

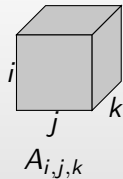
# TTC: A Compiler for Tensor Transpositions

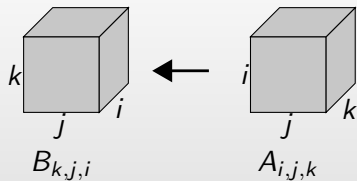
**Paul Springer** and Paolo Bientinesi

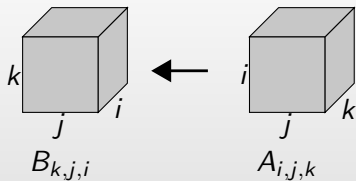
Aachen Institute for Advanced Study in  
Computational Engineering Science

Paris, 14.04.16 — SIAMPP16

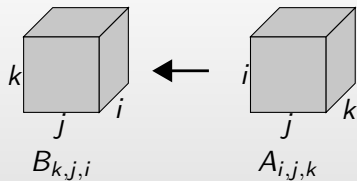




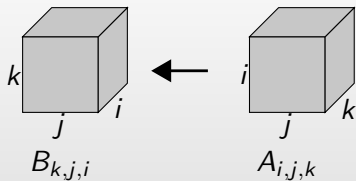




- Building block for tensor contractions
  - Map tensor contraction to GEMM



- Building block for tensor contractions
  - Map tensor contraction to GEMM
- Packing routines for dense linear algebra kernels
  - e.g., GEMM



- Building block for tensor contractions
  - Map tensor contraction to GEMM
- Packing routines for dense linear algebra kernels
  - e.g., GEMM
- Huge search space of viable implementations
  - Implementations vary greatly in performance

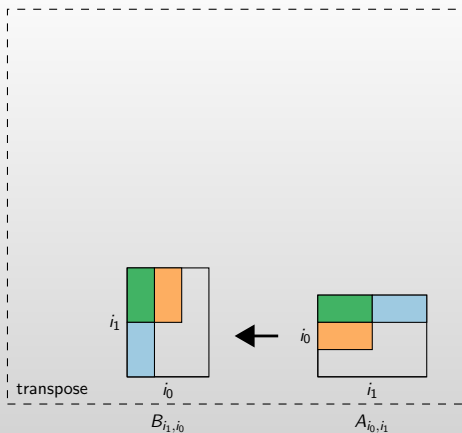
- Domain-specific compiler
  - Multidimensional tensor transpositions

- Domain-specific compiler
  - Multidimensional tensor transpositions
- Generates high-performance C++/CUDA C code
  - Parallelized
  - Vectorized

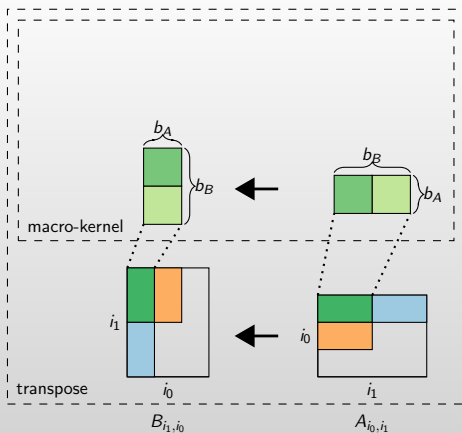


- Domain-specific compiler
  - Multidimensional tensor transpositions
- Generates high-performance C++/CUDA C code
  - Parallelized
  - Vectorized
- Support for multiple architectures
  - NVIDIA GPUs, KNC (Xeon Phi), AVX-enabled CPUs

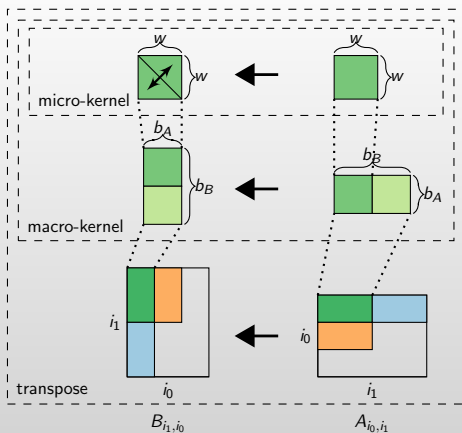
- Domain-specific compiler
  - Multidimensional tensor transpositions
- Generates high-performance C++/CUDA C code
  - Parallelized
  - Vectorized
- Support for multiple architectures
  - NVIDIA GPUs, KNC (Xeon Phi), AVX-enabled CPUs
- Support for common numerical data types
  - single, double, single-complex, double-complex
  - mixed precision



