# Distributed Machine Learning with Zero ETL

Yury Babak

Head of development, GridGain
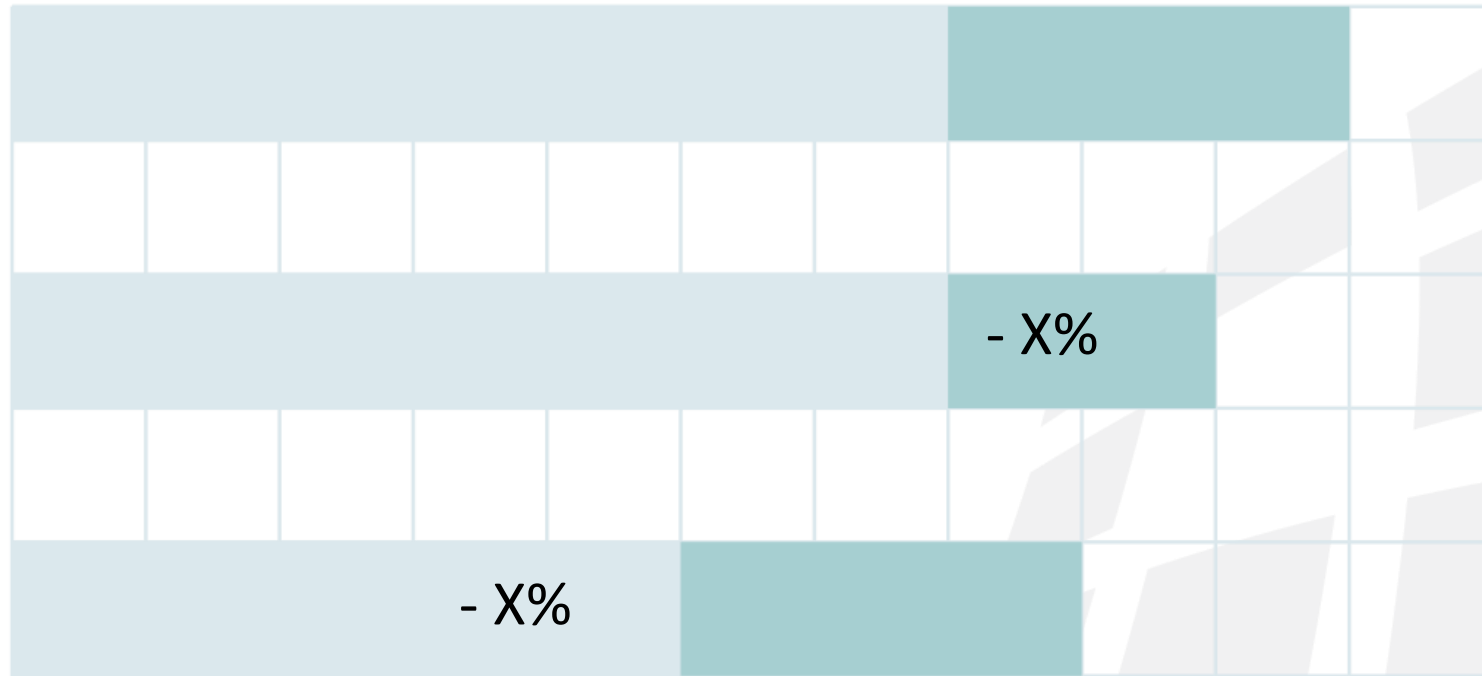
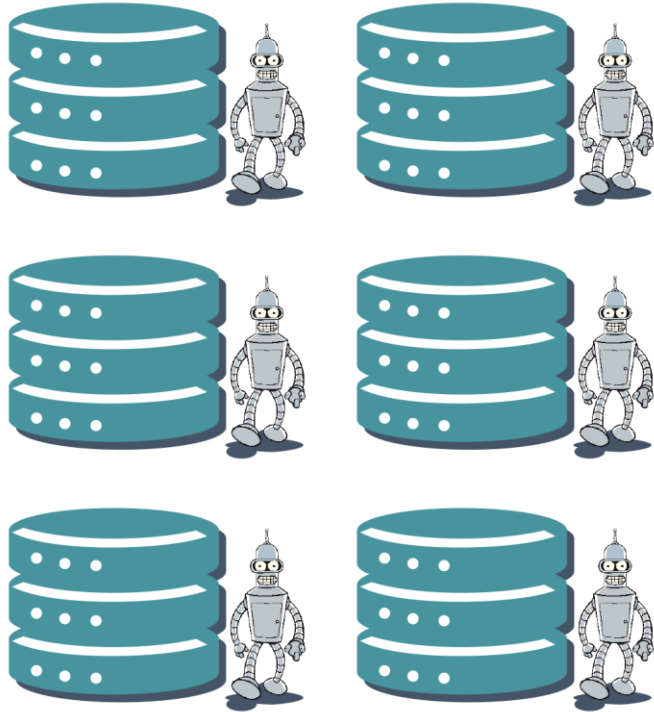# Long ETL

