# Theory-Based Evaluation in Practice

2022 NEC Conference
Jos Vaessen, PhD

/undp_evaluation     /Indep. Evaluation Office

# Introduction

# Objectives of the workshop

After this workshop, participants have developed an initial sound understanding of the role of program theory in evaluation and how to apply theory-based evaluation in practice. More specifically, participants will have a greater understanding of:

- Different purposes and uses of program theory in evaluation

- Principles for reconstructing a program theory

- Applications of theory-based evaluation in practice

# Outline

- 9.00 – 10.30: Principles of Theory-Based Evaluation

coffee/tea

- 11.15 – 12.30: Reconstructing a Program Theory (exercise)

lunch

- 14.00 – 15.15: Reconstructing a Program Theory (continued and plenary discussion)

coffee/tea

- 15.45 – 17.00: Using Program Theory as a framework for evaluation

# Principles of Theory-Based Evaluation

# Definitions

"[Program theory] is a set of hypotheses upon which people build their program plans" (Weiss, 1998:55).

"[TBE] consists of an explicit theory or model of how the program causes the intended or observed outcomes and an evaluation that is at least partly guided by this model" (Rogers et al., 2000:5).

Program theory cannot be simply 'observed' but most be reconstructed.

# GEI Theory of Change

## Scale

The GEI brand enhances the partnership's convening power, supports the establishment of effective collaborations with external partners, and helps to attract new partners and funding, all of which contribute to achieving economies of scale in ECD

## Quality

GEI's quality of work benefits from common standards, expertise and knowledge sharing among partners, and joint work to optimize partners' comparative advantages
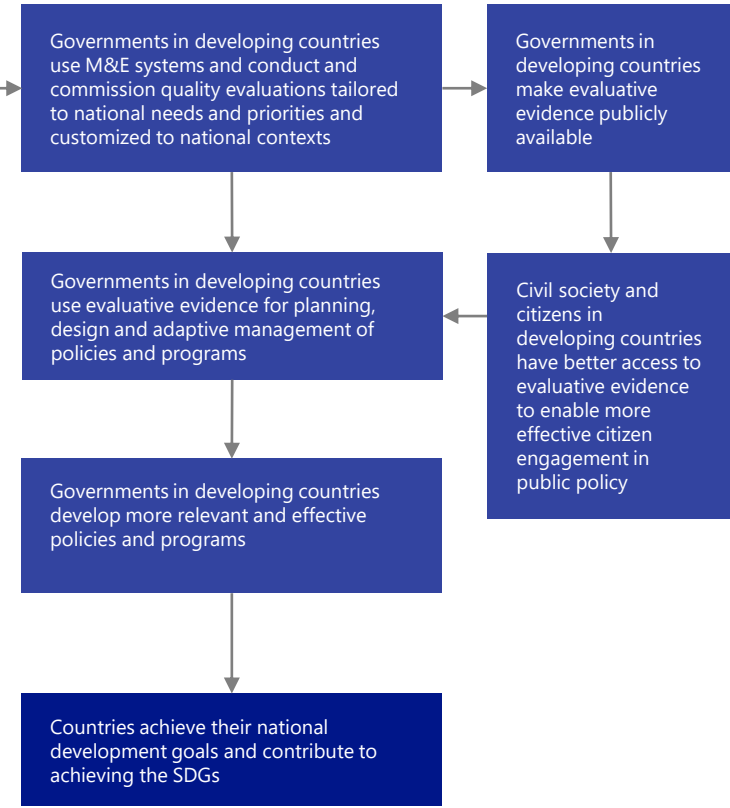
## Strategic Orientation

The strategic orientation of GEI's work benefits from a shared work program and an efficient division of labor based on partners' comparative advantages, as well as strategic collaborations that leverage key synergies

### Develop a culture of evidence-informed decision making in developing countries

› GEI contributes to bringing together national and international stakeholders to better coordinate evaluation plans and initiatives to strengthen M&E systems and capacities in governments in (selected) developing countries
› GEI engages in awareness-raising activities on the role of M&E among governments and other stakeholders in developing countries
› GEI provides TA and advisory services to governments in (selected) developing countries:
  • To strengthen the enabling environment (understanding of the role of M&E in learning and accountability; legislation; policies)
  • To develop and support organizational processes and systems

### Strengthen a cadre of evaluators, M&E specialists, and other evaluation stakeholders in developing countries (especially in priority M&E areas: gender, environmental sustainability and inclusion)

› GEI provides tailored global, regional, national M&E trainings to evaluation stakeholders from developing countries
› GEI provides institution-specific training on M&E issues to governments in (selected) developing countries
› GEI establishes a scholarship scheme to support training M&E professionals, prioritizing (E)FDEs
› GEI establishes an internship program for (emerging) evaluators and M&E specialists in developing countries
› GEI develops, applies and shares good practices and international standards for M&E training
› GEI helps develop quality M&E curricula and competencies in (selected) academic institutions in developing countries

### Generate M&E knowledge (especially in priority M&E areas: gender, environmental sustainability and inclusion)

› GEI collects and curates knowledge and lessons learned from internal and external sources
› GEI (co-) conducts research and (co-) generates knowledge on M&E related themes, practices, processes, systems and methods

