# Faculty of Informatics

Stefano Gonçalves Simao
Student
goncast@usi.ch

Prof. Rolf Krause
Advisor
rolf.krause@usi.ch

Dr. Alena Kopaničáková
Co-Advisor
alena.kopanicakova@usi.ch

# Implementation of a hybrid data parallel algorithm for deep neural network training with reduced communication targeted to GPU-based supercomputers
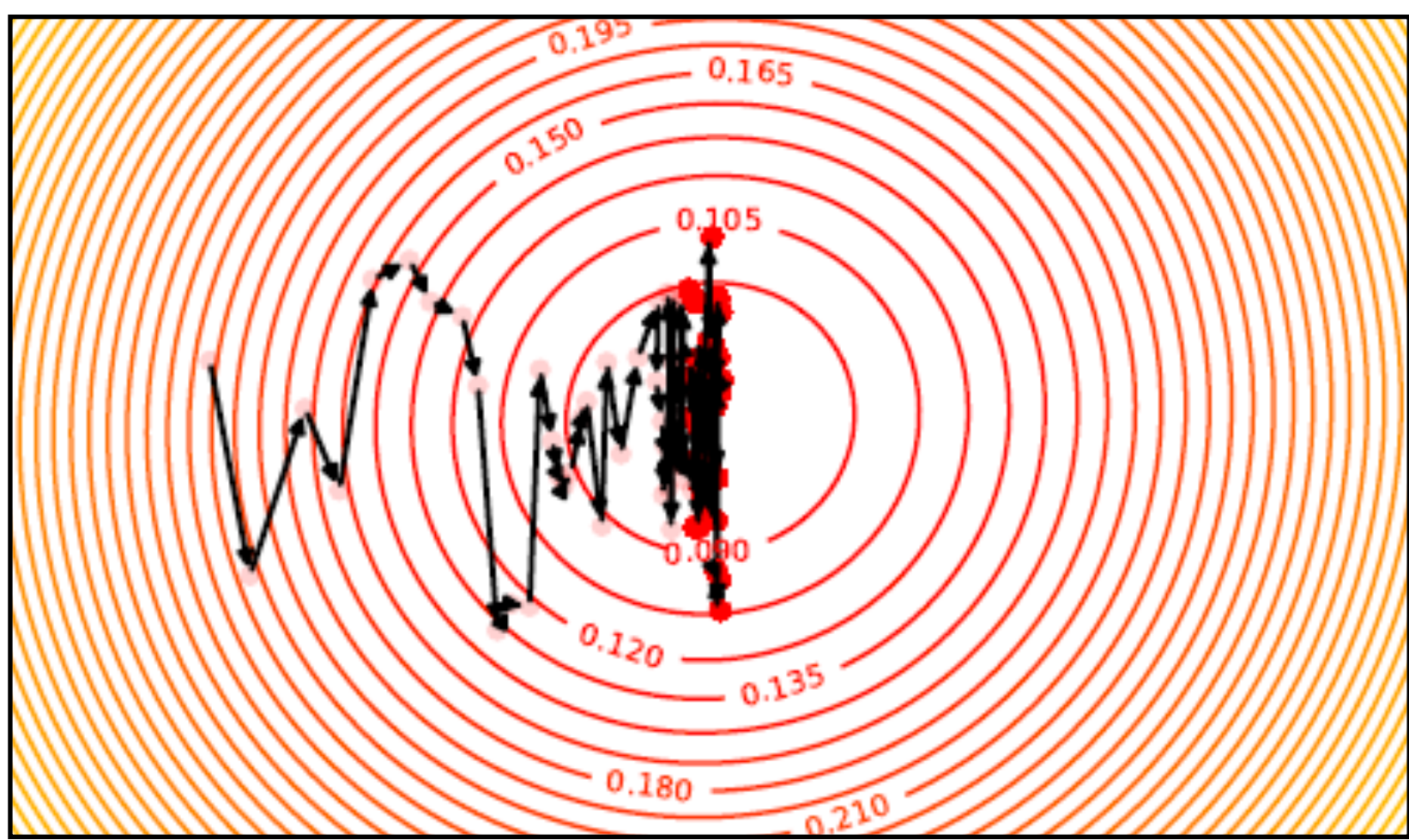
## GOAL

The goal of this project is to investigate existing Parallel Programming strategies to distribute the work of Machine Learning algorithms for training Deep Neural Networks and propose a novel algorithm that reduces communication complexity.

