, social service delivery and health promotion. Within a cluster of projects under evaluation, implementation among interventions may vary widely but single interventions are still linked in terms of common strategies, target populations or problems that are addressed (Worthen and Schmitz, 1997).

The approach was developed by the Kellogg Foundation in the nineties and since then has been taken up by other institutions. Four elements characterize cluster evaluation (Kellogg Foundation, 1991):

- it focuses on a group of projects in order to identify common issues and patterns;
- it focuses on 'what happened' as well as 'why did it happen';
- it is based on a collaborative process involving all relevant actors, including evaluators and individual project staff;
- project-specific information is confidential and not reported to the higher level; evaluators only report aggregate findings; this type of confidentiality between evaluators and project staff induces a more open and collaborative environment.

Cluster evaluation is typically applied during program implementation (or during the planning stage) in close collaboration with stakeholders from all levels. Its purpose is, on the one hand, formative as evaluators in close collaboration with stakeholders at project level try to explore common issues as well as variations between sites. At the program level the evaluation's purpose can both be formative in terms of supporting planning processes as well as summative, i.e. judging what went wrong and why. A common question at the

program level would be for example to explore the factors that in the different sites are associated with positive impacts. In general, the objective of cluster evaluation is not so much to prove as to improve, based on a shared understanding of why things are happening the way they do (Worthen and Schmitz, 1997). It should be noted that not only cluster evaluation but also multi-site evaluation are applicable to homogenous programs with little variation in terms of implementation and context among single interventions.

## Appendix 9. Where to find reviews and synthesis studies that report about mechanisms underlying processes of change

#### Books on social mechanisms

Authors like Elster(1989; 2007), Farnsworth (2007), Hedström and Swedberg (1998), Swedberg (2005), Bunge (2004) and Mayntz (2004) have summarized and synthesized the research literature on different (types of) social mechanisms. Elster's 'Explaining social behaviour' (2007) summarizes insights from neurosciences to economics and political science and discusses 20-plus mechanisms. They range from 'motivations', 'emotions' and 'self interest' to 'rational choice, games and behavior and collective decisionmaking'.

Farnsworth (2007) takes legal arrangements like laws and contracts as a starting point and dissects which (types of) mechanisms play a role when one wants to understand why laws sometimes (do) (not) work. He combines insights from psychology, economics and sociology and discusses mechanisms such as the 'slippery slope', the endowment effect, framing effects and public goods production.

#### Review journals

Since the 1970s there review journals trying have been developed to address important developments within a discipline. An example is 'Annual Reviews', which publishes analytic reviews in 37 disciplines within the Biomedical, Life, Physical, and Social Sciences.

#### Knowledge repositories

Hansen and Rieper (2009) have inventoried a number of second-order evidence-producing organizations within the social (and behavioral) sciences. In recent years the production of systematic reviews has been institutionalized in these institutions. There are two main international organizations: the Cochrane Society working within the health field; and the Campbell Collaboration working within the fields of social welfare, education and criminology. Both organizations subscribe to the idea of producing globally valid knowledge about the effects of interventions, if possible through synthesizing the results of primary studies designed as RCTs and using meta analysis as the form of syntheses. In many (western) countries second-order knowledge-producing organizations have been established at the national level that not all are based on findings from RCTs. Hansen and Rieper (2009) present information about some 15 of them, including web addresses.

#### Knowledge repositories and development intervention impact

The report 'When will we ever learn' documents several 'knowledge repositories'. We quote from the report and refer to some of them (CGD, 2006: 58 and further).

The Coalition for Evidence-Based Policy offers "Social Programs That Work," a Web site providing policymakers and practitioners with clear, actionable information on what works in social policy, as demonstrated in scientifically valid studies. [www.evidencebasedprograms.org/].

The International Organization for Cooperation in Evaluation (IOCE), a loose alliance of regional and national evaluation organizations from around the world, builds evaluation leadership and capacity in developing countries, fosters the cross-fertilization of evaluation theory and practice around the world, addresses international challenges in evaluation, and assists the evaluation professionals to take a more global approach to identifying and solving problems. It offers links to other evaluation organizations; forums that network evaluators internationally; news of events and important initiatives; and opportunities to exchange ideas, practices, and insights with evaluation associations, societies, and networks. [http://ioce.net].

The Abdul Latif Jameel Poverty Action Lab (J-PAL) fights poverty by ensuring that policy decisions are based on scientific evidence. Located in the Economics Department at the Massachusetts Institute of Technology, J-PAL brings together a network of researchers at several universities who work on randomized evaluations. It works with governments, aid agencies, bilateral donors, and NGOs to evaluate the effectiveness of antipoverty programs using randomized evaluations, disseminate findings and policy implications, and promote the use of randomized evaluations, including by training practitioners to carry them out. [www.povertyactionlab.com/].

The Development Impact Evaluation Initiative (DIME). The DIME initiative is a World Bank-led effort involving thematic networks and regional units under the guidance of the Bank's Chief Economist. Its objectives are:

- To increase the number of Bank projects with impact evaluation components;
- To increase staff capacity to design and carry out such evaluations;
- To build a process of systematic learning based on effective development interventions with lessons learned from completed evaluations.

#### Appendix 10. Further information on review and synthesis approaches in IE

#### Realist synthesis

This approach is different from the systematic research reviews. It conceptualizes *interventions*, *programs and policies* as theories and it collects earlier research findings by interpreting the specific policy instrument that is evaluated, as an example or specimen of *more generic instruments and tools* (of governments). Next it describes the intervention in terms of its context, mechanisms (what makes the program work) and outcomes (the deliverables).

In stead of synthesizing results from evaluations and other studies *per intervention* or *per program*, realist evaluators first open up the black box of an interventions, and synthesize knowledge about social and behavioral mechanisms. Examples are Pawson's study of incentives (Pawson, 2002), on naming and shaming and Megan's law (Pawson, 2006) and Kruisbergen's (2005) on fear-arousal communication campaigns trying to reduce the smuggling of cocaine.

Contrary to producers of systematic research reviews, realist evaluators do not use a hierarchy of research designs. For realists an impact study using the RCT design is not necessarily better than a comparative case study design or a process evaluation. The problem (of an evaluation) that needs to be adressed is crucial in selecting the design or method and not vice versa.

#### Combining Different Meta Approaches

In a study on the question which public policy programs designed to reduce and/or prevent violence in the public arena work best, Van der Knaap et al. (2008) have shown the relevance of *combining* the *systematic research review* and the *realist synthesis*. Both perspectives have something to offer each other. Opening up the black box of an intervention under review will be helpful for experimental evaluators if they want to understand why interventions have (no) effects and/or side effects. Realists are confronted with the problem of the selection of studies to be taken into account; ranging from opinion surveys, oral history, and newspaper content analysis to results based on more sophisticated methodologies. As the methodological quality of evaluations can be and sometimes is a problem, particularly with regard to the measurement of the impact of a program, realists can benefit from a *stricter methodology and protocol,* like the one used by the Campbell Collaboration, when doing a synthesis. For, knowledge that is to be generalized should be credible and valid.

In order to combine Campbell standards and the Realist Evaluation approach Van der Knaap et al. (2008) first conducted a systematic review according to the Campbell standards. The research questions were formulated, and next the inclusion and exclusion criteria were determined. This included a number of questions. What types of interventions are included? Which participants should interventions be aimed at? What kinds of outcome data should be reported? At this stage, criteria were also formulated for inclusion and exclusion of study designs and methodological quality. As a third step, the search for potential studies was explicitly described. Once potentially relevant studies hade been identified, they were screened for eligibility according to the inclusion and exclusion criteria.

After selecting the relevant studies, the quality of these studies had to be determined. Van der Knaap et al (2008) used the Maryland Scientific Methods Scale (MSMS) (Sherman at al., 1998; Welsh and Farrington, 2006). This is a five-point scale that enables researchers to draw conclusions on methodological quality of outcome evaluations in terms of the internal validity. Using a scale of 1 to 5, the MSMS is applied to rate the strength of scientific evidence, with 1 being the weakest and 5 the strongest scientific evidence needed for inferring cause and effect.

Based on the MSMS-scores, the authors then classified each of the 36 interventions that were inventoried by analyzing some 450 English, German, French and Dutch written articles and papers, into the following categories: 1) effective, 2) potentially effective, 3) potentially ineffective, and 4) ineffective. However, not all studies could be grouped in one of the four categories: in sixteen cases the quality of the study design was not good enough to decide on the effectiveness of a measure. The (remaining) nine interventions were labeled effective and the (final) six were labeled potentially effective. Four interventions were labeled potentially ineffective and one was labeled ineffective in preventing violence in the public and semi-public domain.

In order to combine Campbell standards and the Realist Evaluation approach, the realist approach was applied *after finishing the Campbell-style systematic review*. This means that only then the underlying mechanisms and contexts as described in the studies included in the review were on the agenda of the evaluator. This was done for the four types of interventions, whether they were measured as being effective, potentially ineffective or ineffective. As a first step, information was collected concerning social and behavioral mechanisms that were assumed to be at work when the program or intervention was implemented. Pawson (2006: 24) refers to this process as "to look beneath the surface [of a program] in order to inspect how they work". One way of doing this is to search articles under review for statements that address the why-question: why will this intervention be working or why has it not worked? Two researchers independently articulated these underlying mechanisms. The focus was on behavioral and social 'cogs and wheels' of the intervention (Elster, 1989; 2007).