### Share M&E knowledge (especially in priority M&E areas: gender, environmental sustainability and inclusion)

› GEI publishes and shares knowledge through publication series, tools and learning events
› GEI (co-) implements a dedicated knowledge platform on M&E issues
› GEI (co-) organizes the National Evaluation Capacities (NEC) Conference
› GEI organizes the gLocal Evaluation Week
› GEI (co-) organizes and participates in other learning and convening events
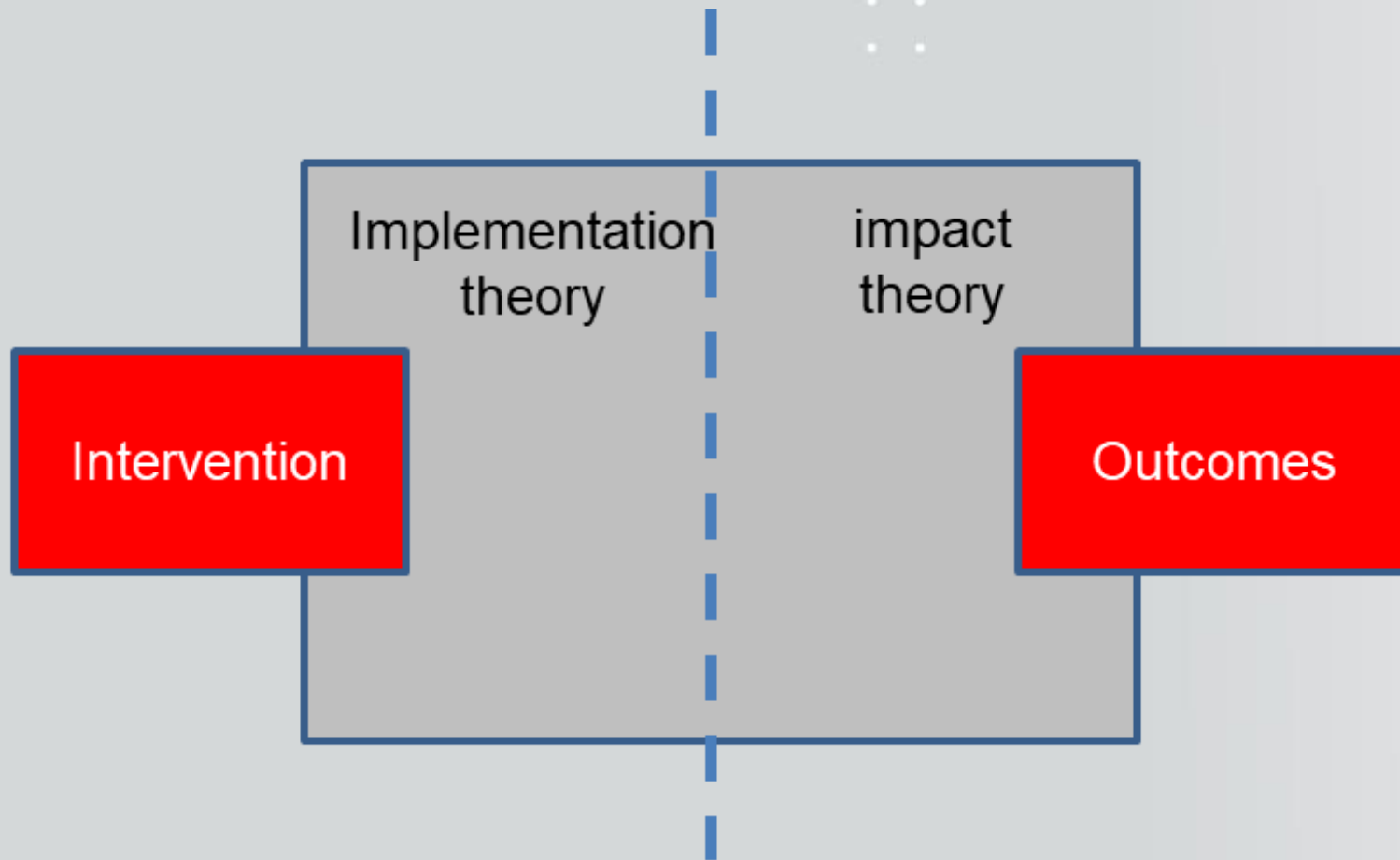› GEI collaborates with ECD partners to facilitate knowledge sharing and dialogue

---

Governments in developing countries are increasingly capable of coordinating evaluations at central government and sectoral levels as well as initiatives to strengthen M&E systems and capacities

Governments and other stakeholders in developing countries better understand the role of M&E in (evidence-informed) decision-making, organizational learning and accountability

Governments in developing countries put in place of improve a regulatory and policy environment that helps practitioners and decision-makers to produce and use evaluative evidence more effectively

Governments in developing countries put in place or improve M&E organizational frameworks, processes and systems to support (evidence-informed) decision-making, organizational learning and accountability

Governments and other stakeholders in developing countries are better capable of conducting evaluations and of managing and using M&E systems for (evidence-informed) decision-making, organizational learning and accountability (especially in priority M&E areas: gender, environmental sustainability and inclusion)

Governments and other stakeholders in developing countries use M&E knowledge products and attend knowledge events

---

Governments in developing countries use M&E systems and conduct and commission quality evaluations tailored to national needs and priorities and customized to national contexts

Governments in developing countries make evaluative evidence publicly available

Governments in developing countries use evaluative evidence for planning, design and adaptive management of policies and programs

Civil society and citizens in developing countries have better access to evaluative evidence to enable more effective citizen engagement in public policy

Governments in developing countries develop more relevant and effective policies and programs

Countries achieve their national development goals and contribute to achieving the SDGs

---

Feedback

Implementation theory

impact theory

Intervention

Outcomes

Theory failure vs. implementation failure

# Purpose of program theory in evaluation

- **Understanding why** interventions do or do not work (implementation versus theory failure)

- **Generating a consensus** on what the intervention is intended to achieve and how (formative use)

- Program theory as an overall **sense-making framework**

- Using program theory as **a basis for data collection and analysis** or M&E system

- Dealing with **causality**

# Exercise – "good" program theory

- You will be given a particular representation of a program theory

- Please respond to the following question:

