

OPTIMIZATION AND BIG DATA



Peter Richtárik

Parallel coordinate descent methods

ALAN COCHRANE TAXIS
38 HARRYSMUIR GARDENS
THANK YOU

TERMINAL ID: ****3577
MERCHANT ID: *****21661

MASTERCARD
*****6875 ICC
PAN SEQ NO: 12
AID: A0000000041010

SALE
AMOUNT £100009.96

*** CUSTOMER COPY ***

TRANSACTION VOID

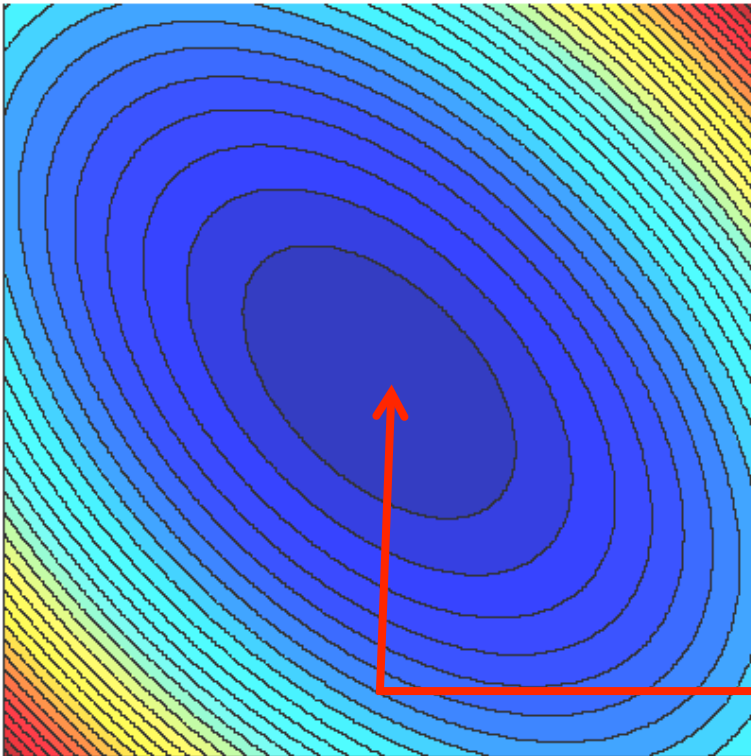
DATE: 05/03/13 TIME: 16:00

Randomized Coordinate Descent in 2D

2D Optimization

Contours of a function

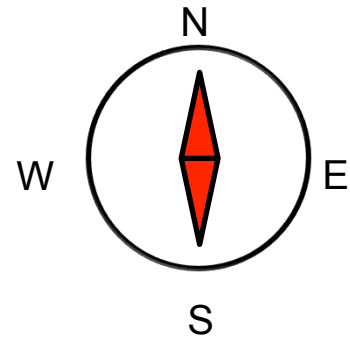
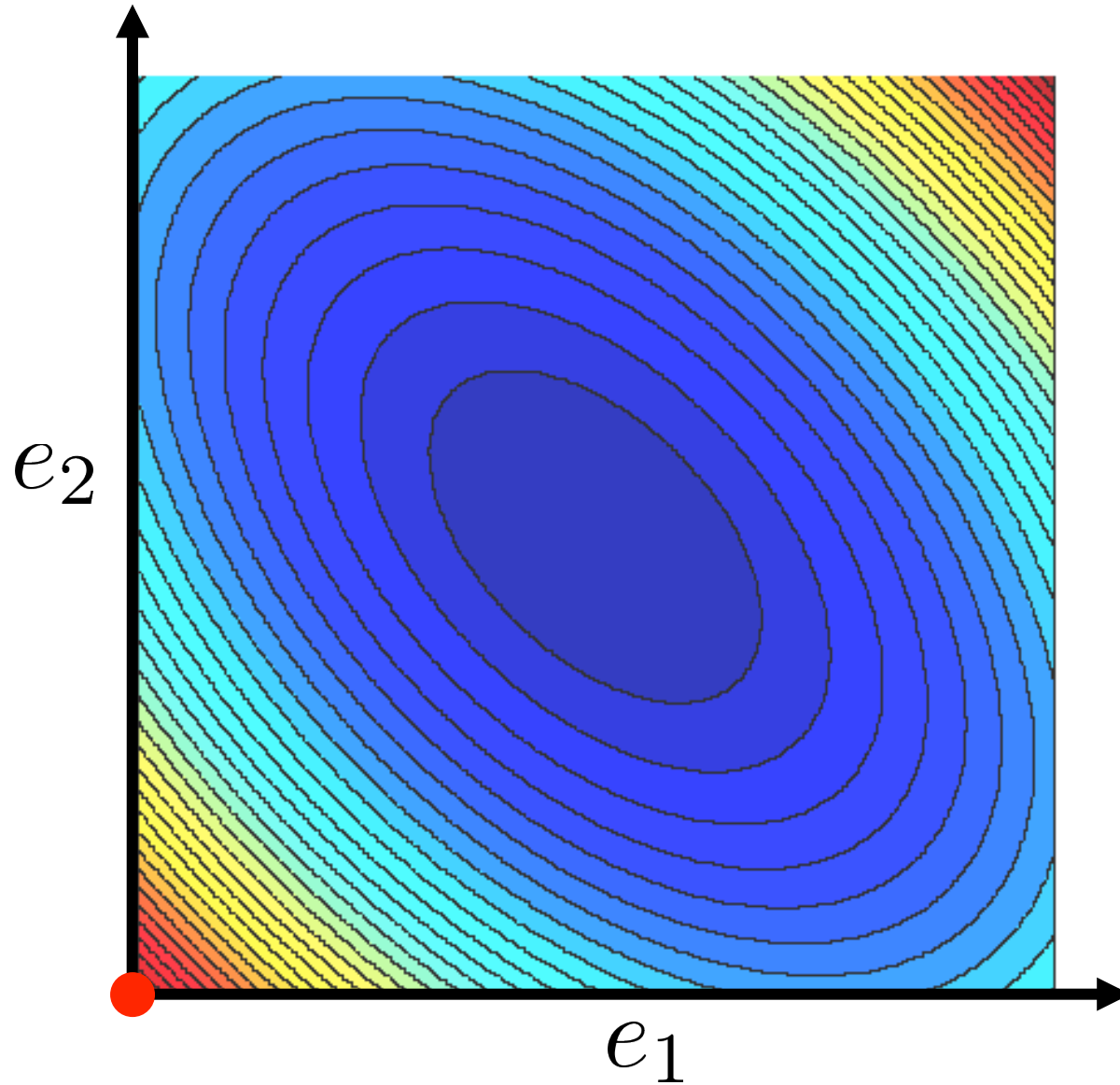
$$F : \mathbb{R}^2 \rightarrow \mathbb{R}$$



