

COCOA

Communication-Efficient Coordinate Ascent

Virginia Smith

*Martin Jaggi, Martin Takáč, Jonathan Terhorst,
Sanjay Krishnan, Thomas Hofmann, & Michael I. Jordan*



LARGE-SCALE OPTIMIZATION

LARGE-SCALE OPTIMIZATION

COCO

LARGE-SCALE OPTIMIZATION

COCOA

RESULTS

LARGE-SCALE OPTIMIZATION

COCOA

RESULTS

Machine Learning with Large Datasets

Machine Learning with Large Datasets

You  Tube

IM  GENET



NETFLIX

illumina

facebook

twitter 

Google

VISA

Machine Learning with Large Datasets

You  Tube

IM  GENET

 amazon.com

NETFLIX

illumina

facebook

twitter 

Google

VISA

image/music/video tagging
document categorization
item recommendation
click-through rate prediction
sequence tagging
protein structure prediction
sensor data prediction
spam classification
fraud detection

Machine Learning Workflow

Machine Learning Workflow



DATA & PROBLEM

*classification, regression,
collaborative filtering, ...*

Machine Learning Workflow



DATA & PROBLEM

*classification, regression,
collaborative filtering, ...*

MACHINE LEARNING MODEL

*logistic regression, lasso,
support vector machines, ...*

Machine Learning Workflow



DATA & PROBLEM

*classification, regression,
collaborative filtering, ...*

MACHINE LEARNING MODEL

*logistic regression, lasso,
support vector machines, ...*

OPTIMIZATION ALGORITHM

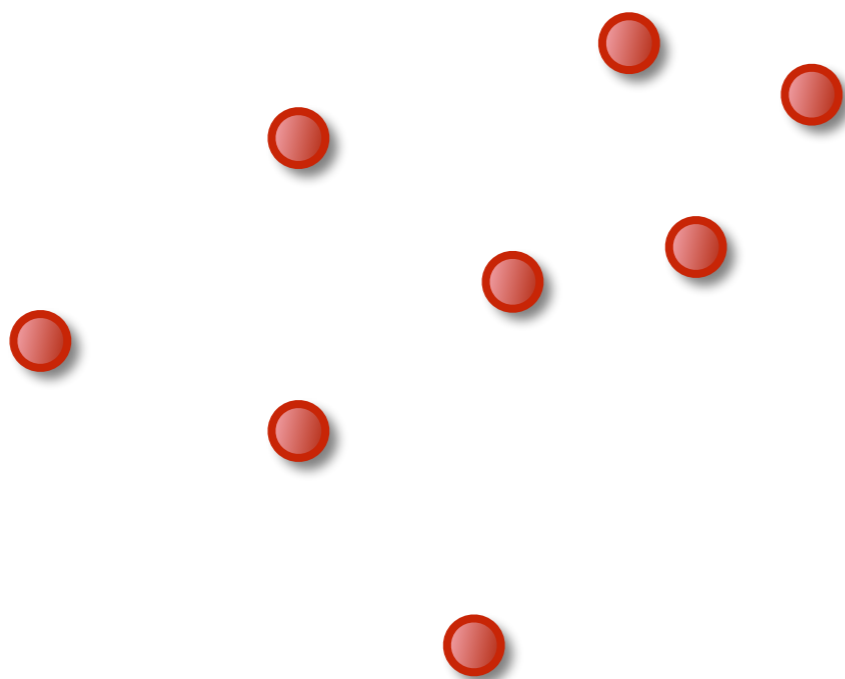
*gradient descent, coordinate
descent, Newton's method, ...*

