# Structure of Deep Learning Frameworks: computational graph, autodiff, and optimizers
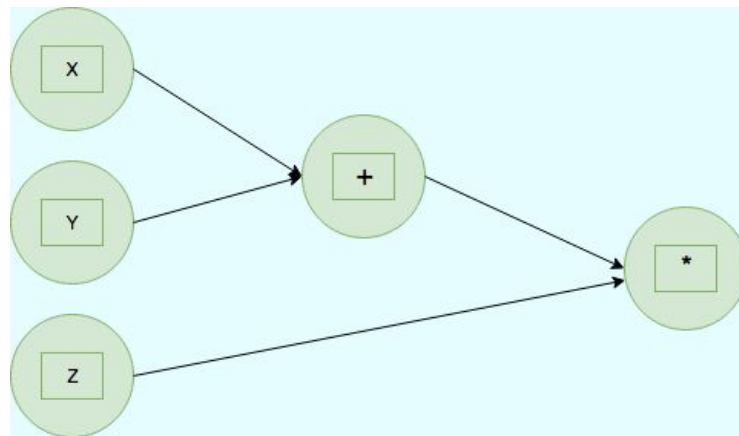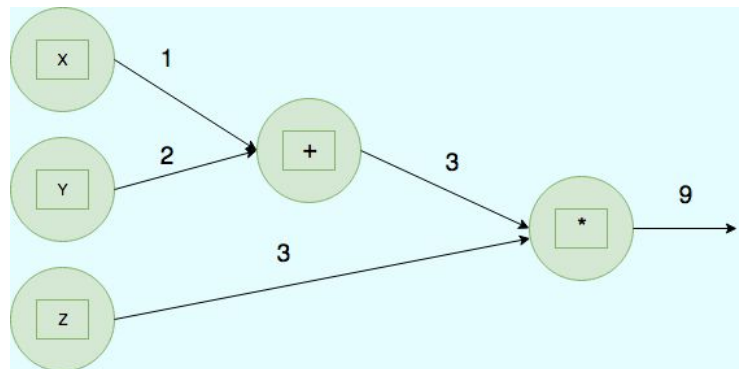
Joris Mollinga
SURF