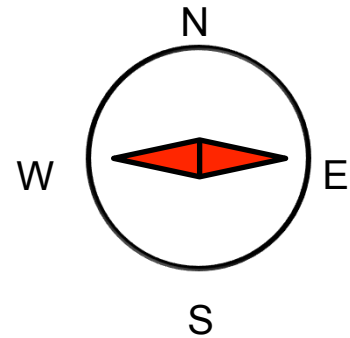
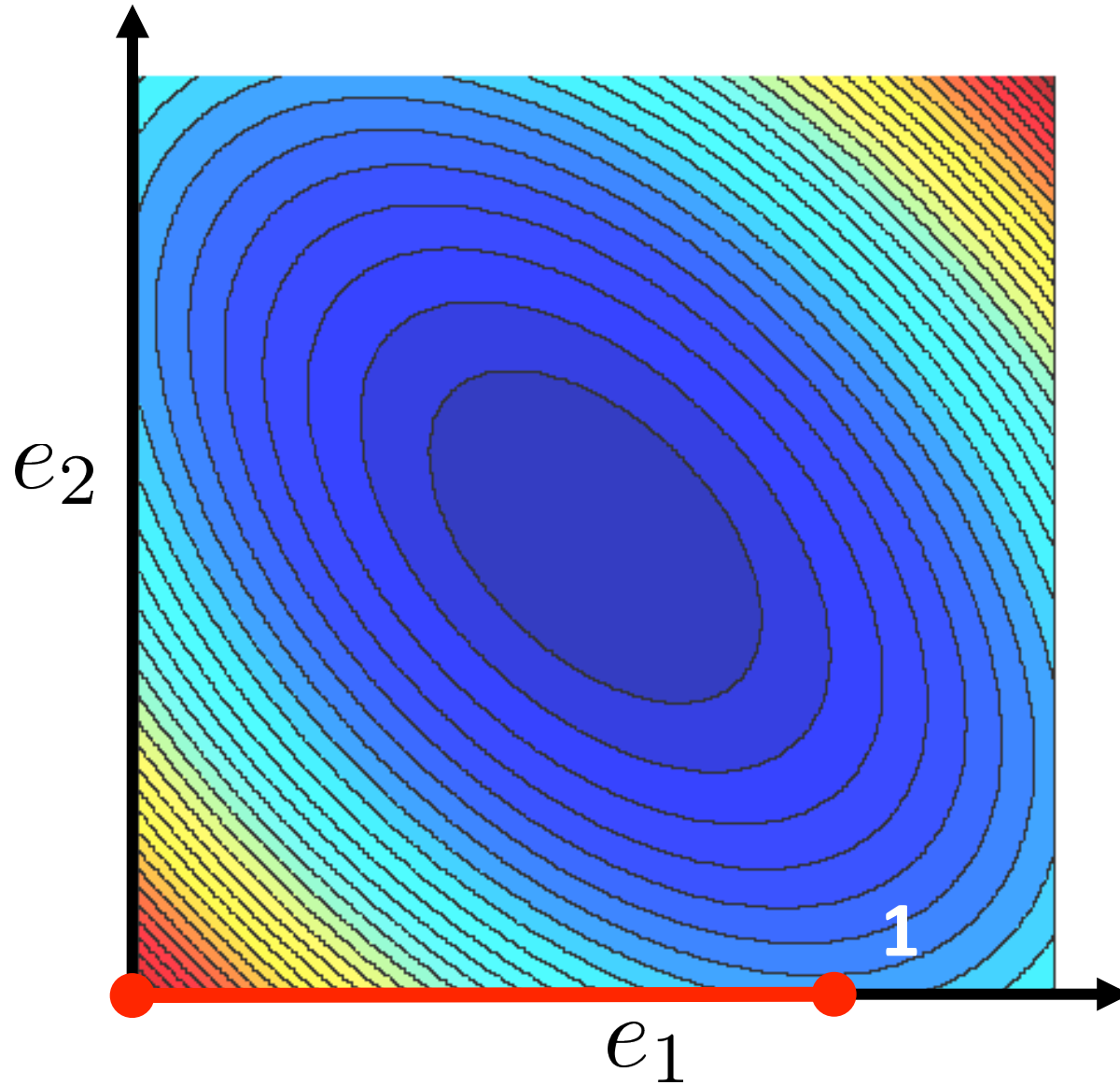
Goal:

