

Towards Observation Lakehouses: Living, Interactive Archives of Software Behavior



Marcus Kessel marcus.kessel@uni-mannheim.de (University of Mannheim, Software Engineering Group)

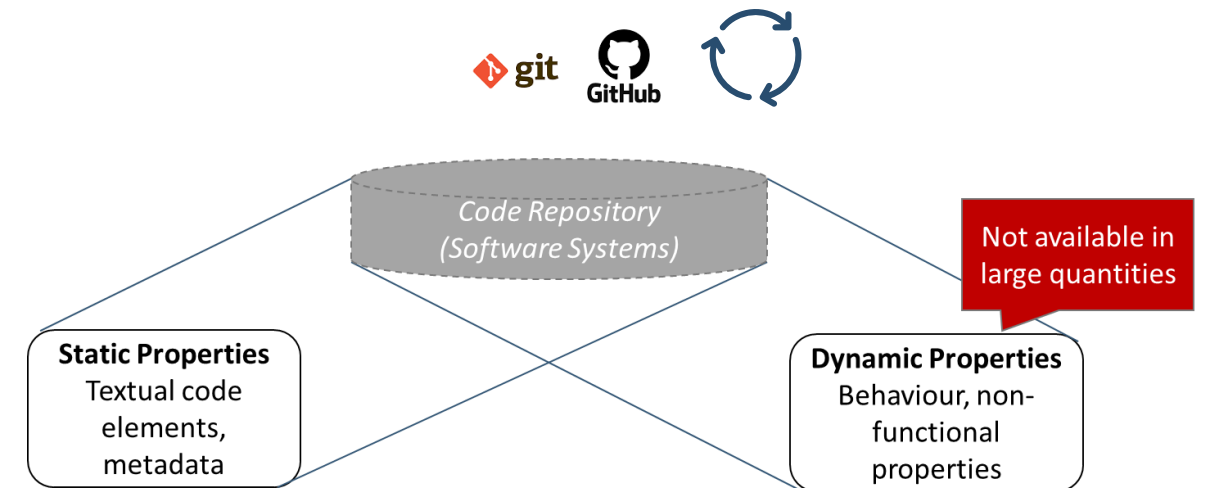


SANER'26 @ Limassol, Cyprus – Tool Demo Track, March 19th, 2026



Motivation – Obtaining Runtime Behavior at Scale

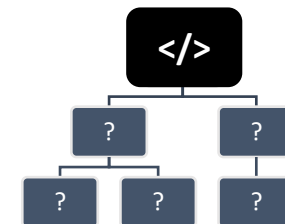
- **Problem:** LLMs (GenAI) typically limited to static artifacts at scale
→ limited availability of runtime behavior
- **Rice's theorem:** semantic (runtime) properties undecidable in general
→ need dynamic observation (software testing)
- **Current Gap**
 - (unit) testing (too little – only pass/fail signal)
 - tracing (too much – expensive)
- **Need: runtime behavior at scale, continuously (!)**
→ living, interactive archives of software behavior **complementing** syntactic artifacts



```
1 import org.junit.Test;
2 import static org.junit.Assert.*;
3
4 public class ArrayStackTest {
5     ...
6
7     @Test
8     public void test_t1() throws Throwable {
9         ArrayStack arrayStack = new ArrayStack();
10        arrayStack.push("hello world");
11        int int0 = arrayStack.size();
12        assertEquals(1, int0);
13    }
14 }
15 ...
16 }
17 }
```



vs.

