

SELF-OPTIMIZING ARCHITECTURES: TOWARD THE TRUE ONE-SIZE-FITS-ALL CPU

Prof. José F. Martínez

<http://csl.cornell.edu/~martinez>

Computer Systems Laboratory – Cornell University

We're Human after All

2

- Moore's Law: transistor density grows exponentially over time — in the billions now for server CPUs

- Architect's approach: Mostly expert intuition
 - "Intuition growth" highly unlikely exponential
 - "Let's just throw more engineers at the problem"
 - Simplistic, rigid, underperforming average-case solutions
 - Multicore by itself hardly the solution
 - "Why throw more engineers at the problem? Let's just throw in more cores!"

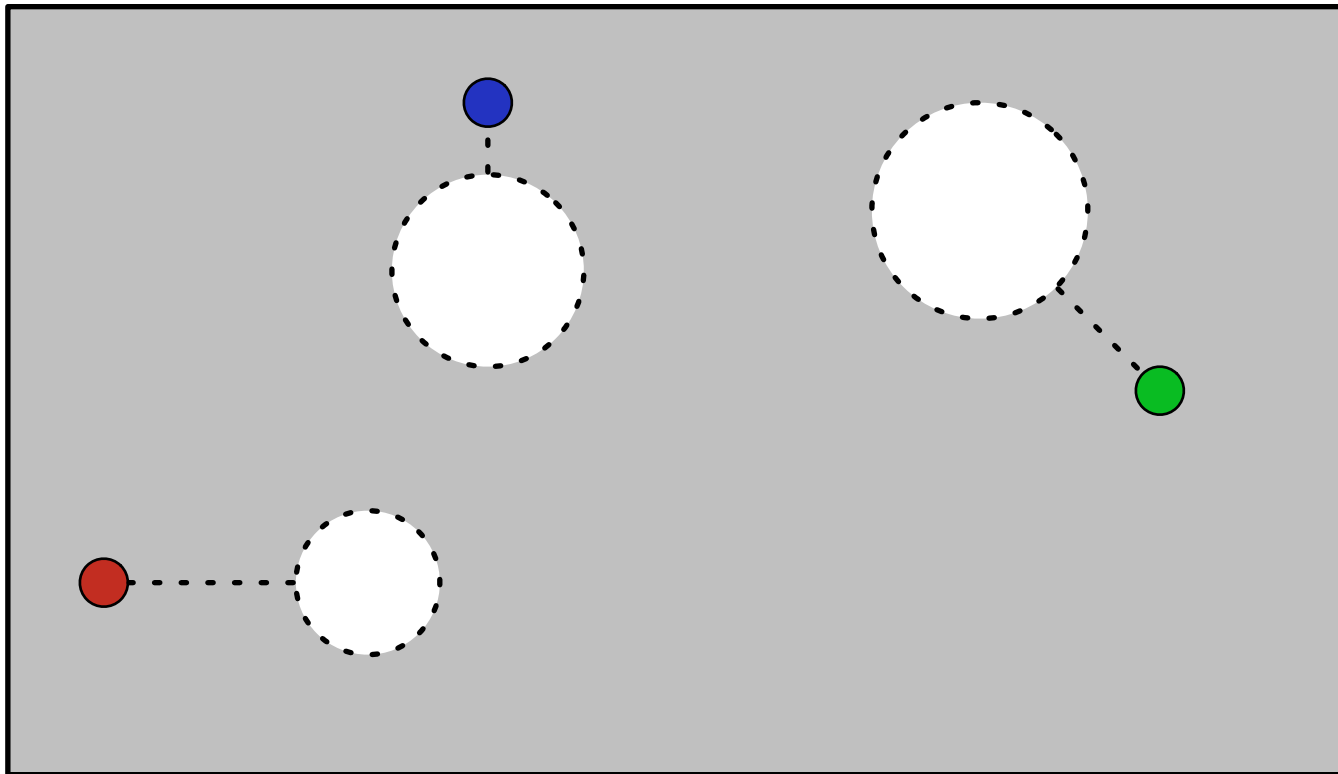
Meanwhile in the Real World...