Find the **minimizer** of F

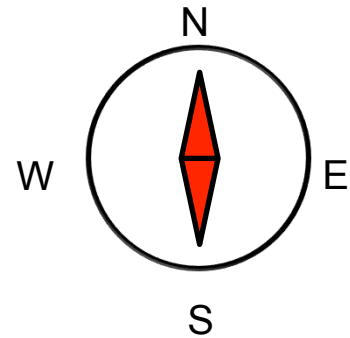
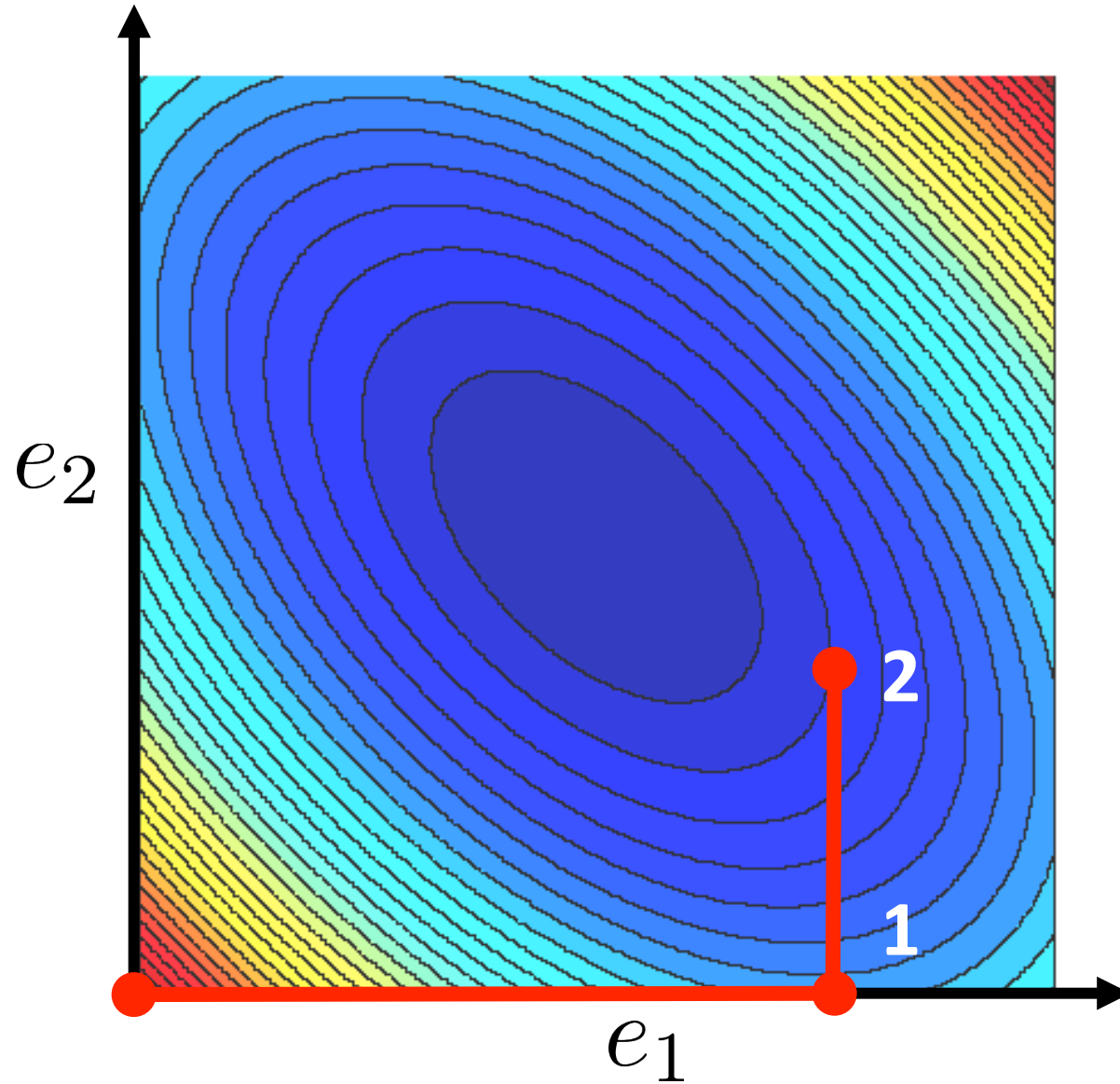
Randomized Coordinate Descent in 2D



