#### TORSTEN HOEFLER, ROBERTO BELLI

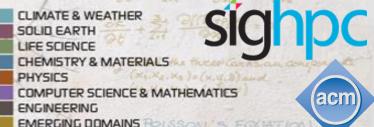
# Scientific Benchmarking of Parallel Computing Systems

Twelve ways to tell the masses when reporting performance results

Performance Engineering for HPC Workshop at ICS'17, Frankfurt, Germany









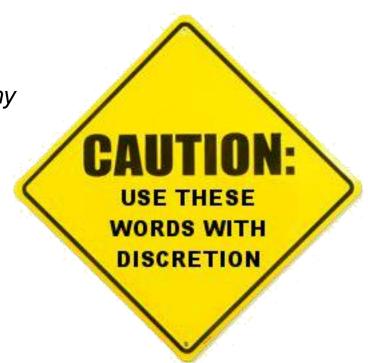


# Disclaimer(s)

- This is an experience talk (published at SC 15 State of the Practice)!
  - Explained in SC15 FAQ:

"generalizable insights as gained from experiences with particular HPC machines/operations/applications/benchmarks, overall analysis of the status quo of a particular metric of the entire field or historical reviews of the progress of the field."

- Don't expect novel insights
   Given the papers I read, much of what I say may be new for many
- My musings shall not offend anybody
  - Everything is (now) anonymized
- Criticism may be rhetorically exaggerated
  - Watch for tropes!
- This talk should be entertaining!

