

# ECE382N.23: Embedded System Design and Modeling

---

## Lecture 5 – ML-Based Predictive Modeling

Andreas Gerstlauer

Electrical and Computer Engineering

University of Texas at Austin

gerstl@ece.utexas.edu



The University of Texas at Austin

Chandra Department of Electrical  
and Computer Engineering

*Cockrell School of Engineering*

---

## Lecture 5: Outline

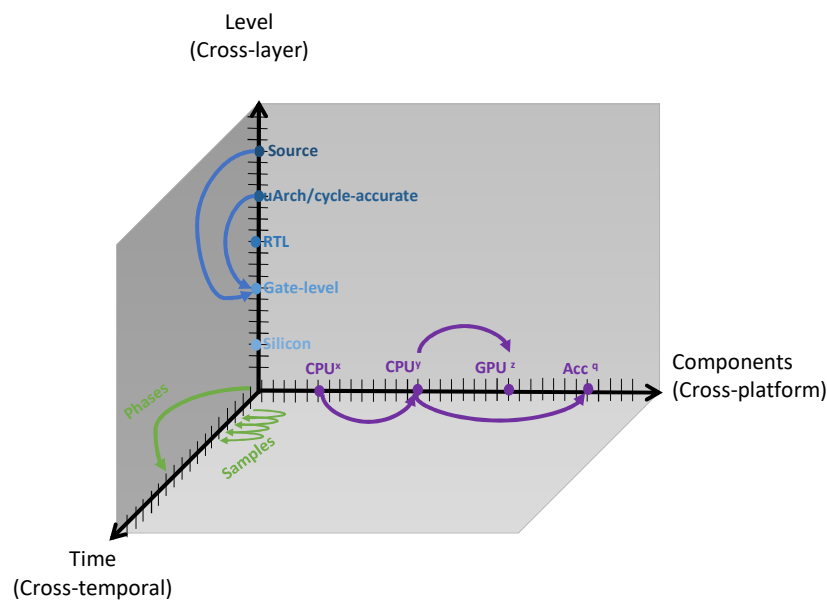
---

- **Overview**
  - Predictive modeling dimensions
- **Cross-layer prediction**
  - Micro-architecture & source level power prediction
- **Cross-platform prediction**
  - CPU to CPU performance and power prediction
- **Cross-temporal prediction**
  - Short-term & long-term workload forecasting

## Machine Learning for Modeling

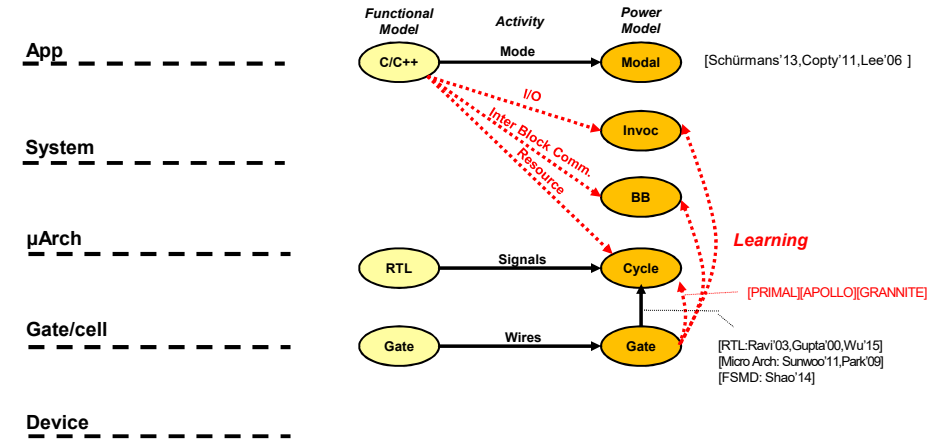
- **Learn rather than construct models**
  - Learn an abstract model from detailed observations
- **Predict rather than simulate**
  - Replace detailed simulations with predictions
- **Supervised learning**
  - Interpolate/extrapolate (complex) behavior
  - From (a few) training samples
  - Exploit inherent correlations
  - Domain-specific feature selection & learning formulations
    - Can not afford deep learning from raw data

