# AI-Assisted Design Space Analysis of High-Performance Arm Processors

18th Nov 2024

**Joseph Moore**
Tom Deakin
Simon McIntosh-Smith

University of Bristol High Performance Computing Group
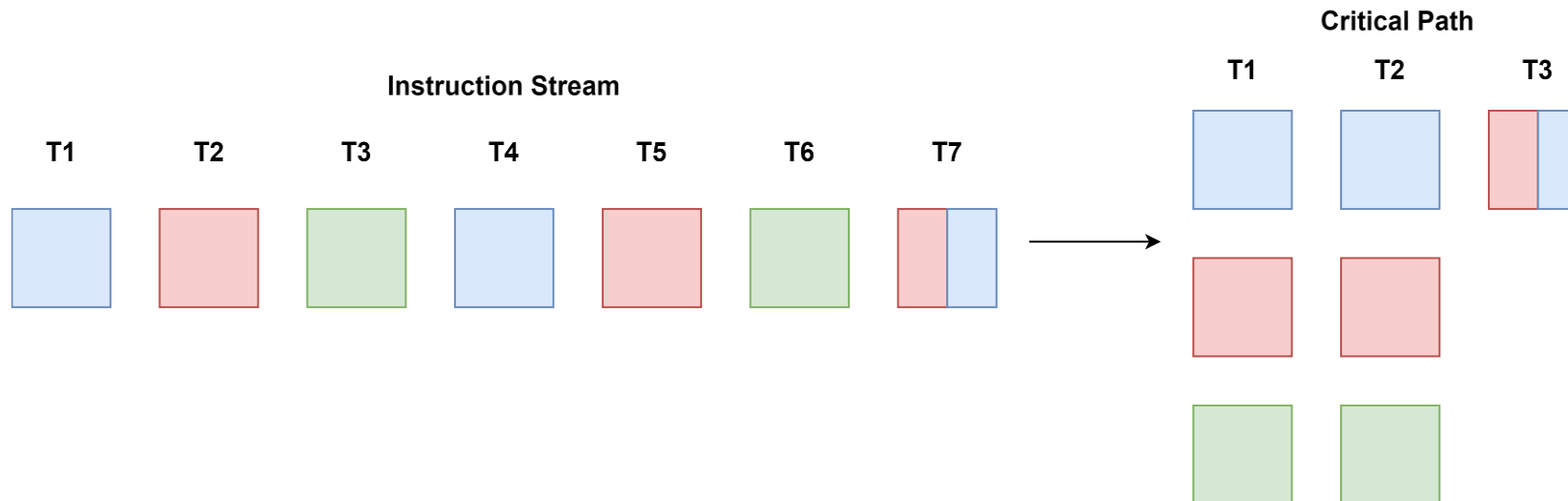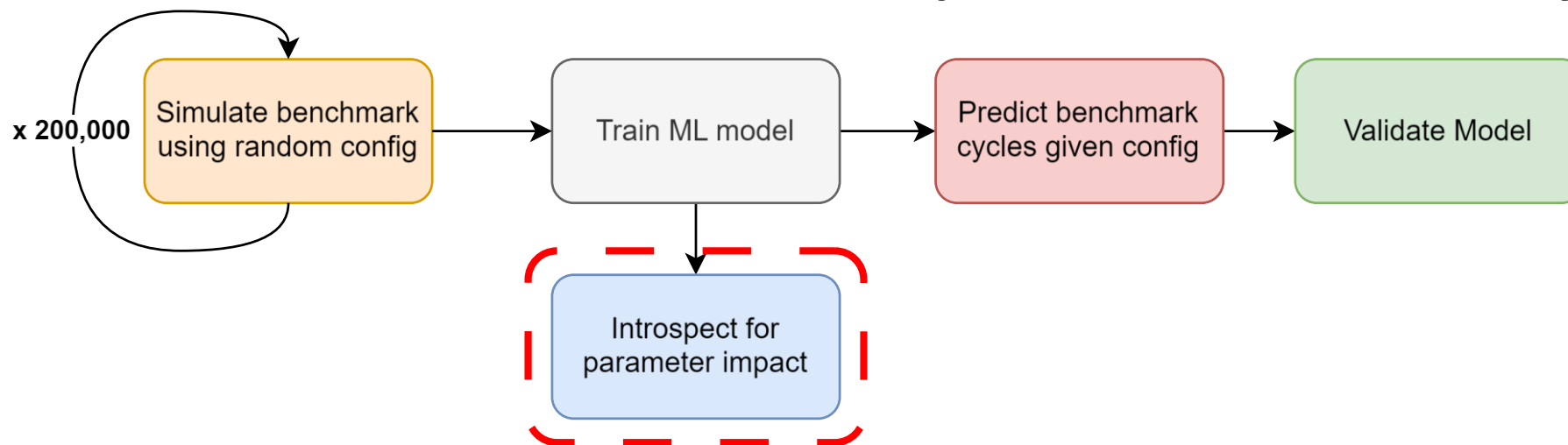
1

# What's the limit of a CPU?

- Consider the trace of executed instructions
  - Critical path = longest chain of data dependent instructions
- Cycles on "perfect" CPU = latency of critical path
  - Bottleneck is the program!
- How can hardware converge to this?

**Instruction Stream**

T1   T2   T3   T4   T5   T6   T7

**Critical Path**

T1   T2   T3

# What this work does

- **Quantifies impact** of single-core bottlenecks

- **Predicts no. cycles** in known HPC applications for different hardware configurations

