

Performance Prediction for Tensor Contractions

Paolo Bientinesi, Edoardo Di Napoli, Diego Fabregat, Elmar Peise

AICES, RWTH Aachen
pauldj@aices.rwth-aachen.de

June 3rd, 2014
PASCConference 14
Zürich, Switzerland



Deutsche
Forschungsgemeinschaft
DFG

MATHEMATICS, PHYSICS → multilinear map
multidimensional array + metric

COMPUTER SCIENCE → multidimensional array

MATHEMATICS, PHYSICS → multilinear map
multidimensional array + metric

COMPUTER SCIENCE → multidimensional array

t dimensional tensors → $S_{\alpha\beta\dots\gamma\delta}$, $S_{\underbrace{\alpha \dots \delta}_{t \text{ indices}}}{\beta \gamma}$, $S_{\alpha\beta \dots \gamma\delta}$, ...

MATHEMATICS, PHYSICS → multilinear map
 multidimensional array + metric

COMPUTER SCIENCE → multidimensional array

t dimensional tensors → $S_{\alpha\beta\dots\gamma\delta}$, $S_{\underbrace{\alpha \dots \delta}_{t \text{ indices}} \beta \gamma}$, $S_{\alpha\beta \dots \gamma\delta}$, ...

Operations →

- low dimensional approximations
- contractions

MATHEMATICS, PHYSICS → multilinear map
 multidimensional array + metric

COMPUTER SCIENCE → multidimensional array

t dimensional tensors → $S_{\alpha\beta\dots\gamma\delta}$, $S_{\alpha \underbrace{\dots}_{t \text{ indices}} \gamma \delta}$, $S_{\alpha\beta \dots \gamma\delta}$, ...

Operations →

- low dimensional approximations
- contractions

Examples: $S_{\alpha ij} T_{ij}$, $S_{\alpha i \beta} M_{ik} T_{k\gamma}$, $S_{\alpha ij} M_{ik} T_{kh} M_{hj}$, ...

MATHEMATICS, PHYSICS → multilinear map
 multidimensional array + metric

COMPUTER SCIENCE → multidimensional array

t dimensional tensors → $S_{\alpha\beta\dots\gamma\delta}$, $S_{\alpha}^{\beta\gamma\delta}$, $S_{\alpha\beta}^{\dots\gamma\delta}$, ...
 t indices

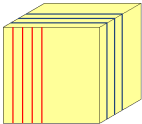
Operations →

- low dimensional approximations
- contractions

Examples: $S_{\alpha ij} T_{ij}$, $S_{\alpha i \beta} M_{ik} T_{k \gamma}$, $S_{\alpha ij} M_{ik} T_{kh} M_{hj}$, ...

Contraction $\rightarrow S_{\alpha ij} T_{j\delta i}$ α, δ **free** indices
 i, j **contracted** indices

$$V_{\alpha\delta} = S_{\alpha ij} T_{j\delta i} \rightarrow \forall \alpha \forall \delta \quad v_{\alpha\delta} = \sum_i \sum_j s_{\alpha ij} t_{j\delta i}$$

Storage $S_{\alpha\beta\gamma\dots} \rightarrow$  $\alpha \rightarrow$ **stride 1**
 $\beta \rightarrow$ **stride $|\alpha|$**
 $\gamma \rightarrow$ **stride $|\alpha||\beta|$**
 \vdots

$$C_{ij} = A_{ik}B_{kj}$$

$$C_{ij} = A_{ik}B_{kj}$$

① **Direct call**

`C := GEMM(A,B)`

$$C_{ij} = A_{ik} B_{kj}$$

0 **Direct call**

`C := GEMM(A,B)`

1 **A is sliced horizontally**

$$C := AB = \begin{bmatrix} a^1 \\ \vdots \\ a^m \end{bmatrix} B = \begin{bmatrix} a^1 B \\ \vdots \\ a^m B \end{bmatrix}$$

for `i=1,...,`
`Ci:=GEMV(Ai,B)`

$$C_{ij} = A_{ik}B_{kj}$$

0 **Direct call**

`C := GEMM(A,B)`

1 **A is sliced horizontally**

$$C := AB = \begin{bmatrix} a^1 \\ \vdots \\ a^m \end{bmatrix} B = \begin{bmatrix} a^1 B \\ \vdots \\ a^m B \end{bmatrix}$$