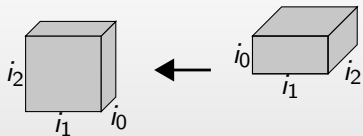
- Decompose a transposition into macro-tiles
  - Parallel over macro-tiles



- Decompose a transposition into macro-tiles
  - Parallel over macro-tiles
- Decompose a macro-tile into micro-tiles
  - Vectorized micro-tiles

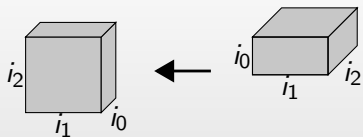


- Decompose a transposition into macro-tiles
  - Parallel over macro-tiles
- Decompose a macro-tile into micro-tiles
  - Vectorized micro-tiles



$B_{i_2, i_1, i_0}$

$A_{i_0, i_1, i_2}$

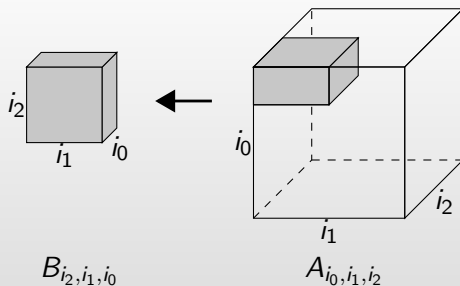


$B_{i_2, i_1, i_0}$

$A_{i_0, i_1, i_2}$

- Input to TTC:

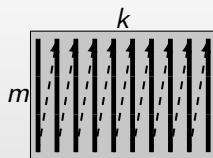
```
ttc --perm=2,1,0 --size=8,16,16
```

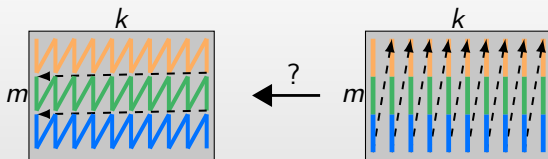


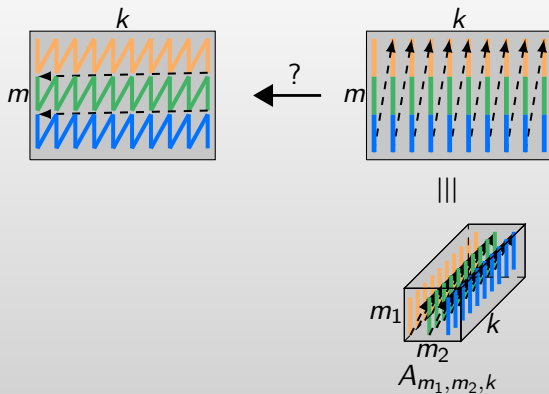
- Input to TTC:

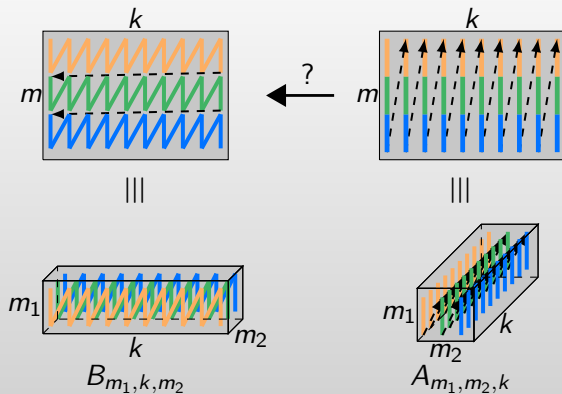
```
t t c  —perm=2,1,0  —size=8,16,16  —lda=32,32,32
```

