# Teaching Survey Research in Software Engineering

**Authors:**
Marcos Kalinowski (Pontifical Catholic University of Rio de Janeiro)
Allysson Allex Araújo (Federal University of Cariri)
Daniel Mendez (Blekinge Institute of Technology, fortiss)

# Agenda

1. Course Syllabus

2. Characteristics and Purpose of Survey Research (LO1)

3. Designing and Evaluating Survey Instruments (LO2)

4. Sampling and Data Collection (LO3)

5. Statistical and Qualitative Analysis (LO4)

# Agenda

**6**    Threats to Validity and Reliability (LO5)

**7**    Ethical Considerations

**8**    Concluding Remarks

# 1) Course Syllabus

# Course syllabus

- The course provides a comprehensive overview of **survey research principles and practices**.
  - How to **design** and **evaluate** survey instruments, focusing on aligning them with research objectives and relevant theories.
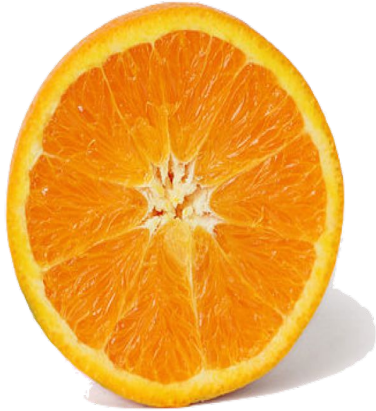
| ID | Learning Objective | Students will be able to ... | Bloom's Taxonomy |
|---|---|---|---|
| LO1 | Understanding the Characteristics and Purposes of Survey Research | ... articulate on the characteristics and purposes of survey research. <br> ... provide survey research application examples. | Remembering & Understanding |
| LO2 | Designing and Evaluating Survey Instruments | ... create survey instruments aligning with specific research objectives and theories. <br> ... critically assess the effectiveness of survey instruments. | Evaluating & Creating |
| LO3 | Mastering Sampling and Data Collection | ... apply best practices in sampling and data collection. <br> ... understand the trade-offs of different sampling and data collection methods. | Understanding & Applying |
| LO4 | Applying Statistical and Qualitative Analysis Methods | ... utilize statistical and qualitative analysis techniques to interpret survey data. | Applying & Analyzing |
| LO5 | Identifying and Addressing Validity and Reliability Threats | ... analyze and address potential threats to the validity and reliability of survey research. | Analyzing & Evaluating |
| LO6 | Understanding Ethical Considerations in Survey Research | ... identify, understand, and apply ethical considerations in survey research. | Understanding & Applying |

**Table 1** Learning Objectives and Bloom's Taxonomy Levels.

# 2) Characteristics and Purpose of Survey Research (LO1)

# Characteristics

**Survey** is an observational method to gather qualitative and/or quantitative data from (a sample of) entities to characterize information, attitudes and/or behaviors from different groups of subjects regarding an object of study.

Surveys
(Cross-sectional)

Case Studies

Experiments
(Case-control)

# Characteristics

- Surveys are probably the **most commonly used research method** worldwide.

- Surveys are conducted when a phenomena (*e.g.*, the use of a technique or tool) **already has taken place** or **before it occurs.**
  - A *survey* provides **no control of the execution or measurement.**
    - *I.e.*, it is not possible to manipulate variables as in the other investigation methods
  - Surveys should aim at obtaining the largest amount of understanding from the fewest number of variables since this reduction also eases the data collection and analysis.

- Surveys are almost never conducted to create an understanding concerning a particular **sample**, the typical focus is on generalizing results to the **population** from which the sample was drawn.
  - Surveys can be **retrospective** (looking back at something that has already happened) or **prospective** (looking ahead to something that is expected to happen)

# Characteristics

Unlike controlled experiments, surveys do not allow for **control over variables** or direct manipulation of the environment.

The observational nature of survey research often leads to challenges in establishing **causality**.

Design surveys to maximize understanding from a **minimal set of variables**.

# Purpose

- General objectives for conducting a survey (Wohlin et al., 2012; Wagner et al., 2020):

| EXPLORATIVE SURVEYS | DESCRIPTIVE SURVEYS | EXPLANATORY SURVEYS |
|---|---|---|
| are used as a pre-study to a more thorough investigation to assure that important issues are not forgotten (e.g., <u>constructs</u> in a theory like requirements elicitation techniques) | can be conducted to enable assertions about some population like the distribution of certain attributes (e.g., usage of requirements elicitation techniques) | aim at making explanatory claims about the population (e.g., why specific requirements elicitation techniques are used in specific contexts) |

# Characteristics and Purpose

"

*Theory building and evaluation can guide the design and analysis of surveys, and surveys can also be applied to test theories.*

*(Wagner et al., 2020)*

"

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. **Challenges in survey research.** In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Examples of Surveys

**Naming the pain in requirements engineering**
**Contemporary problems, causes, and effects in practice**

D. Méndez Fernández[1] · S. Wagner[2] · M. Kalinowski[3] · M. Felderer[4] ·
P. Mafra[3] · A. Vetrò[5] · T. Conte[6] · M.-T. Christiansson[7] · D. Greer[8] ·
C. Lassenius[9] · T. Männistö[10] · M. Nayabi[11] · M. Oivo[12] · B. Penzenstadler[13] ·
D. Pfahl[14] · R. Prikladnicki[15] · G. Ruhe[11] · A. Schekelmann[16] · S. Sen[17] ·
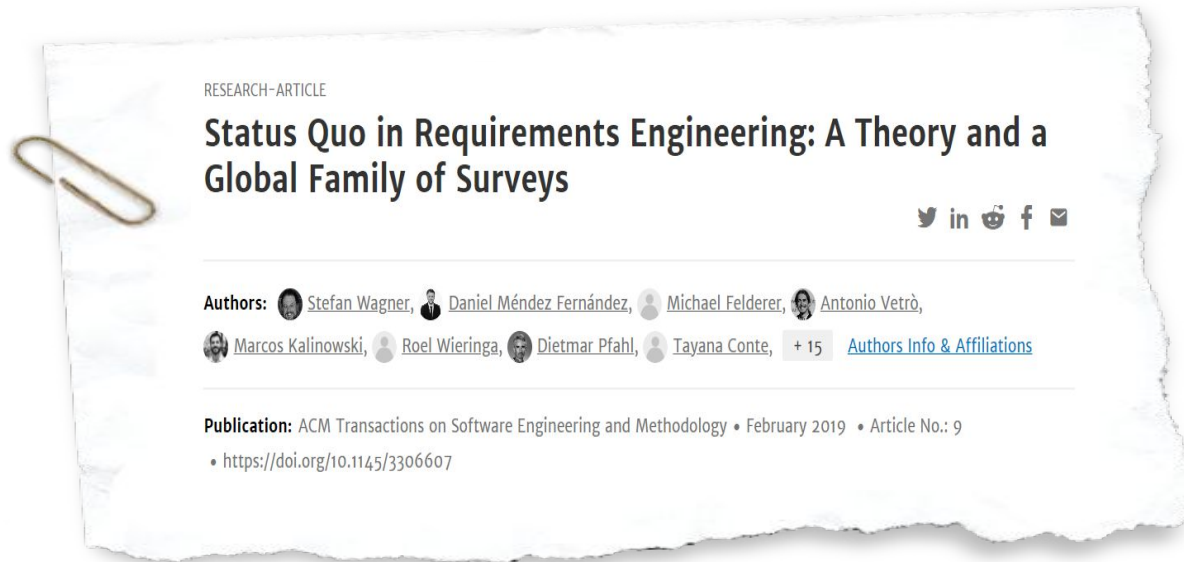R. Spinola[18,19] · A. Tuzcu[20] · J. L. de la Vara[21] · R. Wieringa[22]

Naming the **P**ain in **R**equirements **E**ngineering
NaPiRE

Fernández, D. M.; Wagner, S.; Kalinowski, M.; Felderer, M.; Mafra, P.; Vetro, A.; Conte, T.; Christiansson, M.; Greer, D.; Lassenius, C.; Männistö, T.; Nayabi, M.; Oivo, M.; Penzenstadler, B.; Pfahl, D.; Prikladnicki, R.; Ruhe, G.; Schekelmann, A.; Sen, S.; Spínola, R. O.; Tuzcu, A.; de la Vara, J. L.; and Wieringa, R. Naming the pain in requirements engineering - Contemporary problems, causes, and effects in practice. Empirical Software Engineering, 22(5): 2298-2338. 2017.

# Examples of Surveys

RESEARCH-ARTICLE

## Status Quo in Requirements Engineering: A Theory and a Global Family of Surveys

Authors: Stefan Wagner, Daniel Méndez Fernández, Michael Felderer, Antonio Vetrò, Marcos Kalinowski, Roel Wieringa, Dietmar Pfahl, Tayana Conte, +15  Authors Info & Affiliations

Wagner, S., Fernández, D. M., Felderer, M., Vetro, A., Kalinowski, M., Wieringa, R., Pfahl, D., Conte, T., Christiansson, M., Greer, D., Lassenius, C., Männistö, T., Nayebi, M., Oivo, M., Penzenstadler, B., Prikladnicki, R., Ruhe, G., Schekelmann, A., Sen, S., Spínola, R.O., Tuzcu, A., de la Vara, J. L., and Winkler, D, Status Quo in Requirements Engineering: A Theory and a Global Family of Surveys.  ACM Transactions on Software Engineering and Methdology, 28(2): 9:1-9:48. 2019.

# Examples of Surveys

Check for updates

## Pandemic programming

### How COVID-19 affects software developers and how their organizations can help

Paul Ralph[1] · Sebastian Baltes[2] · Gianisa Adisaputri[1] · Richard Torkar[3,4] ·
Vladimir Kovalenko[5] · Marcos Kalinowski[6] · Nicole Novielli[7] · Shin Yoo[8] ·
Xavier Devroey[9] · Xin Tan[10] · Minghui Zhou[10] · Burak Turhan[11,12] · Rashina Hoda[11] ·
Hideaki Hata[13] · Gregorio Robles[14] · Amin Milani Fard[15] · Rana Alkadhi[16]

Ralph, P., Baltes, S., Adisaputri, G., Torkar, R., Kovalenko, V., Kalinowski, M., Novielli, N., Yoo, S., Devroey, X., Tan, X., Zhou, M., Turhan, B., Hoda, R., Hata, H., Robles, G., Fard, A. M., and Alkadhi, R, Pandemic Programming How COVID-19 affects software developers and how their organizations can help. Empirical Software Engineering (2020), 25: 4927-4961. 2020.

# Examples of Surveys

## What Makes Agile Software Development Agile?

Marco Kuhrmann, Paolo Tell, Regina Hebig, Jil Klünder, Jürgen Münch, Oliver Linssen, Dietmar Pfahl, Michael Felderer, Christian R. Prause, Stephen G. MacDonell, Joyce Nakatumba-Nabende, David Raffo, Sarah Beecham, Eray Tüzün, Gustavo López, Nicolas Paez, Diego Fontdevila, Sherlock A. Licorish, Steffen Küpper, Günther Ruhe, Eric Knauss, Özden Özcan-Top, Paul Clarke, Fergal McCaffery, Marcela Genero, Aurora Vizcaino, Mario Piattini, Marcos Kalinowski, Tayana Conte, Rafael Prikladnicki, Stephan Krusche, Ahmet Coşkunçay, Ezequiel Scott, Fabio Calefato, Svetlana Pimonova, Rolf-Helge Pfeiffer, Ulrik Pagh Schultz, Rogardt Heldal, Masud Fazal-Baqaie, Craig Anslow, Maleknaz Nayebi, Kurt Schneider, Stefan Sauer, Dietmar Winkler, Stefan Biffl, Maria Cecilia Bastarrica, and Ita Richardson

Kuhrmann, M., Tell, P., Hebig, R. et al. What Makes Agile Software Development Agile? Submitted to Transactions on Software Engineering (2021).

# Key Takeaways on Characteristics and Purpose of Survey Research

Characteristics of survey research methods, including strengths and limitations.

