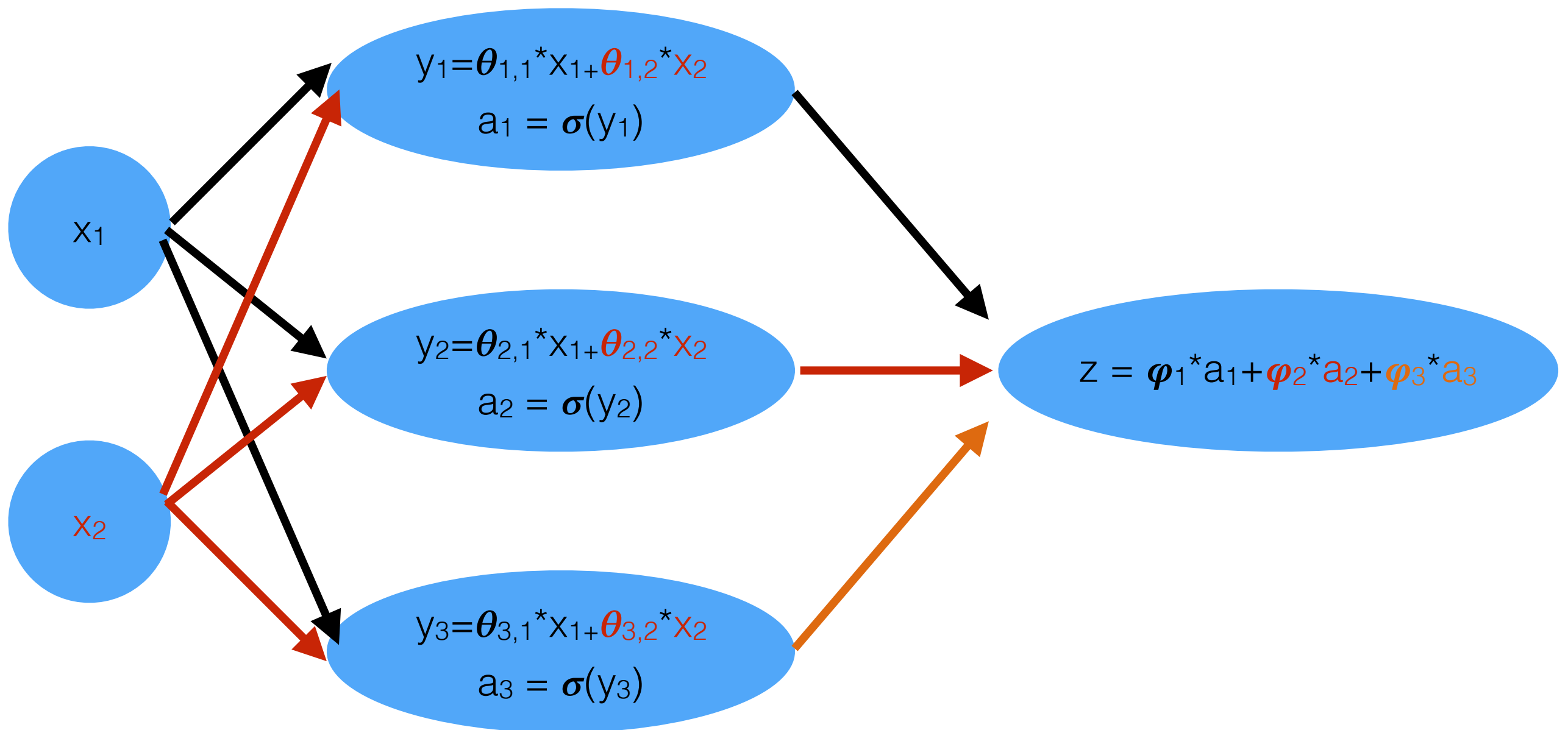# Back-propagation and Parallel Implementation

