# Evidence over Opinion:
# An Empirical Approach to Software Engineering

## Prof. Dr. Sebastian Baltes

empirical-software.engineering
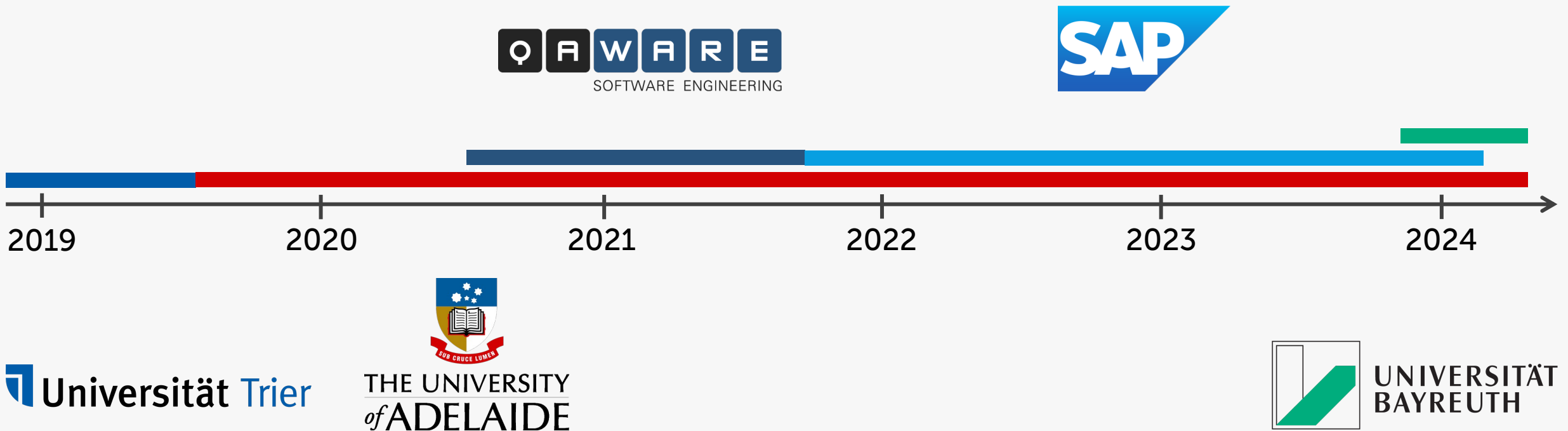
UNIVERSITÄT BAYREUTH

BayLDS

Lecture Series Digital Sciences at UBT
February 2024

# Hi, my name is Sebastian and I'm a pracademic*



* https://en.wikipedia.org/wiki/Pracademic

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

# A Brief History of Software Engineering

# Origin of the term "Software Engineering"



Margaret Hamilton
(1965-1969)

Anthony G Oettinger
(1966)

NATO Software Engineering Conference
(1968)

# The 1968 NATO SE Conference Report



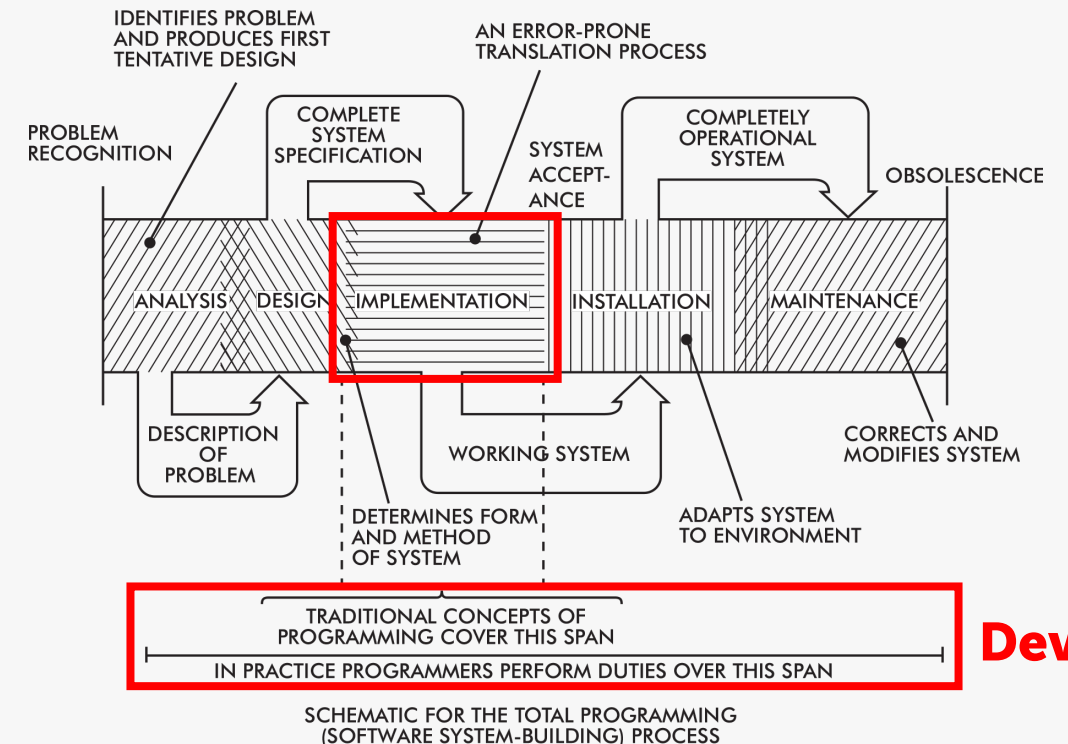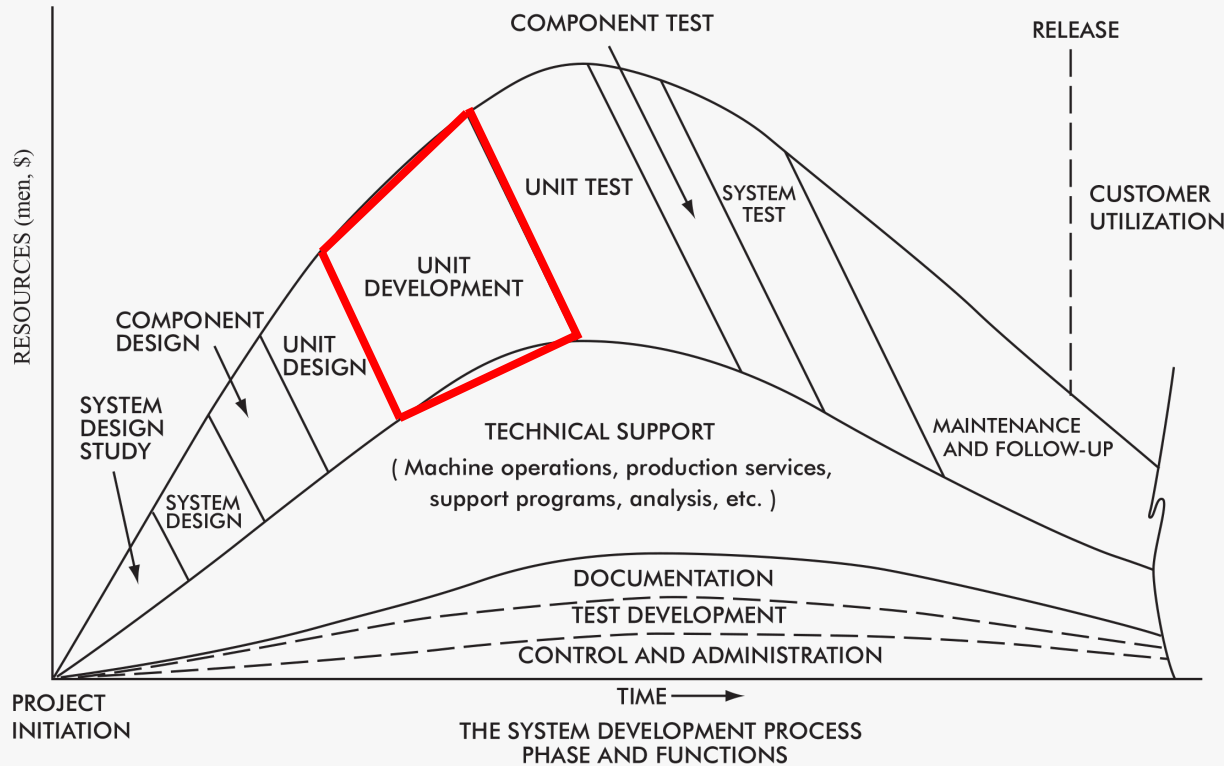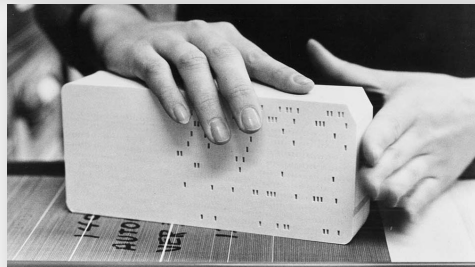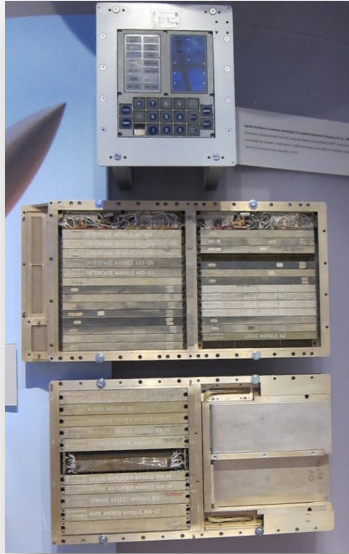Figure 1. From Nash: Some problems in the production of large-scale software systems.