General objectives that surveys can fulfill.

# 3) Designing and Evaluating Survey Instruments (LO2)

# Survey Design

Basics of Survey Design

Goal-Question-Metric-Driven Design

Theory-Driven Design

Issues When Assessing Psychological Constructs

Survey Instrument Evaluation

# Basics of Survey Design

## QUESTIONNAIRE TYPES

✔ Self-administered questionnaire
✔ Interviewer-administered questionnaire

## QUESTION TYPES

✔ Open-ended
✔ Closed-ended
✔ Hybrid questions

## QUESTION CATEGORIES

✔ Demographic questions
✔ Substantive questions
✔ Filter questions
✔ Sensitive questions

# Basics of Survey Design

## Measurement scales

**More Information** (downward arrow)

**Nominal**
- Values can be counted

**Ordinal**
- Values can be counted and ordered

**Interval**
- Values can be counted and ordered
- Distance between values can be interpreted

**Ratio**
- Values can be counted and ordered
- Distance between values can be interpreted
- Radio between values can be interpreted

## Conditions that must be fulfilled to get appropriate responses

Questions must be understandable by the target population

Respondents must have sufficient knowledge to answer

Participants must be motivated and willing to participate

# Basics of Survey Design

**Suggestions to avoid common question <u>wording problems</u> (adapted from Kitchenham and Pfleeger, 2008)**

- ✔ Using appropriate and simple language
- ✔ Avoiding technical terms
- ✔ Keeping questions short
- ✔ Avoiding vague sentences
- ✔ Avoiding sensitive questions
- ✔ Avoiding too demanding questions
- ✔ Avoiding double-barreled questions
- ✔ Avoiding double negatives
- ✔ Avoid asking about long gone events

In a survey, we can either ask for the **opinions** of the participants on topics or for specific **facts** that they experienced.

# Basics of Survey Design

A very simplified process for survey research:

**Survey Planning**

| |
|---|
| Characterising Target Population |
| Sampling |
| Questionnaire Design |
| Recruiting & Measuring |

**Survey Execution**

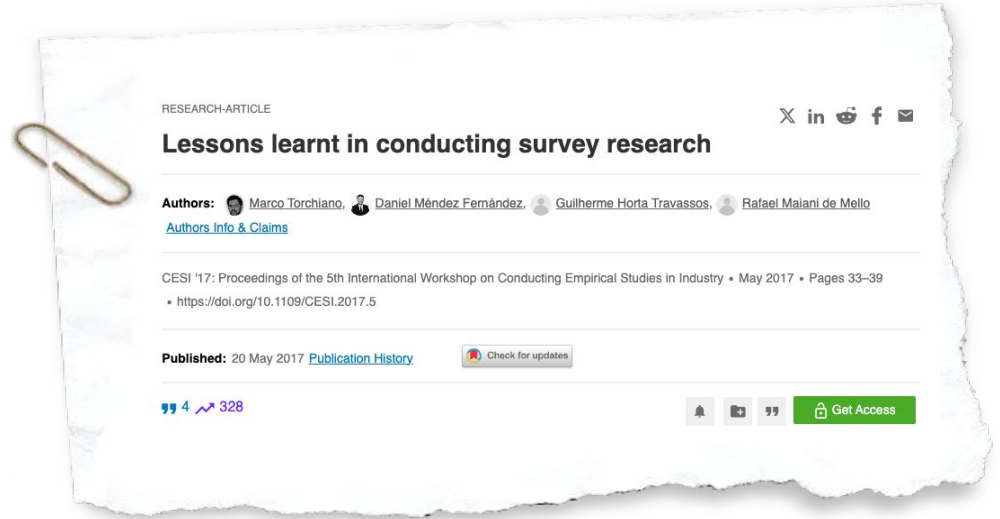| |
|---|
| Data Coding & Editing |
| Post-survey Adjustments |
| Data Analysis & Interpretation |

**Packaging & Reporting**

| |
|---|
| Data Curation & Disclosure |

# Survey Instrument Evaluation Methods

There are too many pitfalls to be handled. For further information, see the work of Torchiano *et al.* about lessons learnt in conducting survey research.
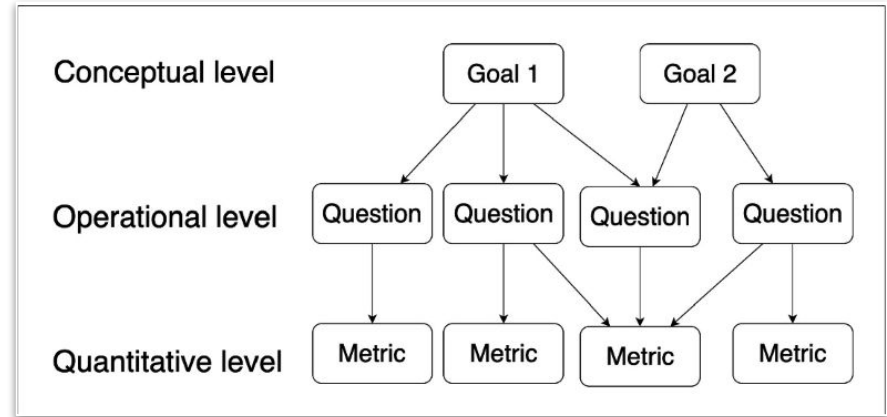
Torchiano, M., Méndez Fernández, D., Travassos, G.H., de Mello, R. M. (2017). Lessons Learnt in Conducting Survey Research. In: Proc. 5th International Workshop on Conducting Empirical Studies in Industry (CESI). ICSE 2017.

# Goal-Question-Metric-Driven Design

Based on the **Goal Question Metric (GQM)** Paradigm *(Basili and Romback, 1988)*

GQM defines a way to plan and execute measurement and analysis activities:

**1** Starts with the declaration of the measurement, **Goals**

**2** From the goals, **Questions** that we would like to answer with the data interpretation are defined

**3** Finally, from the questions, the **Metrics** and the data to be collected are defined

Basili, V.R. and Rombach, H.D., 1988. The TAME project: Towards improvement-oriented software environments. IEEE Transactions on software engineering, 14(6), pp.758-773.

# Goal Definition Template

Analyze **<object of study>**

with the purpose of **<goal>**

with respect to **<quality focus>**

from the point of view of the **<perspective>**

in the context of **<context>**

Measurement activities need <u>clear goals</u>
*GQM: characterize, understand, evaluate, predict, improve.*

Basili, V.R. and Rombach, H.D., 1988. The TAME project: Towards improvement-oriented software environments. IEEE Transactions on software engineering, 14(6), pp.758-773.

# Goal Definition Template (Example)

Analyze **the profile of software development organizations**

with the purpose of **characterizing**

with respect to **the organizations' current profile, satisfaction degree regarding the MPS model, variation of presence in international markets, variation of exportation volume, and variation concerning cost, estimation accuracy, productivity, quality, user satisfaction, and return of investment (ROI)**

from the point of view **the software development organizations**

in the context of s**oftware development organizations with unexpired MPS-SW assessments published in the SOFTEX portal.**

Kalinowski, M., Weber, K.C. and Travassos, G.H., 2008, October. iMPS: an experimentation based investigation of a nationwide software development reference model. In Proceedings of the Second ACM-IEEE international symposium on Empirical Software Engineering and Measurement (ESEM).

# Further Goal-Question-Metric-Driven Design Examples

"Analyze **Social BPM** with the purpose of **characterizing** with respect to **adoption of its practices and technologies during the BPM lifecycle** from the point of view of **BPM participants or managers** In the context of **Brazilian organizations**."

Batista, M., Magdaleno, A. and Kalinowski, M., 2017, May. A Survey on the use of Social BPM in Practice in Brazilian Organizations. In Anais do XIII Simpósio Brasileiro de Sistemas de Informação (SBSI) (pp. 436-443). SBC.

"Analyze **V&V methods** with the purpose of **characterization** with respect to their **suitability for addressing ISO 25010 software quality characteristics** from the point of view of **experts in the area of V&V** in the context of the **software engineering research community**."

Mendoza, I., Kalinowski, M., Souza, U. and Felderer, M., 2019, January. Relating verification and validation methods to software product quality characteristics: results of an expert survey. In Proc. of the Software Quality Days Conference (SWQD) (pp. 33-44).

# Goal Definition Template (Example)

## GOAL

Analyze **software development organizations**

with the purpose of **characterizing**

with respect to **the organizations' current profile, satisfaction degree regarding the MPS model, variation of presence in international markets, variation of exportation volume, and variation concerning cost, estimation accuracy, productivity, quality, user satisfaction, and return of investment (ROI)**

from the point of view **the software development organizations**

in the context of **software development organizations with unexpired MPS-SW assessments published in the SOFTEX portal**

Kalinowski, M., Weber, K.C. and Travassos, G.H., 2008, October. iMPS: an experimentation based investigation of a nationwide software development reference model. In Proceedings of the Second ACM-IEEE international symposium on Empirical Software Engineering and Measurement (ESEM).

# Goal Definition Template (Example)

## QUESTION

**Q1: What is the organization's estimation accuracy?**

## METRICS

**M1.1**: Average Project Duration = Average duration of projects conducted within the last 12 months, measured in months.

**M1.2**: Average Project Estimated Duration = Average estimated duration of projects conducted within the last 12 months, measured in months.

**M1.3**: Estimation Accuracy = 1 - |((Average Project Duration – Average Project Estimated Duration) / Average Project Duration)|

Kalinowski, M., Weber, K.C. and Travassos, G.H., 2008, October. iMPS: an experimentation based investigation of a nationwide software development reference model. In Proceedings of the Second ACM-IEEE international symposium on Empirical Software Engineering and Measurement (ESEM).

# Goal Definition Template (Example)

## QUESTION

**Q2: What is the organization's Return of Investment (ROI) of adopting MPS-SW?**

## METRICS

**M2.1**: Variation in net sales = Percentage of variation in net sales.

**M2.2**: Investment in implementing MPS = Percentage of net sales invested in implementing MPS

**M2.3**: Investment in assessing MPS = Percentage of net sales invested in the MPS assessment

**M2.4**: ROI = (Variation in net sales / (Investment in implementing MPS + Investment in assessing MPS)) * 100

Kalinowski, M., Weber, K.C. and Travassos, G.H., 2008, October. iMPS: an experimentation based investigation of a nationwide software development reference model. In Proceedings of the Second ACM-IEEE international symposium on Empirical Software Engineering and Measurement (ESEM).

