# A High-fidelity Flow for High-Performance RISC-V CPU Design

Luke Yen, Yuanbo Fan, Wei-Han Lien

tenstorrent

# Challenges in high-performance RISC-V CPU Design

## Tedious Process

- Years of efforts
  - Design, implementation & tape-out
- Many turn-arounds
  - uArch tuning based on PPA goals
- IP progress tracking
  - Internal verification & regression
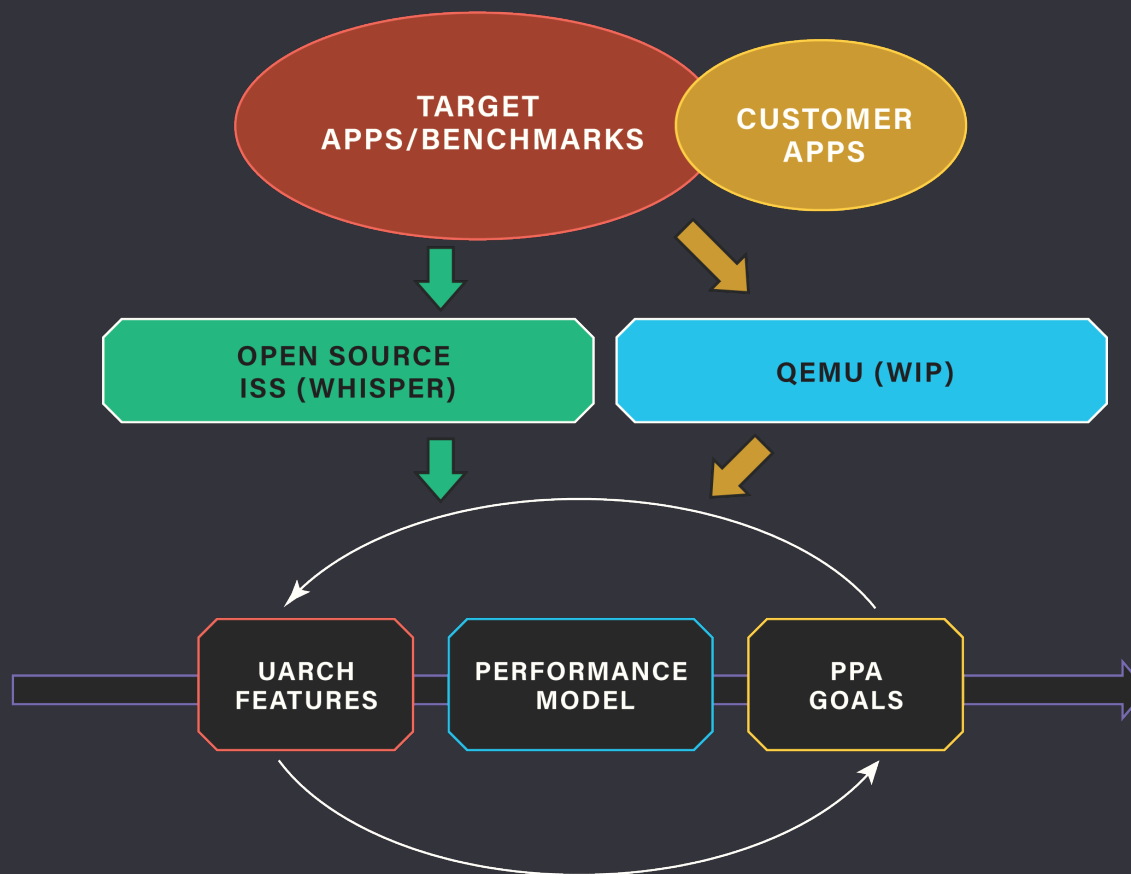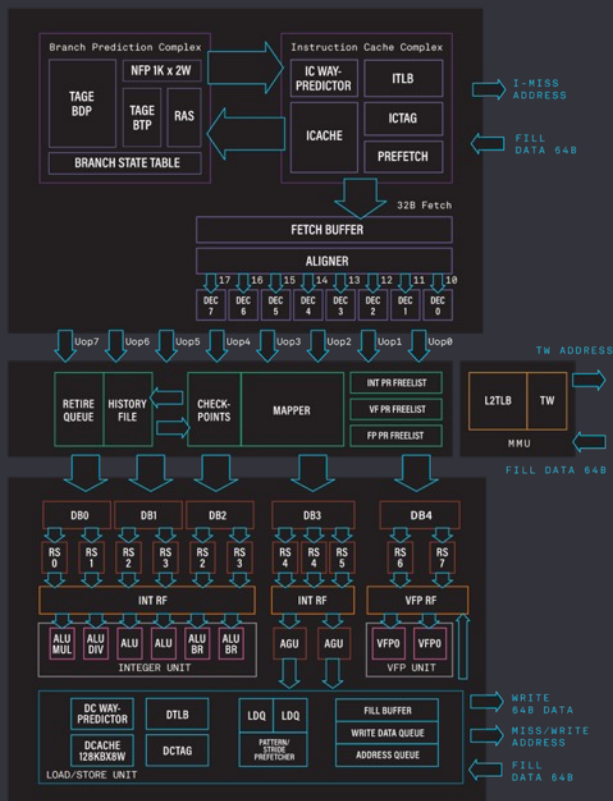  - External visibility & validation

