

# Promoting Open Science in Test-driven Software Experiments



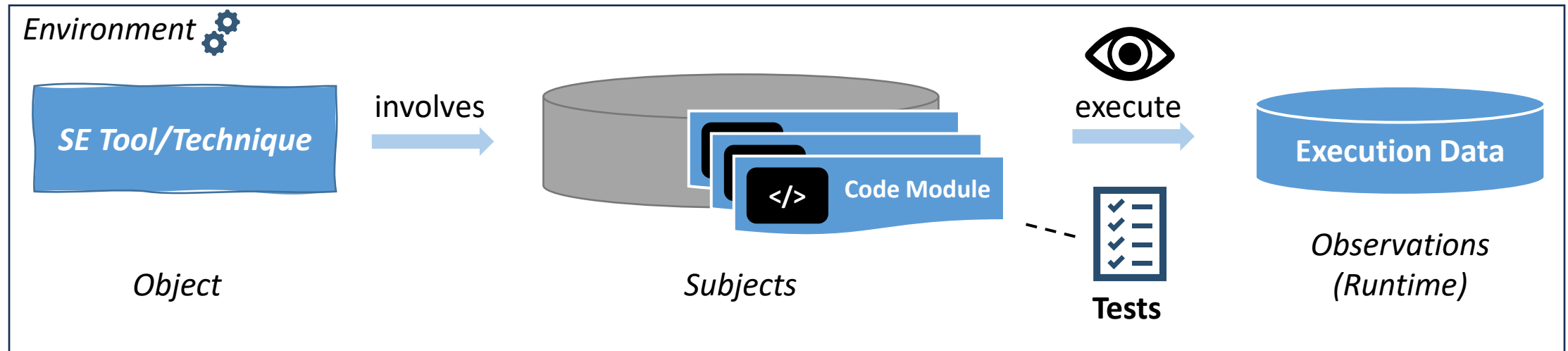
Marcus Kessel [marcus.kessel@uni-mannheim.de](mailto:marcus.kessel@uni-mannheim.de), Colin Atkinson [colin.atkinson@uni-mannheim.de](mailto:colin.atkinson@uni-mannheim.de)



*ICST 2025 @ Naples, Italy – Journal First (Journal of Systems and Software)*

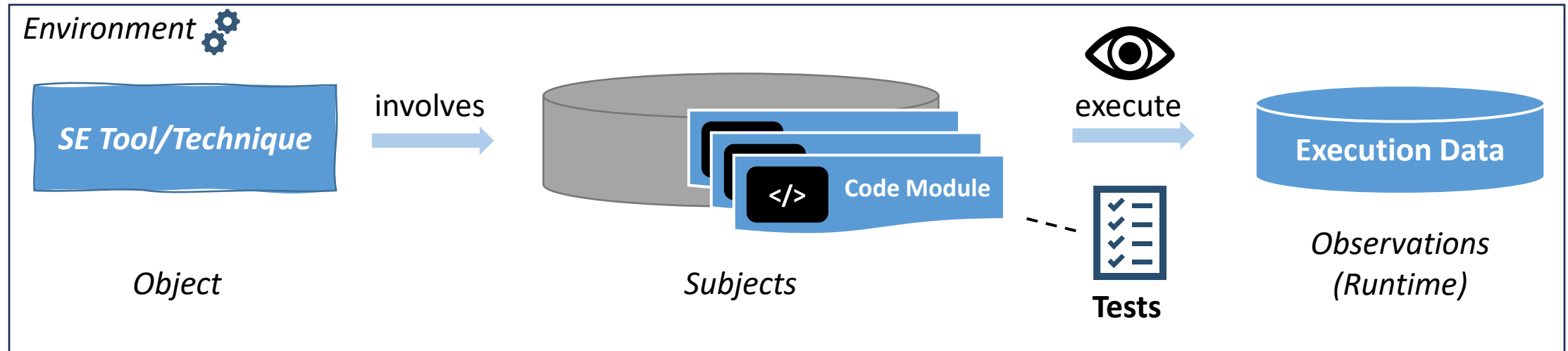


# Test-Driven Software Experiments (TDSEs)



→ TDSEs are experiments that involve controlled testing of software subjects (i.e., code modules) under various conditions, revealing important properties of the code's run-time behavior that cannot be predicted solely through static analysis

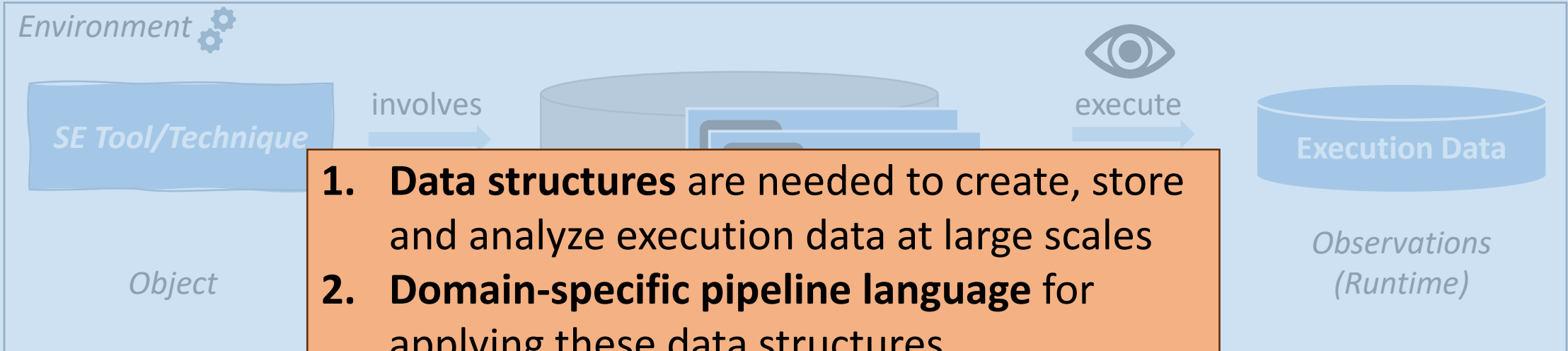
# Open Science Challenges in TDSEs



Reproducibility of TDSEs is hindered by

- ad hoc scripting and custom code (general-purpose languages)
- ad hoc representation of tests and execution data (beyond pass/fail)
- lack of executable corpora
- complex execution logistics

