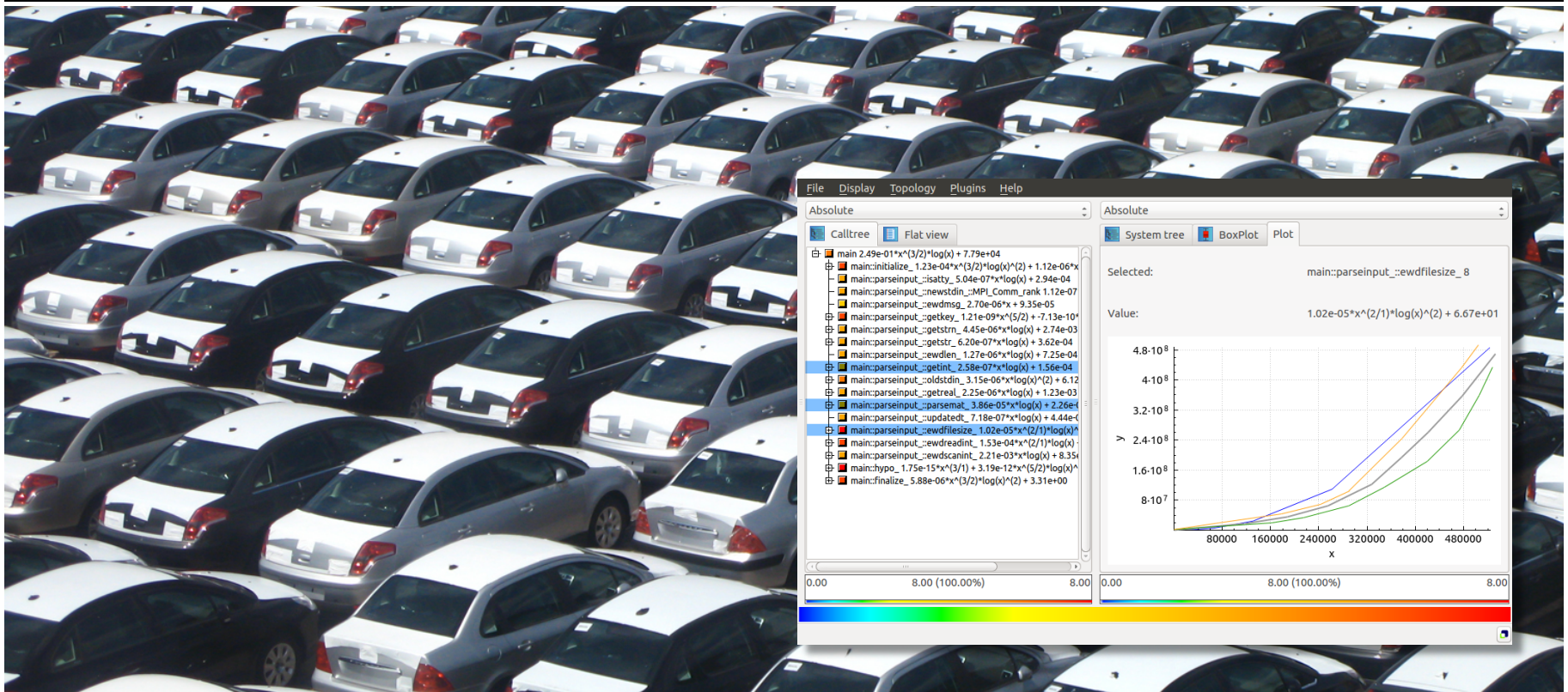


Is your software ready for exascale? – How the next generation of performance tools can give you the answer



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Felix Wolf, TU Darmstadt



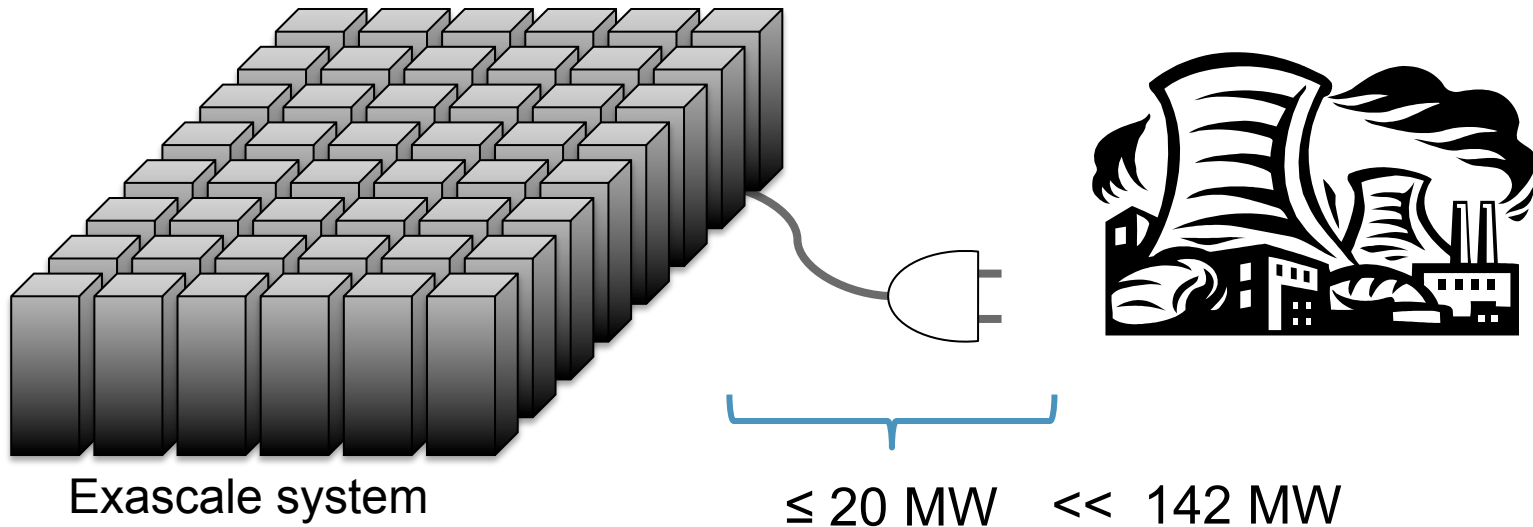
Acknowledgement

- Alexandru Calotoiu (TU Darmstadt)
- Torsten Höfler (ETH Zurich)
- Sergei Shudler (TU Darmstadt)
- Alexandre Strube (Jülich Supercomputing Centre)
- Andreas Vogel (GU Frankfurt)

- Marius Poke (RWTH Aachen)
- Paul Wiedeking (RWTH Aachen)



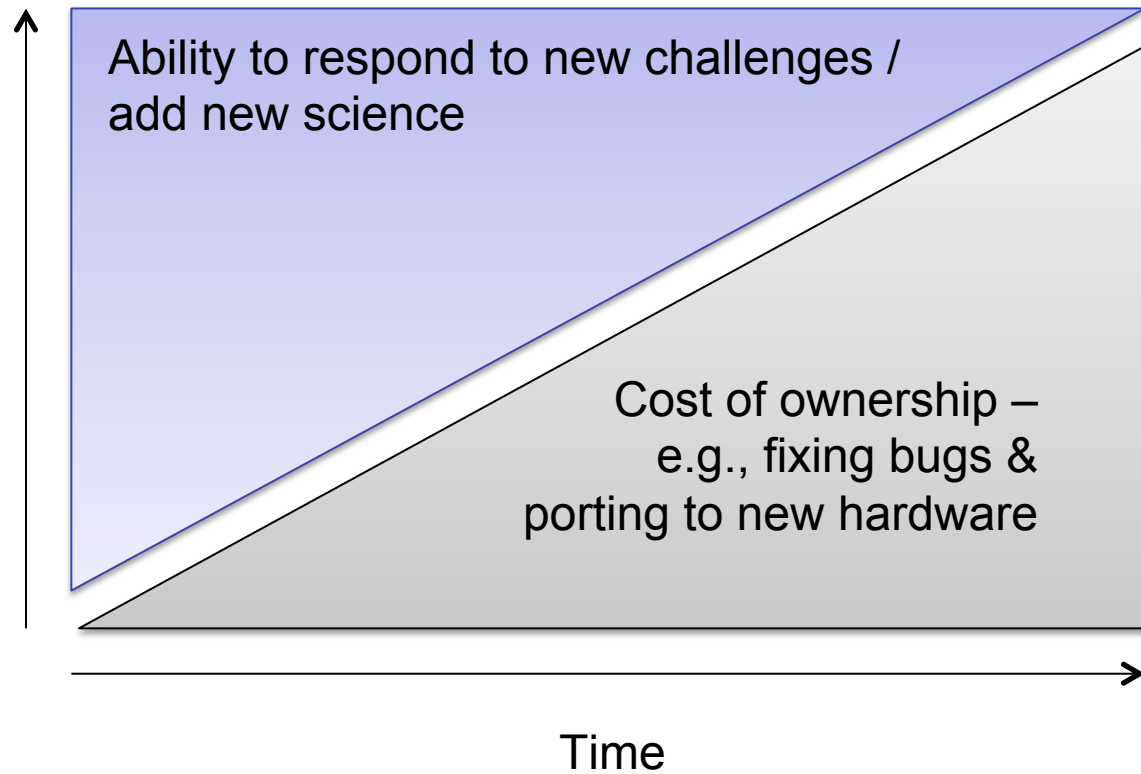
Power envelope of hardware



Green 500 # 1	
GFLOPS/W	7.03
Site	RIKEN
Computer	Shoubu - ExaScaler
Total Power (kW)	50.32

Manpower envelope of software

Available
manpower



Electrical power vs. manpower

Energy costs	
Power	20 MW
Price per kWh	0.1 €
Hours per year	5,000 h
Energy costs per year	10 M€



200 FTEs
(50k€ per FTE)

To amortize the investment, one FTE needs to tune the workload by **0.5%**

Electrical power vs. manpower

Energy costs	
Power (DARPA estimate)	60 MW
Price per kWh	0.1 €
Hours per year	5,000 h
Energy costs per year	30 M€



