# IN5140: Smart processes and agile methods in software engineering

## Empirical Research Methods in Software Engineering

Associate professor
Gunnar Bergersen

gunnab@ifi.uio.no

25.10.2023

UiO : **Department of Informatics**
University of Oslo

# Q: Why use «empirical methods» ?

① Why empirical methods?

not empirical methods math logic

rigour
data - observations
reproducable result
useful knowledge

truth

⇒ incremental?     morality
⇒ falsifiable?
transparent

# Truth

"Pale blue dot"
https://www.youtube.com/watch?v=GO5FwsblpT8



"We can judge our progress by the courage of our questions and the depth of our answers, our willingness to embrace what is true rather than what feels good."
—Carl Sagan

"We don't need any scientific report that tells us …"

Experience trumps research?

# strawman

**Misrepresenting someone's argument to make it easier to attack.**

By exaggerating, misrepresenting, or just completely fabricating someone's argument, it's much easier to present your own position as being reasonable, but this kind of dishonesty serves to undermine rational debate.

After Will said that we should put more money into health and education, Warren responded by saying that he was surprised that Will hates our country so much that he wants to leave it defenceless by cutting military spending.

# false cause

**Presuming that a real or perceived relationship between things means that one is the cause of the other.**

Many people confuse correlation (things happening together or in sequence) for causation (that one thing actually causes the other to happen). Sometimes correlation is coincidental, or it may be attributable to a common cause.

Pointing to a fancy chart, Roger shows how temperatures have been rising over the past few centuries, whilst at the same time the numbers of pirates have been decreasing; thus pirates cool the world and global warming is a hoax.

# slippery slope

**Asserting that if we allow A to happen, then Z will consequently happen too, therefore A should not happen.**

The problem with this reasoning is that it avoids engaging with the issue at hand, and instead shifts attention to baseless extreme hypotheticals. The merits of the original argument are then tainted by unsubstantiated conjecture.

Colin Closet asserts that if we allow same-sex couples to marry, then the next thing we know we'll be allowing people to marry their parents, their cars and even monkeys.

# ad hominem

**Attacking your opponent's character or personal traits in an attempt to undermine their argument.**

Ad hominem attacks can take the form of overtly attacking somebody, or casting doubt on their character. The result of an ad hom attack can be to undermine someone without actually engaging with the substance of their argument.

After Sally presents an eloquent and compelling case for a more equitable taxation system, Sam asks the audience whether we should believe anything from a woman who isn't married, was once arrested, and smells a bit weird.

# special pleading

**Moving the goalposts or making up exceptions when a claim is shown to be false.**

Humans are funny creatures and have a foolish aversion to being wrong. Rather than appreciate the benefits of being able to change one's mind through better understanding, many will invent ways to cling to old beliefs.

Edward Johns claimed to be psychic, but when his abilities were tested under proper scientific conditions, they magically disappeared. Edward explained this saying that one had to have faith in his abilities for them to work.

# loaded question

**Asking a question that has an assumption built into it so that it can't be answered without appearing guilty.**

Loaded question fallacies are particularly effective at derailing rational debates because of their inflammatory nature - recipients of a loaded question are compelled to defend themselves and may appear flustered or on the back foot.

Grace and Helen were both romantically interested in Brad. One day, with Brad sitting within earshot, Grace asked in an inquisitive tone whether Helen was having any problems with a fungal infection.

# the gambler's fallacy

**Believing that 'runs' occur to statistically independent phenomena such as roulette wheel spins.**

This commonly believed fallacy can be said to have helped create a city in the desert of Nevada USA. Though the overall odds of a 'big run' happening may be low, each spin of the wheel is itself entirely independent from the last.

Red had come up six times in a row on the roulette wheel, so Greg knew that it was close to certain that black would be next up. Suffering an economic form of natural selection with this thinking, he soon lost all of his savings.

# bandwagon

**Appealing to popularity or the fact that many people do something as an attempted form of validation.**

The flaw in this argument is that the popularity of an idea has absolutely no bearing on its validity. If it did, then the Earth would have made itself flat for most of history to accommodate this popular belief.

Shamus pointed a drunken finger at Sean and asked him to explain how so many people could believe in leprechauns if they're only a silly old superstition. Sean, however, had had a few too many Guinness himself and fell off his chair.

# black-or-white

**Where two alternative states are presented as the only possibilities, when in fact more possibilities exist.**

Also known as the false dilemma, this insidious tactic has the appearance of forming a logical argument, but under closer scrutiny it becomes evident that there are more possibilities than the either/or choice that is presented.

Whilst rallying support for his plan to fundamentally undermine citizens' rights, the Supreme Leader told the people they were either on his side, or on the side of the enemy.

# begging the question

**A circular argument in which the conclusion is included in the premise.**

This logically incoherent argument often arises in situations where people have an assumption that is very ingrained, and therefore taken in their minds as a given. Circular reasoning is bad mostly because it's not very good.

The word of Zorbo the Great is flawless and perfect. We know this because it says so in The Great and Infallible Book of Zorbo's Best and Most Truest Things that are Definitely True and Should Not Ever Be Questioned.

# appeal to emotion

**Manipulating an emotional response in place of a valid or compelling argument.**

Appeals to emotion include appeals to fear, envy, hatred, pity, guilt, and more. Though a valid, and reasoned, argument may sometimes have an emotional aspect, one must be careful that emotion doesn't obscure or replace reason.

Luke didn't want to eat his sheep's brains with chopped liver and brussels sprouts, but his father told him to think about the poor, starving children in a third world country who weren't fortunate enough to have any food at all.

# tu quoque

**Avoiding having to engage with criticism by turning it back on the accuser - answering criticism with criticism.**

Literally translating as 'you too' this fallacy is commonly employed as an effective red herring because it takes the heat off the accused having to defend themselves and shifts the focus back onto the accuser themselves.

Nicole identified that Hannah had committed a logical fallacy, but instead of addressing the substance of her claim, Hannah accused Nicole of committing a fallacy earlier on in the conversation.

# burden of proof

**Saying that the burden of proof lies not with the person making the claim, but with someone else to disprove.**

The burden of proof lies with someone who is making a claim, and is not upon anyone else to disprove. The inability, or disinclination, to disprove a claim does not make it valid (however we must always go by the best available evidence).