# Open Science Challenges in TDSEs



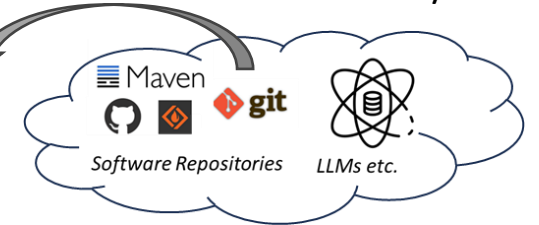
- 1. Data structures** are needed to create, store and analyze execution data at large scales
- 2. Domain-specific pipeline language** for applying these data structures
- 3. A dedicated platform** to execute these pipelines (automation) and store code/data

## Reproducibility

- ad hoc scripting and custom code (general-purpose languages)
- ad hoc representation of tests and execution data (beyond pass/fail)
- lack of executable corpora
- complex execution logistics

# Tabular Data Structures and DSL Languages for TDSEs

code modules/tests



**Problem:** GCD (Greatest Common Divisor)  
Interface: `gcd(int, int)->int`

`test(p1:Class, p2:int, p3:int)`

	A	B	C	D	E
1		create	?p1		
2		gcd	A1	?p2	?p3

output    operation    inputs

*N versions (candidate implementations)*

	V1	V2	V3	...	V <sub>n</sub>
test1	<code>test(V1,03,07)</code>	<code>test(V2, 03,07)</code>	<code>test(V3, 03,07)</code>		<code>test(V<sub>n</sub>, 03,07)</code>
test2	<code>test(V1,10,15)</code>	<code>test(V2,10,15)</code>	<code>test(V3,10,15)</code>		<code>test(V<sub>n</sub>,10,15)</code>
test3	<code>test(V1,49,14)</code>	<code>test(V2,49,14)</code>	<code>test(V3,49,14)</code>		<code>test(V<sub>n</sub>,49,14)</code>
...					
test <sub>m</sub>	<code>test<sub>m</sub>(V1, x, y)</code>	<code>test<sub>m</sub>(V2, x, y)</code>	<code>test<sub>m</sub>(V3, x, y)</code>		<code>test<sub>m</sub>(V<sub>n</sub>, x, y)</code>

Tests

**Stimulus Sheet**

**Stimulus Matrix (SM)**



**Actuation Sheet**

**Stimulus Response Matrix (SRM)**



DSL "LSL"



Test Driver "Arena"

	A	B	C	D	E
1		create	V1		
2	1	gcd	A1	3	7

	V1	V2	V3	...	V <sub>n</sub>
test1	<code>test(V1,03,07)</code> 1	<code>test(V2,03,07)</code> 7	<code>test(V3,03,07)</code> -7		<code>test(V<sub>n</sub>,03,07)</code> ...
test2	<code>test(V1,10,15)</code> 5	<code>test(V2,10,15)</code> 15	<code>test(V3,10,15)</code> -15		<code>test(V<sub>n</sub>,10,15)</code> ...
test3	<code>test(V1,49,14)</code> 7	<code>test(V2,49,14)</code> 14	<code>test(V3,49,14)</code> 28		<code>test(V<sub>n</sub>,49,14)</code> ...
...					
test <sub>m</sub>	<code>test<sub>m</sub>(V1, x, y)</code> ...	<code>test<sub>m</sub>(V2, x, y)</code> ...	<code>test<sub>m</sub>(V3, x, y)</code> ...		<code>test<sub>m</sub>(V<sub>n</sub>, x, y)</code> ...

Differential Testing (identify discrepancies)

**Platform**



LASSO Platform

# Pipelines – LASSO Scripting Language (LSL)



```

dataSource 'lasso_quickstart'
study(name: 'OpenAI-DGAI') {

  profile('java17Profile') {
    scope('class') { type = 'class' }
    environment('java17') {
      image = 'maven:3.9-eclipse-temurin-17'
    }
  }

  def humanEval = loadBenchmark("humaneval-java-reworded")

  action(name: "createStimulusMatrices") {
    execute {
      def myProblems = [humanEval.abstractions['HumanEval_13_greatest_common_divisor']]
      myProblems.each { problem ->
        stimulusMatrix(problem.id, problem.lq, ["impls"/], problem.tests) // id, interface, impls, tests
      }
    }
  }

  action(name: 'generateCodeGpt', type: 'GenerateCodeOpenAI') {
    dependsOn 'createStimulusMatrices'
    include '*'
    profile('java17Profile')

    apiKey = "demo"
    model = "gpt-4o-mini"
    samples = 1

    prompt { stimulusMatrix ->
      def prompt = [:] // create prompt model
      prompt.promptContent = """implement a java class with the following interface specification, but do not inherit a java interface:
      """ + stimulusMatrix.lq + """ Only output the java class and nothing else."""
      return [prompt] // list of prompts
    }
  }

  action(name: 'generateTestsGpt', type: 'GenerateTestsOpenAI') {
    dependsOn 'generateCodeGpt'
    include '*'
    profile('java17Profile')

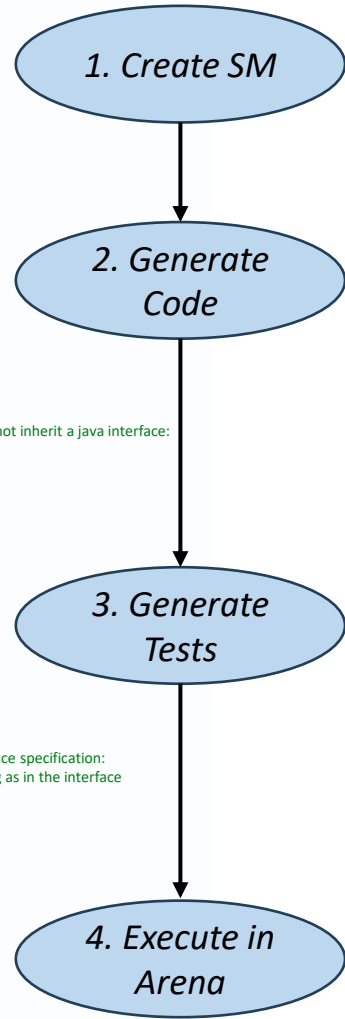
    apiKey = "demo"
    model = "gpt-4o-mini"
    samples = 1

    prompt { stimulusMatrix ->
      def prompt = [:] // create prompt model
      prompt.promptContent = """generate a junit test class to test the functionality of the following interface specification:
      """ + stimulusMatrix.lq + """ Assume that the specification is encapsulated in a class that uses the same naming as in the interface
      specification. Only output the JUnit test class and nothing else."""
      return [prompt] // list of prompts
    }
  }

  action(name: 'execute', type: 'Arena') {
    features = ["cc"] // coverage: code coverage, mutation testing etc.

    dependsOn 'generateTestsGpt'
    include '*'
    profile('java17Profile')
  }
}

