# Software Developers' Work Habits and Expertise

**Sebastian Baltes**
@s_baltes

empirical-software.engineering

THE UNIVERSITY
*of* ADELAIDE

# Interaction

# My Background

# Evidence-based Practice through Practice-based Evidence

# Studying Developers' Work Habits

Observe
Describe
Explain
} Software Developers' Work Habits $\Longrightarrow$ Expand knowledge

→ Derive requirements for better tool support

→ Identify possible process improvements

→ Communicate findings back to practitioners

# Habits?

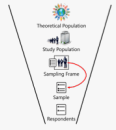A habit is a „**settled tendency or usual manner of behavior**"
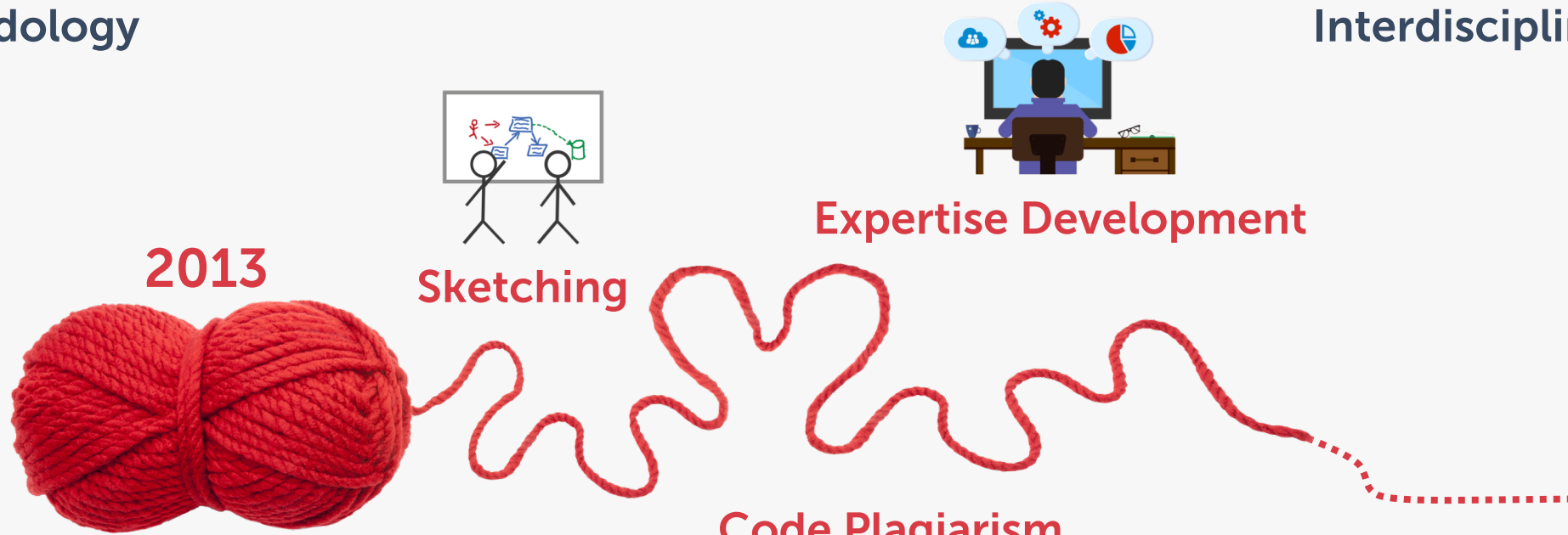
*Personal habits*

*Work habits*

# Studied Habits

Issues in Sampling
Software Developers
**Methodology**

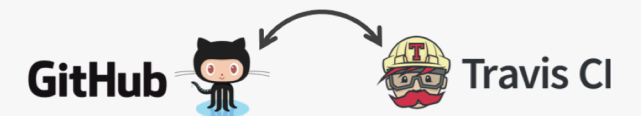Constructing Urban
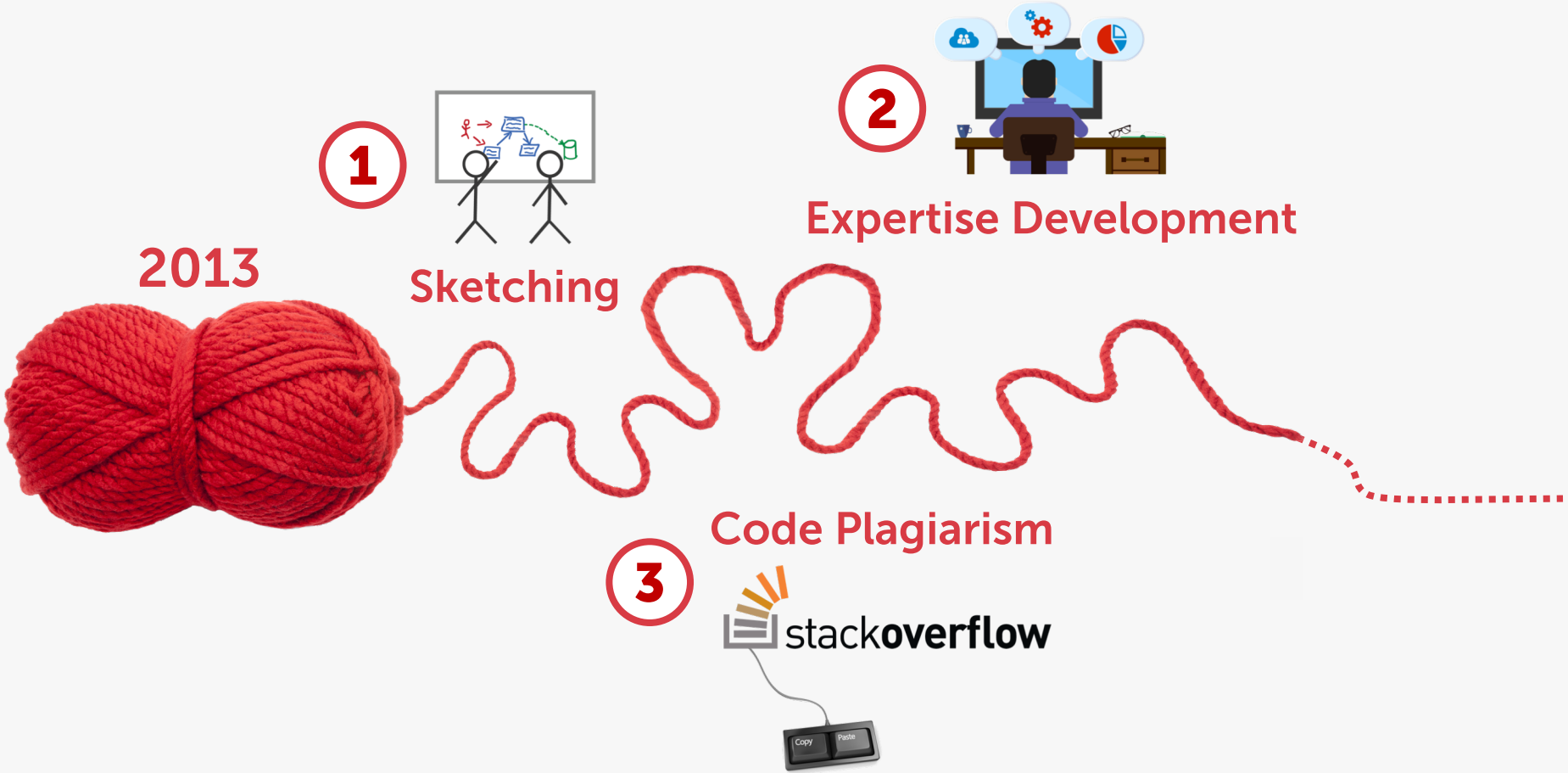Tourism Space Digitally
**Interdisciplinary Research**