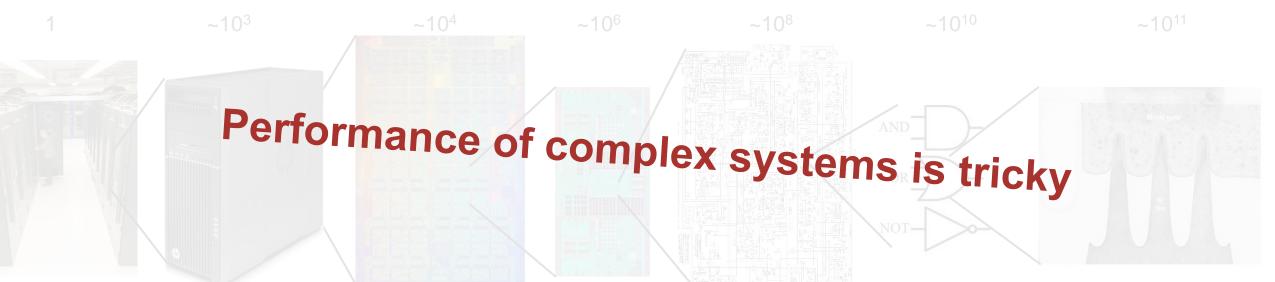




# **High Performance Computing**







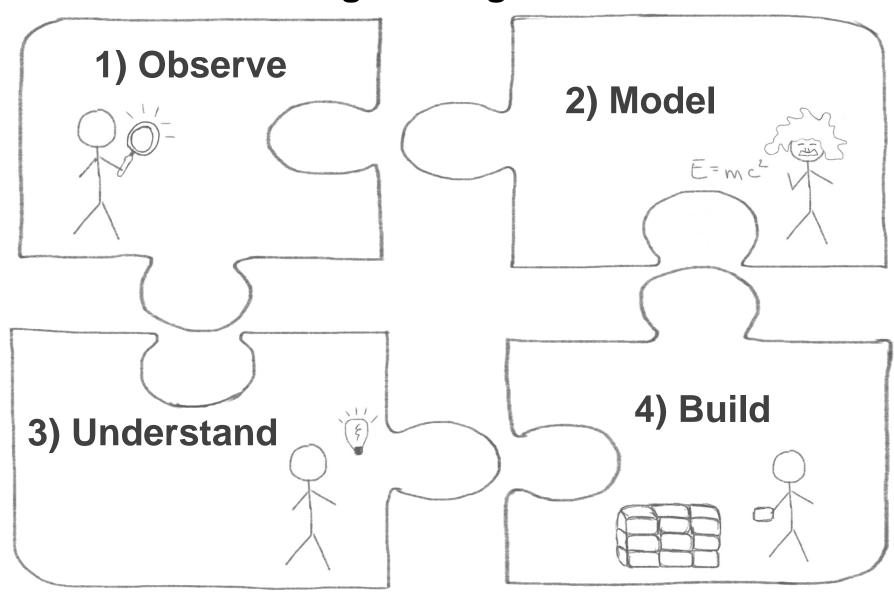


# **HPC** is used to solve complex problems!

Treat performance-centric programming and system design like physical systems

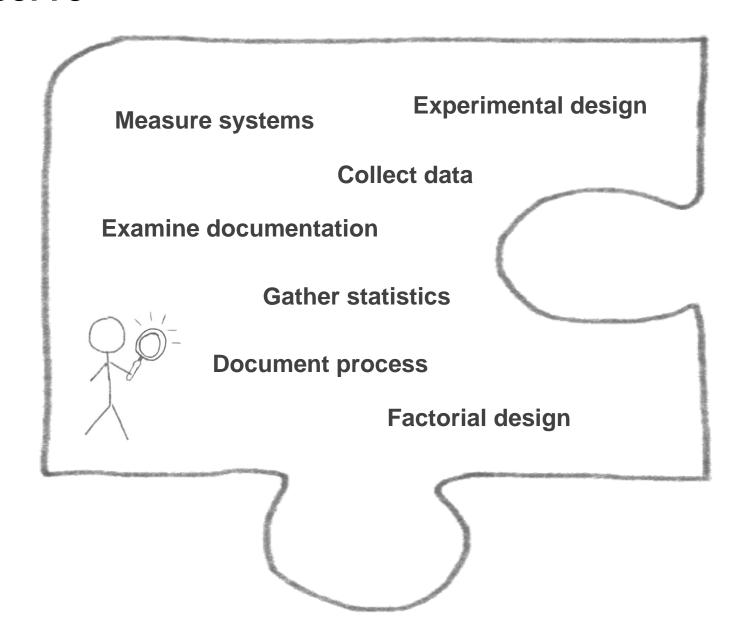


## **Scientific Performance Engineering**



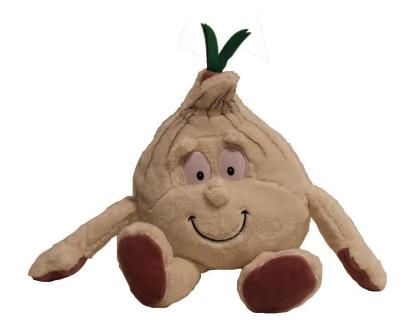


#### **Part I: Observe**



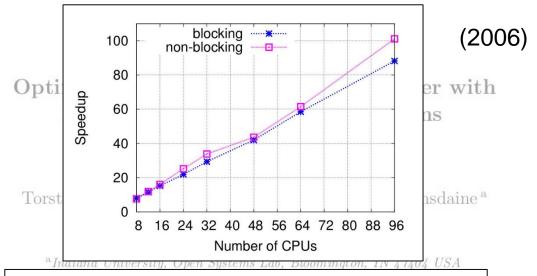
#### How does Garth measure and report performance?

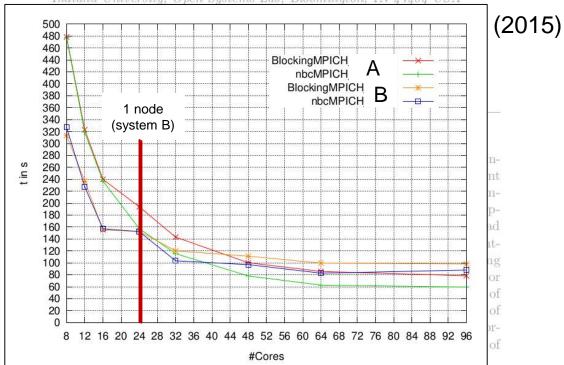
- We may be interested in High Performance Computing
  - We (want to) see it as a science reproducing experiments is a major pillar of the scientific method
- When measuring performance, important questions are
  - "How many iterations do I have to run per measurement?"
  - "How many measurements should I run?"
  - "Once I have all data, how do I summarize it into a single number?"
  - "How do I compare the performance of different systems?"
  - "How do I measure time in a parallel system?"
  - ...
- How are they answered in the field today?
  - Let me start with a little anecdote ... a reaction to this paper ©











#### Original findings:

- If carefully tuned, NBC speed up a 3D solver
   Full code published
- 800³ domain 4 GB (distributed) array
   1 process per node, 8-96 nodes
   Opteron 246 (old even in 2006, retired now)
- Super-linear speedup for 96 nodes
   ~5% better than linear

#### ■ 9 years later: attempt to reproduce ©!

System A: 28 quad-core nodes, Xeon E5520 System B: 4 nodes, dual Opteron 6274

"Neither the experiment in A nor the one in B could reproduce the results presented in the original paper, where the usage of the NBC library resulted in a performance gain for practically all node counts, reaching a superlinear speedup for 96 cores (explained as being due to cache effects in the inner part of the matrix vector product)."



#### **State of the Practice in HPC**

- Stratified random sample of three top-conferences over four years
  - HPDC, PPoPP, SC (years: 2011, 2012, 2013, 2014)
  - 10 random papers from each (10-50% of population)
  - 120 total papers, 20% (25) did not report performance (were excluded)



#### Main results:

- Most papers report details about the hardware but fail to describe the software environment.
   Important details for reproducibility missing
- 2. The average paper's results are hard to interpret and easy to question Measurements and data not well explained
- 3. No statistically significant evidence for improvement over the years  $\odot$
- Our main thesis:

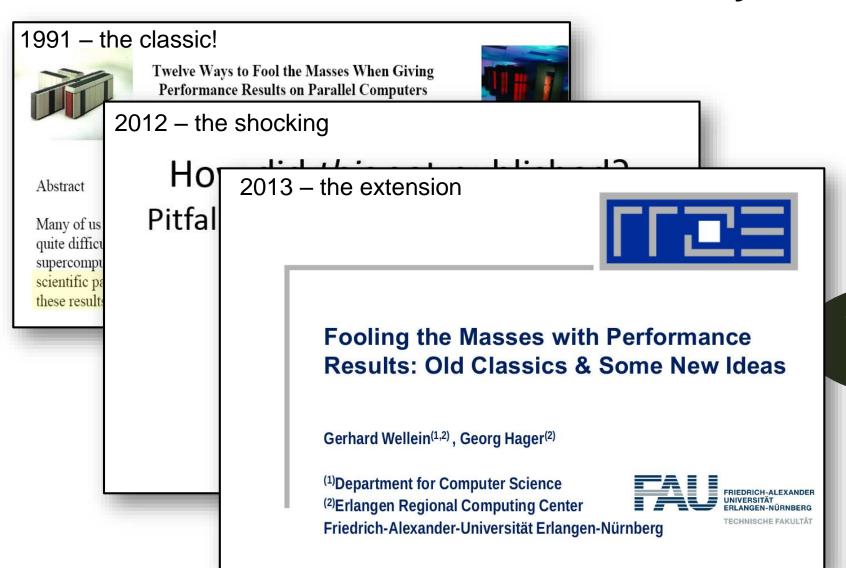
Performance results are often nearly impossible to reproduce! Thus, we need to provide enough information to allow scientists to understand the experiment, draw own conclusions, assess their certainty, and possibly generalize results.

This is especially important for HPC conferences and activities such as the Gordon Bell award





Well, we all know this - but do we really know how to fix it?







