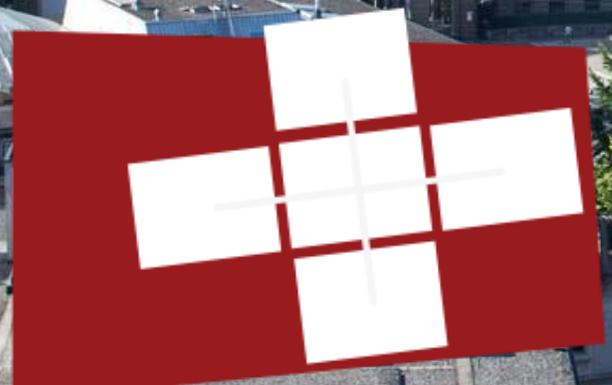


TORSTEN HOEFLER

Theory and practice in HPC: Modeling, Programming, and Networking



Systems@ETH zürich

2017 Platform for Advanced Scientific Computing Conference
Lugano Switzerland | 26-28 June 2017

- CLIMATE & WEATHER
- SOLID EARTH
- LIFE SCIENCE
- CHEMISTRY & MATERIALS
- PHYSICS
- COMPUTER SCIENCE & MATHEMATICS
- ENGINEERING
- EMERGING DOMAINS

Background of the banner features mathematical equations and diagrams, including the Poisson equation: $\Delta u = f$.



24 years ago



Technische Universität Chemnitz, Saxony, Germany



November 28th - December 1st

Next Cluster Series Conference
CLUSTER2001 in Newport Beach, CA



CLUSTER2000 was sponsored by



TFCC
IEEE Task Force on Cluster Computing



- Home
- Attendees
- Impressions
- Forum
- FAQ
- Study
- Cluster
- Sponsors
- Partners
- Venue

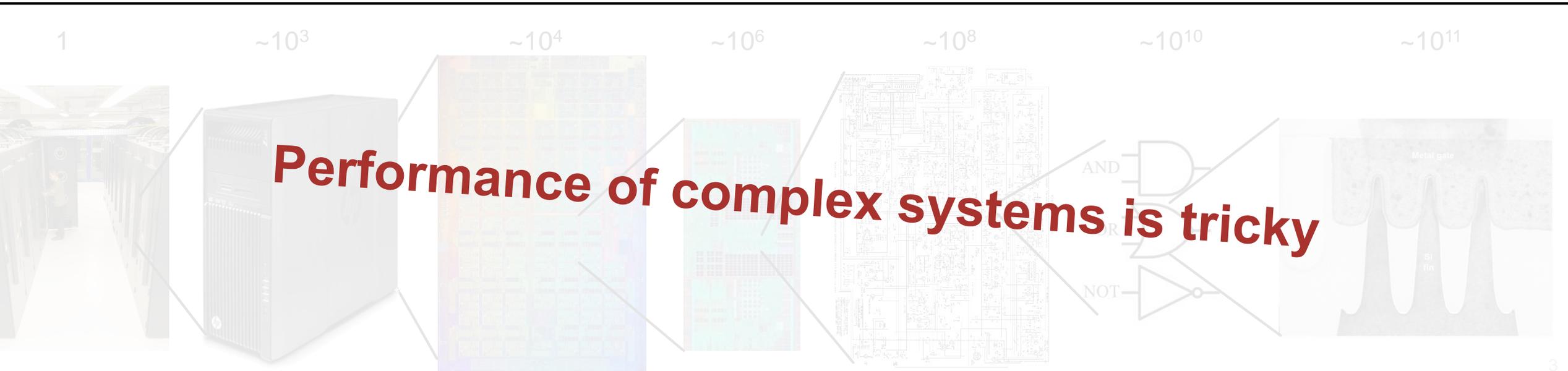


High Performance Cluster Computing

```
dgemm("N", "N", 50, 50, 50, 1.0, A, 50, B, 50, 1.0, C, 50);
```



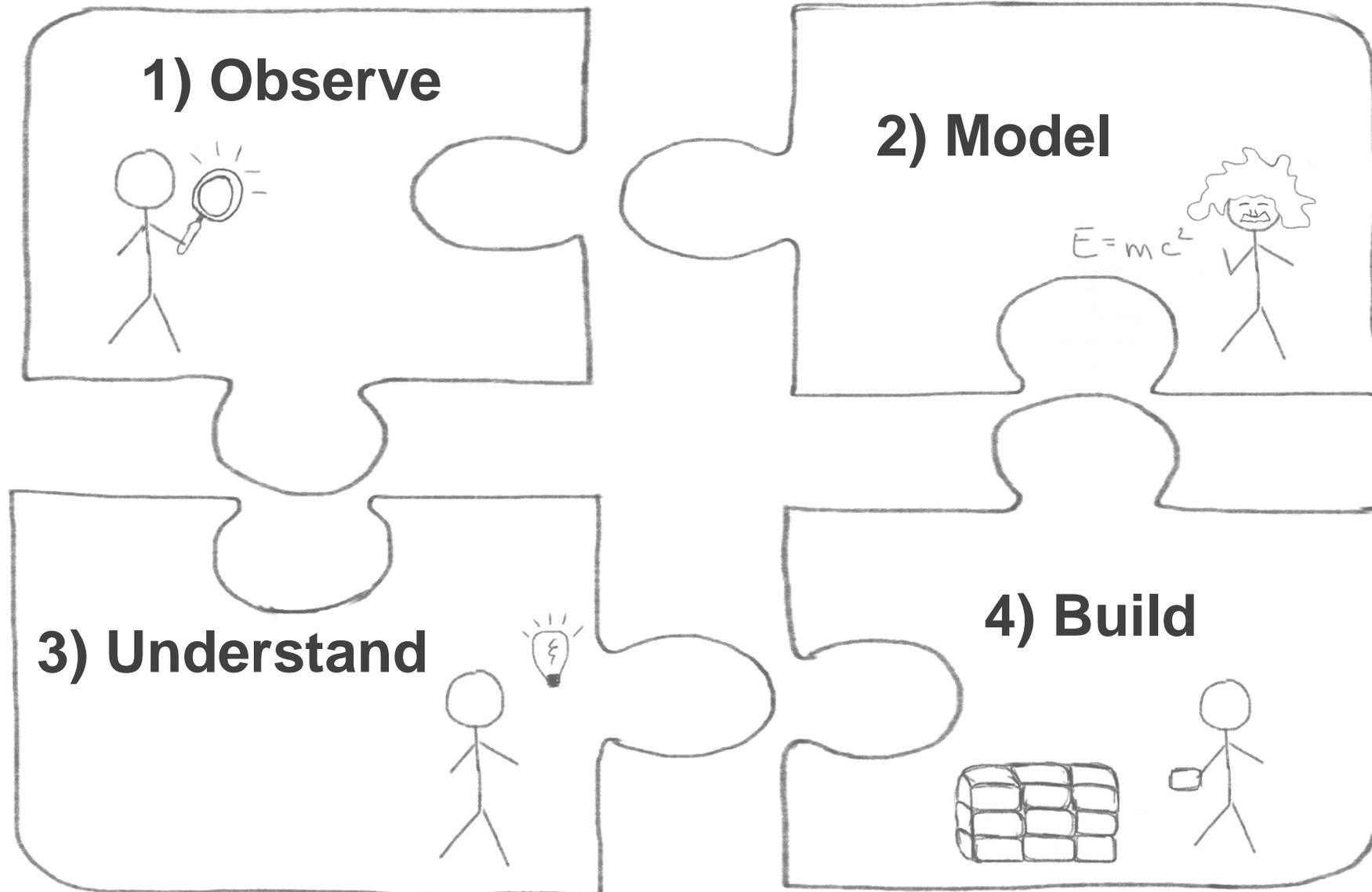
Performance is nondeterministic and not modular



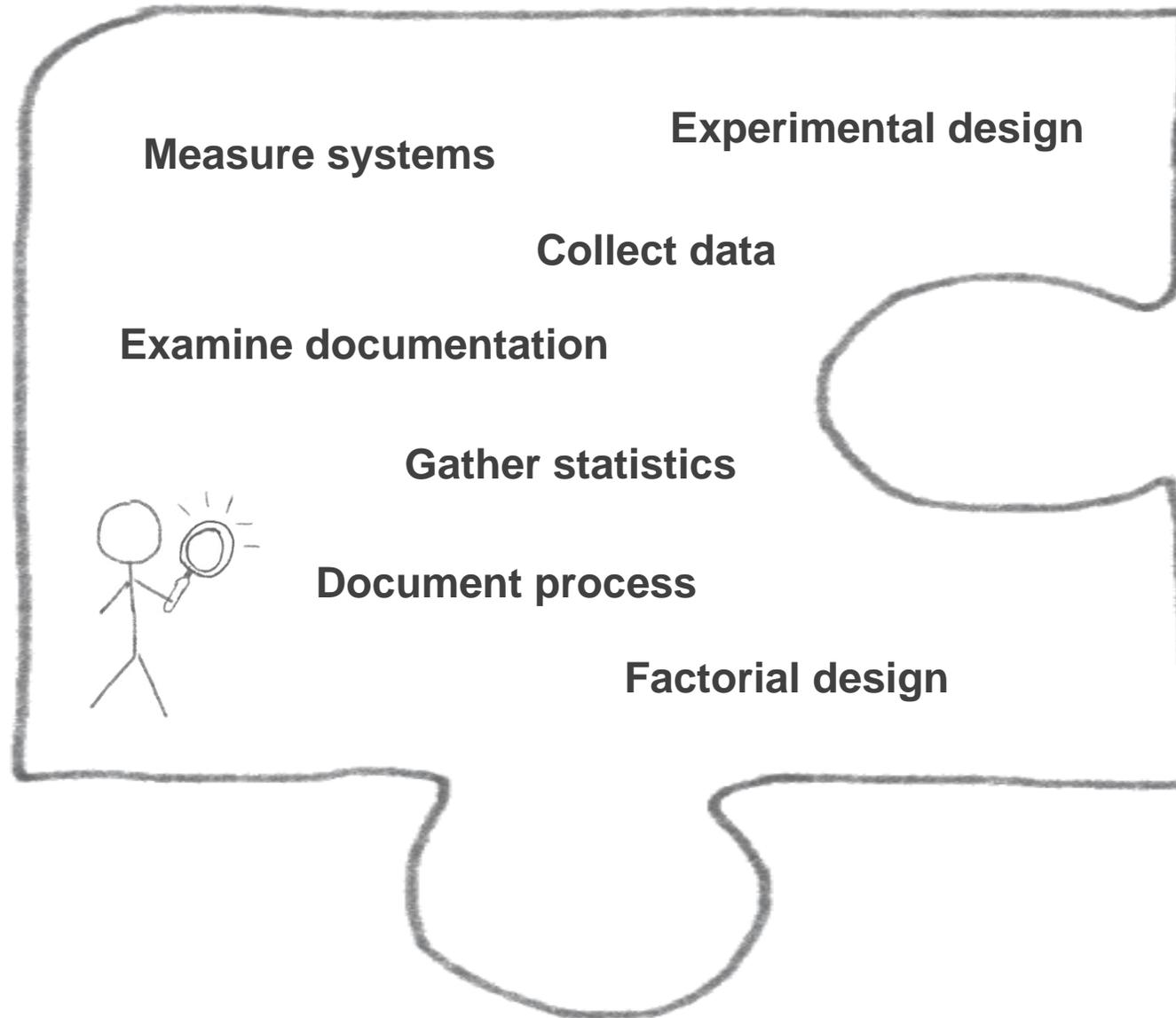
HPC is used to solve complex problems!

Treat performance-centric programming and system design like physical systems

Scientific **Performance** Engineering



Part I: Observe



Observe the state of the art in performance measurement



- Stratified random sample of three top HPC conferences for 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)

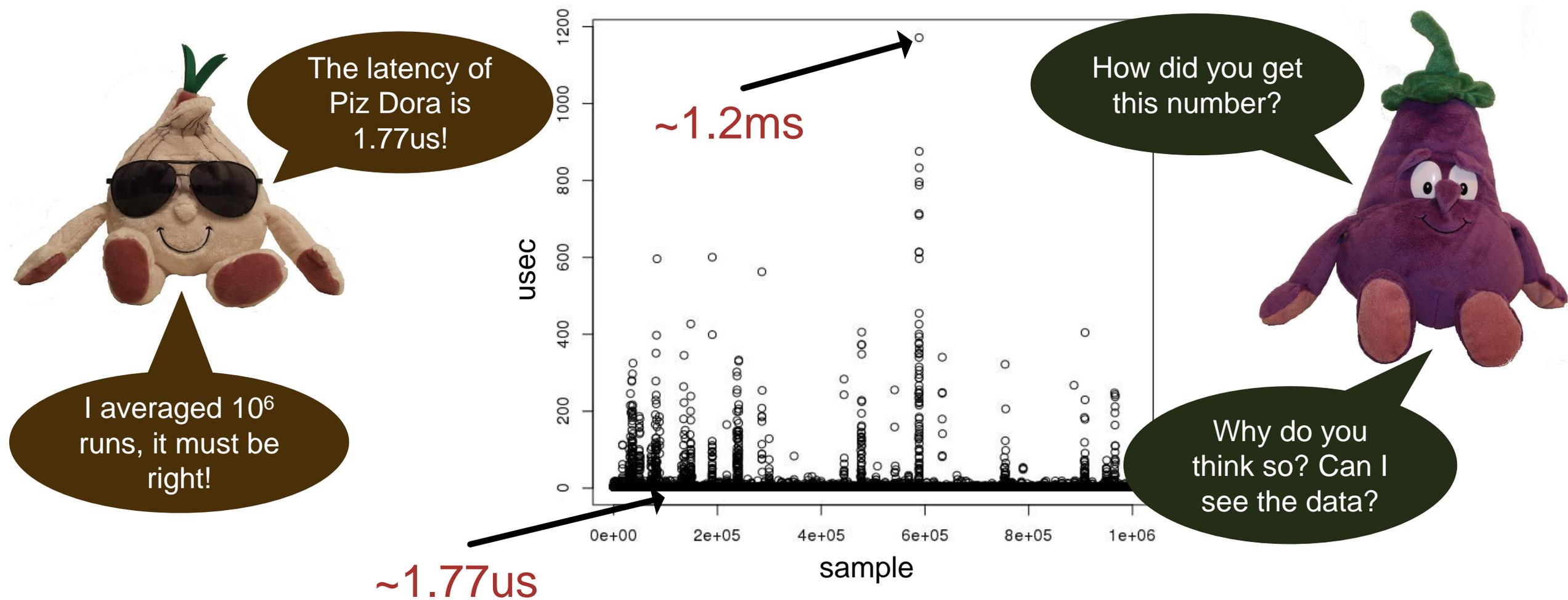
	ConfA	ConfB	ConfC	Tot ✓
Experimental Design	2013	2014	2011	(79/95)
Hardware				(26/95)
Software				(60/95)
Configuration				(35/95)
Data Analysis				(20/95)
Results				(12/95)
				(48/95)
				(30/95)
				(7/95)
				(51/95)
				(13/95)
				(9/95)
				(17/95)

Experimental Design

→

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.

Example: Simple ping-pong latency benchmark

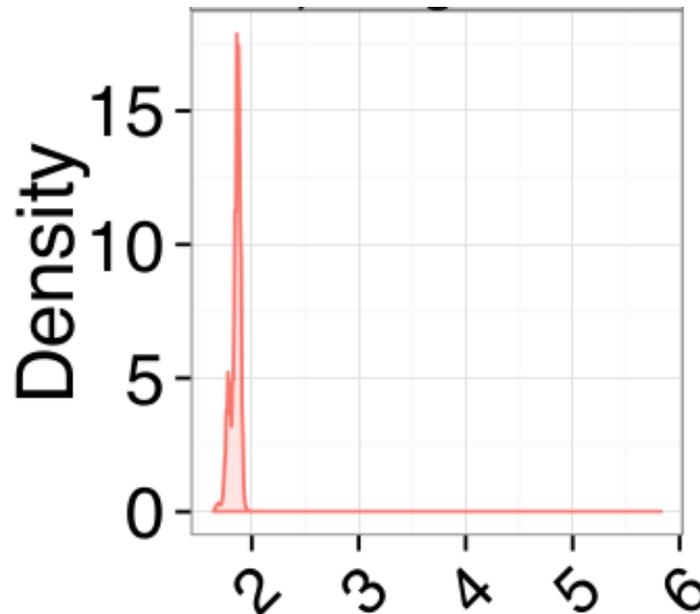


Dealing with variation

The 99.9% confidence interval is 1.765us to 1.775us



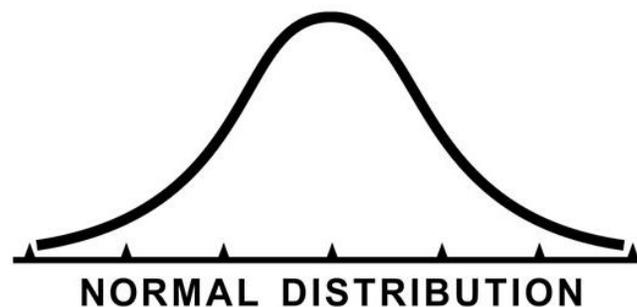
Ugs, the data is not normal at all. The nonparametric 99.9% CI is much wider: 1.6us to 1.9us!



Did you assume normality?



Can we test for normality?



Looking at the data in detail

This CI makes me nervous. Let's check!

Clearly, the mean/median are not sufficient!

Try quantile regression!