In a second step the studies under review were searched for information on *contexts* (schools, streets, banks etc., but also types of offenders and victims and type of crime) and *outcomes*. This completed the C[ontext], M[echanism] and O[utcome]- approach that characterizes realist evaluations. However, not every original evaluation study described which mechanisms are assumed to be at work when the program is implemented. The same goes for contexts and outcomes. This meant that in most cases missing links in or between different statements in the evaluation study had to be identified through *argumentational analysis*.

Based on the evaluations analyzed, Van der Knaap et al. (2008) traced the following three mechanisms to be at work in programs that had demonstrated their impact or the very-likely-to-come-impact:

- The first is of a cognitive nature, focusing on learning, teaching and training.
- The second (overarching) mechanism concerns the way in which the (social) environment is rewarding or punishing behavior (through bonding, community development and the targeting of police activities).
- The third mechanism is risk reduction, for instance by promoting protective factors.

## Concluding remarks on review and synthesis approaches

Given the 'fleets' (Weiss, 1998) and the 'streams of studies' (Rist and Stame, 2006) in the world of evaluation, it is not recommendable to start an impact evaluation of a specific program, intervention or 'tool of government' without making use of the accumulated evidence to be found in systematic reviews and other types of meta-studies. One reason concerns the efficiency of the investments: what has been sorted out, does not need (always) to be sorted out again. If over and over again it has been found that awareness-raising leads to behavior changes only under specific conditions, than it is wise to have that knowledge ready before designing a similar program or the evaluation. A second reason is that by using results from synthesis studies the test of an intervention theory can be done with more rigour. The larger the discrepancy between what is known about mechanisms a policy or programs believes to be 'at work' and what the policy in fact tries to set into motion, the smaller the chances of an effective intervention.

Different approaches in the world of (impact) evaluation are a wise thing to have, but (continuous) paradigm wars ('randomistas vs. relativistas'; realists versus experimentalists) run the risk of developing into intellectual ostracism. Wars also run the risk to vest an image of evaluations as a `helter-skelter mishmash [and] a stew of hit-or-miss procedures" (Perloff, 2003), which is not the best perspective to live with. Combining perspectives and paradigms should therefore be stimulated.

## Appendix 11. Evaluations based on qualitative and quantitative descriptive methods

# Case 1: Combining qualitative and quantitative descriptive methods - Ex-Post Impact Study of the Noakhali Rural Development Project in Bangladesh<sup>39</sup>

### 1. Summary

The evaluation examined the intended and unintended socio-economic impacts, with particular attention to the impact for women and to the sustainability and sustainment of these impacts. The evaluation drew on a wide range of existing evidence and also used mixed methods to generate additional evidence; because the evaluation was conducted 9 years after the project had ended, it was possible to directly investigate the extent to which impacts had been sustained. Careful attention was paid to differential impacts in different contexts in order to interpret the significance of before/after and with/without comparisons; the intervention was only successful in contexts which provided the other necessary 'ingredients' for success. The evaluation had significant resources and was preceded by considerable planning and review of existing evidence.

#### 2. Summary of the intervention; its main characteristics

The Noakhali Rural Development Project (NRDP) was an Integrated Rural Development Projects (IRDP) in Bangladesh, funded for DKK 389 million by Danida. It was implemented in two phases over a period of 14 years, 1978-92, in the greater Noakhali district, one of the poorest regions of Bangladesh, which had a population of approximately 4 million. More than 60 long-term expatriate advisers – most of them Danish – worked 2-3 years each on the project together with a Bangladeshi staff of up to 1,000 (at the peak). During NRDP-I the project comprised activities in 14 different areas grouped under four headings:

- Infrastructure (roads, canals, market places, public facilities);
- Agriculture (credit, cooperatives, irrigation, extension, marketing);
- Other productive activities (livestock, fish ponds, cottage industries);
- Social sector (health & family planning, education).

The overarching objective of NRDP-I was: to promote economic growth and social progress in particular aiming at the poorer sections of the population. The poorer sections were to be reached in particular through the creation of temporary employment in construction activities (infrastructure) and engaging them in income generating activities (other productive activities). There was also an aim to create more employment in agriculture for landless laborers through intensification. Almost all the major activities started under NRDP-I continued under NRDP-II, albeit with some modifications and additions. The overarching objective was kept with one notable addition: to promote economic growth and social progress in particular aiming at the poorer segments of the population including women". A special focus on women was thus included, based on the experience that so far most of the benefits of the project had accrued to men.

## 3. The purpose, intended use and key evaluation questions

This ex-post impact study of the Noakhali Rural Development Project (NRDP) was carried out nine years after the project was terminated. At the time of implementation NRDP was one of the largest projects funded by Danida, and it was considered an excellent example of integrated rural development, which was a common type of support during the seventies and eighties. In spite of the potential lessons to be learned from the project, it was not evaluated upon completion in 1992. This fact and an interest in the sustainability factor in Danish development assistance led to the commission of the study. What type of impact could still be traced in Noakhali nine years after Danida terminated its support to the project?

While the study dealt with aspects of the project implementation, its main focus was on the project's socioeconomic impact in the Noakhali region. The study aimed to identify the intended as well as unintended impact of the project, in particular whether it had stimulated economic growth and social development and improved the livelihoods of the poor, including women, such as the project had set out to do. The evaluation focused on the following questions:

- What has been the short- and long-term intended as well as unintended impact of the project?
- Has the project stimulated economic growth and social development in the area?

<sup>&</sup>lt;sup>39</sup> This case study is drawn from the 2002 report published by the Ministry of Foreign Affairs, Denmark.

- Has the project contributed to improving the livelihoods of the poorest section of the population, including women?
- Have the institutional and capacity-building activities engendered or reinforced by the project produced sustainable results?

## 4. Concise description of the evaluation, focusing on the approach, the rationale for the choice of approach and methods - linked to the four key tasks described in this document

#### Identifying impacts of interest

This study focuses on the impact of NRDP, in particular the long-term impact (i.e. nine years after). But impact cannot be understood in isolation from implementation and hence the study analyses various elements and problems in the way the project was designed and executed. Impact can also not be understood isolated from the context, both the natural/physical and in particular the societal – social, cultural, economic, political – context. In comparison with ordinary evaluations this study puts a lot more emphasis on understanding the national and in particular the local context.

#### Gathering evidence of impacts

One of the distinguishing features of this impact study, compared to normal evaluations, is the order and kind of fieldwork. The fieldwork lasted four months and involved a team of eight researchers (three European and five Bangladeshi) and 15 assistants. The researchers spent 11/2-31/2 months in the field, the assistants 2-4 months.

The following is a list of the methods used:

- Documentary study (project documents, research reports etc.)
- Archival work (in the Danish embassy, Dhaka)
- o Questionnaire survey with former advisers and Danida staff members
- o Stakeholder interviews (Danida staff, former advisers, Bangladeshi staff etc.)
- o Quantitative analysis of project monitoring data
- o Key informant interviews
- Compilation and analysis of material about context (statistics, articles, reports etc.)
- Institutional mapping (particularly NGOs in the area)
- Representative surveys of project components
- Assessment of buildings, roads and irrigation canals (function, maintenance etc.)
- o Questionnaire-based interviews with beneficiaries and non-beneficiaries
- Extensive and intensive village studies (surveys, interviews etc.)
- o Observation
- Focus group interviews
- o In-depth interviews (issue-based and life stories)

In the history of Danish development cooperation no other project has been subject to so many studies and reports, not to speak of the vast number of newspaper articles. Most important for the impact study have been the appraisal reports and the evaluations plus the final project completion report. But in addition to this there exists an enormous amount of reports on all aspects of the project. A catalogue from 1993 lists more than 1,500 reports produced by and for the NRDP. Both the project and the local context were, moreover, intensively studied in a research project carried out in cooperation between the Centre for Development Research (CDR) and Bangladesh Institute of Development Studies (BIDS).

A special effort was made to solicit the views of a number of key actors (or stakeholders) in the project and other key informants. This included numerous former NRDP and BRDB officers, expatriate former advisers as well as former key Danida staff, both based in the Danish Embassy in Dhaka and in the Ministry of Foreign Affairs in Copenhagen. They were asked about their views on strengths and weaknesses of the project and the components they know best, about their own involvement and about their judgment regarding likely impact. A questionnaire survey was carried out among the around 60 former expatriate long-term advisers and 25 former key staff members in the Danish embassy, Danida and other key informants. In both cases about half returned the filled-in questionnaires. This was followed up by a number of individual interviews.

The main method in four of the five component studies was surveys with interviews, based on standardized questionnaires, with a random – or at least reasonably representative – sample of beneficiaries (of course combined with documentary evidence, key informant interviews etc.). A great deal of effort was taken in ensuring that the survey samples are reasonably representative.

The infrastructure component was studied by partly different methods, since in this case the beneficiaries were less well defined. It was decided to make a survey of all the buildings that were constructed during the first phase of the project in order to assess their current use, maintenance standard and benefits. In this phase the emphasis was on construction; in the second phase it shifted to maintenance. Moreover, a number of roads were selected for study, both of their current maintenance standard, their use etc., but also the employment the road construction and maintenance generated, particularly for groups of destitute women. The study also attempted to assess socio-economic impact of the roads on different groups (poor/better-off, men/women etc.).