for `i=1,...,`
`Ci:=GEMV(Ai,B)`

2 **B is sliced vertically**

$$C := AB = A[b_1|b_2|\dots|b_n] = [Ab_1|Ab_2|\dots|Ab_n]$$

for `i=1,...,`
`Ci:=GEMV(A,Bi)`

$$C_{ij} = A_{ik}B_{kj}$$

- 8 **A is sliced vertically and B horizontally**

$$[a_1 | \dots | a_k] \begin{bmatrix} b^1 \\ \vdots \\ b^k \end{bmatrix} = a_1 b^1 + a_2 b^2 + \dots + a_k b^k$$

for $i=1, \dots,$
 $C+=GER(A_i, B_i)$

$$C_{ij} = A_{ik}B_{kj}$$

3 **A is sliced vertically and B horizontally**

$$[a_1 | \dots | a_k] \begin{bmatrix} b^1 \\ \vdots \\ b^k \end{bmatrix} = a_1 b^1 + a_2 b^2 + \dots + a_k b^k$$

for $i=1, \dots,$
 $C += \text{GER}(A_i, B_i)$

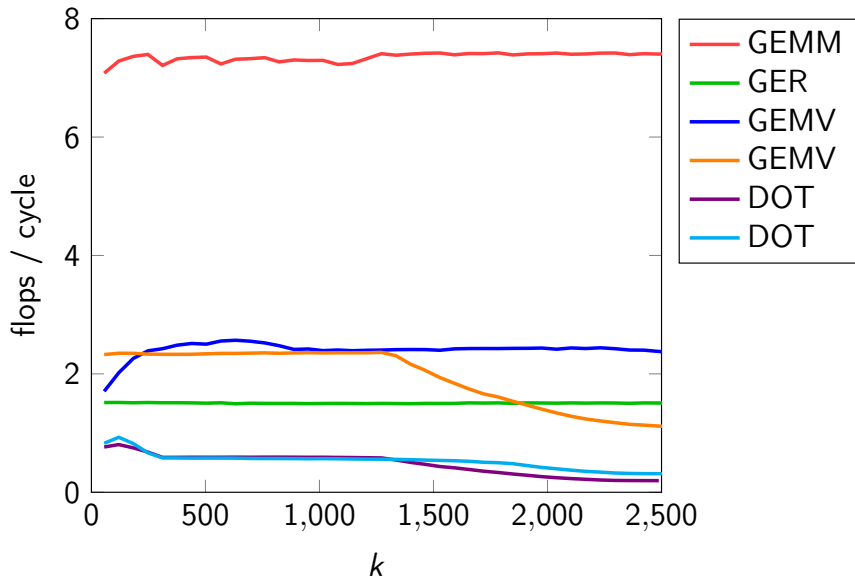
4 **A is sliced horizontally and B vertically**

$$\begin{bmatrix} a^1 \\ \vdots \\ a^m \end{bmatrix} [b_1 | \dots | b_n] = \begin{bmatrix} a^1 b_1 & \dots & a^1 b_n \\ \vdots & \ddots & \vdots \\ a^m b_1 & \dots & a^m b_n \end{bmatrix}$$

for $i=1, \dots,$
 for $j=1, \dots,$
 $C_{ij} := \text{DOT}(A_i, B_j)$

Mathematically equivalent, but . . .

All experiments: OpenBLAS 0.2.8, Intel IvyBridge_EP E5-2680v2



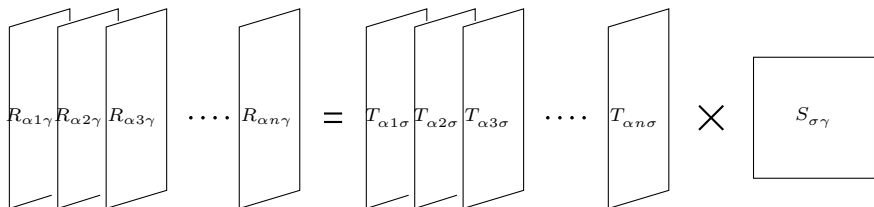
How to use BLAS for contractions?

$$R_{\alpha\beta\gamma} := T_{\alpha\beta\sigma} S_{\sigma\gamma}$$

How to use BLAS for contractions?