Example: SVM Classification

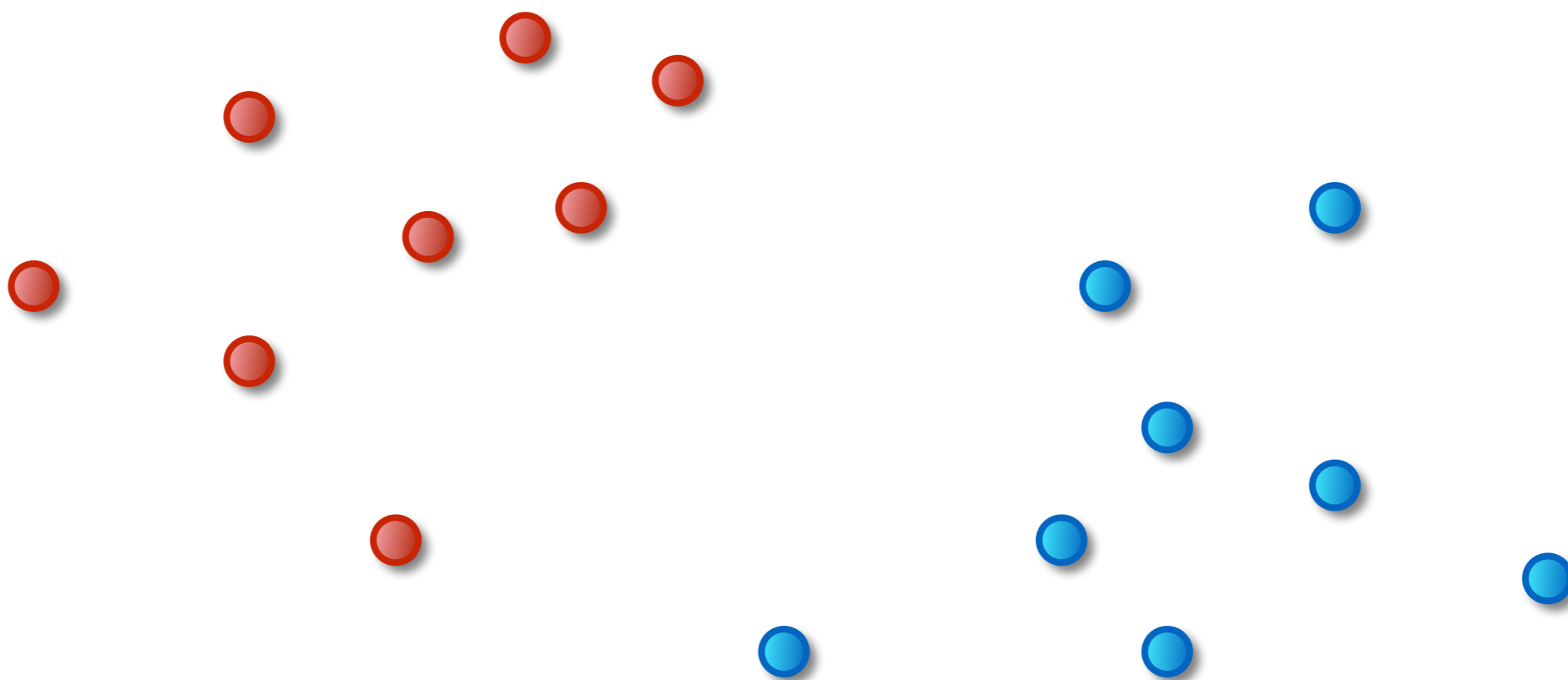
Example: SVM Classification



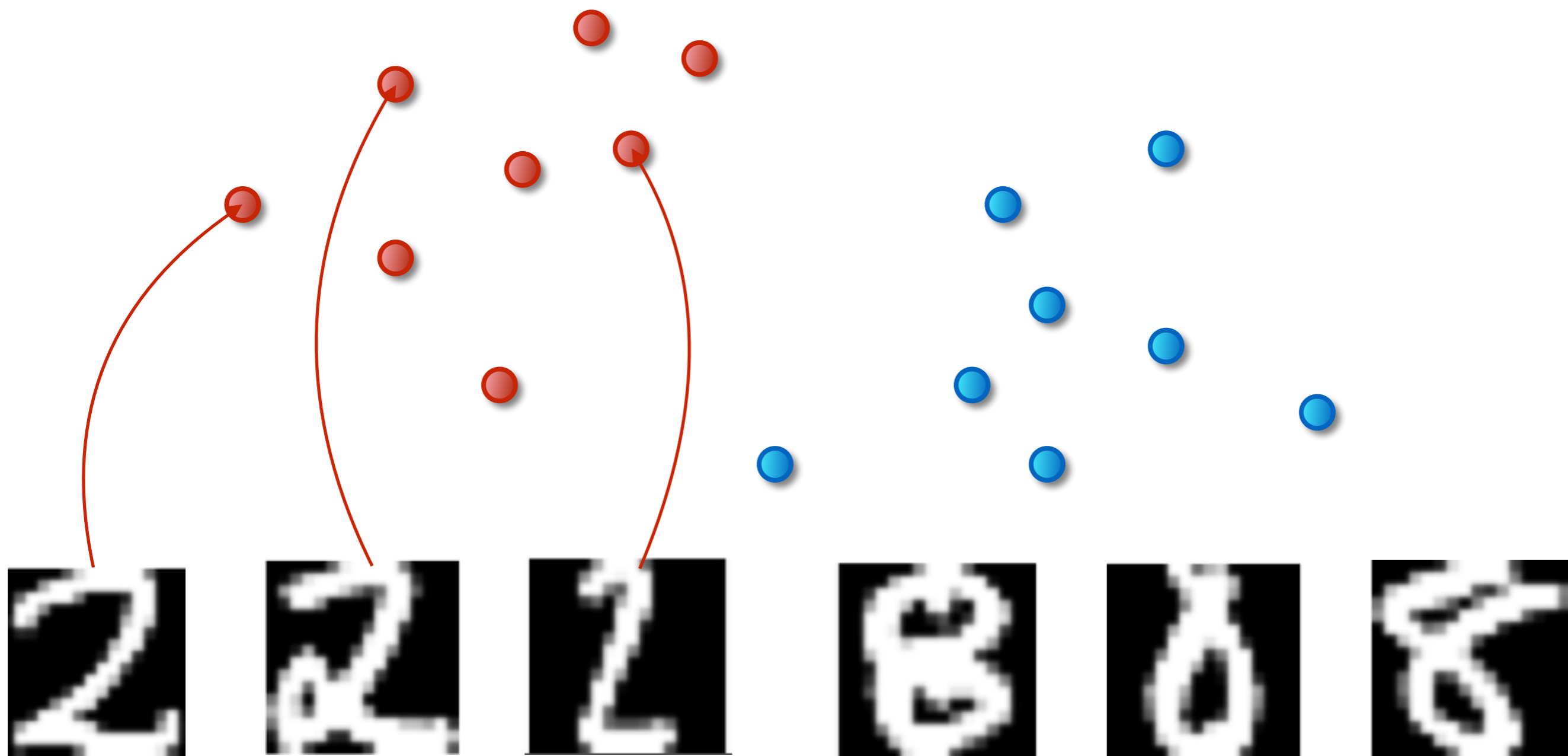
Example: SVM Classification



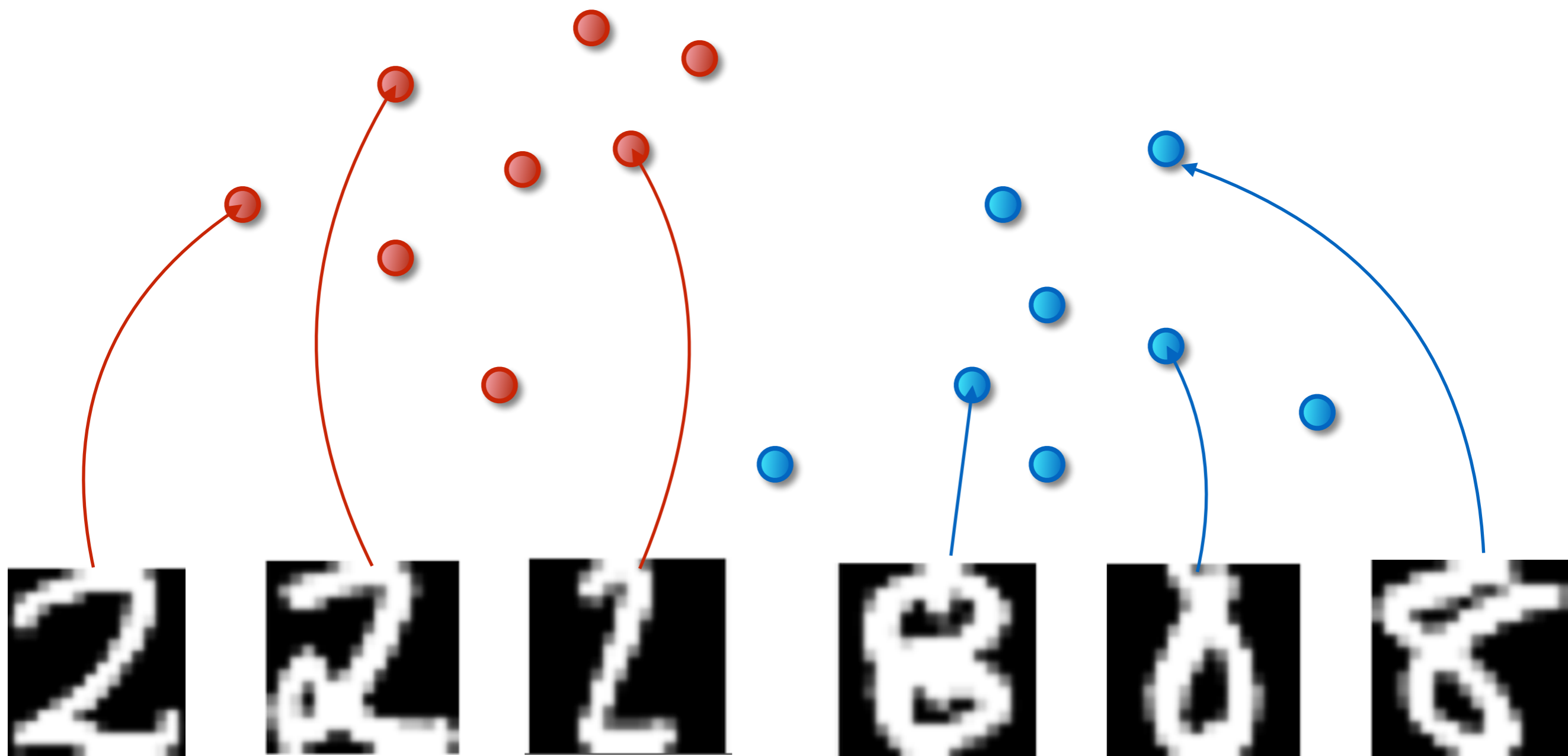
Example: SVM Classification



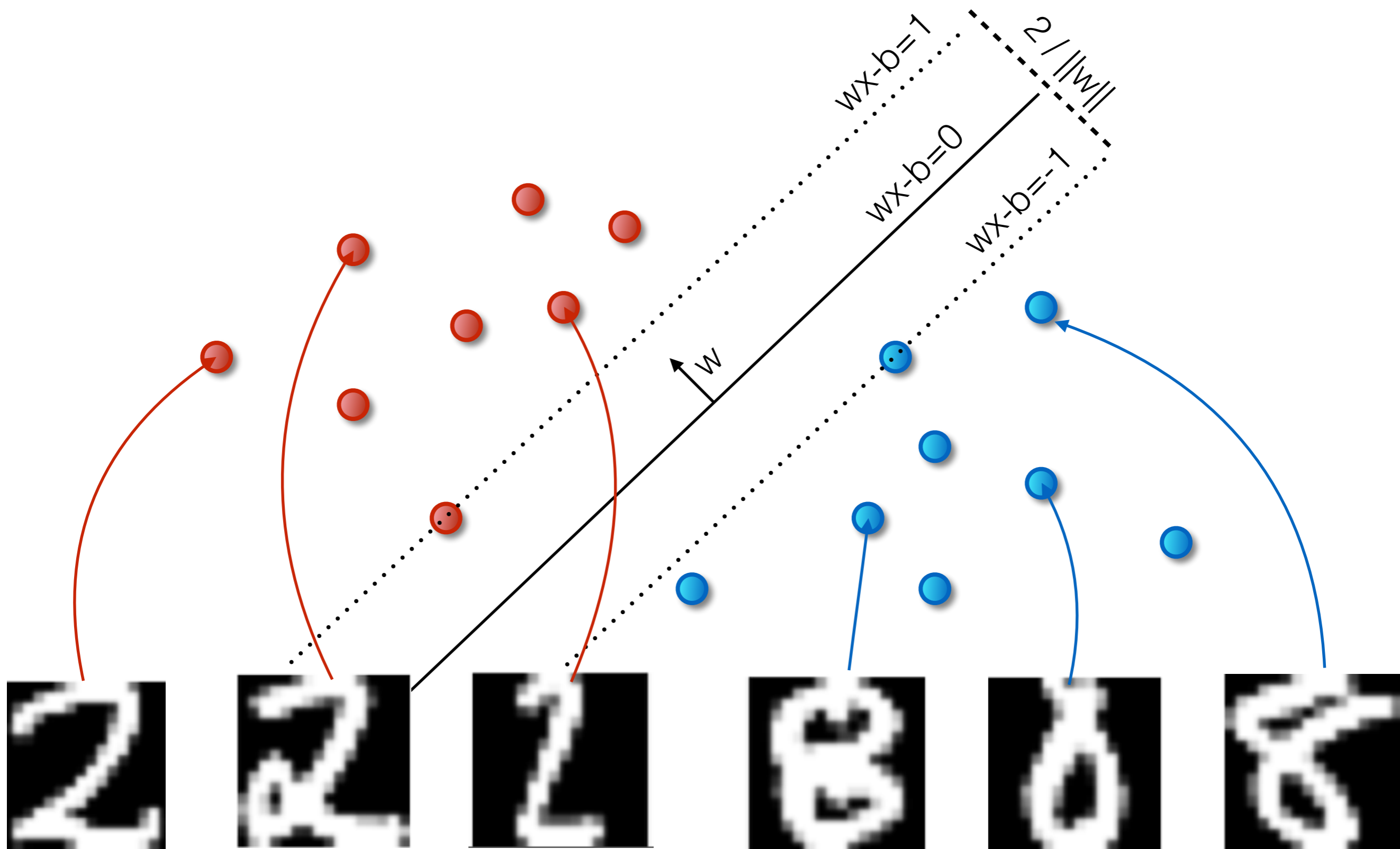
Example: SVM Classification



Example: SVM Classification

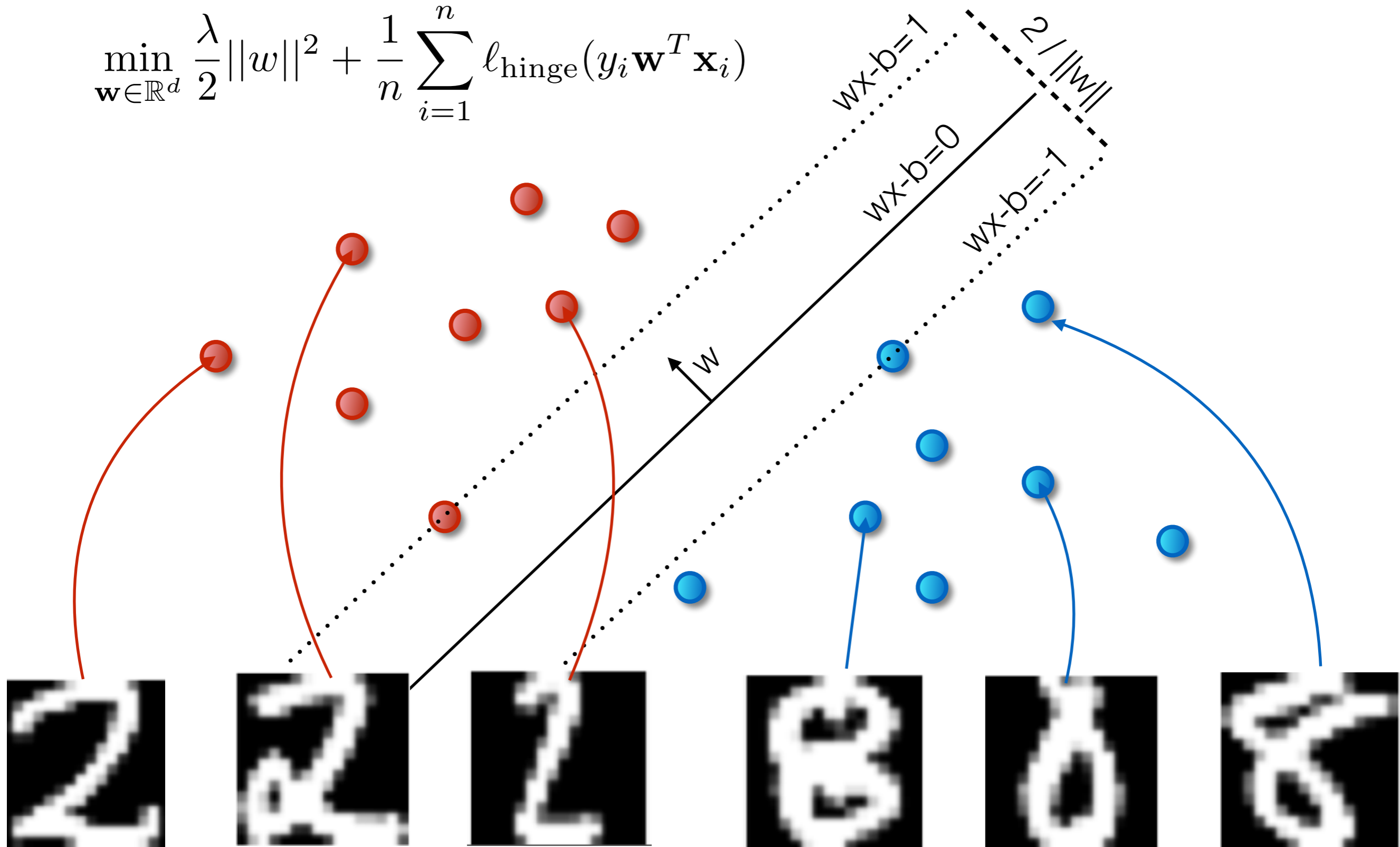


Example: SVM Classification



Example: SVM Classification

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y_i \mathbf{w}^T \mathbf{x}_i)$$



Example: SVM Classification

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_{\text{hinge}}(y_i \mathbf{w}^T \mathbf{x}_i)$$

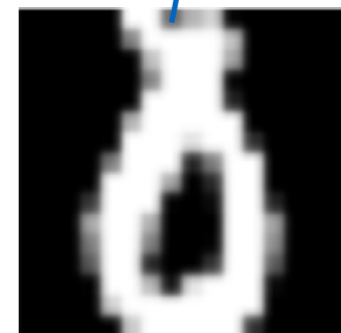
$$\begin{aligned} & \text{--- } x \cdot b = 1 \\ & \dots \dots \dots \\ & \text{--- } 2 / \|\mathbf{w}\| \end{aligned}$$

Descent algorithms and line search methods
Acceleration, momentum, and conjugate gradients
Newton and Quasi-Newton methods
Coordinate descent
Stochastic and incremental gradient methods

SMO

SVM^{light}

LIBLINEAR



Linear Regularized Loss Minimization

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression
- ▶ lasso regression

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression
- ▶ lasso regression
- ▶ ridge regression

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression
- ▶ lasso regression
- ▶ ridge regression
- ▶ etc...

Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression
- ▶ lasso regression
- ▶ ridge regression
- ▶ etc...

image/music/video tagging
document categorization
item recommendation
click-through rate prediction
sequence tagging
protein structure prediction
sensor data prediction
spam classification
fraud detection

Machine Learning Workflow



DATA & PROBLEM

*classification, regression,
collaborative filtering, ...*

MACHINE LEARNING MODEL

*logistic regression, lasso,
support vector machines, ...*

OPTIMIZATION ALGORITHM

*gradient descent, coordinate
descent, Newton's method, ...*

Machine Learning Workflow



DATA & PROBLEM

*classification, regression,
collaborative filtering, ...*

MACHINE LEARNING MODEL

*logistic regression, lasso,
support vector machines, ...*

SYSTEMS SETTING

*multi-core, cluster, cloud,
supercomputer, ...*

OPTIMIZATION ALGORITHM

*gradient descent, coordinate
descent, Newton's method, ...*

Machine Learning Workflow

classification regression

Open Problem:

efficiently solving objective
when data is **distributed**

OPTIMIZATION ALGORITHM

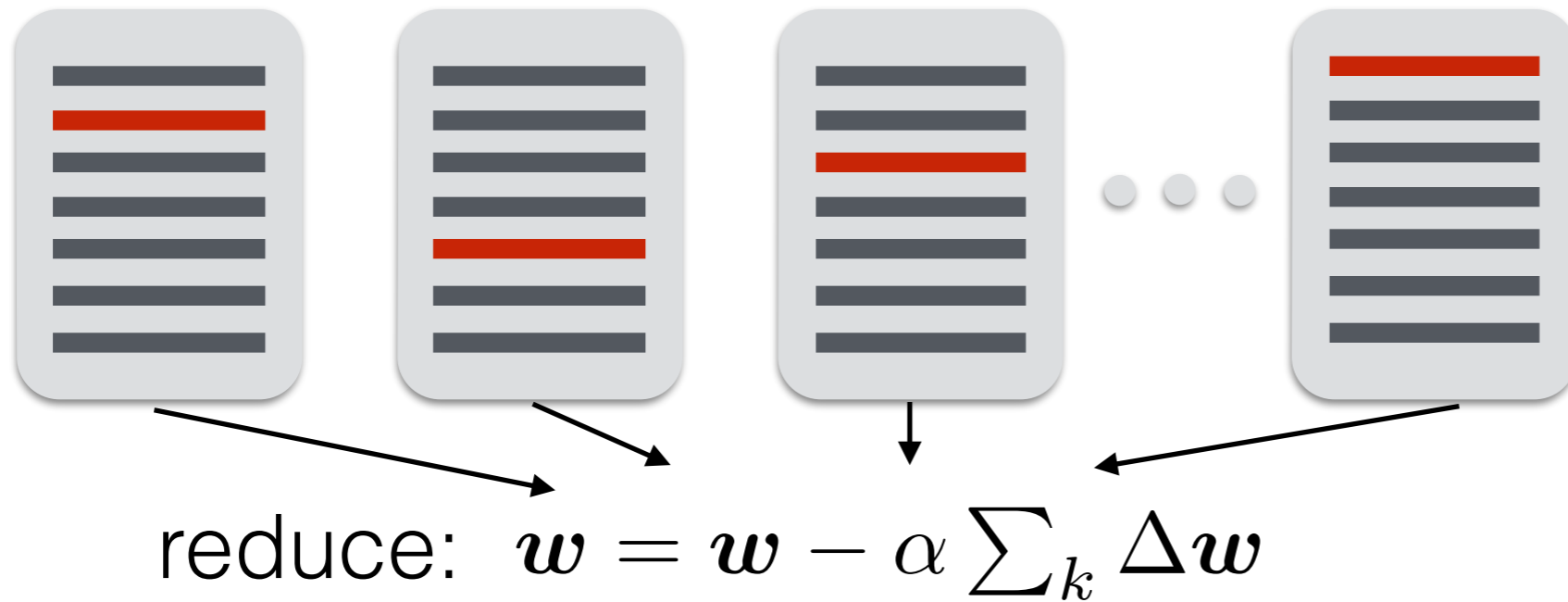
*gradient descent, coordinate
descent, Newton's method, ...*