#### Assessing causal contribution

The impact of a development intervention is a result of the interplay of the intervention and the context. It is the matching of what the project has to offer and people's needs and capabilities that produces the outcome and impact. Moreover, the development processes engendered unfold in a setting, which is often characterized by inequalities, structural constraints and power relations. This certainly has been the case in Noakhali. As a consequence there will be differential impact, varying between individuals and according to gender, socio-economic group and political leverage.

In addition to the documentary studies, interviews and questionnaire survey the actual fieldwork has employed a range of both quantitative and qualitative methods. The approach can be characterized as a contextualized, tailor-made ex-post impact study. There is considerable emphasis on uncovering elements of the societal context in which the project has been implemented. This covers both the national context and the local context. The approach is tailor-made in the sense that it will be made to fit the study design outlined above and apply an appropriate mix of methods.

An element in the method is the incorporation in the study of both before/after and with/without perspectives. These, however, are not seen as the ultimate test of impact (success or failure), but interpreted cautiously, bearing in mind that the area's development has also been influenced by a range of other factors (market forces, changing government policies, other development interventions etc.), both during the 14 years the project was implemented and during the nine years that have lapsed since its termination.

Considerable weight was accorded to studying what has happened in the villages that have previously been studied and for which some comparable data exist. Four villages were studied intensively in 1979 and briefly restudied in 1988 and 1994. These studies – together with a thorough restudy in the year 2001 – provide a unique opportunity to compare the situation before, during and after the project. Moreover, 10 villages were monitored under the project's 'Village-wise Impact Monitoring System' in the years 1988-90, some of these being 'with' (+NRDP) and some (largely) 'without' (-NRDP) the project. Analysis of the monitoring data combined with a restudy of a sample of these villages illuminates the impact of the project in relation to other factors. It was decided to study a total of 16 villages, 3 intensively (all +NRDP, about 3 weeks each), 12 extensively (9 +NRDP, 3 – NRDP, 3-5 days each). As a matter of principle, this part of the study looks at impact in terms of the project as a whole. It brings in focus the project benefits as perceived by different groups and individuals and tries to study how the project has impinged on economic and social processes of development and change. At the same time it provides a picture of the considerable variety found in the local context.

In the evaluation of the Mass Education Program, the problem of attribution was dealt with as carefully as possible. Firstly, a parallel comparison has been made between the MEP beneficiaries on the one hand and non-beneficiaries on the other, in order to identify (if any) the changes directly or indirectly related to the program. Such comparison was vital due to the absence of any reliable and comparable baseline data. Secondly, specific queries were made in relation to the impact of the program as perceived by the beneficiaries and other stakeholders of MEP, assuming that they would be able to perceive the impact of the program intervention on their own lives in a way that would not be possible for others. And finally, views of non-beneficiaries and non-stakeholders were sought in order to have less affected opinion from people who

do not have any valid reason for either understating or overstating the impact of MEP. It was through such a cautious approach that the question of attribution was addressed. Arguably, elements of subjectivity may still have remained in the conclusions and assumptions, but that is unavoidable in a study that seeks to uncover the impact of an education project.

#### Managing the impact evaluation

The impact study was commissioned by Danida and carried out by Centre for Development Research, who also co-funded the study as a component of its Aid Impact Research Program. The research team comprised independent researchers from Bangladesh, Denmark and the UK. A reference group of nine persons (former advisers, Danida officers and researchers) followed the study from the beginning to the end. It discussed the approach paper in an initial meeting and the draft reports in a final meeting. In between it received three progress reports from the team leader and took up discussions by e-mail correspondence. The study was prepared during the year 2000 and fieldwork carried out in the period January-May 2001. The study consists of a main report and 7 topical reports.

The first step in establishing a study design was the elaboration of an approach paper (study outline) by the team leader. This was followed by a two weeks' reconnaissance visit to Dhaka and the greater Noakhali area. During this visit Bangladeshi researchers and assistants were recruited to the team, and more detailed plans for the subsequent fieldwork were drafted. Moreover, a background paper by Hasnat Abdul Hye, former Director General of BRDB and Secretary, Ministry of Local Government, was commissioned.

The fieldwork was preceded by a two days methodology-cum-planning workshop in Dhaka. The actual fieldwork lasted four months – from mid-January to mid-May 2001. The study team comprised 23 persons, five Bangladeshi researchers (two men, three women), three European researchers (two men, one woman), six research assistants (all men) and nine field assistants (including two women, all from Bangladesh). The researchers spent 11/2-31/2 months in the field, the assistants 2-4 months. Most of the time the team worked 60-70 hours a week. So it takes a good deal of resources to accomplish such a big and complex impact study.

## Case 2: Combining qualitative and quantitative descriptive methods - Mixed method impact evaluation of IFAD projects in Gambia, Ghana and Morocco

#### 1. Summary

The evaluation included intended and unintended impacts, and examined the magnitude, coverage and targeting of changes. It used mixed methods to gather evidence of impacts and the quality of processes with cross-checking between sources. With regard to *assessing causal contribution*, it must be noted that no baseline data were available. Instead a comparison group was constructed, and analysis of other contributing factors was made to ensure appropriate comparisons. The evaluation was undertaken within significant resource constraints and was carried out by an interdisciplinary team.

#### 2. Introduction and Background

Evaluations of rural development projects and country programs are routinely conducted by the Office of Evaluation of IFAD. The ultimate objectives of these evaluations is to (i) set a basis for accountability by assessing the development results and (ii) contribute to learning and improvement of design and implementation by providing lessons learned and practical recommendations. These evaluations follow a standardized methodology and a set of evaluation questions including the following: (i) project performance (relevance, effectiveness, and efficiency), (ii) project impact, (iii) overarching factors (sustainability, innovation and replication) and (iv) the performance of the partners. As can be seen, impact is but one the key evaluation questions and the resources allocated to the evaluation (budget, specialists and time) have to be shared for the entirety of the evaluation question.

As such, these evaluations are to be conducted under resource constraints. In addition, very limited data are available on socio-economic changes taking place in the project area that can be ascribed to an impact definition. IFAD adopts an impact definition which is similar to the DAC definition. The key feature of IFAD evaluation is that they are conducted just before or immediately after project conclusion: the effects can be observed after 4-7 years of operations and the future evolution can be estimated through an educated guess on sustainability perspectives. Several impact domains are considered including: (i) household income and assets, (ii) human capital, (iii) social capital, (iv) food security, (v) environment and (vi) institutions.

## 3. Sequencing of the process and choice of methods

This short case study is based on evaluations conducted in Gambia, Ghana and Morocco between 2004 and 2006. As explained above, evaluations had multiple questions to answer and impact assessment was but one of them. Moreover, impact domains were quite diverse. This meant that some questions and domains required quantitative evidence (for example in the case of household income and assets) while a more qualitative assessment would be in order for other domains (for example social capital). In many instances, however, more than one method would have to be used to answer the same questions, in order to cross-check the validity of findings, identify discrepancy and formulate hypotheses on the explanation of apparent inconsistencies.

As the final objective of the evaluation was not only to assess results but also provide future intervention designers with adequate knowledge and insights, the evaluation design could not be confined to addressing a dichotomy between "significant impact has been observed" and "no significant impact has been observed". Findings would need to be rich enough and grounded in field experience in order to provide a plausible explanation that would lead, when suitable, to a solution to identified problems and to recommendations to improve the design and the execution of the operations.

Countries and projects considered in this case study were ostensibly diverse. In all cases however, the first step in the evaluation consisted of a desk review of the project documentation. This allowed the evaluation team to understand or reconstruct the intervention theory (often implicit) and the logical framework. In turn, this would help to identify a set of hypotheses on changes that may be observed in the field as well as on intermediary steps that would lead to those changes.

In particular, the preliminary desk analysis highlighted that the results assessment would have to be supplemented with some analysis of implementation performance. The latter would include some insight in

the business processes (for example the management and resource allocation made by the project implementation unit) and the quality of service rendered (for example the topics and the communication quality of an extension service or the construction quality of a feeder road or of a drinking water scheme).

The second step was to conduct a preparatory mission. Among other purposes, this mission was instrumental in fine tuning our hypotheses on project results and designing the methods and instruments. Given the special emphasis of the IFAD interventions on the rural poor, impact evaluation would need to shed light, to the extent possible, on the following dimensions of impact: (i) magnitude of changes, (ii) coverage (that is the number of persons or households served by the projects and (iii) targeting, that is gauging the distribution of project benefits according to social, ethnic or gender grouping.

As pointed out before, a major concern was the absence of a baseline survey which could be used as a reference for impact assessment. This required reconstructing the "before project" situation. By the same token, it was clear that the observed results could not simply attributed to the evaluated interventions. In addition to exogenous factors such as weather changes, other important factors were at play, for example changes in government strategies and policies (such as the increased support to grassroots associations by Moroccan public agencies) or operations supported by other development organizations in the same or in adjacent zones. This meant that the evaluated interventions would interplay with existing dynamics and interact with other interventions. Understanding synergies or conflicts between parallel dynamics could not be done simply through inferential statistic instruments but required interaction with a wider range of stakeholders.