Bertrand declares that a teapot is, at this very moment, in orbit around the Sun between the Earth and Mars, and that because no one can prove him wrong his claim is therefore a valid one.

# composition /division

**Assuming that what's true about one part of something has to be applied to all, or other, parts of it.**

Often when something is true for the part it does also apply to the whole, but because this isn't always the case it can't be presumed to be true. We must show evidence for why a consistency will exist.

Daniel was a precocious child and had a liking for logic. He reasoned that atoms are invisible, and that he was made of atoms and therefore invisible too. Unfortunately, despite his thinky skills, he lost the game of hide and go seek.

# no true scotsman

**Making what could be called an appeal to purity as a way to dismiss relevant criticisms or flaws of an argument.**

This fallacy is often employed as a measure of last resort when a point has been lost. Seeing that a criticism is valid, yet not wanting to admit it, new criteria are invoked to dissociate oneself or one's argument.

Angus declares that Scotsmen do not put sugar on their porridge, to which Lachlan points out that he is a Scotsman and puts sugar on his porridge. Furious, like a true Scot, Angus yells that no **true** Scotsman sugars his porridge.

# the fallacy fallacy

**Presuming a claim to be necessarily wrong because a fallacy has been committed.**

It is entirely possible to make a claim that is false yet argue with logical coherency for that claim, just as it is possible to make a claim that is true and justify it with various fallacies and poor arguments.

Recognising that Amanda had committed a fallacy in arguing that we should eat healthy food because a nutritionist said it was popular, Alyse said we should therefore eat bacon double cheeseburgers every day.

# personal incredulity

**Saying that because one finds something difficult to understand, it's therefore not true.**

Subjects such as biological evolution via the process of natural selection require a good amount of understanding before one is able to properly grasp them; this fallacy is usually used in place of that understanding.

Kirk drew a picture of a fish and a human and with effusive disdain asked Richard if he really thought we were stupid enough to believe that a fish turned into a human through just, like, random things happening over time.

# ambiguity

**Using double meanings or ambiguities of language to mislead or misrepresent the truth.**

Politicians are often guilty of using ambiguity to mislead and will later point to how they were technically not outright lying if they come under scrutiny. It's a particularly tricky and premeditated fallacy to commit.

When the judge asked the defendant why he hadn't paid his parking fines, he said that he shouldn't have to pay them because the sign said 'Fine for parking here' and so he naturally presumed that it would be fine to park there.

# genetic

**Judging something good or bad on the basis of where it comes from, or from whom it comes.**

To appeal to prejudices surrounding something's origin is another red herring fallacy. This fallacy has the same function as an ad hominem, but applies instead to perceptions surrounding sources, or source of context.

Accused on the 6 o'clock news of corruption and taking bribes, the senator said that we should all be very wary of the things we hear in the media, because we all know how very unreliable the media can be.

# appeal to authority

**Saying that because an authority thinks something, it must therefore be true.**

It's important to note that this fallacy should not be used to dismiss the claims of experts, or scientific consensus. Appeals to authority are not valid arguments, but nor is it reasonable to disregard the claims of experts who have a demonstrated depth of knowledge unless one has a similar level of understanding.

Not able to defend his position that evolution 'isn't true' Bob says that he knows a scientist who also questions evolution (and presumably isn't himself a primate).

# appeal to nature

**Making the argument that because something is 'natural' it is therefore valid, justified, inevitable, good, or ideal.**

Many 'natural' things are also considered 'good', and this can bias our thinking; but naturalness itself doesn't make something good or bad. For instance murder could be seen as very natural, but that doesn't mean it's justifiable.

The medicine man rolled into town on his bandwagon offering various natural remedies, such as very special plain water. He said that it was only natural that people should be wary of 'artificial' medicines like antibiotics.

# anecdotal

**Using personal experience or an isolated example instead of a valid argument, especially to dismiss statistics.**

It's often much easier for people to believe someone's testimony as opposed to understanding variation across a continuum. Scientific and statistical measures are almost always more accurate than individual perceptions and experiences.

Jason said that that was all cool and everything, but his grandfather smoked, like, 30 cigarettes a day and lived until 97 - so don't believe everything you read about meta analyses of sound studies showing proven causal relationships.

# the texas sharpshooter

**Cherry-picking data clusters to suit an argument, or finding a pattern to fit a presumption.**

This 'false cause' fallacy is coined after a marksman shooting at barns and then painting a bullseye target around the spot where the most bullet holes appear. Clusters naturally appear by chance, and don't necessarily indicate causation.

The makers of Sugarette Candy Drinks point to research showing that of the five countries where Sugarette drinks sell the most units, three of them are in the top ten healthiest countries on Earth, therefore Sugarette drinks are healthy.

# middle ground

**Saying that a compromise, or middle point, between two extremes must be the truth.**

Much of the time the truth does indeed lie between two extreme points, but this can bias our thinking: sometimes a thing is simply untrue and a compromise of it is also untrue. Half way between truth and a lie, is still a lie.

Holly said that vaccinations caused autism in children, but her scientifically well-read friend Caleb said that this claim had been debunked and proven false. Their friend Alice offered a compromise that vaccinations cause some autism.

PLATO    SOCRATES    ARISTOTLE

# thou shalt not commit logical fallacies

A logical fallacy is a flaw in reasoning. Strong arguments are void of logical fallacies, whilst arguments that are weak tend to use logical fallacies to appear stronger than they are. They're like tricks or illusions of thought, and they're often very sneakily used by politicians, the media, and others to fool people.