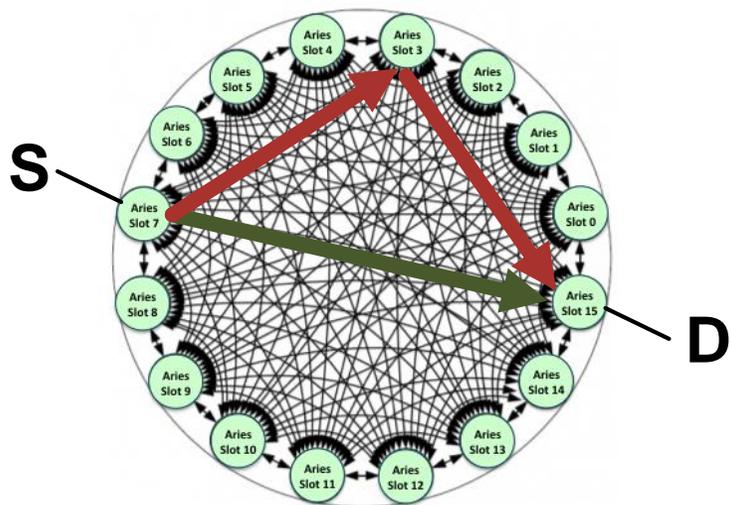
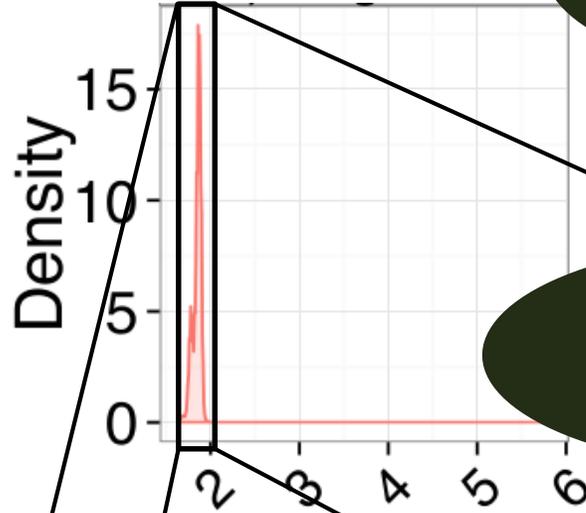
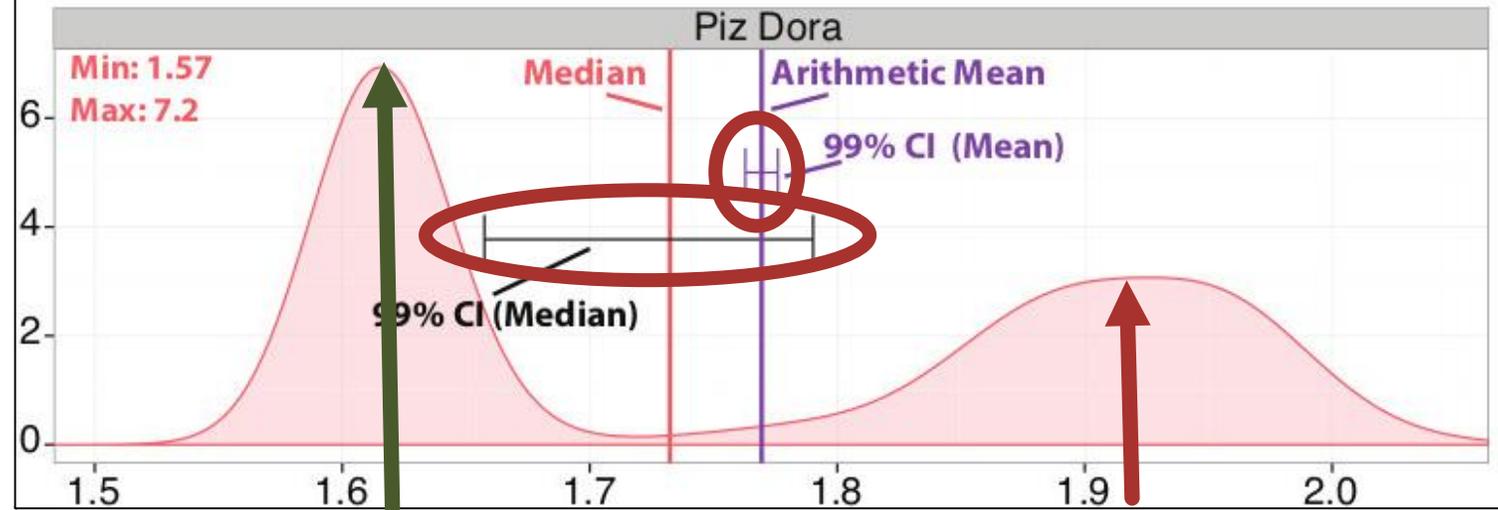


Image credit: nersc.gov



Scientific benchmarking of parallel computing systems



ACM/IEEE Supercomputing 2015 (SC15)

Scientific Benchmarking of Parallel Computing Systems

Twelve ways to tell the masses when reporting performance results

Torsten Hoefler
Dept. of Computer Science
ETH Zurich
Zurich, Switzerland
htor@inf.ethz.ch

Roberto Belli
Dept. of Computer Science
ETH Zurich
Zurich, Switzerland
bellir@inf.ethz.ch

ABSTRACT

Measuring and reporting performance of parallel computers constitutes the basis for scientific advancement of high-performance computing (HPC). Most scientific reports show performance improvements of new techniques and are thus obliged to ensure reproducibility or at least interpretability. Our investigation of a stratified sample of 120 papers across three top conferences in the field shows that the state of the practice is lacking. For example, it is often unclear if reported improvements are deterministic or observed by chance. In addition to distilling best practices from existing work, we propose statistically sound analysis and reporting techniques and simple guidelines for experimental design in parallel computing and codify them in a portable benchmarking library. We

Reproducing experiments is one of the main principles of the scientific method. It is well known that the performance of a computer program depends on the application, the input, the compiler, the runtime environment, the machine, and the measurement methodology [20, 43]. If a single one of these aspects of *experimental design* is not appropriately motivated and described, presented results can hardly be reproduced and may even be misleading or incorrect.

The complexity and uniqueness of many supercomputers makes reproducibility a hard task. For example, it is practically impossible to recreate most hero-runs that utilize the world's largest machines because these machines are often unique and their software configurations changes regularly. We introduce the notion of *interpretability*, which is weaker than reproducibility. We call an ex-

Interpret the
by lines if
valid.

Simplifying Measuring and Reporting: LibSciBench



```
#include <mpi.h>
#include <liblsb.h>
#include <stdlib.h>

#define N 1024
#define RUNS 10

int main(int argc, char *argv[]){
    int i, j, rank, buffer[N];

    MPI_Init(&argc, &argv);
    LSB_Init("test_bcast", 0);

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    /* Output the info (i.e., rank, runs) in the results file */
    LSB_Set_Rparam_int("rank", rank);
    LSB_Set_Rparam_int("runs", RUNS);

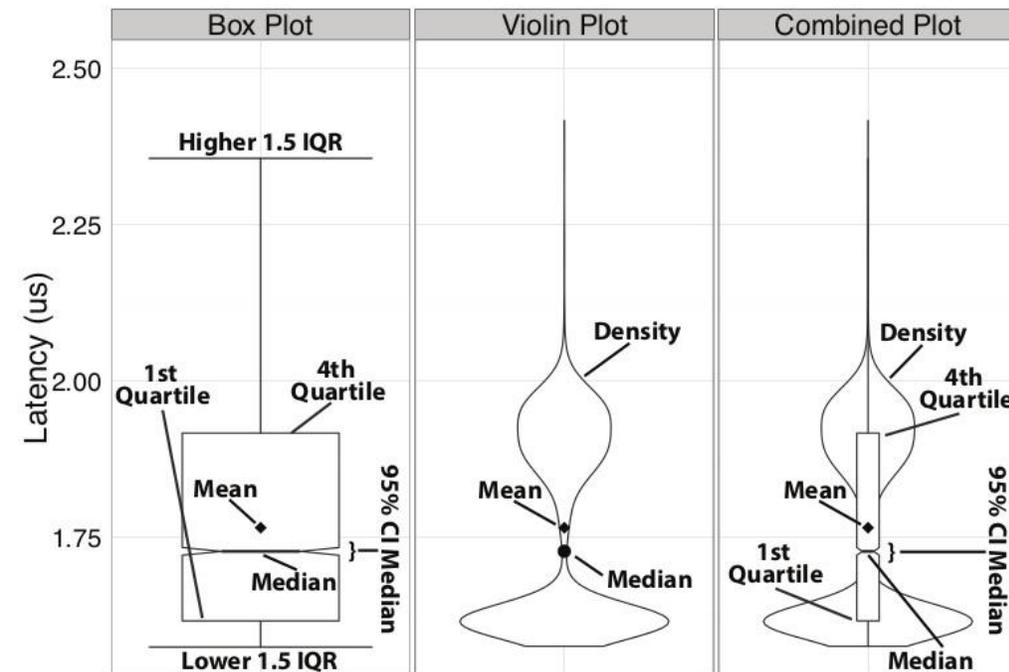
    for (sz=1; sz<=N; sz*=2){
        for (j=0; j<RUNS; j++){
            /* Reset the counters */
            LSB_Res();

            /* Perform the operation */
            MPI_Bcast(buffer, sz, MPI_INT, 0, MPI_COMM_WORLD);

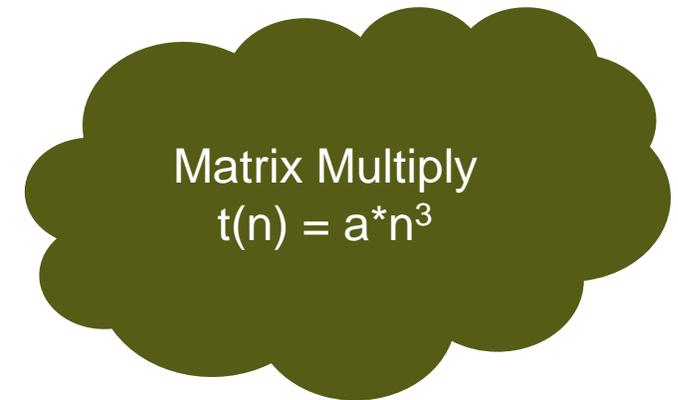
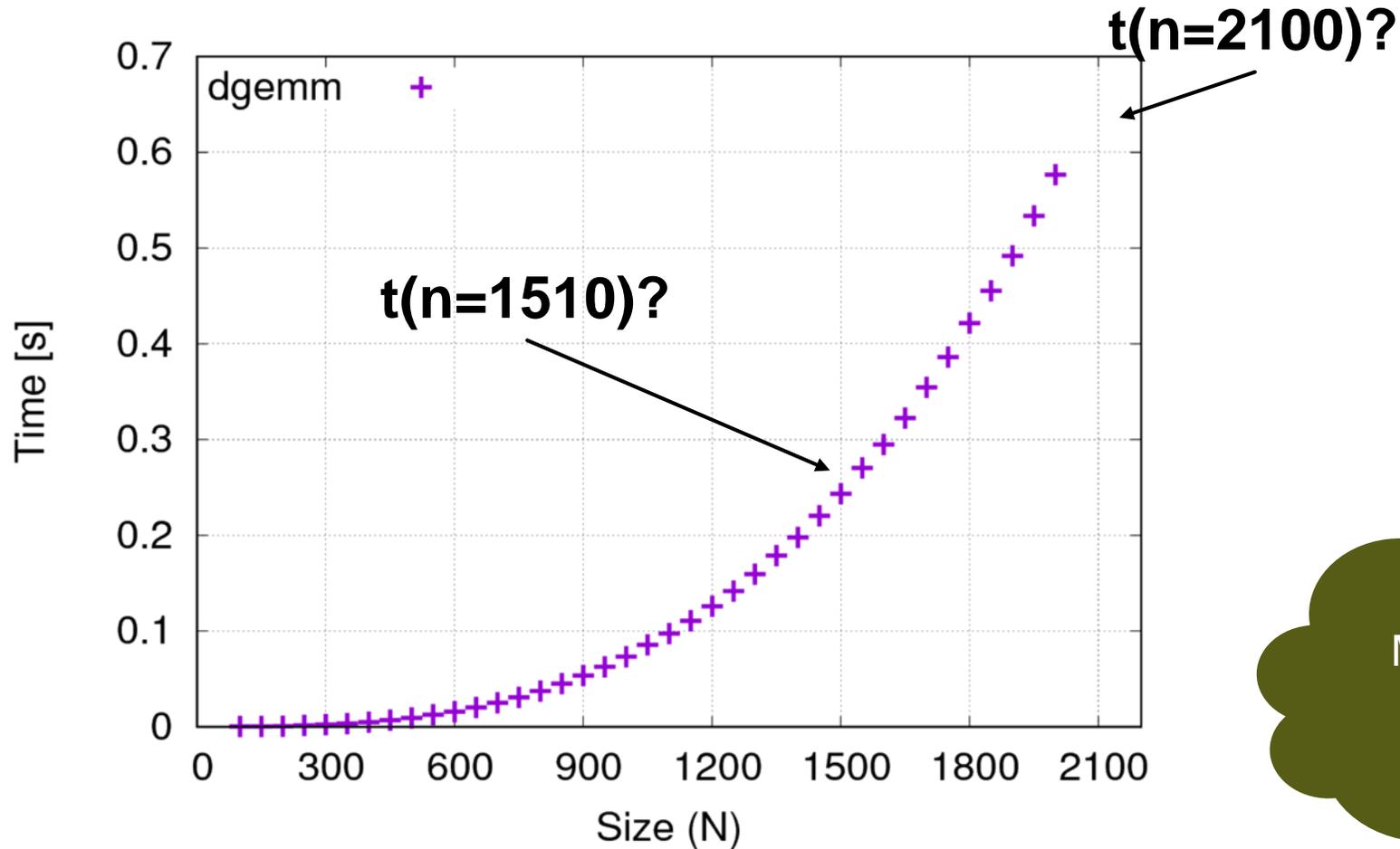
            /* Register the j-th measurement of size sz */
            LSB_Rec(sz);
        }
    }

    LSB_Finalize();
    MPI_Finalize();
    return 0;
}
```

- Simple MPI-like C/C+ interface
- High-resolution timers
- Flexible data collection
- Controlled by environment variables
- Tested up to 512k ranks
- Parallel timer synchronization
- R scripts for data analysis and visualization

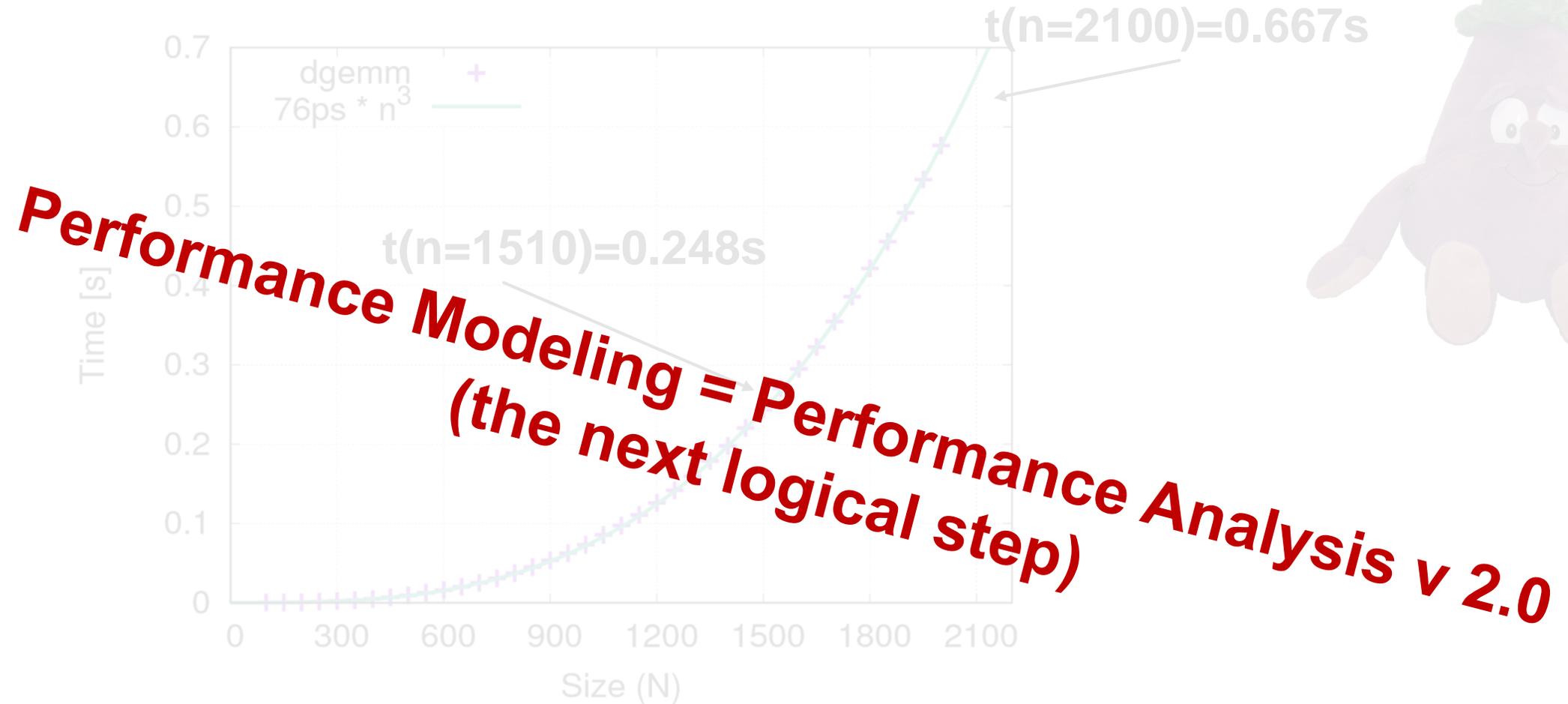


We have the (statistically sound) data, now what?



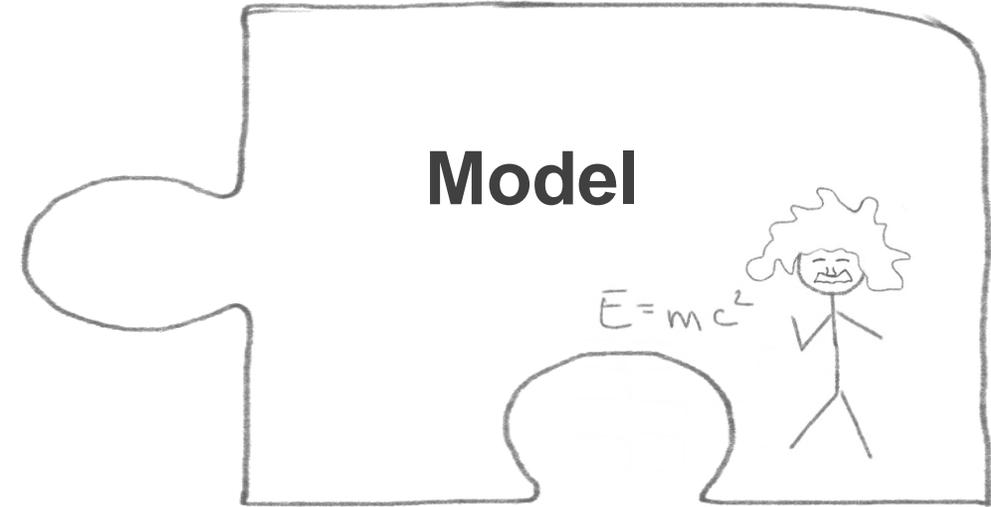
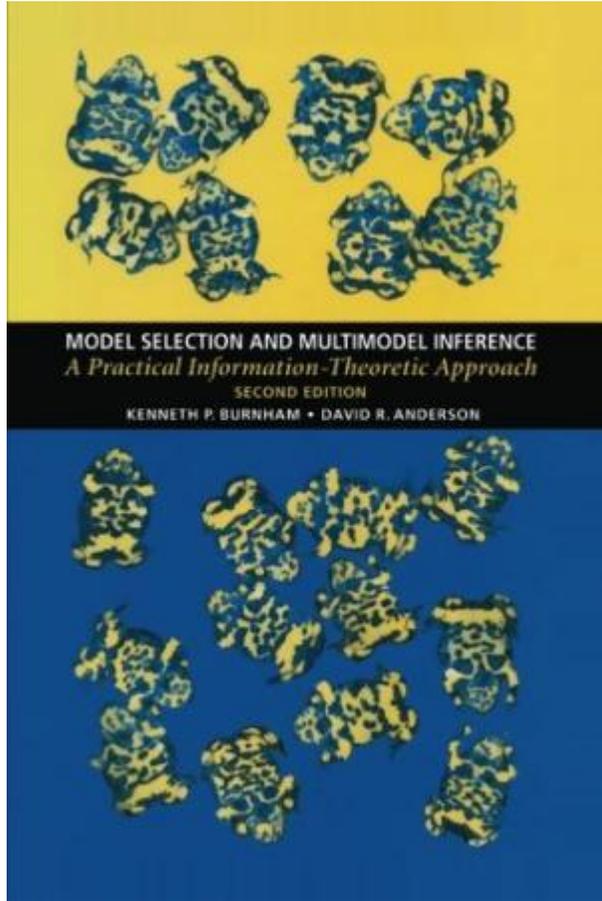
The 99% confidence interval is within 1% of the reported median.