Figure 2. From Selig: Documentation for service and users. Originally due to Constantine.

http://homepages.cs.ncl.ac.uk/brian.randell/NATO/nato1968.PDF

DevOps

# A lot has changed since the 1960s



AWS Lambda

Amazon Aurora

**1966:** Apollo guidance computer

**2024:** Embedded systems still relevant, but way more application domains and abstractions of hardware and software.

While many of the **fundamental concepts** introduced in the first years of software engineering are **still valid**, the application domains, stakeholders, tools, infrastructure, and processes today are **more diverse** than ever before.

Disciplinary Boundaries of Software Engineering

# Disciplinary Boundaries of Software Engineering

With a **traditional view** emphasizing software engineering's **roots in computer and systems engineering** many questions of modern software development **cannot be answered**.
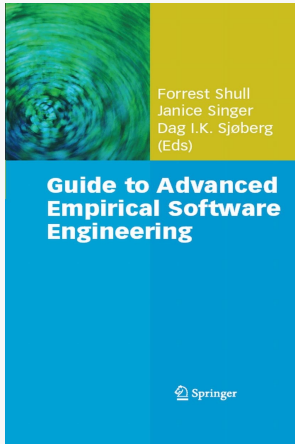
Examples:

- *How can we develop visual programming environments without knowledge of cognition?*

- *How can we study pair or mob programming without a deep understanding of verbal and non-verbal communication?*

- *How can we fully grasp the implications of AI-generated code without understanding copyright legislation and software licenses?*
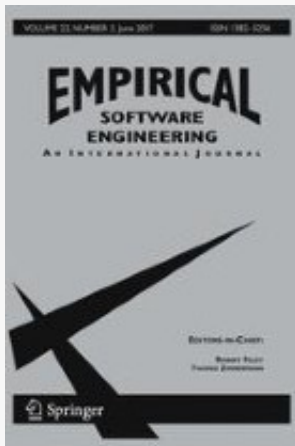
# Personal Observations

1. Many problems **relevant** in the **software industry** are rooted in software engineering but often have an **interdisciplinary** angle.

2. To impact industry, academia needs to provide **actionable recommendations** addressing **problems rooted in practitioners' actual needs**.

3. **Empirical research methods are essential** for identifying such problems (*problem space*) and corroborating recommendations/proposed solutions with empirical evidence (*solution space*).

# Empirical Software Engineering

2008: "Ten years ago, it was **rare to see a conference or journal article about a software development tool or process that had empirical data to back up the claims**. Today, in contrast, it is becoming **more and more common** […]."

https://link.springer.com/book/10.1007/978-1-84800-044-5

2020: "[…] **it has become clear that empirical studies are a fundamental component of software engineering research and practice** […]."

https://www.springer.com/journal/10664/aims-and-scope

# Institutionalized Boundaries

# Institutionalized Boundaries



"If you were using **MDSE\*** for building your mobile app, you'd see huge quality improvements, see papers x, y, z."

"Have you heard about things like **time-to-market** and quickly responding to customer feedback? We're not building safety-critical software."
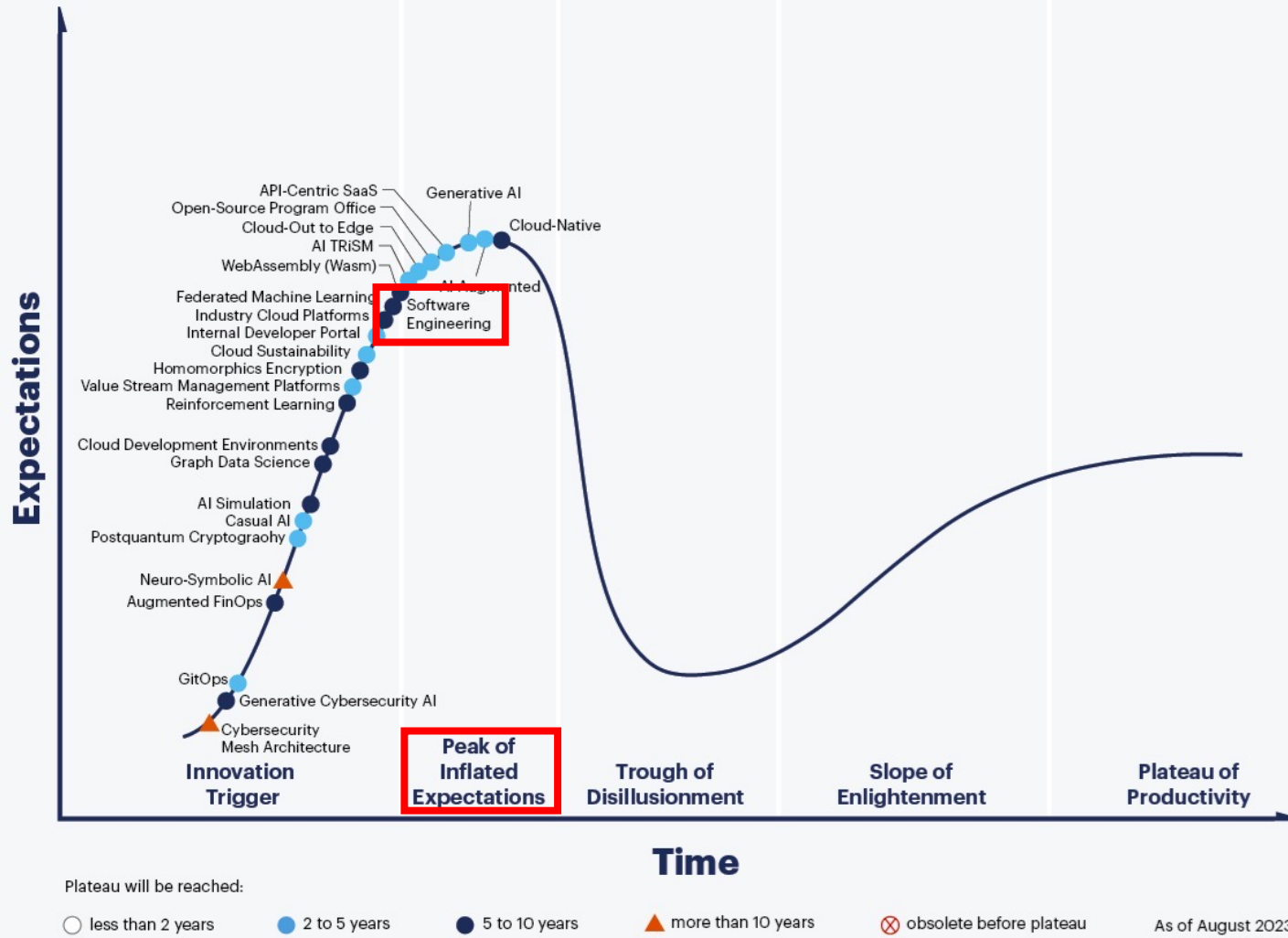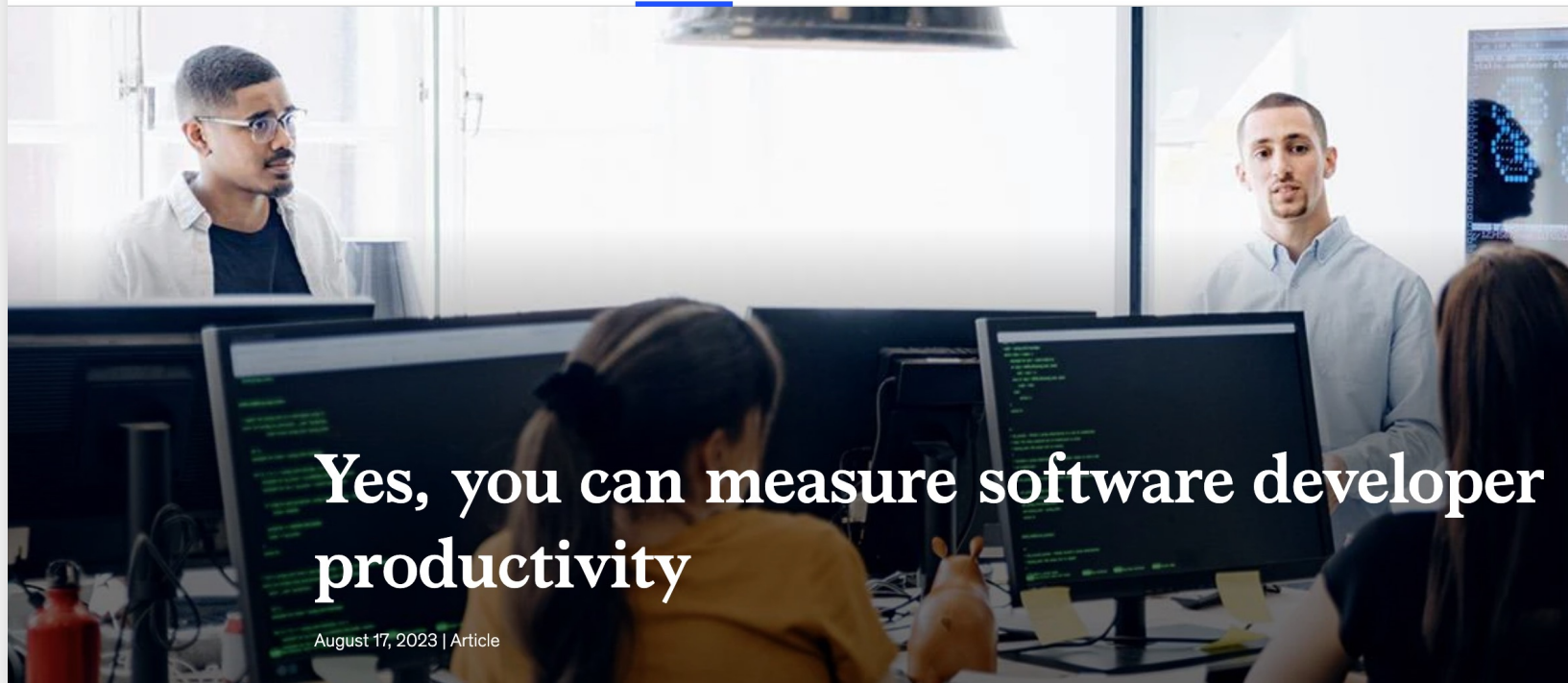
**Research**

