

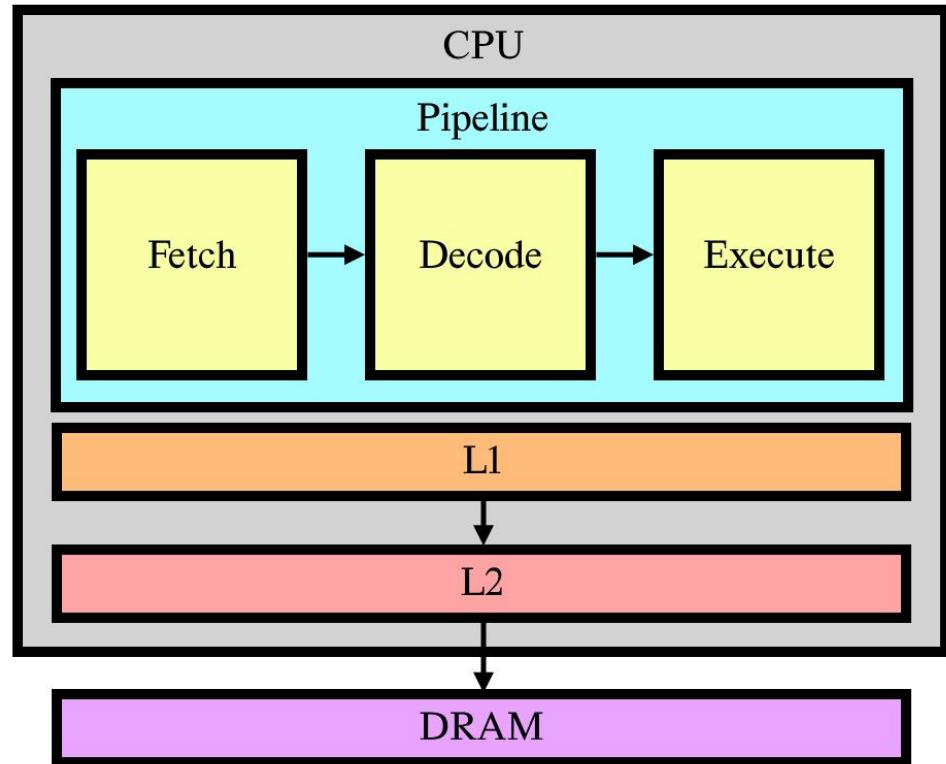
Hardware and Software Optimizations for Deep Learning Workloads on Graphics Processing Units

Justin Garrigus



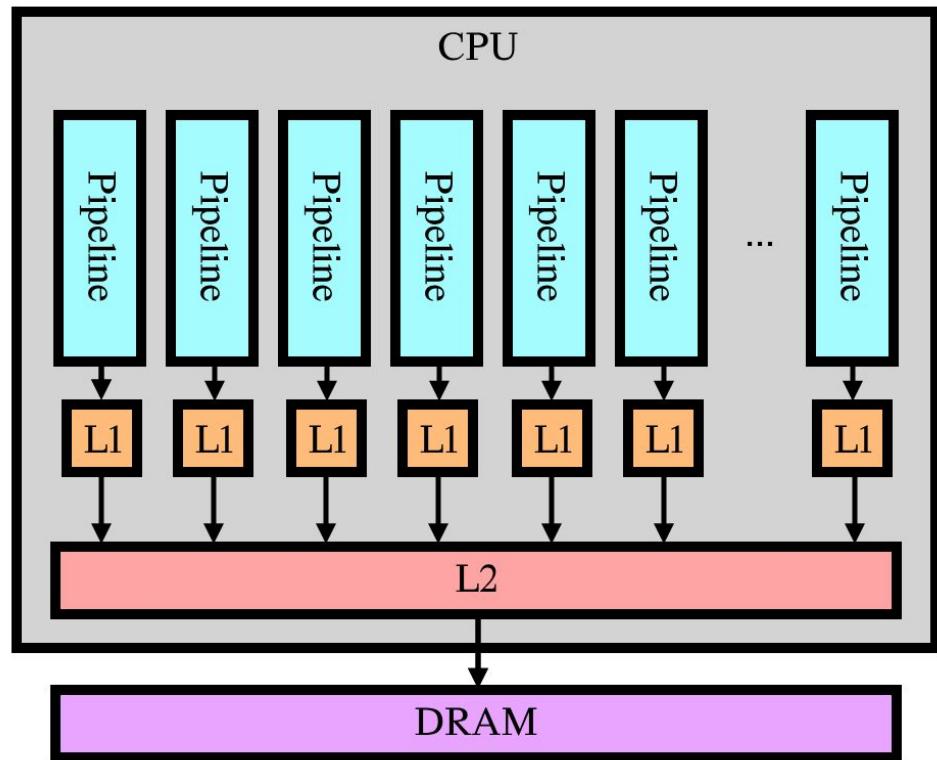
Single-core Processor

- Single pipelines
- One instruction type executed at a time



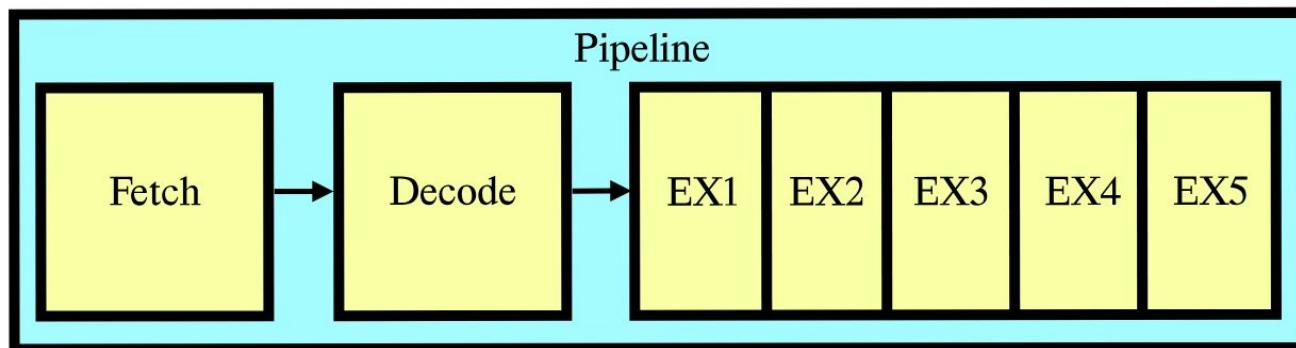
Multi-core Processor

- Multiple pipelines
- Multiple instruction types executed at a time



Vector Processor: Pipelined Execution Unit

- Sub-pipelines
- Execution phase split
- 1 CPI on average for vectors



$MUL[0] = 6$ cycles

$MUL[1] = 7$ cycles

$MUL[2] = 8$ cycles

$MUL[3] = 9$ cycles

$MUL[4] = 10$ cycles

$MUL[5] = 11$ cycles

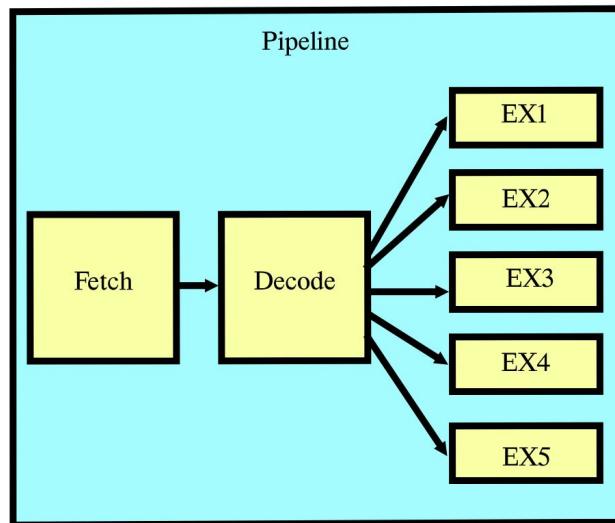
...

$MUL[99] = 105$ cycles

$105 / 99 \approx 1$ CPI

Vector Processor: Duplicated Execution Units

- No sub-pipelines
- Duplicated pipelines, equal to startup time
- 1 CPI on average for vectors

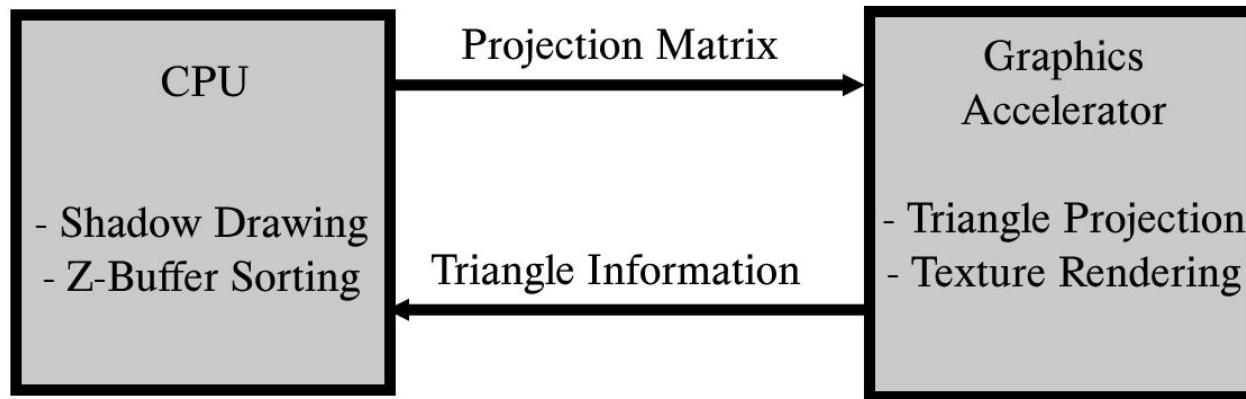


$\text{MUL}[0] = 6 \text{ cycles}$
 $\text{MUL}[1] = 6 \text{ cycles}$
 $\text{MUL}[2] = 6 \text{ cycles}$
 $\text{MUL}[3] = 6 \text{ cycles}$
 $\text{MUL}[4] = 6 \text{ cycles}$
 $\text{MUL}[5] = 12 \text{ cycles}$
 $\text{MUL}[6] = 12 \text{ cycles}$
 $\text{MUL}[7] = 12 \text{ cycles}$
 ...
 $\text{MUL}[99] = 120 \text{ cycles}$

$$120 / 99 \approx 1 \text{ CPI}$$

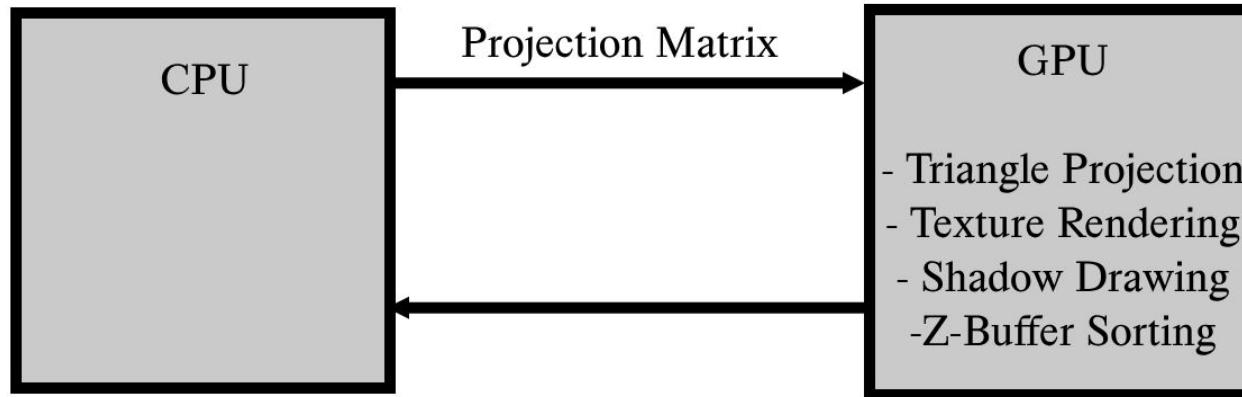
3D Graphics Accelerator

- Polygon rendering
 - Vector translator
- Shadows and textures still on host



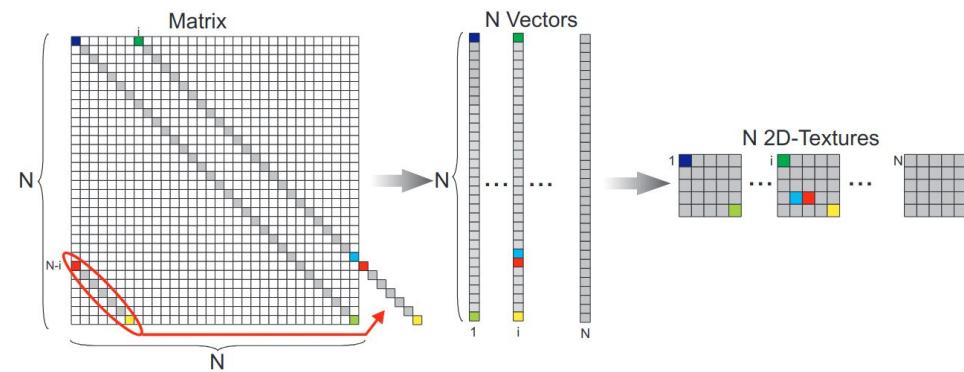
NVIDIA GeForce 256: First “Graphics Processing Unit”