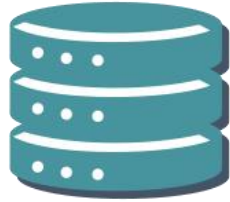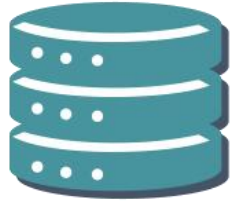# Long ETL



ETL      ML

- X%

- X%

GridGain

# Distributed Training

# Node Crash

# Apache Ignite

GridGain

# Apache Ignite: Replicated Caches



Client

Server Node 1

Server Node 2

Server Node 3

Server Node 4

# Map Reduce

GridGain

# Iterative Optimization Algorithm

GridGain

# Partition Based Data Set

# Restoration of partitions after a failure

# Recovering calculations after failure

# OLS sample

Loss function
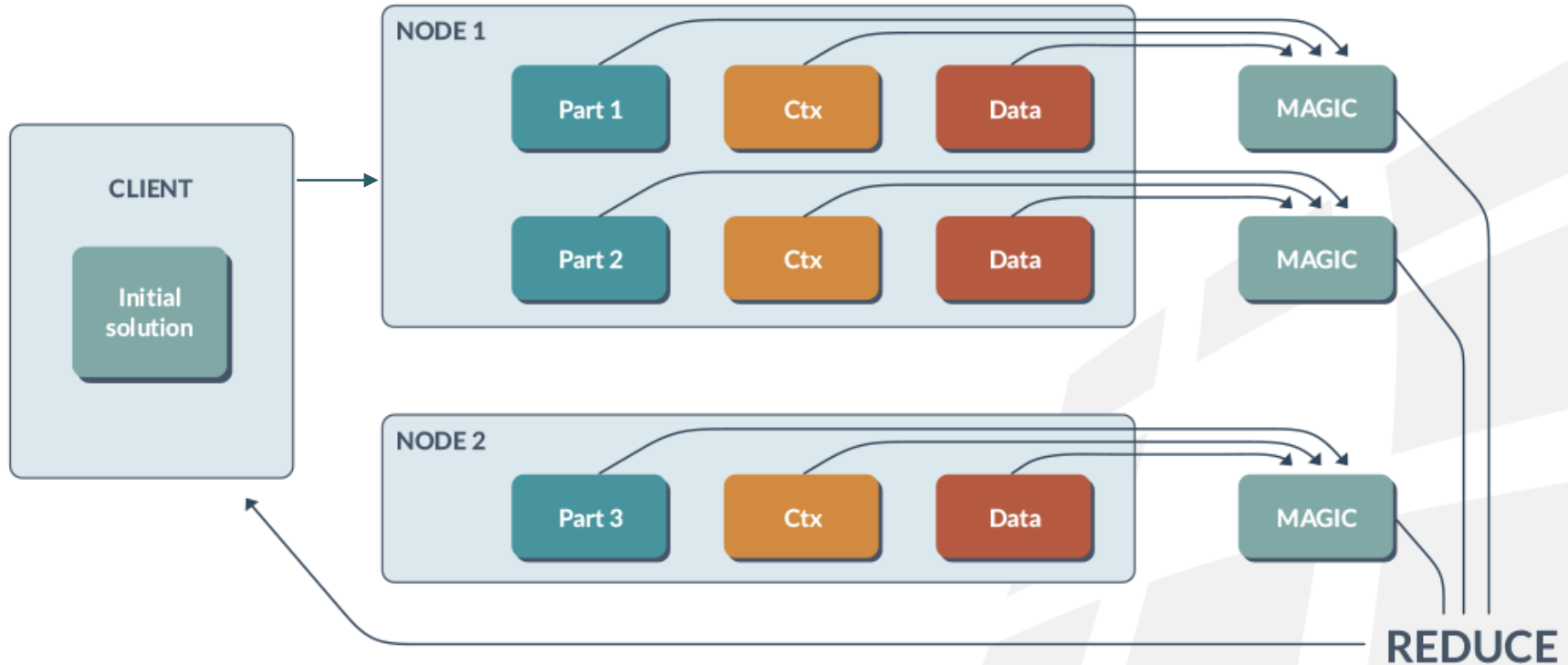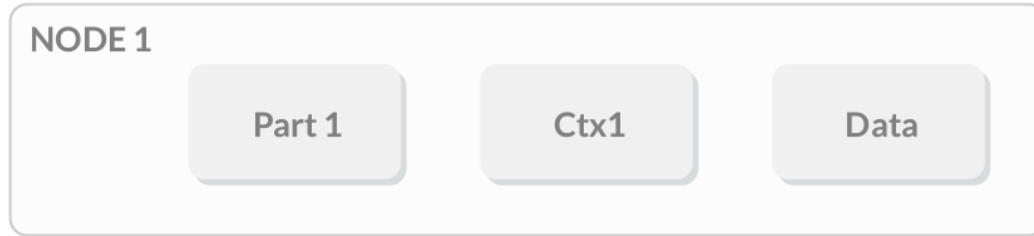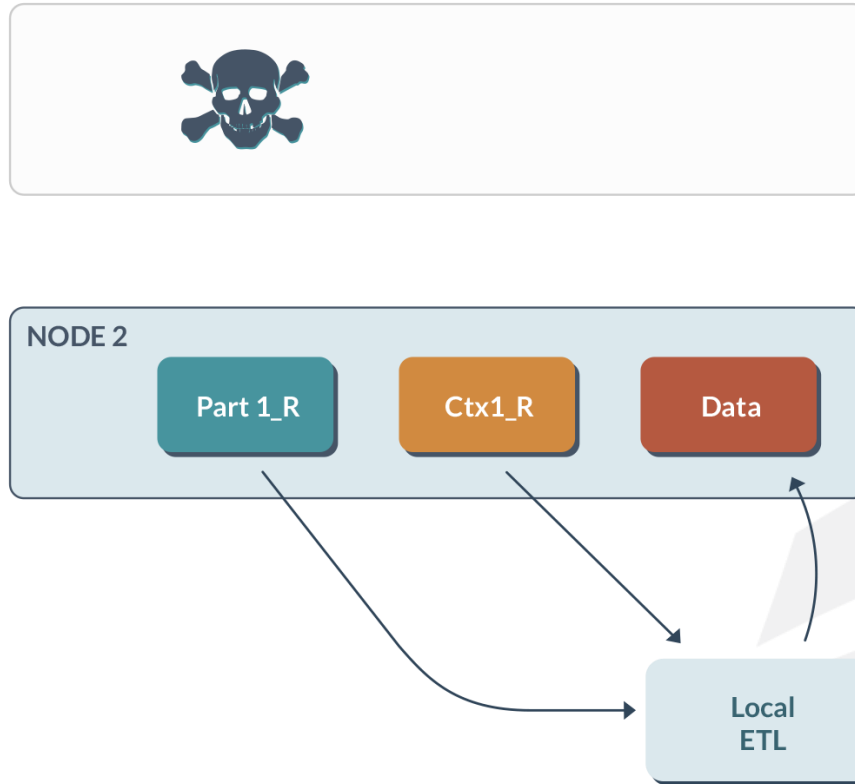
$$(f(x_1) - y_1)^2 + (f(x_2) - y_2)^2 + \cdots + (f(x_n) - y_n)^2$$

Gradient of loss function

$$2(f(x_1) - y_1) * f'(x_1) + 2(f(x_2) - y_2) * f'(x_2) + \cdots + 2(f(x_n) - y_n) * f'(x_n)$$