# Computational Graphs

- Represents math in graph format
- Nodes and edges

# Computational Graphs

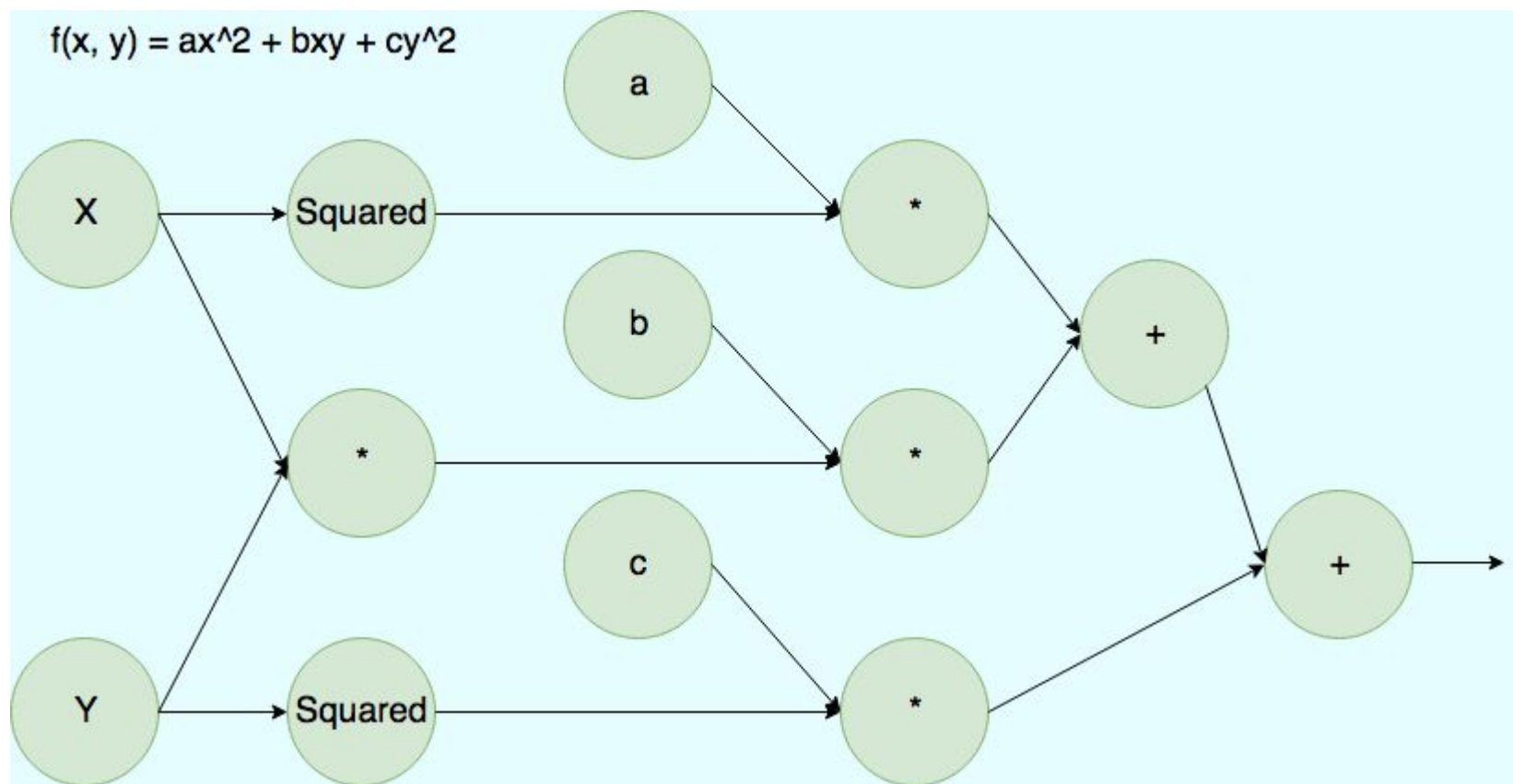- Represents math in graph format
- Nodes and edges