- Marketing term
- Shaders, lighting, and more on device



GPUs for Linear Algebra

- Vector arithmetic
- Matrix-vector product
- Sparse matrices
- General numerical computation
 - Conjugate gradient method
 - Solving sparse $Ax = b$
 - Gauss-Seidel solver
 - Solving $x^i = Lx^i + (D + U)x^{(i-1)}$



CUDA

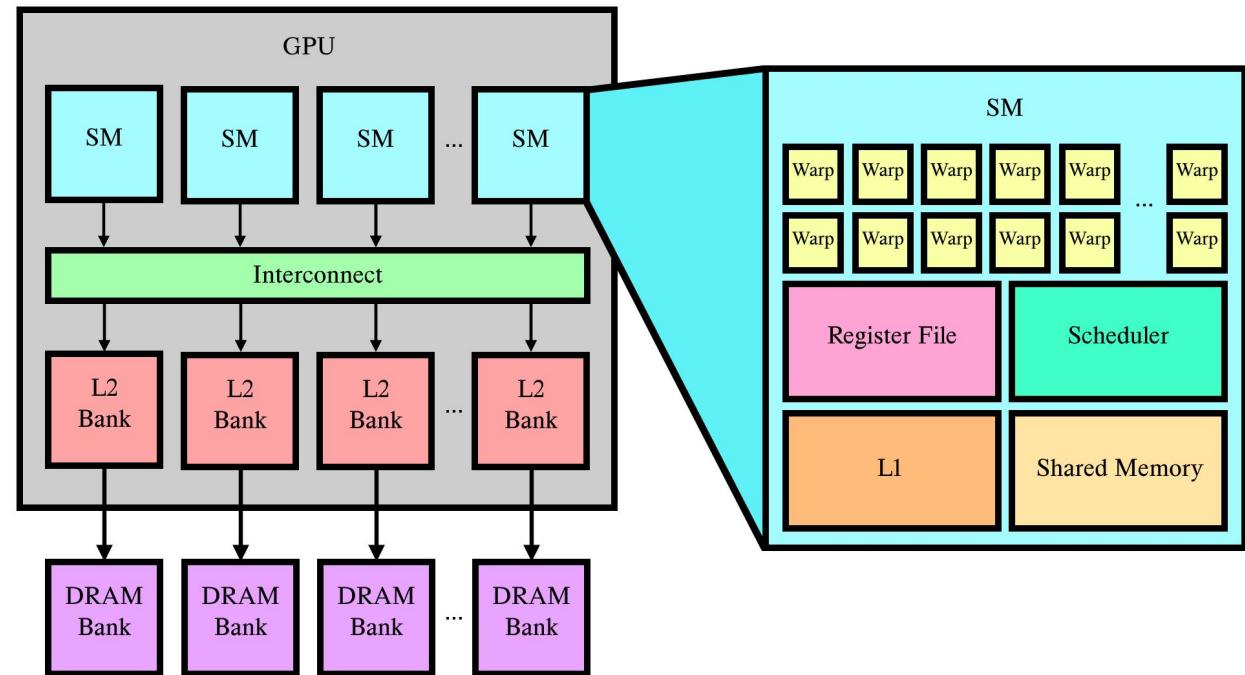
- Host (CPU) code and device (GPU) code
- Host data and device data
 - cudaMemcpy, send data across wire
- vector_add executed 10,000 times
- CUDA → PTX → SASS

```

1  __global__ void vector_add(int* result, int* op1, int* op2, int length) {
2      int idx = threadIdx.x + blockDim.x * blockIdx.x;
3      if (idx < length)
4          result[idx] = op1[idx] + op2[idx];
5  }
6
7  void main() {
8      int length = 10000;
9      int *h_vec1 = create_and_fill_array(length);
10     int *h_vec2 = create_and_fill_array(length);
11
12     int *d_vec1;
13     cudaMalloc(&d_vec1, sizeof(int) * length);
14     cudaMemcpy(d_vec1, h_vec1, sizeof(int) * length, cudaMemcpyHostToDevice);
15
16     int *d_vec2;
17     cudaMalloc(&d_vec2, sizeof(int) * length);
18     cudaMemcpy(d_vec2, h_vec2, sizeof(int) * length, cudaMemcpyHostToDevice);
19
20     int *d_result;
21     cudaMalloc(&d_result, sizeof(int) * length);
22
23     vector_add<<<length / 1024, 1024>>>(d_result, d_vec1, d_vec2, length);
24     cudaDeviceSynchronize();
25
26     int *h_result = (int*)malloc(sizeof(int) * length);
27     cudaMemcpy(h_result, d_result, sizeof(int) * length, cudaMemcpyDeviceToHost);
28 }
```

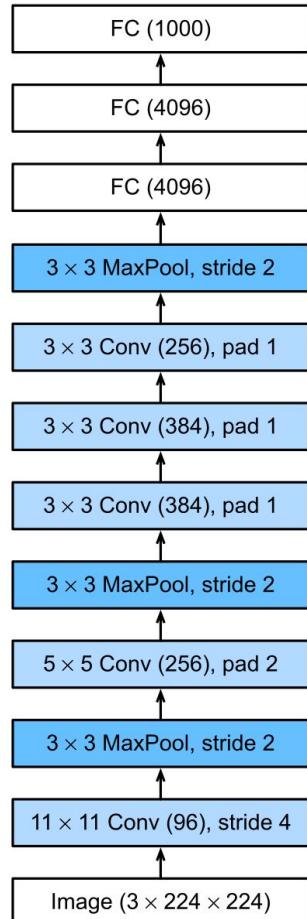
Graphics Processing Units

- Warps
 - Multiple lanes
- Latency hiding
 - Scheduler
 - Banked memory
 - MSHRs
- Tensor cores

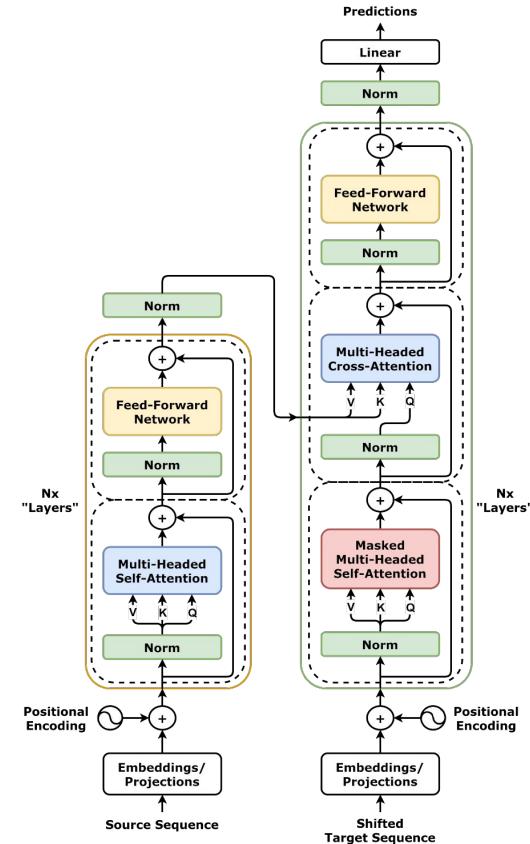


Deep Learning

- Deep neural networks
 - 1990: LeNet
 - 2012: AlexNet
 - 2014: VGG
 - 2015: ResNet
 - 2017: Transformer
- Compute-heavy
- Lots of composite functions



AlexNet



Transformer

Deep Learning Compilers

- Find optimized software representation
- Nonlinear problem
- Direct and indirect optimizations

```

for (n=0; n<N; n++)
    for (co=0; co<Co; co++)
        for (ci=0; ci<Ci; ci++)
            for (ho=0; ho<Ho; ho++)
                for (wo=0; wo<Wo; wo++)
                    for (hk=0; hk<Hk; hk++)
                        for (wk=0; wk<Wk; wk++)
                            Ofm[n][co][ho][wo] +=
                                Ifm[n][ci][ho+hk][wo+wk] * W[co][ci][hk][wk]

```

```

for (n=0; n<N; n++)
    for (coe=0; coe<Co; coe+=TCo)
        for (cie=0; cie<Ci; cie+=TCi)
            for (hoe=0; hoe<Ho; hoe+=THo)
                for (woe=0; woe<Wo; woe+=TWo)

```

Off-chip memory access pattern

```

for (co=coe; co<coe+TCo; co++)
    for (ci=cie; ci<cie+TCi; ci++)
        for (ho=hoe; ho<hoe+THo; ho++)
            for (wo=woe; wo<woe+TWo; wo++)
                for (hk=0; hk<Hk; hk++)
                    for (wk=0; wk<Wk; wk++)
                        Ofm[n][co][ho][wo] +=
                            Ifm[n][ci][ho+hk][wo+wk] * W[co][ci][hk][wk]

```

On-chip computation

Direct Optimizations

- Optimizations applied without profiling
- “Best guess”
 - Value propagation
 - Operator inlining
 - Strength reduction
 - Loop restructuring

```

1 TransformedIR:
2 /* define loop extents as variables for code brevity *:
3 extent_ko, extent_ki = (C_k / TB_tile_k), (TB_tile_k / Warp_tile_k)
4 /* Declare buffer size. */
5 alloc A_shared[3][...]
6 alloc A_reg[2][...]

7 /* Prologue for A_shared and A_reg */
8 for ko in 0 .. 2:
9   /* load into shared memory buffer (same as Line 15-17) */
10  for ki in 0 .. 1:
11    /* load into reg. buffer (same as Line 24-27) */

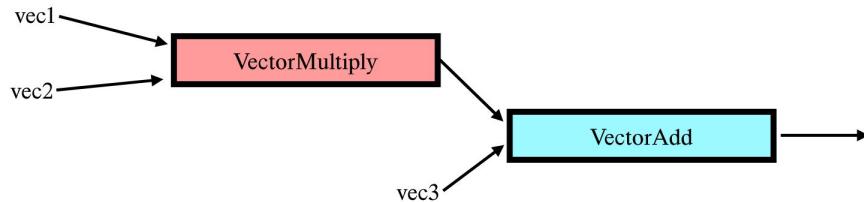
12 for ko in 0 .. extent_ko:
13   /* load into shared memory buffer */
14   /* guard data copy with producer primitives at Line 15 and Line 17 */
15   A_shared.producer_acquire()
16   async_memcpy(A_shared [ (ko + 2) % 3 ][...], A[...], (ko + 2) % extent_ko )
17   A_shared.producer_commit()

18 /* compute with data in shared memory buffer */
19 /* guard data usage with consumer primitives at Line 22 and Line 30 */
20 for ki in 0 .. extent_ki:
21   if ((ki + 1) % extent_ki )== 0:
22     A_shared.consumer_wait()
23     /* load into register buffer */
24     async_memcpy(
25       A_reg [ (ki + 1) % 2 ][...],
26       A_shared [ (ko + ((ki+1) / extent_ki) ) % 3 ][..., (ki + 1) % extent_ki ]
27     )

28 /* tensor-core compute with data in register buffer */
29 wmma(A_reg [ (ki % 2) ][...], ...)
30 A_shared.consumer_release()

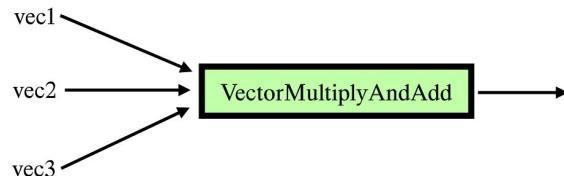
```