600 FTEs

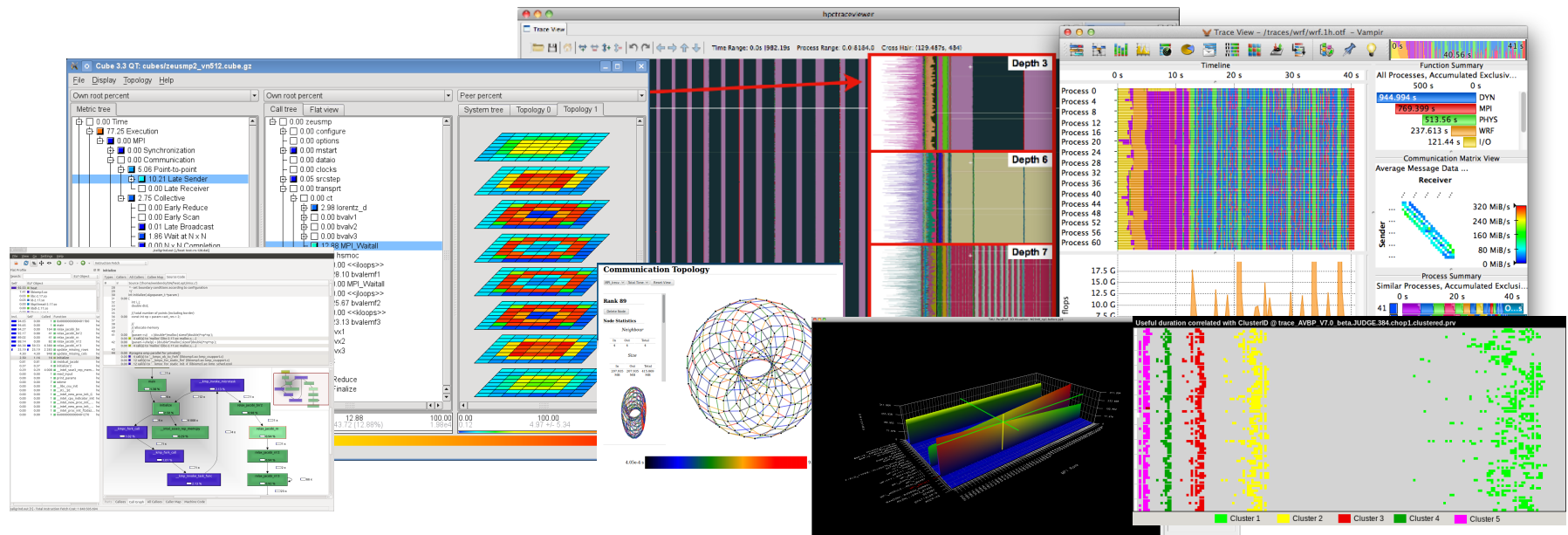
To amortize the investment, one FTE needs to tune the workload by **0.15%**

- Potential in trading hardware for brainware*
- Productivity of staff can be further increased through performance tools
- Early resolution of performance issues maximizes benefit

* C. Bischof et al.: Brainware for green HPC, Computer Science-Research and Development, Springer

Traditional performance tools

Provide insight into measured performance behavior

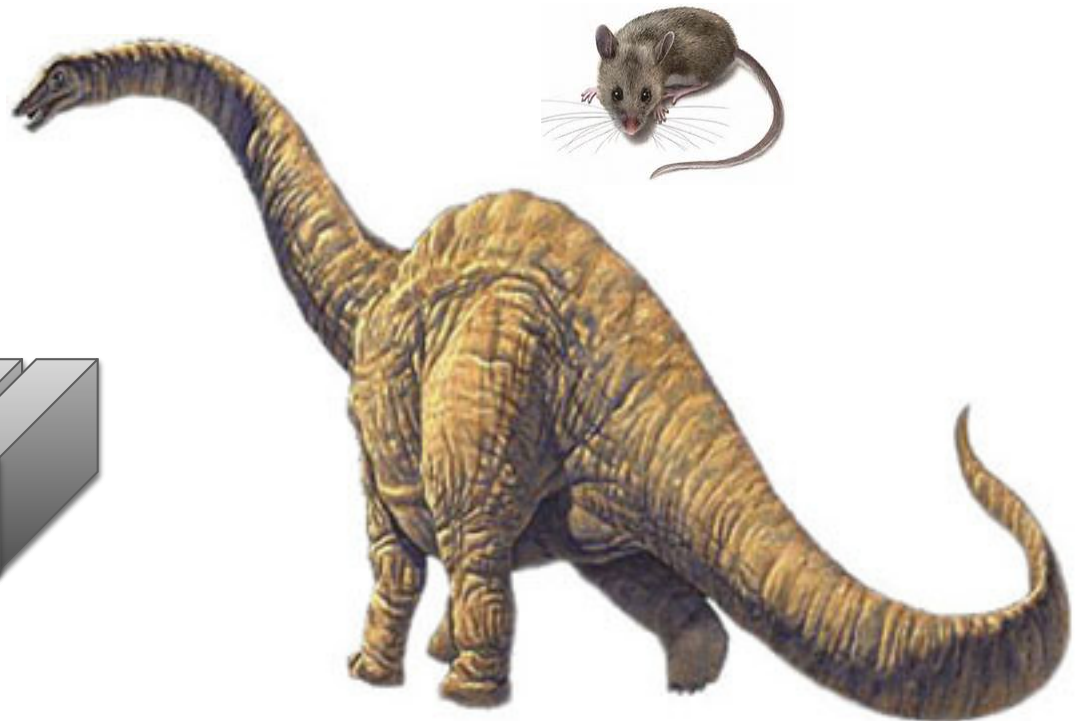
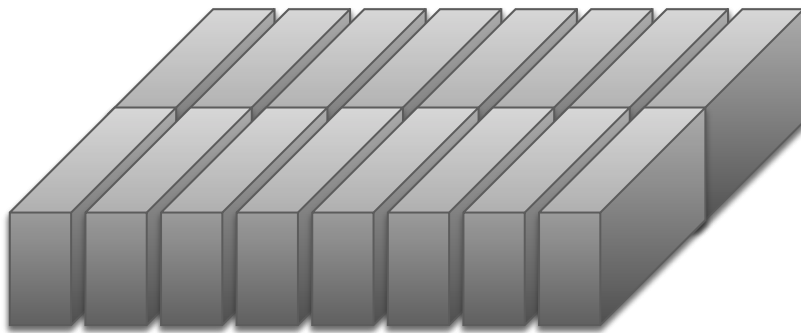
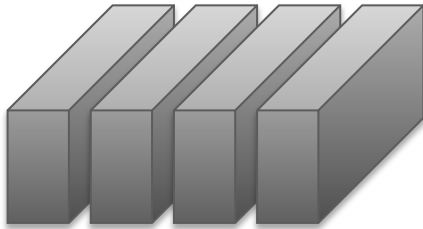