## Predictive System Modeling

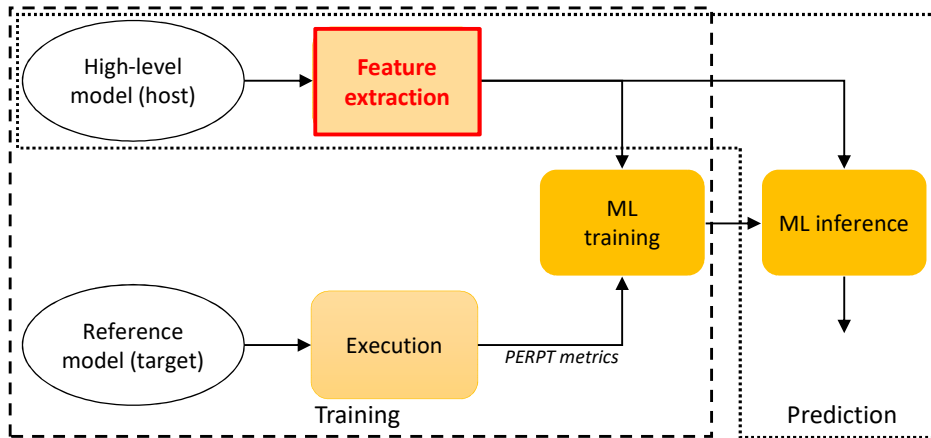


## Cross-Layer Prediction

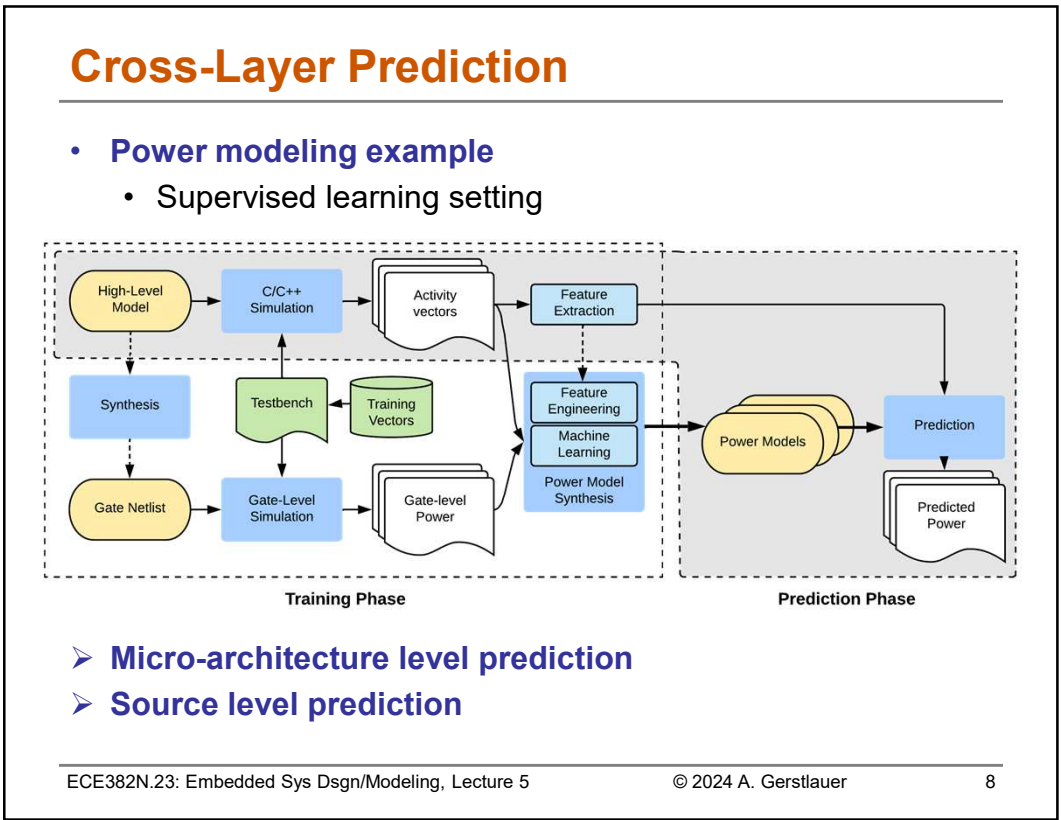
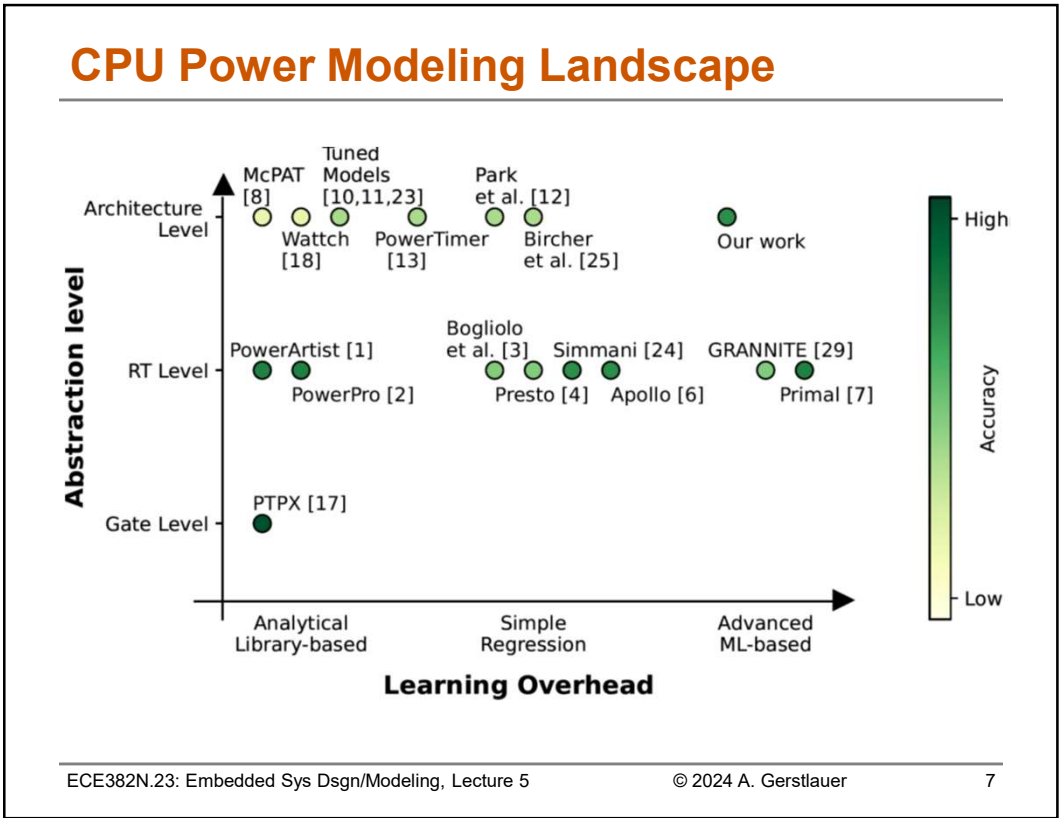
- Power modeling example