- Does so through **Machine Learning**

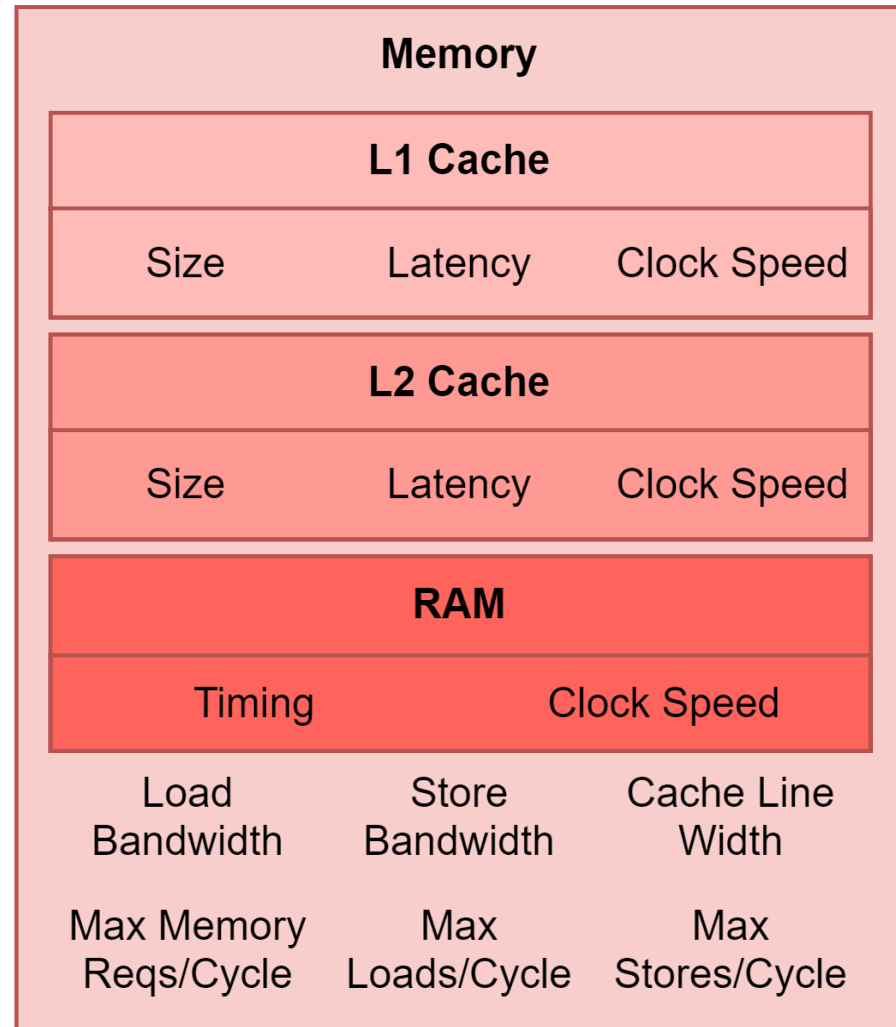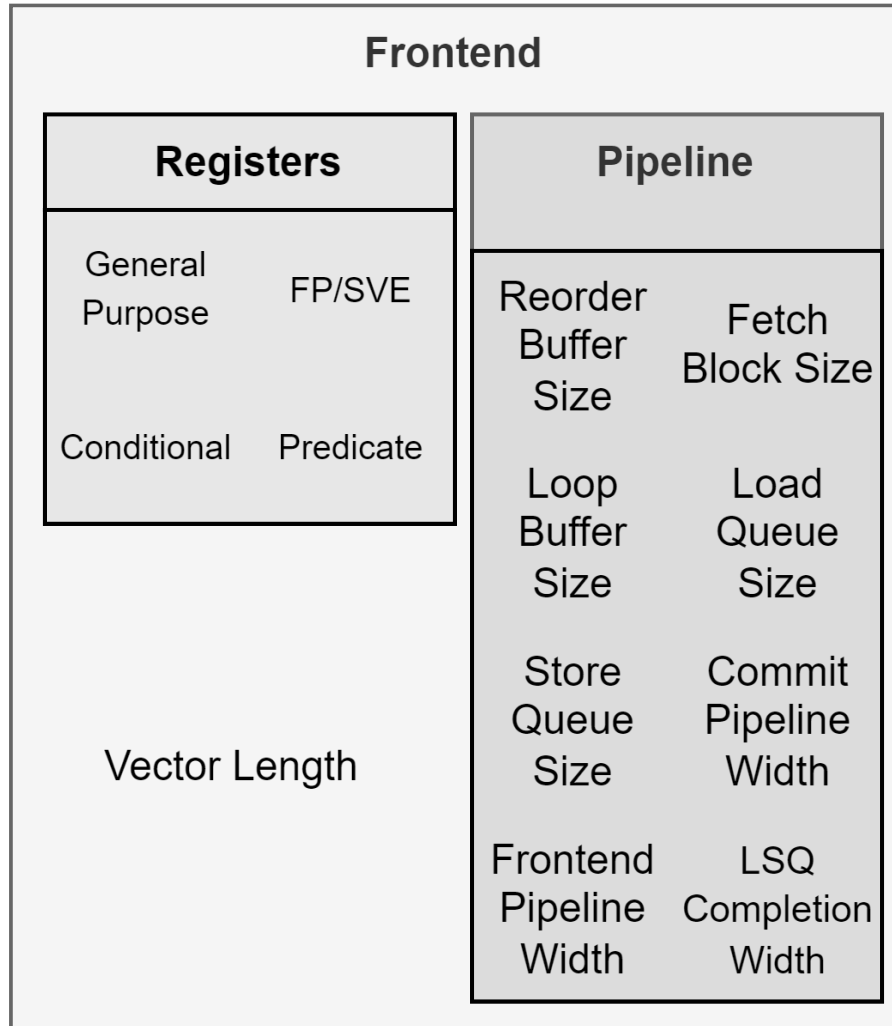- Learn what the model learns – **how do parameters influence cycles**?

# Related work

- 2006-2007 "Golden Age"
  - P.J.Joseph et al, Lee et al, Dubach et al
- Lots of work using traditional ML/AI for parameter searches on few parameters
- Significant jumps in computer architecture since
  - Vectors i.e. Scalable Vector Extension are now commonly used!
- Most work since is focused or models power, space etc.
  - Gap in updated, broad view of core architecture

# SimEng – Our Simulation Framework

- Cycle-approximate Out-Of-Order CPU simulator

- Allows simulation of every stage of the pipeline

- SST integration for memory model

- ~1 MIPS on moderate hardware

- Easy to use - Simple YAML to define CPU properties

https://uob-hpc.github.io/SimEng

# What we are modelling

## Frontend

### Registers

| General Purpose | FP/SVE |
|---|---|
| Conditional | Predicate |

Vector Length

### Pipeline

| Reorder Buffer Size | Fetch Block Size |
|---|---|
| Loop Buffer Size | Load Queue Size |
| Store Queue Size | Commit Pipeline Width |
| Frontend Pipeline Width | LSQ Completion Width |

## Memory

### L1 Cache

| Size | Latency | Clock Speed |
|---|---|---|

### L2 Cache

| Size | Latency | Clock Speed |
|---|---|---|

### RAM

| Timing | Clock Speed |
|---|---|

| Load Bandwidth | Store Bandwidth | Cache Line Width |
|---|---|---|
| Max Memory Reqs/Cycle | Max Loads/Cycle | Max Stores/Cycle |

# What we are **not** modelling

- Reservation Stations
- Execution Units
- No. Cores (just 1)
- Instruction Cache
- L3 Cache – just L1+L2+RAM
- Instruction Set Architecture – fixed to Armv8.4-a+sve

