

Neural Network Preliminary Tensor Computing

Weijie Zhao

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

Neural Network Preliminary Tensor Computing

Weijie Zhao

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

Neural Network Preliminary

~~Tensor~~ Computing

Matrix

Weijie Zhao

- Matrix multiplication
- Non-linear activation
- Gradient descent

•Scalar

•Vector

•Matrix

•Tensor

•Rank

•Dimension

Neural Network Preliminary

~~Tensor~~ Computing

Matrix

Weijie Zhao

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

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

Neural Network Preliminary

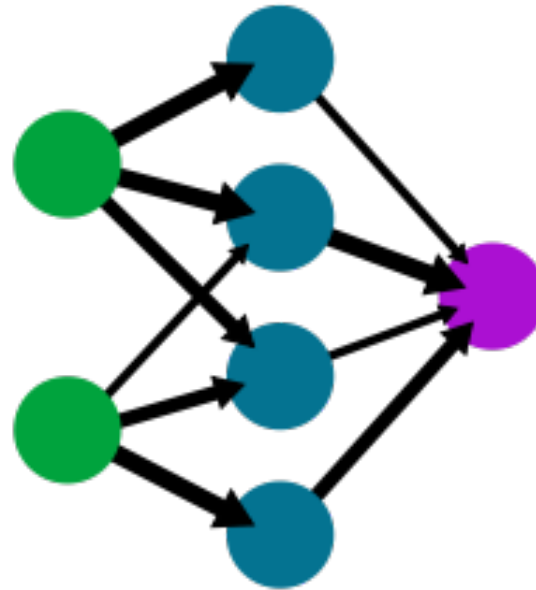
~~Tensor~~ Computing Matrix

Weijie Zhao

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?