3

- Data centers backbone of cloud computing
 - Must service wide variety of workloads
 - Continuously changing demands
- Cost is king
 - One-size-fits-all CPU?
 - Works for Intel!
- These seem contradictory
 - They are if using traditional designs

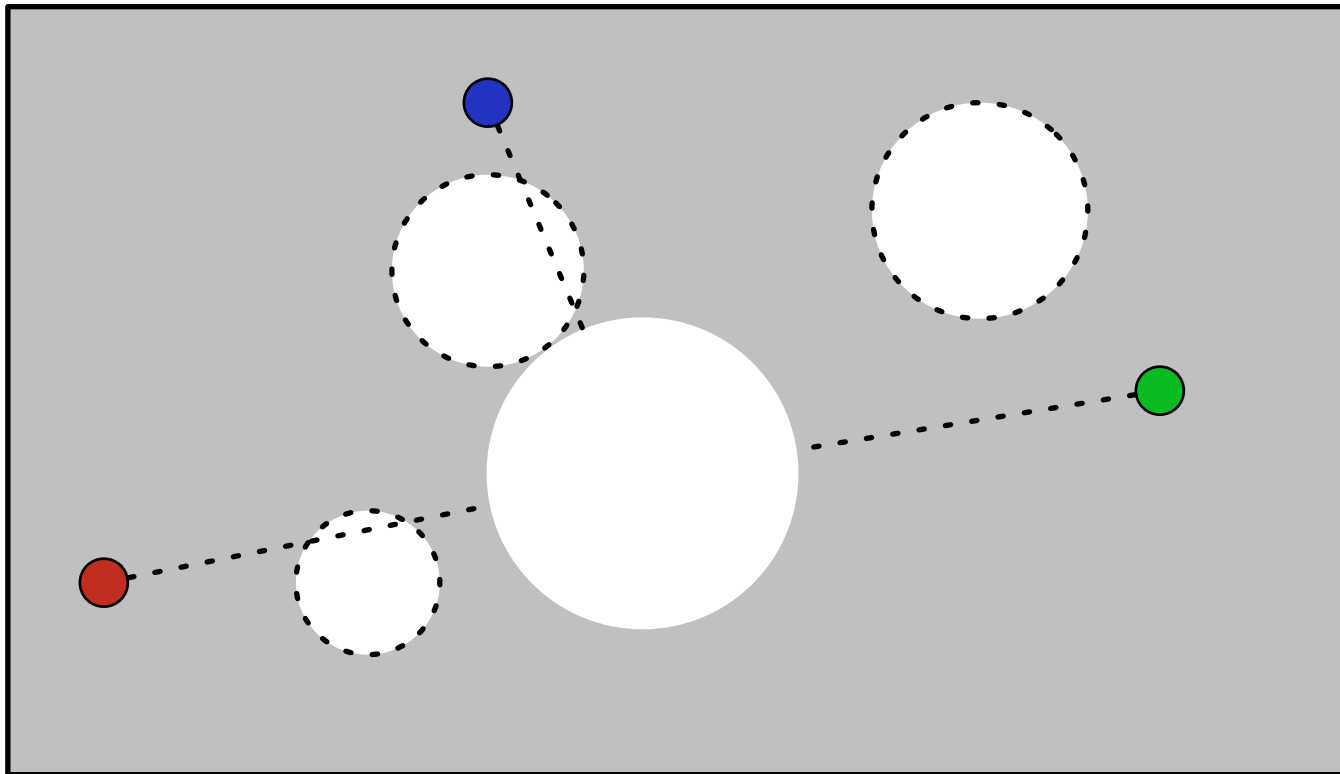
Application-Architecture Space

4



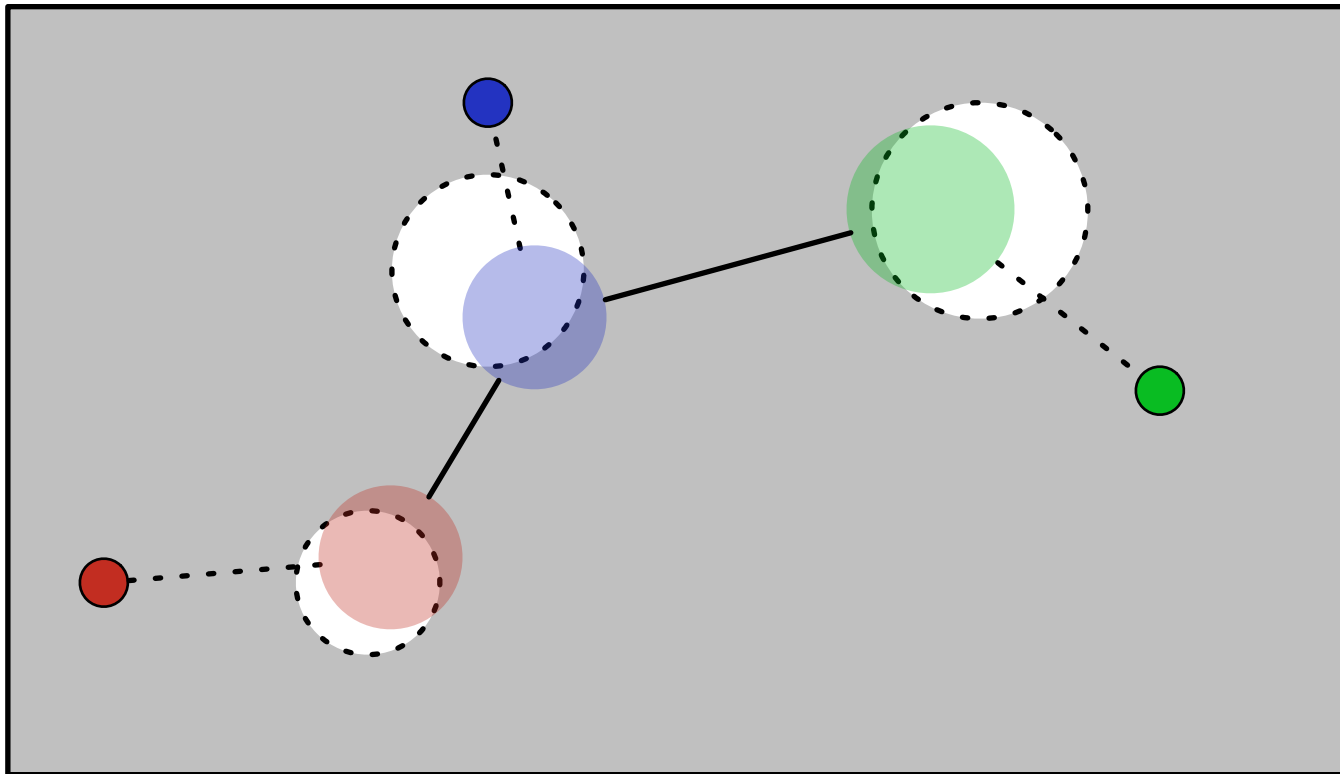
Traditional General-Purpose

5



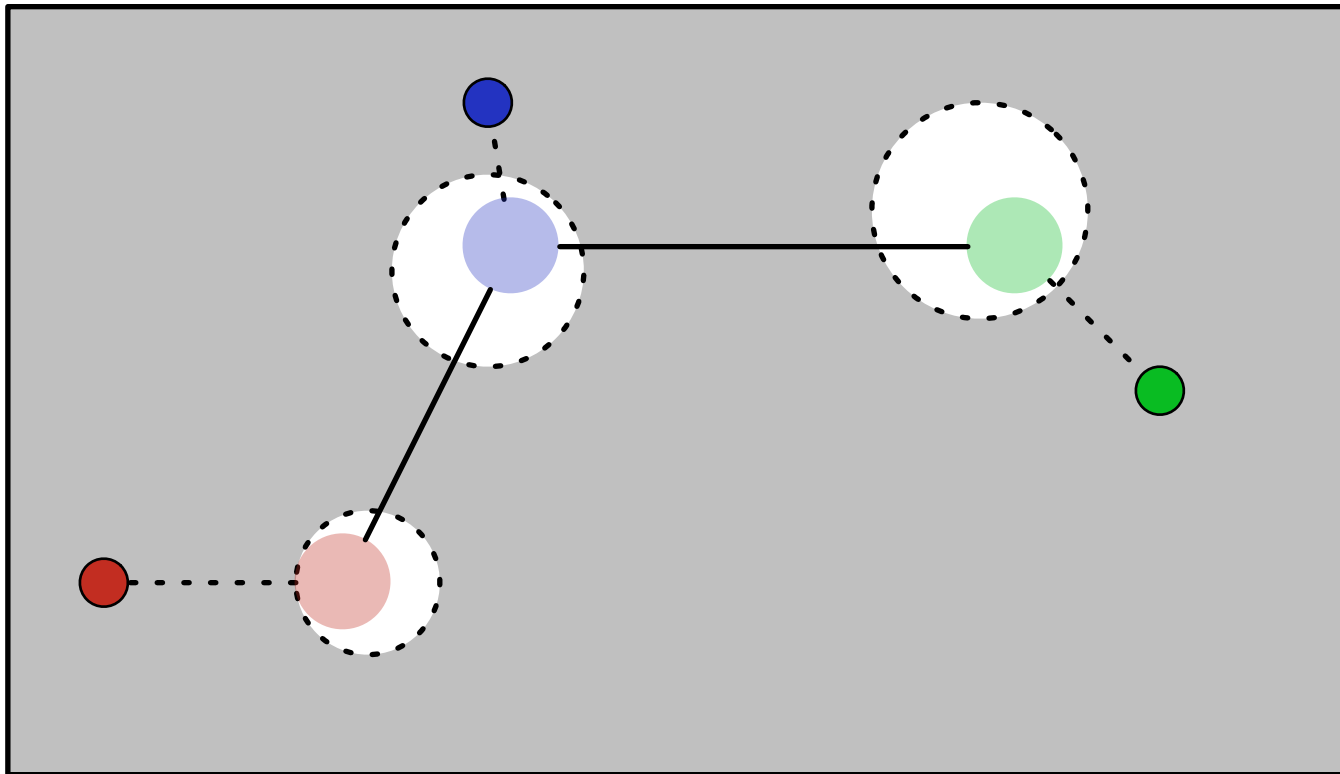
Reconfiguration

6



Self-optimization

