# BenchPAL

# Systematic Benchmarking Test Case Generation

### Motivation

- Proof assistants are critical tools in formal verification used to ensure the correctness of mathematical theorems, software, and hardware systems.
- Lean, Idris, Agda, and Rocq are functional programming languages commonly used as proof assistants.
- Current user experience suggests that there are many low-level deficiencies.
- BenchPAL systematically tests these proof assistant languages (PALs) to highlight areas for improving their performance.









# System Design

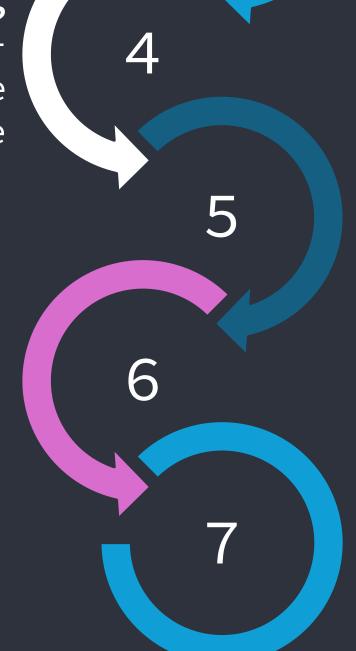
Test generator Uses MHPG to dynamically generate scalable test modules based on size input.

### Tracking

Tests are typechecked, and both time and memory usage are recorded.

### Visualization

Graphs are plotted using Matplotlib from JSON file and served in Flask backend.



### CI Pipeline

Allows users to run the system utilizing Docker and Vercel.

### Translators

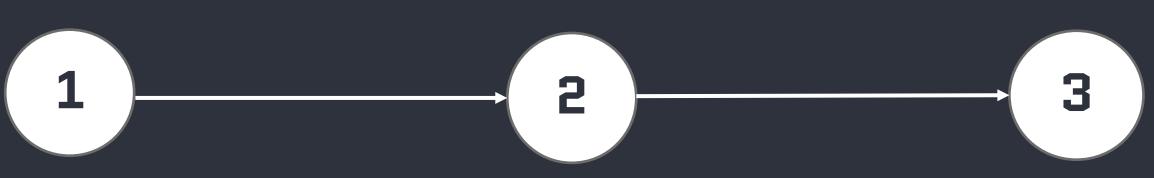
Translates the generated tests into each PAL.

### Calculations

Startup times subtracted from runtimes, and log values are calculated if log interval selected.

### Future development Test suite can be expanded by adding tests to Tests module (Tests.hs).

### Structure of BenchPAL



# [GitHub Actions]

The user interacts with our project through a CI implemented in GitHub Actions or locally.

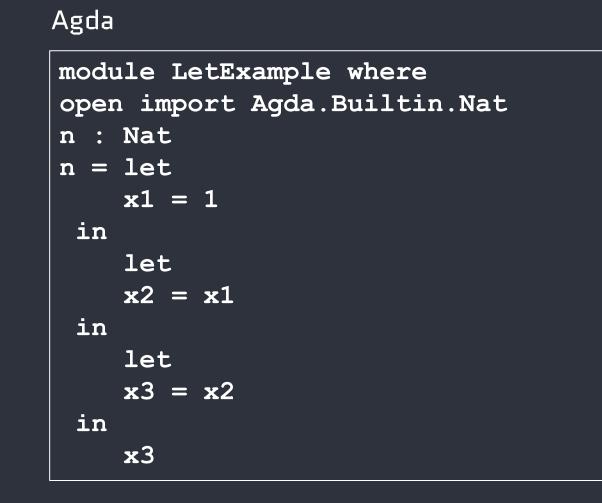
### Automated Code Generator

Automated code generator creates tests of increasing size and translates them to the four proof assistant languages.

### Webpage Generation

Test complexities visualized using Python's matplotlib in a Flask backend and hosted on Vercel to display results.