**Expertise Development**

**2013**

**Sketching**

**Code Plagiarism**

**stackoverflow**

**Regular Expressions**

**RegViz**

**Continuous Integration**

GitHub      Travis CI

# Overview of this Talk

**2013**

**1** Sketching

**2** Expertise Development

**3** Code Plagiarism

stack**overflow**

# Overview of this Talk



2013

**1** Sketching

**2** Expertise Development

**3** Code Plagiarism

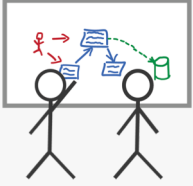stack**overflow**

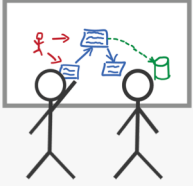# Sketching

# Research Questions

**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**?

**Approach:**

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

# Sketches and Diagrams in Practice

FSE 2014

Sebastian Baltes
Computer Science
University of Trier
Trier, Germany
s.baltes@uni-trier.de

Stephan Diehl
Computer Science
University of Trier
Trier, Germany
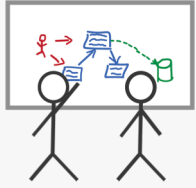diehl@uni-trier.de

## ABSTRACT

Sketches and diagrams play an important role in the daily work of software developers. In this paper, we investigate the use of sketches and diagrams in software engineering practice. To this end, we used both quantitative and qualitative methods. We present the results of an exploratory study in three companies and an online survey with 394 participants. Our participants included software developers, software architects, project managers, consultants, as well as researchers. They worked in different countries and on projects from a wide range of application areas. Most questions in the survey were related to the last sketch or diagram that the participants had created. Contrary to our expectations and previous work, the majority of sketches and
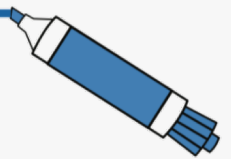
## 1. INTRODUCTION

Over the past years, studies have shown the importance of sketches and diagrams in software development [6,11,43]. Most of these visual artifacts do not follow formal conventions like the *Unified Modeling Language* (UML), but have an informal, ad-hoc nature [6,11,23,25]. Sketches and diagrams are important because they depict parts of the mental model developers build to understand a software project [21]. They may contain different views, levels of abstraction, formal and informal notations, pictures, or generated parts [6, 11,41,42]. Developers create sketches and diagrams mainly to understand, to design, and to communicate [6]. Media for sketch creation include whiteboards, engineering notebooks, scrap papers, but also software tools like Photoshop
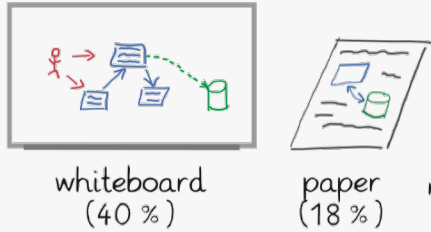
https://empirical-software.engineering/projects/sketches/

Sketching

SketchLink

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

# Overview of this Talk



**2013**

① Sketching

② **Expertise Development**

③ Code Plagiarism

stack**overflow**

# Expertise Development

# Towards a Theory of Software Development Expertise

Sebastian Baltes
University of Trier
Trier, Germany
research@sbaltes.com

ESEC/FSE 2018

Stephan Diehl
University of Trier
Trier, Germany
diehl@uni-trier.de

## ABSTRACT

Software development includes diverse tasks such as implementing new features, analyzing requirements, and fixing bugs. Being an expert in those tasks requires a certain set of skills, knowledge, and experience. Several studies investigated individual aspects of software development expertise, but what is missing is a comprehensive theory. We present a first conceptual theory of software development expertise that is grounded in data from a mixed-methods survey with 335 software developers and in literature on expertise and expert performance. Our theory currently focuses on programming, but already provides valuable insights for researchers, developers, and employers. The theory describes important properties of software development expertise and which factors foster or hinder its formation, including how developers' performance may decline over time. Moreover, our quantitative results show that developers' expertise self-assessments are context-dependent and that experience is not necessarily related to expertise.

expert performance [78]. Bergersen et al. proposed an instrument to measure programming skill [9], but their approach may suffer from learning effects because it is based on a fixed set of programming tasks. Furthermore, aside from programming, software development involves many other tasks such as requirements engineering, testing, and debugging [62, 96, 100], in which a software development expert is expected to be good at.

