

Estimating the Efficiency of BLAS-based Tensor Contractions

2nd Annual Report

Elmar Peise

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

November 6th 2014



Tensor Contractions?

$$\alpha := x^T y \qquad \alpha := \overline{x^T} \mid y \qquad \checkmark \text{ (dot)}$$

$$y := Ax \qquad y := \overline{\boxed{A}} \mid x \qquad \checkmark \text{ (gemv)}$$

$$C := AB \qquad \boxed{C} := \boxed{A} \boxed{B} \qquad \checkmark \text{ (gemm)}$$

$$??? \qquad \text{3D box} := \boxed{\text{2D box}} \text{3D box} \qquad \times$$

Tensor Contractions!

$$C := \overline{A_i} \overline{B_i}$$

$$C := \sum_i A[i] B[i]$$

$$C_a := \boxed{A_{ai}} \overline{B_i}$$

$$C[a] := \sum_i A[a, i] B[i]$$

$$\boxed{C_{ab}} := \boxed{A_{ai}} \boxed{B_{ib}}$$

$$C[a, b] := \sum_i A[a, i] B[i, b]$$

$$\boxed{A_{abc}} := \boxed{B_{ai}} \boxed{B_{ibc}}$$

$$C[a, b, c] := \sum_i A[a, i] B[i, b, c]$$

$$C_{abc} := A_{ija} B_{jbc}$$

$$C[a, b, c] := \sum_{i, j} A[i, j, a] B[j, b, i, c]$$

free indices

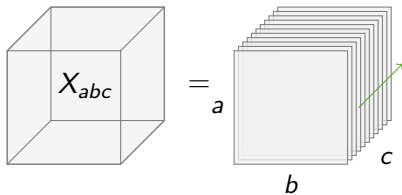
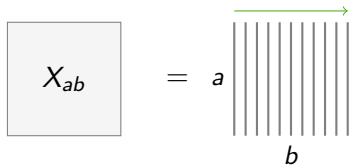
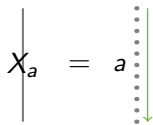
contracted indices

Tensor Storage

- Contiguous
- First index varies the fastest (Fortran style)

Tensor Storage

- Contiguous
- First index varies the fastest (Fortran style)



Outline

① Algorithm Generation

② Performance Prediction

Repeated Execution

Cache Setup

Prefetching

Prefetching Failures

First Iterations

③ Results

$C_{abc} := A_{ai}B_{ibc}$ — Setup 2

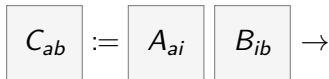
$C_a := A_{iaj}B_{ji}$

$C_{abc} := A_{aij}B_{jbic}$

Multithreading

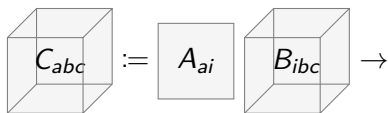
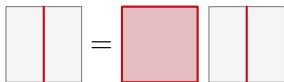
Efficiency

Casting Computation in Terms of BLAS



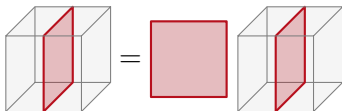
$C_{ab} := A_{ai} B_{ib}$ *b-gemv*

```
for b = 1:b  
  C[:,b] = A[:,:] B[:,b]
```

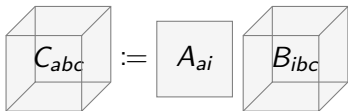


$C_{abc} := A_{ai} B_{ibc}$ *b-gemm*

```
for b = 1:b  
  C[:,b,:] = A[:,:] B[:,b,:]
```

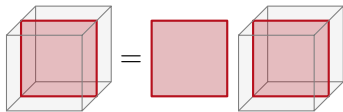


More Examples: $C_{abc} := A_{ai} B_{ibc}$



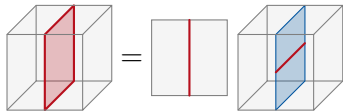
$C_{abc} := A_{ai} B_{ibc}$ *c-gemm*

```
for c = 1:c  
  C[:, :, c] = A[:, :] B[:, :, c]
```



