

Accurate & Efficient Regression Modeling for Microarchitectural Performance & Power Prediction

Benjamin C. Lee, David M. Brooks

{bclee,dbrooks}@eecs.harvard.edu
Division of Engineering and Applied Sciences
Harvard University

24 October 2006



Outline

Motivation & Background

- Simulation Challenges
- Simulation Paradigm
- Regression Theory

Model Derivation

- Experimental Methodology
- Derivation Overview
- Model Specification

Model Evaluation

- Performance
- Power

Conclusion



Outline

Motivation & Background

Simulation Challenges

Simulation Paradigm

Regression Theory

Model Derivation

Experimental Methodology

Derivation Overview

Model Specification

Model Evaluation

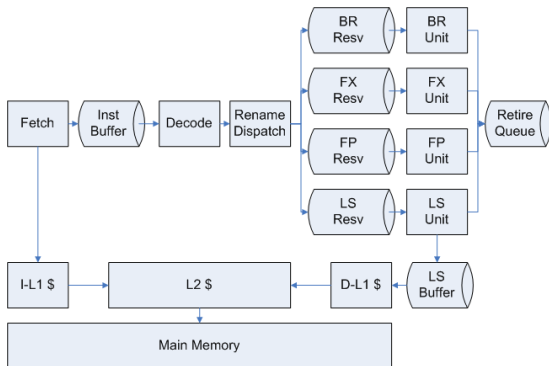
Performance

Power

Conclusion



Microarchitectural Design Space



- Increasing diversity of interesting, viable designs
- Examples :: Power 4, Pentium 4, UltraSPARC T1
- Tractably quantify trends across comprehensive design space



Microarchitectural Simulation Challenges

● Cycle-Accurate Simulation

- Accurately identifies trends in design space
- Tracks instructions' progress through microprocessor
- Estimates performance, power, temperature, . . .

● Simulation Costs

- Long simulation times (minutes, hours per design)
- Number of potential simulations scale exponentially (m^p)
 - p :: parameter count
 - m :: parameter resolution



Microarchitectural Sampling

● Temporal Sampling

- Sample from instruction traces in time domain
- Reduce simulation costs via size of inputs
- Synthetic traces from profiled workloads ¹
- Sampled traces from phase analysis ²

● Spatial Sampling

- Sample from design space
- Reduce simulation costs via number of simulations

¹Eeckhout [ISPASS'00]

²Sherwood [ASPLOS'02], Wunderlich [ISCA'03]



Simulation Paradigm

- **Comprehensively understand design space**
 - Specify large, high-resolution design space
 - Consider all design parameter simultaneously
- **Selectively simulate modest number of designs**
 - Sample points randomly from design space for simulation
 - Decouple resolution of design space and simulation
- **Efficiently leverage simulation data with inference**
 - Reveal trends, trade-offs from sparse sampling
 - Enable predictions for metrics of interest



Regression Theory

● Statistical Inference

- Models approximate solutions to intractable problems
- Requires initial data to train, formulate model
- Leverages correlations from initial data for prediction

● Regression Models

- Low formulation costs (1K samples from 1B designs)
- Accurate inference (4 – 7% median error)
- Efficient computation (100's of predictions per second)



Model Formulation

● Notation

- n observations \triangleright {*simulated design samples*}
- Response $:: \vec{y} = y_1, \dots, y_n$ \triangleright {*e.g., performance, power*}
- Predictor $:: \vec{x}_i = x_{i,1}, \dots, x_{i,p}$ \triangleright {*e.g., depth, cache*}
- Regression Coefficients $:: \vec{\beta} = \beta_0, \dots, \beta_p$
- Random Error $:: \vec{e} = e_1, \dots, e_n$ where $e_i \sim N(0, \sigma^2)$
- Transformations $:: f, \vec{g} = g_1, \dots, g_p$