Scale of insight = scale of experiment

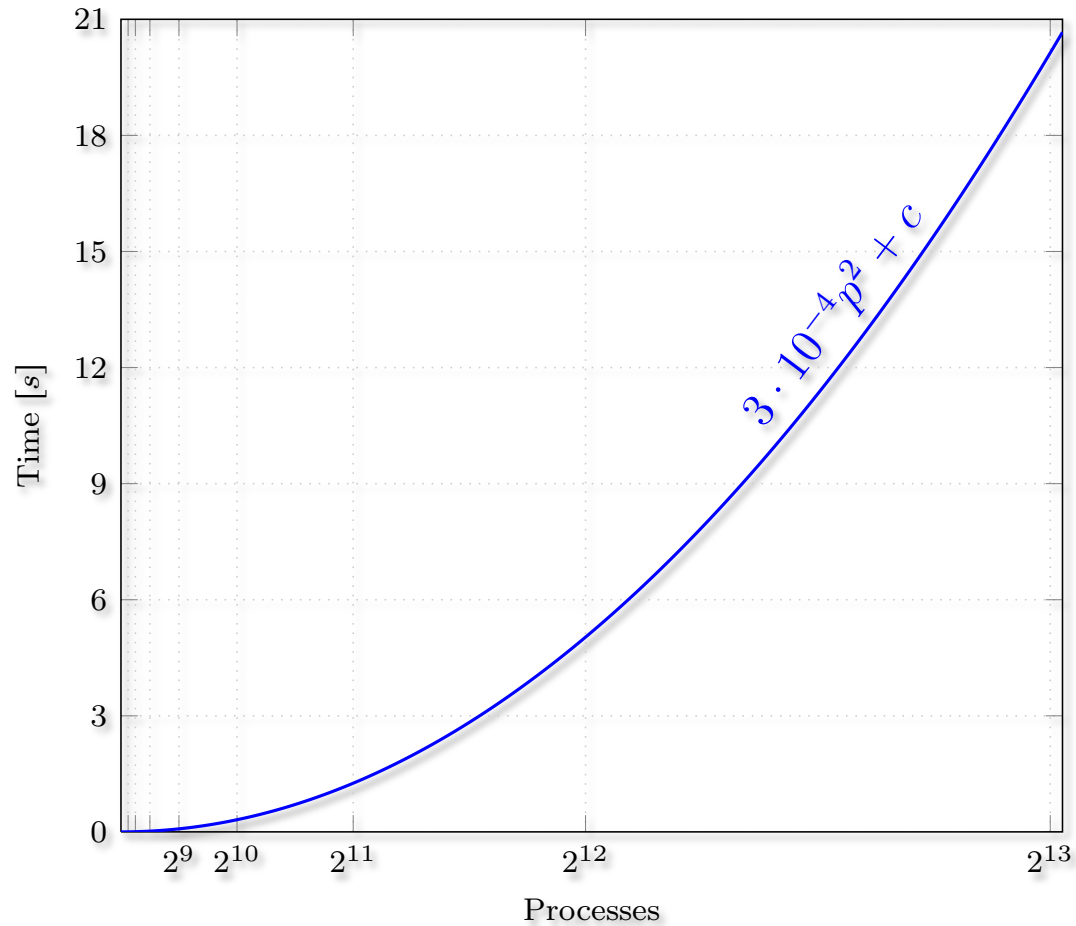
Latent scalability bugs

System size

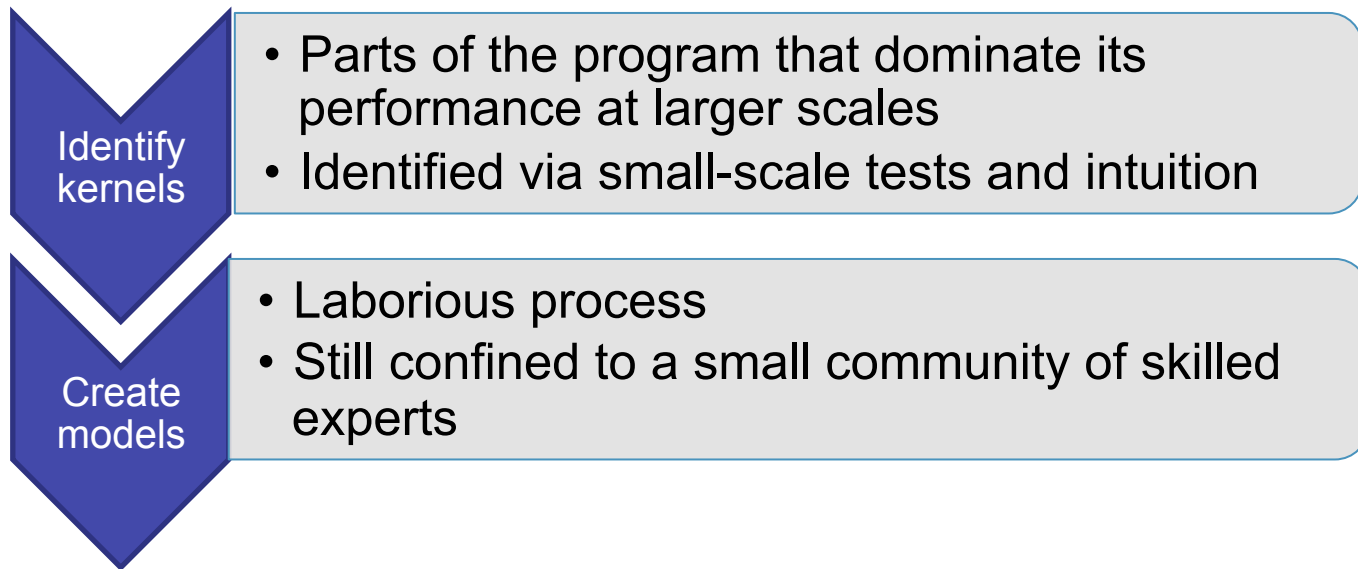
Execution time



Scalability model



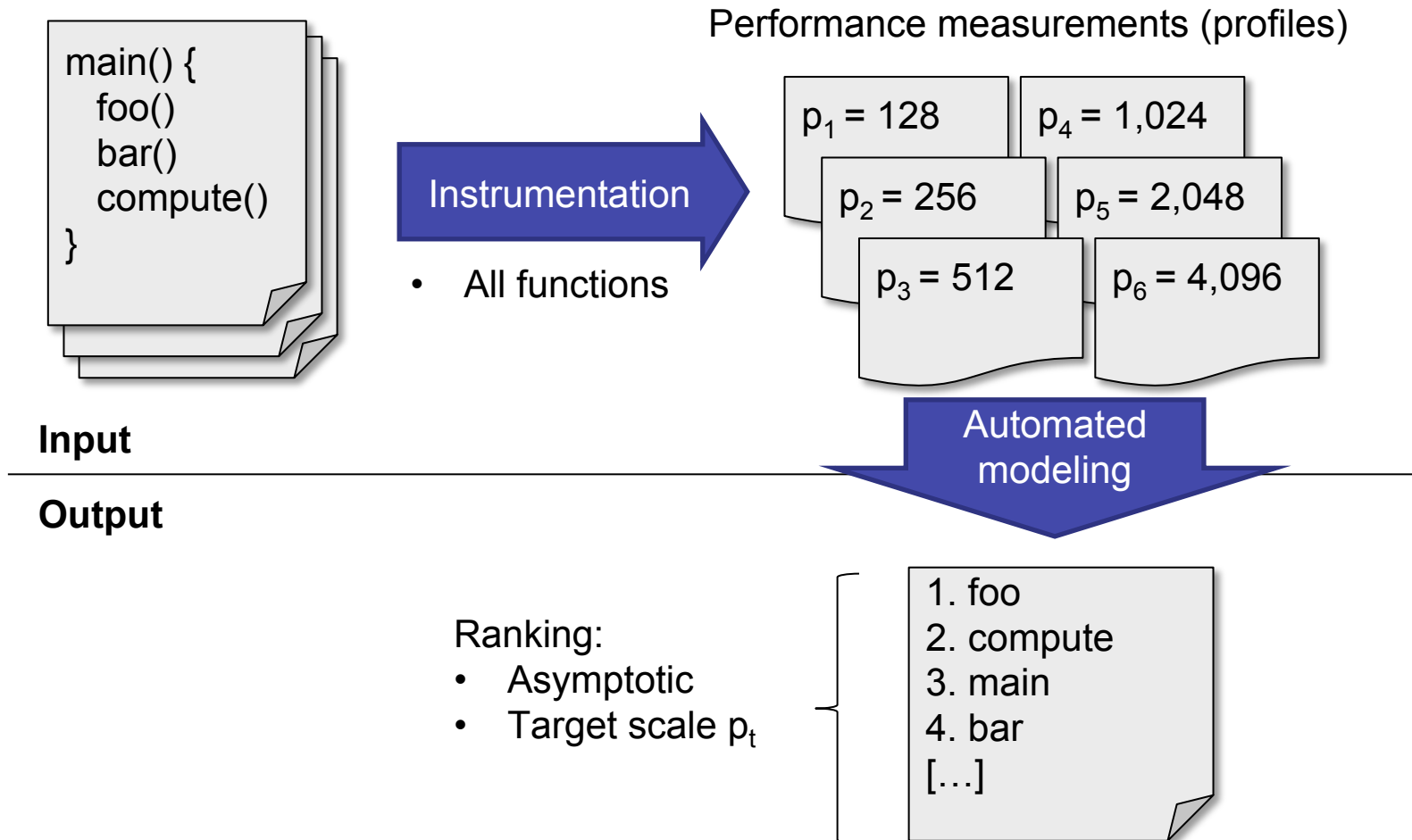
Analytical scalability modeling



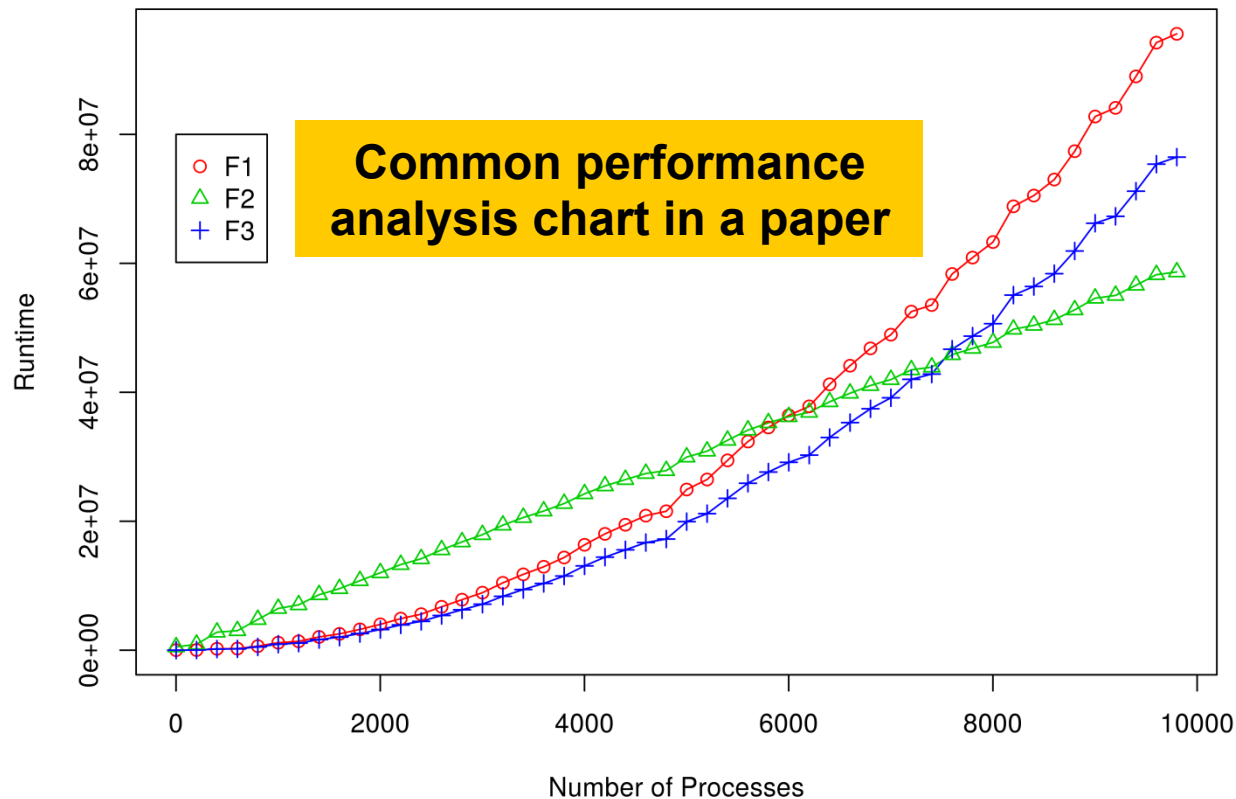
Disadvantages

- Time consuming
- Danger of overlooking unscalable code

Automated empirical modeling (2)



