# Insight, not (random) numbers

**Tom Conte**

**NC State University**

---

# The Quotation

**"The purpose of computing is insight, not numbers"**

**- Richard Hamming (1915-98)**

*from Numerical Methods for Scientists and Engineers, 1962*

# What did Hamming mean ca. 1962?
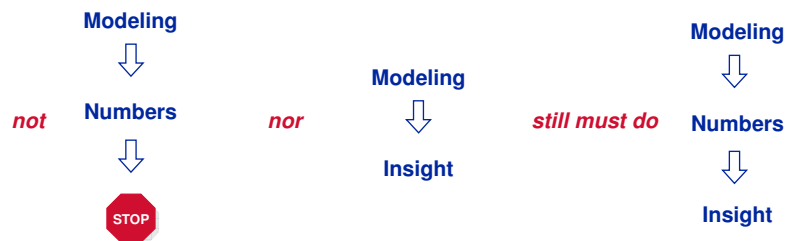
○ **Context: Originally, computing was this:**

Computation

"Computing" = Number crunching
Why number crunch?
*Simulation-based modeling*

---

# Some quotation flowcharting…

**The purpose of  modeling  is insight, not numbers**

*not*   Modeling
⇩
Numbers
⇩
**STOP**

*nor*   Modeling
⇩
Insight

*still must do*   Modeling
⇩
Numbers
⇩
Insight

○ **Hamming was saying two things then:**

1. *Develop a method to gain insight from numbers*

2. *And,* **guarantee the** *quality* **of the numbers so you have a hope of gaining insight!**

# The current state of affairs: On insight from numbers

- ○ **We have "modeled" this car and determined:**
  - ◆ 0-60 time in 4.9 sec, 500 horsepower and 383 lb.-ft. of torque. Engineered to rev with a redline of 8,250, Top speed of 205 mph

**2005 BMW M6**



**(It's $160,000)**

**(plus tax)**

**So??? Does it fit my needs?**

**What *are* my needs? (and who's 'me', for that matter– driver, designer or salesman?)**

---

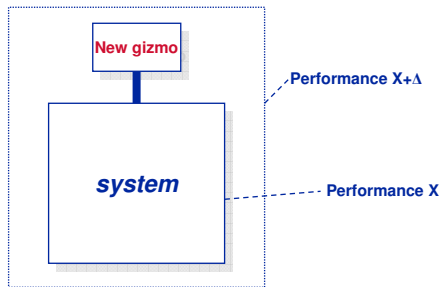# On the road to insight… who consumes the numbers?

- ○ **In the case of a Bimmer, "number consumers" are:**
1. **The buyer (the obvious)**
2. **The marketer/salesman (the parasitic)**
3. **The car designer (the noble?)**

- ○ **And in the case of computers? The same three suspects: buyer, marketer, architect**
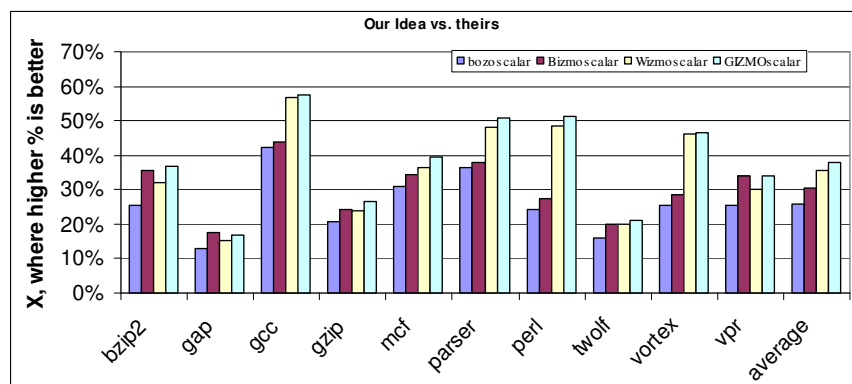- ○ **Let's take these on one at a time (…in reverse order)**

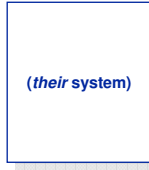# What are we trying to do with modeling?