## Supervised Learning Formulation

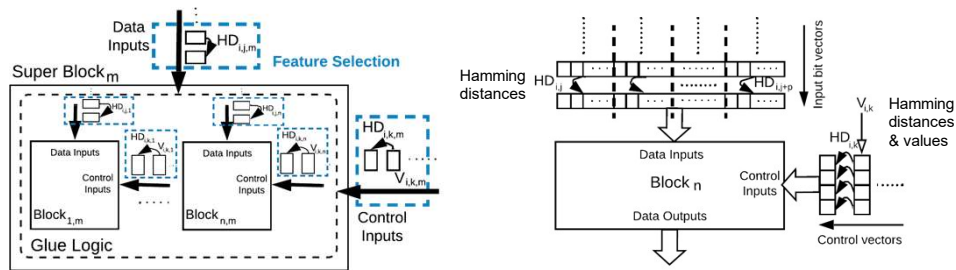


- **Manually selected/engineered vs. learned features**
  - Traditional ML vs. deep learning



## Micro-Architecture Level Prediction

- **Predict gate-level power from high-level CPU simulations**
  - Activity features from micro-architectural simulations
    - Block-level I/O for key data vs. control signals
    - Hamming distances & actual values
  - Hierarchically compose micro-architecture block models
    - Blocks -> Pipeline stages -> Core
    - Include glue logic model at the whole core level
- Predict cycle- and block-level power up to complete CPU



Source: A. K. Ananda Kumar, et al., "Machine Learning-Based Microarchitecture-Level Power Modeling of CPUs," IEEE TC, 2022.

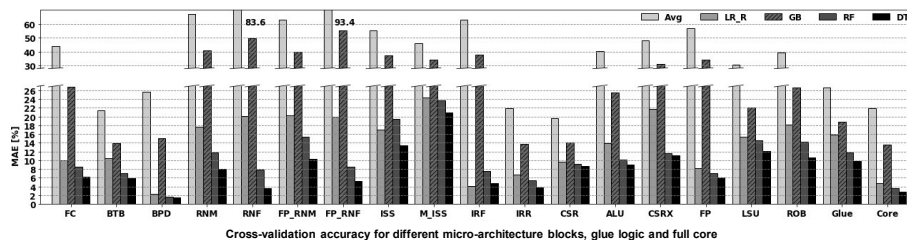
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

9

## Learning Setup

- **Apply advanced machine learning regressors**
  - Linear (LR), gradient boosting (GB)
  - Random forest (RF), decision tree (DT)
- **In-order and out-of-order RISC-V cores (RI5CY & BOOM)**
  - Training & cross-validation using micro-benchmarks
  - Testing using full applications (CoreMark & FFT)
  - Cycle-by-cycle power error vs. gate-level simulations



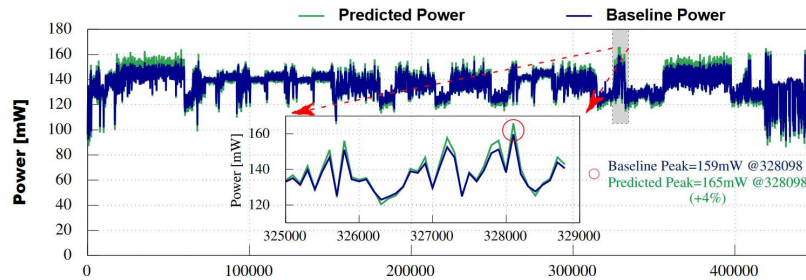
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

