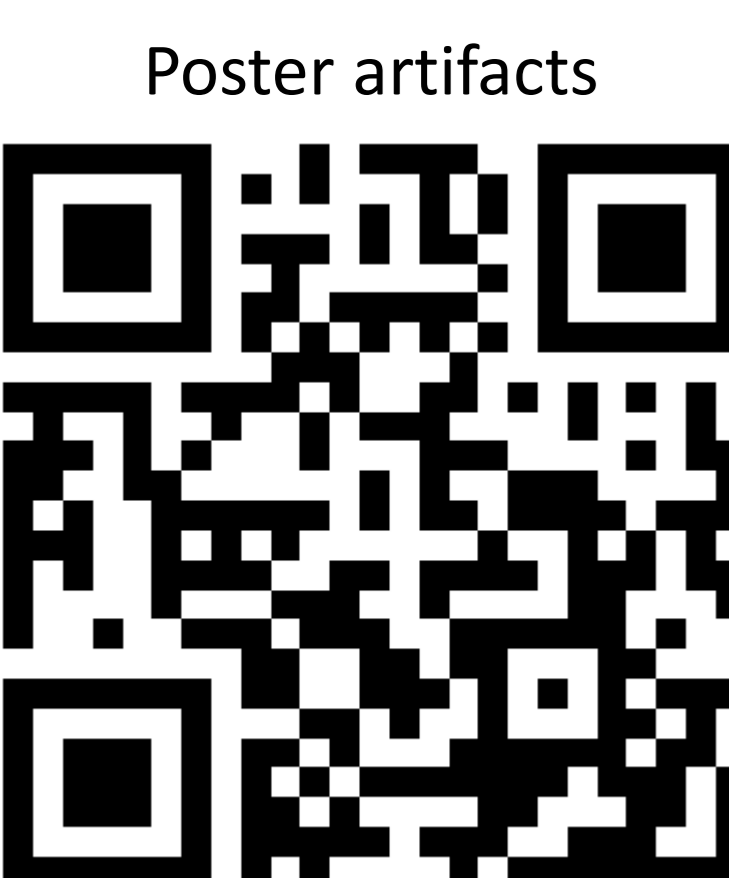
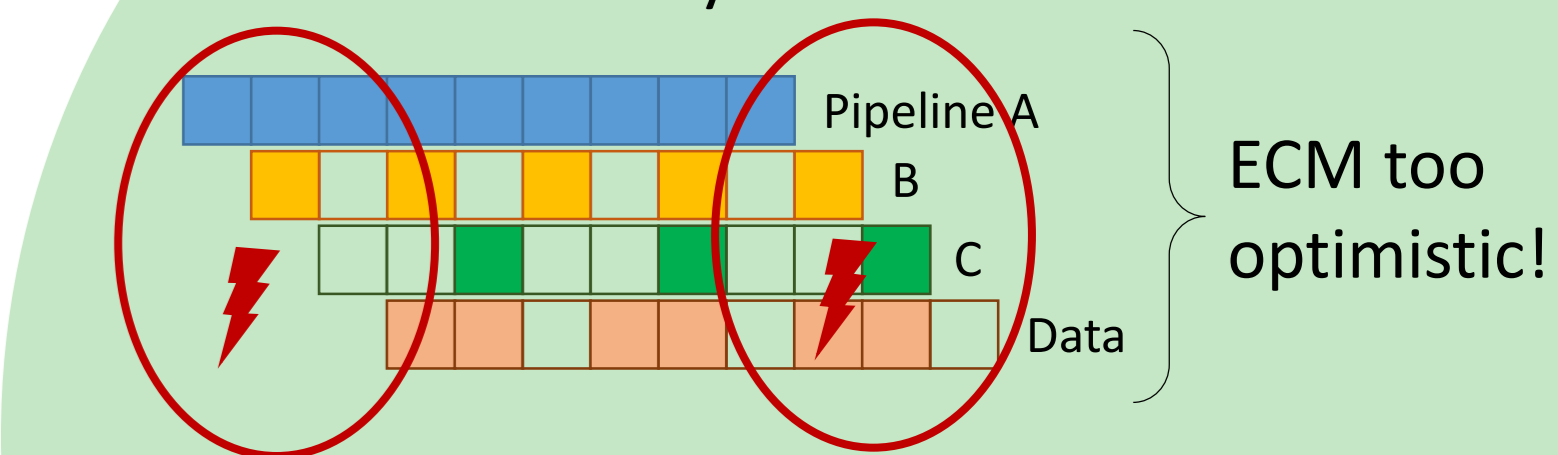