○ **View #1: *The architect***



**New gizmo**

Performance X+Δ

*system*

Performance X

**If Δ > 0, then my idea is a good idea…**

---

# Gizmoscalar *The future of Computer Architecture*



**Our Idea vs. theirs**

■ bozoscalar ■ Bizmoscalar □ Wizmoscalar □ GIZMOscalar

X, where higher % is better — y-axis: 0% to 70%

x-axis: bzip2, gap, gcc, gzip, mcf, parser, perl, twolf, vortex, vpr, average

4

# What are we trying to do with modeling?

○ **View #2:** *Marketing*

| | |
|---|---|
| (*their* system) | *Our system* |
| **Performance X** | **Performance X+ Δ** |

**If Δ > 0 Buy our system**

**Well, not really...**

---

5

## Slide 1

DELL HOME DESKTOPS
Features Comparison

All Dell Dimension™ desktops are customizable and can be configured to contain the following features. Features indicated with an "X" are recommended as being "great for" each system.

Click on the name of the desktop below for product details.

| | Choose Dimension 8250 | Choose Dimension 4550 | Choose Dimension 2350 |
|---|---|---|---|
| E-mail/Internet | X | X | X |
| Word Processing | X | X | X |
| Microsoft® Excel/Powerpoint | X | X | X |
| Educational Software | X | X | X |
| Digital Imaging/Photography | X | X | X |
| Digital Video Editing | X | X | |
| Digital Audio | X | X | |
| Intense Gaming | X | X | |
| Broadband Internet | X | X | |
| Voice Recognition | X | | |
| Advanced Multimedia | X | | |
| Cutting Edge Graphics | X | | |

HOME & HOME OFFICE

Buy Online or Call
1-800-915-3355

Computers | Software & Peripherals | Service & Support | Learning Center

**Insight for free
No numbers required!**

---

## Slide 2

# What are we trying to do with modeling?

**Which one?**

⦾ **View #3: *The (smart) users***

*My favorite application*

**Problem is: Who will run his favorite application on Systems #1-4?**

| System #1 | System #2 | System #3 | System #4 |
|---|---|---|---|

**Performance $X_1$**     **Performance $X_2$**     **Performance $X_3$**     **Performance $X_4$**

**Pick max of $X_i$ …**

6

# How are these done?
## _Benchmarking_

- **The purpose of benchmarking then depends on who you talk to:**
  - ◆ **The architect: Prove _my_ gizmo is _great!_**
  - ◆ **Marketing: Make us look good to _sell $$_ and _crush our competition, get enough commission to buy the red bimmer …_**
  - ◆ **The users: Be <u>our proxy</u>, <u>run our applications</u> on new systems so we don't waste our money or our time**
- **…For the purposes of this talk, we can safely ignore the marketing purpose**

---

# Updated road map



- **If benchmarks are good proxies**
- **And the numbers match the benchmarks**
- **Then…**
  - ◆ _No magic required!_

## Part I: "If benchmarks are good proxies" When are they good proxies?

- Which benchmark do I believe?
- Answer: the one that is closest to what I do
- Question: Which one is that?
- Answer: Read the descriptions

## What's in SPEC…

| Benchmark | Description |
|---|---|
| 164.gzip | Lempel-Ziv data compression algorithm |
| 175.vpr | FPGA place and route tool (combinatorial optimization) |
| 176.gcc | GNU C compiler |
| 181.mcf | Single-depot vehicle scheduling solver (combinatorial optimization) |
| 186.crafty | Computer chess game |
| 197.parser | Link Grammar Parser (word processing) |
| 252.eon | Probabilistic ray tracer (computer visualization) |
| 253.perlbmk | Perl programming language |
| 254.gap | Language and library implementation for group computing (group theory) |
| 255.vortex | Single-user object-oriented database |
| 256.bzip2 | Seward compression algorithm, occurs entirely in memory |

- Say what you do is:
1. Surf the web
2. Database accesses
3. Logic simulation
4. CAD synthesis

- Which benchmark is the right one to listen to?

