

On the Performance Prediction of BLAS-based Tensor Contractions

Elmar Peise Diego Fabregat-Traver Paolo Bientinesi

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

November 16th 2014 — PMBS14



Tensor Contractions?

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

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

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

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

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

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

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

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

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

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

free indices

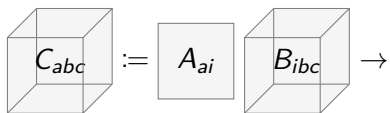
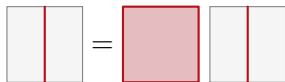
contracted indices

Use BLAS!



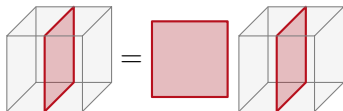
$C_{ab} := A_{ai} B_{ib}$ *b-gemv*

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

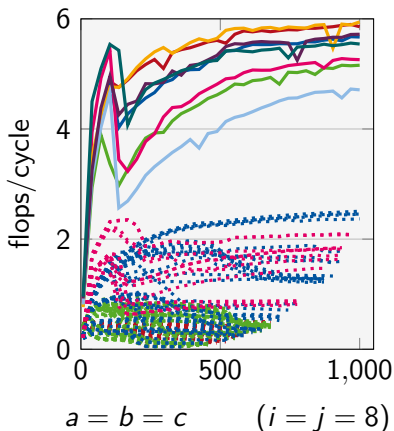
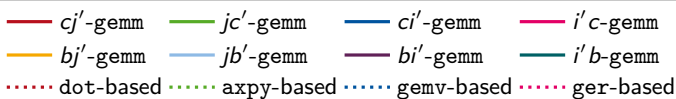


$C_{abc} := A_{ai} B_{ibc}$ *b-gemm*

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



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



Total: 176 Algorithms!

Goals

- Generate algorithms
- Predict their performance

Outline

1 Algorithm Generation

2 Performance Prediction

1. Repeated Execution
2. Cache Setup
3. Prefetching
4. Prefetching Failures
5. First Iterations

3 Results

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

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

Multithreading

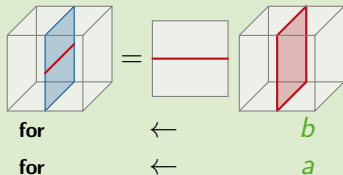
Efficiency

Algorithm Generation

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

Example

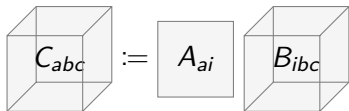
$$\begin{array}{ccc} C_{\alpha} := A_{\alpha l} B_l & \rightarrow & C_{abc} = A_{ai} B_{ibc} \\ l & \rightarrow & i \\ \alpha & \rightarrow & c \end{array}$$



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

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

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



• BLAS-1

- 6 dot-based: $(C := A_l B_l)$

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

- 18 axpy-based: $(C_\alpha := AB_\alpha)$

*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*

• BLAS-2

- 6 gemv-based: $(C_\alpha := A_{\alpha l} B_l)$

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

- 4 ger-based: $(C_{\alpha\beta} := A_\alpha B_\beta)$

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

• BLAS-3

- 2 gemm-based: $(C_{\alpha\beta} := A_{\alpha l} B_{l\beta})$

c-gemm b-gemm

Outline

① Algorithm Generation

② Performance Prediction

1. Repeated Execution
2. Cache Setup
3. Prefetching
4. Prefetching Failures
5. First Iterations

③ Results

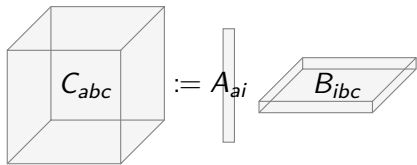
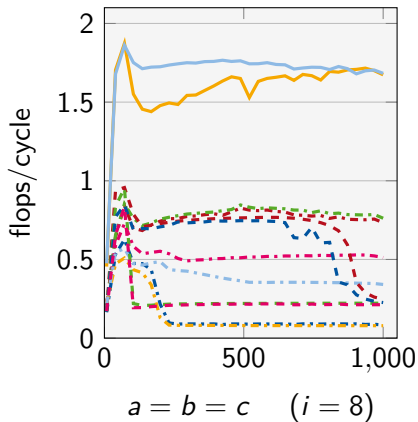
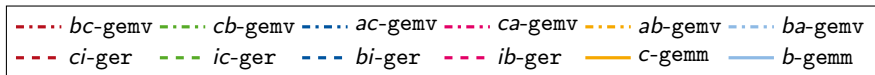
$$C_a := A_{iaj}B_{ji}$$

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

Multithreading

Efficiency

What are we predicting?