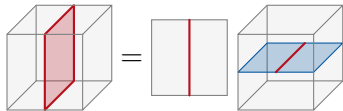
$C_{abc} := A_{ai} B_{ibc}$ *bi-ger*

```
for b = 1:b  
  for i = 1:i  
    C[:, b, :] = A[:, i] B[i, b, :]
```



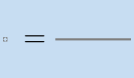
$C_{abc} := A_{ai} B_{ibc}$ *ib-ger*

```
for i = 1:i  
  for b = 1:b  
    C[:, b, :] = A[:, i] B[i, b, :]
```



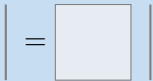
BLAS Kernels

- BLAS-1:

$$C := A_l B_l \quad \text{dot}$$


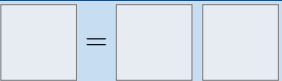
$$C_\alpha := AB_\alpha \quad \text{axpy}$$


- BLAS-2:

$$C_\alpha := A_{\alpha l} B_l \quad \text{gemv}$$


$$C_{\alpha\beta} := A_\alpha B_\beta \quad \text{ger}$$


- BLAS-3:

$$C_{\alpha\beta} := A_{\alpha l} B_{l\beta} \quad \text{gemm}$$


free indices

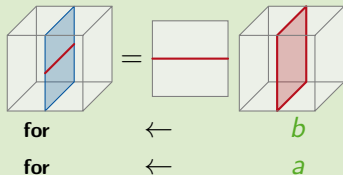
contracted indices

Algorithm Generation

- Select kernel
- Match kernel indices to tensor indices
- Cast remaining tensor indices as for-loops
- Assemble algorithm (AST, C-code)

Example

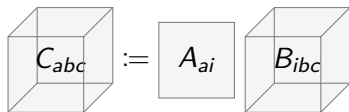
$$\begin{array}{ccc} C_{\alpha} := A_{\alpha l} B_l & \rightarrow & C_{abc} = A_{ai} B_{ibc} \\ l & \rightarrow & i \\ \alpha & \rightarrow & c \end{array}$$



$$C_{abc} := A_{ai} B_{ibc} \quad \text{ba-gemv}$$

```
for b = 1:b
  for a = 1:a
    C[a,b,:] = A[a,:] B[:,b,:]
```

Algorithms for $C_{abc} := A_{ai} B_{ibc}$



36 algorithms:

- 6 dot-based:

abc-dot acb-dot bac-dot bca-dot cab-dot cba-dot

- 18 axpy-based:

ibc-axpy icb-axpy bic-axpy bci-axpy cib-axpy cbi-axpy

iac-axpy ica-axpy aic-axpy aci-axpy cia-axpy cai-axpy

iab-axpy iba-axpy aib-axpy abi-axpy bia-axpy bai-axpy

- 6 gemv-based:

bc-gemv cb-gemv ac-gemv ca-gemv ab-gemv ba-gemv

- 4 ger-based:

ic-ger ci-ger ib-ger bi-ger

- 2 gemm-based:

c-gemm b-gemm

Micro-Benchmarks

```
 $C_{abc} := A_{ai} B_{ibc}$  ba-gemv  
for b = 1:b  
  for a = 1:a  
    C[a,b,:] = A[a,:] B[:,b,:]
```

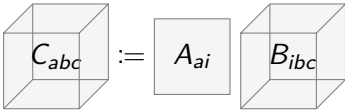
```
Micro-benchmark  
C[a,b,:] = A[a,:] B[:,b,:]
```