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



The 99% confidence interval is within 1% of the reported median.
The adjusted R² of the model fit is 0.99

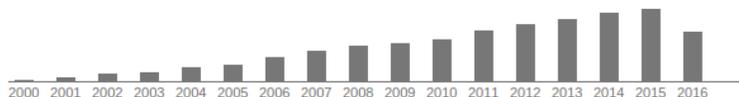
Part II: Model



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.”*

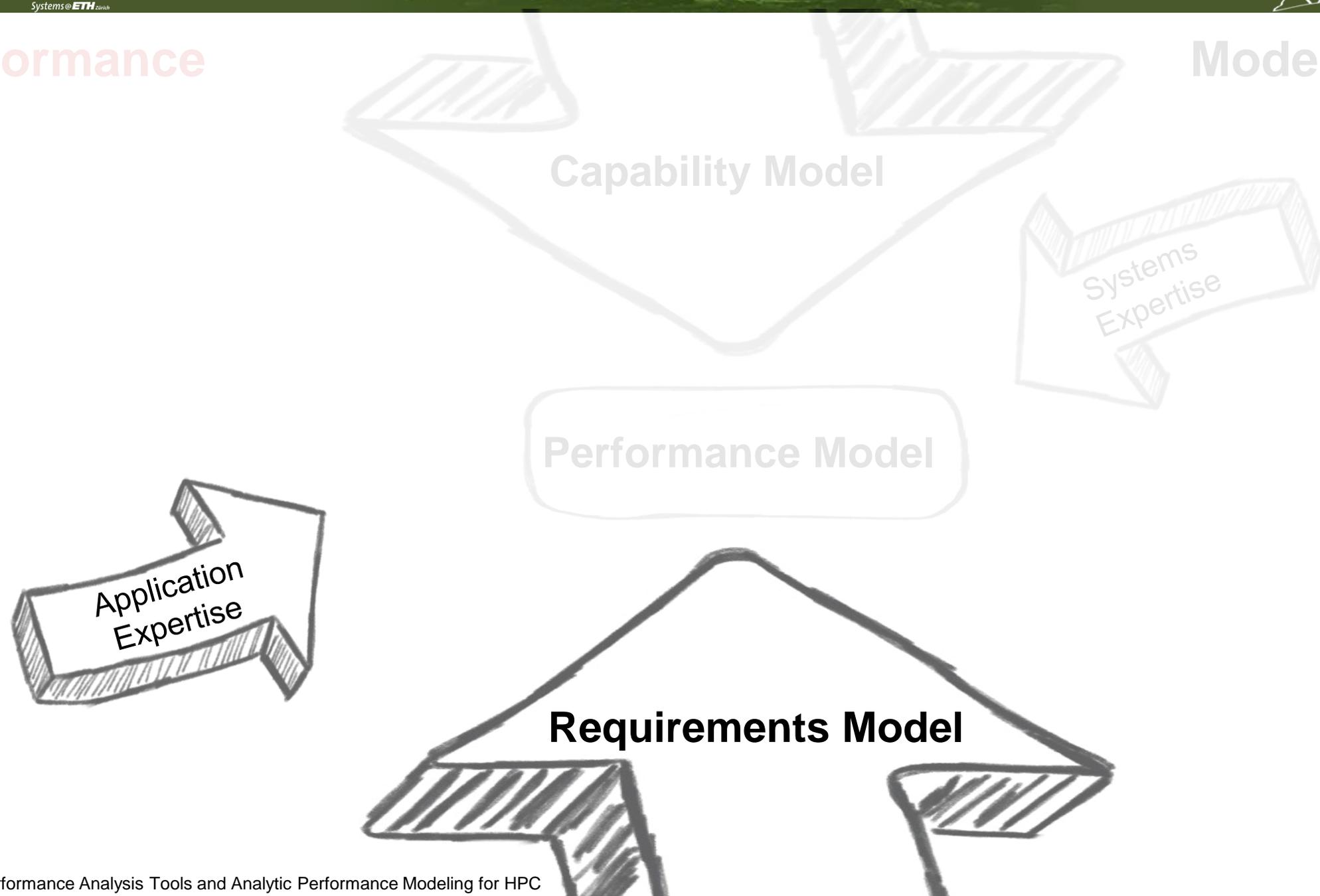
This is generally true for all kinds of modeling.
We focus on **performance modeling** in the following!

Cited by 33599

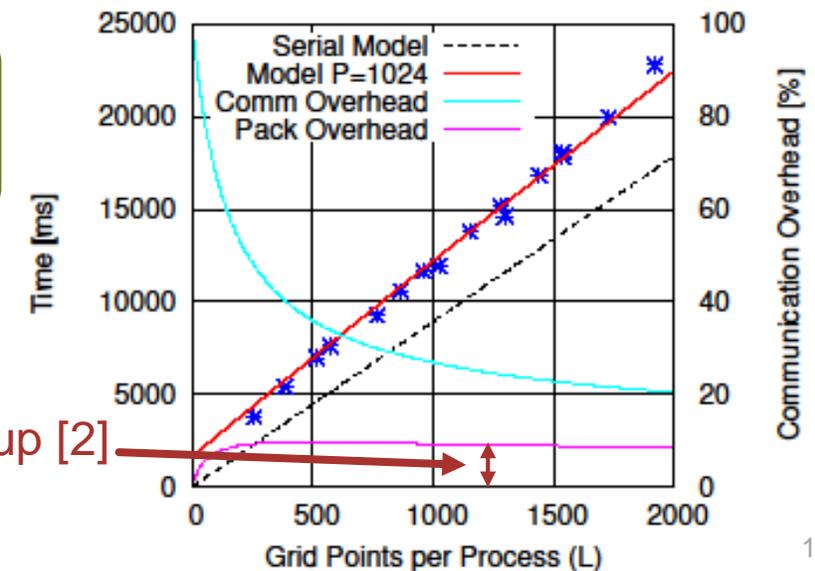
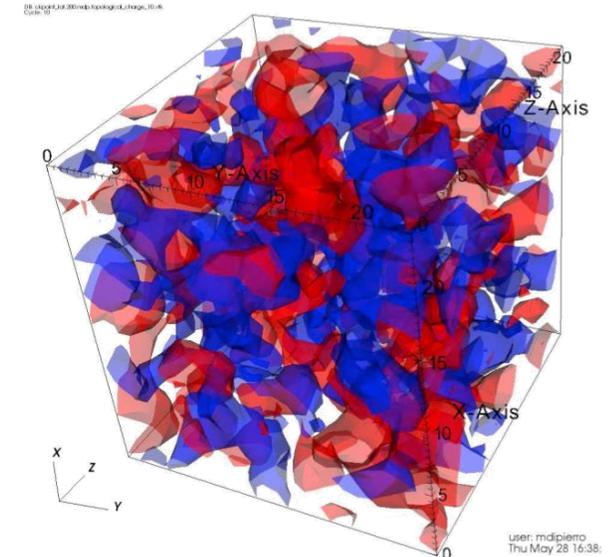
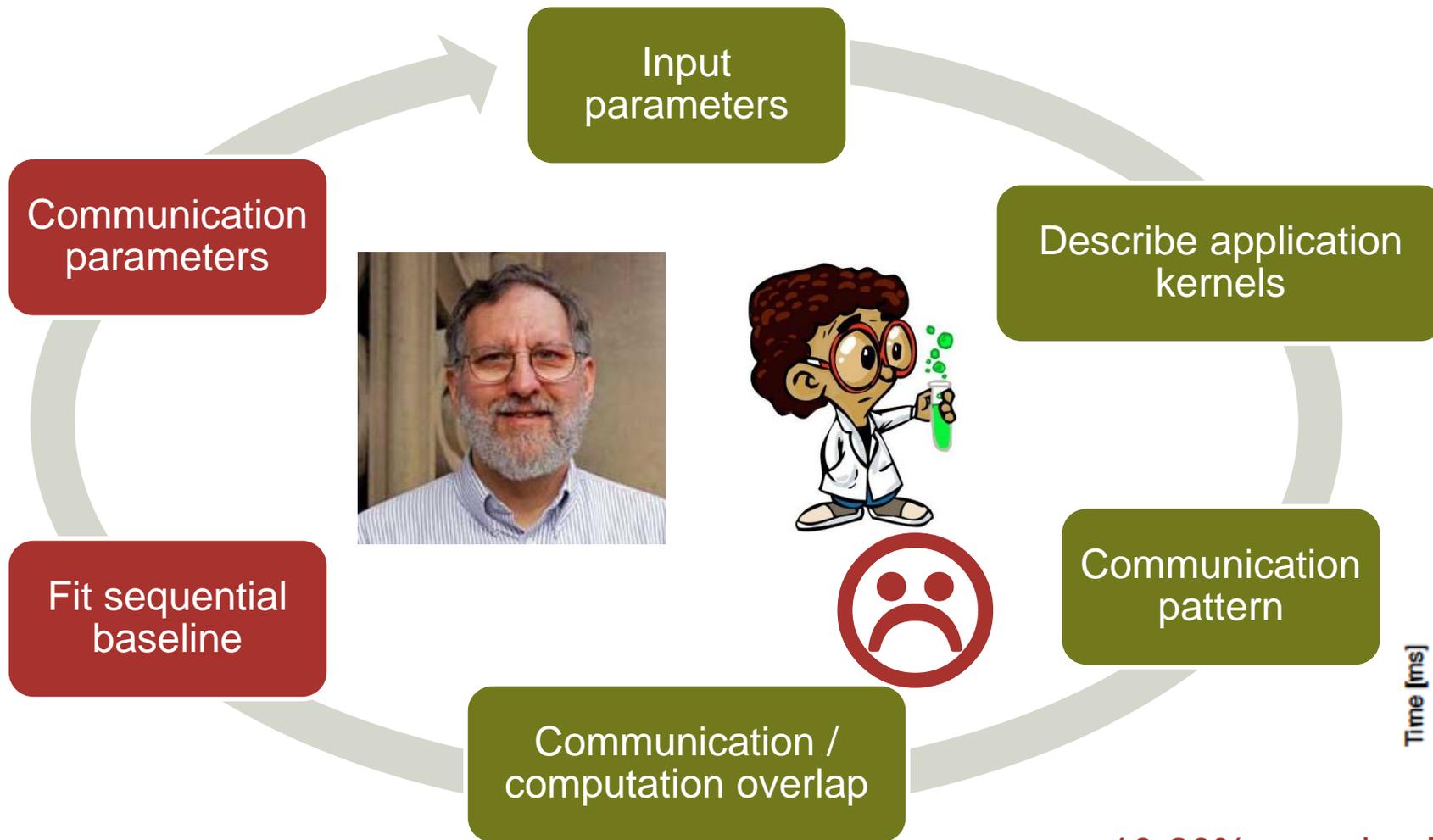


Performance

Modeling



Requirements modeling I: Six-step performance modeling



[1] TH, W. Gropp, M. Snir and W. Kramer: Performance Modeling for Systematic Performance Tuning, SC11
 [2] TH and S. Gottlieb: Parallel Zero-Copy Algorithms for Fast Fourier Transform and Conjugate Gradient using MPI Datatypes, EuroMPI'10

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

Number of processes

$$f(p) = \sum_{k=1}^n c_k \cdot p^{i_k} \cdot \log_2^{j_k}(p)$$

number of terms

(model) constant

$$\begin{aligned} n &\hat{=} \mathbb{N} \\ i_k &\hat{=} I \\ j_k &\hat{=} J \\ I, J &\hat{=} \mathbb{Q} \end{aligned}$$

$$\begin{aligned} n &= 1 \\ I &= \{0, 1, 2\} \\ J &= \{0, 1\} \end{aligned}$$

c_1	$c_1 \times \log(p)$
$c_1 \times p$	$c_1 \times p \times \log(p)$
$c_1 \times p^2$	$c_1 \times p^2 \times \log(p)$

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

$$f(p) = \prod_{k=1}^n c_k \times p^{i_k} \times \log_2^{j_k}(p)$$

$$n = 2$$

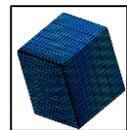
$$I = \{0, 1, 2\}$$

$$J = \{0, 1\}$$

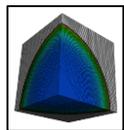
$$\begin{aligned} n &\hat{=} \mathbb{N} \\ i_k &\hat{=} I \\ j_k &\hat{=} J \\ I, J &\hat{=} \mathbb{Q} \end{aligned}$$

$c_1 \cdot \log(p) + c_2 \cdot p$
 $c_1 \cdot \log(p) + c_2 \cdot p \cdot \log(p)$
 $c_1 \cdot \log(p) + c_2 \cdot p^2$
 $c_1 + c_2 \times p$
 $c_1 + c_2 \times p^2$
 $c_1 + c_2 \times \log(p)$
 $c_1 + c_2 \times p \times \log(p)$
 $c_1 + c_2 \times p^2 \times \log(p)$
 $c_1 \cdot p + c_2 \cdot p^2$
 $c_1 \cdot p + c_2 \cdot p^2 \cdot \log(p)$
 $c_1 \cdot p + c_2 \cdot p \cdot \log(p)$
 $c_1 \cdot p \cdot \log(p) + c_2 \cdot p^2$
 $c_1 \cdot p \cdot \log(p) + c_2 \cdot p^2 \cdot \log(p)$
 $c_1 \cdot p^2 + c_2 \cdot p^2 \cdot \log(p)$

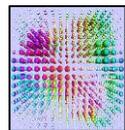
Tool support: Extra-P for automated best-fit modeling [1]



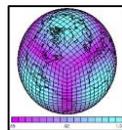
Sweep3d



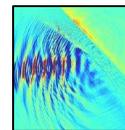
Lulesh



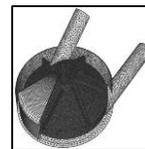
Milc



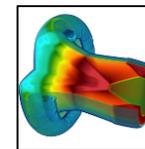
HOMME



JUSPIC



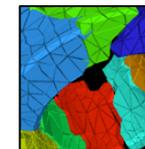
XNS



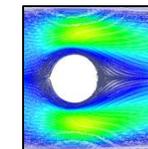
BLAST



NEST



UG4



MP2C

Talk: **Fast Multi-Parameter Performance Modeling**

A. Calotoiu, et al.

Tomorrow!!

10:30am

Room: Grand Hall



Tutorial: **Insightful Automatic Performance Modeling**

A. Calotoiu, F. Wolf, TH, M. Schulz

Sunday, November 13th

1:30pm - 5pm

Room 355-C



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

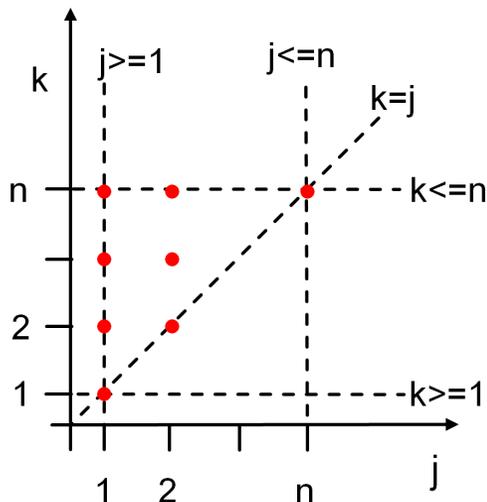
[2] A. Calotoiu, D. Beckingsale, C. W. Earl TH, I. Karlin, M. Schulz, F. Wolf: Fast Multi-Parameter Performance Modeling, IEEE Cluster 2016

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);
```



$$N = (n + 1) \log_2 n - n + 2$$

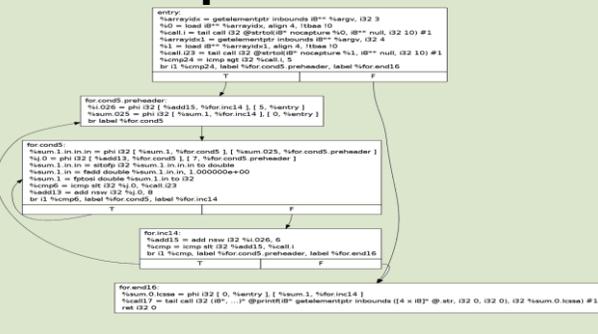
Parallel program

