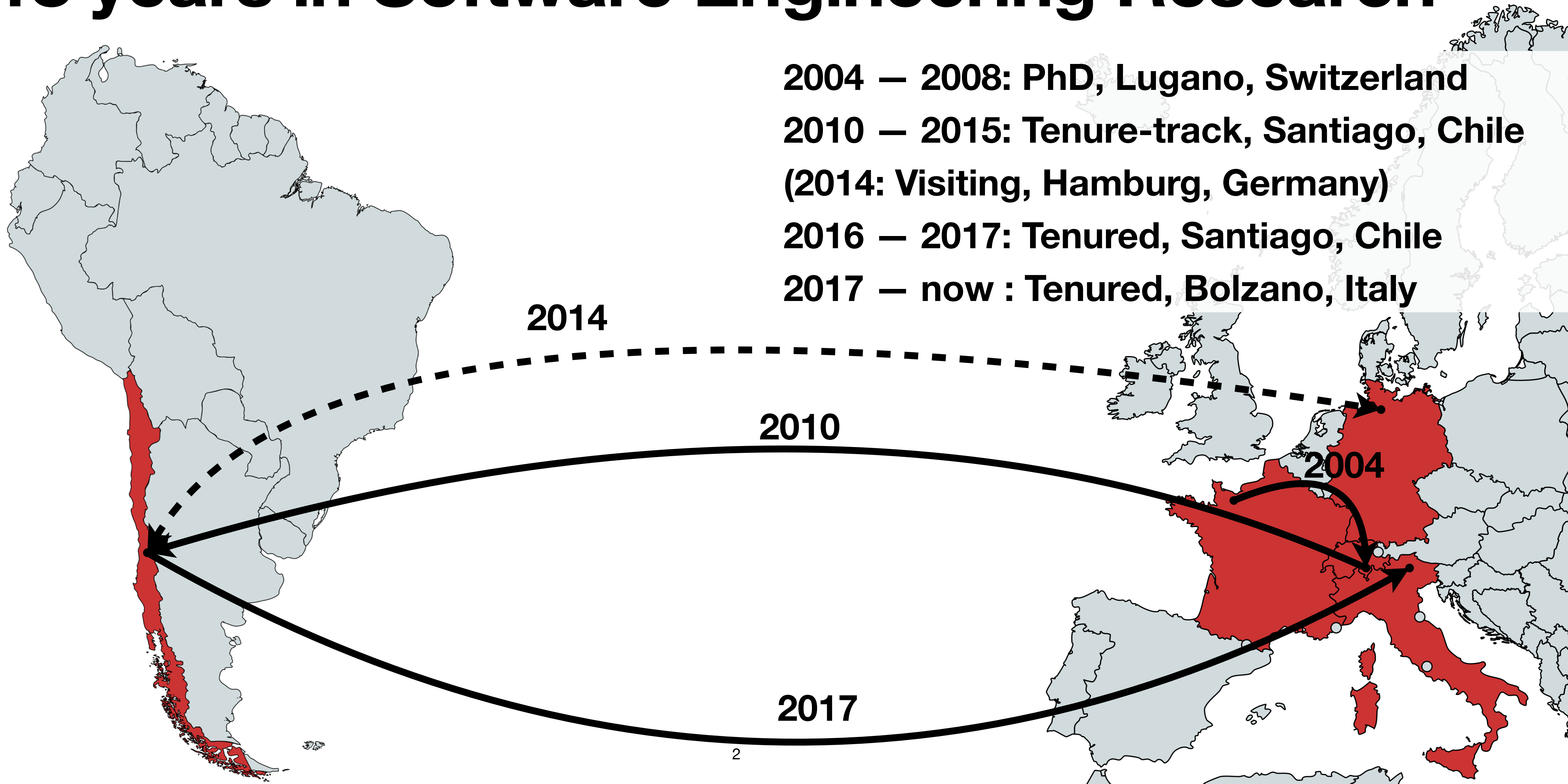# Machine Learning for Software Engineering with Scarce Data

