- A Neural Network (NN)

weights
↓
layer.
weights.

How to compute $\nabla F$
" the gradient function"?.

depth. $\longrightarrow$

$L_n \Rightarrow L_2 \Rightarrow ---$ Loss function.

## Backpropagation.

A rule from calculus, the chain rule: $z$ depend on $y$ and $y$ depends on $x$.

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$\underbrace{\quad}_{(1)} \quad \underbrace{\quad}_{(2)}$

$$z = (y)^2 \text{ (1)}, \quad y = 2x \quad \Rightarrow z = (2x)^2$$
$$= y x^2$$

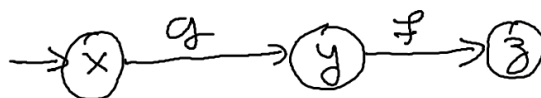$$\frac{dz}{dy} = 2y, \quad \frac{dy}{dx} = 2 \quad \Rightarrow \frac{dz}{dx} = 2y \times 2 = 4y.$$
$$= 8x$$

$$z' = 4 \times (2x) = 8x.$$

How does it translate to multivariate function or having vector as parameters.

$y = (y_1, ---, y_m).$

$(x_1, x_2, ---, x_m) = \underline{x} \in \mathbb{R}^m, \quad y \in \mathbb{R}^n, \quad z \in \mathbb{R}:$ variables

$\underline{g}: \mathbb{R}^m \to \mathbb{R}^n, \quad \underline{f}: \mathbb{R}^n \to \underline{R} \quad :$ functions.

$\Rightarrow z = f(y), \quad y = g(x)$

Loss.

$\longrightarrow \overset{g}{\underset{x}{\bigcirc}} \longrightarrow \overset{f}{\underset{y}{\bigcirc}} \longrightarrow \overset{}{\underset{z}{\bigcirc}}$

$$\frac{\partial \mathfrak{z}}{\partial x_i} = \sum_{\mathfrak{j}} \frac{\partial \mathfrak{z}}{\partial y_{\mathfrak{j}}} \cdot \frac{\partial y_{\mathfrak{j}}}{\partial x_i} \qquad \text{(chain rule)}.$$

$\underbrace{\qquad\qquad\qquad\qquad}_{N \text{ components}}$   (generalized to vectors).

Can be written compactly as:    $\swarrow$ transpose

$$\nabla_x (\mathfrak{z}) = \underbrace{\left( \frac{\partial y}{\partial x} \right)^T}_{\substack{\text{Correspond} \\ \text{to (2)}}} \underbrace{\nabla_y (\mathfrak{z})}_{\text{correspond to (1)}}$$

$\frac{\partial y}{\partial x}$ is $n \times m$ matrix : Jacobian of $\mathfrak{z}$ or $y$.

$$J = \frac{\partial y}{\partial x} = \quad n \left[ \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_m} \\ & \ddots & \\ \frac{\partial y_n}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_m} \end{array} \right]$$

$\underbrace{\qquad\qquad}_{m}$ — matrix of all possible partial derivative I can take.

$$\left( \frac{\partial y}{\partial x} \right)^T = \left( \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & - & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_m} & - \cdots - & \frac{\partial y_n}{\partial x_m} \end{array} \right) \Bigg] m \qquad m \times n \text{ matrix.}$$

$\underbrace{\qquad\qquad}_{n}$

$$\left( \frac{\partial y}{\partial x} \right)^T \nabla_y (\mathfrak{z}) = \left( \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & - & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_m} & - \cdots - & \frac{\partial y_n}{\partial x_m} \end{array} \right) \left( \begin{array}{c} \frac{\partial \mathfrak{z}}{\partial y_1} \\ \vdots \\ \frac{\partial \mathfrak{z}}{\partial y_n} \end{array} \right)$$