The third step in the process was the fielding of a data collection survey (after pre-testing the instruments) which would help the evaluation cope with the dearth of impact data. The selected techniques for data collection included: (i) a quantitative survey with a range of 200 – 300 households (including both project and control groups) and (ii) a more reduced set of focus group discussion with groups of project users and "control groups" stratified based on the economic activities in which they had engaged and the area where they were leaving.

In the quantitative survey standardized questionnaires were administered to final project users (mostly farmers or herders) as well as to non-project groups (control observations) on the situation before (recall methods) and after the project. Recall methods were adopted to make up for the absence of a baseline.

In the course of focus group interviews, open-ended discussion guidelines were adopted and results were mostly of qualitative nature. Some of the focus group facilitators had also been involved in the quantitative survey and could refer the discussion to observations previously made. After the completion of data collection and analysis, a first cross-checking of results could be made between the results of quantitative and qualitative analysis.

As a fourth step, an interdisciplinary evaluation team would be fielded. Results from the preliminary data collection exercise were made available to the evaluation team. The data collection coordinator was a member of the evaluation team or in a position to advise its members. The evaluation would conduct field visits and conduct a further validate survey and focus group data through participant observations and interviews with key informants (and further focus group discussions if necessary). The team would also spend adequate time with project management units in order to gather a better insight of implementation and business processes.

The final impact assessment would be made by means of triangulation of evidence captured from the (scarce) existing documentation, the preliminary data collection exercise and the main interdisciplinary mission (Figure 1).

Figure 1



#### 4. Constraints in data gathering and analysis

Threats to the validity of recall methods. According to the available literature sources<sup>40</sup> and our own experience, the reliability of recall methods may be questionable for monetary indicators (e.g. income) but higher for easier-to-remember facts (e.g. household appliances, approximate herd size). Focus group discussions helped identify possible sources of bias in the quantitative survey and ways to address them.

Finding "equivalent" samples for with and without-project observations. One of the challenges was to extract a control sample that would be "similar" in the salient characteristics to the project sample. In other words, problems of sampling bias and endogeneity should have been controlled for (e.g. more entrepreneurial people are more likely to participate in a rural finance intervention). In sampling control observations serious attempts were made to match project and non-project households based on similarity of main economic activities, agro-ecological environment, household size and resource endowment. In some instances, household that had just started to be served by the projects ("new entries") were considered as control groups, on the grounds that they would broadly satisfy the same eligibility criteria at entry as "older" project clients. However, no statistical technique (e.g. instrumental variables, Heckman's procedure or propensity score) was adopted to test for sampling bias, due to limited time and resources.

Coping with linguistic gaps. Given the broad scope of the evaluations, a team of international sector specialists was required. However, international experts were not necessarily the best suited for data collection analysis which calls for fluency in the local vernacular, knowledge of local practices and skills to obtain the most possible information within a limited time frame. Staggering the process in several phases was a viable solution. The preliminary data collection exercise was conducted by a team of local specialists, with university students or local teachers or literate nurses serving as enumerators.

## 5. Main value added of mixed methods and opportunities for improvement

The choice of methods was made taking into account the objectives of the evaluations and the resource constraints (time, budget and expertise) in conducting the exercise. The combination of multiple methods allowed us to cross-check the evidence and understand, for example, when survey questions were likely to be misinterpreted or generate over or under-reporting. On the other hand, quantitative evidence allowed us to shed light on the prevalence of certain phenomena highlighted during the focus group discussion. Finally the interactions with key informants and project managers and staff helped us better understand the reasons for under or over-achievements and come up with more practical recommendations.

The findings, together with the main conclusions and recommendations in the report were adopted in order to design new projects or a new country strategy. Also there was interest from the concerned project implementation agencies to adopt the format of the survey in order to conduct future impact assessments on their own. Due to time constraints, only inferential analysis was conducted on the quantitative survey

<sup>&</sup>lt;sup>4°</sup> Typical problems with recall methods are that of: (i) incorrect recalling and (ii) telescoping, i.e. projecting backward or forward an event: for example the purchase of a durable good which took place 7 years ago (before the project started) could be projected to four years ago, during project implementation. See, for example, Bamberger et al. (2004).

data. A full-fledged econometric analysis would have been desirable. By the same token, further analysis of focus group discussion outcomes would be desirable in principle.

## 6. A few highlights on the management of the process

The overall process design, as well as the choice of methods and the design of the data collection instruments was made by the Lead Evaluator in the Office of Evaluation of IFAD, in consultation with international sectoral specialists and the local survey coordinator. The pre-mission data collection exercise was coordinated by a local rural sociologist, with the help of a statistician for the design of the sampling framework and data analysis.

Time required conducting the survey and focus groups:

- Develop draft questionnaire and sampling frame, identify enumerators: 3 weeks;
- Conduct a quick trip on the ground, contact project authorities and pre-test questionnaires: 3 days;
- Train enumerators' and coders' team: 3 days;
- Survey administering: depending on the length of the questionnaire, on average an enumerator will be able to fill no more than 3-5 questionnaires per day. In addition time needs to be allowed for travel, rest. With a team of 6 enumerators, in 9-10 working days up to 200 questionnaires can be filled in, in the absence of major transportation problems;
- Data coding: it may vary depending on the length and complexity of the questionnaire. It is safe to assume 5-7 days;
- Time for conducting focus groups discussions: 7 days based on the hypothesis that around 10 FGD would be conducted by 2 teams;
- Data analysis. Depending on the analysis requirement it will require 1-2 weeks only to generate the tables and summary of focus group discussions;
- Drafting survey report: 2 weeks.

Note: as some of the above tasks can be conducted simultaneously, the total time for conducting a preliminary data collection exercise may be lower than the sum of its parts.

## Case 3: Combining qualitative and quantitative descriptive methods - Impact Evaluation: Agricultural Development Projects in Guinea

#### 1. Summary

The evaluation focused on impact in terms of poverty alleviation; the distribution of benefits was of particular interest, not just the mean effect. All data gathering was conducted after the intervention had been completed; mixed methods were used, including attention to describing the different implementation contexts. Assessing causal contribution is the major focus of the case study. A counter-factual was created by creating a comparison group, taking into account the endogenous and exogenous factors affecting impacts. Modeling was used to develop an estimate of the impact. With regard to the management of the impact evaluation, it should be noted that the study was undertaken as part of a PhD; the stakeholder engagement and subsequent use of it was limited.

This impact evaluation concerned two types of agricultural projects based in the Kpèlè region, in Guinea. The first one<sup>41</sup> was the Guinean Oil Palms and Rubber Company (SOGUIPAH). It was founded in 1987 by the Guinean government to take charge of developing palm oil and rubber production at the national level. With the support of several donors, SOGUIPAH quickly set up a program of industrial plantations<sup>42</sup> by negotiating the ownership of 22,830 ha with villagers. In addition, several successive programs were implemented between 1989 and 1998 with SOGUIPAH to establish contractual plantations<sup>43</sup> on farmers' own land and at the request of the farmers (1552 ha of palm trees and 1396 ha of rubber trees) and to improve 1093 ha of lowland areas for irrigated rice production.

The impact evaluation took place in a context of policy debates between different rural stakeholders at a regional level: two seminars had been held in 2002 and 2003 between the farmers' syndicates, the state administration, private sector and development partners (donors, NGOs) to discuss a regional strategy for agricultural development. These two seminars revealed that there was little evidence of what should be done to alleviate rural poverty, despite a long history of development projects. The impact of these projects on farmers' income seemed to be particularly relevant to assess, notably in order to compare the projects efficiency.

This question was investigated through a PhD thesis which was entirely managed by the AGROPARISTECH<sup>44</sup>. It was financed by AFD, one of the main donors in the rural sector in Guinea. This thesis proposed a new method, the systemic impact evaluation, aiming at quantifying impact using a qualitative approach. It enables to understand both the process through which impact materializes and to rigorously quantify the impact of agricultural development projects on the farmers' income, using a counterfactual. The analysis is notably based on the comprehension of the agrarian dynamics and the farmers' strategies, and permits the quantification of ex-post impact but also to devise a model of ex- ante evolution for the following years.

## 2. Gathering evidence of impact

The data collection was carried out entirely ex post. Several types of surveys and interviews were used to collect evidences of impact.

First, a contextual analysis realized all along the research work with key informants was necessary to describe the project implementation scheme, the contemporaneous events and the existing agrarian dynamics. It was also used to assess qualitatively whether those dynamics were attributable or not to the project. A series of surveys and historical interviews (focused on the pre-project situation) were notably

<sup>&</sup>lt;sup>41</sup> The second project was inland valley development for irrigated rice cultivation and will not be presented here.

<sup>&</sup>lt;sup>42</sup> Industrial plantations are the property of SOGUIPAH and are worked by salaried employees.

<sup>&</sup>lt;sup>43</sup> A contract between SOGUIPAH and the farmer binds the farmer to reimburse the cost of the plantation and deliver his production to SOGUIPAH.

<sup>&</sup>lt;sup>44</sup> AGROPARISTECH is a member of the Paris Institute of Technology which is a consortium of 10 of the foremost French Graduate Institutes in Science and Engineering. AGROPARISTECH is a leader Institute in Life Sciences and Engineering.

conducted in order to establish the most reliable baseline possible. An area considered "witness" to the agrarian dynamic that would have existed in the project's absence was identified.

