

GEMM-like Tensor-Tensor Contraction (GETT)

Paul Springer and Paolo Bientinesi

Introduction

Tensor contractions (TC) are a performance critical component in numerous scientific computations. Despite the close connection between matrix-matrix products (GEMM) and TCs, the performance of the latter is in general vastly inferior to that of an optimized GEMM. To close such a gap, we propose a novel approach: **GEMM-like Tensor-Tensor multiplication (GETT)** [1].

Tensor Contraction

- Tensor contractions are the generalization of matrix-matrix multiplications to higher dimensions.

Let \mathcal{A} , \mathcal{B} , and \mathcal{C} be multi-dimensional tensors, a tensor contraction is:

$$\mathcal{C}_{l_c} \leftarrow \alpha \times \mathcal{A}_{l_A} \times \mathcal{B}_{l_B} + \beta \times \mathcal{C}_{l_c},$$

where l_A , l_B , and l_C denote symbolic index sets. We further

define three important index sets:

- Free indices of \mathcal{A} , $l_m := l_A \cap l_c$
- Free indices of \mathcal{B} , $l_n := l_B \cap l_c$
- Contracted indices, $l_k := l_A \cap l_B$.

All indices of l_k are implicitly summed.

Example: Matrix-Matrix Multiplication.

$$C_{mn} \leftarrow A_{mk} \times B_{kn}$$

with $l_m = \{m\}$, $l_n = \{n\}$, and $l_k = \{k\}$.

Example: Tensor Contraction.

$$C_{m_1 n_1 m_2 n_2} \leftarrow A_{m_2 k_1 m_1 k_2} \times B_{n_1 n_2 k_2 k_1}$$

with $l_m = \{m_1, m_2\}$, $l_n = \{n_1, n_2\}$, and $l_k = \{k_1, k_2\}$.

Prior Approaches

- Transpose-Transpose-GEMM-Transpose (TTGT)**
 - Flatten tensors into matrices
 - Contract via GEMM
 - Fold result back into a tensor
 - Adopters: Cyclops Tensor Framework, Tensor Toolbox, LibTensor
- Loop-over-GEMM (LoG)**
 - Identify 2D slices of the tensors
 - Contract all slices via GEMM

Both approaches exhibit significant disadvantages:

- Transpose-Transpose-GEMM-Transpose**
 - Suboptimal I/O cost
 - Flattening accounts for pure overhead
 - Auxiliary memory for flattened tensors
- Loop-over-GEMM (LoG)**
 - Suboptimal I/O cost
 - May require strided GEMM

Classification

Tensor contractions can be mapped to BLAS kernels:

- GEMM: $l_m \neq \emptyset$ and $l_n \neq \emptyset$ and $l_k \neq \emptyset$
- GEMV: $(l_m = \emptyset$ or $l_n = \emptyset)$ and $l_k \neq \emptyset$
- GER: $l_m \neq \emptyset$ and $l_n \neq \emptyset$ and $l_k = \emptyset$
- AXPY: $(l_m = \emptyset$ or $l_n = \emptyset)$ and $l_k = \emptyset$
- DOT: else.

GETT's Key Features

GETT constitutes a *native* tensor contraction algorithm that does not rely on an available GEMM implementation. The key features are:

- GEMM-like design
- Multi-dimensional packing routines
 - Packing via tensor transpositions [2]
 - Preserve stride-1 index
 - Explicitly vectorized