## Performance Projections

- Starting from pre-silicon stage
- Target applications/benchmarks

## Collaboration

- Customized features
  - Branch predictor & prefetcher
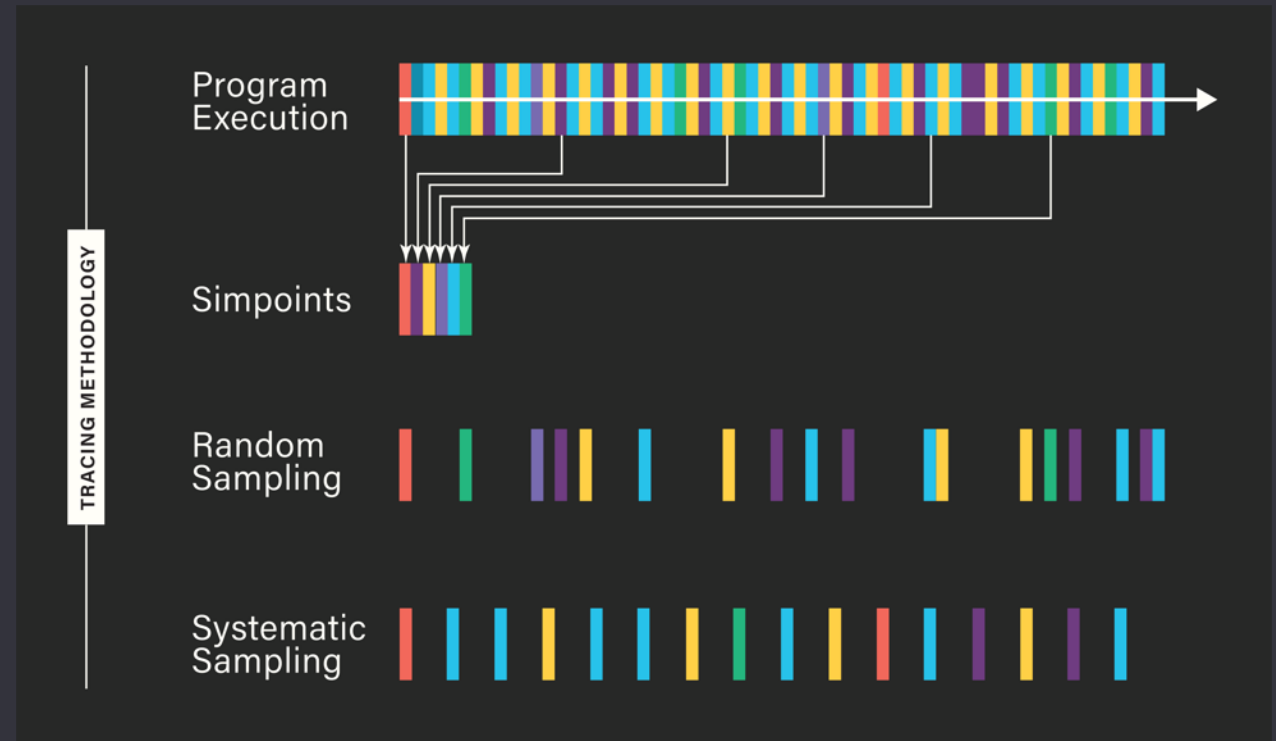- Sub-system
  - Cache & memory
- Vector unit design
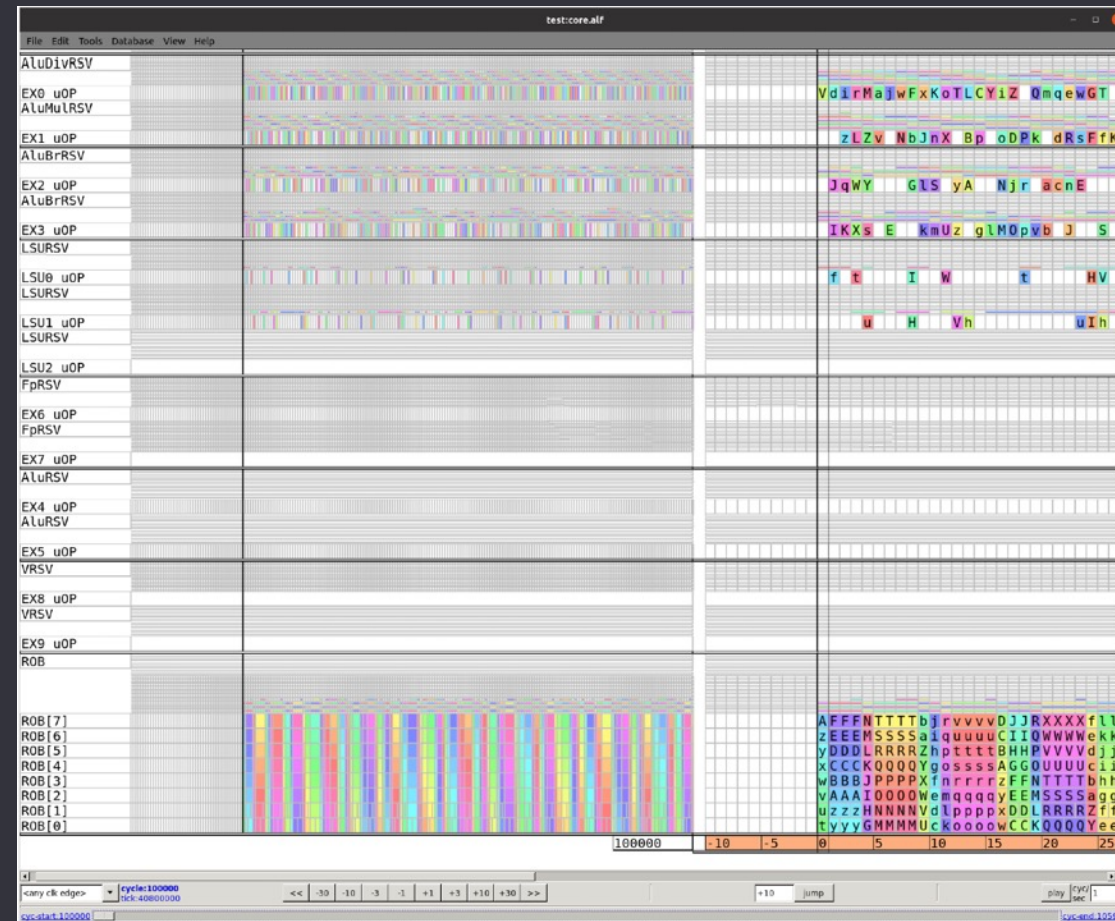
tenstorrent     Confidential

