

Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal,

LLVM Workshop at CGO 2018 February 24, 2018





- Zach DeVito, William S. Moses, Sven Verdoolaege, Andrew Adams, Albert Cohen





The Tensor Comprehensions Team

























cudnnConvolutionForward()



• • •

• • •

dnnConvolutionCreateForwa rd_F32()

* TF also can compile via XLA, discussed later









Someone has a clever idea

- Suppose a ML researcher invents a new layer: hconv
- He/she can implements it two ways:
 - Inefficiently cobbling together existing operators [slow]
 - Write optimized GPU/CPU kernel [difficult, time-consuming]
- Even when the operator exists, it often misses peakperformance, lacking cross-operator-optimization and data-shape/size tuning [1]

[1] Fast Convolutional Nets With fbfft : A GPU Performance Evaluation, ICLR, 2015



"Abstraction without regret"

- To make development efficient, we need abstractions that provide productivity without sacrificing performance
- Given the enormous number of potential kernels, suggests a dynamic-code-generation approach



Prior work

- "Direct generation" such as active library [2] or built-to-order (BTO) [3] provide usability, but miss optimization
- DSLs such as Halide [4] provide usability, and permit scheduling transformations, though manually specify.
- Compilers like XLA [5] or Latte [6] optimize and fuse operators, though performance lacking as the language can't represent complex schedules crucial to GPU/others.

[2] Run-time code generation in C++ as a foundation for domain-specific optimization, 2003 [3] Automating the generation of composed linear algebra kernels, 2009 [4] Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines, 2013 [5] Xla:domain-specific compiler for linear algebra to optimizes tensorflow computations, 2017 [6] Latte: A language, compiler, and runtime for elegant and efficient deep neural networks, 2016



Tensor Comprehensions

- High-level DSL to express tensor computations by extending Einstein-notation.
- End-to-End compilation flow capable of lowering tensor comprehensions to efficient GPU code (CPU in progress)
- Collection of polyhedral compilation algorithms with a specific domain and target orientation
- Autotuning framework built off JIT compilation and caching Integration into ML Frameworks (Caffe2, Pytorch)











Tensor Comprehensions In Progress Polyhedral Transformations Tapir/LLVM Polly Cilk/OpenMP





Optimizations at the appropriate time

- High-level polyhedral: broader scheduling optimizations (mapping, tiling, fusion, etc)
- Halide: Expression simplification / optimizations
- Tapir/LLVM [7] <-> Polly [8]: Runtime-level optimization/scheduling (coarsening, vectorization), instruction-level optimization (i.e. LICM, fuse instructions)
- [7] Tapir: Embedding Fork-Join Parallelism into LLVM's Intermediate Representation, 2017 <u>https://github.com/wsmoses/Tapir-LLVM</u>
- [8] Polly Performing polyhedral optimizations on a low-level intermediate representation, 2012 <u>https://github.com/wsmoses/Tapir-Polly</u>



TC language Concise, emits 1000's of optimized LOC

```
def mv(float(M,K) A, float(K) x) -> (C) {
 C(i) +=! A(i,k) * x(k)
}
```

01(n, o, h, w) +=! I(n, c, h + kh, w + kw) * W1(o, c, kh, kw) 01(n, o, h, w) = fmax(01(n, o, h, w), 0) // relu 02(n, d, h, w) +=! 01(n, d, h + kh, w + kw) * W2(d, o, kh, kw) O2(n, d, h, w) = fmax(O2(n, d, h, w), 0)03(n, e, h, w) +=! 02(n, c, h + kh, w + kw) * W3(e, d, kh, kw) 03(n, e, h, w) = fmax(03(n, e, h, w), 0)





Polyhedral + TC

- High Level Polyhedral IR (ISL) => Easy Transformations
- Schedule heuristic folds into a single kernel
- Schedule tiled to facilitate the mapping and reuse of memory hierarchy of GPU/CPU
- GPU mapping borrows from PPCG, with extensions for more complex/imperfectly nested structures
- Memory promotion into shared cache







ISL scheduling

def sgemm(float a, float b float(N,M) A, float(M,K) B) -> (C) { C(i,j) = b//S(i,j) C(i,j) += a * A(i,k) * B(k,j) //T(i,j,k)

Domain
$$\begin{bmatrix} \{S(i,j) \mid 0 \le i < N \land 0 \le j < K\} \\ \{T(i,j,k) \mid 0 \le i < N \\ \land 0 \le j < K \land 0 \le k < M\} \end{bmatrix}$$

Sequence
Filter{ $S(i,j)$ }
Band{ $S(i,j) \rightarrow (i,j)$ }
Filter{ $T(i,j,k)$ }
Band{ $T(i,j,k) \rightarrow (i,j,k)$ }
Band{ $T(i,j,k) \rightarrow (i,j,k)$ }
Band node: (partial) exect