Don't be fooled! This poster has been designed to help you identify and call out dodgy logic wherever it may raise its ugly, incoherent head. If you see someone committing a logical fallacy online, link them to the relevant fallacy to school them in thinkiness e.g. yourlogicalfallacyis.com/strawman

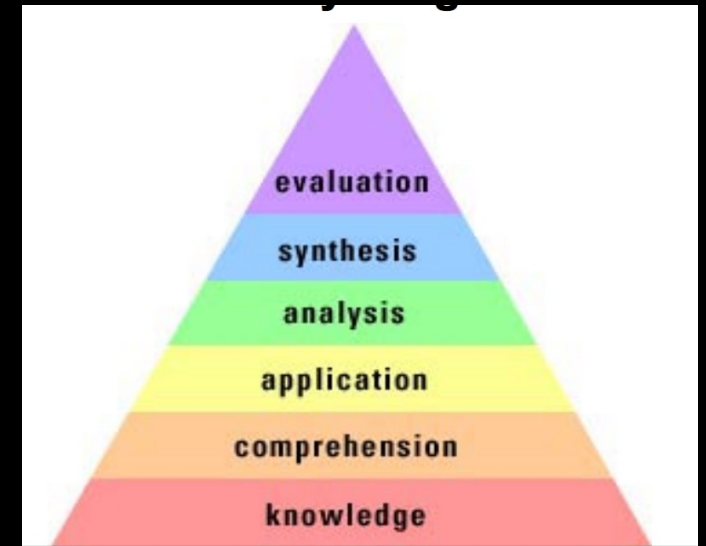You can download this poster for free at yourlogicalfallacyis.com /poster

# About me

- Current:
  - Associate professor at Software Engineering
- Education
  - MSc (2001) and PhD (2015) from UiO
  - PhD thesis: "Measuring programming skill"
- Prior work experience
  - Programmer
  - IT Project leader two companies
  - CEO three companies
  - Startup based on my PhD

Golonka et al (2023): the construct of cuteness https://doi.org/10.3389/fpsyg.2023.1068373

# Learning objectives



After this lecture, you should be able to …

- *Describe* the central elements of e**mpirical research** and
  - *Explain* the steps involved in **Evidence-based software engineering** and provide *critique* of claims that based on use of theory and empirical results
  - *Discuss* the strengths and weaknesses of different **empirical research methods** and *suggest* what method(s) to use for a specific situation
  - *Identify* how changes in **context** may affect the answer to a research question

# Structure

- Empirical research
  - Evidence-based software engineering
  - Empirical research methods
    - Validity
  - The importance of context

**Note: I only present a subset of the slides you will have to study**

# Context



Understand Software Processes → Process Modeling / **Process models**

Understand Software Processes → Process Improvement

Implement (Agile) software processes → Lean Software Engineering

Implement (Agile) software processes → Large-Scale Agile Projects
- Large-Scale Agile and Architecture
- Large Scale Agile Transformation
- Global Software Engineering
- Managing Technical-, Social-, Process- Debt

Implement (Agile) software processes → Quality Assurance Processes

Agile Practices and Teamwork

Assess Software Processes → Empirical Methods in Software Engineering

# Empirical research

- Empirical research concerns the acquisition of knowledge by empirical methods

- Empirical research seeks to explore, describe, predict, and explain natural, social, or cognitive phenomena by using evidence based on observation or experience

- What constitutes knowledge, and the methods for acquiring it, rests on basic assumptions regarding:
  - Ontology, i.e., what we believe to exist,
  - Epistemology i.e., how beliefs are acquired and what justifies them,
  - Methodology, e.g., the inductive or the hypothetico-deductive method.

**Figure 1: A pragmatic-realist view of measurement**

# Empirical evidence



- Empirical evidence is the data on which a conclusion or judgment may be based.

- Interpreting and judging such evidence depends on the "eye of the observer".

- Much research applies to groups of individuals or populations, and are not always relevant or valid for other situations.

- Accurate prediction or absolute proof of causality applicable to individuals or to real-life settings are virtually impossible.

- The contributions of empirical research to any situation depend on the context, judgment and values, understanding of probability, and tolerance for uncertainty.

# Falsifiability (and opening up for the possibility of being wrong, c.f. "ontology")



**Prof. Feynman** @ProfFeynman · 7h

SCIENCE:

If you don't make mistakes, you're doing it wrong.

If you don't correct those mistakes, you're doing it really wrong.

If you can't accept that you're mistaken, you're not doing it at all.

💬 39          🔁 1 565          ♡ 6 740          ⬆

# What is SE practice based on?

- Mostly, the SE discipline is based on a combination of <span style="color:red">human authority</span> and <span style="color:red">anecdotal experience</span>:*

    – We know that a particular technique is good because John Doe, who is an authority in the field, says that it is good (human authority); and that

    – John Doe knows that it is good because it worked for him (anecdotal experience).

*C. Michael Holloway, Software Engineering and Epistemology, *Software Engineering Notes*, 1995, 20(2): 20-21.

# Current problems of software development

- The prevalence of fads more typical of the fashion industry than an engineering discipline.
- The lack of a sound, widely accepted theoretical basis.
- The huge number of methods and method variants, with differences little understood and artificially magnified.
- The lack of credible empirical evaluation and validation.
- The split between (software) industry and academia.

Jacobson, Ng, McMahon, Spence, and Lidman *The Essence of Software Engineering*, Addison-Wesley, 2013.

Demetri @PhDemetri · 7h

Undergrad: Here is the problem, find the solution

**Demetri** @PhDemetri · 7h

Undergrad: Here is the problem, find the solution

Masters: Here is part of the problem. How does the solution change when the problem changes?

**Demetri** @PhDemetri · 7h

Undergrad: Here is the problem, find the solution

Masters: Here is part of the problem. How does the solution change when the problem changes?

PhD: What is the problem and the solution?

Industry: Here is the solution, find the problem

💬 7     🔁 65     ♡ 423     ⬆️

# Delegated vs. centralized control – expert opinions

- **The Delegated Control Style:**
  - *Rebecca Wirfs-Brock*: A delegated control style ideally has clusters of well defined responsibilities distributed among a number of objects. To me, a delegated control architecture feels like object design at its best…
  - *Alistair Cockburn*: [The delegated coffee-machine design] is, I am happy to see, robust with respect to change, and it is a much more reasonable "model of the world."

- **The Centralized Control Style:**
  - *Rebecca Wirfs-Brock*: A centralized control style is characterized by single points of control interacting with many simple objects. To me, centralized control feels like a "procedural solution" cloaked in objects…
  - *Alistair Cockburn*: Any oversight in the "mainframe" object (even a typo!) [in the centralized coffee-machine design] means potential damage to many modules, with endless testing and unpredictable bugs.