Value Propagation and Operator Inlining



```

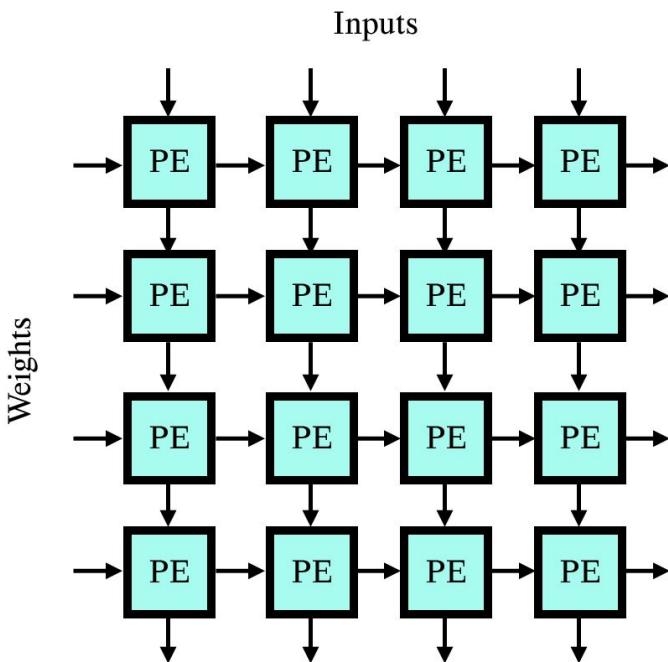
1 __global__ void VectorMultiply(int* result, int length, int* vec1, int* vec2) {
2     int idx = threadIdx.x + blockDim.x * blockIdx.x;
3     if (idx < length)
4         result[idx] = vec1[idx] * vec2[idx];
5 }
6
7 __global__ void VectorAdd(int* result, int length, int* vec1, int* vec2) {
8     int idx = threadIdx.x + blockDim.x * blockIdx.x;
9     if (idx < length)
10        result[idx] = vec1[idx] + vec2[idx];
11 }
  
```



```

1 __global__ void VectorMultiplyAndAdd(int* result, int* vec1, int* vec2, int* vec3, int length) {
2     int idx = threadIdx.x + blockDim.x * blockIdx.x;
3     if (idx < length)
4         result[idx] = vec1[idx] * vec2[idx] + vec3[idx];
5 }
  
```

Strength Reduction with Tensor Cores



```

1   for (int i = 0; i < n; i++)
2       for (int j = 0; j < m; j++)
3           for (int r = 0; r < k; r++)
4               C[i, j] += A[i, r] * B[r, j];

```

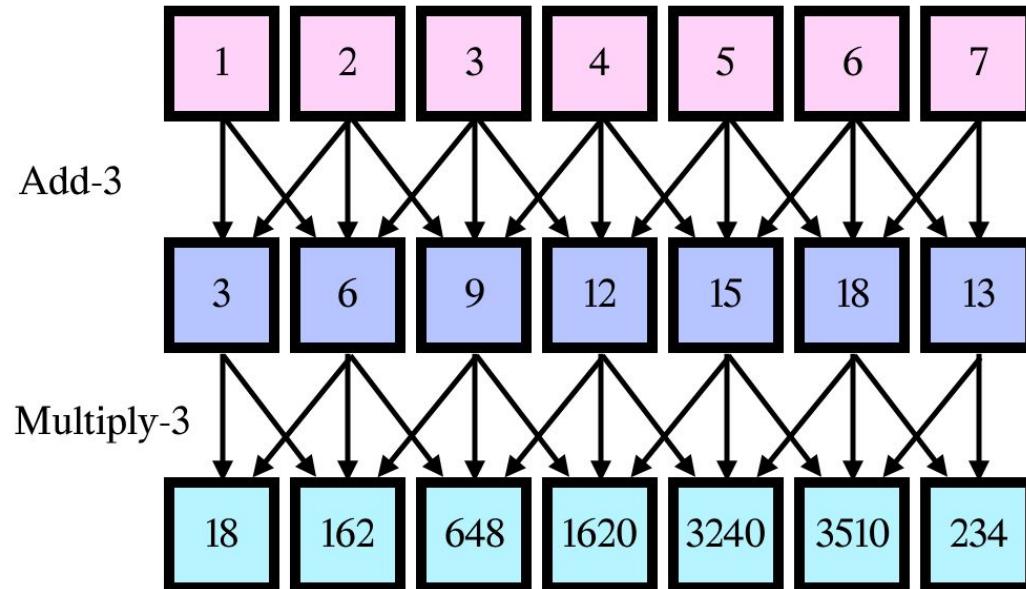
```

1   Buffer<fp16, 16, 16> tensorA;
2   Buffer<fp16, 16, 16> tensorB;
3   Buffer<fp32, 16, 16> tensorC;
4   for (int i = 0; i < n; i += 16) {
5       for (int j = 0; j < m; j += 16) {
6           for (int r = 0; r < k; r += 16) {
7               tensorA = Load(A[i:16, r:16]);
8               tensorB = Load(B[r:16, j:16]);
9               tensorC += TensorCore(tensorA, tensorB);
10          }
11      Store(C[i:16, j:16], tensorC);
12  }
13 }

```

Loop Restructuring

- Loop unrolling
- Vectorization
- Loop tiling



Polyhedral Model

- Describe loop with affine inequalities
 - Iteration domain forms lattice
 - Dependency graph links iterations
- Represent with mathematical expression
 - “Scattering” function describes statement instances
 - Create new scattering by solving equations
 - Different possible solutions
 - Scattering creates new code

```
for (i=0; i<N; i++)
  for (j=0; j<N; j++)
    S1: A[i,j] = A[i,j]+u1[i]*v1[j] + u2[i]*v2[j];

for (k=0; k<N; k++)
  for (l=0; l<N; l++)
    S2: x[k] = x[k]+A[l,k]*y[l];
    Original code
```

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -1 & 0 & 1 & -1 \\ 0 & -1 & 1 & -1 \end{pmatrix} \begin{pmatrix} i \\ j \\ l \\ N \\ 1 \end{pmatrix} \geq \vec{0}$$

Domain of S1



Data dependence graph

	<i>S1</i>	<i>const</i>	<i>S2</i>	<i>const</i>	
<i>c</i> ₁	0	1	0	1	0
<i>c</i> ₂	1	0	0	0	1
<i>c</i> ₃	0	0	0	0	1

Statement-wise transformation

```
for (c1=0; c1<N; c1++)
  for (c2=0; c2<N; c2++)
    A[c2,c1] = A[c2,c1]+u[c2]*v[c1];
    x[c1] = x[c1]+A[c2,c1]*y[c2];
```

Transformed code

Dependence polyhedron for S1→S2 edge

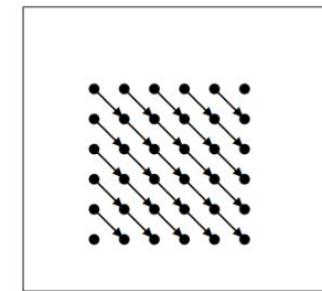
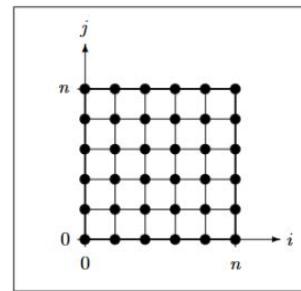
Algorithm 1 SCoP pseudocode

Require: $n, a[n], b[n], c[n]$

```

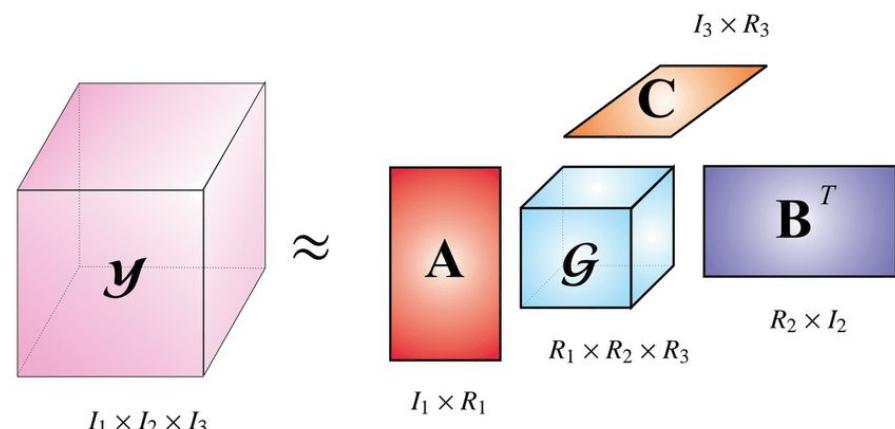
1: for  $i \leftarrow 0, n$  do
2:   for  $j \leftarrow n, 0$  do
3:     if  $i == 0 \parallel j == n$  then
4:        $c[i + j] \leftarrow a[i] \times b[j]$            ▷ S1
5:     else
6:        $c[i + j] \leftarrow c[i + j] + a[i] \times b[j]$    ▷ S2
7:     end if
8:   end for
9: end for

```



Tensor Operations

- Overparameterized networks
 - Tucker decomposition
 - “Core” tensor and matrices
 - Canonical polyadic decomposition
 - List of vectors
 - Dimensionality reduction
 - Clustering
 - Noise reduction
- Some operators are naturally tensors
 - Convolution
 - Multi-headed attention
 - Remove “flatten” operator



Indirect Optimizations

- An optimization may cause adverse effects
- Optimization space very large for many deep-learning programs
 - Local (loop-wide, function-wide) vs. global (program-wide)

```

1  for (int i = 0; i < 256; i++) {           // Loop "L1"
2      a[i] = 0;
3      for (int j = 0; j < 256; j++)        // Loop "L2"
4          a[i] += b[j + i] + c[j + i];
5  }

```

Program to optimize: L1 (outer), L2 (inner)

```

1  for (int i = 0; i < 256; i += 4) {
2      a[i] = 0;
3      a[i+1] = 0;
4      a[i+2] = 0;
5      a[i+3] = 0;
6      for (int j = 0; j < 256; j++)
7          a[i] += b[j + i] + c[j + i];
8      for (int j = 0; j < 256; j++)
9          a[i+1] += b[j + i + 1] + c[j + i + 1];
10     for (int j = 0; j < 256; j++)
11         a[i+2] += b[j + i + 2] + c[j + i + 2];
12     for (int j = 0; j < 256; j++)
13         a[i+3] += b[j + i + 3] + c[j + i + 3];
14 }