● Model

$$f(y) = \beta_0 + \sum_{j=1}^p \beta_j g_j(x_j) + e$$



Predictor Interaction

● Modeling Interaction

- Suppose effects of predictors x_1, x_2 cannot be separated
- Construct predictor $x_3 = x_1x_2$

$$y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1x_2 + e_i$$

● Example

- Let x_1 be pipeline depth, x_2 be L2 cache size
- Performance impact of pipelining affected by cache size

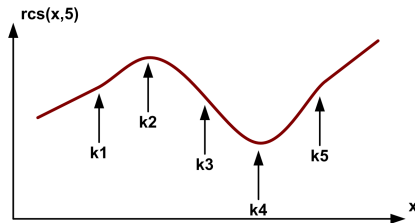
$$Speedup = \frac{Depth}{1 + Stalls/Inst}$$



Predictor Non-Linearity I

● Restricted Cubic Splines

- Divide predictor domain into intervals separated by knots
- Piecewise cubic polynomials joined at knots
- Higher order polynomials provide better fits³



³Stone [SS'86]



Predictor Non-Linearity II

- **Location of Knots**

- Location of knots less important than number of knots ⁴
- Place knots at fixed predictor quantiles

- **Number of Knots**

- Flexibility, risk of over-fitting increases with knot count
- 5 knots or fewer are often sufficient
- 4 knots balances flexibility, risk of over-fitting

⁴Harrell [Springer'01]



Prediction

- **Expected Response**

- β are known from least squares
- $x_{i,1}, \dots, x_{i,p}$ are known for a given query i
- Expected response is weighted sum of predictor values

$$\begin{aligned} E[y] &= E\left[\beta_0 + \sum_{j=1}^p \beta_j x_j\right] + E[e] \\ &= \beta_0 + \sum_{j=1}^p \beta_j x_j \end{aligned}$$



Outline

Motivation & Background

Simulation Challenges

Simulation Paradigm

Regression Theory

Model Derivation

Experimental Methodology

Derivation Overview

Model Specification

Model Evaluation

Performance

Power

Conclusion



Tools and Benchmarks

● Simulation Framework

- Turandot :: a cycle-accurate trace driven simulator
- PowerTimer :: power models derived from circuit analyses
- Baseline simulator models POWER4/POWER5 architecture

● Benchmarks

- SPEC2kCPU :: compute-intensive benchmarks
- SPECjbb :: Java server benchmark

● Statistical Framework

- R :: software environment for statistical computing
- Hmisc and Design packages⁵

⁵Harrell [Springer,'01]



Spatial Sampling

● Design Space

- S_i :: set of values for parameter x_i , $i \in [1, p]$
- $S = \prod_{i=1}^p S_i$:: design space
- B :: set of benchmarks
- $|S| \approx 10^9$ and $|B| = 22$

● Sampling Uniformly at Random (UAR)

- Sample $n = 4,000$ designs and benchmarks for simulation
- Decouple resolution of design space and simulation
- Unbiased observations from full range of parameter values



Predictors :: Microarchitecture

	Set	Parameters	Measure	Range	$ S_i $
S_1	Depth	depth	FO4	9::3::36	10
S_2	Width	width L/S reorder queue store queue functional units	insn b/w entries entries count	4,8,16 15::15::45 14::14::42 1,2,4	3
S_3	Physical Registers	general purpose (GP) floating-point (FP) special purpose (SP)	count count count	40::10::130 40::8::112 42::6::96	10
S_4	Reservation Stations	branch fixed-point/memory floating-point	entries entries entries	6::1::15 10::2::28 5::1::14	10
S_5	I-L1 Cache	i-L1 cache size	$\log_2(\text{entries})$	7::1::11	5
S_6	D-L1 Cache	d-L1 cache size	$\log_2(\text{entries})$	6::1::10	5
S_7	L2 Cache	L2 cache size L2 cache latency	$\log_2(\text{entries})$ cycles	11::1::15 6::2::14	5
S_8	Control Latency	branch latency	cycles	1,2	2
S_9	FX Latency	ALU latency FX-multiply latency FX-divide latency	cycles cycles cycles	1::1::5 4::1::8 35::5::55	5
S_{10}	FP Latency	FPU latency FP-divide latency	cycles cycles	5::1::9 25::5::45	5
S_{11}	L/S Latency	Load/Store latency	cycles	3::1::7	5
S_{12}	Memory Latency	Main memory latency	cycles	70::5::115	10