Distributed Optimization

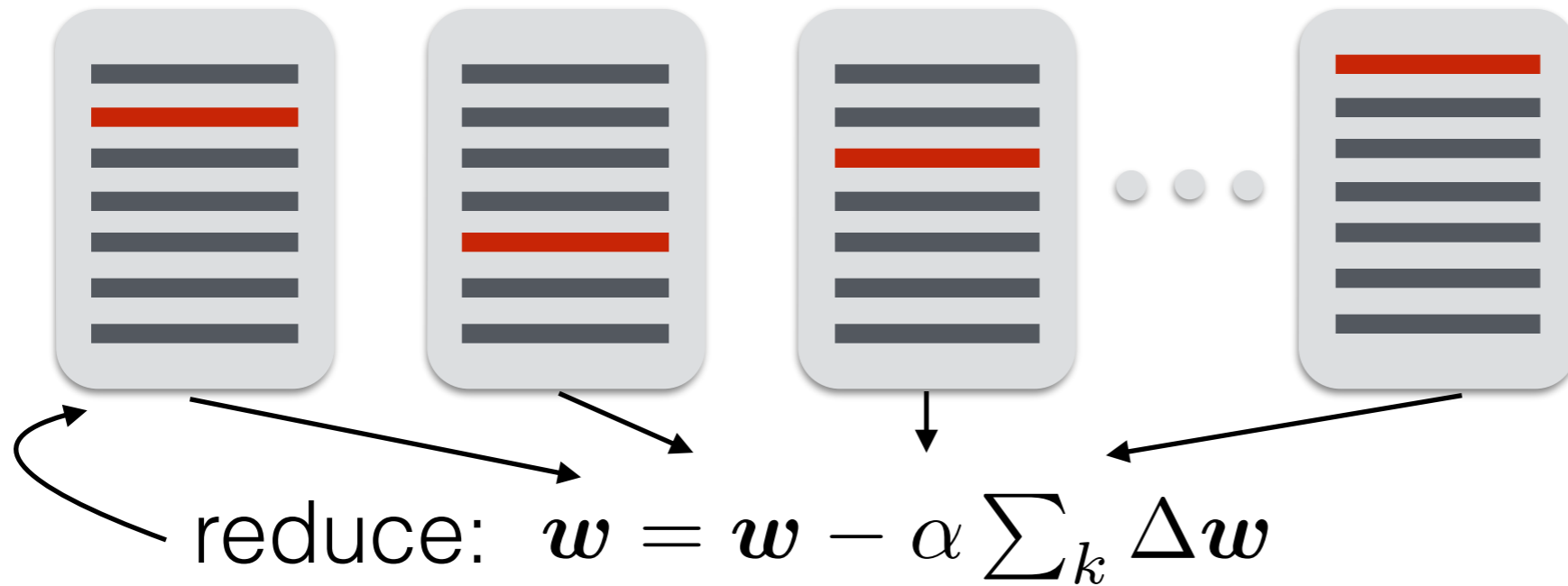
Distributed Optimization



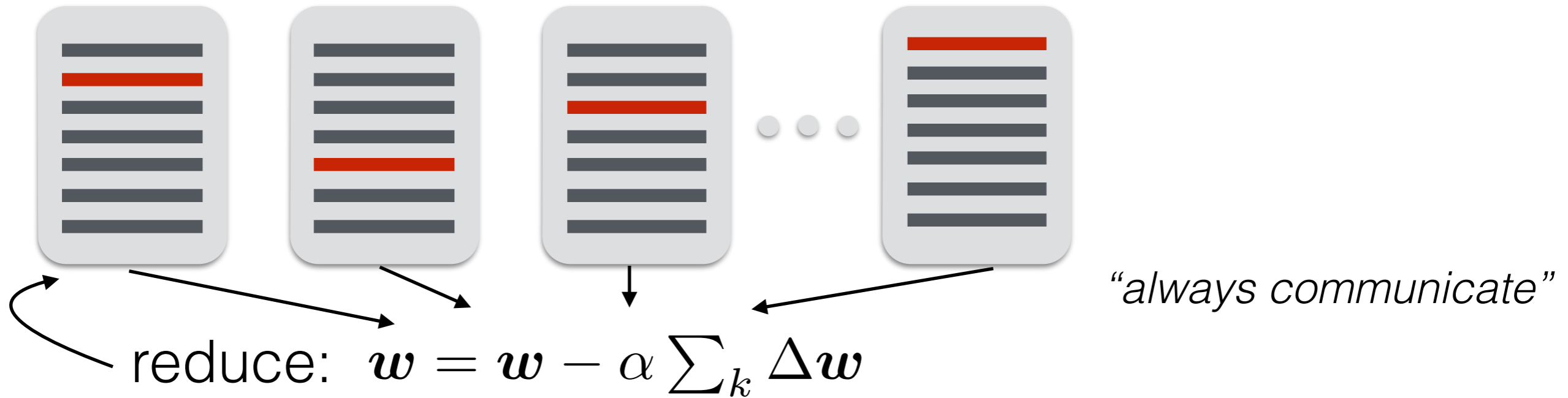
Distributed Optimization



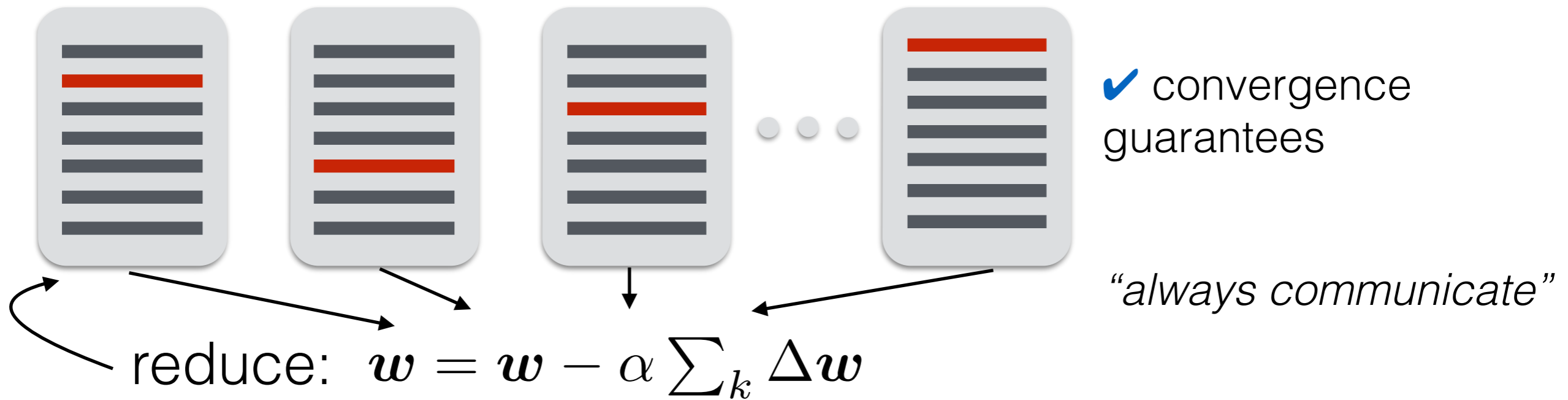
Distributed Optimization



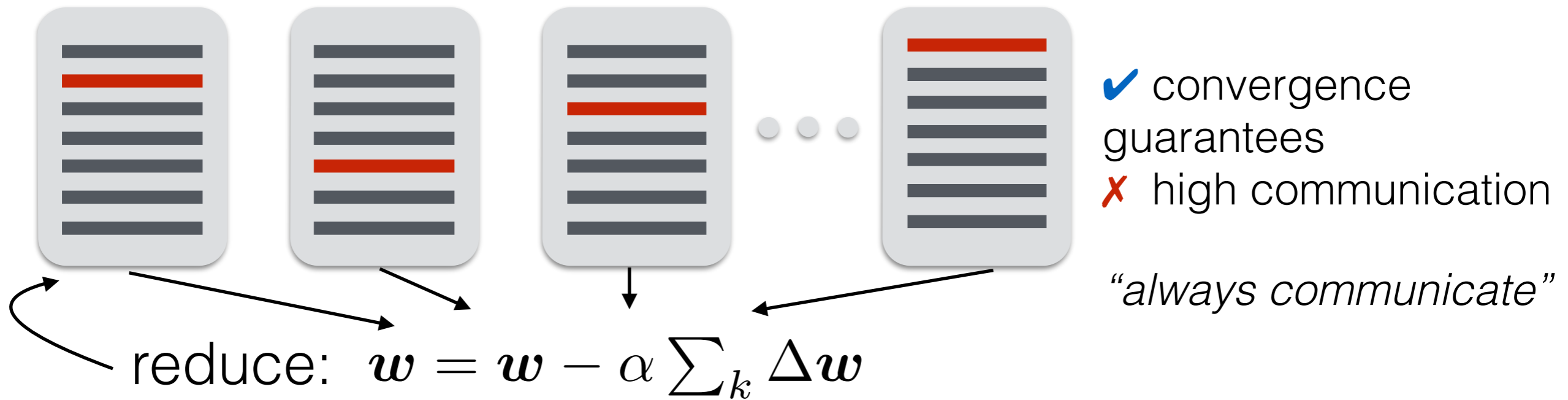
Distributed Optimization



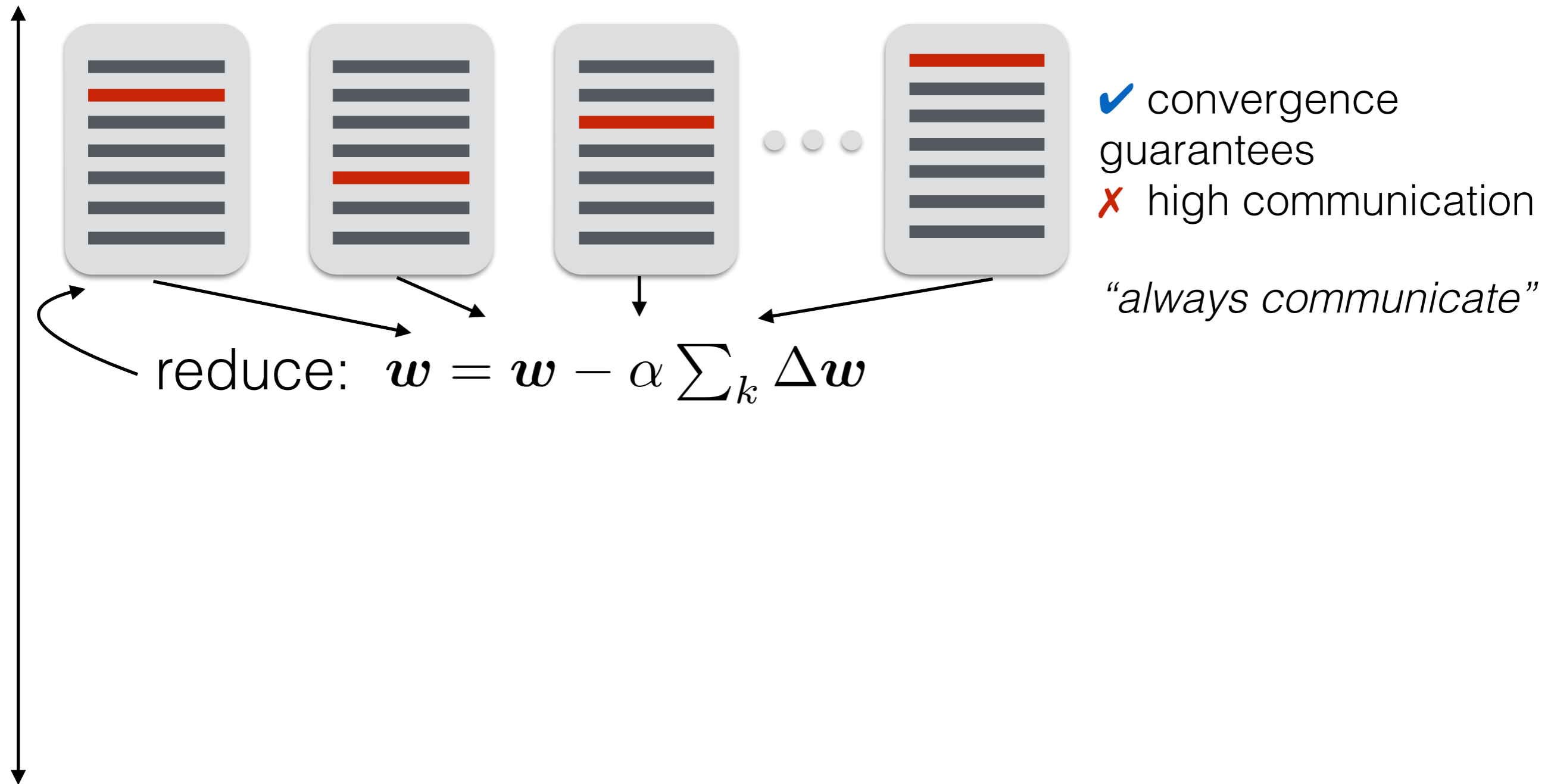
Distributed Optimization



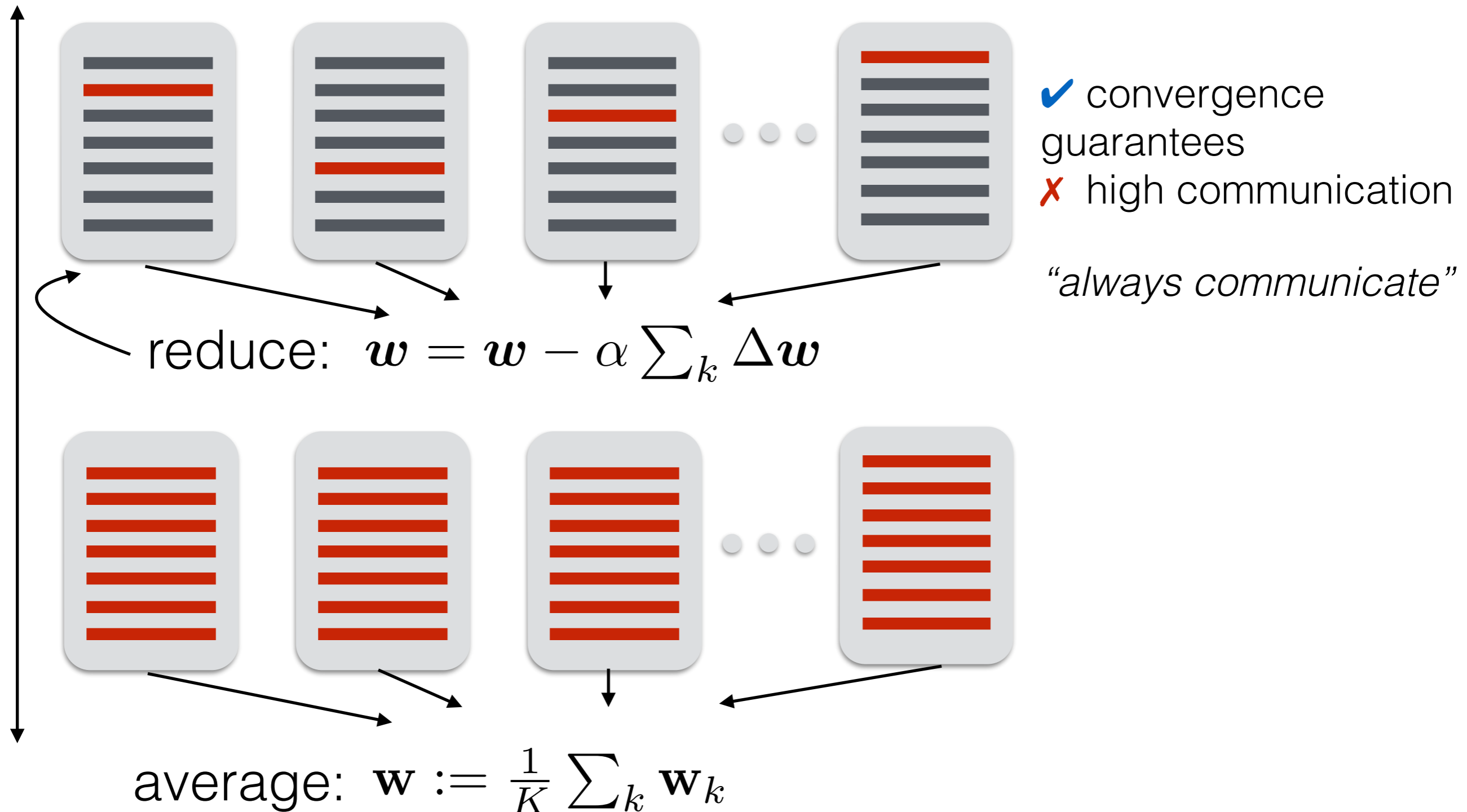
Distributed Optimization



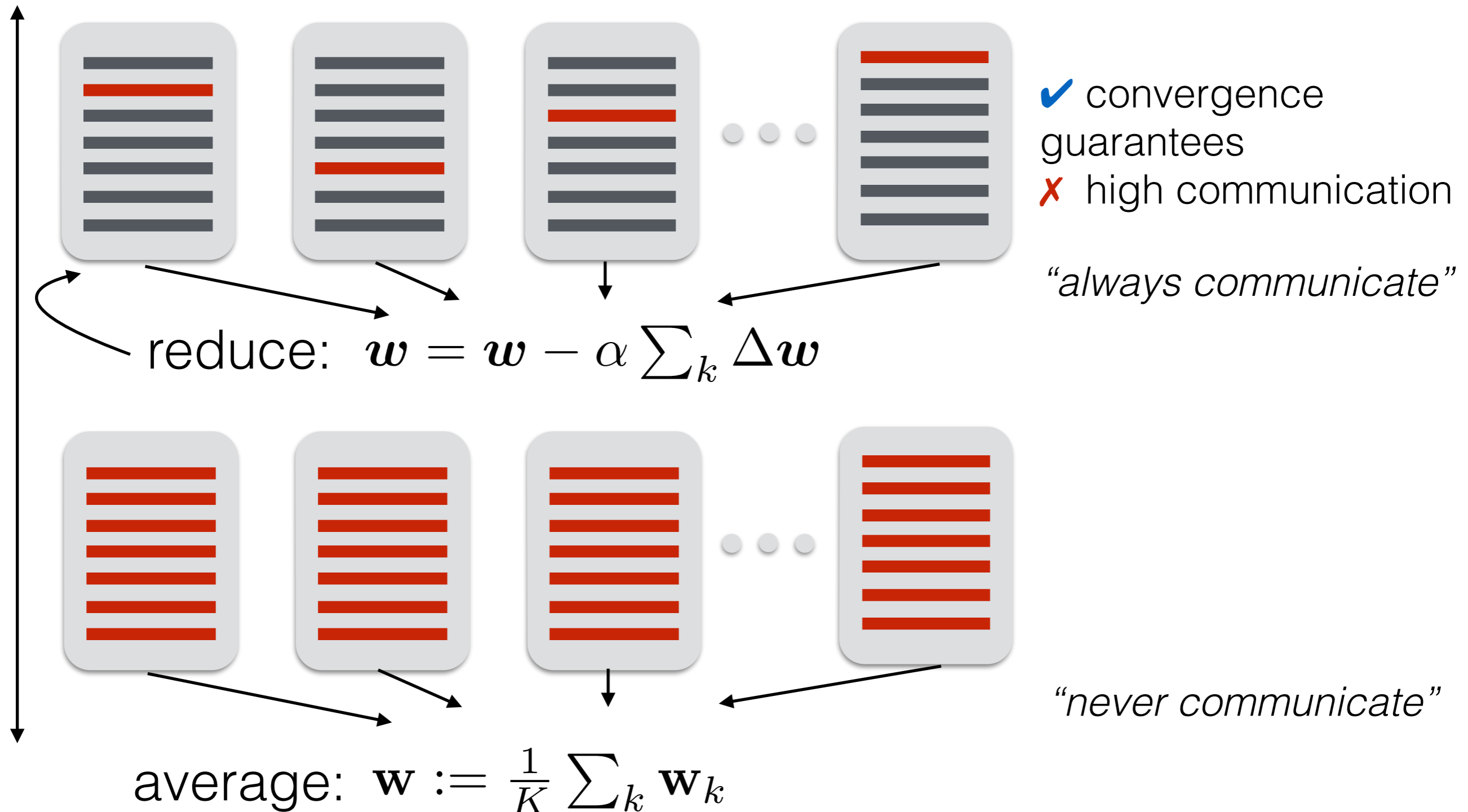
