

Tensor Computing in Deep Learning

Weijie Zhao

02/20/2025

- Scalar
- Vector
- Matrix
- Tensor
 - Rank
 - Dimension

Tensor Computing in Deep Learning

Weijie Zhao

02/20/2025

- Scalar
- Vector
- Matrix
- Tensor
 - Rank
 - Dimension

Matrix ~~Tensor~~ Computing in Deep Learning

Weijie Zhao

02/20/2025

- Scalar
- Vector
- Matrix
- Tensor
 - Rank
 - Dimension

Matrix ~~Tensor~~ Computing in Deep Learning

Weijie Zhao

02/20/2025

- Matrix multiplication
- Non-linear activation
- Gradient descent

- Scalar
- Vector
- Matrix
- Tensor
 - Rank
 - Dimension

Matrix ~~Tensor~~ Computing in Deep Learning

Weijie Zhao

02/20/2025

- Matrix multiplication
- Non-linear activation
- ~~Gradient descent~~
Graduate student descent

Tensor Operations

- Element-wise add
- Element-wise plus
- Element-wise division
- Hadamard product
- Matrix multiplication
- Batched matrix multiplication
- More linear algebra operations...
- Collect, Scatter, Reduce...

Libraries

- Numpy
- Blas
- cuBlas
- cuSparse
- MKL
- TensorFlow
- PyTorch
- MXNet
- ...

Lazy Evaluation and Code Generation

$c = a + b$

$d = c * 2$

for $i = 1$ to n do

$c[i] = a[i] + b[i]$

for $i = 1$ to n do

$d[i] = c[i] * 2$

for $i = 1$ to n do

$d[i] = (a[i] + b[i]) * 2$

Graph Optimizer

- Graph minimization and canonicalization
 - Constant Folding
 - Common subexpression elimination
 - Remove unnecessary operations
- Algebraic simplification and reassociation
- Copy propagation

Meta Optimizer

```
i = 0
while i < config.meta_optimizer_iterations (default=2):
    Pruning()          # Remove nodes not in fanin of outputs, unused functions
    Function()         # Function specialization & inlining, symbolic gradient inlining
    DebugStripper() *
    ConstFold()        # Constant folding and materialization
    Shape()            # Symbolic shape arithmetic
    Remapper()         # Op fusion
    Arithmetic()       # Node deduping (CSE) & arithmetic simplification
    if i==0: Layout() # Layout optimization for GPU
    if i==0: Memory() # Swap-out/Swap-in, Recompute*, split large nodes
    Loop()             # Loop Invariant Node Motion*, Stack Push & Dead Node Elimination
    Dependency()       # Prune/optimize control edges, NoOp/Identity node pruning
    Custom()           # Run registered custom optimizers (e.g. TensorRT)
    i += 1
```

Constant Folding Optimizer

```
do:  
    InferShapesStatically() # Fixed-point iteration with symbolic shapes  
    graph_changed = MaterializeConstants() # grad broadcast, reduction dims  
    q = NodesWithKnownInputs()  
    while not q.empty():  
        node = q.pop()  
        graph_changed |= FoldGraph(node, &q) # Evaluate node on host  
        graph_changed |= SimplifyGraph()  
    while graph_changed
```

Constant Folding Optimizer: SimplifyGraph()

- Removes trivial ops, e.g. identity Reshape, Transpose of 1-d tensors, $\text{Slice}(x) = x$, etc.
- Rewrites that enable further constant folding
- Arithmetic rewrites that rely on known shapes or inputs, e.g.
 - Constant push-down:
 - $\text{Add}(c1, \text{Add}(x, c2)) \Rightarrow \text{Add}(x, c1 + c2)$
 - $\text{ConvND}(c1 * x, c2) \Rightarrow \text{ConvND}(x, c1 * c2)$
 - Partial constfold:
 - $\text{AddN}(c1, x, c2, y) \Rightarrow \text{AddN}(c1 + c2, x, y)$,
 - $\text{Concat}([x, c1, c2, y]) = \text{Concat}([x, \text{Concat}([c1, c2]), y])$
 - Operations with neutral & absorbing elements:
 - $x * \text{Ones}(s) \Rightarrow \text{Identity}(x)$, if $\text{shape}(x) == \text{output_shape}$
 - $x * \text{Ones}(s) \Rightarrow \text{BroadcastTo}(x, \text{Shape}(s))$, if $\text{shape}(s) == \text{output_shape}$
 - Same for $x + \text{Zeros}(s)$, $x / \text{Ones}(s)$, $x * \text{Zeros}(s)$ etc.
 - $\text{Zeros}(s) - y \Rightarrow \text{Neg}(y)$, if $\text{shape}(y) == \text{output_shape}$
 - $\text{Ones}(s) / y \Rightarrow \text{Recip}(y)$ if $\text{shape}(y) == \text{output_shape}$

Arithmetic Optimizer

```
DedupComputations () :
```

```
    do :
        stop = true
        UniqueNodes reps
        for node in graph.nodes():
            rep = reps.FindOrInsert (node, IsCommutative (node) )
            if rep == node or !SafeToDedup (node, rep):
                continue
            for fanout in node.fanout():
                ReplaceInputs (fanout, node, rep)
            stop = false
    while !stop
```

Arithmetic Optimizer

- Arithmetic simplifications
 - Flattening: $a+b+c+d \Rightarrow \text{AddN}(a, b, c, d)$
 - Hoisting: $\text{AddN}(x * a, b * x, x * c) \Rightarrow x * \text{AddN}(a+b+c)$
 - Simplification to reduce number of nodes:
 - Numeric: $x+x+x \Rightarrow 3*x$
 - Logic: $!(x > y) \Rightarrow x \leq y$
- Broadcast minimization
 - Example: $(\text{matrix1} + \text{scalar1}) + (\text{matrix2} + \text{scalar2}) \Rightarrow (\text{matrix1} + \text{matrix2}) + (\text{scalar1} + \text{scalar2})$
- Better use of intrinsics
 - $\text{Matmul}(\text{Transpose}(x), y) \Rightarrow \text{Matmul}(x, y, \text{transpose_x=True})$
- Remove redundant ops or op pairs
 - $\text{Transpose}(\text{Transpose}(x, \text{perm}), \text{inverse_perm})$
 - $\text{BitCast}(\text{BitCast}(x, \text{dtype1}), \text{dtype2}) \Rightarrow \text{BitCast}(x, \text{dtype2})$
 - Pairs of elementwise involutions $f(f(x)) \Rightarrow x$ (Neg, Conj, Reciprocal, LogicalNot)
 - Repeated Idempotent ops $f(f(x)) \Rightarrow f(x)$ (DeepCopy, Identity, CheckNumerics...)
- Hoist chains of unary ops at Concat/Split/SplitV
 - $\text{Concat}([\text{Exp}(\text{Cos}(x)), \text{Exp}(\text{Cos}(y)), \text{Exp}(\text{Cos}(z))]) \Rightarrow \text{Exp}(\text{Cos}(\text{Concat}([x, y, z])))$
 - $[\text{Exp}(\text{Cos}(y)) \text{ for } y \text{ in } \text{Split}(x)] \Rightarrow \text{Split}(\text{Exp}(\text{Cos}(x)), \text{num_splits})$

Differentiation

- Numerical differentiation
- Symbolic differentiation
 - Chain rules
- Forward mode auto differentiation
- Reverse mode auto differentiation