



# GOSSIP COMMUNICATIONS FOR DEEP LEARNING DISTRIBUTED OPTIMIZATION

18 mars 2017

Michael Blot<sup>1</sup>, David Picard<sup>2</sup>, Matthieu Cord<sup>1</sup>, Nicolas Thome<sup>1</sup>

(1) Université Pierre et Marie Curie, UPMC–Sorbonne Universités, LIP6,  
Paris, France

# Outline

- 1 Context
- 2 Distributed framework
- 3 Gossip Stochastic Gradient Descent : GoSGD
- 4 Experiment

# Outline

- 1 Context
- 2 Distributed framework
- 3 Gossip Stochastic Gradient Descent : GoSGD
- 4 Experiment

# Expected Loss Optimization

## Objective

We want to find  $x^*$  that minimize

$$L(x) = \mathbb{E}_Y[\ell(x, Y)]$$

- $Y \sim \mathcal{I}$  is a random couple (input, label) variable following **natural images** distribution
- $x$  is the network parameters to optimize
- $\ell$  the error function or loss

# Stochastic gradient descent

## Gradient descent

$$x \leftarrow x - \eta \nabla L(x)$$

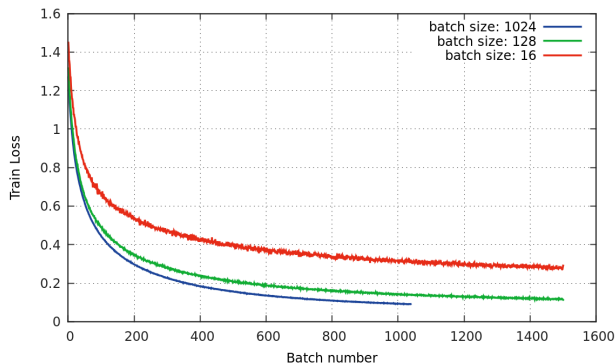
## Stochastic estimation of the gradient

Monte Carlo method :

$$\begin{aligned}\nabla L(x) &= \nabla_x \left( \mathbb{E}_Y[\ell(x, Y)] \right) \\ &= \mathbb{E}_Y[\nabla_x \ell(x, Y)] \\ &\approx \frac{1}{N} \sum_{i=1}^N \nabla_x \ell(x, Y_i)\end{aligned}$$

# Convergence improvement

Comparison of convergence at same number of gradient descent iteration for different batch size : 16, 128, 1024

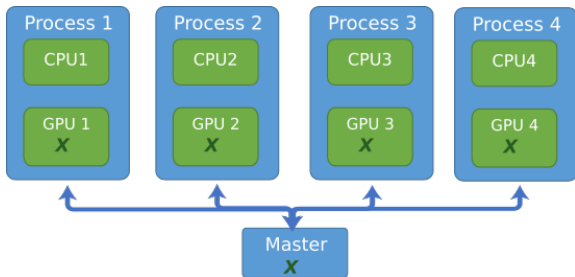


CIFAR10 : Image size  $32 \times 32$  RGB pixels

# Outline

- 1 Context
- 2 Distributed framework**
- 3 Gossip Stochastic Gradient Descent : GoSGD
- 4 Experiment

# Map Reduce Framework



M threads :

Synchronous aggregation of the gradient

$$\widehat{\nabla} L_N(x) = \frac{1}{M} \sum_{m=1}^M \widehat{\nabla} L_{N'}^m(x)$$



# Distributed gradient descent

## Synchronous framework variable

$$x^{t+1} = x^0 - \sum_{\tau=0}^t \eta^\tau \sum_{m=1}^M \frac{1}{M} \widehat{\nabla} L_{N'}^m(x^\tau)$$

## General framework

$$x^{t+1} = x^0 - \sum_{\tau=0}^t \eta^\tau \sum_{m=1}^M a_m^\tau \widehat{\nabla} L_{N'}^m(x_m^\tau)$$

With  $\forall \tau \sum_{m=1}^M a_m^\tau = 1$

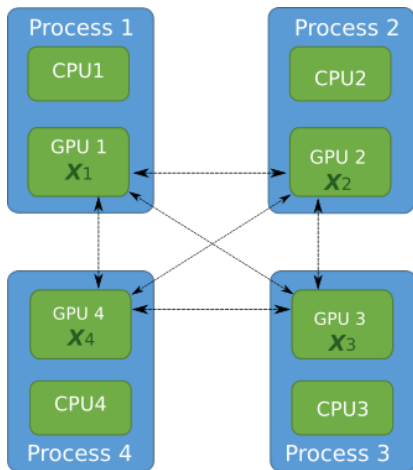
Difficulty : how far are the gradient estimates :

$$\|\widehat{\nabla} L_{N'}^m(x_m^\tau) - \widehat{\nabla} L_{N'}^m(x^\tau)\| < \epsilon$$

# Outline

- 1 Context
- 2 Distributed framework
- 3 Gossip Stochastic Gradient Descent : GoSGD**
- 4 Experiment

# Gossip communication



## GoSGD pseudo-code

**Input** :  $p$  : probability of exchange,  $M$  : number of threads,  $\eta$  : learning rate

**Initialize** :  $x$  is initialized randomly,  $x^i = x$ ,  $\alpha_i = \frac{1}{M}$

**repeat**

$$x^i = x^i - \nu g(x^i, \text{batch})$$

**if**  $S \sim B(p)$  **then**

$$j = \text{Random}(M)$$

$$\alpha^j = \frac{\alpha^j}{2}$$

$$x^i = \frac{\alpha^i}{\alpha^i + \alpha^j} x^i + \frac{\alpha^j}{\alpha^i + \alpha^j} x^j$$

$$\alpha_i = \alpha_i + \alpha_j$$

**end if**

**until** Forever

**get**  $\frac{1}{M} \sum_{m=1}^M x_m$

# Outline

- 1 Context
- 2 Distributed framework
- 3 Gossip Stochastic Gradient Descent : GoSGD
- 4 Experiment**

# Illustration : Gossip vs synchronous vs single thread

- Database : CIFAR10 (50000 images of  $32 \times 32$  RGB pixels)
- Batch size : 20 images, learning rate : of 0.01
- Networks : 286154 parameters, 8 layers
- Number of thread (M) : 8