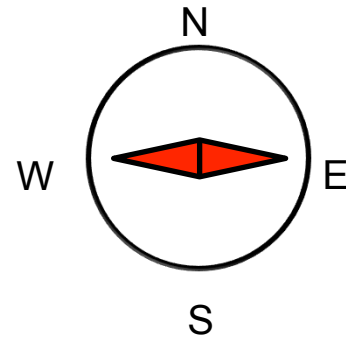
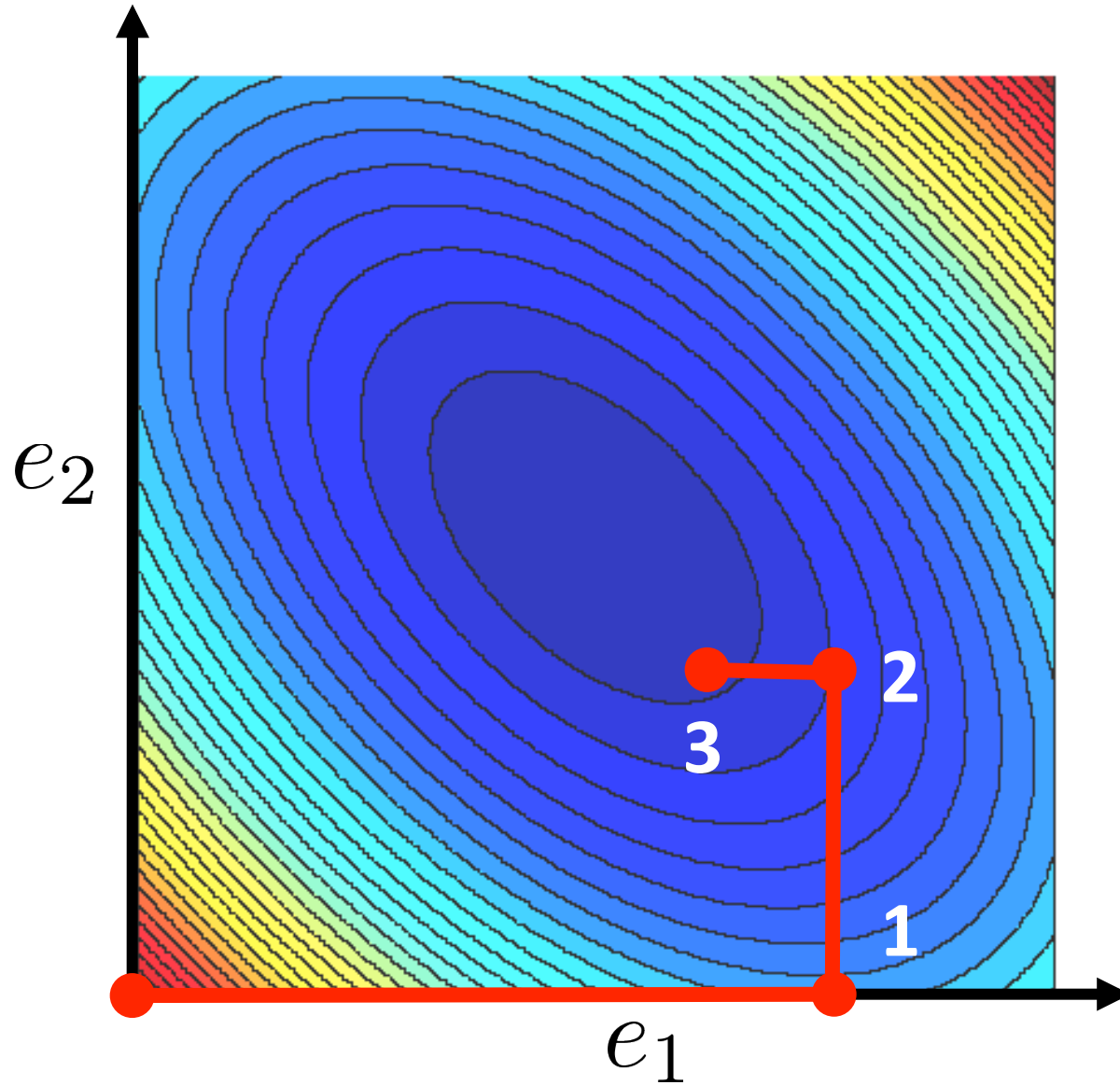
Randomized Coordinate Descent in 2D



