

Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

Building Brains - Parallelisation Strategies of Large-Scale Deep Learning Neural Networks on Parallel Scale Out Architectures Like Apache Spark

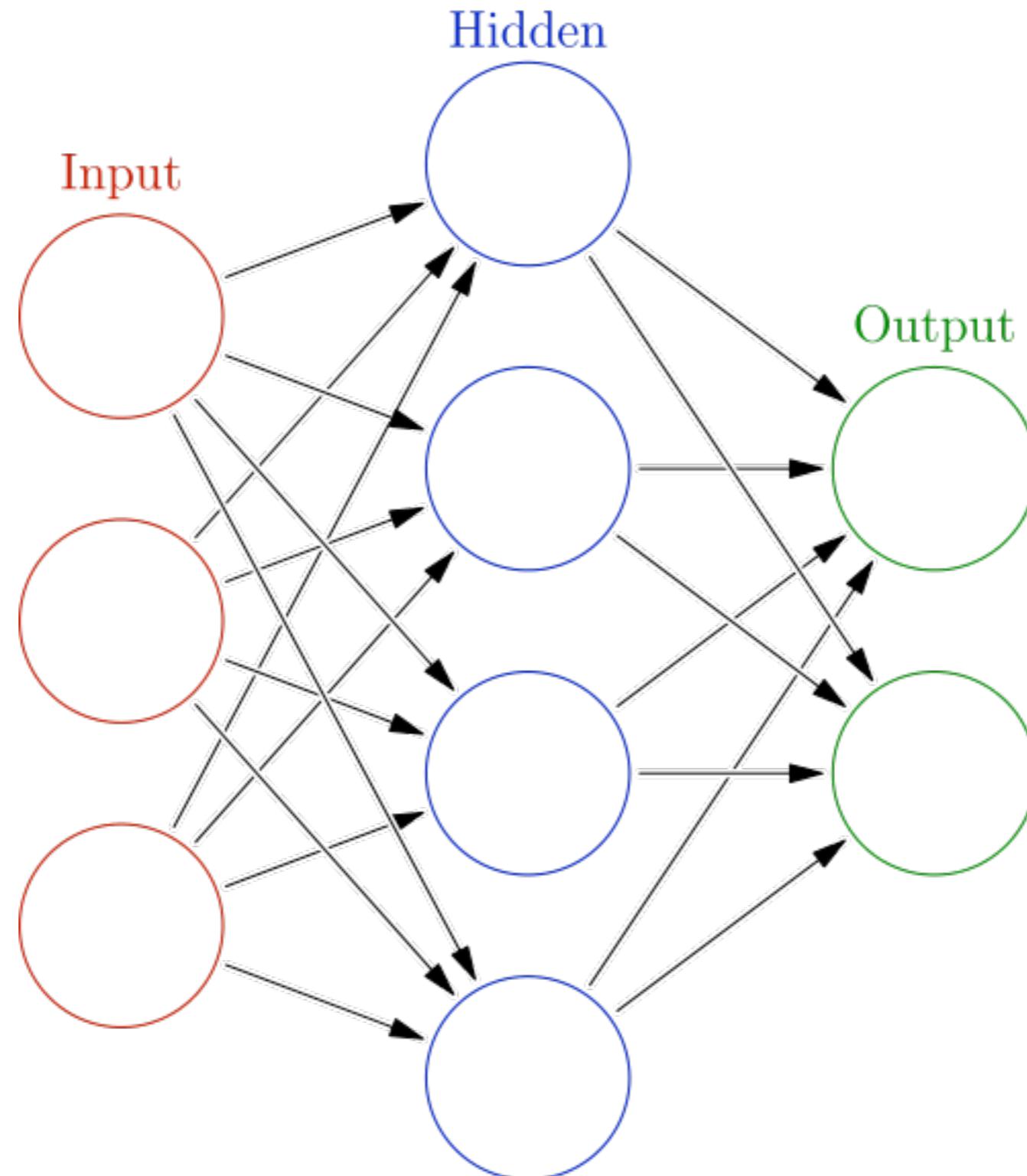


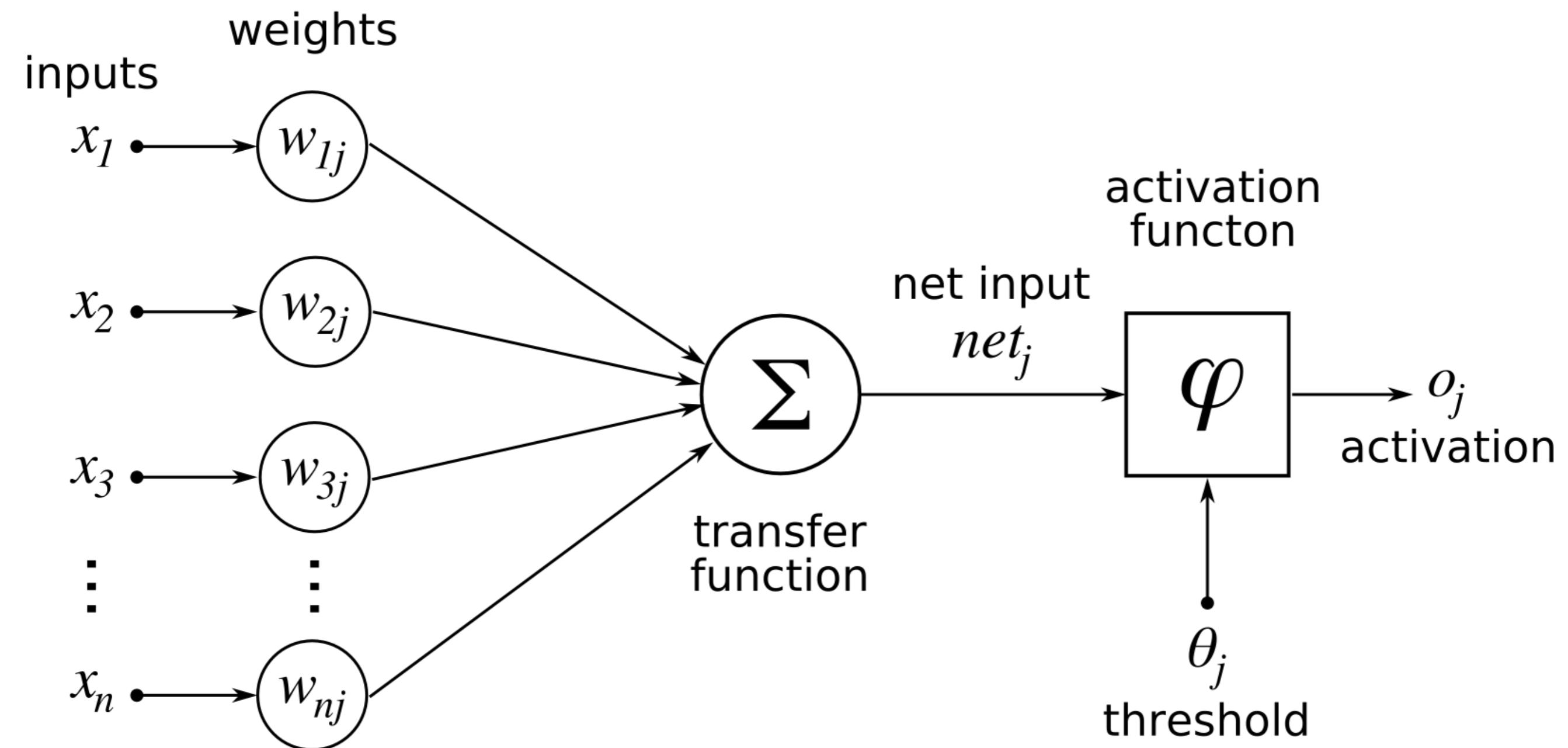
Romeo Kienzler, IBM Watson IoT

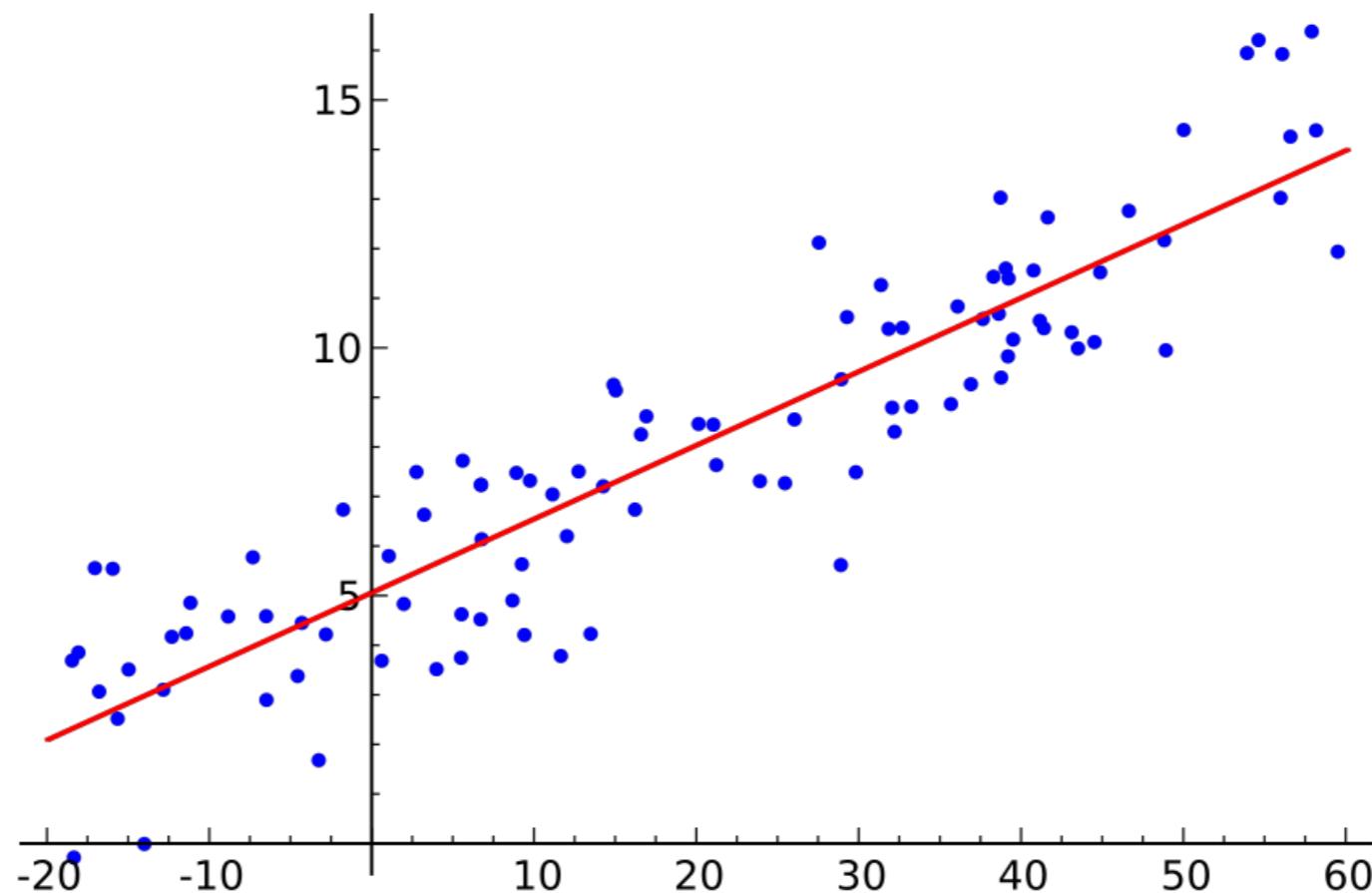
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk

Romeo Kienzler

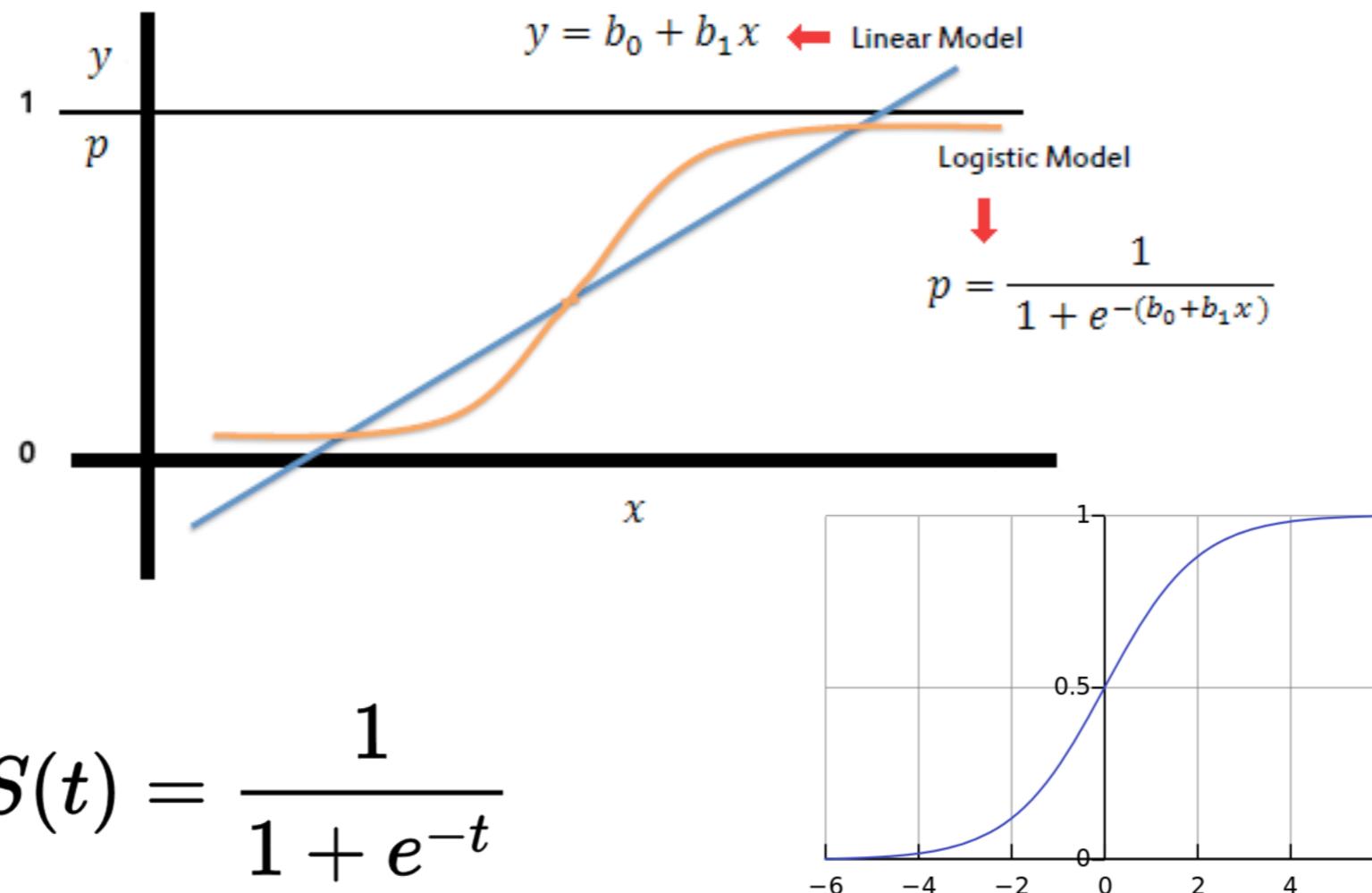
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me