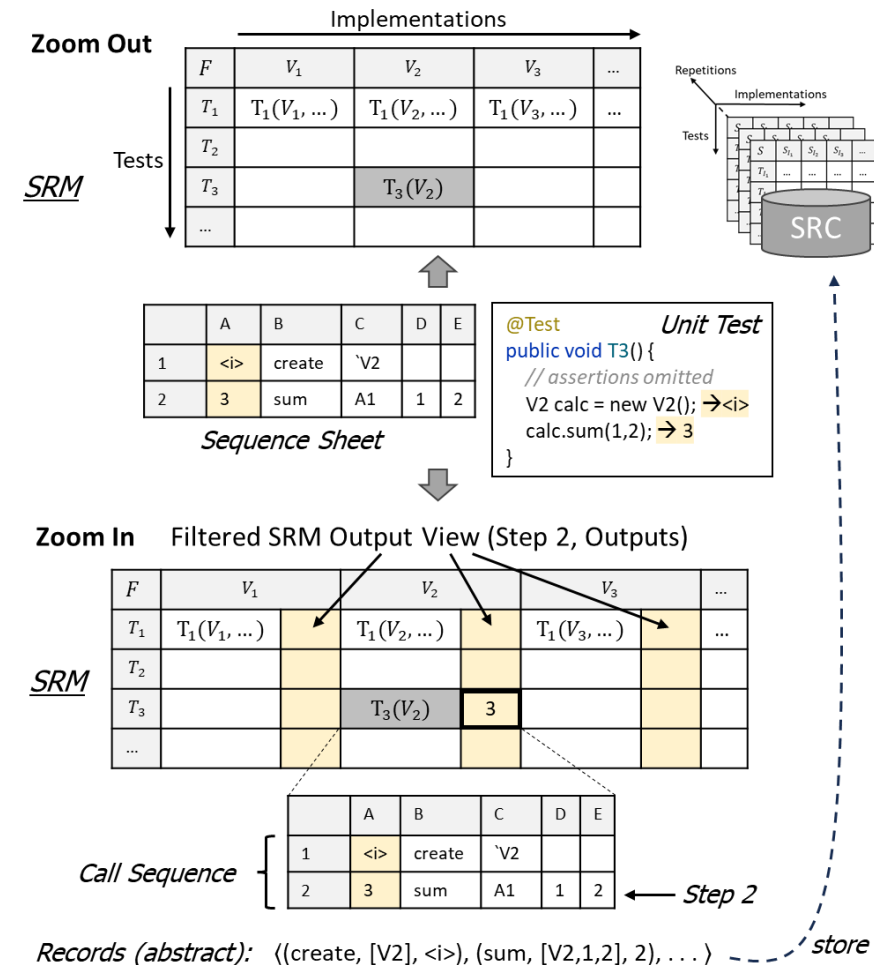


xUnit
... too little
(pass/fail signals)

Tracing
... too much
(detailed tracing)

Continual SRC – Tabular Representations

- **Sequence Sheets (Test Sequences)**
 - stimulus/response encoding
 - invocation-centric: $\langle (\text{op1}, \text{in1}, \text{out1}), (\text{op2}, \text{in2}, \text{out2}), \dots \rangle$
- **SRMs (Stimulus-Response-Matrix)**
 - tests \times implementations comparison matrix
- **SRCs (Stimulus-Response-Cube)**
 - multi-dimensional SRMs across contexts (time, repetitions, envs etc.)
- **“Living” (Continual) property**
 - add rows/columns dynamically



Observation Lakehouse – Architecture

- **Observation Lakehouse: Local store**
 - “lake”: “cheap” storage, flexible schema
 - “warehouse”: SQL analytics (OLAP), ACID
 - Impl.: *Iceberg+Arrow+DuckDB+Parquet*
- **Invocation-centric: append-only table**
 - Stores invocations of test sequences
 - $\langle (op1, in1, out1), (op2, in2, out2), \dots \rangle$
 - Full context awareness (environment, ...)
- **Partitioning** avoids full table scans
 - Benchmarks: *data set / problem*
 - Practitioners: *repository / fully-qualified code unit (e.g., class)*

```

| observations/
| | data/ # partitioned by data_set_id and problem_id
| | | data_set_id=HumanEval/
| | | | problem_id=HumanEval_0_has_close_elements/
| | | | | *.parquet # SRM for problem
| | | | ...
| | | ...
| | ...
| ...
  
```

Table	Contextual Keys	Core Columns & Description
code_ implementations	implementation_id	Stores source code and static metrics (e.g., LOC) for each implementation once.
tests	test_id	Stores the stimulus definition (e.g., a sequence sheet or a mined unit-test method) once.
observations	impl_id, test_id, step_id	The “fact” table storing each invocation step of a call sequence: operation, inputs, outputs and context.

data_set_id	problem_id	implementation_id	test_id	step_id	operation	inputs	output
HumanEval	HumanEval_13_greatest_common_divisor	0dc12996310a4e46-872b-3c3fa7c8636f_original_0	test3()	1	public static long qwen25latest.Problem.great...	144,60	12
HumanEval	HumanEval_13_greatest_common_divisor	e4fba3ef01b641d2-8caf-f99f58685f19_original_0	test3()	1	public static long qwen25latest.Problem.great...	144,60	12
HumanEval	HumanEval_13_greatest_common_divisor	f60dd7ce-aa18-407b-913b-6023f9021d0c_original_0	test3()	1	public static long llama31latest.Problem.great...	144,60	12