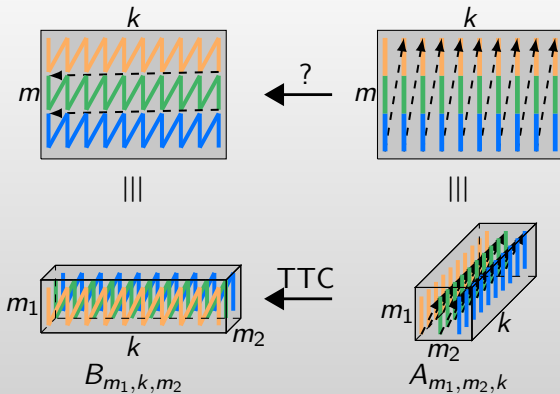












- Input to TTC:

```
ttc --perm=0,2,1 --size= $S_{m_1}, S_{m_2}, S_k$ 
```

## Speedup

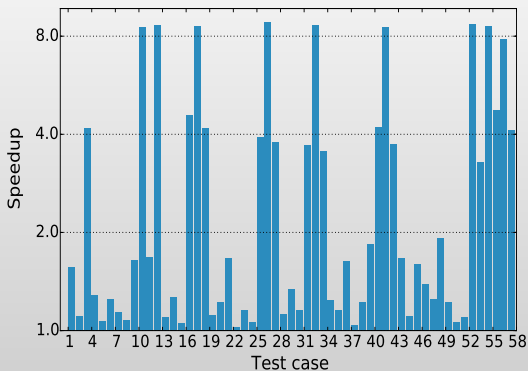


Figure: Intel Haswell architecture. 2x Intel Xeon E5-2680v3 @ 24 threads, ICPC 16.0.1.

## Bandwidth

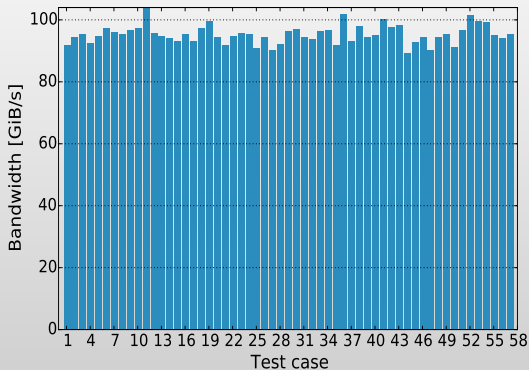


Figure: Intel Haswell architecture. 2x Intel Xeon E5-2680v3 @ 24 threads, ICPC 16.0.1.

## Bandwidth

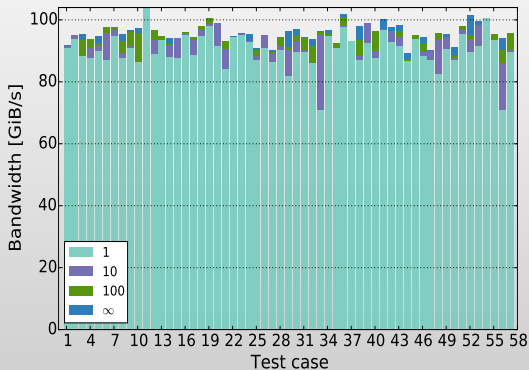
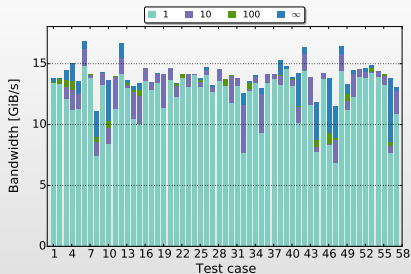
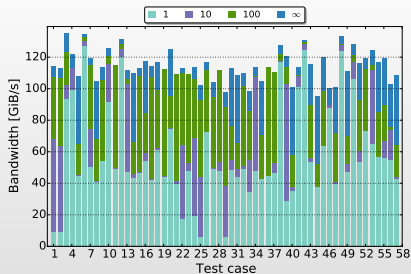


Figure: Intel Haswell architecture. 2x Intel Xeon E5-2680v3 @ 24 threads, ICPC 16.0.1.