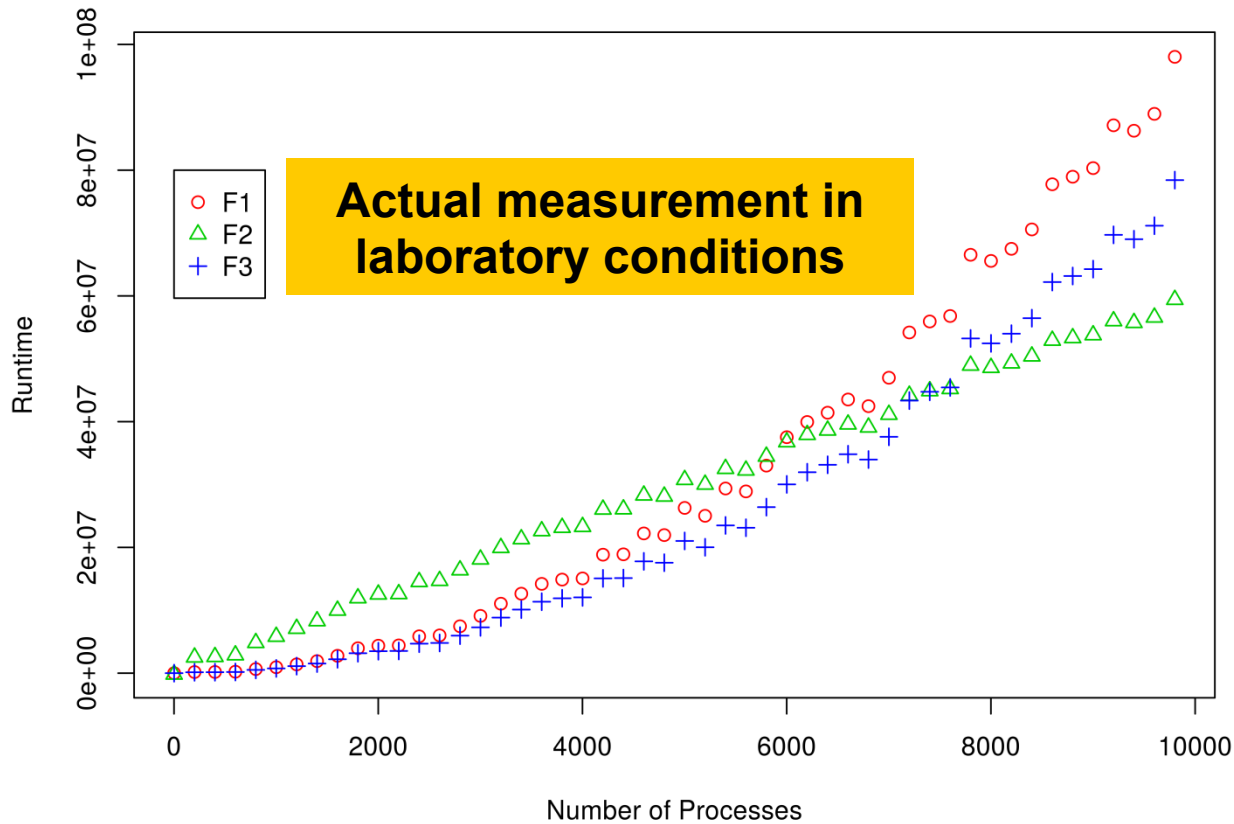
Primary focus on scaling trend



Our ranking

1. F_1
2. F_3
3. F_2

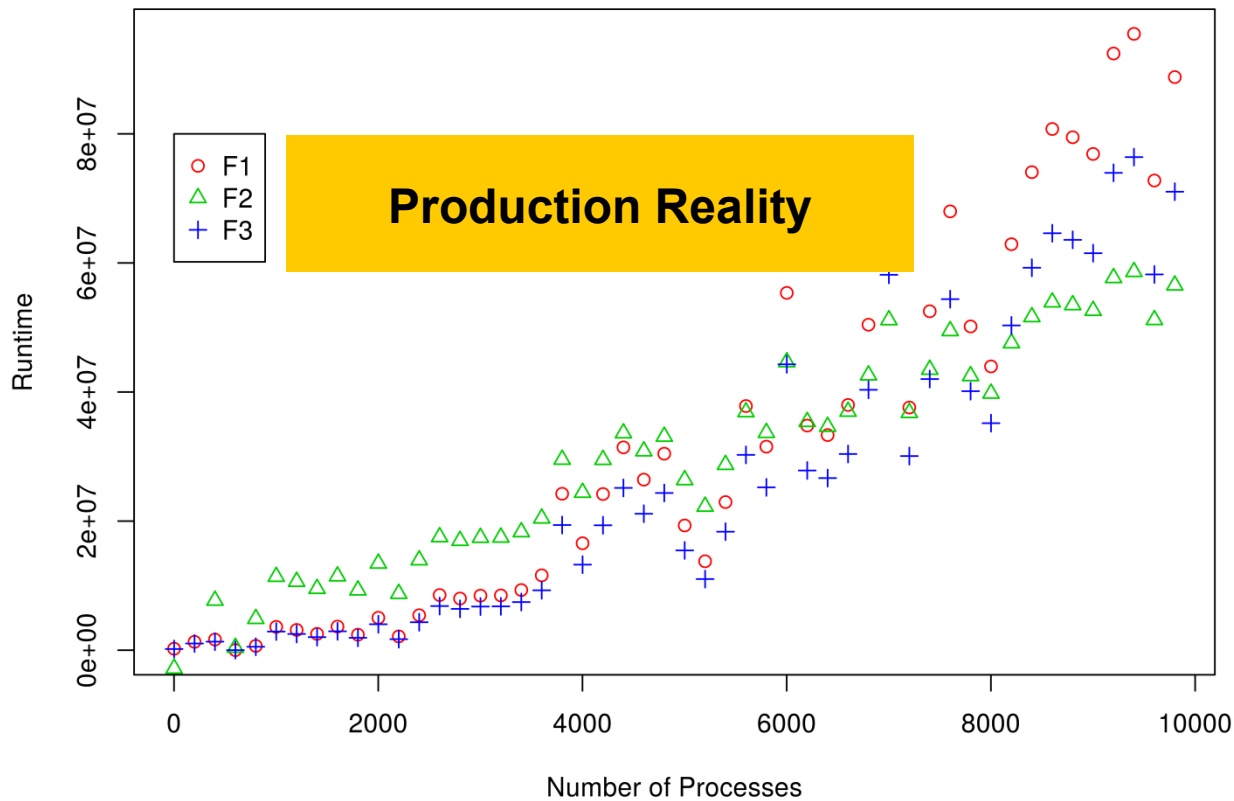
Primary focus on scaling trend



Our ranking

1. F_1
2. F_3
3. F_2

Primary focus on scaling trend



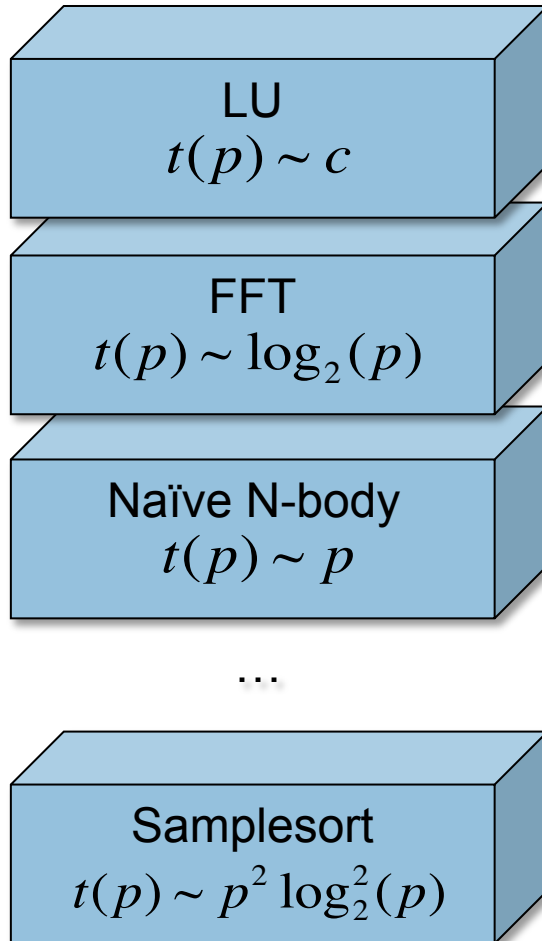
Our ranking

1. F_1
2. F_3
3. F_2

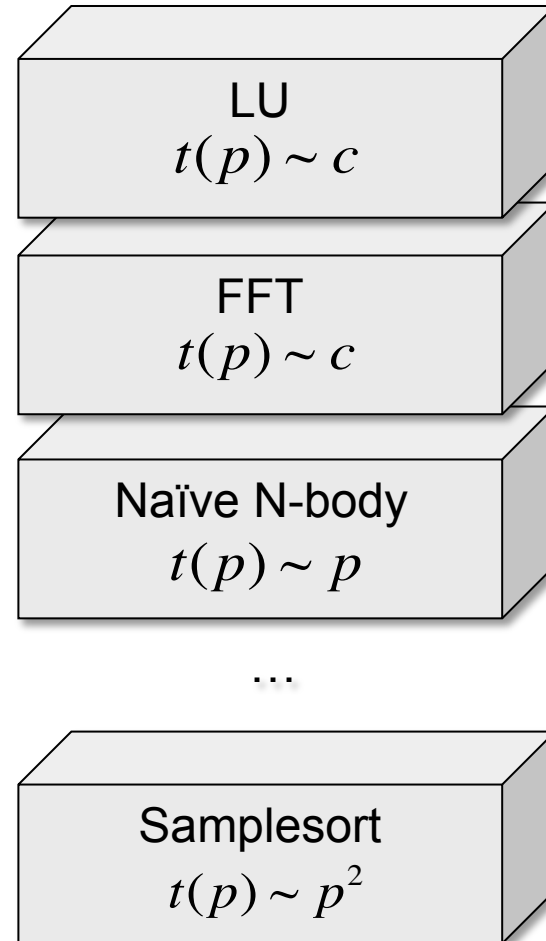
Model building blocks



Computation



Communication



Performance model normal form



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

$$n \in \mathbb{N}$$

$$i_k \in I$$

$$j_k \in J$$

$$I, J \subset \mathbb{Q}$$

$$n = 1$$

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

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

$$c_1$$

$$c_1 \cdot \log(p)$$

$$c_1 \cdot p$$

$$c_1 \cdot p \cdot \log(p)$$

$$c_1 \cdot p^2$$

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

Performance model normal form



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

$n \in \mathbb{N}$

$n = 2$

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

$J = \{0, 1\}$

$$c_1 + c_2 \cdot p$$

$$c_1 + c_2 \cdot p^2$$

$$c_1 + c_2 \cdot \log(p)$$

$$c_1 + c_2 \cdot p \cdot \log(p)$$

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

$$c_1 \cdot \log(p) + c_2 \cdot p$$

$$c_1 \cdot \log(p) + c_2 \cdot p \cdot \log(p)$$

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

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

$$c_1 \cdot p + c_2 \cdot p \cdot \log(p)$$

$$c_1 \cdot p + c_2 \cdot p^2$$

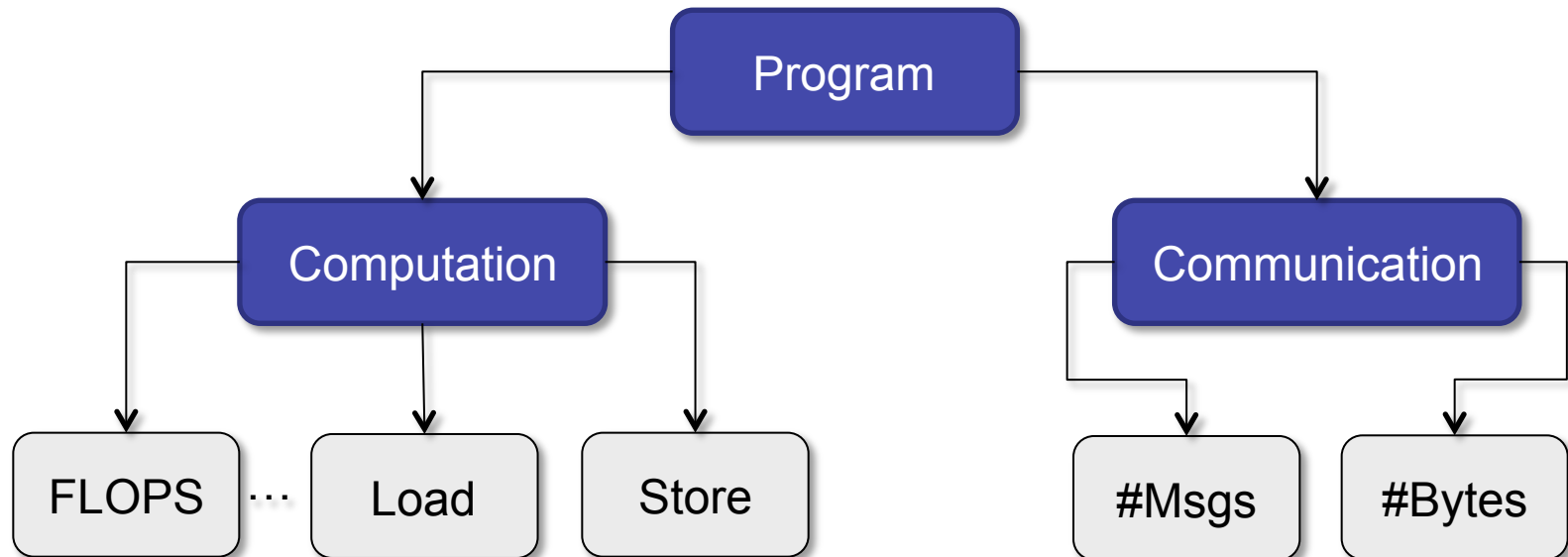
$$c_1 \cdot p + c_2 \cdot p^2 \cdot \log(p)$$