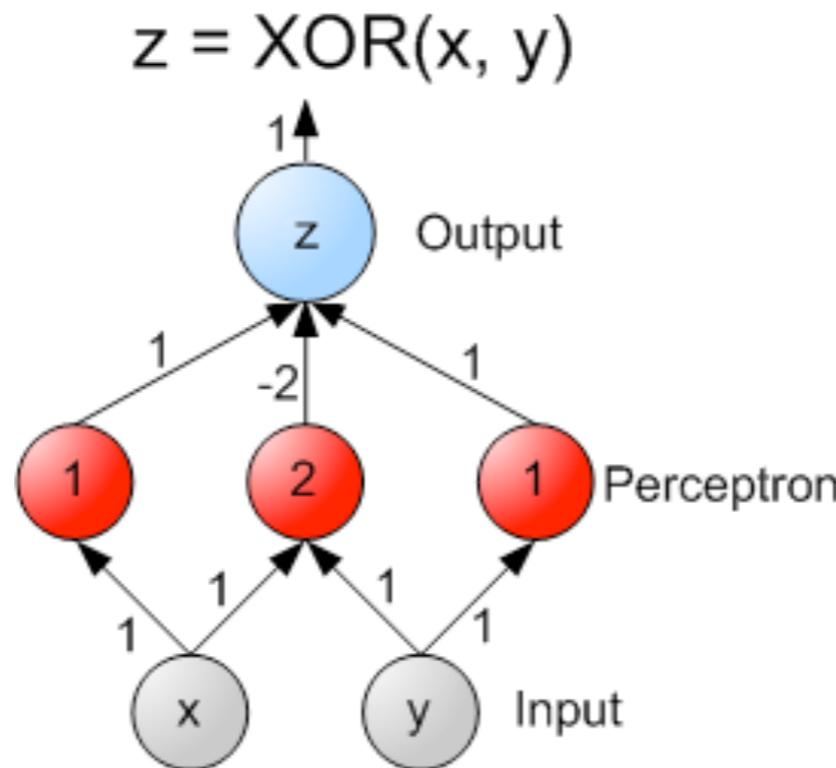




$$y_i = \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n$$

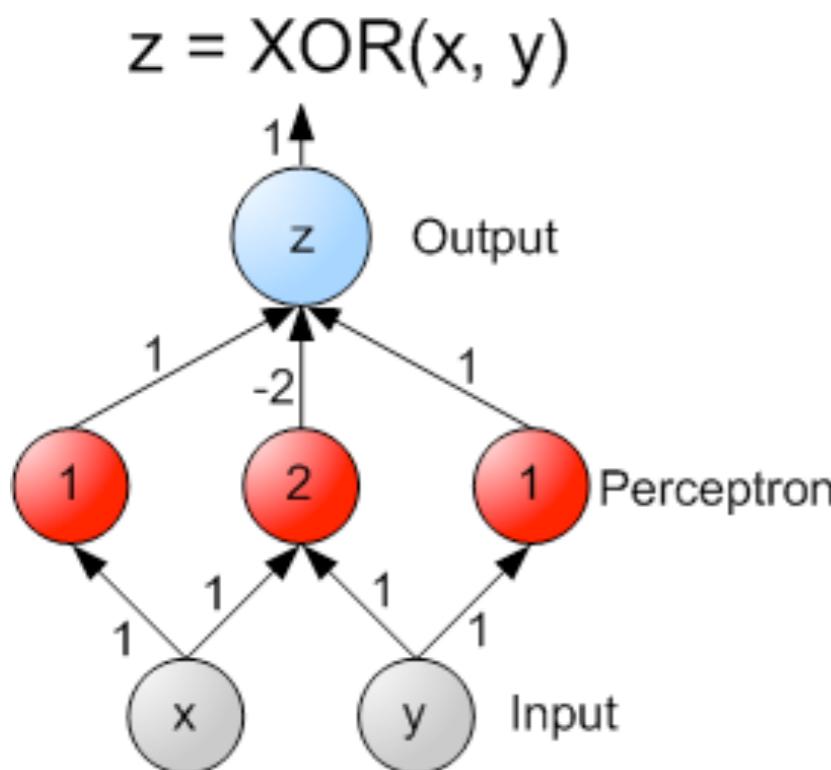


the forward pass



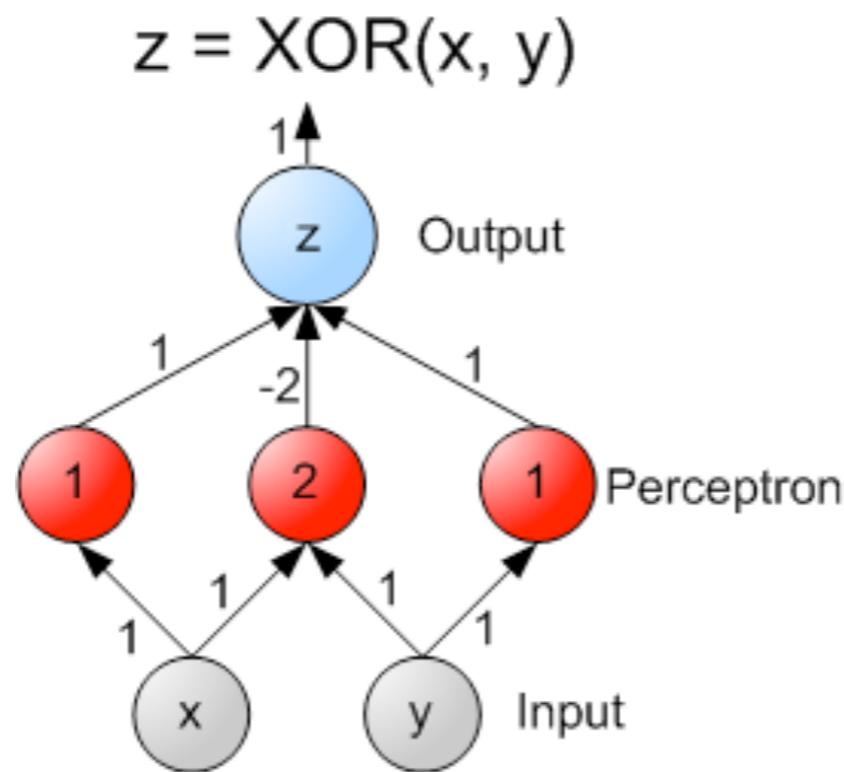
$$o = f(\sum_{k=1}^n i_k W_k)$$

back propagation



$$z^{(2)} = XW^{(1)} \quad (1)$$
$$a^{(2)} = f(z^{(2)}) \quad (2)$$
$$z^{(3)} = a^{(2)}W^{(2)} \quad (3)$$
$$\hat{y} = f(z^{(3)}) \quad (4)$$
$$J = \sum \frac{1}{2} (y - \hat{y})^2 \quad (5)$$

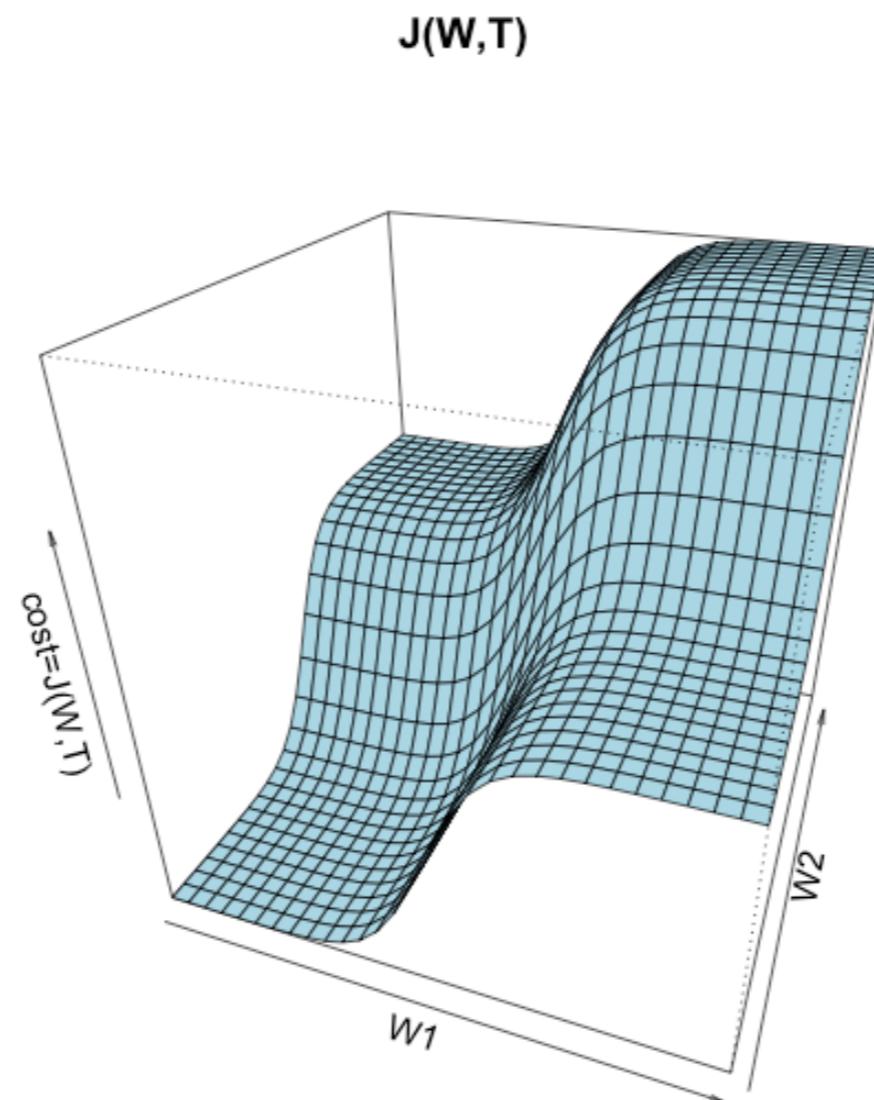
back propagation



$$J = \sum \frac{1}{2} (y - f(f(XW^{(1)}))W^{(2)})^2$$

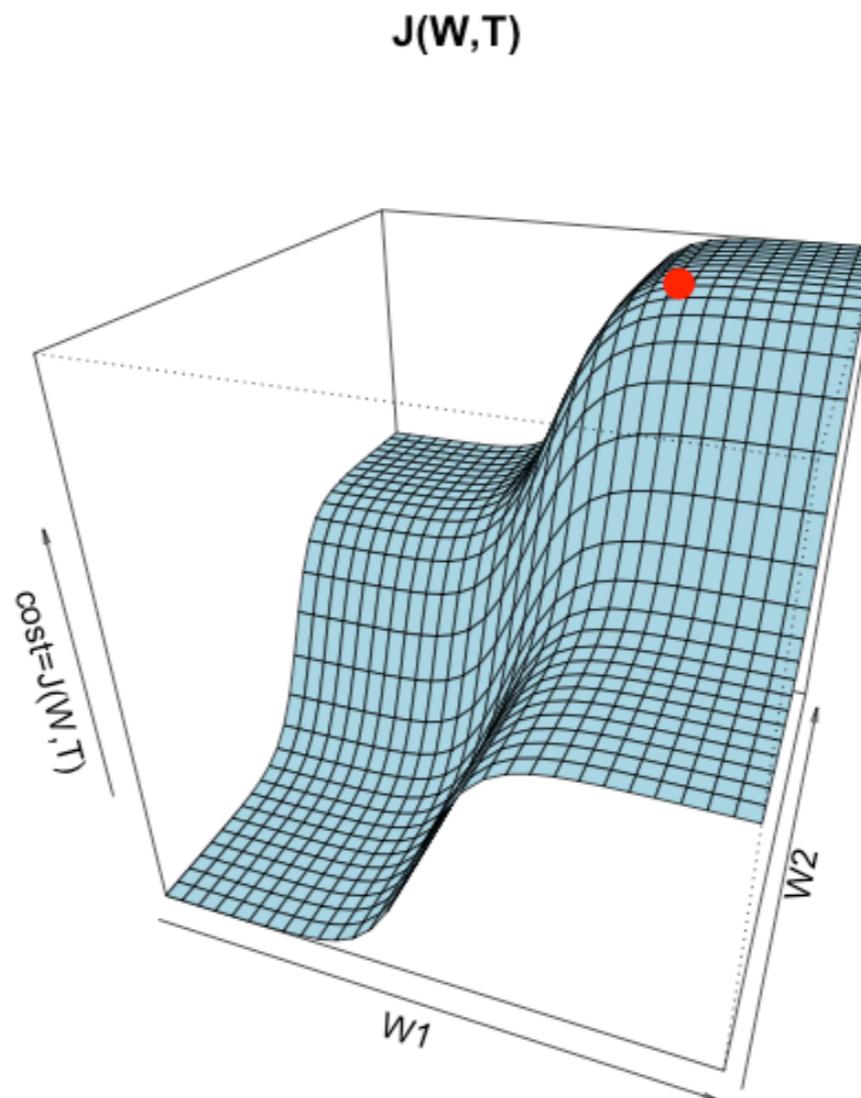
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



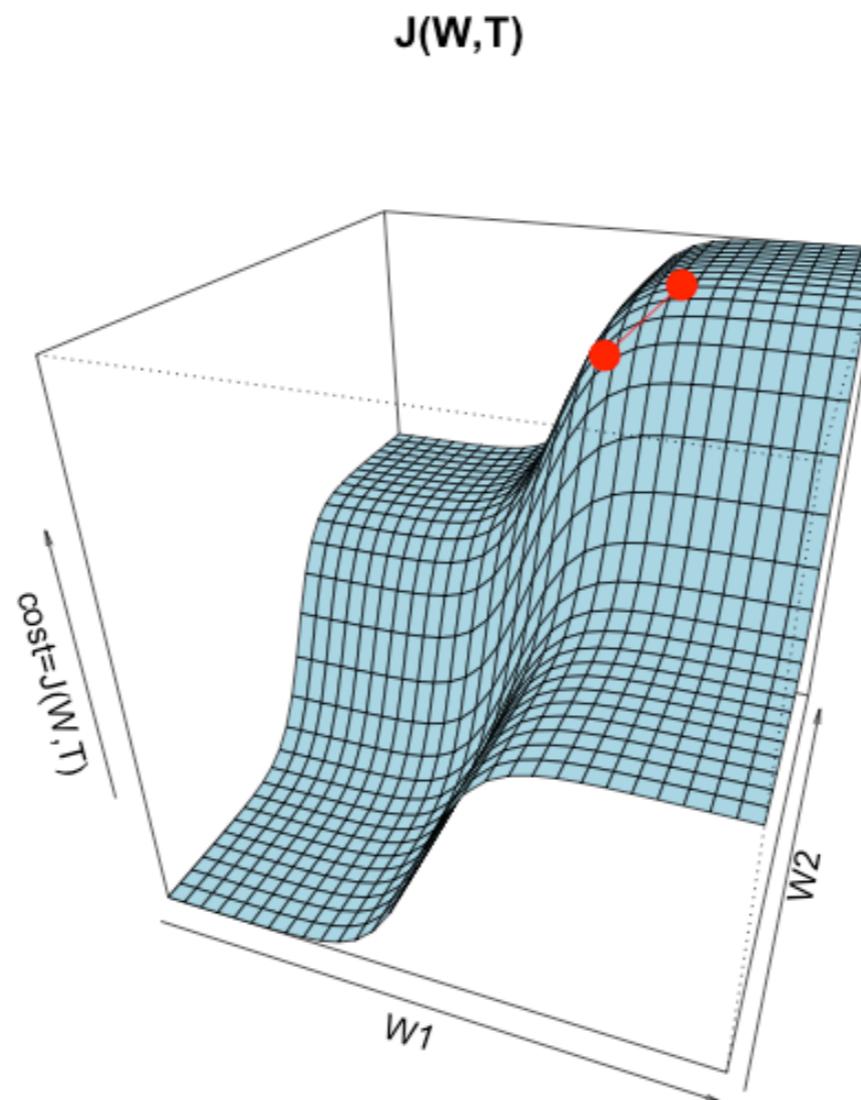
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



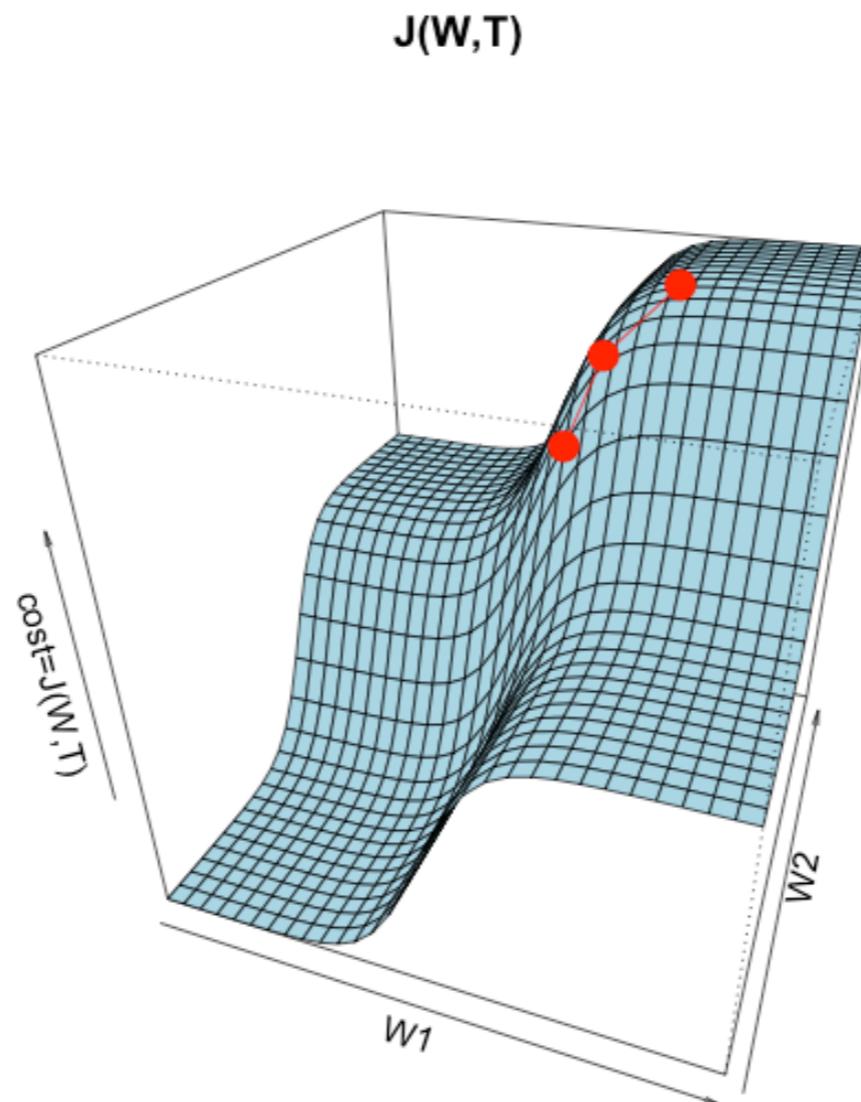
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



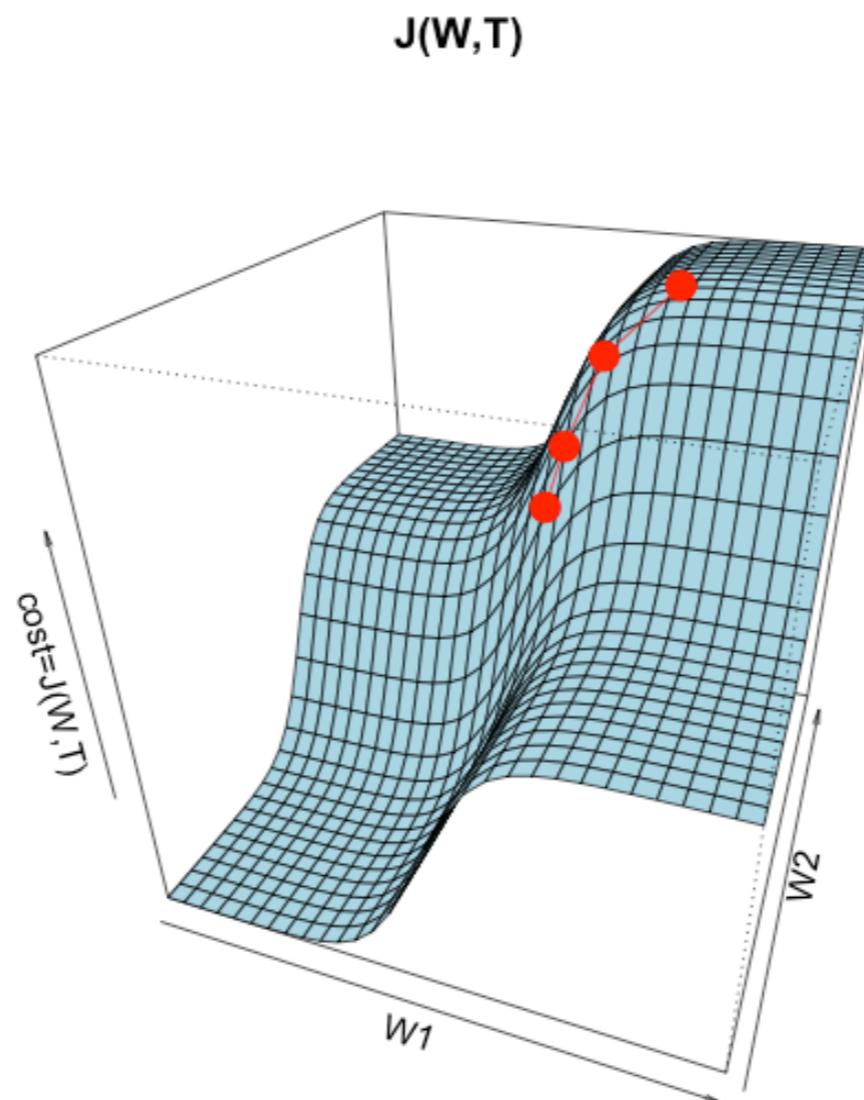
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