7




Self-optimizing Architectures

8

- **Let silicon itself handle the complexity**

- Remember: Transistors *are* plentiful
- Designer focuses on “what”
- Hardware focuses on “how”



Much better use of human ingenuity!

- **Reconfigurable architectures**

- Optimize “baseline”

- **“Gray-box” optimizers** based on machine learning

- CPUs are the “big data” of silicon
 - $2.6 (i+d)/\text{cycle} \times 64\text{b} \times 1.2 \text{ GHz} \approx 200 \text{ Gbps}$

Three Quick Examples

9

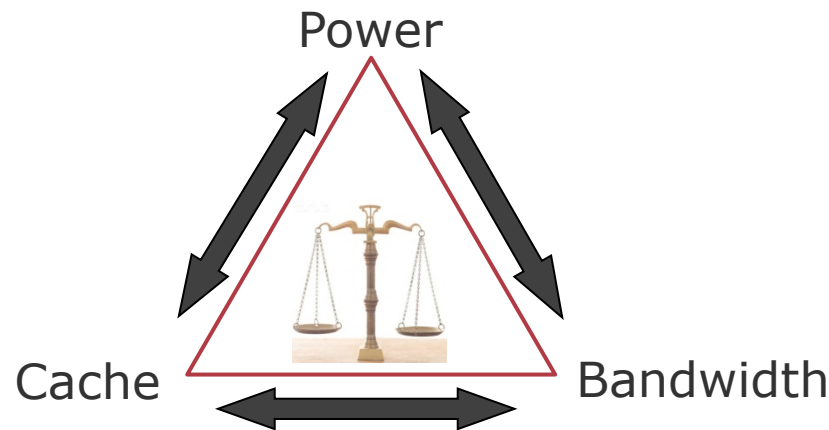
- **Core Fusion:** On-the-fly multicore configuration
 - ▣ Coarse- vs. fine-grain problem goes away [ISCA'07]
- **“Smart” multi-objective memory scheduling:** On-the-fly optimization of memory operations
 - ▣ Employs *genetic programming* and *reinforcement learning* [ISCA'08,HPCA'12]
- **“Smart” multi-resource allocation:** On-the-fly optimization of shared multicore resource utilization
 - ▣ Employs *supervised learning*