# Overview

- Back-propagation review
- MPI Implementation
- Parameter Server Implementation
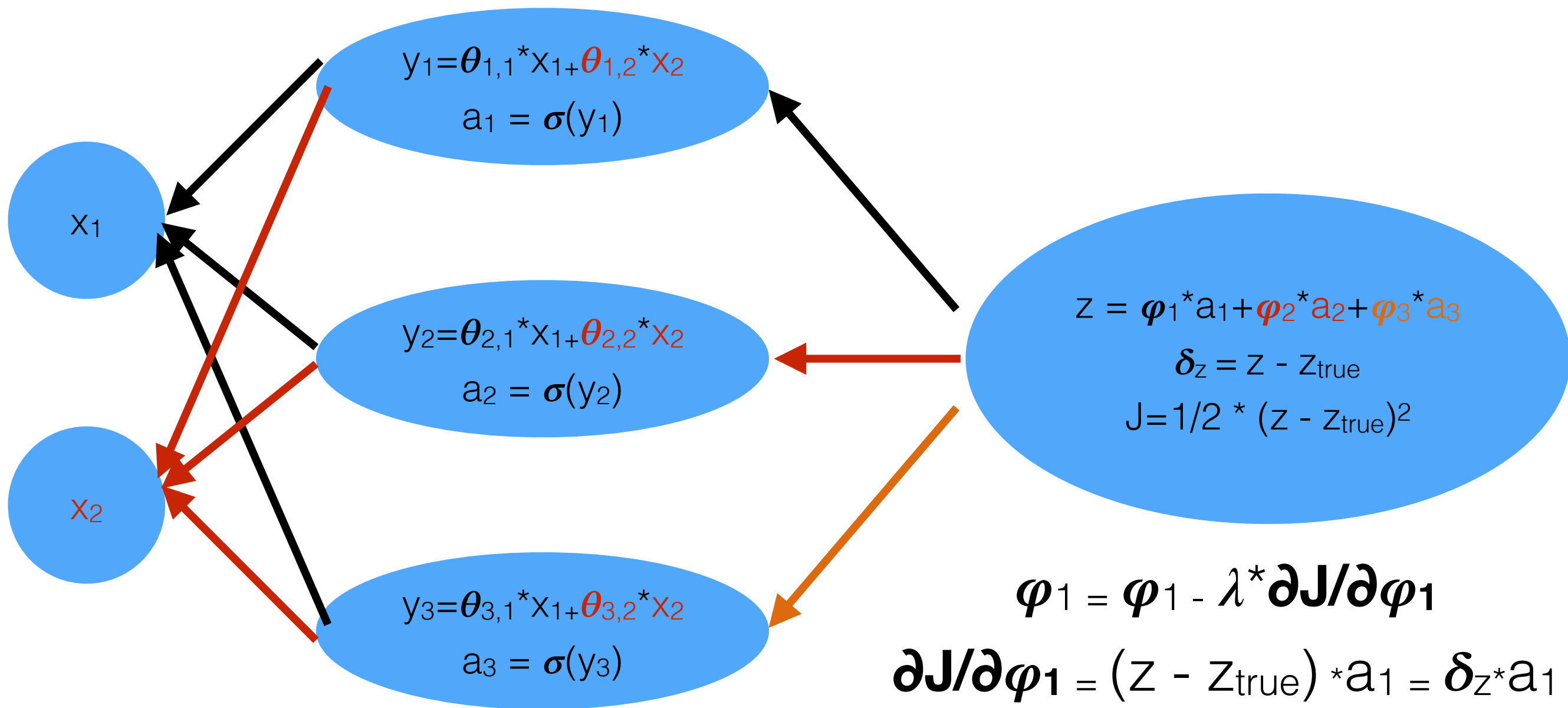
# Back-propagation Review



$$y_1 = \theta_{1,1}*x_1 + \theta_{1,2}*x_2$$
$$a_1 = \sigma(y_1)$$

$$y_2 = \theta_{2,1}*x_1 + \theta_{2,2}*x_2$$
$$a_2 = \sigma(y_2)$$

$$y_3 = \theta_{3,1}*x_1 + \theta_{3,2}*x_2$$
$$a_3 = \sigma(y_3)$$

$$z = \varphi_1*a_1 + \varphi_2*a_2 + \varphi_3*a_3$$

$x_1$

$x_2$

$$\theta = \begin{array}{|c|c|} \hline \theta_{1,1} & \theta_{1,2} \\ \hline \theta_{2,1} & \theta_{2,2} \\ \hline \theta_{3,1} & \theta_{3,2} \\ \hline \end{array}$$

$$\varphi = \begin{array}{|c|} \hline \varphi_1 \\ \hline \varphi_2 \\ \hline \varphi_3 \\ \hline \end{array}$$