# Applying the Execution-Cache-Memory Performance Model: Current State of Practice

G. Hager<sup>1</sup>, J. Eitzinger<sup>1</sup>, J. Hornich<sup>1</sup>, F. Cremonesi<sup>2</sup>, C. A. Alappat<sup>1</sup>, T. Röhl<sup>1</sup>, G. Wellein<sup>1</sup>

<sup>1</sup> Erlangen Regional Computing Center (RRZE), Erlangen, Germany <sup>2</sup> École Polytechnique Fédérale de Lausanne (EPFL), Geneva, Switzerland

**Problem:**  
Non-steady-state, latencies  
Main ECM model prerequisite:  
Steady-state execution



https://tiny.cc/ECM-SC18-AD

J. Hofmann et al.: On the accuracy and usefulness of analytic energy models for contemporary multicore processors. DOI: 10.1007/978-3-319-92040-5\_2

J. Hammer et al.: Kerncraft: A Tool for Analytic Performance Modeling of Loop Kernels. DOI: 10.1007/978-3-319-56702-0\_1

G. Hager et al.: Exploring performance and power properties of modern multicore chips via simple machine models. DOI: 10.1002/cpe.3180

H. Stengel et al.: Quantifying performance bottlenecks of stencil computations using the Execution-Cache-Memory model. DOI: 10.1145/2751205.2751240

T. M. Malas et al.: Optimization of an electromagnetics code with multicore waveform diamond blocking and multi-dimensional intra-tile parallelization. DOI: 10.1109/IPDPS.2016.87

T. Ewart et al.: Neuromapp: A Mini-application Framework to Improve Neural Simulators. DOI: 10.1007/978-3-319-58667-0\_10

Layer Condition Calculator: [tiny.cc/LayerConditions](http://tiny.cc/LayerConditions)

Kerncraft: [tiny.cc/kerncraft](http://tiny.cc/kerncraft)

Pycachesim: [tiny.cc/pycachesim](http://tiny.cc/pycachesim)

IACA: [tiny.cc/IACA](http://tiny.cc/IACA)

OSACA: [tiny.cc/OSACA](http://tiny.cc/OSACA)

LIKWID tools: [tiny.cc/LIKWID](http://tiny.cc/LIKWID)

