When 1 + 1 > 2The Power of Interdisciplinary Research

Paolo Bientinesi + Edoardo di Napoli

Aachen Institute for Advanced Study in Computational Engineering Science RWTH Aachen University

SLQM Opening workshop, JSC April 4, 2017





1 Historical perspective

2 Eigensolvers

3 Tensors



Good old times





- Separation of concerns
- Specialization
- High-performance
- Standardization
- Layering

$$y = X\beta + Zu + \epsilon$$

$$\min_{\mathbf{x}} \|A\mathbf{x} - \mathbf{b}\|^2 + \|\Gamma\mathbf{x}\|^2$$

LINEAR MIXED MODELS

$$V_{LJ} = 4\varepsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^{6} \right]$$

Lennard-Jones potential

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{r},t) = \left[\frac{-2\hbar}{2\mu} \nabla^2 + V(\mathbf{r},t)\right] \Psi(\mathbf{r},t)$$

SCHRÖDINGER EQN.



But...

$$\mathbf{y} = X\boldsymbol{\beta} + Z\mathbf{u} + \boldsymbol{\epsilon}$$
$$\min_{x} \|A\mathbf{x} - \mathbf{b}\|^{2} + \|\Gamma\mathbf{x}\|^{2}$$
$$\text{LINEAR MIXED MODELS}$$
$$V_{LJ} = 4\varepsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^{6} \right]$$

LENNARD-JONES POTENTIAL

$$i\hbar\frac{\partial}{\partial t}\Psi(\mathbf{r},t) = \left[\frac{-2\hbar}{2\mu}\nabla^2 + V(\mathbf{r},t)\right]\Psi(\mathbf{r},t)$$

SCHRÖDINGER EQN.







But...

Paolo Bientinesi + Edoardo di Napoli | When 1 + 1 > 2 The Power of Interdisciplinary Research

A successful collaboration



Computer Science Duke University



Physics U. of North Carolina

A successful collaboration



Computer Science Duke University



Physics U. of North Carolina

- Eigensolvers
- Tensor contractions
- Tensor calculations

- Development in isolation
- "There is a matrix at the door"
- Expert in numerics, expert in performance
- Disconnect between HPC and users

Accuracy vs. Time

Mixed-precision MR3-based symmetric eigensolver



"Improved Accuracy and Parallelism for MRRR-based Eigensolvers", M. Petschow, P. Bientinesi. SIAM Journal of Scientific Computing, 36(2), 240–263, 2014.

- Eigensolver as part of a process
- Where are the matrices coming from?
- What are the known properties?
- What are the eigenvectors used for?

Eigenvectors



"Correlations in Sequences of Generalized Eigenproblems Arising in Density Functional Theory" E. Di Napoli, S. Bluegel, P. Bientinesi Comp. Phys. Comm., 183(8), 1674–1682, 2012.

A practical Density Functional Theory example

Sequences of eigenvalue problems arising in FLAPW



Sequences of Eigenproblems

Adjacent iteration cycles



Sequences of Eigenproblems

Adjacent iteration cycles



An alternative solving strategy: the ChASE algorithm

Adjacent cycles



Speed-up



The core of the algorithm: Chebyshev filter

In practice

Three-terms recurrence relation

 $C_{m+1}(t) = 2xC_m(t) - C_{m-1}(t); \qquad m \in \mathbb{N}, \quad C_0(t) = 1, \quad C_1(t) = x$

$$Z_m \doteq p_m(\tilde{H}) Z_0$$
 with $\tilde{H} = H - cI_n$

For: $i = 1 \rightarrow \text{Deg} - 1$



END FOR.

Efficiency and portability

 $Au_{98}Ag_{10} - n = 8,970 - 32$ cores.



"An optimized and scalable eigensolver for sequences of eigenvalue problems" M. Berljafa, D. Wortmann and E. Di Napoli. Conc. Comp.: Pract. Exper. 2015; 27:905–922

Tensors

Tensors

- Hot, emerging topic in HPC
- · Lost in translation: tensor vs. multi-dimensional array
- Tensor contraction = generalization of matrix product ⇒ can I use GEMM? Yes/no? When?

$$V_{h_1h_2\ldots}:=S_{i_1i_2\ldots}T_{j_1j_2\ldots}$$

Definition: $\Delta(X) = \#$ of free indices of X

Class 1:
$$\Delta(S) = 0 \land \Delta(T) = 0$$

X BLAS3 X BLAS2 VBLAS1

Class 2:
$$\Delta(S) \ge 1 \land \Delta(T) = 0$$
 or $\Delta(S) = 0 \land \Delta(T) \ge 1$
× BLAS3 × BLAS2 (+ transp) × BLAS1

Class 3:
$$\Delta(S) \ge 1 \land \Delta(T) \ge 1$$

 \checkmark BLAS3 (+ transp) \checkmark BLAS2 \checkmark BLAS1

"Towards an Efficient Use of the BLAS Library for Multilinear Tensor Contractions",

E. Di Napoli, D. Fabregat-Traver, G. Quintana-Orti, P. Bientinesi, Applied Mathematics and Computation, 235, 454–468, 2014.