# Our constructive approach: provide a set of (12) rules

- Attempt to emphasize interpretability of performance experiments
- The set is not complete
  - And probably never will be
  - Intended to serve as a solid start
  - Call to the community to extend it
- I will illustrate the 12 rules now see with Performance
  - Using real-world examples All anonymized!
  - Garth and Eddie will represent the bad/good scientist

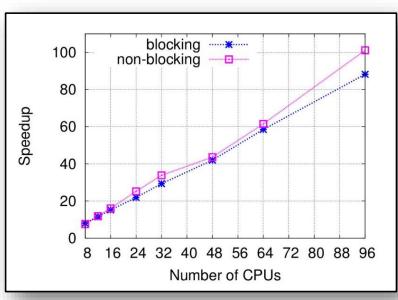
Yes, this is a garlic press!





#### The most common issue: speedup plots







#### Most common and oldest-known issue

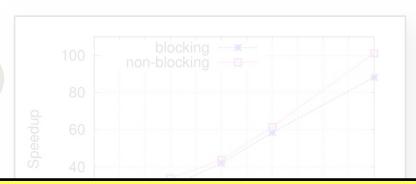
- First seen 1988 also included in Bailey's 12 ways
- 39 papers reported speedups 15 (38%) did not specify the base-performance 🗵
- Recently rediscovered in the "big data" universe
  - A. Rowstron et al.: Nobody ever got fired for using Hadoop on a cluster, HotCDP 2012
  - F. McSherry et al.: Scalability! but at what cost?, HotOS 2015







### The most common issue: speedup plots



Rule 1: When publishing parallel speedup, report if the base case is a single parallel process or best serial execution, as well as the absolute execution performance of the base case.

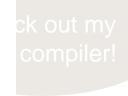
- - A simple generalization of this rule implies that one should never report ratios without absolute values.

    15 (38%) did not specify the base-performance ©





### Garth's new compiler optimization



Performance in Gflop/s

How d perform and B

**Rule 2**: Specify the reason for only reporting subsets of standard benchmarks or applications or not using all system resources.

Well, GarthCC

NAS CG

NAS LU

VAS EP

and was 2. This implies: Show results even if your code/approach stops scaling!







The mean parts of means - or how to summarize data

**Rule 3**: Use the arithmetic mean only for summarizing costs.

Use the harmonic mean for summarizing rates.

Rule 4: Avoid summarizing ratios; summarize the costs or rates that the ratios base on instead. Only if these are not available use the geometric mean for summarizing ratios.

Ah, true, the

- 51 papers use means to summarize data, only four (!) specify which mean was used
  - A single paper correctly specifies the use of the harmonic mean
  - Two use geometric means, without reason
  - Similar issues in other communities (PLDI, CGO, LCTES) see N. Amaral's report ine of
- harmonic mean ≤ geometric mean ≤ arithmetic mean



# **Dealing with variation**

The latency of Piz Dora is

How did you get to this?

**Rule 5**: Report if the measurement values are deterministic. For nondeterministic data, report confidence intervals of the measurement.

- Most papers report nondeterministic measurement results
  - Only 15 mention some measure of variance
  - Only two (!) report confidence intervals

Why do you think so? Can see the data?

- Cls allow us to compute the number of required measurements!
- Can be very simple, e.g., single sentence in evaluation:

"We collected measurements until the 99% confidence interval was within 5% of our reported means."



# **Dealing with variation**

The confidence interval is 1.765us to 1.775us

Rule 6: Do not assume normality of collected data (e.g., based on the number of samples) without diagnostic checking.

- Most events will slow down performance
  - Heavy right-tailed distributions
- The Central Limit Theorem only applies asymptotically
  - Some papers/textbook mention "30-40 samples", don't trust them!

normal at all! The real

Two papers used CIs around the mean without testing for normality

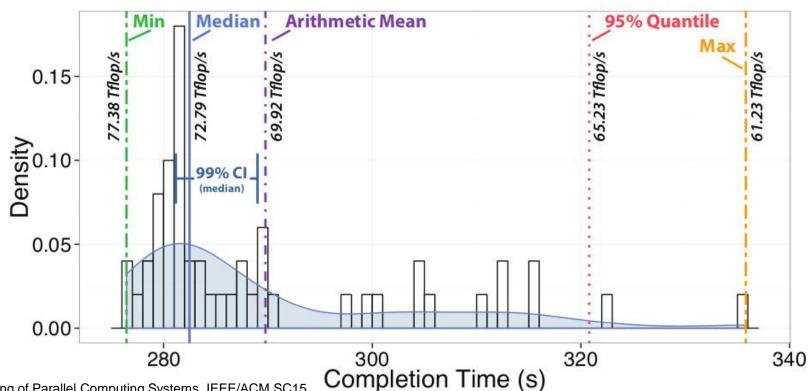
Can we test for normality?





# Dealing with non-normal data – nonparametric statistics

- Rank-based measures (no assumption about distribution)
  - Essentially always better than assuming normality
- Example: median (50<sup>th</sup> percentile) vs. mean for HPL
  - Rather stable statistic for expectation
  - Other percentiles (usually 25<sup>th</sup> and 75<sup>th</sup>) are also useful