```



test1
test2
test3

New Stimulus Matrix (SM) from coding problem

	V1	V2	...	V <sub>n</sub>
test1				
test2				
test3				

Columns – Add implementations (versions)

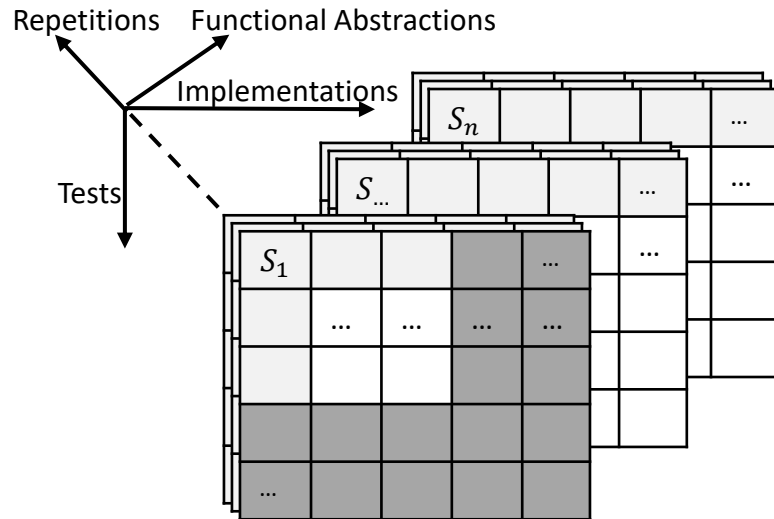
	V1	V2	...	V <sub>n</sub>
test1				
test2				
...				
test <sub>m</sub>				

Rows – Add tests

	V1	V2	...	V <sub>n</sub>
test1	test(V1,03,07) 1	test(V2,03,07) 7		test(V <sub>n</sub> ,03,07) ...
test2	test(V1,10,15) 5	test(V2,10,15) 15		test(V <sub>n</sub> ,10,15) ...
...				
test <sub>m</sub>	test <sub>m</sub> (V1, x, y) ...	test <sub>m</sub> (V2, x, y) ...		test <sub>m</sub> (V <sub>n</sub> , x, y) ...

Observe/store outputs, metrics ...

# Data-driven Analysis of SRMs

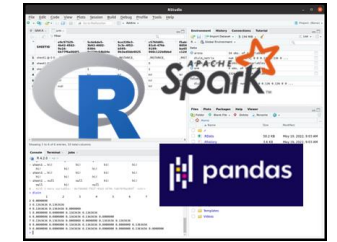


*Collections of SRMs  
(navigation and dimensions)*

*Oracle values or reference implementation*

SM

	Oracle	V1
test1	test(V1,03,07)	test(V1,03,07)
test2	test(V1,10,15)	test(V1,10,15)
test3	test(V1,49,14)	test(V1,49,14)



execute

SRM

	Oracle		V1		...
test1	test(V1,03,07)	1	test(V1,03,07)	1	
test2	test(V1,10,15)	5	test(V1,10,15)	5	
test3	test(V1,49,14)	7	test(V1,49,14)	11	

*Compare observed outputs (columns)*

*Data-Driven Analysis  
(incl. differential testing and  
judging functional equivalence)*

# Demonstration – Reproducing TDSEs with LASSO



Purpose/ Design		Same Purpose (i.e., RQs)		Different Purpose (i.e., RQs)
		Same Design	Different Design	
Team	Original Team	<b>Repetition:</b> <ul style="list-style-type: none"> <li>• same purpose</li> <li>• original team</li> <li>• same (logical) design</li> </ul>	<b>Reproduction:</b> <ul style="list-style-type: none"> <li>• same purpose</li> <li>• original or independent team</li> <li>• different (logical) design</li> </ul>	<b>Repurposing:</b> <ul style="list-style-type: none"> <li>• different purpose (i.e., RQs)</li> <li>• original or independent team</li> <li>• different (logical) design</li> <li>• performed in a different way</li> </ul>
	Different Team	<b>Replication:</b> <ul style="list-style-type: none"> <li>• same purpose</li> <li>• independent team</li> <li>• same (logical) design</li> </ul>		