Distributed Optimization



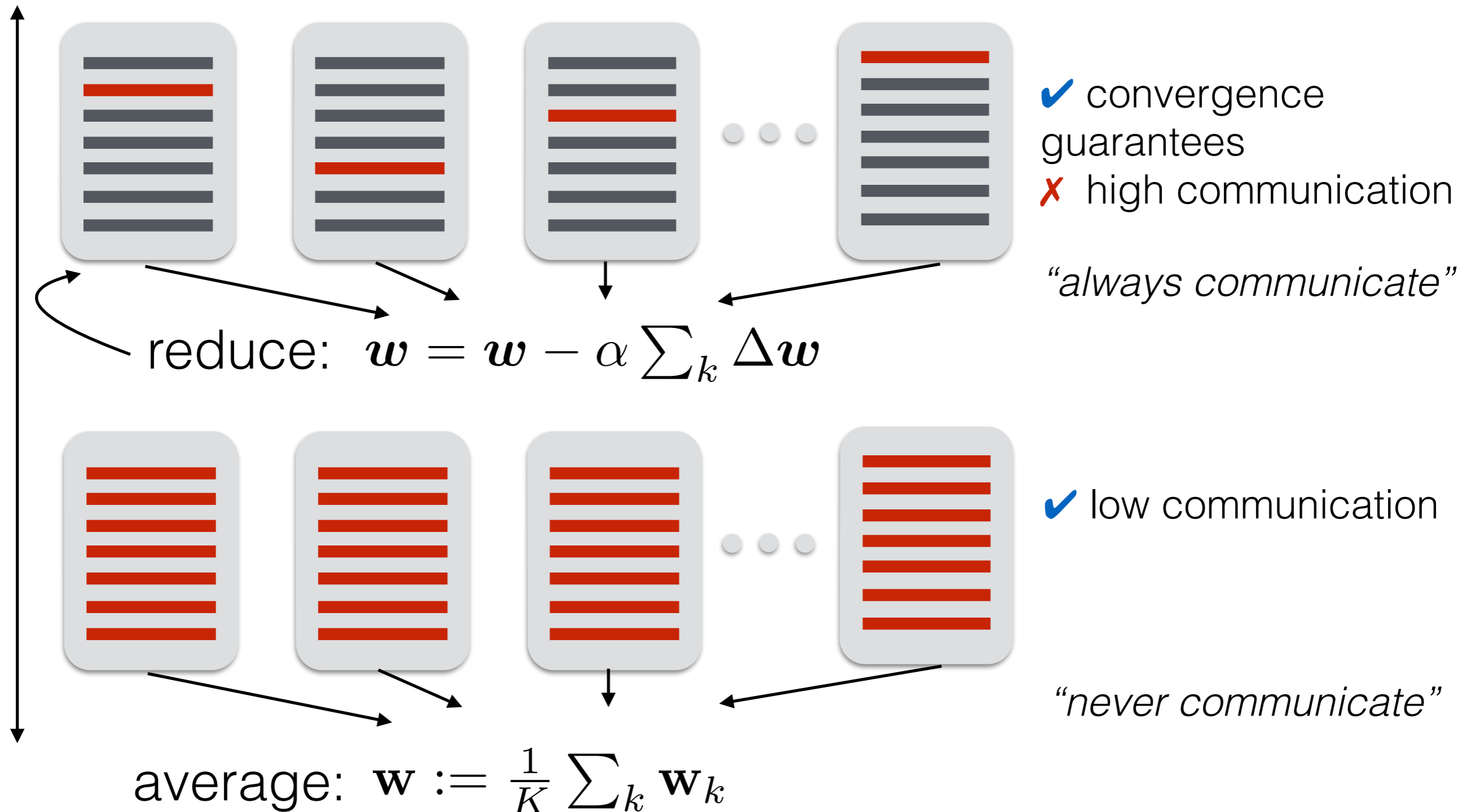
Distributed Optimization



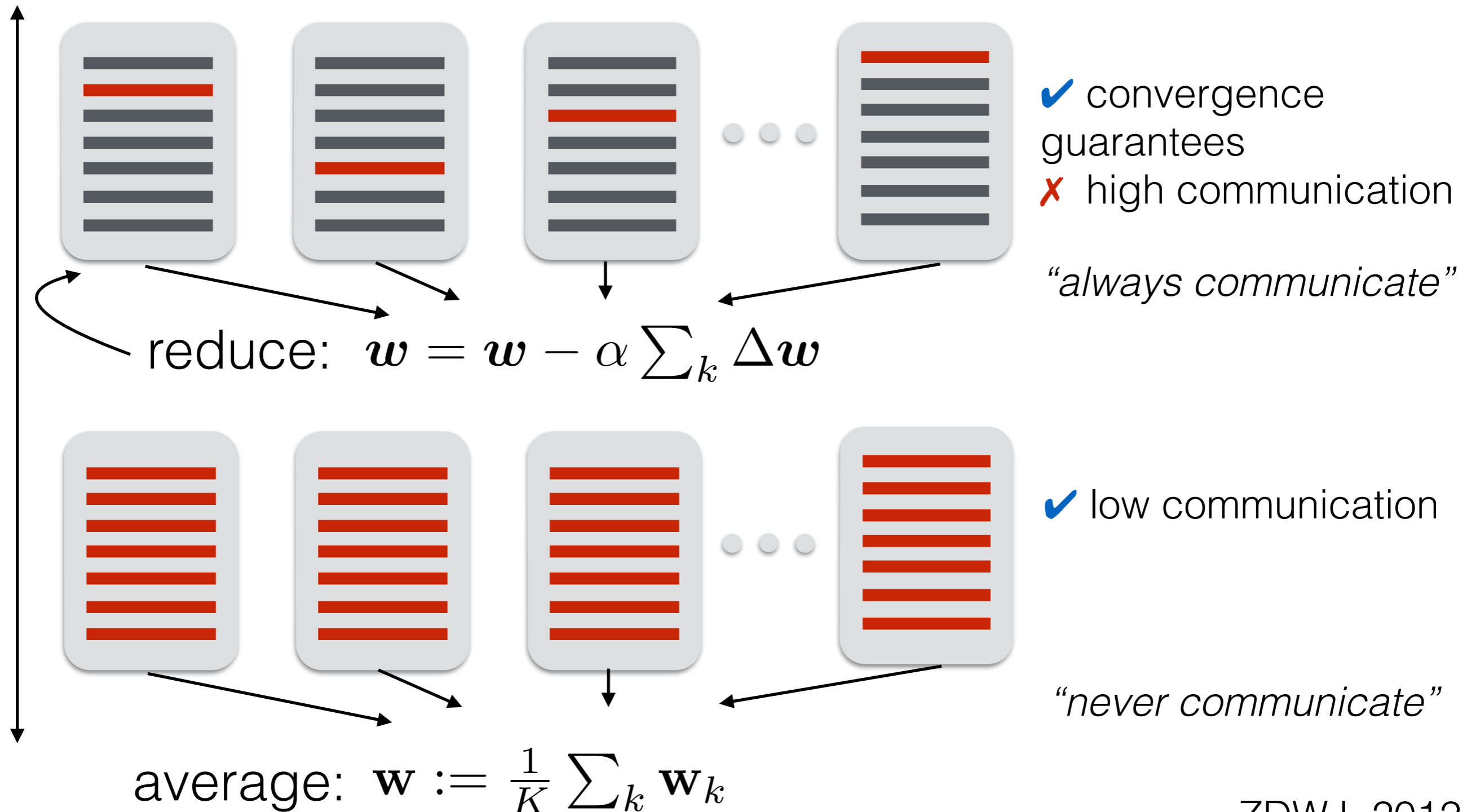
Distributed Optimization



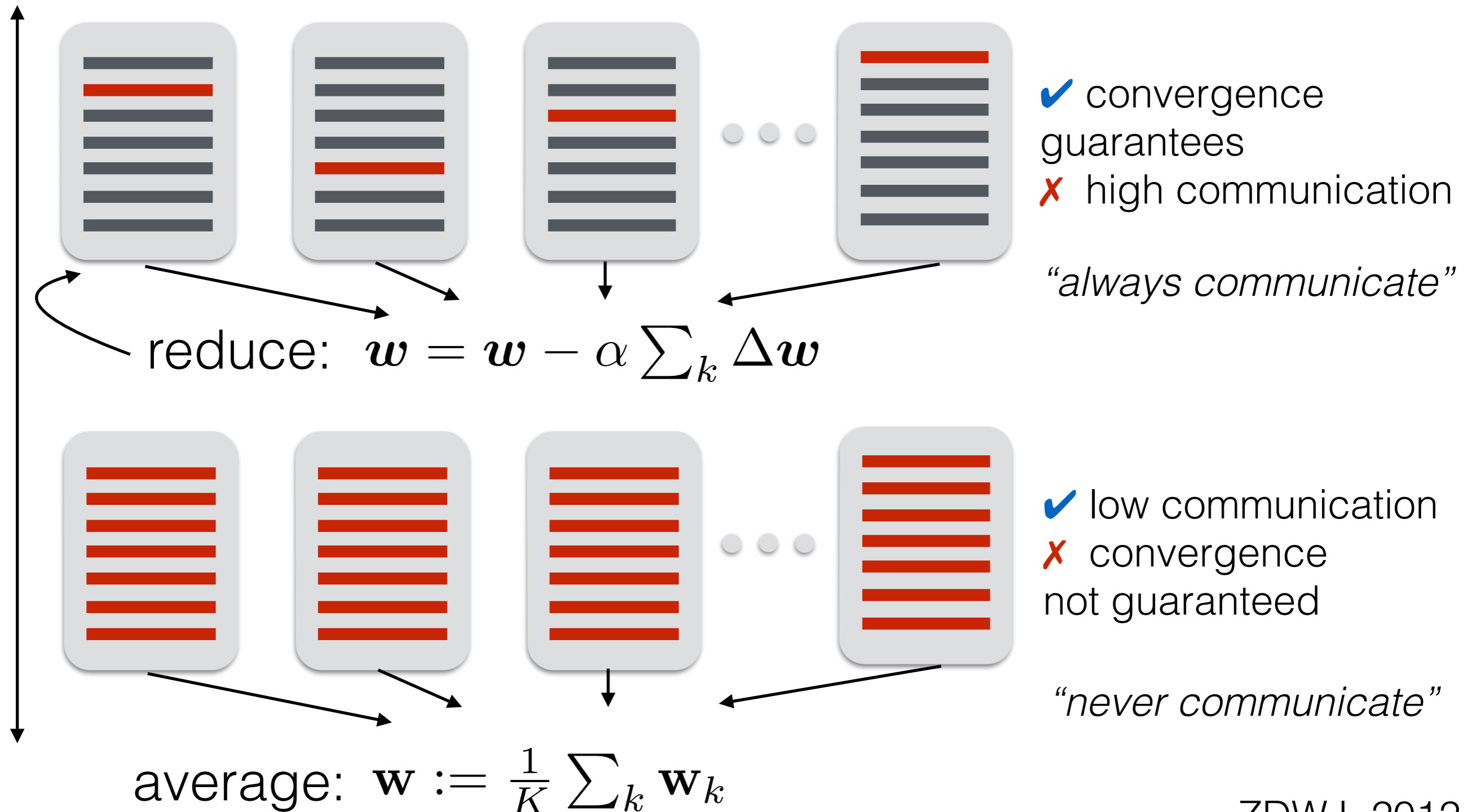
Distributed Optimization



Distributed Optimization



Distributed Optimization



Distributed Optimization

average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



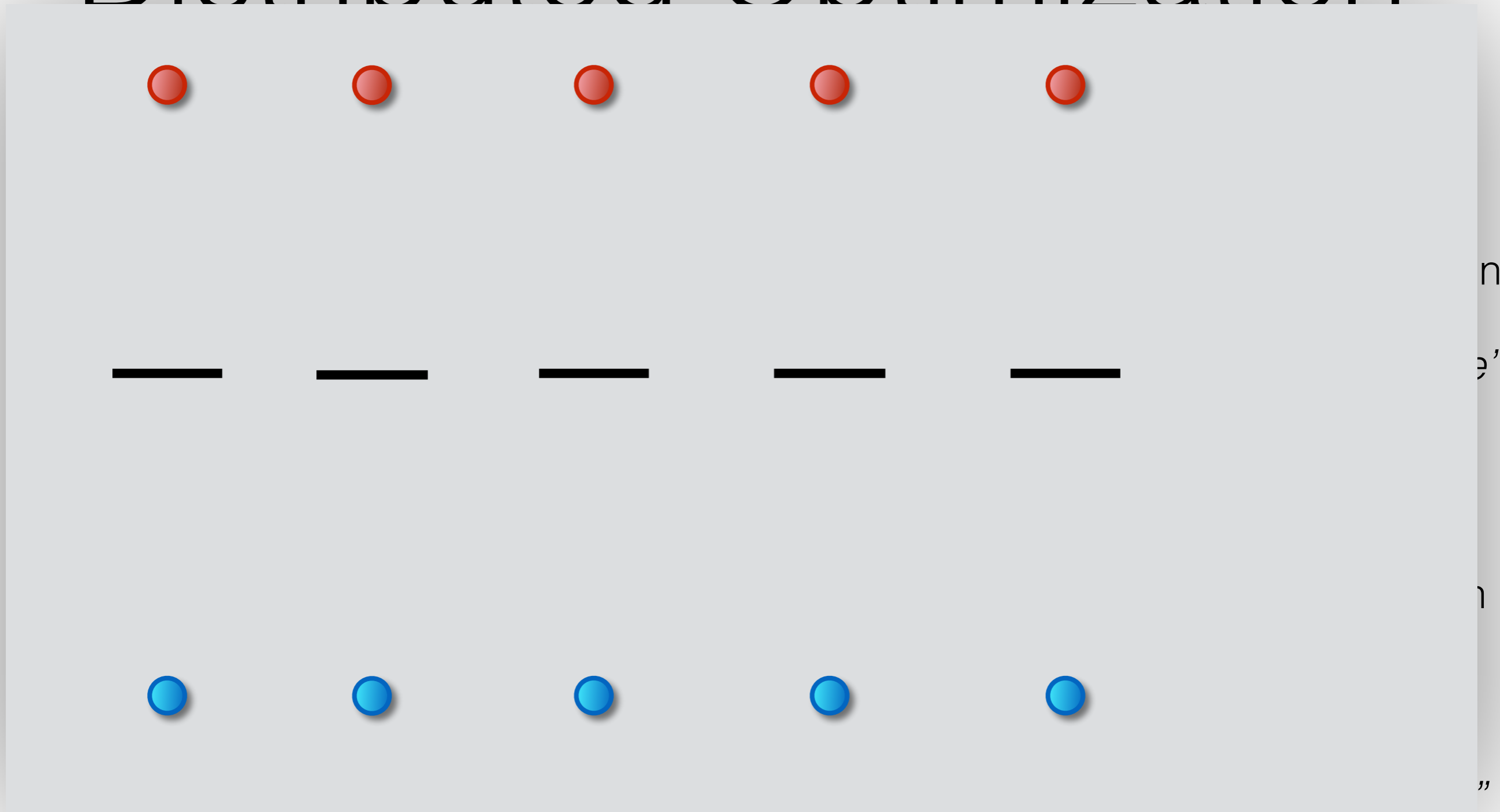
average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