1. Do you find this a convincing program theory?

2. Which purpose(s) of PT in evaluation would this theory support? (e.g. mention 1 or 2)?

3. What do you consider to be strong aspects of this program theory?

4. What do you consider to be weak aspects of this program theory?

# 1 – school inspection



Quality asssssment → Strengths, weaknesses, and the legal requirements schools fail to meet are made clear

Proportional inspection → Schools develop quality assurance

Proportional inspection → Inspection uses its capacity more efficiently → More frequent and intense inspection of weak schools → Schools start improving → Schools improve their educational quality → Schools attain satisfactory levels of educational quality / Schools offer more added value

Inspection results are public → Schools account for actions and results → Parents know about school quality → Parents choose schools

Inspection results are public → Parents take note of results → Parents know about school quality → Parents assess the quality of schools → Parents address schools (or their administrators)

risk aversity
farmers

capital
constraints

opportunity costs
of labour

farmers will feel more capable of
changing their situation in a
positive direction

THEN 3

IF farmers: THEN 1 farmers participate in THEN 2 farmers will better THEN 4 farmers will apply the THEN 8 farmers will change THEN 9 the full benefits associated
- are motivated      the programme          understand their          imparted knowledge on their   their farming system       with the different practices
- receive            - attend courses        farming system            own farms                     from conventional          and integrated application
incentives           - receive technical     and the problems                                        to organic farming         of practices will emerge
                     assistance              they face                 - physical
                                                                       practices                                                 - higher yields      higher
selection            distance and                                     - cultural practices                                       - less reliance on   household
mechanisms           travelling costs                                                                                            external inputs      income
                                                                                                                                 - better soils
                                             THEN 5                                                                              - soil and water
                                                                                                                                 conservation
                                                                                                                                 - conservation and
                                             farmers will share the                                                             manifestation of the
                                    THEN 6   knowledge with their                                                               Mayan
characteristics                              neighbours                                                                         farming tradition
social structure
                                             farmers will organize
                                             themselves to exchange ideas
                                             and experiences

                                             THEN 7

                                             IRDP will be able to
                                             work with farmer groups

# 3 – trade facilitation



**OUTCOMES**

| ACTIVITIES | OUTPUT | Immediate | Intermediate | Final |

**Change in behavior** · **Change in trade related costs** · **Change in international trade flows and achievement of public policy goals**

**ACTIVITIES**

Government facing **POLICY ADVICE**
(IFC, T&C)

Public and Private **LENDING and INVESTMENT**
(T&C, AGR, IFC, MIGA)

Public and Private **CAPACITY BUILDING**
(IFC, T&C, AGR)

**PUBLIC GOODS** Knowledge Generation; Convening Power

**OUTPUT**

- Assessments on trade facilitation simplification completed
- Review of strategies completed
- Regulatory framework developed
- Policy actions recommended

- Regulations are harmonized
- Trade requirements are made more transparent
- PPD system is facilitated or established
- IT system is reengineered
- Custom governance are introduced or enhanced
- Trade services are enhanced
- Single Window is established

- Public/Private trade related agencies staff is trained
- Trade related agencies coordination is enhanced
- Inspection are revised

- Standards are developed
- Knowledge products and indicators are generated
- Global agreements are supported
- Global policy actors are aligned

- Trade related public health and safety or environmental regulations/procedures are introduced or enhanced
- Risk management procedures are adopted or enhanced

**Immediate**

- Policies are enacted
- Regulations are adopted
- IT systems are operational
- Single Window is functioning
- Customs and other agencies apply new or enhanced rules, procedures, or systems
- Public and private operators increase collaboration
- Logistic services are more efficient
- Data is used to identify policy priorities
- Trade approaches are harmonized

- Risk management techniques are operational
- Health, safety, envir. regulations are implemented

**Intermediate**

**DIRECT COSTS***
- Simplified documentation
*(Potential indicators: Fees; number of documents; frequency of inspections)*

**INDIRECT COSTS***
- Enhanced procedures
- Improved IT systems
- Reduced inventory holding
- Efficient functioning or coordination of relevant agencies
*(Potential indicators: Cost to export/import; No of days to clear customs; Time to complete bureaucratic procedures for export/import; Share of cargo physically inspected; quality of logistic services)*

**HIDDEN COSTS***
- Corruption
- Smuggling
*(Potential indicator: irregular payments in export/import)*

**PUBLIC HEALTH, SAFETY, ENVIRONMENTAL COSTS (AND BENEFITS)**

**POSITIVE SPILLOVERS**

**NEGATIVE SPILLOVERS**

**Final**

**Trade Flows** Facilitate export/import flows

**Public Policy Goals** Reduced risk to public safety, health, environment, corruption, smuggling, informality, public revenues

*Costs are measured with monetary and time indicators

EXTERNAL FACTORS: Macro environment, Infrastructure and domestic logistics, Financial markets, etc.
COUNTRY CHARACTERISTICS: Closeness to markets, Membership to trade agreements, other multilateral support, etc..