# TDSE Replication – Benchmarking LLMs for Code Generation



IEEE TRANSACTIONS ON SOFTWARE ENGINEERING, VOL. 49, NO. 7, JULY 2023 3675

## MultiPL-E: A Scalable and Polyglot Approach to Benchmarking Neural Code Generation

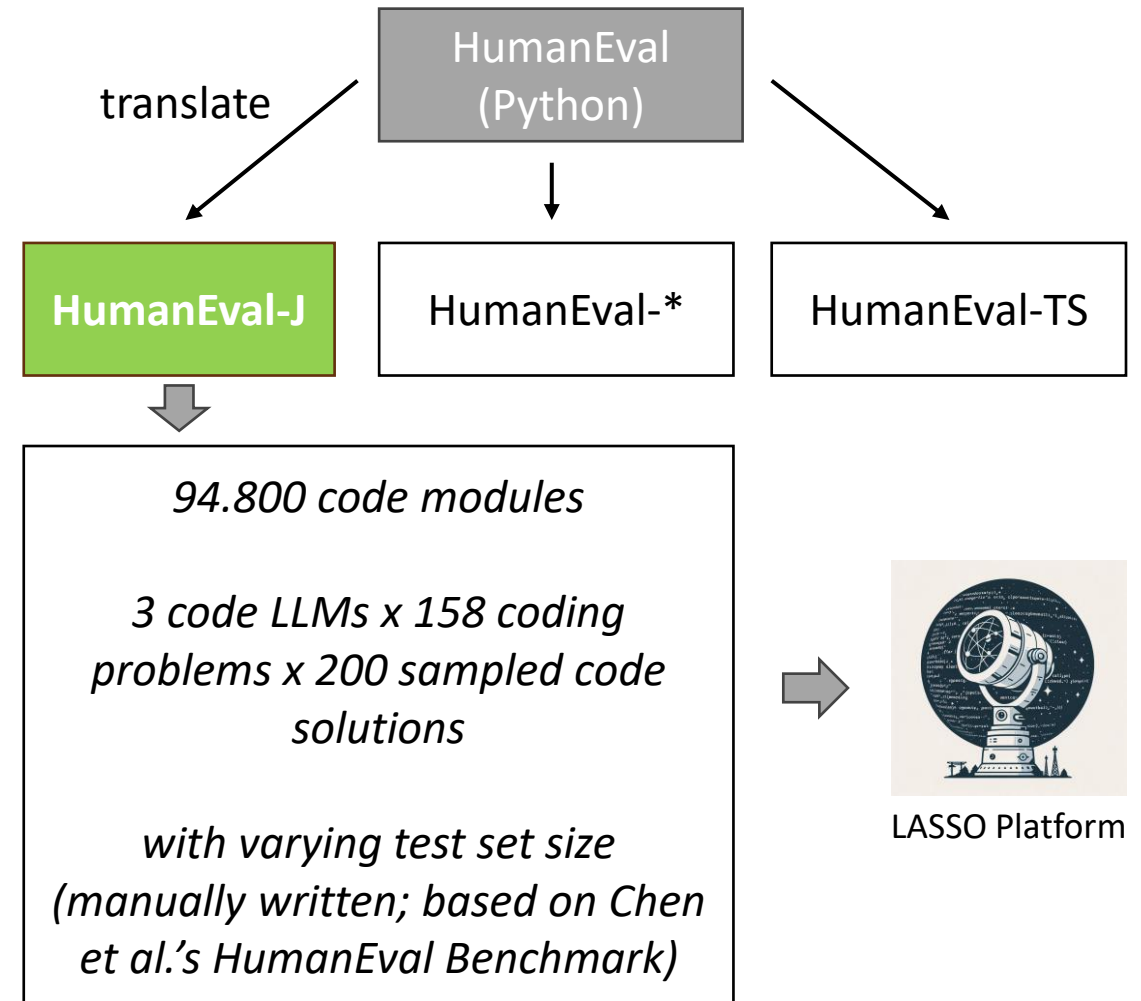
Federico Cassano<sup>1</sup>, John Gouwar<sup>2</sup>, Daniel Nguyen<sup>3</sup>, Sydney Nguyen<sup>3</sup>, Luna Phipps-Costin<sup>4</sup>, Donald Pinckney<sup>5</sup>, Ming-Ho Yee<sup>6</sup>, Yangtian Zi<sup>7</sup>, Carolyn Jane Anderson<sup>8</sup>, Molly Q Feldman<sup>9</sup>, Arjun Guha<sup>10</sup>, Michael Greenberg<sup>11</sup>, and Abhinav Jangda<sup>12</sup>

**Abstract**—Large language models have demonstrated the ability to generate both natural language and programming language text. Although contemporary code generation models are trained on corpora with several programming languages, they are tested using benchmarks that are typically monolingual. The most widely used code generation benchmarks only target Python, so there is little quantitative evidence of how code generation models perform on other programming languages. We propose MultiPL-E, a system for translating unit test-driven code generation benchmarks to new languages. We create the first massively multilingual code generation benchmark by using MultiPL-E to translate two popular Python code generation benchmarks to 18 additional programming languages. We use MultiPL-E to extend the HumanEval benchmark (Chen et al., 2021) and MBPP benchmark (Austin et al., 2021) to 18 languages that encompass a range of programming paradigms and popularity. Using these new parallel benchmarks, we evaluate the multi-language performance of three state-of-the-art code generation models: Codex (Chen et al., 2021), CodeGen (Nijkamp et al., 2022) and InCoder (Fried et al., 2022). We find that Codex matches or even exceeds its performance on Python for several other languages. The range of programming languages represented in MultiPL-E allow us to explore the impact of language frequency and language features on model performance. Finally, the MultiPL-E approach of compiling code generation benchmarks to new programming languages is both scalable and extensible, making it straightforward to evaluate new models, benchmarks, and languages.