# Computational Graphs

- Represents math in graph format
- Nodes and edges

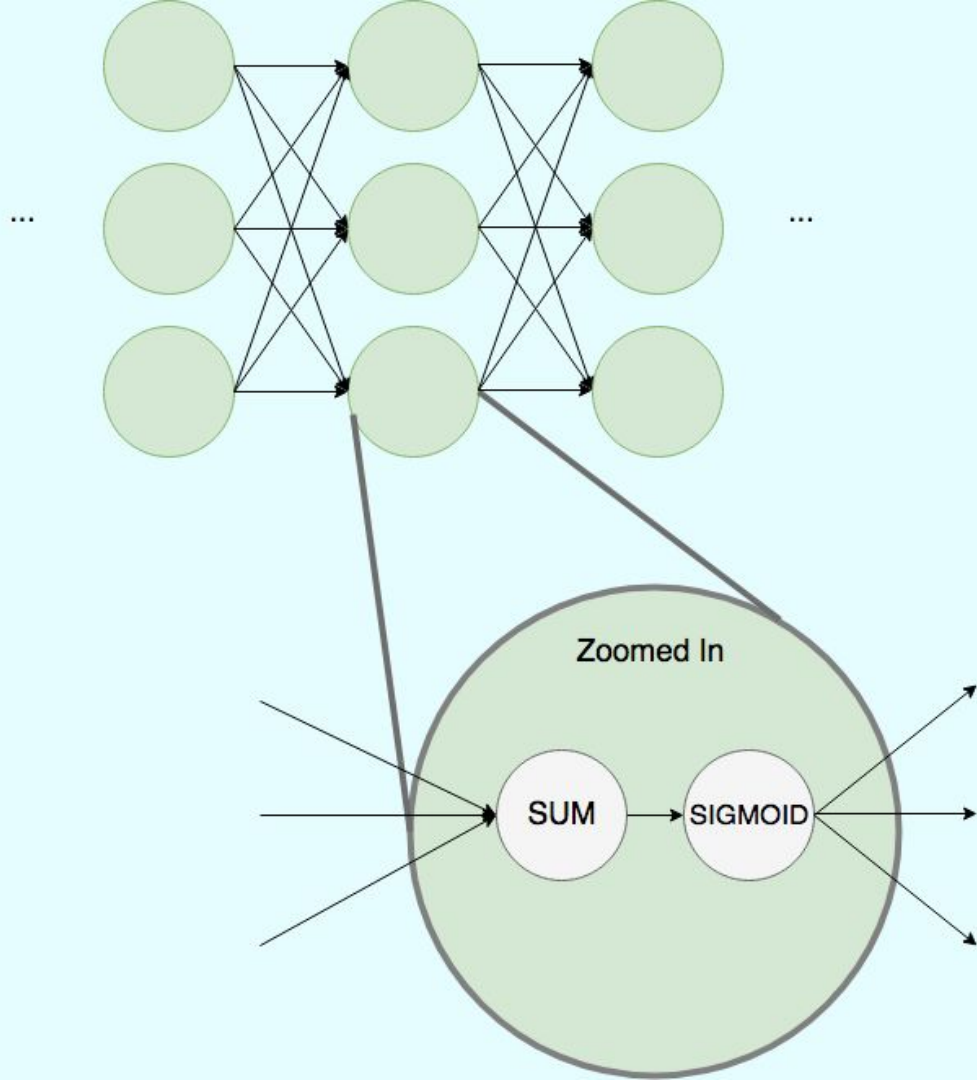f(x, y) = ax^2 + bxy + cy^2

Zoomed In

SUM → SIGMOID

W0

Output from previous layer

*

W1

Output from previous layer

*

W2

Output from previous layer

*

Zoomed In

Bias

SUM

SIGMOID

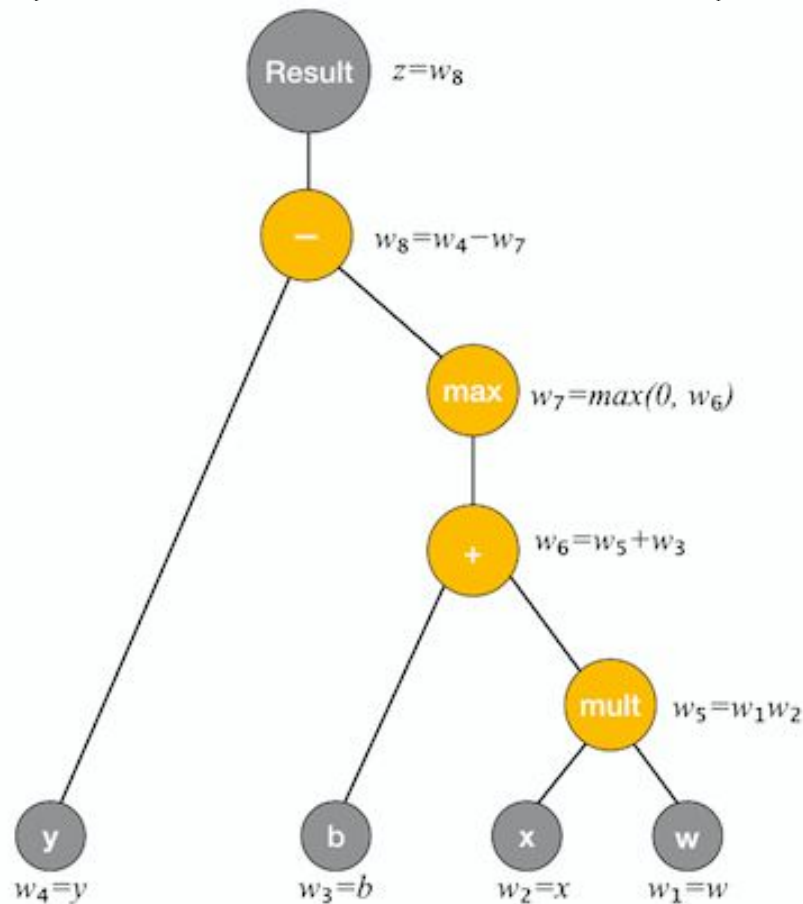**Autodiff**

$$C(y, w, x, b) = y - max(0, w \cdot x + b)$$

$$C(y, w, x, b) = y - max(0, w \cdot x + b)$$

| Node | Expression | Value |
|------|------------|-------|
| $w_1$ | $w$ | 2 |
| $w_2$ | $x$ | 1 |
| $w_3$ | $b$ | 1 |
| $w_4$ | $y$ | 5 |
| $w_5$ | $w_1 \cdot w_2$ | 2 |
| $w_6$ | $w_5 + w_3$ | 3 |
| $w_7$ | $\max(0, w_6)$ | 3 |
| $w_8$ | $w_4 - w_7$ | 2 |
| $z$ | $w_8$ | 2 |