Predictors :: Application-Specific

- **Application Characteristics**

- Collect program characteristics on baseline architecture
- Instruction throughput
- Cache access patterns
- Branch patterns
- Sources of pipeline stalls

- **Application Effects**

- Significant interactions with microarchitectural predictors
- Example :: Impact of d-L1 cache affected by access rates



Derivation Overview

- **Hierarchical Clustering**
- **Performance Associations and Correlations**
 - qualitative scatterplots, quantitative ρ^2
- **Model Specification**
 - predictor interaction, non-linearity
- **Assessing Fit**
 - R^2 statistic
- **Residual Analysis**
 - normality (quantile-quantile), randomness (scatterplots)
- **Significance Testing**
 - hypothesis testing, F-statistic, p-values

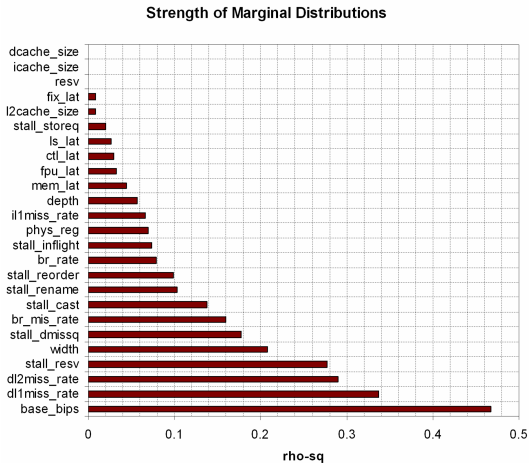


Derivation Overview

- Hierarchical Clustering
- **Performance Associations and Correlations**
 - qualitative scatterplots, quantitative ρ^2
- **Model Specification**
 - predictor interaction, non-linearity
- Assessing Fit
 - R^2 statistic
- Residual Analysis
 - normality (quantile-quantile), randomness (scatterplots)
- Significance Testing
 - hypothesis testing, F-statistic, p-values



Performance Correlations



Model Specification

● Interactions

- Pipeline width/depth interact with
 - instruction bandwidth structures (queues, register file)
 - cache hierarchy
- Cache hierarchy sizes interact with
 - adjacent levels in hierarchy
 - application-specific access rates
- Baseline performance interacts with resource sizings

● Restricted Cubic Splines

- Weaker relationships (latencies, caches, queues) :: 3 knots
- Stronger relationships (depth, registers) :: 4 knots
- Baseline performance :: 5 knots



Outline

Motivation & Background

Simulation Challenges
Simulation Paradigm
Regression Theory

Model Derivation

Experimental Methodology
Derivation Overview
Model Specification

Model Evaluation

Performance
Power

Conclusion



Validation Approach

● Framework

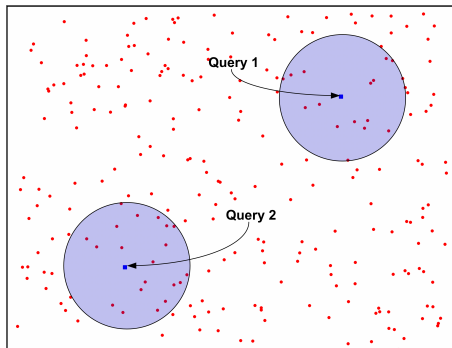
- Formulate models with $n < 4,000$ samples
- Obtain 100 additional random samples for validation
- Quantify percentage error, $100 * |\hat{y}_i - y_i|/y_i$