**Practice**

# Institutionalized Boundaries

*"If you were using **MDSE\*** for building your mobile app, you'd see huge quality improvements, see papers x, y, z."*

*"Have you heard about things like **time-to-market** and quickly responding to customer feedback? We're not building safety-critical software."*

## Research

## Practice

**Issue in practice:**
Hype-driven Software Engineering

Prof. Dr. Sebastian Baltes – Evidence over Opinion: An Empirical Approach to Software Engineering

**McKinsey & Company**

**Technology, Media & Telecommunications**

How We Help Clients    Our Insights    Our People    Contact Us

# Yes, you can measure software developer productivity

August 17, 2023 | Article

https://www.mckinsey.com/industries/technology-media-and-telecommunications/our-insights/yes-you-can-measure-software-developer-productivity

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

Yes, you can
productivit

August 17, 2023 | Article

https://www.mckinsey.com/indus
our-insights/yes-you-can-measur

WS 2023/24  /  Seminar  /  Software Development Productivity

# Software Development Productivity

Der Begriff "Software Development Productivity" bezieht sich auf die Effizienz und Effektivität, mit der einzelne Softwareentwickler/-innen oder ganze Teams Softwaresysteme erstellen und warten. Die Relevanz dieses Themenbereichs in einem Umfeld von sich schnell wandelnden Anforderungen und Marktbedingungen und einem gleichzeitigen Mangel an Softwareentwickler/-innen ist offensichtlich. Allerdings ist Produktivität in der Softwareentwicklung nur sehr schwer allumfassend quantifizierbar, da das reine "Produzieren" von Quellcode einen relativ kleinen Teil der Softwareentwicklung ausmacht. Daher sind traditionelle Ansätze wie das Messen der Anzahl geschriebener Codezeilen (LoC, Lines of Code) pro Zeiteinheit unzureichend, da sie die tatsächliche Qualität und den Nutzen der erstellten Software nicht angemessen erfassen. Insbesondere im unternehmerischen Umfeld sorgt diese Tatsache für Unmut, ein Umstand den das Beratungsunternehmen McKinsey aktuell für sich auszunutzen versucht. Zwei Antworten auf McKinseys Behauptung man könne Produktivität in der Softwareentwicklung messen, bieten einen guten Einstieg in das Themenfeld:

https://ubt-se.github.io/teaching/semesters/23-24_ws/seminar/productivity.html

# Institutionalized Boundaries



*"If you were using **MDSE\*** for building your mobile app, you'd see huge quality improvements, see papers x, y, z."*
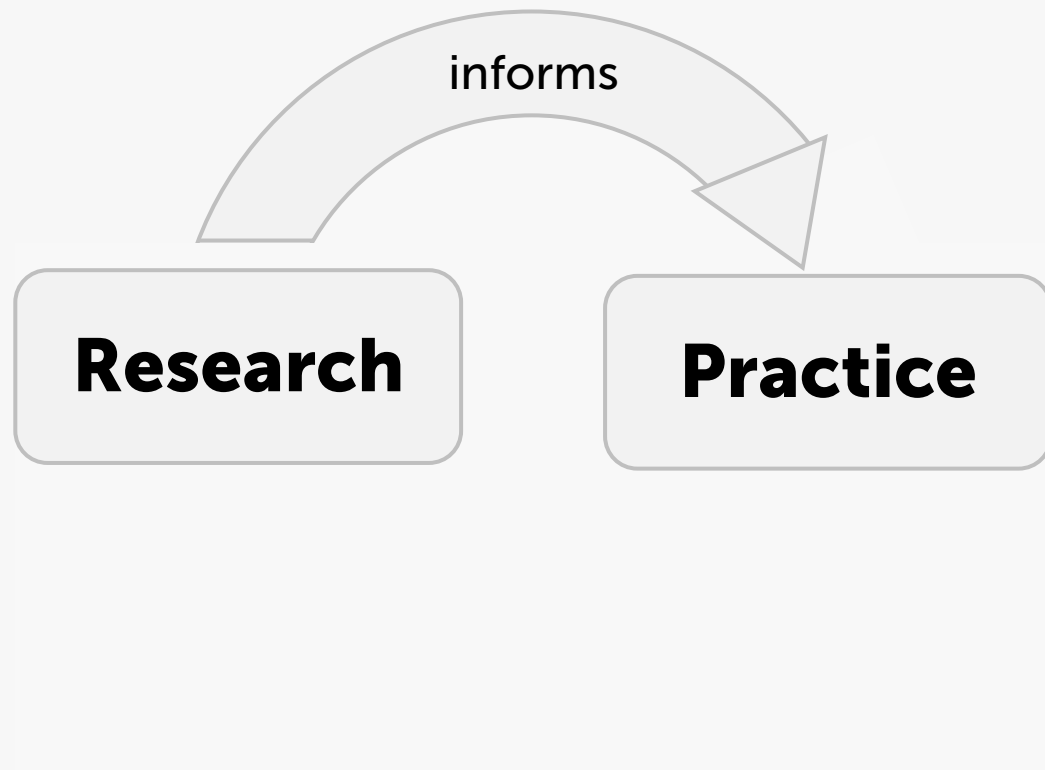
*"Have you heard about things like **time-to-market** and quickly responding to customer feedback? We're not building safety-critical software."*
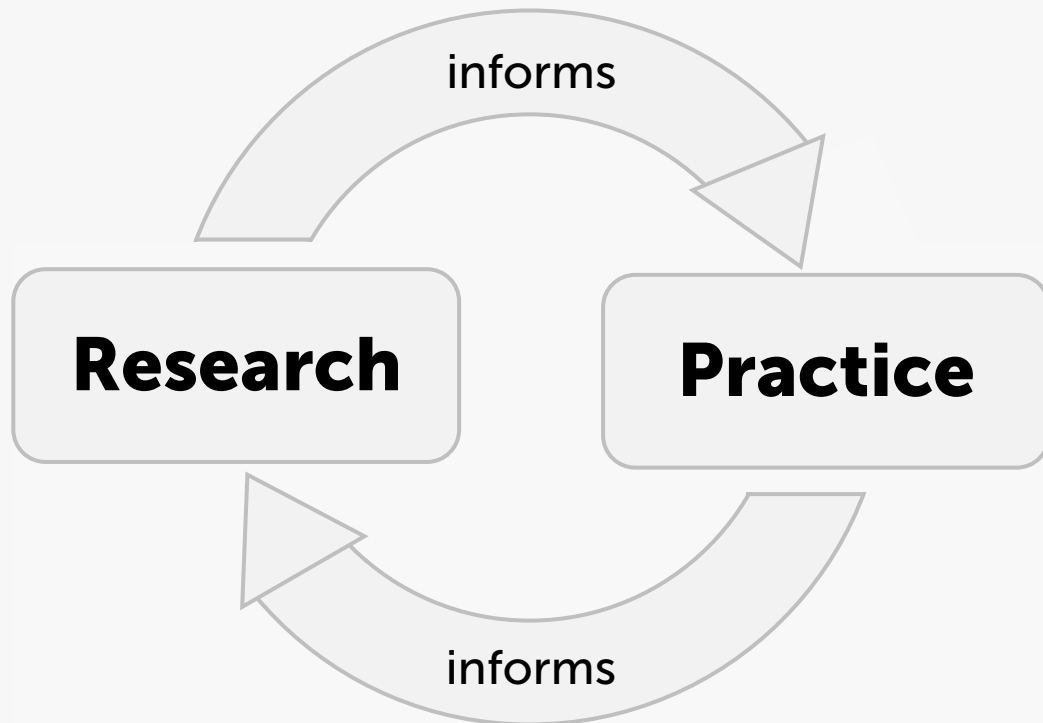
**Research**

**Practice**

# Institutionalized Boundaries

# Institutionalized Boundaries



*Implications for researchers:*

1) Strong understanding of **state of practice** is essential.

2) To reach this understanding, we need to utilize **diverse empirical research methods** and **learn from other disciplines.**

3) To advance evidence-based practice, we need to invest effort into **communicating findings/potential solutions back** to practitioners.