$\underbrace{\qquad\qquad}_{n}$

$(m \times n) \times (n \times 1) = m \times 1$

$$= \left( \sum_{\mathfrak{j}} \frac{\partial \mathfrak{z}}{\partial y_{\mathfrak{j}}} \frac{\partial y_{\mathfrak{j}}}{\partial x_1} - - - - - - \sum_{\mathfrak{j}} \frac{\partial \mathfrak{z}}{\partial y_{\mathfrak{j}}} \frac{\partial y_{\mathfrak{j}}}{\partial x_m} \right)^T$$

$$= \nabla_x \mathfrak{z} \qquad (m \times 1) \text{ matrix aka a vector.}$$

$$\boxed{ \nabla_x (\mathfrak{z}) = \left( \frac{\partial \mathfrak{z}}{\partial y} \right)^T \nabla_y (\mathfrak{z}) } \quad \begin{array}{l} (3) \\ \text{(can be recursively} \\ \text{called )}. \end{array}$$

A derivative can be computed recursively in the reverse order. ( Let's say for each function we implent a jacobian, then we muliply as in (1) to obtain the gradient of the loss functioba) .

---

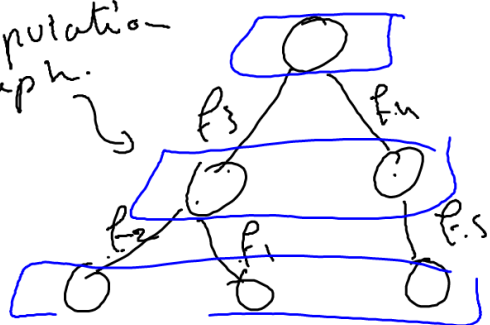Another example ( More simple ). (why we need a forward pass for backprop).

$$x \xrightarrow{g} \boxed{y} \xrightarrow{f} \boxed{3} \xrightarrow{e} \boxed{u}$$

input →

g(x) over first arrow, f(y)·f(g(x)) over second arrow

$$\frac{\partial u}{\partial x} = \frac{\partial u}{\partial 3} \times \frac{\partial 3}{\partial y} \times \frac{\partial y}{\partial x} , \qquad u = e(3) = e'(3)$$

$$\frac{\partial u}{\partial x} = e'(f(g(x))) \times f'(g(x)) \times g'(x)$$

How to obtain : $f(g(x))$ , $g(x)$

I use a forward pass .

Computation complexity :

Computation graph.



$f_3$   $f_n$
$f_2$   $f_1$   $f_s$

forward pass.
$O(n)$
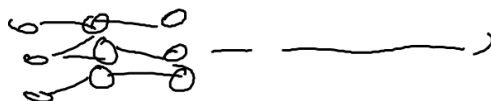n : # of nodes.

• Assuming that.
function evaluations have
cost $n$   $f(\cdot)$
$\approx$   $\approx$   $\frac{\partial f(\cdot)}{\partial x}$ nees

cost ?.

m : # of edges.

Back propagations: $O(n) + O(m) = O(n^2)$

In a fully connected graph. $O(n^2)$

Because NNs have chain structure

practically backprop. ~ $O(n)$ complexity.

Two types of derivates used in practice.

Symbol-to-member: Torch, caffe.

Symbol-to-symbol: Theano, Tensorflow.

↳ adds extra nodes to the computational graph, and the derivatives are computed by traversing the graph, where as symbol to-muber are computed within nodes.

## Techniques used in ① NW training.

### Surrogate Loss functions and Early stopping.

what is a surrogate Loss function is a gradient friendly function ( friendly = I can take the derivate ), for example,

$(0-1)$ loss. $\mathbb{1}\{y_i \neq \hat{y}_i\}$ insted to work with this directly, I can work (NLL) negative.Log-Likeli hood.

. As the optimization is executed. I use a criterion.

on the validation set based on the true loss (0-1 loss).
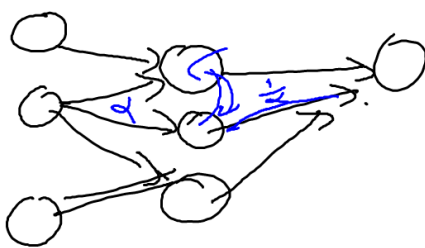(to avoid overfitting).

I stop when it's satisfied. ⟹
. example ( loss is not changing or low enough).
the gradient may still be big.

$\neq$ Classical optimization
$\nabla F \approx 0$.

# NN Landscapes & local minimas

Landscape

I can swap nodes and I obtain an equivalent model.

(Model Identifiability).

- Weight space symmetry.

- scaling of the · weights } equivalent models.
  $\frac{1}{\alpha}, \alpha$.

  can
  $\Rightarrow$ a NN have uncountably infinitly many local minimas.

