# Concorde: Fast and Accurate CPU Performance Modeling with Compositional Analytical-ML Fusion

Arash Nasr-Esfahany, Mohammad Alizadeh, Victor Lee
Hanna Alam, Brett W. Coon, David Culler, Vidushi Dadu
Martin Dixon, Henry M. Levy, Santosh Pandey
Parthasarathy Ranganathan, Amir Yazdanbakhsh

Motivation

## 1. CPU Simulation

**Microarchitecture simulation is a key tool in design and exploration, but we lack fast and accurate performance models.**
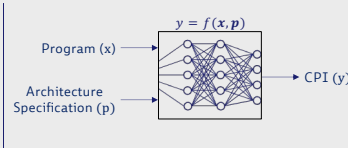
Cycle-level simulators are slow.
Analytical Models are fast, but not accurate.

Motivation

## 2. Prior Work

Ignores problem structure
❌ High sample complexity
❌ Bulky neural networks
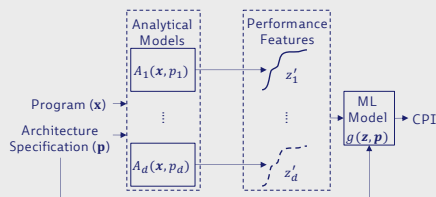❌ Slow training and inference
❌ $\mathcal{O}(\#instrs)$



**SimNet**: Accurate and High-performance Architecture Simulation using Deep Learning, ACM SIGMETRICS/IFIP PERFORMANCE '22
**TAO**: Re-Thinking DL-based Microarchitecture Simulation, ACM SIGMETRICS/IFIP PERFORMANCE '24

Design

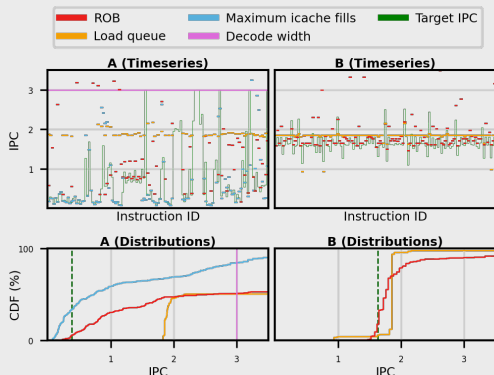## 3. Compositional Analytical-ML Fusion

Multiple lightweight models work together to progressively achieve high fidelity with low computational complexity.
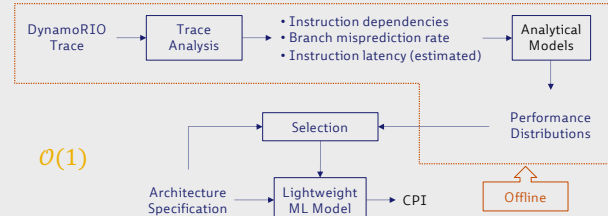


Design

## 4. Performance Features

Per-resource analytical modeling produces a rich performance characterization of a program.
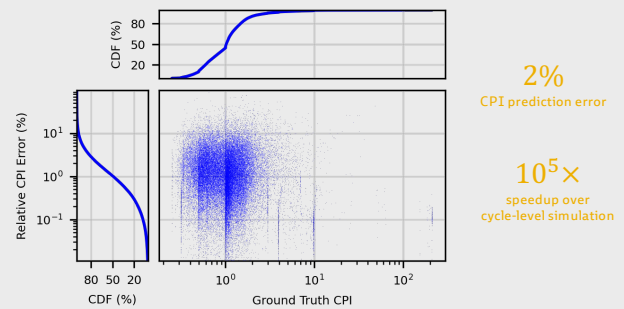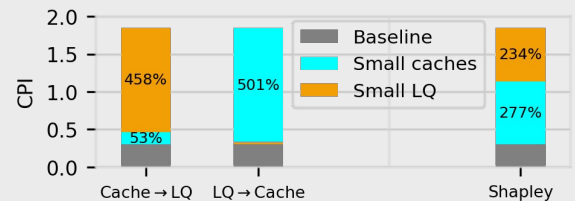


Design

## 5. Concorde



$\mathcal{O}(1)$

Evaluation

## 6. Concorde is fast and accurate



**2%**
CPI prediction error

$10^5 \times$
speedup over
cycle-level simulation

Application

## 7. Fine-Grained Performance Attribution

Shapley Value: A fair, order-independent attribution



Case Study

## 8. Large-Scale Sensitivity Analysis

143M 100k-instruction segment CPI evaluations, in just an hour!