#### ERLANGEN REGIONAL COMPUTING CENTER



#### **Making Sense of Performance Numbers**

Georg Hager Erlangen Regional Computing Center (RRZE) Friedrich-Alexander-Universität Erlangen-Nürnberg

OpenMPCon 2018 Barcelona, 2018-09-25



FRIEDRICH-ALEXANDER UNIVERSITÄT ERLANGEN-NÜRNBERG

#### Agenda

- Performance Engineering
- Analytic performance models
- A simple example: Sparse matrix-vector multiplication
- Another example: Sparse matrix-transpose vector multiplication
- An advanced chip-level model: ECM
- Yet another example: Composite modeling of a (P)CG solver





## **Performance Engineering**

- Performance Engineering (PE): a structured process based on analytic (white-/gray-box) models to optimize/parallelize codes
- Basic questions addressed by analytic performance models
  - What is the bottleneck? → optimization technique
  - What is the next bottleneck? → performance potential of the optimization
  - Am I done? What about other hardware?
  - Impact of processor frequency and socket scalability
     Appropriate execution parameters, energy-optimized operating point
- Engineering for Performance in High Performance Computing (Bill Gropp; PASC 2015): <u>https://www.youtube.com/watch?v=sadfSARXSC0</u>





## **Motivation for analytic modeling**

- Advantages of analytic models
  - Identification of universality
  - Identification of governing mechanisms
  - Insight via model nature
  - Insight via model failure

- Performance models
  - Microarchitecture analysis: Determine bottlenecks and influencing factors
  - Design space exploration: What would happen if resource X were improved?





#### NODE LEVEL PERFORMANCE MODEL FOR SPARSE MATRIX (TRANSPOSE) VECTOR MULTIPLY





FRIEDRICH-ALEXANDER UNIVERSITÄT RLANGEN-NÜRNBERG

FF 🔤 🖃

## SpMV – Roofline Model (Comp. Intensity)

Sparse MVM in double precision with CRS data storage: (N<sub>nzr</sub> : avg. non-zeros per row)

do i = 1, 
$$N_r$$
 do j = row\_ptr(i), row\_ptr(i+1) - 1
 C(i) = C(i) + val(j) \* B(col\_idx(j))
 enddo
enddo
enddo

2

 $8 + 4 + 8\alpha + 20/N_{nzr}$ 

#### Double precision computational intensity

- α quantifies traffic for loading RHS (B)
  - >  $\alpha = 0 \rightarrow \text{RHS}$  is in cache
  - →  $\alpha = 1/N_{nzr}$  → RHS loaded once
  - >  $\alpha = 1 \rightarrow$  no cache
  - >  $\alpha > 1 \rightarrow$  Houston, we have a problem!
- "Expected" performance = b<sub>S</sub> × I<sub>CRS</sub> (Roofline model)

#### See

W. D. Gropp, D. K. Kaushik, D. E. Keyes, and B. F. Smith, "Towards realistic performance bounds for implicit CFD codes," in Proceedingsof Parallel CFD99. Elsevier, 1999, pp. 233–240.

 $I_{CRS}^{DP}$ 

- G. Schubert et al. Parallel Processing Letters 21(3), 339-358 (2011). DOI: 10.1142/S0129626411000254
- M. Kreutzer et al.: SIAM Journal on Scientific Computing 36(5), C401–C423 (2014). DOI: 10.1137/130930352,



flops

byte

### **SpMV** – Quantifying RHS impact ( $\alpha$ )

$$I_{CRS}^{DP} = \frac{2}{8+4+8\alpha+20/N_{nzr}} \frac{\text{flops}}{\text{byte}} = \frac{N_{nz} \cdot 2 \text{ flops}}{V_{meas}}$$

Vmeas is measured overall memory data traffic (using, e.g., likwid)

$$\alpha = \frac{1}{4} \left( \frac{V_{meas}}{N_{nz} \cdot 2 \text{ bytes}} - 6 - \frac{10}{N_{nzr}} \right)$$

Example: kkt\_power matrix from the UoF collection

Intel Xeon Sandy Bridge (8c; 20 MB L3)

• 
$$N_{nz} = 14.6 \cdot 10^6$$
,  $N_{nzr} = 7.1$ 

- $V_{meas} \approx 258 \text{ MB} \rightarrow \alpha = 0.36, \alpha N_{nzr} = 2.5$
- → RHS is loaded 2.5 times from memory

$$\frac{I_{CRS}^{DP}(1/N_{nzr})}{I_{CRS}^{DP}(\alpha)} = 1.11$$





## **SpMV – understanding performance**

#### Intel Xeon E5-2680 (**8c**@2.7 GHz; 20 MB L3) Assumption: $\alpha = \frac{1}{N_{nzr}} \rightarrow P = I_{CRS}^{DP} \times b$

#	Test case	N	$N_{ m nz}$	$N_{ m nzr}$
1	RM07R	381,689	37,464,962	98.16
<b>2</b>	kkt_power	2,063,494	14,612,663	7.08
3	Hamrle3	1,447,360	5,514,242	3.81
4	ML_Geer	$1,\!504,\!002$	$110,\!879,\!972$	73.72
5	pwtk	217,918	11,634,424	53.39
6	shipsec1	140,874	7,813,404	55.46
7	$\operatorname{consph}$	83,334	6,010,480	72.13
8	pdb1HYS	36,417	4,344,765	119.31
9	cant	62,451	4,007,383	64.17
10	cop20k_A	121,192	2,624,331	21.65
11	rma10	46,835	2,374,001	50.69
12	mc2depi	525,825	2,100,225	3.99
13	qcd5_4	49,152	1,916,928	39.00
14	mac_econ_fwd500	206,500	1,273,389	6.17
15	scircuit	170,998	958,936	5.61
16	rail4284	$^{4,284 imes}_{1,092,610}$	11,279,748	2,632.99
17	dense2	2,000	4,000,000	2,000.00
18	webbase-1M	1,000,005	3,105,536	3.11