**Audition DR2 CNRS**

**CRIStAL Lille (UMR 9189)**
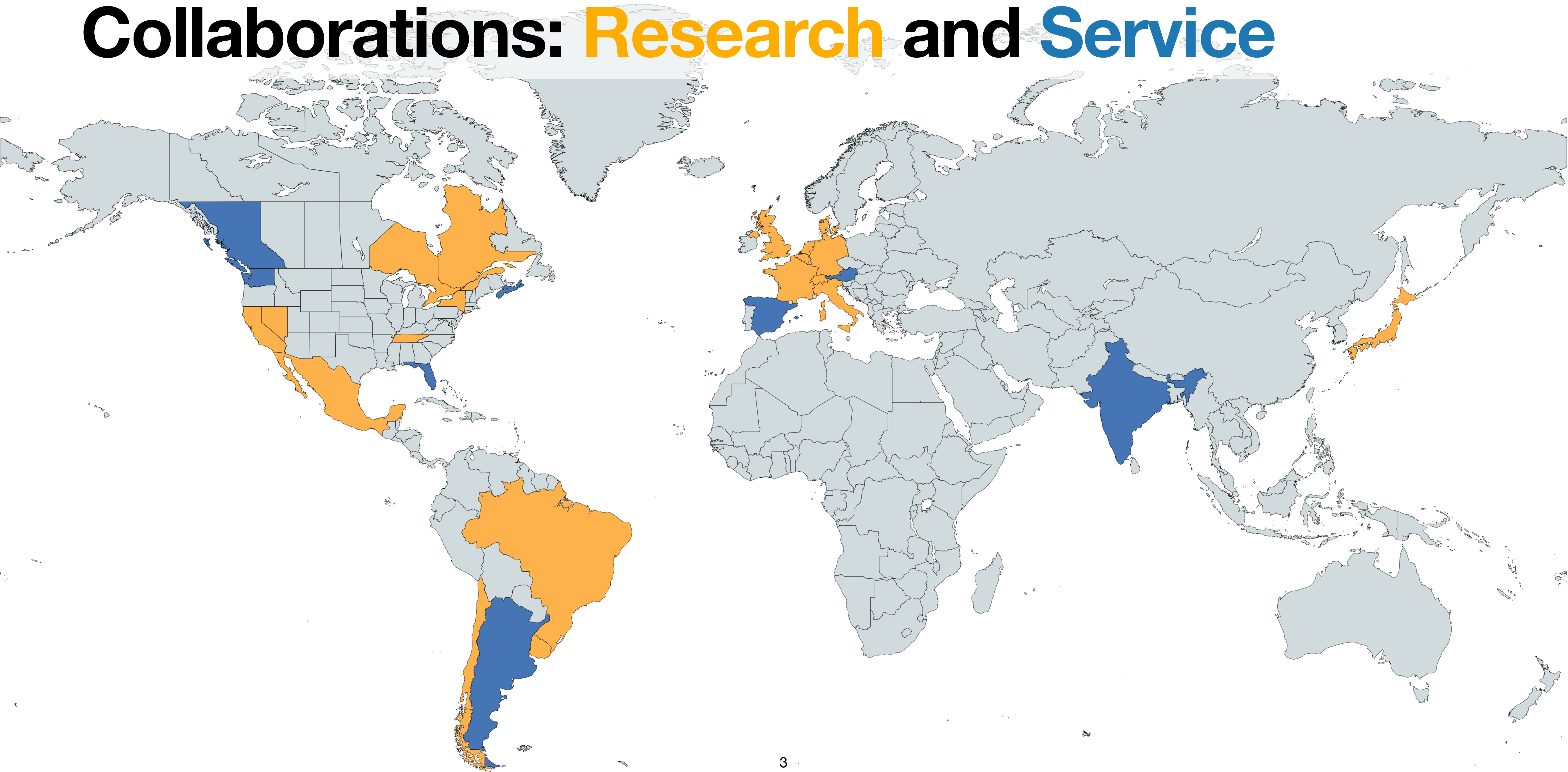
**Romain Robbes, 30/03/2022**

# 18 years in Software Engineering Research

2004 — 2008: PhD, Lugano, Switzerland

2010 — 2015: Tenure-track, Santiago, Chile

(2014: Visiting, Hamburg, Germany)