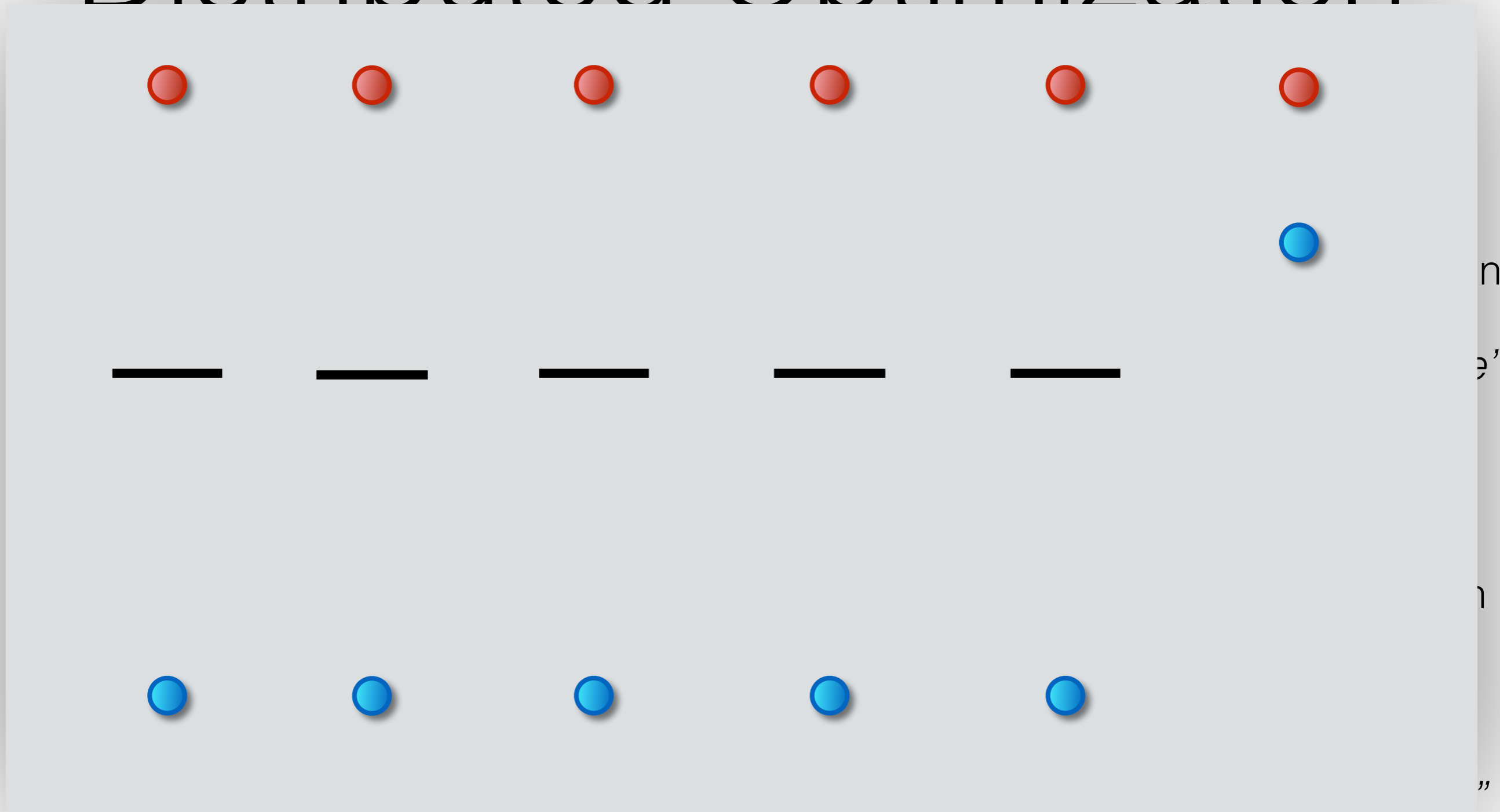
average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



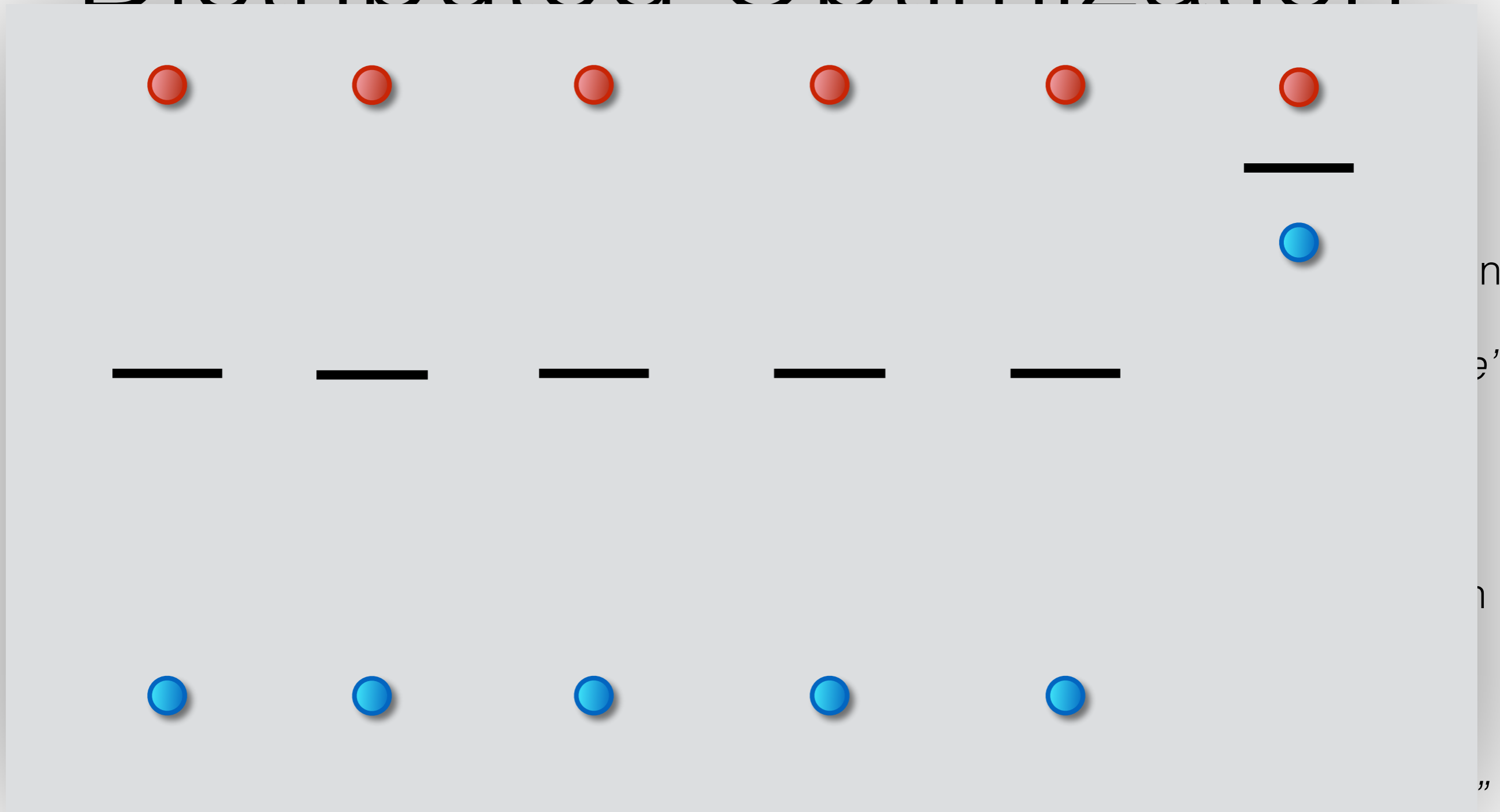
average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



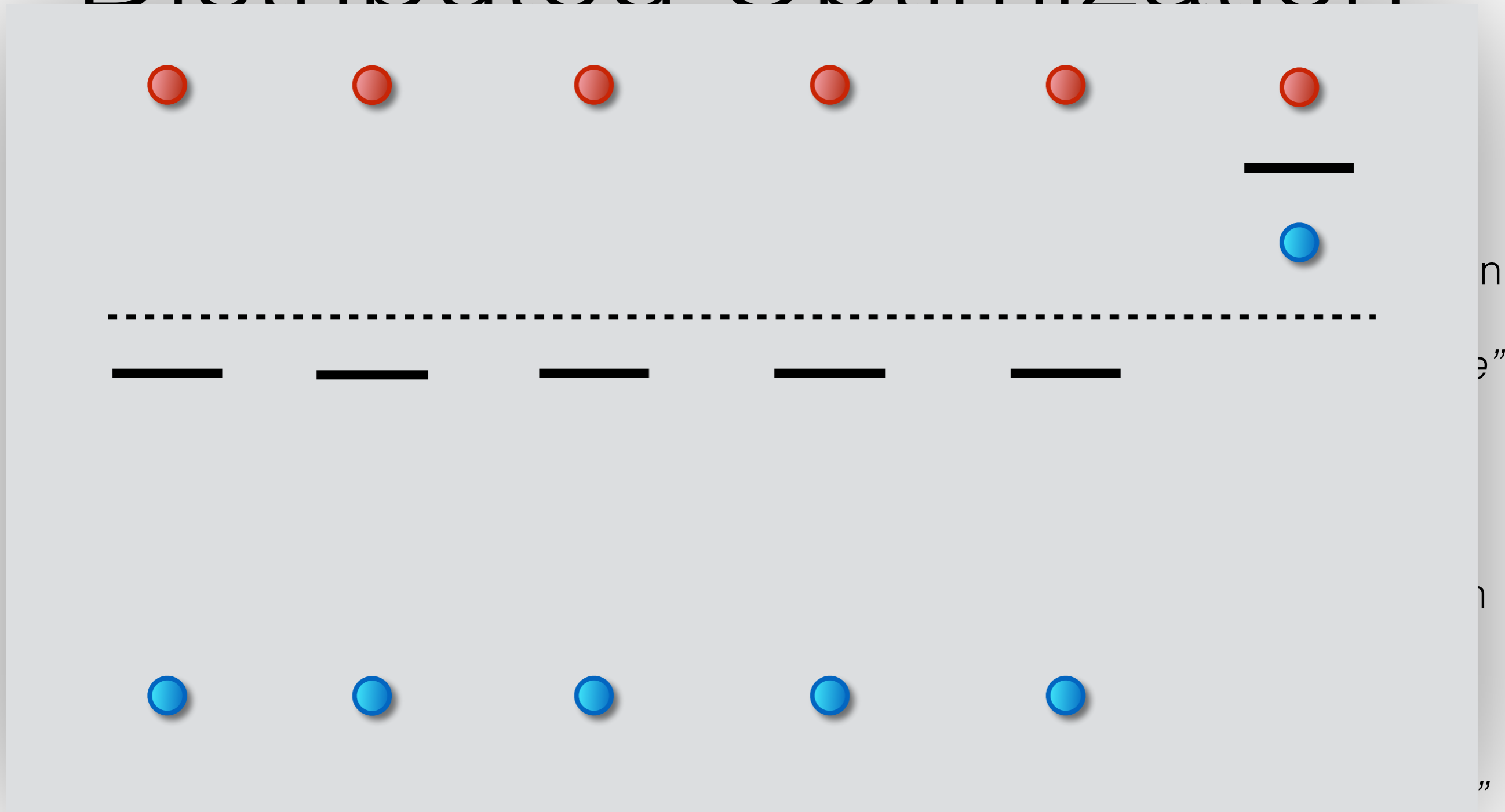
average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Distributed Optimization



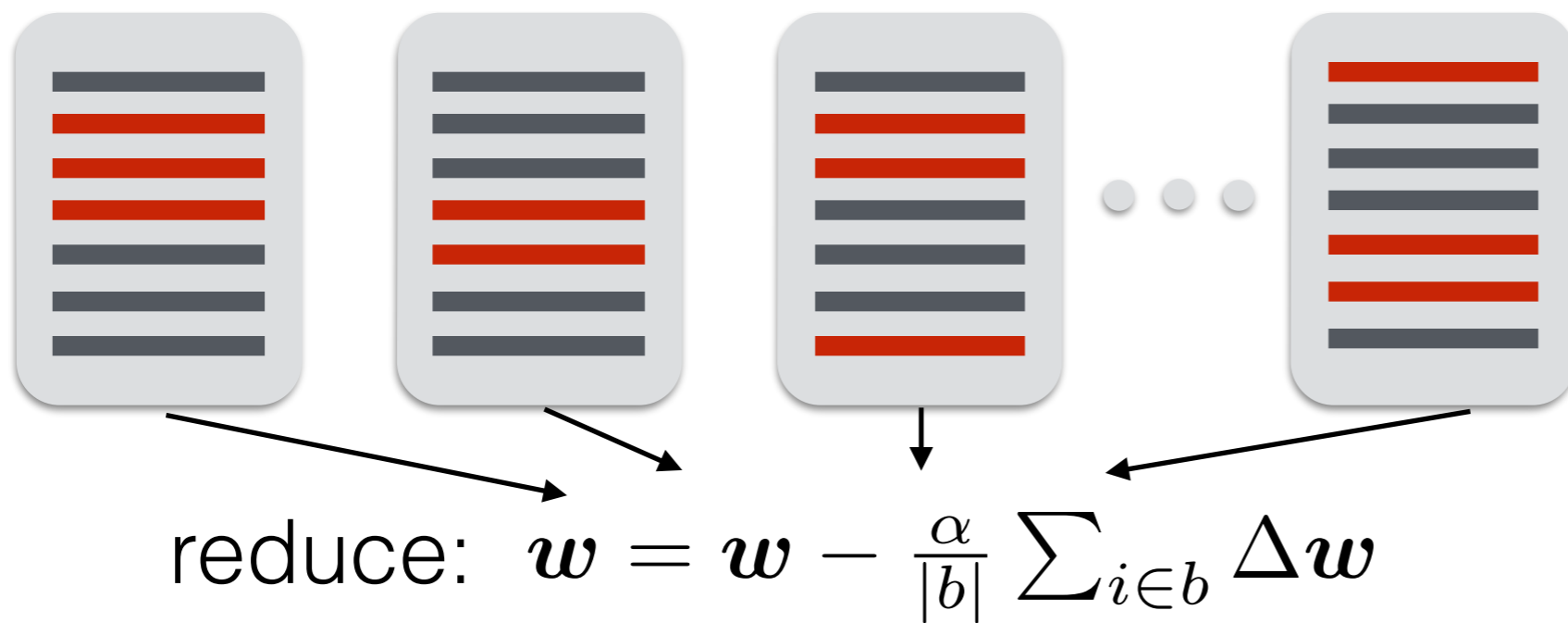
average: $\mathbf{w} := \frac{1}{K} \sum_k \mathbf{w}_k$