- In practice, one suspects local minima values aren't much different than global minima.

## Other critical points (Maximas not much of a problem)

saddle point — because it's unstable, usually SGD manges.
  - because of randomness.
  - Saddle points Impirically gradient meth
    (problem for
    (second order
    methods) ← they are attracted by saddle points.
  , 2nd order information,
  , makes you attrackted).

- Cliffs causes. $\Rightarrow$ exploding gradients.    big.

  valley

  clipping.

  ← huge gradient. } heuristic: to clip the
  gradients only take the direction information.

  method that is not proven to have optimality prop.

  $\theta$

  m

  exploding → { - numerical instability.

$L$ - driven away from minimas.

## Adaptive Learning rates modifications for SGD

### Adaptive gradient.

ADAGRAD:

$$G_k = \sum_{k'=1}^{k} \underset{\tilde{g}}{g(w_{k'}, \xi_{k'})} \; g(w_{k'}, \xi_{k'})^{\top}$$

$$(G_k)_{jj} = \left( \sum_{k'=1}^{k} g \, g^{\top} \right)_{jj} \qquad \begin{pmatrix} a_1 & \\ & a_q \end{pmatrix}^{-\frac{1}{2}} = \begin{pmatrix} \frac{1}{\sqrt{a_1}} & \\ & \frac{1}{\sqrt{a_q}} \end{pmatrix}$$

Hadamard product.

$$w_{k+1} = w_k - \eta \, \text{diag}(G_k)^{-\frac{1}{2}} \odot g$$

$\uparrow$ fixed

$$(w_{k+1})_j = (w_k)_j - \frac{\eta}{\sqrt{G_{jj}}} \left( g(w_k, \xi_k) \right)_j$$

$\sqrt{G_{jj}} \simeq$ estimate of the Lipchitz constat

inspired by online learning techniques.

(ADAM) ( Adaptive moment estimation).

Walking average,

$$\begin{cases} m_{k+1} = \beta_1 \, \underline{m_k} + (1 - \beta_1) \, g(w_k, \xi_k) \leftarrow \\ v_{k+1} = \beta_2 \, \underline{v_k} + (1 - \beta_2) \left( g(w_k, \xi_k) \right)^2 \end{cases}$$

past   new value.

$\beta_1, \beta_2 =$ forgetting factors.

$\rightsquigarrow \uparrow$ forget faster.

$$\tilde{m} = \frac{m_{k+1}}{1 - (\beta_1)^{k+1}} \qquad \tilde{v} = \frac{v_{k+1}}{1 - (\beta_2)^{k+1}}$$

as $k$ becomes larger I give more importance

to recent values of $v$ and $m$.

$$w_{t+1} = w_t \ominus \eta \times \frac{\widetilde{m}}{\sqrt{\widetilde{v}} + \varepsilon}$$

gradient

to avoid division over zero.



← average.
← walking average.

RMS PROP ( Root mean square propagation).

$$v(w_t) = \gamma \, v(w_{t-1}) + (1-\gamma)(\nabla g)^2.$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v(w,t)}} g \cdots$$

# Parameter Initialization

- Break symmetry ⊙ something that we want.

- Randomly initialize the weights.

- If I make the weight large ⇒

Break better the symmetry.

- ✓ ✓ ✓ large ⇒ exploding gradients.

give $m \times n$ stage

- 2 general rule of thumbs ⊙ $w_{ij} \sim U\left(-\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}}\right)$

- $w_{ij} \sim U\left(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}}\right)$.

⇒ compromize between constant activation variance and constant gradient variance.

# Batch normalization

Reparametrize the n.n. Let $H$ be a minibatch
of activations at a hidden layer (we t the
a subset).

$$\underline{H'} = \frac{H - \mu}{\delta} \quad \Big] \quad \mu = \frac{1}{M} \sum_{i=1}^{M} H_i$$

$$\delta = \sqrt{\underbrace{\delta + \frac{1}{M} \sum_i (H - \mu)_i^2}_{> 0}}$$
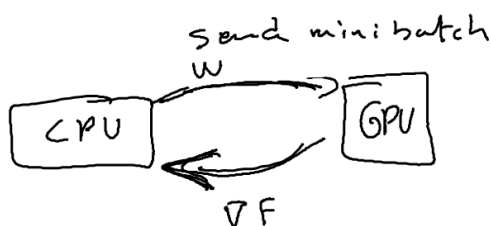
to have $\delta > 0$.
to avoid division over $0$.