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

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

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

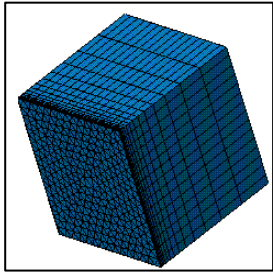
Modeling operations vs. time



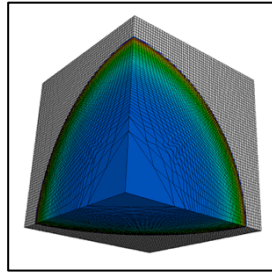
Disagreement may be indicative of wait states

Time

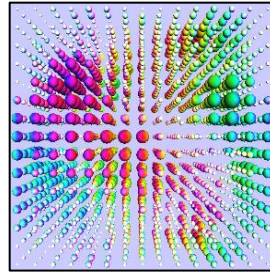
Case studies



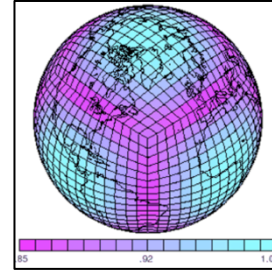
Sweep3d



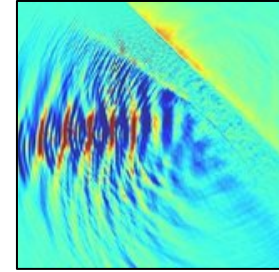
Lulesh



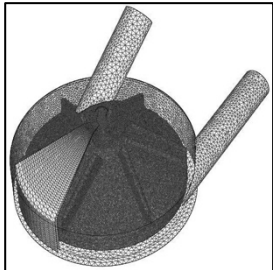
Milc



HOMME



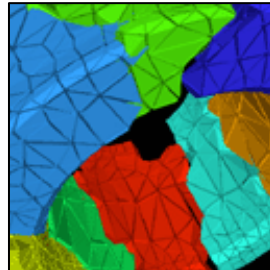
JUSPIC



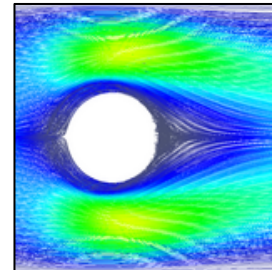
XNS



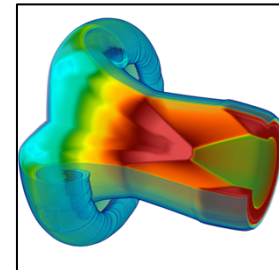
NEST



UG4



MP2C



BLAST

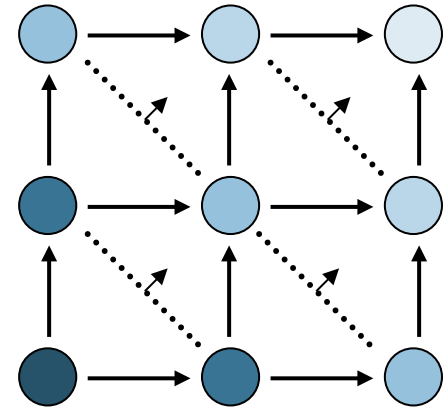


Sweep3D - Neutron transport simulation

- LogGP model for communication developed by Hoisie et al.

$$t^{comm} = [2(p_x + p_y - 2) + 4(n_{sweep} - 1)] \cdot t_{msg}$$

$$t^{comm} = c \cdot \sqrt{p}$$



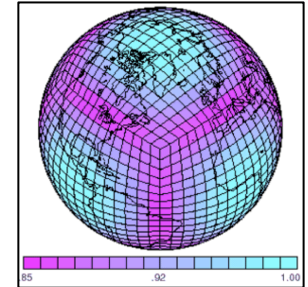
Kernel [2 of 40]	Model [s] $t = f(p)$	Predictive error [%] $p_t=262k$
sweep \rightarrow MPI_Recv	$4.03\sqrt{p}$	5.10
sweep	582.19	01

#bytes \approx const.
#msg \approx const.

$$p_i \leq 8k$$

Core of the Community Atmospheric Model (CAM)

- Spectral element dynamical core on a cubed sphere grid

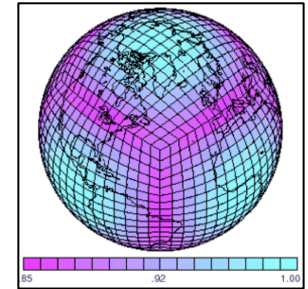


Kernel [3 of 194]	Model [s] $t = f(p)$	Predictive error [%] $p_t = 130k$
box_rearrange → MPI_Reduce	$0.026 + 2.53 \cdot 10^{-6} p \cdot \sqrt{p} + 1.24 \cdot 10^{-12} p^3$	57.02
vlaplace_sphere_vk	49.53	99.32
compute_and_apply_rhs	48.68	1.65

$$p_i \leq 15k$$

Core of the Community Atmospheric Model (CAM)

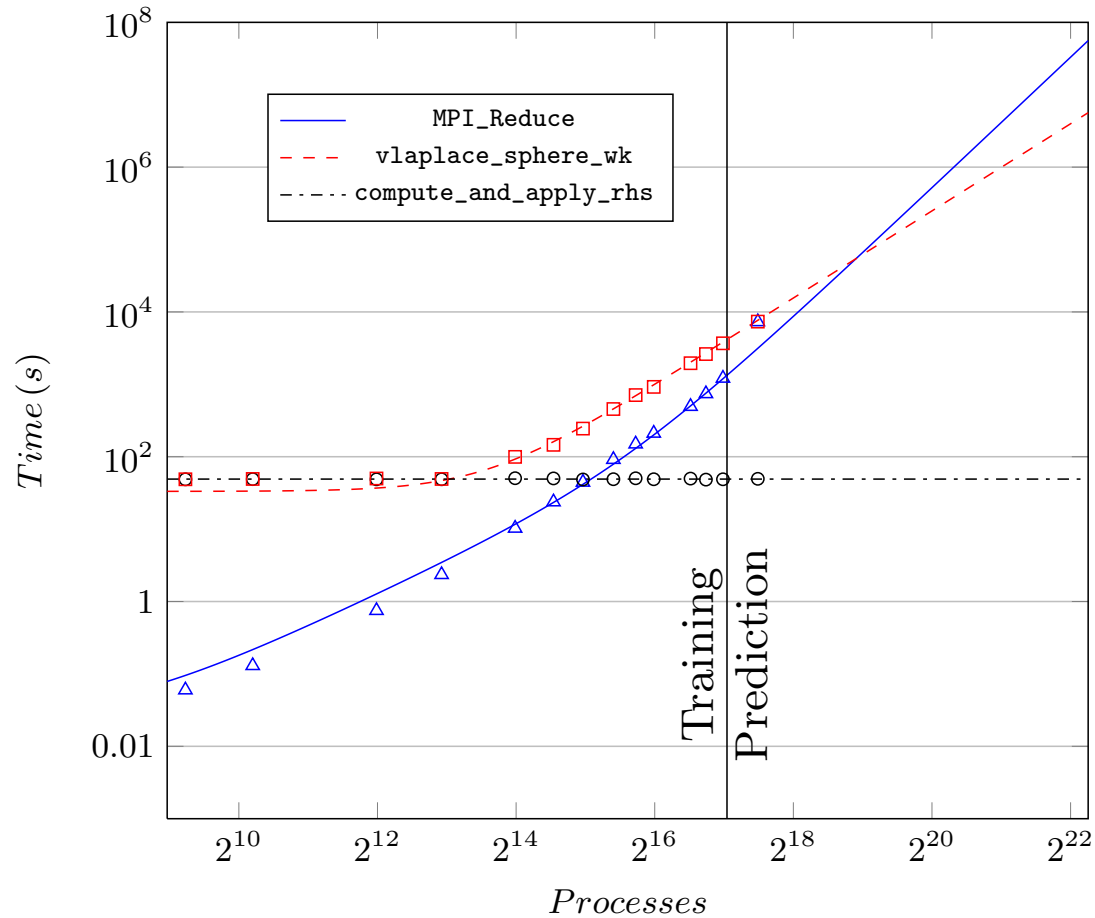
- Spectral element dynamical core on a cubed sphere grid



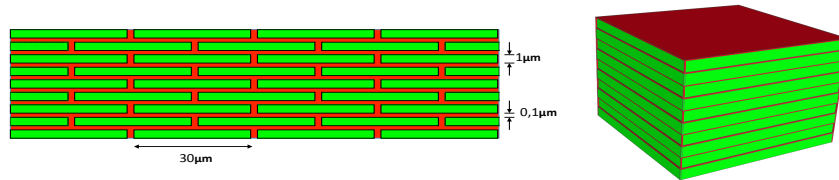
Kernel [3 of 194]	Model [s] $t = f(p)$	Predictive error [%] $p_t = 130k$
box_rearrange → MPI_Reduce	$3.63 \cdot 10^{-6} p \cdot \sqrt{p} + 7.21 \cdot 10^{-13} p^3$	30.34
vlaplace_sphere_vk	$24.44 + 2.26 \cdot 10^{-7} p^2$	4.28
compute_and_apply_rhs	49.09	0.83

$$p_i \leq 43k$$

HOMME – Climate (2)



UG4

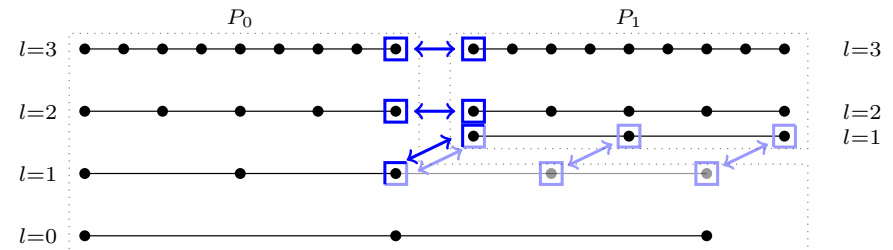


- Numerical framework for grid-based solution of partial differential equations (~500,000 lines of C++ code, 2,000 kernels)
 - Application: drug diffusion through the human skin
- In general, all kernels scale well
 - Multigrid solver kernel (MGM) scales logarithmically
 - Number of iterations needed by the unpreconditioned conjugate gradient (CG) method depends on the mesh size
 - Increases by factor of two with each refinement
 - Will therefore suffer from iteration count increase in weak scaling

Kernel	Model (time [s])
CG	$0.227 + 0.31 * p^{0.5}$
MGM	$0.219 + 0.0006 * \log^2(p)$

Issue with MPI communicator group creation

- Create MPI communicator groups for each level of multigrid hierarchy
- Exclude processes that do not own a grid part on that level
- *Before*: Membership info communicated using MPI_Allreduce with array of length p - non-scalable $p * O(\text{MPI_Allreduce})$ complexity
- *Now*: MPI_Allreduce replaced by MPI_Comm_split - enhanced algorithms of which are known to have $O(\log^2 p)$ complexity