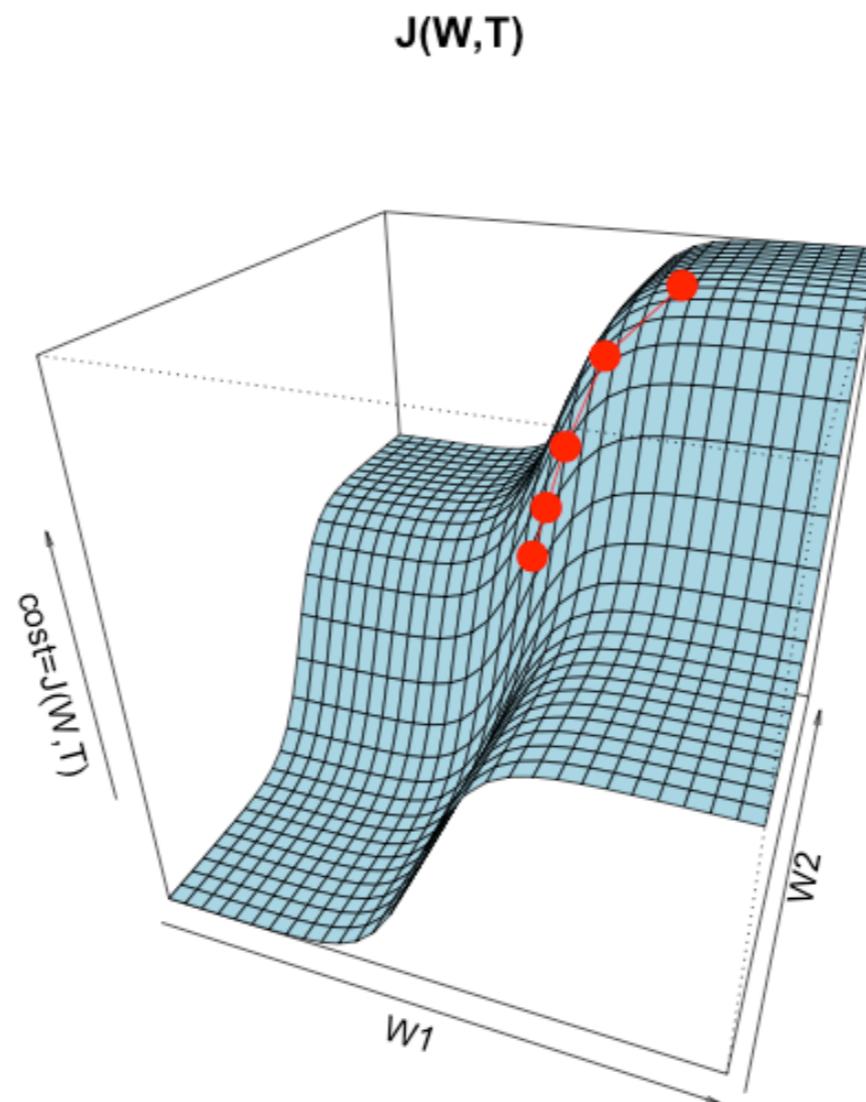
Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me

gradient descent



parallelisation

- ▶ inter-model parallelism
- ▶ data parallelism
- ▶ intra-model parallelism
- ▶ pipelined parallelism

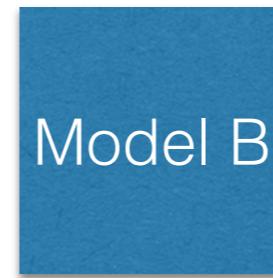
inter-model parallelism

aka. hyper parameter space exploration / tuning



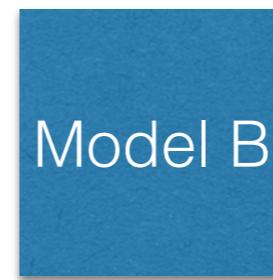
inter-model parallelism

aka. hyper parameter space exploration / tuning



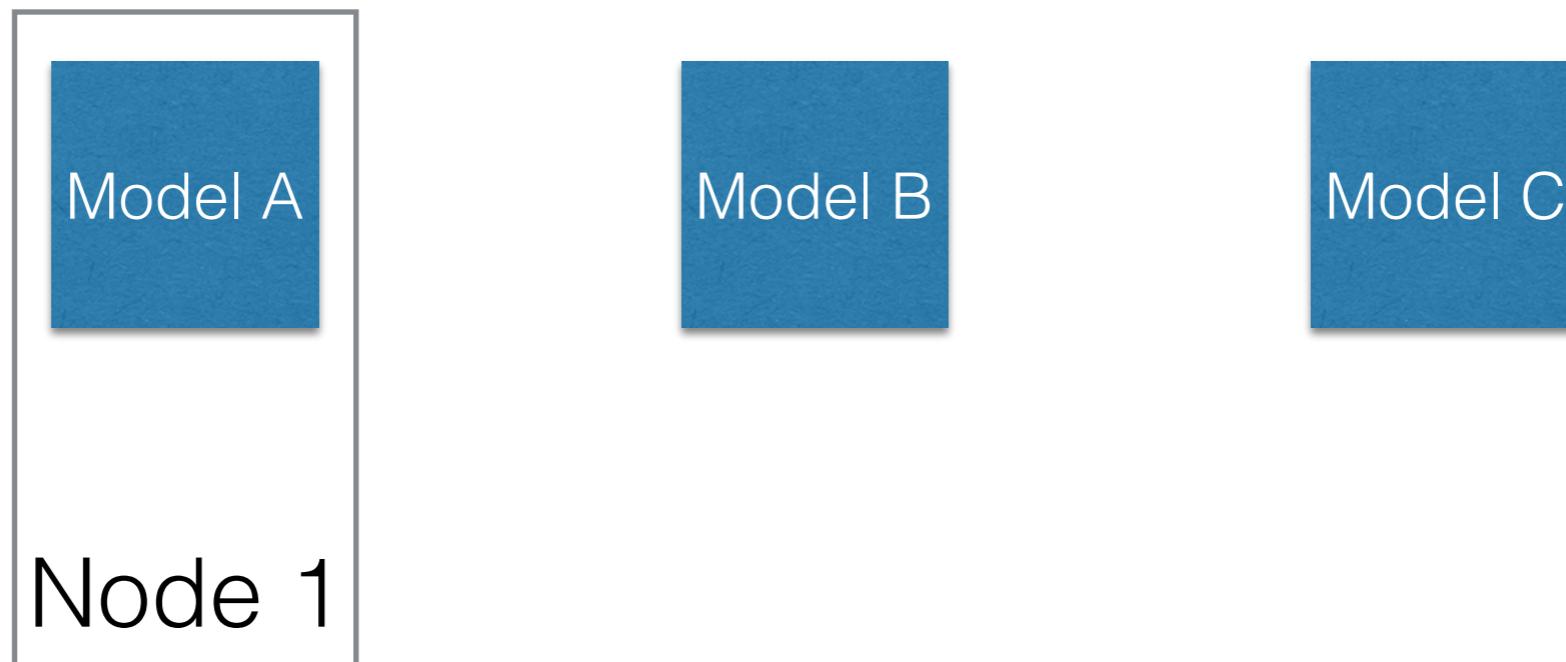
inter-model parallelism

aka. hyper parameter space exploration / tuning



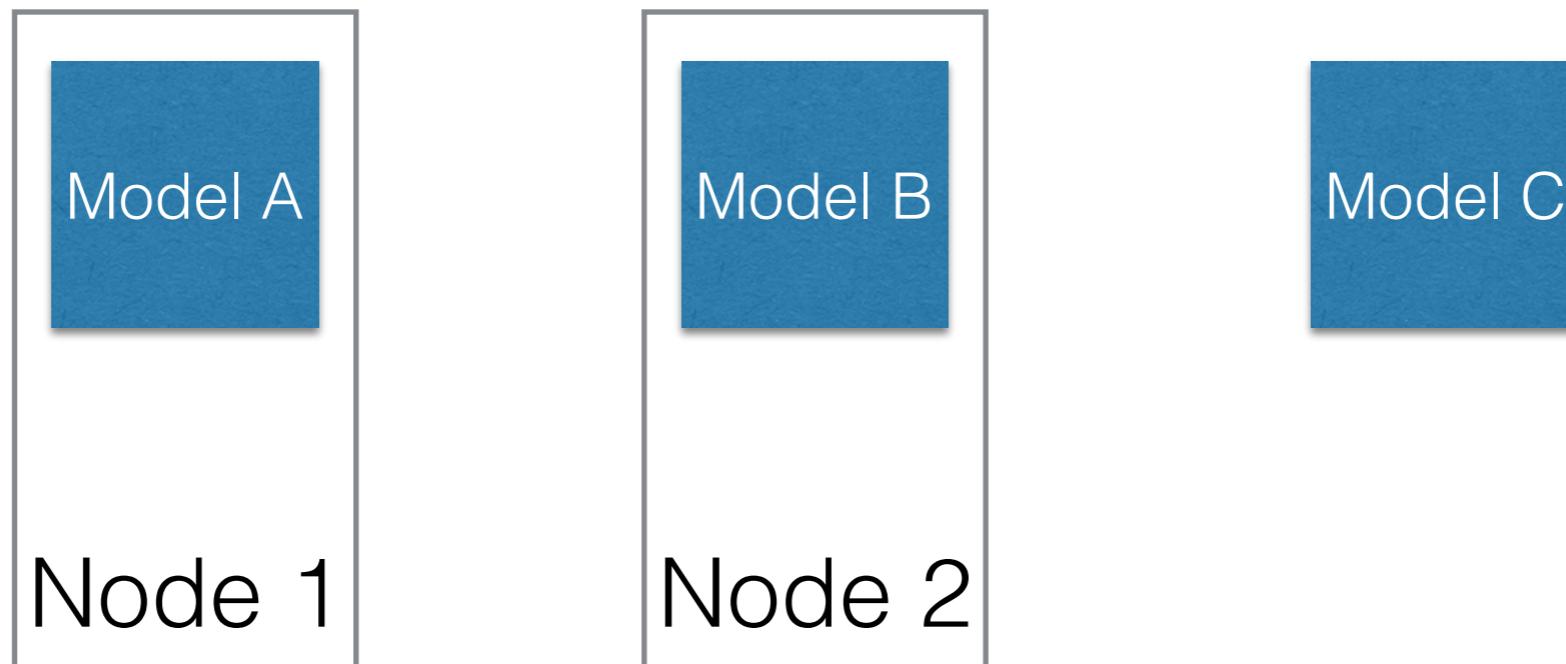
inter-model parallelism

aka. hyper parameter space exploration / tuning



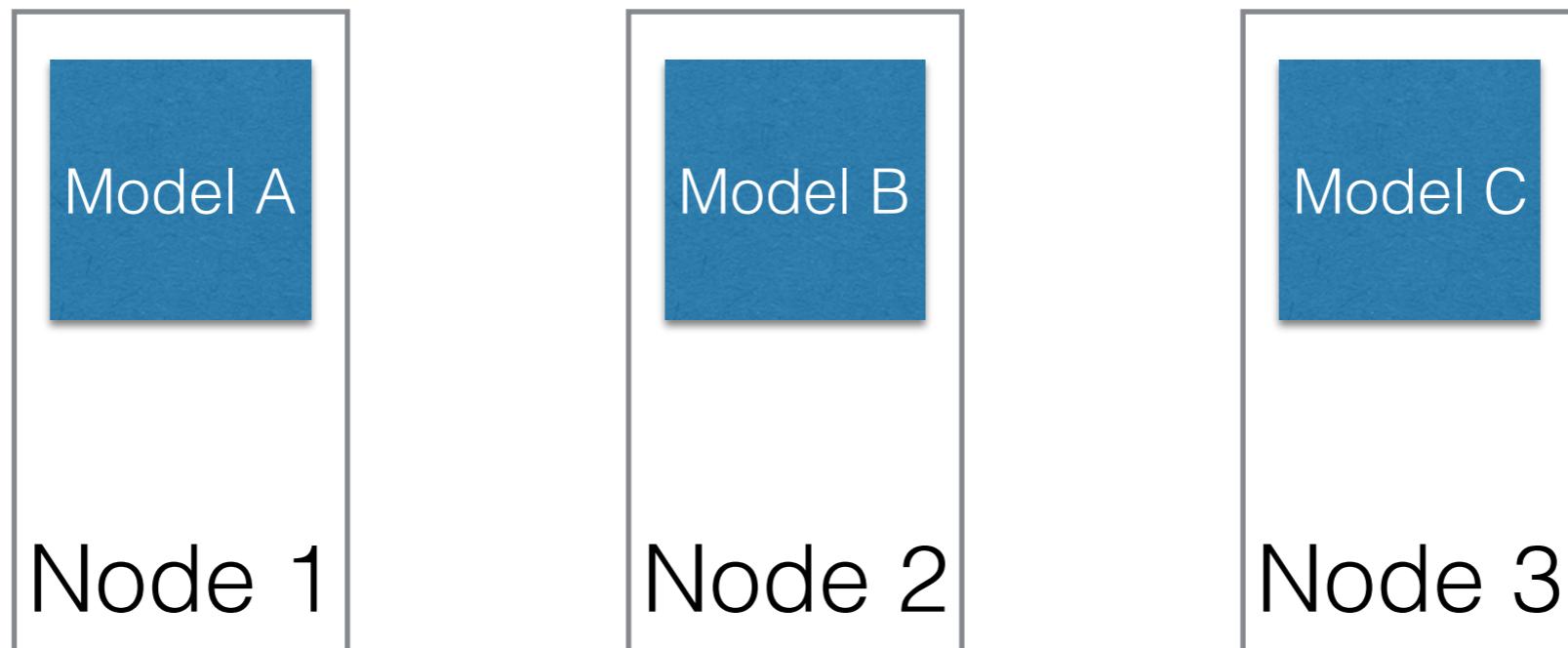
inter-model parallelism

aka. hyper parameter space exploration / tuning



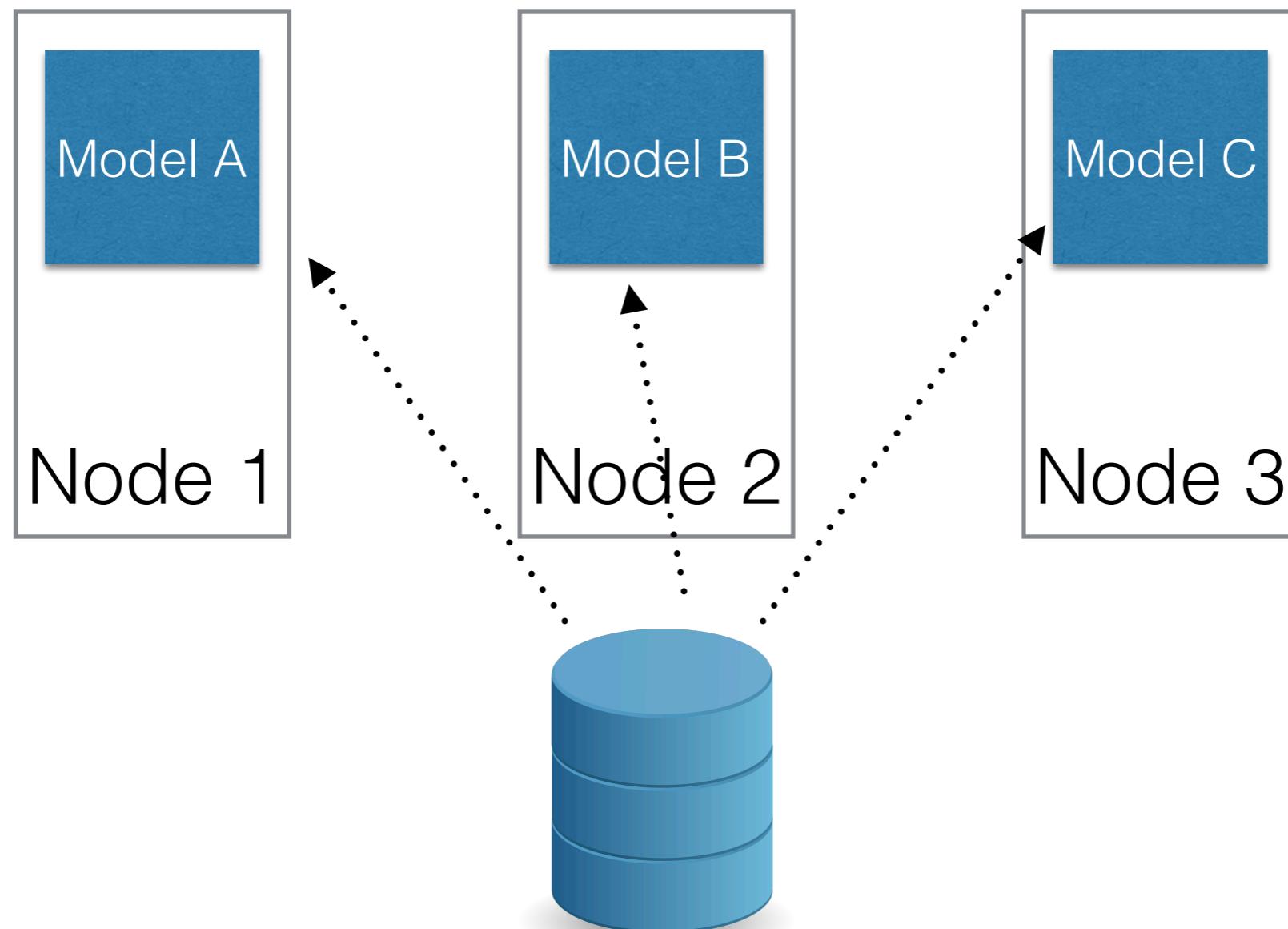
inter-model parallelism

aka. hyper parameter space exploration / tuning



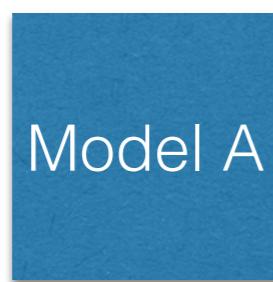
inter-model parallelism

aka. hyper parameter space exploration / tuning



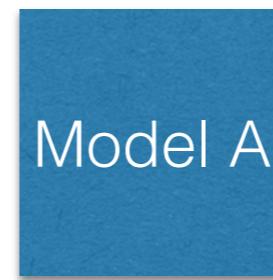
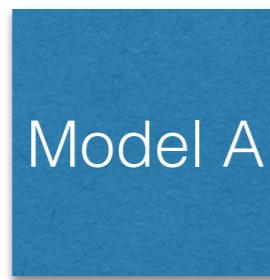
data parallelism