# Benchmarks used

- STREAM – Memory Bound
  - Sustained memory bandwidth benchmark
- MiniBude – Compute Bound
  - Drug Screening Mini-app
- TeaLeaf – Memory Bound
  - Heat Conduction Mini-app
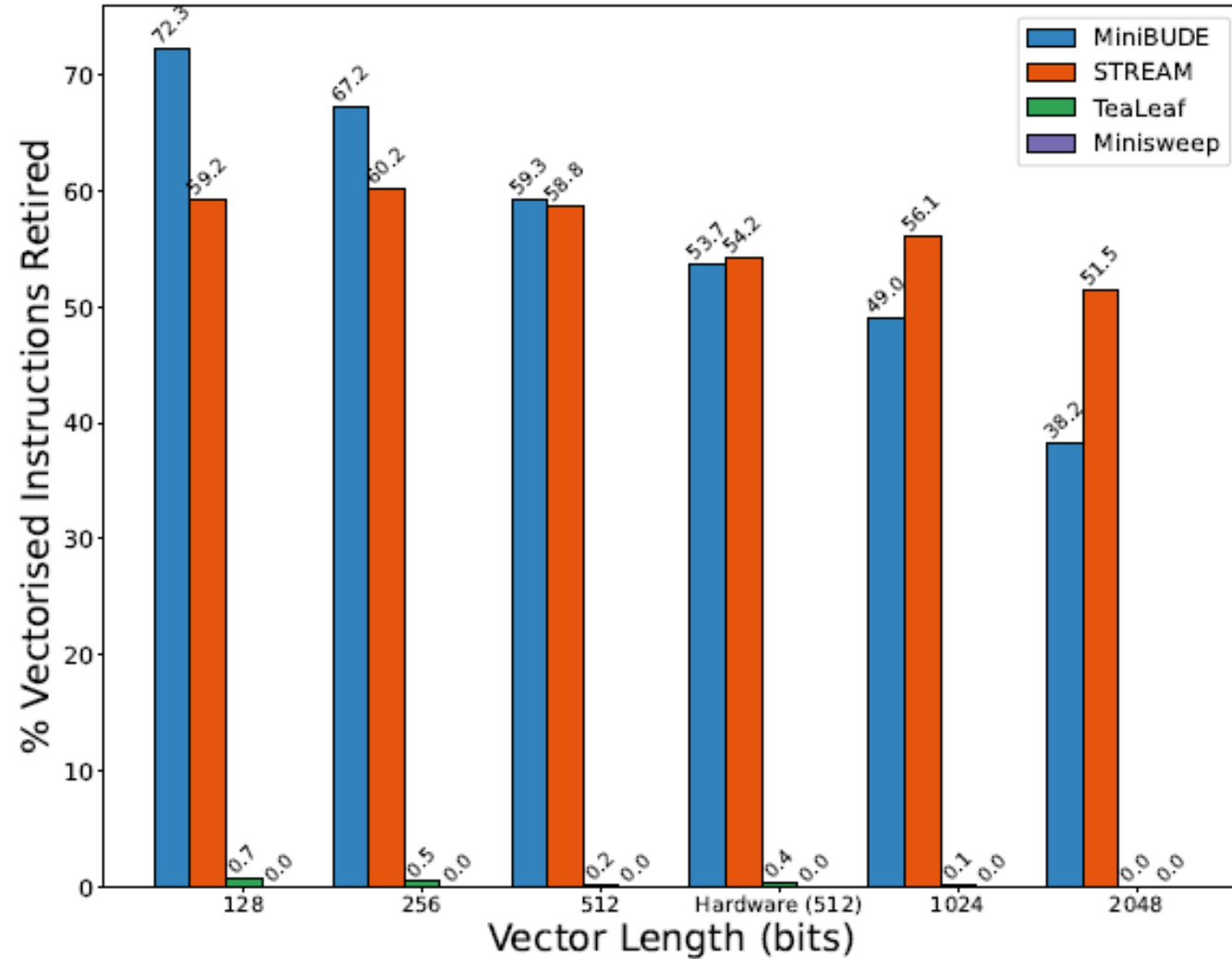- MiniSweep – Compute bound for single core
  - Sn Radiation Transport Mini-app

# Remarks on the benchmarks

- All problems *mostly* fit into L1 or L2 cache (larger takes too much time)
  - For example, STREAM = ~600KiB
- All compiled with Arm Compiler for Linux v23.04.1
  - Compiled statically with –O3, OpenMP (single threaded), and no MPI
  - SVE Vector Length set to "scalable"

SIMULATED SINGLE-CORE CYCLES COMPARED TO HARDWARE CYCLES ON MARVELL'S THUNDERX2 FOR OUR CHOSEN APPLICATIONS IN SIMENG WITH SST
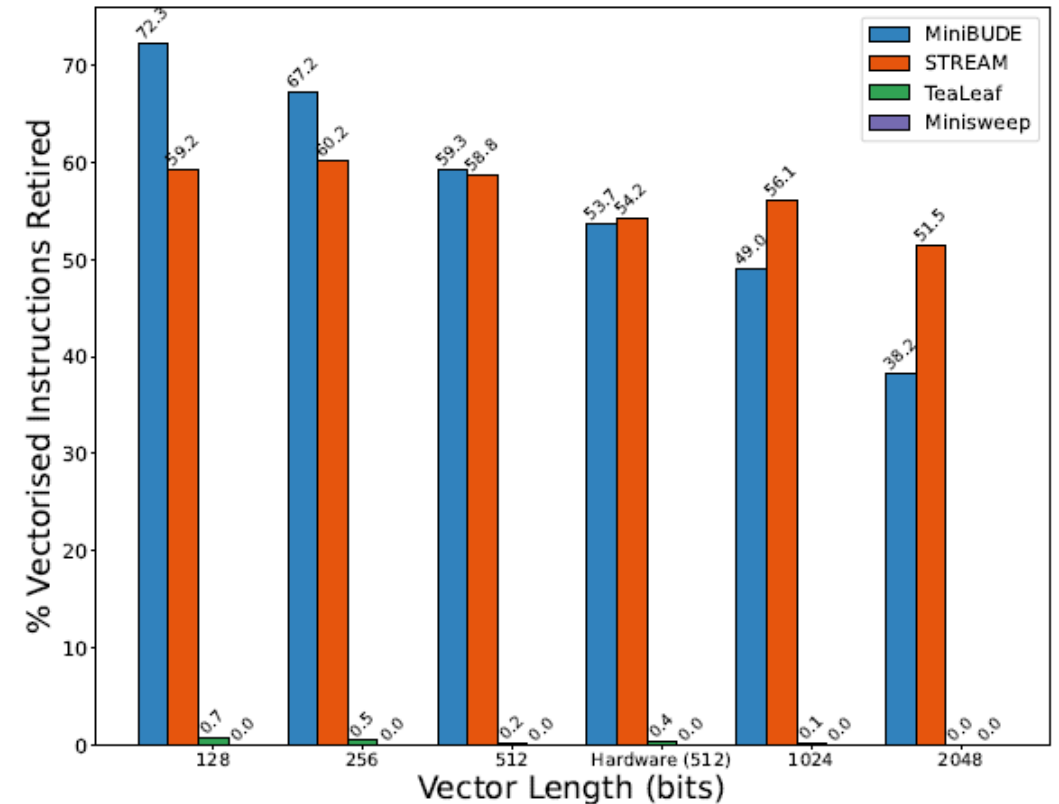
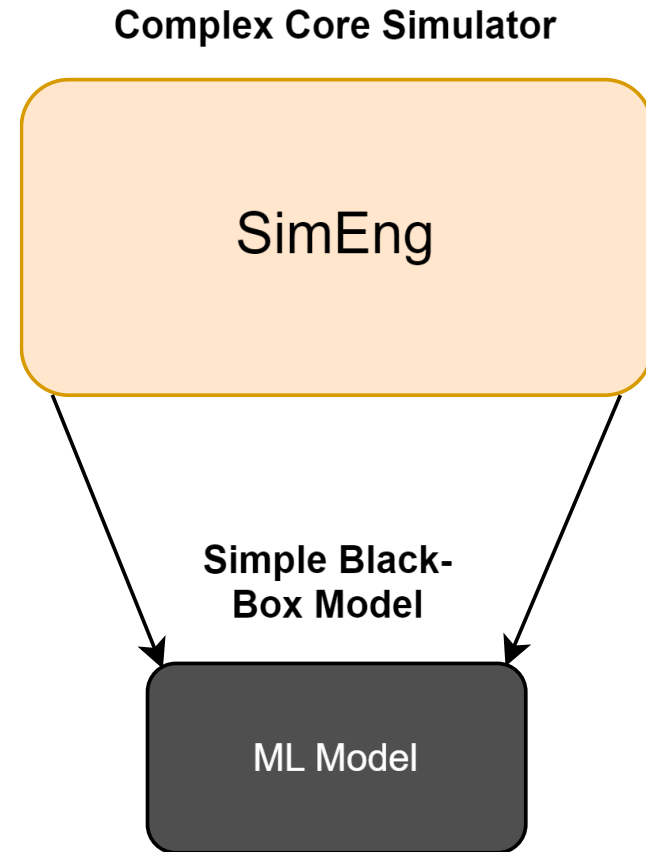| | Simulated Cycles | Hardware Cycles | % Difference |
|---|---|---|---|
| STREAM | 25,078,088 | 26,665,221 | 5.95% |
| MiniBude | 42,436,227 | 48,778,524 | 13.05% |
| TeaLeaf | 19,966,725 | 14,607,184 | 36.69% |
| MiniSweep | 6,529,912 | 10,374,617 | 37.05% |