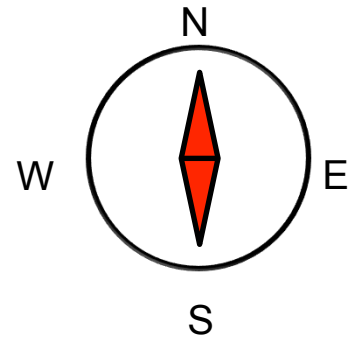
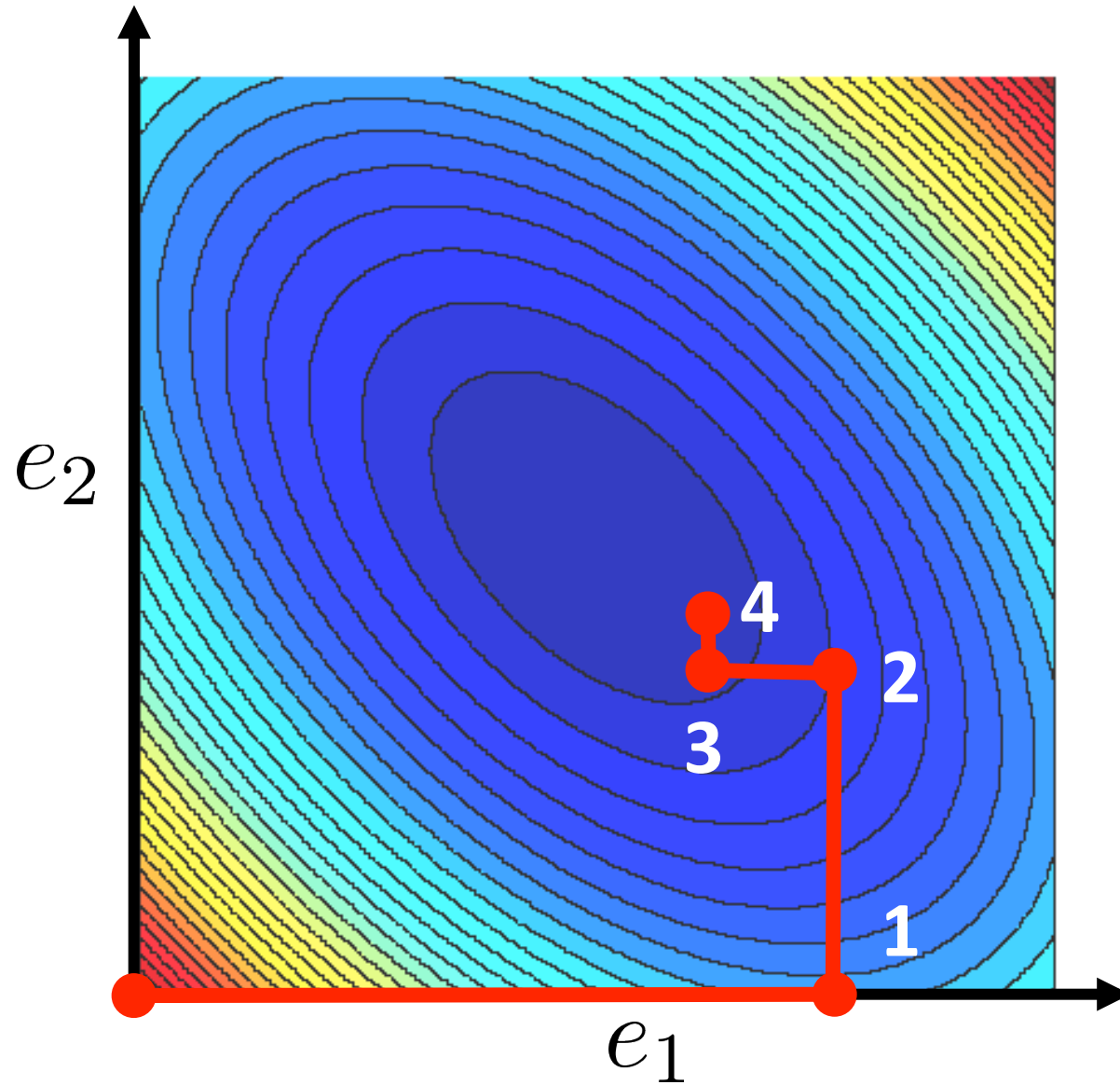
Randomized Coordinate Descent in 2D