- Hot, emerging topic in HPC
- · Lost in translation: tensor vs. multi-dimensional array
- Tensor contraction = generalization of matrix product ⇒ can I use GEMM? Yes/no? When?
- Contractions implemented as explicit loops
 Pros: Symmetries, simplifications
 Cons: Efficiency, parallelism, portability!!

- Hot, emerging topic in HPC
- Lost in translation: tensor vs. multi-dimensional array
- Tensor contraction = generalization of matrix product ⇒ can I use GEMM? Yes/no? When?
- Contractions implemented as explicit loops
 Pros: Symmetries, simplifications
 Cons: Efficiency, parallelism, portability!!
- Our approach:
 - Expose properties
 - Map to linear algebra libraries! (and develop tensor libraries)
 - Pros: Inherited efficiency, parallelism, portability
 - Cons: Not all symmetries are exploited



"TTC: A high-performance Compiler for Tensor Transpositions", P. Springer, J.R. Hammond, P. Bientinesi, ACM TOMS, 2017.

"Design of a high-performance GEMM-like Tensor-Tensor Multiplication", P. Springer, P. Bientinesi, ACM TOMS, 2017.

Operatorial form

$$(H)_{G',G} = \sum_{a} \iint \varphi_{G'}^*(\mathbf{r}) \hat{H}_{\mathrm{KS}} \varphi_G(\mathbf{r}) \mathrm{d}\mathbf{r}, \quad (S)_{G',G} = \sum_{a} \iint \varphi_{G'}^*(\mathbf{r}) \varphi_G(\mathbf{r}) \mathrm{d}\mathbf{r}.$$

Entrywise form

$$(S)_{G',G} = \sum_{a} \sum_{L=(l,m)} \left(A_L^{a,G'} \right)^* A_L^{a,G} + \left(B_L^{a,G'} \right)^* B_L^{a,G} \|\dot{u}_{l,a}\|^2$$

$$(H)_{G',G} = \sum_{a} \sum_{L',L} \left(\left(A_{L'}^{a,G'} \right)^* T_{L',L;a}^{[AA]} A_L^{a,G} \right) + \left(\left(A_{L'}^{a,G'} \right)^* T_{L',L;a}^{[AB]} B_L^{a,G} \right) \\ + \left(\left(B_{L'}^{a,G'} \right)^* T_{L',L;a}^{[BA]} A_L^{a,G} \right) + \left(\left(B_{L'}^{a,G'} \right)^* T_{L',L;a}^{[BB]} B_L^{a,G} \right) .$$

Tensor form

$$H = \sum_{a=1}^{N_A} \underbrace{A_a^H T_a^{[AA]} A_a}_{H_{AA}} + \underbrace{A_a^H T_a^{[AB]} B_a + B_a^H T_a^{[BA]} A_a + B_a^H T_a^{[BB]} B_a}_{H_{AB+BA+BB}}$$

$$S = \underbrace{\sum_{a=1}^{N_A} A_a^H A_a}_{S_{AA}} + \underbrace{\sum_{a=1}^{N_A} B_a^H \dot{U}_a^H \dot{U}_a B_a}_{S_{BB}}$$

Constructing SAA

An example of memory layout re-structuring

$$S_{AA} = \sum_{a=1}^{N_A} A_a^H A_a.$$

1: for $a := 1 \rightarrow N_A$ do 2: $S_{AA} = A_a^H A_a$ 3: end for

 \triangleright (zherk: $4N_LN_G^2$ Flops)



1:
$$S_{AA} = A^H_{\star}A_{\star}$$

 \triangleright (zherk: $4N_AN_LN_G^2$ Flops)

Constructing $H_{AB+BA+BB}$

An example of algorithm re-structuring

$$H_{AB+BA+BB} = \sum_{a=1}^{N_A} B_a^H (T_a^{[BA]} A_a) + (A_a^H T_a^{[AB]}) B_a + \frac{1}{2} B_a^H (T_a^{[BB]} B_a) + \frac{1}{2} (B_a^H T_a^{[BB]}) B_a$$
$$= \sum_{a=1}^{N_A} B_a^H (T_a^{[BA]} A_a + \frac{1}{2} T_a^{[BB]} B_a) + (A_a^H T_a^{[AB]} + \frac{1}{2} B_a^H T_a^{[BB]}) B_a$$