8

# What's in Mediabench…

| Benchmark | Description |
|---|---|
| JPEG | Jpeg compression/decompression |
| MPEG | Decoding mpeg-1 and mpeg-2 video streams |
| GSM | Speech transcoding using RPE/LTP coding at 13kbits/s |
| ADPCM | Adaptive Differential Pulse Code Modulation algorithm for speech compression/decompression |
| G.721 | CCITT G.711, G.721, and G.723 voice compressions |
| PGP | Public key encryption and authentication |
| PEGWIT | Public key encryption and authentication |
| SPHERE | Read and format NIST-formatted speech waveforms |
| RASTA | Filtering for speech recognition |
| Ghostscript | Postscript language interpreter, postscript graphics generation, PDF |
| Mesa | 3D graphics library |
| EPIC | Image compression |

○ **And which one here matches what you do with your cellphone?**

---

# A better way: Quantitative Benchmark Characteristics

**Some examples:**

○ **IPC (with large memory system)**

○ **Branch predictability (for gshare)**

○ **Preferred L1 instruction cache size**

○ **Preferred L1 data cache size**

○ **Preferred L2 unified cache size**

○ **Total virtual memory requirements (4KB page size)**

○ **Others:**

   ◆ **TLB requirements**

   ◆ **Instruction frequency by type**

   ◆ **System Call usage**

   ◆ **…**

# Let's try it out…

**Consider these benchmark sets:**
- **MediaBench (UCLA)**
- **NetBench (UCLA and NWU)**
- **SPEC CPU CINT2000**

---

# Kiviat view

# Kiviat view

**Heavy memory usage
Moderate branchiness**

---

# Kiviat view

**Moderate memory usage
Low branchiness
High parallelism**

11

**ADPCM**

**MPEG2 ENCODE**

*Hard to predict branches*
*Moderate IPC*
*Very small cache needs*

**GZIP**

**MCF**

**JPEG**

**UNEPIC**

**TWOLF**

**G721 DECODE**

**G721 ENCODE**

*Hard to predict branches*
*Moderate-high IPC*
*Small cache needs*

12

## Slide 25 (BZIP, VORTEX, PARSER, GSM TOAST, GCC)

**BZIP**
IPC *10
100
25.4
10
Memory
22
Branch misprediction ratio
4.71
12
10
10
L2 unified Cache FA
L1 Icache FA
L1 Dcache FA

**VORTEX**
IPC *10
100
24.2
10
Memory
20
Branch misprediction ratio
3.93
7
15
11
11
L1 Icache FA
L1 Dcache FA

**PARSER**
IPC *10
100
18
10
Memory
22
Branch misprediction ratio
3.16
13
L2 unified Cache FA
12
11
L1 Icache FA
L1 Dcache FA

**GSM TOAST**
IPC *10
100
34.1
10
Memory
18
Branch misprediction ratio
3.59
14
12
L2 unified Cache FA
8
L1 Dcache FA

**GCC**
IPC *10
100
24
10
Memory
20
Branch misprediction ratio
2.86
17
FA
11
14
L1 Icache FA
L1 Dcache FA

*Somewhat hard branches*
*Moderate IPC*
*Small cache needs*

## Slide 26 (EPIC, MPEG2 DECODE, GAP, VPR)

**EPIC**
IPC *10
100
30.3
10
Memory
18
Branch misprediction ratio
2.8
10
L2 unified Cache FA
10
L1 Icache FA
L1 Dcache FA

**MPEG2 DECODE**
IPC *10
100
31
10
Memory
18
Branch misprediction ratio
2.38
11
L2 unified Cache FA
10
L1 Icache FA
L1 Dcache FA

*Somewhat hard branches*
*High IPC*
*Very Small cache needs*

**GAP**
IPC *10
100
32.3
10
Memory
19
Branch misprediction ratio
1.58
14
9
L2 unified Cache FA
8
L1 Icache FA
L1 Dcache FA

**VPR**
IPC *10
100
19.1
10
Memory
21
Branch misprediction ratio
1.5
14
12
L2 unified Cache FA
11
L1 Icache FA
L1 Dcache FA