10

## Micro-Architecture Prediction Results

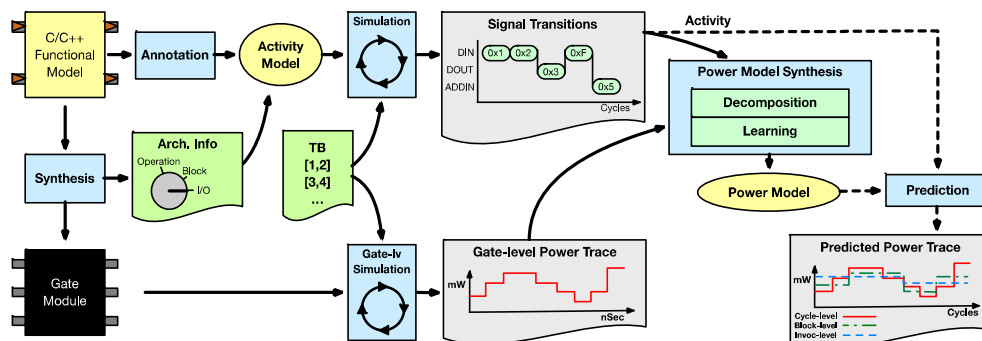
- Cycle-by-cycle / avg. power with 3.6% / less than 1% error
  - MAE decreases exponentially with increasing averaging



- Very fast learning rate and low prediction overhead
  - Models trained in less than 20k samples
  - Models predict with a throughput of 4Mcycles/s

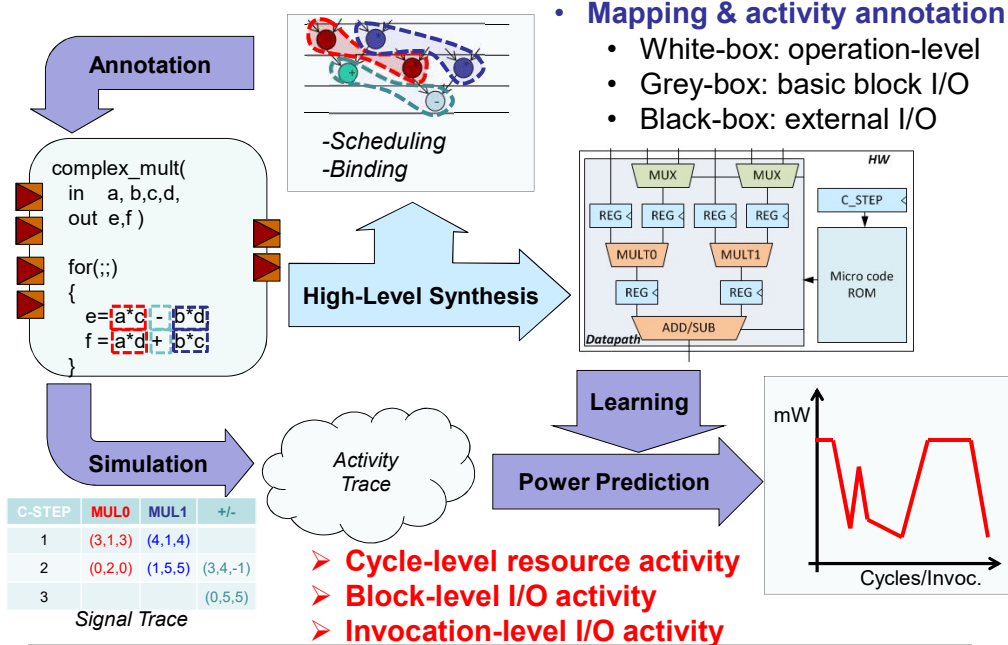
## Source-Level Prediction

- Custom hardware accelerator power modeling
  - White / grey / black box hardware models
  - Operation / block / I/O activity from C/C++ simulation
  - Predict gate-level trace at cycle / block / invocation granularity



➤ Data-dependent, Fast

## Source-Level Prediction Flow



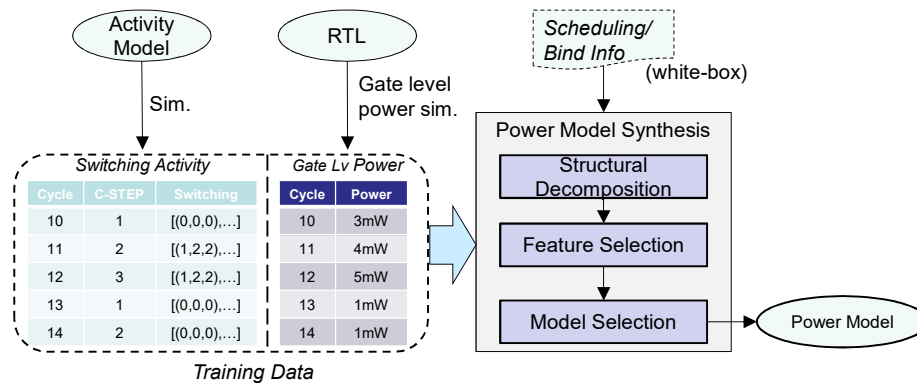
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

13

## Power Model Synthesis

- **Learning- vs. library-based power model**
  - More accurate than library-based models (glue logic, etc.)
  - One-time training overhead using gate-level simulation



- **Apply state-of-the-art machine learning techniques**
  - Structural model decomposition & feature selection

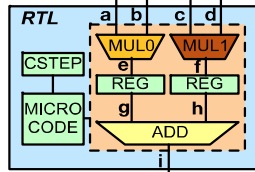
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

