

# Meta-programming for Cross-Domain Tensor Optimizations

Adilla Susungi<sup>1</sup>, Norman A. Rink<sup>2</sup>, Albert Cohen<sup>3</sup>, Jerónimo Castrillón<sup>2</sup>, Claude Tadonki<sup>1</sup>

<sup>1</sup>MINES ParisTech, PSL Research University

<sup>2</sup>Chair for Compiler Construction, Technische Universität Dresden

<sup>3</sup>Inria, Ecole normale supérieure

17th International Conference on Generative Programming: Concepts & Experiences (GPCE'18)  
Boston, USA  
November 5, 2018



# Tensor optimizations and frameworks

## Tensors

- ▶ Fundamental algebraic structure with applications to many domains
- ▶ Operations on multi-dimensional and computationally intense loop nests
- ▶ Involves multiple optimization strategies: loop, data layout, algebraic transformations, mapping decisions, etc.

## Existing optimizing frameworks

- ▶ Built-in strategies do not always generalize well
- ▶ Lack of flexibility in composing finely tuned, target-specific optimizations

# Transformation meta-languages

## Meta-languages offering transformation heuristics as first-class citizens

```
# Double fusion of the three nests
motion(enclose(C2L1_2_1_2_1),TARGET_2_1_2_1)
motion(enclose(C1L1_2_1_2_1),C2L1_2_1_2_1)
motion(enclose(C3L1_2_1_2_1),C1L1_2_1_2_1)

# Register blocking and unrolling (factor 2)
time_stripmine(enclose(C3L1_2_1_2_1,2),2,2)
time_stripmine(enclose(C3L1_2_1_2_1,1),4,2)
interchange(enclose(C3L1_2_1_2_1,2))
time_peel(enclose(C3L1_2_1_2_1,3),4,'2')
time_peel(enclose(C3L1_2_1_2_1,2),3),4,'N-2')
time_peel(enclose(C3L1_2_1_2_1_2_1,1),5,'2')
time_peel(enclose(C3L1_2_1_2_1_2_1,1),5,'M-2')
fullunroll(enclose(C3L1_2_1_2_1_2_1,2))
fullunroll(enclose(C3L1_2_1_2_1_2_1,1))
```

URUK (Cohen et al.,  
ICS'05)

```
reorder((, (0,2,1,3,4,5,6))
fuse_next((0))
fuse_next((0))
fuse_next((0))
fuse_next((0, 2))
```

Clay (Bagnères et al.,  
CGO'16)

```
permute([1,3,2])
tile(0,3,TK)
split(0,2,[d3≥d1+TK])
tile(0,3,TI,2)
tile(0,3,TJ,2)
datacopy(0,3,2)
datacopy(0,4,3,[1])
unroll(0,4,UJ1)
unroll(0,5,UI1)
datacopy(1,2,3,[1])
unroll(1,2,UJ2)
unroll(1,3,UI2)
```

CHILL (Chen et al., 2008)

```
partialDot(x: [float]N, y: [float]N) =
  (join ◦ mapWrg0 ◦
  (join ◦ toGlobal(mapLcl0(mapSeq(id))) ◦ sp
  iterate0 ◦ join ◦
  mapLcl0( toLocal(mapSeq(id)) ◦
  reduceSeq(add, 0) )
  split2 ) ◦
  join ◦ mapLcl0( toLocal(mapSeq(id)) ◦
  reduceSeq(multAndSumUp, 0) ) ◦ spl
  ) ◦ split128)( zip(x, y) )
```

Lift (Steuer et al.,  
CGO'17)

```
Func blur_3x3(Func input) {
  Func blur_x, blur_y;
  Var x, y, xi, yi;

  // The algorithm - no storage or order
  blur_x(x, y) = (input(x-1, y) + input(x, y) + inp
  blur_y(x, y) = (blur_x(x, y-1) + blur_x(x, y) + b

  // The schedule - defines order, locality; imple
  blur_y.tile(x, y, xi, yi, 256, 32)
  .vectorize(xi, 8).parallel(y);
  blur_x.compute_at(blur_y, x).vectorize(x, 8);

  return blur_y;
}
```