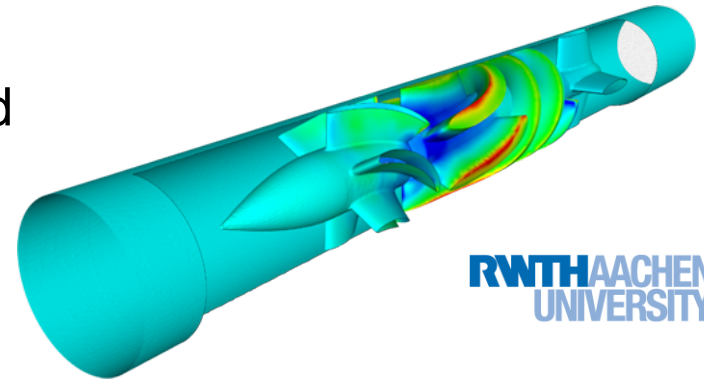
(C. Siebert, F. Wolf: Parallel sorting with minimal data. Recent Advances in the Message Passing Interface, 2011)

Weak vs. strong scaling

- Wall-clock time not necessarily monotonically increasing under strong scaling
 - Harder to capture model automatically
 - Different invariants require different reductions across processes

	Weak scaling	Strong scaling
Invariant	Problem size per process	Overall problem size
Model target	Wall-clock time	Accumulated time
Reduction	Maximum / average	Sum

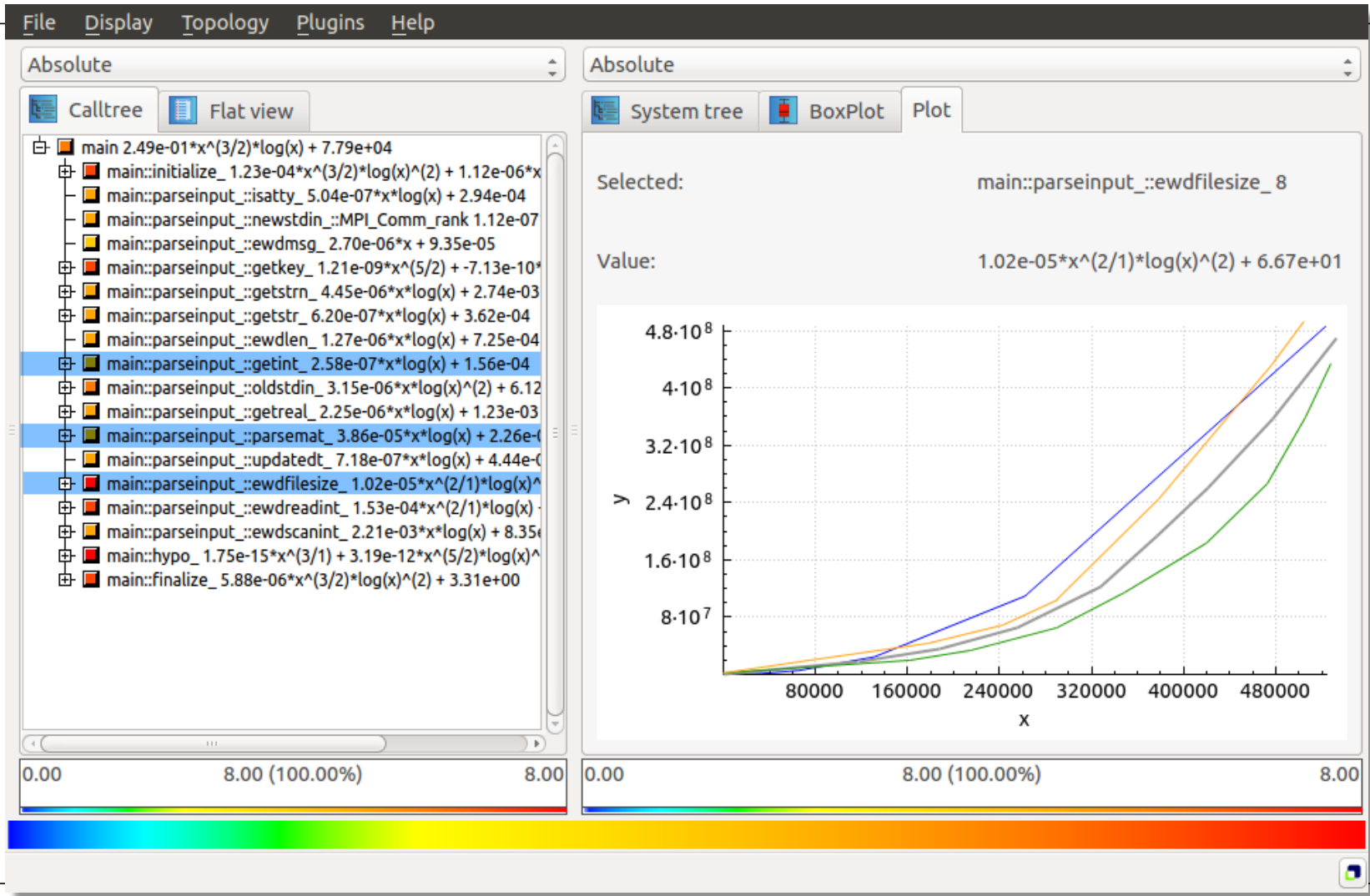
- Finite element flow simulation
- Strong scaling analysis using accumulated time across processes as metric



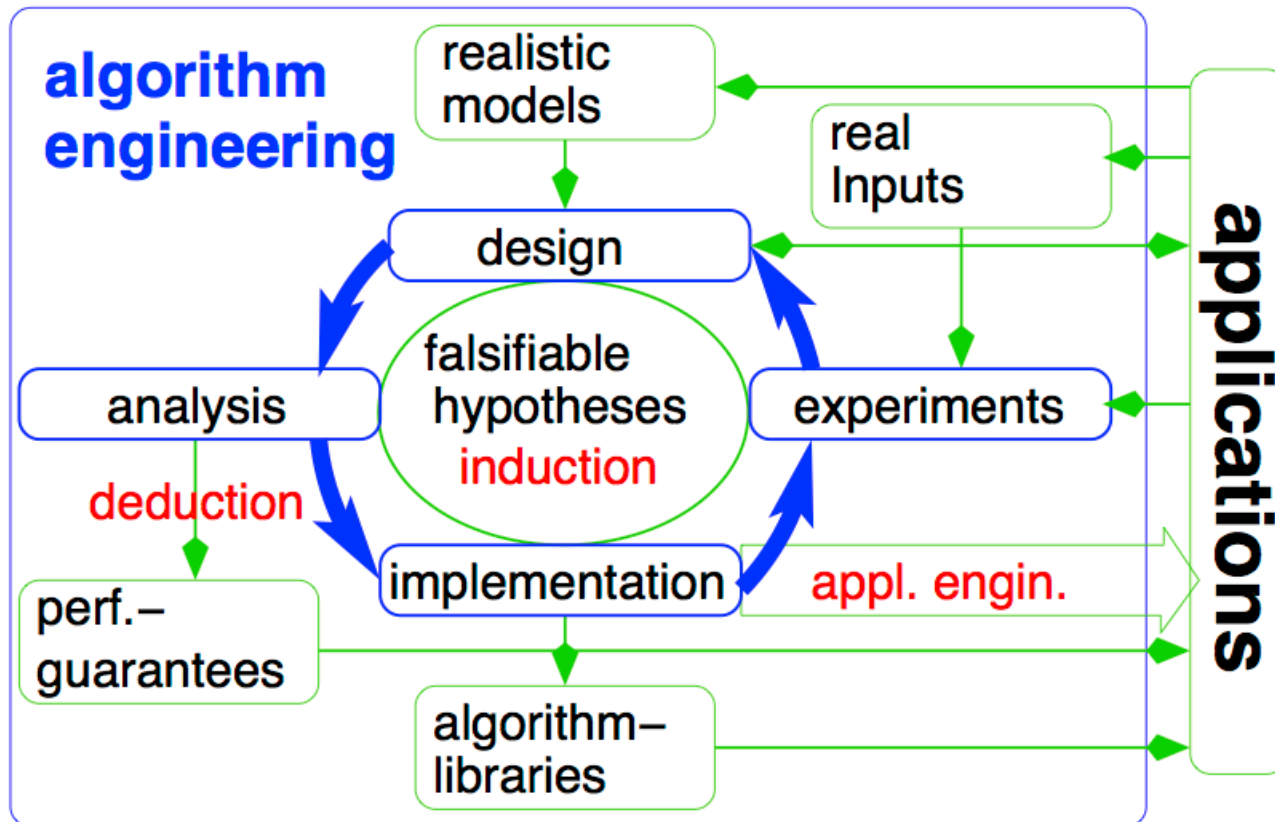
RWTH AACHEN
UNIVERSITY

Kernel	Runtime[%] p=128	Runtime[%] p=4096	Model [s] t = f(p)
ewdgennprm->MPI_Recv	0.46	51.46	$0.029 \cdot p^2$
ewddot	44.78	5.04	$p \cdot \log(p)$