**Index Terms**—B.2.3 reliability, testing, and fault-tolerance, I.5.1.D neural nets.

### I. INTRODUCTION

CODE generation models, also known as large language models (LLMs) of code, are deep neural networks trained on massive corpora of source code. Over the past few years, code generation models have demonstrated their utility on a wide variety of software engineering tasks, including test generation, documentation generation, and even synthesizing working programs from natural language descriptions [1], [3], [4], [5], [6]. New products such as GitHub Copilot<sup>1</sup>, Amazon CodeWhisperer<sup>2</sup>, and Tabnine<sup>3</sup> built on code generation models are growing in popularity with developers [7]. Although several code generation models are trained on multiple programming languages, they are typically only evaluated on a single programming language: Python. Machine learning researchers are familiar with Python: they have painstakingly constructed several Python code generation benchmarks [1], [2], [8], [9] and it is the best represented language in training datasets [1], [2], [9], [10]. However, we should also evaluate code generation models with other languages to support a wider variety of programmers. There is prior work on multi-language evaluation [6], but it uses perplexity as a proxy for performance, instead of benchmarks that check correctness.

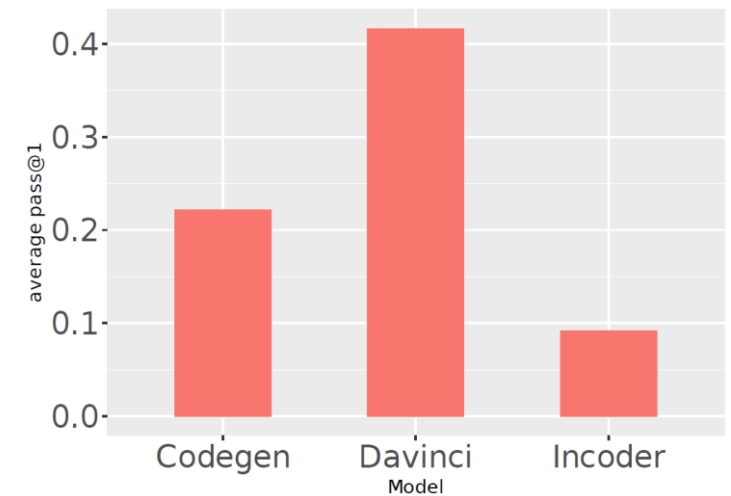


# Replication – Judging Functional Correctness

**RQ1:** What is the Java code completion performance of the three code models (with temperature set to 0.2) as judged by functional correctness? (pass@k metric)

	$O$	$S_{1LLM_1}$	...	$S_{kLLM_1}$	...	$S_{1LLM_n}$	...	$S_{kLLM_n}$
$T_1$	$T_1(O, \dots)$	$T_1(S_{1LLM_1}, \dots)$	...	$T_1(S_{kLLM_1}, \dots)$	...	$T_1(S_{1LLM_n}, \dots)$	...	$T_1(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...	...
$T_j$	$T_j(O, \dots)$	$T_j(S_{1LLM_1}, \dots)$	...	$T_j(S_{kLLM_1}, \dots)$	...	$T_j(S_{1LLM_n}, \dots)$	...	$T_j(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...	...
$T_m$	$T_m(O, \dots)$	$T_m(S_{1LLM_1}, \dots)$	...	$T_m(S_{kLLM_1}, \dots)$	...	$T_m(S_{1LLM_n}, \dots)$	...	$T_m(S_{kLLM_n}, \dots)$

Same design, same RQ



Test@Invocation	Oracle	CodeGen <sub>1</sub>	CodeGen <sub>2</sub>	CodeGen <sub>3</sub>	CodeGen <sub>4</sub>
$T_1@1$	1	1	-7	1	1
$T_2@1$	5	5	-15	5	5
$T_3@1$	7	7	28	7	7
$T_4@1$	12	12	60	12	12

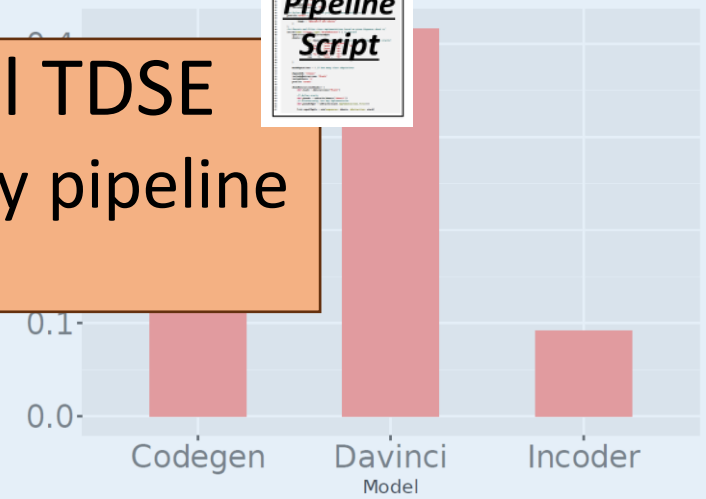
Judging functional correctness using SRMs

# Replication – Judging Functional Correctness

RQ: What is the Java code completion performance of the three code models (with temperature set to 0.2) as judged by functional correctness? (pass@k metric)

	$O$	$S_{1LL}$	$S_{2LL}$	$S_{3LL}$	$S_{4LL}$	$S_{5LL}$	$S_{6LL}$	$S_{7LL}$	$S_{8LL}$
$T_1$	$T_1(O, \dots)$	$T_1(S_{1LL}, \dots)$							
...	...	...							
$T_j$	$T_j(O, \dots)$	$T_j(S_{1LL}, \dots)$							
...	...	...							
$T_m$	$T_m(O, \dots)$	$T_m(S_{1LLM_1}, \dots)$	...	$T_m(S_{kLLM_1}, \dots)$	...	$T_m(S_{1LLM_n}, \dots)$	...	$T_m(S_{kLLM_n}, \dots)$	