---

## <span style="color:red">Distributed optimization</span>

Usually it's benificial to split the work load if
I have more hardware.

- I delegate tasks to Processing Units
  (PUs) : CPUs, GPUs. TPUs.
  
  $\underbrace{\text{Tensor PUs}}$.

---

Assume I am delegating tasks by sending minibatches
then I receive gradient. nmb : size of minibatche

send mini batch



$t_c$ : compute time in GPU per data

compute time $= nmb \times t_c$.

$$d_1 + \underbrace{nmb \, t_T}_{\text{time to transmit}} + \underbrace{nmb \, t_c}_{\text{time to compute}} + d_2 \leftarrow \text{time to prepare info}$$

time
to give data to GPU
$\uparrow$
iteration time.

transmission.

$d_1 = d_2 = d$

Time to converge to $\varepsilon$ precision.

$$T\varepsilon \sim A \left( \frac{1}{\varepsilon} \times \frac{1}{nmb} \right) \times \left( 2d + nmb \times (tc + t_T) \right)$$
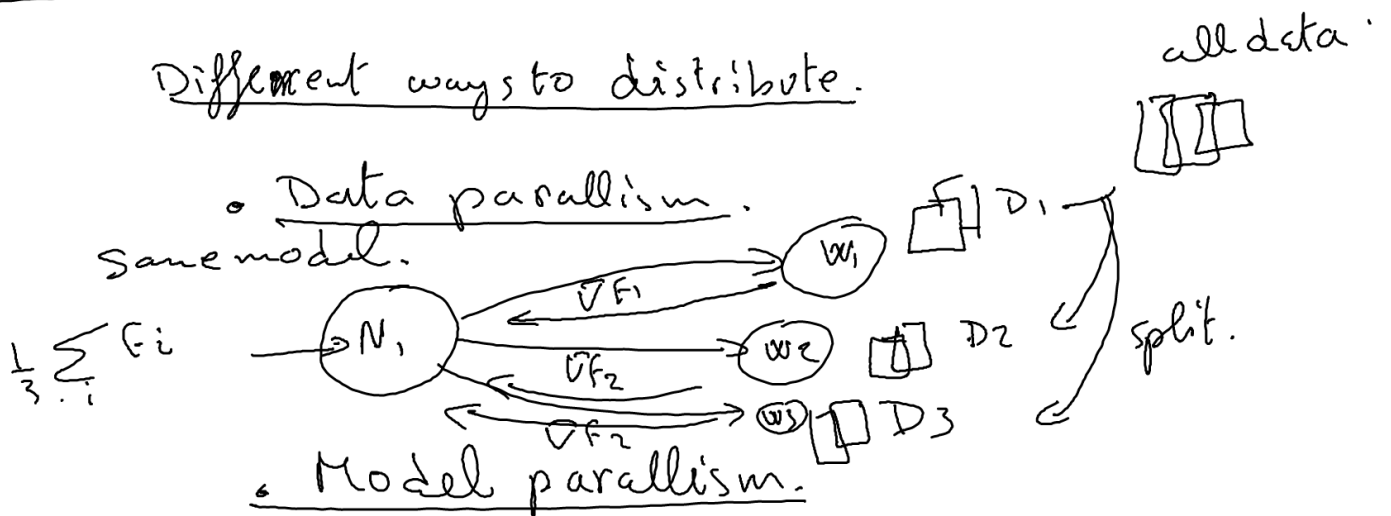
Real time.