2016 — 2017: Tenured, Santiago, Chile

2017 — now : Tenured, Bolzano, Italy

2014

2010

2004

2017

# Collaborations: Research and Service

# Publications and service in quality venues

**12** Full papers in top SE & PL venues:
ICSE, FSE, ASE, ECOOP, OOPSLA
**1 award**
**+10** short papers

**13** Articles in top journals:
EMSE, TSE, JSS, CSUR
**+11** other journals

**17** Full papers in specialist venues:
ICSME, MSR, ICPC, SANER/WCRE
**4 awards**

**3** Track co-chair:
ICSE * 2, OOPSLA
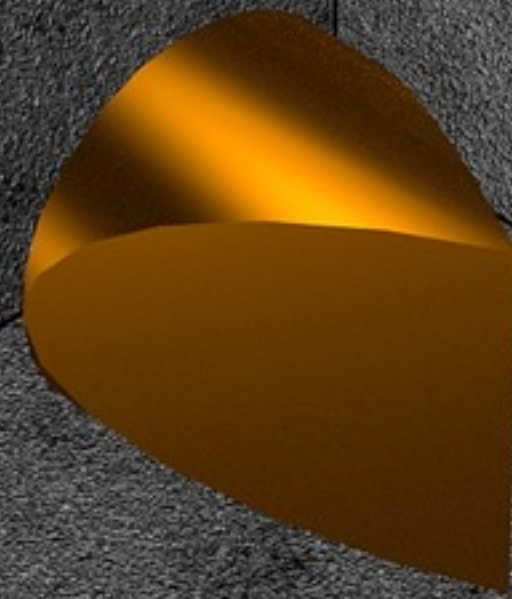
**2** Editorial boards:
EMSE, JSS

**3** PC co-chair:
MSR * 2, WCRE

**Software systems can be very complex …**

**… and still need to continually change**
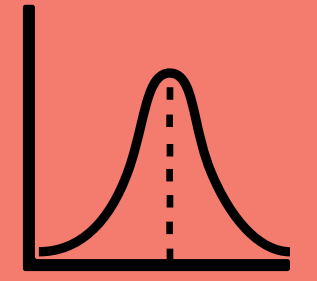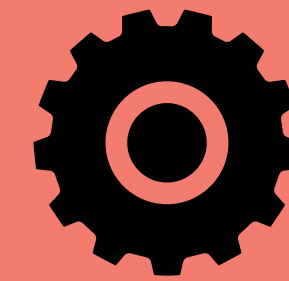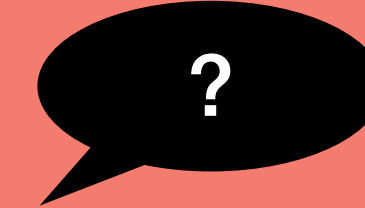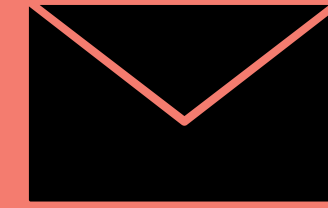
# Software engineering is a complex topic

Human Studies

Mining Software Repositories (MSR)

Machine Learning for Software Engineering (ML4SE)

# Major Research contributions

## Mining Software Repositories

**MSR brings multiple perspectives on Software Systems**

## Human Studies

**Bring back the human perspective:**

- **Do we ask the right questions?**
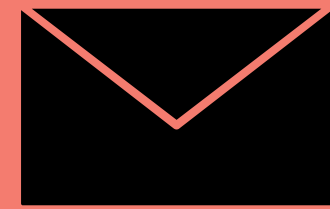- **Do we have the right answers?**

## ML4SE

**Machine Learning for Software Engineering**

**can integrate these multiple perspectives**

# SE artefacts mix structure and language

Hi Alice,

Can you help me figure out this bug?

```
// updating credentials
User u = DB.getUser(id)
u.setPassword(username)
DB.update(id, u)
```

Hope you can help, it's urgent. Thanks in advance!