# An Example from Industry: Tech Stack Harmonization

Paraphrased excerpt from a status report to the SAP CTO:

- **Context:** SAP has a wide range of languages, tools, and infrastructure ("technology stacks") being used to build, test, and deploy products and services.

- **Problem:** Maintaining too heterogenous tech stacks is expensive and inefficient, it also increases the company's attack surface.

- **Progress:**
  - Derived a layered tech stack model based on interviews and a survey.
  - Used that model to build a data mining pipeline and dashboard for visualizing core tech stack aspects.
  - Discussed findings with product teams after a live demo.
  - Developed proposal for <CONFIDENTIAL>.

1. Strong understanding of state of practice.

2. Diverse empirical research methods.

3. Communicating findings back.

4. Solution based on research results.

# An Example from Industry: Tech Stack Harmonization

# An Example from Industry: Service Dependencies



Visually Analyzing Company-wide Software Service Dependencies: An Industrial Case Study

Fig. 1. SAP-managed service dependencies of a large organizational unit. Node color encodes native cloud environment of a service. Two large clusters correspond to the older (yellow) and newer (green, blue) environments.
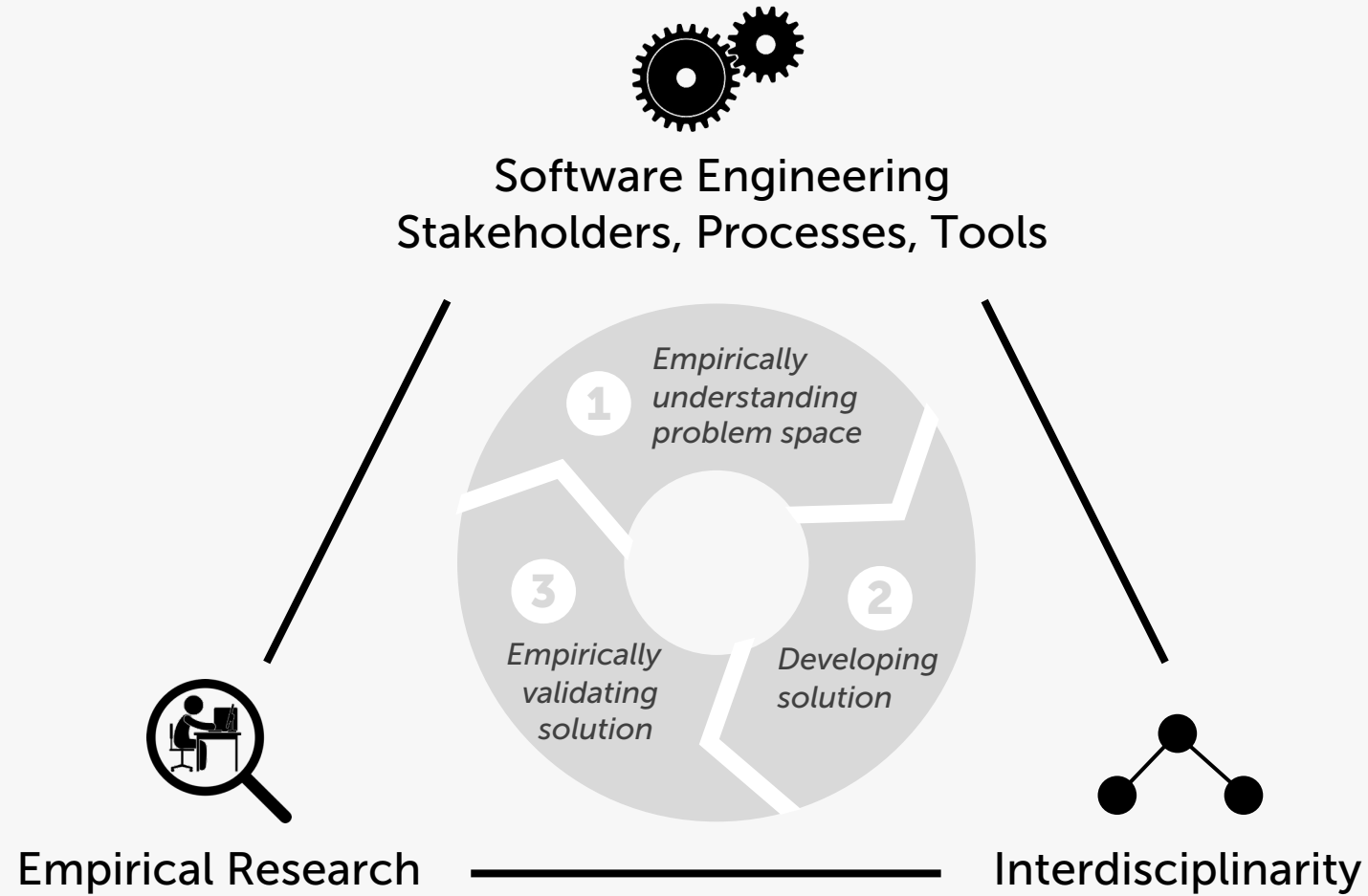
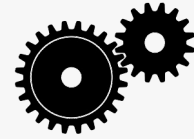https://empirical-software.engineering/publications/#vissoft23-service-dependency-viz
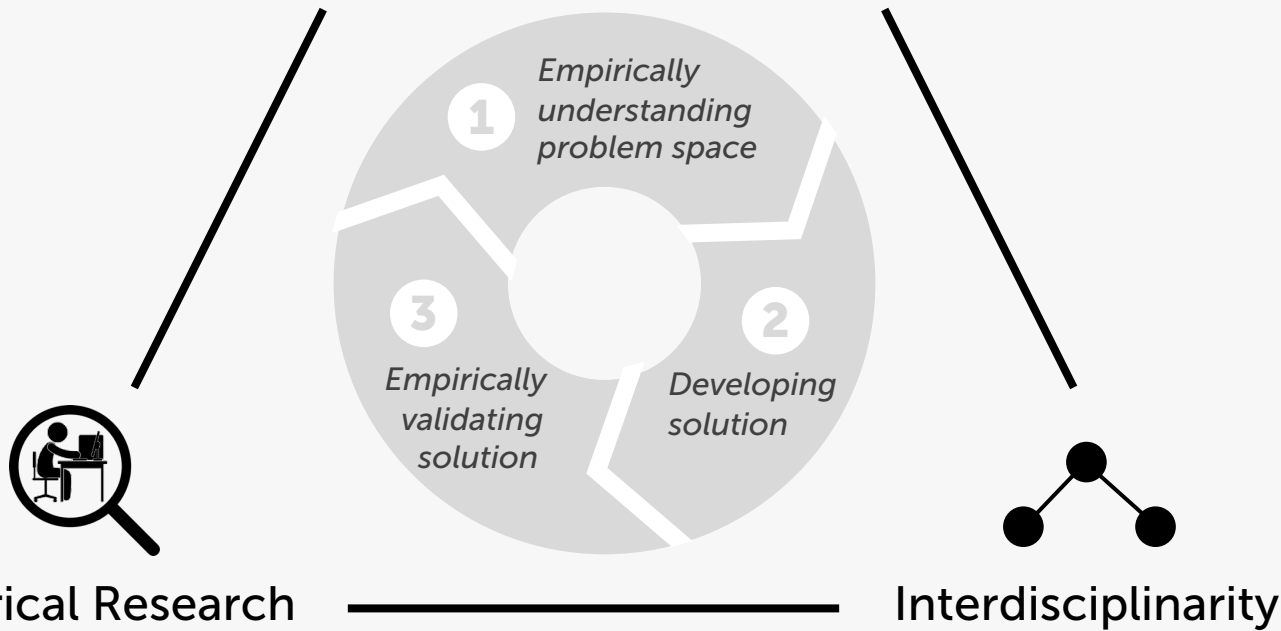
My framework for
Empirical Software Engineering Research
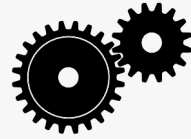
# Empirical Software Engineering



