

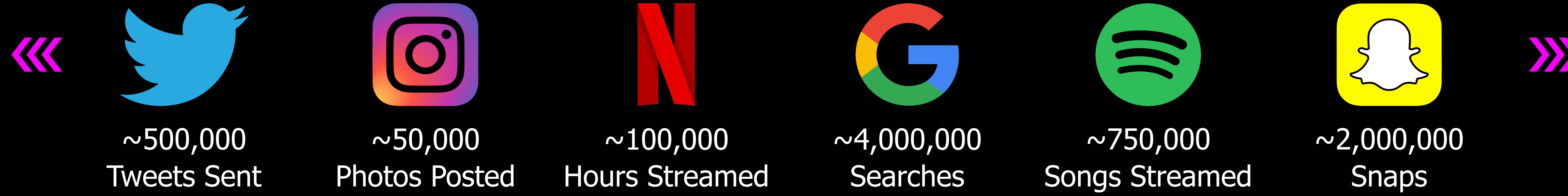
DINGO: Distributed Newton-Type Method for Gradient-Norm Optimization



Rixon Crane*, Fred Roosta*

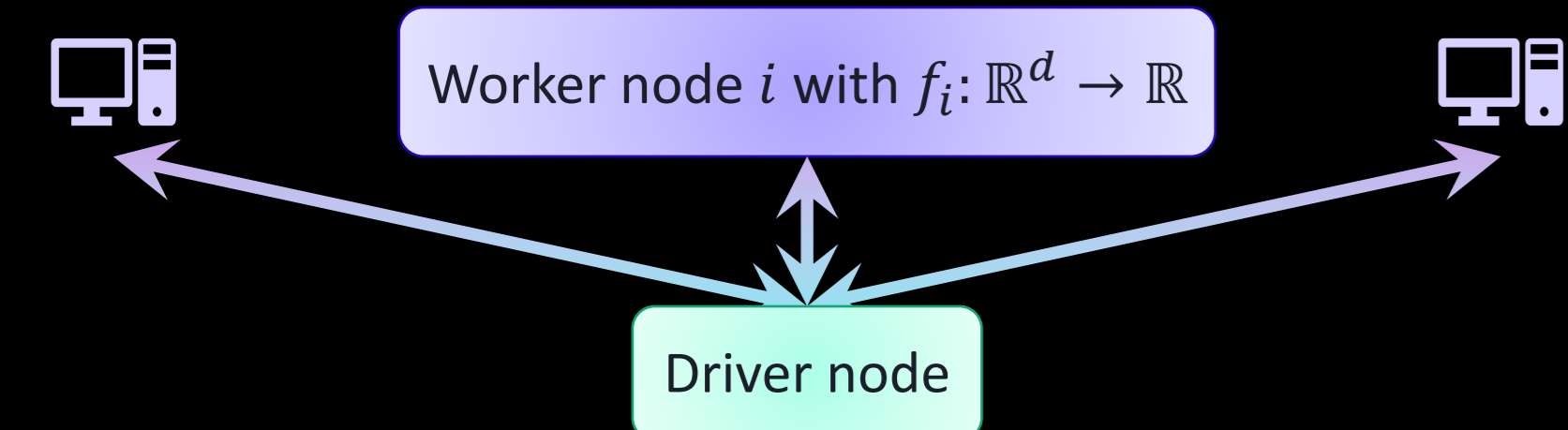
*School of Mathematics and Physics, University of Queensland, Australia.

Every Minute There Are:



The Problem

Centralized Computing Environment

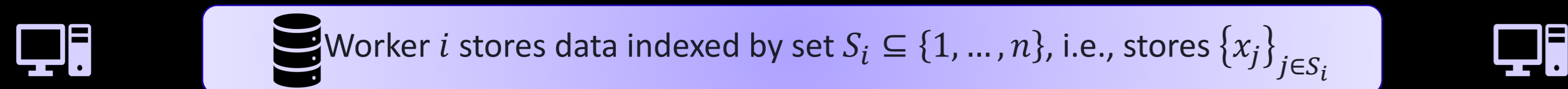


Goal
Over m workers, solve:

$$\min_{w \in \mathbb{R}^d} \left\{ f(w) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{i=1}^m f_i(w) \right\}.$$

Use Case: Big Data Regimes

Distributively Working With a Very Large Dataset $\{x_i\}_{i=1}^n$



Why use Second-Order Methods?