Halide (Ragan-Kelley et  
al., PLDI'14)

```
k = tvm.reduce_axis((0, K), 'k')
A = tvm.placeholder((M, K), name='A')
B = tvm.placeholder((K, N), name='B')
C = tvm.compute((M, N), lambda x, y: A[x, k] * B
s = tvm.create_schedule(C.op)
func = tvm.build(s, [A, B, C], target=target, name
xo, yo, xi, yi = s[C].tile(C.op.axis[0], C.op.axis
k, = s[C].op.reduce_axis
ko, ki = s[C].split(k, factor=4)
```

TVM (Chen et al., 2018)

# Transformation meta-languages

**We are interested in meta-languages for program transformation, because**

- ▶ They help increasing expert productivity when hand-writing optimizations
- ▶ They ease the composition and cancellation of transformations
- ▶ They make the optimization paths explicit and future-proof

**Strong allies for building adaptive, portable and efficient compiler infrastructures to face the complexity of parallel architectures**

# Contributions outline

## Keys to

- ▶ Widen optimization search space
- ▶ Enhance the ability to flexibly compose optimization paths
- ▶ Formally characterize their semantics

**Design and semantics of a tensor optimizations meta-language (TeML)**

# TeML overview

<code>&lt;program&gt;</code>	<code>::= &lt;stmt&gt; &lt;program&gt;</code> <code>  <math>\epsilon</math></code>
<code>&lt;stmt&gt;</code>	<code>::= &lt;id&gt; = &lt;expression&gt;</code> <code>  &lt;id&gt; = @&lt;id&gt; : &lt;expression&gt;</code> <code>  codegen (&lt;ids&gt;)</code> <code>  init (...)</code>
<code>&lt;expression&gt;</code>	<code>::= &lt;Texpression&gt;</code> <code>  &lt;Lexpression&gt;</code>
<code>&lt;Texpression&gt;</code>	<code>::= scalar ()</code> <code>  tensor ([&lt;ints&gt;])</code> <code>  eq (&lt;id&gt;, &lt;iters&gt;? <math>\rightarrow</math> &lt;iters&gt;)</code> <code>  vop (&lt;id&gt;, &lt;id&gt;, [&lt;iters&gt;?], &lt;iters&gt;?)</code> <code>  op (&lt;id&gt;, &lt;id&gt;, [&lt;iters&gt;?], &lt;iters&gt;?) <math>\rightarrow</math> &lt;iters&gt;</code>
<code>&lt;Lexpression&gt;</code>	<code>::= build (&lt;id&gt;)</code> <code>  stripmine (&lt;id&gt;, &lt;int&gt;, &lt;int&gt;)</code> <code>  interchange (&lt;id&gt;, &lt;int&gt;, &lt;int&gt;)</code> <code>  fuse (&lt;id&gt;, &lt;id&gt;, &lt;int&gt;)</code> <code>  unroll (&lt;id&gt;, &lt;int&gt;)</code>
<code>&lt;iters&gt;</code>	<code>::= [&lt;ids&gt;]</code>
<code>&lt;ids&gt;</code>	<code>::= &lt;id&gt; (, &lt;id&gt;)*</code>
<code>&lt;ints&gt;</code>	<code>::= &lt;int&gt; (, &lt;int&gt;)*</code>

## Every function returns either

- ▶ Tensors
- ▶ Loops

## Operations on tensors

- ▶ Computation specification
- ▶ Layout transformations
- ▶ Data initialization, mapping

## Operations on loops

- ▶ Expansion from tensor computation
- ▶ Transformation

# TeML overview

Raising the level of abstraction

## A contraction chain

$$v_{ijk} = \sum_{l,m,n} A_{kn} \cdot A_{jm} \cdot A_{il} \cdot u_{lmn}$$

## Control the evaluation order

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot (A_{jm} \cdot (A_{il} \cdot u_{lmn})))$$