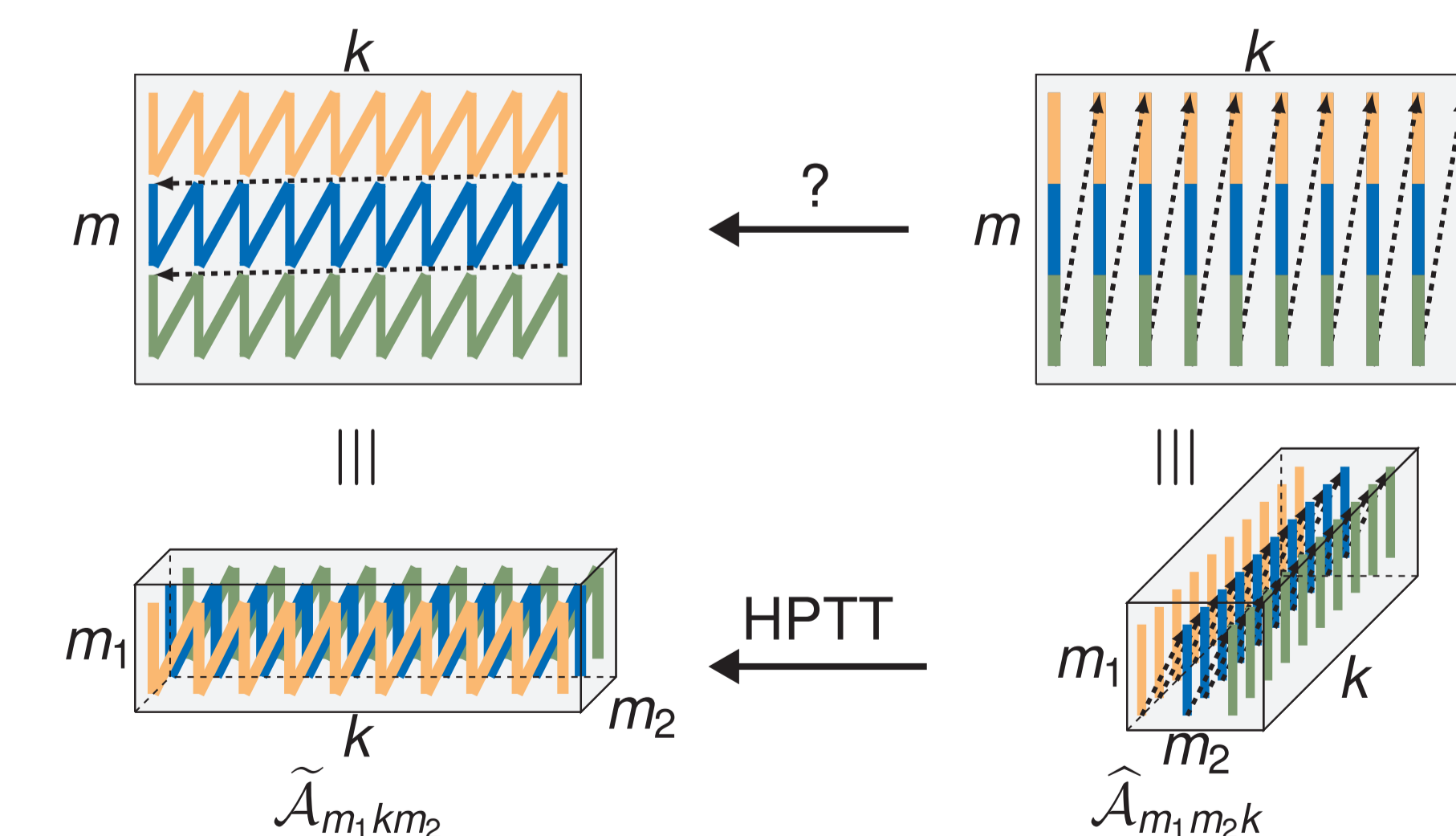
Design

- Similar to a high-performance GEMM
- Complexity offloaded into packing routines

```

// N-Loop
for n = 1 : n_C : S_n
    // K-Loop (contracted)
    for k = 1 : k_C : S_k
        B = identify_subtensor(B, n, k)
        // pack B into B-tilde
        B-tilde = packB(B)
        // M-Loop
        for m = 1 : m_C : S_m
            A = identify_subtensor(A, m, k)
            // pack A into A-tilde
            A-tilde = packA(A)
            C = identify_subtensor(C, m, n)
            // matrix-matrix product: C-tilde ← αA-tilde × B-tilde + βC
            macroKernel(A-tilde, B-tilde, C, α, β)
    end
end
    
```