#bytes = ~p
#msg = ~p

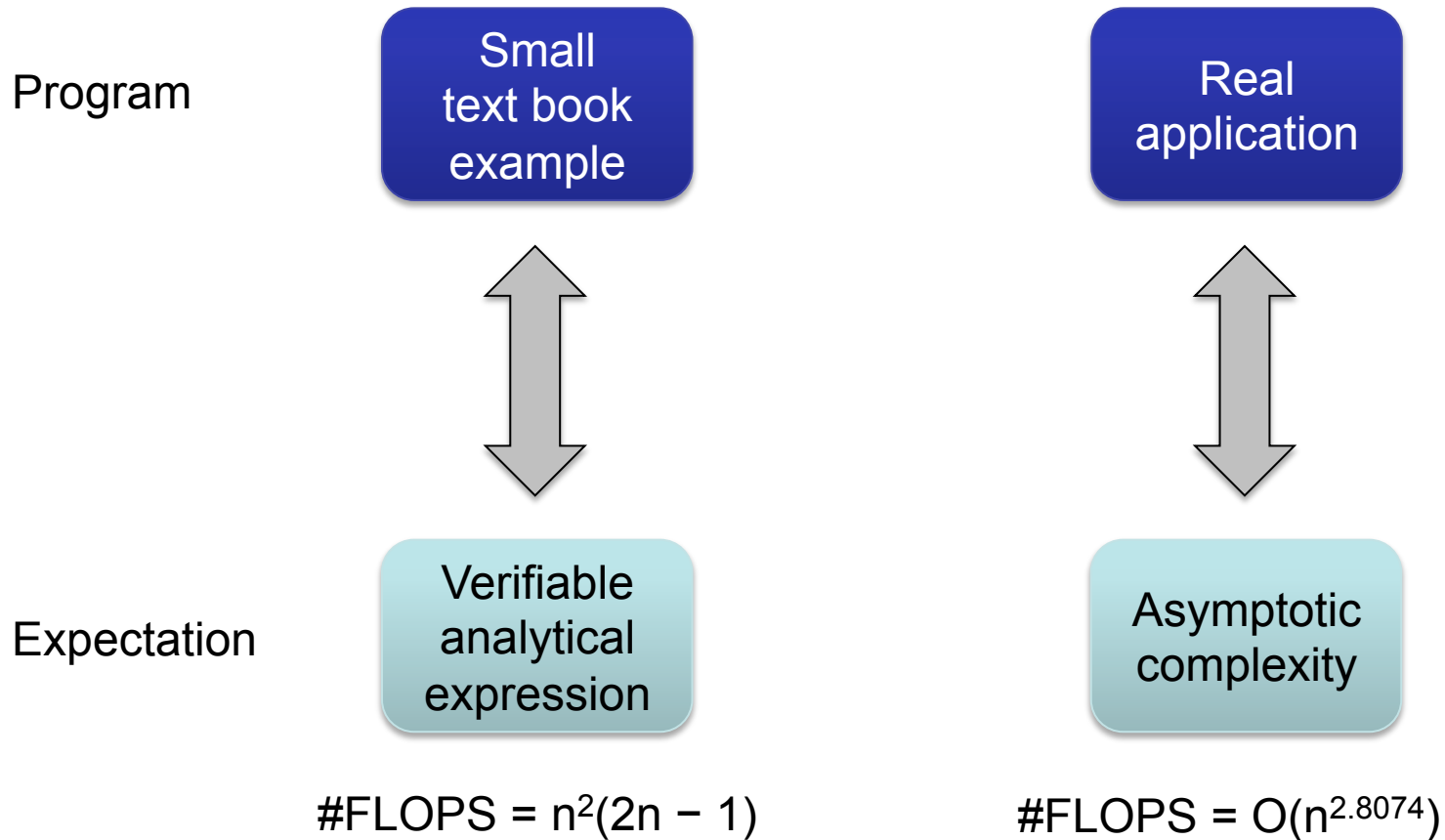


Algorithm engineering



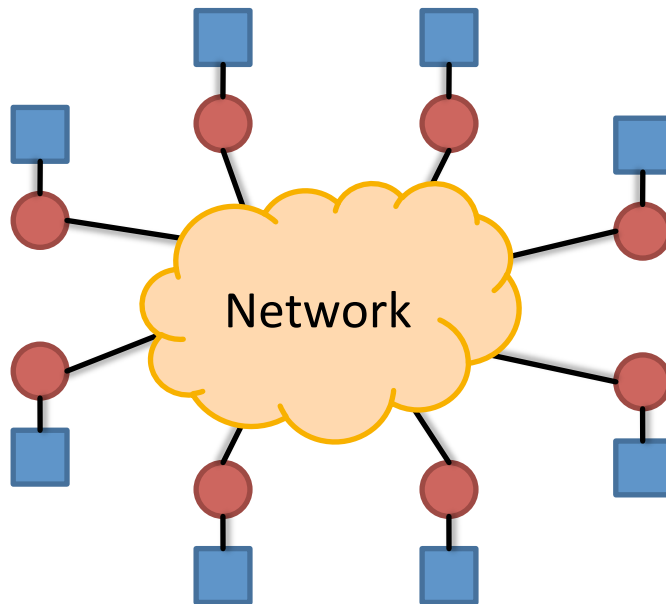
Courtesy of Peter Sanders, KIT

How to validate scalability in practice?



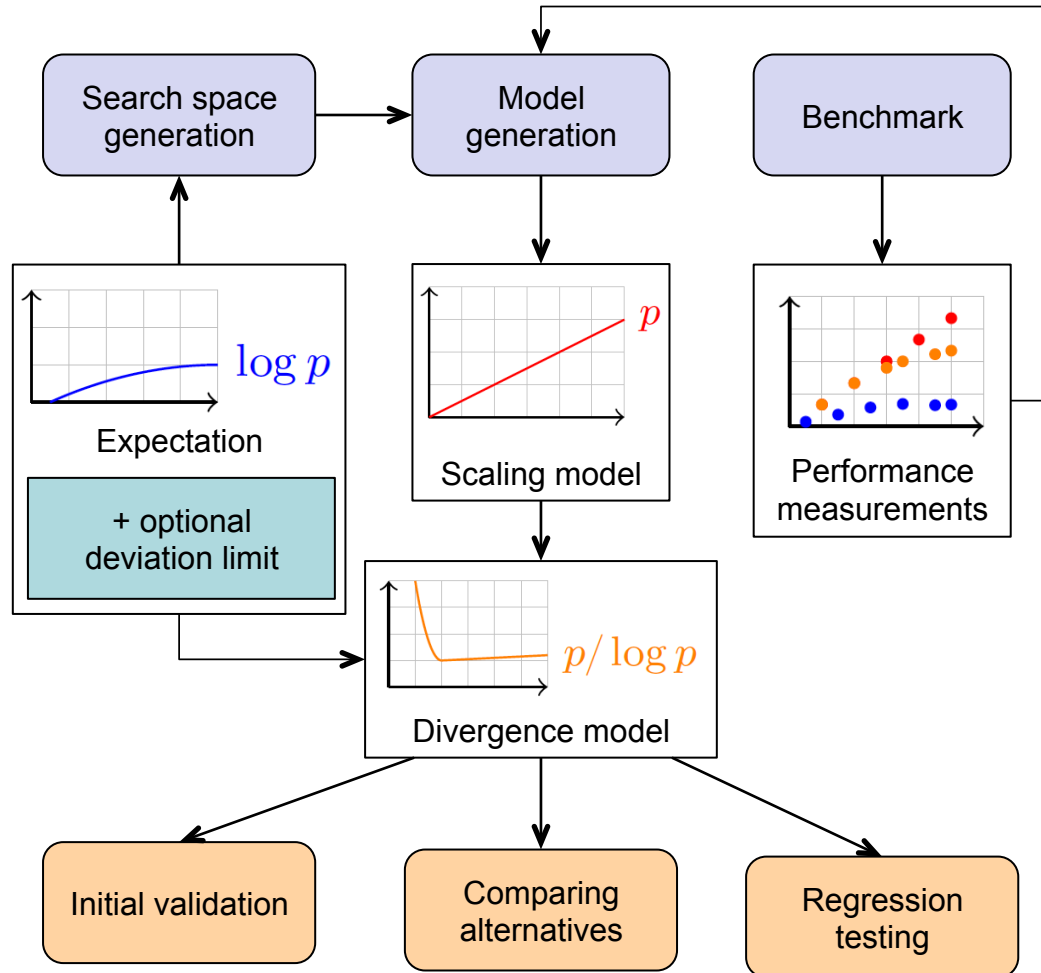
HPC libraries

- Focus on algorithms rather than applications
- Theoretical expectations more common
- Reuse factor makes scalability even more important



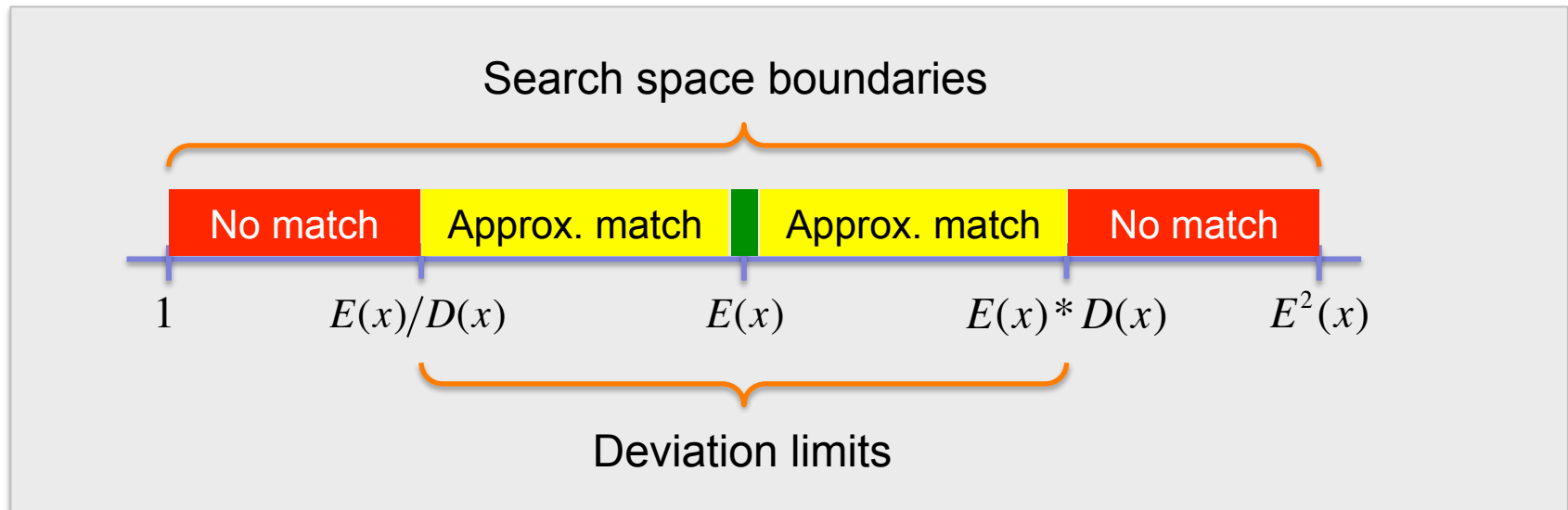
Example:
MPI communication library

Scalability evaluation framework



Customized search space

- Constructed around expectation
- Supports wider range of model functions than original PMNF



Test systems



Platform	Topology	Nodes	CPU	Cores*	Memory*	MPI
Juqueen (BG/Q)	5D torus	28,672	PPC A2	16	16 GB	PAMI
Juropa (Intel)	Fat tree w/ IB	3,288	X5570	8	24 GB	ParTec
Piz Daint (Cray-XC30)	Dragonfly	5,727	E5-2670	8	32 GB	Cray