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot A_{jm}) \cdot (A_{il} \cdot u_{lmn})$$

$$v_{ijk} = \sum_{l,m,n} (A_{kn} \cdot ((A_{jm} \cdot A_{il}) \cdot u_{lmn}))$$

- ▶ The evaluation order may dramatically impact execution time
- ▶ May be combined with other transformation heuristics

# TeML overview

## Raising the level of abstraction

### Tensor-algebraic transformations are essential some applications

- ▶ Out of the scope of polyhedral-based meta-languages
- ▶ Or requires additional analyses to (re)discover algebraic tensor properties

```
# -- Begin program specification
w = tensor(double, [13])
u = tensor(double, [13, 13, 13])
L = tensor(double, [13, 13])
M_ = outerproduct([w, w, w])
Lh = div(L, w, [[i1, i2], [i2]] ->
↪ [i1, i2])
M = entrywise_mul(M_, u)
r1 = contract(Lh, M, [[2, 1]])
r2 = contract(Lh, M, [[2, 2]])
r3 = contract(Lh, M, [[2, 3]])
# -- End program specification
```

- ▶ We want such characterizations to be native to the language
- ▶ Provides room for encoding algebraic properties



# TeML overview

By example: facilitating transformation composition

- ▶ Existing meta-languages are either fully imperative or mix a functional specification of the computation with an imperative transformation sequence
- ▶ We use a functional style for both program stages

```
# -- Begin program specification
w = tensor(double, [13])
u = tensor(double, [13, 13, 13])
L = tensor(double, [13, 13])
M_ = outerproduct([w, w, w])
Lh = div(L, w, [[i1, i2], [i2]] ->
↳ [i1, i2])
M = entrywise_mul(M_, u)
r1 = contract(Lh, M, [[2, 1]])
r2 = contract(Lh, M, [[2, 2]])
r3 = contract(Lh, M, [[2, 3]])
# -- End program specification

# Generate loops
11 = build(M_)
12 = build(Lh)
13 = build(M)
14 = build(r1)
15 = build(r2)
16 = build(r3)
```

```
# Code generation without
↳ transformations
codegen([11, 12, 13, 14, 15, 16])

17 = fuse(14, 15, 3)
18 = fuse(17, 16, 3)
# Code generation with loop fusions
↳ only
codegen([11, 12, 13, 17])

19 = parallelize(11, 1, None)
110 = parallelize(12, 1, None)
111 = parallelize(13, 1, None)
112 = parallelize(18, 1, None)
113 = vectorize(19, 3)
114 = vectorize(110, 2)
115 = vectorize(111, 3)

# Code generation with fusion,
↳ parallelism and vectorization
codegen([113, 114, 115, 112])
```

# Denotational semantics

## Domains of trees for tensors (**T**) and loops (**L**)

### State

- ▶ A state in a `TEML` meta-program maps identifiers to trees representing either tensors or loops

$$\mathbf{S} = \textit{identifier} \rightarrow (\mathbf{T} + \mathbf{L})$$

$$\sigma : \textit{identifier} \rightarrow (\mathbf{T} + \mathbf{L})$$

### Valuation functions

- ▶ Different manipulations of a state  $\sigma$  for each syntactic entity

$$\mathcal{P}_{prog} : \textit{program} \rightarrow (\mathbf{S} \rightarrow \mathbf{S})$$

$$\mathcal{P}_{stmt} : \textit{stmt} \rightarrow (\mathbf{S} \rightarrow \mathbf{S})$$

$$\mathcal{E}_t : \textit{Texpression} \rightarrow (\mathbf{S} \rightarrow \mathbf{T})$$

$$\mathcal{E}_l : \textit{Lexpression} \rightarrow (\mathbf{S} \rightarrow \mathbf{L})$$