# 4 – leather shoe industry

# Program theory as a sense-making framework

# Nested theories



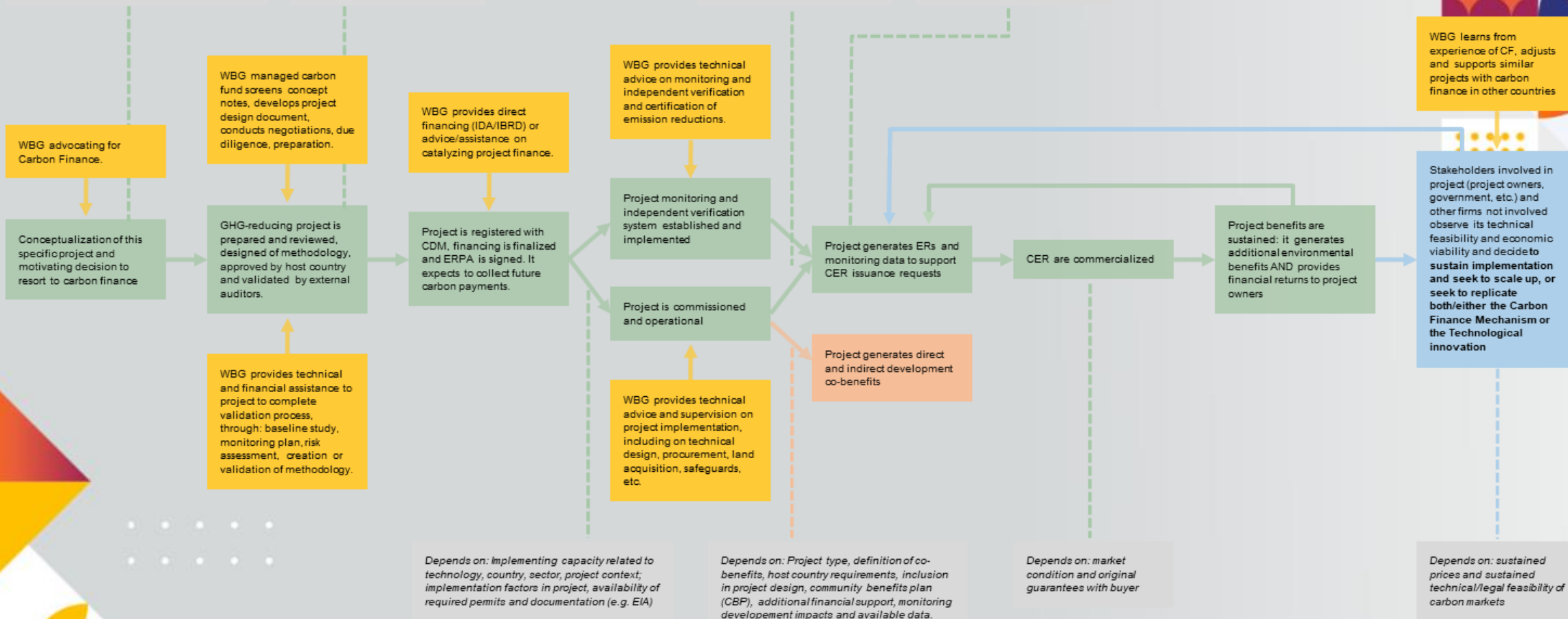Synthetic 'high-level' ToC



Nested 'detailed' ToC

# *Testable program theory*



Additionality assumptions: carbon finance addresses financial (investment rate of return) or non-financial barriers (technological barrier, others) or there is benefits to using carbon finance as a mechanisms.

Depends on: Type (risk) of projects, expected credit delivery timeline, technology (depending on the buyer), country risk

Depends on: actual project performance (e.g. actual power produced/saved, methane captured, etc.)

Depends on: actual credit issuance, abatement costs, and original methodology to compute additionality

WBG learns from experience of CF, adjusts and supports similar projects with carbon finance in other countries

WBG managed carbon fund screens concept notes, develops project design document, conducts negotiations, due diligence, preparation.

WBG provides direct financing (IDA/IBRD) or advice/assistance on catalyzing project finance.

WBG provides technical advice on monitoring and independent verification and certification of emission reductions.

WBG advocating for Carbon Finance.

Stakeholders involved in project (project owners, government, etc.) and other firms not involved observe its technical feasibility and economic viability and decide to **sustain implementation and seek to scale up, or seek to replicate both/either the Carbon Finance Mechanism or the Technological innovation**

Conceptualization of this specific project and motivating decision to resort to carbon finance

GHG-reducing project is prepared and reviewed, designed of methodology, approved by host country and validated by external auditors.

Project is registered with CDM, financing is finalized and ERPA is signed. It expects to collect future carbon payments.

Project monitoring and independent verification system established and implemented

Project generates ERs and monitoring data to support CER issuance requests

CER are commercialized

Project benefits are sustained: it generates additional environmental benefits AND provides financial returns to project owners

WBG provides technical and financial assistance to project to complete validation process, through: baseline study, monitoring plan, risk assessment, creation or validation of methodology.

Project is commissioned and operational

Project generates direct and indirect development co-benefits

WBG provides technical advice and supervision on project implementation, including on technical design, procurement, land acquisition, safeguards, etc.

Depends on: Implementing capacity related to technology, country, sector, project context; implementation factors in project, availability of required permits and documentation (e.g. EIA)

Depends on: Project type, definition of co-benefits, host country requirements, inclusion in project design, community benefits plan (CBP), additional financial support, monitoring developement impacts and available data.

Depends on: market condition and original guarantees with buyer

Depends on: sustained prices and sustained technical/legal feasibility of carbon markets

# "Good" program theory

- What is "good" program theory depends on the purpose of the program theory in the evaluation

- Good principles of a *testable* program theory in evaluation:
  - Be specific
  - Be consistent in formulations
  - Think about the warrants (i.e. is it logical to expect that a contributes to b)
  - Think about the underlying assumptions (i.e. under what conditions is a likely to contribute to b)

# Two broad strategies for reconstructing program theories

- 'Right to left': from objectives (or "problems") to underlying causes to activities/outputs

- 'Left to right': from activities/outputs to direct outcomes to indirect outcomes (objectives)

- Or combination

# Program theory reconstruction for evaluation (main sources)

- Intervention-related documents (policy, strategy, project ; design, monitoring, supervision, research,…..)

- Interviews with stakeholders (funders, implementing agencies, beneficiaries,…..)

- Existing knowledge (documentation) about similar (types of) interventions (broader literature ; policy/grey literature, academic literature,…..)

# Frameworks for reconstructing program theory

- **Policy instruments:** sticks, carrots, sermons (Bemelmans-Videc et al., 2003)

- **Behavioral mechanisms**: social norms, profit-seeking behavior, demonstration and copying behavior, peer pressure, etc.