Cheers,
Bob

**Deep Learning** learns **vector representations** of the meaning of words in context

| Bug | .8 | .2 | .7 | .1 | .9 | .5 | .7 | .3 | .1 | .3 | .4 | .9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Defect | .9 | .2 | .8 | .2 | .8 | .4 | .6 | .3 | .3 | .4 | .4 | .8 |

**Homogeneous** vector representations support **multiple types of artefacts**

Some models support **structure** such as **trees** or **graphs**

9

# Code: machine and human perspectives

**Program analysis ignores identifiers and code comments**

```
// … … … … … …
i272.f428(i823);
```

**Identifiers and comments are critical for program understanding**

```
// updating credentials
u.setPassword(username);
```

**Merging machine and human perspectives can detect new kinds of bugs (e.g., code-comment incoherence)**

# Example: Code Completion

# Every programmer uses code completion

alphabeticallySorted
exhaustive
identifier
listWhichIsQuiteLong
manyWhoAre
not
remotelyCloseTo
whatYouWant

➡️

whatYouWant
shortList

**2008** — **We define code completion as a task** and develop optimized algorithms

**2010** — **We implement and evaluate a code completion tool, still used in practice**

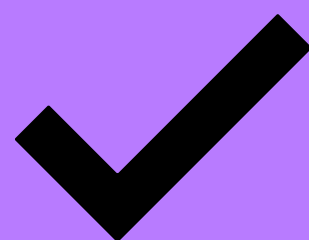2017 — Others use Neural Language Models for code completion, but **struggle**

**We successfully apply Neural Language Models to code completion**

**2020**

# Source Code Neural Language Models (NLMs)

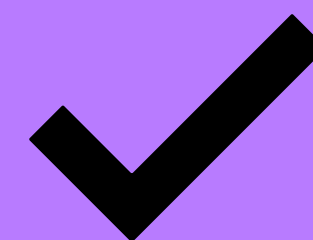**NLMs are pre-trained on source code by predicting tokens based on context**

```
public static void main(String [ ] args) {
    String commandName = args [0] ;
    String uncommonFlag = args [1] ;
```
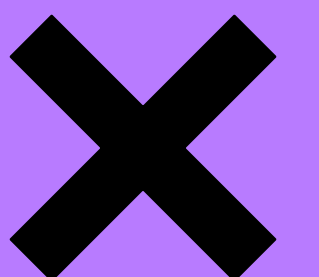
**Used "out of the box" for code completion** ✔

**Fine-tunable on many other tasks**

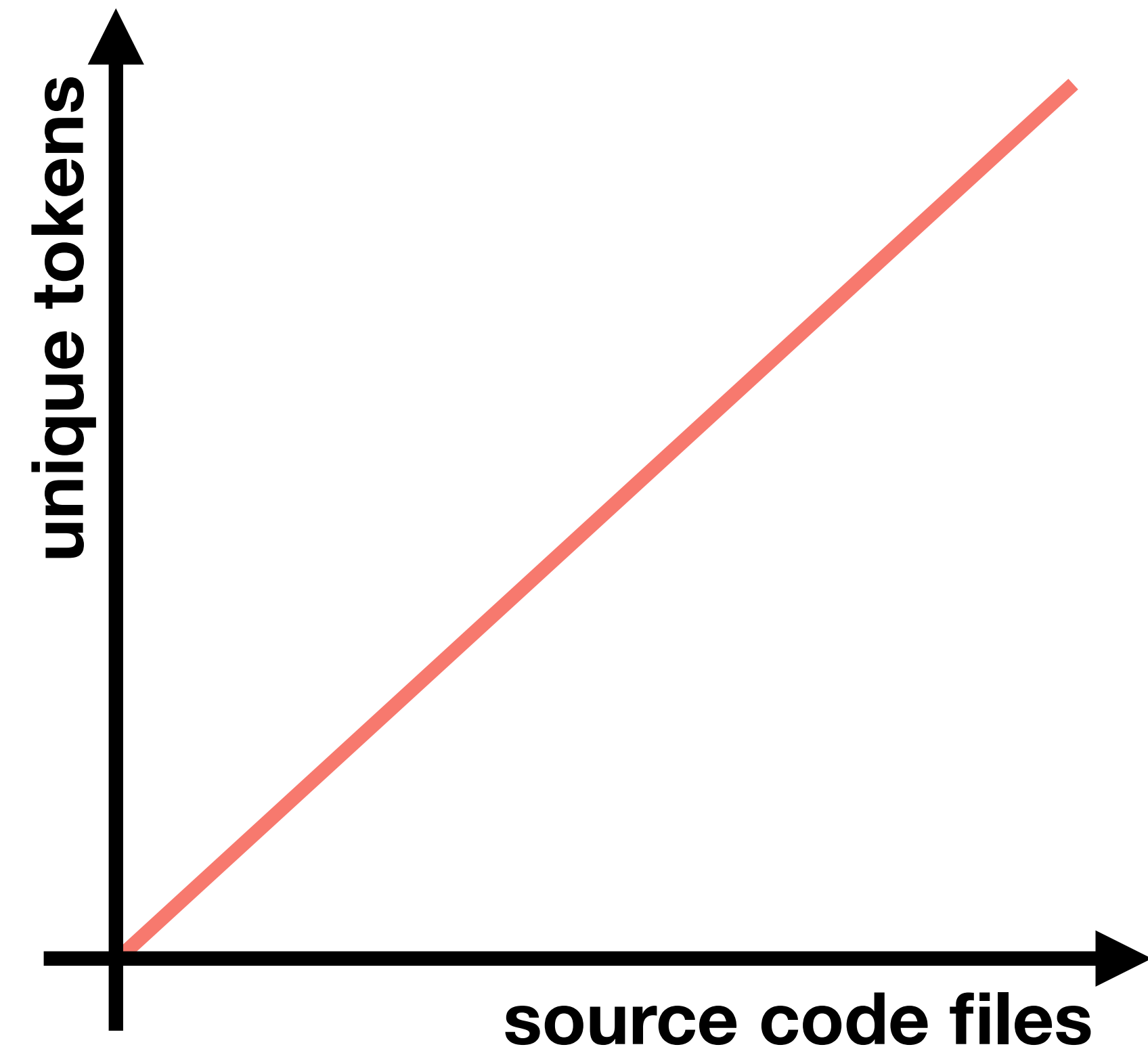Buggy code? { Yes: 0.02
             No: 0.98 ✔