# Back-propagation Review



$y_1 = \theta_{1,1} * x_1 + \theta_{1,2} * x_2$

$a_1 = \sigma(y_1)$

$y_2 = \theta_{2,1} * x_1 + \theta_{2,2} * x_2$

$a_2 = \sigma(y_2)$

$y_3 = \theta_{3,1} * x_1 + \theta_{3,2} * x_2$

$a_3 = \sigma(y_3)$

$x_1$

$x_2$

$z = \varphi_1 * a_1 + \varphi_2 * a_2 + \varphi_3 * a_3$

$\delta_z = z - z_{true}$

$J = 1/2 * (z - z_{true})^2$

$\varphi_1 = \varphi_1 - \lambda * \partial J / \partial \varphi_1$

$\partial J / \partial \varphi_1 = (z - z_{true}) * a_1 = \delta_z * a_1$

$$\theta = \begin{array}{|c|c|} \hline \theta_{1,1} & \theta_{1,2} \\ \hline \theta_{2,1} & \theta_{2,2} \\ \hline \theta_{3,1} & \theta_{3,2} \\ \hline \end{array}$$
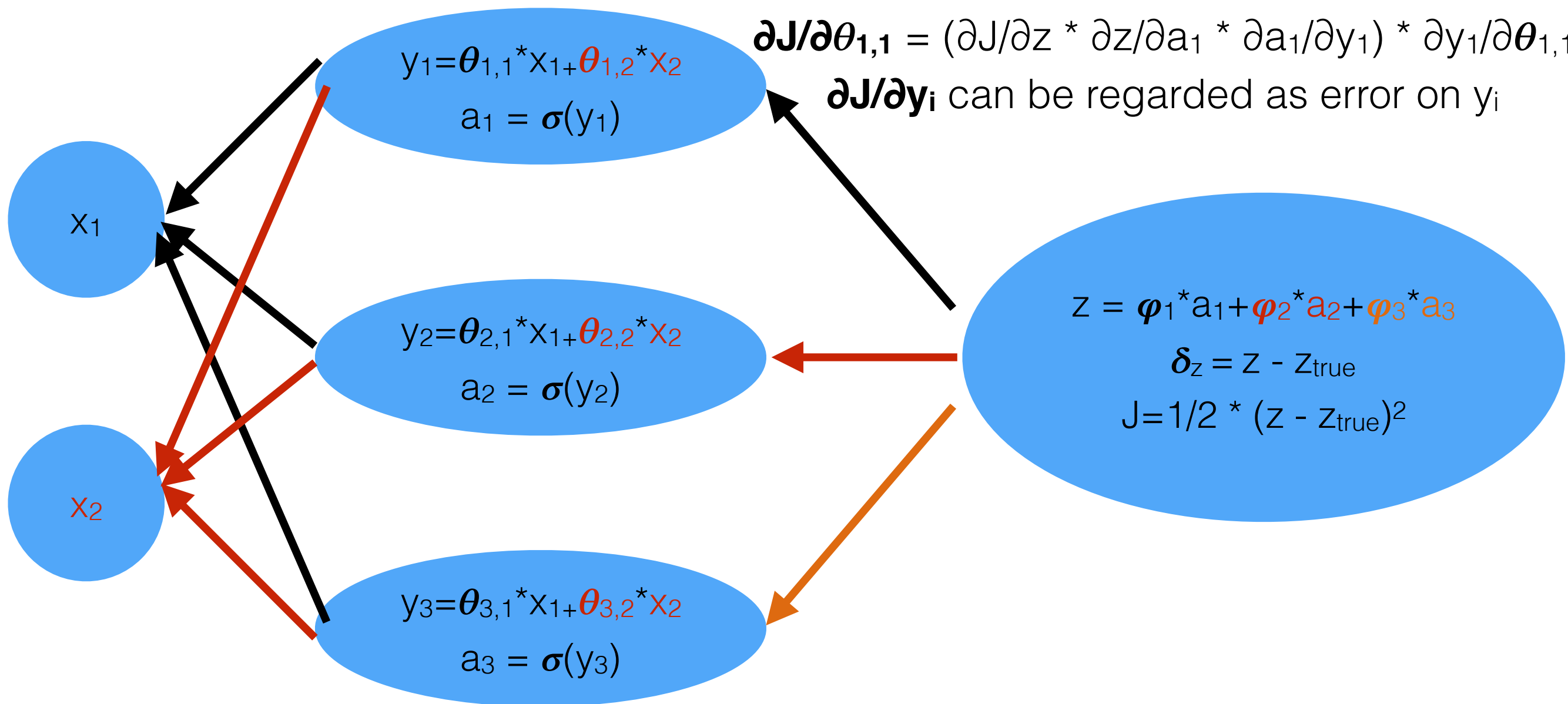
$$\varphi = \begin{array}{|c|c|c|} \hline \varphi_1 & \varphi_2 & \varphi_3 \\ \hline \end{array}$$

# Back-propagation Review

$$\theta_{1,1} = \theta_{1,1} - \lambda * \partial J/\partial\theta_{1,1}$$

$$\partial J/\partial\theta_{1,1} = (\partial J/\partial z * \partial z/\partial a_1 * \partial a_1/\partial y_1) * \partial y_1/\partial\theta_{1,1}$$
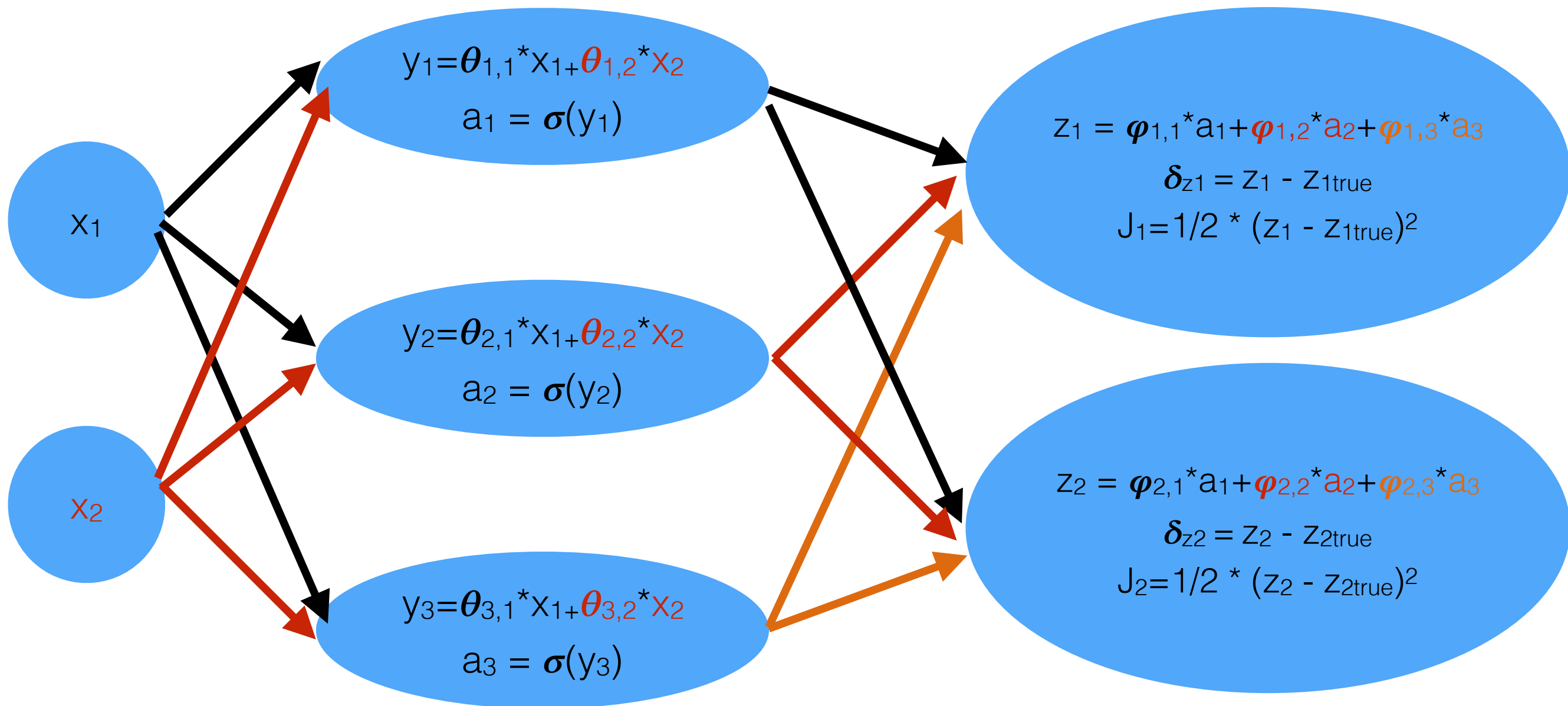
$\partial J/\partial y_i$ can be regarded as error on $y_i$