Responsibility Driven Design

Role Modelling

Object1
Message3
Message1
Message2
Object5
Object2
Message4
Message6
Message5
Object4
Object3

Analytical

**Evaluating the effect of a delegated vs. centralized control style on the maintainability of object-oriented software**



"Assuming that it is not only highly skilled experts who are going to maintain an object-oriented system, a viable conclusion from the controlled experiment reported in this paper is that a design with a centralized control style may be more maintainable than is a design with a delegated control style."

Erik Arisholm and Dag Sjøberg, *IEEE Transactions on Software Engineering*, vol. 30, no. 8, August 2004, pp. 521-534.

FTD Stasjon
JAN MAYEN

TEORI ER NÅR MAN FORSTÅR ALT
MEN INGEN TING VIRKER

PRAKSIS ER NÅR ALT VIRKER
MEN INGEN FORSTÅR HVORFOR

PÅ DENNE STASJONEN FORENER VI TEORI OG PRAKSIS
SLIK AT INGEN TING VIRKER OG INGEN FORSTÅR HVORFOR

Theory is when one understands
everything, but nothing works

Practice is when everything works,
but no one understand why

At [Jan Mayen] we unite theory and
practice so that nothing works and
no one understands why
(my translation)

UiO : **Department of Informatics**
University of Oslo

# Theory and practice

"The research I have available claims the opposite … "

# Structure

- Empirical research
- Evidence-based software engineering
    - The steps of EBSE
    - Research synthesis
- Empirical research methods
    - Controlled experiments
    - Case studies
    - Surveys
    - Action research
    - Validity
- The importance of context

# Evidence-based software engineering

- Adapted from Evidence-Based Medicine
  - To provide the means by which <span style="color:red">current best evidence</span> from research can be integrated with <span style="color:red">practical experience</span> and human values in the decision making process regarding the development and maintenance of software

- EBSE sets requirements on practitioners and researchers:
  - Practitioners need to track down and use best evidence in context of practice
  - Researchers need to provide best evidence

**Archie Cochrane**

"It is surely a great criticism of our profession that we have not organised a critical summary, by specialty or subspecialty, adapted periodically, of all relevant randomised controlled trials."

# Software engineering challenges

- No comparable (to medicine) research infrastructure.

- No agreed standards for empirical studies

- Few software engineering guidelines based on empirical evidence.

- Challenges in addressing software engineering specifics
  - The skill factor
  - The lifecycle issue
  - The context dependences

# The five steps of EBSE:

1. Converting a relevant problem or information need into an answerable question.

2. Searching the literature for the best available evidence to answer the question.

3. Critically appraising the evidence for its validity, impact, and applicability.

4. Integrating the appraised evidence with practical experience and the values and circumstances of the customer to make decisions about practice.

5. Evaluating performance and seeking ways to improve it.

REPORT

SOCIAL NETWORKS

# A causal test of the strength of weak ties

Karthik Rajkumar[1], Guillaume Saint-Jacques[1], Iavor Bojinov[2], Erik Brynjolfsson[3,4], Sinan Aral[5]*

The authors analyzed data from multiple large-scale randomized experiments on LinkedIn's People You May Know algorithm, which recommends new connections to LinkedIn members, to test the extent to which weak ties increased job mobility in the world's largest professional social network. The experiments randomly varied the prevalence of weak ties in the networks of over 20 million people over a 5-year period, during which 2 billion new ties and 600,000 new jobs were created. The results provided experimental causal evidence supporting the strength of weak ties and suggested three revisions to the theory. First, the strength of weak ties was nonlinear. Statistical analysis found an inverted U-shaped relationship between tie strength and job transmission such that weaker ties increased job transmission but only to a point, after which there were diminishing marginal returns to tie weakness. Second, weak ties measured by interaction intensity and the number of mutual connections displayed varying effects. Moderately weak ties (measured by mutual connections) and the weakest ties (measured by interaction intensity) created the most job mobility. Third, the strength of weak ties varied by industry. Whereas weak ties increased job mobility in more digital industries, strong ties increased job mobility in less digital industries.

# Step 1: Asking an answerable question

- The first step in EBSE is to convert a relevant problem or information need into an answerable question.

- Typical questions ask for specific knowledge about how to appraise and apply methods, tools, and techniques in practice.

- Well formulated questions usually have three components:
  - The main intervention or action we are interested in.
  - The context or specific situations of interest.
  - The main outcomes or effects of interest.

- Example:
  - "Does the use of pair programming lead to improved code quality when practiced by professional software developers?"

# Step 2: Finding the best evidence

- Finding an answer to our question includes selecting an appropriate information resource and executing a search strategy.

- The main source of research-based evidence is articles published in scientific journals. Examples of databases that index published articles include:
  - IEEE *Xplore*, http://ieeexplore.ieee.org
  - ACM Digital Library, http://www.acm.org/dl
  - ISI Web of Science, http://isiknowledge.com
  - Google scholar

- Often, reading important magazines such as the *Communications of the ACM*, *IEEE Computer*, *IEEE Software*, and *IT Professional* would probably be enough to get a general overview of the latest developments within software engineering.

# Step 3: Critically appraising the evidence

- Unfortunately, published <span style="color:red">research isn't always of good quality</span>; the problem under study might be unrelated to practice or the research method could have weaknesses so that the results cannot be trusted.

- To assess whether research is of good quality and can be applied to practice, we must be able to <span style="color:red">critically appraise</span> the evidence.
  - Is there any vested interest?
  - Is the evidence valid?
  - Is the evidence important?
  - Can the evidence be used in practice?
  - Is the evidence in this study consistent with the evidence in other available studies?

# Step 4: Applying the evidence

- Active use of new knowledge is characterized by applying or adapting specific evidence to a specific situation in practice.

- Therefore, in order to practice EBSE, the individual software developer must commit him or herself to actively engage in a learning process, combining the externally transmitted evidence with prior knowledge and experience.

- Thus, it is at this point that EBSE needs to be integrated with process improvement.

- EBSE should provide the scientific basis for undertaking specific process changes while SPI should manage the process of introducing a new technology.