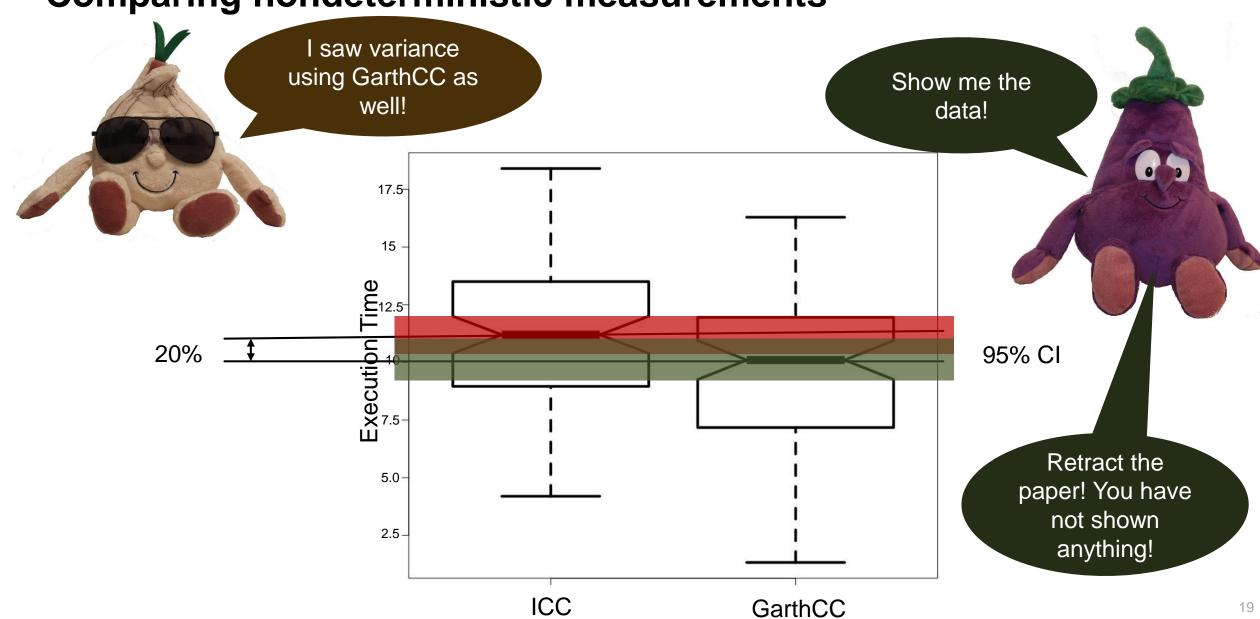






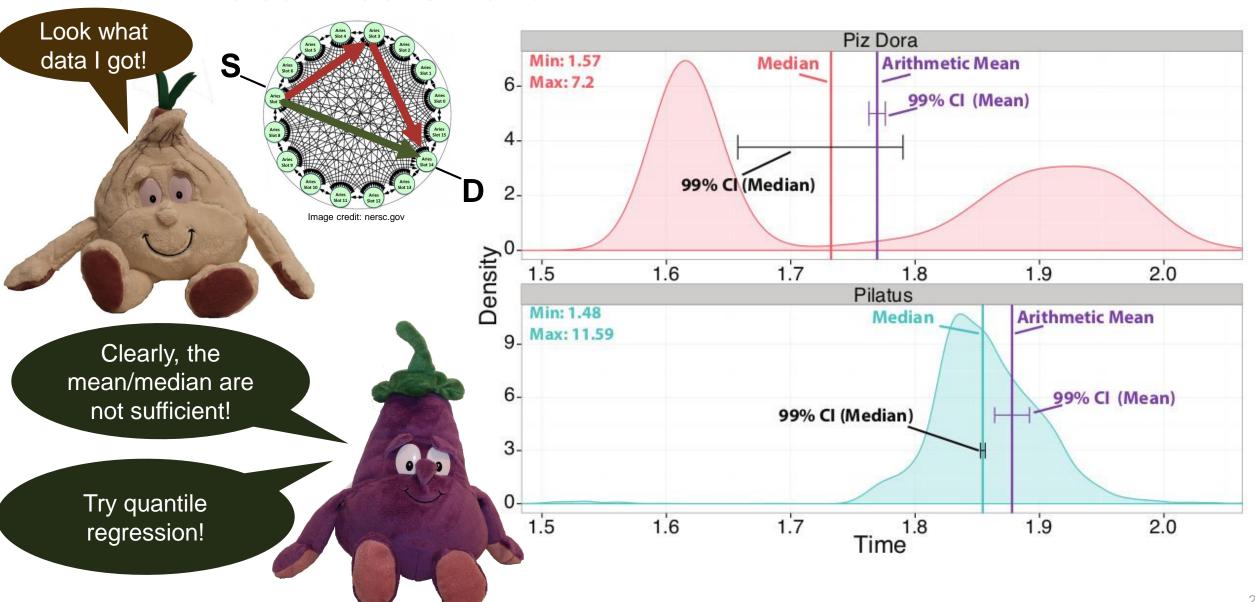
### **Comparing nondeterministic measurements**

\*\*SPCL





#### What if the data looks weird!?







Wow, so Pilatus is better for (worstcase) latency-critical workloads even though Dora is expected to be faster

Rule 8: Carefully investigate if measures of central tendency such as mean or median are useful to report. Some problems, such as worst-case latency, may require other percentiles.

 Check Oliveira et al. "Why you should care about quantile regression". SIGARCH Computer Architecture News, 2013.



# How many measurements are needed?

- Measurements can be expensive!
  - Yet necessary to reach certain confidence
- How to determine the minimal number of measurements?
  - Measure until the confidence interval has a certain acceptable width
  - For example, measure until the 95% CI is within 5% of the mean/median
  - Can be computed analytically assuming normal data
  - Compute iteratively for nonparametric statistics
- Often heard: "we cannot afford more than a single measurement"
  - E.g., Gordon Bell runs
  - Well, then one cannot say anything about the variance
     Even 3-4 measurement can provide very tight CI (assuming normality)
     Can also exploit repetitive nature of many applications





## **Experimental design**

I don't believe you, try other numbers of processes!

MPI Reduce

Rule 9: Document all varying factors and their levels as well as the complete experimental setup (e.g., software, hardware, techniques) to facilitate reproducibility and provide interpretability.

- We recommend factorial design
- Consider parameters such as node allocation, process-to-node mapping, network or node contention
  - If they cannot be controlled easily, use randomization and model them as random variable
- This is hard in practice and not easy to capture in rules

That's nonsense!