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

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,:]
```

estimate

Micro-benchmark

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

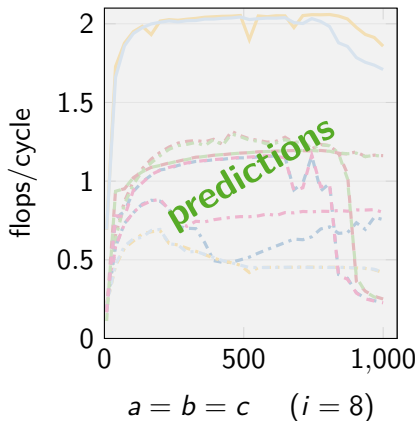
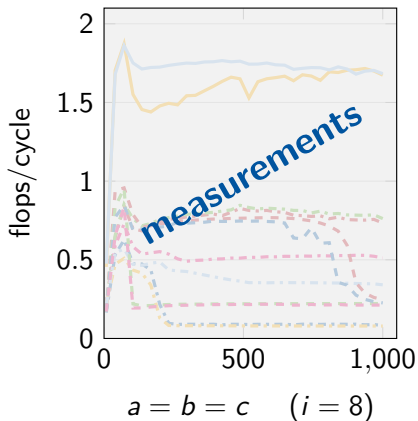
time 10×

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

1. Repeated Execution

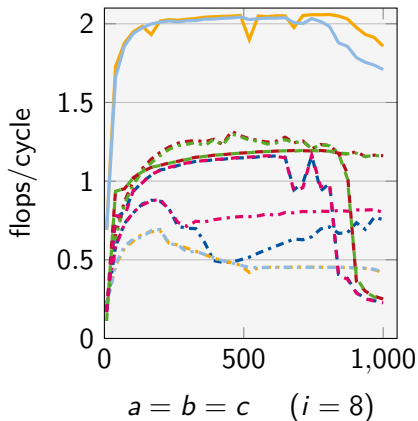
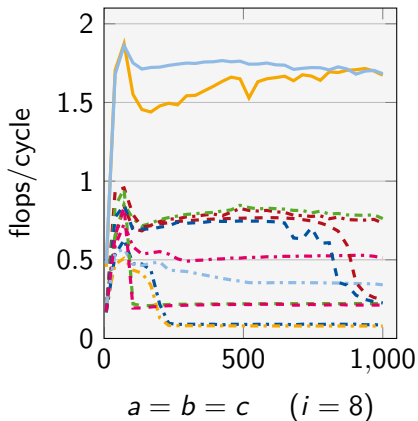
Legend for the plots:

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



1. Repeated Execution

--- bc -gemv -.- cb -gemv -.- ac -gemv -.- ca -gemv -.- ab -gemv -.- ba -gemv
-.- ci -ger -.- ic -ger -.- bi -ger -.- ib -ger — c -gemm — b -gemm



Problem: general overestimation

2. Caching!

Problem: general overestimation

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

Solution Approach

Recreate cache precondition
within micro-benchmark

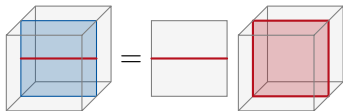
Assumption: Fully associative Least Recently Used (LRU) cache replacement policy

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

Access Distance

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

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

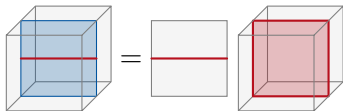


Access Distance $d(M) = \text{size}(\text{all data used since the last access to } M)$

Access Distance

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

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



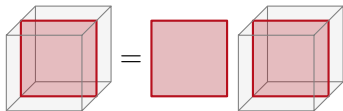
Access Distance $d(M) = \text{size}(\text{all data used since the last access to } M)$

- $B[:, :, c]$ **doesn't vary** in for a
 $d(B[:, :, c]) = 0$ doubles
- $A[a, :]$
- $C[a, :, c]$

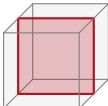

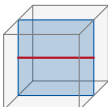
Access Distance

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

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



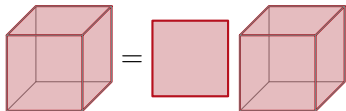
Access Distance $d(M) = \text{size}(\text{all data used since the last access to } M)$

-  $B[:, :, c]$ **doesn't vary** in for a
 $d(B[:, :, c]) = 0$ doubles
-  $A[a, :]$ **varies** in for a; **doesn't vary** in for c