Multi-Resource Allocation in Multicore

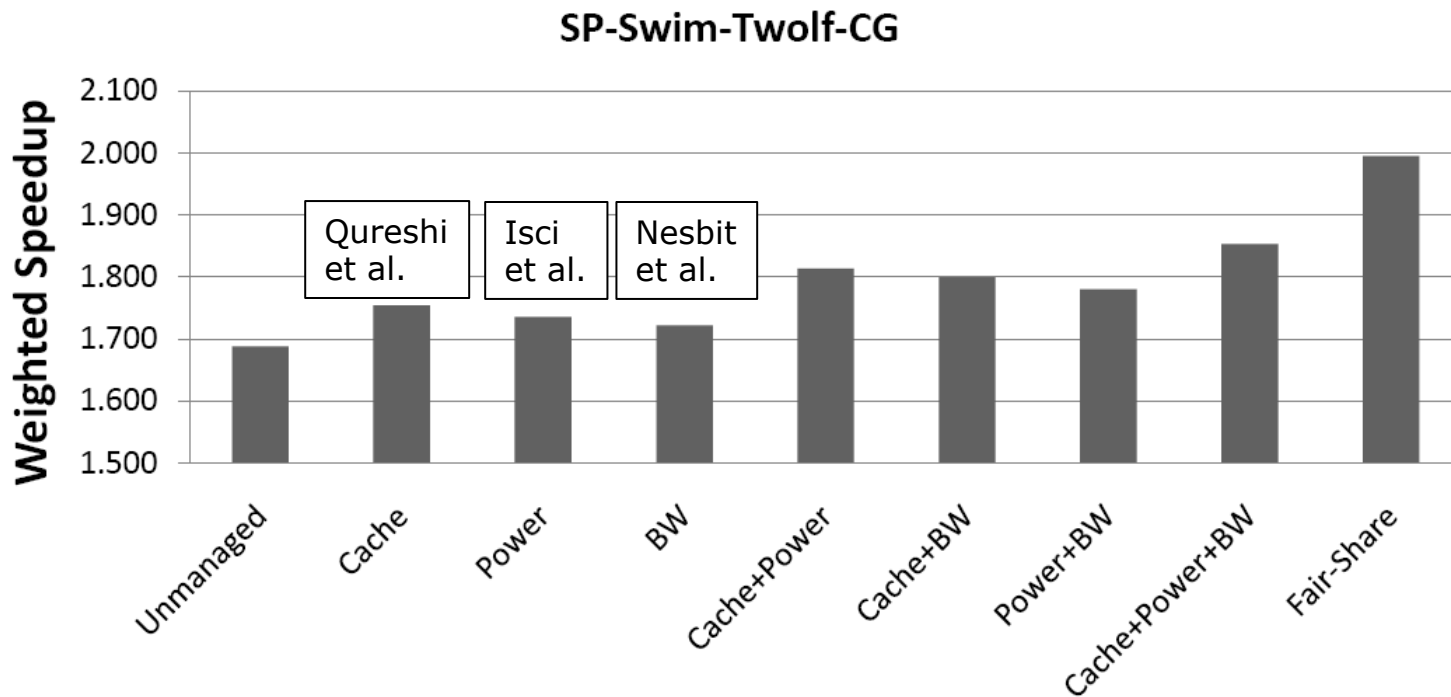
10

- Unrestricted resource sharing can lead to destructive interference (monopoly)
- Several proposals that manage power, cache, or off-chip bandwidth exist but lack complete control over them
- **Need to take resource interactions into account** for optimal performance