Supported by

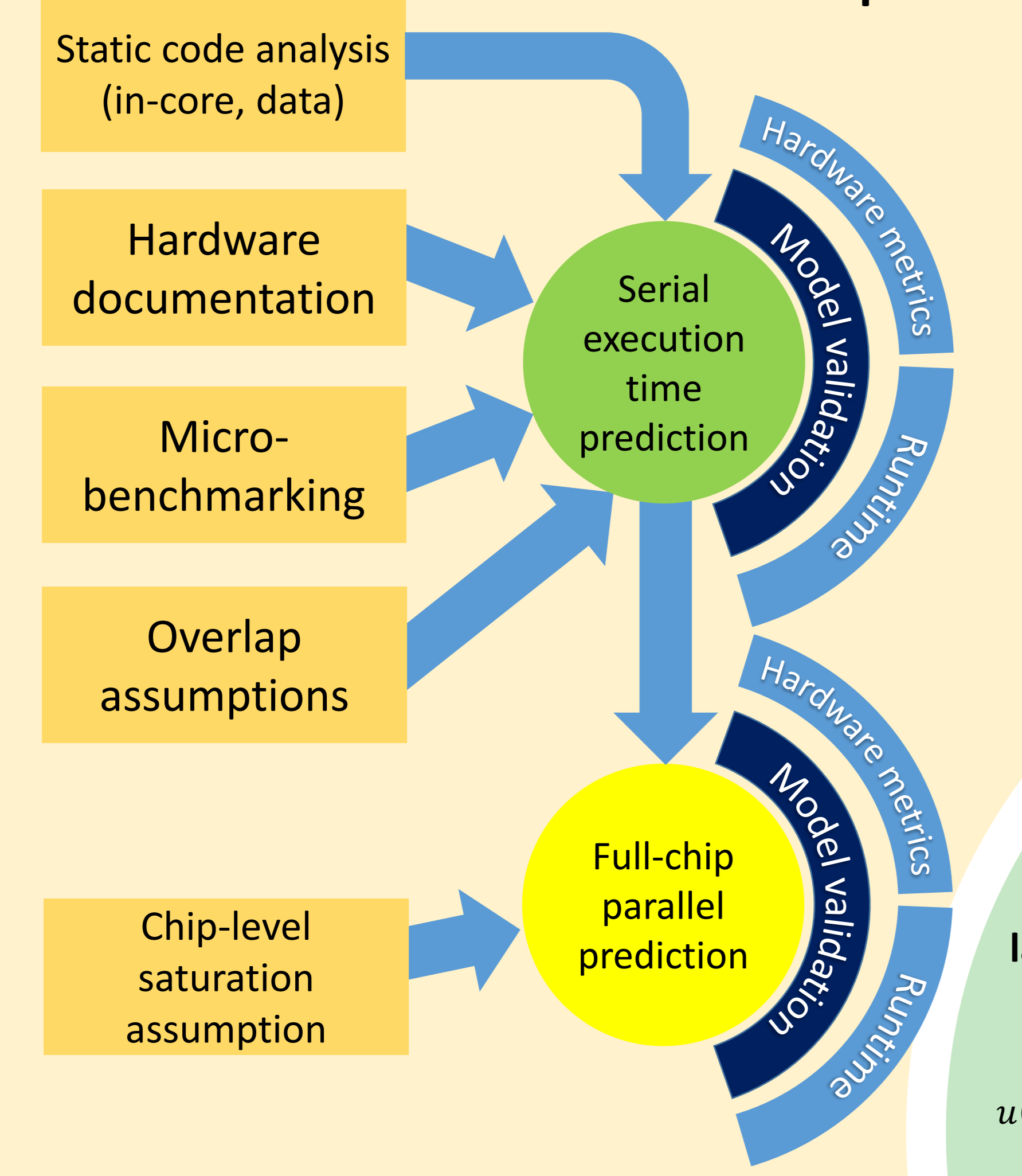


Embedded MM content



https://tiny.cc/ECM-SC18

## An analytic runtime model for loop code on the core and chip level



## Cache architectures and transfer times

Optimistic transfer times:  $T_i = \frac{V_i}{b_i}$ . Example: STREAM triad  $a(:) = b(:) + s * c(:)$  on Xeon Broadwell.2.3 GHz in CoD mode: inclusive caches,  $b_{L3Mem} = 32$  GB/s per NUMA domain (saturated). Short notation:  $\{T_{L1L2} | T_{L2L3} | T_{L3Mem}\} = \{4 | 8 | 18.4\}$  cy/8 iter. Analysis tools: Pen & paper, Layer conditions, pycachesim, Kerncraft. Traffic validation tools: Likwid-perfctr, PAPI, Linux perf.

**Problem: Saturation**  
Plain ECM model too optimistic @ saturation for tight kernels with small  $T_{OL}$ . Refinement: Adaptive latency penalty, depends on bus utilization  $u(n)$ :  $u(1) = \frac{T_{L3Mem}}{T_{ECM}^{Mem}}$ ,  $u(n) = \frac{T_{L3Mem}}{T_{ECM}^{Mem} + (n-1)u(n-1)p_0}$ . Hofmann et al., ISC 2018.

## Putting the model together: Overlap assumptions

Notation for model contributions:  $\{T_{OL} | T_{nOL} | T_{L1L2} | T_{L2L3} | T_{L3Mem}\} = \{7 | 2 | 4 | 8 | 18.4\}$  cy/8 iter. Most pessimistic assumption: no overlap of data-related contributions.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL} + T_{L1L2} + T_{L2L3} + T_{L3Mem})$ . Any deviation from non-overlap behavior in the hardware makes model non-optimistic!

Most optimistic assumption: full overlap of data-related contributions.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL}, T_{L1L2}, T_{L2L3}, T_{L3Mem})$ . Fully optimistic (light speed) model, but not the same as Roofline:  $T_i = \frac{V_i}{b_{i,max}}$  for  $i \in \{MemReg, \dots\}$ .