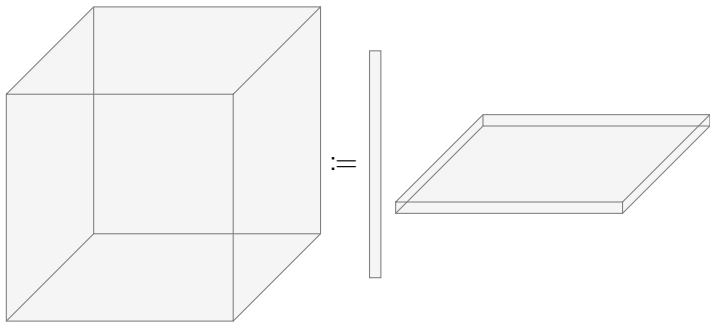
estimate

time 10×

$a \cdot b \cdot (\text{median time})$

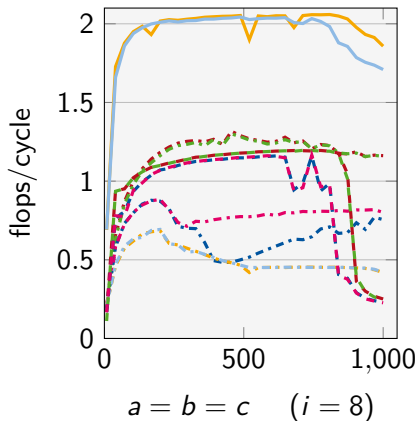
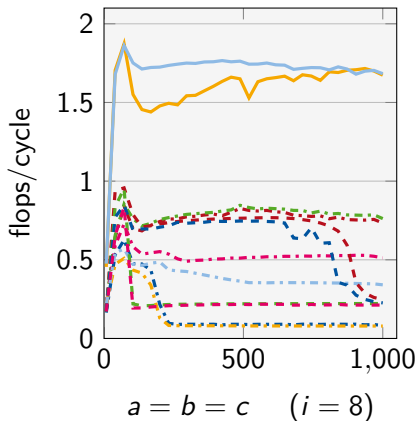
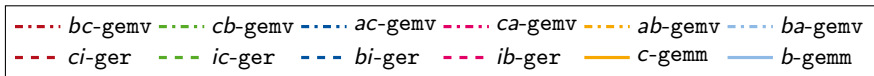
Test Setup

-  $C_{abc} := A_{ai} B_{ibc}$
- $i = 8, a = b = c = 8 \dots 1,000$



- Intel Penryn E5450 (Harpertown)
- Single-threaded OPENBLAS

Repeated Execution



Problem: general overestimation

Problem: general overestimation

Cause: Cache locality not accounted for
(micro-benchmark works in-cache)

Solution Approach

Recreate cache precondition
within micro-benchmark

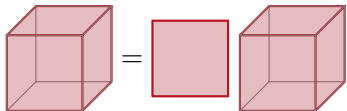
Assumption: Fully associative Least Recently Used (LRU) cache replacement policy

⇒ Cache state is defined by the order of memory accesses.

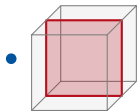
Access Distance

$C_{abc} := A_{ai} B_{ibc}$ *ca-gemv*

```
for c = 1:c
  for a = 1:a
    C[a, :, c] = A[a, :] B[:, :, c]
```



Access Distance $d(M)$: How much data was loaded since M 's last access.



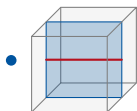
$B[:, :, c]$ **doesn't depend** on for a

$$d(B[:, :, c]) = 0$$



$A[a, :]$ **depends** on for a; **doesn't depend** on for c

$$\begin{aligned} d(A[a, :]) &= \text{size}(\text{join}_a(C[a, :, c], A[a, :], B[:, :, c])) \\ &= \text{size}(C[:, :, c], A[:, :], B[:, :, c]) \\ &= a \cdot b + a \cdot i + i \cdot b \end{aligned}$$



$C[a, :, c]$ **depends** on for a; **depends** on for c

$$\begin{aligned} d(C[a, :, c]) &= \text{size}(\text{join}_{ac}(C[a, :, c], A[a, :], B[:, :, c])) \\ &= \text{size}(C[:, :, :], A[:, :], B[:, :, :]) \\ &= a \cdot b \cdot c + a \cdot i + i \cdot b \cdot c \end{aligned}$$

Setup

Access distances

$$\begin{aligned}d(B[:, :, c]) &= 0 &&= 0 \\d(A[a, :]) &= a \cdot b + a \cdot i + i \cdot b &&= 166,400 \\d(C[a, :, c]) &= a \cdot b \cdot c + a \cdot i + i \cdot b \cdot c &&= 65,283,200\end{aligned}$$

Sizes: $a = b = c = 400$, $i = 8$.