$$R_{\alpha\beta\gamma} := T_{\alpha\beta\sigma} S_{\sigma\gamma}$$

1) Slicing along β

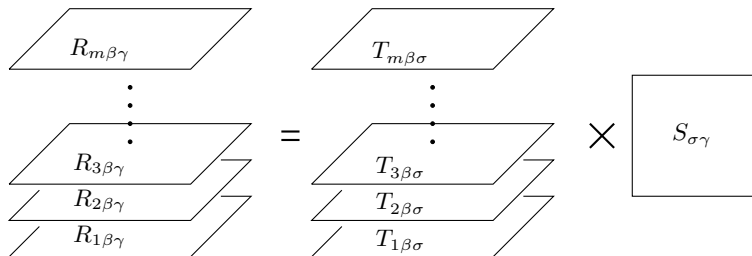


GEMM

How to use BLAS for contractions?

$$R_{\alpha\beta\gamma} := T_{\alpha\beta\sigma} S_{\sigma\gamma}$$

2) Slicing along α



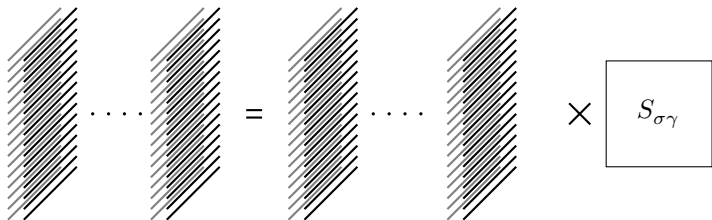
~~X~~ GEMM

✓ Transposition + GEMM

How to use BLAS for contractions?

$$R_{\alpha\beta\gamma} := T_{\alpha\beta\sigma} S_{\sigma\gamma}$$

3) Slicing along α and β



GEMV

$$V_{h_1 h_2 \dots} := S_{i_1 i_2 \dots} T_{j_1 j_2 \dots}$$

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

$$V_{h_1 h_2 \dots} := S_{i_1 i_2 \dots} T_{j_1 j_2 \dots}$$

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

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

 BLAS3

 BLAS2

 BLAS1

$$V_{h_1 h_2 \dots} := S_{i_1 i_2 \dots} T_{j_1 j_2 \dots}$$

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

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


 BLAS3

 BLAS2

 BLAS1

Class 2: $\Delta(S) \geq 1 \wedge \Delta(T) = 0$ or $\Delta(S) = 0 \wedge \Delta(T) \geq 1$

 BLAS3

 BLAS2 (+ transp)

 BLAS1

$$V_{h_1 h_2 \dots} := S_{i_1 i_2 \dots} T_{j_1 j_2 \dots}$$

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

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


 BLAS3

 BLAS2

 BLAS1


Class 2: $\Delta(S) \geq 1 \wedge \Delta(T) = 0$ or $\Delta(S) = 0 \wedge \Delta(T) \geq 1$

 BLAS3

 BLAS2 (+ transp)

 BLAS1

Class 3: $\Delta(S) \geq 1 \wedge \Delta(T) \geq 1$

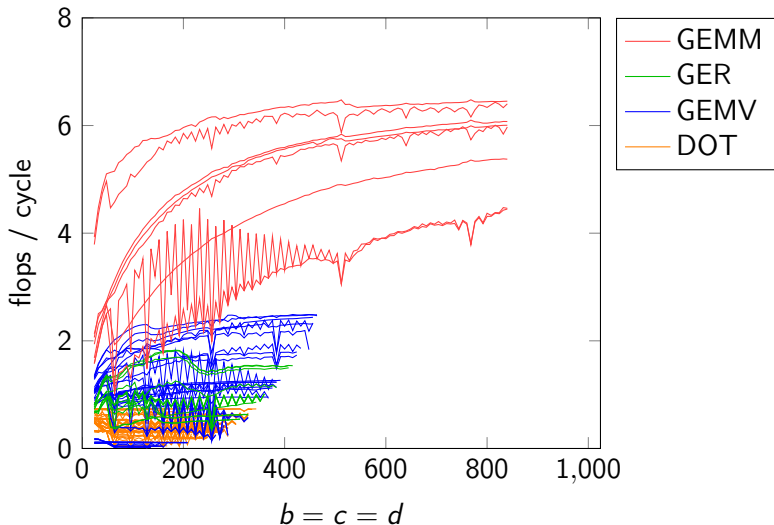
 BLAS3 (+ transp)