**NLMs have a closed vocabulary of pre-defined tokens** ✘
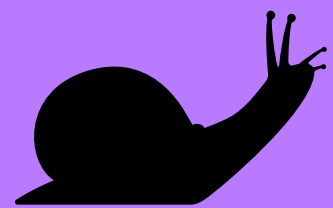
# Developers create new identifiers at will

```
foo
Foo
Bar
FooBar
FOO_BAR
a234
a282
a9095
my_taylor_is_rich
MeinTaylorIstReich
mon-tailleur-est-riche
FooBarManagerController
ArbitrarilyLongAndComplexIdentifierIsAllowed
FooBarManagerControllerRuntimeInstantiationException
```

unique tokens

source code files

**New code contains new tokens, leading to a Vocabulary Explosion**

# Large vocabulary leads to three issues

**Slow training**

Does not scale to millions of tokens

**Infrequent tokens**

| ? | ? | ? | ? | ? | ? |
|---|---|---|---|---|---|

Don't have good representations

**Unseen tokens**

\<UNK\>

New code has unknown tokens

**Only ~100 projects**
NLM < n-grams

## Are Deep Neural Networks the Best Choice for Modeling Source Code?

Vincent J. Hellendoorn
Computer Science Dept., UC Davis

Premkumar Devanbu
Computer Science Dept., UC Davis

# We went back to pre-processing

**Exhaustive preprocessing** on 10,000 Java projects (also C, Python)

| | # Unique Tokens | Corpus Size | Frequency (%< 10) | | OOV (% @ 75K) | |
|---|---|---|---|---|---|---|
| **No preprocessing**<br>`FooBarManagerController` | 11.6M | 1.0 | | 83% | | 79% |
| **Coding Conventions**<br>`Foo Bar Manager Controller` | 1.3M | 1.6 | | 81% | | 20% |
| **Best standard approach**<br>(Out of a dozen) | 0.5M | 1.9 | | 70% | | 9% |

**HUGE vocabulary**

**Many rare tokens**

**Many OOV tokens**

**This works, but … this is not enough**

# We looked further

**Exhaustive preprocessing** on 10,000 Java projects (also C, Python)

| | # Unique Tokens | Corpus Size | Frequency (%< 10) | OOV (% @ 75K) |
|---|---|---|---|---|
| **No preprocessing** *FooBarManagerController* | 11.6M | 1.0 | 83% | 79% |
| **Best standard approach** (Out of a dozen) | 0.5M | 1.9 | 70% | 9% |