# Step 5: Evaluating performance

- We need to consider how well we perform each step in the EBSE process and how we might improve our use of EBSE.
    - In particular, we should ask ourselves how well we are integrating evidence with practical experience, customer requirements, and our knowledge of the specific circumstances.

- Following SPI practice, we also need to assess whether process change has been effective.
    - This might include After Action Reviews, Postmortem Analyses, and organization-wide measurement programs.

# What is research synthesis?

- Collective term for a family of methods for summarizing, integrating, combining, and comparing the findings of different studies on a topic or research question.

- Embodies the idea that individual studies or pieces of evidence are combined to produce a coherent whole, in the form of an argument, theory, or conclusions.

- It can provide conclusions with increased accuracy and less uncertainty compared to individual studies.

- A guiding principle is to be as rigorous and as transparent as possible.

**Confidence in research synthesis depends on body of evidence strength and quality in primary studies and synthesis**

- The confidence we can place in the conclusions and recommendations arising from a research synthesis depends on three issues:
    - The quality of the primary studies
    - The quality of the synthesis itself
    - The strength of the total body of evidence

# Synthesis Types



**Evidence**
- Qualitative
- Quantitative

**Qualitative Synthesis**
- **Interpreting** evidence to reach a new theoretical or conceptual level of understanding.

**Quantitative Synthesis**
- **Aggregating** evidence to reach statistical conclusions.

**Synthesis Results**
- Decision Support
- Knowledge Building

46

# What if we have weak evidence in?

# Systematic bias in publications – funnel plot (previous slide)

# What if we have a weak process of synthesis?

**Research Synthesis**

**Strong Evidence**

**Qualitative**

**Quantitative**

Qualitative Approaches:
    Narrative Synthesis
    Grounded Theory
    Meta-ethnography
    Thematic Synthesis
    …
Quantitative Approaches
    Content Analysis
    Case Survey
    Comparative Analysis
    Meta-Analysis
    …

**Synthesis Results**

THIS IS GARBAGE. DON'T RECYCLE IT.

# Empirical Research Methods

Bent Hamer: "Kitchen Stories," 2003. **Trailer**

# Structure

- Empirical research
- Evidence-based software engineering
  - The steps of EBSE
  - Research synthesis
- Empirical research methods
  - Controlled experiments
  - Case studies
  - Surveys
  - Action research
  - Validity
- The importance of context

# An alternative, supporting approach to study SE practice



Lars Mathiassen, Collaborative Practice Research, *Information Technology & People*, Vol. 15 No. 4, 2002, pp. 321-345.

# Controlled experiments



- An experiment is a study in which an intervention is deliberately introduced to observe its effects:
    - The identification of causal relations provides an explanation of *why* a phenomenon occurred.

    - The identification of casual processes yields an account of *how* a phenomenon occurred.

- Experiments are conducted when the investigator wants control over the situation, with direct, precise, and systematic manipulation of the behavior of the phenomenon to be studied.

- All experiments involve at least a treatment, an outcome measure, units of assignment, and some comparison from which change can be inferred and (hopefully) attributed to the treatment.

# Classical experimental design

- **Randomized experiment**
    - An experiment in which units are assigned to receive the treatment or an alternative condition by a random process

- **Quasi-experiment**
    - An experiment in which units are not assigned to conditions randomly



O1: Pre-intervention data collection point
O2: Post-intervention data collection point

# Advantages

- They are a well established strategy, seen by many as the 'scientific' and therefore most acceptable approach
- They are the only research strategy that can prove causal relationships
- Laboratory experiments permit high levels of precision in measuring outcomes and in analyzing data

# Disadvantages

- Laboratory experiments (e.g., with students at the university) often create artificial situations, which are not comparable with real-world situations
- It is often difficult or impossible to control all the relevant variables
- It is often difficult to recruit a representative sample of participants
- It may be necessary to conceal from the participants the purpose of the research so they do not skew the results

# Experiment – example

## What?

Research Question:

- What is best – Pair Programming or Solo Programming?

## Why?

Many studies with contradicting results – mostly conducted with students (not with professional developers).

Source:

E. Arisholm, H. Gallis, T. Dybå, and D. Sjøberg, "Evaluating Pair Programming with Respect to System Complexity and Programmer Expertise," *IEEE Transactions on Software Engineering*, 2007, 33(2): 65-86.

## Who, where and when?

- 295 junior, intermediate and senior professional Java consultants from 29 companies were paid to participate (one work day)
- Norway, Sweden, UK; 2001-2005
- 99 individuals, 98 pairs
- The pairs and individuals performed the same Java maintenance tasks on either:
    - a "simple" system (centralized control style), or
    - a "complex" system (delegated control style)
- We measured:
    - duration (elapsed time)
    - effort (cost)
    - quality (correctness) of their solutions

# Case study research

Case study research is an empirical inquiry that:

- Investigates a contemporary phenomenon within its real-life context, especially when
- the boundaries between phenomenon and context are not clearly evident.

Types of case studies:

- Singlecase, multicase
- Exploratory, descriptive, explanatory
- Holistic, embedded
- Qualitative, quantitative
- Positivist, interpretative, critical

# Advantages

- It can deal with complex situations where it is difficult to study a single factor in isolation
- It is appropriate for situations where the researcher has little or no control over events
- It is suitable for both theory building and theory testing
- It allows the researcher to show complexities and to explore alternative meanings and explanations
- It produces data that is close to people's experience

# Disadvantages

- It is sometimes seen as lacking rigor and leading to generalizations with poor credibility
- It can be difficult and time-consuming to negotiate access to the necessary settings, people and documents
- The presence of the researcher can affect how people behave
- There aren't really any rules to follow

# Case study – example

## What?

Research Question:

- What are the challenges of shared decision-making in agile software development teams?

## Why?

Agile software development changes the nature of collaboration, coordination, and communication in software projects

Source:

N.B. Moe, A. Aurum, and T. Dybå, "Challenges of Shared Decision-Making: A Multiple Case Study of Agile Software Development," Information and Software Technology, 2010, 54(8): 853-865.