 BLAS2

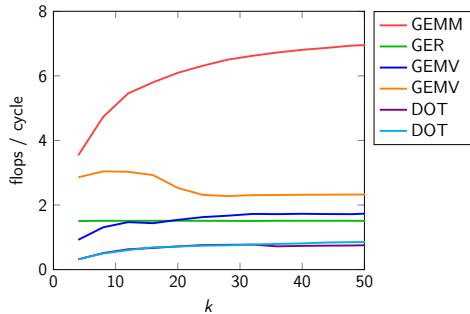
 BLAS1

$$V_{bcd} := S_{ijb} T_{icjd}$$

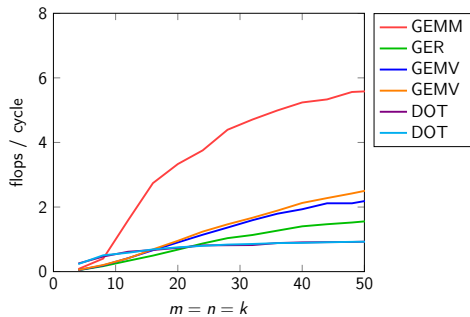
$$V_{bcd} := S_{ijb} T_{icjd}$$



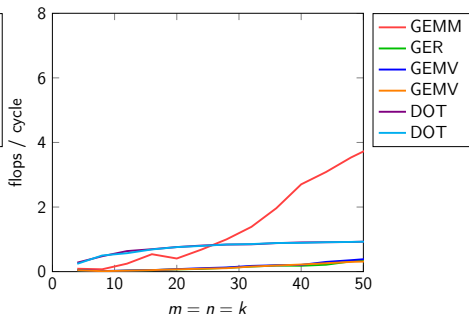
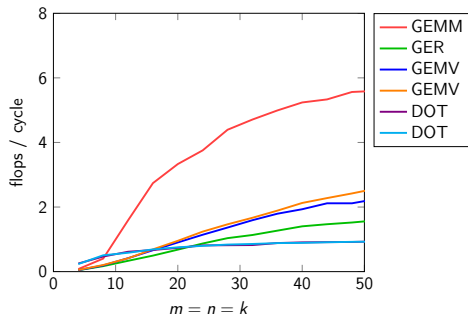
Small dimensions...



Small dimensions...



Small dimensions...