# Code Vectorisation

# Poor Vectorisation?

- Compiler dependent, not the fault of the hardware!

- Some discrepancies between simulation vs hardware counting

- Huge performance implications

- Not the fault of *-march* flags etc

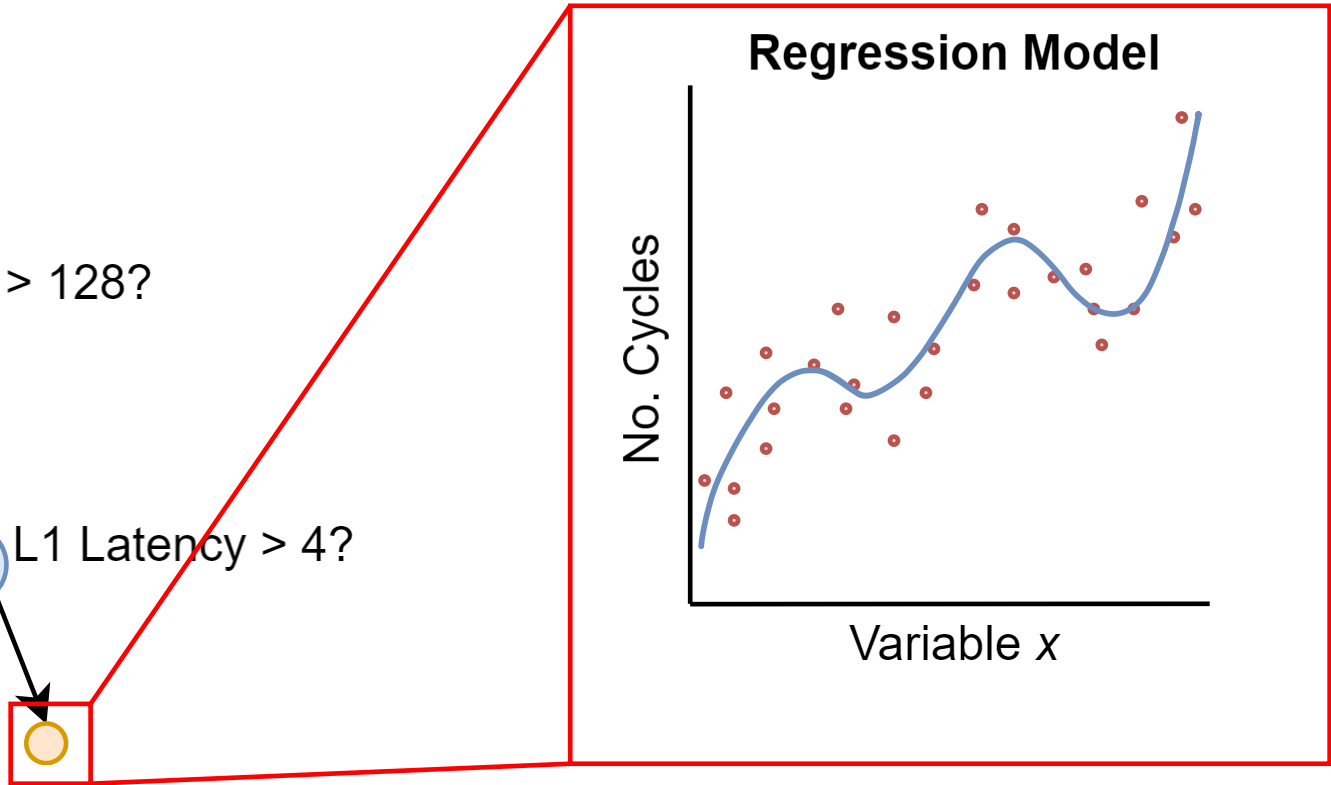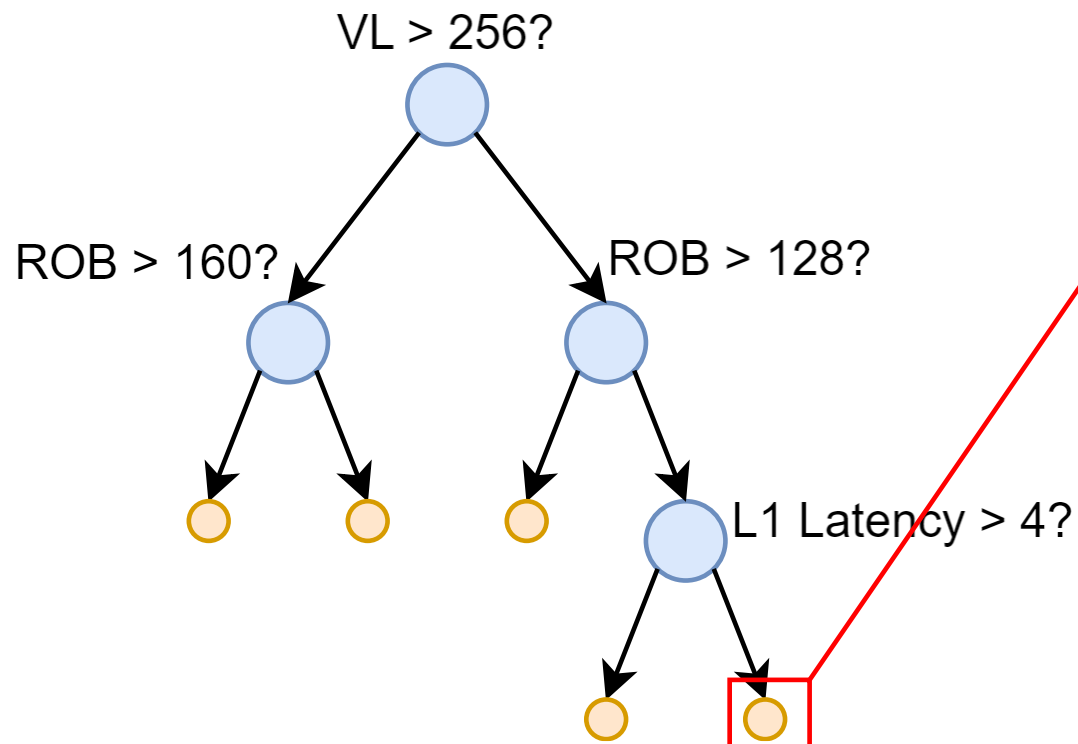- Interesting to consider both well vs poorly vectorized performance