Second, a preliminary structured survey (on about 240 households) was implemented, using recall to collect data on the farmers' situation in the pre-intervention period and during the project. It was the basis of a judgment sample to realize in depth interviews (see bellow), which aimed at describing the farming systems and quantifying rigorously the farmers' income.

## 3. Assessing causal attribution

By conducting an early contextual analysis, the evaluator was able to identify a typology of farming systems which existed before the project. In order to set up a sound counterfactual, a judgment sample was realized amongst the 240 households surveyed, by choosing 100 production units which had belonged to the same initial types of farming system and which had evolved with (in the project area) or without the project (in the witness area).

In-depth understanding of the endogenous and exogenous factors influencing the evolution and possible trajectories of farming systems enabled the evaluator to rigorously identify the individuals whose evolution with or without the project were comparable. This phase of judgment sample was followed by in-depth interviews with the hundred farmers. The evaluator's direct involvement in data collection was then essential, hence the importance of a small sample. It would not have been possible to gather reliable data on yields, modifications to production structures over time and producers' strategies from a large survey sample in a rural context.

Then, based on the understanding of the way the project proceeded and of the trajectories of these farmers, with or without the project, it was possible to build a quantitative model, based on Gittinger's method of economic analysis of development projects (Gittinger, 1982). As the initial diversity of production units was well identified before sampling, this model was constructed for each type of farming system existing before the project. Understanding the possible evolutions for each farming system with and without the project allowed for the estimation of the differential created by the project on farmers' income, that is its impact.

## 4. Ensuring rigor and quality

Although the objective differences between each production unit studied appear to leave room for the researcher's subjectivity when constructing the typology and sample, the rationale behind the farming system concept made it possible to transcend this possible arbitrariness. What underlies this methodological jump from a small number of interviews to a model is the demonstration that a finite number of types of farming systems exists in reality.

Moreover, (i) the use of a comparison group, (ii) the triangulation of most data collected by in-depth interviews through direct observation and contextual analysis and (iii) the constant implication of the principal researcher, were key factors to ensure rigor and quality.

## 5. Key findings

The large survey realized by interviewers on 240 households allowed identifying 11 trajectories related to the implementation of the project. Once each trajectory and each impact was characterized and quantified through in-depth interviews and modeling, this survey permitted as well quantifying a mean impact of the project, on the basis of the weight of each type in the population. The mean impact was only 24  $\epsilon$ /year/household in one village poorly served by the project, due to its enclosed situation, whereas it was 200  $\epsilon$ /year/household in a central village.

Despite a positive mean impact there were also highly differentiated impacts that existed, depending on the original farming system and the various trajectories with and without the project, which could not be ignored. Whereas former civil servants or traditional landlords beneficiated large contractual plantations, other villagers were deprived of their land for the needs of the project or received surfaces of plantations too limited to improve their economic situation.

Therefore, it seems important that the impact evaluation of a complex development project includes an analysis of the diversity of cases created by the intervention, directly or indirectly.

The primary interest of this new method was to give the opportunity to build a credible impact assessment entirely ex post. Second, it gave an estimate of the impact on different types of farming systems, making explicit the existing inequalities in the distribution of the projects' benefits. Third, it permitted a subtle understanding of the reasons why the desired impacts materialized or not.

## 6. Influence

The results from this first impact assessment were available after four years of field work and data treatment. They were presented to the Guinean authorities and to the local representatives of the main donors in the rural sector. In the field, the results were delivered to the local communities interviewed and to the farmers' syndicates. The Minister of Agriculture declared that he would try to foster more impact evaluations on agricultural development projects. Unfortunately, there is little hope that the conclusions of this research will change the national policy about these types of projects, in the absence of an institutionalized forum for discussing it between the different stakeholders.

## Case 4: A theory-based approach with qualitative methods - GEF Impact Evaluation 2007<sup>45</sup>

#### Evaluation of Three GEF Protected Area Projects in East Africa

#### 1. Description of Evaluation

The objectives of this evaluation included:

- To test evaluation methodologies that can assess the impact of GEF interventions. The key activity of the GEF is "providing new and additional grant and concessional funding to meet the agreed incremental costs of measures to achieve agreed global environmental benefits"<sup>46</sup>. The emphasis of this evaluation was therefore on verifying the achievement of agreed global environmental benefits.
- Specifically, to test a *Theory of Change approach* to evaluation in GEF's biodiversity focal area, and assess its potential for broader application within GEF evaluations.
- To assess the sustainability and replication of the benefits of GEF support, and extract lessons. It evaluated whether and how project benefits have continued, and will continue, after project closure.

#### Primary users

The primary users of the evaluation are GEF entities. They include: the GEF Council, which requested the evaluation; GEF Secretariat, which will approve future protected area projects, Implementing Agencies (such as the World Bank, UN agencies and regional Development Banks) and national stakeholders who will implement future protected area projects.

#### 2. Evaluation design

#### Factors driving selection of evaluation design

The Approach Paper to the Impact Evaluation<sup>47</sup> considered the overall GEF portfolio in order to develop an entry-point which could provide a good opportunity to develop and refine effective and implementable impact evaluation methodologies. Themes and projects that are relatively straightforward to evaluate were emphasized. The EO adopted the DAC definition of impact, which determined that closed projects would be evaluated to assess the sustainability of GEF interventions.

#### Biodiversity and protected areas:

The biodiversity focal area has the largest number of projects within the GEF portfolio of currently active and completed projects. In addition, biodiversity has developed more environmental indicators and global data sets than other focal areas, both within the GEF and in the broader international arena. The Evaluation Office chose protected areas as the central theme for this phase of the Impact Evaluation because: protected areas are one of the primary approaches supported by the GEF biodiversity focal area and its implementing agencies, and the GEF is the largest supporter of protected areas globally; previous evaluations have noted that an evaluation of the GEF support for protected areas are based on a set of explicit change theories, not just in the GEF, but in the broader conservation community; in many protected area projects, substantial field research has been undertaken, and some have usable baseline data on key factors to be changed by the intervention ; a protected areas strategy can be addressed at both a thematic and regional cluster level (as in East Africa, the region chosen for the study); the biodiversity focal area team

<sup>&</sup>lt;sup>45</sup> The GEF Evaluation Office section of the GEF website contains the 11 papers produced by the Impact Evaluation 2007, under the heading of "ongoing evaluations."

<sup>&</sup>lt;sup>46</sup> Instrument for the Establishment of the Restructured Global Environment Facility

<sup>&</sup>lt;sup>47</sup> GEF EO, "Approach Paper to Impact Evaluation", February 2006.

has made considerable progress in identifying appropriate indicators for protected areas through its "Managing for Results" system.

## The choice of projects

Lessons from a set of related interventions (or projects) are more compelling than those from an isolated study of an individual project. To test the potential for aggregation of project results, enable comparisons across projects and ease logistics, it was decided to adopt a sub-regional focus and select a set of projects that are geographically close to each other. East Africa is the subregion with the largest number of complete and active projects in the GEF portfolio with a protected area component, utilizing large GEF and cofinancing expenditure.

The following three projects were selected for evaluation:

- Bwindi Impenetrable National Park and Mgahinga Gorilla National Park Conservation Project, Uganda (World Bank);
- Lewa Wildlife Conservancy, Kenya (World Bank);
- Reducing Biodiversity Loss at Cross-Border Sites in East Africa, Regional: Kenya, Tanzania, Uganda (UNDP).

These projects were implemented on behalf of the GEF by the World Bank and UNDP. They have a variety of biodiversity targets, some of which are relatively easy to monitor (gorillas, zebras, rhinos). Also, these projects were evaluated positively by terminal and other evaluations and the continuance of long term results was predicted. The Bwindi Impenetrable National Park and Mgahinga Gorilla National Park Conservation Project is a \$6.7 million full-size-project and the first GEF-sponsored trust fund in Africa. The Lewa Wildlife Conservancy is a medium-sized-project, within a private wildlife conservation company. The Reducing Biodiversity Loss at Cross-Border Sites in East Africa Cross project is a \$12 million project, implemented at field level by Government agencies, that aims to foster an enabling environment for the sustainable use of biodiversity.

#### The advantages of a Theory of Change approach



some of which are outside the purview of the intervention (e.g. actions of

of several

are then

Box 1). The

actions at



exterior actors at the local, national or global levels or, change in political situations, regional conflicts and natural disasters). Subsequently, an intervention may have different levels of achievement in its component parts, giving mixed results towards its objectives.