$x_1$

$x_2$

$y_1 = \theta_{1,1} * x_1 + \theta_{1,2} * x_2$

$a_1 = \sigma(y_1)$

$y_2 = \theta_{2,1} * x_1 + \theta_{2,2} * x_2$

$a_2 = \sigma(y_2)$

$y_3 = \theta_{3,1} * x_1 + \theta_{3,2} * x_2$

$a_3 = \sigma(y_3)$

$z = \varphi_1 * a_1 + \varphi_2 * a_2 + \varphi_3 * a_3$

$\delta_z = z - z_{true}$

$J = 1/2 * (z - z_{true})^2$

$$\theta = \begin{array}{|c|c|} \hline \theta_{1,1} & \theta_{1,2} \\ \hline \theta_{2,1} & \theta_{2,2} \\ \hline \theta_{3,1} & \theta_{3,2} \\ \hline \end{array}$$

$$\varphi = \begin{array}{|c|c|c|} \hline \varphi_1 & \varphi_2 & \varphi_3 \\ \hline \end{array}$$

5

# Back-propagation Review

# Back-propagation Review



$y_1 = \boldsymbol{\theta}_{1,1} * x_1 + \boldsymbol{\theta}_{1,2} * x_2$

$a_1 = \boldsymbol{\sigma}(y_1)$

$z_1 = \boldsymbol{\varphi}_{1,1} * a_1 + \boldsymbol{\varphi}_{1,2} * a_2 + \boldsymbol{\varphi}_{1,3} * a_3$

$\boldsymbol{\theta}_{1,1} = \boldsymbol{\theta}_{1,1} - \lambda * \partial J / \partial \boldsymbol{\theta}_{1,1}$

$\partial J / \partial \boldsymbol{\theta}_{1,1} = (\partial J_1 / \partial y_1) * \partial y_1 / \partial \boldsymbol{\theta}_{1,1} + (\partial J_2 / \partial y_1) * \partial y_1 / \partial \boldsymbol{\theta}_{1,1}$

$x_1$

$y_2 = \boldsymbol{\theta}_{2,1} * x_1 + \boldsymbol{\theta}_{2,2} * x_2$

$a_2 = \boldsymbol{\sigma}(y_2)$

$x_2$

$z_2 = \boldsymbol{\varphi}_{2,1} * a_1 + \boldsymbol{\varphi}_{2,2} * a_2 + \boldsymbol{\varphi}_{2,3} * a_3$

$\boldsymbol{\delta}_{z2} = z_2 - z_{2true}$

$J_2 = 1/2 * (z_2 - z_{2true})^2$

$y_3 = \boldsymbol{\theta}_{3,1} * x_1 + \boldsymbol{\theta}_{3,2} * x_2$

$a_3 = \boldsymbol{\sigma}(y_3)$

$\theta =$

| | |
|---|---|
| $\boldsymbol{\theta}_{1,1}$ | $\boldsymbol{\theta}_{1,2}$ |
| $\boldsymbol{\theta}_{2,1}$ | $\boldsymbol{\theta}_{2,2}$ |
| $\boldsymbol{\theta}_{3,1}$ | $\boldsymbol{\theta}_{3,2}$ |