# Semantics of tensor expressions

## Subtleties

```
for (int i1 = 0; i1 < (N-1); i1++)  
  for (int i2 = 0; i2 < (N-1); i2++)  
    E[i1][i2] = C[i1][i2] * (A[i1][i2] + B[i1][i2]);
```

```
A = tensor([N, N])  
B = tensor([N, N])  
C = tensor([N, N])  
D = vadd(A, B, [[i1, i2], [i1, i2]])  
E = mul(C, D, [[i1, i2], ] -> [i1, i2])
```

- ▶ We use virtual operators (vops) to compose beyond 3-address expressions
- ▶ Tensors returned by vops only hold subexpressions eventually expanded recursively at instances of ops
- ▶ Tensors returned by vops do not have shapes of their own
- ▶ Others have their shape inferred, as well as their loop domains

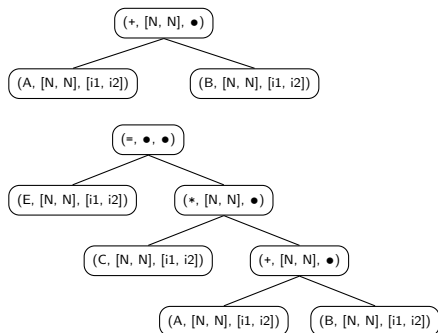
# Semantics of tensor expressions

## Low-level operations

### Essential informations to capture

- ▶ Shape
- ▶ Expression tree
- ▶ Associated list of iterators

```
A = tensor([N, N])
B = tensor([N, N])
C = tensor([N, N])
D = vadd(A, B, [[i1, i2], [i1, i2]])
E = mul(C, D, [[i1, i2], ] -> [i1, i2])
```



$$\sigma_1 = \mathcal{P}_{stmt} \llbracket A = \text{tensor}([N, N]) \rrbracket \emptyset$$
$$= \{ A \mapsto \langle (A, [N, N], \epsilon), [] \rangle \}$$

$$\sigma_2 = \mathcal{P}_{stmt} \llbracket B = \text{tensor}([N, N]) \rrbracket \sigma_1$$
$$= \{ A \mapsto \langle (A, [N, N], \epsilon), [] \rangle, B \mapsto \langle (B, [N, N], \epsilon), [] \rangle \}$$

$$\sigma_3 = \dots$$

# Semantics of tensor expressions

## High-level operations

### The example of tensor contraction

$$\mathcal{P}_{stmt} \llbracket t' = \text{contract}(t_0, t_1, [r_0, r_1]) \rrbracket = \mathcal{P}_{prog} \left[ \left[ \begin{array}{l} t_2 = \text{vmul}(t_0, t_1, [I, J]) \\ t' = \text{add}(t', t_2, [I', \epsilon] \rightarrow I') \end{array} \right] \right]$$

where

$$\begin{aligned} I &= [i_0, \dots, i_{(r_0-1)}, \mathbf{k}, i_{(r_0+1)}, \dots, i_{s_0}], \\ J &= [j_0, \dots, j_{(r_1-1)}, \mathbf{k}, j_{(r_1+1)}, \dots, j_{s_1}], \\ I' &= (I \setminus \{\mathbf{k}\}) \parallel (J \setminus \{\mathbf{k}\}). \end{aligned}$$

# Semantics of loop transformations

- ▶ Principles of loop transformations are quite well understood.
- ▶ The polyhedral model is a rich formalism to abstracts the effects of loop transformations
- ▶ The idea here is to formalize such principles in a meta-language context

## Example

```
for (int i1 = 0; i1 <= (N-1); i1++) {  
  C[i1] = A[i1] - B[i1]; // tC  
  for (int i2 = 0; i2 <= (N-1); i2++) {  
    E[i1][i2] = D[i2] * C[i1]; // tE  
    F[i1][i2] = E[i1][i2]; // tF  
  }  
  for (int i3 = 0; i3 <= (N-1); i3++) {  
    G[i1] = G[i1] + F[i1][i3] // tG  
  }  
}
```