In the past, researchers investigated certain aspects of software development expertise (SDExp) such as the influence of programming experience [95], desired attributes of software engineers [63], or the time it takes for developers to become "fluent" in software projects [117]. However, there is currently no theory combining those individual aspects. Such a theory could help structuring existing knowledge about SDExp in a concise and precise way and hence facilitate its communication [44]. Despite many arguments in favor of developing and using theories [46, 56, 85, 109], theory-driven research is not very common in software engineering [97].

https://empirical-software.engineering/projects/expertise/

# Software Development Expertise?

Implementing new features
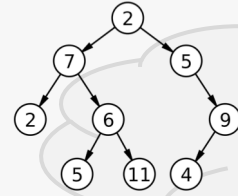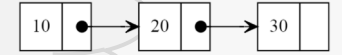
Algorithms & Data structures

Testing

Communication

Debugging

# Software Development Expertise?

Implementing new features

Algorithms & Data structures

JUnit Testing jbehave

Communication

Debugging

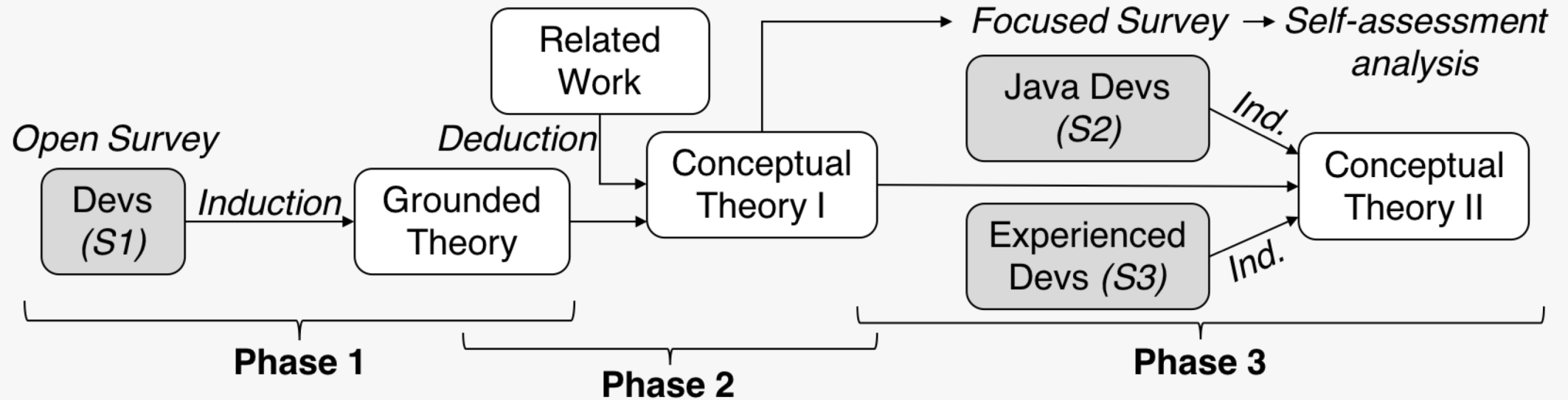**Expertise Development**

## Questions:

How to **structure** all those expertise-related aspects? Which factors influence **expertise development** over time?

## Approach:
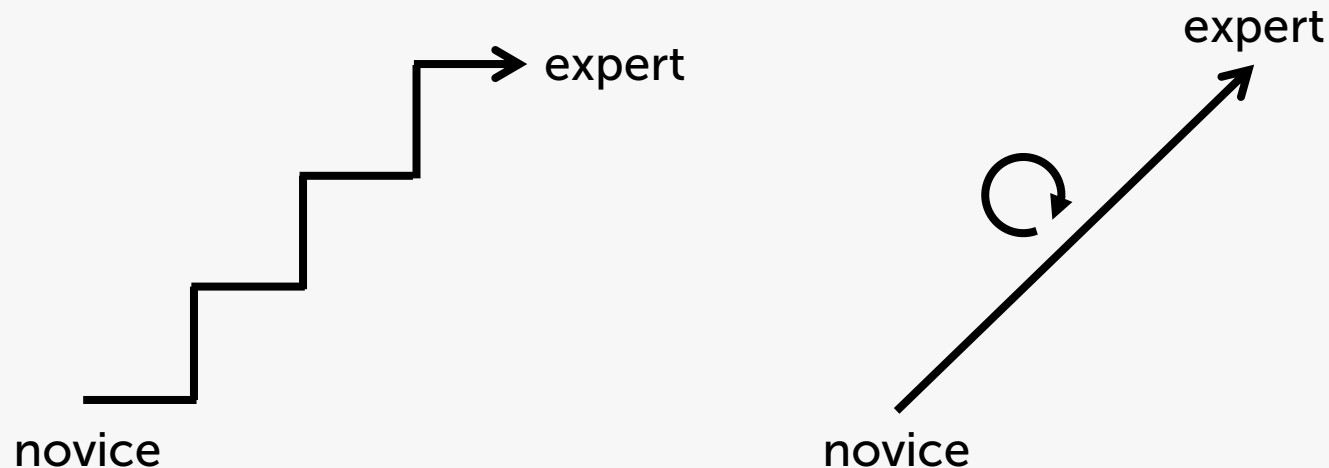
Iterative theory building

# Research Design



- **Induction:** 335 online survey participants in total
- **Deduction:** Main source *"Cambridge Handbook of Expertise and Expert Performance"*

# Our Expertise Model

- **Task-specific** (e.g., writing code, debugging, testing)
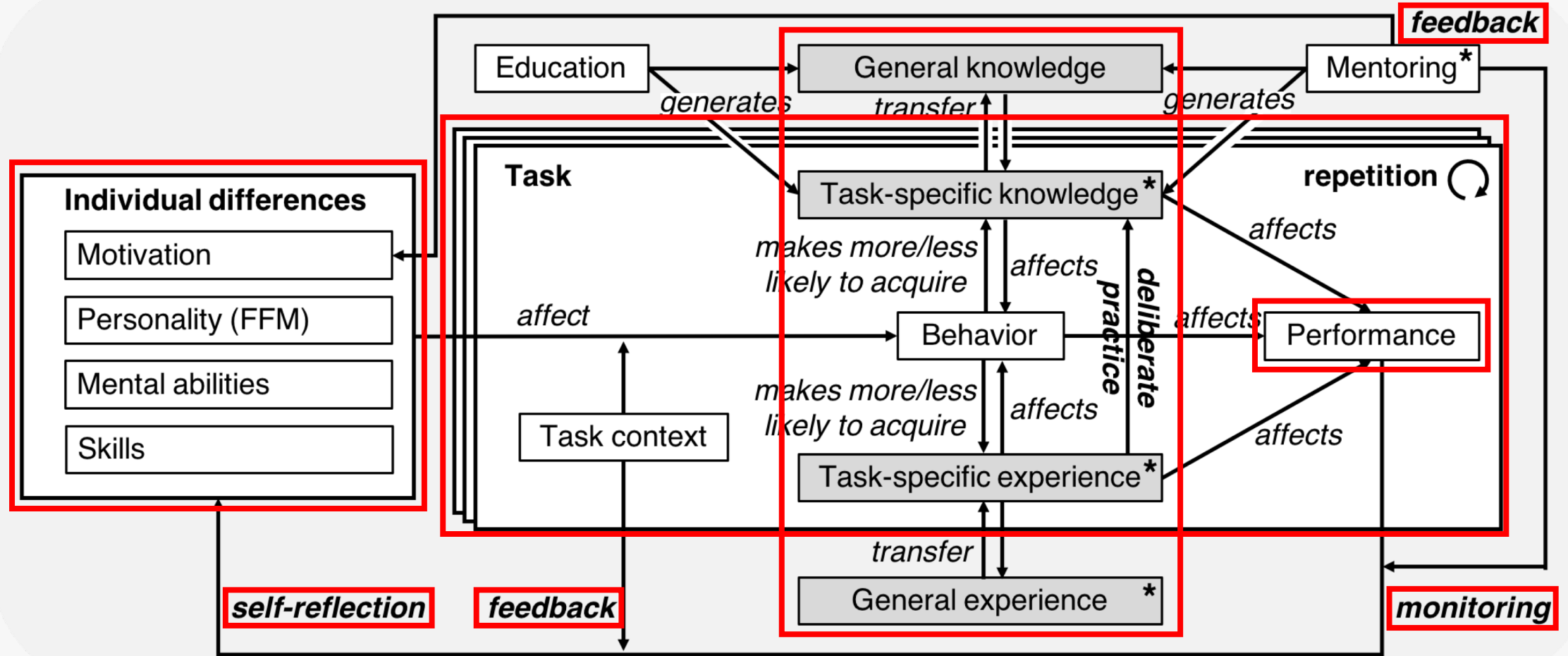- Focuses on **individual developers**
- **Process** view (repetition of tasks)
- Notion of **transferable knowledge and experience** from related fields or tasks
- **Continuum** instead of discrete expertise steps