## Who, where and when?

- Multiple case study of four projects in two software product companies
- Norway; 2007-2010
- Both companies recently adopted Scrum
  - One company introduced Scrum in the middle of two 3-year projects
  - One company introduced Scrum at the beginning of two 9-12 month projects
- We collected data in semi-structured interviews, through participant observations, and from process artifacts
- Data collected over a period of 11-12 months in all four projects

# Action Research

- Simultaneously contribute to the practical concerns in a concrete situation and to the goals of science.

- Dual commitment to study a system and concurrently to collaborate with members of the system in changing it.

- Active collaboration between researchers and practitioners underlines the importance of co-learning as a primary aspect of the research process.



64

**Action Research attempts to provide**

# practical value

**to the client organization while simultaneously contributing**

**to the acquisition of**

# new theoretical knowledge

# Criticisms of Action Research

- Action Research has been criticized for:

  - its lack of methodological rigor
  - its lack of distinction from consulting, and
  - its tendency to produce either

**'research with little action or action with little research'**

# Action research – example

**What?**

Research Question:

- What benefits and challenges can arise from introducing knowledge redundancy interventions based on job rotation in software development?

**Why?**

Establish a formalized support service and contribute to improved flexibility in project staffing by knowledge redundancy

Source:

T.E. Fægri, T. Dybå, and T. Dingsøyr (2010) "Introducing Knowledge Redundancy Practice in a Small Software Organization: Experiences with Job Rotation in Support Work," *Information and Software Technology*, 52(10): 1118-1132.

**Who, where and when?**

- Action research in one company to integrate organizational change with scientific inquiry.
- Norway; 2008
- The practical objectives were:
  - to establish customer support as a legitimate organizational function that would shield developers from support enquiries, and
  - to contribute to improved flexibility in project staffing by enabling overlapping product experience among developers.
- During a period of 18 weeks, nine developers rotated to customer support.
- We collected data in meetings, from comprehensive interviews, and from customer support work logs.

# Survey research

- A survey is useful for studying a large number of variables using a large sample size and rigorous statistical analysis.

- They are used when:
    - control of the independent and dependent variables is not possible or not desirable,
    - when the phenomena of interest must be studied in their natural setting, and
    - when the phenomena of interest occur in current time or the recent past.

# Example

**Teamwork Quality and Project Success in Agile Software Development: A Survey of Agile Development Teams**

Yngve Lindsjørn[a], Dag I.K Sjøberg[a,b], Torgeir Dingsøyr[b], Gunnar R. Bergersen[a], Tore Dybå[b,a]

[a] Department of Informatics, University of Oslo, Norway {ynglin, dagsj, gunnab} @ifi.uio.no
[b] SINTEF, Trondheim, Norway {torgeir.dingsoyr, tore.dyba}@sintef.no

*Table 9 - Items in Questionnaire*

| Construct (no of Items) | Items (Questions) |
|---|---|
| **Teamwork Quality (38)**<br><br>Communication (10) | 1. There is frequent communication within the team<br>2. The team members communicate often in spontaneous meetings, phone conversations, etc.<br>3. The team members communicate mostly directly and personally with each other<br>4. There are mediators through whom much communication is conducted (*)<br>5. Relevant ideas and information relating to the teamwork is shared openly by all team members<br>6. Important information is kept away from other team members in certain situations (*)<br>7. In the team there are conflicts regarding the openness of the information flow (*)<br>8. The team members are happy with the timeliness in which they receive information from other team members<br>9. The team members are happy with the precision of the information they receive from other team members<br>10. The team members are happy with the usefulness of the information they receive from other team members |
| Coordination (4) | 11. The work done on subtasks within the team is closely harmonized<br>12. There are clear and fully comprehended goals for subtasks within our team<br>13. The goals for subtasks are accepted by all team members<br>14. There are conflicting interests in our team regarding subtasks/subgoals (*) |
| Mutual Support (7) | 15. The team members help and support each other as best they can<br>16. If conflicts come up, they are easily and quickly resolved<br>17. Discussions and controversies are conducted constructively<br>18. Suggestions and contributions of team members are respected<br>19. Suggestions and contributions of team members are discussed and further developed<br>20. The team is able to reach consensus regarding important issues<br>21. The team cooperate well |
| Effort (4) | 22. Every team member fully pushes the teamwork<br>23. Every team member makes the teamwork their highest priority<br>24. The team put(s) much effort into the teamwork<br>25. There are conflicts regarding the effort that team members put into the teamwork (*) |
| Cohesion (10) | 26. The teamwork is important to the team<br>27. It is important to team members to be part of the team<br>28. The team does not see anything special in this teamwork (*)<br>29. The team members are strongly attached to the team<br>30. All team members are fully integrated in the team<br>31. There were many personal conflicts in the team (*)<br>32. There is mutual sympathy between the members of the team<br>33. The team sticks together<br>34. The members of the team feel proud to be part of the team<br>35. Every team member feels responsible for maintaining and protecting the team |
| Balance of member Contribution (3) | 36. The team recognizes the specific characteristics (strengths and weaknesses) of the individual team members<br>37. The team members contribute to the achievement of the team's goals in accordance with their specific potential<br>38. Imbalance of member contributions cause conflicts in our team (*) |

# Common in society

- Requires relatively *few resources* to include many people

- Create statistics and test hypotheses over characteristics of the target group (the population being investigated)

- Obtain information about people's *opinion* about what, how much, how many, how and why or what people *say* they do

  - As opposed to experiments, one does *not control* independent and dependent variables

  - As opposed case studies and ethnography, one does *not observe*

# Types of surveys

- **Cross-sectional surveys** are used to gather information on a population at a single point in time.

- **Longitudinal surveys** gather data over a period of time. The researcher may then analyze changes in the population and attempt to describe and/or explain them.

  – *Trend studies* focus on a particular population, which is sampled and scrutinized repeatedly. While samples are of the same population, they are typically not composed of the same people.

  – *Cohort studies* also focus on a particular population, sampled and studied more than once. A cohort study would sample the same group of people, every time.

  – *Panel studies* allow the researcher to find out why changes in the population are occurring, since they use the same sample of people every time.