$\langle i1, [tC,$   
 $\quad \langle i2, [tE, tF]\rangle,$   
 $\quad \langle i3, [tG]\rangle$   
 $\quad \rangle \rangle$

# Loop creation from tensor expressions

- ▶ The semantics of build

```
A = tensor([N, N])
B = tensor([N, N])
C = tensor([N, N])
D = vadd(A, B, [[i1, i2], [i1, i2]])
E = mul(C, D, [[i1, i2], ] -> [i1, i2])
```

$$\underline{\mathcal{E}_l[\text{build}(E)]\sigma_5 = \langle i1, [\langle i2, [\sigma_5(E)] \rangle] \rangle :}$$

```
for (int i1 = 0; i1 <= (N-1); i1++)
  for (int i2 = 0; i2 <= (N-1); i2++)
    E[i1][i2] = C[i1][i2] * (A[i1][i2] + B[i1][i2]);
```

# Semantics of loop expressions

## Stripmining

- ▶ Divides an iteration space into smaller blocks

```
A = tensor([N, N])
B = tensor([N, N])
C = tensor([N, N])
D = vadd(A, B, [[i1, i2], [i1, i2]])
E = mul(C, D, [[i1, i2], ] -> [i1, i2])
L = build(E)
S = stripmine(L, 1, 32)
```

$$\begin{aligned}\sigma_n &= \mathcal{P}_{stmt} \llbracket L = \text{build}(E) \rrbracket \sigma_{n-1} \\ &= \{ L \mapsto \langle i1, [\langle i2, [\sigma_{n-1}(E)] \rangle] \rangle \}\end{aligned}$$

$$\begin{aligned}\sigma_{n+1} &= \mathcal{P}_{stmt} \llbracket S = \text{stripmine}(L, 1, 32) \rrbracket \sigma_n \\ &= \{ L \mapsto \langle i1, [\langle i2, [\sigma_n(E)] \rangle] \rangle, S \mapsto \langle t1, [\langle i1, [\langle i2, [\sigma_n(E)] \rangle] \rangle] \rangle \}\end{aligned}$$

$$\mathcal{E}_t \llbracket \text{stripmine}(L, 1, 32) \rrbracket \sigma_n = \langle t1, [\langle i1, [\langle i2, [\sigma_n(E)] \rangle] \rangle] \rangle :$$

```
for (int t1 = 0; t1 <= (N-1)/32; t1++)
  for (int i1 = 32 * t1; i1 <= min((N-1), 32 * t1 + 31); i1++)
    for (int i2 = 0; i2 <= (N-1); i2++)
      E[i1][i2] = C[i1][i2] * (A[i1][i2] + B[i1][i2]);
```



# Semantics of loop expressions

## Interchange

- ▶ Swaps dimensions of a loop nest

```
A = tensor([N, N])
B = tensor([N, N])
C = tensor([N, N])
D = vadd(A, B, [[i1, i2], [i1, i2]])
E = mul(C, D, [[i1, i2], ] -> [i1, i2])
L = build(E)
I = interchange(L, [1, 2])
```

$$\begin{aligned}\sigma_n &= \mathcal{P}_{stmt} \llbracket L = \text{build}(E) \rrbracket \sigma_{n-1} \\ &= \{ L \mapsto \langle i1, [\langle i2, [\sigma_{n-1}(E)] \rangle] \rangle \}\end{aligned}$$

$$\begin{aligned}\sigma_{n+1} &= \mathcal{P}_{stmt} \llbracket I = \text{interchange}(L, [1, 2]) \rrbracket \sigma_n \\ &= \{ L \mapsto \langle i1, [\langle i2, [\sigma_n(E)] \rangle] \rangle, I \mapsto \langle i2, [\langle i1, [\sigma_n(E)] \rangle] \rangle \}\end{aligned}$$

$$\underline{\mathcal{E}_l \llbracket \text{interchange}(L, [1, 2]) \rrbracket \sigma_n = \langle i2, [\langle i1, [\sigma_n(E)] \rangle] \rangle :}$$

```
for (int i2 = 0; i2 <= (N-1); i2++)
  for (int i1 = 0; i1 <= (N-1); i1++)
    E[i1][i2] = C[i1][i2] * (A[i1][i2] + B[i1][i2]);
```