 $d(A[a, :]) = \text{size}(\text{all operands in for a})$
 $= \text{size}(C[:, :, c]) + \text{size}(A[:, :]) + \text{size}(B[:, :, c])$
 $= a \cdot b + a \cdot i + i \cdot b$
-  $C[a, :, c]$

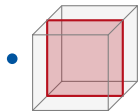
Access Distance

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

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



Access Distance $d(M) = \text{size}(\text{all data used since the last access to } M)$



$B[:, :, c]$ **doesn't vary** in for a

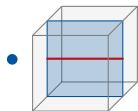
$d(B[:, :, c]) = 0$ doubles



$A[a, :]$ **varies** in for a; **doesn't vary** in for c

$d(A[a, :]) = \text{size}(\text{all operands in for a})$

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



$C[a, :, c]$ **varies** in for a; **varies** in for c

$d(C[a, :, c]) = \text{size}(\text{all operands in for c})$

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

Setup

Access Distances

$d(B[:, :, c])$	$= 0$	$= 0$	doubles
$d(A[a, :])$	$= a \cdot b + a \cdot i + i \cdot b$	$= 166,400$	doubles
$d(C[a, :, c])$	$= a \cdot b \cdot c + a \cdot i + i \cdot b \cdot c$	$= 65,283,200$	doubles

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

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

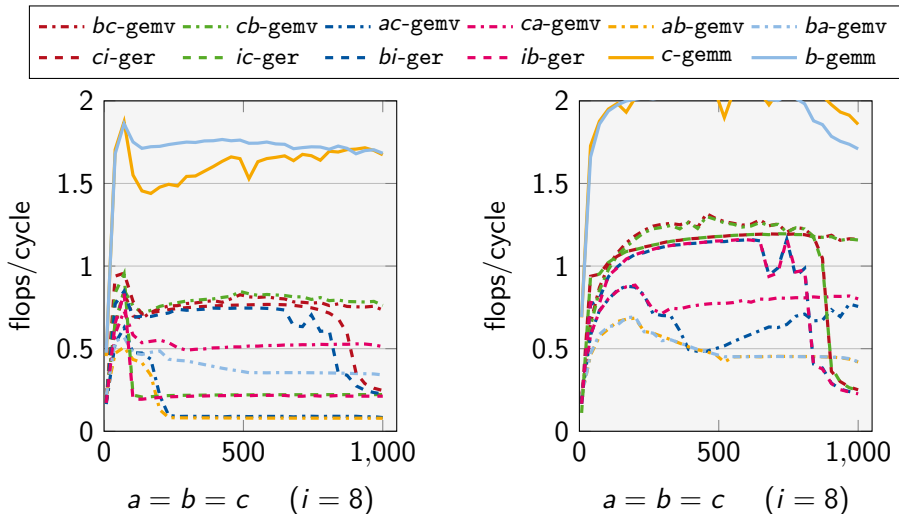
```
for c = 1:c
  for a = 1:a
    C[a, :, c] = A[a, :] B[:, :, c]
```

Micro-benchmark

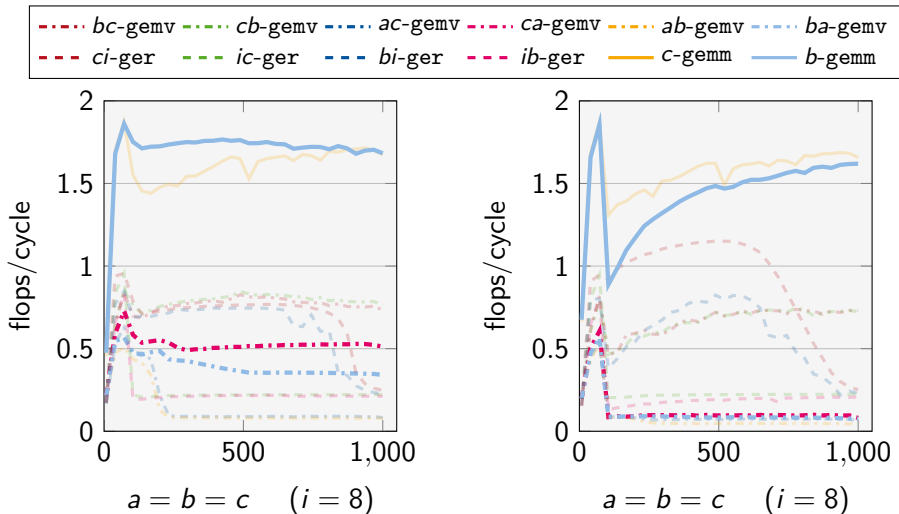
```
flush 816,632
touch A[a, :]
flush 163,200
touch B[:, :, c]
C[a, :, c] = A[a, :] B[:, :, c]
```

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

2. Estimates with Cache Setup



2. Estimates with Cache Setup



Problem: underestimation

3. Prefetching!

Problem: selective underestimation

Cause: Prefetching not accounted for

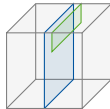
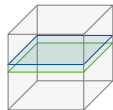
Solution Approach

Mimick prefetching

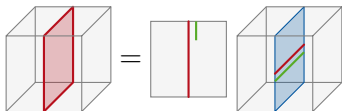
Prefetching:

access

prefetch



Mimicking the Prefetching


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

bi-ger

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

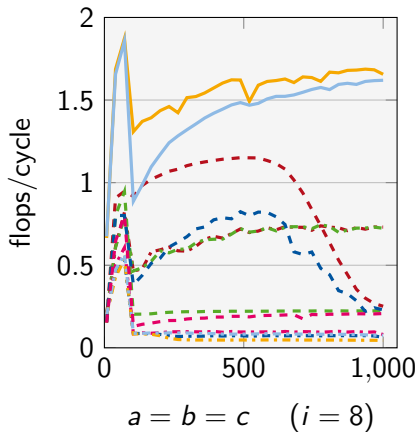
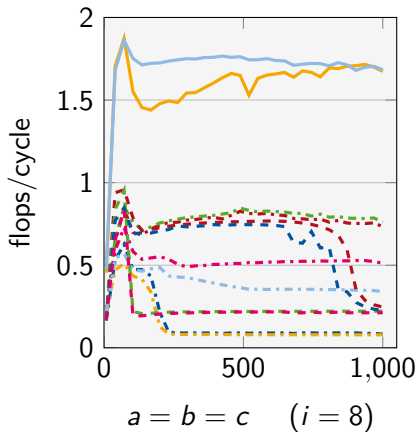
Micro-benchmark