```

Unroll factor: L1 = 4, L2 = 1

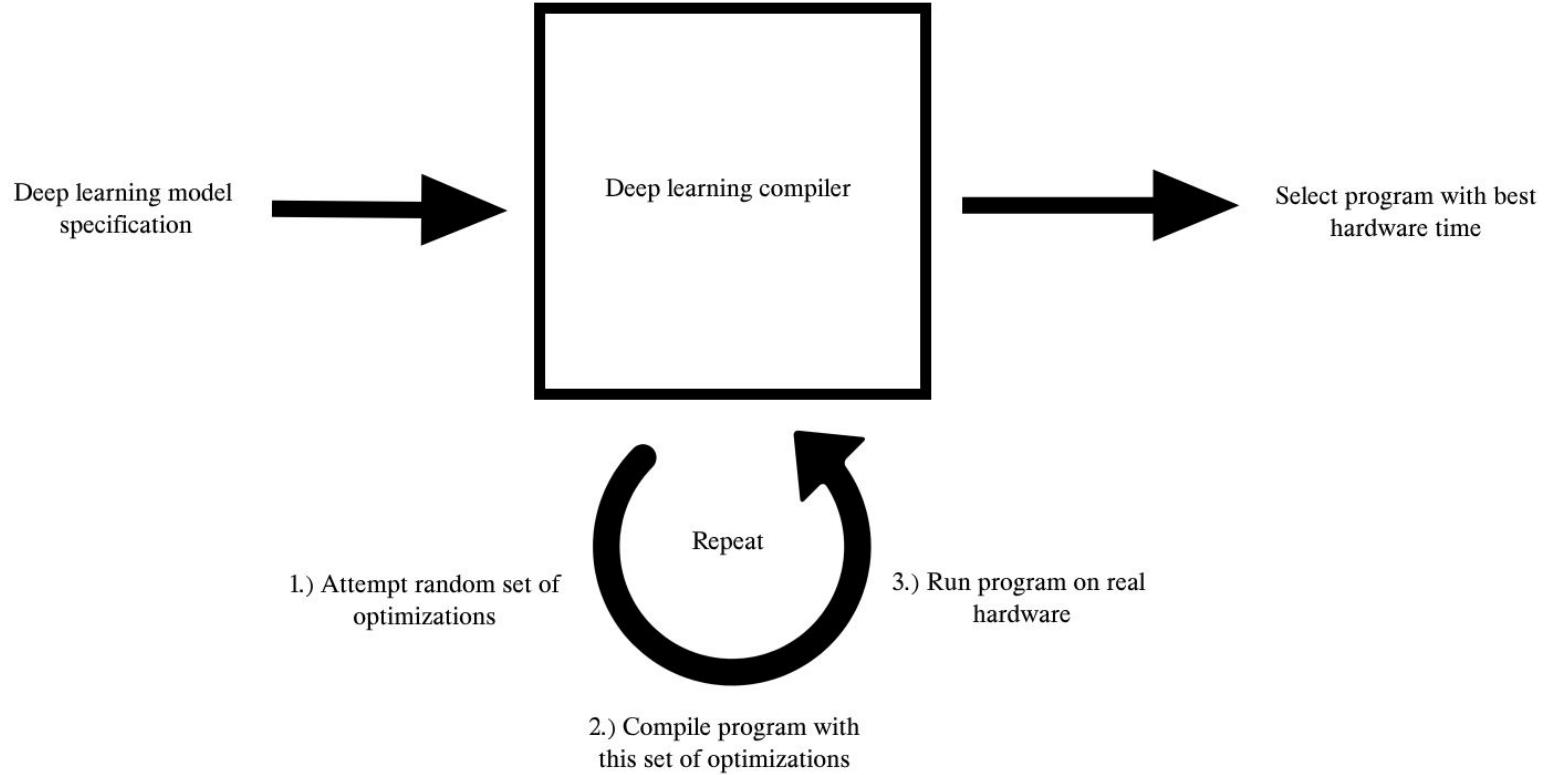
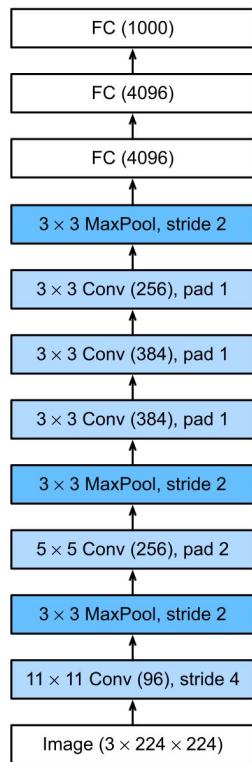
```

1  for (int i = 0; i < 256; i++) {
2      a[i] = 0;
3      for (int j = 0; j < 256; j += 8) {
4          a[i] += b[j + i] + c[j + i] +
5              b[j + 1 + i] + c[j + 1 + i] +
6              b[j + 2 + i] + c[j + 2 + i] +
7              b[j + 3 + i] + c[j + 3 + i] +
8              b[j + 4 + i] + c[j + 4 + i] +
9              b[j + 5 + i] + c[j + 5 + i] +
10             b[j + 6 + i] + c[j + 6 + i] +
11             b[j + 7 + i] + c[j + 7 + i];
12     }
13 }

```

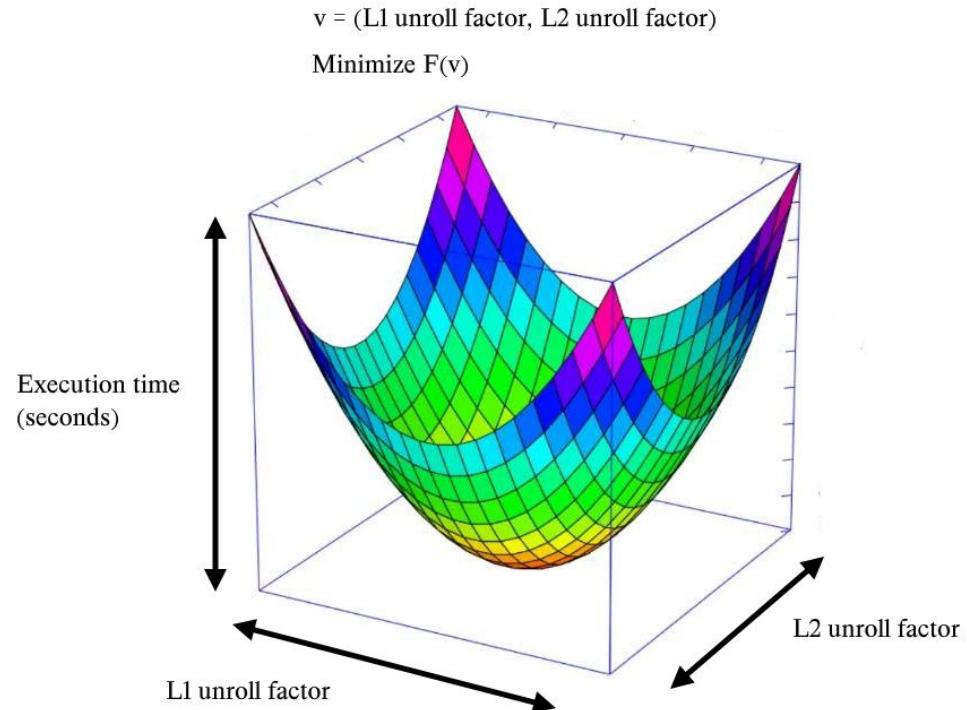
Unroll factor: L1 = 1, L2 = 8

Indirect Optimizations



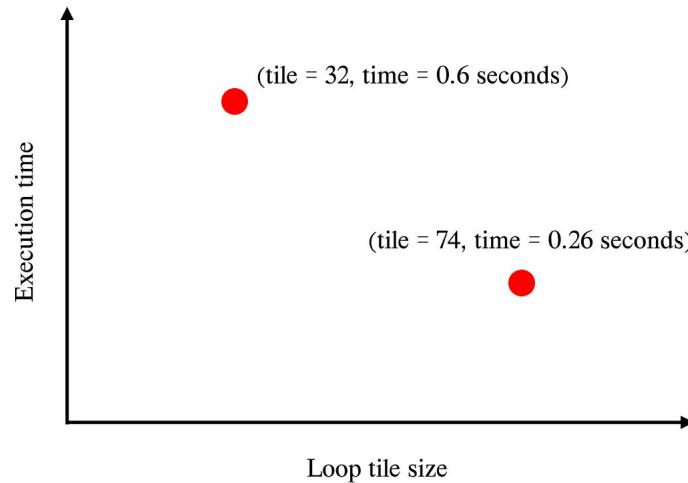
Optimization Pool

- Create pool of possible optimizations
- Form into vector
 - $|v_1 - v_2| < \epsilon$, optimizations are similar
- On-device measurement
 - Time in seconds = $F(v)$
- Approximate minimum of function
 - Random search
 - Grid search
 - Genetic algorithms
 - Simulated annealing
 - Deep learning



Genetic Algorithms and Simulated Annealing

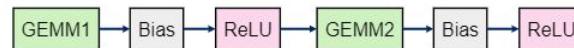
- Requires constraint that nearby points have nearby costs
- Device measurements are expensive
- Simulated annealing
 - Single-variable genetic
- Genetic algorithms
 - Selection
 - Crossover
 - Mutation



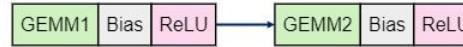
BOLT: Auto-tuner for Template Library Parameters

- Use template libraries (CUTLASS)
 - Instead of generating low-level CUDA code
- Design parameters:
 - Thread/block counts
 - Operator fusion
 - Data alignment
 - Synchronization

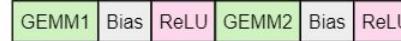
Non-fused:



Epilogue fusion:



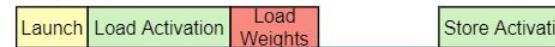
Persistent kernel fusion:



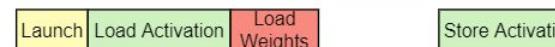
(a) The graph view of persistent kernel fusion

Non-fused:

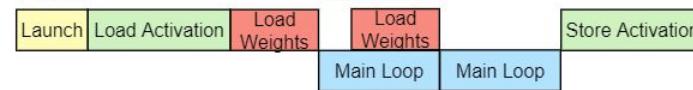
GEMM/conv 0



GEMM/conv 1



Fused:



(b) The kernel view of persistent kernel fusion

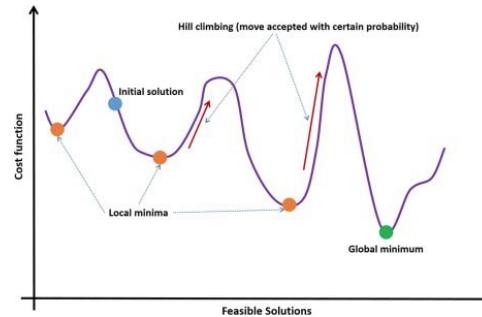
Halide: Graph Optimizer and Auto-Tuner for Image Pipelines

- Memory-limited stencils interconnected in a graph
 - Sharpen, blur, etc.
- Performance limited by interleaving allocation, execution, and communication
 - Separate algorithm (high-level operator) from schedule (low-level implementation) in domain-specific language
 - Automatically search

$\begin{aligned} \text{blurx}(x,y) &= \text{in}(x-1,y) + \text{in}(x,y) + \text{in}(x+1,y) \\ \text{out}(x,y) &= \text{blurx}(x,y-1) + \text{blurx}(x,y) + \text{blurx}(x,y+1) \end{aligned}$ <p>(Baseline) A 3x3 two-stage blur algorithm.</p>	
<pre>alloc blurx[2048][3072] for each y in 0..2048: for each x in 0..3072: blurx[y][x] = in[y][x-1] + in[y][x] + in[y][x+1] alloc out[2046][3072] for each y in 1..2047: for each x in 1..3072: out[y][x] = blurx[y-1][x] + blurx[y][x] + blurx[y+1][x]</pre> <p>(1) Compute and store every point in <code>blurx</code> before doing anything in the output. High parallelism, but little producer-consumer locality.</p>	<pre>alloc out[2046][3072] for each y in 1..2047: for each x in 0..3072: alloc blurx[-1..1] for each i in -1..1: blurx[i] = in[y+i][x-1]+in[y+i][x]+in[y+i][x+1] out[y][x] = blurx[0] + blurx[1] + blurx[2]</pre> <p>(2) Could compute each point in <code>blurx</code> immediately before the point that uses it. High parallelism and maximum locality, but lots of redundant work.</p>
<pre>alloc out[2046][3072] alloc blurx[3][3072] for each y in -1..2047: for each x in 0..3072: blurx[(y+1)%3][x]=in[y+1][x-1]+in[y+1][x]+in[y+1][x+1] if y < 1: continue out[y][x] = blurx[(y-1)%3][x] + blurx[y % 3][x] + blurx[(y+1)%3][x]</pre> <p>(3) Could interleave the two stages over a sliding window, wasting no work but no opportunities for parallelization.</p>	<pre>alloc out[2046][3072] for each ty in 0..2048/8: alloc blurx[-1..1][3072] for y in -2..8: for x in 0..3072: blurx[(y+1)%3][x] = in[ty*8+y+1][tx*32+x-1] + in[ty*8+y+1][tx*32+x] + in[ty*8+y+1][tx*32+x+1] if y < 0: continue for x in 0..3072: out[ty*8+y][x] = blurx[(y-1)%3][x] + blurx[y % 3][x] + blurx[(y+1)%3][x]</pre> <p>(4) Each version 1-3 had a pitfall. The optimal solution is somewhere in the middle, like sliding a window over strips of independent scanlines. Hardware independent, so this is searched for automatically.</p>