Mini-batch

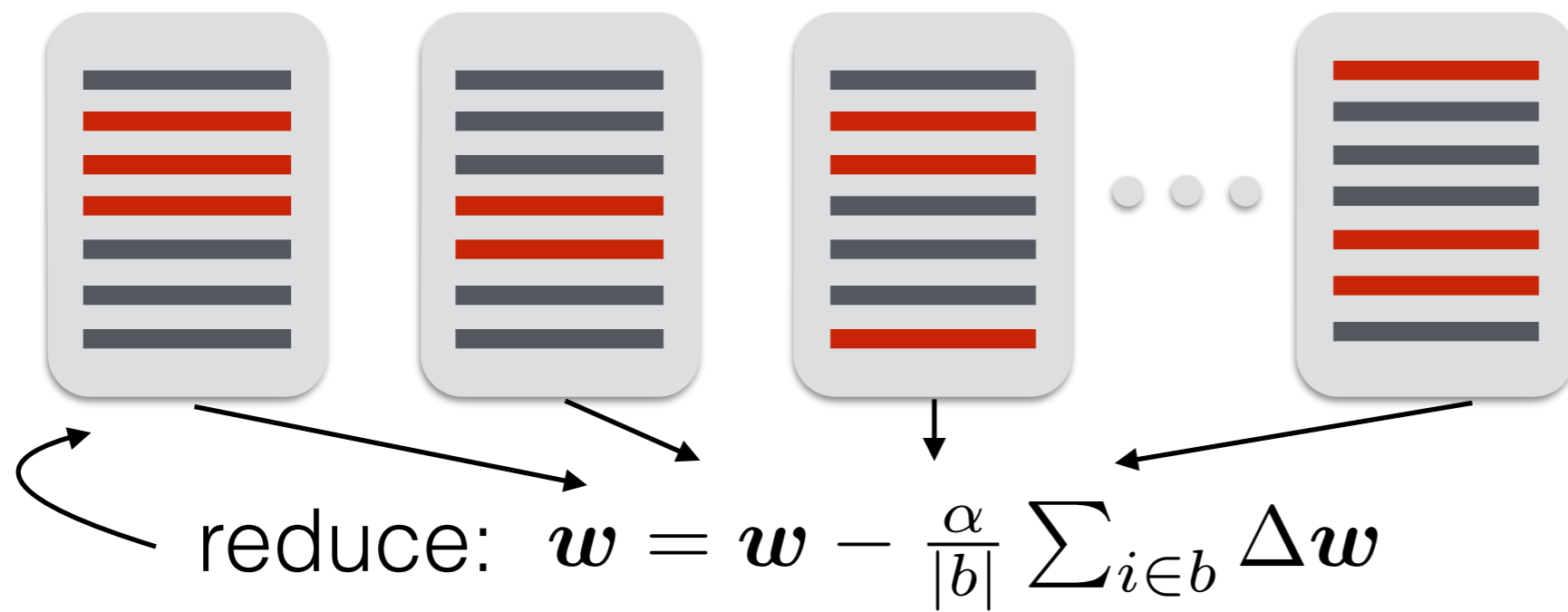
Mini-batch



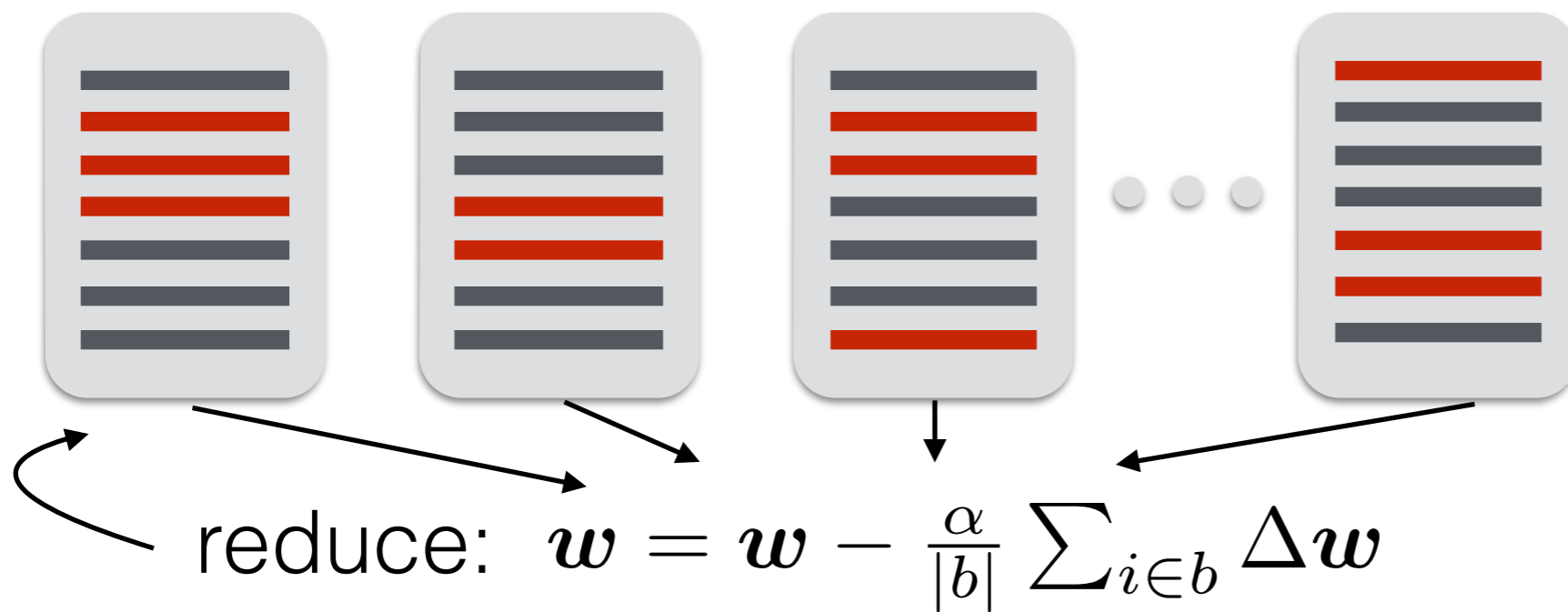
Mini-batch



Mini-batch

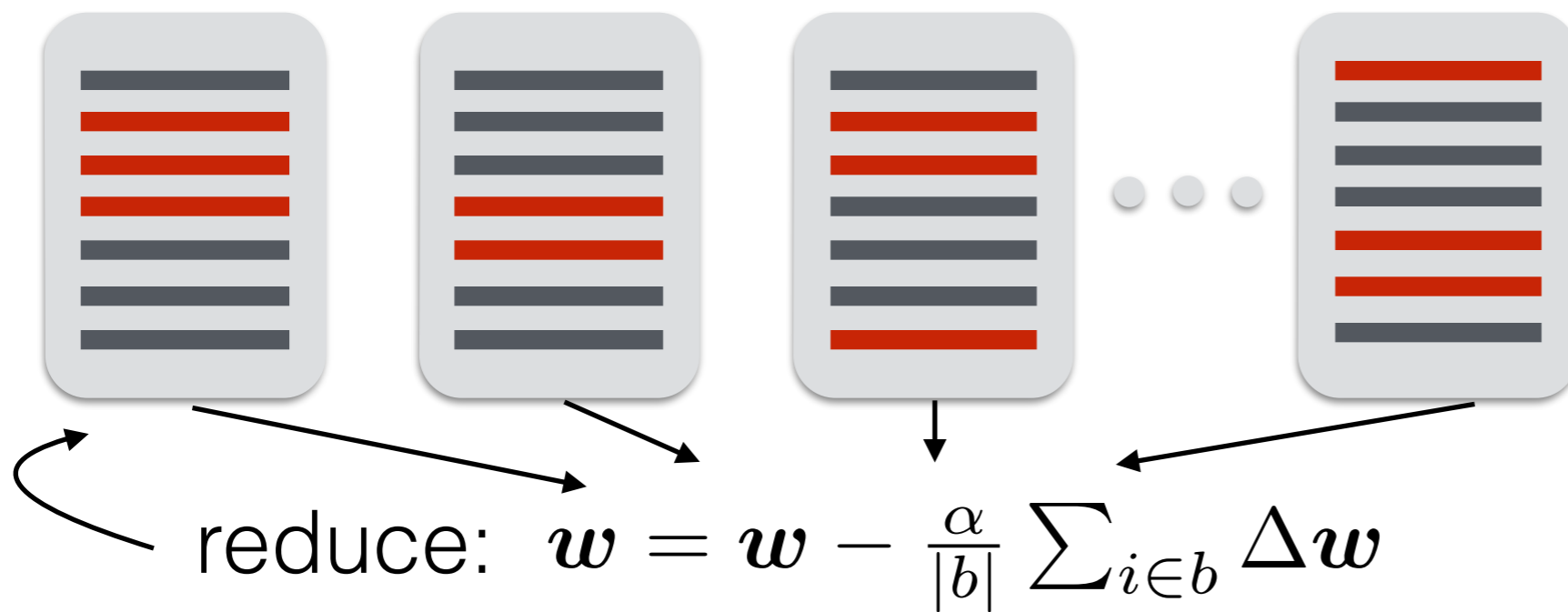


Mini-batch



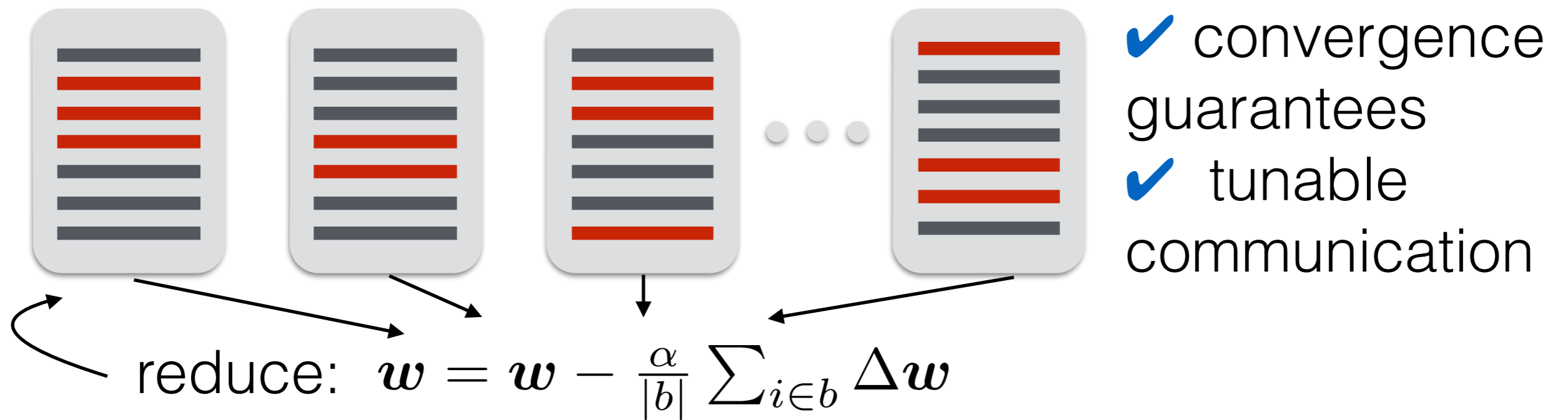
✓ convergence guarantees

Mini-batch



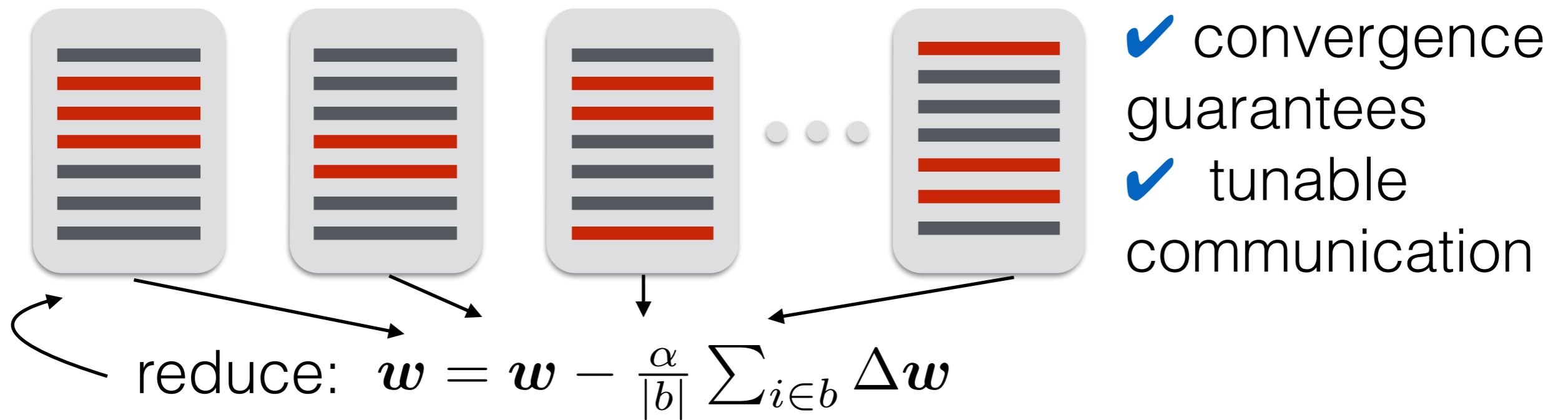
- ✓ convergence guarantees
- ✓ tunable communication