# Conceptual Theory

# Conceptual Theory

# Summary

**Researchers** can...

- Use our theory to **design studies** on expertise development
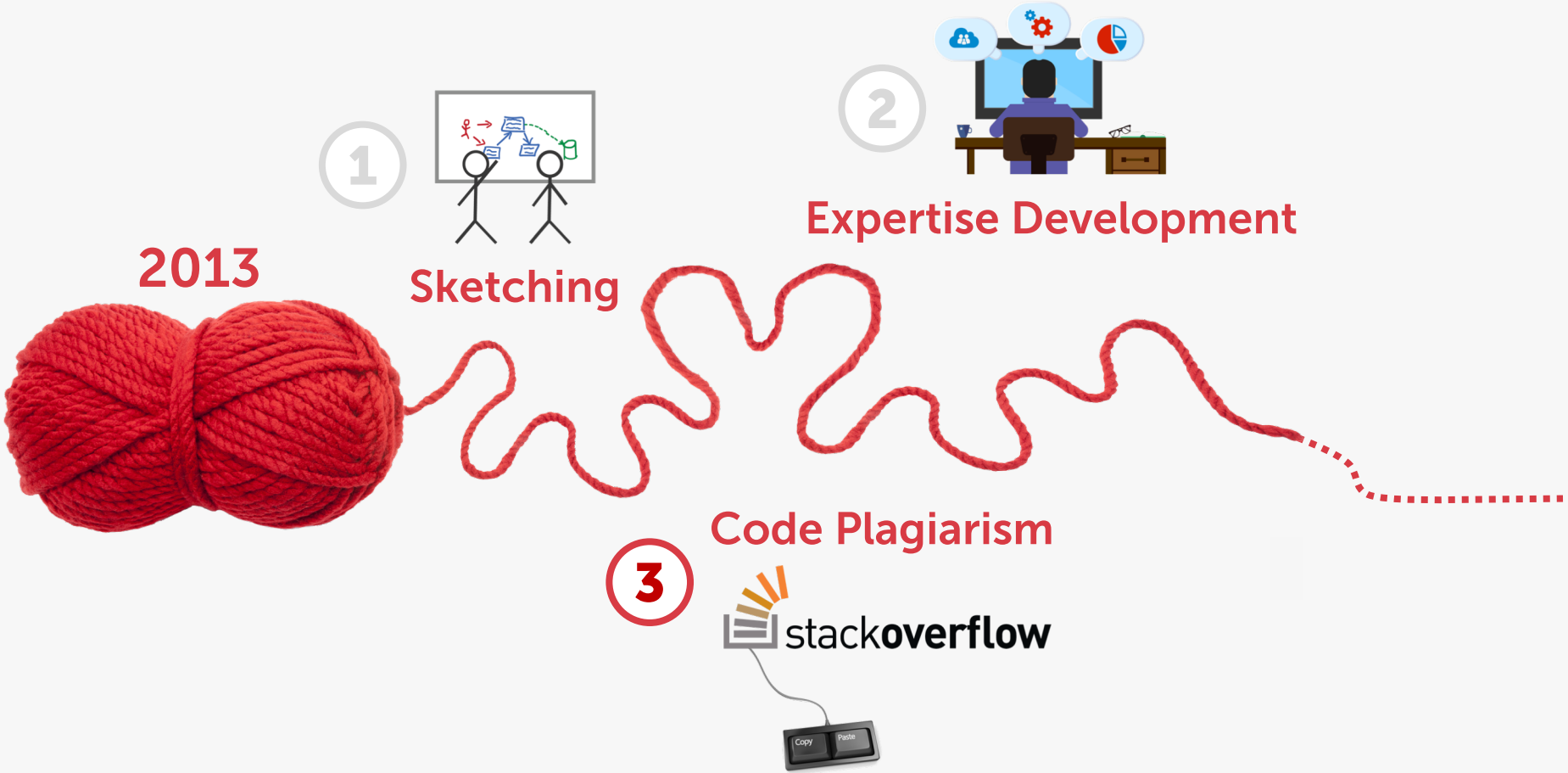- Adopt our **theory building** approach

**Developers** can...

- Learn what other developers expect from **experts/mentors**
- Learn which **behaviors** may lead to becoming an expert

**Employers** can...

- Learn what **(de)motivates** employees and thus fosters or hinders expertise development
- Reflect on ideas to build a work environment **supporting self-improvement** of their staff

# Overview of this Talk



2013

① Sketching

② Expertise Development

③ Code Plagiarism

**stackoverflow**

# Code Plagiarism

# Code Plagiarism

**stackoverflow**

CrossMark

# Usage and attribution of Stack Overflow code snippets in GitHub projects

Sebastian Baltes[1] · Stephan Diehl[1]