Apache TVM

- Relay
 - High-level IR
 - Computation graph
 - Optimizations allow graph rewriting
- Tensor Expression
 - Low-level IR
 - “Template” defined from placeholders, knobs, tuners, parallelism, vectorization
 - Search space + cost model
 - AutoTVM
 - Simulated annealing



e compute expression

```
A = t.placeholder((1024, 1024))
B = t.placeholder((1024, 1024))
k = t.reduce_axis((0, 1024))
C = t.compute((1024, 1024),
  lambda y, x:
    t.sum(A[k, y] * B[k, x], axis=k))
```

x₀ default code

```
for y in range(1024):
    for x in range(1024):
        C[y][x] = 0
        for k in range(1024):
            C[y][x] += A[k][y] * B[k][x]
```

s₁ loop tiling

```
yo, xo, yi, xi, s1 = s[C].title(y, x, ty, tx)
s[C].reorder(yo, xo, k, yi, xi)
```

x₁ = *g*(*e*, *s₁*)

```
for yo in range(1024 / ty):
    for xo in range(1024 / tx):
        C[yo*ty:yo*ty+ty][xo*tx:xo*tx+tx] = 0
        for k in range(1024):
            for yi in range(ty):
                for xi in range(tx):
                    C[yo*ty+yi][xo*tx+xi] +=
                        A[k][yo*ty+yi] * B[k][xo*tx+xi]
```

s₂ tiling, map to micro kernel intrinsics

```
yo,xo,ko,yi,xi,ki = s[C].title(y,x,k,8,8,8)
s[C].tensorize(yi, intrin.gemm8x8)
```

x₂ = *g*(*e*, *s₂*)

```
for yo in range(128):
    for xo in range(128):
        intrin.fill_zero(C[yo*8:yo*8+8][xo*8:xo*8+8])
        for ko in range(128):
            intrin.fused_gemm8x8_add(
                C[yo*8:yo*8+8][xo*8:xo*8+8],
                A[ko*8:ko*8+8][yo*8:yo*8+8],
                B[ko*8:ko*8+8][yo*8:yo*8+8])
```

Anstor: Hierarchical Optimizer

- Rule-based kernel expansion
 - Input: mathematical notation
 - Evolutionary search
 - Rules generate sketches (high-level)
 - Don't specify tile dimensions, parallelization, or unrolling
 - Randomly label loops as parallel or vectorized; randomly assign tile dimensions; genetic algorithm iterates

Example Input 1:

* The mathematical expression:

$$C[i, j] = \sum_k A[i, k] \times B[k, j]$$

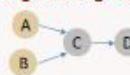
$$D[i, j] = \max(C[i, j], 0.0)$$

where $0 \leq i, j, k < 512$

* The corresponding naive program:

```
for i in range(512):
    for j in range(512):
        for k in range(512):
            C[i, j] += A[i, k] * B[k, j]
for i in range(512):
    for j in range(512):
        D[i, j] = max(C[i, j], 0.0)
```

* The corresponding DAG:



Generated sketch 1

```
for i.0 in range(TILE_I0):
    for j.0 in range(TILE_J0):
        for i.1 in range(TILE_I1):
            for j.1 in range(TILE_J1):
                for k.0 in range(TILE_K0):
                    for i.2 in range(TILE_I2):
                        for j.2 in range(TILE_J2):
                            for k.1 in range(TILE_I1):
                                for i.3 in range(TILE_I3):
                                    for j.3 in range(TILE_J3):
                                        C[...] += A[...] * B[...]
for i.4 in range(TILE_I2 * TILE_I3):
    for j.4 in range(TILE_J2 * TILE_J3):
        D[...] = max(C[...], 0.0)
```

No	Rule Name
1	Skip
2	Always Inline
3	Multi-level Tiling
4	Multi-level Tiling with Fusion
5	Add Cache Stage
6	Reduction Factorization
...	User Defined Rule

Sampled program 1

```
parallel i.0@j.0@i.1@j.1 in range(256):
    for k.0 in range(32):
        for i.2 in range(16):
            unroll k.1 in range(16):
                unroll i.3 in range(4):
                    vectorize j.3 in range(16):
                        C[...] += A[...] * B[...]
    for i.4 in range(64):
        vectorize j.4 in range(16):
            D[...] = max(C[...], 0.0)
```

Layout Structure Organization

- TASO: substitute subgraphs using mathematical properties of operators
 - Theorem prover
- DNNFusion: based on rules and categories
 - Different fusion/codegen rules depending on label (one-to-one, etc.)
 - Check associativity/distributivity/commutativity of each operator
- Apollo: rule-based under polyhedral model

Table 3. Mapping type analysis. The first column and the first row (both without color) show the mapping types of first and second operators, respectively, before fusion, and the colored cells show the mapping type of the operator after fusion. Green implies that these fusion combinations can be fused directly (i.e., they are profitable). Red implies that these fusions are unprofitable. Yellow implies that further profiling is required to determine profitability.

Second op First op	One-to-One	One-to-Many	Many-to-Many	Reorganize	Shuffle
One-to-One	One-to-One	One-to-Many	Many-to-Many	Reorganize	Shuffle
One-to-Many	One-to-One	One-to-Many	X	One-to-Many	One-to-Many
Many-to-Many	Many-to-Many	Many-to-Many	X	Many-to-Many	Many-to-Many
Reorganize	Reorganize	One-to-Many	Many-to-Many	Reorganize	Reorganize
Shuffle	Shuffle	One-to-Many	Many-to-Many	Reorganize	Shuffle

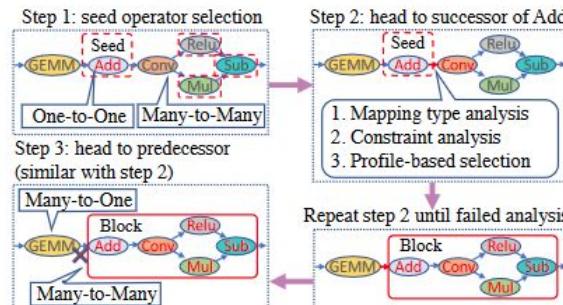
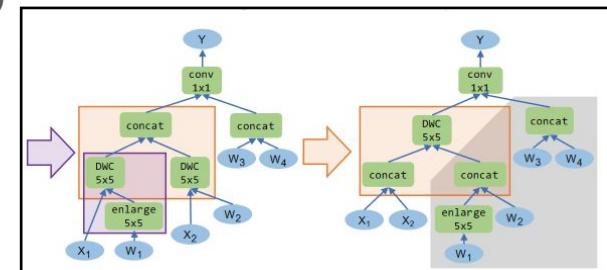
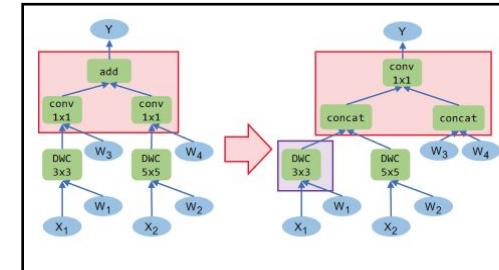


Figure 3. An example of fusion plan exploration. Assume Add, Conv, Relu, Mul, and Sub have identical output shape and IRS size.



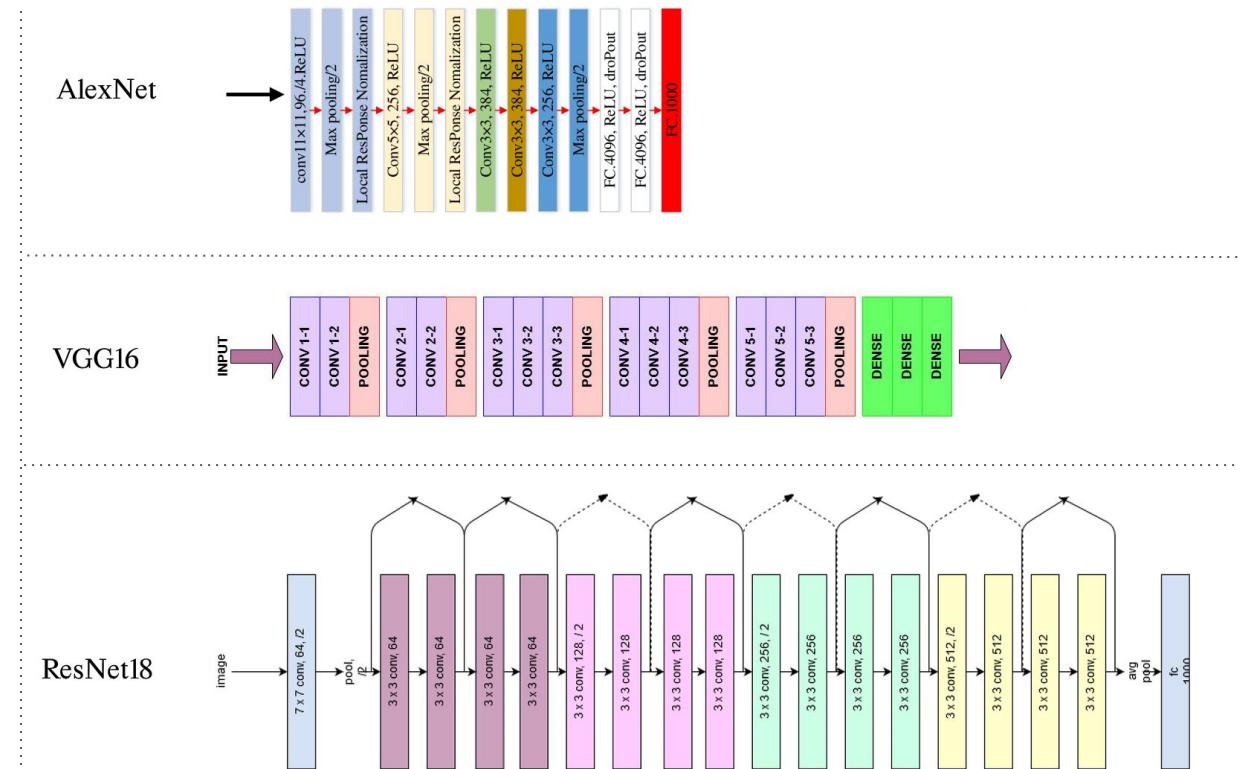
$$\forall s, p, x, y, z. \text{ conv}(s, p, \text{A}_{\text{none}}, \text{ewadd}(x, y), z) = \text{ewadd}(\text{conv}(s, p, \text{A}_{\text{none}}, x, z), \text{conv}(s, p, \text{A}_{\text{none}}, y, z))$$

Example mathematical property, asserts convolution without activation is linear in first argument.

- Z. Jia, O. Padon, J. Thomas, T. Warszawski, M. Zaharia, and A. Aiken, "TASO: Optimizing deep learning computation with automatic generation of graph substitutions," in Proceedings of the 27th ACM Symposium on Operating System Principles (SOSP), 2019, pages 47–62.
- J. Zhao, X. Gao, R. Xia, Z. Zhang, D. Chen, L. Chen, R. Zhang, Z. Geng, B. Cheng, and X. Jin, "Apollo: Automatic partition-based operator fusion through layer by layer optimization," in Proceedings of the 5th MLSys Conference (MLSys), 2022.
- W. Niu, J. Guan, Y. Wang, G. Agrawal, and B. Ren, "DNNFusion: Accelerating deep neural networks execution with advanced operator fusion," in Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI), 2021.