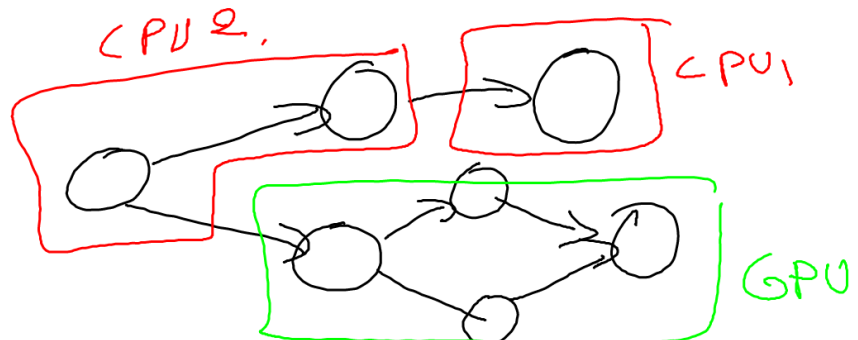if $d = 0$.   $nmb \leftarrow$ doesn't matter.

if $d > 0$   $nmb \gg 1$.

when $nmb \gg 1$

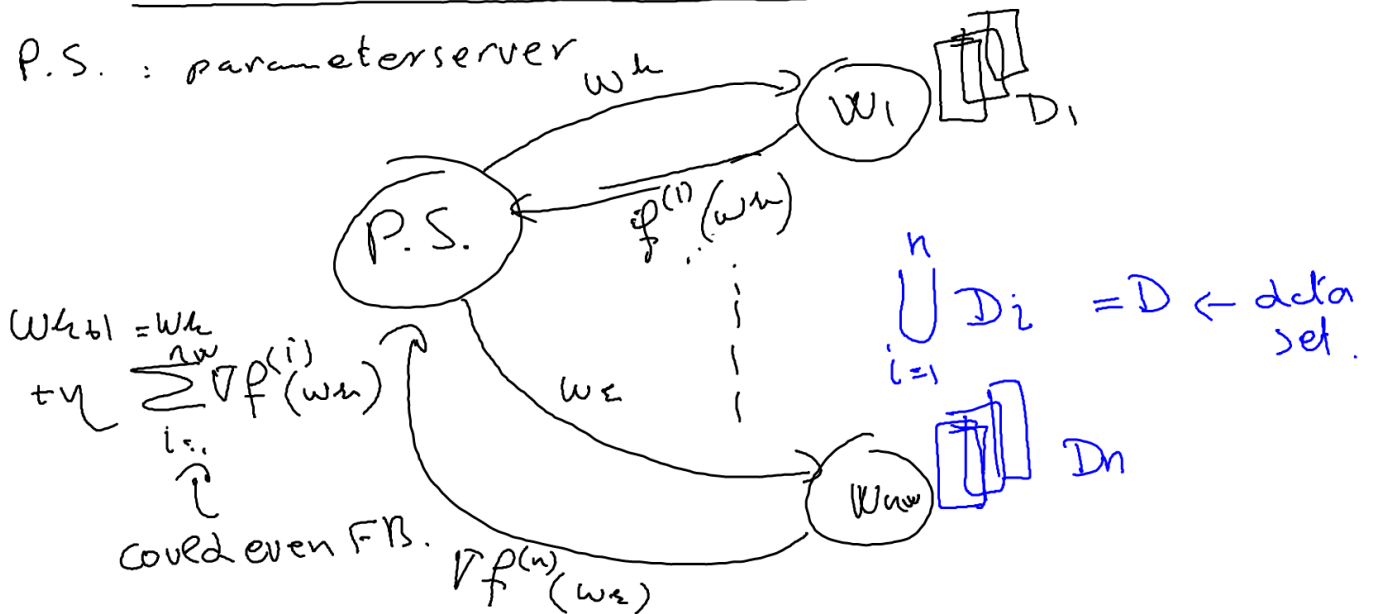$$\boxed{T\varepsilon = \frac{1}{\varepsilon}(tc + t_T)}$$

---

Different ways to distribute.

all data.


• Data parallism.
Same model.



$\frac{1}{3}\sum_i f_i$

$\nabla F_1$
$\nabla F_2$
$\nabla f_2$

$D_1$
$D_2$  split.
$D_3$

• Model parallism.

Instead to lean $(\underbrace{w_1, w_2}_{w_1}, \underbrace{- - -}_{w_2}, \underbrace{- w_d}_{w_{10}})$ in one

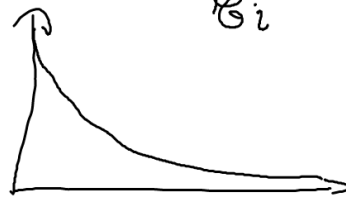machine, I can split the Learning to different

machines.