# Sampling and generalization



Who do you want to generalize to? → The Theoretical Population

What population can you get access to? → The Study Population

How can you get access to them? → The Sampling Frame

Who is in your study? → The Sample

W.M. Trochim, The Research Methods Knowledge Base, http://www.socialresearchmethods.net/kb/

# Types of questions

- All researchers must make two basic decisions when designing a survey – they must decide:
  1. whether they are going to employ an oral, written, or electronic method, and
  2. whether they are going to choose questions that are open or close-ended.
- We will focus on **written** and **close-ended** methods.

# The importance of wording …

Two catholic priests wondered if one is allowed to smoke when one prays? They both sent a letter to the Pope:

**P1: "Is it allowed to smoke when one prays?"**
Answer: **NO** – the pray should get full attention

**P2: "Is it allowed to pray when one smokes?"**
Answer: **YES** – it is always a good thing to pray

# Question formats

- Classification of objects or individuals.

  "Are you: Male___  Female ___ Other ___?"

- Ranking of items in order to reflect the relative ordering of phenomena.

  "Please rank the following factors in order of importance (1-4)"

- Pairwise comparison

  "Which do you prefer"

- Rating of characteristics

  – Simple, single-item scales, e.g.,

    "Programming is a terrific course (check one)"

    "Strongly agree__, agree__, neither__, disagree__, strongly disagree__"

# Likert type scales

- **Evaluation-type**

  **Example:**

  – "Familiarity with and comprehension of the software development environment"

  - ❑ Little
  - ❑ Unsatisfactory
  - ❑ Neutral
  - ❑ Satisfactory
  - ❑ Excellent

- **Frequency-type**

  **Example:**

  – "Customers provide information to the project team about the requirements"

  - ❑ Never
  - ❑ Rarely
  - ❑ Neutral
  - ❑ Occasionally
  - ❑ Most of the time

- **Agreement-type**

  **Example:**

  – "The tasks supported by the software at the customer site change frequently"

  - ❑ Strongly Agree
  - ❑ Agree
  - ❑ Neutral
  - ❑ Disagree
  - ❑ Strongly Disagree

Dag Sjøberg

**UNIVERSITETE' I OSLO**

# Advantages

- They provide a wide an inclusive coverage of people or events
- They can be administered from remote locations using mail, email or telephone
- They can provide a lot of data in a short time at a reasonable cost
- They lend themselves to quantitative analysis
- They can be replicated
- Usually, high reliability is easy to obtain

# Disadvantages

- They lack depth
- They tend to focus on what can be counted or measured
- They do not establish cause and effect
- They cannot judge the accuracy or honesty of people's responses by observing their body language

**Important dimensions of empirical methods:**
- obtrusiveness,
- generality,
- artificiality, and
- point of maximum concern



Obtrusive Research Operations

B

Laboratory Experiments

Experimental Simulations

Judgment Tasks

II  II

III  I

Field Experiments

Sample Surveys

III  I

IV  IV

Field Studies

Formal Theory

Computer Simulations

A

Unobtrusive Research Operations

Universal Behavior Systems

Particular Behavior Systems

I. Settings in natural systems.
II. Contrived and created settings.
III. Behavior not setting dependent.
IV. No observation of behavior required.

A. Point of maximum concern with generality over actors.
B. Point of maximum concern with precision of measurement of behavior.
C. Point of maximum concern with system character of context.

Philip J. Runkel & Joseph E. McGrath, Research on human behavior: Systematic guide to method. New York: Holt, Rinehart and Wilson, 1972, p. 85.

# Selecting the research method

| Research question | Controlled experiment | Longitudinal survey | Cross-sectional survey | Case study Action research |
|---|---|---|---|---|
| *Effectiveness*: Does it work? Does method A work better than method B? | ++ | + | - | -- |
| *Explanation*: How does it work? Why does it work? | -- | - | + | ++ |
| *Context*: In what circumstances does it work, for whom? | -- | - | + | ++ |
| *Safety*: Will it do more good than harm? | ++ | + | + | + |
| *Acceptability*: Will the target group accept the new method of working? | -- | - | + | ++ |
| *Prevalence*: How often is this method/ technique applied/implemented? | -- | -- | ++ | -- |
| *Appropriateness*: Is this the right process/method for this target group? | -- | - | + | ++ |

Adapted from cebma.org

"The purpose of computing is insight, not numbers."

Richard Hamming

# Structure

- Empirical research
- Evidence-based software engineering
  - The steps of EBSE
  - Research synthesis
- Empirical research methods
  - Controlled experiments
  - Case studies
  - Surveys
  - Action research
  - Validity
- The importance of context

# Reliability and <u>validity</u> of empirical studies

- **Reliability**
  - Can the study can be repeated (i.e., by other researchers) and yield the same results?

- **Statistical conclusion validity**
  - Is the statistical inference valid?

- **Internal validity**
  - Does the observed covariation between *A* (the presumed treatment) and *B* (the presumed outcome) reflects a causal relationship from *A* to *B?*

- **Construct validity**
  - Do the measures in the study represent the (abstract, possibly theoretical) constructs they are intended to measure?

- **External validity**
  - Does the cause–effect relationship hold over variations in persons, settings, treatment variables, and measurement variables?

# The four validities

See Sjøberg & Bergersen (2023) - Construct validity in SE, figure 4, for a better picture
https://ieeexplore.ieee.org/document/9780058

# The quality of empirical studies

Three methodological features have been shown to influence the results of primary studies:

- **Randomization** can avoid selection bias by making sure that each subject in the study has an equal chance of getting into each treatment group.

- **Blinding** of study participants and personnel may reduce the risk that knowledge of which treatment was received, rather than the treatment itself, affects outcomes and outcome measurements.

- **Missing outcome data**, due to attrition (withdrawal and dropout) during the study or exclusions from the analysis, raise the possibility that the observed effect estimate is biased.

# Context: What is best? Bicycle or helicopter?

# Structure

- Empirical research
- Evidence-based software engineering
  - The steps of EBSE
  - Research synthesis
- Empirical research methods
  - Controlled experiments
  - Case studies
  - Surveys
  - Action research
  - Validity
- The importance of context

Understand which technologies that cause which outcomes in which situations, e.g.:

– When is technique X more efficient than technique Y?

– What resources are needed to use method X in a given situation?