- Coleman's **Theory of Social Action** (1986)
  - Situational mechanisms
  - Action-formation mechanisms
  - Transformational mechanisms

- **Intervention-specific** templates for program theory

# Looking at interventions across…

Portfolio-level: GEF-funded activities directed at rural landowners

# Focus on behavioral mechanisms

- There is no grand theory of social change, only **patterns of regularity** (Merton, 1967; Elster, 1989; Pawson and Tilley, 1997; Hedström and Swedberg,1998; Astbury and Leeuw, 2010)

- Describing patterns of change in terms of **mechanisms, contextual variables and outcomes**

- **Generative causality:** under what conditions does an intervention trigger a response (mechanism) that results in particular outcome

# Impact theory – microcredit

Based on Coleman (1986, 1990); Hedström and Swedberg (1998), see also Leeuw (2008)

# Intervention-specific templates for program theory



## Simplified Theory of Change Training

E.g. Incentives & Resources

E.g. Staff Turnover

Identification of Training Needs → Design Training → Implementation Training → Participants Learn → Participants Change Behavior → Organization Benefits

**Issues**

Relevance of Training …

Content and Quality of Curriculum …

Quality of Delivery …

Knowledge uptake …

Practices of Planning, Design, Research, Operations, Management …

Quality of Planning, Design, Research, Operations, Management …

Number and Type of Participants …

**Methods**

Specific methods and data sources differ according to causal step or underlying assumption

# Whose theory?

- Government, implementing organizations, beneficiaries (etc.) may have different expectations and assumptions regarding how an intervention is intended to work and what it may achieve

- Reconstructing different stakeholder theories is helpful in understanding the different views and beliefs of stakeholder groups

- Generating consensus on how an intervention is intended to work can be helpful in improving stakeholder relations and may benefit the intervention implementation process and subsequent benefits

# How you reconstruct program theory also depends on:

- The purpose of the evaluation (and the theory *of* evaluation)

  ➢ Goal-oriented (objectives-based) evaluation

    versus

  ➢ Goal-free evaluation

# Reconstructing a Program Theory (exercise and plenary)

# Group exercise

Read the case on the health sector intervention. You are then requested to work with your group on two tasks:

- Develop a program theory of the intervention.

- Identify to the extent possible (measurable) indicators relating to the different causal steps in the program theory.


- The necessary elements for the program theory are in the text. Indicators are not mentioned in the text but logically flow from the causal steps in the program theory. After the group work there will be a plenary discussion where each group will present its findings.

**Assumptions:** Socio-economic/economic/demographic factors (age, income, education, social integration, religion, family size); geographical factors; attitudes and beliefs (fear, aversity to risk); gender dimensions; political affiliation; sexual behavior (age of first sexual encounter, number of sex partners); influence of partners; seasonal influence.

► # coordination meetings; documentary evidence of improved coordination

MoH strengthens the coordination between Mobile Clinics and (regular, static) health care centers and hospitals

► # women and men who express change in SRH-related behavior

Women and men change their SRH-related behavior

► duration of campaigns; # and geographic coverage of campaigns

► # women and men who express awareness of SRH issues

► # women and men who express change in attitude toward using health services

► # women (and men) using SRH services

MoH organizes awareness campaigns on SRH

Women and men become more aware of SRH issues

Women and men change their health-seeking attitudes

Women (and men) use SRH services through the (regular, static) public and private health care system

Women are correctly treated for cervical cancer and other SRH issues

Women's health is improved (reduction in cervical cancer-related morbidity and mortality)

► # women correctly treated

MoH deploys Mobile Clinics

Women (and men) visit Mobile Clinics

Women (and men) receive more information on SRH issues, risks and treatments (including cervical cancer)

► # women (and men) provided with information sessions/materials

► morbidity rate; mortality rate

► # and geographic coverage of mobile clinics

► # women (and men) visiting mobile clinics

Women are examined for cervical cancer

► # women screened/examined

Women are correctly diagnosed and (where needed) referred for treatment

► # women correctly diagnosed and referred

► Indicator

**Assumptions:** Socio-economic/economic/demographic factors (age, income, education, social integration, religion, family size); geographical factors; attitudes and beliefs (fear, aversity to risk); gender dimensions; political affiliation; sexual behavior (age of first sexual encounter, number of sex partners); influence of partners; seasonal influence.

# Using Program Theory as a framework for evaluation

# Using program theory as a framework for evaluation

- Program theory is not 'method-specific'

- Program theory as a framework for particular assumptions being tested / refined, using:

  - (Quasi-)experimental techniques
  - Regression-based techniques
  - Descriptive and inferential statistical techniques
  - (Advanced) modelling approaches
  - Participatory techniques
  - Semi-structured interviews, open interviews, focus group interviews, discourse analysis, unobtrusive measures, etc.
  - Etc. etc.

# Evaluation of training in organic agriculture

# Evaluation of training in organic agriculture

- EU-supported rural development projects in 8 provinces
- In each of the provinces a national NGO provided training in organic agriculture
- In-depth evaluation (case study) of 1 out of 8 provinces
- Objective: assess implementation (participation), delivery of trainings and TA to farmers and outcomes