Software Engineering
Stakeholders, Processes, Tools

1 Empirically understanding problem space

2 Developing solution

3 Empirically validating solution

Empirical Research

Interdisciplinarity

# Remainder of this talk



Software Engineering
Stakeholders, Processes, Tools

1 Empirically understanding problem space

2 Developing solution

3 Empirically validating solution

Empirical Research

Interdisciplinarity

# Remainder of this talk



Sketching

*FSE '14, ESEM '15, VISSOFT '17*

Software Engineering
Stakeholders, Processes, **Tools**

1 *Empirically understanding problem space*

2 *Developing solution*

3 *Empirically validating solution*

Empirical Research

Interdisciplinarity

# Remainder of this talk

**Software Engineering**
Stakeholders, **Processes**, Tools



**Empirical Research**

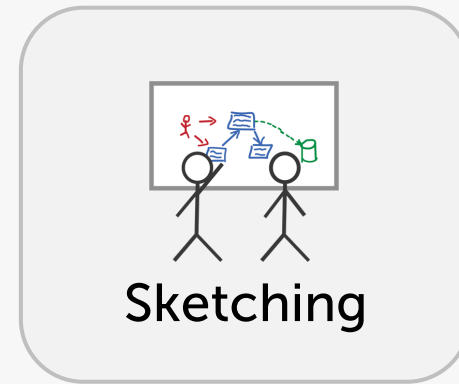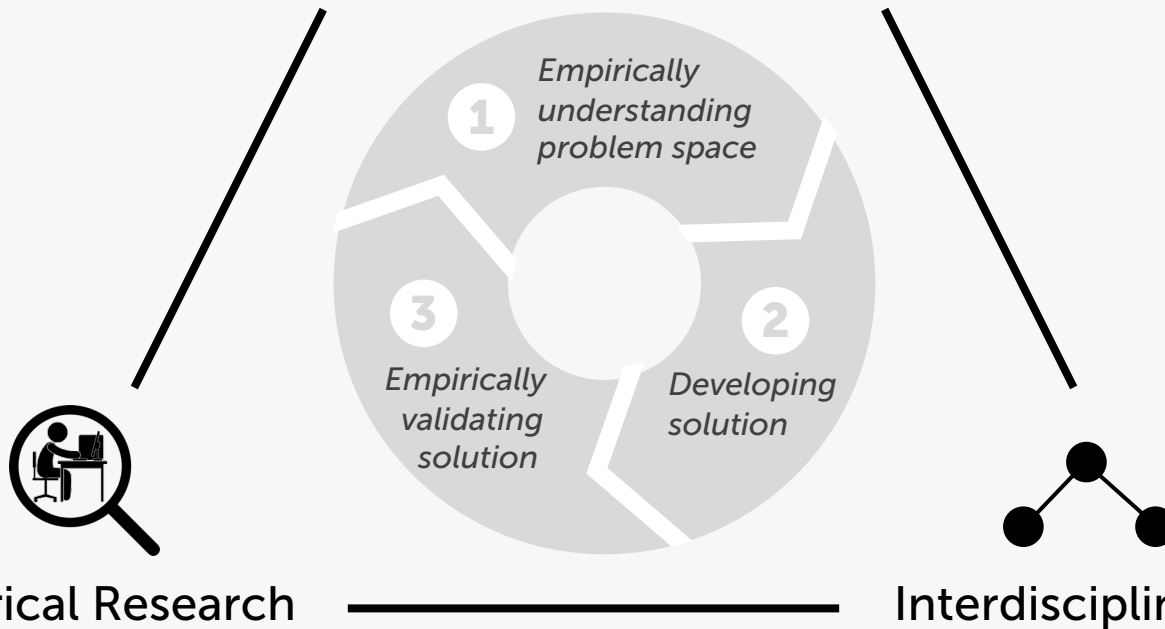1 *Empirically understanding problem space*

2 *Developing solution*

3 *Empirically validating solution*

**Interdisciplinarity**

**Sketching**
*FSE '14, ESEM '15, VISSOFT '17*

**stackoverflow**
**Code Plagiarism**
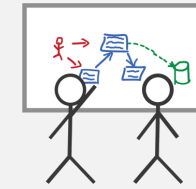*EMSE '18, MSR '18, MSR '19, ICSE '20*

# Remainder of this talk



Software Engineering
**Stakeholders**, Processes, Tools

1 Empirically understanding problem space

2 Developing solution

3 Empirically validating solution

Empirical Research

Interdisciplinarity

Sketching

*FSE '14, ESEM '15, VISSOFT '17*

Code Plagiarism

*EMSE '18, MSR '18, MSR '19, ICSE '20*

**Pandemic Programming**

*EMSE '20*

# Pandemic Programming

# Impact beyond SE

**Philosophical Psychology** ›
Latest Articles

| 76 | 0 | 2 |
| Views | CrossRef citations to date | Altmetric |

Research Article

## Perceived threat of COVID-19, self-assessment of physical health and mental resilience

Analyzing the impact of agile mindset adoption on software development teams productivity during COVID-19

Chaitanya Arun Sathe, Chetan Panse ▼

Journal of Advances in Management Research

ISSN: 0972-7981
Article publication date: 31 October 2022

DOWNLOADS
50

Reprints & Permissions

### Working in the digital economy: A systematic review of the impact of work from home arrangements on personal and organizational performance and productivity

Amy Hackney, Marcus Yung, Kumara G. Somasundram, Behdin Nowrouzi-Kia, Jodi Oakman, Amin Yazdani

Published: October 12, 2022 • https://doi.org/10.1371/journal.pone.0274728

**Applied Ergonomics**
Volume 102, July 2022, 103749

ELSEVIER

## The effect of training and workstation adjustability on teleworker discomfort during the COVID-19 pandemic

Megan J. McAllister [a], Patrick A. Costigan [a], Joshua P. Davies [a], Tara L. Diesbourg [b]

### "In the office nine to five, five days a week... those days are gone": qualitative exploration of diplomatic personnel's experiences of remote working during the COVID-19 pandemic

Samantha K. Brooks, Charlotte E. Hall, Dipti Patel & Neil Greenberg

BMC Psychology 10, Article number: 272 (2022) | Cite this article

708 Accesses | 6 Altmetric | Metrics

Occup Environ Med. 202
1.1097/JOM.00000000

## Digital Talent Management pp 29–

## Cultural Adaptation and Validation of the Health a Performance Questionnaire in German

Golz, Christoph MScN; Gerlach, Maisa MScN; Kilcher, Gablu MD, MPH; Peter, Karin Anne PhD
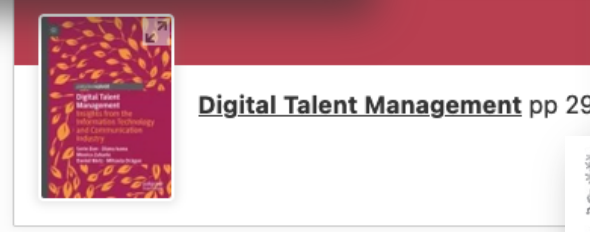
Author Information ☑

## Remote Work in a Changing World: A Nod to Personal Space, Regulation and Other Health and Wellness Strategies

Sybil Geldart ☑

Psychology Program, Wilfrid Laurier University Brantford Campus, Brantford, ON N3T 2Y3, Cana

t. J. Environ. Res. Public Health 2022, 19(8), 4873; https://doi.org/10.3390/ijerph19084873

eceived: 7 March 2022 / Revised: 29 March 2022 / Accepted: 15 April 2022 / Published: 17 A

ith Performan

oph Golz [1], Maisa Ger

**European Economic Review**
Volume 151, January 2023, 104323

ELSEVIER

## Does working from home work? A natural

**Decision Analytics Journal**
Volume 3, June 2022, 100037

ELSEVIER

## A decision framework for software startups to succeed in COVID-19 environment

**Self and Identity** ›
Volume 21, 2022 - Issue 8

Open access

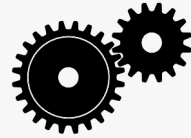| 1,516 | 2 | 3 |
| Views | CrossRef citations to date | Altmetric |

Listen ▶

Research Article

## Identity integration matters: The case of parents working fr during the COVID-19 health emergency

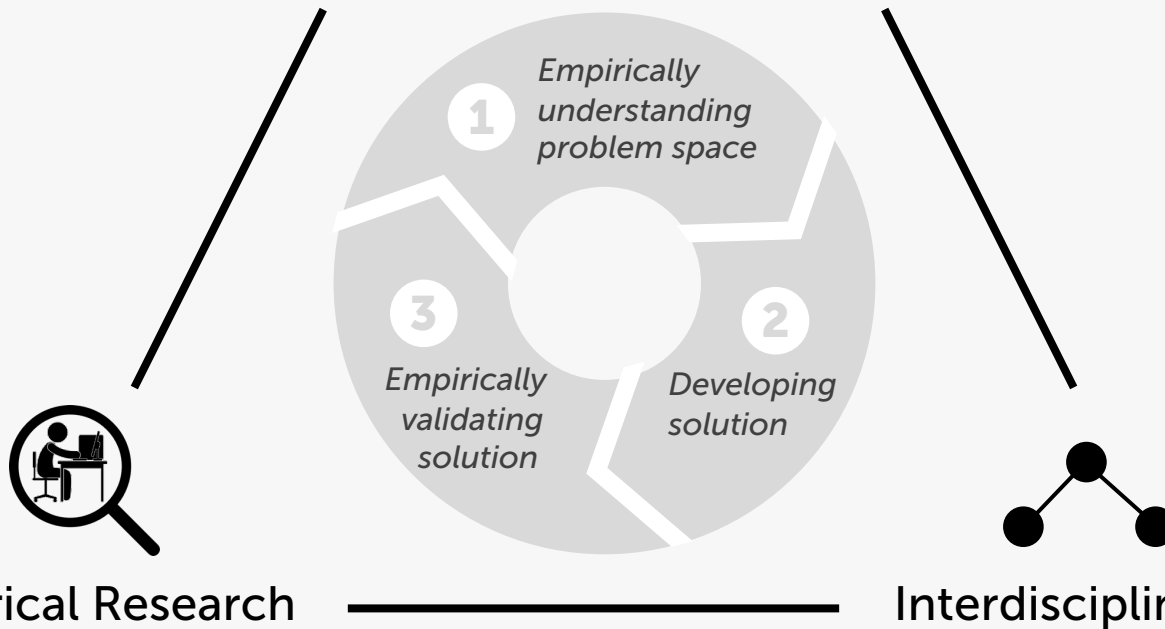Claudia Manzi, Yasin Koc ☑, Verónica Benet-Martínez & Eleonora Reverberi

Pages 914-938 | Received 23 Dec 2020, Accepted 04 Nov 2021, Published online: 02 Dec 2021

Download citation | https://doi.org/10.1080/15298868.2021.2004217

Check for updates

# Remainder of this talk

Software Engineering
**Stakeholders**, Processes, Tools

**Empirical Research**

1. *Empirically understanding problem space*
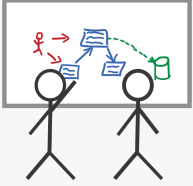2. *Developing solution*
3. *Empirically validating solution*

**Interdisciplinarity**

**Sketching**

*FSE '14, ESEM '15, VISSOFT '17*

**Code Plagiarism**

stack**overflow**

*EMSE '18, MSR '18, MSR '19, ICSE '20*

**Pandemic Programming**

*EMSE '20*

**SAP** HANA

**Timeout Flakiness**

*ICSE '24*

# Sketching

# Research Design

**Questions:**

**How** and **why** do software practitioners use sketches and diagrams?
How are they related to **source code**?
How can we provide better **tool support**?