## KEY WORDS

### Supervised Learning Algorithm

### Stochastic Gradient Descent

$$\min_{\theta} L := \sum_{i=1}^{n} l(g(x_i; \theta), y_i),$$
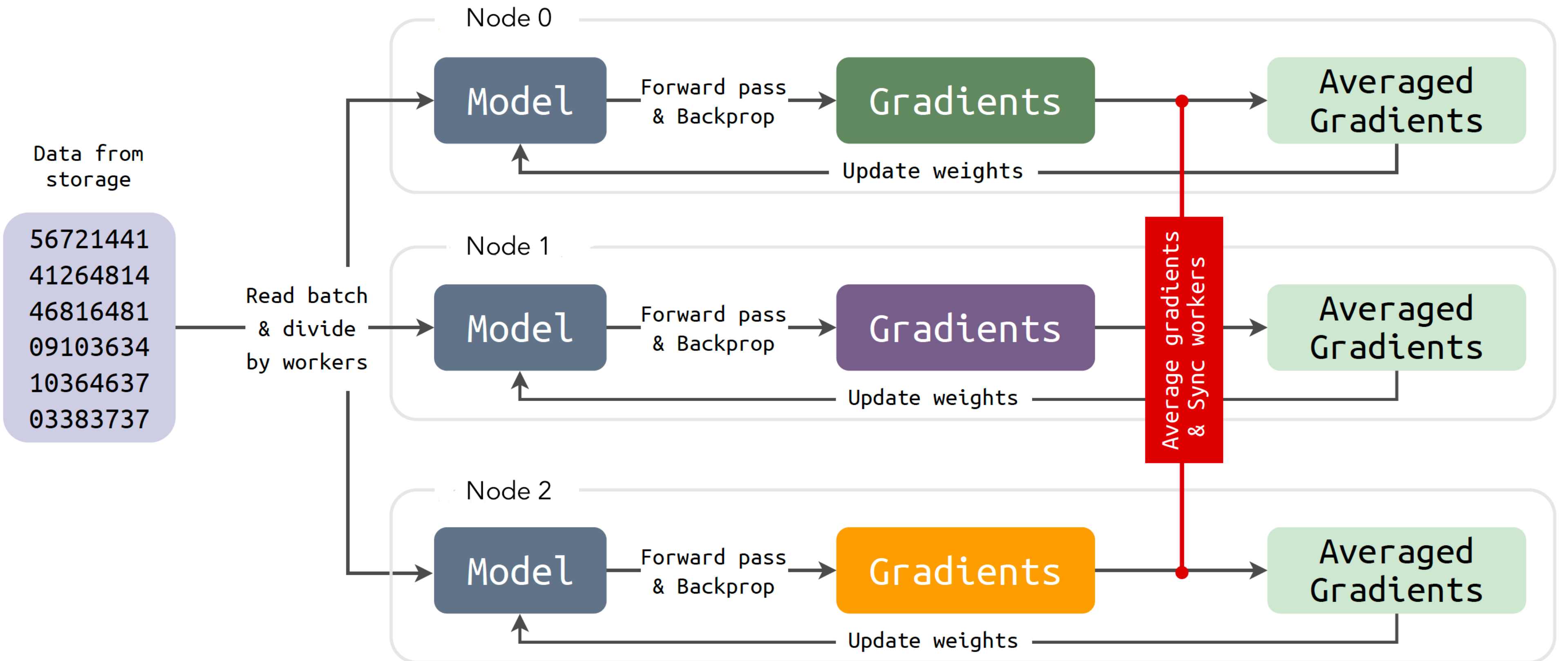