Node 1

Node 2

Node M

**GridGain**

# Sample 2 LSQR

## Golub-Kahan-Lanczos Bidiagonalization Procedure

Choose $v_1 = $ unit 2-norm vector and set $\beta_0 = 0$,

for $k = 1, 2..n$

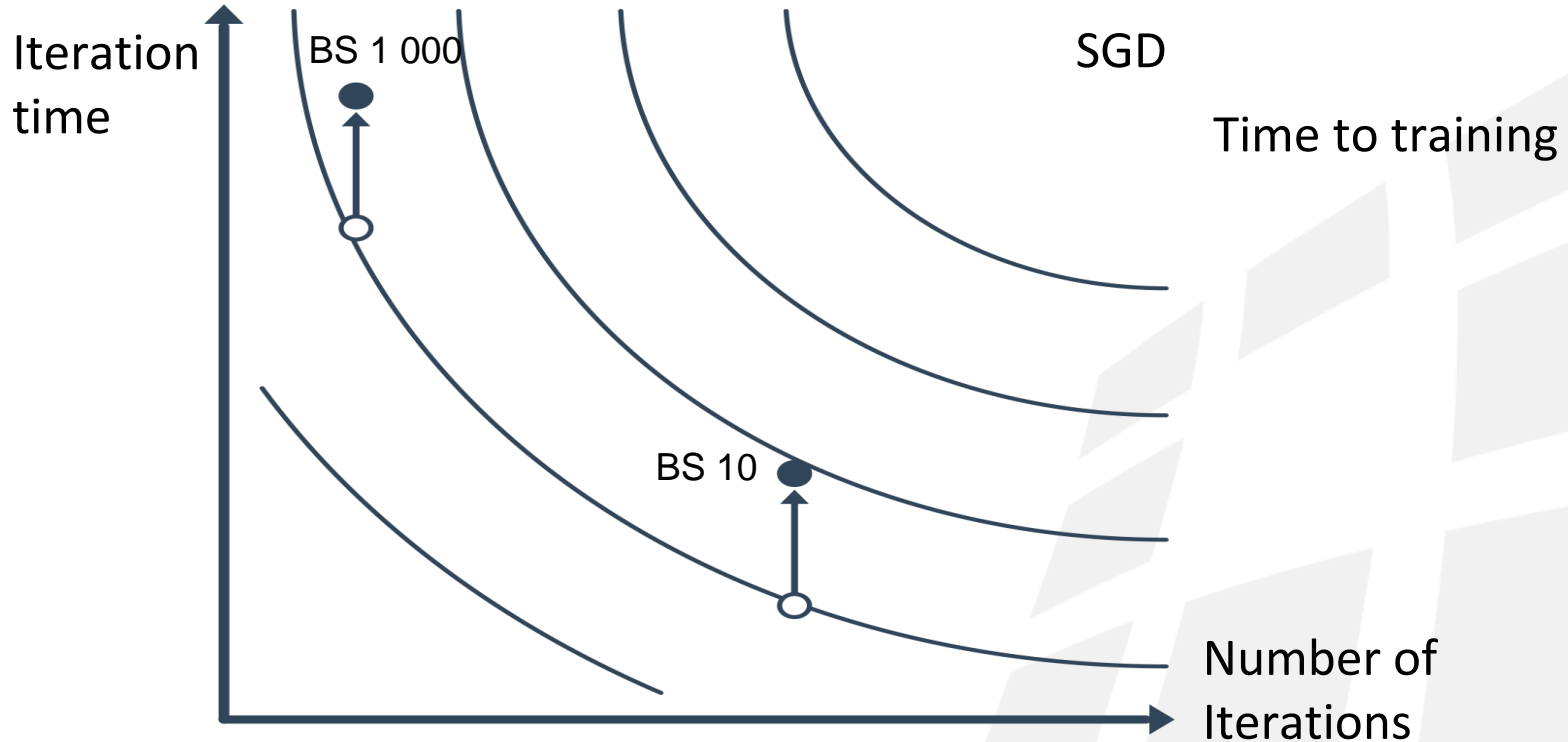$\quad v_k = Av_k - \beta_{k-1} u_{k-1}$

$\quad a_k = ||u_k||_2$

$\quad u_k = u_k / a_k$

$\quad v_{k+1} = A^* u_k - a_k v_k$

$\quad \beta_k = ||v_{k+1}||_2$

$\quad v_{k+1} = v_{k+1} / \beta_k$

end.

GridGain

# Limitations of Applicability

# Want to learn more?

https://ignite.apache.org

https://apacheignite.readme.io/docs

https://github.com/apache/ignite

ybabak@gridgain.com

**GridGain**