$\varphi =$

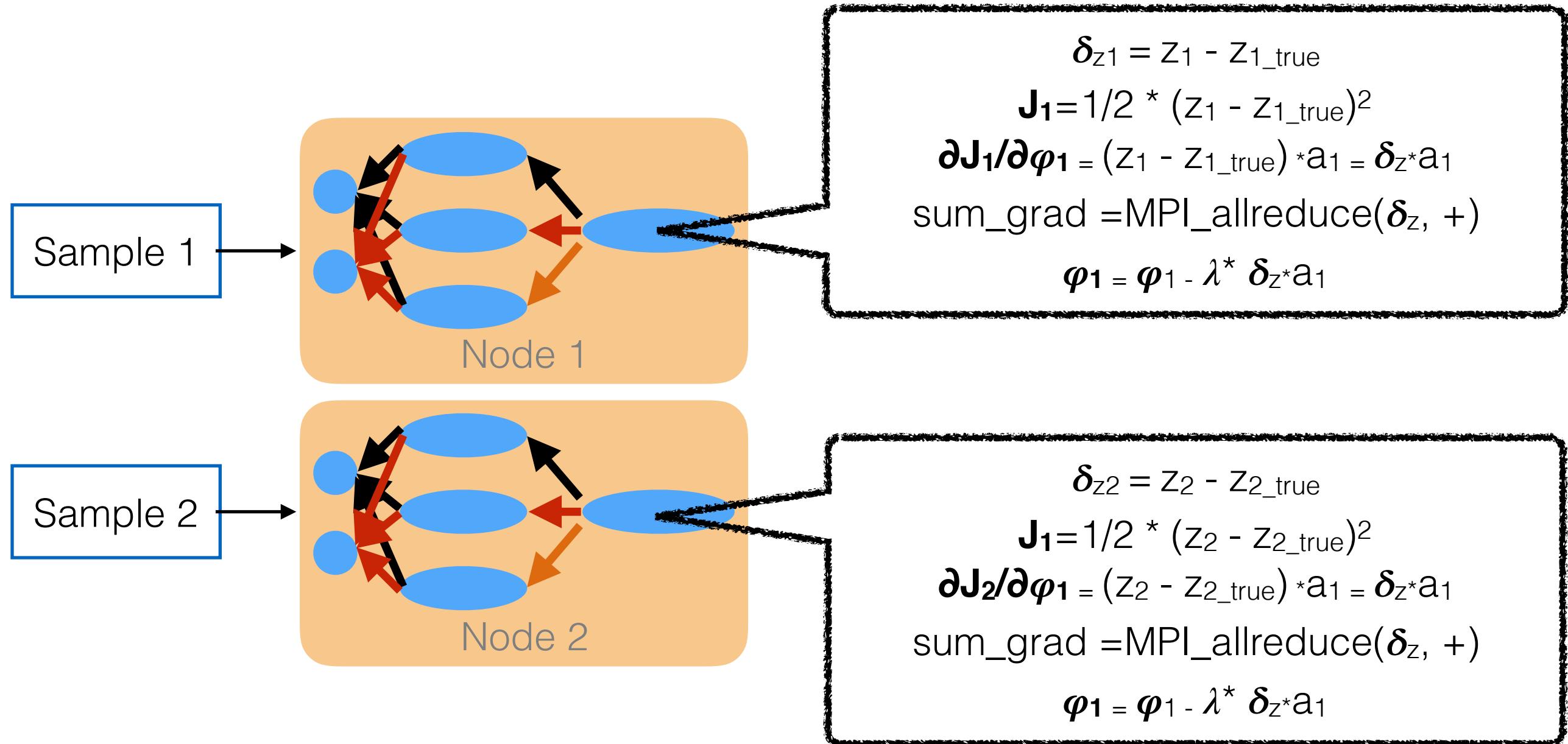| | | |
|---|---|---|
| $\boldsymbol{\varphi}_{1,1}$ | $\boldsymbol{\varphi}_{1,2}$ | $\boldsymbol{\varphi}_{1,3}$ |
| $\boldsymbol{\varphi}_{2,1}$ | $\boldsymbol{\varphi}_{2,2}$ | $\boldsymbol{\varphi}_{2,3}$ |

# Intel MLSL Implementation

- MPI based library MLSL
- Supports Caffe, Theano
- Data parallelism based
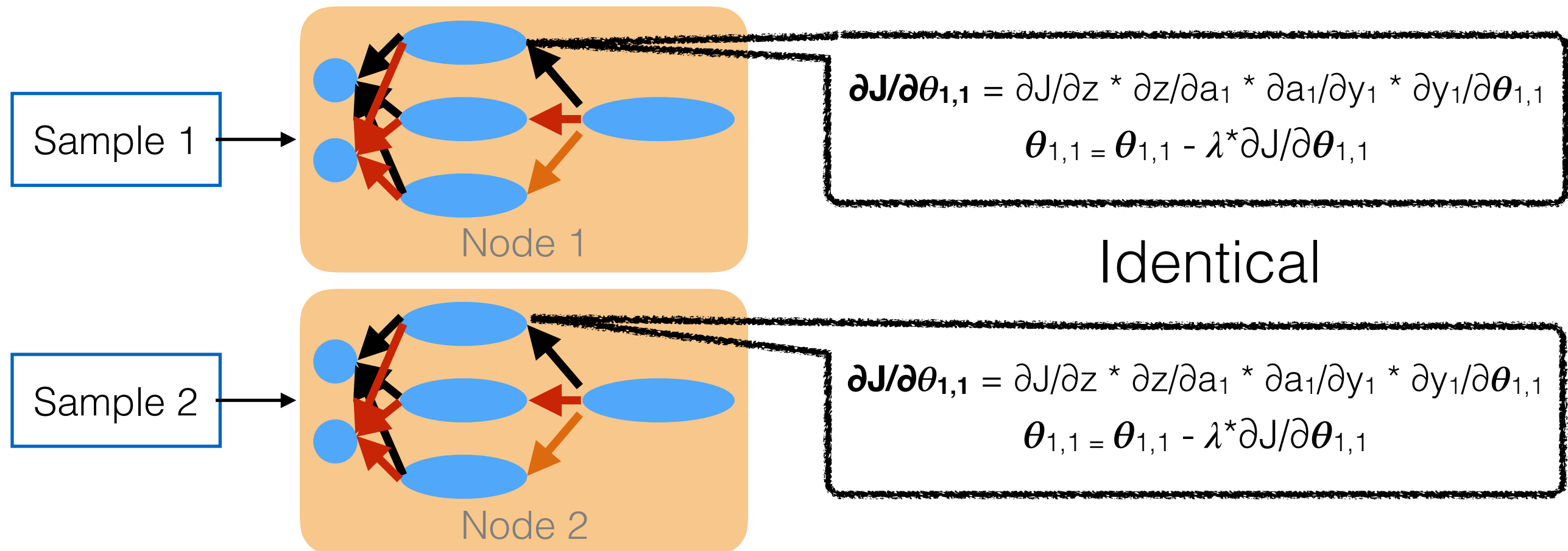- It is as simple as an MPI_allreduce between forward propagation and back propagation