**Successful replication of original TDSE**  
 → Replication resulted in one LSL study pipeline script (LOC < 50)



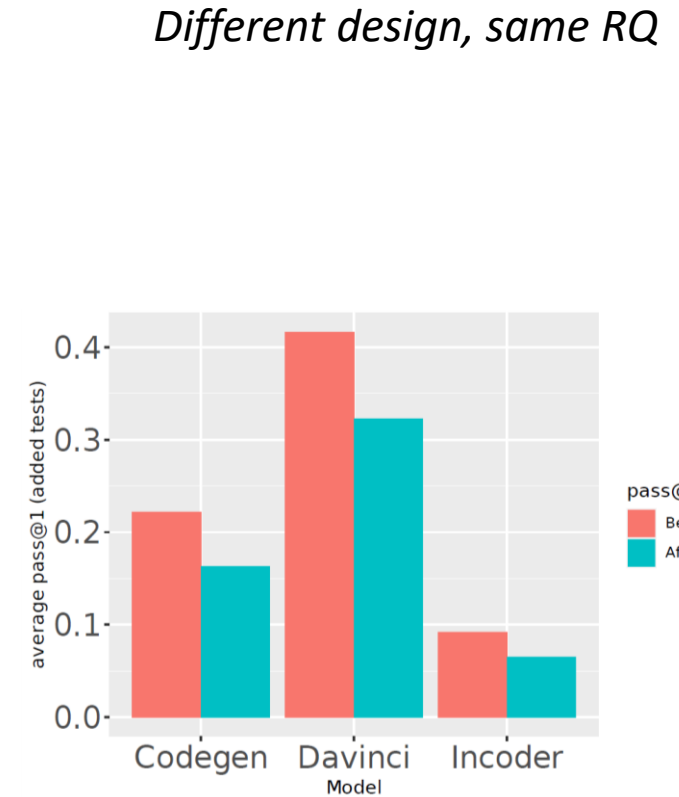
Test@Invocation	Oracle	CodeGen <sub>1</sub>	CodeGen <sub>2</sub>	CodeGen <sub>3</sub>	CodeGen <sub>4</sub>
$T_1@1$	1	1	-7	1	1
$T_2@1$	5	5	-15	5	5
$T_3@1$	7	7	28	7	7
$T_4@1$	12	12	60	12	12

*Judging functional correctness using SRMs*

# Reproduction – Adding more Tests (EvoSuite)

**RQ1:** What is the Java code completion performance of the three code models (with temperature set to 0.2) as judged by functional correctness ? (pass@k metric)

	$S_{1LLM_1}$	...	$S_{kLLM_1}$	...	$S_{1LLM_n}$	...	$S_{kLLM_n}$
$T_1$	$T_1(S_{1LLM_1}, \dots)$	...	$T_1(S_{kLLM_1}, \dots)$	...	$T_1(S_{1LLM_n}, \dots)$	...	$T_1(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...
$T_j$	$T_j(S_{1LLM_1}, \dots)$	...	$T_j(S_{kLLM_1}, \dots)$	...	$T_j(S_{1LLM_n}, \dots)$	...	$T_j(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...
$T_m$	$T_m(S_{1LLM_1}, \dots)$	...	$T_m(S_{kLLM_1}, \dots)$	...	$T_m(S_{1LLM_n}, \dots)$	...	$T_m(S_{kLLM_n}, \dots)$
$E_1$	$E_1(S_{1LLM_1}, \dots)$	...	$E_1(S_{kLLM_1}, \dots)$	...	$E_1(S_{1LLM_n}, \dots)$	...	$E_1(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...
$E_n$	$E_n(S_{1LLM_1}, \dots)$	...	$E_n(S_{kLLM_1}, \dots)$	...	$E_n(S_{1LLM_n}, \dots)$	...	$E_n(S_{kLLM_n}, \dots)$



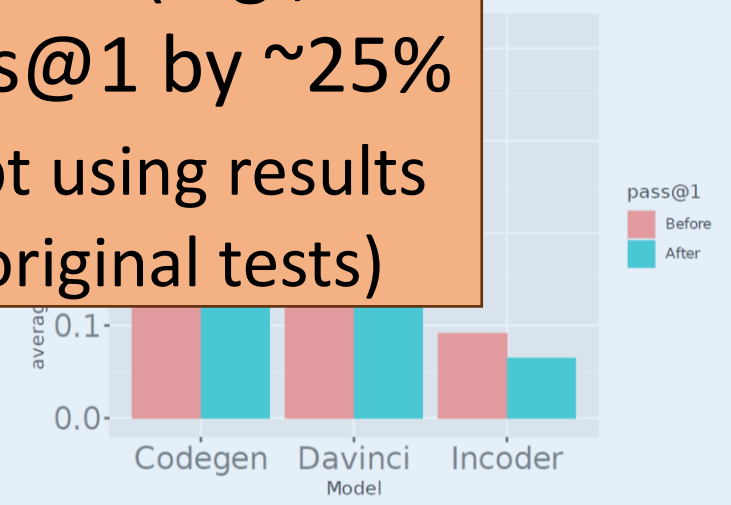
Adding more **tests** to SRMs using *EvoSuite* (Fraser and Arcuri)

# Reproduction – Adding more Tests (EvoSuite)

RQ: What is the Java code completion performance of the three code models (with temperature set to 0.2) as judged by functional correctness ? (pass@k metric)

Aligns with findings of related TDSEs (e.g., *EvalPlus*; Liu et al.) – decrease of pass@1 by ~25%  
 → Reproduction resulted in one LSL script using results from previous replication (no re-run of original tests)

$T_1$	$T_1$							
...	...							
$T_j$	$T_j$							
...	...							
$T_m$	$T_m$							
}	$E_1$	$E_1(S_{1LLM_1}, \dots)$	...	$E_1(S_{kLLM_1}, \dots)$	...	$E_1(S_{1LLM_n}, \dots)$	...	$E_1(S_{kLLM_n}, \dots)$
	...	...	...	...	...	...	...	...
	$E_n$	$E_n(S_{1LLM_1}, \dots)$	...	$E_n(S_{kLLM_1}, \dots)$	...	$E_n(S_{1LLM_n}, \dots)$	...	$E_n(S_{kLLM_n}, \dots)$