# **Multi-method approach**

- Review of project implementation reports

- In-depth interviews with EU project staff, NGO staff, farmers

- Review of training curriculum

- Observation of training sessions

- Farms visits to inspect land use practices

- Quasi-experimental design based on baseline and ex post survey

# Program theory



**Where do the data fit into the theory?**

# Addressing the attribution issues: a quasi-experiment

start | participants

CHANGE

end | participants ←→ control group

DIFFERENCE

# Data – outcomes

| practice | participants start | participants end | control group end |
|---|---|---|---|
| burning crop residues (%) | 27 % ** | 2 % | 29 % ** |
| applying green material (%) | 25 % ** | 63 % | 18 % ** |
| 'chemical' fertilizers (%) | 96 % * | 79 % | 97 % * |
| 'organic' fertilizers (%) | 79 % [a] | 83 % | 18 % ** |
| ditches (%) | 56 % [a] | 73 % | 24 % ** |
| barriers (%) | 44 % [a] | 58 % | 21 % ** |
| minimum tillage (%) | nihil [b] | 54 % | nihil [b] |
| latrines (%) | 15 % ** | 56 % | 8 % ** |
| furnaces (%) | 60 % | 69 % | 34 % ** |
| pig sties (%) | 42 % | 60 % | 45 % |
| nurseries (%) | 33 % | 44 % | 3 % ** |
| medicinal plants (no. plants) | 3.2 (5.3) ** | 8.7 (7.0) | 3.2 (3.5) ** |
| crop diversity (no. crops) | 4.3 (1.7) * | 4.9 (2.4) | 3.2 (1.4) ** |
| fruit tree diversity (no. trees) | 4.8 (2.9) * | 6.2 (3.2) | 4.6 (2.3) ** |

# Evaluation of police literacy training

```
institutional development MoIA
        ↓
4 senior master trainers
        ↓
20 provincial master trainers
        ↓
500 facilitators
        ↓
10,000+ patrol men and women  ←  certification by MoE
```

# Evaluation focus

- Initial focus: effectiveness (of different facilitator incentives on quality )

- Initial (local) purpose: to inform donor, to inform process of harmonization and improvement of the effectiveness and sustainability of police literacy training in MoIA

- Revised focus: 'impact' (of participation in literacy training on literacy levels)

- Revised (local) purpose: to inform donor, to inform other offices with ongoing or potential projects on police literacy training

- Slightly different stakeholder audience

# APLS: 'piggy-backing' on recently collected data

- Construction of indices of numeracy, writing and reading skills based on the results of several tests administered to respondents
  - Numeracy: counting, number recognition, basic maths
  - Writing: dictation, form filling
  - Reading: ability to read, comprehension, speed

- Independent variables: literacy training, prior education, other individual and regional characteristics

- Insufficient explanatory variables for statistical matching → imperfect explanatory regression model: multinomial logistic regression model

APLS impact analysis

Illustration of results: multinomial logit regression – dependent variable reading literacy

| | Variables | P(Y=Low) | P(Y=Medium) | P(Y=High) |
|---|---|---|---|---|
| | Age | -0.001 | -0.006 | -0.009 |
| | | (0.006) | (0.010) | (0.009) |
| | Male | 0.347* | 0.884** | 0.863*** |
| | | (0.208) | (0.354) | (0.316) |
| | Single | 0.004 | 0.016 | 0.087 |
| | | (0.076) | (0.112) | (0.103) |
| | Years of schooling | 0.140*** | 0.309*** | 0.439*** |
| | | (0.010) | (0.011) | (0.011) |
| | Household size | -0.011* | -0.037*** | -0.023*** |
| | | (0.006) | (0.010) | (0.008) |
| | Rural | -0.428*** | -0.330** | -0.558*** |
| | | (0.104) | (0.148) | (0.149) |
| Job-related factors | Uniform police | 0.032 | -0.199** | -0.287*** |
| | | (0.059) | (0.088) | (0.084) |
| | Job duration (yrs) | -0.027** | 0.020 | 0.009 |
| | | (0.012) | (0.018) | (0.017) |
| Literacy-related factors | Attendance (months) | 0.047*** | 0.074*** | 0.075*** |
| | | (0.006) | (0.007) | (0.008) |
| | Literacy training | 0.697*** | 1.060*** | 0.284*** |
| | | (0.073) | (0.113) | (0.100) |
| Lang (Dari/Pashto) | Both | -0.239 | -0.080 | -0.628** |
| | | (0.200) | (0.331) | (0.289) |
| | Only 1 of them | -0.323 | 0.031 | -0.554* |
| | | (0.197) | (0.326) | (0.284) |
| | One of them + other | -0.292 | -0.261 | -0.943*** |
| | | (0.204) | (0.337) | (0.298) |
| | | -1.143*** | -3.564*** | -2.689*** |
| | | (0.357) | (0.582) | (0.526) |
| | Observations | 8083 | 8083 | 8083 |
| | Log likelihood | -8165 | -8165 | -8165 |
| | Pseudo R$^2$ | 0.188 | 0.188 | 0.188 |

Standard errors in parentheses
Statistical Significance : *** p<0.01, ** p<0.05, * p<0.1

APLS impact analysis

Illustration of results: multinomial logit regression – dependent variable reading literacy

| | Variables | Total Sample | | |
|---|---|---|---|---|
| | | (1) P(Y=Low) | (2) P(Y=Medium) | (3) P(Y=High) |
| | Age | -0.001 | -0.006 | -0.009 |
| | | (0.006) | (0.010) | (0.009) |
| | Male | 0.347* | 0.884** | 0.863*** |
| | | (0.208) | (0.354) | (0.316) |
| | Single | 0.004 | 0.016 | 0.087 |
| | | (0.076) | (0.112) | (0.103) |
| | Years of schooling | 0.140*** | 0.309*** | 0.439*** |
| | | (0.010) | (0.011) | (0.011) |
| | Household size | -0.011* | -0.037*** | -0.023*** |
| | | (0.006) | (0.010) | (0.008) |
| | Rural | -0.428*** | -0.330** | -0.558*** |
| | | (0.104) | (0.148) | (0.149) |
| Job-related factors | Uniform police | 0.032 | -0.199** | -0.287*** |
| | | (0.059) | (0.088) | (0.084) |
| | Job duration (yrs) | -0.027** | 0.020 | 0.009 |
| | | (0.012) | (0.018) | (0.017) |
| Literacy-related factors | Attendance (months) | 0.047*** | 0.074*** | 0.075*** |
| | | (0.006) | (0.007) | (0.008) |
| | Literacy training | 0.697*** | 1.060*** | 0.284*** |
| | | (0.073) | (0.113) | (0.100) |
| Lang (Dari/Pashto) | Both | -0.239 | -0.080 | -0.628** |
| | | (0.200) | (0.331) | (0.289) |
| | Only 1 of them | -0.323 | 0.031 | -0.554* |
| | | (0.197) | (0.326) | (0.284) |
| | One of them + other | -0.292 | -0.261 | -0.943*** |
| | | (0.204) | (0.337) | (0.298) |
| | | -1.143*** | -3.564*** | -2.689*** |
| | | (0.357) | (0.582) | (0.526) |
| | Observations | 8083 | 8083 | 8083 |
| | Log likelihood | -8165 | -8165 | -8165 |
| | Pseudo R$^2$ | 0.188 | 0.188 | 0.188 |