**Methods:**

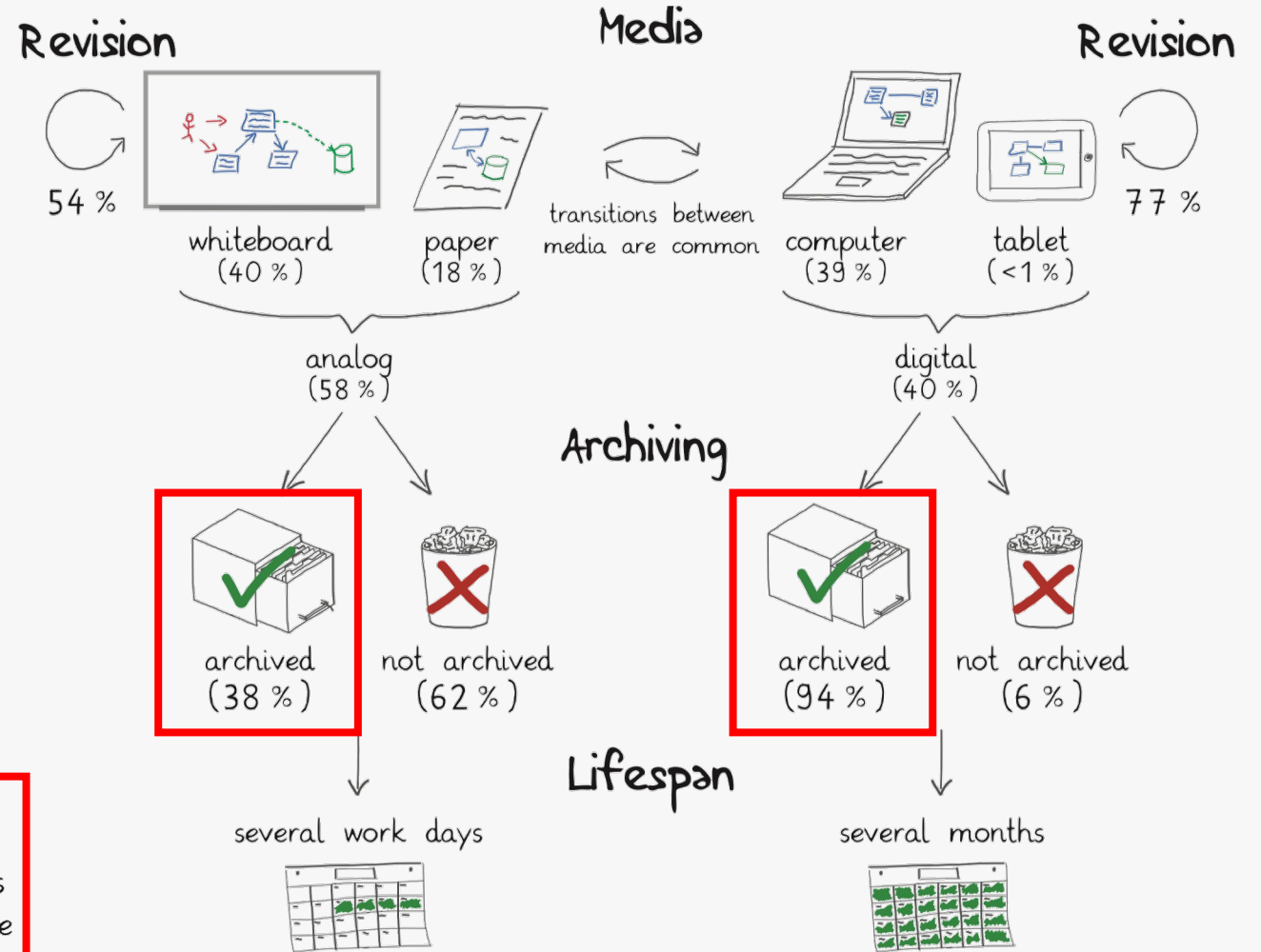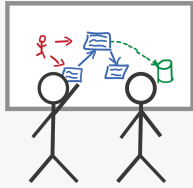Field study, interviews, lab study, online survey, formative tool evaluations.

# Results



**Sketching**

**Online survey with 394 participants** from 32 countries, asking for last sketch/diagram created.

Revision — 54 % — whiteboard (40 %) — paper (18 %)

Media — transitions between media are common

Revision — 77 % — computer (39 %) — tablet (<1 %)

analog (58 %) — digital (40 %)

Archiving

archived (38 %) — not archived (62 %)

archived (94 %) — not archived (6 %)

Lifespan

several work days — several months

**Relation to Source Code**

47 % of the sketches are rated as helpful for others to understand the related source code artifacts.

Sketching

SketchLink

https://www.youtube.com/watch?v=mG6xCiQpS80

# A startup recently gaining traction went a step further

# Code Plagiarism

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

Cutting corners to meet arbitrary management deadlines

Essential

# Copying and Pasting from Stack Overflow

O'REILLY®

The Practical Developer
@ThePracticalDev

## The full stackoverflow developer

*Friday, July 17th, 2015 at 1:04 pm*

In a few talks and interviews I lamented about a phenomenon in our market that's always been around, but seems to be rampant by now: the one of **the full stackoverflow developer**. Prompted by Stephen Hay on Twitter, I shall now talk a bit about what this means.

**Full Stack Overflow developers work almost entirely by copying and pasting code from Stack Overflow instead of understanding what they are doing. Instead of researching a topic, they go there first to ask a question hoping people will just give them the result.**

https://christianheilmann.com/2015/07/17/the-full-stackoverflow-developer/

https://twitter.com/ThePracticalDev/status/705825638851149824

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

# Research Design

**Question:**

How **frequently** is code from Stack Overflow posts used in public GitHub projects **without** the required **attribution**?

**Method:**

**Triangulation** of an estimate for the attribution ratio using three different **data mining** approaches.

# **Background**

*"Well, but these snippets are rather trivial and not protected by copyright."*

- **Not all** snippets on Stack Overflow copyrightable, but some experts argue that the **threshold is low.**
[Engelfriet 2016]

- No *"international standard for originality".*
[Creative Commons 2017b]

- CC BY-SA is a **viral copyleft license**, affecting all modifications and derived works.

http://theconversation.com/why-universities-cant-be-expected-to-police-copyright-infringement-82677

# Triangulated Attribution Ratio

**stackoverflow**

*Question:* **How frequently** is code from Stack Overflow posts **used** in public GitHub projects **without the required attribution**?

1. **Exploratory study**

2. **Code clone detector study**

3. **Exact matches study**

$$\overline{r}_{\mathrm{attr}}$$

We used **popularity** and **length** of the snippets as a **proxy for originality** and checked **external availability**.

# Attribution

## Code Plagiarism

*Attribution ratio:*

- Method 1 (regular expressions): $\bar{r}_{\mathrm{attr}} = 23\%$
- Method 2 (code clone detector): $\bar{r}_{\mathrm{attr}} = 24\%$
- Method 3 (exact matches): $\bar{r}_{\mathrm{attr}} = 8\%$

*Conservative estimate:* $\boxed{\bar{r}_{\mathrm{attr}} \leq 25\%}$

**Code Plagiarism**

stack**overflow**

Copy Paste

# Share-alike

**SA**

Only **2%** of all analyzed repositories (methods 1-3) containing code from Stack Overflow **attributed** its source and used a **compatible license.**

# Reaching out to Developers

- **Contacted owners** of GitHub repositories containing copies of Stack Overflow snippets.

- **75% not aware** of CC BY-SA licensing.

- Many thankful responses.

# Stack Overflow Code in the OpenJDK

JDK / JDK-8170860

**Get rid of the humanReadableByteCount() method in openjdk/hotspot**

## Details

| | | | |
|---|---|---|---|
| Type: | Bug | Status: | RESOLVED |
| Priority: | P2 | Resolution: | Fixed |
| Affects Version/s: | 9 | Fix Version/s: | 9 |
| Component/s: | hotspot | | |

implement the method humanReadableByteCount which body was copied from the Stack Overflow site: https://stackoverflow.com/a/3758880

It's just a few lines of code, but it could cause legal issues. The method should be either re-implemented or removed.