```
do i = 1, procCols
  call mpi_irecv( buff, 2, dp_type, reduce_exch_proc(i),
    > i, mpi_comm_world, request, ierr )
  > call mpi_send( buff2, 2, dp_type, reduce_exch_proc(i),
    i, mpi_comm_world, ierr )
  > call mpi_wait( request, status, ierr )
enddo

do i = id * n/p, ( id + 1 ) * n/p
  do j = 1, nSize
    call compute
```

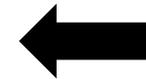


Loop extraction



Number of iterations

$$N = \sum_{i_1=0}^{n_1(x_0,1)} \sum_{i_2=0}^{n_2(x_0,2)} \dots \sum_{i_{r-1}=0}^{n_{r-1}(x_0,r-1)} n_r(x_0,r)$$



Requirements Models

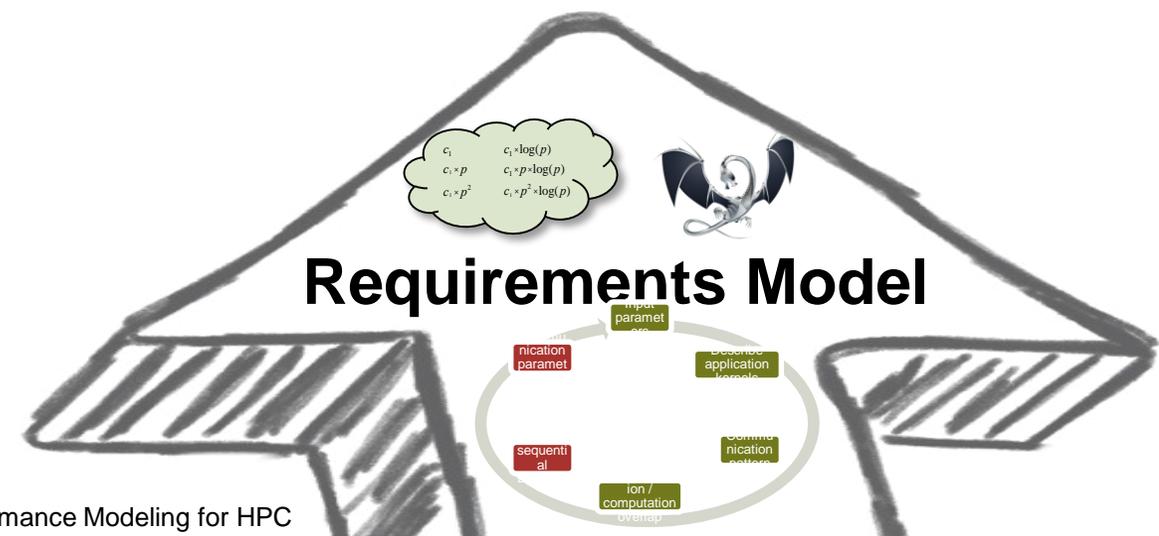
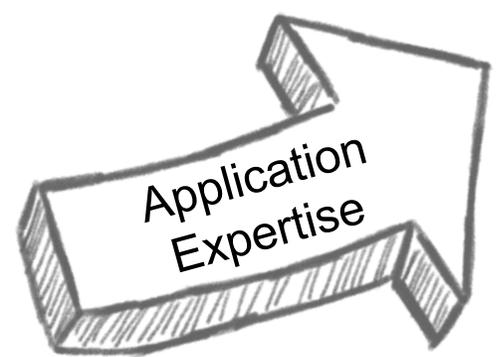
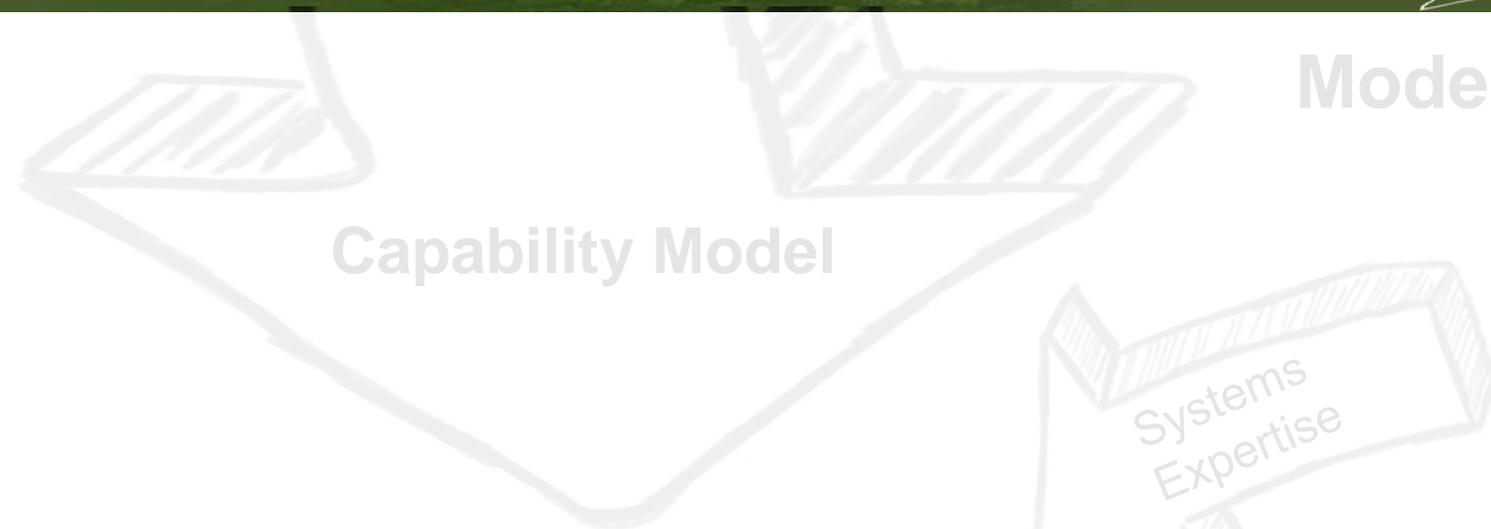
$$W = N \Big|_{p=1}$$

$$D = N \Big|_{p \rightarrow \infty}$$

[1]: TH, G. Kwasniewski: Automatic Complexity Analysis of Explicitly Parallel Programs, ACM SPAA'14

Performance

Modeling



Performance

Modeling

Capability Model

Systems Expertise

Performance Model

Application Expertise

Requirements Model

$$\begin{matrix} c_1 & c_1 \cdot \log(p) \\ c_1 \cdot p & c_1 \cdot p \cdot \log(p) \\ c_1 \cdot p^2 & c_1 \cdot p^2 \cdot \log(p) \end{matrix}$$



Capability models for network communication

The LogP model family and the LogGOPS model [1]



LogP

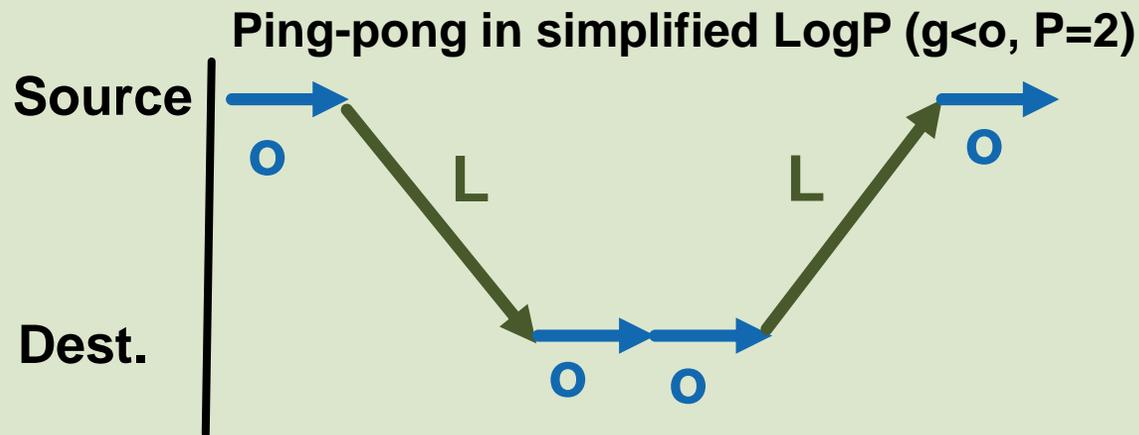
A PRACTICAL MODEL of PARALLEL COMPUTATION

OUR 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 real machine.

PRAM consists of a collection of processors which compute synchronously in parallel and communicate with a global random access memory.

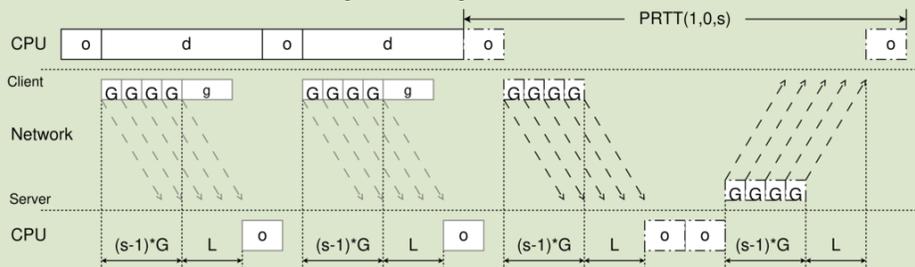
A new parallel machine model reflects the critical technology trends underlying parallel computers

David E. Culler, Richard M. Karp, David Patterson, Abhijit Sahay, Eunice E. Santos, Klaus Erik Schauer, Ramesh Subramonian, and Thorsten von Eicken



Finding LogGOPS parameters

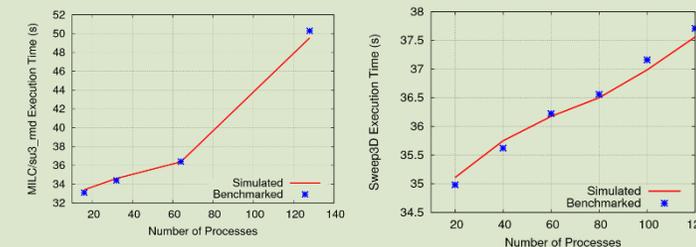
Netgauge [2], model from first principles, fit to data using special kernels



Large scale LogGOPS Simulation

LogGOPSim [1], simulates LogGOPS with 10 million MPI ranks

<5% error



[1]: TH, T. Schneider and A. Lumsdaine: LogGOPSIm - Simulating Large-Scale Applications in the LogGOPS Model, LSAP 2010, <https://spcl.inf.ethz.ch/Research/Performance/LogGOPSIm/>

[2]: TH, T. Mehlman, A. Lumsdaine and W. Rehm: Netgauge: A Network Performance Measurement Framework, HPC 2007, <https://spcl.inf.ethz.ch/Research/Performance/Netgauge/>

Performance

Modeling



Capability Model



Performance Model



Requirements Model



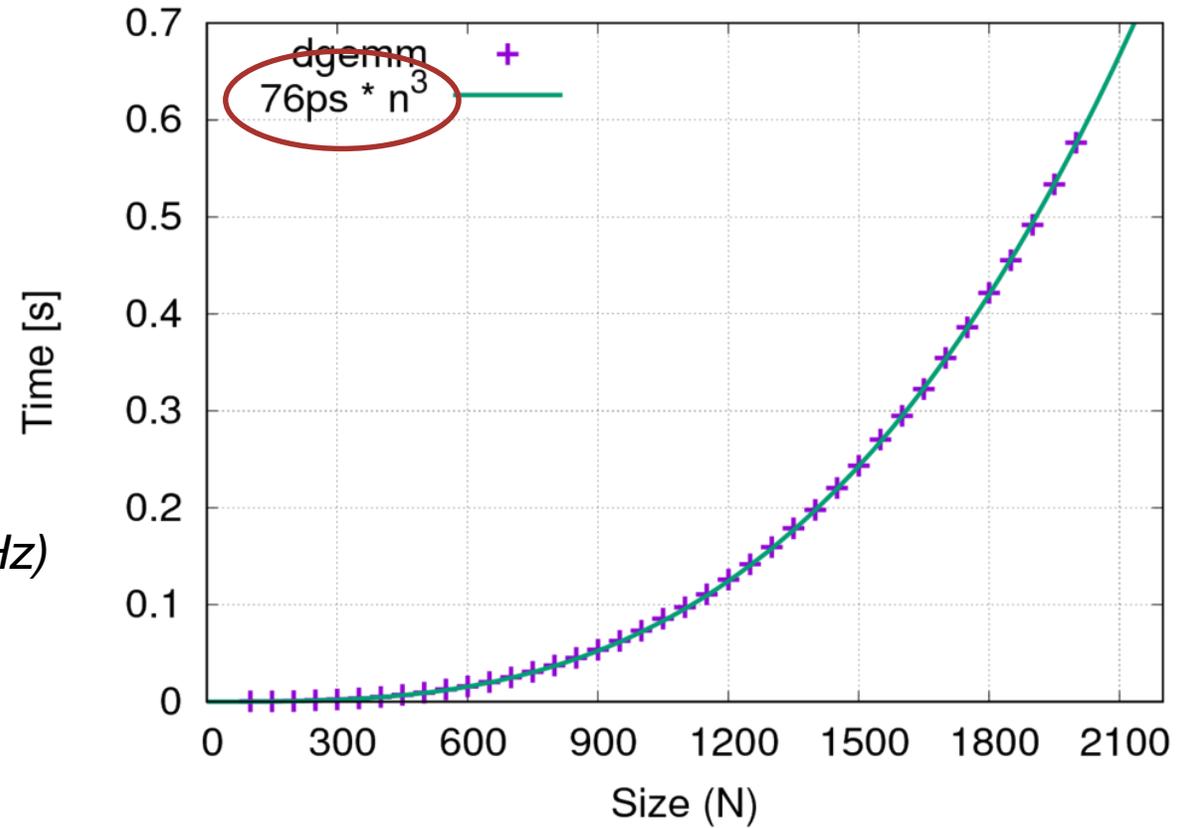
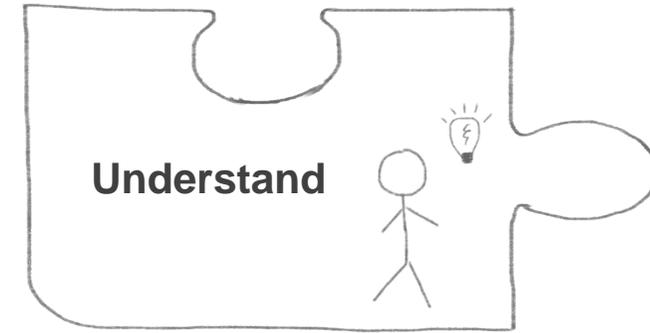
Part III: Understand

- **Use models to**
 1. Proof optimality of real implementations
 - *Stop optimizing, step back to algorithm level*
 2. Design optimal algorithms or systems in the model
 - *Can lead to non-intuitive designs*