● Model Variants

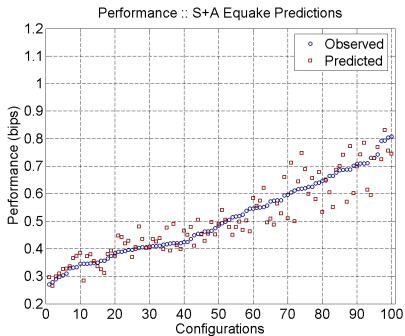
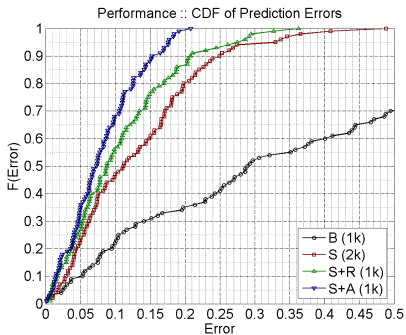
- **Baseline (B)** :: Non-transformed response
- **Stabilized (S)** :: Square-root of response
- **Regional (S+R)** :: Per query with similar samples
- **Application (S+A)** :: Per benchmark with similar samples



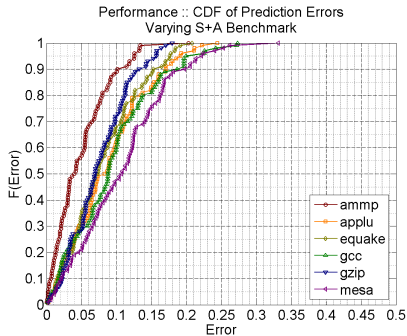
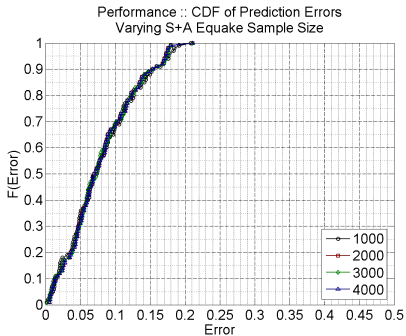
Regional Sampling



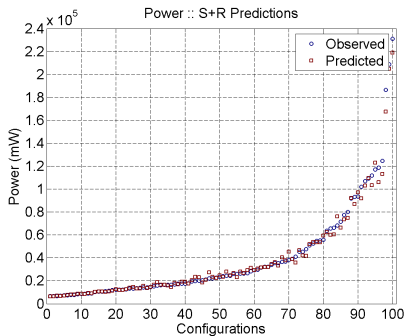
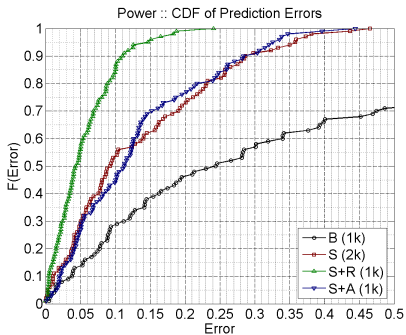
Performance Prediction



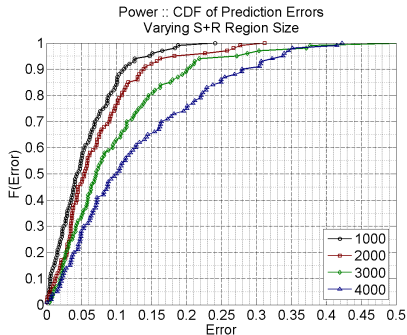
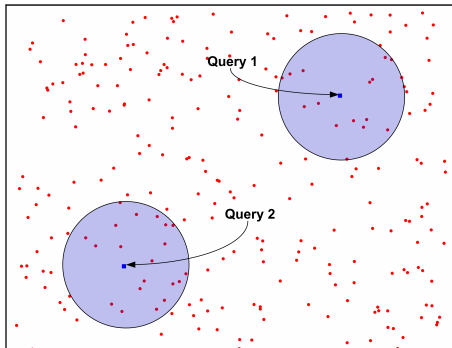
Performance Sensitivity :: S+A



Power Prediction



Power Sensitivity :: S+R Region



Outline

Motivation & Background

Simulation Challenges
Simulation Paradigm
Regression Theory

Model Derivation

Experimental Methodology
Derivation Overview
Model Specification

Model Evaluation

Performance
Power

Conclusion



Conclusion

- **Simulation Paradigm**

- Comprehensively understand design space
- Selectively simulate modest number of designs
- Efficiently leverage simulation data with inference