13

**GSM UNTOAST**

**OSDEMO**

*Easy branches*
*High IPC*
*Very Small cache needs*

**TEXGEN**

**URL**

---

**DRR**

**PERL**

*High branches*
*High IPC*
*High cache needs*

**ROUTING**

**AES**

14

# What's it all mean?

1. **Benchmarks have distinct characteristics**
2. **Some benchmarks are similar…**
   - **Across different benchmark suites**
   - **Across different application domains**
- **Bottom line: *There's hope* (!) that characteristics can be used to guarantee a benchmark is a good proxy**

---

# The missing piece:
# *How to make a million dollars*

- **Create a tool that runs unobtrusively**
- **The tool collects statistics about *usage* characteristics**
- **So you know then which benchmark to choose as your proxy …**

- **Taking it a step further, the tool finds the benchmarks that have those characteristics**



Hey, you need Mediabench "GSM Untoast"

# But what about benchmark *suites?*

- Who creates benchmark suites today?
  - ◆ Mostly industry
  - ◆ Why? *Marketing!*
- Who speaks for the users?
  - ◆ They do. *Trust them.*
- A modest proposal:
  - ◆ Poll what users care about
  - ◆ Create benchmarks for them
  - ◆ Have an impartial panel select among these based on quantitative characteristics
  - ◆ Use this to create a benchmark suite
  - ◆ Rigorously review the suite every year
  - ◆ IMHO, better suited to academia than industry

---

# Back to the road map

Benchmarks ← = My proxy
⇩
*do they match?*
Modeling
⇩
Numbers
⇩
Insight

- "If benchmarks are good proxies"
  - ◆ How? Use benchmark characterization
- "<u>And the numbers match the benchmarks</u>"
- Then... *insight*

**Part II:** **"And the numbers match the benchmarks"**
**Easy!** *Just Simulate The Benchmarks!*

- **Problem…Simulation takes** *time*
- **We can't quickly simulate and get an accurate number!**
- **Solution: Don't simulate the entire benchmark**
- **How?**
  - **who cares.**

---

# Shame on us

- **Skip 100million, simulate 1 billion**
- **Skip 1billion, simulate 100million**
- **Skip 1billion, simulate 1billion**
- **Change the inputs**
- **Change the benchmarks**
- **Use only benchmarks that show my gizmo shines**
- **And my favorite… Skip benchmarks that crash or don't compile …**
- ***How good are these numbers? How much can you trust them?***

## Why don't we include error bars?
## I think I know why:
## Gizmoscalar revisited

**Our Idea vs. theirs**



*(by the way, this is from one of my/my students' papers)*

---

# Better way: Sampling

○ **How to predict who will be the next president of the US:**
   ◆ **Solution #1: Ask all Americans**
      ◇ **Takes too long**
   ◆ **Solution #2: Ask random Americans**
      ◇ **Which ones?  Be careful!  (e.g., not just TX… not just Austin, TX)**
○ **Saves a lot of work …**
   ◆ **Pick random pieces of a benchmark *trace* and simulate only those**
○ **Great idea!**
   ◆ **History of sampling for fast architecture simulation:**
   1. **Original credit due to Laha, Patel, lyre in 1988**
   2. **Early work for cache sampling only – Kessler, Fu, …**
   3. **Processor sampling work ca. 1992 and onwards – Menezes, Poursepanj, …**
   4. **Latest work on whole system sampling – cast of thousands**

# But can you trust sampling?

- How accurate is your sampling?
  - ◆ Silly question!
  1. Run the sampling trace, say get $X_{sample}$
  2. Run the full trace, say get metric $X_{true}$
  3. Error is just $(X_{sample} - X_{true})/ X_{true}$
  - ◆ Simple!
- Not really, of course
  - ◆ To get $X_{true}$ to calculate error, *you didn't save any work*, it requires a full simulation! …
- Or just use the published error to find your error bars

---

# Getting error bars

**Benchmarks**

**Published sampling regimen** → *My system simulation*

**published error**      **Numbers +/- published error**

- "Trust my error"
- But we can do better
  - ◆ Sampling theory allows calculation of error (confidence intervals) a priori using Student-t statistics