Image Recognition Networks

- VGG-16 spends 93.1% of execution time running convolution



2D Convolution

Input Image

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

$$\begin{aligned}
 & 1 * 5 + 2 * 3 + 3 * 1 + \\
 & 6 * 6 + 7 * 2 + 8 * 8 + \\
 & 11 * 7 + 12 * 3 + 9 * 13 = 358
 \end{aligned}$$

Step 1

Filter

5	3	1
6	2	8
7	3	9

*

Step 2

Input Image

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

Filter

5	3	1
6	2	8
7	3	9

*

$$\begin{aligned}
 & 2 * 5 + 3 * 3 + 4 * 1 + \\
 & 7 * 6 + 8 * 2 + 9 * 8 + \\
 & 12 * 7 + 13 * 3 + 14 * 9 = 402
 \end{aligned}$$

Lowering

Input Image

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

Filter

5	3	1
6	2	8
7	3	9

*

1	2	3	6	7	8	11	12	13
---	---	---	---	---	---	----	----	----

*

5	3	1	6	2	8	7	3	9
---	---	---	---	---	---	---	---	---

Input Image (first iteration)

Filter (first channel)

Lowering and Duplication

1	2	3	6	7	8	11	12	13
2	3	4	7	8	9	12	13	14
3	4	5	8	9	10	13	14	15
6	7	8	11	12	13	16	17	18
7	8	9	12	13	14	17	18	19
8	9	10	13	14	15	18	19	20
11	12	13	16	17	18	21	22	23
12	13	14	17	18	19	22	23	24
13	14	15	18	19	20	23	24	25

*

5
3
1
6
2
8
7
3
9

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

5	3	1
6	2	8
7	3	9

- ResNet18 is 14.5x faster, but has 8.4x more memory transactions

Duplo: Reducing Redundant Memory Accesses

- Repetition occurs in predictable locations, if we know convolution parameters
- Register renaming for WMMA registers
 - Load history buffer
 - Proceed to L2 if not exists in L1, else rename

$$dprod_{total} = dprod_{row} \times dprod_{col} \cdot input_{channels}$$

$$dprod_{row} = 1 + \frac{input_{rows} - filter_{rows}}{stride_{rows}}$$

$$dprod_{col} = 1 + \frac{input_{col} - filter_{cols}}{stride_{cols}}$$

$$dprod_{total} = dprod_{row} \times dprod_{col}$$

$$workspace_{index} = \frac{address - input_{start\ address}}{size_{data\ type}}$$

$$workspace_{row} = workspace_{index} / dprod_{size}$$

$$workspace_{col} = workspace_{index} \bmod dprod_{size}$$

$$element_id = F(address)$$

$$filter_{row} = (workspace_{row} / dprod_{col}) \times stride_{rows}$$

$$filter_{col} = (workspace_{row} \bmod dprod_{col}) \times stride_{cols}$$

$$relative_{row} = (workspace_{col} / filter_{cols}) \bmod filter_{rows}$$

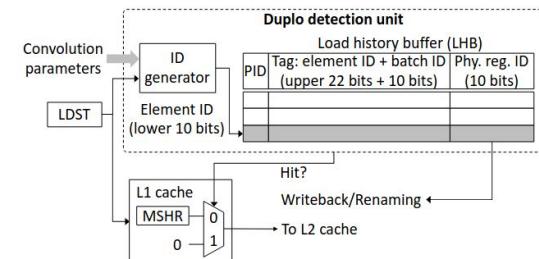
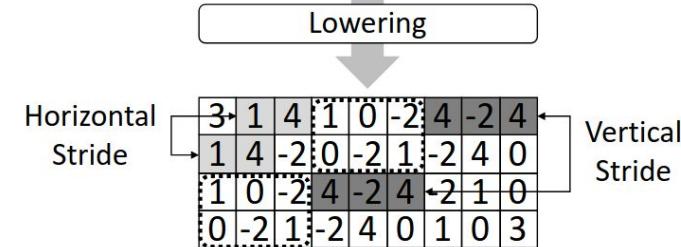
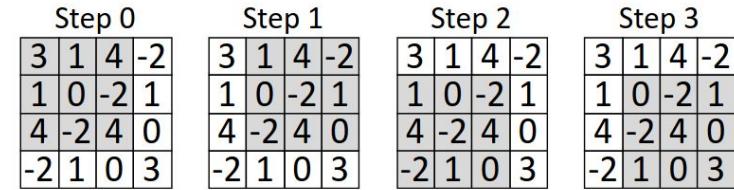
$$relative_{col} = workspace_{col} \bmod filter_{cols}$$

$$channel = workspace_{col} / (filter_{rows} \times filter_{cols})$$

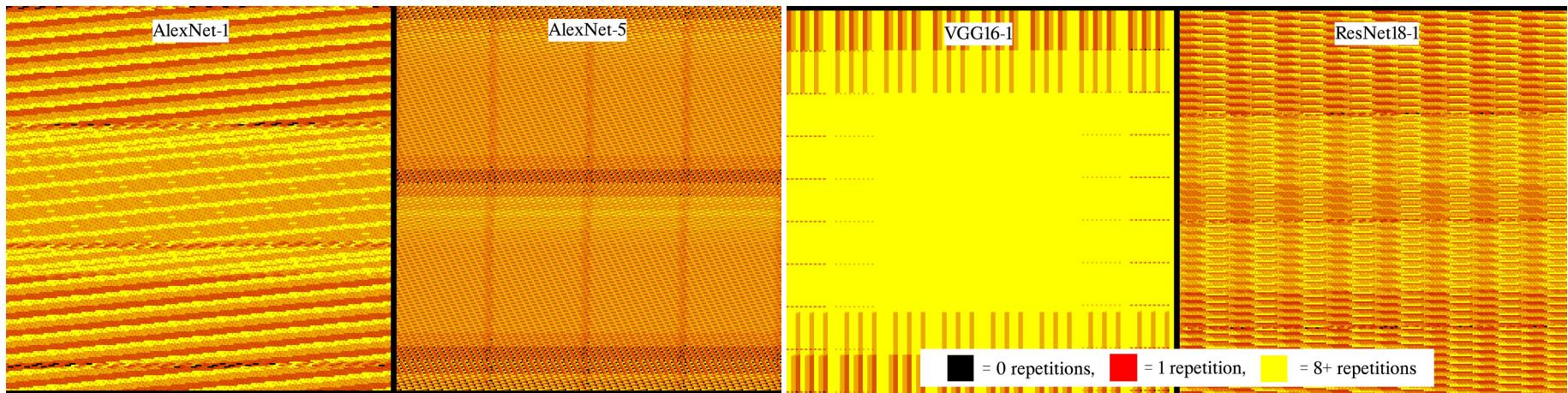
$$row = filter_{row} + relative_{row}$$

$$col = filter_{col} + relative_{col}$$

$$id = channel \times (input_{rows} \times input_{cols}) + (row \times input_{cols}) + col$$



Effect on the L2 Cache



Channels: 3

Filter size: 11x11

Stride: 4

Channels: 256

Filter size: 3x3

Stride: 1

Channels: 3

Filter size: 3x3

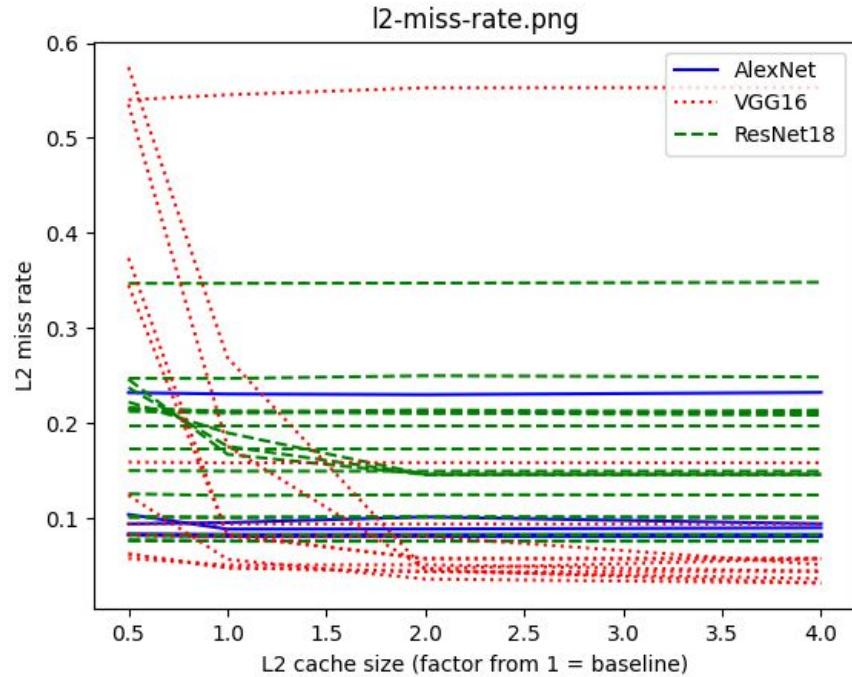
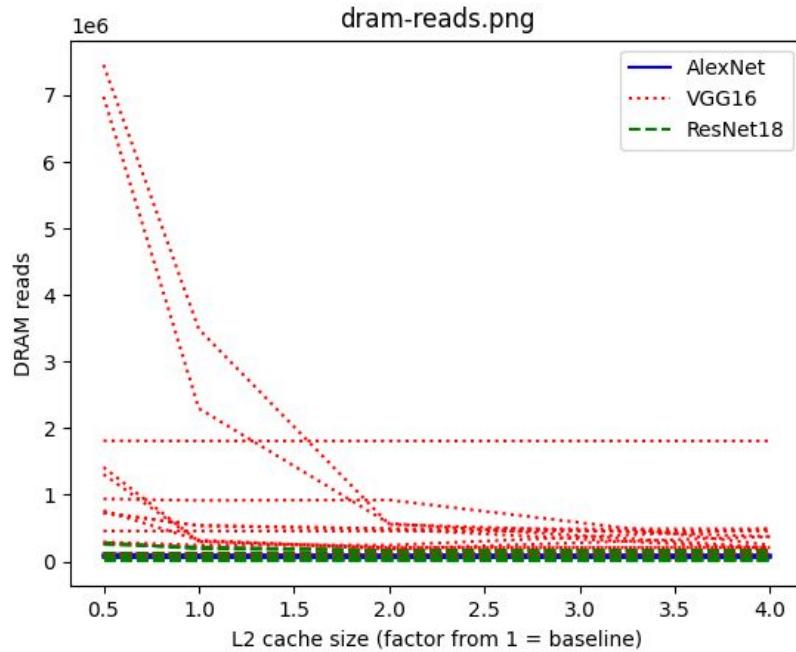
Stride: 1

Channels: 3

Filter size: 7x7

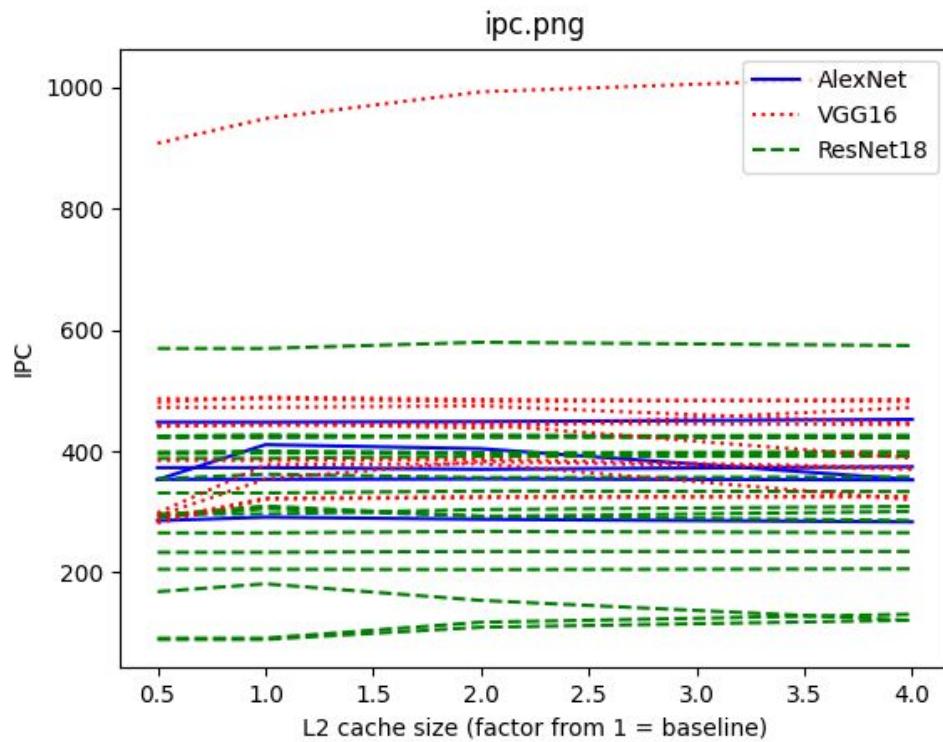
Stride: 2

Effect on the L2 Cache



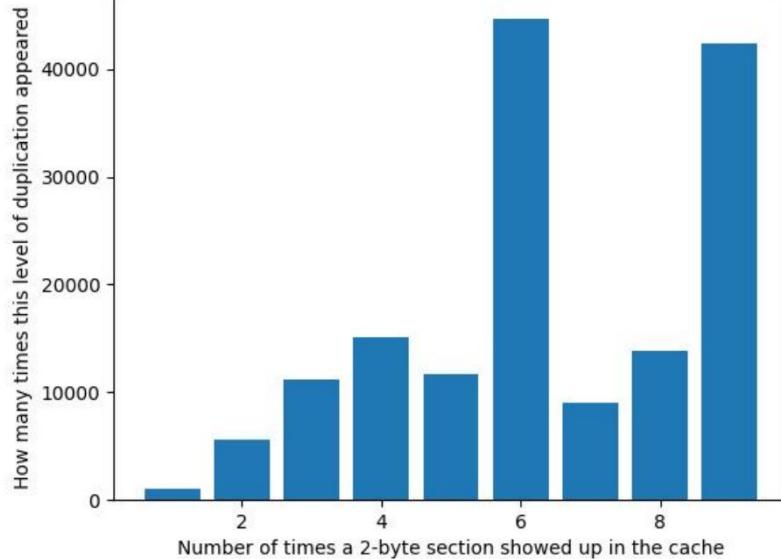
Effect on System Performance

- Increasing cache size will not improve performance
- Instead, halve the cache size
 - Prevent a decrease in IPC