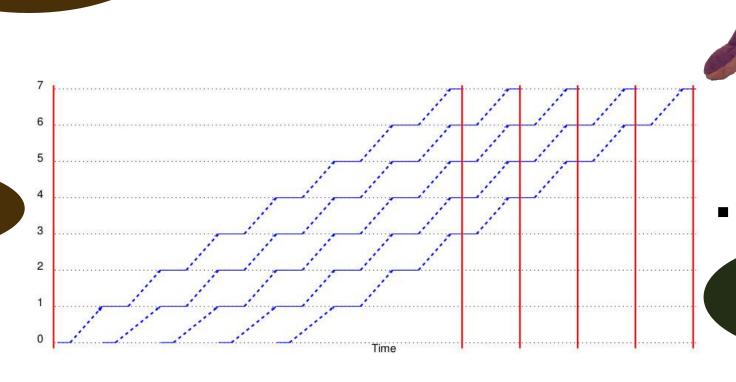


# Time in parallel systems

My simple broadcast takes only one latency!

But I measured it so it must be true!

```
t = -MPI_Wtime();
for(i=0; i<1000; i++) {
    MPI_Bcast(...);
}
t += MPI_Wtime();
t /= 1000;</pre>
```







# Summarizing times in parallel systems!

Come on, show me the data!

My new reduce

Rule 10: For parallel time measurements, report all measurement, (optional) synchronization, and summarization techniques.

- Measure events separately
  - Use high-precision timers
  - Synchronize processes
- Summarize across processes:
  - Min/max (unstable), average, median depends on use-case



# Give times a meaning!

have no clue.

compute 10<sup>10</sup>

Rule 11: If possible, show upper performance bounds to facilitate interpretability of the measured results.

#### Model computer system as k-dimensional space

- Each dimension represents a capability
   Floating point, Integer, memory bandwidth, cache bandwidth, etc.
- Dk T Features are typical rates
  - Determine maximum rate for each dimension
     E.g., from documentation or benchmarks
- Can be used to proof optimality of implementation
  - If the requirements of the bottleneck dimension are minimal

Can you provide?
Ideal speedup

Amdahl's speedup

Parallel overheads



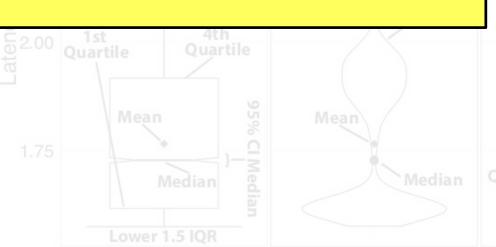
### Plot as much information as possible!

My most common request was "show me the data"



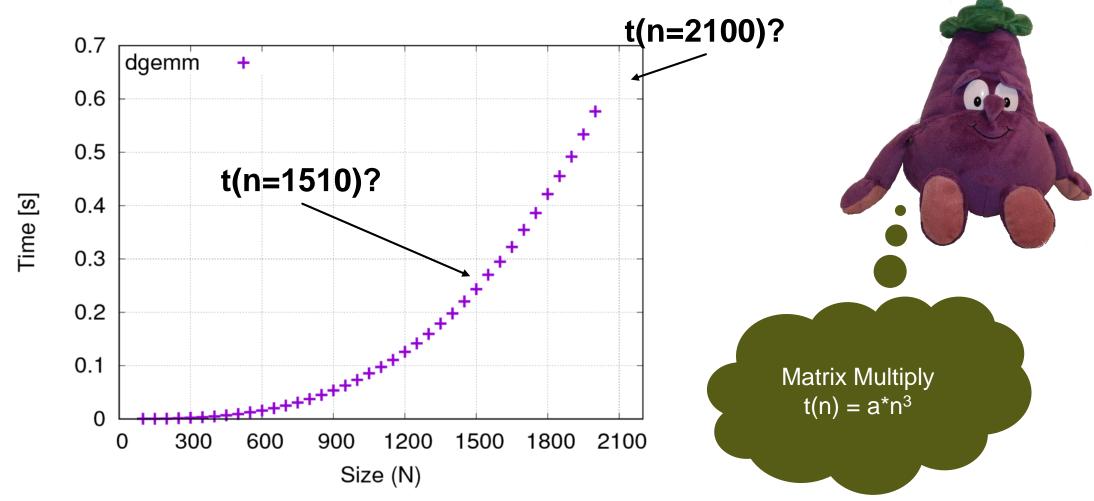
Rule 12: Plot as much information as needed to interpret the experimental results. Only connect measurements by lines if they indicate trends and the interpolation is valid.

This is how I should have presented the Dora results.





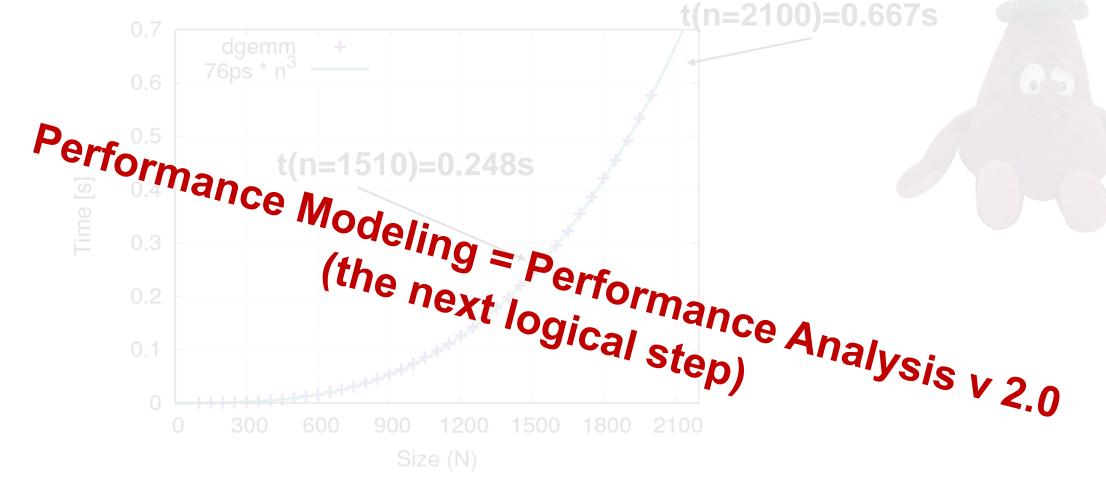
# We have the (statistically sound) data, now what?