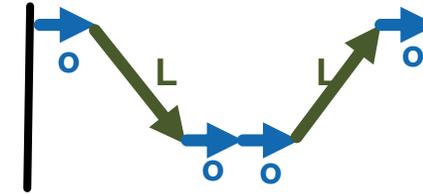
- **Proof optimality of matrix multiplication**
 - Intuition: flop rate is the bottleneck
 - $t(n) = 76ps * n^3$
 - **Flop rate** $R = 2flop * n^3 / (76ps * n^3) = 27.78 \text{ Gflop/s}$
 - **Flop peak:** $3.864 \text{ GHz} * 8 \text{ flops} = 30.912 \text{ Gflop/s}$
Achieved ~90% of peak (IBM Power 7 IH @3.864GHz)



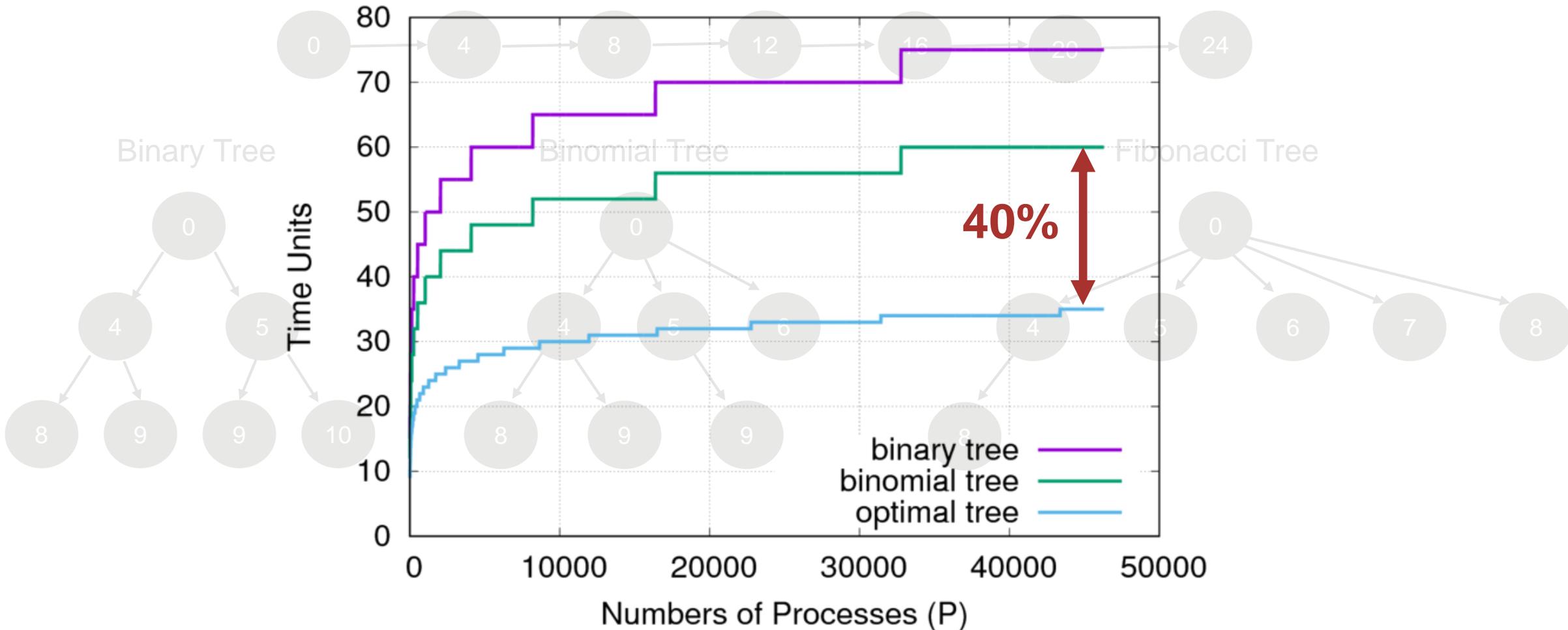
- **Gets more complex quickly**
 - Imagine sparse matrix-vector



2) Design optimal algorithms – small broadcast in LogP

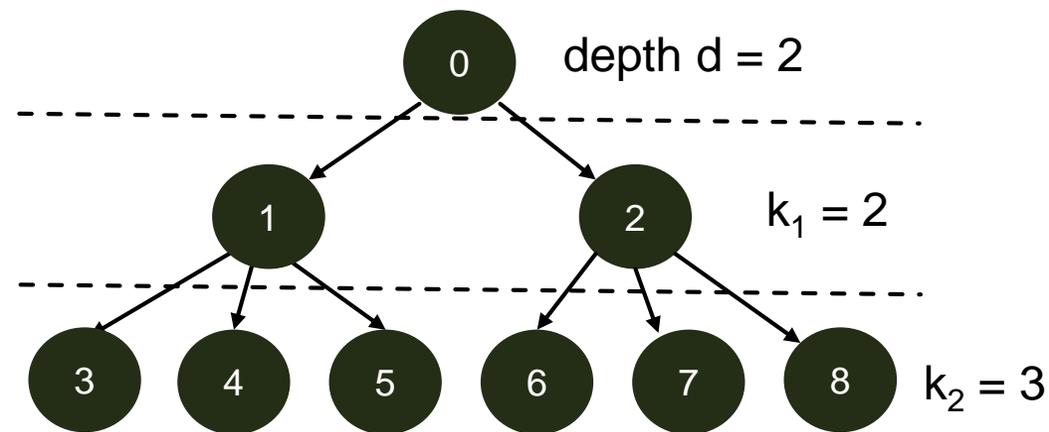


$L=2, o=1, P=7$



Design algorithms – bcast in cache-to-cache model

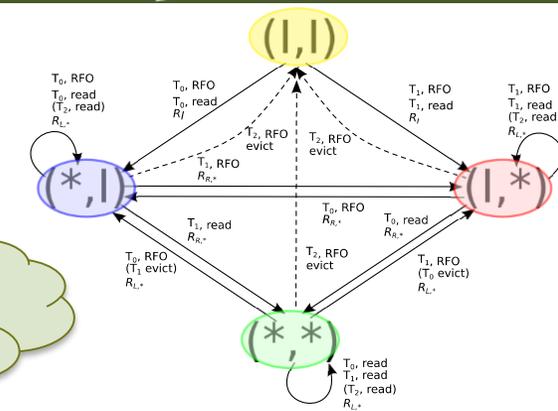
Multi-ary tree example



Tree depth

Level size

Tree cost



$$\mathcal{T}_{tree} = \sum_{i=1}^d \mathcal{T}_C(k_i) = \sum_{i=1}^d (c \cdot k_i + b)$$

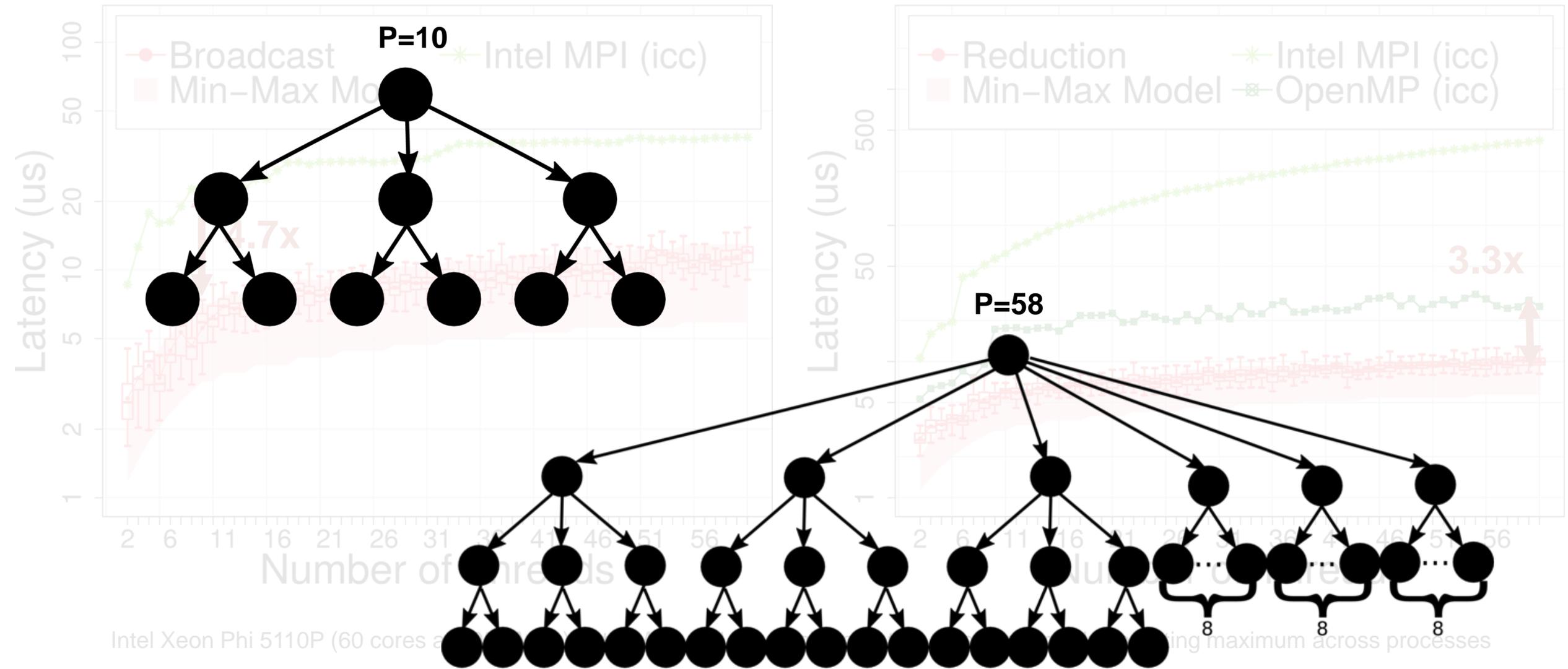
$$= \sum_{i=1}^d (R_R + R_L + c \cdot (k_i - 1))$$

$$\mathcal{T}_{sbcast} = \min_{d, k_i} \left(\mathcal{T}_{fw} + \sum_{i=1}^d (c \cdot k_i + b) + \sum_{i=1}^d \mathcal{T}_{nb}(k_i + 1) \right)$$

Reached threads

$$N \leq 1 + \sum_{i=1}^d \prod_{j=1}^i k_j, \quad \forall i < j, k_i \leq k_j$$

Measured results – small broadcast and reduction

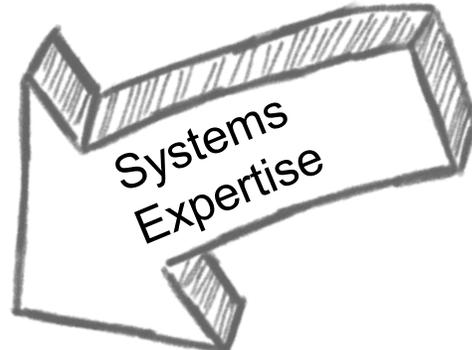
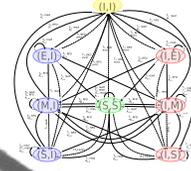


Performance

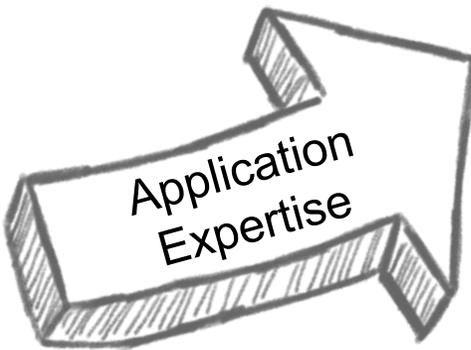
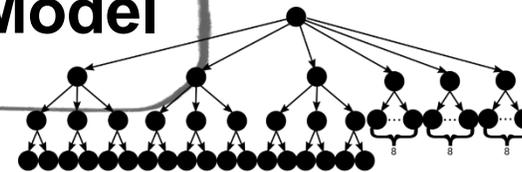
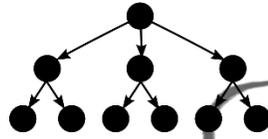
Modeling



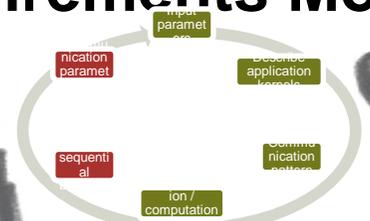
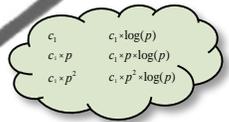
Capability Model



Performance Model



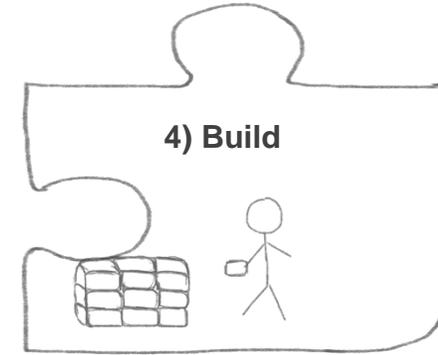
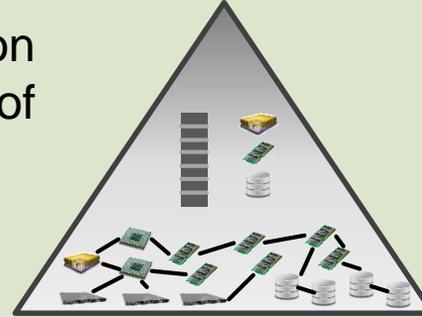
Requirements Model



Part IV: Build

Abstraction is Key

- Enables to focus on essential aspects of a system

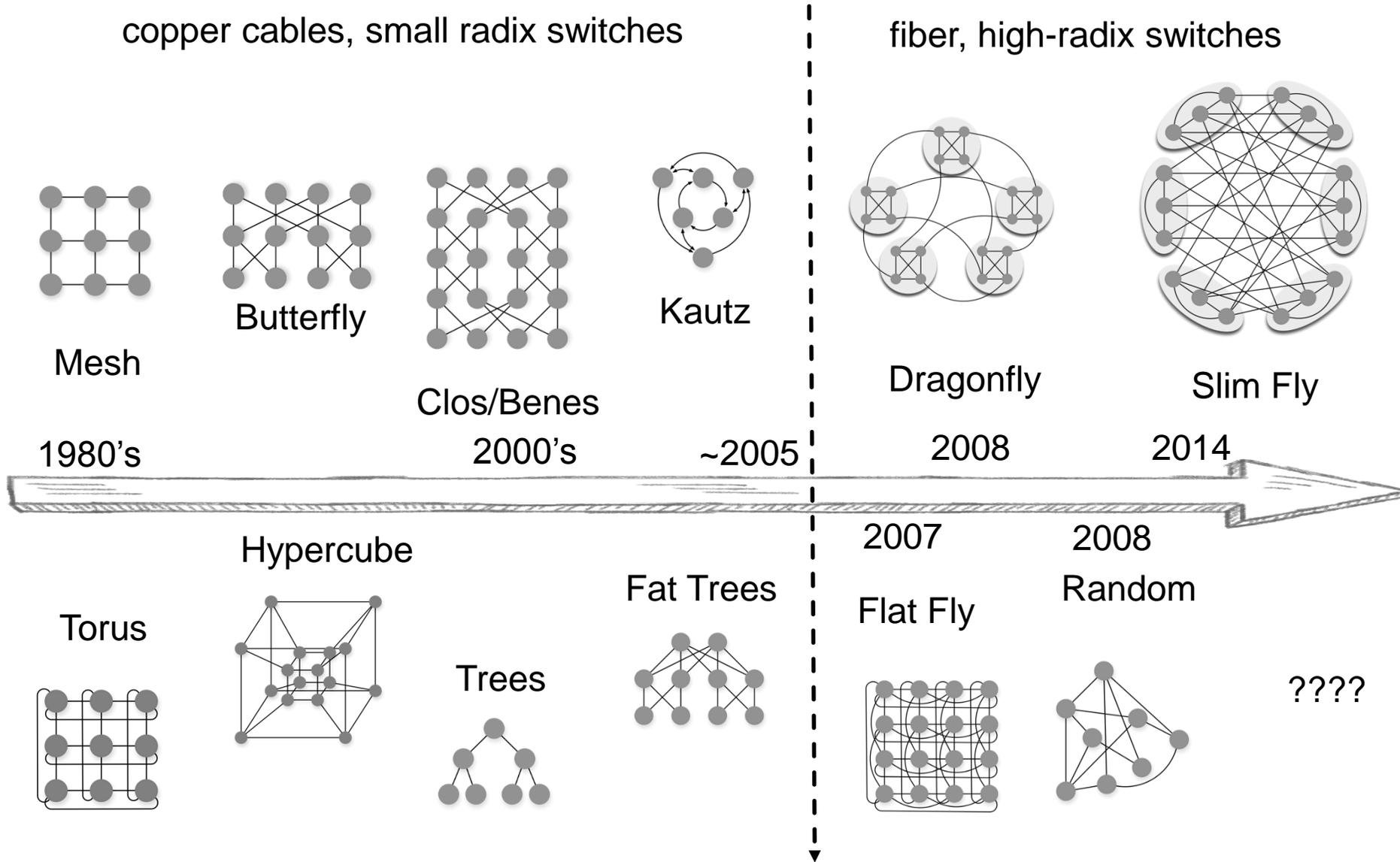


Case study: Network Topologies

- **Observe:** optimize for cost, maintain performance:
 - router radix, number of cables, number of routers → cost
 - number of endpoints, latency, global bandwidth → capabilities
- **Model:** system as graph
- **Understand:** degree-diameter graphs
- **Build:** Slim Fly topology
- Result: non-trivial topology that outperforms all existing ones



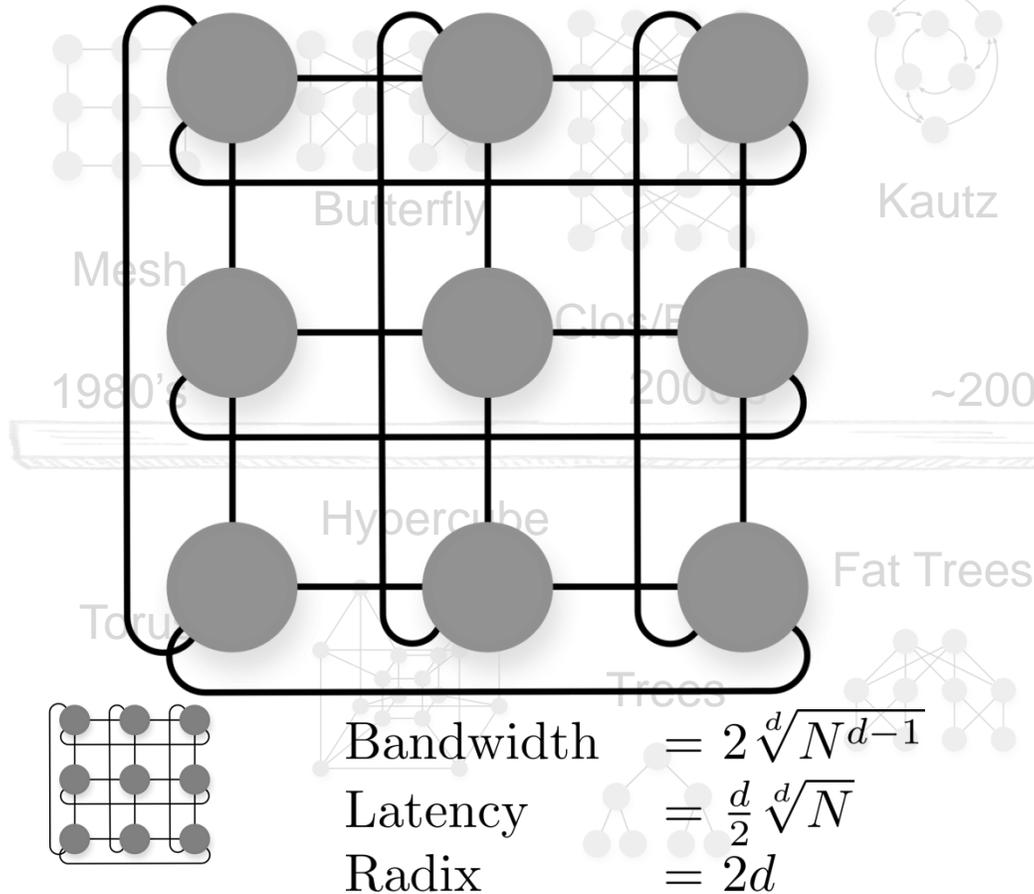
A BRIEF HISTORY OF NETWORK TOPOLOGIES



A BRIEF HISTORY OF NETWORK TOPOLOGIES

copper cables, small radix switches

fiber, high-radix switches



2008

2014

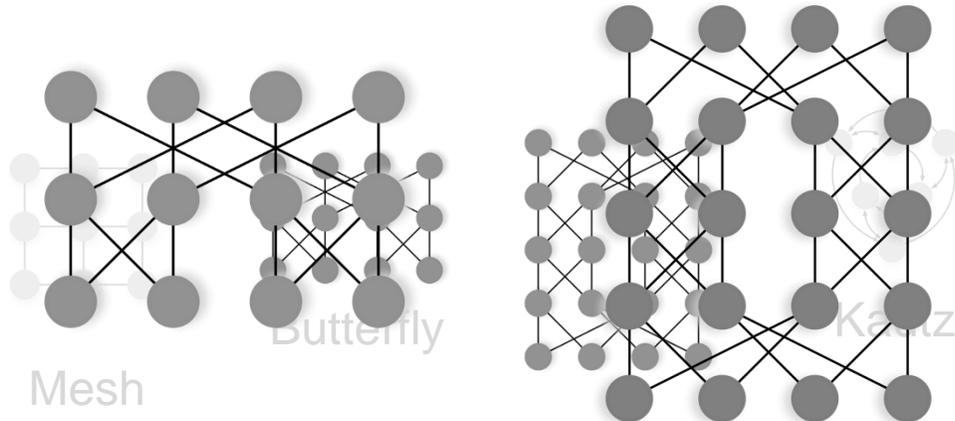


2010

Flattened

A BRIEF HISTORY OF NETWORK TOPOLOGIES

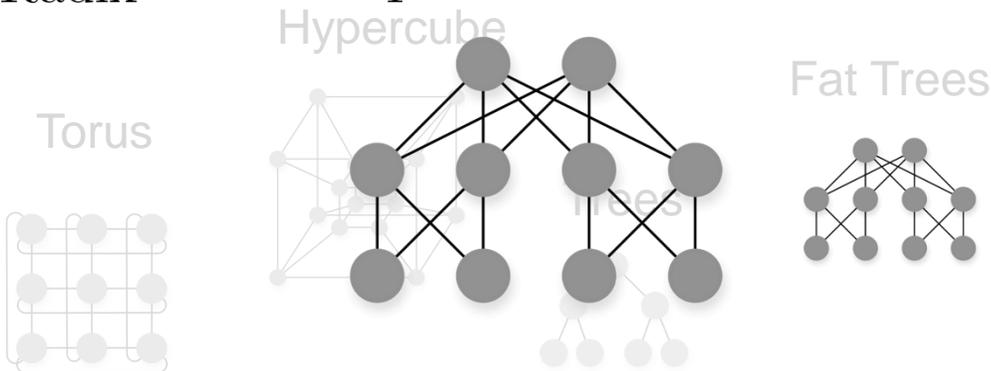
copper cables, small radix switches



Bandwidth = $\frac{N}{2}$

Latency = $2 \log_2 N$

Radix = 4



Dragonfly Slim Fly

2008

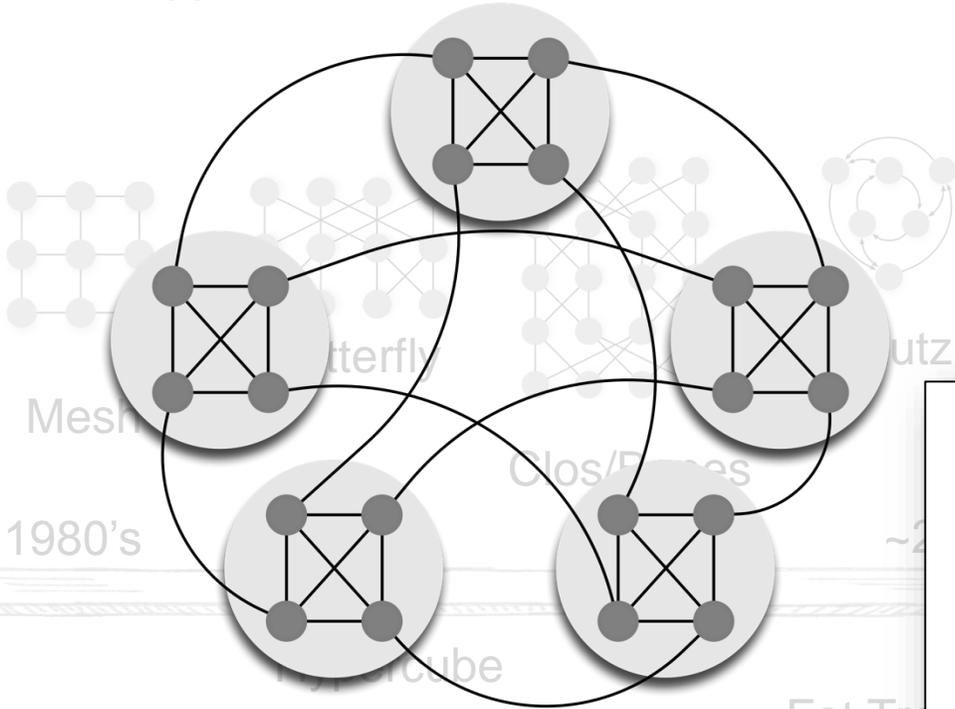
2014



A BRIEF HISTORY OF NETWORK TOPOLOGIES

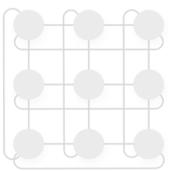
copper cables, small radix switches

fiber



1980's