Mixed model: partial overlap of data-related contributions. Example: no overlap at L1, full overlap of all other contributions.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL} + T_{L1L2}, T_{L2L3}, T_{L3Mem})$ .

## Notation for model predictions

Example: no-overlap model.  $\{T_{L1}^{ECM} | T_{L2}^{ECM} | T_{L3}^{ECM} | T_{Mem}^{ECM}\}$ . Example: AVX DP sum reduction w/ single accumulator on Broadwell EP CoD 2.3 GHz ( $b_{L3Mem} = 32$  GB/s).  $\{6 | 6 | 6 | 7.6\}$  cy/8 iter.  $\rightarrow$  Constant performance up to L3.

## Multicore scaling and saturation

Optimistic assumption: Performance scaling is linear until a bandwidth bottleneck (e.g.,  $b_{Mem}$ ) is saturated  $\rightarrow$  memory transfer time as lower limit. Runtime vs. cores (memory bottleneck): Number of cores at saturation:  $T_{ECM}(n) = \max(\frac{T_{ECM}^{Mem}}{n}, T_{L3Mem}) \Rightarrow n_s = \frac{T_{ECM}^{Mem}}{T_{L3Mem}}$ . Example: AVX DP sum reduction from above.  $\{6 | 6 | 6 | 7.6\}$  cy/8 iter.  $\Rightarrow n_s = \frac{7.6}{4.6} = 2$ .

## In-core execution model

Example: 2D DP 5-pt stencil on Xeon Skylake Gold 6148, AVX-512. `do i=1,N do j=1,N do k=1,N y(i,j) = 0.25d0 * (x(i-1,j) + x(i+1,j) + x(i,j-1) + x(i,j+1))` enddo. Throughput / critical path analysis. Best case: max throughput. Worst case: critical path.  $T_{core}^{min} = \max(T_{nOL}, T_{OL})$ ,  $T_{core}^{max} = T^{CP}$ .  $T_{nOL}$  interacts with cache hierarchy,  $T_{OL}$  does not.

## From time to performance

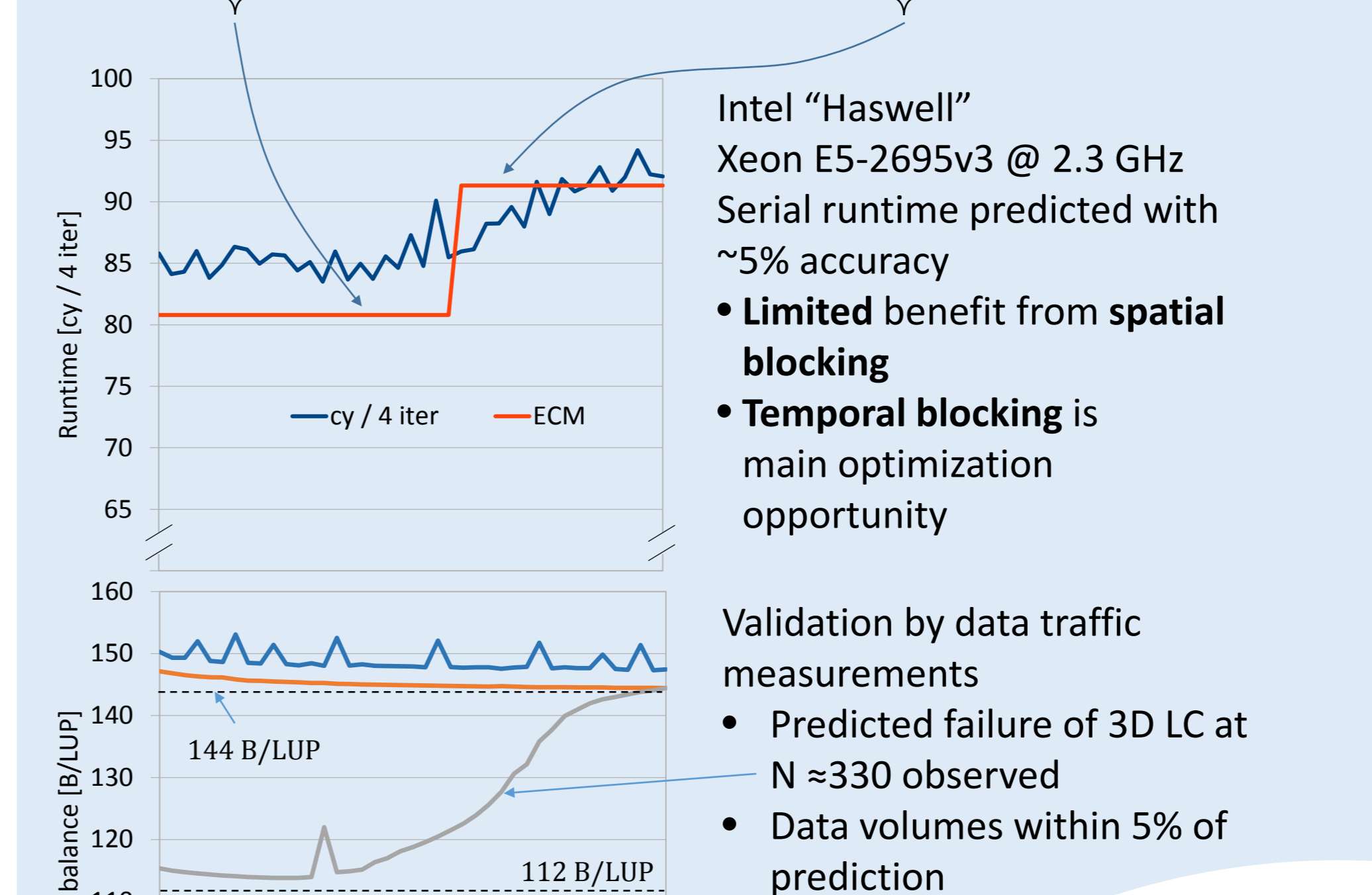
$f$  CPU frequency,  $f_0$  CPU base frequency,  $W$  Amount of work.  $P = \frac{W}{T_{ECM}(f, f_0)}$ .