# Theory-Driven Survey Design

A theory provides **explanations** and **understanding** in terms of basic **constructs** and underlying **mechanisms**, which constitute an important counterpart to knowledge of passing trends and their manifestation *(Hannay et al. 2007)*:

- From the practical perspective, theories should **be useful and explain or predict phenomena** that occur in software engineering

- From a scientific perspective, theories should **guide and support further research** in software engineering

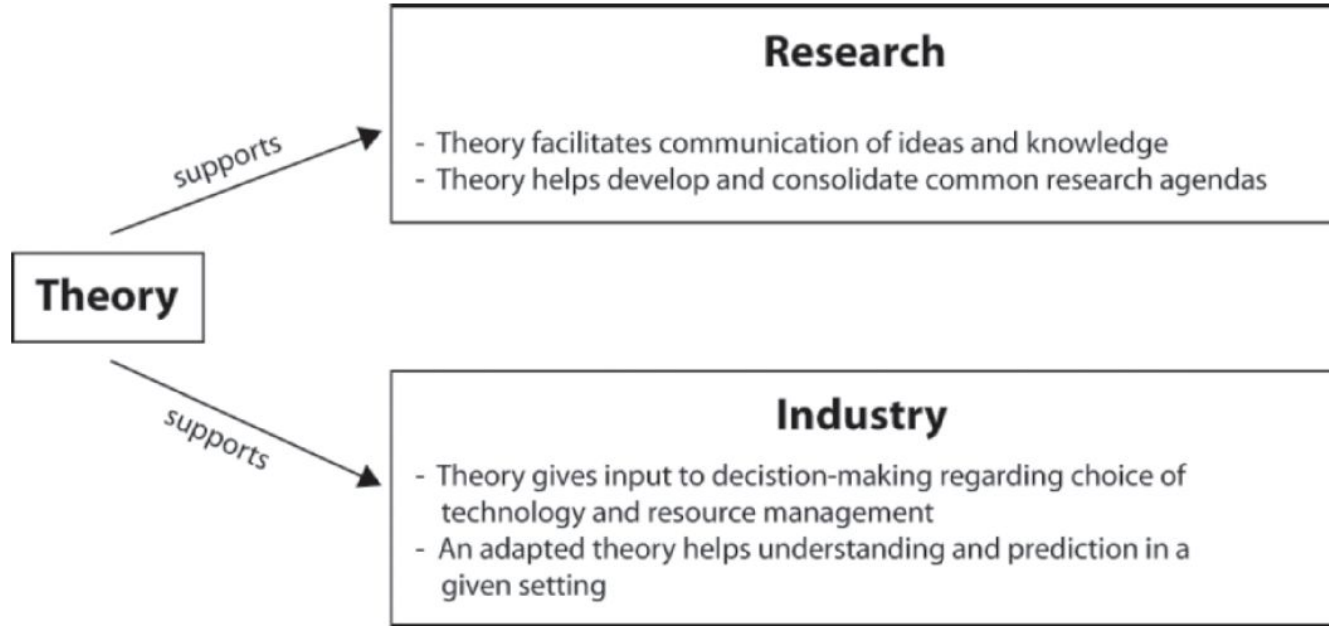**THEORY BUILDING BLOCKS**
*(Sjøberg et al., 2008)*

| Constructs |
| --- |
| **Propositions** |
| **Explanations** |
| **Scope** |

Basili, V.R. and Rombach, H.D., 1988. The TAME project: Towards improvement-oriented software environments. IEEE Transactions on software engineering, 14(6), pp.758-773.

# Theory-Driven Survey Design

Sjøberg, D.I., Dybå, T., Anda, B.C. and Hannay, J.E., 2008. Building theories in software engineering. In Guide to advanced empirical software engineering (pp. 312-336). Springer, London.

# Theory-Driven Survey Design

- **Theory building** and **survey research** are strongly interrelated;

- Initial theories can be drawn from **observations** and available **literature;**

- An initial theory may be a **taxonomy of constructs** or a **set of statements relating constructs:**

  - For NaPiRE, a set of constructs and propositions was elaborated based on available literature and expert knowledge,

  - For Pandemic Programming, a theoretical model was designed based on related work

  - The surveys, in both cases, were designed to test the theory (and to potentially extend it)
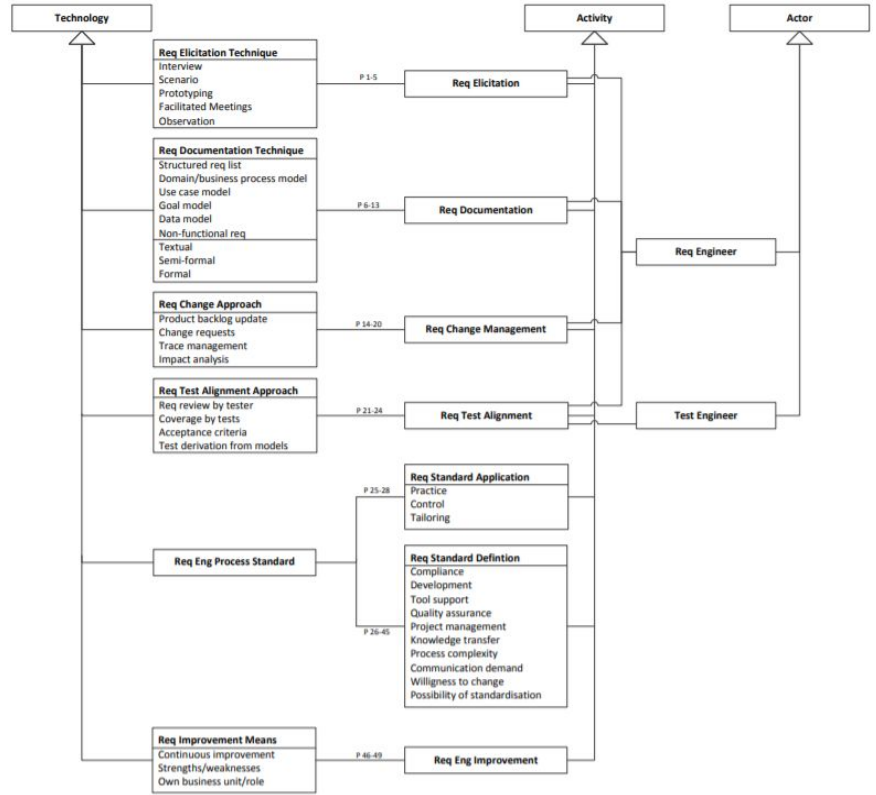
# Theory-Driven Survey Design: NaPIRE

**INITIAL THEORY**

| Constructs | | Type |
|---|---|---|
| C 1 | Requirements Elicitation | Activity |
| C 2 | Requirements Documentation | Activity |
| C 3 | Requirements Change Management | Activity |
| C 4 | Requirements Test Alignment | Activity |
| C 5 | Requirements Standard Application | Activity |
| C 6 | Requirements Standard Definition | Activity |
| C 7 | Requirements Engineering Improvement | Activity |
| C 8 | Requirements Engineer | Actor |
| C 9 | Test Engineer | Actor |
| C 10 | Requirements Elicitation Technique | Technology |
| C 11 | Requirements Documentation Technique | Technology |
| C 12 | Requirements Change Approach | Technology |
| C 13 | Requirements Test Alignment Approach | Technology |
| C 14 | Requirements Engineering Process Standard | Technology |
| C 15 | Requirements Improvement Means | Technology |

**Scope**

The theory is supposed to be applicable to contemporary requirements engineering in practice world-wide. There could be differences in different regions of the world because of cultural differences or different economic environments as well as differences in different application domains.



Wagner, S. et al. Status Quo in Requirements Engineering: A Theory and a Global Family of Surveys. ACM Transactions on Software Engineering and Methodology, 28(2): 9:1-9:48. 2019.

# Theory-Driven Survey Design: NaPIRE

| No. | Propositions |
|-----|--------------|
| P 1 | Requirements are elicited via interviews |
| P 2 | Requirements are elicited via scenarios |
| P 3 | Requirements are elicited via prototyping |
| P 4 | Requirements are elicited via facilitated meetings (including workshops) |
| P 5 | Requirements are elicited via observation |

| No. | Explanations | Propositions |
|-----|--------------|--------------|
| E 1 | Interviews, scenarios, prototyping, facilitated meetings, and observations allow the requirements engineers to include many different viewpoints including those from nontechnical stakeholders | P1–P5 |
| E 2 | Prototypes and scenarios promote a shared understanding of the requirements among stakeholders | P2, P3 |

# Theory-Driven Survey Design: NaPIRE

**DESIGNED QUESTIONNAIRE**

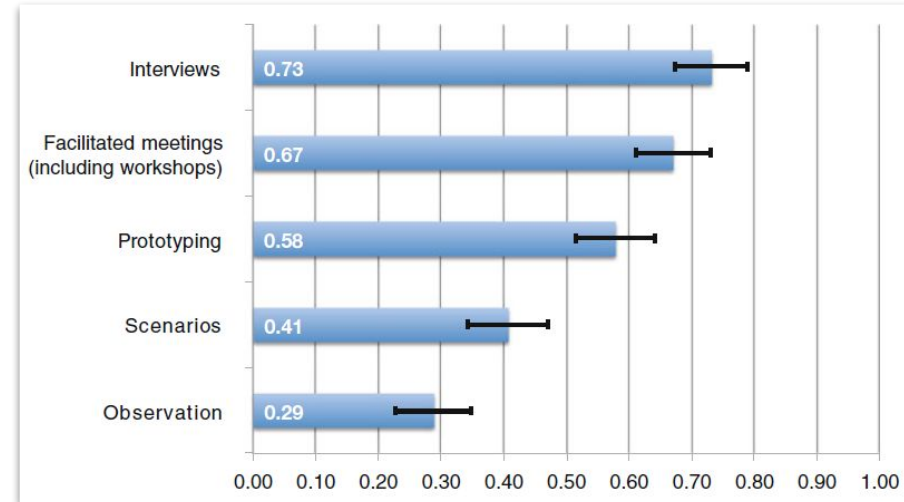| RQ | No. | Question | Type |
|---|---|---|---|
| – | Q 1 | What is the size of your company? | Closed(SC) |
| | Q 2 | Please describe the main business area and application domain. | Open |
| | Q 3 | Does your company participate in globally distributed projects? | Closed(SC) |
| | Q 4 | In which country are you personally located? | Open |
| | Q 5 | To which project role are you most frequently assigned? | Closed(SC) |
| | Q 6 | How do you rate your experience in this role? | Closed(SC) |
| | Q 7 | Which organisational role does your company take most frequently in your projects? | Closed(SC) |
| | Q 8 | Which process model do you follow (or a variation of it)? | Closed(MC) |
| RQ 1 | Q 9 | How do you elicit requirements? | Closed(MC) |
| | Q 10 | How do you document functional requirements? | Closed(MC) |
| | Q 11 | How do you document non-functional requirements? | Closed(SC) |
| RQ 2 | Q 21 | How do you perform change management in your requirements engineering? | Closed(MC) |
| | Q 12 | How do you deal with changing requirements after the initial release? | Closed(SC) |
| | Q 13 | Which traces do you explicitly manage? | Closed(MC) |
| | Q 14 | How do you analyse the effect of changes to requirements? | Closed(MC) |
| | Q 15 | How do you align the software test with the requirements? | Closed(MC) |
| RQ 3 | Q 16 | What RE standard have you established at your company? | Closed(MC) |
| | Q 17 | Which reasons do you agree with as a motivation to define a company standard for RE in your company? | Likert |
| | Q 18 | Which reasons do you see as a barrier to define a company standard for RE in your company? | Likert |
| | Q 19 | Is the requirements engineering standard mandatory and practised? | Closed(SC) |
| | Q 20 | How do you check the application of your requirements engineering standard? | Closed(MC) |
| | Q 22 | How is your RE standard applied (tailored) in your regular projects? | Closed(MC) |
| RQ 4 | Q 23 | Is your RE continuously improved? | Closed(SC) |
| | Q 24 | Why do you continuously improve your requirements engineering? | Closed(MC) |

**RQ 1** How are requirements elicited and documented?

**RQ 2** How are requirements changed and aligned with tests?

**RQ 3** How are RE standards applied and tailored?

**RQ 4** How is RE improved?

Wagner, S. et al. Status Quo in Requirements Engineering: A Theory and a Global Family of Surveys. ACM Transactions on Software Engineering and Methodology, 28(2): 9:1-9:48. 2019.