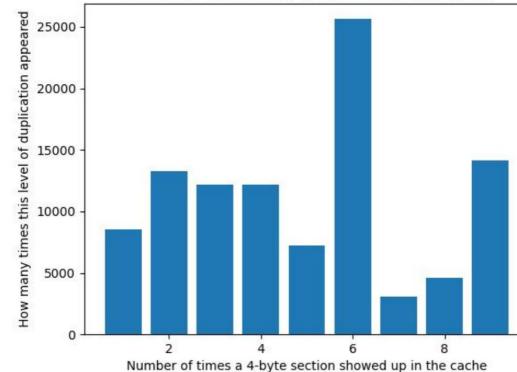


Subsequent Duplication

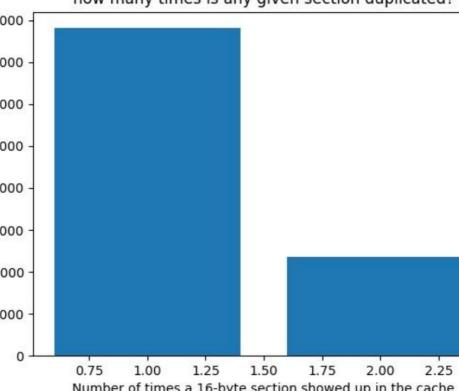
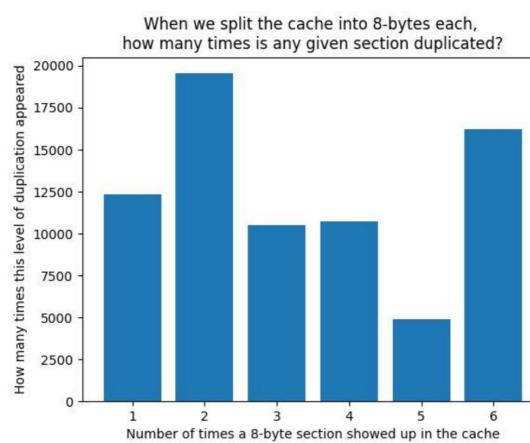
When we split the cache into 2-bytes each,
how many times is any given section duplicated?



When we split the cache into 4-bytes each,
how many times is any given section duplicated?

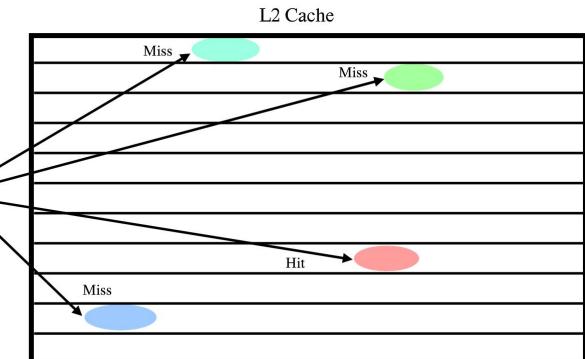


When we split the cache into 16-bytes each,
how many times is any given section duplicated?



Duplication Patterns

- Invert element-ID function
 - Normal: $\text{element_id} = F(\text{address})$
 - Inverse: $\text{address}[] = F(\text{element_id})$



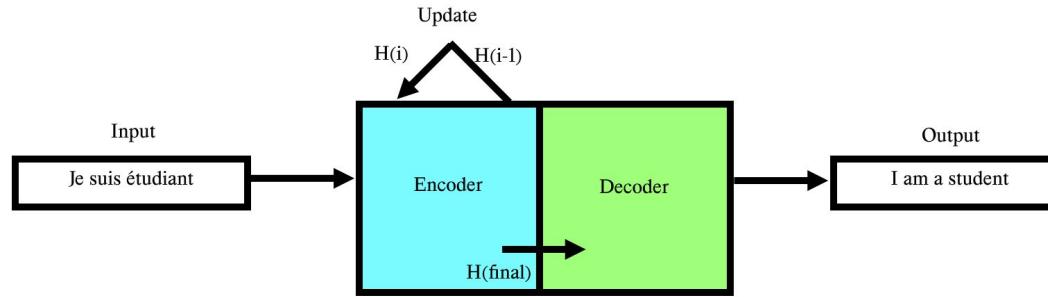
```

185 * Performs the inverse of the "g_workspace_indices" function. Given an
186 * element_id pointing somewhere in the lowered matrix, find an address from
187 * the lowered matrix (between g_B_start_ptr and g_B_end_ptr) that would map
188 * to the same element_id if it were put through "g_workspace_indices" again.
189 */
190 bool g_inverse_workspace(int element_id, new_addr_type** addrs, int* length) {
191     // -1 is an indicator of a value in the lowered matrix that points to some
192     // padded area, so it's not something we can return from this function.
193     if (element_id == -1)
194         return 0;
195
196     // We can consider the 0 padding as a part of the matrix in this step. This
197     // should be considered "padding" in the same way the 16-element-padding
198     // is considered "padding"; thus, it should ideally have an element_id equal
199     // to -1.
200     int conv_input_rows = g_conv_input_rows + g_conv_padding_rows * 2;
201     int conv_input_cols = g_conv_input_cols + g_conv_padding_cols * 2;
202
203     // Before we do stuff specifically related to this function, obtain the
204     // general info for workspace calculations. This part is the same as the
205     // main "g_workspace_indices" function.
206     // How many elements are in one dot product.
207     int dprod_size = g_conv_filter_rows * g_conv_filter_cols *
208         g_conv_input_channels;
209
210     // How many dot products are in a row/column/filter.
211     int dprods_per_row = (conv_input_rows - g_conv_filter_rows) /
212         g_conv_stride_rows + 1;
213     int dprods_per_col = (conv_input_cols - g_conv_filter_cols) /
214         g_conv_stride_cols + 1;
215     int num_dprods = dprods_per_row * dprods_per_col;
216
217     // The array is padded by a multiple of 16, while the actual lowered matrix
218     // is not. While we don't need to know if "element_id" points to somewhere
219     // outside of the workspace (since we already know the element_id is not -1),
220     // we still need to know the padding for address calculations.
221     int pad_amount_row = (16 - dprod_size % 16) % 16;
222     int row_size = dprod_size + pad_amount_row;
223     int pad_amount_col = (16 - num_dprods % 16) % 16;
224     int col_size = num_dprods + pad_amount_col;
225
226     // Obtain the original components of the 3d vector that make up the position.
227     // These three values denote a single location in the convolution tensor.
228     int ichab = element_id / (conv_input_rows * conv_input_cols);
229     element_id = element_id % (conv_input_rows * conv_input_cols);
230     int rowab = element_id / conv_input_cols;
231     element_id = element_id % conv_input_cols;
232     int colab = element_id;
233
234     // There are multiple workspace rows/columns that will contain this position,
235     // so locate the first one (the dot-product operation that would have
236     // occurred first, or the one with the smallest row index). Each counter here
237     // represents how many dot-products in the same row/column would need to
238     // occur before we have the first dot-product that hits the element. For
239     // example, if row counter = 5 and col counter = 7, then the dot product
240     // happening on iteration (col_counter * dprods_per_row + row_counter)
241     // would contain this element.
242     int col_counter = max(
243         colab - g_conv_filter_cols + g_conv_stride_cols /
244             g_conv_stride_cols, 0);
245     int row_counter = max(
246         rowab - g_conv_filter_rows + g_conv_stride_rows /
247             g_conv_stride_rows, 0);
248     int dprod_counter = row_counter * dprods_per_row + col_counter;
249
250     // We need to identify the position within the filter the first occurrence
251     // of the element_id is within the dot products.
252     int crel = colab - col_counter * g_conv_stride_cols;
253     int rrel = rowab - row_counter * g_conv_stride_rows;
254
255     // Now we can construct workspace rows/columns from these values. The first
256     // workspace row is the dprod_counter, and the workspace_col is the flattened
257     // relative position.
258     int workspace_row_0 = dprod_counter;
259     int workspace_col_0 = rrel * g_conv_filter_cols + crel;
260
261     // Now we can construct workspace rows/columns from these values. The first
262     // workspace row is the dprod_counter, and the workspace_col is the flattened
263     // relative position.
264     int workspace_row_0 = dprod_counter;
265     int workspace_col_0 = rrel * g_conv_filter_cols + crel;
266
267     // Now we can construct workspace rows/columns from these values. The first
268     // workspace row is the dprod_counter, and the workspace_col is the flattened
269     // relative position.
270     int workspace_row_0 = dprod_counter;
271     int workspace_col_0 = rrel * g_conv_filter_cols + crel;
272
273     // Each dot-product is put into the lowered B matrix in order, so
274     // contains this element. We can calculate how many times dot-products will
275     // contain this element in the same row, and again for the same column.
276     int repeat_col = ceil((float)(crel + 1) / g_conv_stride_cols);
277     int repeat_row = ceil((float)(rrel + 1) / g_conv_stride_rows);
278
279     // Finally, the list of workspace rows/columns that contain this element is
280     // equal to (repeat_row * repeat_col). We can iterate across each of them
281     // and convert their workspace positions into real byte positions easily.
282     *length = repeat_row * repeat_col;
283     *addr = (new_addr_type*)malloc(sizeof(int) * length);
284     int workspace_row_current = workspace_row_0;
285     int workspace_col_current = workspace_col_0;
286     for (int r = 0; r < row_counter; r++) {
287         for (int c = 0; c < col_counter; c++) {
288             // The variables "workspace_row.col_current" map to a position in the
289             // workspace that contains the same element_id as the parameter to this
290             // function. First, add back in the 16-element padding that got added
291             // to each dprod to make it a multiple of 16. Then, turn the idx (points
292             // to elements) into a byte offset (each element is a half, so two bytes)
293             // and add it to the base address (start of the B matrix). The result is
294             // the final address.
295             int widx = workspace_row_current * dprod_size + workspace_col_current;
296             int idx = widx + workspace_row_current * pad_amount_row;
297             new_addr_type byte_offset = idx * 2;
298             new_addr_type addr = byte_offset + g_B_start_ptr;
299             (*addr)[c * col_counter + c] = addr;
300
301             // Within a row (or, across columns), move the filter to the right by the
302             // stride. This makes the workspace_col decrease by the stride. The
303             // next dot-product will contain this element_id still, so increase the
304             // workspace_row by 1.
305             workspace_col_current -= g_conv_stride_cols;
306             workspace_row_current++;
307         }
308
309     }
310
311     // Once we've reached the end of the dot-products in this row that contain
312     // the element, we move on to the next row. The number of dot-products we
313     // pass between these points is equal to (dprods_per_row - repeat_row).
314     // The new position within the filter is equal to the original (first)
315     // filter position moved upwards by the stride amount.
316     workspace_row_current += dprods_per_row - repeat_row;
317     workspace_col_current = workspace_col_0 -
318         g_conv_stride_rows * g_conv_filter_cols * (1 + row_counter);
319
320     return 0;
321 }

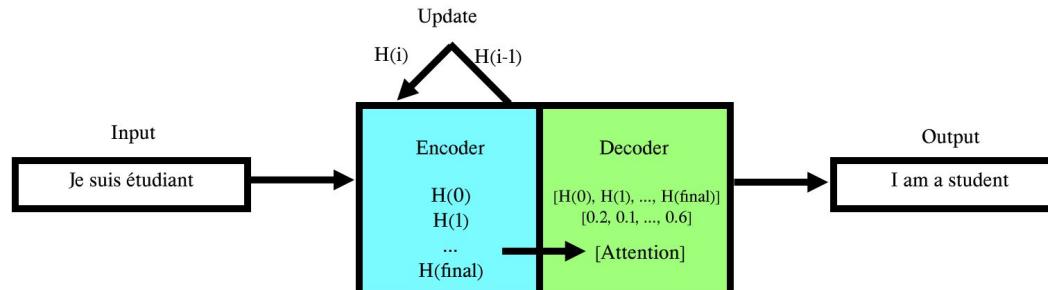
```