Besides the potential legal issues, duplicating a code is not a good practice.

https://bugs.openjdk.java.net/browse/JDK-8170860

# Reaching out to Developers

Code Plagiarism

stackoverflow

Microsoft / ApplicationInsights-Home · Watch 133 · Star 172 · Fork 155
Microsoft / rDSN · Watch 165 · Star 845 · Fork 203
Microsoft / Windows-universal-samples · Watch 1,064 · Star 6,540 · Fork 6,589

<> Code · ⊙ Issues 42 · Pull requests 55 · Projects 0 · Wiki · Insights

**Unclear licensing situation for code in BindableFlyout.cs**
#1070

⊙ Open · sbaltes opened this issue a day ago · 1 comment

## Why Code Snippets From Stack Overflow Can Break Your Project

You'll be surprised how many of the most common solutions contain security vulnerabilities

Mahdhi Rezvi  Follow
Jun 5, 2020 · 6 min read ★

**Stack Overflow Code Snippets**
⏱ 12 minute read

stackoverflow  Code Snippets
in GitHub  Projects

Dr. Sebastian Baltes
Software Engineering
Empirical Research
Interdisciplinarity

## The most copied StackOverflow Java code snippet contains a bug

Nine years later, developer corrects code snippet.

ZDNet

By Catalin Cimpanu for Zero Day | December 5, 2019 -- 00:09 GMT
(00:09 GMT) | Topic: Developer

Posted by u/fhoffa 3 years ago
33  Finding Stack Overflow Code Snippets in GitHub Projects
sbaltes.github.io/blog/s... ☐

**Y Hacker News** new | past | comments | ask | show | jobs | submit

▲ The most copied StackOverflow snippet of all time is flawed (programming.guide)
216 points by chris_wot on Dec 4, 2019 | hide | past | favorite | 88 comments

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

# Summary

**Code Plagiarism**

stack**overflow**

**Quantification** of code plagiarism in open-source projects, **feasibility of detection**, outreach.

Software Engineering
Stakeholders, Processes, Tools

**Triangulation** using three data mining approaches, online survey, (qualit. analysis).

1 Empirically understanding problem space

3 Empirically validating solution

2 Developing solution

Research on worldwide copyright and licensing **legislation**, exemplary court cases.

Empirical Research

Interdisciplinarity

# Industry Relevance

- Initial idea: **building a tool** integrated into CI/CD.

- Nowadays: **Snippet scanning** is supported by commercial tools, but usually deactivated due to false positives.

- Especially since the US Executive Order 14028 in 2021, companies have invested in automatically creating **Software Bill of Materials** (SBOMs).

- However, SBOMs still only cover **reuse on the level of libraries** and frameworks.

# One Project Became Two

# Evolution

15 — Rollback to Revision 13 - Edit approval overridden by post owner or moderator
source  link
edited Sep 2 '17 at 1:27

14 — null pointer exception fix as recommended in comments - also fixed Scanner stream not being closed
source  link
edit approved Aug 23 '17 at 17:52

13 — change "stupid scanner tricks" url to its new home on oracle; hat tip to @eng-samer-t
source  link
edited Jul 26 '17 at 19:33

12 — replaced http://stackoverflow.com/
source  link

11 — fix broken (first link)/old (java 6 has
source  link

10 — Nixed irrelevant access modifier.
source  link

9 — Completed a chopped word
source  link

8 — A little note about the stream not be
source  link

7 — Changing to public static.
source  link

6 — Simplify!
source  link

5 — deleted 7 characters in body
source  link

4 — added 4 characters in body
source  link

3 — Minor tweak to code, so it works rig
source  link

2 — Made the function more robust when handling an empty input stream
source  link
edited Feb 19 '12 at 5:49

1 — source  link
answered Mar 26 '11 at 20:40

## Revision 2 detail

2 — Made the function more robust when handling an empty input stream
source  link
edited Feb 19 '12 at 5:49
Pavel Repin

inline    side-by-side    side-by-side markdown

Here's a way using only standard Java library.

**Text block (local id 1)**

```java
import java.util.Scanner;
import java.util.NoSuchElementException;

public String convertStreamToString(InputStream is) {
    try {
        return new Scanner(is).useDelimiter("\\A").next();
    } catch (NoSuchElementException e) {
        return "";
    }
}
```

**Code block (local id 2)**

I learned this ~~one-liner~~trick from "Stupid Scanner tricks" article. The reason it works is because Scanner iterates over tokens in the stream, and in this case we separate tokens using "beginning of the input boundary" (\A) thus giving us only one token for the entire contents of the stream.

Note, if you need to be specific about the input stream's encoding, you can provide the second argument to `Scanner` ctor that indicates what charset to use (e.g. "UTF-8").

Hat tip goes also to Jacob, who once pointed me to the said article.

**EDITED:** Thanks to a suggestion from Patrick, made the function more robust when handling an empty input stream.

**Text block (local id 3)**

https://stackoverflow.com/posts/5445161/revisions

# Post Block Matching: Example



I can authenticate fine with other clients including SleekXMPP and Strophe. Using Prosody 0.8.2 on Ubuntu 12.04 and latest master HEAD of jaxl (2518a44b9dfeb9ec947922f078cf4f8663497712). from client:

```
<body xmlns="http://jabber.org/protocol/httpbind"
content="text/xml; charset=utf-8" to="localhost"
route="xmpp:localhost:5222" secure="true" xml:lang="en"
xmpp:version="1.0" xmlns:xmpp="urn:xmpp:xbosh" hold="1"
wait="30" rid="3937" ver="1.10" from="yang@localhost">
```

from server:

```
<body authid='72604504-a5be-4ab6-aba0-9686cca478f3' xmpp:version='1.0'
xmlns:stream='http://etherx.jabber.org/streams'
xmlns:xmpp='urn:xmpp:xbosh' inactivity='60' wait='30' polling='5'
secure='true' hold='1' from='localhost' ver='1.6'
sid='72604504-a5be-4ab6-aba0-9686cca478f3' requests='2'
xmlns='http://jabber.org/protocol/httpbind'>
```

from client:

```
<body sid="72604504-a5be-4ab6-aba0-9686cca478f3" rid="3938"
xmlns="http://jabber.org/protocol/httpbind">
```

from server:

```
<body xmlns='http://jabber.org/protocol/httpbind'
sid='72604504-a5be-4ab6-aba0-9686cca478f3' xmlns:stream =
'http://etherx.jabber.org/streams'>
```

from client:

```
<body xmlns="http://jabber.org/protocol/httpbind"
sid="72604504-a5be-4ab6-aba0-9686cca478f3" rid="3939">
```

I can authenticate fine with other clients including SleekXMPP and Strophe. Using Prosody 0.8.2 on Ubuntu 12.04 and latest master HEAD of jaxl (2518a44b9dfeb9ec947922f078cf4f8663497712). from client:
The code:

```
require 'JAXL/jaxl.php';
$cli = new JAXL(array(
    'jid' => 'yang@localhost',
    'pass' => 'asdf',
    'bosh_url' => 'http://localhost/chat/candy/example/http-bind/'
));
$cli->add_cb('on_auth_success', function() {
    print 'yay';
});
$cli->start();
```

from client:

```
<body xmlns="http://jabber.org/protocol/httpbind"
content="text/xml; charset=utf-8" to="localhost"
route="xmpp:localhost:5222" secure="true" xml:lang="en"
xmpp:version="1.0" xmlns:xmpp="urn:xmpp:xbosh" hold="1"
wait="30" rid="3937" ver="1.10" from="yang@localhost">
```

from server:

```
<body authid='72604504-a5be-4ab6-aba0-9686cca478f3' xmpp:version='1.0'
xmlns:stream='http://etherx.jabber.org/streams'
xmlns:xmpp='urn:xmpp:xbosh' inactivity='60' wait='30' polling='5'
secure='true' hold='1' from='localhost' ver='1.6'
sid='72604504-a5be-4ab6-aba0-9686cca478f3' requests='2'
xmlns='http://jabber.org/protocol/httpbind'>
```

from client:

●●●

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

# Post Block Matching: Example

# Post Block Matching: Predecessor Selection Strategy

**Algorithm 2** Revised Matching Strategy

```
for all p_{2≤i≤n} do
    // set predecessors where only one candidate exists
    for all b^τ_{(i,1≤j≤|p_i|)} do
        if |Pred(b^τ_{(i,j)})| = 1 then
            Let pred be the equal or similar predecessor
            if available(pred) then // new
                if |Succ(pred)| = 1 then
                    Set pred as predecessor of b^τ_{(i,j)}
                    continue
                end if
            else
                setPredPositionRunnerUp(p_i)
            end if
        end if
    end for
    // set predecessors using context
    predSet = true
    while predSet do
        predSet = setPredContext(p_i, BOTH)
    end while
    while predSet do
        predSet = setPredContext(p_i, BELOW)
    end while
    while predSet do
        predSet = setPredContext(p_i, ABOVE)
    end while
    // set predecessors using position
    setPredPosition(p_i)
    // set runner-up predecessors for the remaining post blocks
    setPredPositionRunnerUp(p_i)
end for
```

- If post block has **only one possible successor** and this possible successor has **only one possible predecessor**, no strategy required

- Otherwise, consider the **context**

- Then, try to set remaining post blocks using **position** (min. local id distance)

- Sometimes exact matches are not the correct predecessor

# Post Block Matching: Performance

**Text:**
*manhattanFourGramNormalized*
Threshold:        0.17
TPR:               0.99
FPR:               0.14
**Matthews Corr.: 0.86**

**Code:**
winnowingFourGramDiceNormalized
Threshold:        0.23
TPR:               0.99
FPR:               0.07
**Matthews Corr.: 0.92**

https://github.com/sotorrent/posthistory-extractor

Prof. Dr. Sebastian Baltes - Evidence over Opinion: An Empirical Approach to Software Engineering

# SOTorrent

- Among other features, the dataset provides the **version history** of Stack Overflow content on the **level of individual text or code blocks**

- Was official **mining challenge** of MSR 2019

## sotorrent.org

*Dataset available on Zenodo*


Open Data

# Timeout Flakiness

SAP
HANA

# Common Flakiness Definition

"A test can be considered flaky when it exhibits both passing and failing results for the same code."