Torus



$$\begin{aligned} \text{Bandwidth} &\approx \frac{N}{4} \\ \text{Latency} &= 3 - 5 \\ \text{Radix} &= 48 - 64 \end{aligned}$$

International Symposium on Computer Architecture

Technology-Driven, Highly-Scalable Dragonfly Topology*

John Kim
Northwestern University
Evanston, IL 60208
jki12@northwestern.edu

William J. Dally
Stanford University

Steve Scott
Cray Inc.

Dennis Abts
Google Inc.

2010 18th IEEE Symposium on High Performance Interconnects

The PERCS High-Performance Interconnect

Baba Arimilli *, Ravi Arimilli *, Vicente Chung *, Scott Clark *, Wolfgang Denzel †, Ben Drerup *, Torsten Hoefler ‡, Jody Joyner *, Jerry Lewis *, Jian Li †, Nan Ni * and Ram Rajamony †

* IBM Systems and Technology Group, 11501 Burnet Road, Austin, TX 78758
† IBM Research (Austin, Zurich), 11501 Burnet Road, Austin, TX 78758
‡ Blue Waters Directorate, NCSA, University of Illinois at Urbana-Champaign, Urbana, IL 61801
E-mail: arimilli@us.ibm.com, rajamony@us.ibm.com, htor@illinois.edu

Abstract—The PERCS system was designed by IBM in response to a DARPA challenge that called for a high-productivity high-performance computing system. A major innovation in the PERCS design is the network that is built using Hub chips that are integrated into the compute nodes. Each Hub chip is about 580 mm² in size, has over 3700 signal I/Os, and is packaged in a module that also contains LGA-attached optical electronic devices.

The Hub module implements five types of high-bandwidth interconnects with multiple links that are fully-connected with a high-performance internal crossbar switch. These links provide over 9 Tbits/second of raw bandwidth and are used to construct a two-level direct connect topology spanning up to tens of thou-

bandwidths do not scale accordingly. For instance, while High Performance Linpack performance [5], [10] shows a steady improvement over time, interconnect-intensive metrics such as G-RandomAccess and G-FFTE [5] show very little improvement.

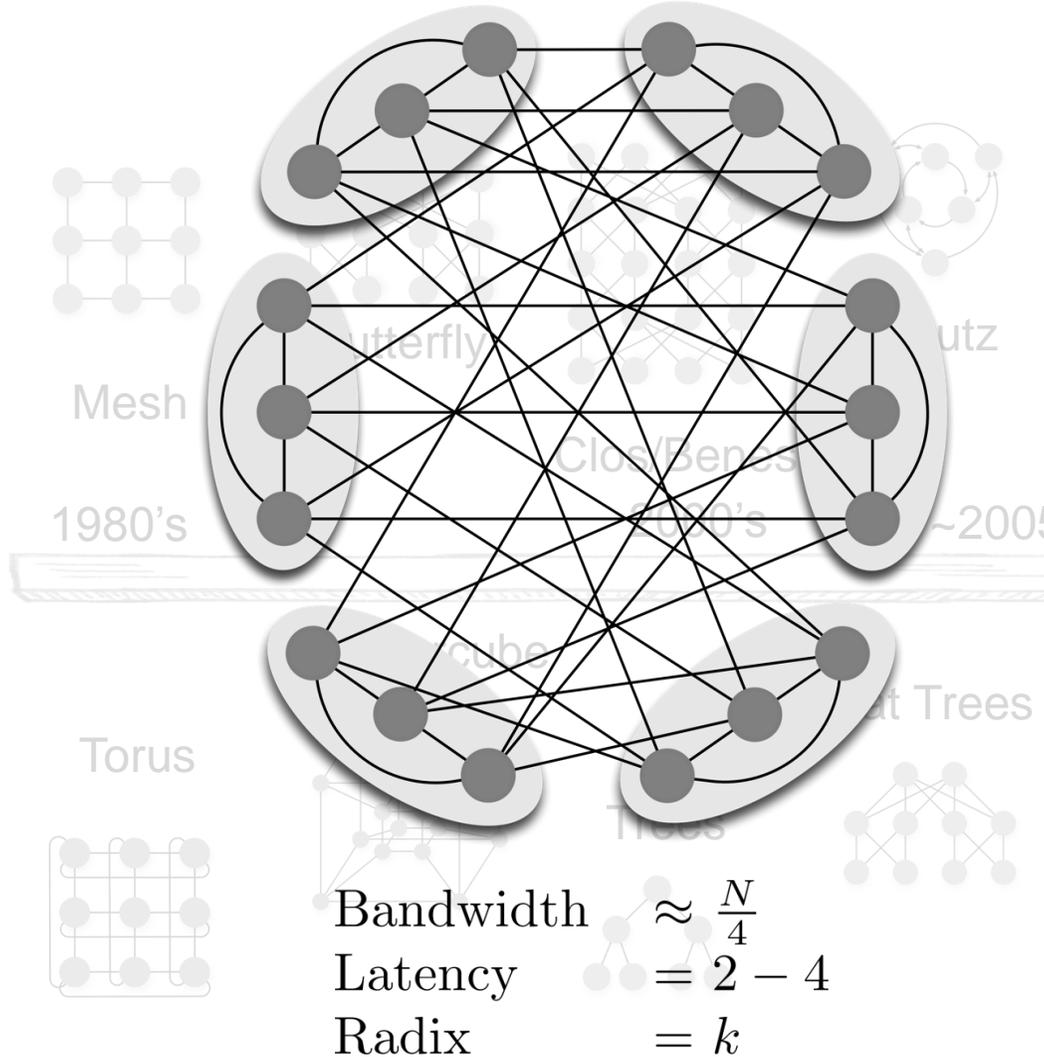
The challenge of building a high-performance, highly productive, multi-Petaflop system forced us to recognize early on that the entire infrastructure had to scale along with the microprocessor's capabilities. A significant component of our scaling solution is a new switchless interconnect with very high fanout organized into a two-level direct connect

A BRIEF HISTORY OF NETWORK TOPOLOGIES



copper cables, small radix switches

fiber, high-radix switches



Key insight:

“It’s the diameter, stupid”

Lower diameter:

- Fewer cables traversed
- Fewer cables needed
- Fewer routers needed

Cost and energy savings:

- Up to 50% over Fat Tree
- Up to 33% over Dragonfly

DESIGNING LOWEST-DIAMETER TOPOLOGIES: A ONE-MINUTE PROOF

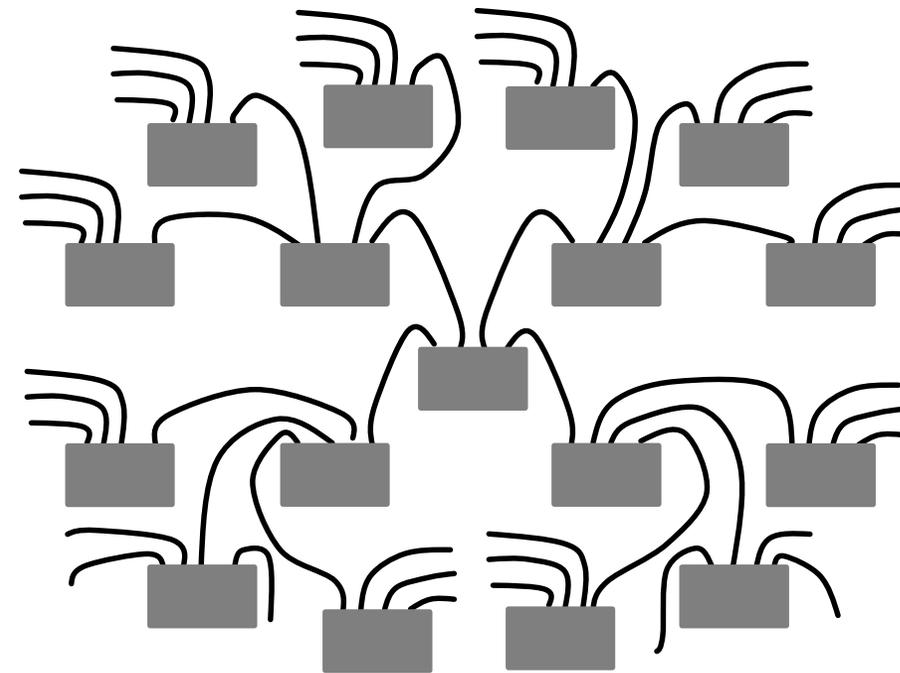


Key method