14

## Power Model Decomposition

### • Single power model



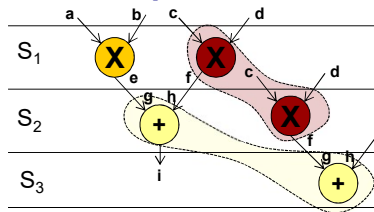
- White-box: cycle-level, signal activity

$$P(t) = f_{cycle}(HD_n(t)), n = a \dots i$$

- Black-box: invocation-level, I/O history

$$P(t) = f_{invoc}(HD_n(t), HD_n(t-1), HD_n(t-2)), n = a \dots d$$

### • Decomposed model



- White-box: decomposed model

$$P_{S1}(t) = f_1(HD_a(t), \dots, HD_f(t))$$

$$P_{S2}(t) = f_2(HD_c(t), HD_d(t), HD_f(t), HD_g(t), HD_h(t), HD_i(t))$$

$$P_{S3}(t) = f_3(HD_g(t), HD_h(t), HD_i(t))$$

- Black-box: ensemble model

$$P(t) = \text{avg}(f_{S(i)}(HD_n(t), HD_n(t-1), HD_n(t-2))), n = a \dots d$$

- Only capture signals (features) utilized in each state
- Apply feature selection to further remove correlated signals

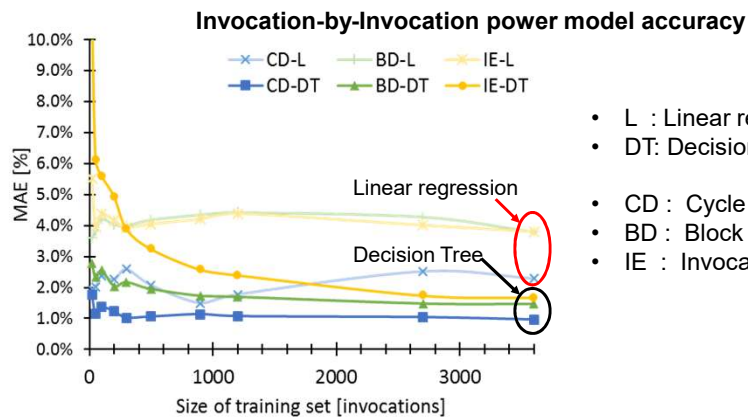
## Learning Formulation

### • Dedicated, domain-specific learning formulations

- Structural model decomposition & feature selection

### • Advanced, non-linear regression models

- *Not* deep learning (small training size)



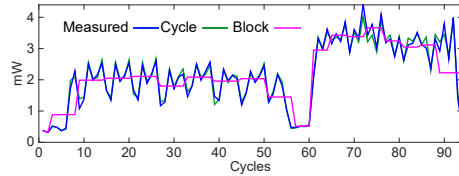
- L : Linear regression
- DT: Decision Tree regression
- CD : Cycle decomposed model
- BD : Block decomposed model
- IE : Invocation ensemble model



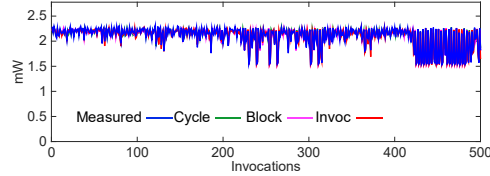
## Source-Level Prediction Results

### • Pipelined 2D-DCT

- Cycle-by-cycle trace

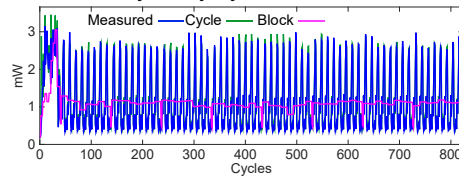


- Invocation-by-invocation trace

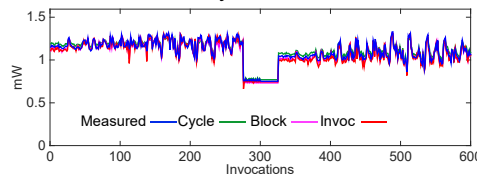


### • Pipelined HDR weight comp.

- Cycle-by-cycle trace



- Invocation-by-invocation trace



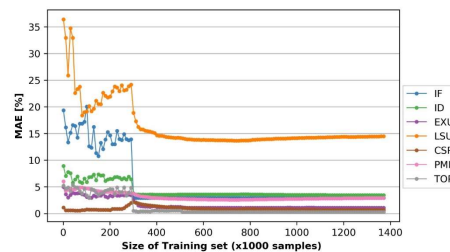
➤ **> 90-97% accuracy @ 1-10Mcycles/s speed**

- 2,000-10,000x faster than gate-level, 100x-500x faster than RTL

## Cross-Layer Prediction Questions

### • Training time and training set

- Time to collect training samples and train model less than time to simulate
- Deep learning not an option!