#### The use of a hybrid evaluation model

During the process of field-testing, it was decided that, given the intensive data requirements of a theory-of-change approach and the intention to examine project impacts, the evaluation would mainly focus on the later elements of each project's theory-of-change, when outcomes are expected to lead to impact. Based on this approach, the evaluation developed a methodology composed of three components (see Box 2):



#### Box 2. Components of Impact Evaluation Framework

- Assessing implementation success and failure: To understand the contributions of the
  project at earlier stages of the results continuum, leading to project outputs and
  outcomes, a Logframe analysis is used. Though the normally complex and iterative process
  of project implementation is not captured by this method, the Logframe provides a means
  of tracing the realization of declared objectives. GEF interventions aim to "assist in the
  protection of the global environment and promote thereby environmentally sound and
  sustainable economic development"<sup>48</sup>.
- Assessing the level of contribution (i.e. impact): To provide a direct measure of project impacts, a <u>Targets-Threats analysis (Threats-Based Analysis</u>) is used to determine whether global environmental benefits have actually been produced and safeguarded<sup>49</sup>. The robustness of global environment benefits identified for each project (or 'targets') is evaluated by collecting information on attributes relating to the targets' biological composition, environmental requirements and ecological interactions. This analysis of targets is complemented by an assessment of the level of 'threat' (e.g., predation, stakeholder attitude and behavior) faced by the global environment benefits. For targets and significant threats, trends over time (at project start, at project close, and currently) and across project and non-project areas are sought, so that a comparison is available to assess levels of change.
- Explanations for observed impact: To unpack the processes by which the project addresses and contributes to impact, an <u>Outcomes-Impacts Theory-of-Change analysis</u> is used. This theory-of-change approach constructs and validates the project logic connecting outcomes and ultimate project impact. It involves a comprehensive assessment of the activities undertaken after project closure, along with their explicit and implicit assumptions. This component enables an assessment of the sustainability and/or catalytic nature of project interventions, and provides a composite qualitative ranking for the achievements of the projects. Elements of the varied aspects of sustainability include behavior change and the effectiveness of capacity-building activities, financial mechanisms, legislative change and institutional development.

<sup>&</sup>lt;sup>48</sup> See the Preamble, "Instrument for the Establishment of the Restructured Global Environment Facility".

<sup>&</sup>lt;sup>49</sup> This is based on Nature Conservancy's 'Conservation Action Planning' methodology.

The model incorporates three different elements that it is suggested are involved in the transformation of project outcomes into impacts. These are as follows, and were each scored for the level of achievement of the project in converting outcomes into impacts.

- Intermediary States. These are conditions that are expected to be produced on the way to delivering the intended impacts.
- Impact Drivers. These are significant factors or conditions that are expected to contribute to the ultimate realization of project impacts. Existence of the Impact Driver in relation to the project being assessed suggests that there is a good likelihood that the intended project impact will have been achieved. Absence of these suggests that the intended impact may not have occurred, or may be diminished.
- External Assumptions. These are potential events or changes in the project environment that would negatively or positively affect the ability of a project outcome to lead to the intended impact, but that are largely beyond the power of the project to influence or address.

#### 3. Data Collection and Constraints:

#### Logical Framework and Theory of Change Model:

The approach built on existing project logical frameworks, implying that a significant part of the Framework could be relatively easily tested through an examination of existing project documentation, terminal evaluation reports and, where available, monitoring data. Where necessary, targeted consultations and additional studies were carried out.

#### Assessing conservation status and threats to Global Environment Benefits:

A data collection framework for assessing the status of the targets and associated threats was developed, identifying indicators for each, along with the potential sources of information. For the Bwindi and Lewa projects, the task of collecting and assessing this information was undertaken by scientists from the Institute of Tropical Forest Conservation, headquartered in Bwindi Impenetrable National Park, and the Lewa Research Department respectively. For the Cross Borders project, this exercise was done by Conservation Development Center, based on the existing project documentation, a field visit to the project site and consultations with key informants. The objective of this exercise was to provide quantitative measures for each indicator from *before the project* (baseline), at the *project close*, and *present day*. Where quantitative data were not available, detailed qualitative data were collected.

#### Improving rigor

#### Internal validity:

The evaluation used a participatory approach with substantial involvement of former project staff in drawing out theories-of-change and subsequently providing data for verification. These data were verified by local independent consultants, via a process of triangulating information from project documentation and external sources. Given that all three projects are now closed, the participation from former project staff enabled a candid and detailed exchange of information (during Workshops in Uganda and Kenya). The participants in return found the process to be empowering, as it clarified and supported the rationale for their actions (by drawing out the logical connections between activities, goals and assumptions) and also enabled them to plan for future interventions.

#### External validity:

Given the small number of projects, their variety and age (approved in varied past GEF replenishment phases), the evaluation did not expect to produce findings, which could be directly aggregated. Nevertheless, given the very detailed analysis of the interventions a few years after project-closure, it did provide a wealth of insights into the functioning of protected area projects, particularly elements of their sustainability after project closure. This allowed limited generalization on key factors associated with achievement of impact, on the basis of different levels of results related to a set of common linkages in the theoretical models. On this basis, the

Evaluation Office recommended that the GEF Secretariat ensure specific monitoring of progress towards institutional continuity of protected areas throughout the life of a project.

#### Weaknesses

Impact evaluations are generally acknowledged to be highly challenging. The objective of this particular study, of examining GEF's impact at a 'global' level in biodiversity, make the study particularly complex. A few concerns include:

- The nature of changes in biodiversity is still under debate. Such changes are often nonlinear, with uncertain time-scales even in the short-run, interactions within and across species, and exogenous factors (like climate change). Evidence regarding the achievement of global environment benefits and their sustainability must therefore be presented with numerous caveats.
- Numerous explanations and assumptions may be identified for each activity that is carried out.
- The approach may not always uncover unexpected outcomes or synergies, unless they are anticipated in the theories or assumptions of the evaluation team. However, fieldwork should be able to discern such outcomes, as was the case in the Bwindi case study, which produced evidence of a number of unexpected negative impacts on local indigenous people.
- The association between activities and outcomes in the Theory of Change approach depends on measuring the level of activities carried out, and then consciously (logically) linking them with impact through a chain of intermediate linkages and outcomes. Information on these intermediate outcomes may be difficult to obtain, unless former project implementers participate fully in the evaluation.

#### 4. Concluding thoughts on the evaluation approach

For biodiversity, GEF's first strategic priority is "Catalyzing Sustainability of Protected Area Systems", which aims for an expected impact whereby "biodiversity [is] conserved and sustainably used in protected area systems."

The advantage of the hybrid evaluation model used was that by focusing towards the end of the results-chain, it examined the combination of mechanisms in place that have led to a project's impacts and ensure sustainability of results. It is during this later stage, after project closure, that the ecological, financial, political, socio-economic and institutional sustainability of the project are tested, along with its catalytic effects. During project conceptualization, given the discounting of time, links from outcome to impact are often weak. Once a project closes, the role of actors, activities and resources is often unclear; this evaluation highlighted these links and assumptions.

Adopting a Theory of Change approach also had the potential to provide a mechanism that helped GEF understand what has worked and what has not worked and allows for predictions regarding the probability of success for similar projects. The Evaluation Office team concluded that the most effective combination for its next round of impact evaluation (Phase-out of Ozone Depleting Substances in Eastern Europe) should seek to combine Theory of Change approaches with opportunistic use of existing data sets, which might provide some level of quantifiable counterfactual information.

## Application: Impact of Lewa Wildlife Conservancy (Kenya)<sup>50</sup>

Context

<sup>50</sup> Full Case Study at:

http://www.thegef.org/uploadedFiles/Evaluation\_Office/Ongoing\_Evaluations/Ongoing\_Evals-Impact-8Case\_Study\_Lewa.pdf

The Lewa GEF Medium-Sized Project provided support for the further development of Lewa Wildlife Conservancy ("Lewa"), a not-for-profit private wildlife conservation company that operates on 62,000 acres of land in Meru District, Kenya. The GEF awarded Lewa a grant of \$0.75 million for the period 2000 to the end of 2003, with co-financing amounting to \$3.193 million.

Since the GEF grant, Lewa has been instrumental in initiating the formation of the Northern Rangelands Trust (NRT) in 2004. NRT is an umbrella local organization with a goal of collectively developing strong community-led institutions as a foundation for investment in community development and wildlife conservation in the Northern Rangelands of Kenya. The NRT membership comprises community conservation conservancies and trusts, local county councils, the Kenya Wildlife Service, the private sector, and NGOs established and working within the broader ecosystem. The establishment and functioning of the NRT has therefore been a very important aspect in understanding and assessing the ultimate achievement of impacts from the original GEF investment.

The Lewa Case study implemented the three elements of the Impact Evaluation Framework which are summarized below.