## Example: Non-overlapping CPU with on-chip clock domain down to L3, $b_{Mem}$ independent of $f$

$P_{L1}^{ECM} \dots P_{Mem}^{ECM} = \frac{W \cdot f}{\{T_{L1}^{ECM} | T_{L2}^{ECM} | T_{L3}^{ECM}\} \max(T_{L3}^{ECM} + T_{L3Mem} \cdot f/f_0, T_{OL})}$ .  $f$ -independent  $P$  component.

## Application: Complex stencils

FD2D code solving Maxwell's equations via THIM for thin-film solar cells, AVX2 Vectorization via C intrinsics, complex arithmetic. Malas et al., IPDPS 2016. LC fulfilled ( $B_{C}^{Mem} = 112$  B/LUP):  $\{10 | 8 | 18 | 18 | 36.8\}$  cy/4 iter. LC broken ( $B_{C}^{Mem} = 144$  B/LUP):  $\{10 | 8 | 18 | 18 | 47.3\}$  cy/4 iter.



## Application: Blue Brain Project kernels

What is the computational cost of a synapse? Study performance of brain simulation code via mini-apps  $\rightarrow$  SSE4.2 vectorization by compiler. Case 1: Synaptic current kernel. Ewart et al., ISC 2017. Challenges: `exp()` function call, some divides, some indirect accesses, integer register spill (many streams, call ABI restricts untouched registers). Streaming kernel, adjust data volume via knowledge about structure of index arrays (4% correction). Strong intra-iteration dependency chain. `exp()` overhead (thruput/latency) measured by microbenchmark, added to ECM.

IVY close to CP prediction, HSW data bound! Still saturating @ 3-5 cores on both CPUs!