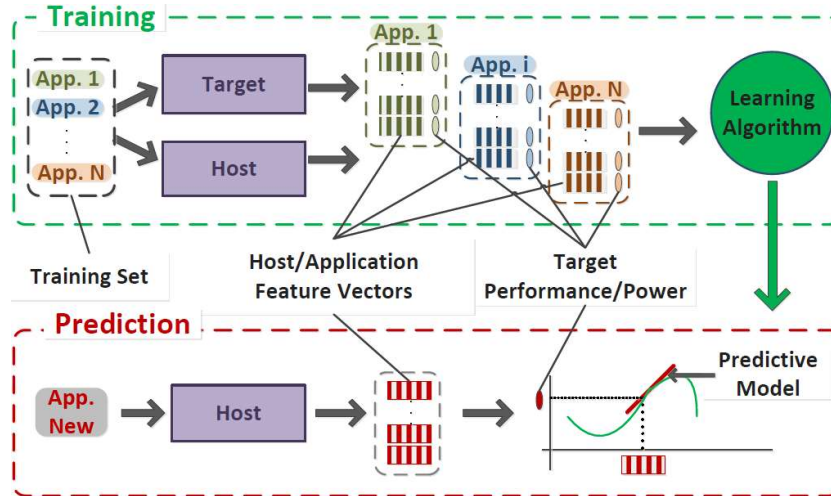


### • Need for accurate reference model

- Modeling effort
- Unless real hardware is available
  - May defeat the purpose of modeling

## Cross-Platform Prediction

- **Power & performance prediction**



- CPU->CPU, GPU->GPU, CPU->FPGA

## CPU-to-CPU Prediction

- **Predict on target CPU while running on host CPU**

- Using hardware counters on host as features
- Predict target performance and power
- At program phase level

- **Instrumentation-based**

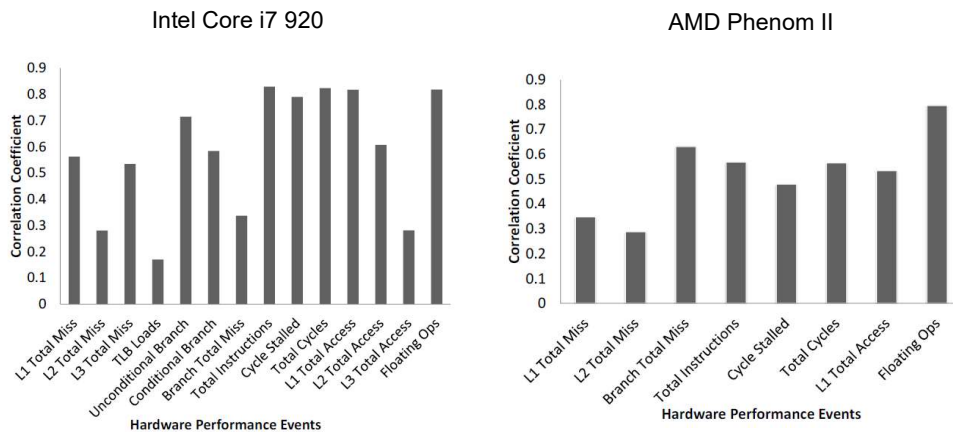
- Compiler-based instrumentation at basic block granularity
- Collect features and train/call model every  $N$  basic blocks

- **Sampling-based**

- Source-oblivious at binary level using timer interrupts
- Sample alignment during training

## Hardware Counters as Features

- **Correlation of host counter events with target timing**
  - ARM target [GEM5] vs. counters on host [PAPI]



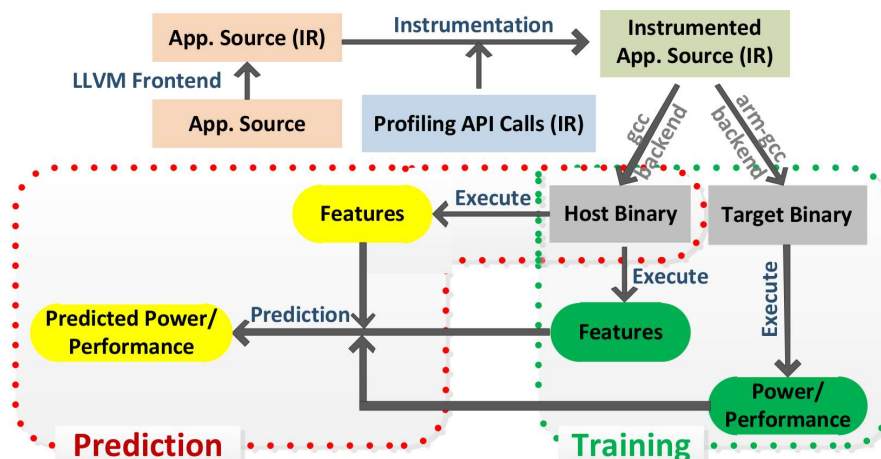
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

21

## Instrumentation-Based Prediction

- **Every  $N$  number of basic blocks (BBs)**
  - Collect host counters and target metrics during training
  - Collect host counters and call model during prediction

Source: X. Zheng, L. John, A. Gerstlauer, "LACross: Learning-Based Analytical Cross-Platform Performance and Power Prediction," *IJPP*, 45(6), 2017.