Standard errors in parentheses
Statistial Significance : *** p<0.01, ** p<0.05, * p<0.1

APLS impact analysis

Illustration of results: multinomial logit regression – dependent variable numeracy literacy

|  | Variables | Total Sample | | |
|---|---|---|---|---|
|  |  | (1) P(Y=Low) | (2) P(Y=Medium) | (3) P(Y=High) |
|  | Age | -0.029*** | 0.000 | -0.001 |
|  |  | (0.008) | (0.008) | (0.009) |
|  | Male | -0.360 | 0.274 | 0.518* |
|  |  | (0.250) | (0.263) | (0.311) |
|  | Single | -0.052 | -0.184** | 0.032 |
|  |  | (0.094) | (0.093) | (0.105) |
|  | Years of schooling | 0.118*** | 0.276*** | 0.463*** |
|  |  | (0.019) | (0.017) | (0.017) |
|  | Household size | -0.006 | 0.001 | 0.006 |
|  |  | (0.008) | (0.007) | (0.008) |
|  | Rural | -0.928*** | -0.489*** | -0.722*** |
|  |  | (0.129) | (0.116) | (0.142) |
| Job-related factors | Uniform police | 0.419*** | 0.036 | -0.354*** |
|  |  | (0.076) | (0.074) | (0.085) |
|  | Job duration (yrs) | 0.011 | -0.025* | 0.007 |
|  |  | (0.015) | (0.015) | (0.017) |
| Literacy-related factors | Attendance (months) | 0.050*** | 0.083*** | 0.098*** |
|  |  | (0.012) | (0.011) | (0.012) |
|  | Literacy training | 0.690*** | 1.142*** | 1.172*** |
|  |  | (0.100) | (0.096) | (0.108) |
| Lang. (Dari/Pashto) | Both | 0.343 | -0.250 | 0.079 |
|  |  | (0.286) | (0.247) | (0.292) |
|  | Only 1 of them | 0.500* | -0.030 | -0.013 |
|  |  | (0.283) | (0.243) | (0.289) |
|  | One of them + other | 0.415 | 0.266 | 0.233 |
|  |  | (0.295) | (0.254) | (0.301) |
|  |  | 0.191 | -0.743* | -2.126*** |
|  |  | (0.465) | (0.442) | (0.521) |
|  | Observations | 8083 | 8083 | 8083 |
|  | Log likelihood | -9079 | -9079 | -9079 |
|  | $X^2$ - test | $X^2_{(33)}$=3894 | | |
|  | Pseudo $R^2$ | 0.175 | 0.184 | 0.184 |

Standard errors in parentheses

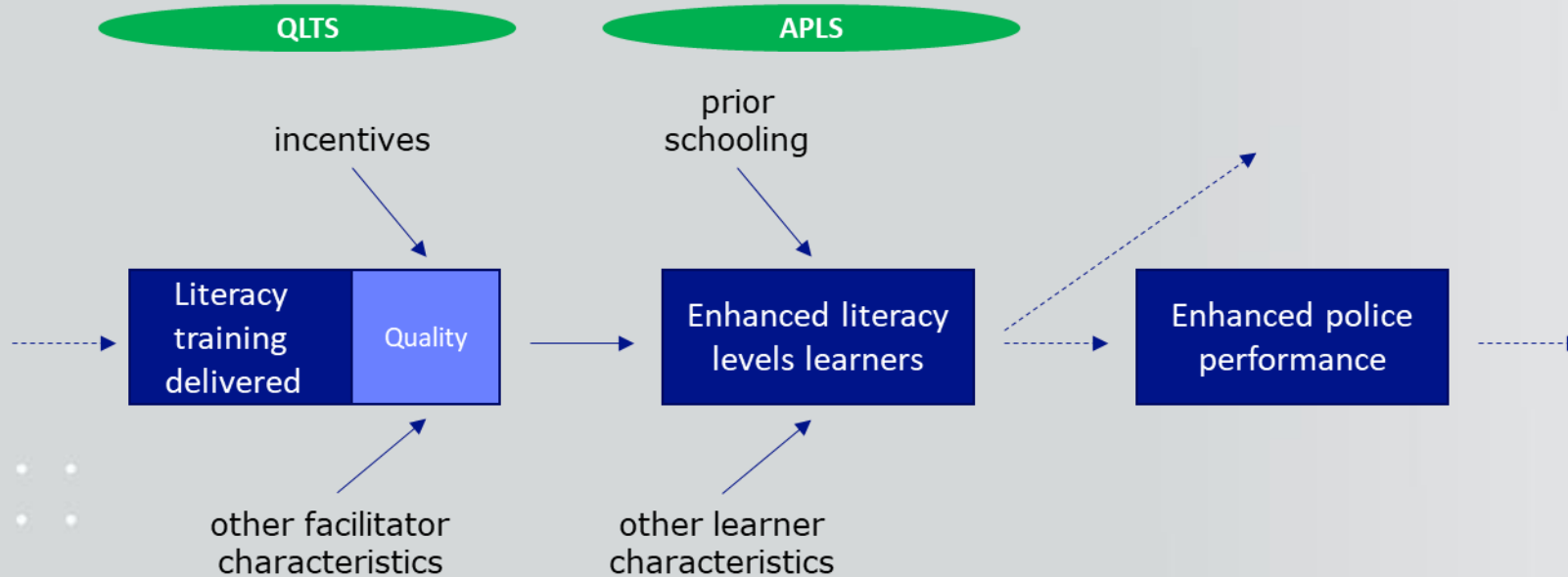Statistical Significance : *** p<0.01, ** p<0.05, * p<0.1