**Byte-Pair Encoding (BPE) was used in Translation, but not yet in NLMs**

**BPE learns an open vocabulary bottom-up, from common character sequences**

# BPE solves all the vocabulary issues

**Exhaustive preprocessing** on 10,000 Java projects (also C, Python)

| | # Unique Tokens | Corpus Size | Frequency (%< 10) | OOV (% @ 75K) |
|---|---|---|---|---|
| **No preprocessing**<br>`FooBarManagerController` | 11.6M | 1.0 | 83% | 79% |
| **Best standard approach**<br>**(Out of a dozen)** | 0.5M | 1.9 | 70% | 9% |
| **Byte-Pair Encoding (BPE)** | 10K | 1.4 | 1% | 0% |

`Foo Bar ManagerController`

**Performs well**

`http client lib`

`Se hr Lan ge Que ll code Ken nun g`

**Degrades gracefully**

**1,000 times smaller**

**Few Rare tokens**

**No OOV**

# Outcomes

We built an **open-vocabulary** NLM (with additional code-specific extensions)

It scales to **10,000+** software projects, **100 X** more than previous work

A thorough evaluation shows **large improvements**



**top 2%** of submitted papers
(distinguished paper)

**90+ citations in 2 years**
data, code, models available

**Needed expertise in both ML & SE, willingness to revisit assumptions**

# Research Program: ML4SE in Scarce Data Scenarios

# Source Code NLMs are getting … larger

**Ours (2020)**

**OpenAI's Codex (2021)**
🐙 **GitHub** Copilot

**DeepMind's AlphaCode (2022)**

**1.8 GB**

**159 GB**

**715 GB**

**Trained on 55 Million
GitHub code repositories
Runs only in the cloud
Closed source**

# Are there 55 million repositories of COBOL?

## MANTIS?

## NCL?

## 4D?

# Legacy systems, in orphan languages

**Legacy Systems:** old, but **critical infrastructure** for important institutions

Complex code, many quality issues