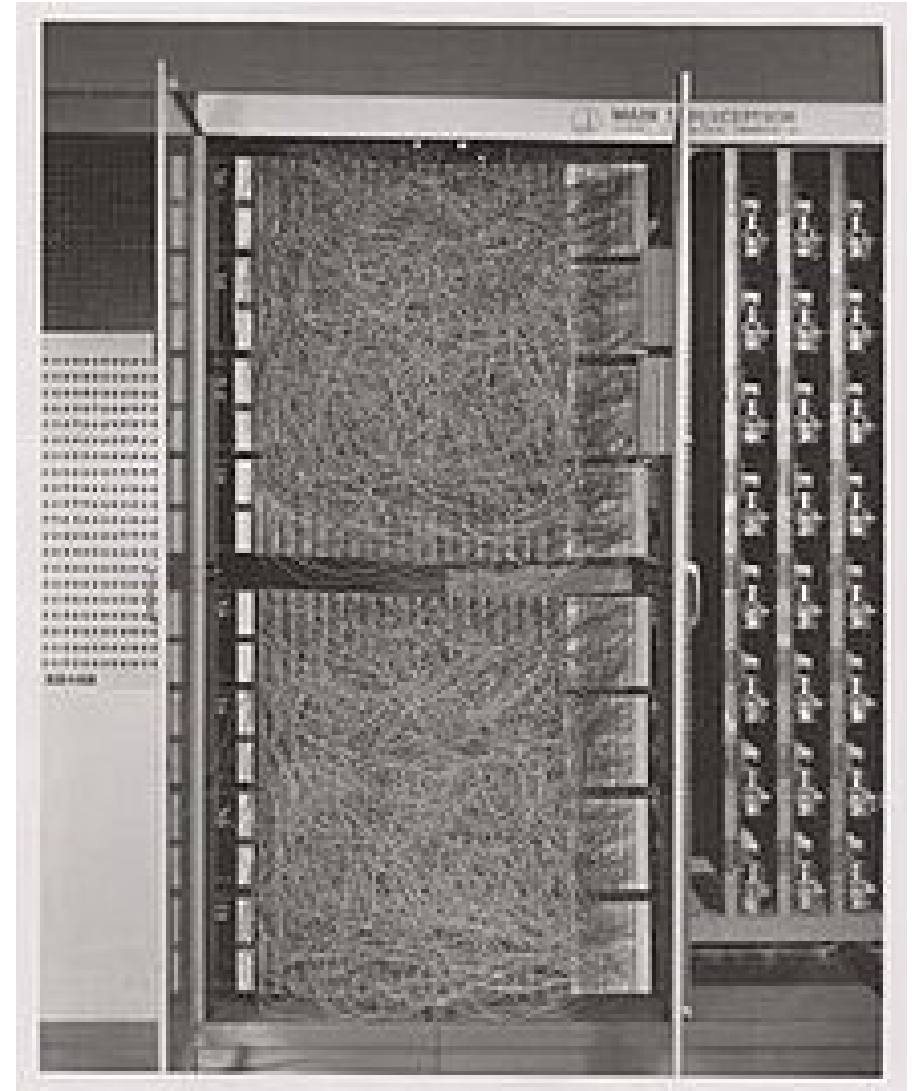
Deep Learning Language

- 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

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

Optimizations

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

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
```

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
 - $\text{Matmul}(\text{Transpose}(x), y) \Rightarrow \text{Matmul}(x, y, \text{transpose_x}=\text{True})$
- 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
 - $\text{Concat}([\text{Exp}(\text{Cos}(x)), \text{Exp}(\text{Cos}(y)), \text{Exp}(\text{Cos}(z))]) \Rightarrow \text{Exp}(\text{Cos}(\text{Concat}([x, y, z])))$
 - $[\text{Exp}(\text{Cos}(y)) \text{ for } y \text{ in Split}(x)] \Rightarrow \text{Split}(\text{Exp}(\text{Cos}(x)), \text{num_splits})$

Tensor

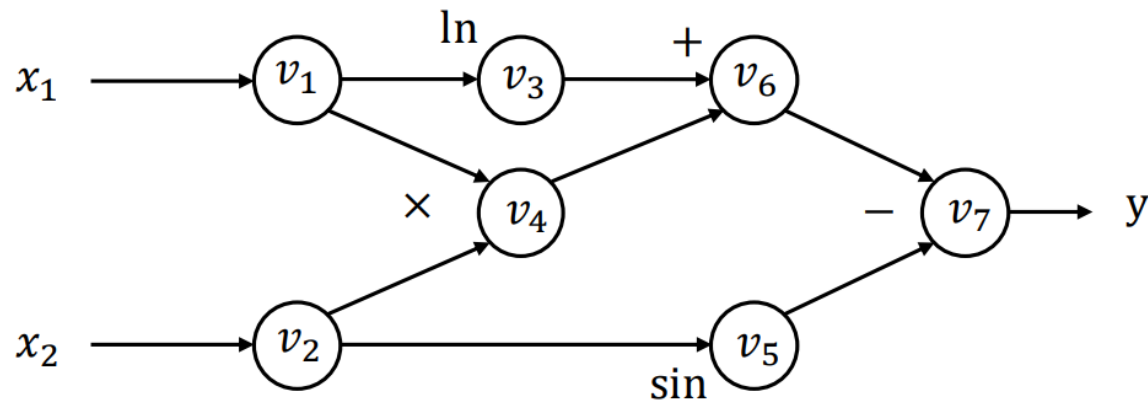
- Dense
 - Column major
 - Row major
 - Stride
- Sparse
 - Compressed representation
 - Set intersection

Differentiation

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

Computational Graph

$$y = f(x_1, x_2) = \ln(x_1) + x_1x_2 - \sin x_2$$

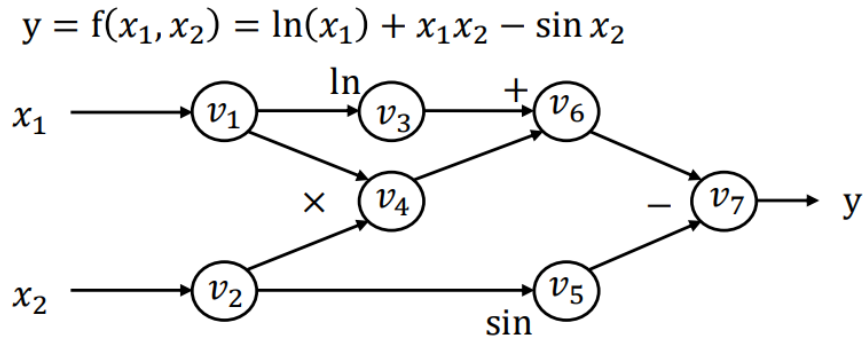


Forward evaluation trace

$$\begin{aligned}v_1 &= x_1 = 2 \\v_2 &= x_2 = 5 \\v_3 &= \ln v_1 = \ln 2 = 0.693 \\v_4 &= v_1 \times v_2 = 10 \\v_5 &= \sin v_2 = \sin 5 = -0.959 \\v_6 &= v_3 + v_4 = 10.693 \\v_7 &= v_6 - v_5 = 10.693 + 0.959 = 11.652 \\y &= v_7 = 11.652\end{aligned}$$

Each node represent an (intermediate) value in the computation. Edges present input output relations.

Forward Mode Auto Differentiation



Forward evaluation trace

$$\begin{aligned}v_1 &= x_1 = 2 \\v_2 &= x_2 = 5 \\v_3 &= \ln v_1 = \ln 2 = 0.693 \\v_4 &= v_1 \times v_2 = 10 \\v_5 &= \sin v_2 = \sin 5 = -0.959 \\v_6 &= v_3 + v_4 = 10.693 \\v_7 &= v_6 - v_5 = 10.693 + 0.959 = 11.652 \\y &= v_7 = 11.652\end{aligned}$$

Define $\dot{v}_i = \frac{\partial v_i}{\partial x_1}$