- **Model Evaluation**

- 7.4%, 4.3% median errors for performance, power
- S+A, S+R more accurate for performance, power

- **Future Directions**





- Demonstrate for comprehensive design studies ⁶
- Expand design space and benchmark suite
- Extend to CMP's and interconnect modeling

⁶Lee [HPCA'07] :: www.deas.harvard.edu/~bclee








Appendix

www.deas.harvard.edu/~bcllee

-  **B.C. Lee and D.M. Brooks.**
Illustrative design space studies with microarchitectural regression models
HPCA-13: International Symposium on High Performance Computer Architecture, Feb 2007.
-  **B.C. Lee and D.M. Brooks.**
Accurate, efficient regression modeling for microarchitectural performance, power prediction.
ASPLOS-XII: International Conference on Architectural Support for Programming Languages and Operating Systems, Oct 2006.
-  **B.C. Lee and D.M. Brooks.**
Statistically rigorous regression modeling for the microprocessor design space.
MoBS-2: Workshop on Modeling, Benchmarking, and Simulation, June 2006.
-  **B.C. Lee and D.M. Brooks.**
Regression modeling strategies for microarchitectural performance and power prediction.
Harvard University Technical Report TR-08-06, March 2006.








References I

-  Y. Li, B.C. Lee, D. Brooks, Z. Hu, K. Skadron.
CMP design space exploration subject to physical constraints.
HPCA-12: International Symposium on High-Performance Computer Architecture, Feb 2006.
-  L. Eeckhout, S. Nussbaum, J. Smith, and K. DeBosschere.
Statistical simulation: Adding efficiency to the computer designer's toolbox.
IEEE Micro, Sept/Oct 2003.
-  R. Liu and K. Asanovic.
Accelerating architectural exploration using canonical instruction segments.
In *International Symposium on Performance Analysis of Systems and Software*, Austin, Texas, March 2006.
-  T. Sherwood, E. Perelman, G. Hamerly, and B. Calder.
Automatically characterizing large scale program behavior.
ASPLOX-X: Architectural Support for Programming Languages and Operating Systems, October 2002.
-  B.C. Lee and D.M. Brooks.
Effects of pipeline complexity on SMT/CMP power-performance efficiency.
ISCA-32: Workshop on Complexity Effective Design, June 2005.



References II

-  C. Stone.
Comment: Generalized additive models.
Statistical Science, 1986.
-  F. Harrell.
Regression modeling strategies.
Springer, New York, NY, 2001.
-  J. Yi, D. Lilja, and D. Hawkins.
Improving computer architecture simulation methodology by adding statistical rigor.
IEEE Computer, Nov 2005.
-  P. Joseph, K. Vaswani, and M. J. Thazhuthaveetil.
Construction and use of linear regression models for processor performance analysis.
In *Proceedings of the 12th Symposium on High Performance Computer Architecture*, Austin, Texas, February 2006.
-  S. Nussbaum and J. Smith.
Modeling superscalar processors via statistical simulation.
In *PACT2001: International Conference on Parallel Architectures and Compilation Techniques*, Barcelona, Sept 2001.



Controlling Simulation Costs

● Hybrid Simulation

- Decouples simulation of microprocessor structures
- Leverages fast, specialized simulators for particular units ⁷

● Trace Sampling/Compression

- Reduces redundant simulation
- Simulate unique, representative instruction segments ⁸

● Synthetic Workloads

- Reduces size of simulator inputs
- Profiles workload to construct smaller, synthetic traces ⁹

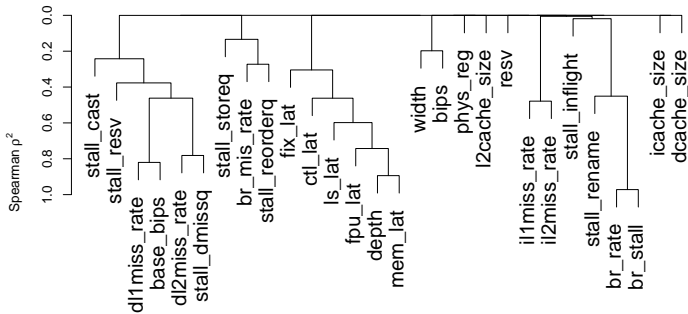
⁷Li, Lee, Brooks, Hu, Skadron [HPCA'06]