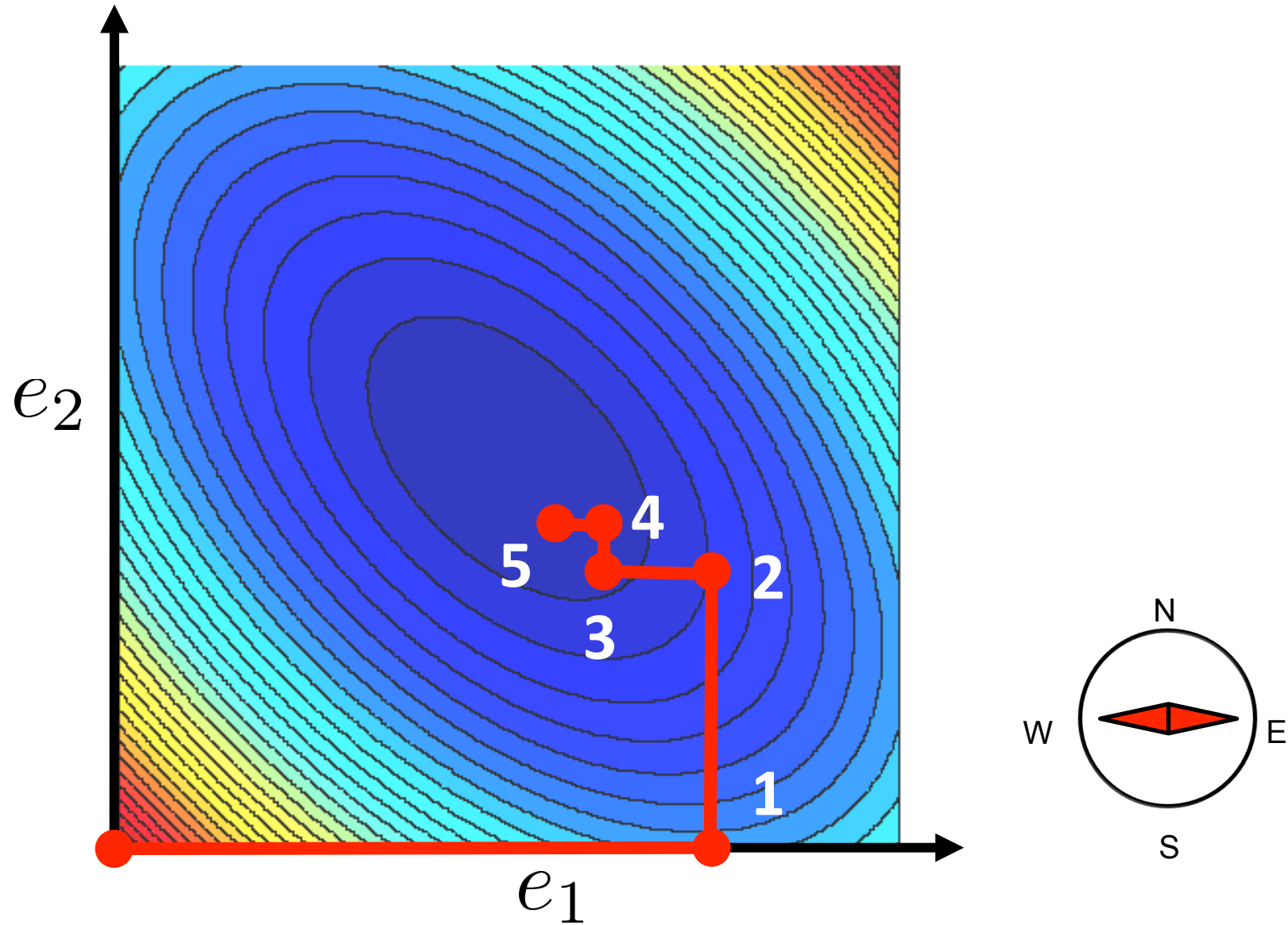
Randomized Coordinate Descent in 2D



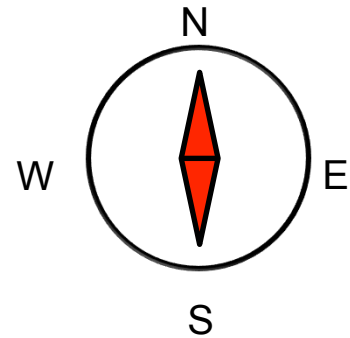
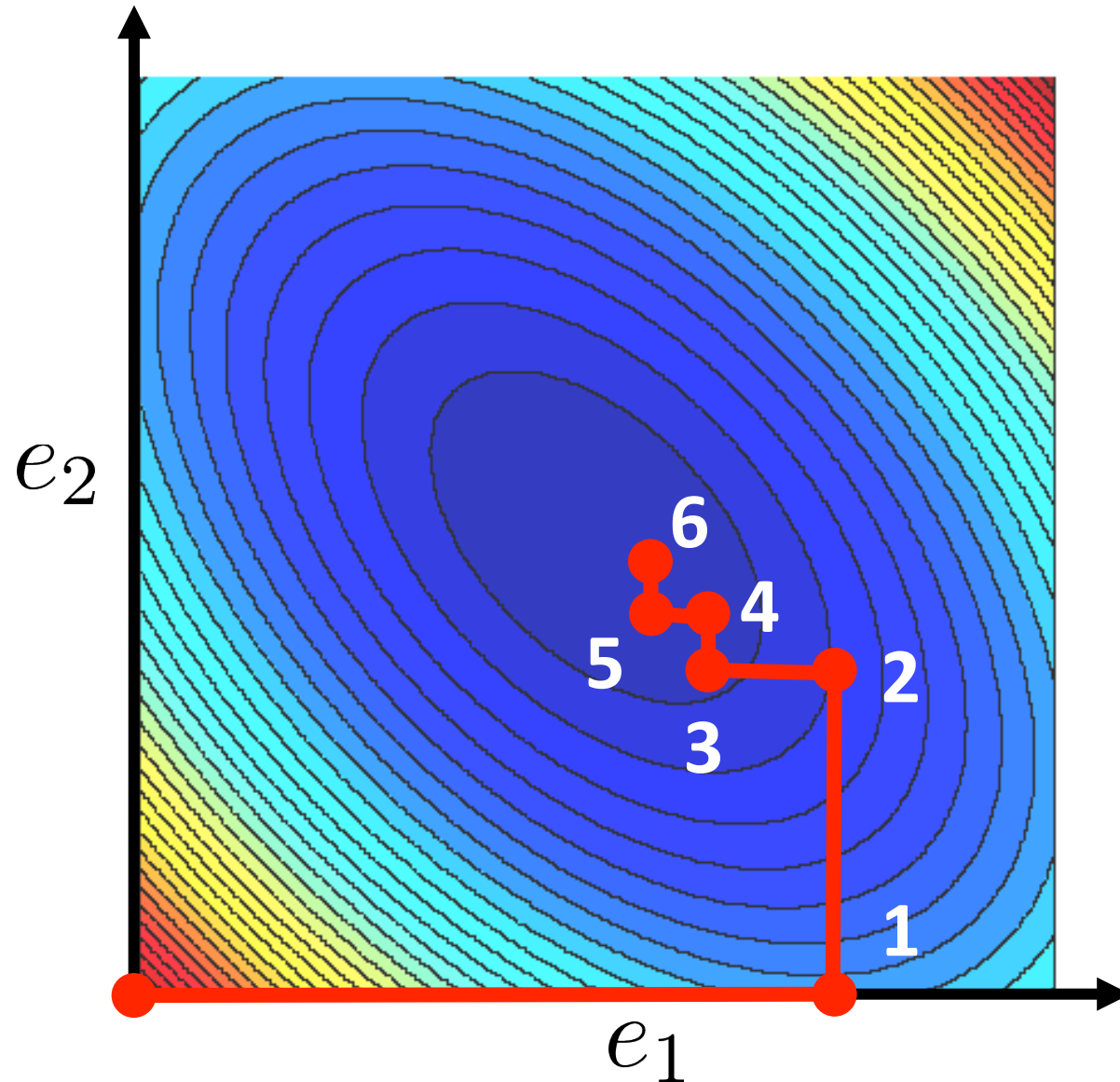
Randomized Coordinate Descent in 2D