Goal Automatic **selection** of the best variants

Idea Performance prediction

Approach

- Kernels execution ✓
- Algorithms execution ✗

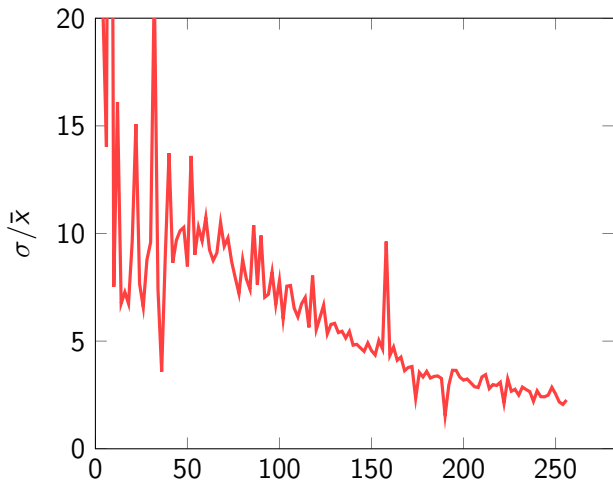
Challenges

Fluctuations – uncertainties
Cache influence

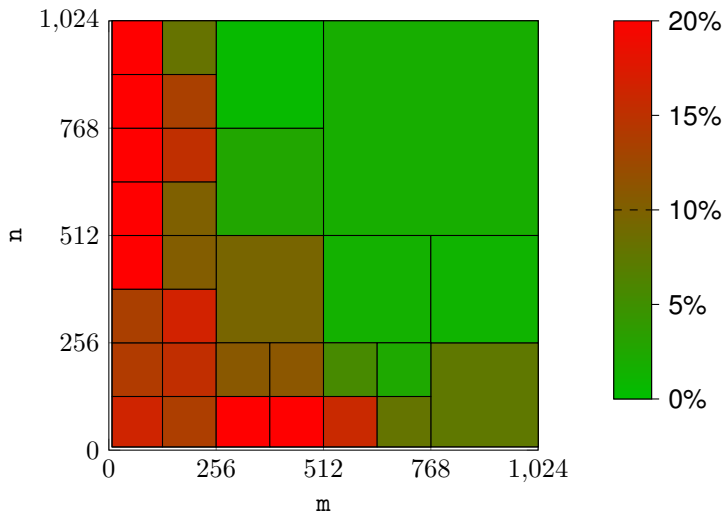
Solution

- Performance models ✗
- Context-aware timings

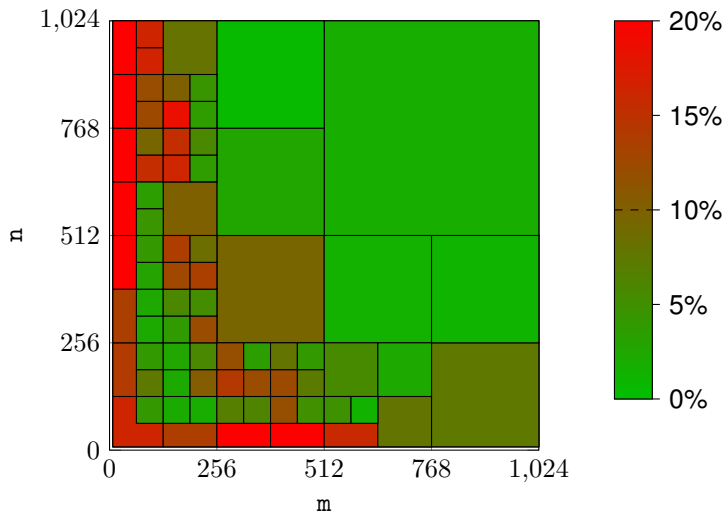
Fluctuations \rightarrow performance models



Fluctuations → performance models



Fluctuations \rightarrow performance models



Observation:

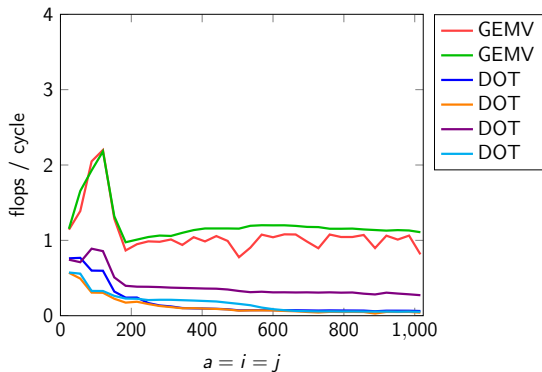
typical linear algebra algorithms → shrinking active region
tensor contractions → identical size slices

Models ... Timings?

Observation:

typical linear algebra algorithms → shrinking active region
tensor contractions → identical size slices

$$V_a := S_{aij}T_{ij}$$

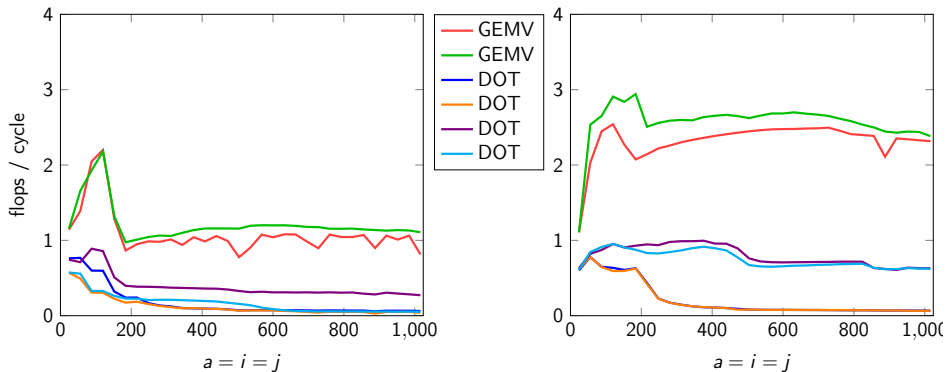


Models ... Timings?