#### Assessed implementation success and failure

Given that no project logical framework or outcomes were defined as such in the original GEF project brief, the GEF Evaluation Office Study of Local Benefits in Lewa (2004), with the participation of senior Lewa staff, identified project outcomes and associated outputs that reflected the various intervention strategies employed by the project and identified missed opportunities in achieving the project goals. The assessment was as follows, and provided an understanding of the project logic used (Box 1) and a review of the fidelity with which the project activities were implemented (Box 2):


| Box 1. (b) Project Outcome  | Assessment         |
|---|--------------------|
| <b>Outcome 1:</b> Long-term institutional and financial capacity of Lewa to provide global and local benefits from wildlife conservation strengthened | Fully achieved (5) |
| <b>Outcome 2:</b> Protection and management of endangered wildlife species in the wider ecosystem strengthened  | Well achieved (4)  |
| <b>Outcome 3:</b> Community-based conservation and natural resource management initiatives strengthened   | Well achieved (4)  |

Assessed the level of contribution (i.e. impact)

A *Targets-Threats analysis* of those ecological features identified as global environment benefits (Black Rhinos and Grevy's Zebra) was undertaken with input from scientists from Lewa and the Northern Rangelands Trust research departments. Box 2. (a) and (b) provide an overview of the variables considered to increase robustness of the understanding of ecological changes that have taken place since before the project started.

|   | Box 2. (a) Change in Key   | Ecological A | Attributes o | ver time       |        |       |
|---|--|--------------|--------------|----------------|--------|-------|
| Vay Faalagiaal                                    |  |              | Co           | nservation S   | tatus  |       |
| Attribute   | Indicator  | Unit         | Baseline     | Project<br>end | Now    | Trend |
| Black Rhino                                       |  |              |              |                |        |       |
| Population size                                   | Total population size of Black<br>rhino on Lewa                    | Number       | 29           | 40             | 54     | Î     |
| Productivity                                      | Annual growth rates at Lewa  | %            | 12           | 13             | 15     |       |
| Suitable secure<br>habitat                        | Size of Lewa rhino sanctuary                                       | Acres        | 55,000       | 55,000         | 62,000 |       |
| Genetic diversity                                 | Degree of genetic variation  | -            | N            | No data availa | ble    |       |
| Grevy's zebra                                     |  |              |              |                |        |       |
| Population size                                   | Total population size of<br>Grevy's zebra on Lewa                  | Number       | 497          | 435            | 430    |       |
| Productivity                                      | Annual foaling rates on Lewa                                       | %            | 11           | 11             | 12     |       |
| Population<br>distribution                        | Number of known sub-<br>populations and connectivity               |              | N            | lo data availa | ble    |       |
| Suitable habitat<br>(grassland & secure<br>water) | Community conservancies set<br>aside for conservation under<br>NRT | Number       | 3            | 4              | 15     |       |
| Genetic diversity                                 | Degree of genetic variation  |              | N            | No data availa | ble    |       |

| Box 2. (b) Current Threats to the G    | lobal environn                        | nent benefits(G                    | EBs)               |
|--|---------------------------------------|------------------------------------|--------------------|
|  | Severity <sup>51</sup><br>Score (1-4) | Scope <sup>52</sup><br>Score (1-4) | Overall<br>ranking |
| Black rhino                            |                                       |                                    |                    |
| Poaching and snaring                   | 3                                     | 3                                  | 3                  |
| Insufficient secure areas              | 2                                     | 3                                  | 2                  |
| Habitat loss (due to elephant density) | 1                                     | 1                                  | 1                  |
| Grevy's zebra                          |                                       |                                    |                    |

<sup>&</sup>lt;sup>51</sup> **Severity** (level of damage): Destroy or eliminate GEBs/Seriously degrade the GEBs/Moderately degrade the GEBs/Slightly impair the GEBs.

<sup>&</sup>lt;sup>52</sup> **Scope** (geographic extent): Very widespread or pervasive/Widespread/Localized/ Very localized.

| Poaching                            | 2 | 2 | 2 |
|-------------------------------------|---|---|---|
| Disease                             | 4 | 2 | 3 |
| Predation                           | 3 | 1 | 2 |
| Habitat loss/ degradation           | 3 | 3 | 3 |
| Insufficient secure areas           | 2 | 2 | 2 |
| Hybridization with Burchell's zebra | 1 | 1 | 1 |

Explanations for observed impact

Theory of Change models were developed for each project Outcome to <u>establish contribution</u>; the framework reflected in Box 6.3(a) was used. This analysis enabled an examination of the links between observed project interventions and observed impact. As per GEF principles, factors that were examined as potentially influencing results included the *appropriateness* of intervention, the *sustainability* of the intervention and its *catalytic effect* – these are referred to as 'impact drivers.' The next step involved the identification of 'intermediary states': examining whether the successful achievement of a specific project outcome would directly lead to the intended impacts and, if not, identifying additional conditions that would need to be met to deliver the impact. Taking cognizance of factors that are 'beyond project control', the final step identified those factors that are necessary for the realization and sustainability of the intermediary state(s) and ultimate impacts, but outside the project's influence.



An illustrative example is provided by a consideration of Outcome 3 that via *Community-based* conservation and natural resource management initiatives strengthened, expected to achieve enhanced conservation of Black Rhinos and Grevy's Zebras. The *theory of change* model linking Outcome 3 to the intended impacts is illustrated below, in Box 6.3(b). The overall logframe assessment of the project's implementation for community-based conservation and natural resource management was *well achieved*. All intermediate factors/impact drivers/external assumptions that were identified received a score of *partially to well achieved*, indicating that together with all its activities, this component was well-conceived and implemented.



## In sum for Lewa

The analysis provided indication that the Black rhino and Grevy's zebra populations on the Lewa Conservancy are very well managed and protected. Perhaps the most notable achievement has been the visionary, catalytic and support role that Lewa has provided for the conservation of these endangered species in the broader ecosystem, beyond Lewa. Lewa has played a significant role in the protection and management of about 40% of Kenya's Black rhino population and is providing leadership in finding innovative ways to increase the coverage of secure sanctuaries for Black rhino. Regarding the conservation of Grevy's zebra, Lewa's role in the establishment of community conservancies, which have added almost one million acres of land set aside for conservation, has been unprecedented in East Africa and is enabling the recovering of Grevy's zebra populations within their natural range. However, the costs and resources required to manage and protect this increasing conservation estate are substantial and unless the continued and increasing financing streams are maintained, it is possible that the substantial gains in the conservation of this ecosystem and its global environmental benefits could eventually be reversed.

### In conclusion

The assessment of project conceptualization and implementation of project activities in Lewa has been favorable, but, this is coupled with indications that threats from poaching, disease and habitat loss in and around Lewa continue to be severe. Moreover, evaluation of the other case studies *Bwindi Impenetrable National Park and Mgahinga Gorilla National Park Conservation Project*, Uganda and *Reducing Biodiversity Loss at Cross-Border Sites in East Africa*, Regional: Kenya, Tanzania, Uganda, confirmed that to achieve long-term results in the generation of global environment benefits the absence of a specific plan for institutionalized continuation would, in particular, reduce results over time – this was the major conclusion of the GEF's pilot impact evaluation.

# Appendix 12. Basic Education in Ghana

# Introduction

In 1986 the Government of Ghana embarked on an ambitious program of educational reform, shortening the length of pre-University education from 17 to 12 years, reducing subsidies at the secondary and tertiary levels, increasing the school day and taking steps to eliminate unqualified teachers from schools. These reforms were supported by four World Bank credits – the Education Sector Adjustment Credits I and II, Primary School Development Project and the Basic Education Sector Improvement Project. An impact study by the World Bank evaluation department, IEG, looked at what had happened to basic education (grades 1 to 9, in primary and junior secondary school) over this period.

### Data and methodology

In 1988/89 Ghana Statistical Service (GSS) undertook the second round of the Ghana Living Standards Survey (GLSS 2). Half of the 170 areas surveyed around the country were chosen at random to have an additional education module, which administered math and English tests to all those aged 9-55 years with at least three years of schooling and surveyed schools in the enumeration areas. Working with both GSS and the Ministry of Education, Youth and Sport (MOEYS), IEG resurveyed these same 85 communities and their schools in 2003, applying the same survey instruments as previously. In the interests of comparability, the same questions were kept, although additional ones were added pertaining to school management, as were two whole new questionnaires – a teacher questionnaire for five teachers at each school and a local language test in addition to the math and English tests. The study thus had a possibly unique data set – not only could children's test scores be linked to both household and school characteristics, but this could be done in a panel of communities over a fifteen year period. The test scores are directly comparable since exactly the same tests were used in 2003 as had been applied fifteen years earlier.

There was no clearly defined 'project' for this study, rather support to the sub-sector through four large operations. The four projects had supported a range of activities, from rehabilitating school buildings to assisting in the formation of community-based school management committees. To identify the impact of these various activities a regression-based approach was adopted which analyzed the determinants of school attainment (years of schooling) and achievement (learning outcomes, i.e. test scores). For some of these determinants – notably books and buildings – the contribution of the World Bank to better learning outcomes could then be quantified. The methodology adopted a theory-based approach to identify the channels through which a diverse range of interventions were having their impact. As discussed below, the qualitative context of the political economy of education reform in Ghana at the time proved to be a vital piece of the story.

### Findings

The first major finding from the study was the factual. Contrary to official statistics, enrolments in basic education have been rising steadily over the period. This discrepancy was readily explained: in the official statistics, both the numerator and denominator were wrong. The numerator was wrong as it relied on the administrative data from the school census, which had incomplete coverage of the public sector and did not cover the rapidly growing private sector. A constant mark up was made to allow for private sector enrolments, but the IEG analysis showed that had gone up fourfold (from 5 to 20% of total enrolments) over the 15 years. The denominator was based on the 1984 census with an assumed rate of growth which turned out to be too high once the 2000 census became available, thus underestimating enrolment growth.

More strikingly still, learning outcomes have improved markedly: 15 years ago nearly twothirds (63 percent) of those who had completed grades 3-6 were, using the English test as a guide, illiterate. By 2003 this figure had fallen to 19 percent. The finding of improved learning outcomes flies in the

face of qualitative data from many, though not all, 'key informant' interviews. But such key informants display a middle class bias which persists against the reforms which were essentially populist in nature.