⁸Liu, Asanovic [ISPASS'06], Sherwood, *et al.*, [ASPLOS'02]

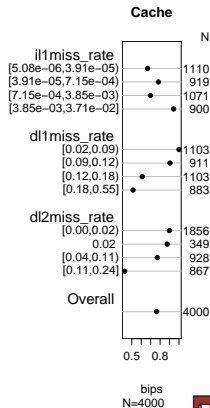
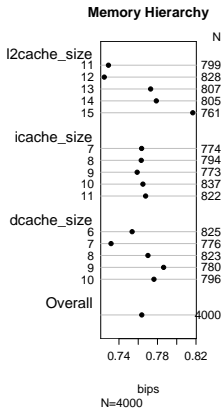
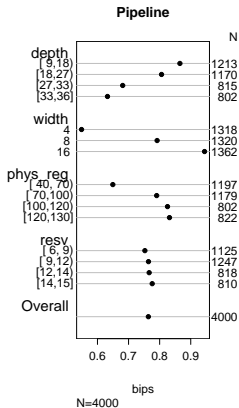
⁹Eeckhout, Nussbaum, Smith, DeBosschre [IEEE Micro'03]



Variable Clustering



Performance Associations



Assessing Fit

● Multiple Correlation Statistic

- R^2 is fraction of response variance captured by predictors
- Large R^2 suggests better fit to observed data
- $R^2 \rightarrow 1$ suggests over-fitting (less likely if $p < n/20$)

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \frac{1}{n} \sum_{i=1}^n y_i)^2}$$

● Residual Distribution Assumptions

- Residuals are normally distributed, $e_i \sim N(0, \sigma^2)$
- No correlation between residuals and response, predictors
- Validate by scatterplots and quantile-quantile plots

$$\hat{e}_i = y_i - \hat{\beta}_0 - \sum_{j=0}^p \hat{\beta}_j x_{ij}$$



Predictor Non-Linearity I

- **Polynomial Transformations**

- Undesirable peaks and valleys
- Differing trends across regions

- **Linear Splines**

- Piecewise linear regions separated by knots
- Inadequate for complex, highly curved relationships

- **Restricted Cubic Splines**

- Higher order polynomials provide better fits
- Continuous at knots
- Linear constraint on tails



Predictor Non-Linearity II

- **Location of Knots**

- Location of knots less important than number of knots
- Place knots at fixed predictor quantiles

- **Number of Knots**

- Flexibility, risk of over-fitting increases with knot count
- 5 knots or fewer are often sufficient ¹⁰
- 4 knots is a good compromise between flexibility, over-fitting
- Fewer knots required for small data sets



¹⁰Stone [SS'86]

Significance Testing I

● Approach

- Given two nested models, hypothesis H_0 states additional predictors in larger model have no response association
- Test H_0 with F-statistics and p-values

● Example

- Predictor interaction requires comparing nested models
- Consider a model $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_1x_2$.
- Test significance of x_1 with null hypothesis $H_0 : \beta_1 = \beta_3 = 0$



Significance Testing II

● F-Statistic

- Compare two nested models using their R^2 and F-statistic
- R^2 is fraction of response variance captured by predictors

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \frac{1}{n} \sum_{i=1}^n y_i)^2}$$

- F-statistic of two nested models follows F distribution

$$F_{k, n-p-1} = \frac{R^2 - R_*^2}{k} \times \frac{n-p-1}{1 - R^2}$$

● P-Values

- Probability F-statistic greater than or equal to observed value would occur under H_0
- Small p-values cast doubt on H_0