* Per node

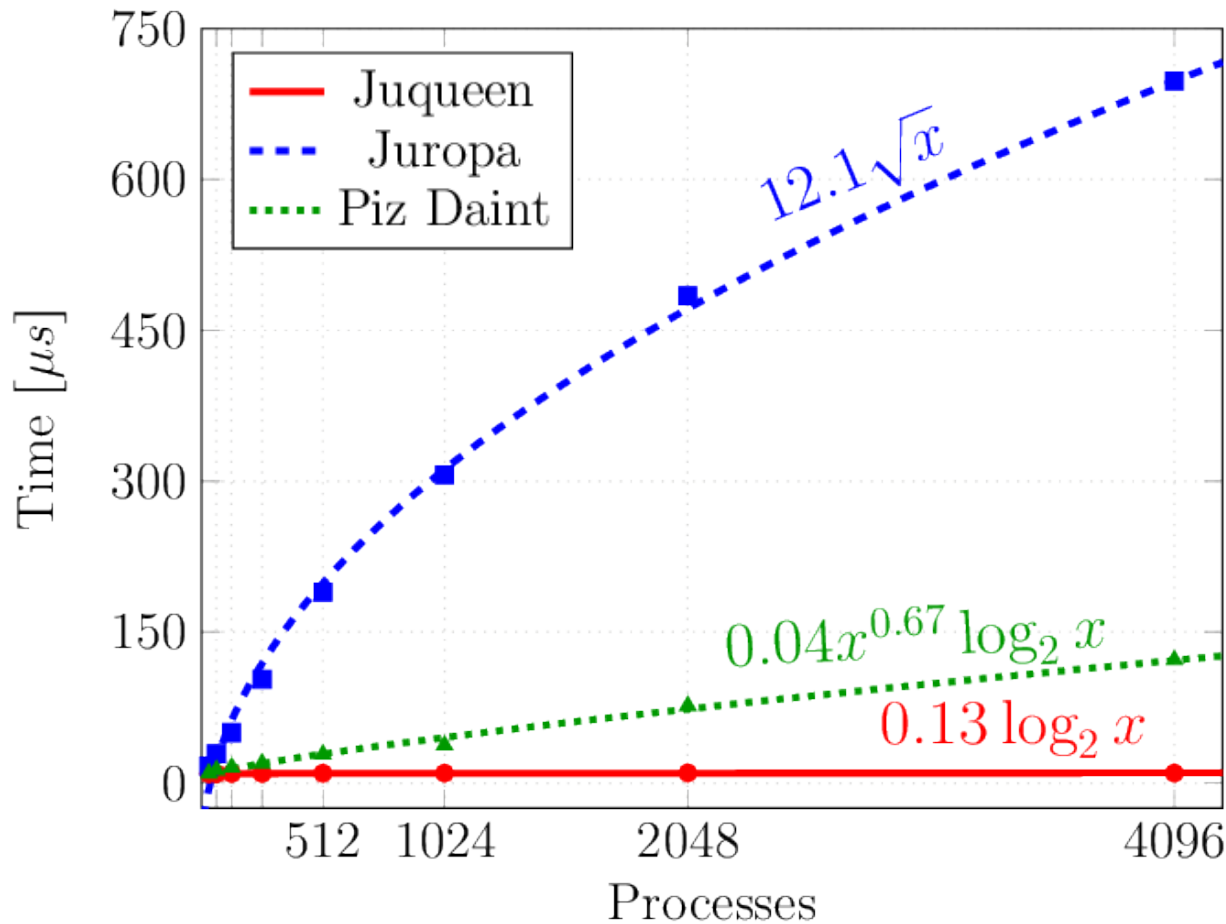


MPI



Platform	Platform	Juqueen	Juropa	Piz Daint	Daint
Barrier [s]					$O(\log p)$
Model	Allreduce [s]		Expectation: $O(\log p)$		
R ²	Model	$O(\log p)$	$O(p^{0.5})$	$O(p^{0.67} \log p)$	$O(\log p)$
Divergence	R ²	0.87	0.99	0.99)
Match					
Bcast [s]	Divergence	0	$O(p^{0.5}/\log p)$	$O(p^{0.67})$	on: $O(p)$
Model	Match	✓	~	✗!)
R ²)
Divergence	Comm_dup [B]		Expectation: $O(1)$		
Match	Model	2.2e5	256	3770 + 18p)
Reduce [s]	R ²	1	1	0.99	on: $O(p)$
Model	Divergence	$O(1)$	$O(1)$	$O(p)$)
R ²	Match	✓	✓	✗)
Divergence)
Match)

Allreduce results



Divergence on Piz
Daint is $O(p^{0.67})$, the
highest of all three

Potential reasons for discrepancies

1. Expectations based on simplified network models (in reality: IB, N-D torus)

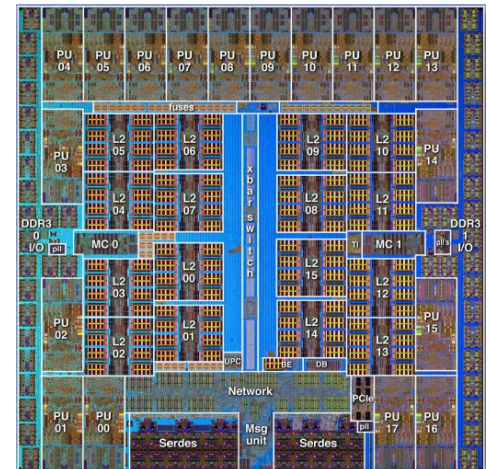
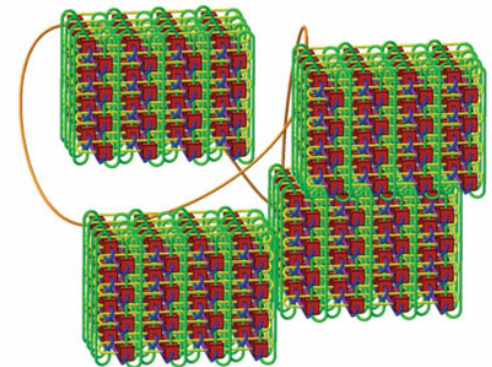
2. Nodes allocation, node's neighborhood

A. Bhatele, K. Mohror, S. H. Langer, and K. E. Isaacs.

There Goes the Neighborhood: Performance Degradation Due to Nearby Jobs, SC' 13

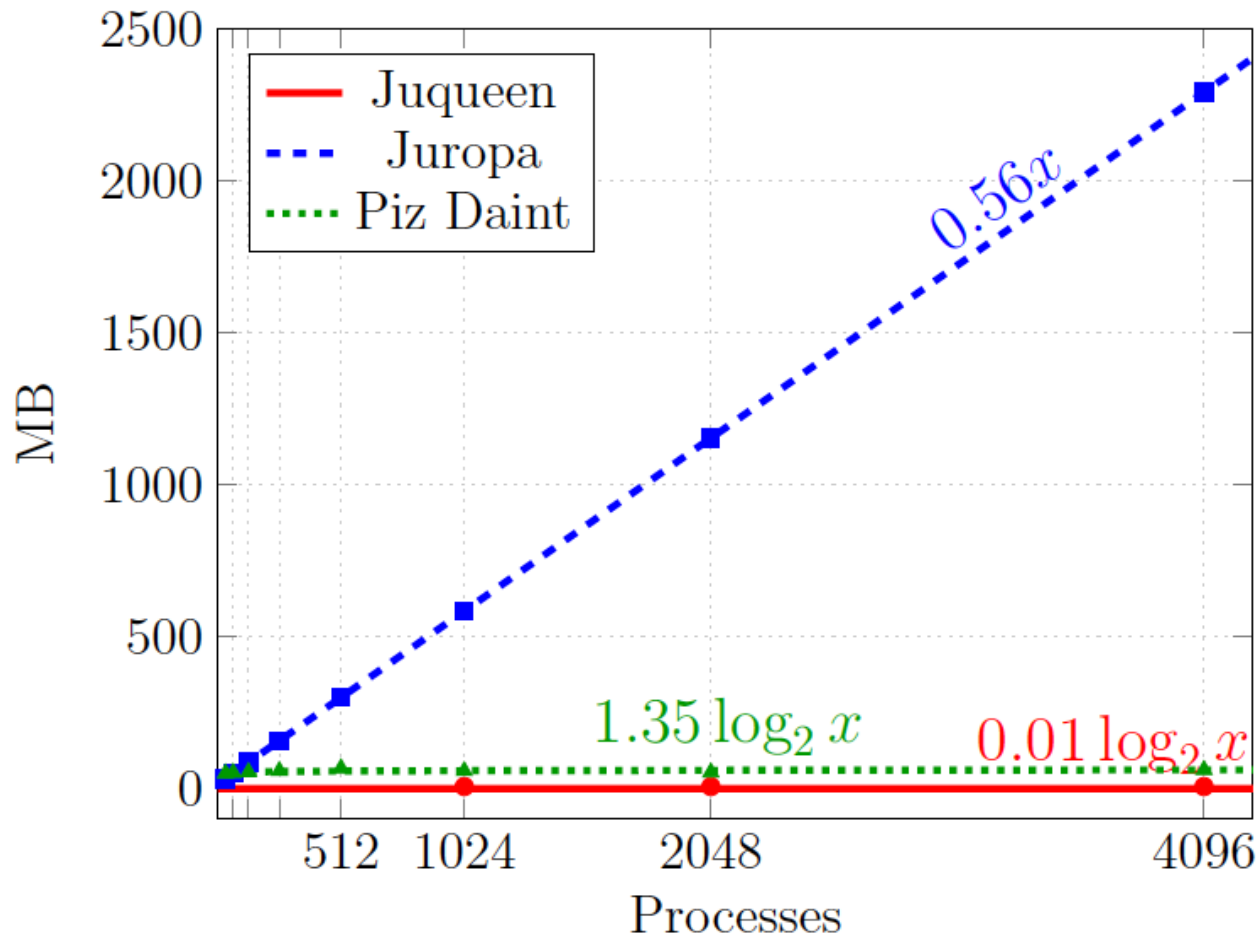
3. Network hardware differences (MU on BG/Q)

T. Hoefler and M. Snir. Generic Topology Mapping Strategies for Large-scale Parallel Architectures, ICS' 11



computing.llnl.gov/tutorials/bgq

MPI memory consumption on all three systems



Linear memory consumption on Juropa

ParaStation MPI uses RC over IB

Sub-space clustering code used in data-mining

- Cluster dimensionality k is the model parameter
- Result: observed behavior matched the expectations

	gen	dedup	pcount	unjoin
Expectation	$O(k^3 2^k)$	$O(k^4 2^k)$	$O(k 2^k)$	$O(k^3 2^k)$
Model	$O(k^4 2^k)$	$O(k^4 2^k)$	$O(k 2^k)$	$O(k^2 2^k)$
Divergence	$O(k)$	$O(1)$	$O(1)$	$O(1/k)$
Match	\sim	✓	✓	\sim

Mass-producing performance models

- 
- Is feasible
 - Offers insight
 - Requires low effort
 - Improves code coverage

References

[1] Christian Iwainsky, Sergei Shudler, Alexandru Calotoiu, Alexandre Strube, Michael Knobloch, Christian Bischof, Felix Wolf: How Many Threads will be too Many? On the Scalability of OpenMP Implementations. In Proc. of the 21st Euro-Par Conference, Vienna, Austria of Lecture Notes in Computer Science, pages 451–463, Springer, August 2015.



[2] Andreas Vogel, Alexandru Calotoiu, Alexandre Strube, Sebastian Reiter, Arne Nägel, Felix Wolf, Gabriel Wittum: 10,000 Performance Models per Minute - Scalability of the UG4 Simulation Framework. In Proc. of the 21st Euro-Par Conference, Vienna, Austria of Lecture Notes in Computer Science, pages 519–531, Springer, August 2015.



[3] Sergei Shudler, Alexandru Calotoiu, Torsten Hoefler, Alexandre Strube, Felix Wolf: Exascaling Your Library: Will Your Implementation Meet Your Expectations?. In *Proc. of the International Conference on Supercomputing (ICS), Newport Beach, CA, USA, pages 1-11, ACM, June 2015*



[4] Alexandru Calotoiu, Torsten Hoefler, Marius Poke, Felix Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes. In Proc. of the ACM/IEEE Conference on Supercomputing (SC13), Denver, CO, USA, pages 1-12, ACM, November 2013.





Thank you!

Upcoming tutorial @ SC15
Insightful Automatic Performance Modeling

Austin, Texas, USA, November 15