# Trace sampling according to sampling theory



- ○ **"Trace sampling" = cluster sampling:**
  - ◆ *n* clusters of *m* execution cycles each
  - ◆ Actually *indirectly sampled*: *n* clusters of *m'* instructions each
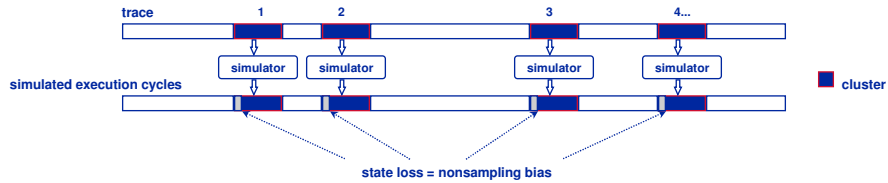- ○ **Error due to three effects:**
  - ◆ sampling bias (e.g., ask 3 people)
  - ◆ sampling variability (e.g., ask only people in TX)
  - ◆ nonsampling bias (e.g., ask people in Canada who their friend will vote for)
    - ◊ If reduced, sampling theory applies, error bars can be calculated!

---

# Nonsampling bias

- ○ **Nonsampling bias due to *indirect sampling***
  - ◆ The measured population is different from the actual
  - ◆ For us: *System state is unknown at start of each cluster simulation …*
- ○ **If eliminated, then**

$$\text{metric}_{\text{true}} = \text{metric}_{\text{sample}} \pm 1.96 \cdot S, \quad \text{for 95\% confidence interval}$$

$$\text{standard error,} \quad S = \frac{s_{\text{metric}}}{\sqrt{N_{\text{cluster}}}}$$

$$\text{standard deviation,} \quad s_{\text{metric}} = \sqrt{\frac{\sum_{i=1}^{N_{\text{cluster}}} \left(\text{metric}_{\text{cluster}}^i - \text{metric}_{\text{sample}}\right)^2}{(N_{\text{cluster}} - 1)}}$$
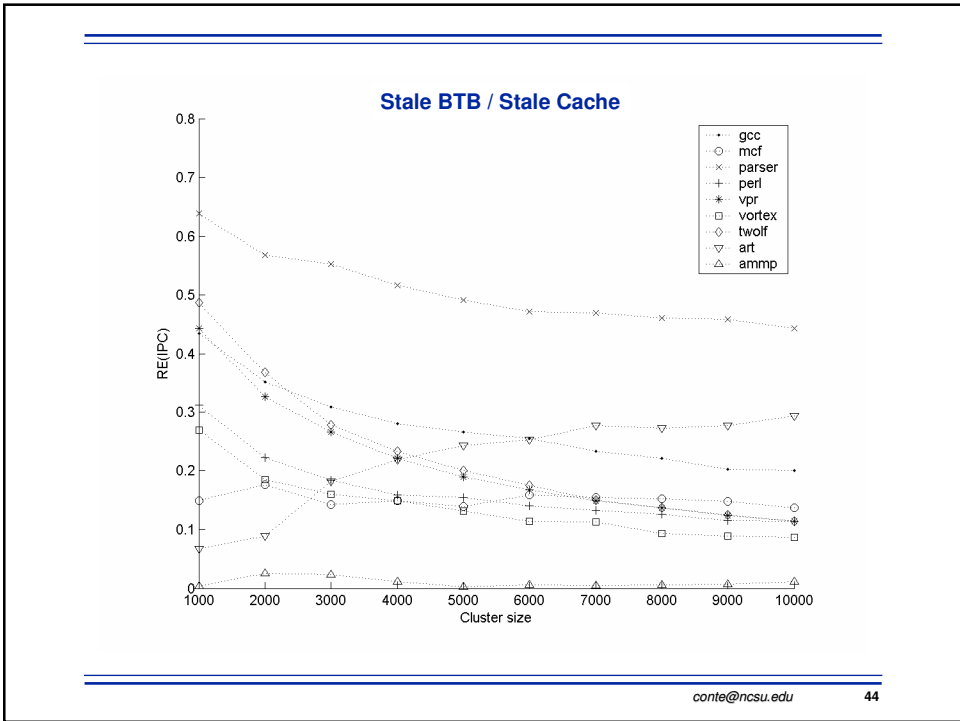
20