$C_{abc} := A_{ai} B_{ibc}$ *ca-gemv*

```
for c = 1:c
  for a = 1:a
    C[a, :, c] = A[a, :] B[:, :, c]
```

Micro-benchmark

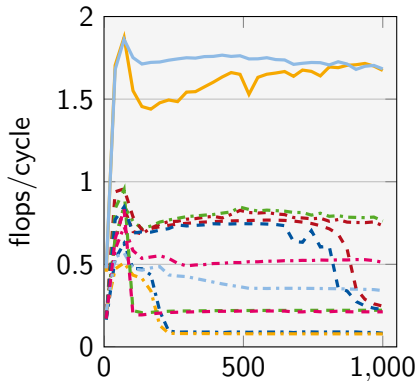
```
touch [816,632]
touch A[a, :]
touch [163,200]
touch B[:, :, c]
C[a, :, c] = A[a, :] B[:, :, c]
```

Limit at $\frac{5}{4} \text{size}(\text{cache}) = \frac{5}{4} \cdot 6\text{MB} = 983,040$ doubles

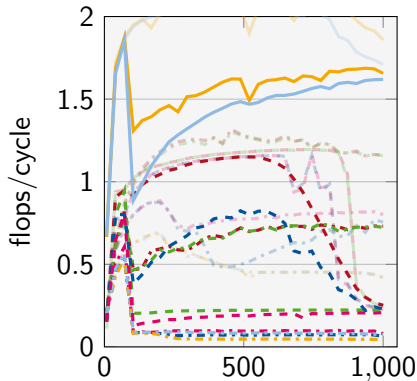
Estimates with Cache Setup

Legend for the plots:

- bc -gemv (red dotted)
- cb -gemv (green dotted)
- ac -gemv (blue dotted)
- ca -gemv (magenta dotted)
- ab -gemv (yellow dotted)
- ba -gemv (light blue dotted)
- ci -ger (red dashed)
- ic -ger (green dashed)
- bi -ger (blue dashed)
- ib -ger (magenta dashed)
- c -gemm (yellow solid)
- b -gemm (light blue solid)



$a = b = c \quad (i = 8)$



$a = b = c \quad (i = 8)$

Problem: underestimation (including ca -gemv (.....))

Prefetching!

Problem: selective underestimation

Cause: Prefetching not accounted for

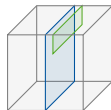
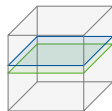
Solution Approach

Imitate prefetching

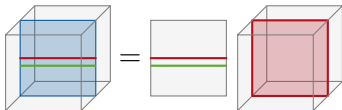
Prefetching:

access

prefetch



Prefetching Imitation



$$C_{abc} := A_{ai} B_{ibc}$$

ca-gemv

```
for c = 1:c
  for a = 1:a
    C[a,:,c] = A[a,:] B[:, :, c]
```

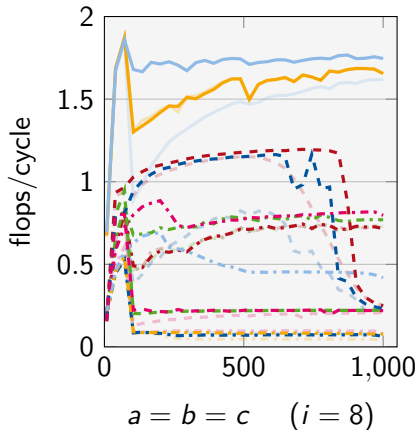
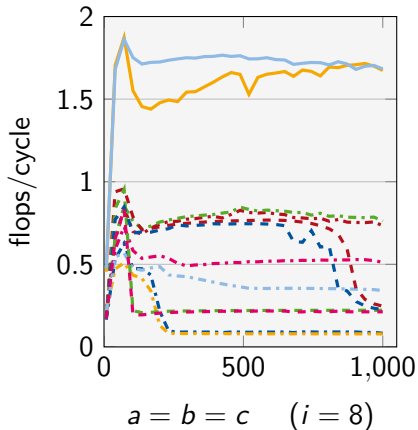
Micro-benchmark

```
C[a,:,c] = A[a,:] B[:, :, c]
```