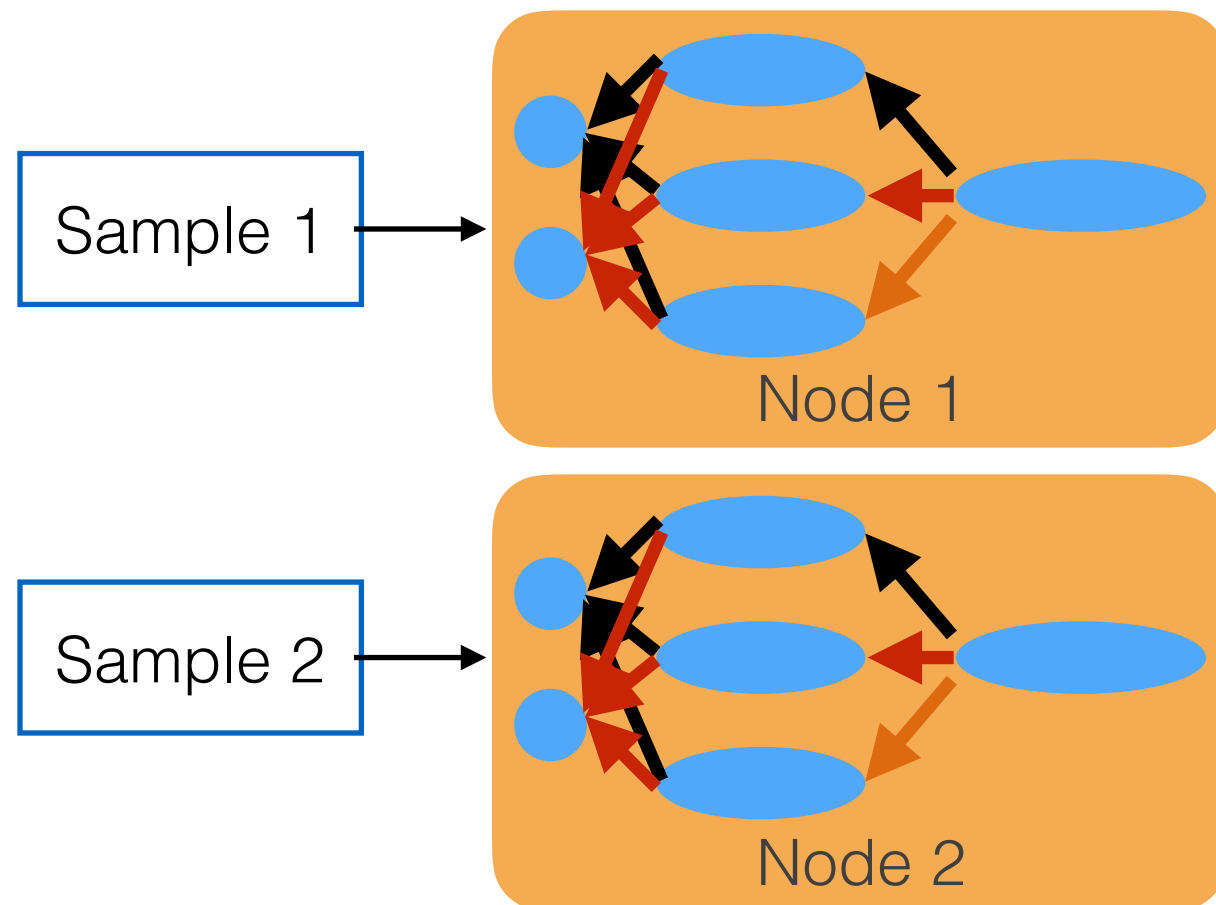# Intel MLSL Implementation



$$\boldsymbol{\delta}_z = z_1 - z_{1\_true}$$
$$\mathbf{J_1} = 1/2 * (z_1 - z_{1\_true})^2$$

$$\boldsymbol{\delta}_z = z_2 - z_{2\_true}$$
$$\mathbf{J_1} = 1/2 * (z_2 - z_{2\_true})^2$$