# Methodology

# Instruction Set Simulator (i.e. Whisper)

- Available: Github Link

- Simulation
  - ➤ Speed: ~100M instructions per second
  - ➤ Static trace generation
  - ➤ Co-simulation (with performance model)
- Tracing Methodology
  - ✓ Random/systematic sampling
  - ✓ Representative phases (i.e. simpoints)
  - ✓ Warmup traces (for branch and I/D cache)
- Limitation
  - ❑ Static linking
  - ❑ Single-thread simulation

# Cycle-Accurate Performance Model

- Cycle-accurate, event-driven Model

  - Highly flexible and configurable

  - Trace-driven simulation

- Performance Metrics

  - Such as IPC, total execution cycles,
    cache misses, TLB misses, branch
    mispredicts

- Visualization & Debug

  - Pipetrace: a visual representation of the
    pipeline execution over time

# PPA-orientated Tuning

## Tradeoff: Projection Accuracy vs. Simulation Speed

- The tracing flow consists of several steps including:
  - ✓ Whisper simulation
  - ✓ Benchmark profiling
  - ✓ Snapshot generation
  - ✓ Trace generation

- Customized microbenchmark
  - C/C++ tests
  - Major uarch components & timing paths
  - Critical perf. metrics (e.g. data cache hit latency, issue bandwidth)

- Calibration & Validation
  - ❖ Tradeoffs: Low projection error vs 10x simulation speedup
  - ❖ End-to-end co-simulation



| CPU2017 INTRATE | I20M,K100 | I20M,K30 |
|---|---|---|
| DEEPSJENG_R | 0.99 | 0.98 |
| EXCHANGE2_R | 1.01 | 1.02 |
| GCC_R | 1.03 | 0.93 |
| LEELA_R | 0.99 | 0.99 |
| MCF_R | 1.00 | 0.98 |
| OMNETPP_R | 0.96 | 0.93 |
| PERLBENCH_R | 0.98 | 0.99 |
| X264_R | 0.97 | 0.97 |
| XALANCBMK_R | 1.00 | 0.96 |
| XZ_R | 1.00 | 1.00 |
| GEOMEAN | 0.99 | 0.97 |

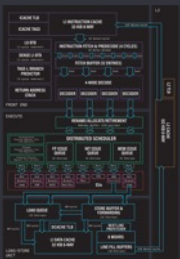Note: normalized by performance projected by 100M simpoints

tenstorrent   Confidential

# Tenstorrent RISC-V O-o-O Processor Family
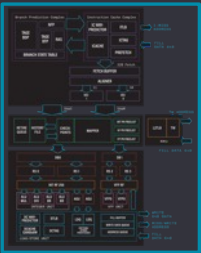
Higher Performance

Performance

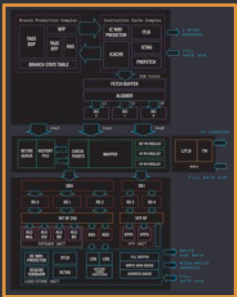**One Design and 5 IPs in a year**



8-Wide Decode
Ascalon
Server, Laptop, and HPC

6-Wide Decode
Alastor
Client and Edge

4-Wide Decode

3-Wide Decode

2-Wide Decode

4-Wide Decode
Sonic Boom with Vector

Decode Width

# Future Collaboration

Instruction Set Simulator (Whisper)

Target applications/benchmarks & sharable traces

Compiler development & optimization

tenstorrent