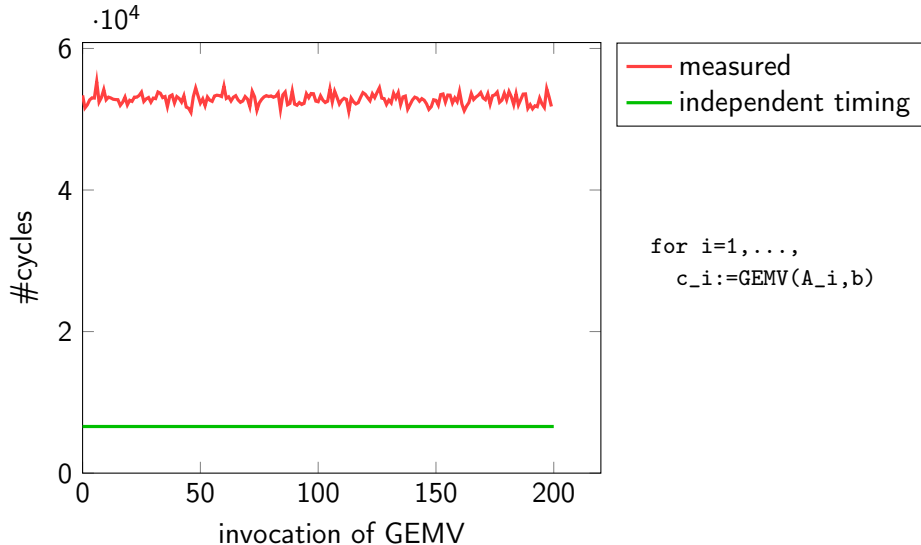
Observation:

typical linear algebra algorithms → shrinking active region
tensor contractions → identical size slices

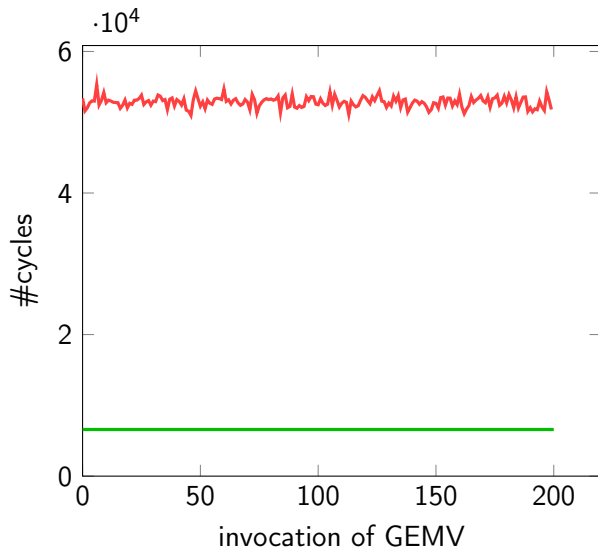
$$V_a := S_{aij}T_{ij}$$



Influence of caching (1/2)



Influence of caching (1/2)

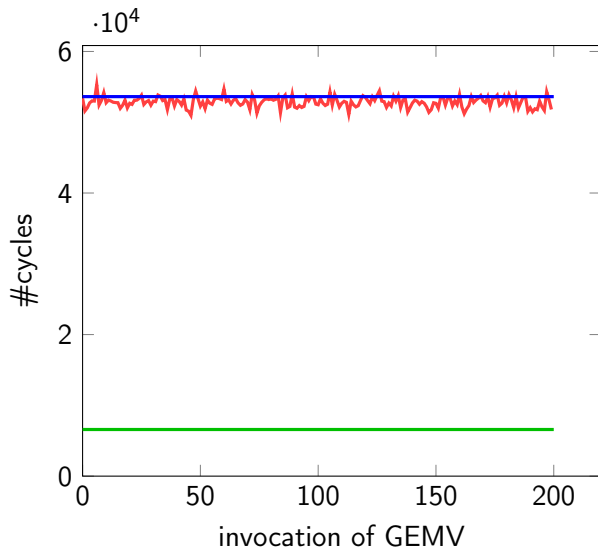


— measured
— independent timing

```
for i=1,...,  
    c_i:=GEMV(A_i,b)
```

Idea: cache setup

Influence of caching (1/2)