```
touch A[:,i]
flush 5,992
touch A[8,i]
touch B[i,b,:]
touch C[:,b,:]
C[:,b,:] = A[:,i] B[i,b,:]
```

3. Estimates with Prefetching

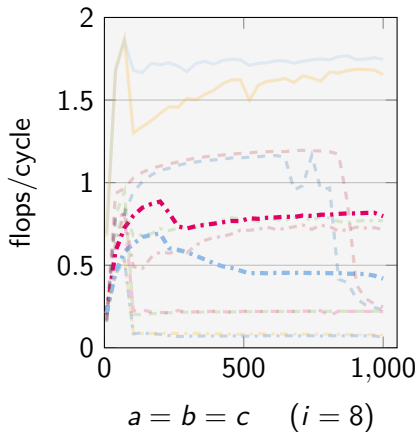
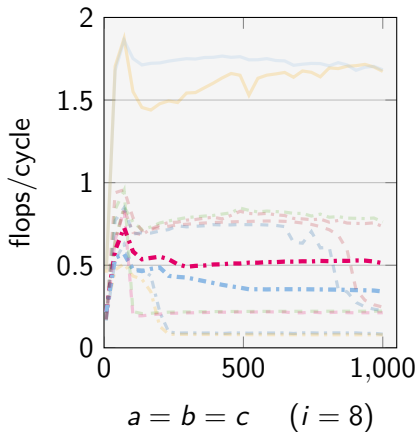
Legend for the plots:

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



3. Estimates with Prefetching

--- *bc-gemv* -.- *cb-gemv* -.- *ac-gemv* -.- *ca-gemv* -.- *ab-gemv* -.- *ba-gemv*
- - *ci-ger* - - *ic-ger* - - *bi-ger* - - *ib-ger* — *c-gemm* — *b-gemm*



Problem: some incorrect prefetching

4. 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)