We can then compute the \dot{v}_i iteratively in the forward topological order of the computational graph

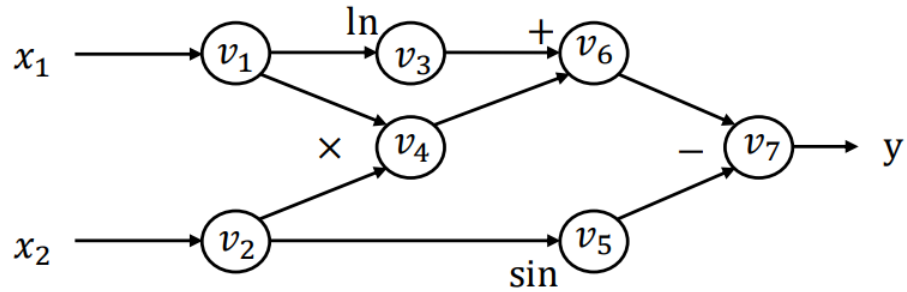
Forward AD trace

$$\begin{aligned}\dot{v}_1 &= 1 \\ \dot{v}_2 &= 0 \\ \dot{v}_3 &= \dot{v}_1/v_1 = 0.5 \\ \dot{v}_4 &= \dot{v}_1 v_2 + \dot{v}_2 v_1 = 1 \times 5 + 0 \times 2 = 5 \\ \dot{v}_5 &= \dot{v}_2 \cos v_2 = 0 \times \cos 5 = 0 \\ \dot{v}_6 &= \dot{v}_3 + \dot{v}_4 = 0.5 + 5 = 5.5 \\ \dot{v}_7 &= \dot{v}_6 - \dot{v}_5 = 5.5 - 0 = 5.5\end{aligned}$$

Now we have $\frac{\partial y}{\partial x_1} = \dot{v}_7 = 5.5$

Reverse Mode Auto Differentiation

$$y = f(x_1, x_2) = \ln(x_1) + x_1 x_2 - \sin x_2$$



Forward evaluation trace

$$\begin{aligned} v_1 &= x_1 = 2 \\ v_2 &= x_2 = 5 \\ v_3 &= \ln v_1 = \ln 2 = 0.693 \\ v_4 &= v_1 \times v_2 = 10 \\ v_5 &= \sin v_2 = \sin 5 = -0.959 \\ v_6 &= v_3 + v_4 = 10.693 \\ v_7 &= v_6 - v_5 = 10.693 + 0.959 = 11.652 \\ y &= v_7 = 11.652 \end{aligned}$$

Define adjoint $\bar{v}_i = \frac{\partial y}{\partial v_i}$

We can then compute the \bar{v}_i iteratively in the **reverse** topological order of the computational graph

Reverse AD evaluation trace

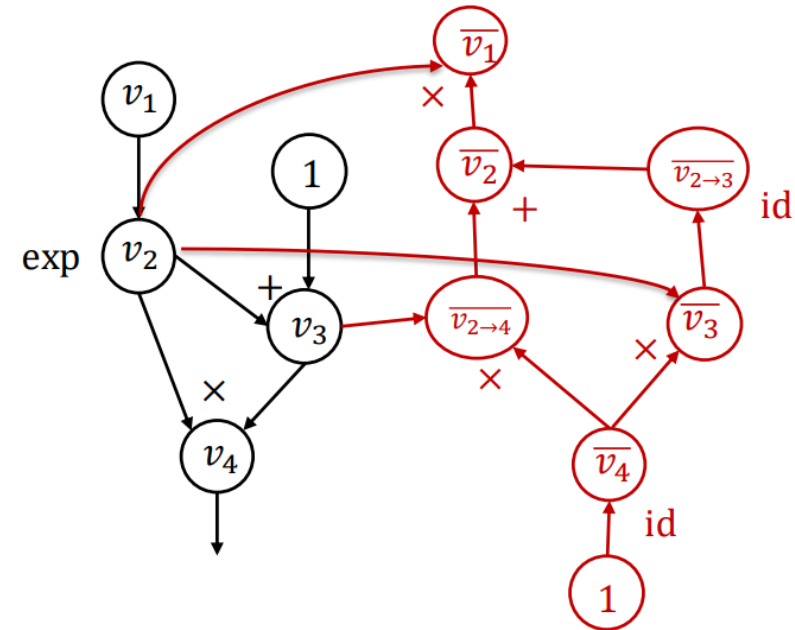
$$\begin{aligned} \bar{v}_7 &= \frac{\partial y}{\partial v_7} = 1 \\ \bar{v}_6 &= \bar{v}_7 \frac{\partial v_7}{\partial v_6} = \bar{v}_7 \times 1 = 1 \\ \bar{v}_5 &= \bar{v}_7 \frac{\partial v_7}{\partial v_5} = \bar{v}_7 \times (-1) = -1 \\ \bar{v}_4 &= \bar{v}_6 \frac{\partial v_6}{\partial v_4} = \bar{v}_6 \times 1 = 1 \\ \bar{v}_3 &= \bar{v}_6 \frac{\partial v_6}{\partial v_3} = \bar{v}_6 \times 1 = 1 \\ \bar{v}_2 &= \bar{v}_5 \frac{\partial v_5}{\partial v_2} + \bar{v}_4 \frac{\partial v_4}{\partial v_2} = \bar{v}_5 \times \cos v_2 + \bar{v}_4 \times v_1 = -0.284 + 2 = 1.716 \\ \bar{v}_1 &= \bar{v}_4 \frac{\partial v_4}{\partial v_1} + \bar{v}_3 \frac{\partial v_3}{\partial v_1} = \bar{v}_4 \times v_2 + \bar{v}_3 \frac{1}{v_1} = 5 + \frac{1}{2} = 5.5 \end{aligned}$$

Reverse Mode Auto Differentiation