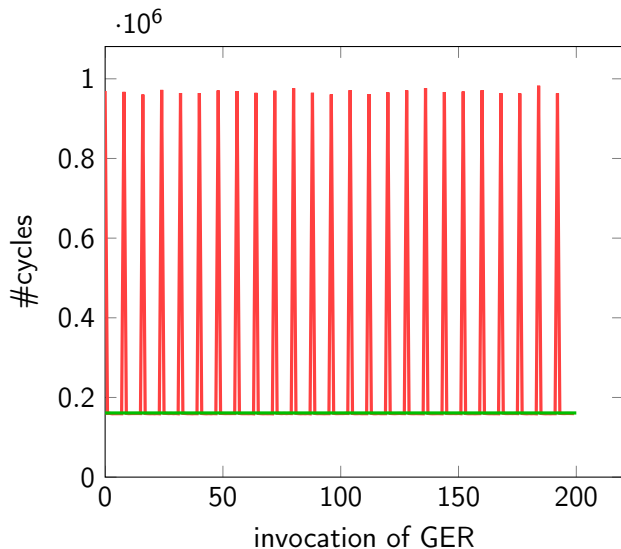


— measured
— independent timing
— cache aware timing

```
for i=1,...,  
    c_i:=GEMV(A_i,b)
```

Idea: cache setup

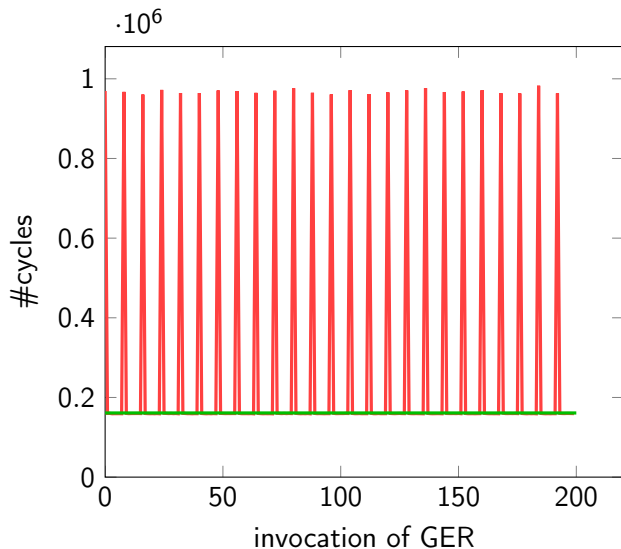
Influence of caching (2/2)



— measured
— cache aware timing

```
for i=1,...,  
  for j=1,...,  
    A_i:=GER(ai,bj)
```


Influence of caching (2/2)

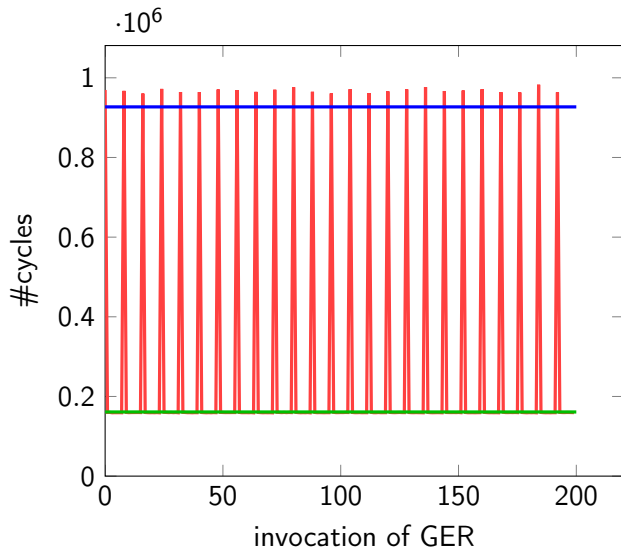


— measured
— cache aware timing

```
for i=1,...,  
  for j=1,...,  
    A_i:=GER(ai,bj)
```

Idea: first iteration

Influence of caching (2/2)