1: for $a := 1 \rightarrow N_A$ do 2: $Z_a = T_a^{[BA]}A_a$ 3: $Z_a = Z_a + \frac{1}{2}T_a^{[BB]}B_a$ 4: Stack Z_a to Z_{\star} and B_a to B_{\star} 5: end for 6: $H = Z_{\star}^H B_{\star} + B_{\star}^H Z_{\star}$

▷ (zgemm: $8N_L^2N_G$ Flops) ▷ (zhemm: $8N_L^2N_G$ Flops)

 \triangleright (zher2k: $8N_AN_LN_G^2$ Flops)

	NaCl ($\mathbf{K}_{max} = 4.0$)					
	lv	yBridge		F	laswell	
	HSDLA	FLEUR	×	HSDLA	FLEUR	×
1 core	31.53	48.31	1.53	19.00	47.41	2.50
2 cores	16.10	24.58	1.53	9.98	24.95	2.50
1 CPU	3.90	6.21	1.59	2.25	5.00	2.22
2 CPUs	2.61	5.20	1.99	1.93	4.03	2.09
	TiO_2 (K _{max} = 3.6)					
		TiC	\mathbf{b}_2 (\mathbf{K}_{m}	$_{\rm nax} = 3.6$)		
	lv	TiC yBridge	\mathbf{b}_2 (\mathbf{K}_{rr}	nax = 3.6) ⊢	laswell	
	lv HSDLA	TiC yBridge FLEUR	$\mathbf{b}_2 (\mathbf{K}_m)$	nax = 3.6) ⊢ HSDLA	laswell FLEUR	×
1 core	lv HSDLA 175.53	TiC yBridge FLEUR 256.15	0 ₂ (K _m × 1.46	ax = 3.6) $HSDLA$ 106.56	laswell FLEUR 259.91	× 2.44
1 core 2 cores	lv HSDLA 175.53 86.68	TiC yBridge FLEUR 256.15 127.90	0 ₂ (K _m × 1.46 1.48	hax = 3.6) HSDLA 106.56 53.48	laswell FLEUR 259.91 131.21	× 2.44 2.45
1 core 2 cores 1 CPU	lv HSDLA 175.53 86.68 19.63	TiC yBridge FLEUR 256.15 127.90 29.35	0 ₂ (K _m × 1.46 1.48 1.50	HSDLA 106.56 53.48 10.63	laswell FLEUR 259.91 131.21 25.95	× 2.44 2.45 2.44

Table: Scalability of HSDLA and FLEUR: execution times in minutes on Haswell (12 cores / CPU) and IvyBridge (10 cores / CPU); speedups of HSDLA over FLEUR in **bold**.

"High-performance generation of the Hamiltonian and Overlap matrices in FLAPW methods." E. Di Napoli, E. Peise, M. Hrywniak and P. Bientinesi. Comp. Phys. Comm. 211 (2017) 61–72

Porting to heterogeneous architectures

Back-of-the-envelope analysis

- 5 lines of the algorithm constitute 97% of flops
- Correspond to BLAS-3 operations (gemm, herk, her2k)
- High arithmetic intensity and should fit GPUs well
- First step: offload these routine calls
- All BLAS kernels. Which library?
 - cuBLAS
 - cuBLAS-XT
 - MAGMA
 - BLASX
- 3x wrappers (zgemm, zherk, zher2k)
- Init and cleanup of cuda runtime and devices
- Allocate data in page-locked memory

Only around 100 lines of additional code

Experimental multi-GPU results

Test case: AuAg $(N_A = 108, N_L = 121, N_G = [3275 - 13379])$



Kmax

- CPU: E5-2650, 2 x 8core, 2.0GHz, 64GBs RAM (Sandy Bridge)
- 2 Nvidia Tesla K20X
- Peak performance: 256 GFs/s + 2 x 1.3 TFs/s

"Hybrid CPU-GPU generation of the Hamiltonian and Overlap matrices in FLAPW methods."

D. Fabregat-Traver, D. Davidović, M. Höhnerbach, E. Di Napoli. LNCS, Vol. 10164, pp. 200-211 .

Interdisciplinary research

- Close collaboration & interactions
- Common language!
- New questions, new challenges
- Lots of opportunities

Interdisciplinary research

- Close collaboration & interactions
- Common language!
- New questions, new challenges
- Lots of opportunities

- ... funding schemes ...
- Education

Thank you for your attention!