Optimize towards the Moore Bound [1]: the upper bound on the *number of vertices* in a graph with given *diameter* D and *radix* k .

$$MB(D, k) = 1 + k + k(k-1) + k(k-1)^2 + \dots$$

$$MB(D, k) = 1 + k \sum_{i=0}^{D-1} (k-1)^i$$

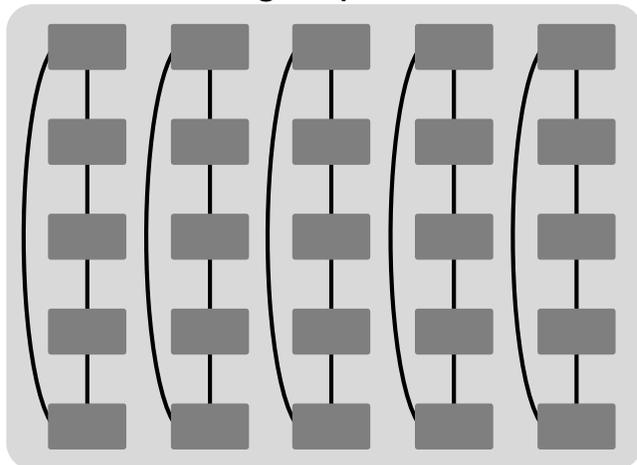


DESIGNING AN EFFICIENT NETWORK TOPOLOGY

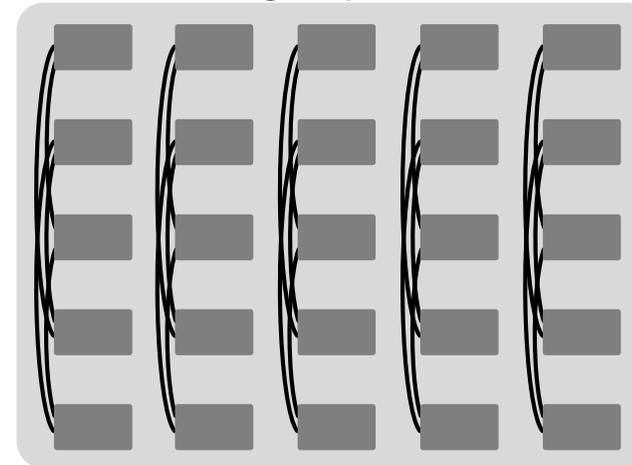
CONNECTING ROUTERS: DIAMETER 2

Example Slim Fly design for $diameter = 2$: *MMS graphs* [1]

A subgraph with
identical groups of routers



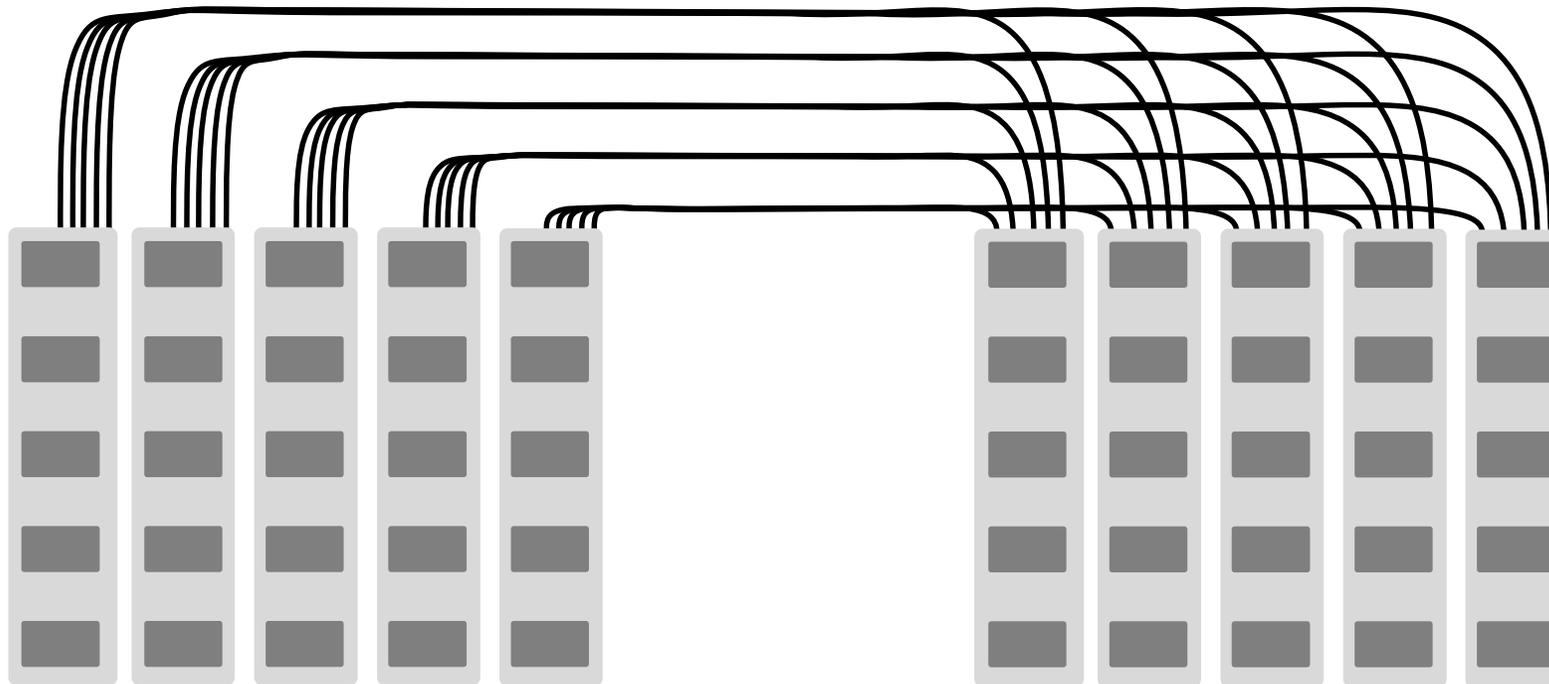
A subgraph with
identical groups of routers



[1] B. D. McKay, M. Miller, and J. Siráň. A note on large graphs of diameter two and given maximum degree. *Journal of Combinatorial Theory, Series B*, 74(1):110 – 118, 1998

DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2



Groups form a fully-connected bipartite graph

DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

1 Select a prime power q

$$q = 4w + \delta;$$

$$w \in \mathbb{N} \quad \delta \in \{-1, 0, 1\},$$

A Slim Fly based on q :
Number of routers: $2q^2$
Network radix: $(3q - \delta)/2$

2 Construct a finite field \mathcal{F}_q .

Assuming q is prime:

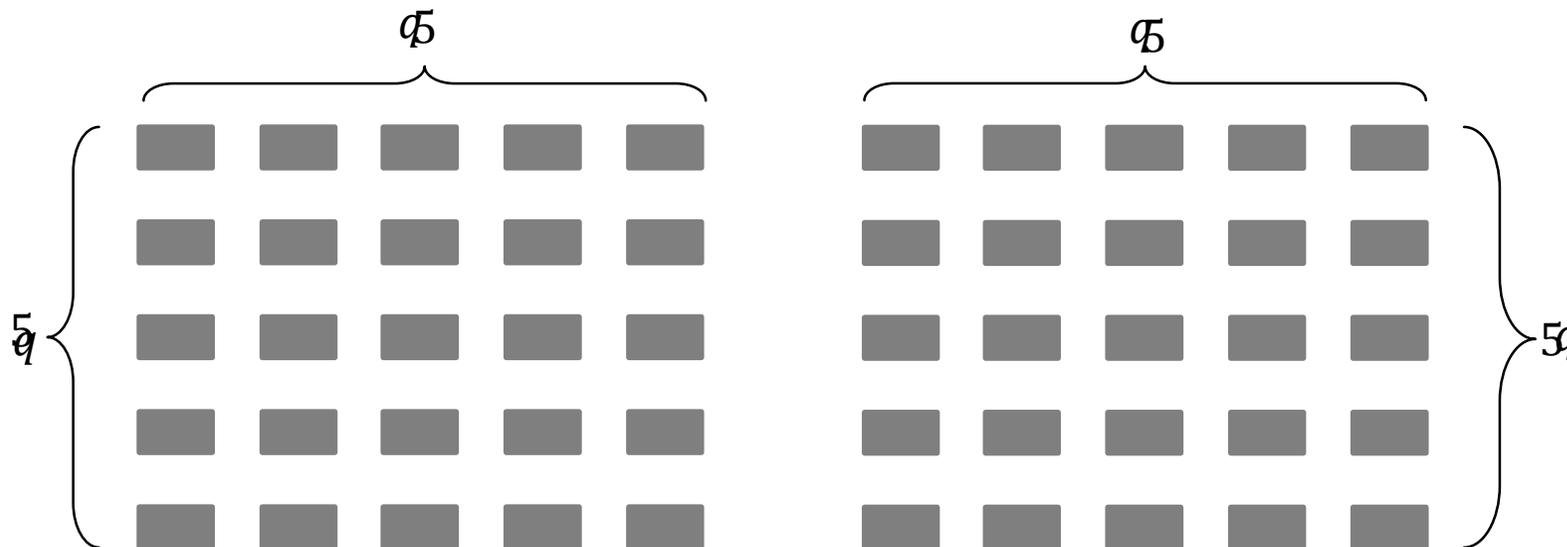
$$\mathcal{F}_q = \mathbb{Z}/q\mathbb{Z} = \{0, 1, \dots, q - 1\}$$

with modular arithmetic.

E Example: $q = 5$

50 routers
network radix: 7

$$\mathcal{F}_5 = \{0, 1, 2, 3, 4\}$$



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

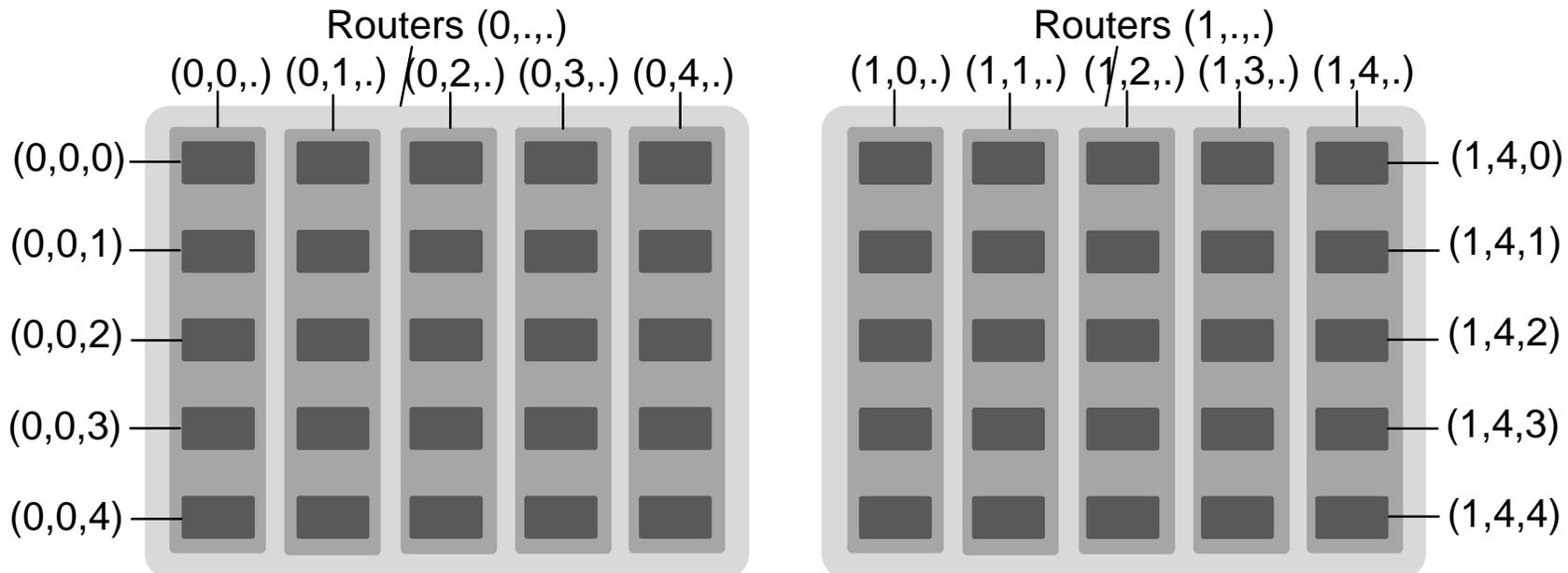
3 Label the routers

Set of routers:

$$\{0,1\} \times \mathcal{F}_q \times \mathcal{F}_q$$

E Example: $q = 5$

...



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

4 Find primitive element ξ

$\xi \in \mathcal{F}_q$ generates \mathcal{F}_q :

All non-zero elements of \mathcal{F}_q can be written as ξ^i ; $i \in \mathbb{N}$

5 Build Generator Sets

$$X = \{1, \xi^2, \dots, \xi^{q-3}\}$$

$$X' = \{\xi, \xi^3, \dots, \xi^{q-2}\}$$

E Example: $q = 5$

$$\mathcal{F}_5 = \{0, 1, 2, 3, 4\}$$

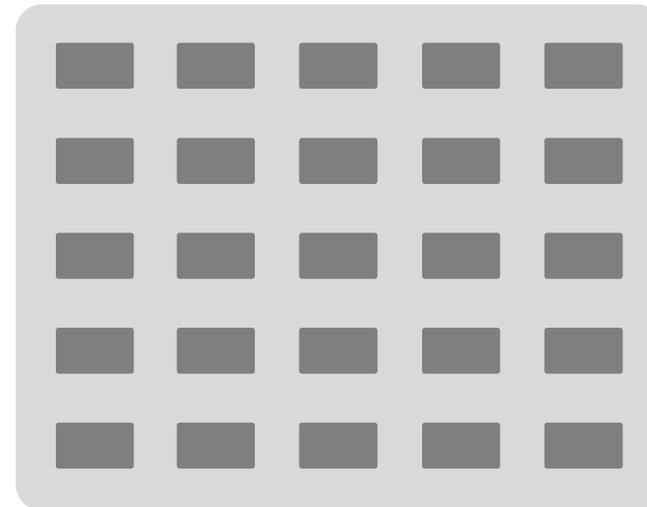
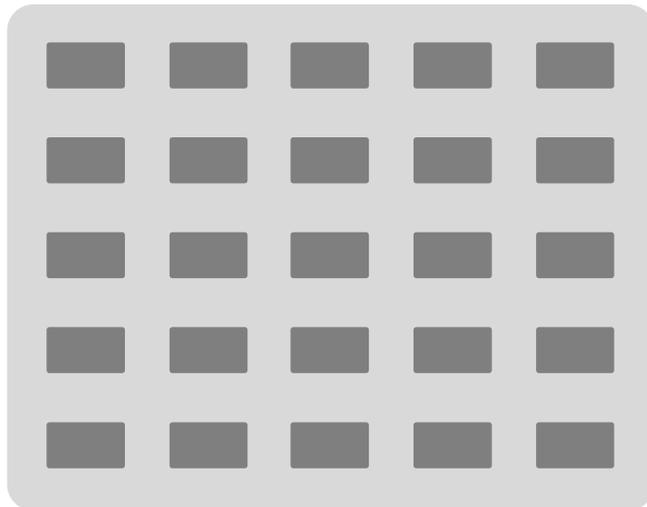
$$\xi = 2$$

$$1 = \xi^4 \text{ mod } 5 =$$

$$2^4 \text{ mod } 5 = 16 \text{ mod } 5$$

$$X = \{1, 4\}$$

$$X' = \{2, 3\}$$



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

6 Intra-group connections

Two routers in one group are connected iff their “vertical Manhattan distance” is an element from:

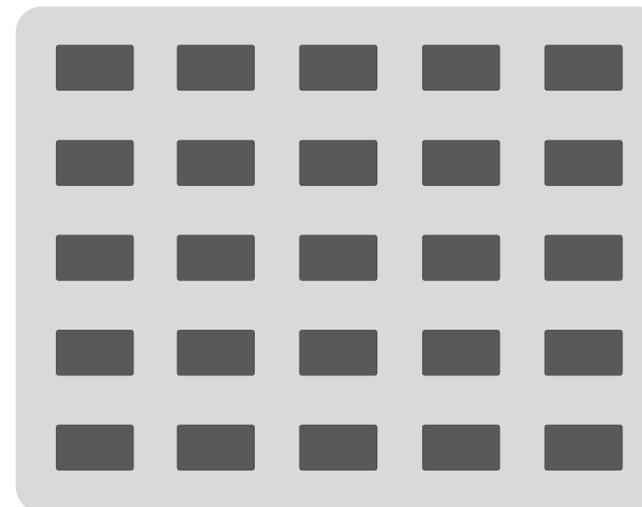
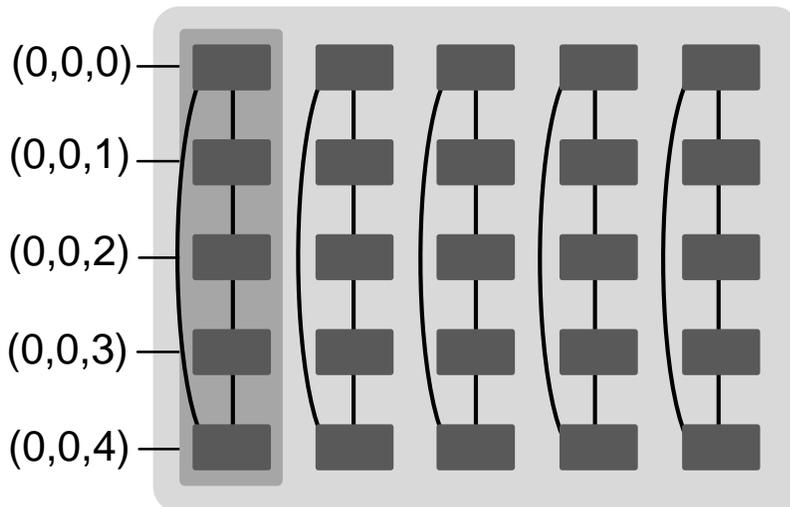
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

$$X' = \{\xi, \xi^3, \dots, \xi^{q-2}\} \text{ (for subgraph 1)}$$

E Example: $q = 5$

Take Routers $(0,0,.)$

$$X = \{1, 4\}$$



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

6 Intra-group connections

Two routers in one group are connected iff their “vertical Manhattan distance” is an element from:

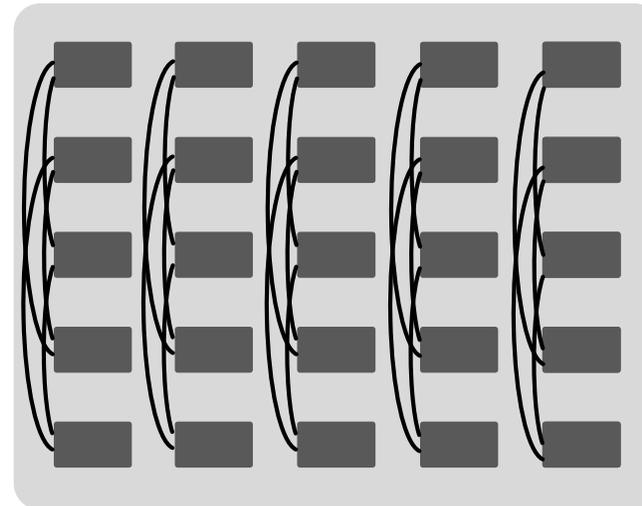
