# Neural Network Preliminary Tensor Computing

Weijie Zhao

10/13/2022

# Neural Network Preliminary Tensor Computing

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

Weijie Zhao

10/13/2022

# Neural Network Preliminary

## ~~Tensor~~ Computing

## <span style="color:red">Matrix</span>

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

Weijie Zhao

10/13/2022

•Matrix multiplication

•Non-linear activation

•Gradient descent

•Scalar

•Vector

•Matrix

•Tensor

  •Rank

  •Dimension

# Neural Network Preliminary

## ~~Tensor~~ Computing

## Matrix

Weijie Zhao

10/13/2022

- Matrix multiplication

- Non-linear activation

- ~~Gradient descent~~  Graduate student descent

- Scalar

- Vector

- Matrix

- Tensor
  - Rank
  - Dimension

# Neural Network Preliminary

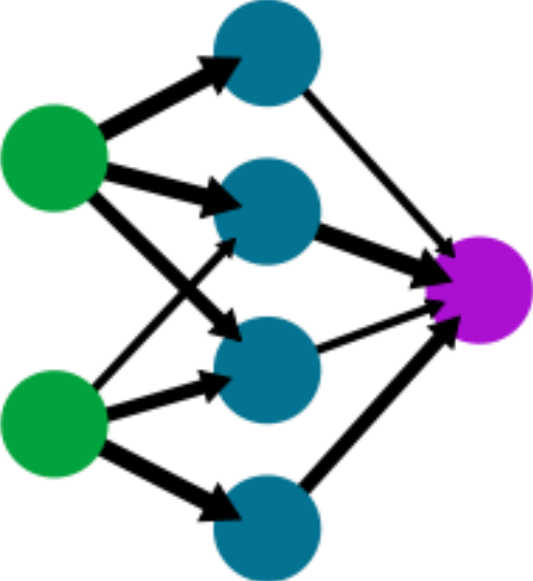## ~~Tensor~~ Computing Matrix

Weijie Zhao

10/13/2022

# Neural Networks



A simple neural network

input layer | hidden layer | output layer

# Deep Learning Framework Implementation

- Knowing the things under the hood

- Deployment

- Deployment on emerging hardware

# Deep Learning Language

- How to represent a deep neural network?

# Deep Learning Language

- How to represent a deep neural network?
  - Abstraction

# Deep Learning Language

- How to represent a deep neural network?
  - Abstraction
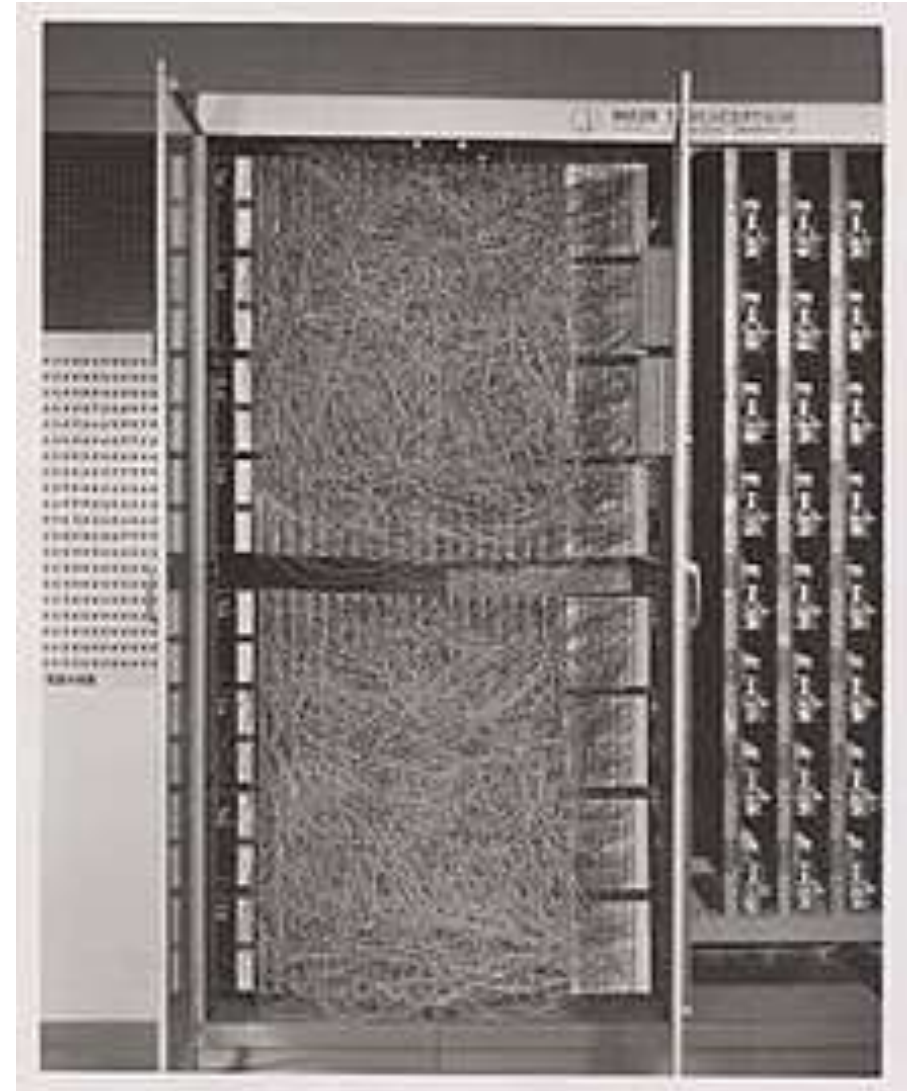- How to implement/deploy the abstraction deep neural network?