# Theory-Driven Survey Design: NaPIRE
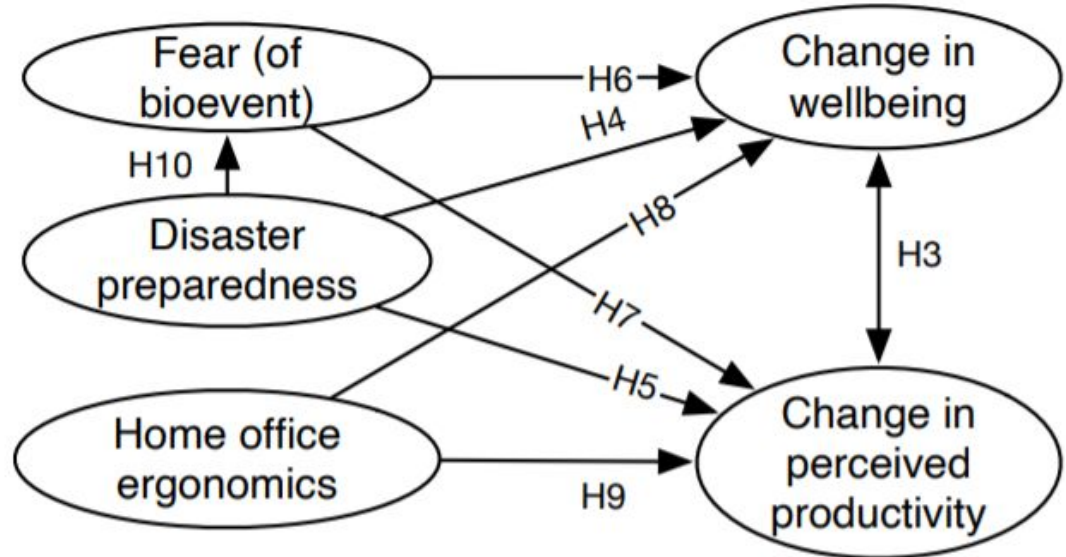
| No. | Propositions |
|---|---|
| P 1 | Requirements are elicited via interviews |
| P 2 | Requirements are elicited via scenarios |
| P 3 | Requirements are elicited via prototyping |
| P 4 | Requirements are elicited via facilitated meetings (including workshops) |
| P 5 | Requirements are elicited via observation |

| No. | Explanations | Propositions |
|---|---|---|
| E 1 | Interviews, scenarios, prototyping, facilitated meetings, and observations allow the requirements engineers to include many different viewpoints including those from nontechnical stakeholders | P1–P5 |
| E 2 | Prototypes and scenarios promote a shared understanding of the requirements among stakeholders | P2, P3 |



Wagner, S. et al. Status Quo in Requirements Engineering: A Theory and a Global Family of Surveys. ACM Transactions on Software Engineering and Methodology, 28(2): 9:1-9:48. 2019.

# Theory-Driven Survey Design: Pandemic Programming

Ralph, P., Baltes, S., Adisaputri, G., Torkar, R., Kovalenko, V., Kalinowski, M., Novielli, N., Yoo, S., Devroey, X., Tan, X., Zhou, M., Turhan, B., Hoda, R., Hata, H., Robles, G., Fard, A. M., and Alkadhi, R, Pandemic Programming How COVID-19 affects software developers and how their organizations can help.  Empirical Software Engineering (2020), 25: 4927-4961. 2020.

# Theory-Driven Survey Design: Pandemic Programming

## SELECTING VALIDATED SCALES FOR THE CONSTRUCTS

**Change in wellbeing**
We used the WHO's five-item wellbeing index (WHO-5)

**Change in perceived productivity**
We used items from the WHO's Health and Work Performance Questionnaire (HPQ)

**Disaster preparedness**
We adapted Yong et al.'s (2017) individual disaster preparedness scale

**Fear (of bioevent)**
We adapted the Bracha-Burkle Fear and Resilience (FR) checklist, a triage tool for assessing patients' reactions to bioevents (including pandemics).

**Home office ergonomics**
We could not find a reasonable scale. Based on our reading of the ergonomics literature, we made a simple six-item, six-point Likert scale concerning distractions, noise, lighting, temperature, chair comfort and overall ergonomics.

Ralph, P., Baltes, S., Adisaputri, G., Torkar, R., Kovalenko, V., Kalinowski, M., Novielli, N., Yoo, S., Devroey, X., Tan, X., Zhou, M., Turhan, B., Hoda, R., Hata, H., Robles, G., Fard, A. M., and Alkadhi, R, Pandemic Programming How COVID-19 affects software developers and how their organizations can help. Empirical Software Engineering (2020), 25: 4927-4961. 2020.

# Theory-Driven Survey Design: Pandemic Programming

**SUPPORTED MODEL**

Ralph, P., Baltes, S., Adisaputri, G., Torkar, R., Kovalenko, V., Kalinowski, M., Novielli, N., Yoo, S., Devroey, X., Tan, X., Zhou, M., Turhan, B., Hoda, R., Hata, H., Robles, G., Fard, A. M., and Alkadhi, R, Pandemic Programming How COVID-19 affects software developers and how their organizations can help. Empirical Software Engineering (2020), 25: 4927-4961. 2020.

# Evaluating Theories

**Table 1** Criteria for evaluating theories

| | |
|---|---|
| Testability | The degree to which a theory is constructed such that empirical refutation is possible |
| Empirical support | The degree to which a theory is supported by empirical studies that confirm its validity |
| Explanatory power | The degree to which a theory accounts for and predicts all known observations within its scope, is simple in that it has few ad hoc assumption, and relates to that which is already well understood |
| Parsimony | The degree to which a theory is economically constructed with a minimum of concepts and propositions |
| Generality | The breadth of the scope of a theory and the degree to which the theory is independent of specific settings |
| Utility | The degree to which a theory supports the relevant areas of the software industry |

Sjøberg, D.I., Dybå, T., Anda, B.C. and Hannay, J.E., 2008. Building theories in software engineering. In Guide to advanced empirical software engineering (pp. 312-336). Springer, London.

# Survey Research and Theory Building

**Key Takeaways** (Wagner et al., 2020):

**1**

Survey research and theory building are strongly interrelated. The exact relationship depends on whether the theory is descriptive, explanatory, or predictive.

**3**

Survey data supports the definition or refinement of constructs, relationships, explanations, and the scope of a theory as well as testing of a theory.

**2**

Theories are of high value to guide the design of surveys.

**4**

Use validated scales as much as possible to improve construct validity.

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. Challenges in survey research. In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Issues When Assessing Psychological Constructs

**Psychological constructs** are theoretical concepts to model and understand human behavior, cognition, affect, and knowledge (Binning, 2016)

Examples include **happiness**, **job satisfaction**, **motivation**, **commitment**, **personality**, **intelligence**, **skills**, and **performance**

These constructs can only be assessed **indirectly**

We need ways to **proxy** our **measurement** of a construct in robust, valid, and reliable ways

▫ This is why, whenever we wish to investigate psychological constructs and their variables, we need to either develop or adopt measurement instruments that are **psychometrically validated**

Scientists have investigated issues of **validity** and **reliability** of psychological tests

Binning JF (2016) Construct. https://www.britannica.com/science/construct

# Issues When Assessing Psychological Constructs

## Validity and Reliability in Psychometrics (AERA et al., 2014)

### VALIDITY

✔ The degree to which evidence and theory support the interpretation of test scores for proposed uses of tests

✔ We need to ensure that any meaning we provide to the values obtained by a measurement instrument needs to be validated

### RELIABILITY

✔ Consistency of a questionnaire score in repeated instances of it; or

✔ Coefficient between scores on two equivalent forms of the same test

AERA, APA, NCME: Standards for educational and psychological testing. American Educational Research Association, Washington, DC (2014)

# Issues When Assessing Psychological Constructs

> " *Software engineering research should favor **psychometric validation** of tests.*
>
> *(Wagner et al., 2020)* "

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. **Challenges in survey research.** In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Issues When Assessing Psychological Constructs

**Key Takeaways** (Wagner et al., 2020):

**1**

Representing and assessing constructs on human behavior, cognition, affect, and knowledge is a difficult problem that requires psychometrically validated measurement instruments.

**2**

Software engineering research should either adopt or develop psychometrically validated questionnaires.

**3**

Adoption or development of psychometrically validated questionnaires should consider psychometric validity and reliability issues, which are diverse and very different from the usual and common validity issues we see in "Threats to Validity" sections.

**4**

Software engineering research should introduce studies on the development and validation of questionnaires.

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. Challenges in survey research. In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Survey Instrument Evaluation Methods

- Used to assess the **validity** and **reliability** of the survey instrument;

- A survey can be evaluated, to avoid **threats to validity and reliability**, using the following methods *(Robson, 2002 apud Linaker et al., 2015)*:

EXPERT REVIEWS  FOCUS GROUPS  PILOT SURVEYS  COGNITIVE INTERVIEWS  EXPERIMENTS

Robson, C., (2002) Real World Research - A Resource for Social Scientists and Practitioner-Researchers, 2nd ed. Malden: Blackwell Publishing.

# Survey Instrument Evaluation Methods

Additionally, the empirically evaluated **checklist for surveys** in software engineering by Molléri et al. [30] can be used as an additional valuable resource for evaluating the survey design (as well as the final survey report).

Molléri, J.S., Petersen, K., Mendes, E.: An empirically evaluated checklist for surveys in software engineering. Information and Software Technology 119, 106240 (2020)

# Data Collection

- Besides all methodological issues… Every survey needs a proper **project plan**:

  1. Plan for methodological challenges

  2. Find a proper project organisation early

  3. Set up a proper project infrastructure

  4. Develop a good project dissemination plan

  5. Organise an efficient data collection

  6. Organise an efficient data curation and analysis

  7. Develop a good packaging and reporting pla

# Key Takeaways on Teaching Designing and Evaluating Survey Instruments

Different types of questionnaires, question types, and question categories, as well as measurement scales and conditions for obtaining accurate responses.

The role of GQM-Driven and Theory-Driven survey design.

Importance of using validated scales to improve construct validity.

Survey instruments may be evaluated using different methods to avoid threats to validity and improve reliability.

# 4) Sampling and Data Collection (LO3)

# Sampling

- At the beginning of any design of survey research, we should clarify what the **target population** is that we try to characterize and generalize to
  - Statistical analysis relies on **systematic sampling** from this target population

- In software engineering surveys, the unit of analysis that defines the granularity of the target population is often (de Mello et al. 2015):

**AN ORGANIZATION**

**A SOFTWARE TEAM OR PROJECT**

**AN INDIVIDUAL**

de Mello RM, da Silva PC, Travassos GH (2015) Investigating probabilistic sampling approaches for large-scale surveys in software engineering. Journal of Software Engineering Research and Development, 3(1):8.

# Sampling

- For common research questions, we are typically interested in **producing results related to all organizations that develop software** in the world or all software developers in the world.
  - We want to find theories that have a scope as wide as possible.

- We have no solid understanding about the **target population**.
  - Which companies are developing software?
  - How many software developers are there in the world?
  - What are the demographics of software engineers in the world?

- We face enormous difficulties to discuss **representativeness of a sample**, the needed size of the sample and, therefore, to what degree we can generalize our results.

# Sampling

- Scientists often rely on **demographic information** published by governmental or other public bodies such as statistical offices
  - These bodies are, so far, rather unhelpful for our task, because they do not provide a good idea about software-developing companies

- There are possibilities to approach the demographics of software engineers
  - Commercial **providers of data** from large surveys such as *Evans Data Corporation*:
    - Estimated number of developers worldwide as of 2018: 23 million
    - Include information on different roles, genders, used development processes and technologies
  - An open alternative is the **Stack Overflow Annual Developer Survey**

# Sampling

- Having demographic information, we can design our survey in a way that we collect **comparable data**.



- Then, we can **compare the distributions** in our survey and the larger surveys to estimate **representativeness**:

  **A** Should be part of the interpretation and discussion of the results;

  **B** Prevents us from overclaiming;

  **C** Gives us more credibility in case we cover the population well.

# Sampling

- A good sample size (n) can be estimated as follows (Yamane, 1973 apud Wagner et al., 2020):

$$n = \frac{N}{1 + Ne^2}$$

$\left\{ \begin{array}{l} \text{n - sample size} \\ \\ N \text{ - population size} \\ \\ e \text{ - level of precision (often set to 0.05 or 0.01)} \end{array} \right.$

- Reasonable sample size for **software developers** (using precision 0.05):

$$n = \frac{23{,}000{,}000}{1 + 23{,}000{,}000 \cdot 0.05^2} = 400$$

# Sampling

**1**

There is **no suitable official data** on the number and properties of software developing companies in the world.

**3**

For individual software engineers, **existing demographic studies** can be used to assess a survey's representativeness.

**2**

**Ethics** needs to be considered before contacting potential survey participants.

**4**

For the estimate of 23 million developers worldwide, a good sample size would be **400 respondents**.