The 99% confidence interval is within 1% of the reported median.





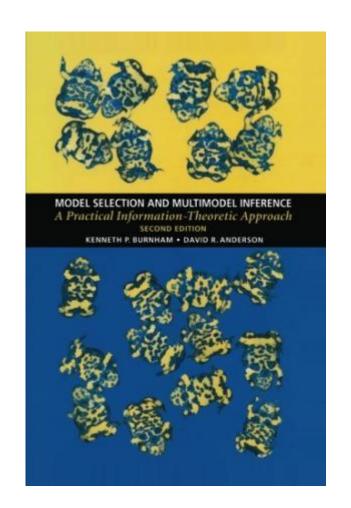


The state of the s

The 99% confidence interval is within 1% of the reported median. The adjusted  $R^2$  of the model fit is 0.99

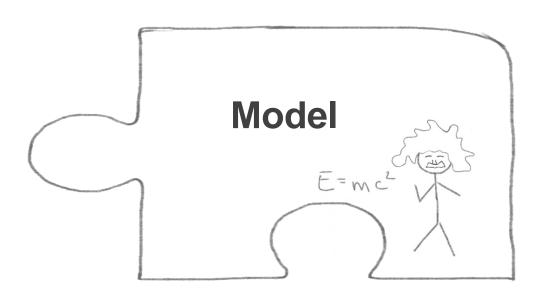


#### Part II: Model









**Burnham, Anderson:** "A model is a simplification or approximation of reality and hence will not reflect all of reality. ... Box noted that "all models are wrong, but some are useful." While a model can never be "truth," a model might be ranked from very useful, to useful, to somewhat useful to, finally, essentially useless."

This is generally true for all kinds of modeling.

We focus on **performance modeling** in the following!





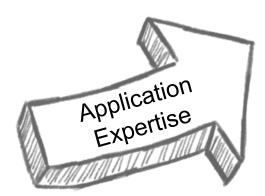
#### **Performance**







**Performance Model** 



**Requirements Model** 



Grid Points per Process (L)



# Requirements modeling I: Six-step performance modeling

Input parameters Communication Describe application parameters kernels 25000 Serial Model Model P=1024 Communication Comm Overhead 20000 Fit sequential Pack Overhead pattern baseline Ime [ms] 15000 10000 Communication / computation overlap 5000 10-20% speedup [2] 1500 2000

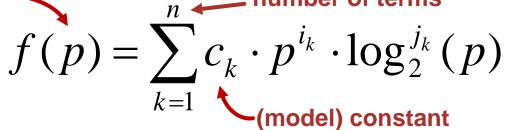


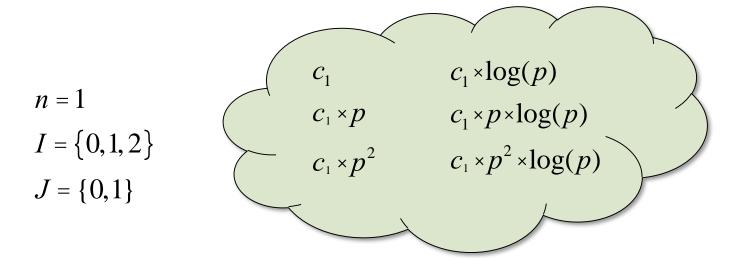


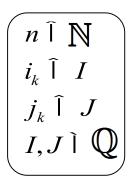
# Requirements modeling II: Automated best-fit modeling

- Manual kernel selection and hypothesis generation is time consuming (boring and tricky)
- Idea: Automatically select best (scalability) model from predefined search space







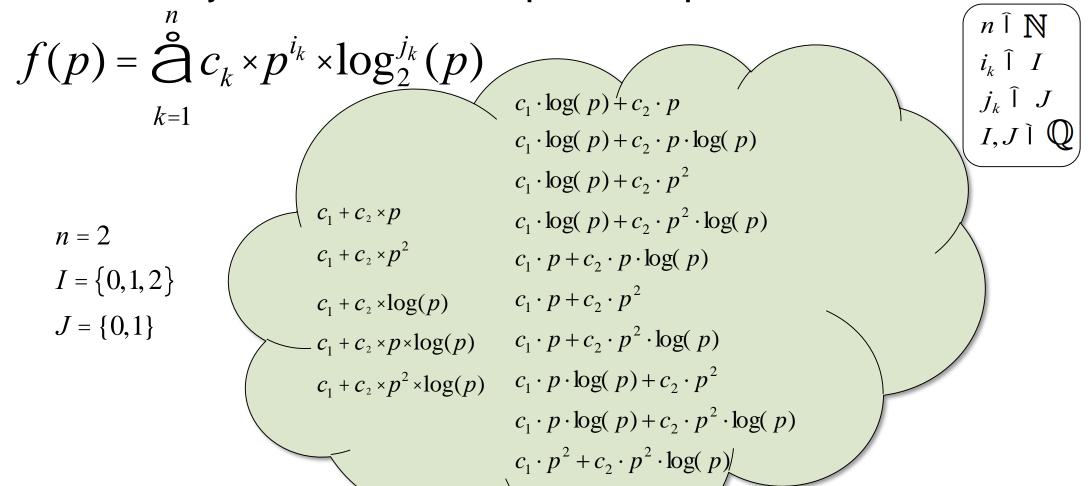






# Requirements modeling II: Automated best-fit modeling

- Manual kernel selection and hypothesis generation is time consuming (and boring)
- Idea: Automatically select best model from predefined space







# Tool support: Extra-P for automated best-fit modeling [1]











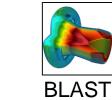


















MP2C

[1] Download Extra-P at: <a href="http://www.scalasca.org/software/extra-p/download.html">http://www.scalasca.org/software/extra-p/download.html</a>