```
touch A[:,i]
flush 5,992
touch A[8,i]
touch B[i,b,:]
touch C[:,b,:]
C[:,b,:] = A[:,i] B[i,b,:]
```

time $10\times$

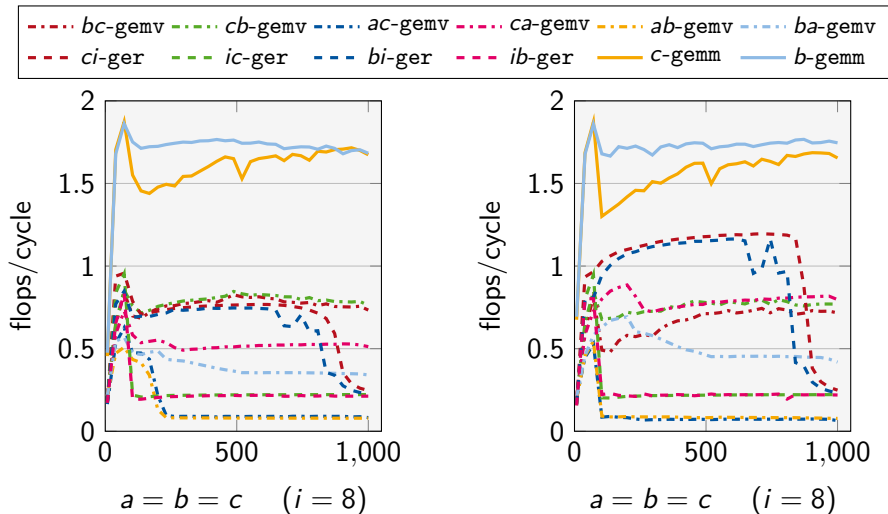
Micro-benchmark (no pre)

```
flush 816,240
touch A[:,i]
flush 5,992
touch A[8,i]
touch C[:,b,:]
C[:,b,:] = A[:,i] B[i,b,:]
```

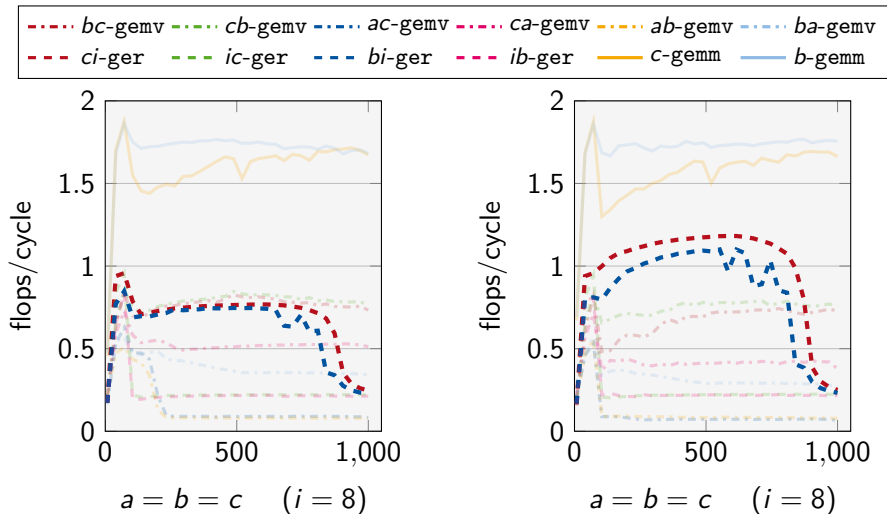
time $10\times$

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

4. Estimates with Prefetching Failures



4. Estimates with Prefetching Failures



Problem: selective overestimation

5. Small Loops' First Iterations

Problem: selective overestimation

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

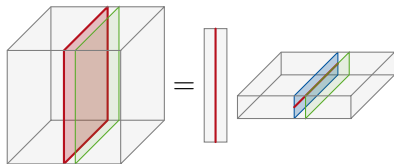
5. Small Loops' First Iterations

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

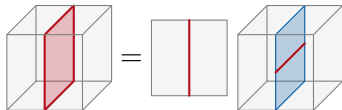
Separate micro-benchmarks for first iteration of small loops

First Iteration Benchmark

- Access Distance:

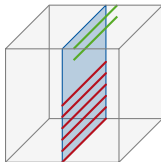
$$C_{abc} := A_{ai} B_{ibc} \quad \text{bi-ger}$$

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

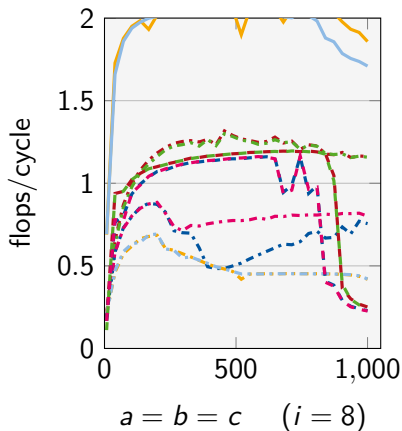
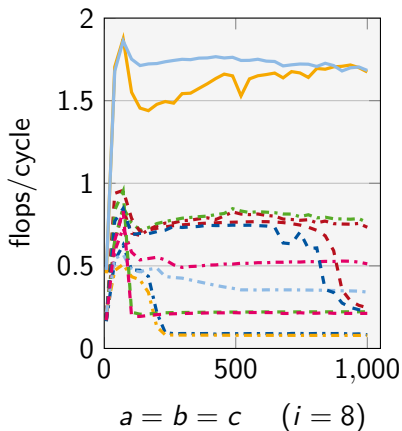
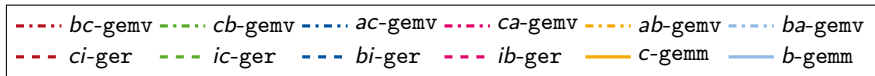


~~Find last access to $A[:,i]$, $B[i,b,:]$, and $C[:,b,:]$ within for i~~
Find last access to $A[:,i]$, $B[i,b,:]$, and $C[:,b,:]$ within for c

- Prefetching:



5. Final Estimates



Outline

① Algorithm Generation

② Performance Prediction

1. Repeated Execution
2. Cache Setup
3. Prefetching
4. Prefetching Failures
5. First Iterations

③ Results

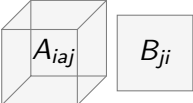
$$C_a := A_{iaj} B_{ji}$$

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

Multithreading

Efficiency

$C_a := A_{iaj} B_{ji}$ — 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, :]
```