## Abstract

Stack Overflow (SO) is the most popular question-and-answer website for software developers, providing a large amount of copyable code snippets. Using those snippets raises maintenance and legal issues. SO's license (CC BY-SA 3.0) requires attribution, i.e., referencing the original question or answer, and requires derived work to adopt a compatible license. While there is a heated debate on SO's license model for code snippets and the

https://empirical-software.engineering/projects/snippets/

# GitHub

- **Hosted version control** platform for (software) projects
- Features include access control, **collaboration features** such as **issue tracking**, wikis, gamification of development activity
- **Public** projects and **private** projects with up to three collaborators are **free**
- As of May 2019: **>37m users** and **>100m projects**

GitHub

# Stack Overflow

- **Question and answer** website for software developers
- Covers a **wide variety** of **programming-related topics**
- Posts can be commented, edited, and up-/down-voted
- **Gamification** through reputation points awarded for different kinds of contributions
- **Jobs** section for advertising employment opportunities
- As of June 2019 **>10.5m** registered users and **>17.7m questions**

# How do I read / convert an InputStream into a String in Java?

**Asked** 10 years, 9 months ago    **Active** 2 days ago    **Viewed** 2.0m times

Ask Question

▲

3775

▼

★

1107

If you have a `java.io.InputStream` object, how should you process that object and produce a `String`?

Suppose I have an `InputStream` that contains text data, and I want to convert it to a `String`, so for example I can write that to a log file.

What is the easiest way to take the `InputStream` and convert it to a `String`?

```
public String convertStreamToString(InputStream is) {
    // ???
}
```

java    string    io    stream    inputstream

share    improve this question

edited Jan 5 at 10:28
**Peter Mortensen**
14.4k ● 19 ● 88 ● 117

asked Nov 21 '08 at 16:47
**Johnny Maelstrom**
19.3k ● 5 ● 17 ● 17

781    Boy, I'm absolutely in love with Java, but this question comes up so often you'd think they'd just figure out that the chaining of streams is somewhat difficult and either make helpers to create various combinations or rethink the whole thing. – **Bill K** Nov 21 '08 at 17:16

29    The answers to this question only work if you want to read the stream's contents *fully* (until it is closed). Since that is not always intended (http requests with a keep-alive connection won't be closed), these method calls block (not giving you the contents). – **f1sh** Jul 14 '10 at 13:32

# Example



## Question

https://stackoverflow.com/q/309424

## Answer

https://stackoverflow.com/a/5445161

Here's a way using only standard Java library (note that the stream is not closed, YMMV).

**2034**

```
static String convertStreamToString(java.io.InputStream is) {
    java.util.Scanner s = new java.util.Scanner(is).useDelimiter("\\A");
    return s.hasNext() ? s.next() : "";
}
```

**Code snippet**

I learned this trick from "Stupid Scanner tricks" article. The reason it works is because Scanner iterates over tokens in ___ase we separate tokens using ___ boundary" (\A) thus giv___ the entire contents of the st___

**Source of snippet**

**Reference to JDK**

**Note, if you need to be specific about the input stream's encoding, you can provide the second argument to** `Scanner` **constructor 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. **One more edit:** nixed try/catch, Patrick's way is more laconic.

share in___er          edited Sep 2___          ___ered Mar 26 '11 at 20:40

**Post edits**          **Reasons for edits**

Pavel Repin
25.3k ● 1 ● 27 ● 36

# Comments



This stuff is clearly a hack.

# Evolution

- Like other **software artifacts**, SO posts **evolve:**
  - Content of **17.3m** posts has been edited
  - **Bugs** in code snippets are fixed
  - **Clarifications** are added in text documenting the code
  - Snippets are **updated** to new language/library versions

- **Evolution of code on SO** differs from regular software projects:
  - **Short** code snippets (12 LOC on average)
  - **No bug tracking** system (just comments and new answers)
  - **No versioning** for individual snippets (just whole posts)

# SO Revisions

**Problems:**

- Version history is only available on the level of whole posts, thus **individual code snippets hard to trace**
- **Comments and edits** are not linked
- Unclear how **external sources** are related



Text block

Code block

Text block

# SOTorrent

- Among other features, the dataset provides the **version history** of Stack Overflow content on the **level of individual text or code blocks**
- Extraction of post blocks and mapping to their predecessors was required, involving an extensive evaluation of similarity metrics

| Type | Metric | | Variants |
|---|---|---|---|
| edit | levenshtein<br>longestCommonSubsequence (LCS) | damerauLevenshtein<br>optimalAlignment (OA) | with/without normalization |
| set | nGram{Jaccard\|Dice\|Overlap}<br>token{Jaccard\|Dice\|Overlap} | nShingle{Jaccard\|Dice\|Overlap} | $n$Gram : $n \in \{2, 3, 4, 5\}$, $n$Shingle : $n \in \{2, 3\}$<br>with/without normalization, padding (nGram) |
| profile | cosineNGram{Bool\|TF\|NormalizedTF}<br>cosineNShingle{Bool\|TF\|NormalizedTF}<br>cosineToken{Bool\|TF\|NormalizedTF} | manhattanNGram<br>manhattanNShingle<br>manhattanToken | $n$Gram : $n \in \{2, 3, 4, 5\}$, $n$Shingle : $n \in \{2, 3\}$<br>with normalization (both) and without (cosine) |
| fingerprint | winnowingNGram{Jaccard\|Dice\|Overlap\|LCS\|OA} | | $n$Gram : $n \in \{2, 3, 4, 5\}$,<br>with/without normalization |
| equal | equal | tokenEqual | with/without normalization |