## Case 2: Sodium ion channel (NaTs2\_t state)

for( $i_{ml} = 0; i_{ml} < cnt_{ml}; ++i_{ml}$ ) { `_nd_idx = _ni[_i_ml]; v = _vec_v[_nd_idx]; mAlpha[_i_ml] = ( 1.0 - ( exp( - ( v + 32.0 ) / 6.0 ) ) ) / ( 1.0 - ( exp( - ( v - 32.0 ) / 6.0 ) ) ) ; mBeta[_i_ml] = ( 0.124 * ( - ( v - 32.0 ) ) / ( 1.0 - ( exp( - ( v - 32.0 ) / 6.0 ) ) ) ) ; mInf[_i_ml] = mAlpha[_i_ml] / ( mAlpha[_i_ml] + mBeta[_i_ml] ) ; mTau[_i_ml] = ( 1.0 / ( mAlpha[_i_ml] + mBeta[_i_ml] ) ) / 2.95 ; hAlpha[_i_ml] = ( - 0.015 * ( v + 60.0 ) ) / ( 1.0 - ( exp( ( v + 60.0 ) / 6.0 ) ) ) ; hBeta[_i_ml] = ( - 0.015 * ( - ( v - 60.0 ) ) / ( 1.0 - ( exp( - ( v - 60.0 ) / 6.0 ) ) ) ) ; hInf[_i_ml] = hAlpha[_i_ml] / ( hAlpha[_i_ml] + hBeta[_i_ml] ) ; hTau[_i_ml] = ( 1.0 / ( hAlpha[_i_ml] + hBeta[_i_ml] ) ) / 2.95 ; m[_i_ml] = m[_i_ml] + ( 1. - exp(dt*( - 1.0 / mTau[_i_ml] ) ) ) * ( - ( mInf[_i_ml] / mTau[_i_ml] ) ) / ( - 1.0 / mTau[_i_ml] ) - m[_i_ml] ; h[_i_ml] = h[_i_ml] + ( 1. - exp(dt*( - 1.0 / hTau[_i_ml] ) ) ) * ( - ( hInf[_i_ml] / hTau[_i_ml] ) ) / ( - 1.0 / hTau[_i_ml] ) - h[_i_ml] ; }`

"Ivy Bridge" E5-2660v2: Full throughput  $T_{OL}$ .  $\{140 | 9.25 | 5.15 | 5.15 | 9\} \rightarrow T_{ECM}^{Mem} \geq 140$  cy. Measurement: 191 cy. Saturation @ 16 cores!

## Overlap assumptions for current architectures

Intel Skylake SP: no overlap in data transfers, victim L3.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL} + T_{L1L2} + T_{L2L3} + T_{L3Mem})$ . Intel Xeon  $\leq$  Broadwell EP: no overlap.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL} + T_{L1L2} + T_{L2L3} + T_{L3Mem})$ . IBM Power8: overlap at L1.  $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL}, T_{L1L2}, T_{L2L3} + T_{L2Mem} + T_{L3Mem})$ . AMD Zen (Epyc): full overlap.  $T_{ECM}^{Mem} = \max(T_{nOL}, T_{L1L2}, T_{L2L3}, T_{L3Mem}, T_{OL})$ .

## Application: Conjugate Gradient solver

while( $\alpha_0 < tol$ ): ECM [cy/8 iter],  $T_{ECM}^{Mem}$  [cy/8 iter],  $n_s$  [cores], Full domain limit [cy/8 iter].  $\vec{v} = A\vec{p}$  { 8 | 4 | 6.7 | 10 | 16.9 } 37.6 3 16.9.  $\lambda = \alpha_0 / (\vec{v}, \vec{p})$  { 2 | 2 | 2.7 | 4 | 9.1 } 17.8 2 9.11.  $\vec{x} = \vec{x} + \lambda\vec{p}$  { 2 | 4 | 6 | 16.9 } 29.0 2 16.9.  $\vec{r} = \vec{r} - \lambda\vec{v}$  { 2 | 4 | 6 | 16.9 } 29.0 2 16.9.  $\alpha_1 = \langle \vec{r}, \vec{r} \rangle$  { 2 | 2 | 1.3 | 2 | 4.6 } 9.90 3 4.56.  $\vec{p} = \vec{r} + \frac{\alpha_1}{\alpha_0}\vec{p}$ ,  $\alpha_0 = \alpha_1$  { 2 | 4 | 6 | 16.9 } 29.0 2 16.9. Sum 152 81.3.

- Multi-loop code well represented
- Single core performance predicted with 5% error
- Saturated performance predicted with < 0.5% error
- Saturation point predicted approximately
- Future work: Include GS preconditioner  $\rightarrow$  HPC modeling