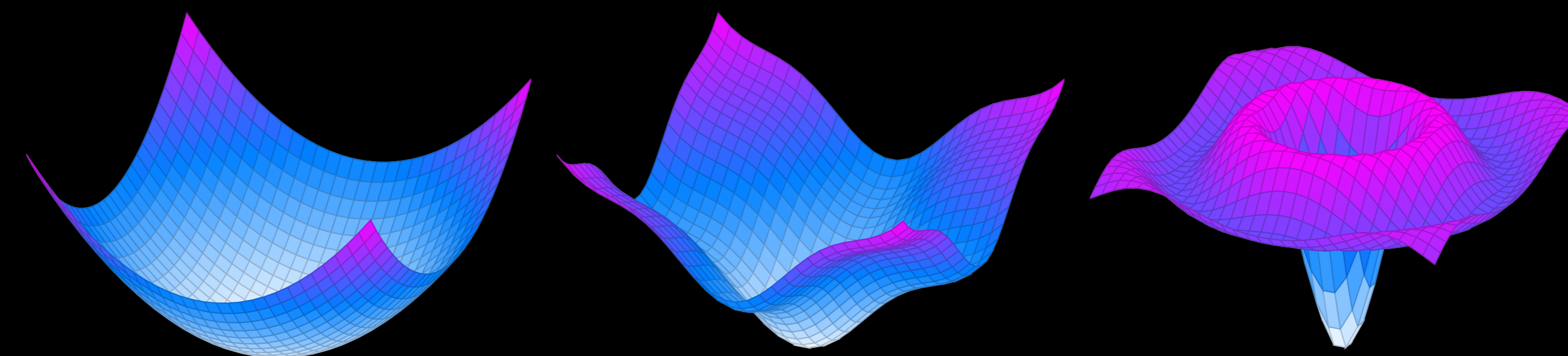
Second-order methods employ curvature (Hessian matrix) information to transform the gradient so that it is a more suitable direction to follow.

Benefits

- Perform more computations per iteration
- May take full advantage of available distributed computational resources
- May require significantly less communication costs
- Often require far fewer iterations to achieve similar results

Related Work

Method	Applicable to Non-Convex Functions	Arbitrary Data Distribution	Arbitrary Form of f_i	Simple Sub-Problems	Not Sensitive to Hyper-Parameters
GIANT	X	X	X	✓	✓
DISCO	X	✓	✓	✓	✓
DANE	✓	✓	✓	X	X
InexactDANE	✓	✓	✓	X	X
AIDE	✓	✓	✓	X	X
DINGO	✓	✓	✓	✓	✓



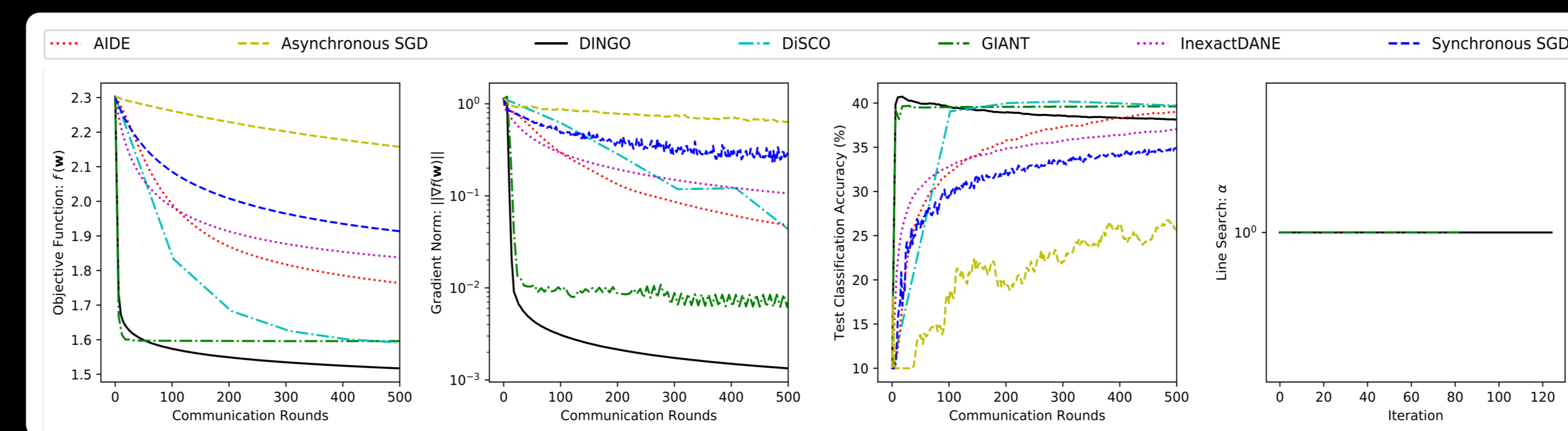
- Convex**
 - Hessian is positive semidefinite.
 - Local minima are global minima.
- Invex**
 - Hessian can be indefinite and singular.
 - Local minima are global minima.
- Non-Convex**
 - Hessian can be indefinite and singular.
 - Not all local minima are global minima.