Large Language Models

- Previous models lost long-term information

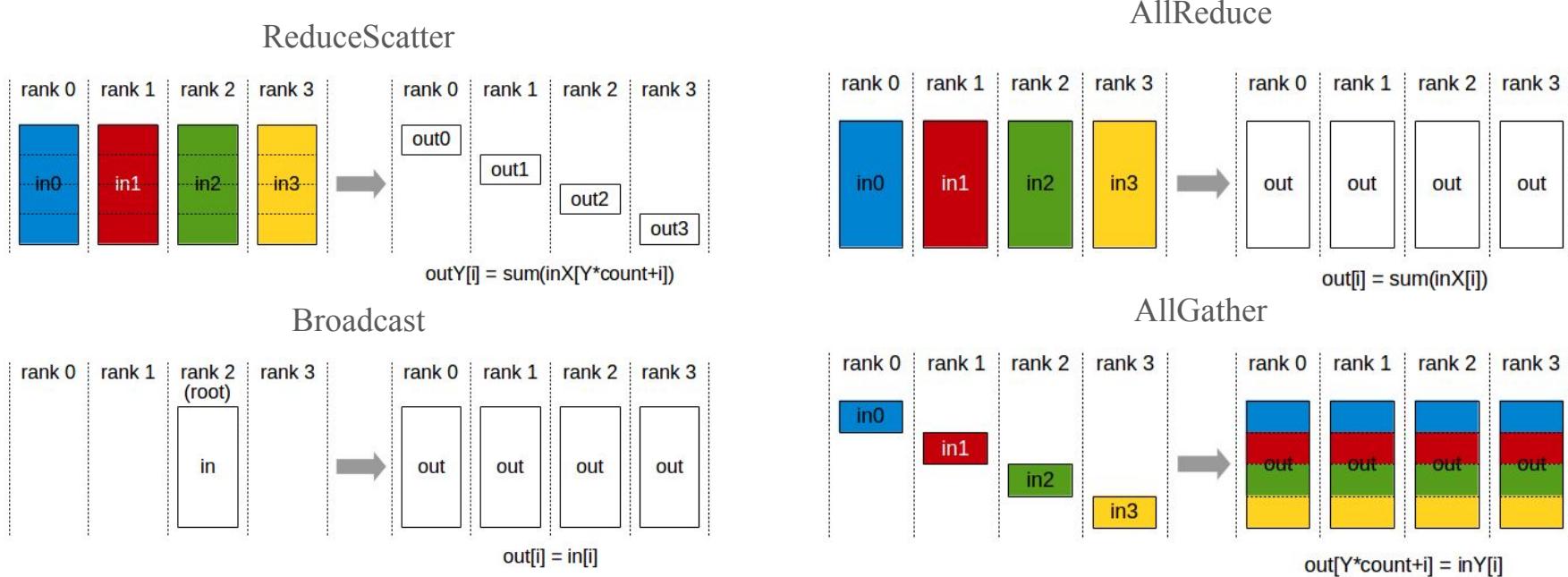


- New “attention” mechanism keeps all intermediate states, less likely to lose long-term information



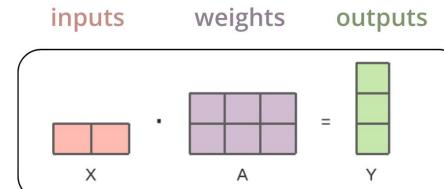
Interconnects

- NVLink: 1.8TB/s bidirectional, direct GPU-to-GPU interconnect
- Allows implementation of communication collectives



Large Language Models

- Megatron-LM
 - 462 billion parameters
 - Trained on 6144 H100 GPUs
- Parallelism modes
 - Data parallel: split batch
 - Tensor parallel: split weights
 - Pipeline parallel: split layers



is equivalent to

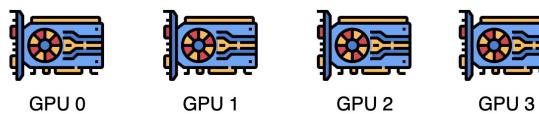
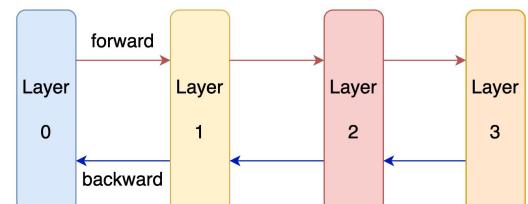
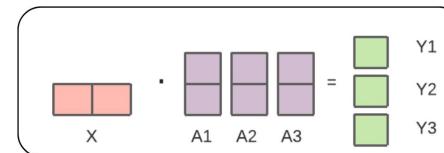
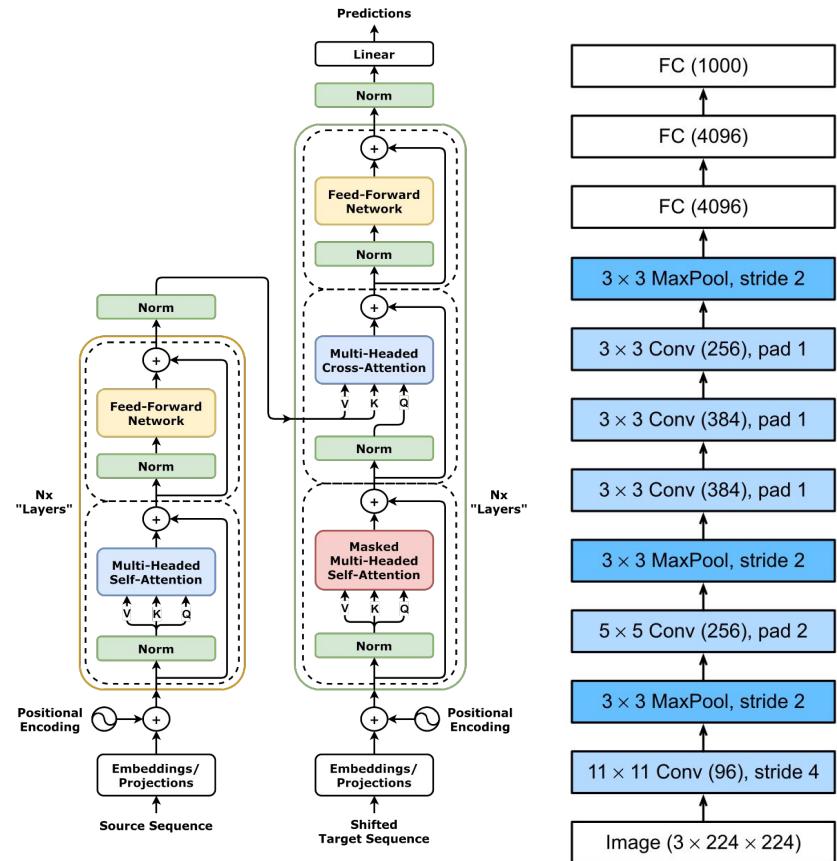


Image Net vs. LLM Comparison

- Significantly more kernels than image networks
 - AlexNet
 - 62 million parameters
 - 11 layers: 3 fully-connected, 5 conv2d
 - BERT
 - 110 million parameters
 - 12 encoders
 - 12 attention heads
 - 64 training examples / 2 epochs
 - 274,825 GEMM kernels (41.6 seconds)
 - 638,408 non-GEMM kernels (31.0 seconds)
 - E.g., data movement, elementwise
 - Time includes profiler



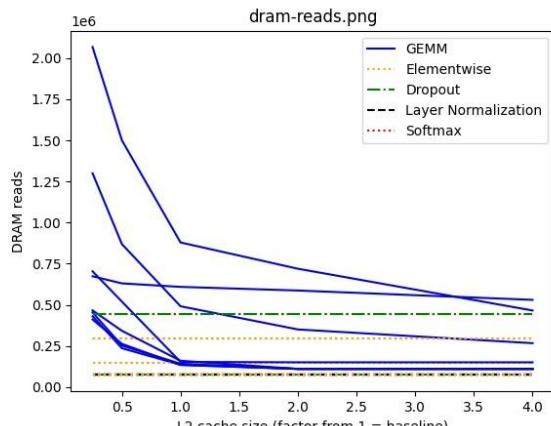
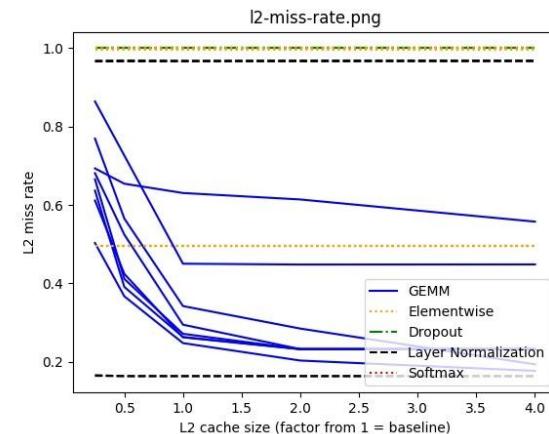
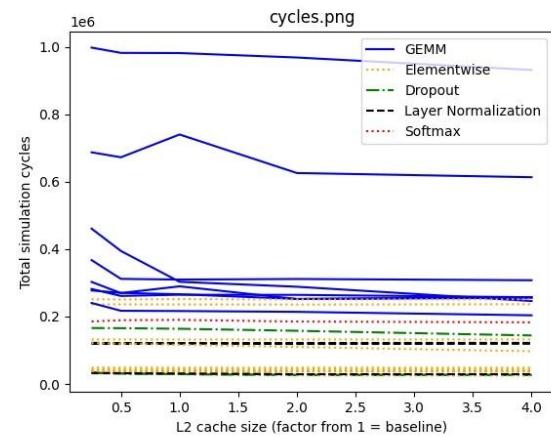
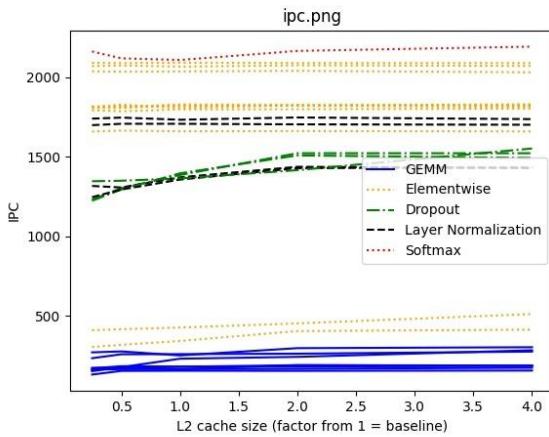
BERT Characteristics

- | | |
|-------------------------|-----------------|
| 1. GEMM | 16. GEMM |
| 2. Elementwise | 17. Elementwise |
| 3. Dropout | 18. GEMM |
| 4. Elementwise | 19. Elementwise |
| 5. Layer Normalization | 20. GEMM |
| 6. Layer Normalization | 21. Elementwise |
| 7. GEMM | 22. GEMM |
| 8. Elementwise | 23. Elementwise |
| 9. Elementwise | 24. Elementwise |
| 10. GEMM | 25. Softmax |
| 11. Elementwise | 26. Dropout |
| 12. Dropout | 27. GEMM |
| 13. Elementwise | 28. Elementwise |
| 14. Layer Normalization | |
| 15. Layer Normalization | |

- Repeats 12 times
- Memory access overlap:
 - GEMM/Elementwise: 39%
 - Elementwise/Dropout: 40%
 - Dropout/Elementwise: 57%
- Average repeat accesses:
 - GEMM: 6.2
 - Elementwise: 1.3
 - Dropout: 1.0
- Average total access count:
 - GEMM: 48,000
 - Elementwise: 6,100
 - Dropout: 7600
- Average unique access count:
 - GEMM: 7600
 - Elementwise: 3000
 - Dropout: 7600

L2 Size

- Does not improve with increasing L2 size
 - Similar to image nets



Conclusion

- Deep learning requires fast hardware and software
 - The unique properties of deep learning allows iterative compilation
 - Compute-heavy operations are good for GPUs
 - Communication libraries make distributed computing across GPUs easy