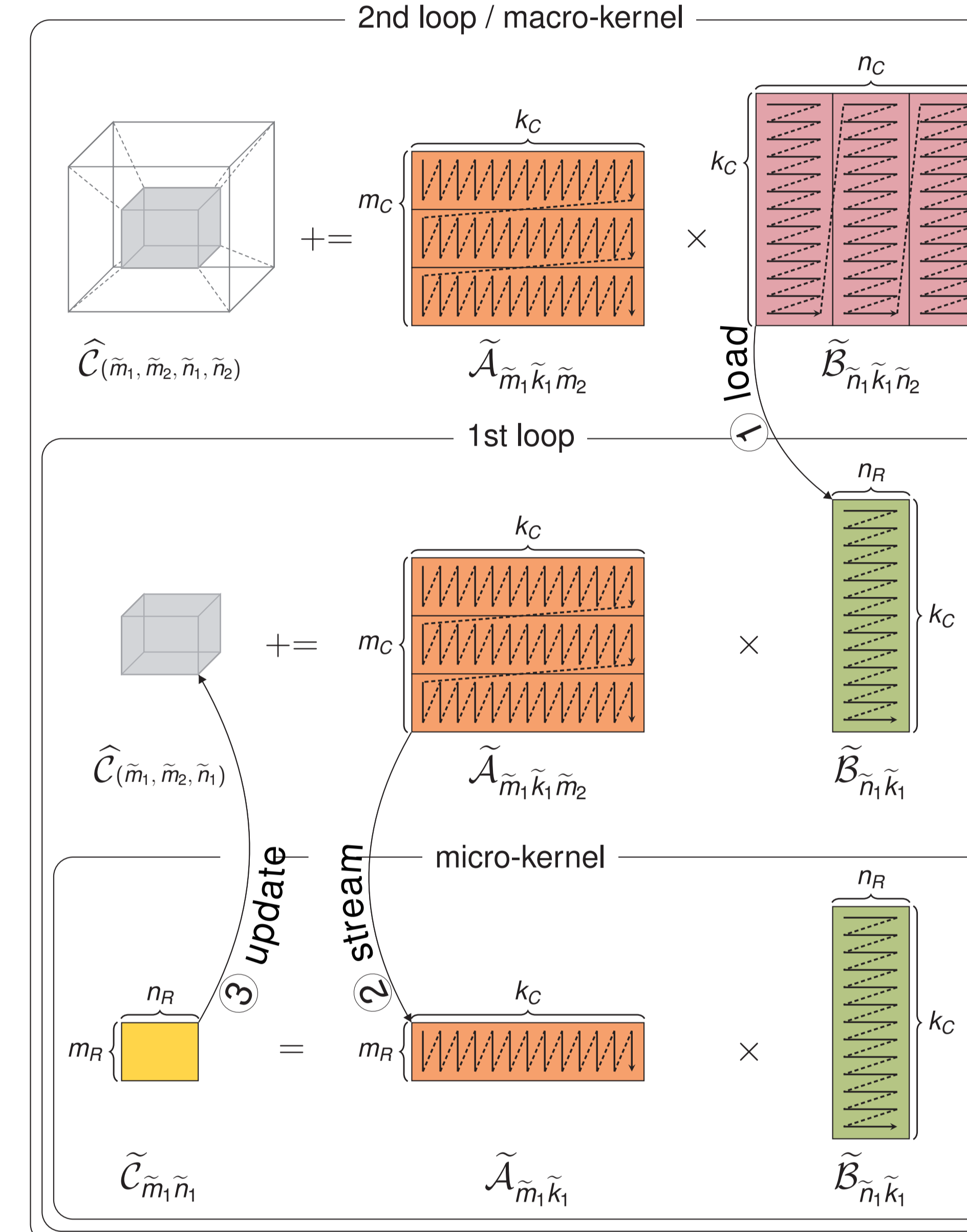
Packing Routines



Macro- & Micro-Kernel

- Blocking for caches
- Blocking for registers
- Explicitly vectorized

main memory L3 L2 L1 registers



Shared-Memory Parallelism

- Parallelism is similar to that of BLIS [4]
 - All loops can be parallelized individually
- GETT also parallelizes along the contracted dimensions
 - Exposes additional parallelism if C is very small
- A promising *parallelism strategy* is chosen according to a cost model:
 - Maximize load balancing
 - Stay within the cache limitations