# Context

- Flaky tests **interfere with CI/CD.**

- SAP HANA is a large industrial DBMS with **long-running system tests.**

- Flaky failures **impede automatic test run assessment** and merging.

- Standard strategy: **restart flaky tests.**

- Configurable **"max duration"** (timeout values) to prevent stuck tests from blocking resources, but this can cause flaky test executions.

- We found that **99% of CI runs are affected** by flaky failures (→ costs, delay).

# Research Questions

**RQ1:** What **level of test flakiness** do we observe in SAP HANA's system tests and what can we identify as a **major contributing factor**?

**RQ2:** What **impact** does **increasing timeout values** have on test flakiness in context of SAP HANA?

**RQ3:** How do **developers commonly adjust timeout values** in the context of SAP HANA?

**RQ4:** To what degree can we **optimize the timeout values** with respect to their average test execution costs?
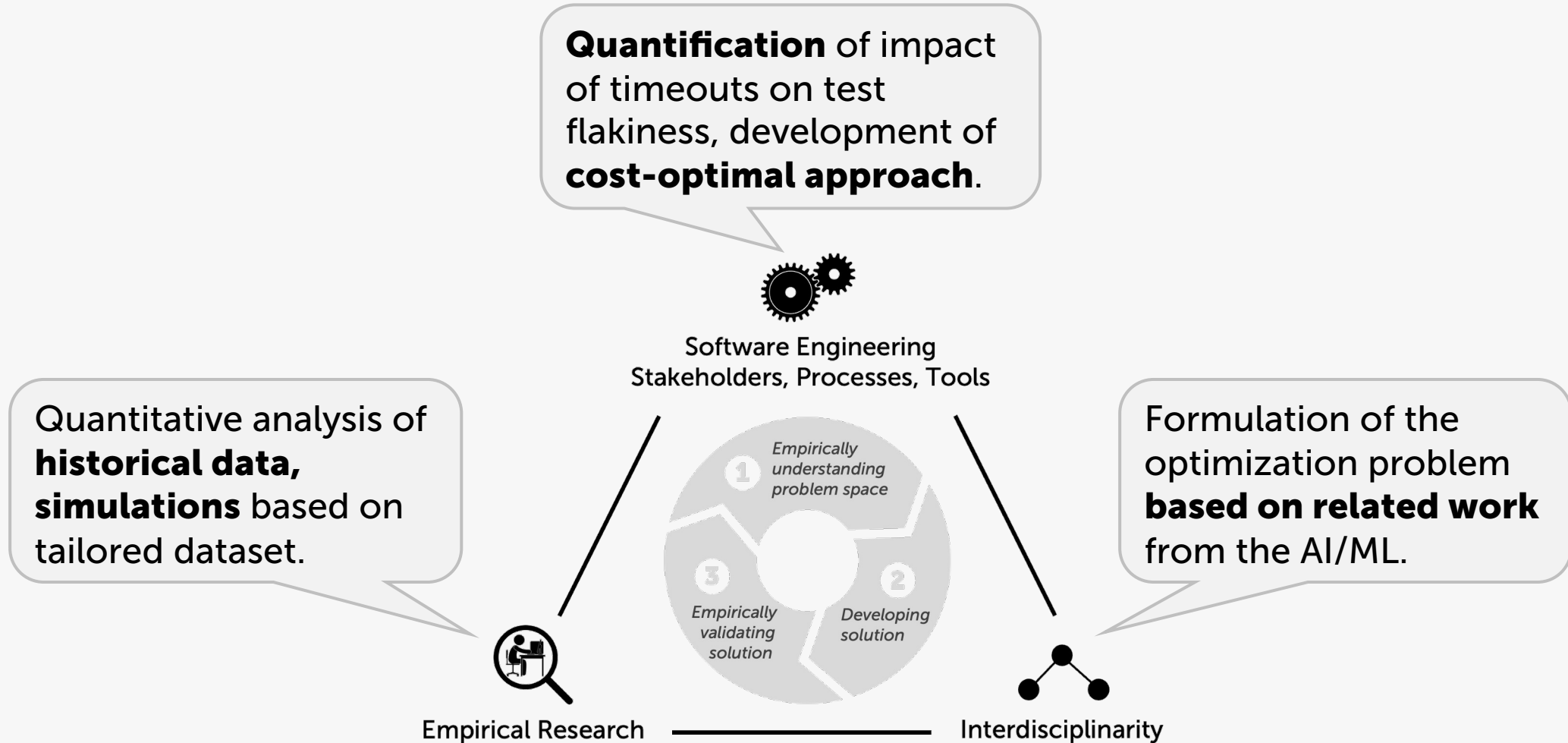
# Conclusion

- Flakiness definition is of **little practical use**, because the test flakiness rate **converges to 1**.

- Timeout values can cause **additional costs.**

- Cost-optimal timeout values **can increase efficiency.**

- Baseline approach with a **fixed global timeout** currently being implemented at SAP.
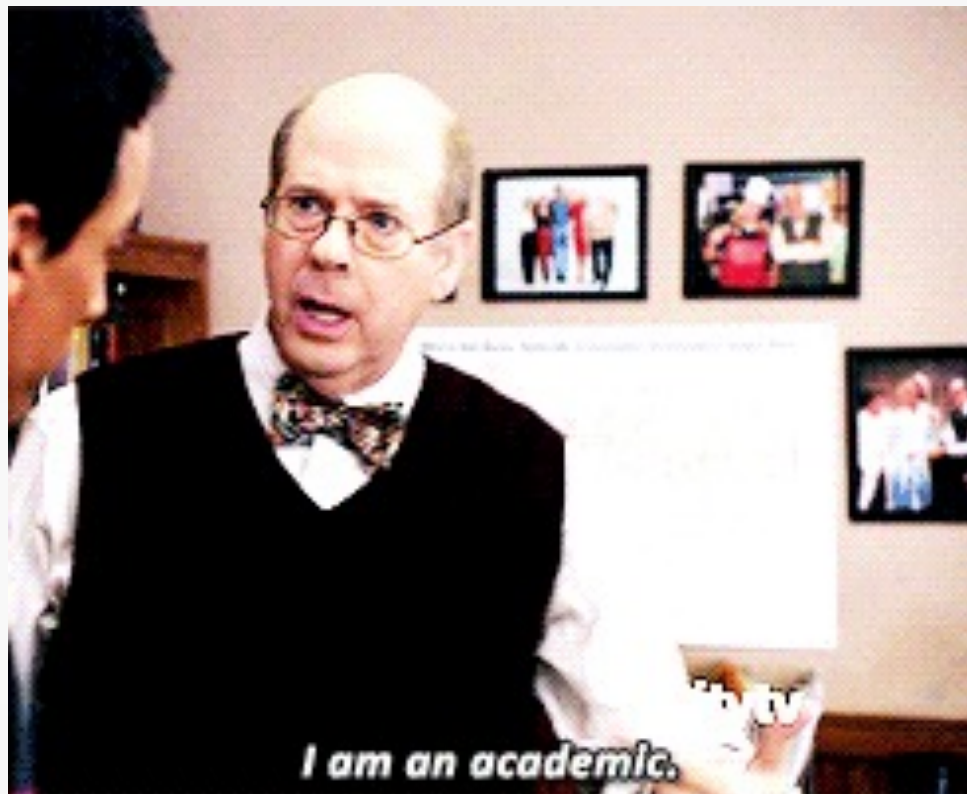
# Back to our previous topic...

Research     VS.     Practice

# The issues of applied SE research in industry

- **Constantly changing (business) priorities**.

- Hard to internally "sell" **projects with potential long-term benefits** requiring short-term investment/exploration.

- Importance of (early) quantification of business impact via KPIs (key performance indicators) makes **foundational research difficult**.

# ...are an opportunity for academia

When doing research in **close collaboration with industry**, academic researchers can work on the topics that:

- Are **promising**, but do not (yet) have an immediate business value.

- Require considerable investment in **exploring the problem space** without an immediate quantifiable business impact.

- Require **deep understanding** of the broad spectrum of empirical software engineering methods.

- Require a **systematic screening of related (academic) work** on a topic to replicate or extend it.

# Software Engineering @ UBT

- Lecture in Summer Semester 2024:
  - *"Software Engineering"* (Bachelor)
- New courses planned for upcoming semesters:
  - *"Advanced Software Engineering"* (Bachelor/Master)
  - *"Software Analytics"* (Master)

- **Interested in software engineering research?**
  - Bachelor/Master theses
  - Openings for PhD positions
  - Please reach out to me in case you're interested!