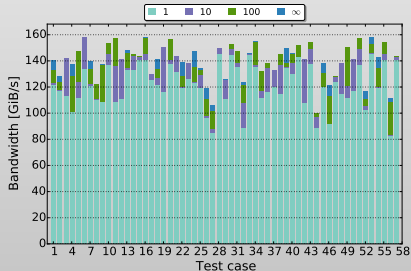




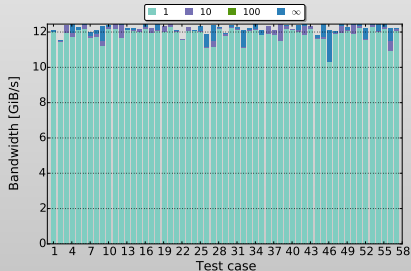
(a) AMD Steamroller, AMD A10-7850K.



(b) Intel KNC, Xeon Phi 5110p.



(c) NVIDIA Kepler, Tesla K40c.

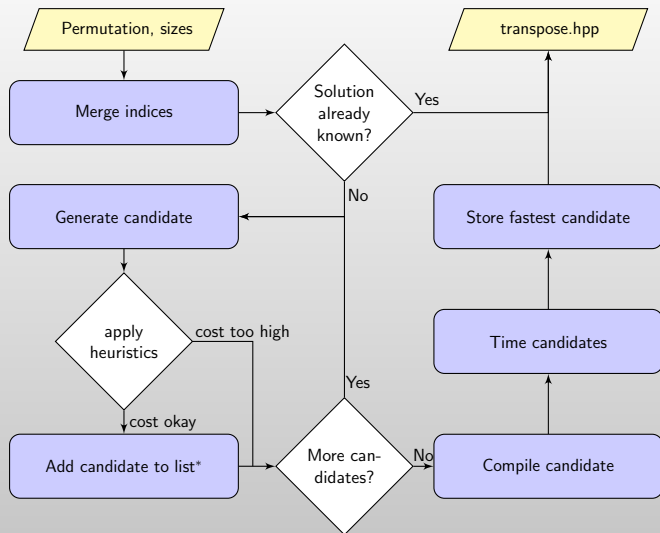


(d) NVIDIA Maxwell, GeForce 840m .

- Domain-specific compiler for multidimensional transpositions
  - Generates high-performance C++/CUDA C code
- Features
  - multiple architectures
  - all numerical data types
  - mixed precision
  - multiple leading dimensions
- Open Source
  - <https://github.com/hpac/ttc>
- Preprint
  - <http://arxiv.org/pdf/1603.02297v1>

- Domain-specific compiler for multidimensional transpositions
  - Generates high-performance C++/CUDA C code
- Features
  - multiple architectures
  - all numerical data types
  - mixed precision
  - multiple leading dimensions
- Open Source
  - <https://github.com/hpac/ttc>
- Preprint
  - <http://arxiv.org/pdf/1603.02297v1>

Thank you for your attention.



\* capacity of the list can be influenced by the user, default capacity: 200.

```
/**
 * B(i2,i1,i0) ← alpha * A(i0,i1,i2)
 *
 * \param[in] A Input tensor
 * \param[out] B Output tensor
 * \param[in] alpha scalar factor of A
 * \param[in] lda leading dimensions of A, can be NULL
 * \param[in] ldb leading dimensions of B, can be NULL
 */
template<int size0, int size1, int size2>
void sTranspose210_8x16x16(const float* A,
                          float* B,
                          float alpha,
                          const int *lda,
                          const int *ldb);
```

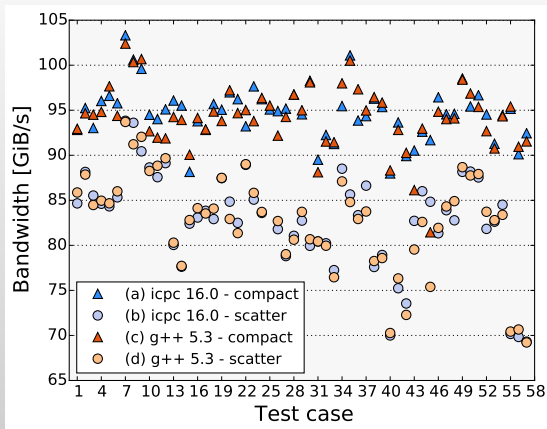


Figure: Influence of thread affinity and compiler on performance.