```
def gradient(out):  
    node_to_grad = {out: [1]}  
    for i in reverse_topo_order(out):  
         $\bar{v}_i = \sum_j \bar{v}_{i \rightarrow j} = \text{sum}(\text{node\_to\_grad}[i])$   
        for k in inputs(i):  
            compute  $\bar{v}_{k \rightarrow i} = \bar{v}_i \frac{\partial v_i}{\partial v_k}$   
            append  $\bar{v}_{k \rightarrow i}$  to node_to_grad[k]  
    return adjoint of input  $\bar{v}_{input}$ 
```



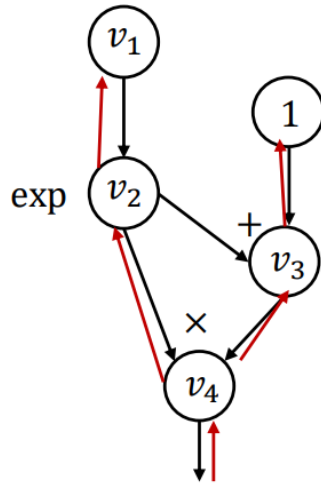
```
i = 2  
node_to_grad: {  
    1: [ $\bar{v}_1$ ]  
    2: [ $\bar{v}_{2 \rightarrow 4}$ ,  $\bar{v}_{2 \rightarrow 3}$ ]  
    3: [ $\bar{v}_3$ ]  
    4: [ $\bar{v}_4$ ]  
}
```



NOTE: id is identity function

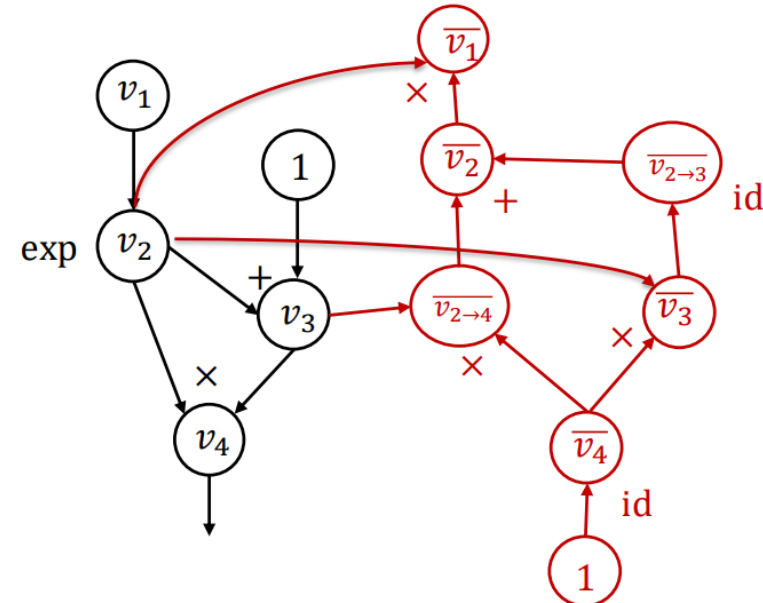
Backprop V.S. Reverse Mode AD

Backprop



- Run backward operations the same forward graph
- Used in first generation deep learning frameworks (caffe, cuda-convnet)

Reverse mode AD by extending computational graph



- Construct separate graph nodes for adjoints
- Used by modern deep learning frameworks