Estimates with Prefetching

Legend for the plots:

- bc -gemv (red dotted)
- cb -gemv (green dotted)
- ac -gemv (blue dotted)
- ca -gemv (magenta dotted)
- ab -gemv (yellow dotted)
- ba -gemv (light blue dotted)
- ci -ger (red dashed)
- ic -ger (green dashed)
- bi -ger (blue dashed)
- ib -ger (magenta dashed)
- c -gemm (yellow solid)
- b -gemm (light blue solid)



Problem: some incorrect prefetching (including ca -gemv (.....))

Prefetching Failures

Problem: selective prefetching failure

Cause: No Prefetching along 1st dimension across cache-lines.
(every 8th iteration is not prefetched)

Solution Approach

Separate micro-benchmarks with and without prefetching

Micro-benchmark (pre)

```
C[a,:,c] = A[a,:] B[:,:,c]
```

time 10×

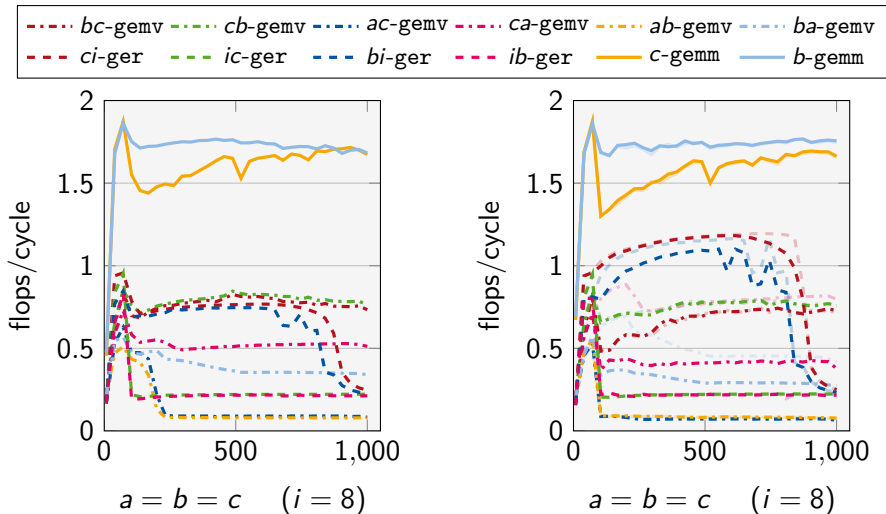
Micro-benchmark (no pre)

```
touch [816,632]  
touch A[a,:]  
touch [163,200]  
touch B[:,:,c]  
C[a,:,c] = A[a,:] B[:,:,c]
```

time 10×

$$\text{estimate} = \frac{1}{8}(7\text{median} + 1\text{median})$$

Estimates with Prefetching Failures



Problem: selective overestimation (including bi -ger (---))

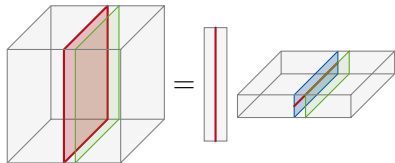
Small Loops' First Iterations

Problem: selective overestimation

Cause: Innermost loop dimension too small
(first iteration of innermost loop differs)

$C_{abc} := A_{ai} B_{ibc}$ *bi-ger*

```
for b = 1:b
  for i = 1:i
    C[:,b,:] = A[:,i] B[i,b,:]
```



Solution Approach

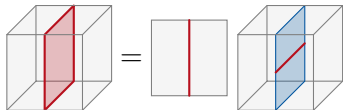
Separate micro-benchmarks for first iteration of small loops

First Iteration Benchmark

- **Access Distance:**

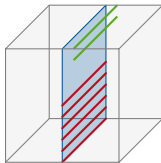
$$C_{abc} := A_{ai} B_{ibc} \quad \text{bi-ger}$$

```
for b = 1:b
  for i = 1:i
    C[:,b,:] = A[:,i] B[i,b,:]
```