aka. “Jeff Dean style” parameter averaging



data parallelism

aka. “Jeff Dean style” parameter averaging



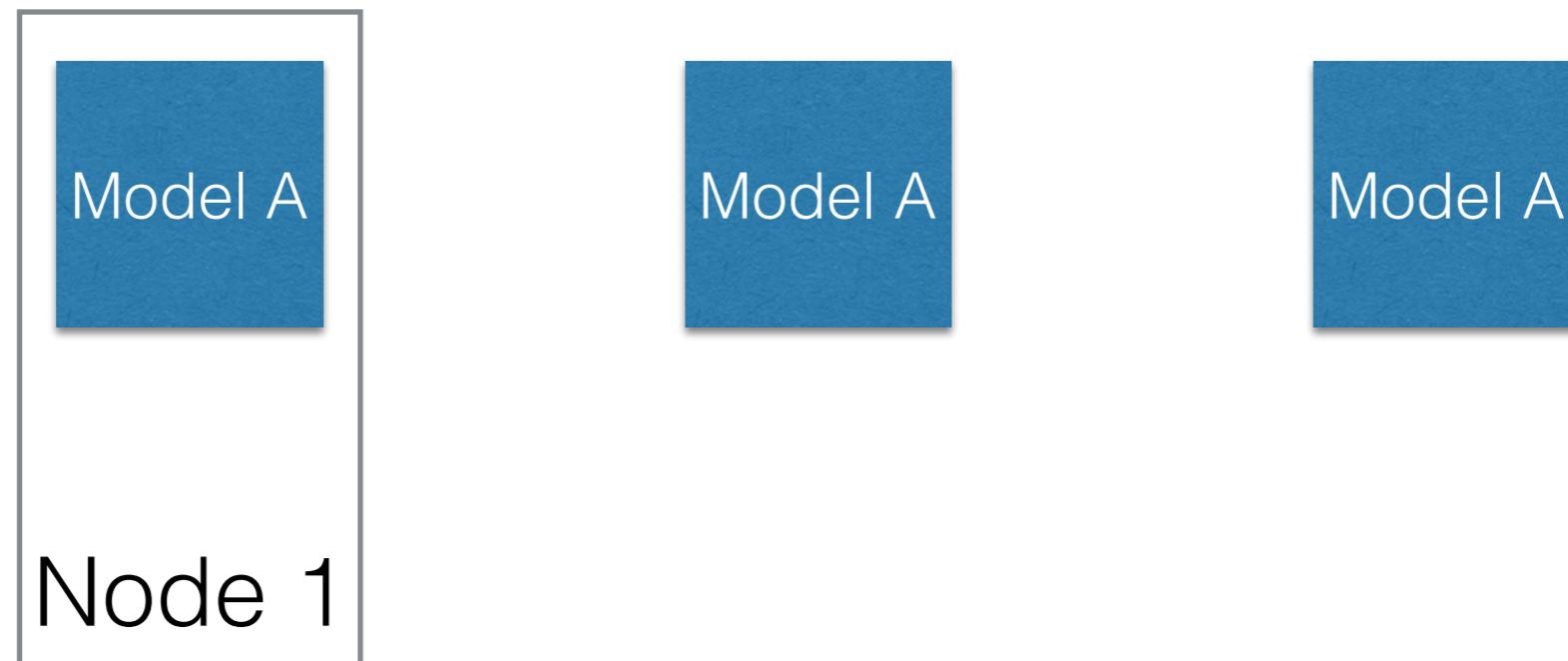
data parallelism

aka. “Jeff Dean style” parameter averaging



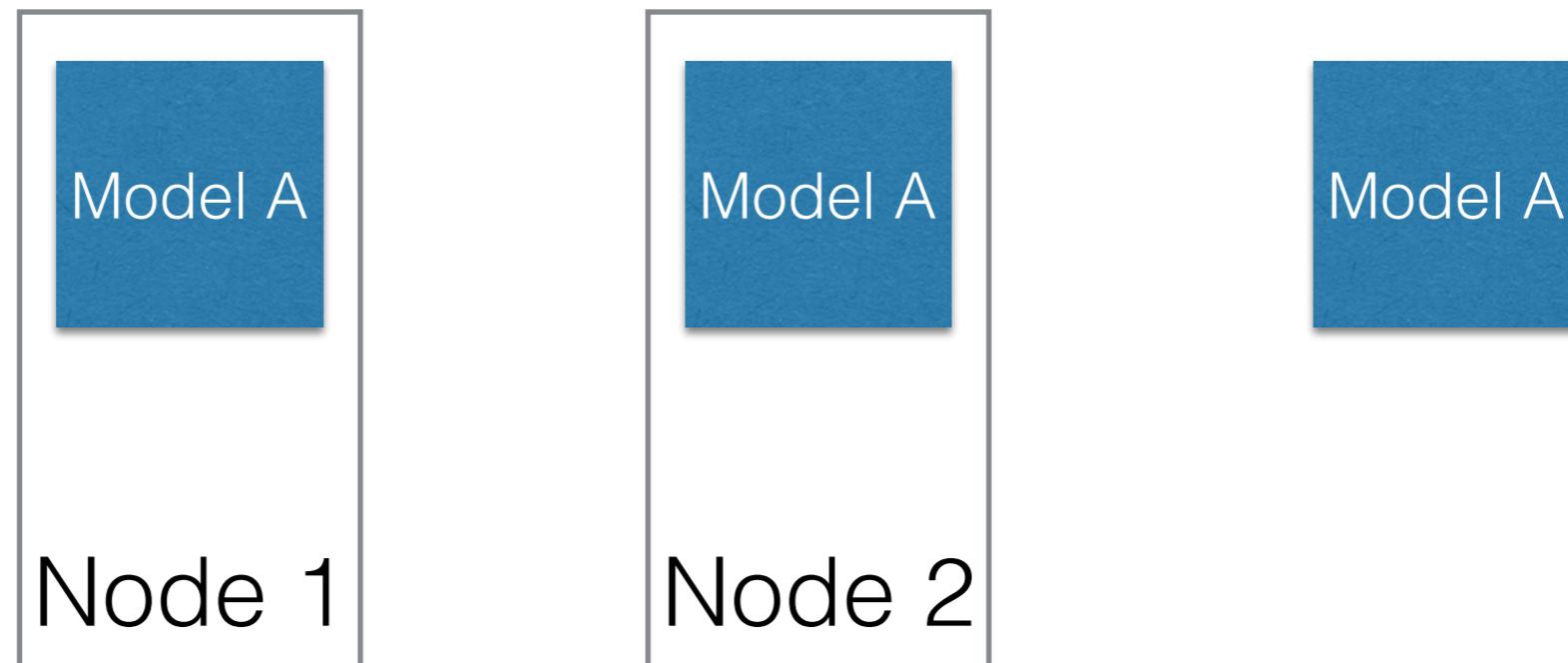
data parallelism

aka. “Jeff Dean style” parameter averaging



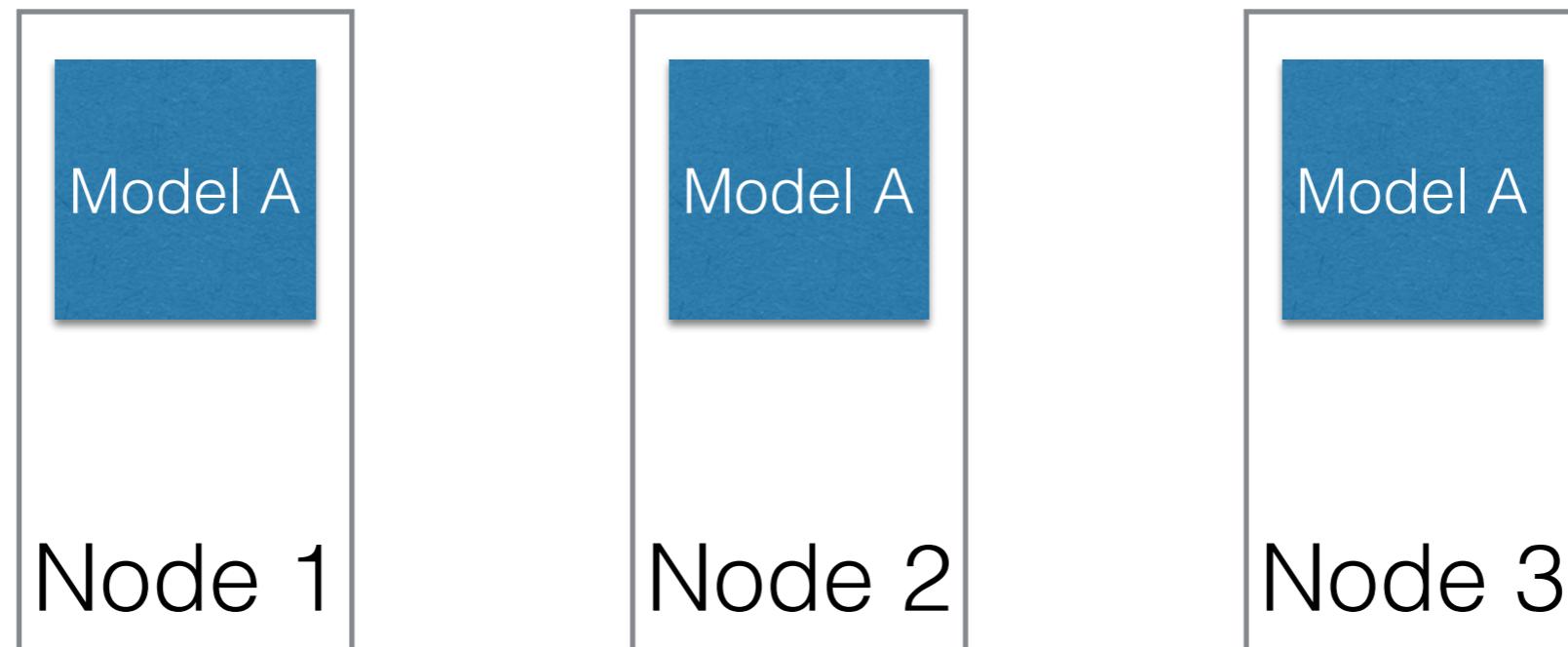
data parallelism

aka. “Jeff Dean style” parameter averaging



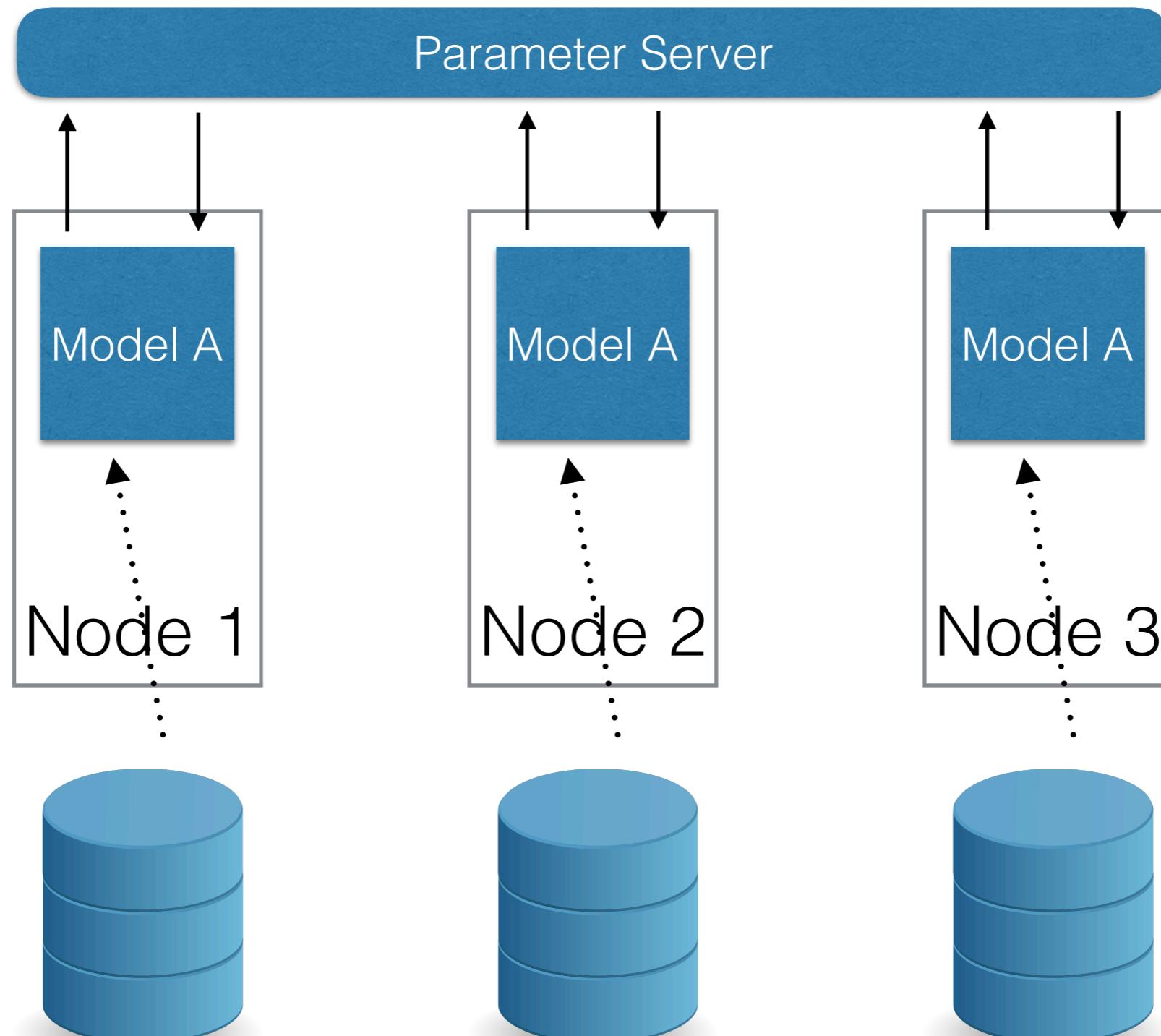
data parallelism

aka. “Jeff Dean style” parameter averaging

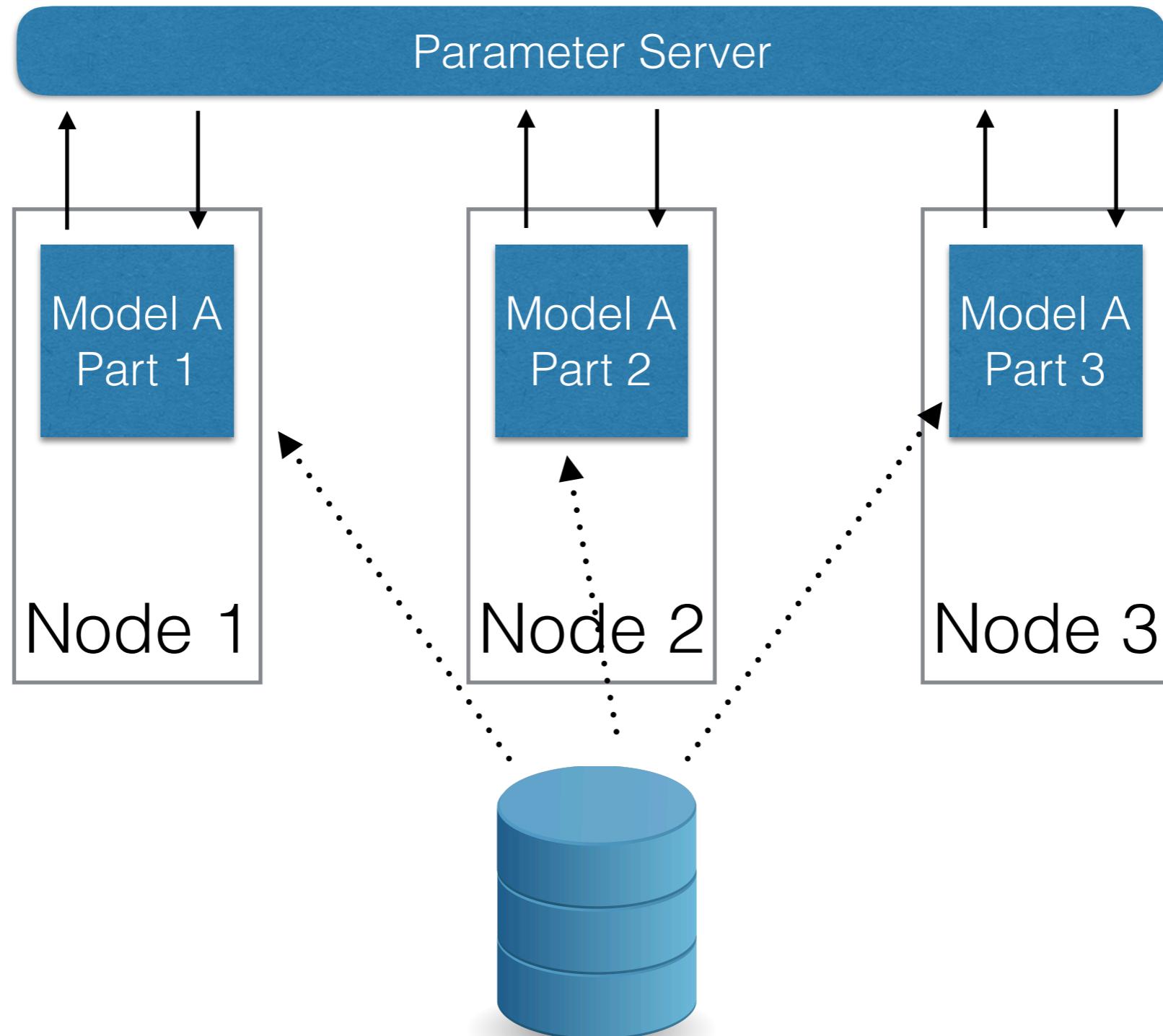


data parallelism

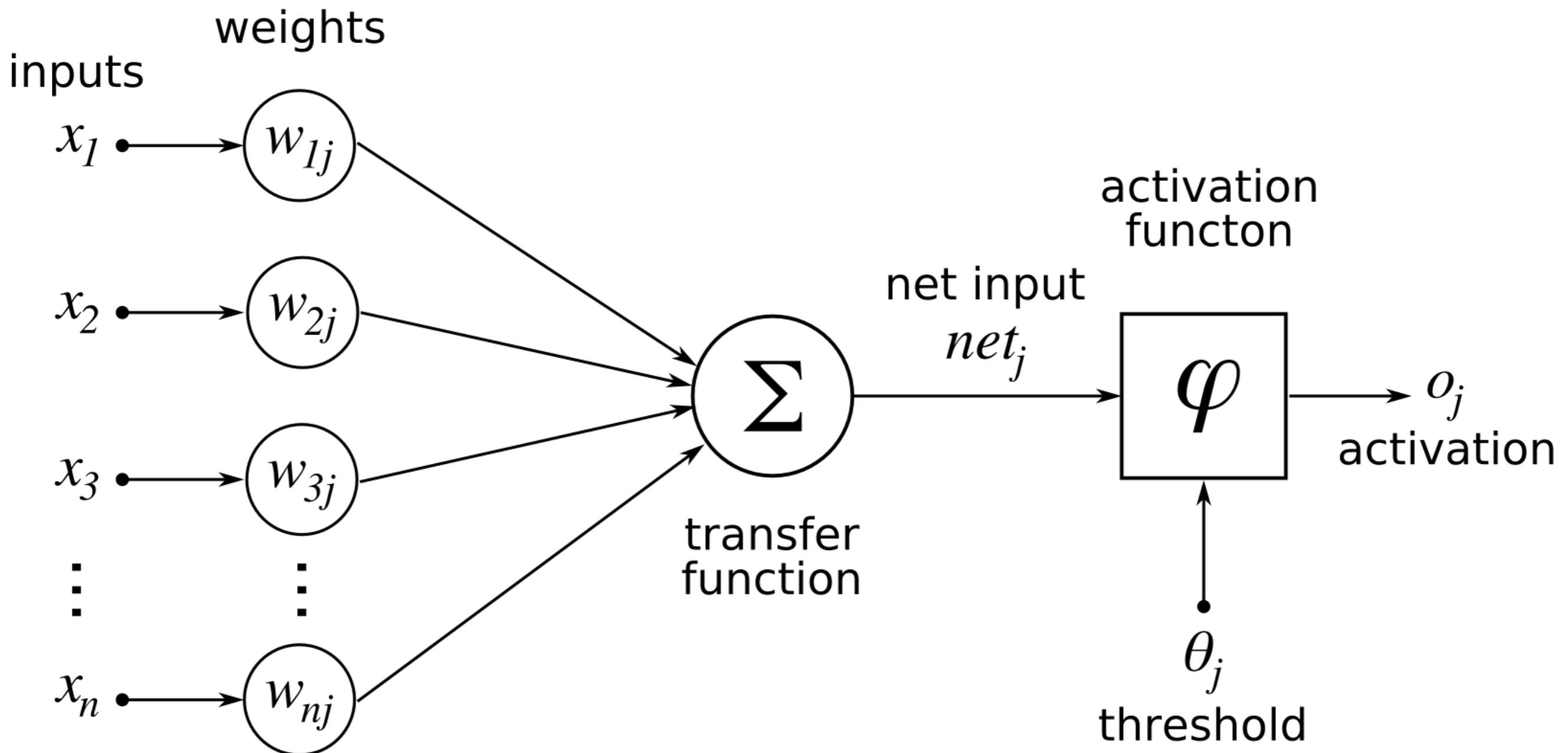
aka. “Jeff Dean style” parameter averaging