# Sampling

- Survey sampling strategies are crucial to understand because they directly impact the validity and generalizability of survey research results
  - Linaker et al. [26] present some common sampling strategies, dividing them into:

| Non-probabilistic | Probabilistic |
|---|---|
| Convenience (Accidental) Sampling | Simple Random Sampling |
| Quota Sampling | Clustered Sampling |
| Purposive (Judgement) Sampling | Stratified Sampling |
| Snowball Sampling | Systematic Sampling |

# Strategies to approach the **population**

### CLOSED INVITATIONS

✔ Approaching *known* groups or individuals to participate per invitation-only;

✔ *Restricting* the survey access to those invited.

### OPEN INVITATIONS

✔ Approaching a broader, often *anonymous* audience via open survey access;

✔ *Anyone* with a link to the survey can participate.

# **Key Takeaways** (Wagner et al., 2020):

**1**

**Both strategies to approach the target population** (closed and open invitations) can be applied, but have distinct implications on the survey design and the recruitment approaches.

**2**

**Closed invitations are suitable in situations in which it is possible to precisely identify and approach a well-defined sample of the target population**. They may also be required in situations where filtering out participants that are not part of the target population would be difficult, harming the sample representativeness.

**3**

**Open invitations allow reaching out for larger samples**. However, they typically require more carefully considering context factors when designing the survey instruments. These context factors can then be used during the analyses to filter out participants that are not representative (e.g., applying the blocking principle to specific context factors).

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. Challenges in survey research. In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Key Takeaways on Teaching Sampling and Data Collection

What are the fundamentals and strategies for sampling and data collection?

What strategies could be explored to approach the target population?

# 5) Statistical and Qualitative Analysis (LO4)

# Statistical Analysis

- **DESCRIPTIVE AND INFERENTIAL STATISTICS**

- **NULL-HYPOTHESIS SIGNIFICANCE TESTING**

- **BOOTSTRAPPING CONFIDENCE INTERVALS**

- **BAYESIAN ANALYSIS**

- **STRUCTURAL EQUATION MODELING**

With the often large number of participants in surveys, we usually aim at a **statistical analysis** of the survey results.

A majority of the questionnaires are typically composed of **closed questions** that have **quantitative results**.

# Statistical Analysis

## Descriptive Statistics

- The goal of descriptive statistics is to **characterize** the answers to one or more questions of our specific sample

- We do not yet talk about generalizing to the population

- Which descriptive statistic is suitable depends on **what we are interested** in most and the **scale** of the data

| Scale | Nominal | Ordinal | Interval | Ratio |
|---|---|---|---|---|
| Values Counting | X | X | X | X |
| Values Ordering | | X | X | X |
| Equidistant Intervals | | | X | X |
| Values Division | | | | X |

# Statistical Analysis

## Descriptive Statistics

- Descriptive statistics for **ordinal scales** (e.g., Likert scales)

  - Frequency counting, mode, median, minimum, maximum, median absolute deviation (MAD), interquartile range (IQR)

  - An interesting alternative is showing the whole distribution of ordinal data in a stacked bar chart.



Generated using the Likert package in R

http://www.labape.com.br/rprimi/statR/T7_plus_likert.html

*Wagner et al., 2019*

# Statistical Analysis

- For **interval** or **ratio** scales we can use all available descriptive statistics, such as mean, variance, and standard deviation.

- Still, we recommend using **boxplots**, to enable eliminating outliers by using the quartile method



**Outliers**

**Maximum Value**

**3rd Quartile**

**Median**

**1st Quartile**

**Minimum Value**

Quartile Method
    Lower Outliers: Q1 - 1.5*IQR
    Upper Outliers: Q3 + 1.5*IQR
    Where IQR = Q3 – Q1.

**Inferential Statistics**

**Descriptive statistics** concern the sample
**Inference statistics** concern the population



Source: https://danawanzer.github.io/stats-with-jamovi/descriptive-vs-inferential-statistics.html

Different possibilities for analyzing quantitative survey results, including:

- null hypothesis significance testing;
- bootstrapping with confidence intervals;
- bayesian analysis;
- structural equation modeling.

# Statistical Analysis

**We need hypotheses to evaluate**
- ✔ A survey should be guided by a theory
- ✔ Propositions can be operationalized into hypotheses to test with the survey data

**In surveys we typically have:**
- ✔ Point estimate hypotheses for answers to single questions
- ✔ Hypotheses on correlations between answers to two questions

# Statistical Analysis

**Null-hypothesis Significance Testing (NHST)**

In general, two **hypotheses** are defined

| Null Hypothesis (H0) | Alternative Hypothesis (H1) |
|---|---|
| Indicates the observed differences are coincidental. It means that this is the hypothesis the researcher would like most to reject with high confidence | Represents the hypothesis indicating some type of effect, that can be accepted, or tested |

# Statistical Analysis

## Types of Errors

| Type I (α) | Type II (β) |
|---|---|
| It happens when the statistical test indicates the existence of a relationship between cause and effect that actually does not exist | It happens when the statistical test does not indicate a relationship between cause and effect that actually does exist |

Statistics tests allow confirming or refuting hypotheses
(according to a previously defined **significance level** - α-value)

# Statistical Analysis

**Null-hypothesis Significance Testing (NHST)**

## Types of Errors

# Statistical Analysis

- **Significance Testing**
    - Shows the likelihood of an type-I error to happen
        - Most common significance level (α): 10%, 5%, 1% and 0.1%
        - We call *p-value* the lowest level of significance that can be used to reject the null hypothesis
        - We say there is statistical significance when the calculated p-value is lower than the adopted significance level (*α-value*)

- Besides significance testing, it is important to also look at effect sizes.
    - *Cohen's d* is defined as the difference between two means divided by a standard deviation for the data:

$$d = \frac{\bar{x}_1 - \bar{x}_2}{s}$$

# Statistical Analysis

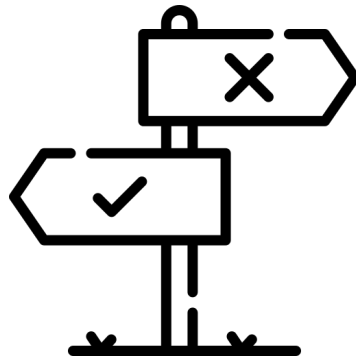## Null-hypothesis Significance Testing (NHST)

- Several statistical significance tests can be applied, with differences in their statistical power (Power= P (H0 rejected | Ho is false))
    - The statistical test with the highest power shall be used to evaluate the hypotheses

# Statistical Analysis

## Null-hypothesis Significance Testing (NHST)

- Several statistical significance tests can be applied, with differences in their statistical power (Power= P (H0 rejected | Ho is false))
    - The statistical test with the highest power shall be used to evaluate the hypotheses

```
                              Normal              Pearson
                              Distribution Data   Linear Regression
    Relationship
    Exploration
                              Non Normal          Spearman
                              Distribution Data   Non-Linear Regression
```

# Statistical Analysis

**Null-hypothesis Significance Testing (NHST)**
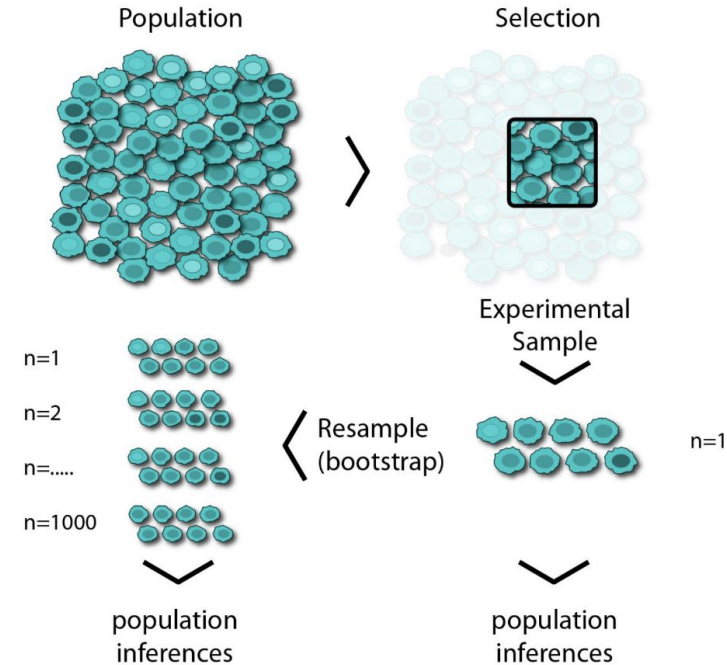
- **Problems with NHST**
    - Dichotomous nature of its results
    - Requires a representative sample  of the population, otherwise it is unclear what NHST actually means

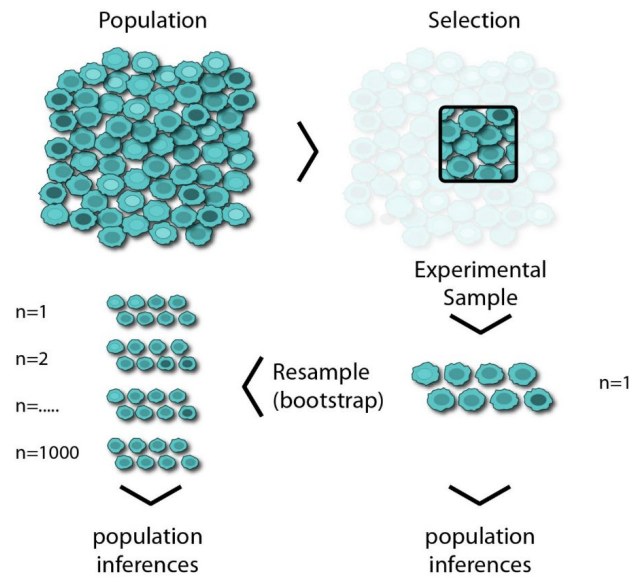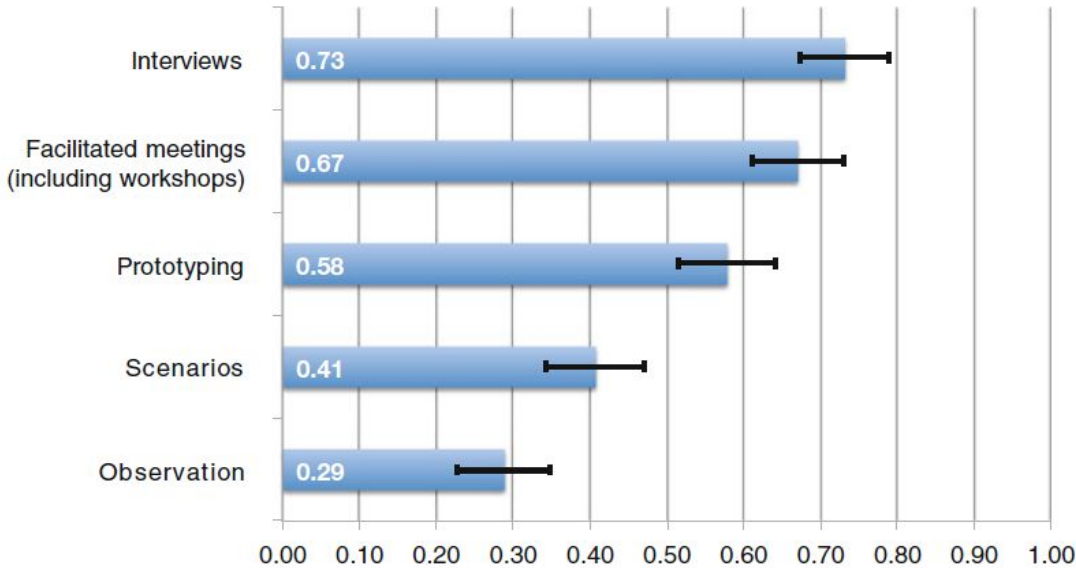- **We need alternatives...**

# Statistical Analysis

## Bootstrapping Confidence Intervals

- Replaces fixed significance **level thresholds**

- Involves **estimating a confidence interval** around a metric we are interested in
  - How large is the confidence interval?
  - How strongly do confidence intervals of methods to compare overlap?