$$\theta_{k+1} = \theta_k - \alpha_k \nabla L_k,$$
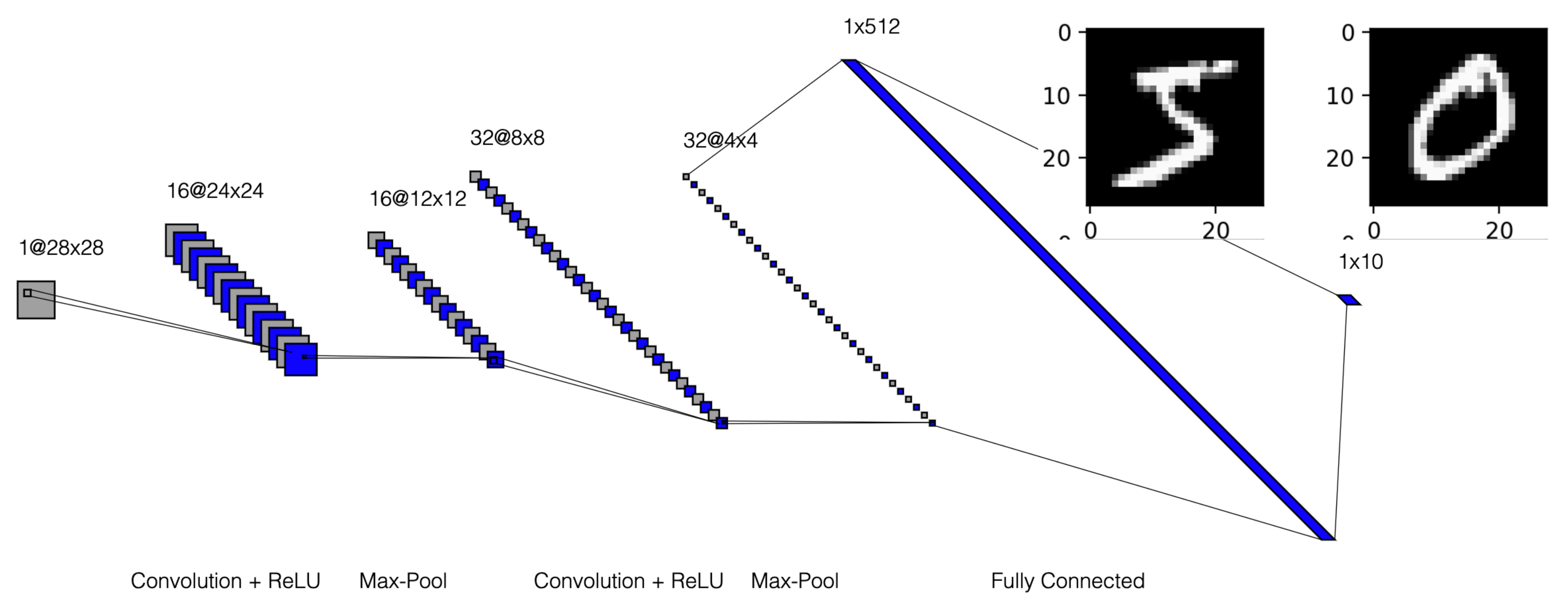
### Data Parallelism

P1
P2
P3



### GPU-based supercomputers
(CSCS - Piz Daint)
5704 x Intel® Xeon® E5-2690 v3 @2.60GHz
(12 cores, 64GB RAM),
NVIDIA® Tesla® P100 16GB

### DistributedDataParallel - NCCL - AllReduce

Node 0
Model — Forward pass & Backprop → Gradients — → Averaged Gradients
Update weights

Data from storage
56721441
41264814
46816481
09103634
10364637
03383737

Read batch & divide by workers

Node 1
Model — Forward pass & Backprop → Gradients — → Averaged Gradients
Update weights

Average gradients & Sync workers

Node 2
Model — Forward pass & Backprop → Gradients — → Averaged Gradients
Update weights



### Convolutional Neural Networks - MNIST



1@28x28    16@24x24    16@12x12    32@8x8    32@4x4    1x512    1x10

Convolution + ReLU    Max-Pool    Convolution + ReLU    Max-Pool    Fully Connected

### Decentralised Synchronous Communication

## Hybrid Parallel Stochastic Gradient Descent

### _Algorithm structure:_

Divide num. epochs by 3

1. Standard Parallel SGD with gradient synchronisation at each iteration
2. Standard Parallel SGD with gradient accumulation
3. Parallel SGD with gradient accumulation without model consistency (independent _influenced_ training)

Model averaging