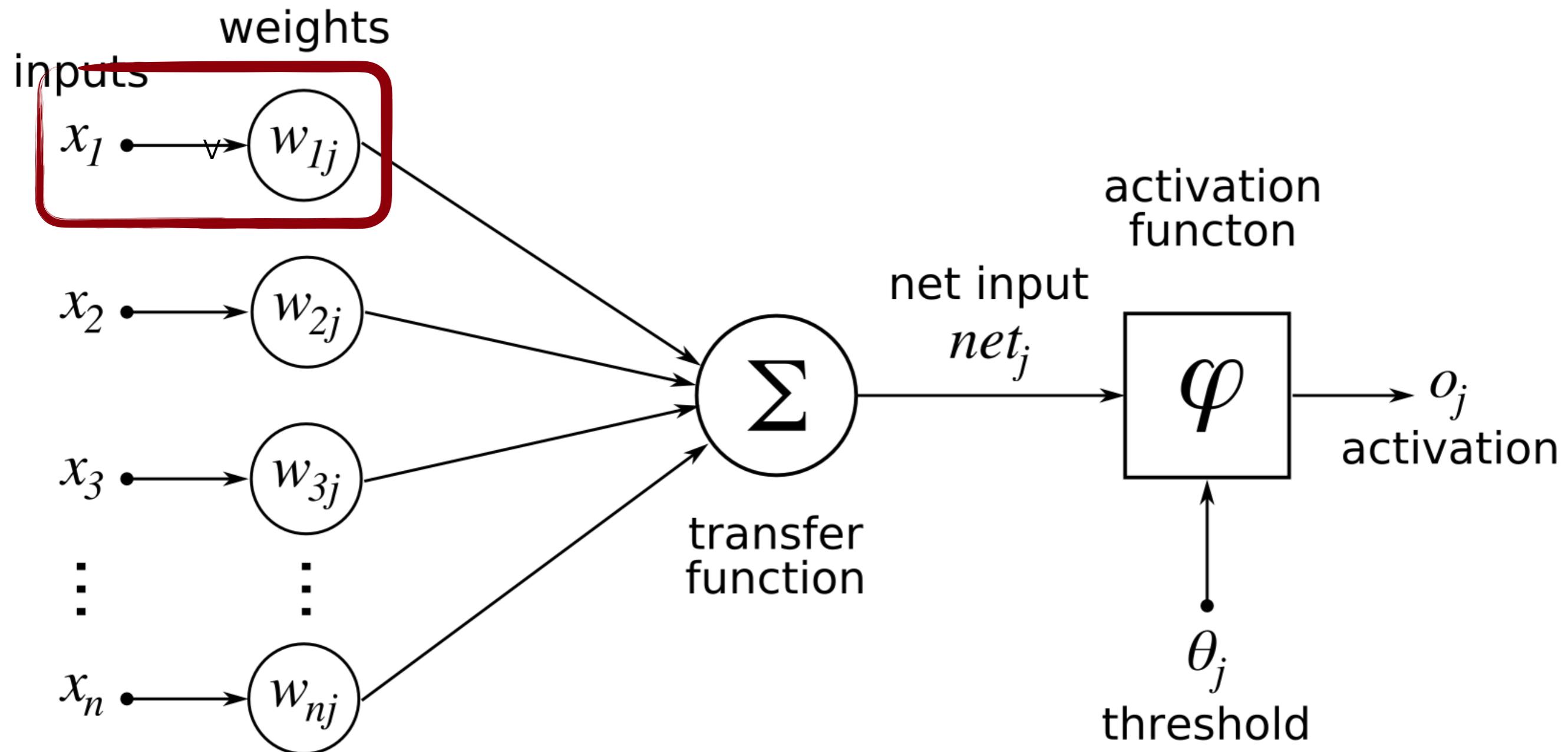
intra-model parallelism



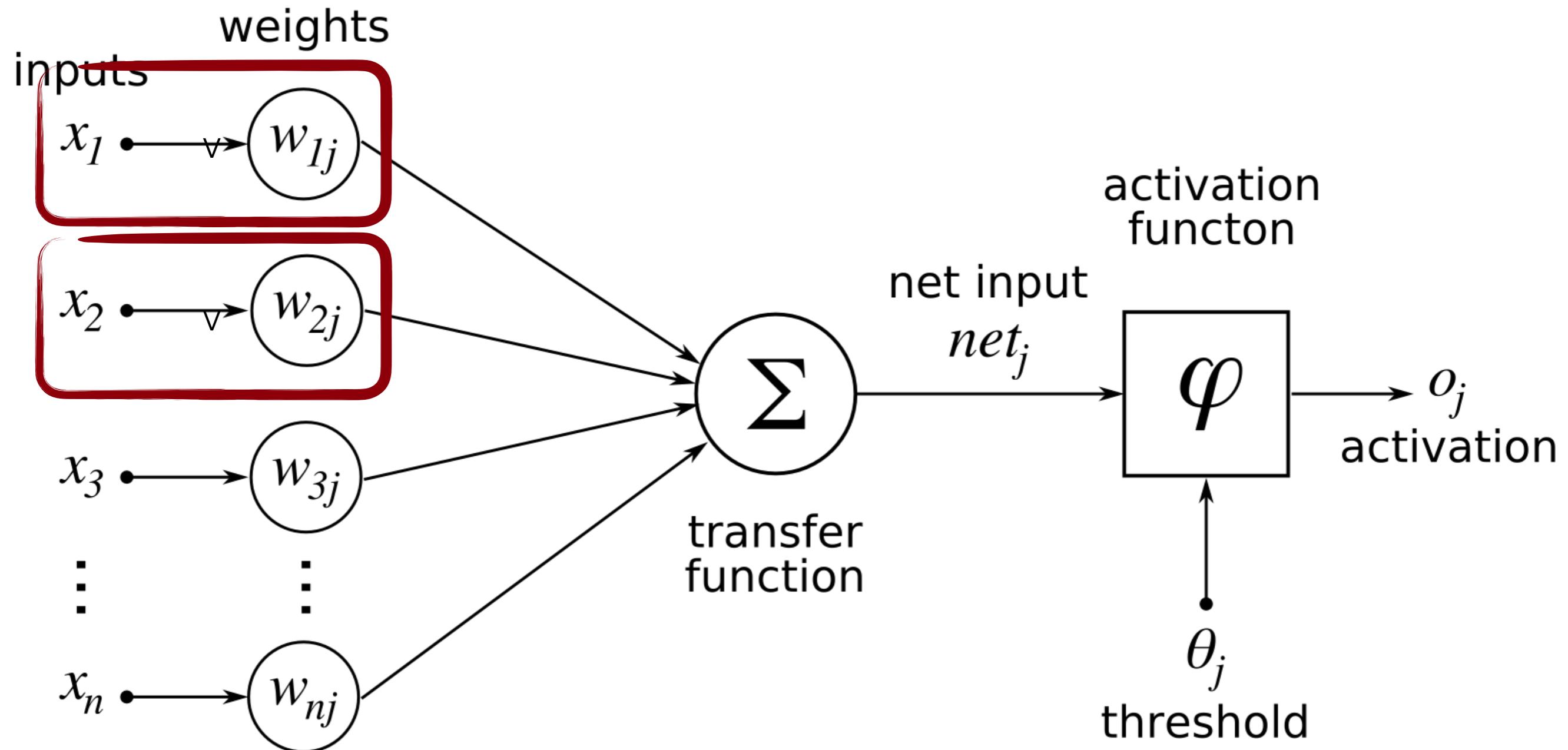
intra-model parallelism



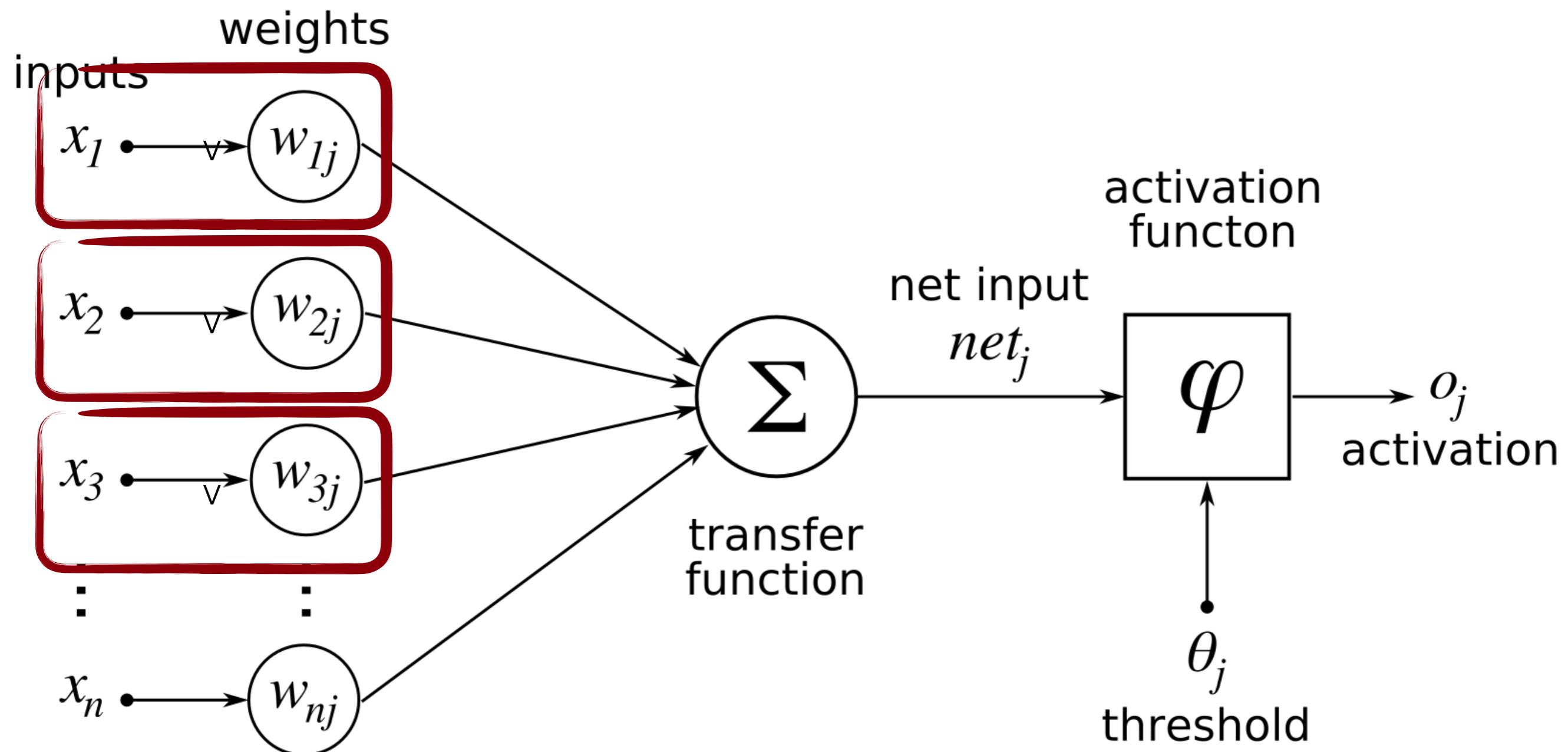
intra-model parallelism



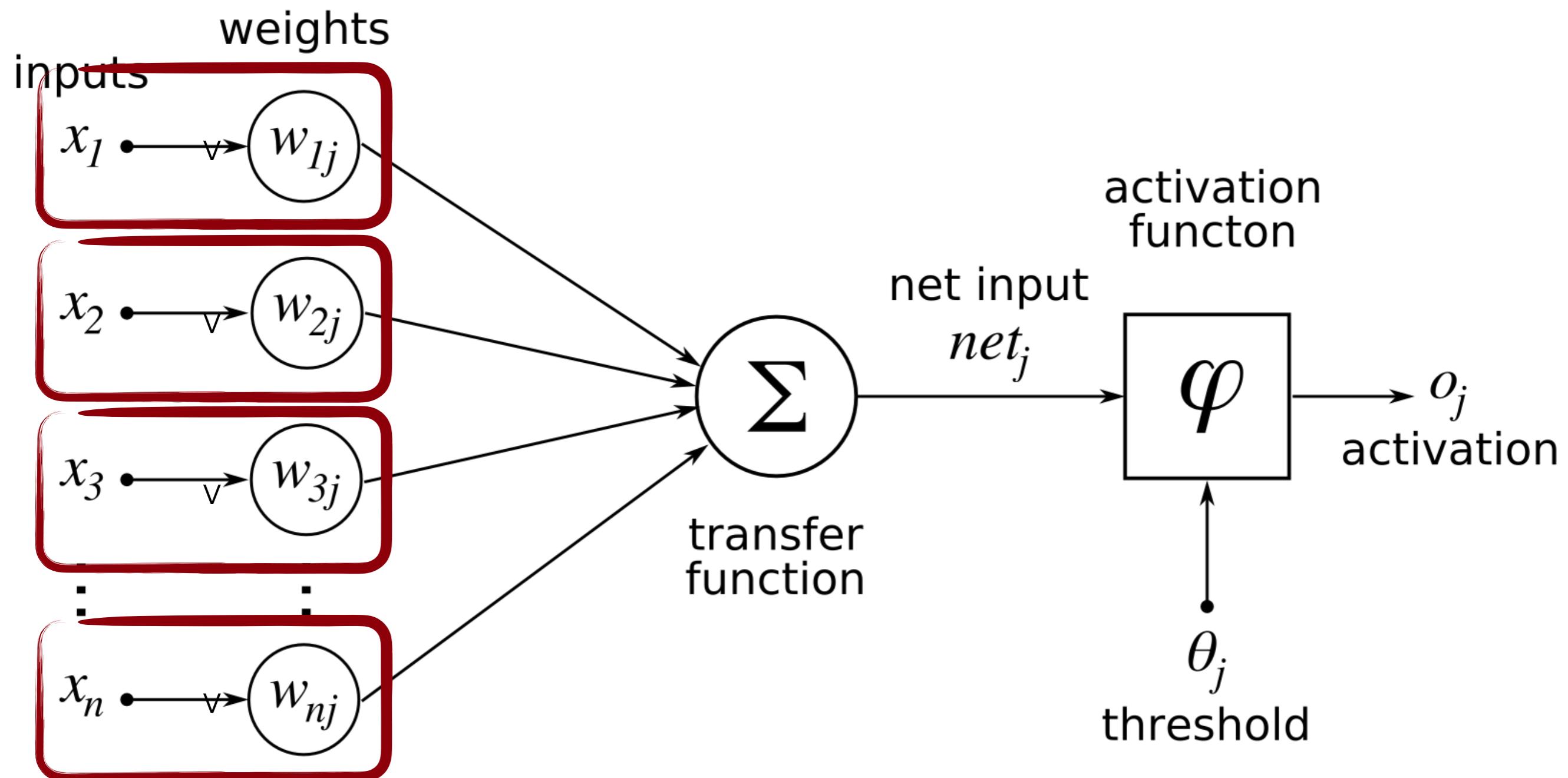
intra-model parallelism



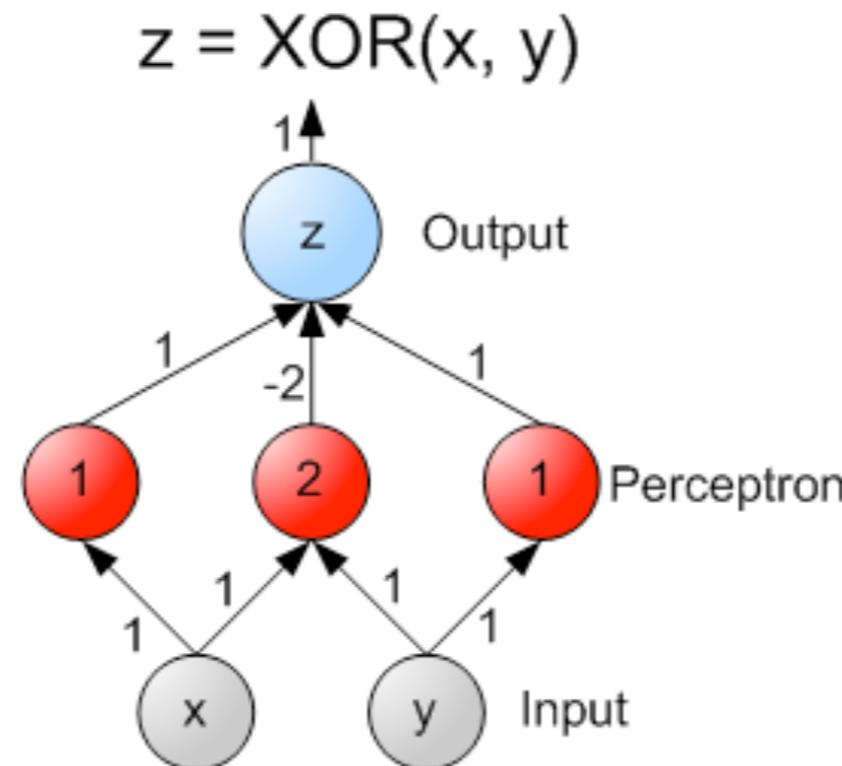
intra-model parallelism



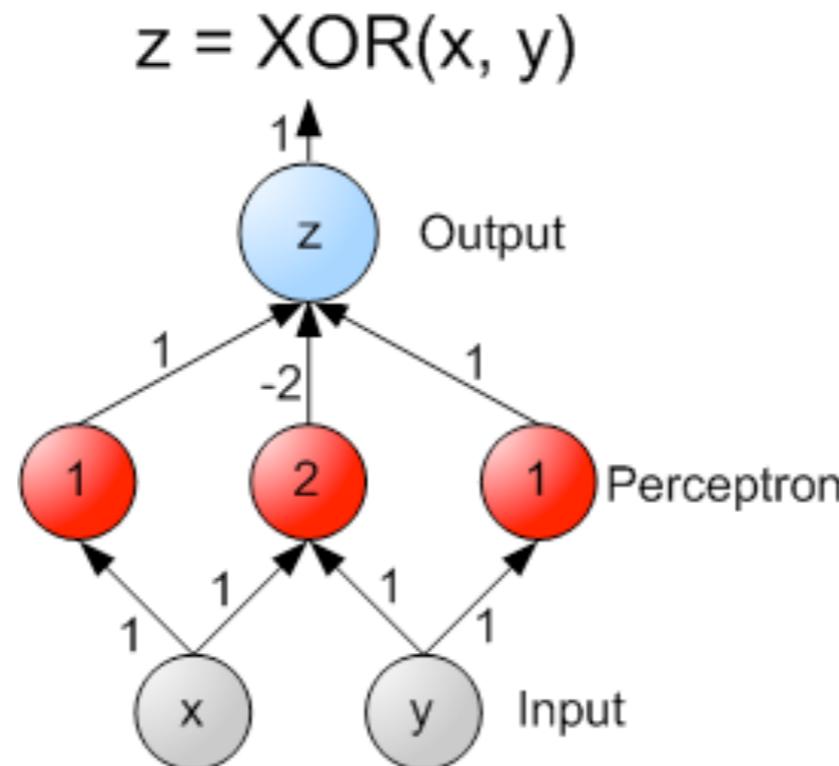
intra-model parallelism



pipelined parallelism

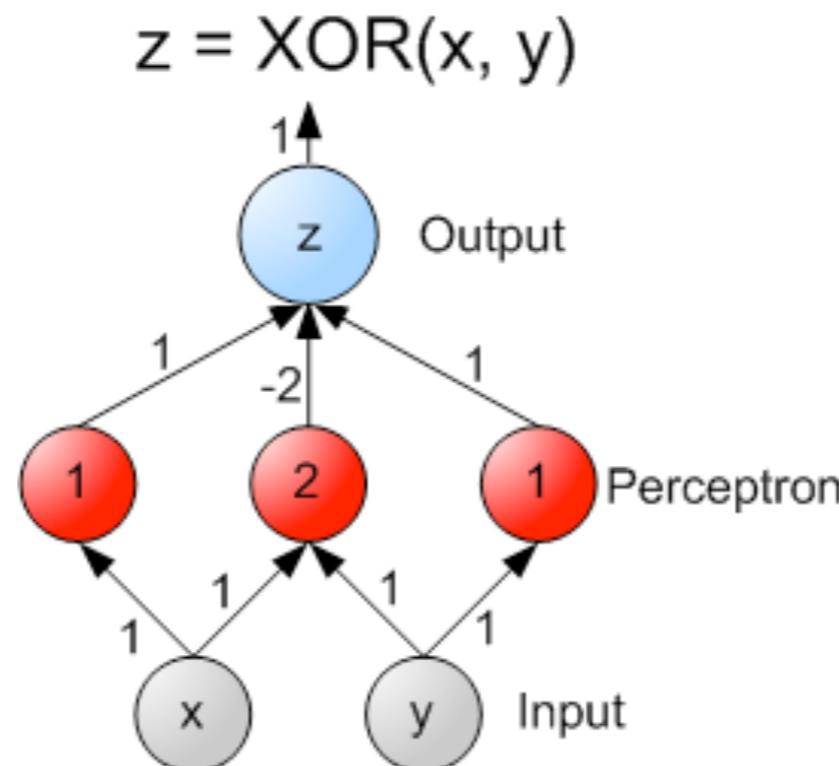


pipelined parallelism



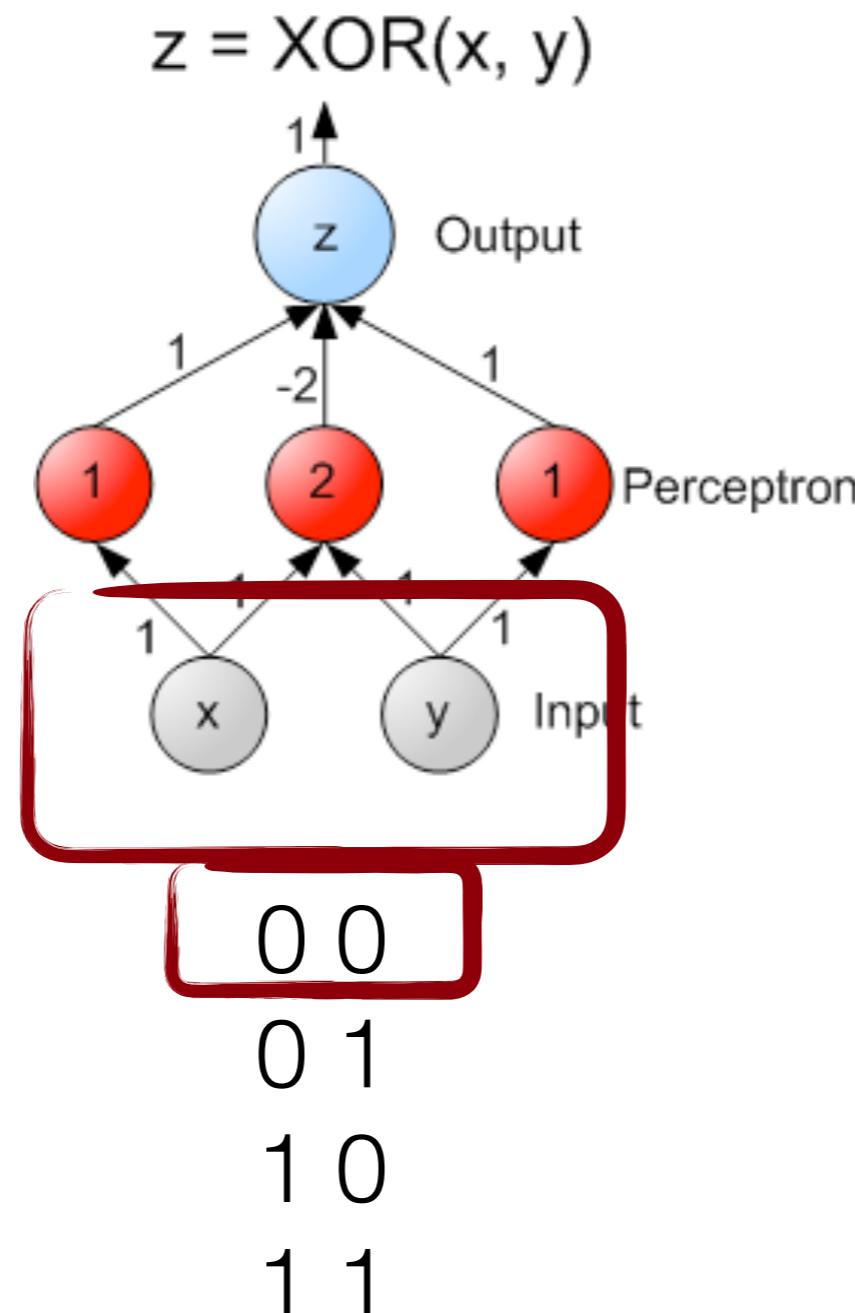
0 0
0 1
1 0
1 1

pipelined parallelism

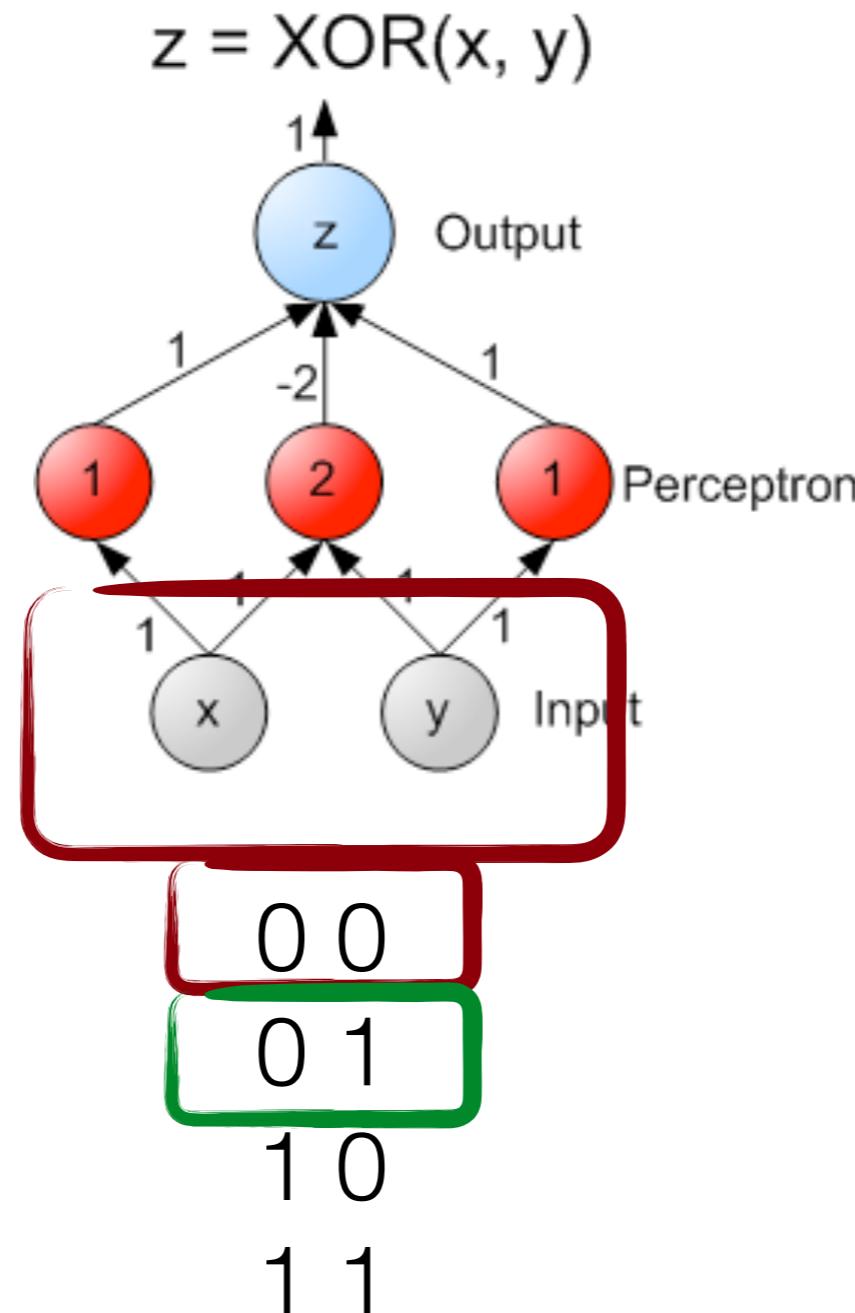