# Sampling bias and variability

- **Sampling bias reduced via three parameters:**
  - ◆ **Cluster size**
  - ◆ **Number of clusters**
  - ◆ **Overall sample size = Cluster size X Number of clusters**
- **Sampling variability improved via random sampling**
  - ◆ **Gaps between clusters are selected using random variable of uniform distribution**
  - ◆ **If you get this wrong, error bars may be too tight**

---

# Example

- **SPEC CPU Cint2000**
- **4-issue, 64-entry window**
- **L1: 32KB, 4-way, L2: 1MB, 8-way**
- **Memory bus contention modeled**
- **BTB: 64k-entry gshare, 1k-entry ret addr stack**
- **Four nonsampling bias removal choices:**
  - ◆ **Leave BTB stale / simulate it during the gap (warm)**
  - ◆ **Leave caches stale / simulate it during the gap (warm)**

**Finding the required number of clusters**

~1000 clusters is reasonable



**Stale BTB / Stale Cache**

22

**Instead, use 10% of the cluster to warm up BTB, caches**

**Using *90%* of the cluster to warm up BTB, caches**

23

**Warm BTB / Stale Cache**

**Stale BTB / Warm Cache**

*Need to keep the caches warm!*

24

Warm BTB / Warm Cache

Best technique for nonsampling bias removal

---

# Finding the right clustersize: Do the statistics predict the actual error?

yes if sampled +/- CI = actual

| Cluster size | 1000 | 2000 | 3000 | 4000 | 5000 | 6000 | 7000 | 8000 | 9000 | 10000 |
|---|---|---|---|---|---|---|---|---|---|---|
| gcc | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| mcf | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| parser | yes | no | yes | no | yes | no | yes | yes | yes | yes |
| perl | yes | yes | yes | yes | no | no | no | no | no | no |
| vortex | no | yes | yes | no | no | no | no | no | no | no |
| vpr | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| twolf | yes | no | yes | yes | yes | yes | yes | yes | yes | yes |
| ammp | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| art | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |

- **Clustersize = 3000 …**
- **(Why is vortex so difficult?)**

25

## Consider the characteristics

**VORTEX**



- ○ **High branch misprediction rate**
- ○ **Moderate memory footprint, but it's enough …**

## But also low cluster variability

- ○ **The lower the variability**
- ○ **…the smaller the standard error**
- ○ **…the tighter the predicted confidence interval**
- ○ **Thus more stress is placed on nonsampling bias removal …**
- ○ **Some benchmarks are tougher than others**

26

# Bracketing the error

| benchmark | True mean $\mu_{IPC}^{true}$ | Estimated mean $\mu_{IPC}^{sample}$ | Standard error $S_{IPC}$ | 95% Error bound CI | Absolute error $\lvert\mu_{IPC}^{true} - \mu_{IPC}^{sample}\rvert$ |
|---|---|---|---|---|---|
| gcc | 0.87314 | 0.89178 | 0.02263 | ±0.04436 | 0.01864 |
| mcf | 0.20854 | 0.22202 | 0.01999 | ±0.03918 | 0.01348 |
| parser | 1.07389 | 1.05273 | 0.01343 | ±0.02632 | 0.02116 |
| perl | 1.28956 | 1.28458 | 0.00761 | ±0.01493 | 0.00498 |
| vpr | 1.18062 | 1.17164 | 0.00601 | ±0.01178 | 0.00898 |
| vortex | 0.92672 | 0.92415 | 0.00487 | ±0.00955 | 0.00257 |
| twolf | 0.97398 | 0.97523 | 0.00599 | ±0.01175 | 0.00125 |
| art | 0.77980 | 0.78220 | 0.01816 | ±0.03560 | 0.00240 |
| ammp | 0.24811 | 0.24390 | 0.02740 | ±0.05371 | 0.00421 |

- **Cluster size = 3000, 1000 clusters**
- **Warm / Warm nonsampling bias removal**
- **The confidence intervals predicted the empirical error!**