Sample 1

Sample 2

Node 1

Node 2

# Intel MLSL Implementation



$\boldsymbol{\delta}_{z1} = z_1 - z_{1\_true}$

$\mathbf{J_1} = 1/2 * (z_1 - z_{1\_true})^2$

$\boldsymbol{\partial J_1/\partial \varphi_1} = (z_1 - z_{1\_true}) * a_1 = \boldsymbol{\delta}_{z} * a_1$

sum_grad = MPI_allreduce($\boldsymbol{\delta}_z$, +)

$\boldsymbol{\varphi_1} = \boldsymbol{\varphi}_1 - \lambda * \boldsymbol{\delta}_{z} * a_1$

$\boldsymbol{\delta}_{z2} = z_2 - z_{2\_true}$

$\mathbf{J_1} = 1/2 * (z_2 - z_{2\_true})^2$

$\boldsymbol{\partial J_2/\partial \varphi_1} = (z_2 - z_{2\_true}) * a_1 = \boldsymbol{\delta}_{z} * a_1$

sum_grad = MPI_allreduce($\boldsymbol{\delta}_z$, +)

$\boldsymbol{\varphi_1} = \boldsymbol{\varphi}_1 - \lambda * \boldsymbol{\delta}_{z} * a_1$

# Intel MLSL Implementation



Sample 1

$$\partial J/\partial \theta_{1,1} = \partial J/\partial z * \partial z/\partial a_1 * \partial a_1/\partial y_1 * \partial y_1/\partial \theta_{1,1}$$

$$\theta_{1,1} = \theta_{1,1} - \lambda * \partial J/\partial \theta_{1,1}$$

Node 1

Identical

Sample 2

$$\partial J/\partial \theta_{1,1} = \partial J/\partial z * \partial z/\partial a_1 * \partial a_1/\partial y_1 * \partial y_1/\partial \theta_{1,1}$$

$$\theta_{1,1} = \theta_{1,1} - \lambda * \partial J/\partial \theta_{1,1}$$

Node 2

# Parameter Server

- Uses a separate set of nodes as key-value stores
- Data parallelism based
- Training workers synchronize via push and pull operations
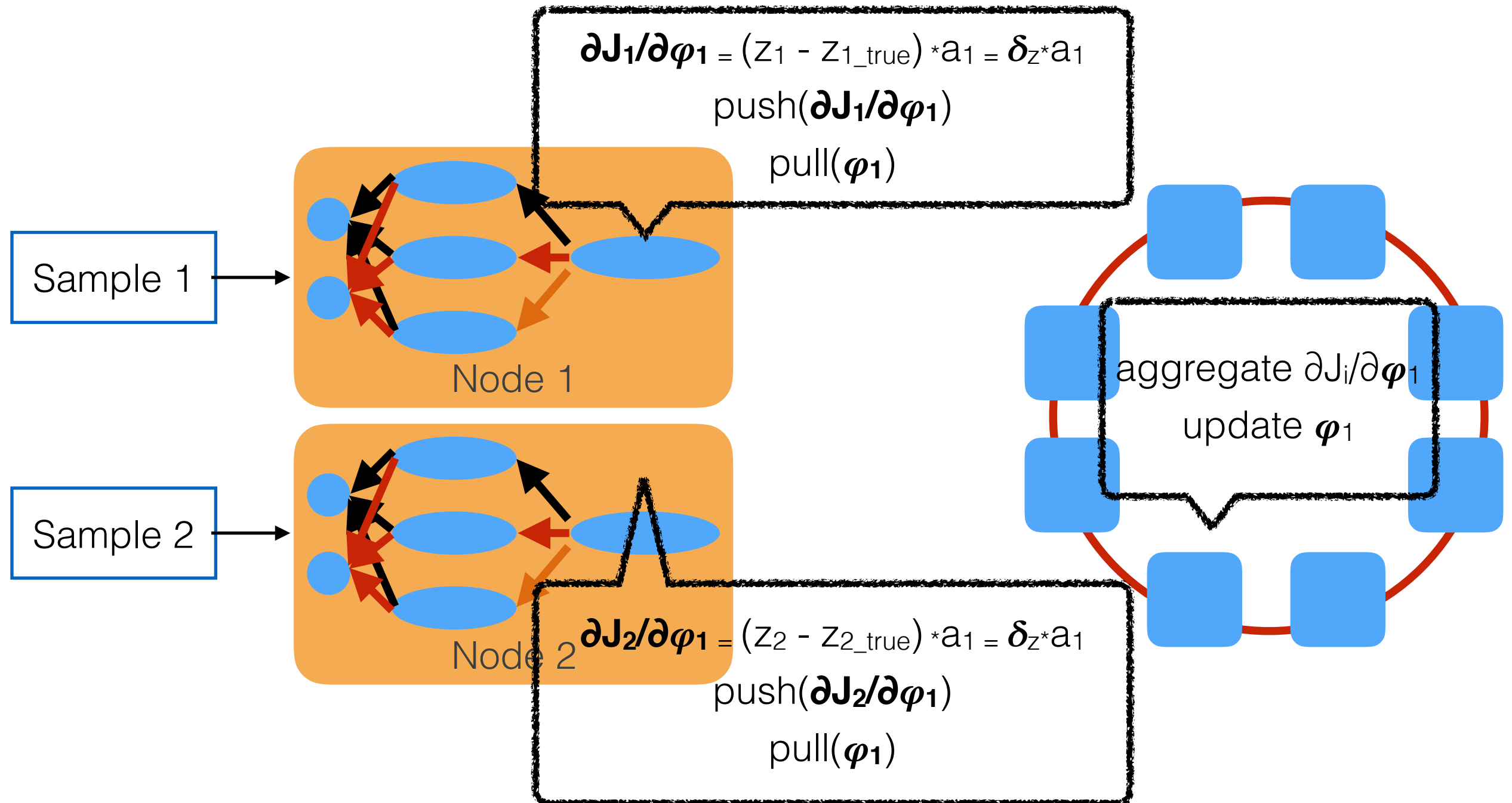- Key innovation: consistency model