# Key conclusion

Despite some limitations in the data, the findings of the APLS impact analysis suggest that police literacy trainings have significantly improved literacy levels among the Afghan Police in all three dimensions of literacy (numeracy, reading, writing), controlling for other factors such as prior education and other individual and regional explanatory variables
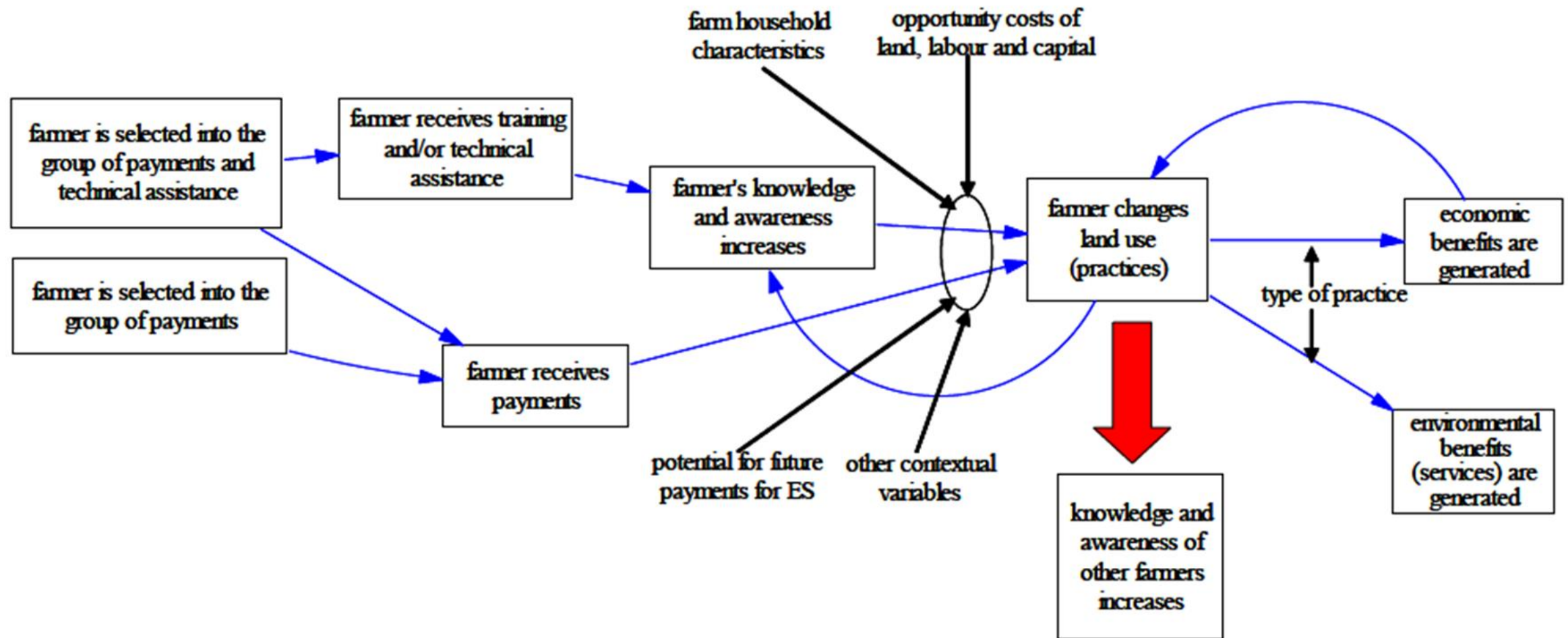
# Is that enough?

# Analysis shows how two complementary inquiries on two steps in the causal chain can enhance validity of causal claim

# Importance of a mixed methods approach: the logic of comparative advantages

- The randomized experiment can test the effectiveness of different incentives (PES and TA) on LU changes (from remote sensing data) and subsequently the environmental (from ES index calculations, remote sensing data) and socio-economic (from survey data) effects of these changes *(internal validity)*

- Survey data ('sub-group') analysis and case studies can tell us how incentives have *different* effects (knowledge, adoption) on *particular types* of farm households *(strengthens internal validity and increases external validity of findings)*

- Direct observation in selected sites, semi-structured interviews and focus group conversations can tell us more about the nature of effects in terms of production, consumption, poverty alleviation, etc. *(internal validity and construct validity)* as well as possible unintended effects (e.g. spillover effects, displacement effects)

# Some results: PES group – control group

# Grounded theory using social network analysis



Knowledge leadership in the Health Sector
in Liberia



Financial Flows in the Health Sector in
Liberia

# Deductive and inductive approach



1. How did outreach evolve? Was there increased outreach among the rural poor?

2. What are the factors that explain outreach/access?

3. What are the implications for poverty alleviation?

# Keep in mind the following:

- Fit for purpose
- Sources of theory
- Principles for developing a *"testable"* program theory
- Objectives-based evaluation and unintended effects
- Intervention-centric bias
- Confirmation bias

# THANK YOU