Missing documentation and tests

Knowledge loss over decades
Constant need for software evolution

**Orphan Languages:** decades-old proprietary languages with **little or no support**

Almost no open source presence:
Not enough data for out of the box ML approaches

# Two challenges: scientific and practical

**Learn with less** **data**          **In practical settings**

**For orphan languages**          **Relevant for legacy systems**

# Learn with less data, but with more …

**"Big data" ➡ Small preprocessing**

source code == **natural text**

**Large but generic NLMs**

**"Small data" ➡ SE-specific preprocessing**

source code != **natural text**
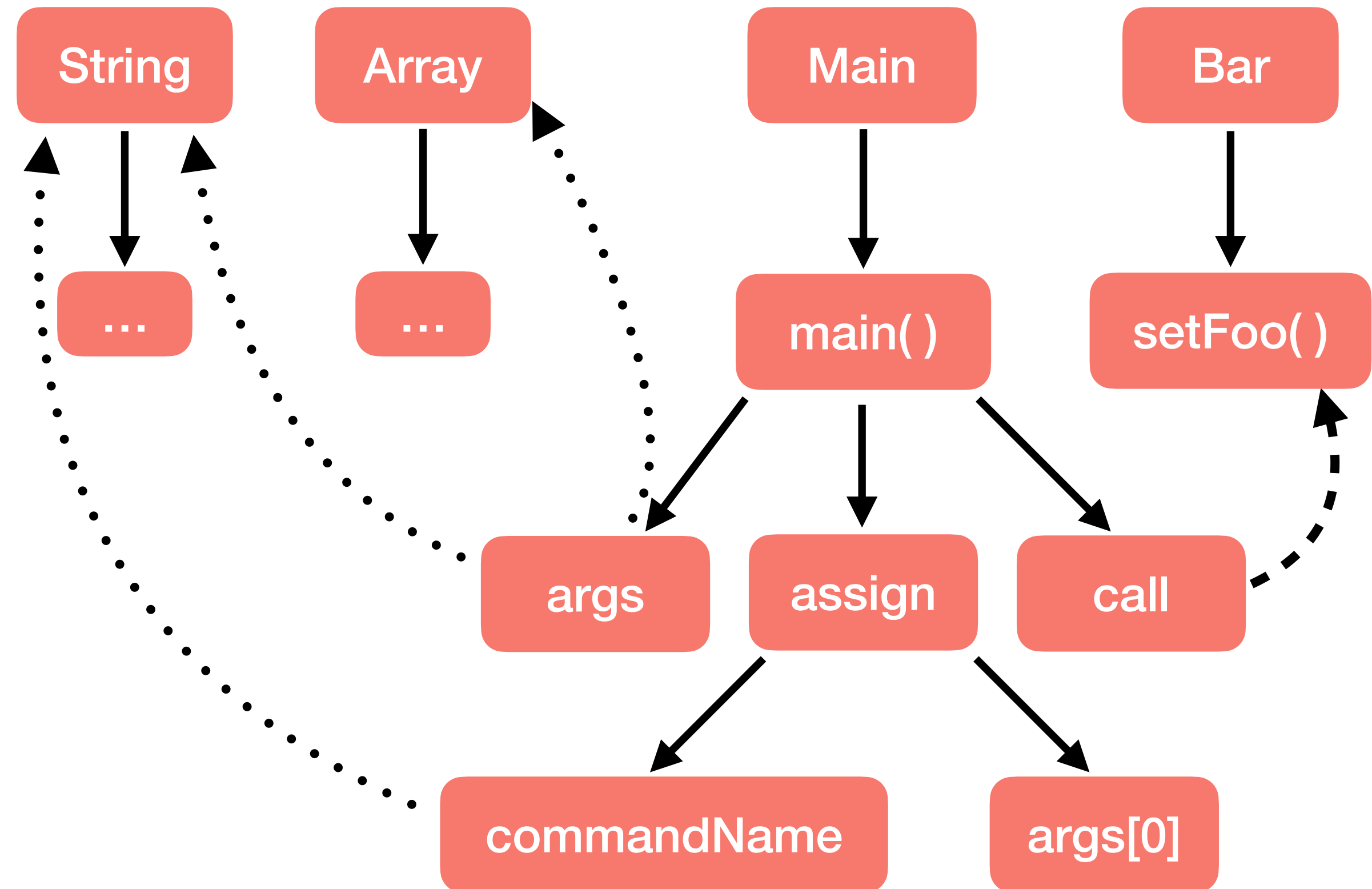
**more structure**

**more information**

**better transfer**

# Learn from less with **more structure**

```
public static void main(String [ ] args) {
    String commandName = args [0] ;
    String uncommonFlag = args [1] ;
    myBar.setFoo(args[2])
```

**Learning implicit information from scratch requires a large model & lots of data**

**Explicit information: smaller model, less data**

String

Array

Main

Bar

...

...

main( )

setFoo( )

args

assign

call

commandName

args[0]

# Most papers focus on the ML perspective

… data work | Collect and process artefacts

Everybody wants to do the model work, not the …

Develop approach

Run on benchmarks

# Practical scenarios require a holistic view

… scoping — **Define the problem**

… data work — **Collect and process artefacts**

Everybody wants to do the model work, not the … — **Develop approach** / **Run on benchmarks**