- Idea of **bootstrapping**:
  - We repeatedly take samples with replacement and calculate the statistic we are interested in
  - This is repeated a large number of times and, thereby, provides us with an understanding of the distribution of the sample



Source: https://medium.com/swlh/bootstrap-sampling-using-pythons-numpy-85822d868977

# Bootstrapping Confidence Intervals: Example

1000 times resampling for **bootstrapping** confidence intervals



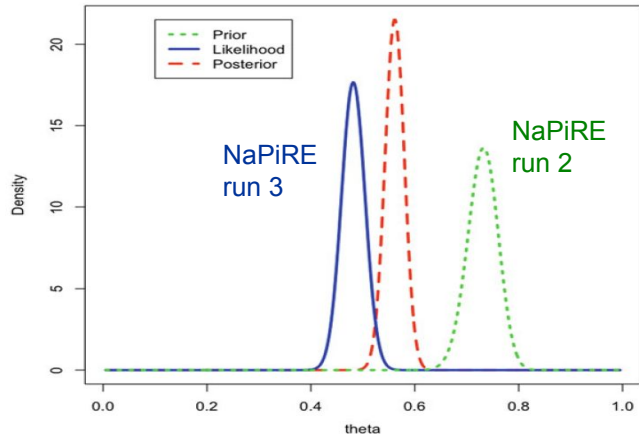Source: https://medium.com/swlh/bootstrap-sampling-using-pythons-numpy-85822d868977

***The Bootstrap Assumption:*** The original sample approximates the population from which it was drawn. So resamples from this sample approximate what we would get if we took many samples from the population. The bootstrap distribution of a statistic, based on many resamples, approximates the sampling distribution of the statistic, based on many samples.

# Statistical Analysis

## Bayesian Analysis

- In Bayesian statistics, probability is understood as a representation of the state of knowledge or belief
    - Acknowledges uncertainty
    - Allows integrating existing evidence and accumulating knowledge



Workshops for eliciting requirements (Wagner et al., 2020)

**Further reading:**

Torkar, R., Feldt, R. and Furia, C.A., 2020. Bayesian Data Analysis in Empirical Software Engineering: The Case of Missing Data. In Contemporary Empirical Methods in Software Engineering (pp. 289-324). Springer, Cham.

# Statistical Analysis

## Structural Equation Modeling

- **Used to test theories involving constructs** (also called latent variables).
    - In our Pandemic Programming survey example fear, disaster preparedness, home office ergonomics, wellbeing and productivity are all constructs

- To design a structural equation model, we first define a **measurement model**, which maps each reflective indicator into its corresponding construct.
    - For example, each of the five items comprising the WHO5 wellbeing scale is modeled as a reflective indicator of wellbeing

- SEM uses **Confirmatory Factor Analysis (CFA)** to estimate each construct as the shared variance of its respective indicators
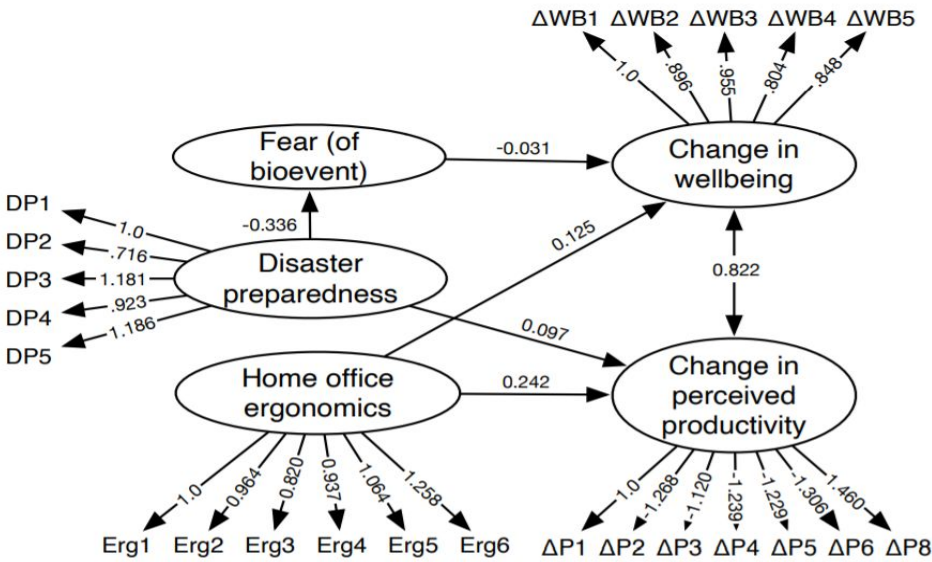
# Statistical Analysis

## Structural Equation Modeling

- Next, we define the **structural model**, which identifies the expected relationships among the constructs
  - The constructs we are attempting to predict are referred to as **endogenous** (dependent variables), while the predictors are **exogenous** (independent variables)

- SEM uses a **path modeling technique** (e.g. regression) to build a model that predicts the endogenous (latent) variables based on the exogenous variables, and to estimate both the strength of each relationship and the overall accuracy of the model.

# Structural Equation Modeling Example: Pandemic Programming

## Supported Model



The **arrows** between the constructs show the supported causal relationships.

The **path coefficients** (the numbers on the arrows) indicate the relative strength and direction of the relationships.

Ralph, P., Baltes, S., Adisaputri, G., Torkar, R., Kovalenko, V., Kalinowski, M., Novielli, N., Yoo, S., Devroey, X., Tan, X., Zhou, M., Turhan, B., Hoda, R., Hata, H., Robles, G., Fard, A. M., and Alkadhi, R, Pandemic Programming How COVID-19 affects software developers and how their organizations can help. Empirical Software Engineering (2020), 25: 4927-4961. 2020.

# Statistical Analysis

**Key Takeaways** (Wagner et al., 2020):

**1**
Always make clear whether you aim at analyzing opinions or facts

**2**
Descriptive statistics are always helpful

**3**
NHST inferential statistics are useful to test theoretical propositions

**4**
Bootstrapping confidence intervals helps to deal with uncertain sampling

**5**
Bayesian analysis allows us to directly integrate prior knowledge

**6**
SEM is a powerful multivariate analysis technique that is widely used in the social sciences and that should be further used in computer science research

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. Challenges in survey research. In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

# Qualitative Analysis

Besides the common focus on statistical analysis, surveys can also be **qualitative** and contain **open questions**

Open questions **do not impose restrictions** on respondents and allow them to more precisely describe the phenomena of interest according to their perspective and perceptions

However, they can lead to a **large amount of qualitative data to analyze**, which is not easy and may require a significant amount of resources
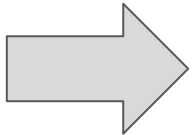
We recommend referring to chapter "Qualitative Data Analysis in Software Engineering: Techniques and Teaching Insights" for further advice on teaching qualitative methods

# Qualitative Analysis

The answers to such **open questions** can help researchers to further understand a phenomenon eventually including causal relations among theory constructs and theoretical explanations
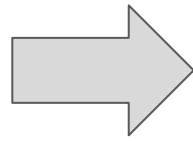
**Open questions can help generating new theories**

A research method commonly employed to support qualitative analyses is **Grounded Theory**.

There are at least three main streams of GT:
✔ Glaser's GT (classic or Glaserian GT) *(Glaser, 1992)*
✔ Corbin and Strauss' GT (Straussian GT) (*Corbin and Strauss, 1990)*
✔ Charmaz's constructivist GT *(Charmaz, 2014)*

Grounded theory, "in theory", involves **inductively** generating theory from data.

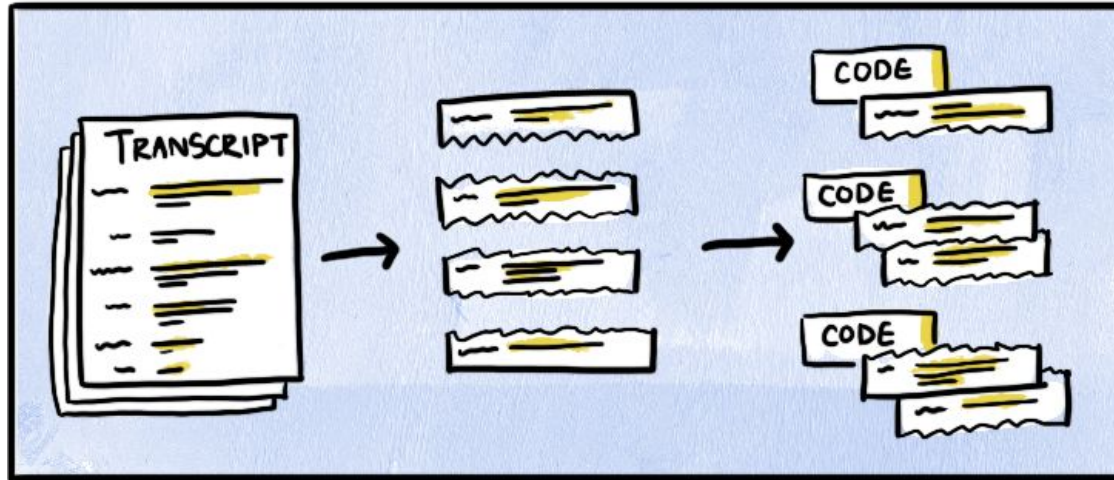**Few "GT" Studies Generate Theory (Stol et al., 2016).**

# Qualitative Analysis

**1** Turn your data into small, discrete components of data

**2** Code each discrete pieces of data with a descriptive label



*Source: https://delvetool.com/blog/openaxialselective*

Corbin, J.M. and Strauss, A., 1990. Grounded theory research: Procedures, canons, and evaluative criteria. Qualitative sociology, 13(1), pp.3-21.
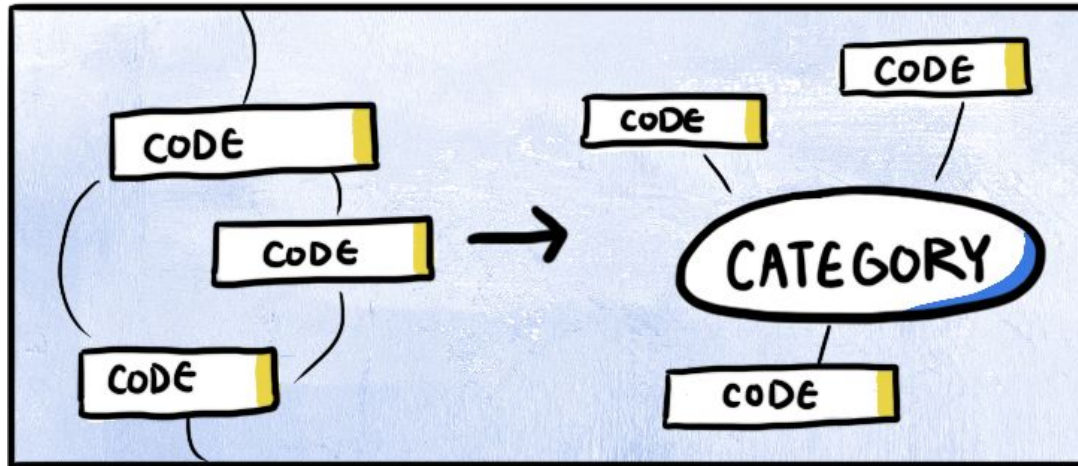
# Qualitative Analysis

Find connections and relationships between codes

4 Aggregate and condense codes into broader categories



*Source: https://delvetool.com/blog/openaxialselective*

Corbin, J.M. and Strauss, A., 1990. Grounded theory research: Procedures, canons, and evaluative criteria. Qualitative sociology, 13(1), pp.3-21.

# Qualitative Analysis

**5** Bring it together with one overarching category

**6** Identify the connections between this overarching category and the rest of your codes and data

**7** Remove categories or codes that don't have enough supporting data

**8** Read the transcript again, and code according to this overarching category



*Source: https://delvetool.com/blog/openaxialselective*

Corbin, J.M. and Strauss, A., 1990. Grounded theory research: Procedures, canons, and evaluative criteria. Qualitative sociology, 13(1), pp.3-21.