Filter node: partition iteration space



uence node: order-dependent collection of nodes

cution



ISL scheduling

- ISL's scheduling algorithm
 - Works by solving a linear program

• Uses affine clustering, computing schedule for each stronglyconnected components then scheduling those together



Extending ISL scheduling

- Extended ISL's scheduler to allow additional constraints Affine constraint added to the LP
- - Supply clustering decision for graph component combining
- Clustering allows for conventional minimum and maximum fusion targets AND maximum fusion that preserves at least three nested parallel loops (i.e. for mapping to CUDA blocks / threads)



Memory promotion

- Cache indirectly accessed arrays

- reused, and >= 1 access without memory coalescing
- coalescing)

O[l+Idx[i][j]][k] => shared_0[l][i][j][k] Only done when 0 and Idx are only read (not written)

• Promote directly accesses if tile of fixed size, elements Promote indirectly accessed arrays in same way (ignore



Autotuning

- Even with heuristics, there's a large space of options
- genetic algorithm with fixed search-time.



Derive schedule (and other parameters) by searching via





How well does it work?

End-to-end benchmarks Baseline CUDA 8.0, CUBLAS 8.0, CUDNN 6.0, CUB recent



8 Pascal nodes with 2 socket, 14 core Intel(R) Xeon(R) CPU E5-2680 v4 @ 2.40GHz, with 8 Tesla P100-SXM2 GPUs and 16GB of memory each. Median runtime out of a batch of 1000 Autotuning time out O(hours)

Pascal GPU Benchmarks



TC overview

"Natural ML math running faster than libraries"

- Productive environment to develop ML
- Comparable or better than hand-coded operators
- Perform true kernel fusion, with optimization
- Specialize to specific architecture and sizes
- Autotuning "unlocks" much of polyhedral benefits



Future work

- Share best implementations, for any architecture
- Port to more architectures & accelerators, leveraging highly optimized primitives
- Implement symbolic automatic differentiation directly
- Allow sparse, vector and mixed- precision types
- Support more dynamic control flow and ML architectures
- Integrate with other frameworks





TC overview

"Natural ML math running faster than libraries"

- Available stand-alone and in Caffe2/PyTorch bindings [public in a few days]
- Open source:

• Paper:

https://arxiv.org/abs/1802.04730

<u>https://github.com/facebookresearch/tensorcomprehensions</u>





Questions?

Backup Slides

TC in Practice

```
import tc
ee = tc.ExecutionEngine()
ee.define("""
 def mm(float(M,K) A,
         float(K,N) B) -> (C) {
   C(m,n) +=! A(m,kk) * B(kk,n)
"""\
```

import torch A = torch.randn(3,4)B = torch.randn(4,5)C = ee.mm(A, B)

Figure 12: JIT compile, tune, or hit the compi-

Figure 11: Build execution engine lation cache, then run

```
def 2LUT(float(E1,D) LUT1, int(B,L1) I1,
        float(E2,D) LUT2, int(B,L2) I2) -> (01,02) {
  01(i,j) +=! LUT1(I1(i,k),j)
 02(i,j) +=! LUT2(I2(i,k),j)
def MLP1(float(B,M) I, float(0,N) W1, float(0) B1) -> (01) {
  01(b,n) = B1(n)
 01(b,n) += I(b,m) * W1(n,m)
 01(b,n) = fmaxf(01(b,n), 0)
def MLP3(float(B,M) I, float(0,N) W2, float(0) B2,
        float(P,0) W3, float(P) B3, float(Q,P) W4,
        float(Q) B4) -> (01,02,03,04) {
  02(b,o) = B2(o)
  02(b,o) += 01(b,n) * W2(o,n)
  02(b,o) = fmaxf(02(b,o), 0)
  03(b,p) = B3(p)
  03(b,p) += 02(b,o) * W3(p,o)
  O3(b,p) = fmaxf(O3(b,p), 0)
  04(b,q) = B4(q)
  04(b,q) += 03(b,p) * W4(q,p)
  04(b,q) = fmaxf(04(b,q), 0)
def prodModel(float(E1,D) LUT1, int(B,L1) I1,
              float(E2,D) LUT2, int(B,L2) I2,
              float(B,WX) I3, float(WY,WX) W,
              float(N,M) W1, float(N) B1,
              float(0,N) W2, float(0) B2,
              float(P,0) W3, float(P) B3,
              float(Q,P) W4, float(Q) B4)
    -> (C1,C2,C3,I,01,02,03,04) {
  (C1, C2)
             = 2LUT(LUT1, I1, LUT2, I2)
           += I3(b,wxx) * W(wy,wxx)
  C3(b,wy)
             = concat(C1, C2, C3) # not implemented yet
  I(b,m)
             = MLP1(I, W1, B1)
  (02,03,04) = MLP3(01,W2,B2,W3,B3,W4,B4)
 # 04 goes out to binary classifier, omitted here
```

Figure 17: Full production model (pseudo-code)