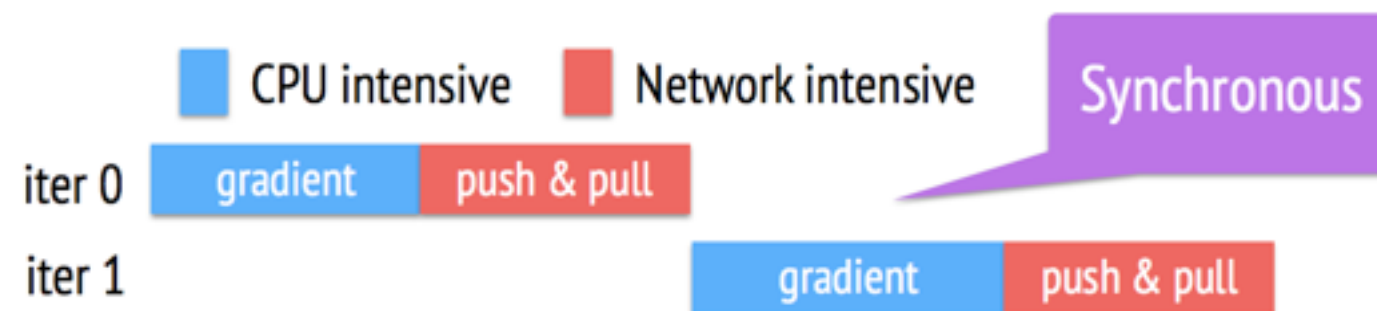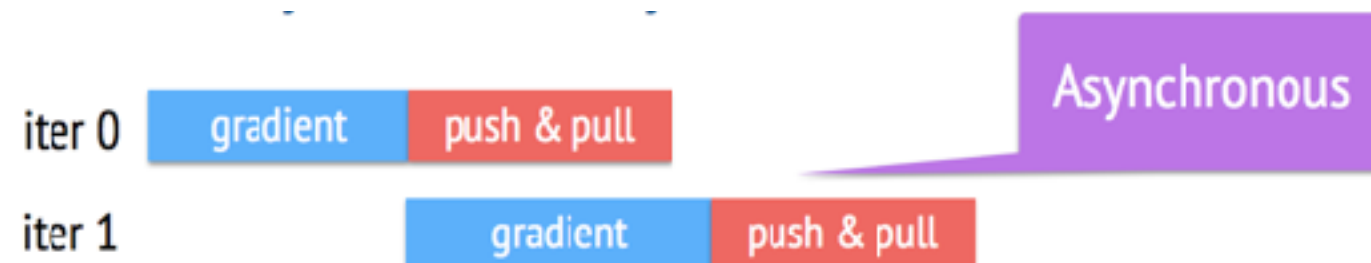# Parameter Server



Key-Value pairs:
(parameter_name, value)

Sample 1

Node 1

Sample 2

Node 2

Parameter Server

# Parameter Server



$\partial J_1/\partial \varphi_1 = (z_1 - z_{1\_true}) * a_1 = \delta_z * a_1$

push($\partial J_1/\partial \varphi_1$)

pull($\varphi_1$)

aggregate $\partial J_i/\partial \varphi_1$

update $\varphi_1$

$\partial J_2/\partial \varphi_1 = (z_2 - z_{2\_true}) * a_1 = \delta_z * a_1$

push($\partial J_2/\partial \varphi_1$)

pull($\varphi_1$)

Sample 1

Sample 2

Node 1

Node 2

# Task

- a push/pull/user defined function (an iteration)
- "execute-after-finished" dependency



- executed asynchronously

# Flexible Iteration Dependency

- Executed asynchronously



- Iteration 1 uses the old parameters as in Iteration 0, and obtains the same gradient.

- It is likely to slow down the convergence progress

- Some algorithms are less sensitive to this type of inconsistency

# Flexible Iteration Dependency

- Sequential

Sequential 

- Eventual

Eventual 

- Bounded Delay

1-bounded delay 

  - $\tau$: maximal delay time. A new iteration will be blocked until all previous tasks $\tau$ time ago have been finished.
  - $\tau = 0$ —> Sequential, $\tau = \infty$ —> Eventual

# Flexible Iteration Dependency

- Bounded Delay consistency model is referred as Stale Synchronous Parallel (SSP) consistency model.

# Practical Scalability

- MPI based solution
  - Intel MLSL — measuring now

- Parameter Server based solution
  - Poseidon — 8 GPU nodes 4-4.5x speedup
  - Common ML algorithm — 5000 workers and 1000 parameter servers