# Qualitative Analysis: Example



Sankey diagram with left-side causes connecting to middle problems and right-side outcomes.

Left column:
- Lack of experience of RE team members
- Lack of time
- Communication flaws between project team and the customer
- Missing direct communication to customer
- Requirements remain too abstract
- Too high team distribution
- Unclear roles and responsibilities at customer side
- Weak qualification of RE team members
- Lack of a well-defined RE process
- Customer does not know what he wants

Middle column:
- Underspecified reqs that are too abstract and allow for various interpretations
- Communication flaws between project team and the customer
- Incomplete and / or hidden requirements
- Communication flaws within the project team
- Inconsistent requirements
- Insufficient support by customer
- Weak access to customer needs and / or (internal) business information
- Time boxing / Not enough time in general
- Moving targets (changing goals, business processes and / or requirements)
- Stakeholders with difficulties in separating reqs from known solution designs

Right column:
- Project Failed
- Project Completed

Fernández, D. M.; Wagner, S.; Kalinowski, M.; Felderer, M.; Mafra, P.; Vetro, A.; Conte, T.; Christiansson, M.; Greer, D.; Lassenius, C.; Männistö, T.; Nayabi, M.; Oivo, M.; Penzenstadler, B.; Pfahl, D.; Prikladnicki, R.; Ruhe, G.; Schekelmann, A.; Sen, S.; Spínola, R. O.; Tuzcu, A.; de la Vara, J. L.; and Wieringa, R. Naming the pain in requirements engineering - Contemporary problems, causes, and effects in practice. Empirical Software Engineering, 22(5): 2298-2338. 2017.
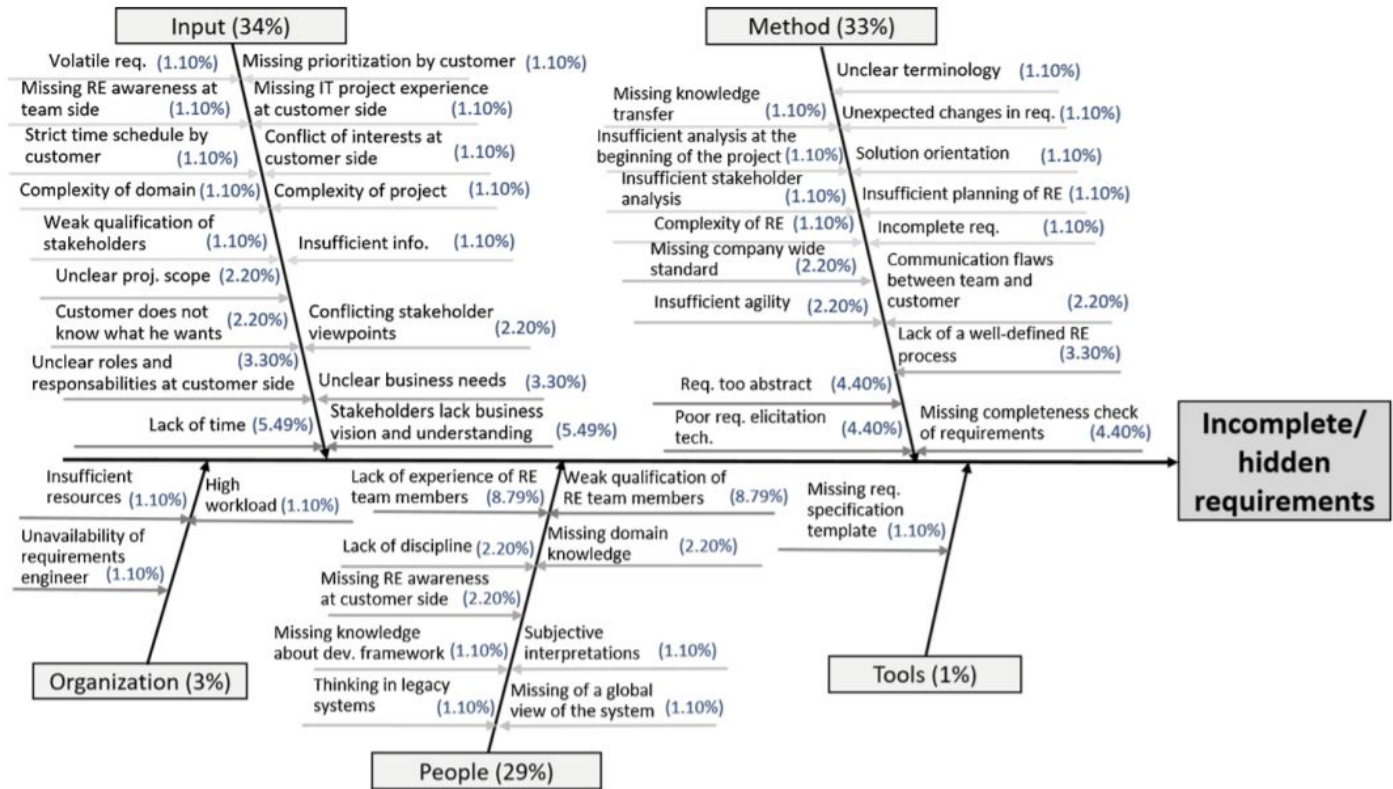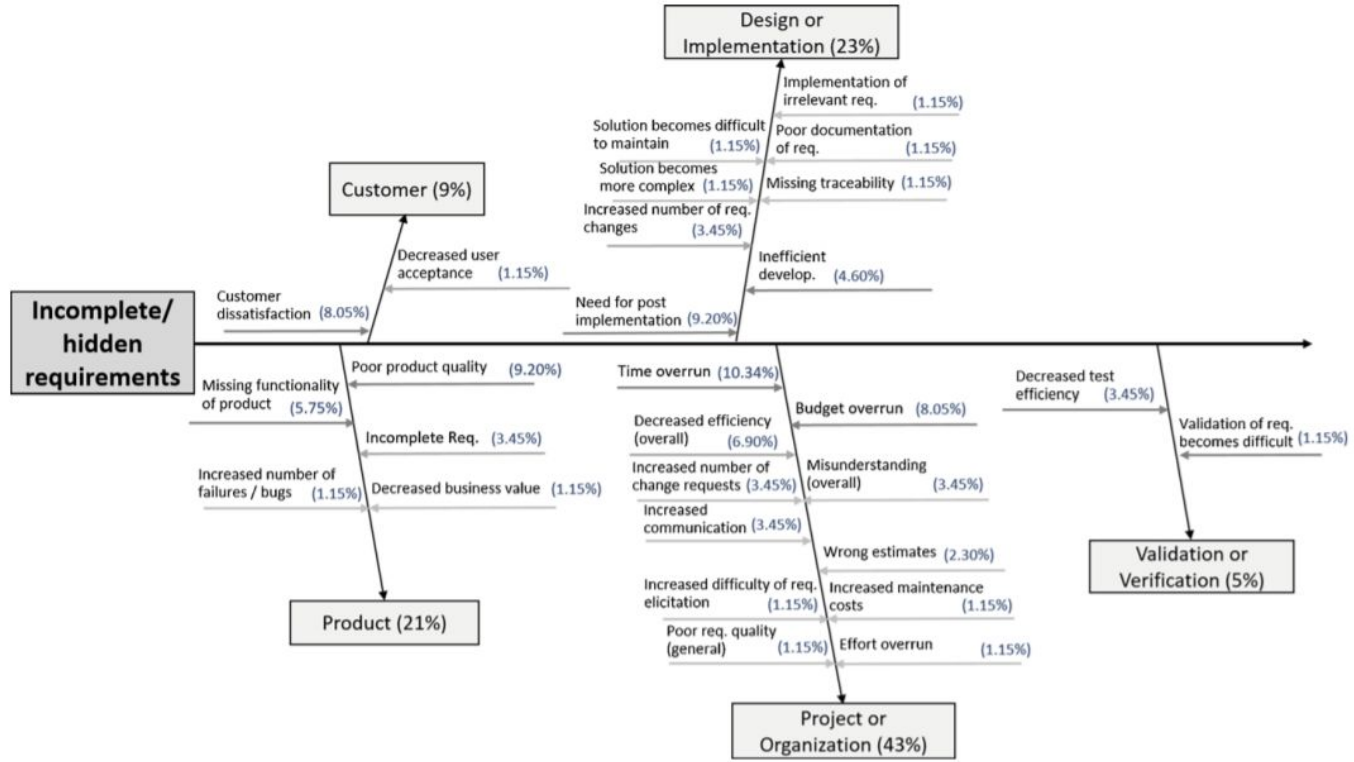
# Qualitative Analysis: Example



Fernández, D. M.; Wagner, S.; Kalinowski, M.; Felderer, M.; Mafra, P.; Vetro, A.; Conte, T.; Christiansson, M.; Greer, D.; Lassenius, C.; Männistö, T.; Nayabi, M.; Oivo, M.; Penzenstadler, B.; Pfahl, D.; Prikladnicki, R.; Ruhe, G.; Schekelmann, A.; Sen, S.; Spínola, R. O.; Tuzcu, A.; de la Vara, J. L.; and Wieringa, R. Naming the pain in requirements engineering - Contemporary problems, causes, and effects in practice. Empirical Software Engineering, 22(5): 2298-2338. 2017.

# Qualitative Analysis: Example

Fernández, D. M.; Wagner, S.; Kalinowski, M.; Felderer, M.; Mafra, P.; Vetro, A.; Conte, T.; Christiansson, M.; Greer, D.; Lassenius, C.; Männistö, T.; Nayabi, M.; Oivo, M.; Penzenstadler, B.; Pfahl, D.; Prikladnicki, R.; Ruhe, G.; Schekelmann, A.; Sen, S.; Spínola, R. O.; Tuzcu, A.; de la Vara, J. L.; and Wieringa, R. Naming the pain in requirements engineering - Contemporary problems, causes, and effects in practice. Empirical Software Engineering, 22(5): 2298-2338. 2017.

# Qualitative Analysis

## **Key Takeaways** (Wagner et al., 2020):

**1**

When preparing your survey, invest effort in avoiding confounding factors that may interfere in having respondents focusing mainly on the survey question when providing their answers (e.g., language issues). Assess the instrument validity.

**3**

When reporting the qualitative analysis of your survey, explicitly state your research method, providing details on eventual deviations.

**2**

Applying coding and analysis techniques from Grounded Theory can help to understand qualitative data gathered through open questions.

**4**

To avoid researcher bias and improve the reliability of the results, qualitative analyses should be conducted in teams and make use of independent validations. Also, ideally the raw and analyzed data should be open to enable other researchers to replicate the analysis procedures.

Wagner, S., Mendez, D., Felderer, M., Graziotin, D. and Kalinowski, M., 2020. Challenges in survey research. In: Contemporary Empirical Methods in Software Engineering (pp. 93-125). Springer, Cham.