Mini-batch



► a natural middle-ground

Mini-batch



► a natural middle-ground

Mini-batch Limitations



Mini-batch Limitations

1. METHODS BEYOND SGD

Mini-batch Limitations

1. METHODS BEYOND SGD
2. STALE UPDATES

Mini-batch Limitations

1. METHODS BEYOND SGD
2. STALE UPDATES
3. AVERAGE OVER BATCH SIZE

LARGE-SCALE OPTIMIZATION

COCO

RESULTS

LARGE-SCALE OPTIMIZATION

CoCoA

RESULTS

Mini-batch Limitations

1. METHODS BEYOND SGD
2. STALE UPDATES
3. AVERAGE OVER BATCH SIZE

Mini-batch Limitations

1. ~~METHODS BEYOND SGD~~

Use Primal-Dual Framework

2. ~~STALE UPDATES~~

Immediately apply local updates

3. ~~AVERAGE OVER BATCH SIZE~~

Average over $K \ll \text{batch size}$

Communication-Efficient Distributed Dual Coordinate Ascent (CoCoA)

1. ~~METHODS BEYOND SGD~~

Use Primal-Dual Framework

2. ~~STALE UPDATES~~

Immediately apply local updates

3. ~~AVERAGE OVER BATCH SIZE~~

Average over $K \ll \text{batch size}$

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[P(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) \right]$$

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[P(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) \right]$$

$$\max_{\alpha \in \mathbb{R}^n} \left[D(\alpha) := -\|A\alpha\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i) \right]$$
$$A_i = \frac{1}{\lambda n} x_i$$

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[P(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) \right]$$

$$\max_{\alpha \in \mathbb{R}^n} \left[D(\alpha) := -\|A\alpha\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i) \right]$$
$$A_i = \frac{1}{\lambda n} x_i$$

- ▶ Stopping criteria given by duality gap

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[P(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) \right]$$

$$\max_{\alpha \in \mathbb{R}^n} \left[D(\alpha) := -\|A\alpha\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i) \right]$$
$$A_i = \frac{1}{\lambda n} x_i$$

- ▶ Stopping criteria given by duality gap
- ▶ Good performance in practice

1. Primal-Dual Framework

PRIMAL

\geq

DUAL

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left[P(\mathbf{w}) := \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i) \right]$$

$$\max_{\alpha \in \mathbb{R}^n} \left[D(\alpha) := -\|A\alpha\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i) \right]$$
$$A_i = \frac{1}{\lambda n} x_i$$

- ▶ Stopping criteria given by duality gap
- ▶ Good performance in practice
- ▶ Default in software packages e.g. liblinear

2. Immediately Apply Updates

2. Immediately Apply Updates

STALE

```
for  $i \in b$   
  |  $\Delta \mathbf{w} \leftarrow \Delta \mathbf{w} - \alpha \nabla_i P(\mathbf{w})$   
end  
 $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$ 
```

2. Immediately Apply Updates

STALE

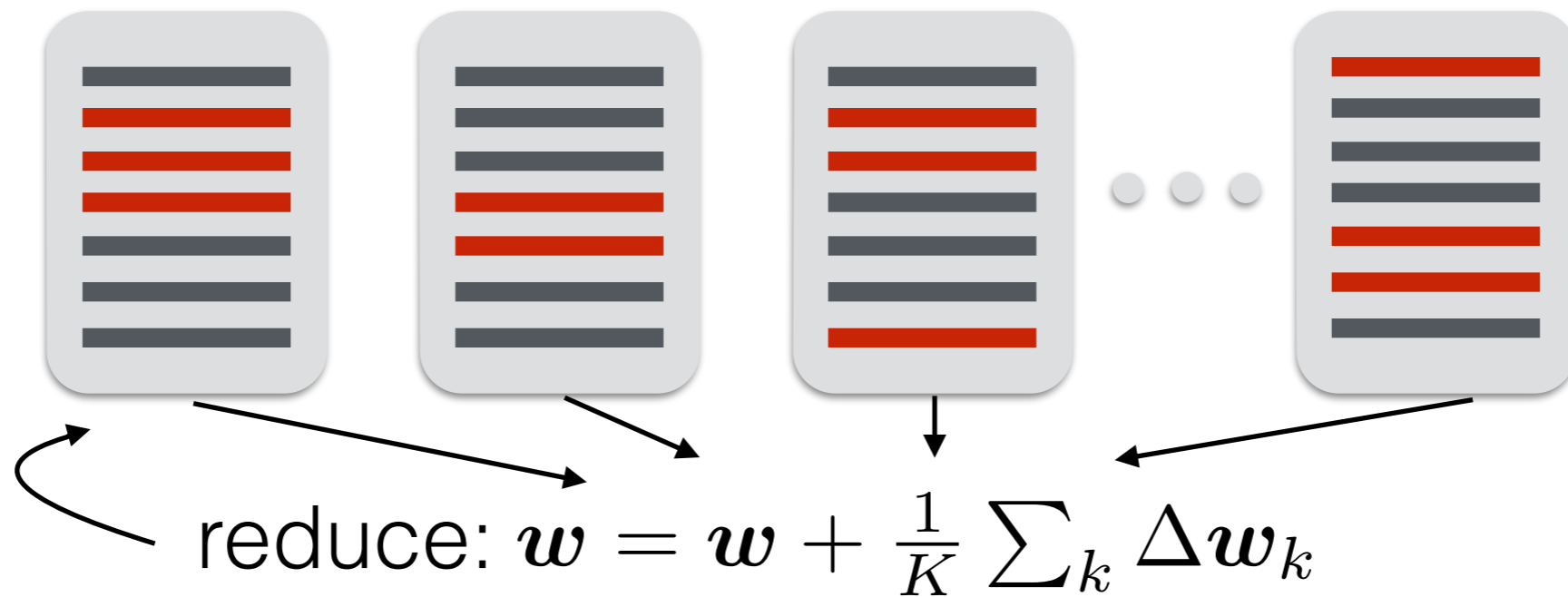
```
for  $i \in b$   
|  $\Delta \mathbf{w} \leftarrow \Delta \mathbf{w} - \alpha \nabla_i P(\mathbf{w})$   
end  
 $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$ 
```

FRESH

```
for  $i \in b$   
|  $\Delta \mathbf{w} \leftarrow \Delta \mathbf{w} - \alpha \nabla_i P(\mathbf{w})$   
|  $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$   
end
```

3. Average over K

3. Average over K



CoCoA

Algorithm 1: CoCoA

Input: $T \geq 1$, scaling parameter $1 \leq \beta_K \leq K$ (default: $\beta_K := 1$).

Data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ distributed over K machines

Initialize: $\alpha_{[k]}^{(0)} \leftarrow 0$ for all machines k , and $\mathbf{w}^{(0)} \leftarrow 0$

for $t = 1, 2, \dots, T$

for all machines $k = 1, 2, \dots, K$ *in parallel*

$(\Delta\alpha_{[k]}, \Delta\mathbf{w}_k) \leftarrow \text{LOCALDUALMETHOD}(\alpha_{[k]}^{(t-1)}, \mathbf{w}^{(t-1)})$

$\alpha_{[k]}^{(t)} \leftarrow \alpha_{[k]}^{(t-1)} + \frac{\beta_K}{K} \Delta\alpha_{[k]}$

end

reduce $\mathbf{w}^{(t)} \leftarrow \mathbf{w}^{(t-1)} + \frac{\beta_K}{K} \sum_{k=1}^K \Delta\mathbf{w}_k$

end

Procedure A: LOCALDUALMETHOD: Dual algorithm on machine k

Input: Local $\alpha_{[k]} \in \mathbb{R}^{n_k}$, and $\mathbf{w} \in \mathbb{R}^d$ consistent with other coordinate blocks of α s.t. $\mathbf{w} = A\alpha$

Data: Local $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$

Output: $\Delta\alpha_{[k]}$ and $\Delta\mathbf{w} := A_{[k]}\Delta\alpha_{[k]}$

CoCoA

✓ <10 lines of code in *Spark*

BLOCKS OF α S.T. $w = A\alpha$

Data: Local $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$

Output: $\Delta\alpha_{[k]}$ and $\Delta w := A_{[k]}\Delta\alpha_{[k]}$

CoCoA

- ✓ <10 lines of code in *Spark*
- ✓ primal-dual framework allows for *any* internal optimization method

BLOCKS OF α S.T. $w = A\alpha$

Data: Local $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$

Output: $\Delta\alpha_{[k]}$ and $\Delta w := A_{[k]}\Delta\alpha_{[k]}$

CoCoA

- ✓ <10 lines of code in *Spark*
- ✓ primal-dual framework allows for *any* internal optimization method
- ✓ local updates applied *immediately*

BLOCKS OF α S.T. $w = A\alpha$

Data: Local $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$

Output: $\Delta\alpha_{[k]}$ and $\Delta w := A_{[k]}\Delta\alpha_{[k]}$

CoCoA

- ✓ <10 lines of code in *Spark*
- ✓ primal-dual framework allows for *any* internal optimization method
- ✓ local updates applied *immediately*
- ✓ average over K

blocks of α s.t. $w = A\alpha$

Data: Local $\{(\mathbf{x}_i, y_i)\}_{i=1}^{n_k}$

Output: $\Delta\alpha_{[k]}$ and $\Delta w := A_{[k]}\Delta\alpha_{[k]}$

Convergence

Convergence

Assumptions:

- ▶ l_i are $1/\gamma$ -smooth
- ▶ LocalDualMethod makes improvement Θ per step

e.g. for SDCA $\Theta = \left(1 - \frac{\lambda n \gamma}{1 + \lambda n \gamma \tilde{n}}\right)^H$

Convergence

Assumptions:

- ▶ l_i are $1/\gamma$ -smooth
- ▶ LocalDualMethod makes improvement Θ per step

e.g. for SDCA $\Theta = \left(1 - \frac{\lambda n \gamma}{1 + \lambda n \gamma} \frac{1}{\tilde{n}}\right)^H$

$$E[D(\boldsymbol{\alpha}^*) - D(\boldsymbol{\alpha}^{(T)})] \leq \left(1 - (1 - \Theta) \frac{1}{K} \frac{\lambda n \gamma}{\sigma + \lambda n \gamma}\right)^T \left(D(\boldsymbol{\alpha}^*) - D(\boldsymbol{\alpha}^{(0)})\right)$$

Convergence

Assumptions:

- ▶ l_i are $1/\gamma$ -smooth
- ▶ LocalDualMethod makes improvement Θ per step

e.g. for SDCA $\Theta = \left(1 - \frac{\lambda n \gamma}{1 + \lambda n \gamma} \frac{1}{\tilde{n}}\right)^H$

$$E[D(\boldsymbol{\alpha}^*) - D(\boldsymbol{\alpha}^{(T)})] \leq \left(1 - (1 - \Theta) \frac{1}{K} \frac{\lambda n \gamma}{\sigma + \lambda n \gamma}\right)^T \left(D(\boldsymbol{\alpha}^*) - D(\boldsymbol{\alpha}^{(0)})\right)$$

- ▶ applies also to duality gap
 $0 \leq \sigma \leq n/K$
- ▶ measure of difficulty of data partition

Convergence

Assumptions:

✓ it converges!

ϵ



$E[D(\alpha^*)$

$(\alpha^{(0)})$

▶ measure of difficulty of data partition

Convergence

Assumptions:

- ▶  it converges!
- ▶  inherits convergence rate of locally used method

$E[D(\alpha^*)$




ep

$(\alpha^{(0)})$

▶ measure of difficulty of data partition

Convergence

Assumptions:

- ▶  it converges!
- ▶  inherits convergence rate of locally used method
- ▶  convergence rate is linear for smooth losses

$E[D(\alpha^*)$

ep

$(\alpha^{(0)})$

▶ measure of difficulty of data partition

LARGE-SCALE OPTIMIZATION

CoCoA

RESULTS!

LARGE-SCALE OPTIMIZATION

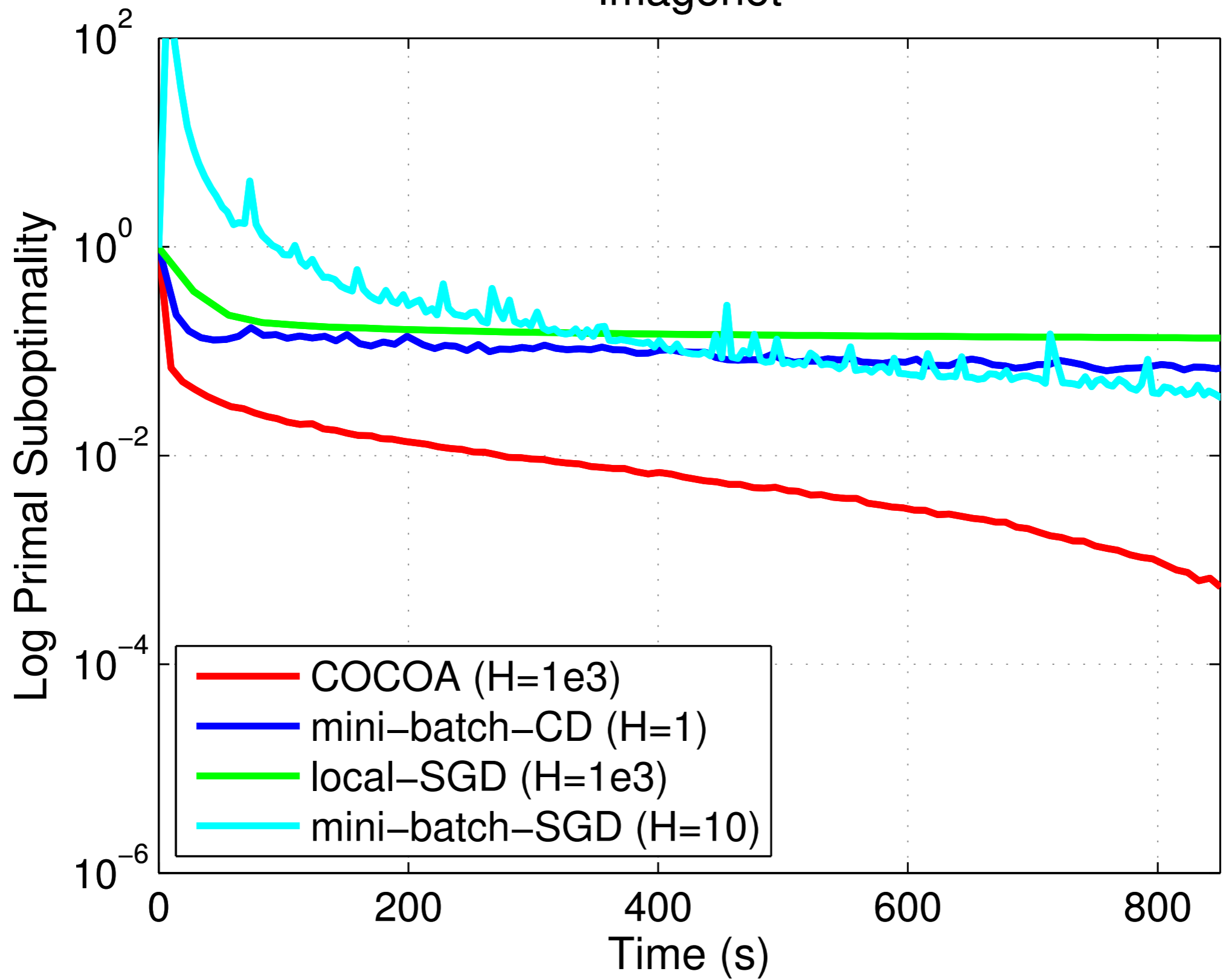
CoCoA

RESULTS!