# Machine Learning Model

- Surrogate model – map simulation to ML
- Model significantly faster but more constrained
- Lots of high dimensional data
- Predicting numerical output – regression
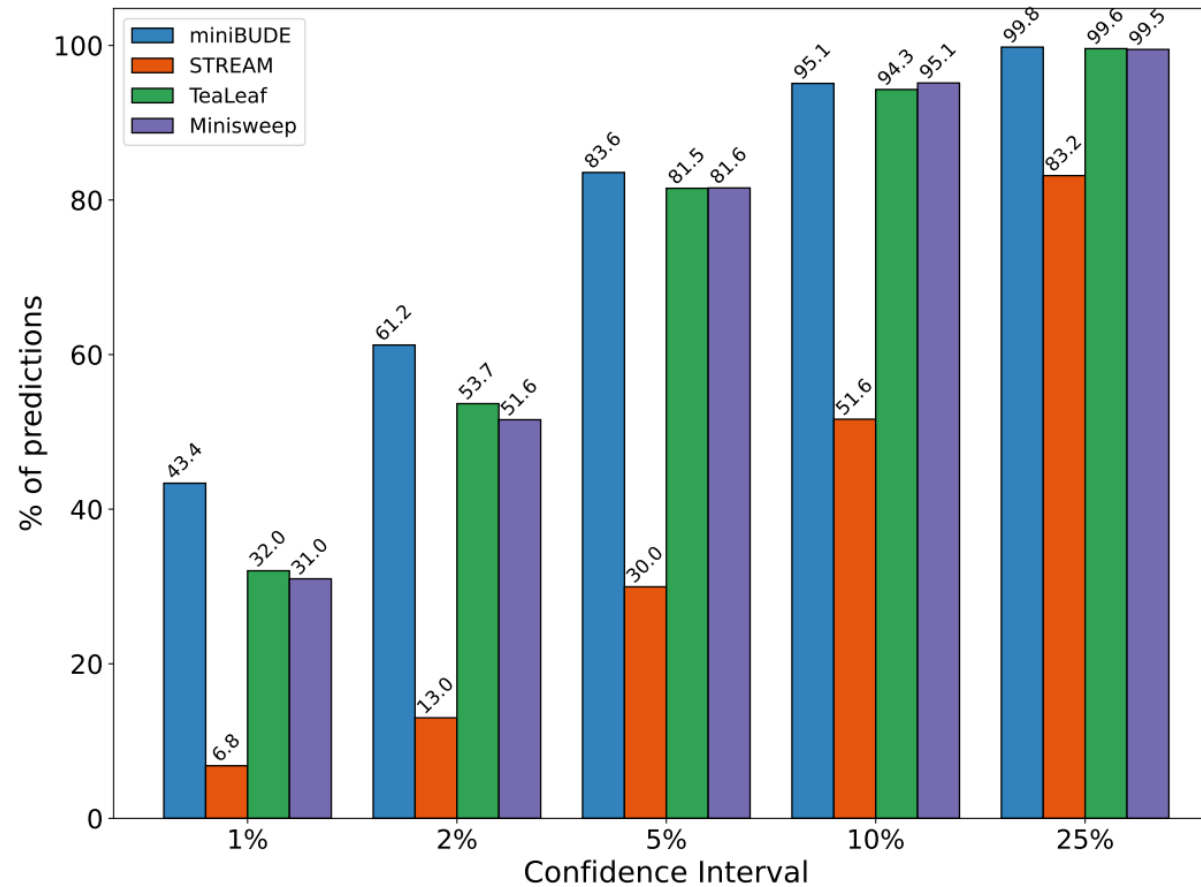- Interested in learned data, not the usage

**Complex Core Simulator**

SimEng

**Simple Black-Box Model**

ML Model

# Decision Tree Regressor

VL > 256?

ROB > 160?

ROB > 128?

L1 Latency > 4?