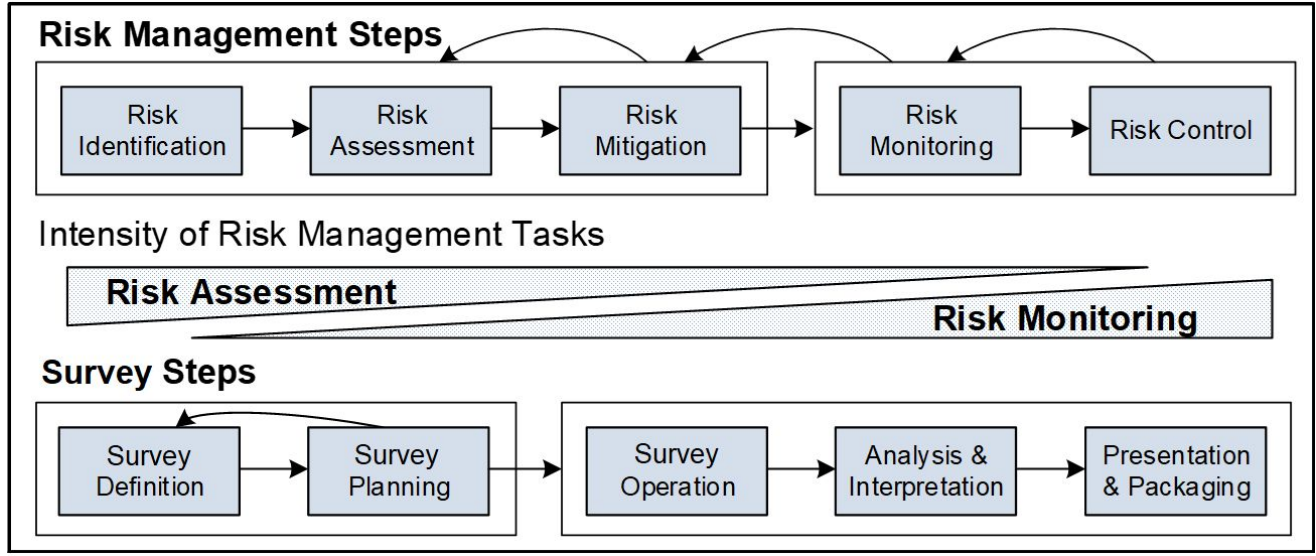# Key Takeaways on Teaching Statistical and Qualitative Analysis

Descriptive statistics provide a foundation for understanding data by summarizing key characteristics, while alternatives to traditional inferential statistics offer robust tools for analyzing data under various conditions and assumptions

Open questions enrich qualitative research by capturing detailed and nuanced responses.

# 6) Threats to Validity and Reliability (LO5)

# Survey Risk Management

**Validity** is a property of inferences and every study faces **Threats to Validity** (Biffl et al., 2014).
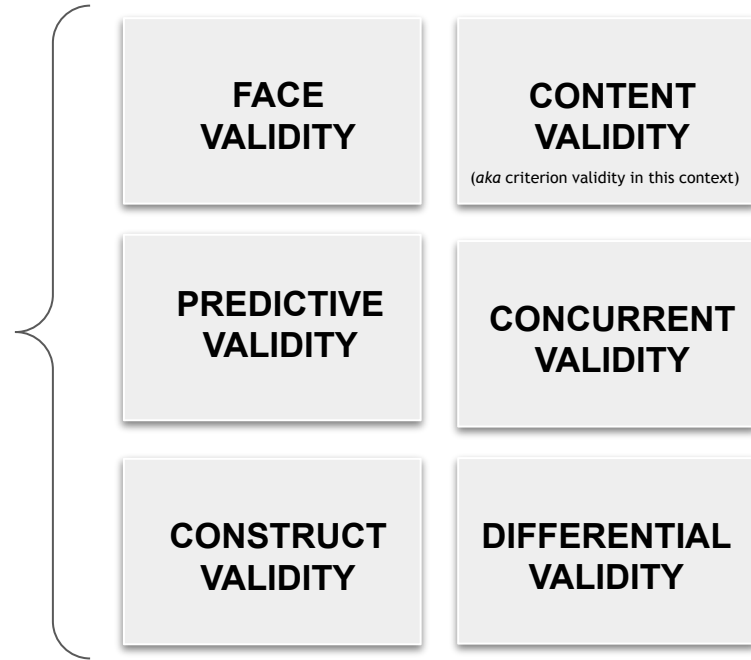
Biffl, S., Kalinowski, M., Ekaputra, F., Neto, A.A., Conte, T. and Winkler, D., 2014, September. Towards a semantic knowledge base on threats to validity and control actions in controlled experiments. In Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (pp. 1-4).

# Validity Assessment

In *psychometrics*, **validity** concerns "the degree to which evidence and theory support the interpretation of test scores for proposed uses of tests" (AERA et al., 2014)

Rust (2009) summarized six **facets of validity** in the context of psychometric tests:

| | |
|---|---|
| **FACE VALIDITY** | **CONTENT VALIDITY** (*aka* criterion validity in this context) |
| **PREDICTIVE VALIDITY** | **CONCURRENT VALIDITY** |
| **CONSTRUCT VALIDITY** | **DIFFERENTIAL VALIDITY** |

AERA, APA, NCME (2014) Standards for educational and psychological testing. American Educational Research Association, Washington.

Rust J (2009) Modern psychometrics: the science of psychological assessment. Routledge, Hove, East Sussex New York.

# Validity Assessment

In *software engineering* we typically aim at assessing whether it is possible to safely conclude that a survey measures what it is supposed to:

The following validity types are discussed in this context (Kitchenham and Pfleeger, 2008 apud Linaker et al., 2015):

| | |
|---|---|
| **FACE VALIDITY** | Typically involves a lightweight review of the questionnaire by randomly chosen respondents |
| **CRITERION VALIDITY** | Refers to how the questionnaire can separate between respondents that belong to different groups. An existing classification and mapping of the different groups in the target population must be in place |
| **CONTENT VALIDITY** | Typically involves having a (focus) group of reviewers evaluating the questionnaire. The group should include subject matter experts and example respondents from the target population |

Kitchenham, B.A. and Pfleeger, S.L., 2008. Personal opinion surveys. In Guide to advanced empirical software engineering (pp. 63-92). Springer, London.

# Validity Assessment

In *software engineering* we typically aim at assessing whether it is possible to safely conclude that a survey measures what it is supposed to:

The following validity types are discussed in this context (Kitchenham and Pfleeger, 2008 apud Linaker et al., 2015):

| FACE VALIDITY | Typically involves a lightweight review of the questionnaire by randomly chosen respondents |
|---|---|
| CRITERION VALIDITY | Refers to how the questionnaire can separate between respondents that belong to different groups. An existing classification and mapping of the different groups in the target population must be in place |
| CONTENT VALIDITY | Typically involves having a (focus) group of reviewers evaluating the questionnaire. The group should include subject matter experts and example respondents from the target population |
| CONSTRUCT VALIDITY | How well the question actually measures the construct it was intended to by the designer |

Kitchenham, B.A. and Pfleeger, S.L., 2008. Personal opinion surveys. In Guide to advanced empirical software engineering (pp. 63-92). Springer, London.

# Reliability Assessment

## **Reliability** (aka External Validity and Generalizability):

### TEST-RETEST RELIABILITY

✔ The same subject responds to the same survey two times, and it is measured whether the subject gives the same answers each time;

✔ Kitchenham and Pfleeger (2008) state that if the correlation between both of the answers is greater than 0.7 the test-retest reliability can be considered good.

### PHRASING / REORDER EFFECT RELIABILITY

✔ Testing whether the phrasing or reordering of questions has any effect on the answers by a respondent (assesses instrument bias on the respondent).

**Reliability** (aka External Validity and Generalizability):

| INTER-OBSERVER RELIABILITY | INTER-OBSERVER RELIABILITY |
|---|---|
| ✔ Assesses observer interview bias in not self-administered surveys;<br><br>✔ Assesses observer analysis bias (e.g., when interpreting and decoding open ended questions);<br><br>✔ Typically addressed by having two or more observers involved in the interview and analysis process | ✔ If conclusions are to be drawn on the whole population, not just on the sample, the reliability needs to be proven and established |

## **Reliability** (aka External Validity and Generalizability):

| Threats | Treatment |
|---------|-----------|
| Face Validity – Bad instrumentation | Revision and evaluation of the questionnaire about the format and formulation of the questions. Questions objectively focused on the 3PDF. Running a pilot study. |
| Content Validity – Inadequate explanation of the constructs | Revision and evaluation of the questionnaire about the format and formulation of the questions. Running a pilot study. Providing a brief explanation on the 3PDF and a link with further details. |
| Criterion Validity – Not surveying the target population. | We identified SE SLR update authors following an explicitly documented and carefully conducted procedure (cf. Section 3). |
| Construct Validity – Inadequate measurement procedures and unreliable results. | We only used frequency counting, which can be safely applied to discrete survey questions concerning the relevance of the 3PDF questions and the agreement with the 3PDF decision drivers. Also, we triangulated the answers with the provided explanations. |
| Reliability – Lack of statistical conclusion validity | This threat strongly depends on the sample size. Unfortunately, while contacting twice the SE SLR update authors we were aware of, our final sample size was still limited. Hence, we focused our results on qualitative analyses and did not make any further claims on conclusion validity. |

Mendes, E., Wohlin, C., Felizardo, K. and Kalinowski, M., 2020. When to update systematic literature reviews in software engineering. Journal of Systems and Software, 167, p.110607.

# Key Takeaways on Teaching Threats to Validity and Reliability

Understanding and ensuring validity and reliability are fundamental for conducting trustworthy and thorough software engineering surveys. Validity ensures that the survey measures what it is intended to measure, while reliability ensures consistent results across different instances of the survey.

Different types of validity are essential in survey research to ensure that the survey accurately reflects the concept being studied.

# 7) Ethical Considerations (LO5)

**Ethical considerations** are paramount in survey research within software engineering, as they ensure **respect** for participants, the **integrity of the data**, and the **credibility** of the research findings.

# Ethical Considerations

- In software engineering, there is yet **no established standard or guidelines** on how to conduct surveys ethically

- The Insight Association provides ethical guidelines that consider **unethical sampling**, among other practices: "*Collection of respondent emails from Websites, portals, Usenet or other bulletin board postings without specifically notifying individuals that they are being 'recruited' for research purposes*".

- We will probably need flexible **rules and guidelines** to keep developers in social media from being spammed by study requests while still allowing research to take place.

- We should all consider thoughtfully **how and whom** we contact for a survey study.

# Ethical Considerations

## INFORMED CONSENT

Participants must be fully informed about the nature of the research, what it involves, the risks and benefits, and their rights to withdraw at any time without penalty.

## PRIVACY AND CONFIDENTIALITY

Researchers must protect the privacy of participants and the confidentiality of their data, using data encryption and anonymization techniques where appropriate.

## INSTITUTIONAL ETHICS REVIEW

Submitting survey research to institutional ethics review boards, as they will ensure the research adheres to ethical standards and protects participant

# Key Takeaways on Teaching Ethical Considerations

Ethics needs to be considered before contacting potential survey participants. Participants must be fully informed about the nature of the research, what it involves, the risks and benefits, and their rights to withdraw at any time without penalty.

Pay attention to the role of the institutional ethics review boards and how to report survey ethics in software engineering publications.

# 8) Concluding Remarks

For further information, see section 3.6 in the chapter.

# Concluding Remarks

- We have explored effective strategies for **survey research**, combining **theoretical foundations** with **practical applications**.

| ID | Learning Objective | Students will be able to ... | Bloom's Taxonomy |
|---|---|---|---|
| LO1 | Understanding the Characteristics and Purposes of Survey Research | ... articulate on the characteristics and purposes of survey research.<br>... provide survey research application examples. | Remembering & Understanding |
| LO2 | Designing and Evaluating Survey Instruments | ... create survey instruments aligning with specific research objectives and theories.<br>... critically assess the effectiveness of survey instruments. | Evaluating & Creating |
| LO3 | Mastering Sampling and Data Collection | ... apply best practices in sampling and data collection.<br>... understand the trade-offs of different sampling and data collection methods. | Understanding & Applying |
| LO4 | Applying Statistical and Qualitative Analysis Methods | ... utilize statistical and qualitative analysis techniques to interpret survey data. | Applying & Analyzing |
| LO5 | Identifying and Addressing Validity and Reliability Threats | ... analyze and address potential threats to the validity and reliability of survey research. | Analyzing & Evaluating |
| LO6 | Understanding Ethical Considerations in Survey Research | ... identify, understand, and apply ethical considerations in survey research. | Understanding & Applying |

**Table 1** Learning Objectives and Bloom's Taxonomy Levels.

# Teaching Survey Research in Software Engineering

**Authors:**
Marcos Kalinowski (Pontifical Catholic University of Rio de Janeiro)
Allysson Allex Araújo (Federal University of Cariri)
Daniel Mendez (Blekinge Institute of Technology)