Also striking are the improvements in school quality revealed by the school-level data: For example:

- In 1988, less than half of schools could use all their classrooms when it was raining, but in 2003 over two-thirds can do so;
- Fifteen years ago over two-thirds of primary schools reported occasional shortages of chalk, only one in 20 do so today, with 86 percent saying there is always enough;
- The percentage of primary schools having at least one English textbook per pupil has risen from 21 percent in 1988 to 72 percent today and for math books in Junior Secondary School (JSS) these figures are 13 and 71 percent, respectively.

School quality has improved across the country, in poor and non-poor communities alike. But there is a growing disparity within the public school sector. Increased reliance on community and district financing has meant that schools in relatively prosperous areas continue to enjoy better facilities than do those in less well off communities.

The IEG study argues that Ghana has been a case of a quality-led quantity expansion in basic education. The education system was in crisis in the seventies; school quality was declining and absolute enrolments falling. But by 2000, over 90 percent of Ghanaians aged 15 and above had attended school compared to 75 percent 20 years earlier. In addition, drop-out rates have fallen, so completion rates have risen: by 2003, 92 percent of those entering grade 1 complete Junior Secondary School (grade 9). Gender disparities have been virtually eliminated in basic enrolments. Primary enrolments have risen in both disadvantaged areas and amongst the lowest income groups. The differential between both the poorest areas and other parts of the country, and between enrolments of the poor and non-poor, have been narrowed but are still present.

Statistical analysis of the survey results showed the importance of building school infrastructure on enrolments. Building a school, and so reducing children's travel time, has a major impact on enrolments. While the majority of children live within 20 minutes of school, some 20 percent do not and school building has increased enrolments among these groups. In one area surveyed, average travel time to the nearest school was cut from nearly an hour to less than 15 minutes with enrolments increasing from 10 to 80 percent. In two other areas average travel time was reduced by nearly 30 minutes and enrolments increased by over 20 percent. Rehabilitating classrooms so that they can be used when it is raining also positively affects enrolments. Complete rehabilitation can increase enrolments by as much as one third. Across the country as a whole, the changes in infrastructure quantity and quality have accounted for a 4 percent increase in enrolments between 1988 and 2003, about one third of the increase over that period. The World Bank has been the main source of finance for these improvements. Before the first World Bank program communities were responsible for building their own schools. The resulting structures collapsed after a few years. The Bank has financed 8,000 school pavilions around the country, providing more permanent structures for the school which can better withstand the weather.

Learning outcomes depend significantly on school quality, including textbook supply. Bankfinanced textbook provision accounts for around one quarter of the observed improvement in test scores. But other major school-level determinants of achievement such as teaching methods and supervision of teachers by the head teacher and circuit supervisor have not been affected by the Bank's interventions. The Bank has not been heavily involved in teacher training and plans to extend in-service training have not been realized. Support to "hardware" has been shown to have made a substantial positive contribution to both attainment and achievement. But when satisfactory levels of inputs are reached — which is still far from the case for the many relatively deprived schools — future improvements could come from focusing on what happens in the classroom. However, the Bank's one main effort to change incentives — providing head teacher housing under the Primary School Development Project in return for the head teacher signing a contract on school management practices — was not a great success. Others, notably DFID and USAID, have made better progress in this direction but with limited coverage.

The policy context, meaning government commitment, was an important factor in making the Bank's contributions work. The government was committed to improving the quality of live in rural areas, through the provision of roads, electricity and schools, as a way of building a political base. Hence there was a desire to make it work. Party loyalists were placed in key positions to keep the reform on track, the army used to distribute textbooks in support of the new curriculum in the early 1990s to make sure they reached schools on time, and efforts made to post teachers to new schools and make sure that they received their pay on time.

Teachers also benefited from the large civil service salary increase in the run up to the 1992 election. Better education leads to better welfare outcomes. Existing studies on Ghana show how education reduces fertility and mortality. Analysis of IEG's survey data shows that education improves nutritional outcomes, with this effect being particularly strong for children of women living in poorer households. Regression analysis shows there is no economic return to primary and JSS education (i.e. average earnings are not higher to children who have attended primary and JSS compared to children who have not), but there is a return to cognitive achievement. Children who attain higher test scores as a result of attending school can expect to enjoy higher income; but children who learn little in school will not reap any economic benefit.

## Some policy implications

The major policy finding from the study relates to the appropriate balance between hard and software in support for education. The latter is now stressed. But the study highlights the importance of hardware: books and buildings. It was also of course important that teachers were in their classrooms: government's own commitment (borne out of a desire to build political support in rural areas) helped ensure this happened.

In the many countries and regions in which educational facilities are inadequate then hardware provision is a necessary step in increasing enrolments and improving learning outcomes. The USAID project in Ghana encourages teachers to arrange children's desks in groups rather than rows – but many of the poorer schools don't have desks. In the words of one teacher, "I'd like to hang posters on my walls but I don't have posters. In fact, as you can see, I don't have any walls".

These same concerns underlie a second policy implication. Central government finances teacher's salaries and little else for basic education. Other resources come from donors, districts or the communities themselves. There is thus a real danger of poorer communities falling behind, as they lack both resources and the connections to access external resources. The reality of this finding was reinforced by both qualitative data – field trips to the best and worst performing schools in a single district in the same day – and the quantitative data, which show the poorer performance of children in these disadvantaged schools. Hence children of poorer communities are left behind and account for the remaining illiterate primary graduates which should be a pressing policy concern.

The study highlighted other areas of concern. First amongst these is low teacher morale, manifested through increased absenteeism. Second is the growing importance of the private sector, which now accounts for 20 percent of primary enrolments compared to 5 percent 15 years earlier. This is a sector which has had limited government involvement and none from the Bank.

| Key<br>T = Time<br>P = Project participants; C = Control group  | Start of<br>project<br>[pre-test] | Project<br>intervention<br>[Process not | Mid-term<br>evaluation | End of<br>project<br>[Post-test] | The stage of the project<br>cycle at which each<br>evaluation design can to be |
|---|-----------------------------------|---|------------------------|----------------------------------|--|
| $P_1$ , $P_2$ , $C_1$ , $C_2$ first and second observations $X = Project$ intervention (a process rather than a discrete event)       |                                   | discrete<br>event]                      |                        |                                  | used.  |
| Quantitative Impact Evaluation Design   | $T_1$                             |   | $T_2$                  | $T_3$                            |  |
| <b>RELATIVELY ROBUST QUASI-EXPERIMENTAL DESIGNS</b>   |                                   |   |                        |                                  |  |
| 1. Pre-test post-test non-equivalent control group design with  | $\mathbf{P}_1$                    | X                                       |                        | $\widetilde{\mathbf{h}}_2$       | Start  |
| statistical matching of the two groups. Participants are either self-<br>selected or are selected by the project implementing agency. | C.                                |   |                        | $C_2$                            |  |
| Statistical techniques (such as propensity score matching), drawing   |                                   |   |                        |                                  |  |
| on high quality secondary data used to match the two groups on a number of relevant variables   |                                   |   |                        |                                  |  |
| 2. Pre-test post-test non-equivalent control group design with judgmental   | P1                                | X                                       |                        | $\mathbf{P}_2$                   | Start  |
| matching of the two groups. Participants are either self-selected or are  | CI                                |   |                        | $C_2$                            |  |
| selected by the project implementing agency Control areas usually selected  |                                   |   |                        |                                  |  |
| judgmentally and subjects are randomly selected from within these areas.  |                                   |   |                        |                                  |  |
| LESS ROBUST QUASI-EXPERIMENTAL DESIGNS  |                                   |   |                        |                                  |  |
| 3. Pre-test/post-test comparison where the baseline study is not conducted  |                                   | Х                                       | $P_1$                  | $\mathbf{P}_2$                   | During project   |
| until the project has been underway for some time (most commonly this is  |                                   |   | $C_1$                  | $C_2$                            | implementation (often at   |
| around the mid-term review).  |                                   |   |                        |                                  | mid-term)  |
| 4. Pipeline control group design. When a project is implemented in phases,  | $\mathbf{P}_1$                    | Х                                       |                        | $\mathbf{P}_2$                   | Start  |
| subjects in Phase 2 (i.e who will not receive benefits until some later point in  | C1                                |   |                        | $C_2$                            |  |
| time) can be used as the control group for Phase 1 subjects.  |                                   |   |                        |                                  |  |
| 5. Pre-test post-test comparison of project group combined with post-test   | $\mathbf{P}_{1}$                  | X                                       |                        | $\mathbf{P}_2$                   | Start  |
| comparison of project and control group.  |                                   |   |                        | $C_2$                            |  |
| 6. Post-test comparison of project and control groups   |                                   | X                                       |                        | P <sub>1</sub>                   | End  |
| NON-EXPERIMENTAL DESIGNS (THE LEAST ROBUST)   |                                   |   |                        |                                  |  |
| 7. Pre-test post-test comparison of project group   | $\mathbf{P}_1$                    | Х                                       |                        | $\mathbf{P}_2$                   | Start  |
| 8. Post-test analysis of project group.   |                                   | Х                                       |                        | $\mathbf{P}_1$                   | End  |

Appendix 13. Hierarchy of quasi-experimental designs

Source: Bamberger et al. (2006)

115