Ignoring Interactions: Bad Idea

11



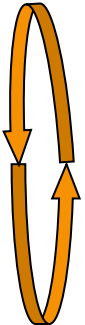
Using Hardware Predictive Models

12

- ❑ Optimizing configuration online via trial-and-error can be costly
 - ❑ Huge search space; learning curve affects performance

- ❑ Proposal: Use “sideline” approach
 - ❑ Monitor applications (sparse sampling)
 - ❑ Learn using on-silicon predictive models (ANNs)
 - ❑ Search for optimum allocation via “on the side” querying
 - Less invasive, much faster
 - ❑ Install resulting allocation periodically

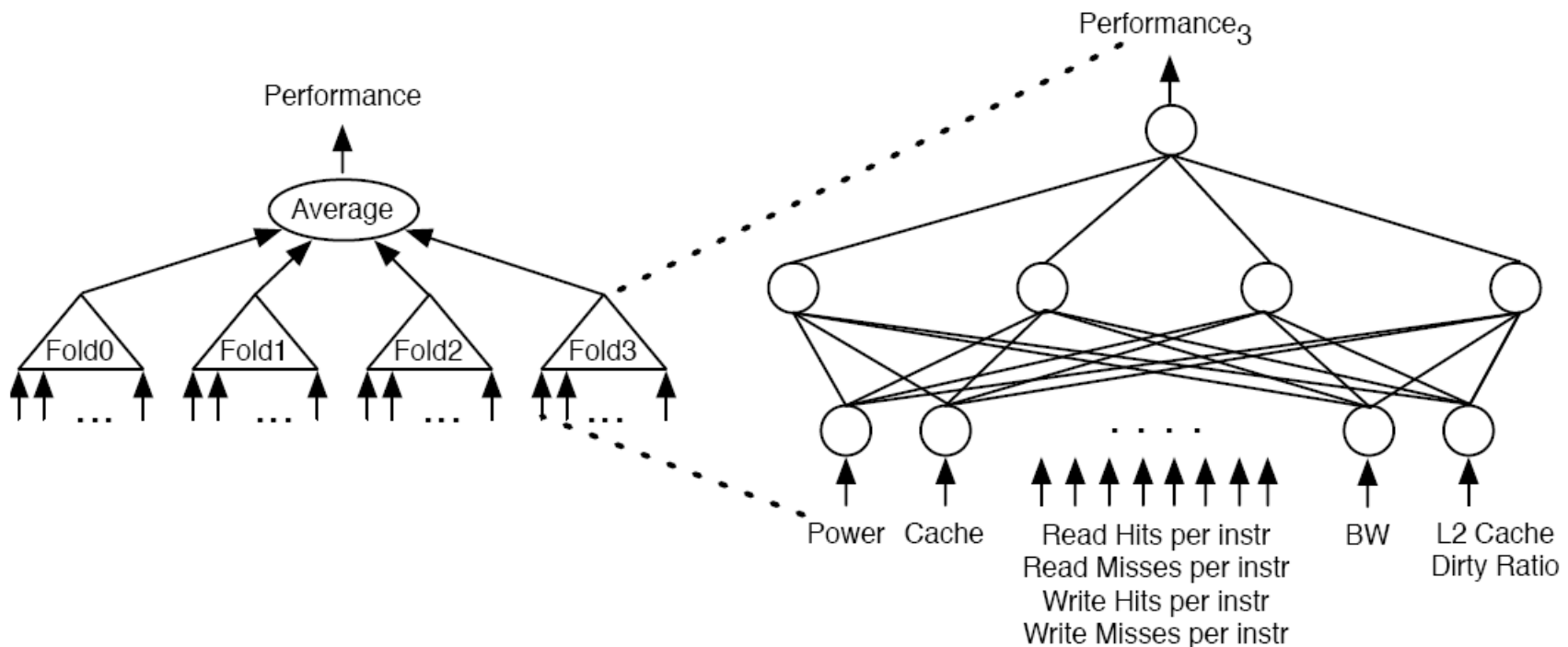
- ❑ **Architect can focus on more productive challenges**