https://github.com/sotorrent/string-similarity

**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// new
                setPredPositionRunnerUp(p_i) // new
            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) // new
end for
```

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

SOTorrent: Reconstructing and Analyzing the Evolution of Stack Overflow Posts

Sebastian Baltes
Lorik Dumani
research@sbaltes.com
dumani@uni-trier.de
University of Trier, German

Christoph Treude
christoph.treude@adelaide.edu.au
University of Adelaide, Australia

Stephan Diehl
diehl@uni-trier.de
University of Trier, Germany

**ABSTRACT**

Stack Overflow (SO) is the most popular
site for software developers, providin
snippets and free-form text on a wide v
software artifacts, questions and answe
for example when bugs in code snippet
to work with a more recent library ver
code snippet is edited for clarity. To be a
on SO evolves, we built *SOTorrent*, an
official SO data dump. *SOTorrent* provid
tory of SO content at the level of whole
code blocks. It connects SO posts to oth
URLs from text blocks and by collecti

SOTorrent: Studying the Origin, Evolution, and Usage of Stack Overflow Code Snippets

Sebastian Baltes
*University of Trier, Germany*
research@sbaltes.com

Christoph Treude
*University of Adelaide, Australia*
christoph.treude@adelaide.edu.au

Stephan Diehl
*University of Trier, Germany*
diehl@uni-trier.de

*Abstract*—Stack Overflow (SO) is the most popular question-and-answer website for software developers, providing a large amount of copyable code snippets. Like other software artifacts, code on SO evolves over time, for example when bugs are fixed or APIs are updated to the most recent version. To be able to analyze how code and the surrounding text on SO evolves, we built *SOTorrent*, an open dataset based on the official SO data dump. *SOTorrent* provides access to the version history of SO content at the level of whole posts and individual text and code blocks. It connects code snippets from SO posts to other platforms by aggregating URLs from surrounding text blocks and comments, and by collecting references from GitHub files to SO posts. Our vision is that researchers will use *SOTorrent* to investigate and understand the evolution and maintenance of code on SO and its relation to other platforms such as GitHub.

dataset [16] that enables researchers to analyze the version history of SO posts at the level of individual text and code blocks (see Figure 1 for exemplary posts). The official SO data dump [1] keeps track of different versions of entire posts, but does not contain information about differences between versions at a more fine-grained level. In particular, extracting different versions of the same code snippet from the history of a post is challenging and required us to develop a complex strategy, involving the evaluation of 134 different string similarity metrics [15]. Beside providing access to the version history, our dataset links SO posts to external resources in two ways: (1) by extracting linked URLs from text blocks of SO posts and from post comments and (2) by providing

**MSR 2018/19**

# sotorrent.org
*Dataset available on Zenodo and BigQuery*
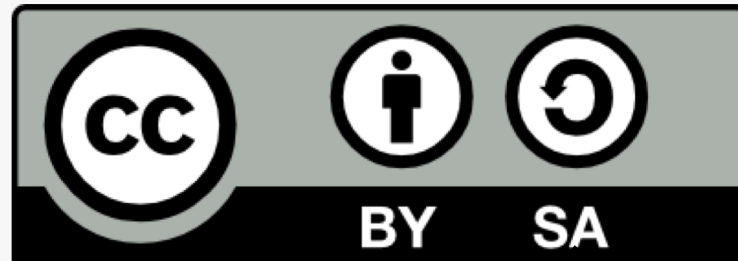
Open Data

# Question for the Audience I

Who admits regularly copying non-trivial code snippets from Stack Overflow?

# Question for the Audience II

Who knew that all content on Stack Overflow is licensed under CC BY-SA?



*"You must give **appropriate credit** [...] and indicate if changes were made."*

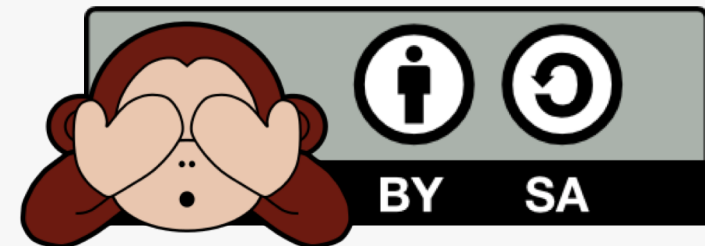*"If you [...] **build upon** the material, you must **distribute your contributions** under the same license as the original."*

**Attribution**          **Share-alike**

# Results from our Online Surveys

- **46%** of the participants admitted copying code from Stack Overflow **without attribution**

- **75%** did **not know** that content on SO is licensed under **CC BY-SA**

- **67%** did **not know** that **attribution is required**

→ **Lack of awareness**

# Background

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

- Not all code snippets on Stack Overflow are copyrightable

- "A snippet that is more than one or two lines of standard function calls would typically be creative enough for copyright" [Engelfriet 2016]

- But no "international standard for originality" [Creative Commons 2017b]

https://stackoverflow.com/a/3145655

https://github.com/pacosal/ownmdm/blob/master/src/com/pacosal/mdm/MyLocation.java

# 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

# … and in Microsoft GitHub Repos

# Implications of Stack Overflow's License

## Permissive Licenses

- Permit using the licensed source code in proprietary software **without publishing changes** or the **derived work**
- *Examples:* MIT, Apache, and BSD license families

## Copyleft Licenses

- Requires either modifications to the licensed content or the complete derived work to be **published under the same or a compatible license** (share-alike)
- *Examples (weak copyleft):* Mozilla/Eclipse Public Licenses
- *Examples (viral copyleft):* GNU General Public Licenses, Creative Commons Share-Alike Licenses (e.g., **CC BY-SA**)

# Enforceability of Copyleft Licenses

- Courts in the US and Europe ruled that open source licenses are **enforceable contracts**

- Authors are able to **sue** when terms such as the share-alike requirement are violated:
  - **Interdict distribution** of derived work
  - **Claim monetary damages**

- USA: DMCA takedown notices for allegedly infringed copyright
  - Example: https://github.com/github/dmca

- Risk in mergers and acquisitions of companies
  - Example: FSF vs. Cisco lawsuit

# Research Question

## Question:

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

## Approach:

Triangulate an estimate for the attribution ratio using three different methods.

# Method 1: Regular Expressions

Google Cloud
GitHub

**209m files in 4.1m projects**

`.java` →

Java

**13m Java files in 336k projects**

`...stackoverflow\.com...` →

stackoverflow

**10 most frequently referenced answers**

Check external availability

**Manually build regular expressions matching code snippets**
(referenced usages as test cases)