**Regression Model**

No. Cycles

Variable *x*

# Training the model

- 180,006 valid data entries
- Data sampled uniformly at random
- Collected across 10 Marvell Thunder-X2 nodes across ~3 days
- "Data" is runtime statistics for all applications + config
- 80/20 Train/Test split
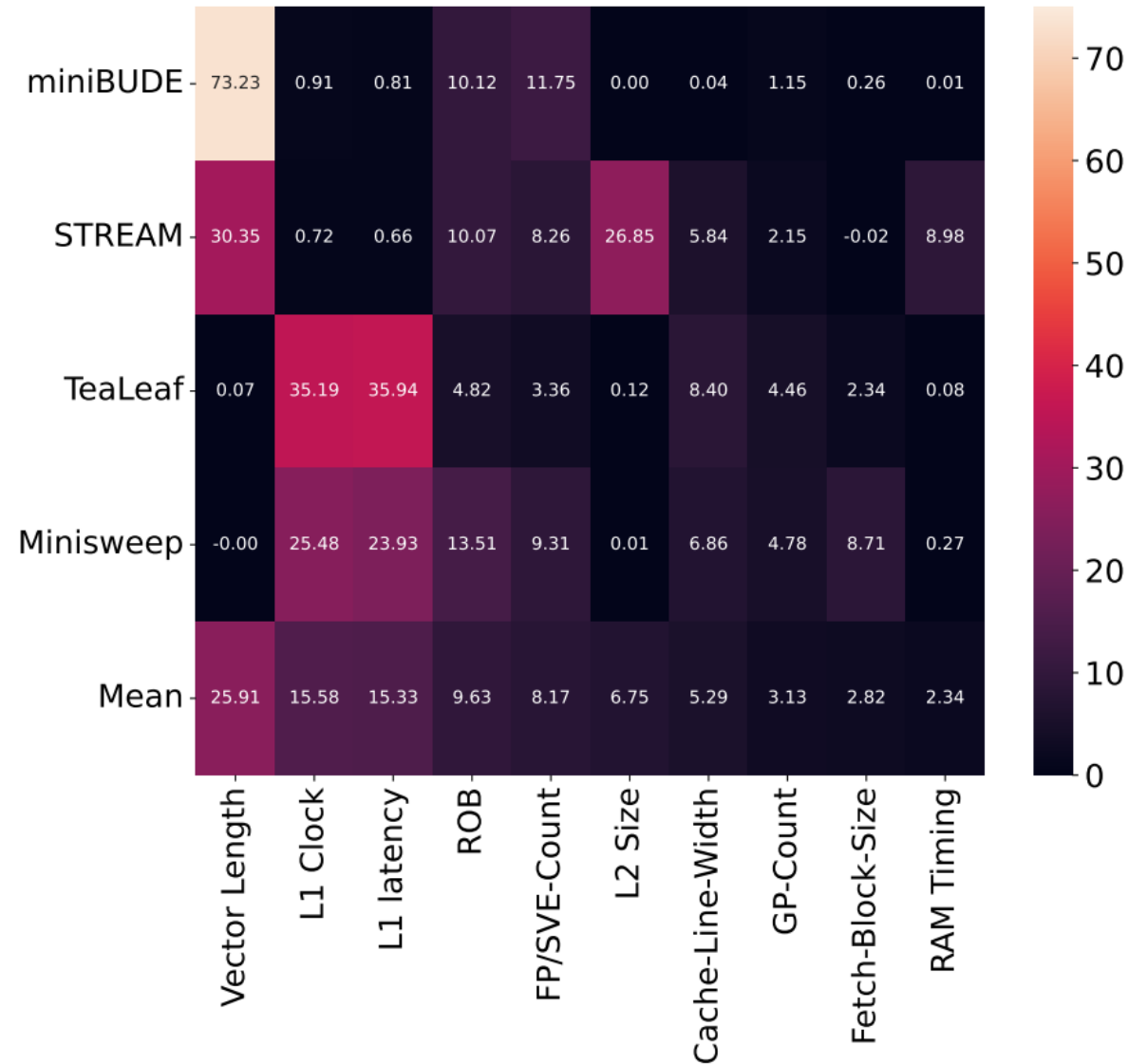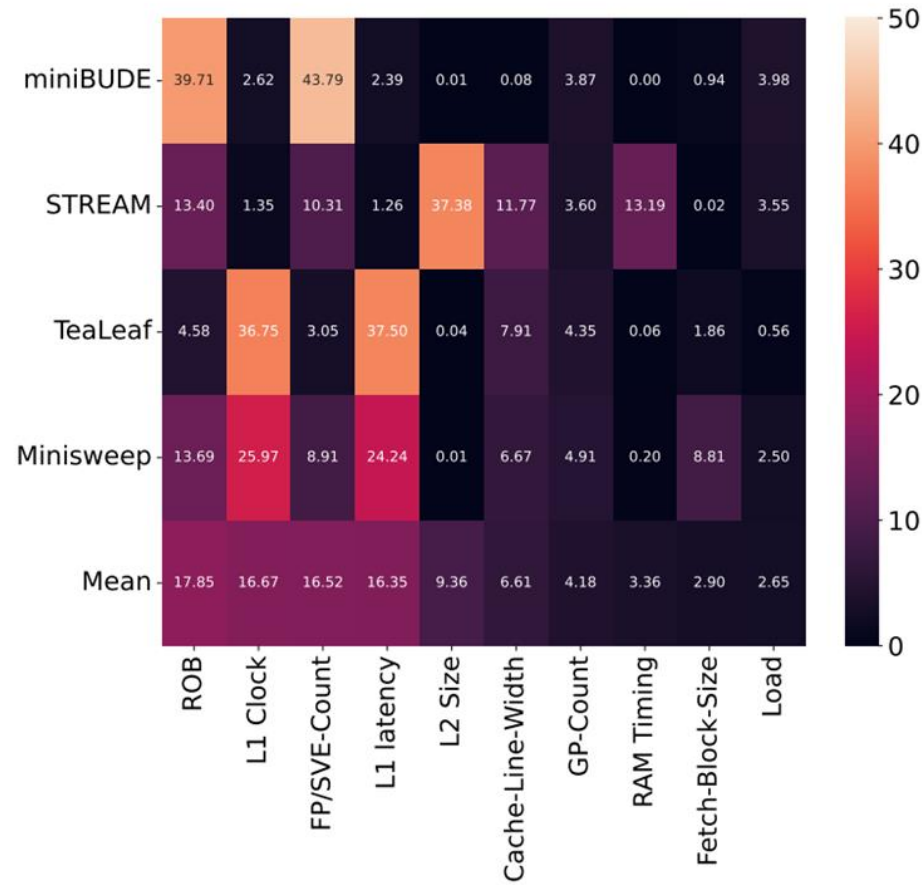- One tree per application

# Model Validation

# Metric of Importance

- Permutation Feature Importance
- Shuffle values of each column and predict
- Measure mean absolute error
- More error caused = more important feature
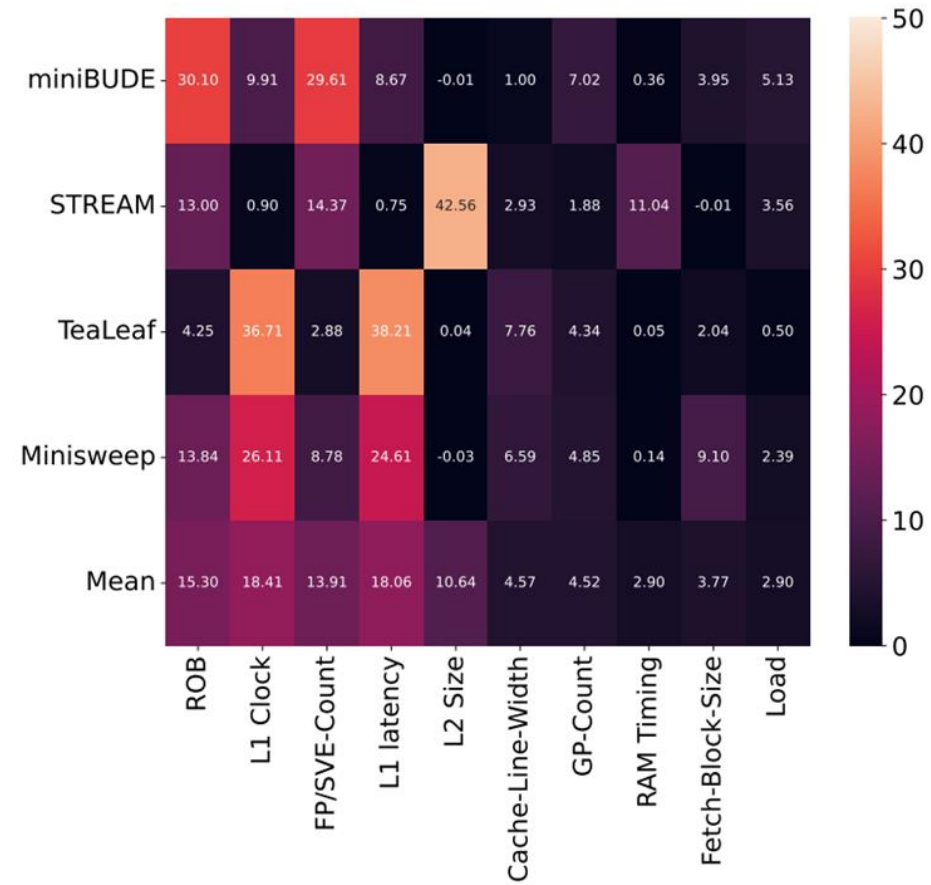- Feature importance = percentage of summed error across all features