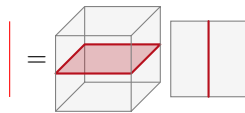
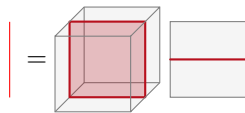
$$C_a := A_{iaj} B_{ji} \quad i'\text{-gemv}$$

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

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

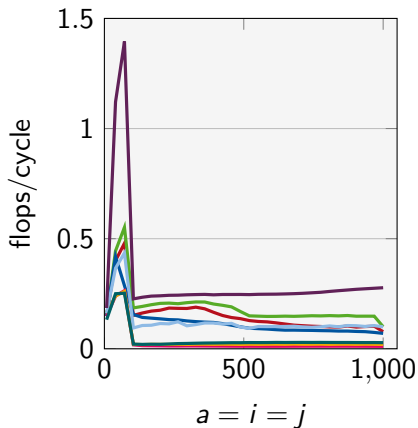
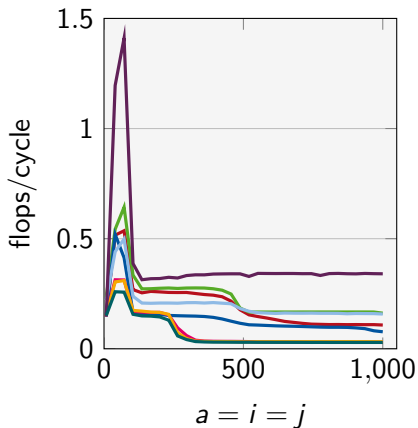
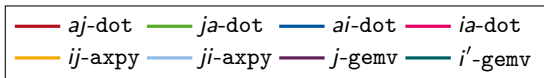
- 2 axpy-based:

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:

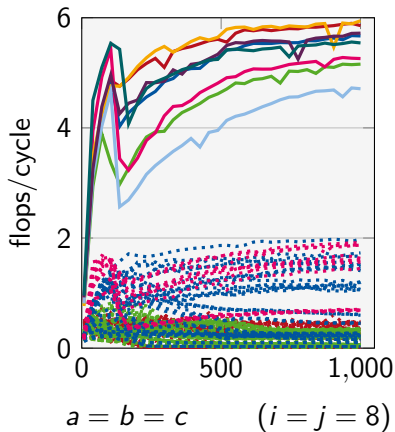
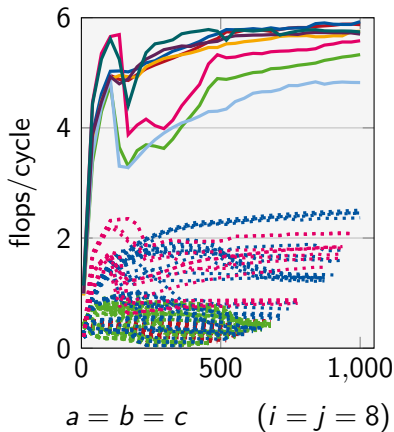
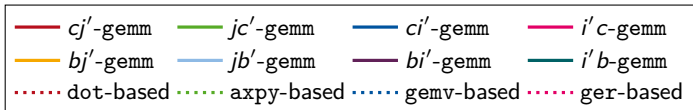
- 48 dot-based
- 72 axpy-based
- 36 gemv-based
- 12 ger-based
- 8 gemm-based:

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

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

Results

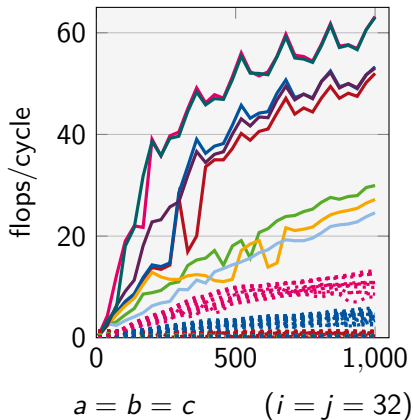
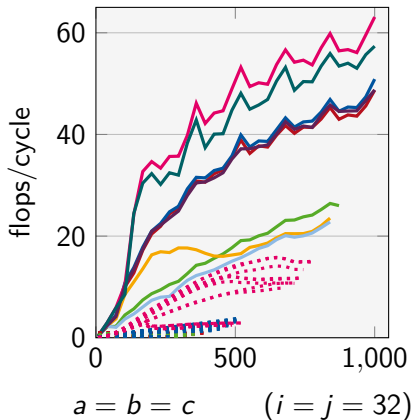
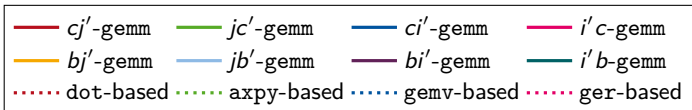


$C_{abc} := A_{aij}B_{jbic}$ — 10 Threads

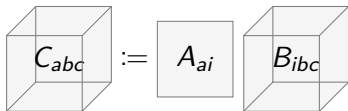
- $C_{abc} := A_{aij}B_{jbic}$
- $i = j = \mathbf{32}$, $a = b = c = 8 \dots 1,000$
- Intel Ivy Bridge E5-2680 v2
- OPENBLAS, 10 threads

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

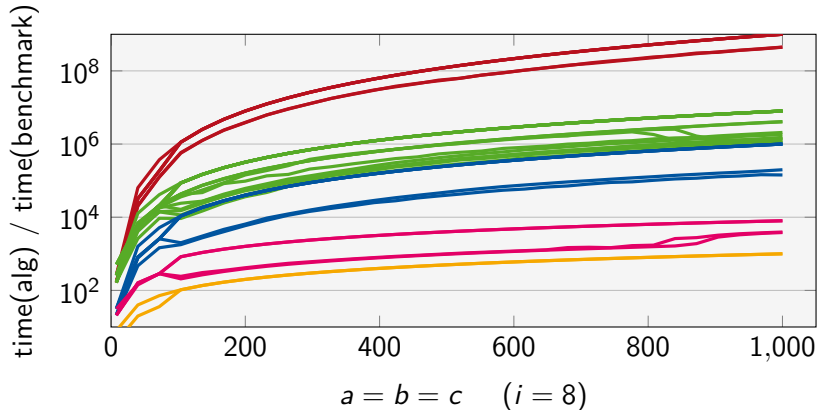
Results



Prediction Efficiency



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



On the Performance Prediction 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