… integration — **Tool integration**

… evaluation — **Practical Evaluations**

# Practical scenarios provide feedback…

… which may lead to **revisiting decisions**

Define the problem

Collect and process artefacts

Develop approach

Run on benchmarks

Tool integration

Practical Evaluations

# The ideal team for this project is RMoD
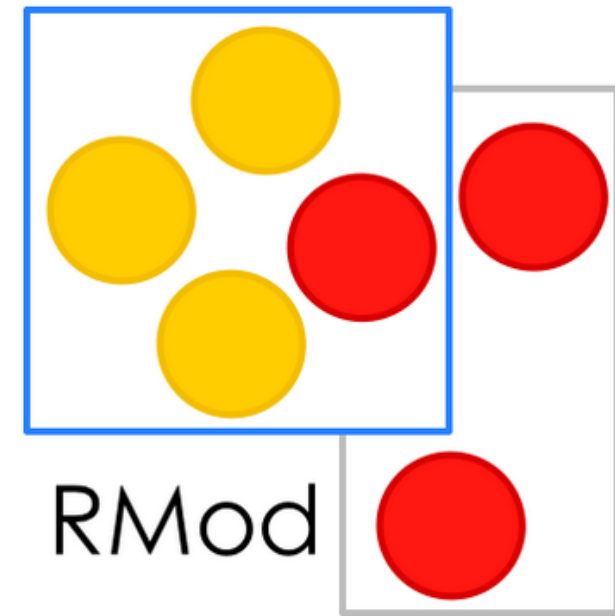
**RMoD** works with **industrial partners** to help them support their **legacy systems**

**Practical problems** include test generation, code migration, business rule extraction, bug triage, & many others
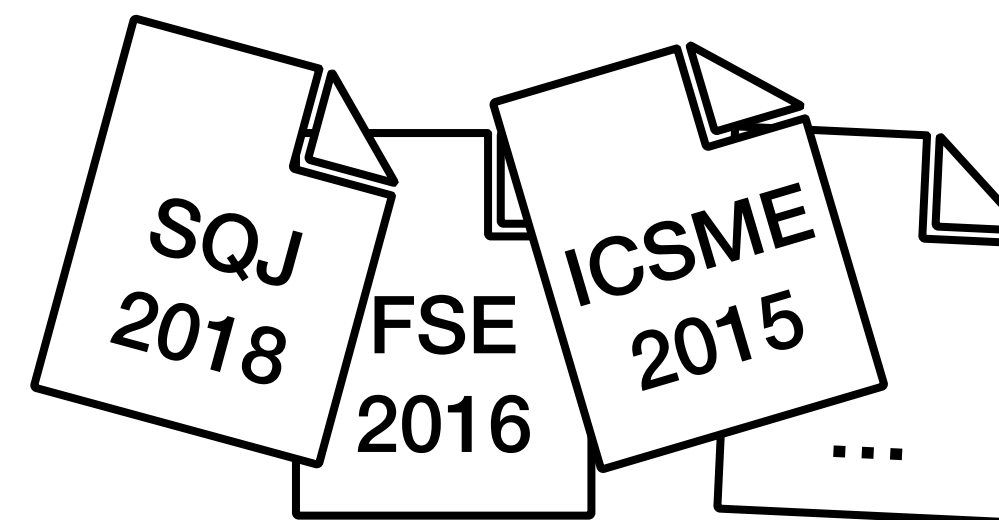
# Integrating with RMoD

**Reinforce RMoD in:** MSR · Human Studies

**RMoD seeks:** ML4SE

**For me, a practical perspective:** Berger Levrault · SIEMENS · arolla · Lifeware

**RMod & I have extensive previous collaborations:** SQJ 2018 · FSE 2016 · ICSME 2015 · … · SATToSE 2019
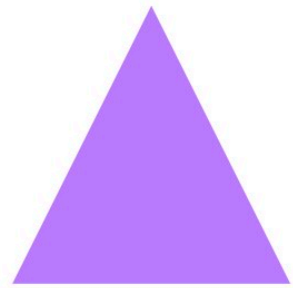
# Major SE Contributions

**MSR** brings **multiple perspectives** on Software Systems

**Human Studies** bring back the **human perspective**

**ML4SE integrates** these multiple perspectives

# Perspectives on code completion

alphabeticallySorted
exhaustive
identifier
listWhichIsQuiteLong
manyWhoAre
not
remotelyCloseTo
whatYouWant

→ whatYouWant
shortList

Defined code completion as a task

Released and evaluated a tool, still used

Applied Deep Learning to Code Completion
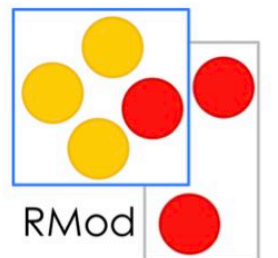
# A research program with two challenges

Learn with less data    In practical settings

For orphan languages    Relevant for legacy systems

# RMoD is the ideal team

Its partners have legacy systems & practical problems

Berger Levrault BL    SIEMENS    arolla    Lifeware

My skills & interests are a great match with RMoD