# Feature Importance

# Feature Importance (Fixed Vector Length)
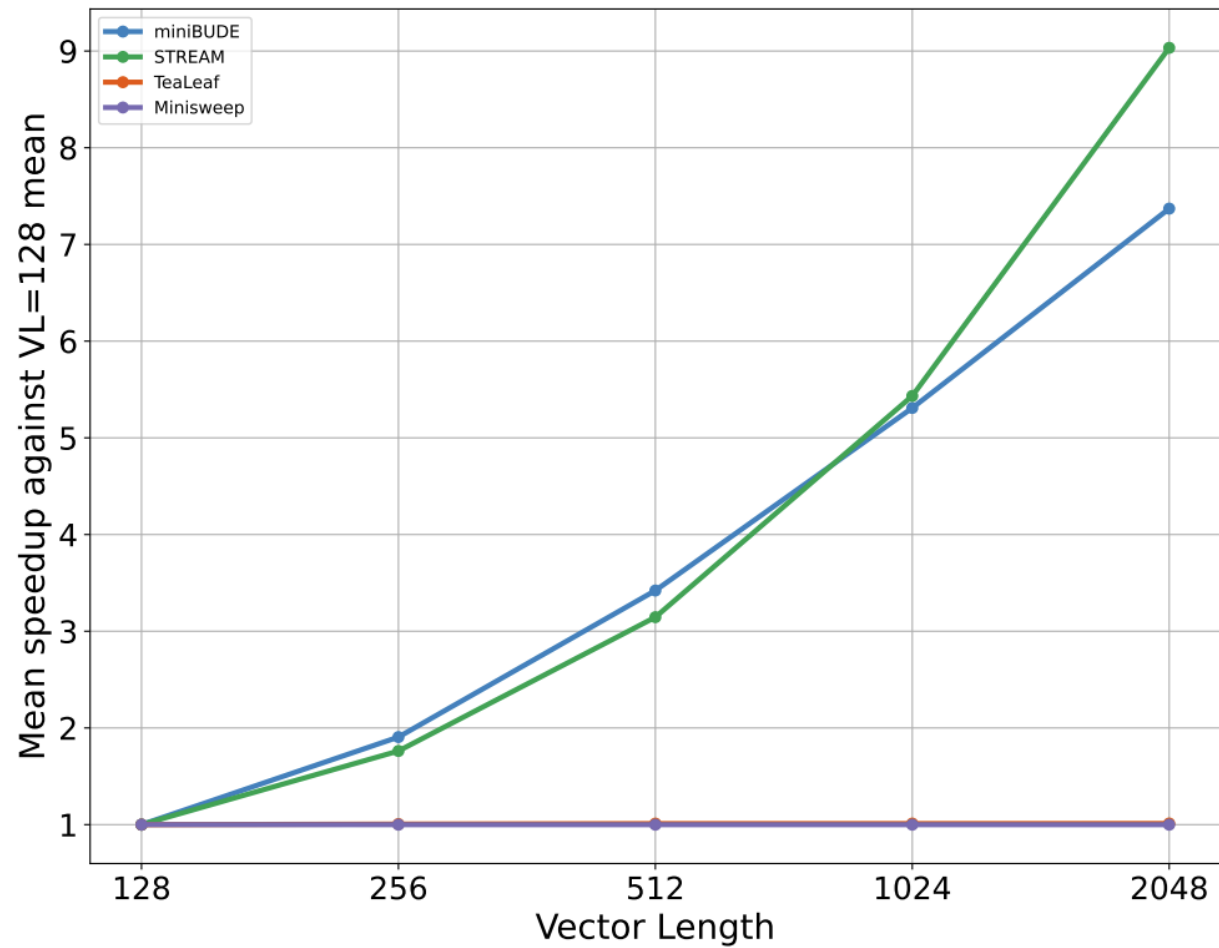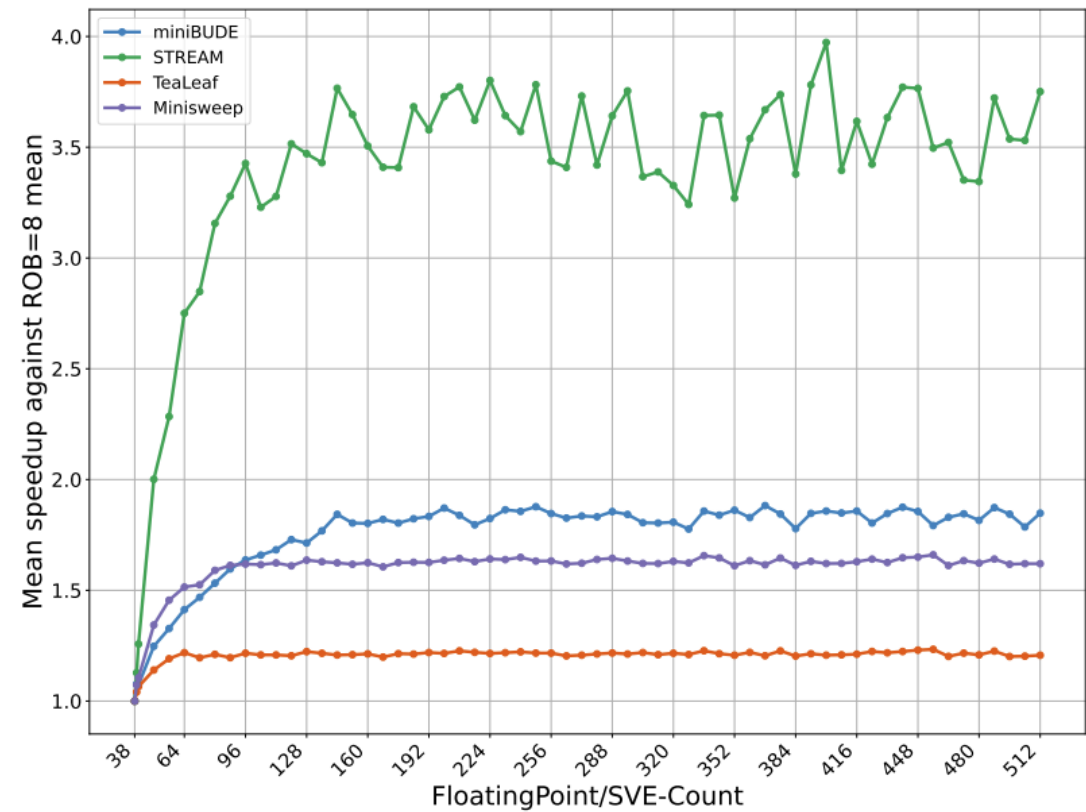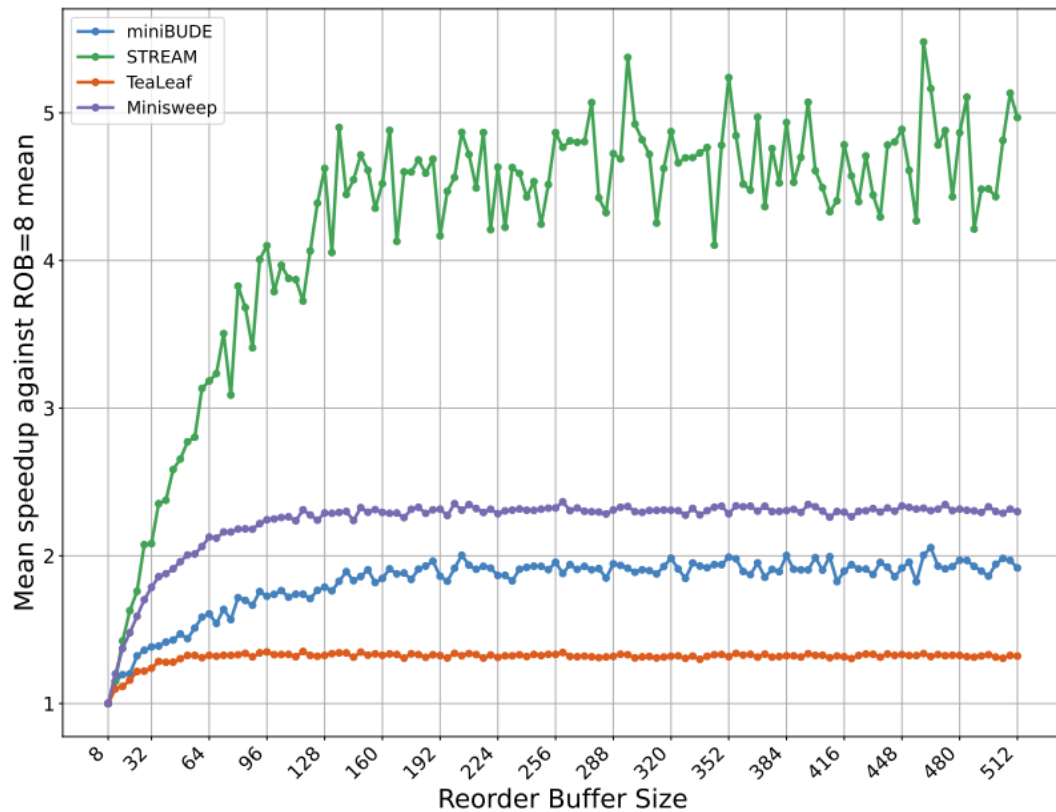