Future Work

- Implement GETT as part of a C++ library

Performance

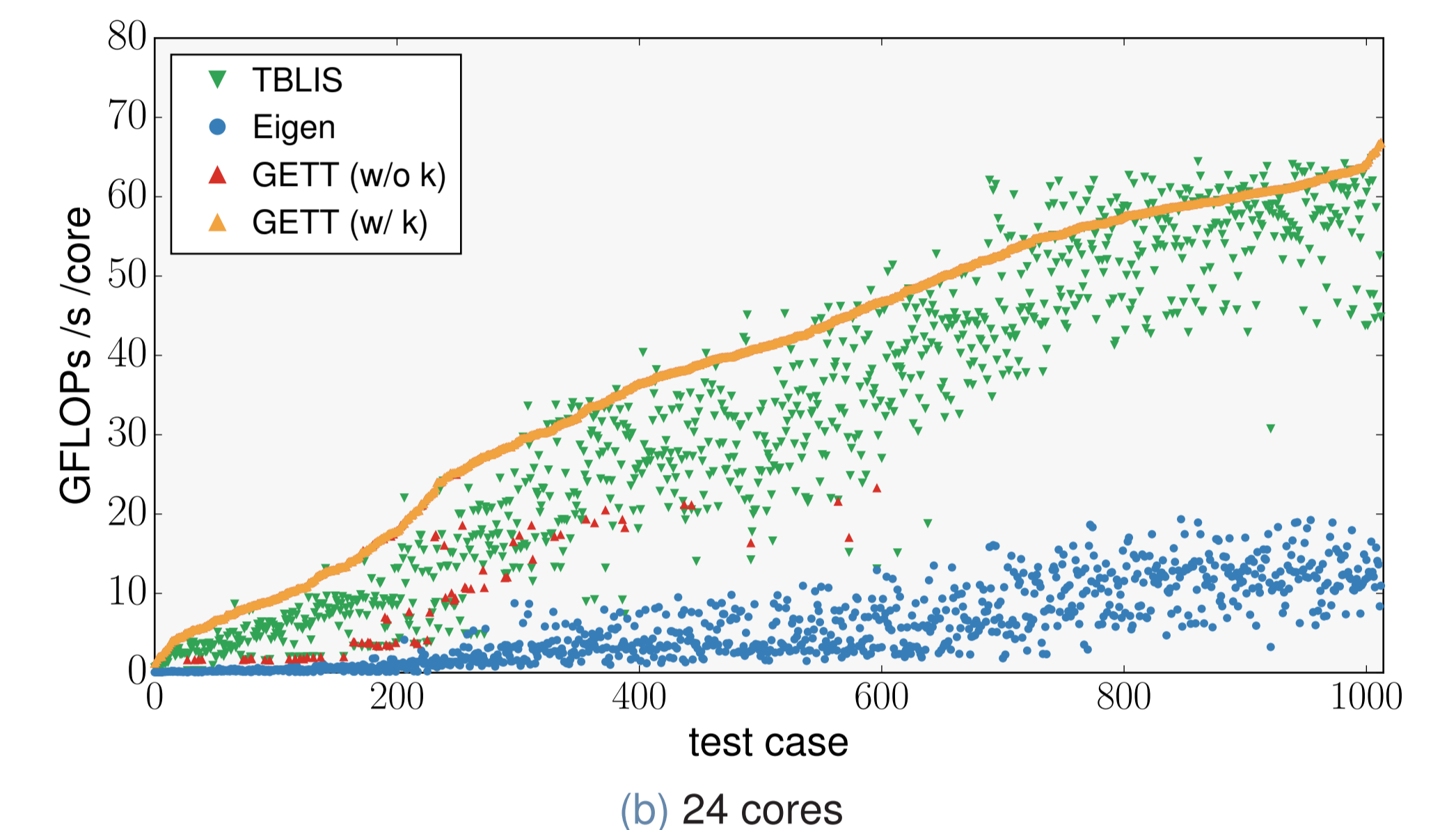
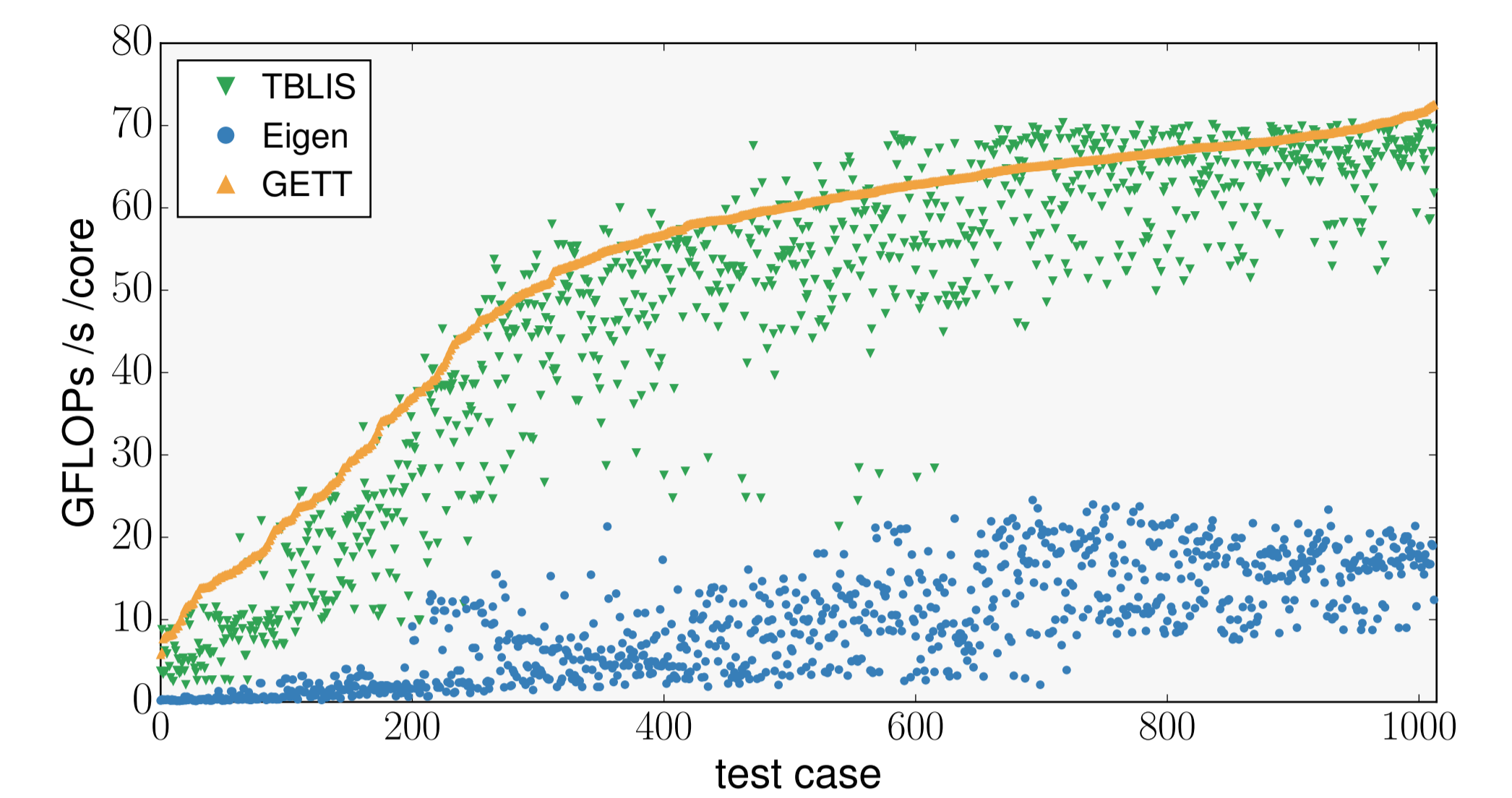