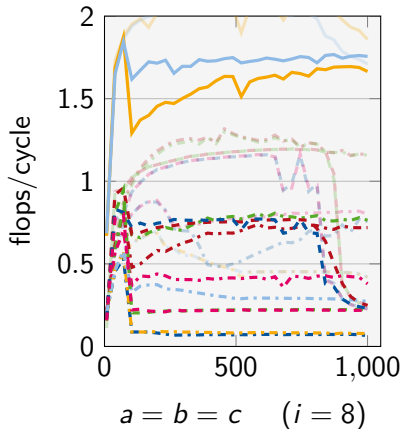
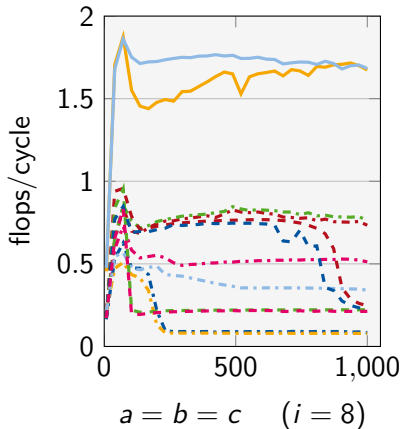
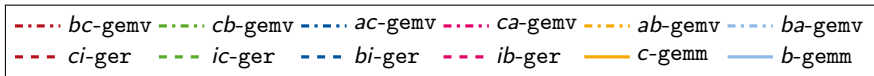


~~Find last access to $A[:,i]$, $B[i,:,c]$, $C[:,:,c]$ within for i~~
Find last access to $A[:,:]$, $B[:,:,c]$, $C[:,:,c]$ within for c

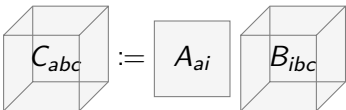
- **Prefetching:**



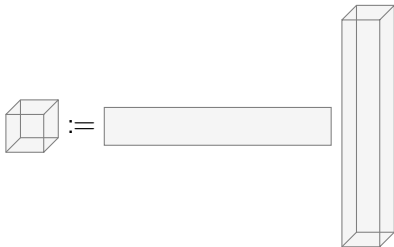
Final Estimates



$C_{abc} := A_{ai} B_{ibc}$ — Setup 2

-  $C_{abc} := A_{ai} B_{ibc}$

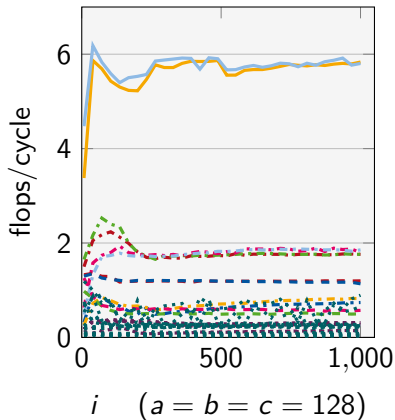
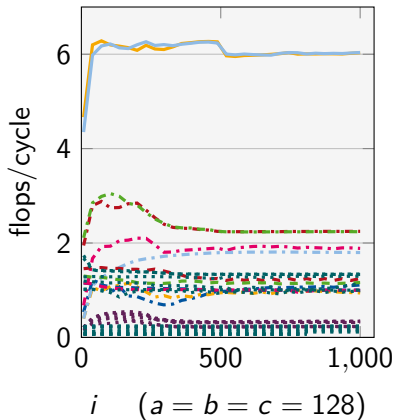
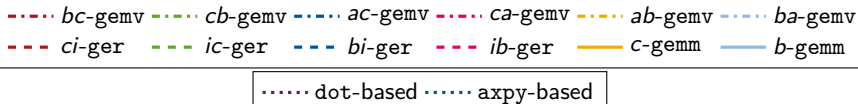
- $a = b = c = 128, i = 8 \dots 1,000$



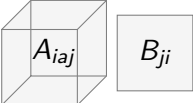
- Intel Ivy Bridge E5-2680 v2
- Single-threaded Intel MKL

$C_{abc} := A_{ai}B_{ibc}$ — Setup 2

Results



$C_a := A_{iaj} B_{ji}$ — Only BLAS-1 and BLAS-2

- $C_a :=$ 

8 algorithms:

- 4 dot-based:

aj-dot *ja-dot* *ai-dot* *ia-dot*

- 2 gemv-based:

$$C_a := A_{iaj} B_{ji} \quad j\text{-gemv}$$

```
for j = 1:j  
  C[:, :] += A[:, :, j] B[j, :]
```

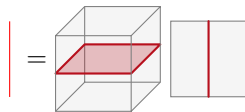
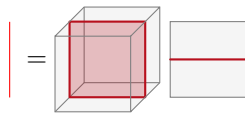
$$C_a := A_{iaj} B_{ji} \quad i'\text{-gemv}$$

```
for i = 1:i  
   $\tilde{A}[:, :] = A[i, :, :]$   
  C[:, :] +=  $\tilde{A}[:, :] B[:, i]$ 
```

- $a = i = j = 8 \dots 1,000$

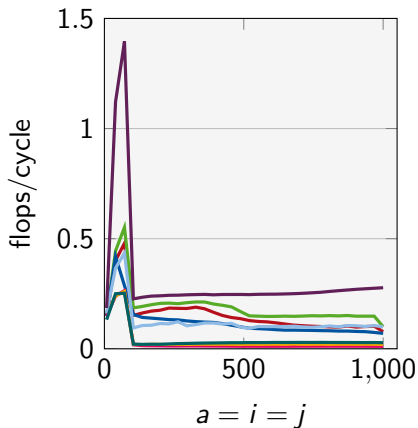
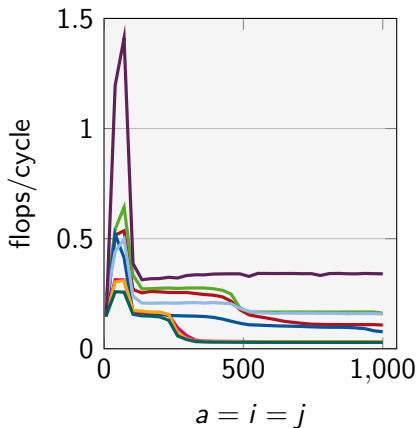
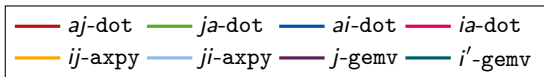
- 2 axpy-based:

ij-axpy *ji-axpy*



$C_a := A_{iaj}B_{ji}$ — Only BLAS-1 and BLAS-2

Results



$C_{abc} := A_{aij}B_{jbic}$ — Challenging Contraction

- $C_{abc} := A_{aij}B_{jbic}$

176 Algorithms:

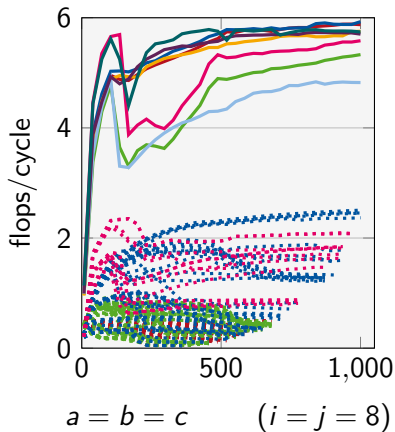
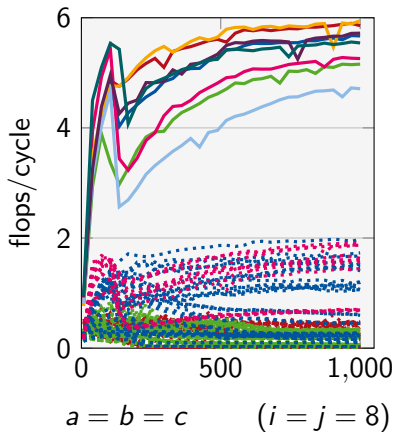
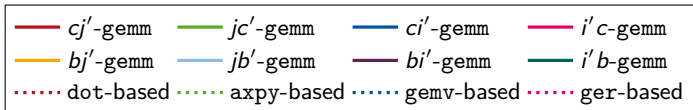
- 48 dot-based
- 72 axpy-based
- 36 gemv-based
- 12 ger-based
- 8 gemm-based:

cj' -gemm	jc' -gemm	ci' -gemm	$i'c$ -gemm
bj' -gemm	jb' -gemm	bi' -gemm	$i'b$ -gemm