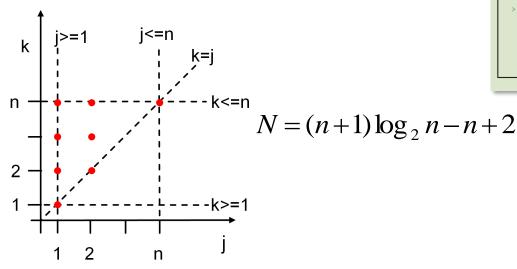




# Requirements modeling III: Source-code analysis [1]

- Extra-P selects model based on best fit to the data
  - What if the data is not sufficient or too noisy?
- Back to first principles
  - The source code describes all possible executions
  - Describing all possibilities is too expensive, focus on counting loop iterations symbolically

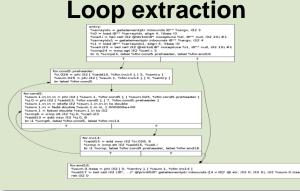
```
for (j = 1; j <= n; j = j*2)
  for (k = j; k <= n; k = k++)
    OperationInBody(j,k);</pre>
```



#### Parallel program





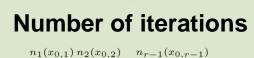


#### **Requirements Models**

$$W = N \big|_{p=1}$$

$$D=N\big|_{p\to\infty}$$





$$N = \sum_{i_1=0}^{n_1(x_0, i_1)} \sum_{i_2=0}^{n_2(x_0, i_2)} \dots \sum_{i_{r-1}=0}^{n_{r-1}(x_0, i_{r-1})} n_r(x_{0,r}).$$





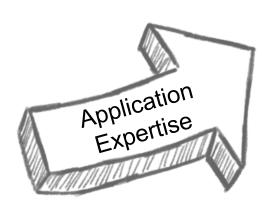


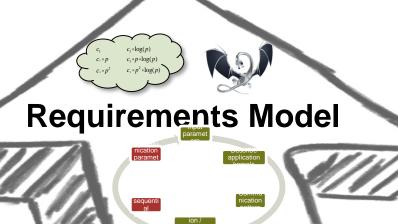
#### **Performance**













#### \*\*\*SPCL

#### **Performance**



**Performance Model** 





LogP





### Capability models for network communication

#### The LogP model family and the LogGOPS model [1]

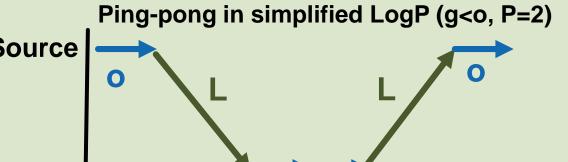
A PRACTICAL MODEL of PARALLEL COMPUTATION

JR GOAL IS TO DEVELOP A MODEL OF PARALLEL COMPUTATION THAT WILL serve as a basis for the design and analysis of fast, portable parallel algorithms, such as algorithms that can be implemented effectively on a wide variety of current and future parallel machines. If we look at the body of parallel algorithms developed under current parallel models, many are impractical because they exploit artificial factors not present in any reaPRAM consists of a collection of processors a global random access

David E. Culler, Richard M. Kari

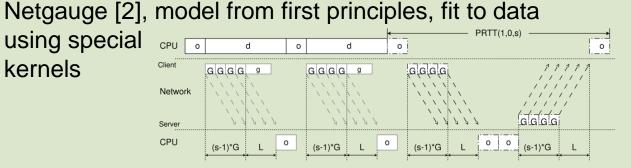
Source

Dest.