Figure: Performance across random TCs with varying arithmetic intensity (sorted w.r.t. GETT). Host: 2x Intel Xeon E5-2680 v3.

	Single-Threaded	Multi-Threaded
TBLIS [3]	1.3x	1.6x
Eigen [5]	13.2x	16.2x
Direct [1]	-	52.3x

Table: GETT's average speedups.

References

- P. Springer and P. Bientinesi. Design of a high-performance GEMM-like Tensor-Tensor Multiplication. 2016. <http://arxiv.org/abs/1607.00145>.
- P. Springer, T. Su, and P. Bientinesi. HPTT: A High-Performance Tensor Transposition C++ Library. *ACM Array Workshop*, 2017. DOI: 10.1145/3091966.3091968.
- Devin A. Matthews. High-Performance Tensor Contraction without BLAS. 2016. <http://arxiv.org/abs/1607.00291>.
- Tyler M. Smith et al. Anatomy of High-Performance Many-Threaded Matrix Multiplication. *IPDPS'14*. DOI: 10.1109/IPDPS.2014.110
- Gaël Guennebaud and Benoît Jacob et al. Eigen v3.3.4. <http://eigen.tuxfamily.org>



Paul Springer
Telefon: +49 241 80 99 141
Schinkelstr. 2, 52062 Aachen, GERMANY
E-Mail: springer@ices.rwth-aachen.de



GitHub

