



From numbers to insight via performance models

Georg Hager

Erlangen National High Performance Computing Center (NHR@FAU)

IACS Seminar, University of Stony Brook

2021-10-14



White-/gray-box performance modeling

Introduction to resource-based modeling

The Execution-Cache-Memory (and Roofline) model

Going beyond the node: The pitfalls of highly parallel modeling

Motivation

Analytic performance modeling:

Constructing a simplified model for the interaction between software and hardware in order to understand lowest-order performance behavior

- Basic questions addressed by analytic performance models
 - What is the bottleneck?
 - What is the next bottleneck after optimization?
 - Impact of hardware features \rightarrow co-design, architectural exploration
- What if the model fails?
 - We learn something
 - We may still be able to use the model in a less predictive way

Getting a little more specific

What data/knowledge can a model be based on?

- Only documented hardware properties + hypotheses
 Purely analytic model
- Hardware properties + measurements + hypotheses
 (Partly) phenomenological model
- Measured performance/speedup data + hypotheses
 - Curve-fitting analytic model





white

gray box

An example from physics



Examples for white-/gray-box models in computing



Models and insights







Resource-based performance models



How much of \$RESOURCE does \$STUFF need on \$HARDWARE, and why?

→ Analytic, resource-based, first-principles models

Mechanistic vs. resource-based modeling

Mechanistic

- Cycle-by-cycle
- Latency

Simulators



Resource based

- Resource utilization
- Data flow
- (Non-)overlapping components
- Simplify to make manageable
- If it doesn't work, refine and iterate



A general view on resource bottlenecks

- What is the maximum performance when limited by a bottleneck?
- Resource bottleneck i delivers resources at maximum rate R_i^{max}
- W_i = needed amount of resources

• Minimum runtime:
$$T_i = \frac{W_i}{R_i^{max}} + \lambda_i$$

- Multiple bottlenecks \rightarrow multiple min. runtimes: $T_{\text{expect}} = f(T_1, \dots, T_n)$
- Overall performance:

$$P_{\text{expect}} = \frac{W}{T_{\text{expect}}}$$

A bottleneck model of computing

Example: two bottlenecks

```
#pragma omp parallel for
for(i=0; i<10<sup>7</sup>; ++i)
        a[i] = a[i] + s * c[i];
```



n-core CPU (1 CMG A64FX 2 GHz)

 $W_{flops} = 2 \times 10^7$ flops $W_{BW} = 3 \times 8 \times 10^7$ bytes

$$T_{flops} = \frac{2 \times 10^7 \text{ flops}}{768 \frac{\text{Gflops}}{\text{s}}} = 26.0 \ \mu s$$

$$T_{BW} = \frac{2.4 \times 10^8 \text{ bytes}}{210 \frac{\text{Gbyte}}{\text{s}}} = 1.14 \text{ ms}$$

Bottleneck models

How do we reconcile the multiple bottlenecks? I.e., what is the functional form of $f(T_1, ..., T_n)$?

→ pessimistic model (no overlap):
→ optimistic model (full overlap):

$$f(T_{1}, ..., T_{n}) = \sum_{i} T_{i}$$

$$f(T_{1}, ..., T_{n}) = \max(T_{1}, ..., T_{n})$$

$$R_{i}^{\text{contine} model} = 1,2008$$

$$R_{i}^{\text{contine} et al., 2008}$$

$$R_{i}^{\text{contine} et al., 2008}$$

Roofline for our example: $T_{\min} = \max(T_{flops}, T_{BW}) = 1.14 \text{ ms}$

Maximum performance ("light speed"): $P_{\text{expect}} = \frac{2 \times 10^7}{1.14 \times 10^{-3}} \frac{\text{flops}}{\text{s}} = 17.5 \text{ Gflop/s}$





The Execution-Cache-Memory (ECM) model





ECM modeling workflow

J. Hofmann, C. L. Alappat, <u>G. H.</u>, D. Fey, G. Wellein, DOI: <u>10.14529/jsfi200204</u>



Automating this workflow is possible in some cases:

J. Hammer, J. Eitzinger, <u>G. H.</u>, G. Wellein, DOI: <u>10.1007/978-3-319-56702-0_1</u> (Kerncraft)

J. Laukemann, J. Hammer, <u>G. H.</u>, G. Wellein, DOI: <u>10.1109/PMBS49563.2019.00006</u> (OSACA)

Overlap assumptions



Model validation (FX1000, large pages)



Multicore (in-memory data set) w/ unrolling

C. Alappat et al., DOI: <u>10.1002/cpe.6512</u>



Sufficient unrolling is crucial (but sometimes it's not enough)

Does it work for "real" code, too?