0 0
0 1
1 0
1 1

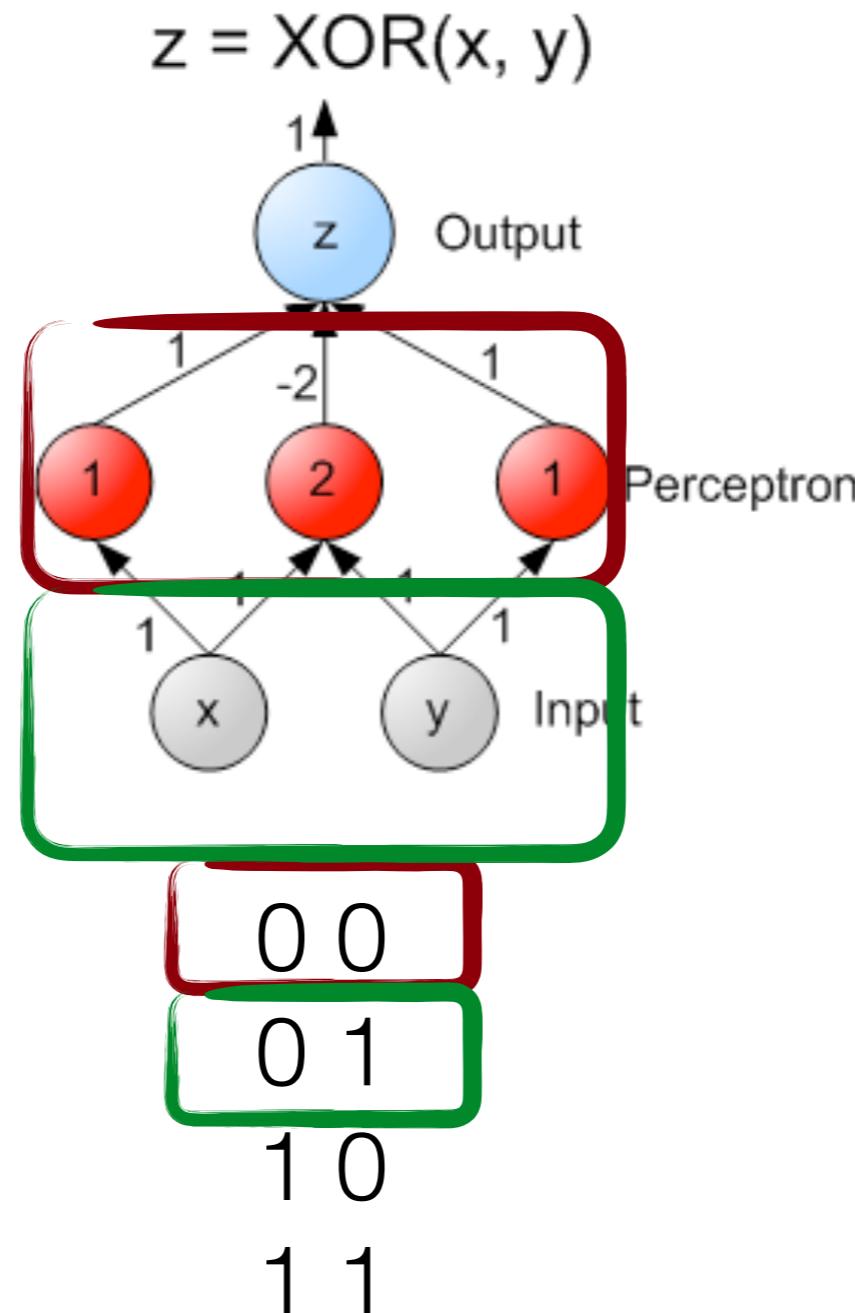
pipelined parallelism



pipelined parallelism



pipelined parallelism



Driver JVM

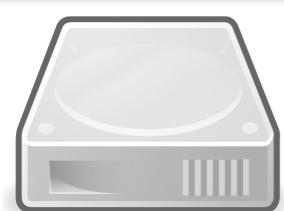
Apache Spark

Compute Node

Executor JVM

Executor JVM

Executor JVM



Compute Node

Executor JVM

Executor JVM

Executor JVM



Compute Node

Executor JVM

Executor JVM

Executor JVM



Compute Node

Executor JVM

Executor JVM

Executor JVM

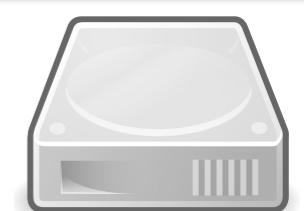


Compute Node

Executor JVM

Executor JVM

Executor JVM



DeepLearning4J

`new MultiLayerNetwork(conf);`

vs.