append-only table

	test_id	step_id	01e02ed1-29a7-454f-8aac-1e30519d9291_original_0	02041a08-238d-4236-9e1f-7d8ce1e68521_original_0	07719ffb-635a-4d68-9596-81a3396e3bfc_original_0
0	gemma3latest0_testSumEmptyArray()	0	SCUT@Problem@0	SCUT@Problem@0	SCUT@Problem@0
1	gemma3latest0_testSumEmptyArray()	1	[]	[]	[]
2	gemma3latest0_testSumEmptyArray()	2	0	0	0
3	gemma3latest0_testSumLargeArray()	0	SCUT@Problem@0	SCUT@Problem@0	SCUT@Problem@0
4	gemma3latest0_testSumLargeArray()	1	[null]	[null]	[null]
...
860	test0()	1	6	6	6
861	test1()	0	SCUT@Problem@0	SCUT@Problem@0	SCUT@Problem@0
862	test1()	1	50	50	50
863	test2()	0	SCUT@Problem@0	SCUT@Problem@0	SCUT@Problem@0
864	test2()	1	3	3	3

SRM view reconstructed (via SQL)

```
SELECT * FROM ...
```

Evaluation / Performance

Dataset



- **~8.6M observations** (= rows)
- **509 coding problems** (from *Cassano et al.'s* MultiPL-E benchmarks)
- **13.384 implementations** (N Java implementations for each problem)
- **95.154 tests** (a combination of human-written and LLM-generated tests)
- obtained with LASSO



Cassano, Federico, et al. "Multipl-e: A scalable and polyglot approach to benchmarking neural code generation." *IEEE Transactions on Software Engineering* 49.7 (2023): 3675-3691.

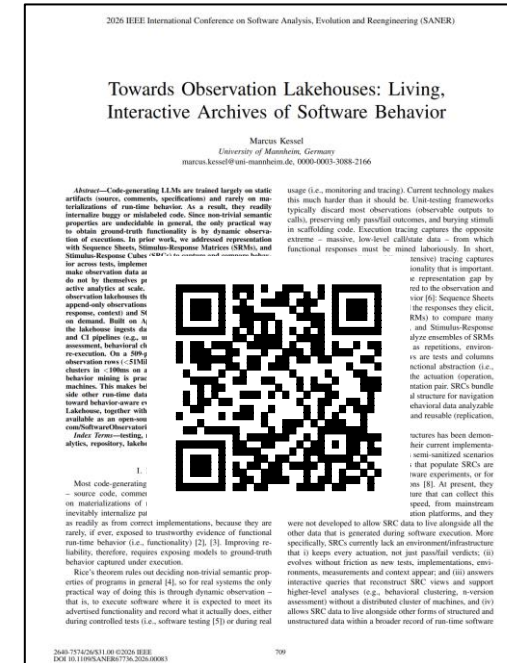
Performance (single laptop)



- **Ingestion**
 - Bulk import (≈ 7 secs)
 - Rate: ≈ 155 K records/second
 - Denormalization strategy for performance
 - Duplicate data deliberately (principles)
- **Storage: 51MiB (observations)**
 - high compression via Parquet format - *Run-length encoding (RLE)* etc.
- **Queries** (\bar{x} , over all problems): < 100 ms
 - **SRM view reconstruction: ≈ 52 ms**
 - **Behavioral clustering: ≈ 29 ms**
 - **Full join (three tables): ≈ 90 ms**

Conclusion/Takeaways

- **Practical realization** of obtaining runtime behavior at scale using a lakehouse stack
 - Initial performance promising (scalability)
 - Data-driven n-version assessment, behavioral clustering, consensus/historical oracles
- **Challenges**
 - Stimulus/response serialization
 - Need for equivalence functions (lossless?)
 - Query-dependent partitioning strategies
- **Ongoing, Future Work**
 - Integration with GenAI, Federation of lakehouses, ...
 - *Mining invocations from unit tests*: Test-driver extensions (Java/JUnit5, python/pytest)
 - *“Replaying”* git repositories
 - *Behavioral Co-Versioning*: behavioral archive alongside git (accepted at FSE’26 IVR track)



Preprint: <https://arxiv.org/abs/2512.02795>

Reproduction package:
<https://github.com/SoftwareObservatorium/observation-lakehouse>





Thank you!



<https://www.wim.uni-mannheim.de/atkinson/team/dr-marcus-kessel/>

<https://www.linkedin.com/in/marcus-kessel-307609312/>