CPU 2.
CPU 1
GPU

P.S. : parameterserver

$W_k$

$W_1$

$D_1$

P.S.

$f^{(1)}(w_k)$

$n$

$\bigcup_{i=1}^{n} D_i = D \leftarrow$ data set.

$w_k$

$D_n$

$W_{nw}$

$W_{k+1} = W_k$
$+ \eta \sum_{i=1}^{nw} \nabla f^{(i)}(w_k)$

$\uparrow$

could even F.B.

$\nabla f^{(n)}(w_k)$

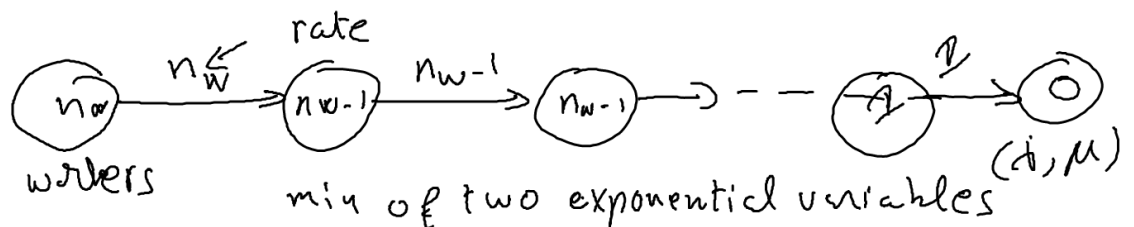Why it is called synchroneous, because I need to wait all the results from each worker to go to the next iteration.

" Bad IDea in practice because of straggler effect ". ( straggler : some one who takes forever to get the job done compared to others).

I am waiting for $n_w$ workers, Each worker gets the job done in random time $\sim Exp(1)$
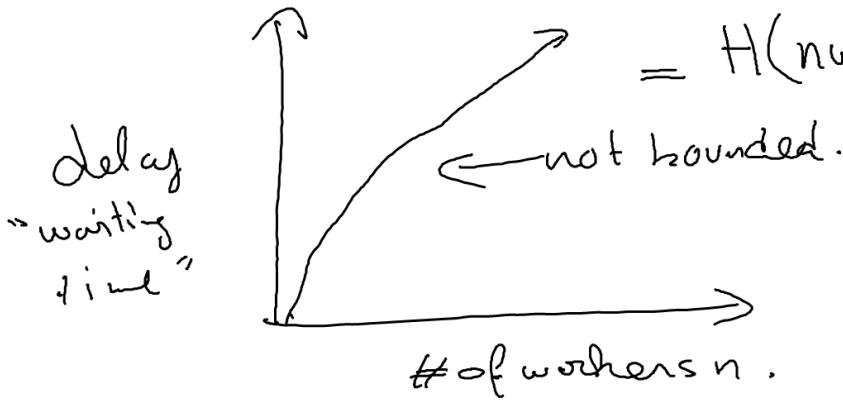
$B_i$

$\mathcal{E}$ : tau.

Expected finish time for a worker is $\mathbb{E}[B_i] = 1$

rate

$n_w$ — $n \overleftarrow{\frac{1}{W}}$ → $n_w-1$ — $n_w-1$ → $n_w-1$ → $-$ $-$ → $1$ $\overset{1}{\to}$ $0$

$(i, \mu)$

workers

min of two exponential variables

is cls exponential with rate $\sigma + \mu$.

$$\mathbb{E}[\max \sigma_i] = \frac{1}{nw} + \frac{1}{nw-1} + \,-\,-\,-\, 1$$

$$= H(nw) \simeq \log(n)$$



← not bounded.

delay
"waiting time"

\# of workers n.

---

In practice it's even worse.



task done 90%

jobs that fail
may refail.

← 5%

3%

≈%

I wait for 100%
to be done.

waiting time.

Synchroneous Parameter server is Bad! ← No one uses t.

---

Ring - All - Reduce Architecture



w

cells = compts of w.

w

to be completed

w

gradients

(1)

(2)

(3)

$e(nw - 1)$ rounds of coms.

split to two types $\begin{cases} \text{aggregation.} (nw-1) \\ \text{exchange} (nw - 1) \\ \text{of} \\ \text{completed} \\ \text{components} \end{cases}$

Ring topology.

dependency is on one worker
rather than $n$.