---

# Is it worth it?  How much speedup?

| benchmark | full sim time (min) | sampled sim time (min) | percentage speedup |
|---|---|---|---|
| gcc | 743 | 46 | 16.2 |
| mcf | 5776 | 66 | 87.5 |
| parser | 675 | 63 | 10.7 |
| perl | 682 | 86 | 7.9 |
| vpr | 613 | 37 | 16.6 |
| vortex | 929 | 113 | 8.2 |
| twolf | 706 | 38 | 18.6 |
| art | 511 | 35 | 14.6 |
| ammp | 3665 | 58 | 63.2 |

- **From 8x to 87x speedup**
- **~1 to 2 hours per benchmark**
- **This would improve with better/more efficient nonsampling bias removal techniques**

# We had this: "Trust my error"

**Benchmarks**

**Published sampling regimen** → **My system simulation**

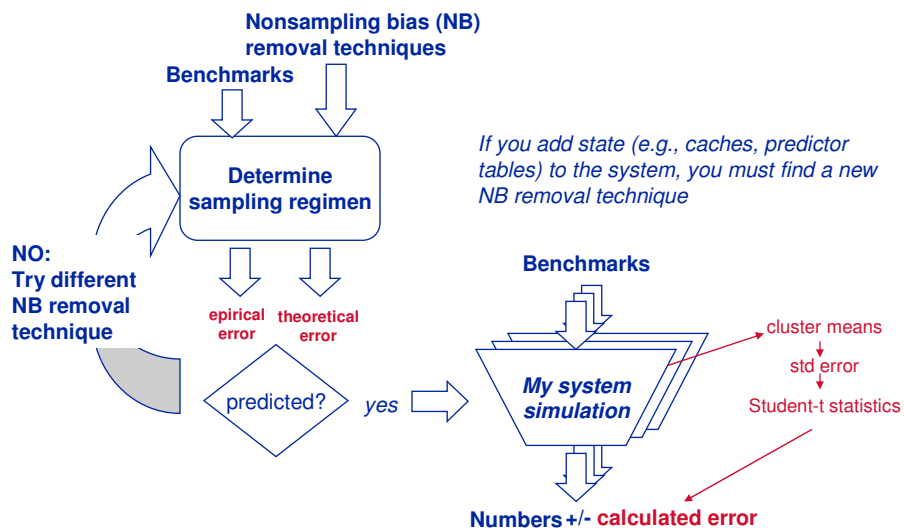**published error**

**Numbers +/- published error**

---

# Now we have this

**Nonsampling bias (NB) removal techniques**

**Benchmarks**

**Determine sampling regimen**

**NO: Try different NB removal technique**

*If you add state (e.g., caches, predictor tables) to the system, you must find a new NB removal technique*

**epirical error**  **theoretical error**

predicted? → *yes*

**Benchmarks**

**My system simulation**

cluster means
std error
Student-t statistics
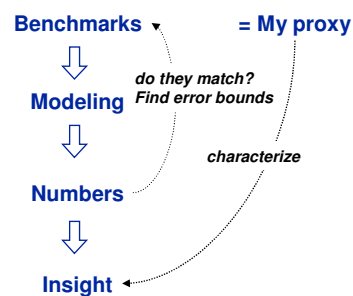
**Numbers +/- calculated error**

28

# Sampling last thoughts

○ **If regimen developed correctly**
- ◆ **can use the derived sampling regimen to calculate confidence intervals**

○ ***You know how much you can trust your numbers***

○ **Much more research is needed into effective nonsampling bias removal techniques**

○ **Should we develop benchmarks just for finding sampling regimens?**

○ **All results should include confidence intervals – even if it makes your gizmo look bad**

---

# The road map

**Benchmarks** → **= My proxy**

⇩

*do they match?*
**Modeling** *Find error bounds*

⇩

*characterize*

**Numbers**

⇩

**Insight** ←

○ **And I think Hamming would be happy with that**

# In closing, one more Hamming quotation

**Mathematicians stand on each other's shoulders while computer scientists stand on each other's toes.**
**- R. Hamming**

30