### Example: Let Bindings Test Case

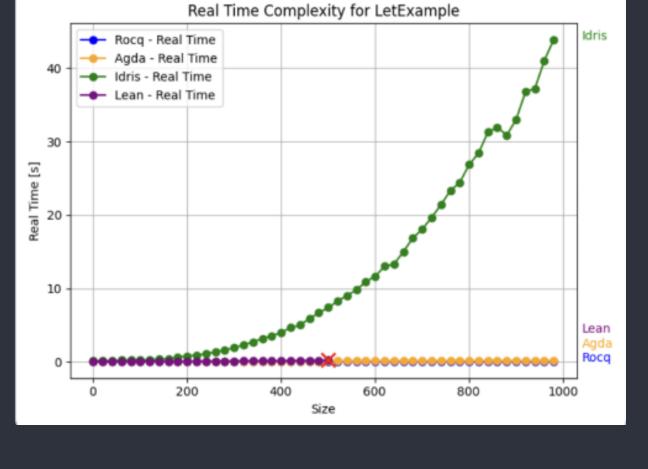


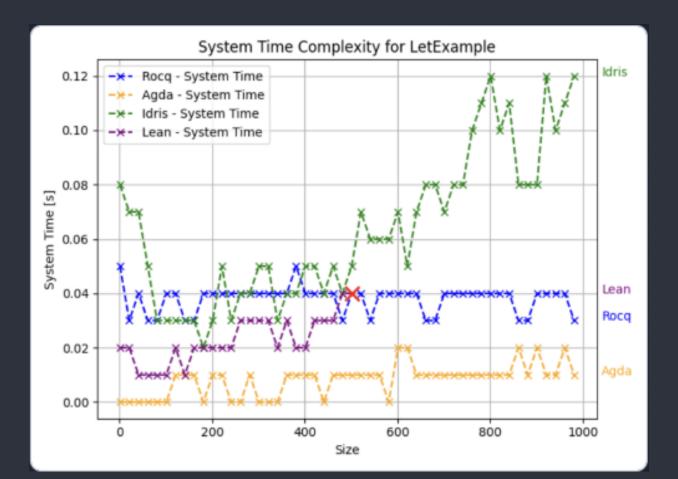
Lean def n : Nat := let x1 := 1 let x2 := x1 let x3 := x2

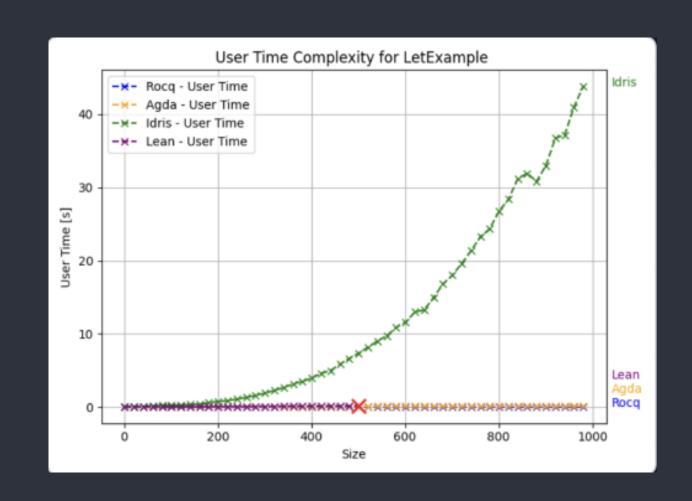
#### module Main n : Nat n = let x1 = 1 in let x2 = x1 in let x3 = x2 in main : IO() main = putStrLn "" Rocq Module LetExample. Definition n : nat := let x1 := 1 in

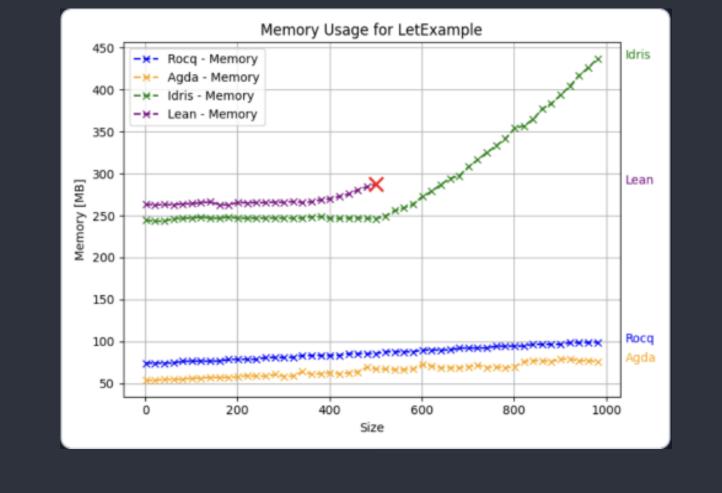
let x2 := x1 in let x3 := x2 in **x**3. End LetExample.

### Results









# Conclusions **Worst Memory** Best time Performance ■ Agda ■ Lean ■ Idris ■ Rocq

Best Time: Rocq Worst Time: Agda Best Memory: Rocq Worst Memory: Lean

# Technology Stack



















### Team

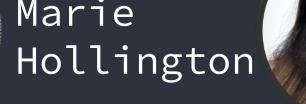




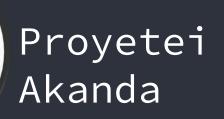
















## Supervisors

Dr. Jacques Carette

Reed Mullanix