- Preconditioned matrix- free conjugate-gradient solver
- Four systems
 - IBM Power9
 - Cavium/Marvell TX2
 - AMD Naples
 - Intel Skylake
- Yes it does.

J. Hofmann et al., DOI: <u>10.14529/jsfi200204</u>







The pitfalls of composite models in the highly parallel case



Composite analytic models

Plausible assumption: $T = T_{exec} + T_{nexec}$



In practice, $T \neq T_{exec} + T_{nexec}$ and it can go in either direction

Initial observation

Two-socket single-core Pentium IV "Prescott" node (2004-ish)

MPI-parallel Lattice-Boltzmann solver timeline view:



Chipset northbridge)

Memory

Markidis et al. (2015)

Simulator-based analysis

Idle waves perceived as "damped linear waves"

Classical wave equation postulated for continuum description



S. Markidis et al.: *Idle waves in high-performance computing*. Phys. Rev. E **91**(1), 013306 (2015). DOI: <u>10.1103/PhysRevE.91.013306</u>

A more modern platform







→ Spontaneous symmetry breaking, "computational wave" Why? Under which conditions?

Research questions

Setting: MPI- or hybrid-parallel bulk-synchronous barrier-free programs

- How do "disturbances" propagate?
 - Injected idle periods
 - Dependence on communication characteristics
- How do idle waves interact with each other, with noise, and with the hardware?
 - Idle wave decay (noise-induced, bottleneck-induced, topology-induced)
- How do computational waves form? Instabilities?
 - Core-bound vs. memory-bound
 - Amplitude of the computational wave?
- Continuum description?







Idle wave propagation and bottleneck-induced decay



Idle waves interact nonlinearly

A wave-like description cannot be based on a linear model

 Basis for noiseinduced decay of idle waves



Noise-induced idle wave decay

- System or application noise "eats away" on the idle wave
- Statistical details do not matter (only integrated noise power)



Formation of computational wavefronts from idle waves

- 2-socket 10-core
- No decay if in non-saturated regime
- Faster decay with stronger saturation



Application: Chebyshev Filter Diagonalization (ChebFD)

DOI: <u>10.1007/978-3-030-50743-5_20</u>

Computes faster in

- Computes inner eigenvalues of a large sparse matrix
- Blocking optimization: M. Kreutzer, <u>G. H.</u>, D. Ernst, H. Fehske, A.R. Bishop, G. Wellein, DOI: <u>10.1007/978-3-319-92040-5_17</u>
- MPI+OpenMP hybrid, topological insulator matrix, Emmy@RRZE

desynchronized state 50 $1T, n_{\rm b} = 2$ $1T, n_b = 32$ $n_{\rm b} = 32$ 400 Performance [Gflop/s] Performance [Gflop/s] 300 405T, $n_{\rm b} = 2$ 5T, $n_{\rm b} = 32$ Rank 40 10T, $n_{\rm b} = 2$ 10T, $n_{\rm b} = 32$ $n_{\rm b} = 2$ 80 30 120200 159 200 20 00Rank 10 100 15 8 10 2 6 8 0 6 $\overline{4}$ 1010Time step Number of nodes Number of node Cores per socket, N

(a) Single socket performance (b) MPI only vs. hybrid communication, $n_b = 2$, $n_b = 32$

(c) Timeline visulization

Current results

- Instability of bulk-synchronous barrier-free programs is bound to the presence of a resource bottleneck
- Desynchronized bottlenecked programs can exhibit automatic communication/execution overlap via formation of computational waves
- Idle waves can be absorbed by fine-grained system noise, and the mechanism behind this is well understood
- Idle waves can decay via topological noise caused by inhomogeneous communication characteristics
- Proof that noise statistics is largely irrelevant for idle wave decay rate
- Analytic model for idle wave velocity w.r.t. communication topology and characteristics
- Experimental evidence that MPI collectives can be transparent to idle waves

Future directions

- Development of a comprehensive, bottleneck-aware simulator framework for message-passing programs
- Analytic description of decaying wave for bottleneck-triggered decay
- Bottlenecks other than memory bandwidth
- Analytic understanding of computational wave amplitude w.r.t. communication characteristics and bottleneck saturation
- Idle wave phenomena in irregular programs
- Physical model for coupled processes (Kuramoto-like)

$$\dot{\theta_i} = \omega_i + \alpha \sum_j T_{ij} V(\theta_j - \theta_i)$$

 Continuum description of parallel system as a nonlinear (dissipative?) medium





Thank you.