Our Method: DINGO

Derived by optimization of the gradient's norm as a surrogate function, i.e.,

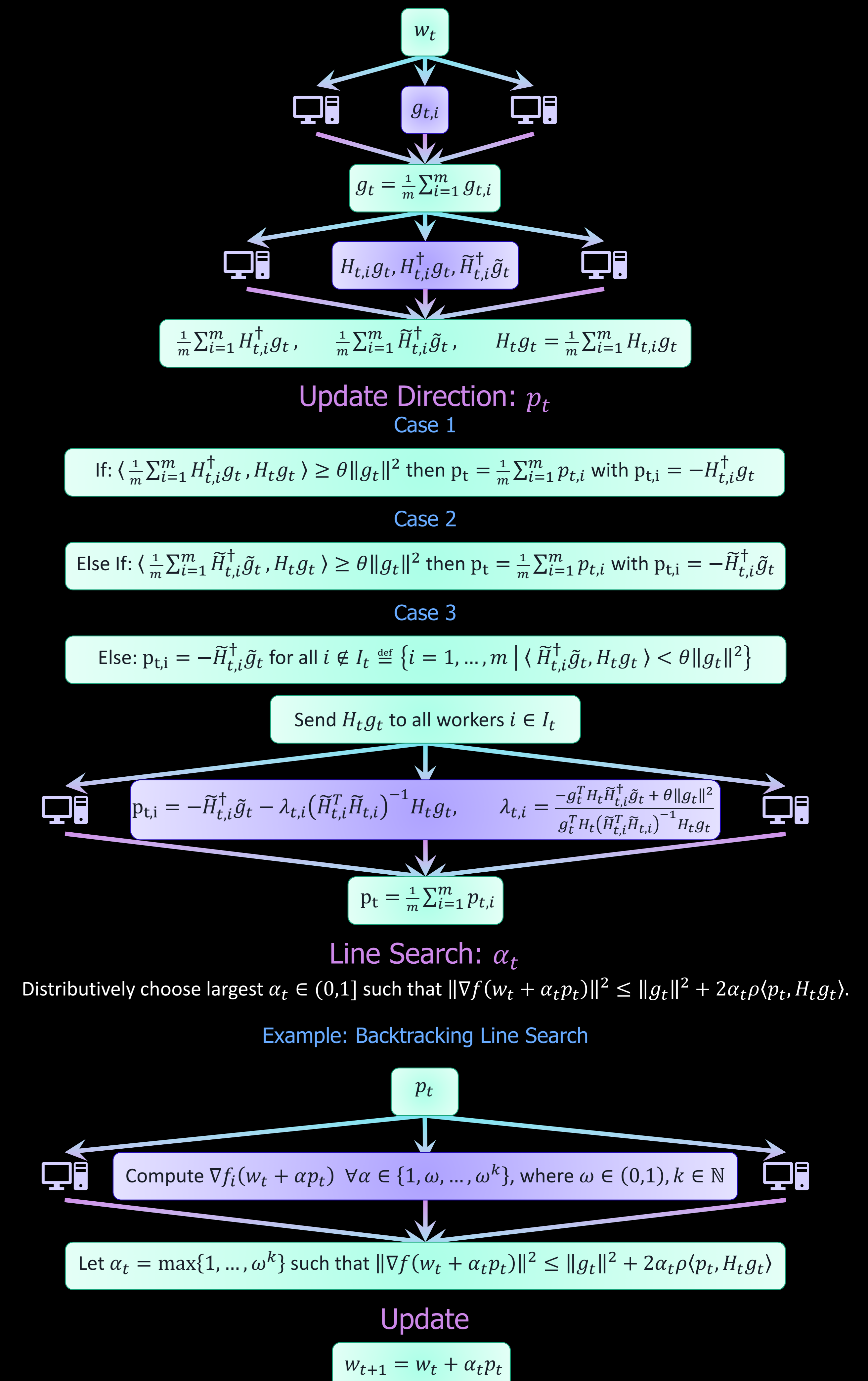
$$\min_{w \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\nabla f(w)\|^2 = \frac{1}{2m^2} \left\| \sum_{i=1}^m \nabla f_i(w) \right\|^2 \right\}.$$

DINGO is for "Distributed Newton-type method for Gradient-norm Optimization". DINGO is particularly suitable for invex objectives. A strict linear-rate reduction in the gradient norm is always guaranteed.



Softmax regression, with regularization, problem on the CIFAR10 dataset.

Each Iteration of DINGO



The constants $\theta, \phi > 0$ and $\rho \in (0,1)$ are hyper-parameters. The vector $w_t \in \mathbb{R}^d$ denotes the point at iteration t . For notational convenience, we denote $g_{t,i} \stackrel{\text{def}}{=} \nabla f_i(w_t)$, $H_{t,i} \stackrel{\text{def}}{=} \nabla^2 f_i(w_t)$, $g_t \stackrel{\text{def}}{=} \nabla f(w_t)$, $H_t \stackrel{\text{def}}{=} \nabla^2 f(w_t)$. We also let

$$\tilde{H}_{t,i} \stackrel{\text{def}}{=} \begin{bmatrix} H_{t,i} \\ \phi I \end{bmatrix} \in \mathbb{R}^{2d \times 2d}, \quad \tilde{g}_t \stackrel{\text{def}}{=} \begin{bmatrix} g_t \\ 0 \end{bmatrix} \in \mathbb{R}^{2d},$$

where I is the identity matrix and 0 is the zero vector. Green and purple rectangles represent the driver node and worker nodes, respectively.