Adding more **tests** to SRMs using EvoSuite (Fraser and Arcuri)

# Repurposing – Behavioral Clustering

**RQ2:** *What is the variability in observable functional behavior exhibited by the software components sampled from the code models?*

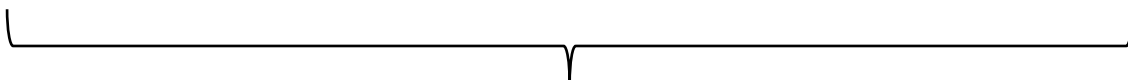
*Same design, different RQ*

	$S_{1LLM_1}$	...	$S_{kLLM_1}$	...	$S_{1LLM_n}$	...	$S_{kLLM_n}$
$T_1$	$T_1(S_{1LLM_1}, \dots)$	...	$T_1(S_{kLLM_1}, \dots)$	...	$T_1(S_{1LLM_n}, \dots)$	...	$T_1(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...
$T_j$	$T_j(S_{1LLM_1}, \dots)$	...	$T_j(S_{kLLM_1}, \dots)$	...	$T_j(S_{1LLM_n}, \dots)$	...	$T_j(S_{kLLM_n}, \dots)$
...	...	...	...	...	...	...	...
$T_m$	$T_m(S_{1LLM_1}, \dots)$	...	$T_m(S_{kLLM_1}, \dots)$	...	$T_m(S_{1LLM_n}, \dots)$	...	$T_m(S_{kLLM_n}, \dots)$



Descriptive statistics of cluster sizes over all functional abstractions (research question RQ 2).

	min	q25	median	q75	max	mean	sd
Total	1	11.00	23.50	49.75	143	32.71	29.25
Codegen	1	3.00	9.00	21.00	80	14.45	16.04
Davinci	1	2.00	4.00	7.00	55	5.83	7.39
InCoder	1	6.00	12.00	24.00	76	17.61	16.05



Cluster implementations (SRM columns)  
by the equivalence of their observed outputs



# Repurposing – Behavioral Clustering

**RQ2:** What is the variability in observable functional behavior exhibited by the software components sampled from the code models?

	$S_{1LLM_1}$	...
$T_1$	$T_1(S_{1LLM_1}, \dots)$	...
...	...	...
$T_j$	$T_j(S_{1LLM_1}, \dots)$	...
...	...	...
$T_m$	$T_m(S_{1LLM_1}, \dots)$	...

**Davinci** provides the least variability in functional behavior (~5.8 clusters on average)  
 → No additional LSL script necessary, since SRMs of previous study were analyzed directly

Cluster implementations (SRM columns)  
by the equivalence of their observed outputs

actions (research question

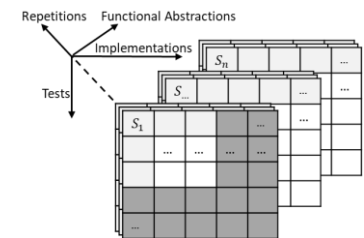
x	mean	sd
	32.71	29.25
	14.45	16.04
	5.83	7.39
	17.61	16.05

# Main Takeaways

- We successfully replicated, reproduced and repurposed a popular TDSE (i.e., benchmarking LLMs for code generation – HumanEval-J)
- Resulting LSL study pipeline scripts
  - encode study designs and assumptions
  - can be shared with others together with SRMs
  - can be repeated by executing it again
  - Increases transparency + fosters open science principles
- Limitations
  - serialization of objects (currently: JSON string representation)
  - equivalence checking in tabular representations

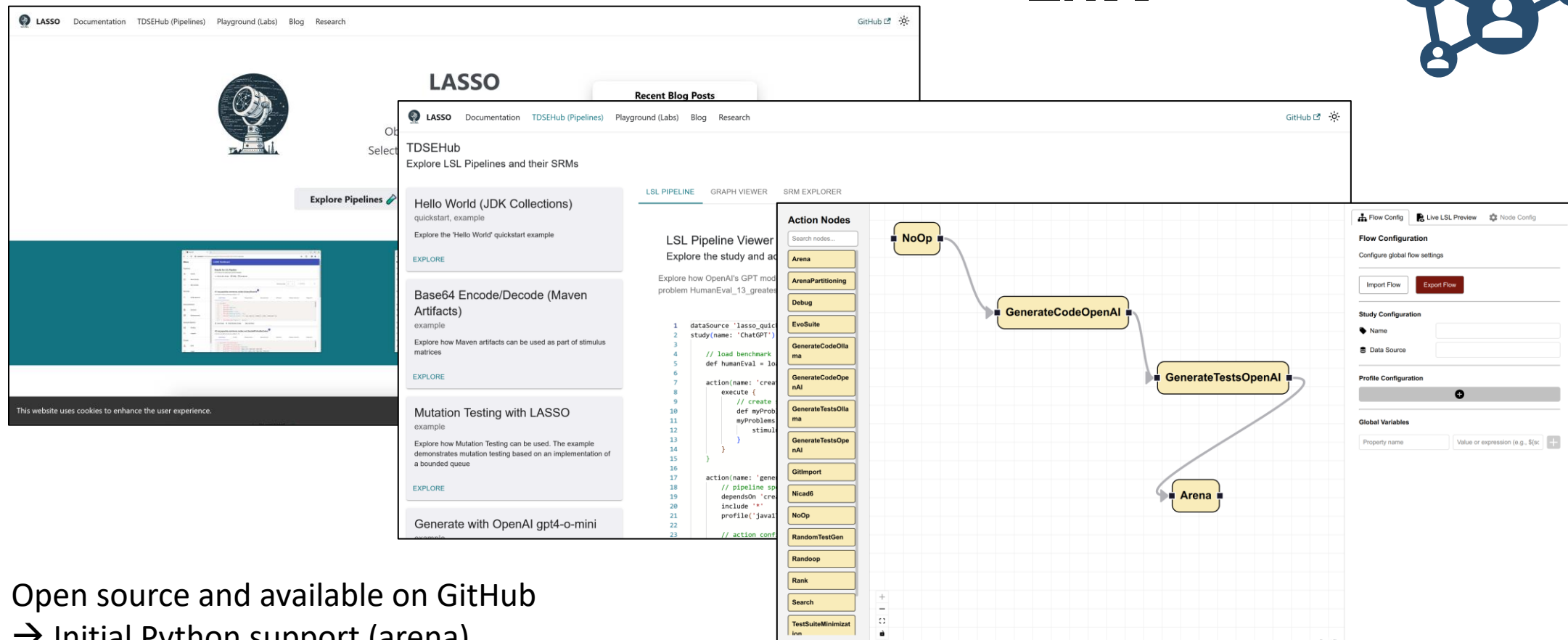
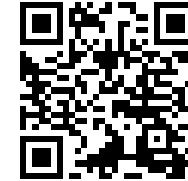