# Semantics of loop expressions

## Loop tiling in denotational semantics

- ▶ Loop tiling is the composition of stripmining and interchange

```
for (int t1 = 0; t1 <= (N-1)/32; t1++)
  for (int t2 = 0; t2 <= (N-1)/32; t2++)
    for (int i1 = 32 * t1; i1 <= min((N-1), 32 * t1 + 31); i1++)
      for (int i2 = 32 * t2; i2 <= min((N-1), 32 * t2 + 31); i2++)
        E[i1][i2] = C[i1][i2] * (A[i1][i2] + B[i1][i2]);
```

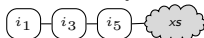
$\mathcal{P}_{stmt} \llbracket l' = \text{tile}(l, v) \rrbracket =$

$$\mathcal{P}_{prog} \left[ \begin{array}{l} l_0 = \text{stripmine}_n(l, d, v) \\ l_1 = \text{interchange}_n(l_0, 2, 2d - 2) \\ l_2 = \text{interchange}_n(l_1, 3, 2d - 3) \\ \dots \\ l_{d-1} = \text{interchange}_n(l_{d-2}, d, d) \\ l' = \text{interchange}_n(l_{d-1}, d + 1, d - 1) \end{array} \right]$$

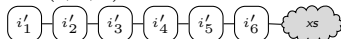
# Semantics of loop expressions

## Loop tiling in denotational semantics

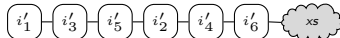
### Initial loop nest



`stripmine_n(-, 3, v)` has introduced  $i'_2$ ,  $i'_4$ , and  $i'_6$



After triple application of `interchange_n`



# TeML evaluation

Expressing tensor computations in comparison to TensorFlow

Application domains: Linear Algebra (LA), Deep Learning (DL), Machine Learning (ML), Data Analytics (DA), Fluid Dynamics (FD), Image Processing (IP).

	Name	Domain	TensorFlow		TeML	
			LOC	Constructs used	LOC	Constructs used
<b>Matrix Multiplication</b>	<i>mm</i>	LA	3	matmul	3	contract
transposed	<i>tmm</i>	DL	3	matmul:transpose=True	4	transpose, contract
batched	<i>bmm</i>		3	einsum	3	mul, add
<b>Grouped Convolutions</b>	<i>gconv</i>		N/A	Not implemented Incompatible with einsum	5	vmul, add
<b>Matricized Tensor Times Khatri-Rao product</b>	<i>mtkrp</i>	DA	4	einsum or tensordot, multiply	5	vcontract, contract
<b>Sampled Dense-Dense Matrix Product</b>	<i>sddmm</i>	ML	4	einsum or tensordot, multiply	6	vcontract, entrywise_mul
<b>Interpolation</b>	<i>interp</i>	FD	3	einsum or tensordot	5	contract
<b>Helmholtz</b>	<i>helm</i>		N/A	Required division not well supported	9	contract, outerproduct, div, entrywise_mul
<b>Blur</b>	<i>blur</i>	IP	N/A	No stencil support.	9	op, vop
<b>Coarsity</b>	<i>coars</i>		6	einsum or multiply, subtract	6	ventrywise_mul, entrywise_sub

# TeML evaluation

## Reproducing optimization paths of Pluto

### Pluto

- ▶ Polyhedral automatic parallelizer
- ▶ Some flexibility in selecting optimizations and their parameters
- ▶ But quite rigid heuristics, mostly “black-box” optimizations