Treatment of Missing Data

- **Missing Completely at Random (MCAR)**

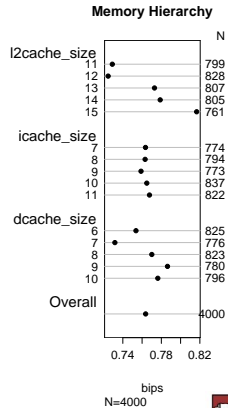
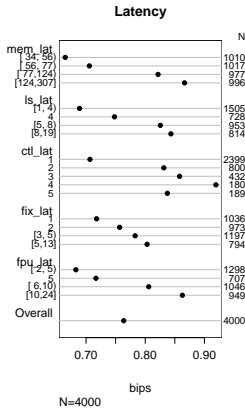
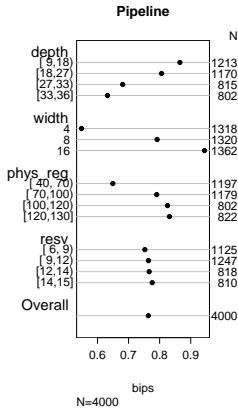
- Treat unobserved design points as missing data
- Sampling UAR ensures observations are MCAR
- Data is missing for reasons unrelated to characteristics or responses of the configuration

- **Informative Missing**

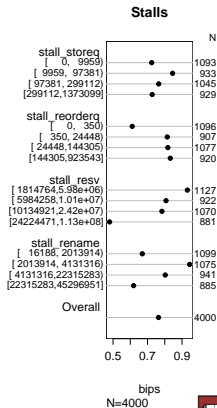
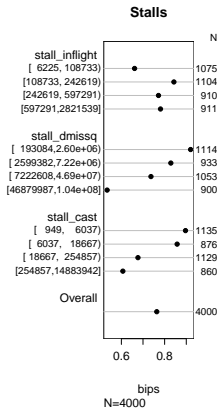
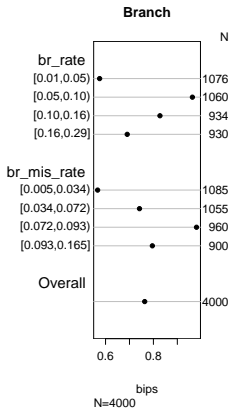
- Data is more likely missing if their responses are systematically higher or lower
- “Missingness” is non-ignorable and must also be modeled
- Sampling UAR avoids such modeling complications



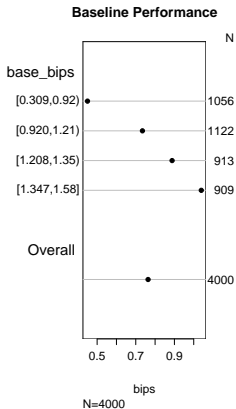
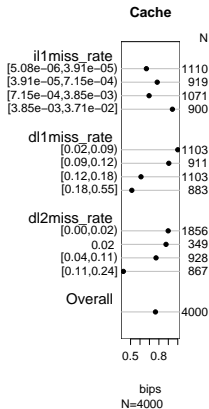
Performance Associations I



Performance Associations II



Performance Associations III



Significance Tests

● Microarchitectural Predictors

- Majority of F-tests imply significance (p -values $< 2.2E - 16$)
- Several predictors were less significant
 - Control latency (p -value = 0.1247)
 - Reservation station size (p -value = 0.1239)
 - L1 instruction cache size (p -value = 0.02941)

● Application-Specific Predictors

- Majority of F-tests imply significance (p -values $< 2.2E - 16$)
- Pipeline stalls classified by structure are less significant
 - Completion and reorder queue stalls (p -values > 0.4)



Related Work

- **Statistical Significance Ranking**

- Yi :: Plackett-Burman, effect rankings
- Joseph :: Stepwise regression, coefficient rankings
- Bound parameter values to improve tractability
- Require simulation for estimation

- **Synthetic Workloads**

- Eeckhout :: Profile workloads to obtain synthetic traces
- Nussbaum :: Superscalar and SMP simulation
- Obtain distribution of instructions and data dependencies
- Require simulation with smaller traces for estimation