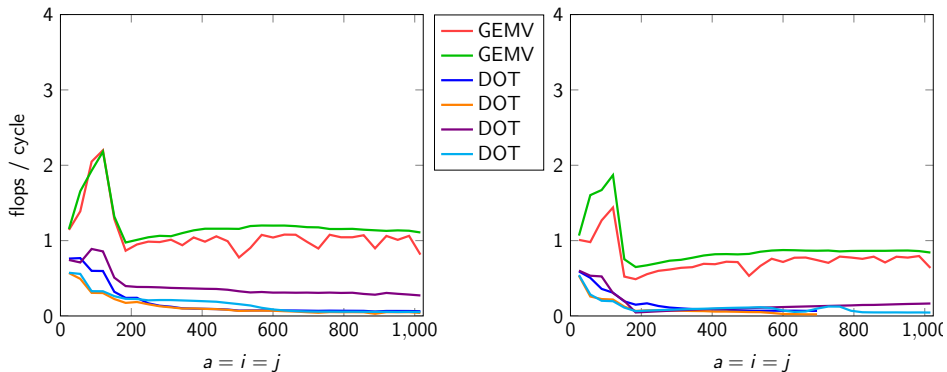
— measured
— cache aware timing
— loop aware timing

```
for i=1,...,  
  for j=1,...,  
    A_i:=GER(ai,bj)
```

Idea: first iteration

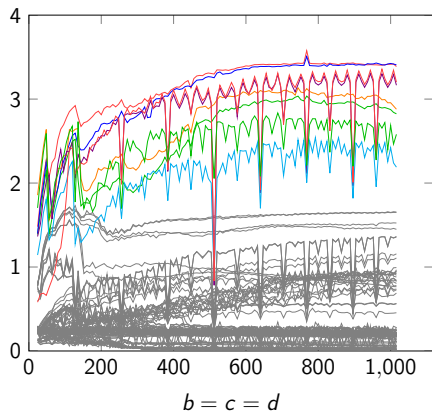
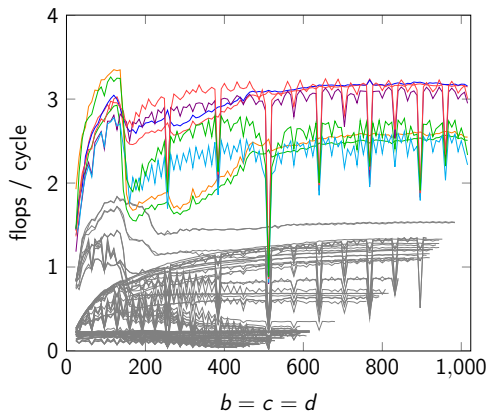
(Nice) results

$$V_a := S_{aij}T_{ij}$$



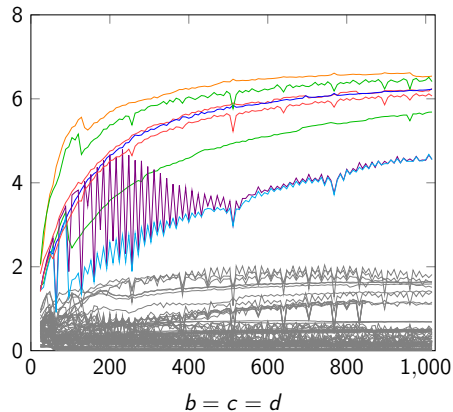
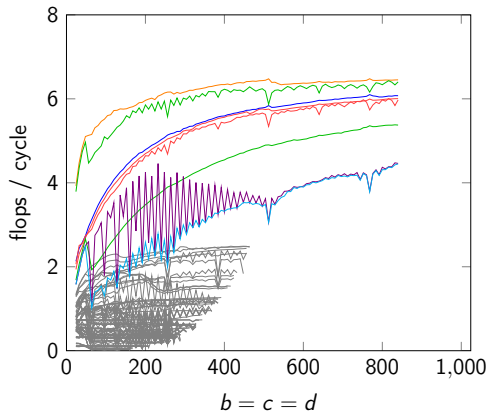
(So so) results

$$V_{bcd} := S_{ijb}T_{icjd}$$



(Awesome) results

$$V_{bcd} := S_{ijb}T_{icjd}$$



Conclusions

Tensor contractions \neq matrix operations
algorithmic space LARGE!
need for BLAS4? maybe
automation? Yes, please!