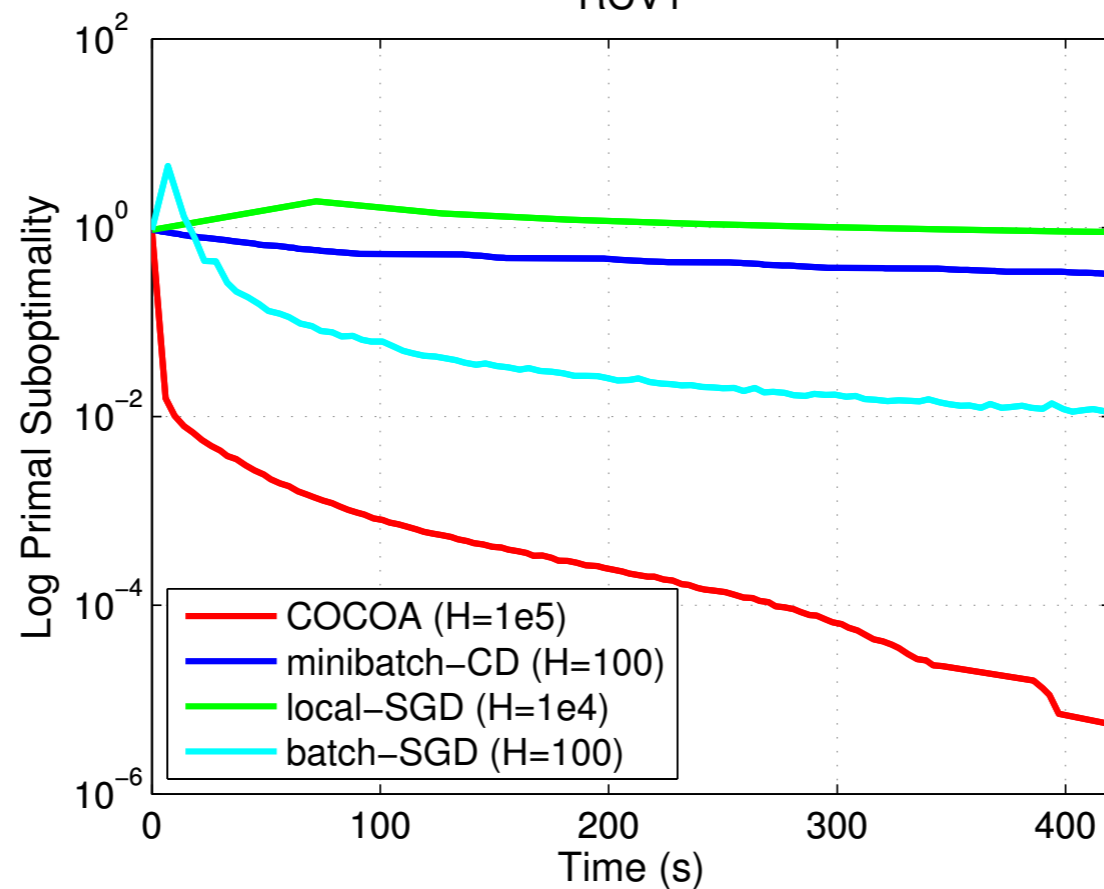
Empirical Results in *Spark*

Dataset	Training (n)	Features (d)	Sparsity	Workers (K)
Cov	522,911	54	22.22%	4
Rcv1	677,399	47,236	0.16%	8
Imagenet	32,751	160,000	100%	32

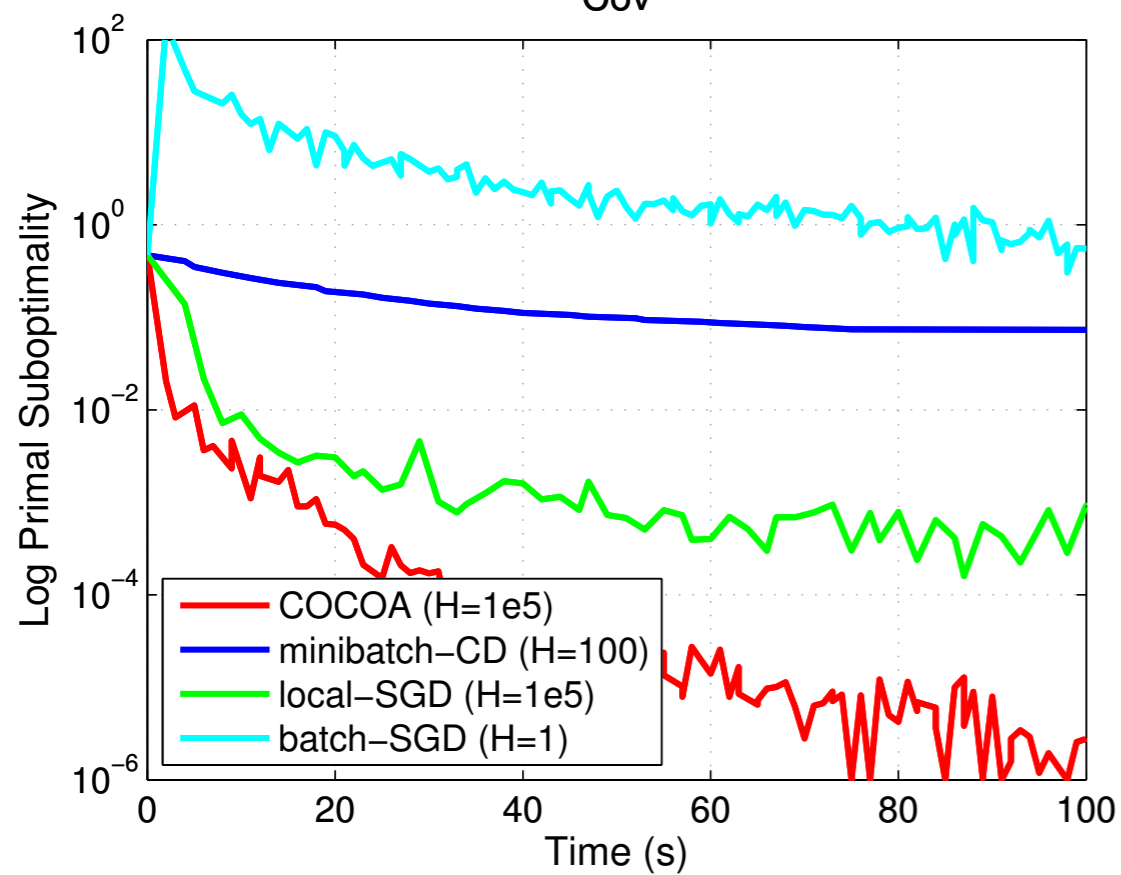
Imagenet



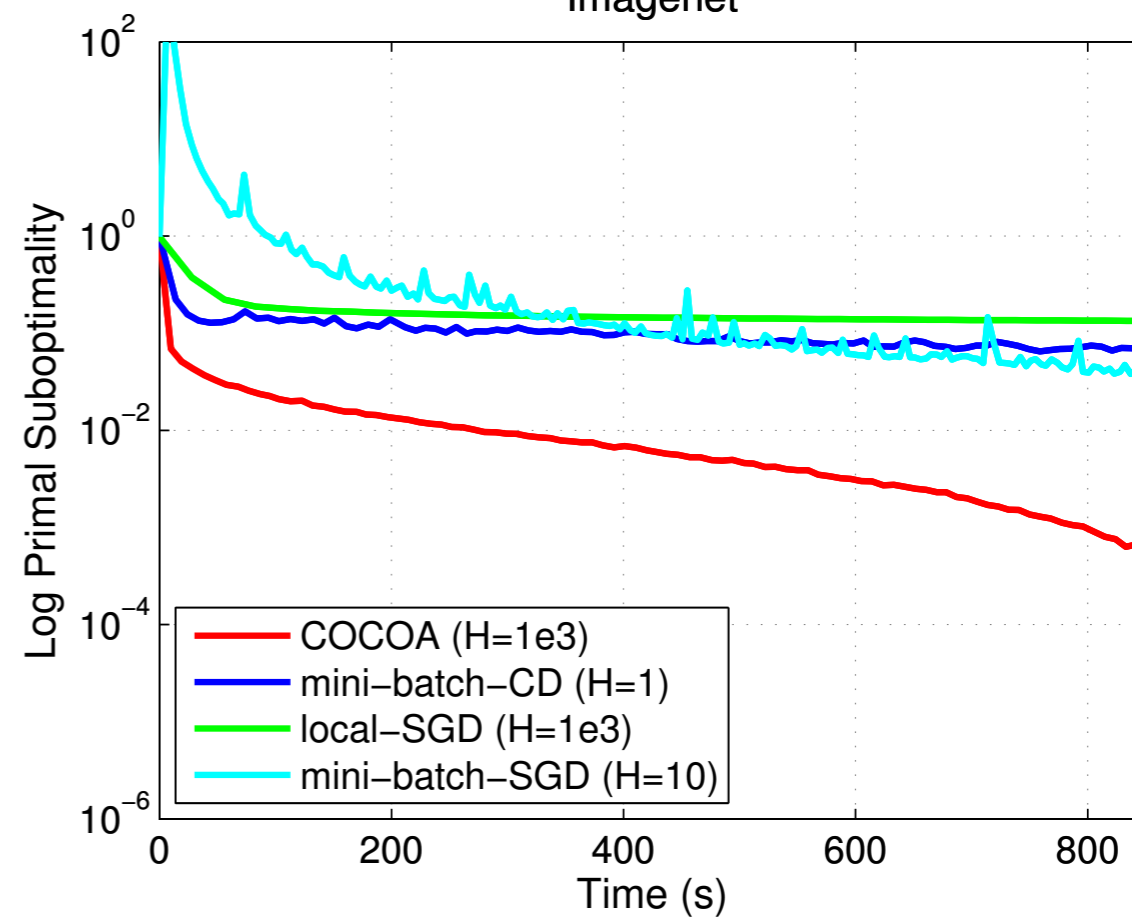
RCV1



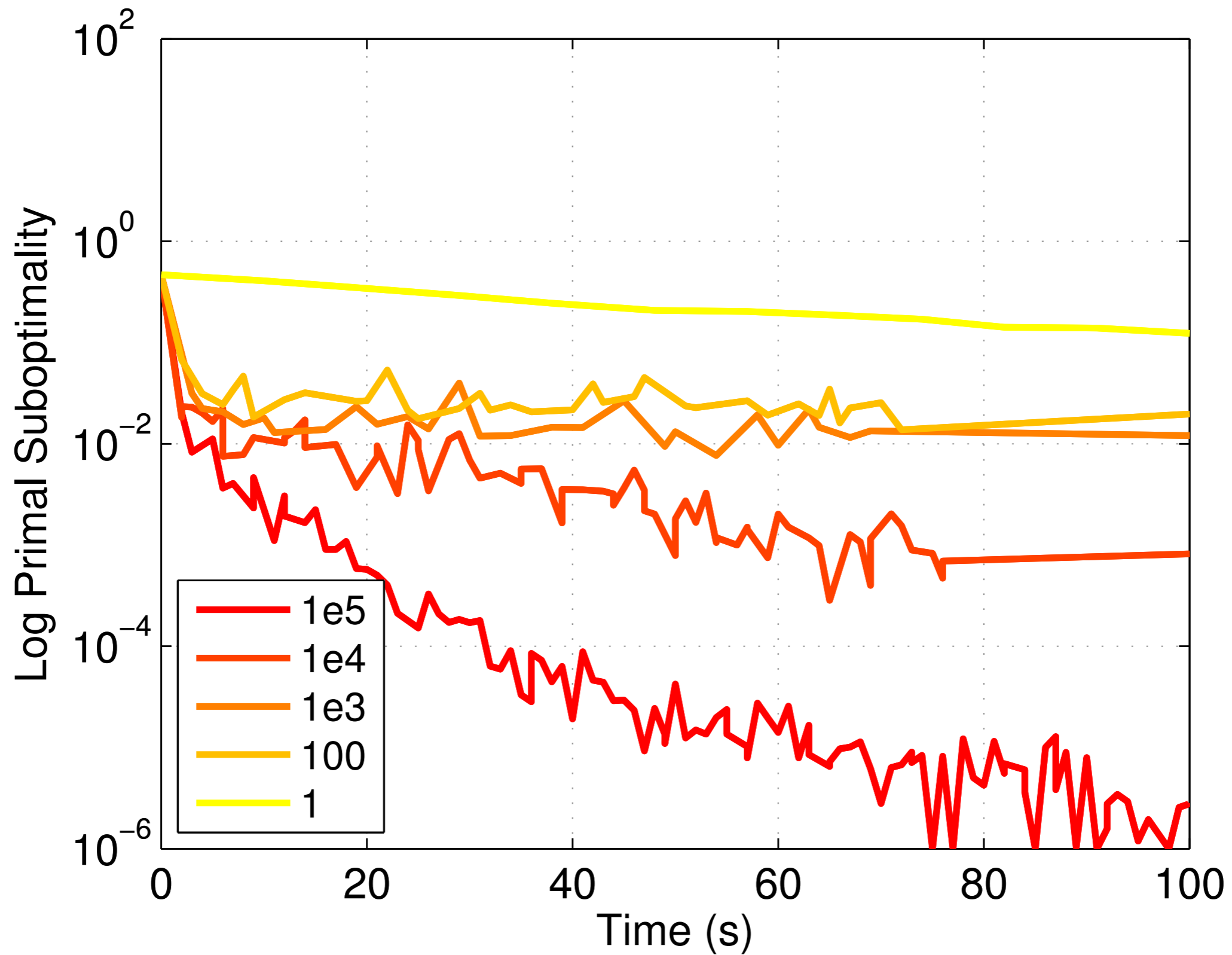
Cov



Imagenet



Effect of H on CoCoA



COCOA Take-Aways

COCOA Take-Aways

- ▶ A framework for distributed optimization

COCO A Take-Aways

- ▶ A framework for distributed optimization
- ▶ Uses state-of-the-art primal-dual methods

COCO A Take-Aways

- ▶ A framework for distributed optimization
- ▶ Uses state-of-the-art primal-dual methods
- ▶ Reduces communication:
 - applies updates immediately
 - averages over # of machines

COCO A Take-Aways

- ▶ A framework for distributed optimization
- ▶ Uses state-of-the-art primal-dual methods
- ▶ Reduces communication:
 - applies updates immediately
 - averages over # of machines
- ▶ Strong convergence guarantees & results in practice

COCO A Take-Aways

- ▶ A framework for **distributed optimization**
- ▶ Uses **state-of-the-art** primal-dual methods
- ▶ Reduces **communication**:
 - applies updates immediately
 - averages over # of machines
- ▶ Strong convergence guarantees & results in practice

COCO A Take-Aways

- ▶ A framework for **distributed optimization**
- ▶ Uses **state-of-the-art** primal-dual methods
- ▶ Reduces **communication**:
 - applies updates immediately
 - averages over # of machines
- ▶ Strong convergence guarantees & results in practice

Future Work

Future Work

- ▶ optimal scaling between adding & averaging

Future Work

- ▶ optimal scaling between adding & averaging
- ▶ similar rates for local-SGD?

Future Work

- ▶ optimal scaling between adding & averaging
- ▶ similar rates for local-SGD?
- ▶ additional experiments: models/datasets

Future Work

- ▶ optimal scaling between adding & averaging
- ▶ similar rates for local-SGD?
- ▶ additional experiments: models/datasets
- ▶ integration into *Spark MLlib*

Thanks!

github.com/gingsmith/cocoa