LSL Script



# LASSO Community

- <https://softwareobservatorium.github.io>



The screenshot displays the LASSO web application interface. The top navigation bar includes links for Documentation, TDSEHub (Pipelines), Playground (Labs), Blog, and Research. The main content area is divided into several sections:

- Recent Blog Posts:** A list of articles including "Hello World (JDK Collections)", "Base64 Encode/Decode (Maven Artifacts)", "Mutation Testing with LASSO", and "Generate with OpenAI gpt4-o-mini".
- LSL Pipeline Viewer:** A section for exploring the study and actions, showing a code snippet for a pipeline configuration.
- Action Nodes:** A list of available actions such as Arena, ArenaPartitioning, Debug, EvoSuite, GenerateCodeOpenAI, GenerateCodeOllama, GenerateTestsOpenAI, GenerateTestsOllama, GenerateTestsOpenAI, GitImport, Nicad6, NoOp, RandomTestGen, Randoop, Rank, Search, and TestSuiteMinimization.
- Flow Configuration:** A section for configuring the global flow settings, including Name, Data Source, Profile Configuration, and Global Variables.

The central part of the interface shows a visual representation of a pipeline flow with nodes: NoOp, GenerateCodeOpenAI, GenerateTestsOpenAI, and Arena, connected by arrows.

Open source and available on GitHub  
→ Initial Python support (arena)



# Thank you!

Feel free to reach out 😊



## Promoting open science in test-driven software experiments<sup>†</sup>

Marcus Kessel<sup>\*</sup>, Colin Atkinson

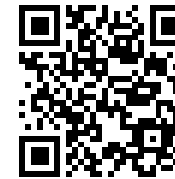
University of Mannheim, 68159 Mannheim, Germany

### ARTICLE INFO

**Keywords:**  
Software  
Engineering  
Empirical  
Experimentation  
Observation  
Behavior  
Reproducibility  
Replication  
Data structures  
Open science  
Large language models  
Machine learning  
Generative artificial intelligence  
Benchmark  
Language-to-code  
Humorful  
Automation  
Measurement

### ABSTRACT

A core principle of open science is the clear, concise and accessible publication of empirical data, including “raw” observational data as well as processed results. However, in empirical software engineering there are no established standards (*de jure* or *de facto*) for representing and “opening” observations collected in test-driven software experiments — that is, experiments involving the execution of software subjects in controlled scenarios. Execution data is therefore usually represented in ad hoc ways, often making it abstract and difficult to access without significant manual effort. In this paper we present new data structures designed to address this issue. We describe the stimuli and responses used to execute their utility, we show how they can be used to support incremental evaluations of AI-based code completion and the incremental expansion of execution data addressing new research questions.



### 1. Introduction

Empirical studies have grown in importance in research over the last few years as leading conferences and journals have placed increased emphasis on research claims to be supported by empirical evidence (e.g., ACM, 2023). Most recently, this trend has evolved to encompass the principles of “Open Science” which aims to facilitate transparency, collaboration, and accessibility in software engineering research by making research data, methods and findings publicly available (Méndez Fernández et al., 2019). Today, open science is actively promoted by most leading software engineering conferences and journals such as EMSE<sup>1</sup>, and has spawned new research communities such as the “Empirical Software Engineering and Measurement” (ESEM) symposium (ESEM, 2023).

The basic motivation for this trend is broadly the same across all scientific and engineering disciplines — namely, making all elements of empirical studies available for other researchers to examine, use and build upon, so that research communities can reach consensus

more easily repeat, replicate and reproduce ACM (ACM, 2023), replication involves the re-execution of an empirical study with the original setup, or a different team reperforming an empirical study and reproduction involves a different team

reperforming the empirical study with a different design and/or setup. The fundamental principles and ingredients of open science are also basically the same across scientific disciplines, even though the terminology and the detailed breakdown of steps may differ. In the social science domain, Misocher et al. (2020) define four fundamental prerequisites for an empirical study to be fully reproducible — “(1) data recoverability, the availability of the involved data and other materials to attempt reproduction of analyses, (2) data usability, the completeness and clarity of data, when available, (3) analytical clarity, the adequacy of published reports for repeating analyses, and (4) agreement of reproduced results with published results”. In their overview of open science in software engineering, Méndez et al. (2020) point out that the first of these, data recoverability, involves the dissemination

<sup>†</sup> Editor: Alexander Serebrenik.

<sup>\*</sup> Corresponding author.

E-mail address: marcus.kessel@uni-mannheim.de (M. Kessel), colin.atkinson@uni-mannheim.de (C. Atkinson).

<sup>1</sup> with its Open Science Initiative (<https://emsejournal.github.io/open-science>)

<https://doi.org/10.1016/j.jss.2024.111971>

Received 2 August 2023; Received in revised form 8 December 2023; Accepted 15 January 2024

Available online 12 March 2024

0164-1212/© 2024 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).