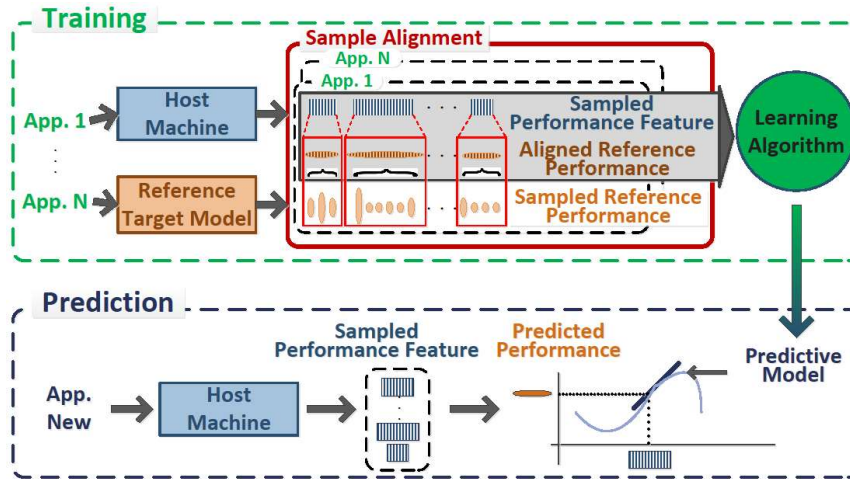
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

22

## Sampling-Based Prediction

- **Source-oblivious, binary-level prediction**
  - Transparent background prediction on timer interrupt
  - Arbitrary library, OS and system code



Source: X. Zheng, et al., "Sampling-Based Binary-Level Cross-Platform Performance Estimation," DATE, 2017.

ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

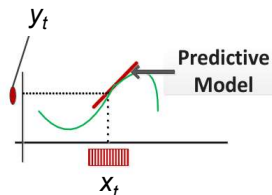
23

## Learning Formulation

- **Given training set  $(x_p, y_i)$** 
  - $x_i \in \mathbb{R}^d$ :  $d$ -dimensional counter feature vector from host
  - $y_i \in \mathbb{R}$ : reference performance/power on target
- **Want to find function  $F(x_i) \approx y_i$** 
  - Fundamentally non-linear
- **Locally linear approximation  $F_t(x_t)$  at input  $x_t$**

$$F_t(x_t) = \theta_t^T x_t$$

- Around neighborhood of  $x_t$
- LASSO regression to solve for  $\theta_t$



ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

24

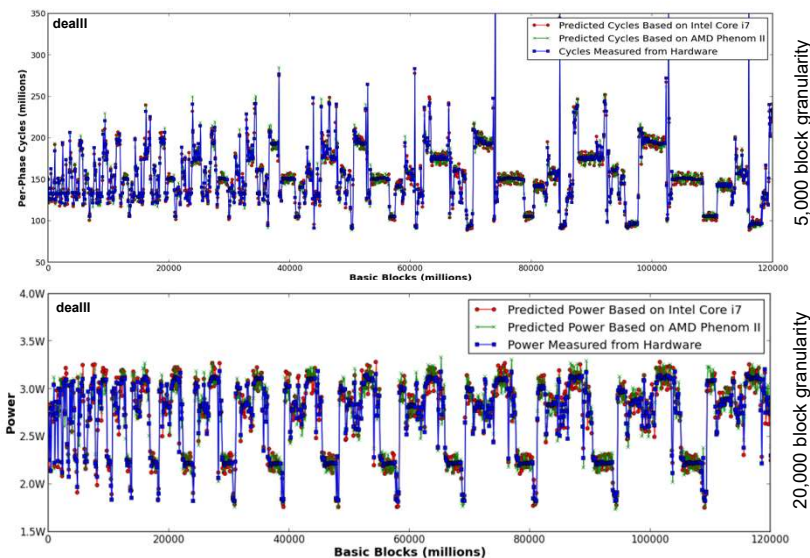
## CPU-to-CPU Prediction Setup

- **Platforms**
  - Target: Samsung ARM A9/A15 Exynos
  - Host: Intel Core i7 / AMD Phenom II
- **Host counters**
  - Instrumentation-based: 14 / 8 counters
  - Sampling-based: 6 counters
- **Learning formulation**
  - Phase-level localized LASSO regression
- **Training set**
  - 157-284 programs of ACM-ICPC competition
- **Test set**
  - 7 programs from MiBench and 8 programs from SD-VBS
  - 19 programs from SPEC CPU 2006
  - 13 Java & Python benchmarks from DaCapo/PyBench

Host Counters
Instructions
Cycles
Total Cache Misses
Total Cache References
Total Branches
Total Branch Misses