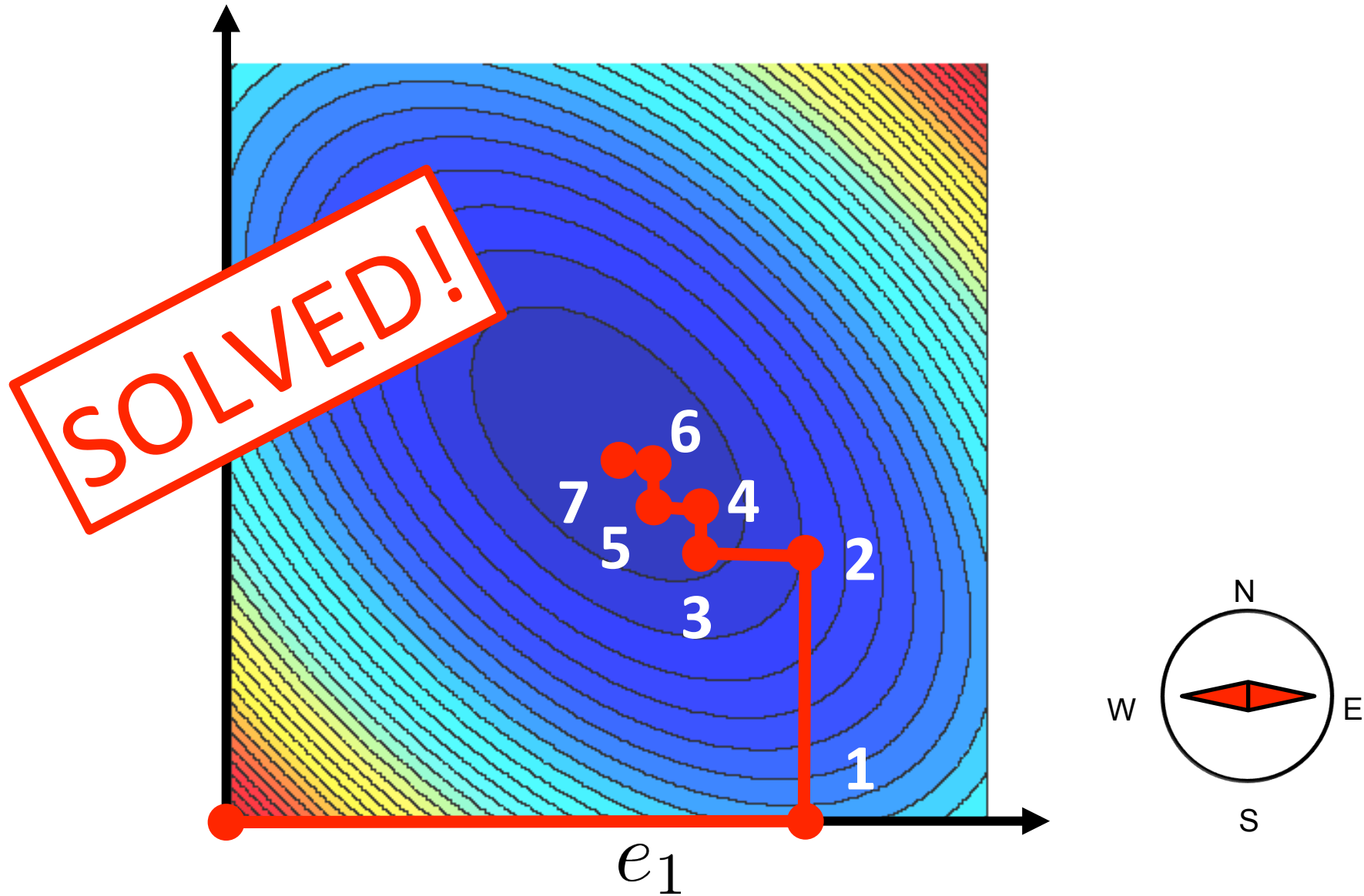
Randomized Coordinate Descent in 2D



Randomized Coordinate Descent in 2D



Randomized Coordinate Descent in 2D



Convergence of Randomized Coordinate Descent

In \mathbb{R}^n , randomized coordinate descent with uniform probabilities needs

$$O(n \times \xi(\epsilon)) \text{ iterations}$$

Strongly convex F

$$\xi(\epsilon) = \frac{1}{\epsilon}$$

Smooth or 'simple' F

Focus on n
(big data = big n)

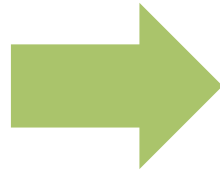
'difficult' nonsmooth F

Parallelization Dream

Serial

(1 coordinate per iteration)

$$O(n \times \xi(\epsilon))$$



Parallel

(τ coordinates per iteration)

$$O\left(\frac{n}{\tau} \times \xi(\epsilon)\right)$$

What do we actually get?

$$O\left(\frac{n\beta}{\tau} \times \xi(\epsilon)\right)$$

WANT

$$\beta = O(1)$$

Depends on to what extent we can **add up** individual updates, which depends on the **properties of F** and the **way coordinates are chosen** at each iteration