```
((?i:String[\s]*\w+\([^\{]*long[^\{]+\)[\s]*\{[\s\S]+if
[\s]*\([^<]+<[^\)]+\)[\s\S]*return[^;]+\+[^;]*\"* B\"
[\s\S]+int[\s][^\=]+\=[\s]*\([\s]+int[\s]+\*\)[\s]*\([\s]*
Math[\s]*\.[\s]*log[\s]*\([^\)]+\)[\s]*\/[\s]*Math[\s]*
\.[\s]*log[\s]*\([^\)]+\)[\s]*\)[\s]*\S]+return[^\}]+
String[\s]*\.[\s]*format[\s]*\([^\}]+\)))
```
←

**4,198 files with matches**

←

**Check if attributed**
(URL to answer or corresponding question)

**Check for false positives**

# Results

| Rank | Matches | | Recall | Attribution | |
| | ALL | DISTINCT | REF | NO-REF | REF/$F_{AQ}$ | REF/DISTINCT | $F_{AQ}$/DIST. |
|---|---|---|---|---|---|---|---|
| 1 | 997 | 448 | 97 | 351 | 79.5% | 21.7% | 27.2% |
| 2 | 1,843 | 913 | 60 | 853 | 60.0% | 6.6% | 11.0% |
| 3 | 2,662 | 902 | 87 | 815 | 80.6% | 9.6% | 12.0% |
| 4 | 420 | 170 | 18 | 152 | 94.7% | 10.6% | 11.2% |
| 5 | 1,492 | 402 | 25 | 377 | 73.5% | 6.2% | 8.5% |
| 6 | 2,642 | 807 | 65 | 742 | 87.8% | 8.1% | 9.2% |
| 7 | 160 | 124 | 12 | 112 | 29.3% | 9.7% | 33.1% |
| 8 | 355 | 174 | 22 | 152 | 61.1% | 12.6% | 20.7% |
| 9 | 295 | 225 | 5 | 220 | 10.6% | 2.2% | 20.9% |
| 10 | 65 | 33 | 11 | 22 | 42.3% | 33.3% | 78.8% |
| All | 10,931 | 4,198 | 402 | 3,796 | $M$ 61.9% | $M$ 12.1% | $M$ 23.2% |

# Method 2: Code Clone Detector

- **Goal:** Use code clone detector to find clones of a sample of Stack Overflow snippets in a sample of GitHub projects

- *Why samples?*
  - Code clone detection is computationally expensive

- *Which snippets and projects to select?*
  - Random samples: Many **toy projects** on GitHub and many **irrelevant snippets** on Stack Overflow
  - Purposive sampling: Limited generalizability

# GitHub Project Sample

- Focus on **popular** GitHub projects
- High precision in selecting "engineered" software projects [Munaiah et al. 2017]
- Greater (potential) impact of licensing issues



Watcher count filter for non-fork Java GH projects (n=925,536)

Sample size:
3,000 / 2,313

# Stack Overflow Snippet Samples

- Non-trivial snippets retrieved from 100 most frequently referenced answers (n=111)

$$\Rightarrow S_{\mathrm{top100}}$$

- Non-trivial snippets retrieved from answers referenced in GitHub projects (n=137)

$$\Rightarrow S_{\mathrm{gh}}$$

- *External sources:* Only three snippets available under a more permissive license than CC BY-SA

# Code Clone Detector Calibration

https://pmd.github.io/

**Comparison of CPD configurations**



Legend:
- Precision (mt) — black dashed
- Precision (mt,ia) — blue dashed
- Precision (mt,ia,ii) — red dashed
- Recall (mt) — black solid
- Recall (mt,ia) — blue solid
- Recall (mt,ia,ii) — red solid

Y-axis: Precision and recall

X-axis: Minimum tokens

# Results

| Set | Snippets | | | | Files | | Repos |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | ALL | MATCHED | ANSWERS | MATCHED | MATCH. | REF | MATCHED |
| $S_{\mathrm{gh}}$ | 137 | 53 (39%) | 102 | 52 (51%) | 163 | 58 (36%) | 124 (5%) |
| $S_{\mathrm{top100}}$ | 111 | 48 (43%) | 85 | 46 (54%) | 173 | 25 (14%) | 125 (5%) |
| $\cup S$ | 222 | 101 (46%) | 169 | 86 (51%) | 297 | 70 (24%) | 199 (9%) |

# Method 3: Exact Matches

- **Goal:** Address shortcomings of Method 1 and 2
  - Increase sample sizes
  - Exclude snippets available on external sources
  - Systematically exclude short snippets

- Select as many projects and snippets as possible and search for (almost) exact matches

# Method 3: Exact Matches

Google Cloud
**GitHub**
**209m files in 4.1m projects**

✓ Project is not a fork, has ≥ 5 Java files and ≥ 1 watcher(s)
✓ File has ending `.java` has ≥ 68 NLOC ($Q_3$)

Java
**1.7m Java files in 64k projects**

Normalization and substring search

**10,358 matches**

Google Cloud
**stackoverflow**
**21m answers**

✓ Question tagged `java` or `android`
✓ Answer score ≥ 10
✓ Code block ≥ 6 NLOC

Java
**29k snippets from 24k answers**

https://iwsc2018.github.io/assets/img/sheep.png

# Details: Filtering of GitHub Projects



**File size filter for GH Java files (n=6,851,022)**

Number of files / Normalized file size (LOC)

Excluded, Excluded, ← 75% quantile



**Watcher count filter for GH Java projects (n=260,498)**

Number of projects / Number of watchers

Excluded, ← 75% quantile



**Fork filter for GH projects containing Java files (n=307,489)**

Number of projects / Fork, No fork

Excluded



**File count filter for GH Java projects (n=260,498)**

Number of projects / Number of Java files

Excluded, ← 25% quantile

# Details: Filtering of Stack Overflow Snippets

**Proxies for originality**

# Method 3: Filtering of Matches

10,358 matches

✓ Use heuristic to detect and exclude matches in mirrors of JDK and Android source code

1,379 matches

✓ Manually analyze answers, exclude snippets that are too trivial, incomplete, or copied from an external source
✓ Use **GitHub** API to remove matches where commit adding snippet is older than answer on Stack Overflow

1,369 matches

Check if attributed
(URL to answer or corresponding question)