VL=128



VL=2048

# Vector Length

# ROB Size / FP/SVE Register Count

# What we found

- Vector Length unlocks huge Data-Level-Parallelism (when it's used)

- Memory speed (and capacity) is key

- Frontend throttles, not accelerates

# More interestingly…

- We can accurately map out known codes across a large search space
- Faster simulators and machines make data collection cheaper
- Decision Tree Regressors work nicely for modelling these high-dimensional relationships
- Relatively easy to map new codes against a defined architecture space
- Reduces one context of simulation down to a faster surrogate

# Future Work

- Multi-core/multi-node modelling to consider communication
- Modelling execution unit design
- Prediction of unseen codes with higher-capacity models
- Improved compiler cost-modelling to fully utilise hardware

# Thank you for listening

Any questions?

*zi23956@bristol.ac.uk*

# Full Search Space

| Parameter | Range | Step |
|---|---|---|
| Vector Length (Bits) | {128-2048} | Powers of 2 |
| Fetch-Block-Size | {4-2048} | Powers of 2 |
| Loop-Buffer-Size | {1-512} | 1 |
| General Purpose (GP) Registers | {38-512} | 8 starting from 40 |
| Floating-Point (FP)/SVE Registers | {38-512} | 8 starting from 40 |
| Predicate Registers | {24-512} | 8 |
| Conditional Registers | {8-512} | 8 |
| Commit Pipeline Width | {1-64} | 1 |
| Frontend Pipeline Width | {1-64} | 1 |
| Load-Store-Queue Completiton Pipeline Width | {1-64} | 1 |
| Reorder Buffer (ROB) Size | {8-512} | 4 |
| Load Queue Size | {4-512} | 4 |
| Store Queue Size | {4-512} | 4 |
| Load Bandwidth (Bytes) | {16-1024} | Powers of 2 |
| Store Bandwidth (Bytes) | {16-1024} | Powers of 2 |
| Permitted Memory Requests Per Cycle | {1-32} | 1 |
| Permitted Memory Loads Per Cycle | {1-32} | 1 |
| Permitted Memory Stores Per Cycle | {1-32} | 1 |

| Parameter | Range | Step |
|---|---|---|
| Cache Line Width (clw) | {32-512} | Powers of 2 |
| L1 Latency (Cycles) | {1-10} | 1 |
| L1 Clock Speed (GHz) | {1-5} | 0.5 |
| L1 Associativity | {1-16} | Powers of 2 |
| L1 Size (KiB) | {16-2048} | Powers of 2 |
| L2 Latency (Cycles) | {6-50} | 1 |
| L2 Clock Speed (GHz) | {1-5} | 0.5 |
| L2 Associativity | {1-16} | Powers of 2 |
| L2 Size (MiB) | {0.25 - 64} | Powers of 2 |
| Ram Timing (ns) | {40-250} | 10 |
| Ram Clock (GHz) | {1-5} | 0.5 |
| Ram Size (GiB) | 8 | N/A |

# Benchmark Parameters

| Application | Input options | Input Values |
|---|---|---|
| STREAM | Programming Model<br>Stream Array Size | OpenMP (single thread)<br>200000 |
| MiniBude | Programming Model<br>Benchmark Name<br>Atoms<br>Poses<br>Iterations | OpenMP (single thread)<br>bm1<br>26<br>64<br>1 |
| TeaLeaf | Programming Model<br>Dimensions<br>Number of cells along {X, Y}<br>Domain {xmin, xmax}, {ymin, ymax}<br>Solver Method<br>Initial Timestep<br>End Step<br>Max Iterations | OpenMP (single thread)<br>2D<br>{32, 32}<br>{0, 10}, {0, 10}<br>Conjugate Gradient<br>0.004<br>5<br>10000 |
| MiniSweep | Programming Model<br>Global number of gridcells along {X, Y, Z}<br>Total number of energy groups<br>Number of angles for each octant direction<br>Sweep Iterations<br>Sweep blocks used to tile the Z dimension | OpenMP (single thread)<br>{4, 4, 4}<br>1<br>32<br>1<br>1 |