#### **Finding LogGOPS parameters**

using special kernels

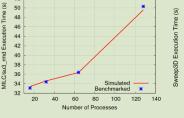


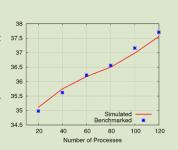
#### **Large scale LogGOPS Simulation**

LogGOPSim [1], simulates LogGOPS with 10

million MPI ranks

<5% error





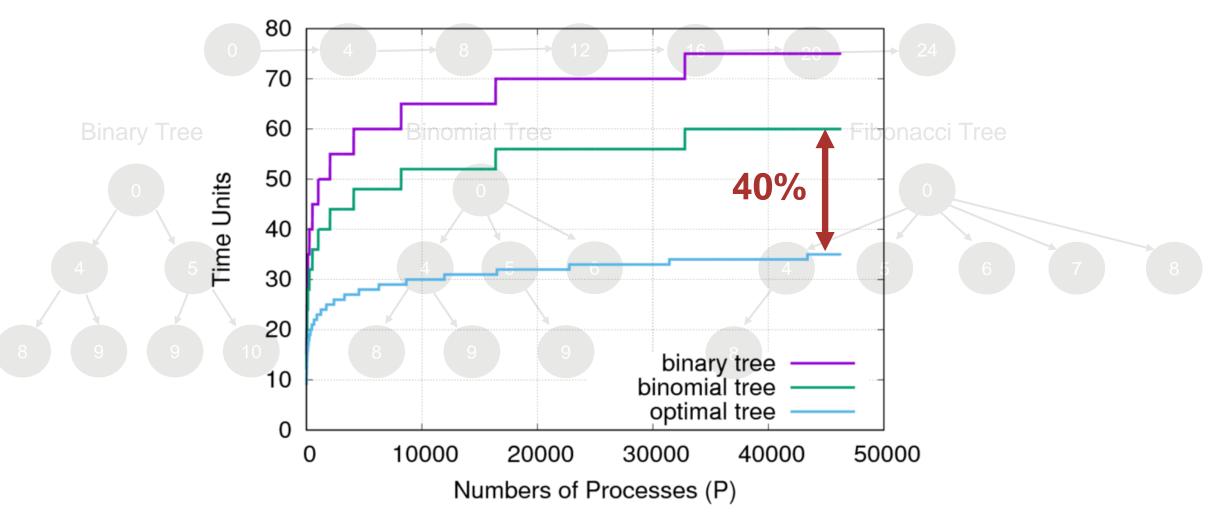




# 2) Design optimal algorithms – small broadcast in LogP 5

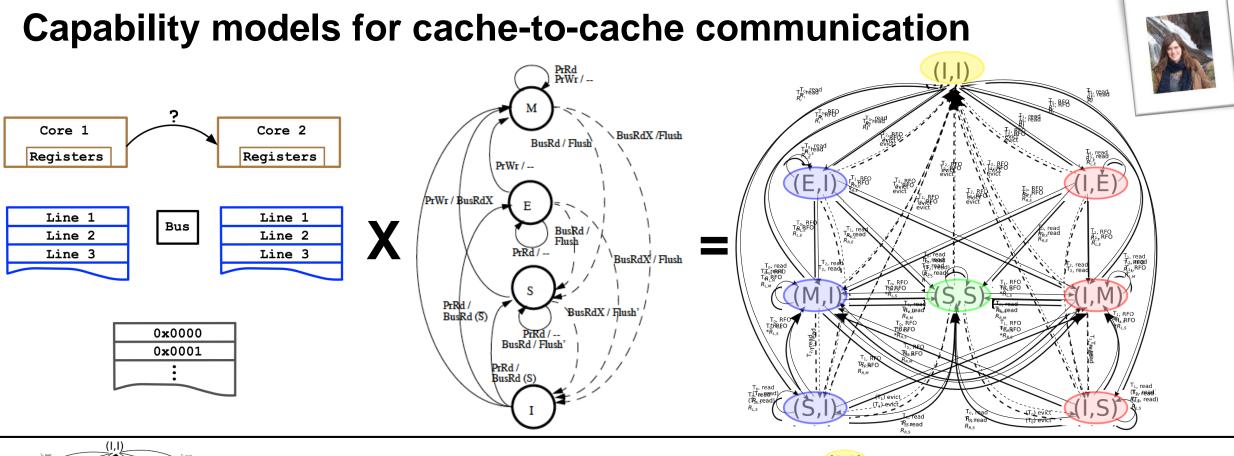


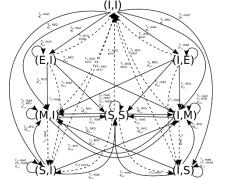
L=2, o=1, P=7



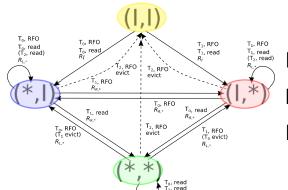


#### \*\*SPCL







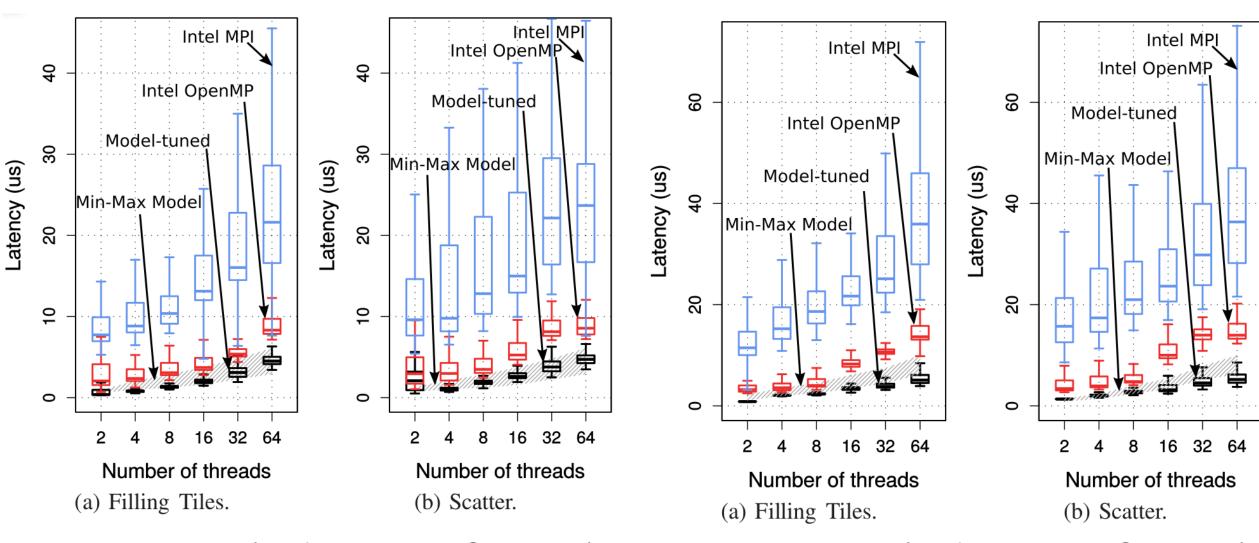


Invalid read  $R_1 \approx 135 \text{ ns}$ Local read:  $R_L = 3.8 \text{ ns}$ Remote read  $R_R \approx 115 \text{ ns}$ 

S. Ramos, TH: "Capability Models for Manycore Memory Systems: A Case-Study with Xeon Phi KNL", IEEE IPDPS'17



### Model-tuned Barrier and Reduce vs. Intel's OpenMP and MPI



Barrier (7x faster than OpenMP)

Reduce (5x faster then OpenMP)







Modeling

**Capability Model** 



**Performance Model** 





**Requirements Model** 











#### Conclusions and call for action

- Performance may not be reproducible
  - At least not for many (important) results
- Interpretability fosters scientific progress
  - Enables to build on results
  - Sounds statistics is the biggest gap today
- We need to foster interpretability
  - Do it ourselves (this is not easy)
  - Teach young students
  - Maybe even enforce in TPCs
- See the 12 rules as a start
  - Need to be extended (or concretized)
  - Much is implemented in LibSciBench [1]



No vegetables were harmed for creating these slides!

#### **Acknowledgments**

- ETH's mathematics department (home of R)
  - Hans Rudolf Künsch, Martin Maechler, and Robert Gantner
- Comments on early drafts
  - David H. Bailey, William T. Kramer, Matthias Hauswirth, Timothy Roscoe, Gustavo Alonso, Georg Hager, Jesper Träff, and Sascha Hunold
- Help with HPL run
  - Gilles Fourestier (CSCS) and Massimiliano Fatica (NVIDIA)