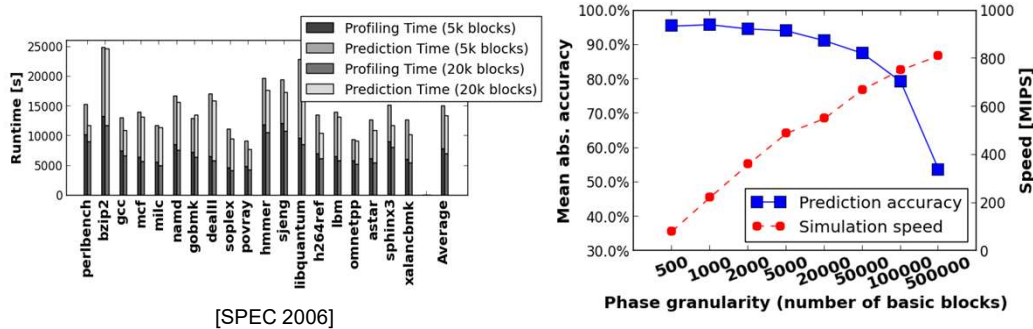
## Instrumentation-Based Prediction Results

- **Performance & power prediction**
  - 90-95% per-phase accuracy @ 500-600 MIPS throughput



## Instrumentation-Based Speed & Accuracy

- **Accuracy & speed vs. phase granularity**
  - Finer granularity requires more prediction overhead
  - But: more & better training data w/ finer granularity
    - Phase similarity: number of unique phases *decreases* linearly
  - Runtime also limited by hardware counter support on host
    - Multiple runs needed to collect all counters



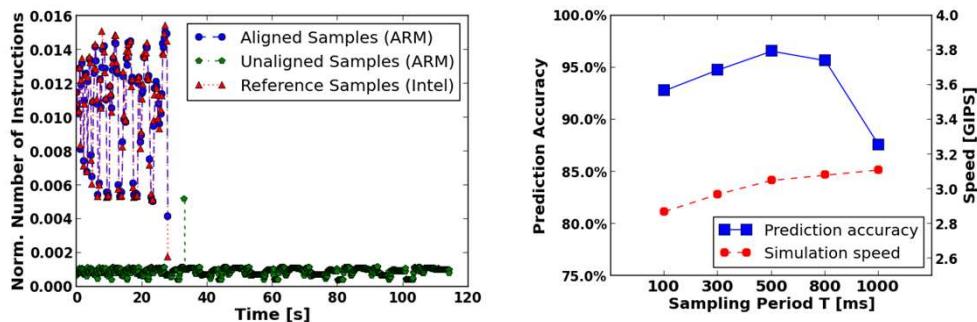
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

27

## Sampling-Based Results

- **Speed & accuracy increase with coarser host sampling  $T$** 
  - Better alignment, until lack of training data ( $T > 500$ ms)



### ➤ 96% accuracy @ 3 GIPS ( $T = 500$ ms)

- No instrumentation overhead (6x faster)
  - Fewer counters, coarser granularity, but requires more training
- 2x faster than running native on ARM target

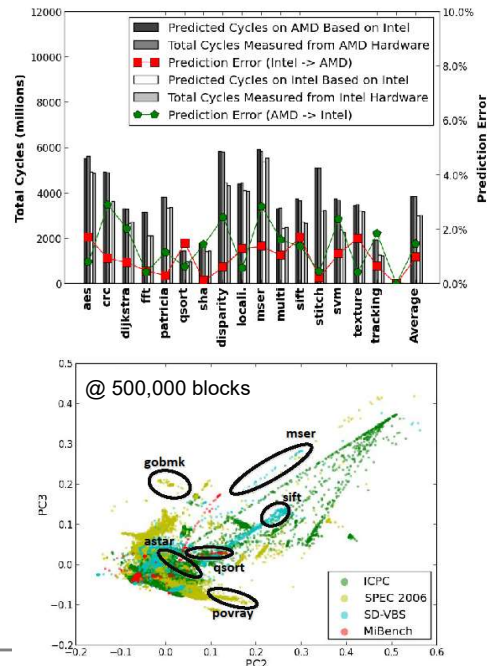
ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

28

## Cross-Platform Prediction Questions

- **Host/target pairs**
  - ARM from x86, x86-to-x86
  - From simple to complex?
- **Prediction features**
  - Which counters?
  - Other information?
- **Training set**
  - Larger granularity requires larger training set
  - Optimal training set?
    - Generate synthetic training set (Genesys) [SAMOS'16]



ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

## Other Cross-Platform Approaches

- **GPU performance models (Intel/UC Riverside, P. Brisk)**
  - GPU-to-GPU prediction using performance counters
  - Commercial GPUs to predict pre-silicon hardware
- **FPGA high-level synthesis models (UC Riverside, P. Brisk)**
  - Predict FPGA performance of code regions of interest
  - Running on host CPU, using hardware counters
- **Heterogeneous ISA models for OSs (UCSD, D. Tullsen)**
  - Predict performance on different CPU cores
  - Use prediction to make OS scheduling decisions
- **CPU benchmark performance models (Harvard, D. Brooks)**
  - Predict benchmark performance from CPU specifications

ECE382N.23: Embedded Sys Dsgn/Modeling, Lecture 5

© 2024 A. Gerstlauer

30

## Lecture 5: Summary

---

- **ML for modeling**
  - Learn, not model
  - Predict, not simulate
- **Long history of learning-based modeling approaches**
  - Various forms of regression
  - Most problems are not linear
- **Advanced machine-learning to capture complex relations**
  - Cross-layer
  - Cross-platform
  - Cross-temporal