Challenge #1: Find Algorithm

13

- Artificial neural network (ANN) ensemble
 - ▣ Feed appropriate system attributes to model apps
 - ▣ Performance prediction will guide allocation search



Challenge #2: Prediction Error

14

- ANNs not always able to predict with sufficient accuracy
- Identify and discard these cases for better performance
- **Simple mechanism based on variance across four ANNs**
 - Intuition: Large variance → Large error

Challenge #3: Building the Model

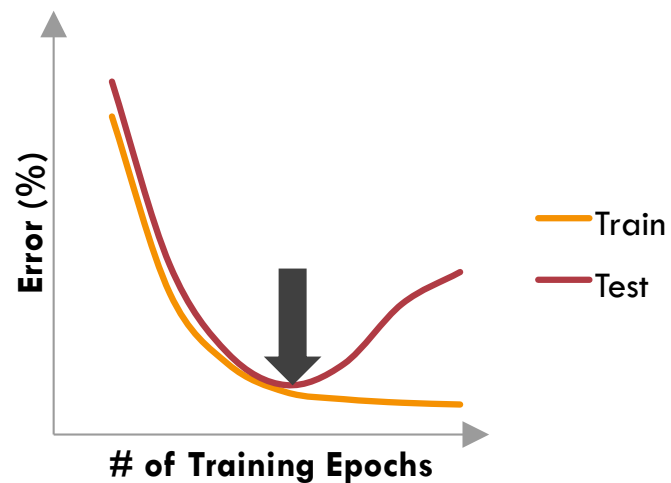
15

- Problem 1: Overfitting
 - Solution: Set aside some data to detect overfitting (**early stopping**)

Data Set:

Train

Test



Challenge #3: Building the Model

16

- Problem 1: Overfitting
 - ▣ Solution: Set aside some data to detect overfitting (**early stopping**)

- Problem 2: Some training data wasted
 - ▣ Solution: **cross validation over the four ANNs**

Data Set:



Data Set:



ANN1:

Train	Train	Train	Test
-------	-------	-------	------

ANN2:

Test	Train	Train	Train
------	-------	-------	-------

ANN3:

Train	Test	Train	Train
-------	------	-------	-------

ANN4:

Train	Train	Test	Train
-------	-------	------	-------

Challenge #4: Search Algorithm

17

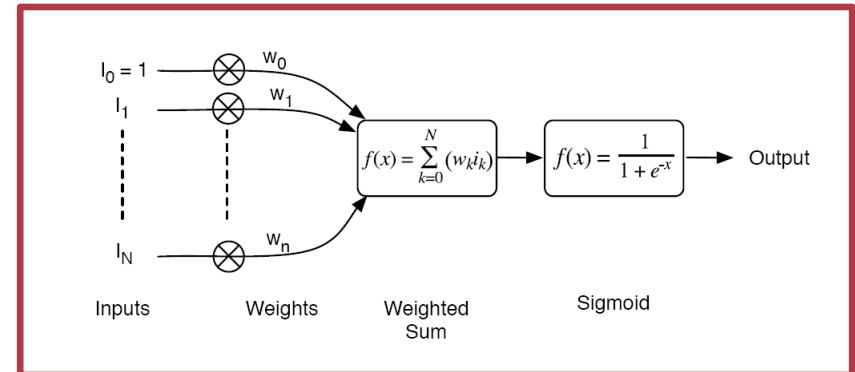
- **Stochastic hill-climbing:** 2,000 steps at start of each interval
 - Estimate performance for all neighbors of current optimum
 - Too many neighbors to try → Move to a better neighbor immediately
 - Try random points with 20% chance
 - Move there if it is better than the best allocation found so far
 - Detect if stuck in local minima
 - Jump to a random point and resume HC

Challenge #5: Implementation

18

- Impractical to build all 16 ANNs in HW
 - ▣ Have 4 applications in our setup
 - ▣ Each application needs 4 ANNs