<i>mttkrp</i> (250*250*250)	<i>sddmm</i> (4096*4096)	<i>bmm</i> (8192*72*26)	<i>gconv</i> (32*32*32*32*7*7)	<i>interp</i> (50000*7*7*7)	<i>helm</i> (5000*13*13*13)	<i>coars</i> (4096*4096)
parallelize(1, 1) interchange(1, 2, 3)	interchange(1, 2, 3), parallelize(1, 1), vectorize(1, 3)	tile(1, 32) interchange(1, 7,8) parallelize(1, 1) vectorize(1, 8)	interchange(11, 4, 5) interchange(11, 5, 6) parallelize(11, 1) vectorize(11, 9) parallelize(12, 1) vectorize(12, 9)	interchange(11, 4, 5), vectorize(11, 5), interchange(12, 4, 5), vectorize(12, 5), parallelize(11, 1), parallelize(12,1), parallelize(13, 1)	fuse_outer(14, 15, 5), fuse_outer(14, 16, 5), parallelize(11, 1), parallelize(12, 1), parallelize(13, 1), parallelize(14, 1), vectorize(11, 2), vectorize(12, 3), vectorize(13, 4)	tile(1, 32) parallelize(1, 1) vectorize(1, 4)

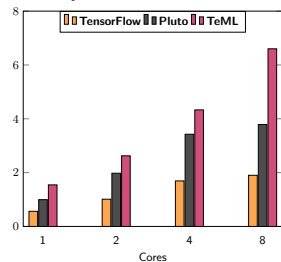
- ▶ Can we outperform Pluto?

# TeML evaluation

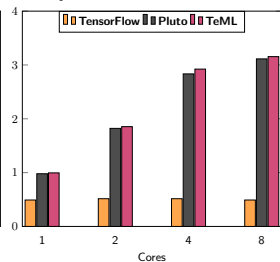
Expressing transformations that outperform Pluto

- ▶ On Intel(R) Core(TM) i7-4910MQ CPU (2.90GHz, 8 hyperthreads, 8192KB of shared L3 cache), Ubuntu 16.04
- ▶ Generated C programs compiled with the Intel C compiler ICC 18.02 (flags: -O3 -xHost -qopenmp)
- ▶ TensorFlow version 1.6 with support for AVX, FMA, SSE, and multi-threading

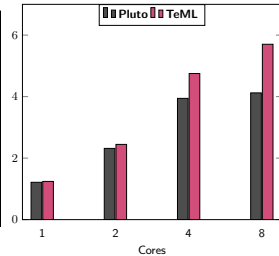
### mttkrp



### interp

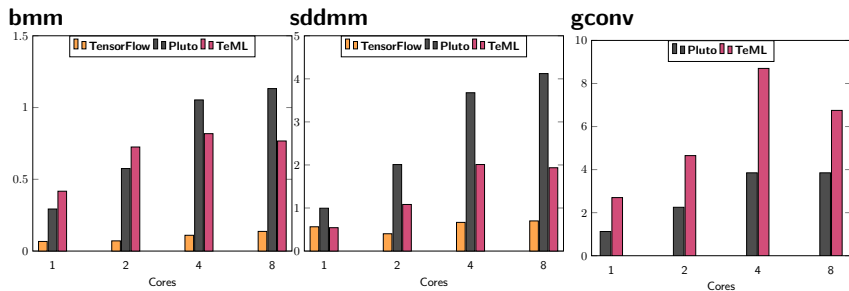


### helm



# TeML evaluation

Expressing transformations that outperform Pluto



- ▶ We are able to express more efficient transformation paths

# Conclusion

## TEML

- ▶ Program construction and transformation phases are both functional
- ▶ Higher-level of abstractions for tensor computations
- ▶ Formal specification of program construction and transformation

## Future work

- ▶ Extensions for parallelism support
- ▶ Abstractions for memory virtualization and corresponding semantics
- ▶ Type system
- ▶ High-level abstractions for stencil patterns, general convolutions, sparse tensors