– How to tailor process X to the actual situation?

# What is best?
## Pair programming or solo programming*

- 295 junior, intermediate and senior professional Java consultants from 29 companies were paid to participate (one work day)

- 99 individuals; 98 pairs

- The pairs and individuals performed the same Java maintenance tasks on either:
  - a "simple" system (centralized control style), or
  - a "complex" system (delegated control style)

- We measured:
  - duration (elapsed time)
  - effort (cost)
  - quality (correctness) of their solutions

*E. Arisholm, H. Gallis, T. Dybå, and D. Sjøberg, "Evaluating Pair Programming with Respect to System Complexity and Programmer Expertise," *IEEE Transactions on Software Engineering*, 2007, 33(2): 65-86.

**Total Effect of PP**

Bar chart titled "Total Effect of PP". The y-axis is labeled "Difference from individuals" ranging from -40 % to 160 %. Three bars: Duration -8 %, Effort 84 %, Correctness 7 %.

**Effect of PP for Juniors**

**Moderating Effect of System Complexity for Juniors**

Legend:
- CC (easy)
- DC (complex)

Y-axis: Difference from individuals (-40 % to 160 %)

Duration: CC 4 %, DC 6 %
Effort: CC 109 %, DC 112 %
Correctness: CC 32 %, DC 149 %

96

**Effect of PP for Seniors**

Difference from individuals

- Duration: -9 %
- Effort: 83 %
- Correctness: -8 %

**Moderating Effect of System Complexity for Seniors**

Legend:
- CC (easy)
- DC (complex)

Y-axis: Difference from individuals

| Category | CC (easy) | DC (complex) |
|---|---|---|
| Duration | -23 % | 8 % |
| Effort | 55 % | 115 % |
| Correctness | -13 % | -2 % |

# So, when should we use PP?

| Programmer Expertise | Task Complexity | Use PP? | Comments |
|---|---|---|---|
| Junior | Easy | Yes | Provided that increased quality is the main goal |
| | Complex | Yes | Provided that increased quality is the main goal |
| Intermediate | Easy | No | |
| | Complex | Yes | Provided that increased quality is the main goal |
| Expert | Easy | No | |
| | Complex | No | Unless you are sure that the task is too complex to be solved satisfactorily even by solo seniors |

**The question of whether PP is beneficial, or not, is meaningless!**

# Important dimensions of SE context

Omnibus context:

◆ What?     ◆ !Who?     ◆ !Where?     ◆ !When?     ◆ !Why?!

'Phenomenon    'Subjects    'Loca1on    'Time    'Ra1onale'

Discrete context:

◆ 'Technical:      ◆ 'Social:      ◆ 'Environmental:!

'Complexity     'Individual'skill     'Uncertainty'
'Technology     'Team autonomy     'Community'
'Task/system     'Organizational structure     'Market'
'Etc…     'Etc…     'Etc…    "

Dybå, T., Sjøberg, D.I.K., and Cruzes, D.S. (2012) "What Works for Whom, Where, When, and Why?
On the Role of Context in Empirical Software Engineering," *Proceedings, ESEM 2012*, pp. 19-28.

# The research process



1. Select topic
2. Define question
3. Design study
4. Collect data
5. Analyze data
6. Interpret data
7. Publish results

| | |
|---|---|
| 1. Research problem<br>   a. Background and rationale<br>   b. Objectives and/or hypotheses | • What is the background of this investigation?<br>• What is the current status of research in this field?<br>• What is the purpose of the study and/or the question being asked? |
| 2. Research context<br>   a. Site selection<br>   b. Personnel<br>   c. Trial period | • What will the site and context of the study be?<br>• What personnel will be needed to conduct the study?<br>• What are their skills and experience?<br>• What is the approximate time schedule for carrying out the study? |
| 3. Study design<br>   a. Variables<br>   b. Design configuration<br>   c. Subject assignment<br>   d. Control of confounding variables | • What are the independent and dependent variables of the study?<br>• How will subjects be assigned to treatments?<br>• How many observations will you have for each treatment?<br>• What confounding variables will be controlled for? |
| 4. Treatment characteristics<br>   a. Description<br>   b. Tasks<br>   c. Duration | • What is the study treatment?<br>• What will you compare it with?<br>• What specific tasks will the subjects perform?<br>• Are they representative of what you want to study?<br>• How will the tasks be ordered? At random? |
| 5. Subject characteristics<br>   a. Selection criteria<br>   b. Representativeness of sample<br>   c. Subject recruitment<br>   d. Subject compliance | • What is the population to be studied?<br>• What steps will you take to ensure that your sample is representative and inclusive?<br>• How will subjects be recruited and selected?<br>• How will you measure their skills and experience? |
| 6. Data collection<br>   a. Scope of data collection<br>   b. Data collection procedure<br>   c. Data collection schedule<br>   d. Data reliability and validity | • What data will be collected?<br>• How and when will it be collected?<br>• How will completeness and accuracy be ensured?<br>• Who will collect the data?<br>• How will the data be stored? |
| 7. Data analysis<br>   a. Data preparation<br>   b. Data presentation<br>   c. Statistical analysis<br>   d. Data synthesis | • What are the expected results?<br>• What will be compared to what?<br>• What sort of analyzes will you do?<br>• How will you perform the analyzes?<br>• How will you aggregate and synthesize the data? |

**A research protocol** is a detailed description of how and why the research will be carried out.

# Selecting the research method

- The choice of method depends among other things on:
  - Suitable study subject (e.g., do participants have enough experience?)
  - Possibility to control the environment
  - The size/scale/cost of the study
  - The need for generality in the results
  - Availability of information/data and other resources
  - What is the purpose of the study? (exploration, prediction, understanding of cause-effect relations, applicability of results in industry, ....)

- Difficult to provide general recommendation with respect to choice of method, however …

# Summary: empirical research methods

- Empirical research is a foundation for modern society
- Synthesis of evidence-based approaches (i.e., good studies) is best
- Different research methods can provide such evidence, albeit with different strengths and weaknesses

Combined, separating "what works" from "what doesn't work" is a goal of both industry and academia.

# EBSE – literature



106

# Empirical research methods – literature