- **Build single ANN and time-share it**
 - ▣ Narrow-width, pipelined multipliers
 - ▣ Lookup tables for sigmoid function approximation

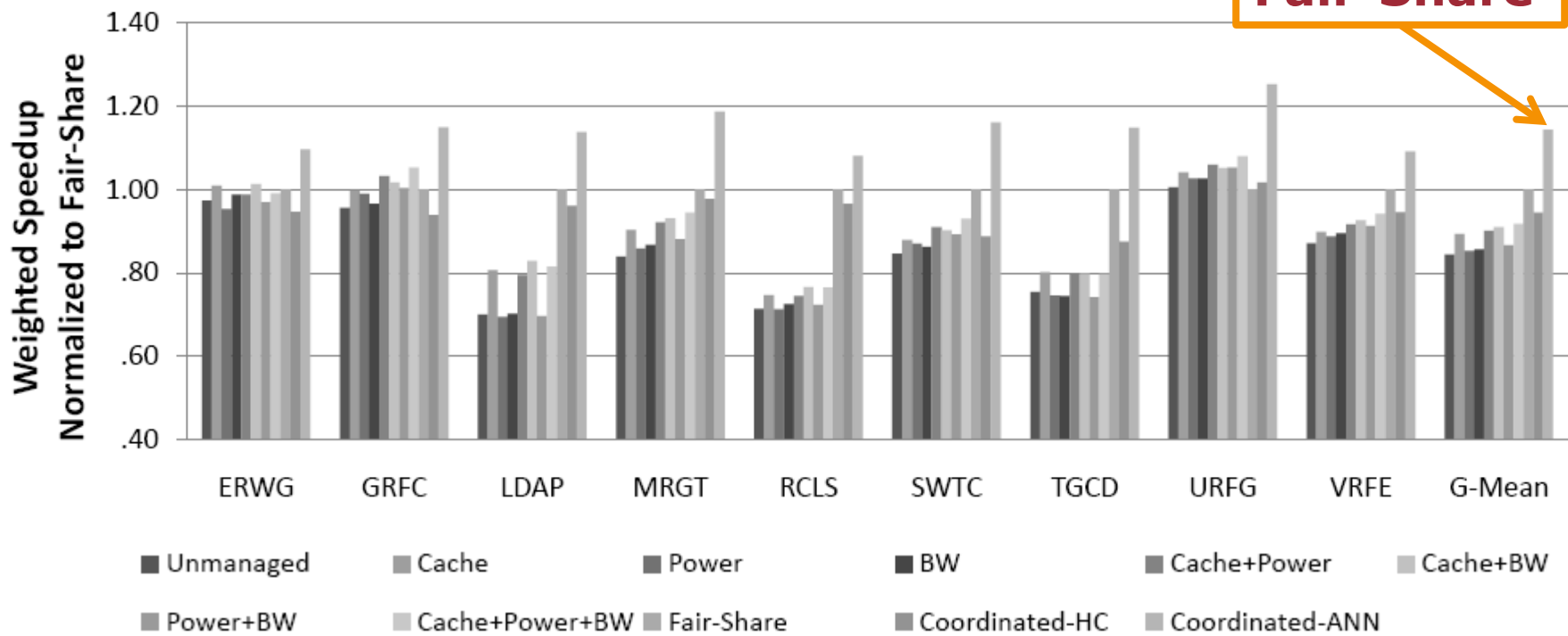


- Area overhead: $\sim 6\text{mm}^2$
- Power overhead: $\sim 3\text{W}$

Performance Results: 4 Cores/Apps

19

**1.14 over
Fair-Share**



Bitirgen et al., MICRO '08

Self-optimization: Takeaway Slide

20

- General-purpose doesn't mean average
 - ▣ Nothing more general-purpose than humans!
 - ▣ Key: ability to specialize over time

- Transistors are there; put them to work *smart*
 - ▣ Give transistors a “higher education”
 - Then watch them do great things on their own
 - ▣ Leverage architect's intuition where it matters
 - I.e., don't re-invent the ANN !