#### 4 corner case matrices from UoF collection

#### Williams collection http://www.nvidia.com/content/NV\_Research/matrices.zip

M. Kreutzer et al.: A Unified Sparse Matrix Data Format for Efficient General Sparse Matrix-Vector Multiplication on Modern Processors with Wide SIMD Units. SIAM Journal on Scientific Computing 2014 36:5, C401-C423

#### **SpMV – understanding performance**

Socket scaling – more recent architectures

PWTK matrix (N<sub>nzr</sub>= 53 &  $\alpha = \frac{1}{N_{nzr}} \rightarrow I_{CRS}^{DP} = \frac{2F}{12.5B} \rightarrow P = I_{CRS}^{DP} \times b$ )



#### **SpMTransposeV – understanding performance**



## **SpMTV – understanding performance**



Publication in preparation

#### THE ECM PERFORMANCE MODEL



#### A quick walk-through





## ECM model components: Data transfer times

- Optimistic transfer times through mem hierarchy
- $T_i = \frac{V_i}{b_i}$
- Transfer time notation for some given loop kernel:

 $\{T_{L1L2} | T_{L2L3} | T_{L3Mem}\} =$   $\{4 | 8 | 18.4\} \text{ cy/8 iter}$ 

- Input:
  - Cache properties
  - Application data transfer prediction







#### ECM model components: In-core execution



#### ECM model components: 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 overlap model: no overlap
  - $T_{ECM}^{Mem} = \max(T_{OL}, T_{nOL} + T_{L1L2} + T_{L2L3} + T_{L3Mem})$  for in-mem data





#### ECM model: Notation for runtime predictions







## **ECM model: (Naive) saturation assuption**

 Performance is assumed to scale across cores until a shared bandwidth bottleneck is hit

$$T_{ECM}(n) = \max\left(\frac{T_{Mem}^{ECM}}{n}, T_{L3Mem}\right) \implies n_{S} = \left[\frac{T_{ECM}^{Mem}}{T_{L3Mem}}\right]$$
Roofline bandwidth ceiling

 This is (sometimes) too optimistic near the saturation point. For improvements see

J. Hofmann, G. Hager, and D. Fey: *On the accuracy and usefulness of analytic energy models for contemporary multicore processors*. Proc. ISC High Performance 2018. DOI: <u>10.1007/978-3-319-92040-5\_2</u>





#### MODELING A CONJUGATE-GRADIENT SOLVER



Building a model from components





#### A matrix-free CG solver

- 2D 5-pt FD Poisson problem
- Dirichlet BCs, matrix-free
- $N_x \times N_y = 40000 \times 1000$  grid
- CPU: Haswell E5-2695v3 CoD mode









#### **Machine characteristics**

- 2.3 GHz (fixed core & Uncore)
  - AVX2, 2 x FMA per cycle
  - 2 load & 1 store per cycle (in practice: 1+1)
- Cache characteristics
  - Inclusive, non-overlapping
  - $b_{L1L2} = 43$  Byte/cy (gray)
  - $b_{L2L3} = 32$  Byte/cy (white)
- Memory bandwidth per ccNUMA domain (saturated)
  - *b<sub>read</sub>* = 32.3 GByte/s = 11.3 Byte/cy
  - *b<sub>copy</sub>* = 26.1 GByte/s = 14.0 Byte/cy



#### **ECM model composition**

#### Naive implementation of all kernels (omp parallel for)

<b>w</b> hile( $\alpha_0 < \text{tol}$ ):	<i>T</i> <sub><i>x</i></sub> [cy/8 iter]	T <sup>ECM</sup> [cy/8 iter]	n <sub>s</sub> [cores]	Full domain limit [cy/8 iter]
$\vec{v} = A\vec{p}$	{ 8    4   6.7   10   16.9 }	37.6	3	16.9
$\lambda = \alpha_0 / \langle \vec{v}, \vec{p} \rangle$	{ 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





## CG performance – 1 core to full socket

- Multi-loop code well
   represented
- Single core performance predicted with 5% error
- Saturated performance predicted with
   < 0.5% error</li>
- Saturation point predicted approximately
  - Can be fixed by improved ECM model







#### **Barrier cost**

- Does the OpenMP barrier impact the performance?
- At which problem size can this be expected?



#### CG with GS preconditioner: Naïve parallelization

Pipeline parallel processing: OpenMP barrier after each wavefront update (ugh!)







# CG with GS preconditioner: additional kernels



- Back substitution does not saturate the memory bandwidth!
  - → full algorithm does not fully saturate
- Impact of barrier still negligible overall, but noticeable in the preconditioner



#### **PCG** measurement

- <2% model error for single threaded and saturated performance
- Expected large impact of barrier at smaller problem sizes in x direction





#### Conclusions

- Analytic modeling is worth the effort
  - Even if it's inaccurate
- Even Roofline can yield amazing insights
- Analytic modeling
  - is not just for "simple kernels"
  - is composable
  - can cover programming model overhead, too
- ... and yes, it is real work.
- No, it does not always work.

http://tiny.cc/kerncraft



Automatic Roofline/ECM modeling tool