- $i = j = 8, a = b = c = 8 \dots 1,000$
- Intel Ivy Bridge E5-2680 v2
- Single-threaded OPENBLAS

$C_{abc} := A_{ajj}B_{jbic}$ — Challenging Contraction

Results

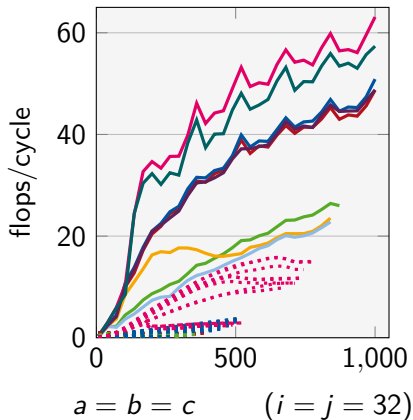
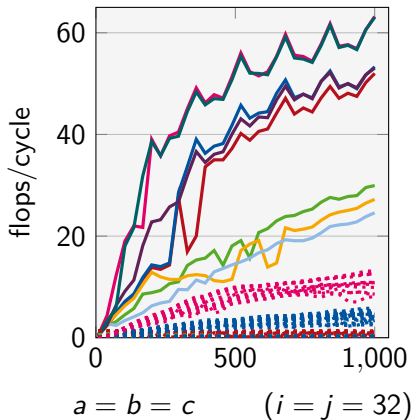
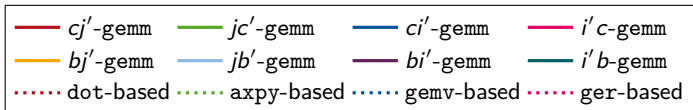


$C_{abc} := A_{aij}B_{jbic}$ — 10 Threads

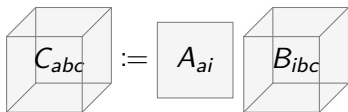
- $C_{abc} := A_{aij}B_{jbic}$
- $i = j = \mathbf{32}$, $a = b = c = 8 \dots 1,000$
- Intel Ivy Bridge E5-2680 v2
- OPENBLAS, 10 threads

$C_{abc} := A_{aij}B_{jbc}$ — 10 Threads

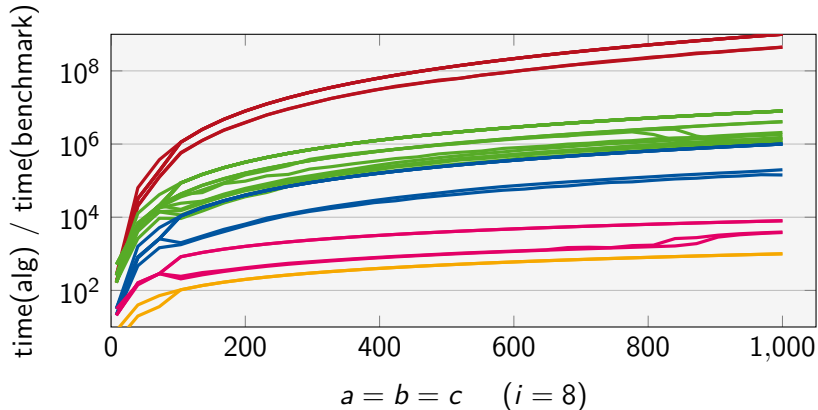
Results



Prediction Efficiency



kernel: — dot — axpy — gemv — ger — gemm



Estimating the Efficiency of BLAS-based Tensor Contractions

- BLAS-based algorithm generation
- Micro-benchmarks with careful cache setup
- Applicable to wide range of challenging scenarios

Funding from Deutsche Forschungsgemeinschaft and Deutsche Telekom Stiftung is gratefully acknowledged.

Deutsche
Forschungsgemeinschaft

DFG

Deutsche
Telekom
Stiftung