# Deep Learning Language

- How to represent a deep neural network?
  - Abstraction

- How to implement/deploy the abstraction deep neural network?
  - Build tensor operations workflow
  - Implement high-performance low-level operations

# Perceptron

- The perceptron was invented in 1943 by McCulloch and Pitts.

- The first implementation was a machine built in 1958 at the Cornell Aeronautical Laboratory by Frank Rosenblatt

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w} \cdot \mathbf{x} + b > 0, \\ 0 & \text{otherwise} \end{cases}$$

# Perceptron

- Linear Layer
    - Matrix multiplication
    - Addition
- ReLU

# 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
- PaddlePaddle
- 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()*       #  Remove assert, print, check_numerics
    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
```

https://web.stanford.edu/class/cs245/slides/TFGraphOptimizationsStanford.pdf

# 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, Slice(x) = x, etc.
- Rewrites that enable further constant folding
- Arithmetic rewrites that rely on known shapes or inputs, e.g.
  - Constant push-down:
    - Add(c1, Add(x, c2)) => Add(x, c1 + c2)
    - ConvND(c1 * x, c2) => ConvND(x, c1 * c2)
  - Partial constfold:
    - AddN(c1, x, c2, y) => AddN(c1 + c2, x, y),
    - Concat([x, c1, c2, y]) = Concat([x, Concat([c1, c2]), y)
  - Operations with neutral & absorbing elements:
    - x * Ones(s) => Identity(x), if shape(x) == output_shape
    - x * Ones(s) => BroadcastTo(x, Shape(s)), if shape(s) == output_shape
    - Same for x + Zeros(s) , x / Ones(s), x * Zeros(s) etc.
    - Zeros(s) - y => Neg(y), if shape(y) == output_shape
    - Ones(s) / y => Recip(y) if shape(y) == 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 => AddN(a, b, c, d)
  - Hoisting: AddN(x * a, b * x, x * c) => x * AddN(a+b+c)
  - Simplification to reduce number of nodes:
    - Numeric: x+x+x => 3*x
    - Logic: !(x > y) => x <= y
- Broadcast minimization
  - Example: (matrix1 + scalar1) + (matrix2 + scalar2) => (matrix1 + matrix2) + (scalar1 + scalar2)
- Better use of intrinsics
  - Matmul(Transpose(x), y) => Matmul(x, y, transpose_x=True)
- Remove redundant ops or op pairs
  - Transpose(Transpose(x, perm), inverse_perm)
  - BitCast(BitCast(x, dtype1), dtype2) => BitCast(x, dtype2)
  - Pairs of elementwise involutions f(f(x)) => x (Neg, Conj, Reciprocal, LogicalNot)
  - Repeated Idempotent ops f(f(x)) => f(x) (DeepCopy, Identity, CheckNumerics...)
- Hoist chains of unary ops at Concat/Split/SplitV
  - Concat([Exp(Cos(x)), Exp(Cos(y)), Exp(Cos(z))]) => Exp(Cos(Concat([x, y, z])))
  - [Exp(Cos(y)) for y in Split(x)] => Split(Exp(Cos(x), num_splits)