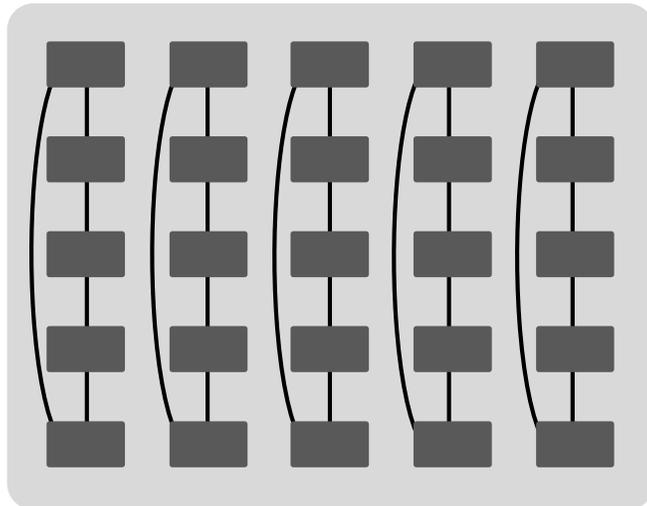
$$X = \{1, \xi^2, \dots, \xi^{q-3}\} \text{ (for subgraph 0)}$$

$$X' = \{\xi, \xi^3, \dots, \xi^{q-2}\} \text{ (for subgraph 1)}$$

E Example: $q = 5$

Take Routers (1,4,.)

$$X' = \{2,3\}$$



DESIGNING AN EFFICIENT NETWORK TOPOLOGY

CONNECTING ROUTERS: DIAMETER 2

7 Inter-group connections

Router $(0, x, y) \leftrightarrow (1, m, c)$

$$\text{iff } y = mx + c$$

E Example: $q = 5$

Take Router $(1,0,0)$

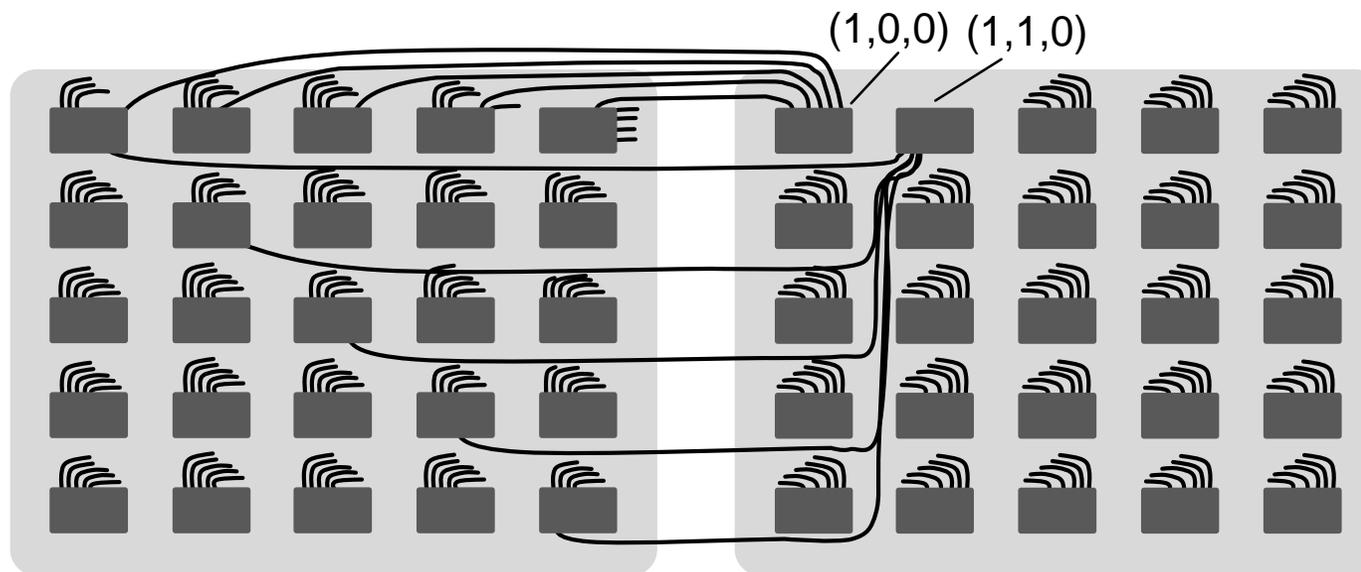
$(1,0,0) \leftrightarrow (0, x, 0)$

$$m = 0, c = 0$$

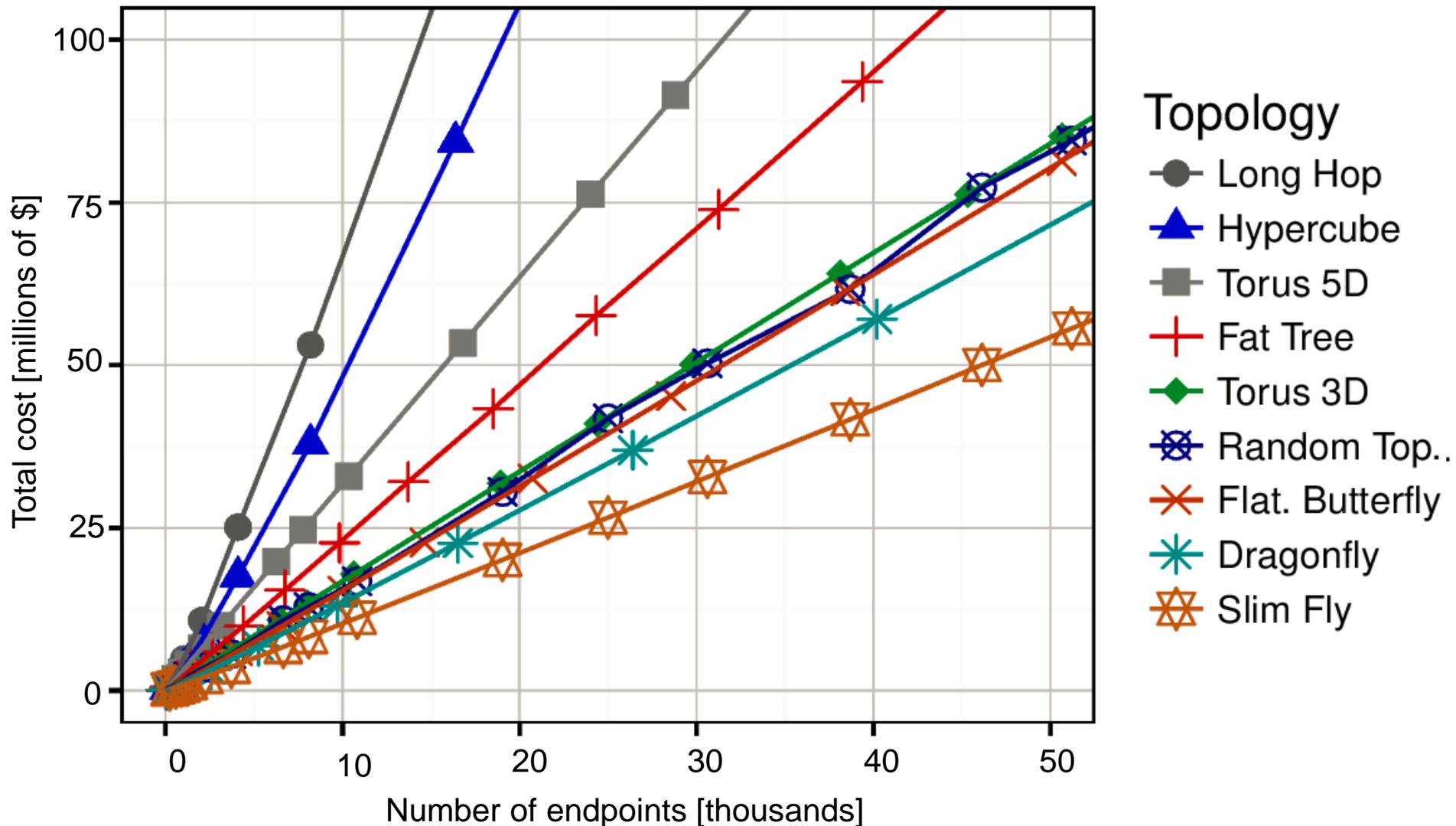
Take Router $(1,1,0)$

$(1,0,0) \leftrightarrow (0, x, x)$

$$m = 1, c = 0$$

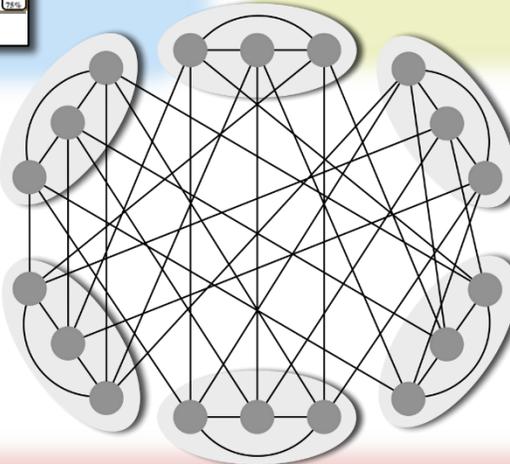
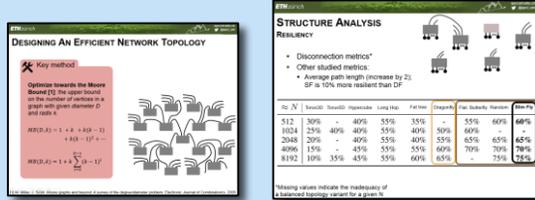


COST COMPARISON



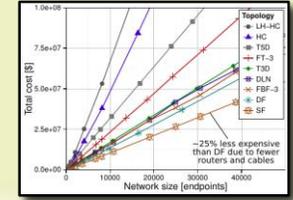
A LOWEST-DIAMETER TOPOLOGY

- Approaching the Moore Bound
- Resilient



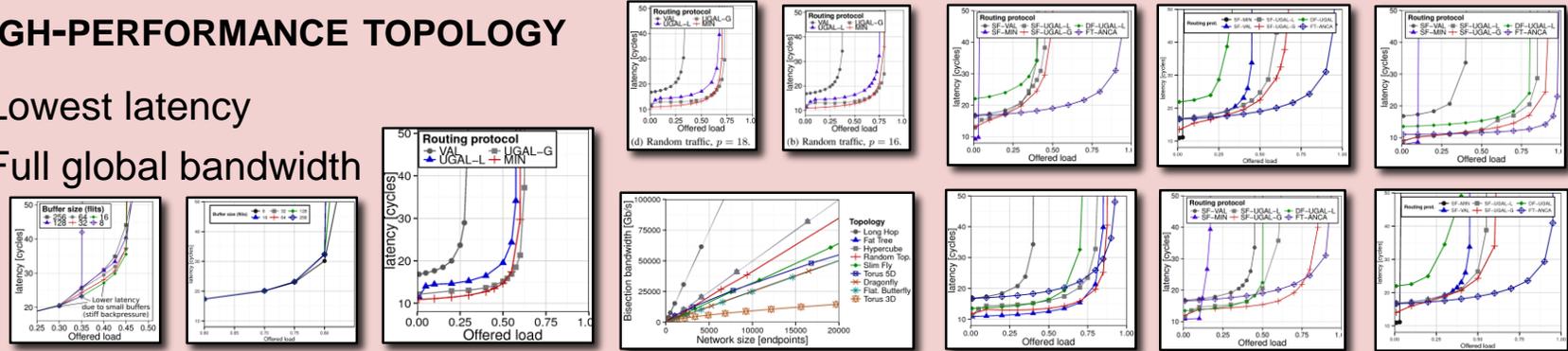
A COST & POWER EFFECTIVE TOPOLOGY

- 25% less expensive than Dragonfly,
- 26% less power-hungry than Dragonfly



A HIGH-PERFORMANCE TOPOLOGY

- Lowest latency
- Full global bandwidth



How to continue from here?

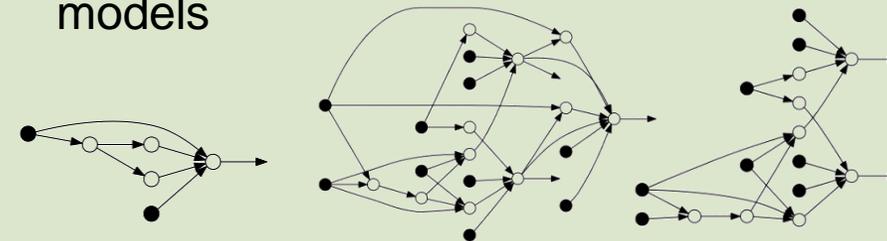
Transformation System

- User-supported, compile- and run-time



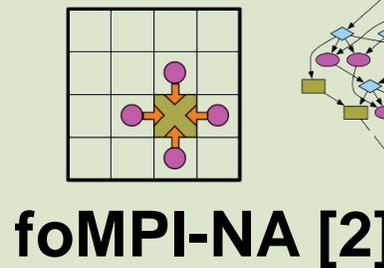
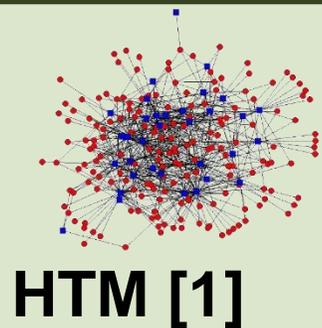
Parallel Language

- Data-centric, explicit requirements models



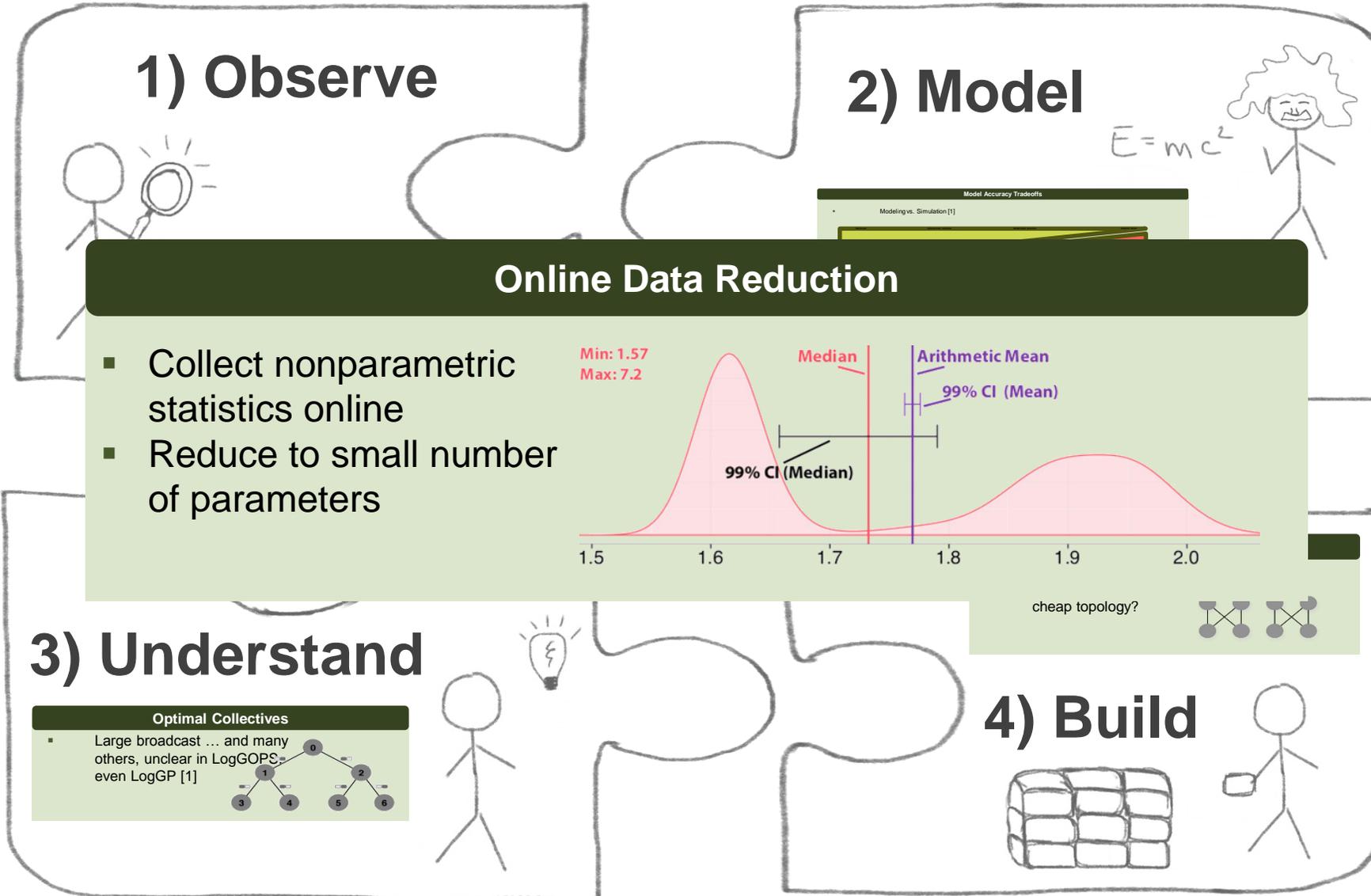
European Research Council
Established by the European Commission
Supporting top researchers
from anywhere in the world

Performance-transparent Platforms



[1]: M. Besta, TH: Accelerating Irregular Computations with Hardware Transactional Memory and Active Messages, ACM HPDC'15
 [2]: R. Belli, TH: Notified Access: Extending Remote Memory Access Programming Models for Producer-Consumer Synchronization, IPDPS'15
 [3]: S. Di Girolamo, P. Jolivet, K. D. Underwood, TH: Exploiting Offload Enabled Network Interfaces, IEEE Micro'16

Advancing with Scientific Performance Engineering

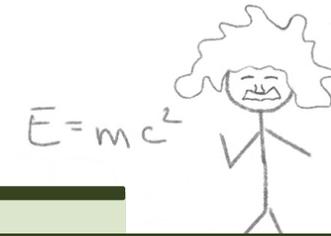


Advancing with Scientific Performance Engineering

1) Observe



2) Model



Model Accuracy Tradeoffs

- Modeling vs. Simulation [1]

Benchmark

Cycle-accurate simulation

Model-based simulation

Analytical Model

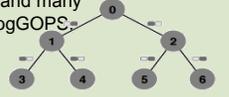
Number of Parameters/Complexity

Model Error

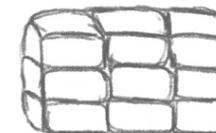
3) Understand

Optimal Collectives

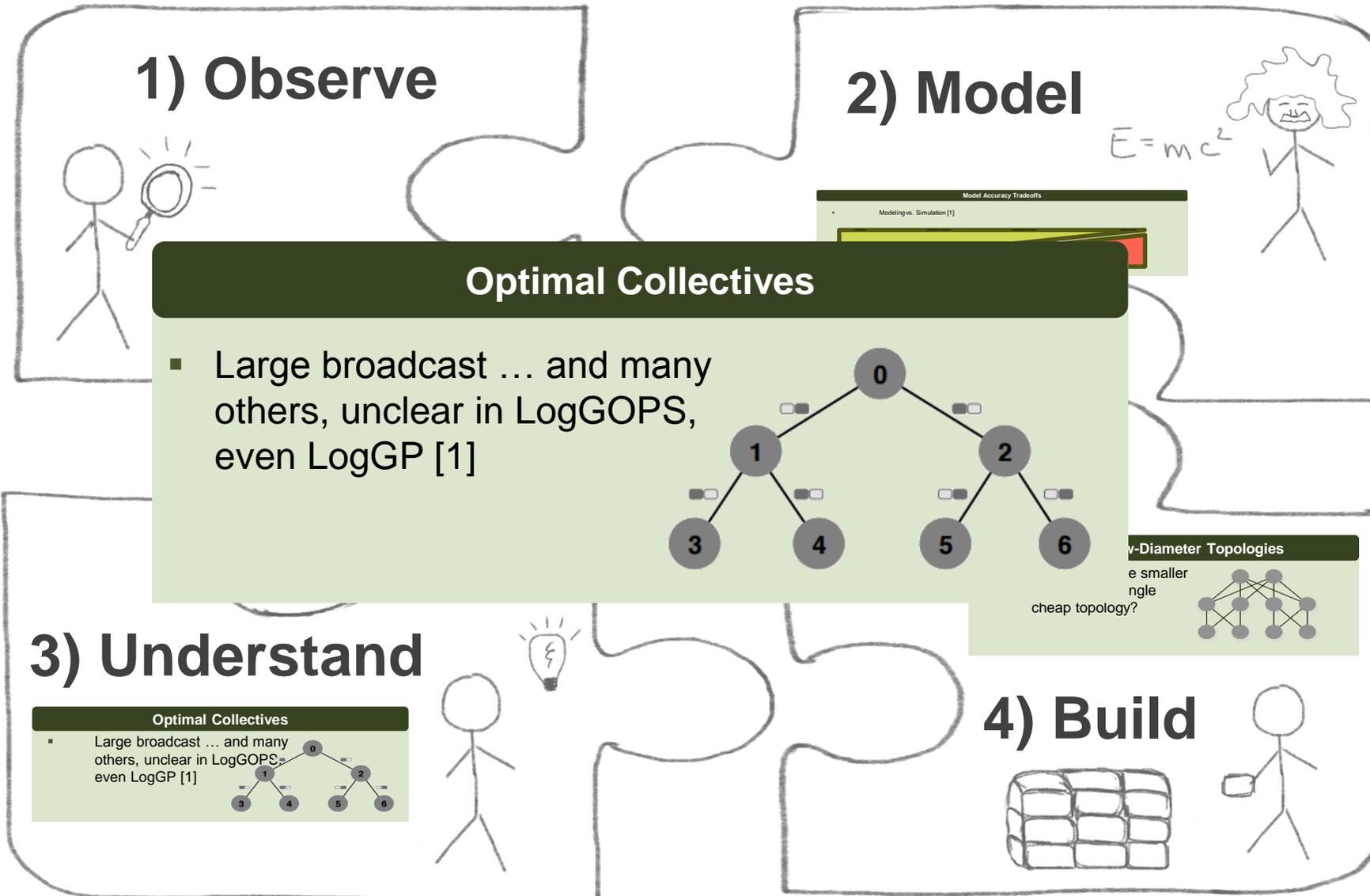
- Large broadcast ... and many others, unclear in LogGOPS, even LogGP [1]



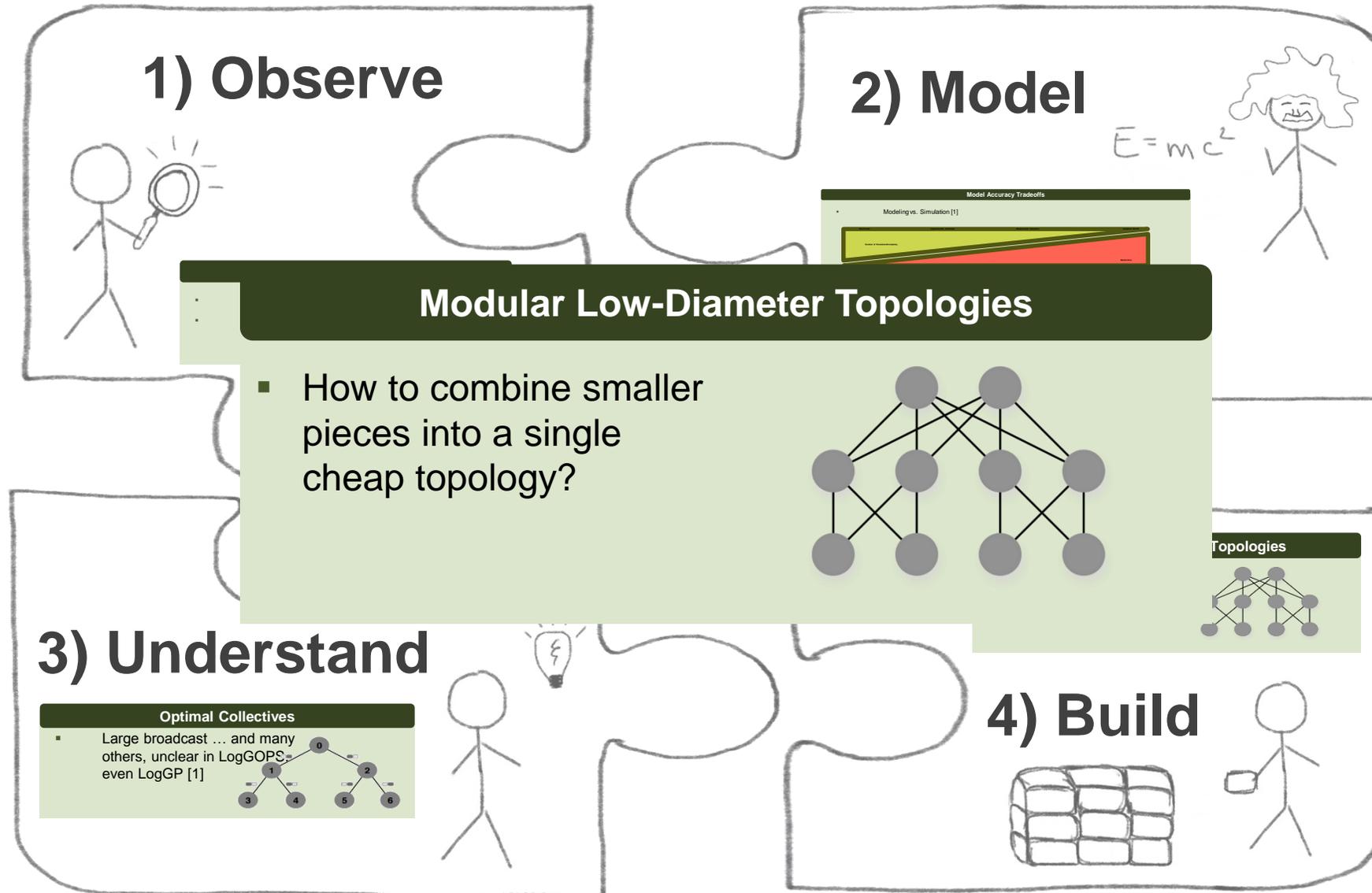
4) Build



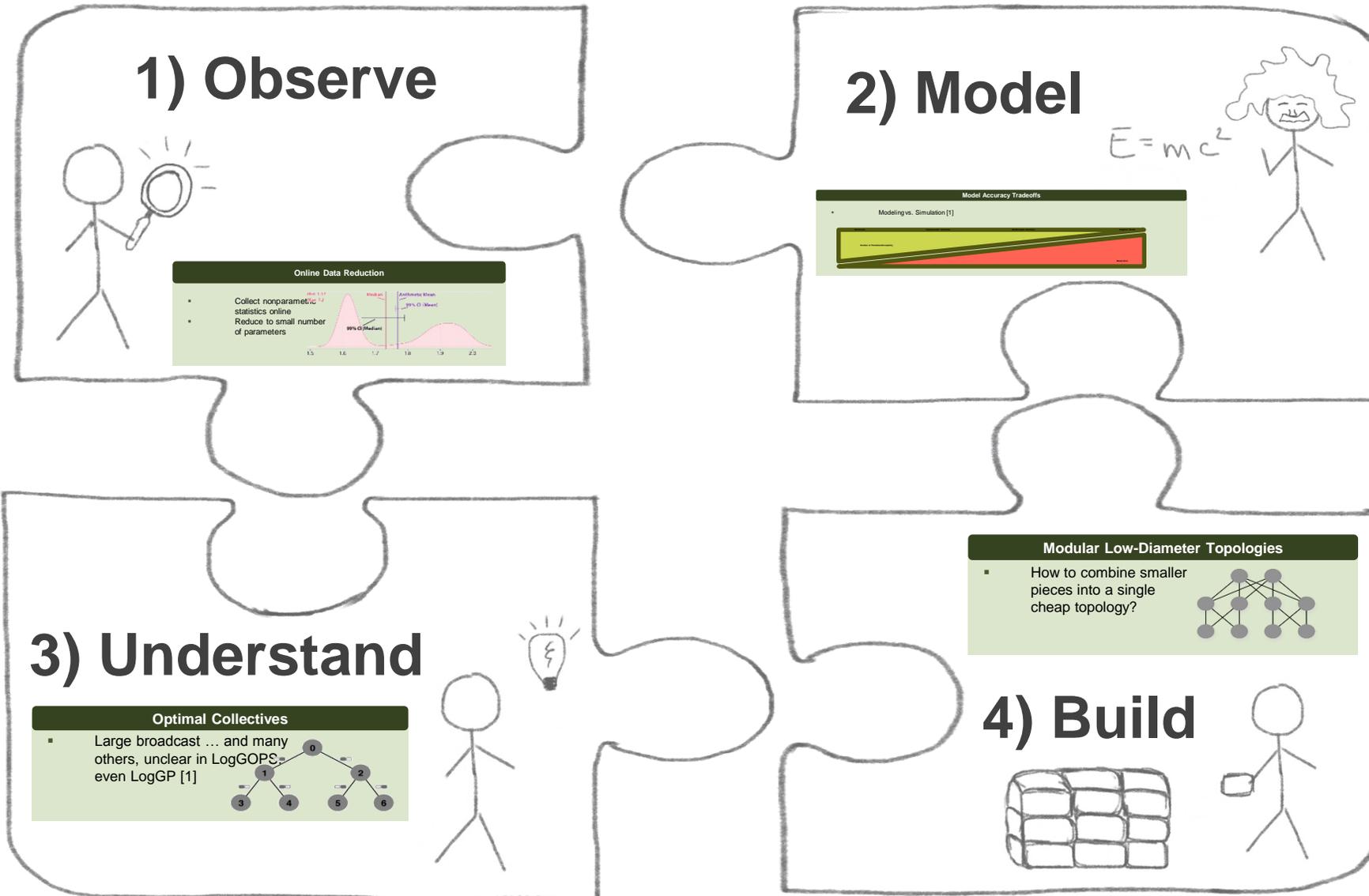
Advancing with Scientific Performance Engineering



Advancing with Scientific Performance Engineering



Advancing with Scientific Performance Engineering



Backup