`new SparkDl4jMultiLayer(sc, conf, tm);`

```
TrainingMaster tm = new ParameterAveragingTrainingMaster.Builder(int dataSet0bjectSize)
    ... (your configuration here)
    .build();
```

DeepLearning4J

`new MultiLayerNetwork(Conf);`

vs.

`new SparkDL4jMultiLayer(sc, conf, tm);`

```
TrainingMaster tm = new ParameterAveragingTrainingMaster.Builder(int dataSetObjectSize)
    ... (your configuration here)
    build()
```

data parallelism

ND4J (DeepLearning4J)

- Tensor support (Linear Buffer + Stride)
- Multiple implementations, one interface
 - vectorized c++ code (JavaCPP), off-heap data storage, BLAS (OpenBLAS, Intel MKL, cuBLAS)
 - GPU (CUDA 8)

ND4J (DeepLearning4J)

- Tensor support (Linear Buffer + Stride)
- Multiple implementations, one interface
 - vectorized c++ code (JavaCPP), off-heap data storage, BLAS (OpenBLAS, Intel MKL, cuBLAS)
 - GPU (CUDA 8)

Apache SystemML

- Custom machine learning algorithms
- Declarative ML
- Transparent distribution on data-parallel framework
 - Scale-up
 - Scale-out
- Cost-based optimiser generates low level execution plans

2007-2008: Multiple projects at IBM Research – Almaden involving machine learning on Hadoop.

2009: We form a dedicated team for scalable ML

2009-2010: Through engagements with customers, we observe how data scientists create **ML solutions**.

2007

2008

2009

2010

Research

2011

2012

2013

2014

June 2015: IBM Announces open-source SystemML

November 2015: SystemML enters Apache incubation

June 2016: Second Apache release (0.10)

September 2015: Code available on Github

February 2016: First release (0.9) of Apache SystemML

2015

2016

Apache SystemML

```
U = rand(nrow(X), r, min = -1.0, max = 1.0);
V = rand(r, ncol(X), min = -1.0, max = 1.0);
while(i < mi) {
    i = i + 1; ii = 1;
    if (is_U)
        G = (W * (U %*% V - X)) %*% t(V) + lambda * U;
    else
        G = t(U) %*% (W * (U %*% V - X)) + lambda * V;
    norm_G2 = sum(G ^ 2); norm_R2 = norm_G2;
    R = -G; S = R;
    while(norm_R2 > 10E-9 * norm_G2 & ii <= mii) {
        if (is_U) {
            HS = (W * (S %*% V)) %*% t(V) + lambda * S;
            alpha = norm_R2 / sum (S * HS);
            U = U + alpha * S;
        } else {
            HS = t(U) %*% (W * (U %*% S)) + lambda * S;
            alpha = norm_R2 / sum (S * HS);
            V = V + alpha * S;
        }
        R = R - alpha * HS;
        old_norm_R2 = norm_R2; norm_R2 = sum(R ^ 2);
        S = R + (norm_R2 / old_norm_R2) * S;
        ii = ii + 1;
    }
    is_U = ! is_U;
}
```

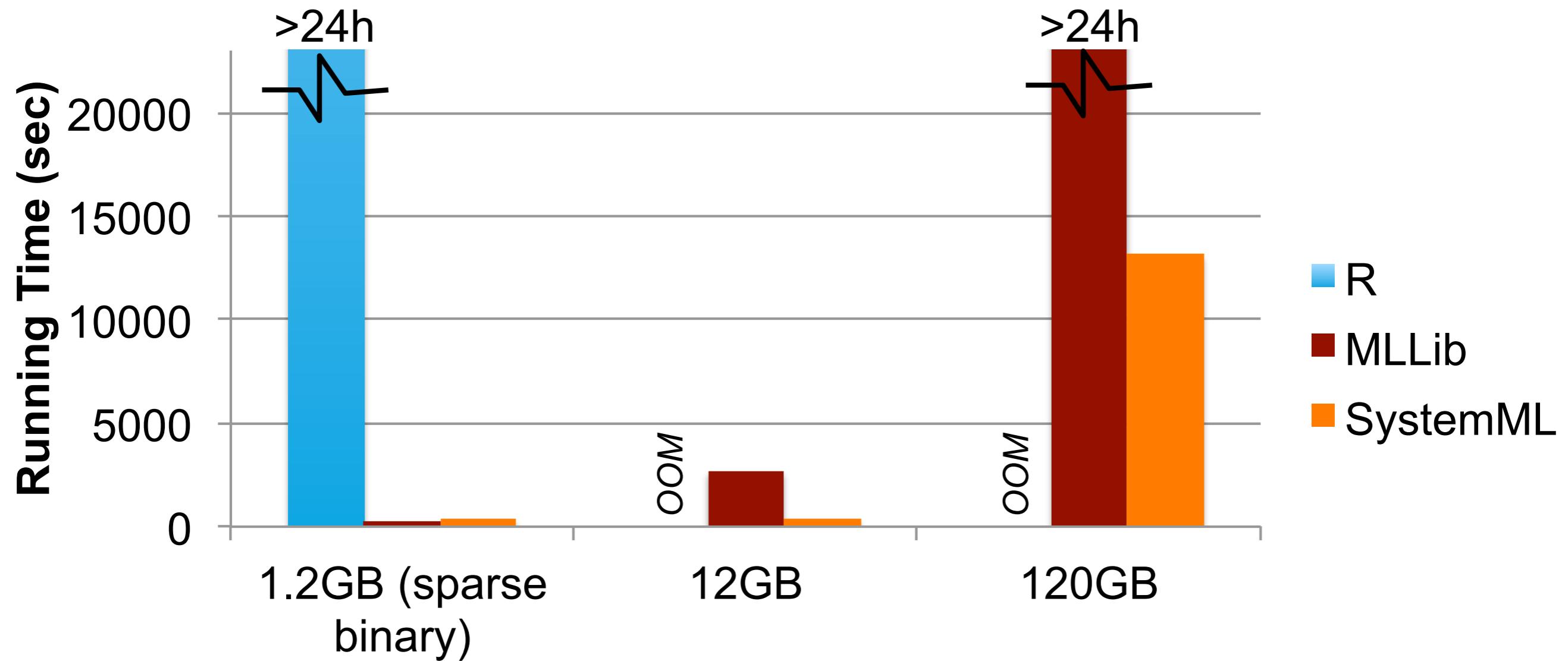
Apache SystemML

```

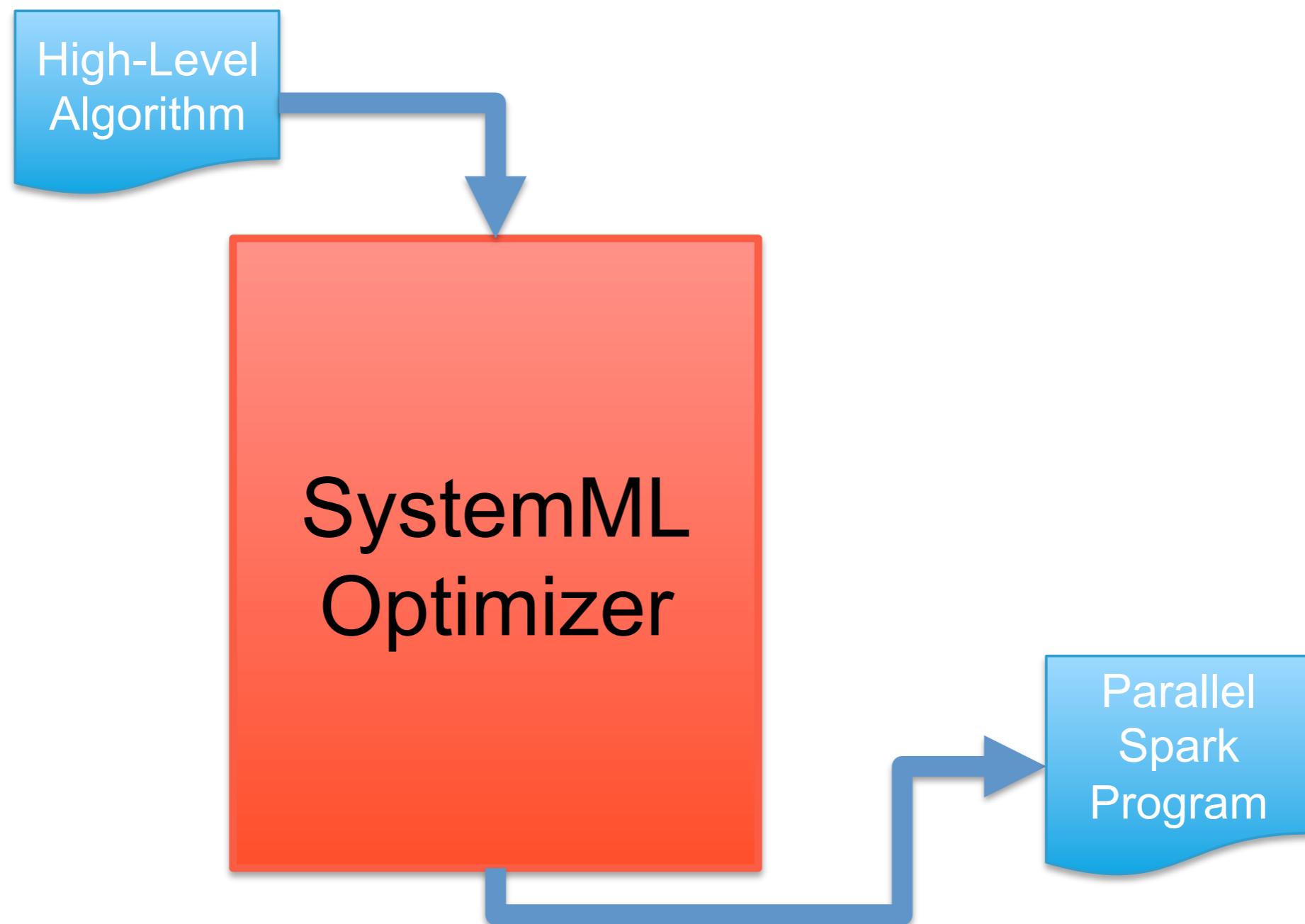
U = rand(nrow(X), r, min = -1.0, max = 1.0);
V = rand(r, ncol(X), min = -1.0, max = 1.0);
while(i < mi) {
  i = i + 1; ii = 1;
  if (is_U)
    G = (W * (U %*% V - X)) %*% t(V) + lambda * U;
  else
    G = t(U) %*% (W * (U %*% V - X)) + lambda * V;
  norm_G2 = sum(G ^ 2); norm_R2 = norm_G2;
  R = -G; S = R;
  while(norm_R2 > 10E-9 * norm_G2 & ii <= mii) {
    if (is_U) {
      HS = (W * (S %*% V)) %*% t(V) + lambda * S;
      alpha = norm_R2 / sum (S * HS);
      U = U + alpha * S;
    } else {
      HS = t(U) %*% (W * (U %*% S)) + lambda * S;
      alpha = norm_R2 / sum (S * HS);
      V = V + alpha * S;
    }
    R = R - alpha * HS;
    old_norm_R2 = norm_R2; norm_R2 = sum(R ^ 2);
    S = R + (norm_R2 / old_norm_R2) * S;
    ii = ii + 1;
  }
  is_U = ! is_U;
}

```

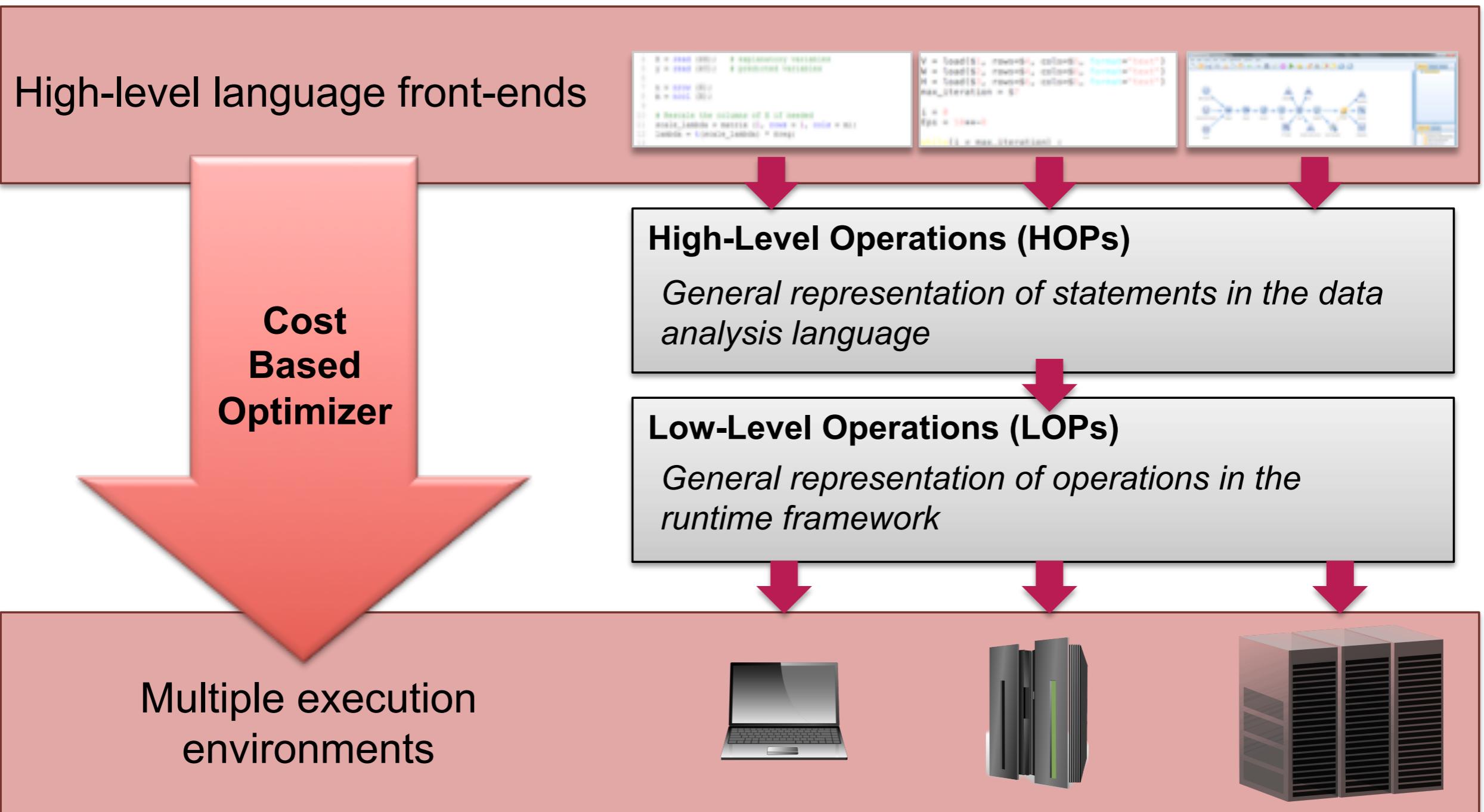
SystemML:
 compile and run at scale
 no performance code needed!