Only **7.6%** attributed

# Attribution

*Attribution ratio:*

- Method 1 (regular expressions): 23 %
- Method 2 (code clone detector): 24 %
- Method 3 (exact matches): 8 %

*Conservative estimate:*

- **Attribution ratio $\leq$ 25%**

# Share-alike

Only **2%** of all analyzed repositories (all methods) containing code from Stack Overflow **attributed** its source and used a **compatible license** (not CC BY-SA, but GPL 3.0).

| SPDX license name | Number of repos containing a SO code snippet clone that was: | |
|---|---|---|
| | unattributed ($n = 2,962$) | attributed ($n = 329$) |
| Apache-2.0 | 921 (31.1%) | 99 (30.1%) |
| MIT | 621 (21.0%) | 72 (21.9%) |
| GPL-3.0 | 435 (14.7%) | 60 (18.2%) |
| GPL-2.0 | 284 (9.6%) | 21 (6.4%) |
| BSD-3-Clause | 82 (2.8%) | 9 (2.7%) |

Method 1

| SPDX license name | Number of repos containing a SO code snippet clone that was: | |
|---|---|---|
| | unattributed ($n = 144$) | attributed ($n = 55$) |
| None | 56 (38.9%) | 18 (32.7%) |
| Apache-2.0 | 33 (22.9%) | 15 (27.3%) |
| GPL-3.0 | 17 (11.8%) | 6 (10.9%) |
| MIT | 6 (4.2%) | 4 (7.3%) |
| GPL-2.0 | 4 (2.8%) | 2 (3.6%) |

Method 2

| SPDX license name | Number of repos containing a SO code snippet clone that was: | |
|---|---|---|
| | unattributed ($n = 1,169$) | attributed ($n = 163$) |
| Apache-2.0 | 353 (30.2%) | 36 (37.4%) |
| MIT | 239 (20.4%) | 25 (15.3%) |
| GPL-3.0 | 211 (18.0%) | 19 (11.7%) |
| None | 153 (13.1%) | 61 (37.4%) |
| GPL-2.0 | 89 (7.61%) | 8 (4.9%) |

Method 3

# Reaching out to Developers

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

- **75% not aware** of CC BY-SA licensing
  (see slide about online surveys)

- Many thankful responses

# Future Work

- *Tool support*: Support **maintainability** of copied snippets by automatically adding links to sources, integration into CI tools

- *Education:* Help developers **understand complex licensing situations** (not only for complete libraries but also for individual snippets)

- *Study:* Analyze links to better understand Stack Overflow's role in the **ecosystem** of documentation resources

# Code Duplication on Stack Overflow

Sebastian Baltes
sebastian.baltes@adelaide.edu.au
The University of Adelaide, Australia

Christoph Treude
christoph.treude@adelaide.edu.au
University of Adelaide, Australia

## ABSTRACT

Despite the unarguable importance of Stack Overflow for the daily work of many software developers and the existing knowledge about the impact of code duplication on software projects, the prevalence and implications of code clones on Stack Overflow have not yet received the attention they deserve. In this paper, we motivate why studies of this aspect are needed and how existing studies on code reuse from Stack Overflow differ from this new research direction. We present similarities and differences between code clones in general and code clones on Stack Overflow and point to open questions that need to be addressed to be able to make data-informed decisions about how to handle clones on this important platform. We present results from a first preliminary investigation indicating that clones on Stack Overflow are common and diverse and conclude with possible directions for future work.
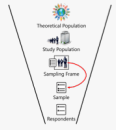
questions rather than supporting the maintenance and evolution of code on Stack Overflow.

Considering the importance that Stack Overflow has today for the daily work of many software developers worldwide and the fact that in many posts, non-trivial code snippets are collected and maintained, it is surprising that Stack Overflow does not have proper code versioning or bug tracking features. Text and code are versioned together as Markdown content [18], making it hard to identify changes to the code snippets in the revision view. [1] Furthermore, there are no language-specific syntax highlighting or error checking in Stack Overflow's online Markdown editor, leading to many snippets being not parseable, compilable, or even runnable [2]. Finally, there is no way to report bugs in Stack Overflow code snippets other than posting a comment or alternative answer.

Despite the above-mentioned challenges, code is maintained and does evolve on Stack Overflow [18]. The purpose of this article is to point the research community to open questions regarding code clones on Stack Overflow and how research in that area may inform significant improvements to the platform. We present a preliminary analysis of code clones within Stack Overflow and point
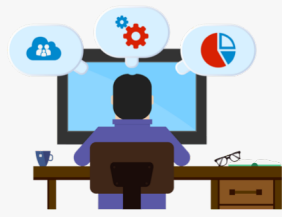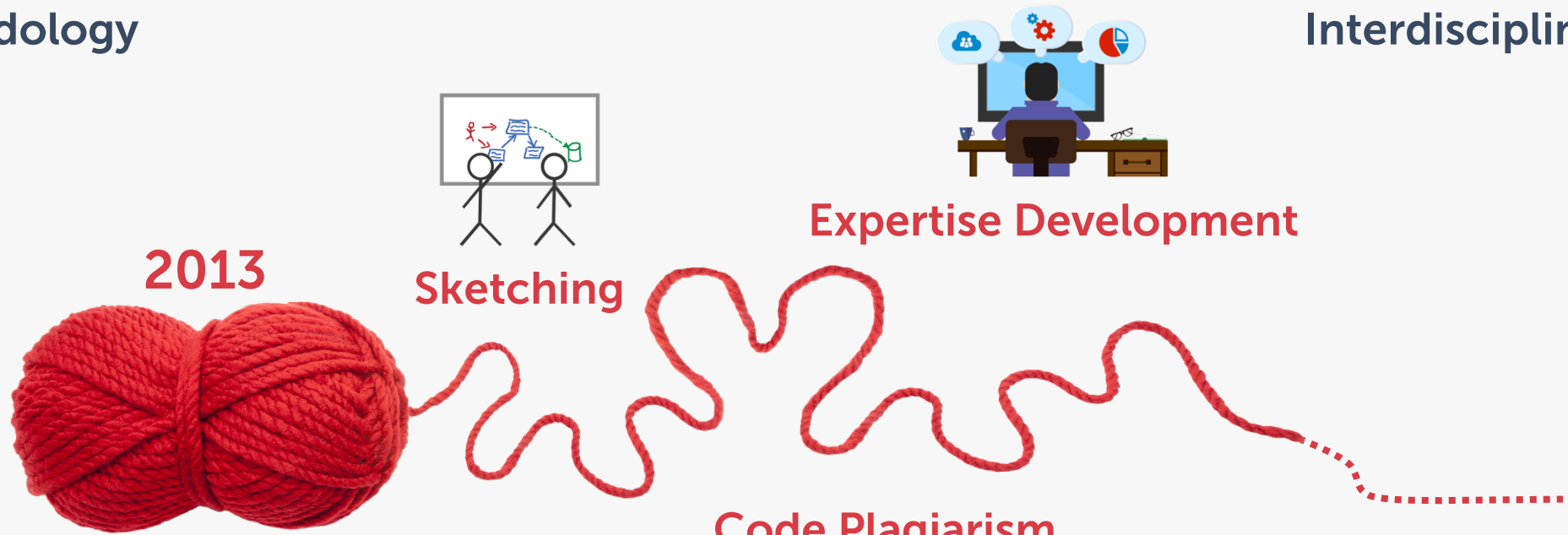
# Studied Habits

Issues in Sampling
Software Developers
**Methodology**

Constructing Urban
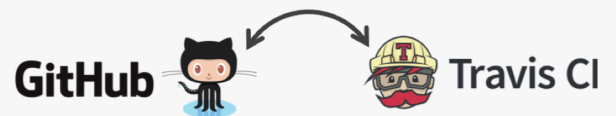Tourism Space Digitally
**Interdisciplinary Research**

**Expertise Development**

**2013**

**Sketching**

**Code Plagiarism**

stack**overflow**

**Regular Expressions**

RegViz

**Continuous Integration**

GitHub    Travis CI