$$C(y, w, x, b) = y - max(0, w \cdot x + b)$$



Result — $z = w_8$

$-$ — $w_8 = w_4 - w_7$

max — $w_7 = max(0, w_6)$

$+$ — $w_6 = w_5 + w_3$

mult — $w_5 = w_1 w_2$

y — $w_4 = y$

b — $w_3 = b$

x — $w_2 = x$

w — $w_1 = w$

| | |
|---|---|
| $\dfrac{\delta w_5}{\delta w_1} = \dfrac{\delta w_1 w_2}{\delta w_1} = w_2$ | $\dfrac{\delta w_5}{\delta w_2} = \dfrac{\delta w_1 w_2}{\delta w_2} = w_1$ |
| $\dfrac{\delta w_6}{\delta w_5} = \dfrac{\delta w_5 + w_3}{\delta w_5} = 1$ | $\dfrac{\delta w_6}{\delta w_3} = \dfrac{\delta w_5 + w_3}{\delta w_3} = 1$ |
| $\dfrac{\delta w_7}{\delta w_6} = \dfrac{\delta max(0, w_6)}{\delta w_6} = \begin{cases} 0, & x < 0 \\ 1, & x > 0 \end{cases}$ | $\dfrac{\delta w_8}{\delta w_7} = \dfrac{\delta w_4 - w_7}{\delta w_7} = -1$ |
| $\dfrac{\delta w_8}{\delta w_4} = \dfrac{\delta w_4 - w_7}{\delta w_4} = 1$ | $\dfrac{\delta z}{\delta w_8} = 1$ |

$$\frac{\delta z}{\delta w_1} = \frac{\delta z}{\delta w_8} \times \frac{\delta w_8}{\delta w_7} \times \frac{\delta w_7}{\delta w_6} \times \frac{\delta w_6}{\delta w_5} \times \frac{\delta w_5}{\delta w_1} = 1 \times (-1) \times \begin{cases} 0, & x < 0 \\ 1, & x > 0 \end{cases} \times 1 \times w_2 = \begin{cases} 0, & x < 0 \\ 1, & x > 0 \end{cases}$$

# Frameworks

- Caffe
- NVCaffe
- IntelCaffe
- PyTorch
- PyTorch Lightning
- Tensorflow
- Horovod (DL distribution framework only)

# PyTorch

- Probably the most popular package nowadays
- Multi-GPU support (more than one node)
- Supports cuDNN and MKL-DNN
- Suppports various communication backends: Gloo, MPI, NCCL

# Tensorflow

- Also very popular
- Multi-GPU in one node through *tf.device*
- Parallelism across nodes with *tf.distribute*
- Support for TPU
- Supports cuDNN and MKL-DNN
- Supports NCCL allreduce
- For CPU, build from source

# Tensorflow optimization tips

- Use *TFRecords* to prevent I/) bottlenecks
- Overlap computation and data preparation using *tf.data.Dataset.prefetch*
- Parallelize data transformation using *tf.data.Dataset.map*
- If data fits in memory, use *tf.data.Dataset.cache*

# Horovod

Is a distribution framework for deep learning (not a deep learning framework itself). Design goals:

- Minimal code changes to make serial program distributed
- High performance distribution

# Horovod

- Support for Keras, MXNet, Tensorflow and PyTorch
- Supports MPI and MCCL as communication backends
- Requires about 6 lines of code changes
- Has it's own profiling ability, making it easy to assess communication overhead.

# PyTorch Lightning

- PyTorch wrapper for high performance AI research
- Supports various distribution strategies (horovod, data parallel, model parallel)
- Has other convenient features too.

# Thank you!