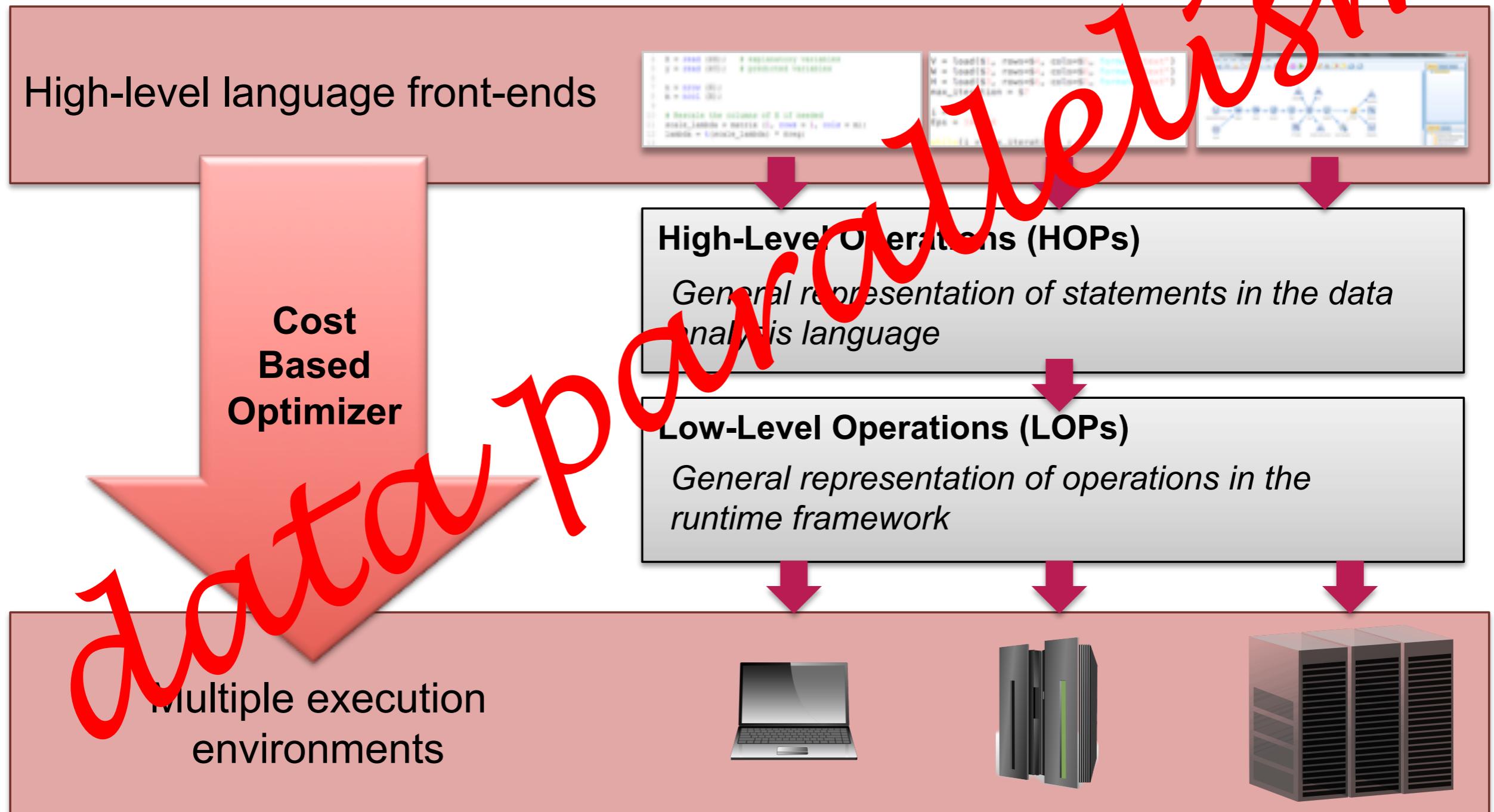
Architecture



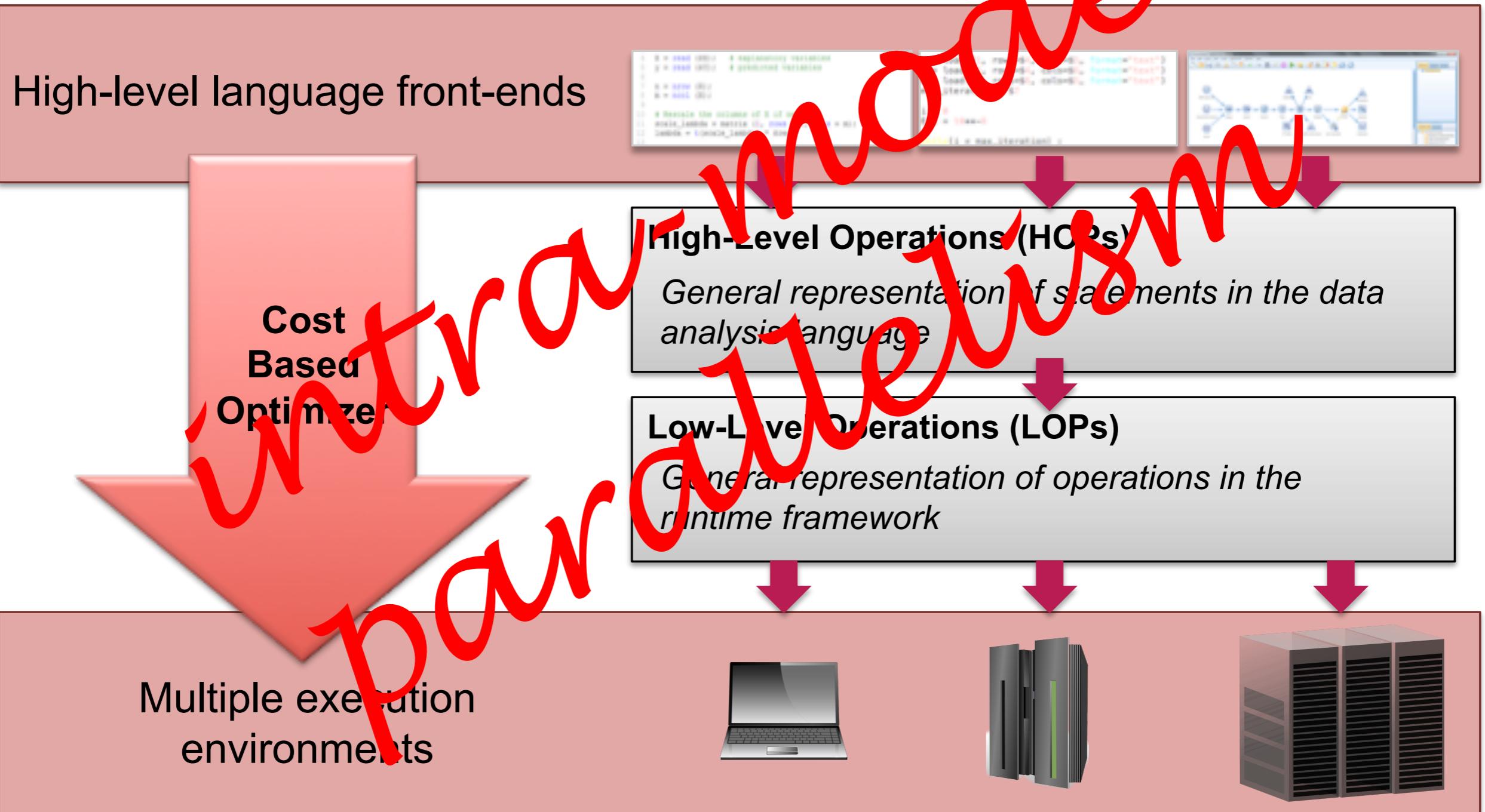
Architecture



Architecture



Architecture



TensorFrames

Compute Node



Compute Node



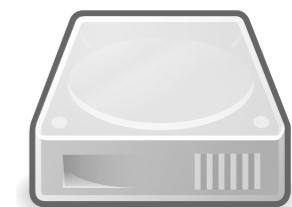
Compute Node



Compute Node

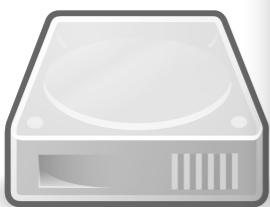


Compute Node



TensorFrames

Compute Node

Executor
JVMExecutor
JVMExecutor
JVM

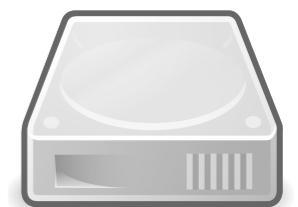
```
# Build a DataFrame of vectors
data = [Row(y=[float(y), float(-y)]) for y in range(10)]
df = sqlContext.createDataFrame(data)
# Because the dataframe contains vectors, we need to analyze it first to find the
# dimensions of the vectors.
df2 = tfs.analyze(df)

# The information gathered by TF can be printed to check the content:
tfs.print_schema(df2)
# TF has inferred that y contains vectors of size 2
# root
# |-- y: array (nullable = false) DoubleType[?,2]

# Let's use the analyzed dataframe to compute the sum and the elementwise minimum
# of all the vectors:
# First, let's make a copy of the 'y' column. This will be very cheap in Spark 2.0+
df3 = df2.select(df2.y, df2.y.alias("z"))

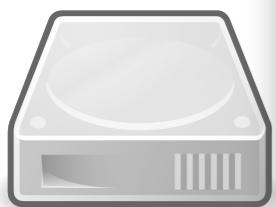
with tf.Graph().as_default() as g:
    # The placeholders. Note the special name that end with '_input':
    y_input = tfs.block(df3, 'y', tf_name="y_input")
    z_input = tfs.block(df3, 'z', tf_name="z_input")
    y = tf.reduce_sum(y_input, [0], name='y')
    z = tf.reduce_min(z_input, [0], name='z')
    # The resulting dataframe
    (data_sum, data_min) = tfs.reduce_blocks([y, z], df3)
```

Compute Node

Executor
JVMExecutor
JVMExecutor
JVM

TensorFrames

Compute Node



```
# Build a DataFrame of vectors
data = [Row(y=[float(y), float(-y)]) for y in range(10)]
df = sqlContext.createDataFrame(data)
# Because the dataframe contains vectors, we need to analyze it first to find the
# dimensions of the vectors.
df2 = tfs.analyze(df)

# The information gathered by TF can be printed to check the content:
tfs.print_schema(df2)
# TF has inferred that y contains vectors of size 2
# root
# |-- y: array(nullable = false) DoubleType[2:2]

# Let's use the analyzed dataframe to compute the sum and the elementwise minimum
# of all the vectors:
# First let's make a copy of the 'y' column. This will be very cheap in Spark 2.0+
df3 = df2.select(df2.y, df2.y.alias("z"))
with tf.Graph().as_default() as g:
    # The placeholders. Note the special name that end with '_input':
    y_input = tfs.block(df3, 'y', tf_name="y_input")
    z_input = tfs.block(df3, 'z', tf_name="z_input")
    y = tf.reduce_sum(y_input, [0], name='y')
    z = tf.reduce_min(z_input, [0], name='z')
    # The resulting dataframe
    (data_sum, data_min) = tfs.reduce_blocks([y, z], df3)
```

parallelism

interoperability

model

Compute Node



Parameter
Server on
Driver



TensorSpark

Compute Node

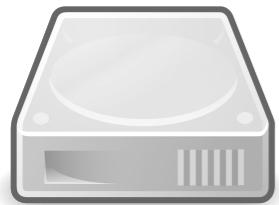
Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

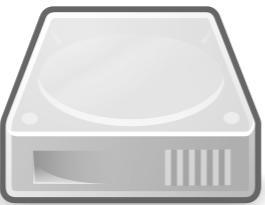
Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

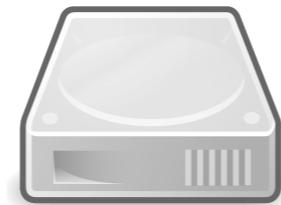
Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

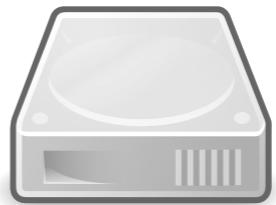
Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

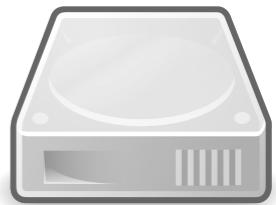
Executor
JVM



Executor
JVM



Executor
JVM



Parameter
Server on
Driver



TensorSpark

Compute Node

Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

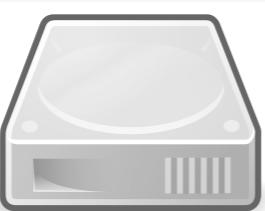
Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

Executor
JVM



Executor
JVM



Executor
JVM



Compute Node

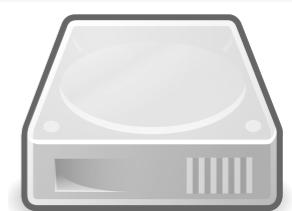
Executor
JVM



Executor
JVM



Executor
JVM

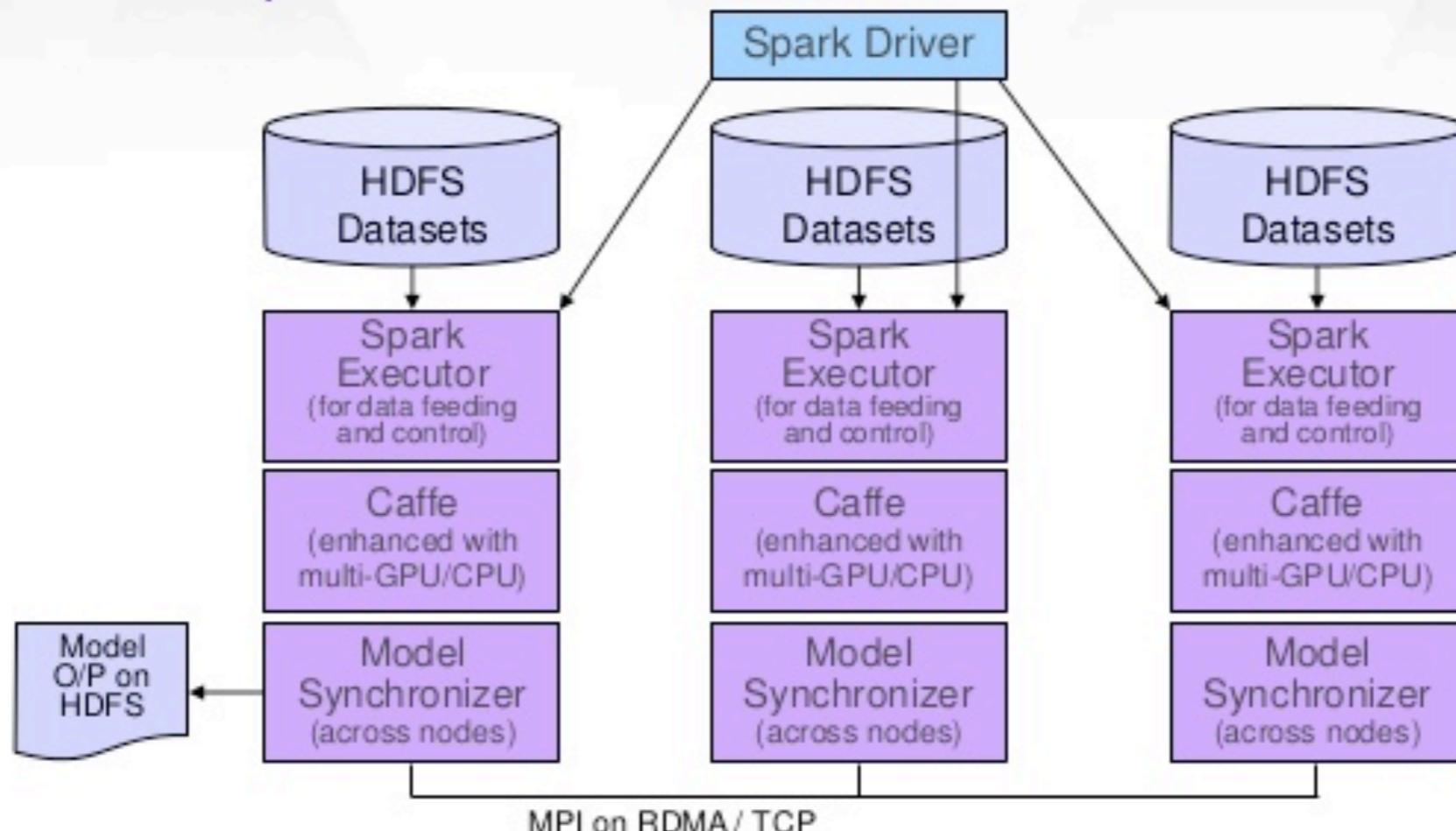


data parallelism

compute parallelism

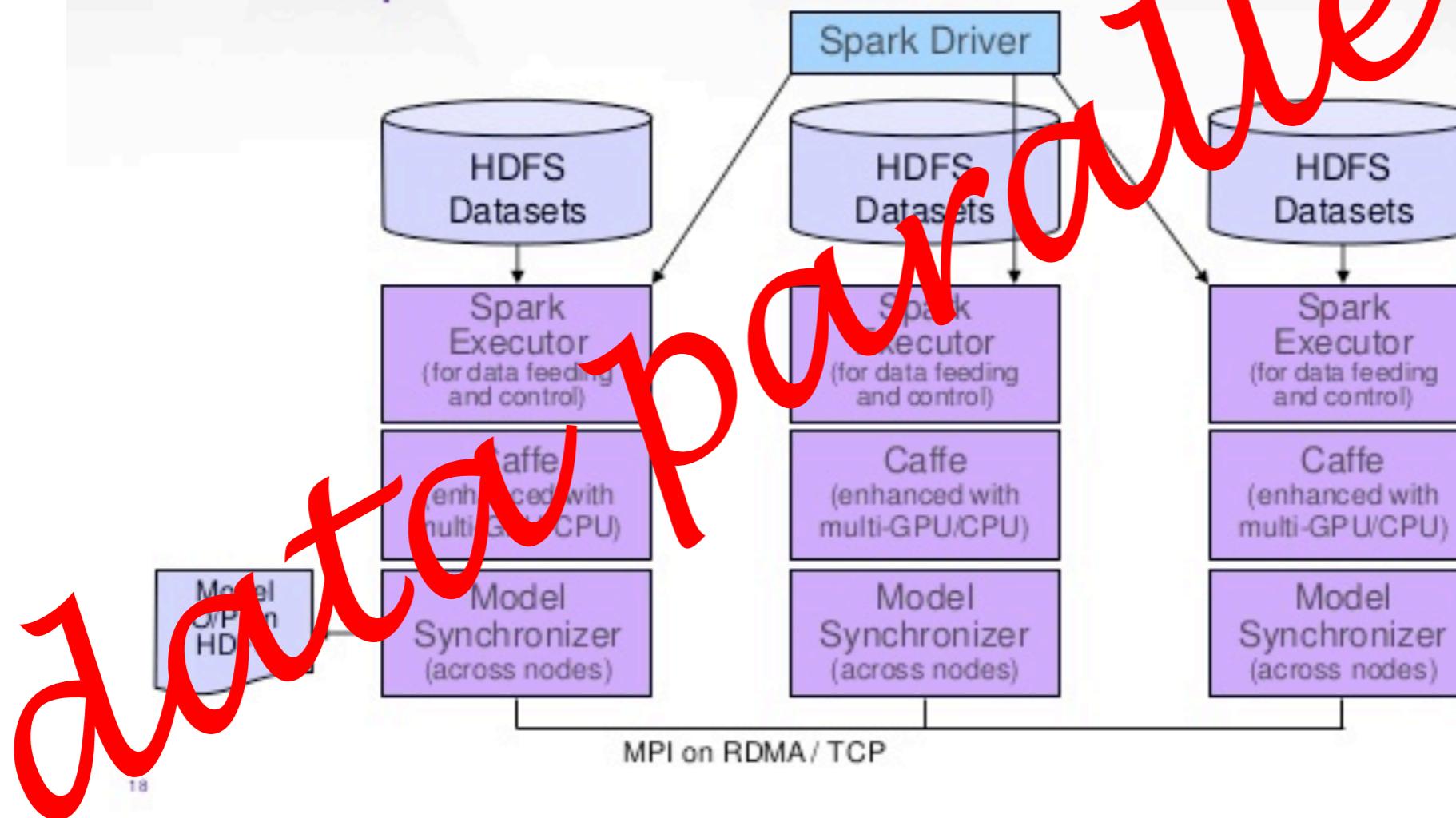
CaffeOnSpark

CaffeOnSpark Architecture – Common Cluster



CaffeOnSpark

CaffeOnSpark Architecture – Common Cluster



Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me



Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk - don't kill me - Disclaimer: This is a beta talk

Romeo Kienzler