How (not) to Parallelize Coordinate Descent

“Naive” parallelization

Do the same thing as before, but

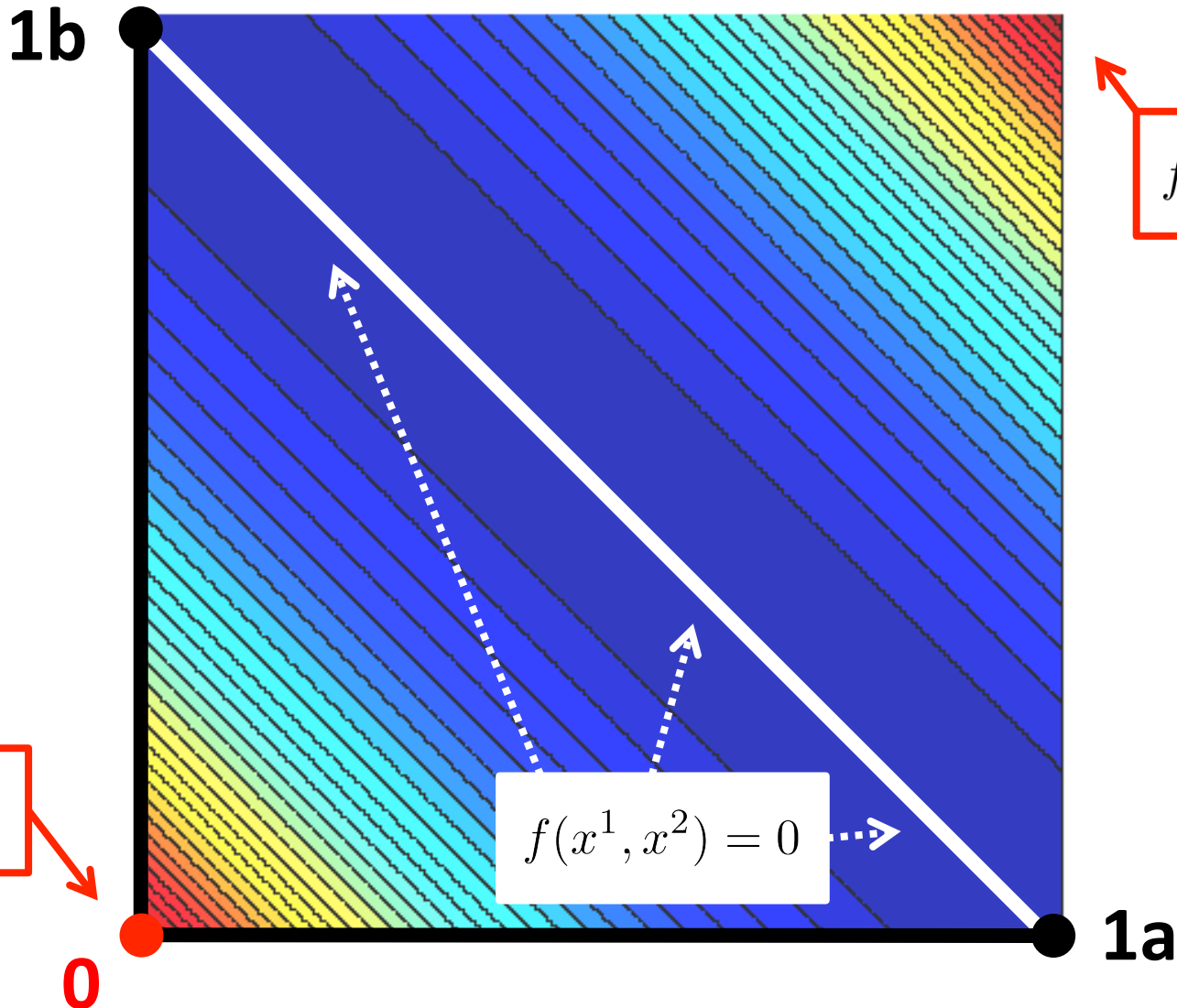
for **MORE or ALL** coordinates

&

ADD UP the updates

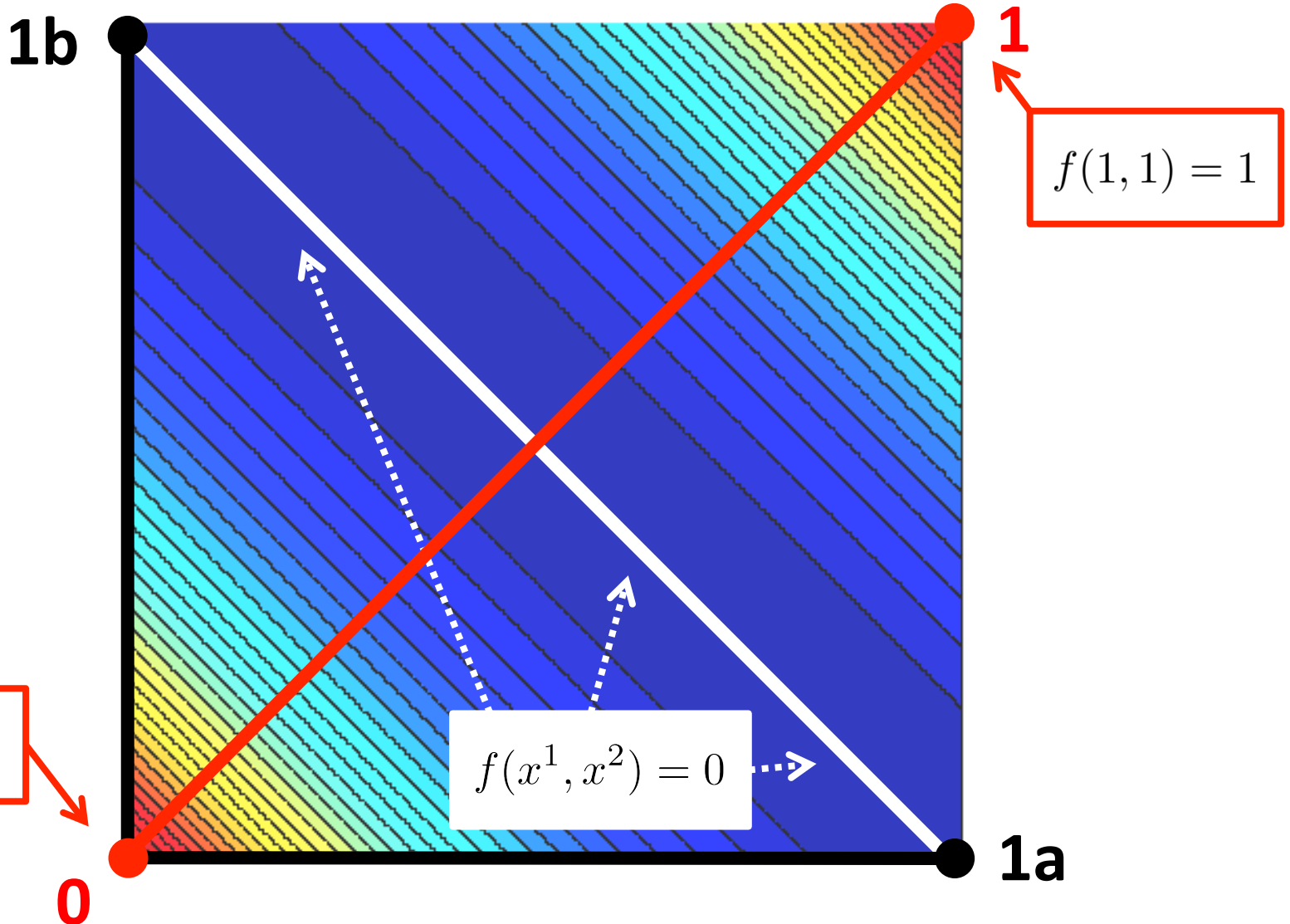
Failure of naive parallelization

$$f(x^1, x^2) = (x^1 + x^2 - 1)^2$$



Failure of naive parallelization

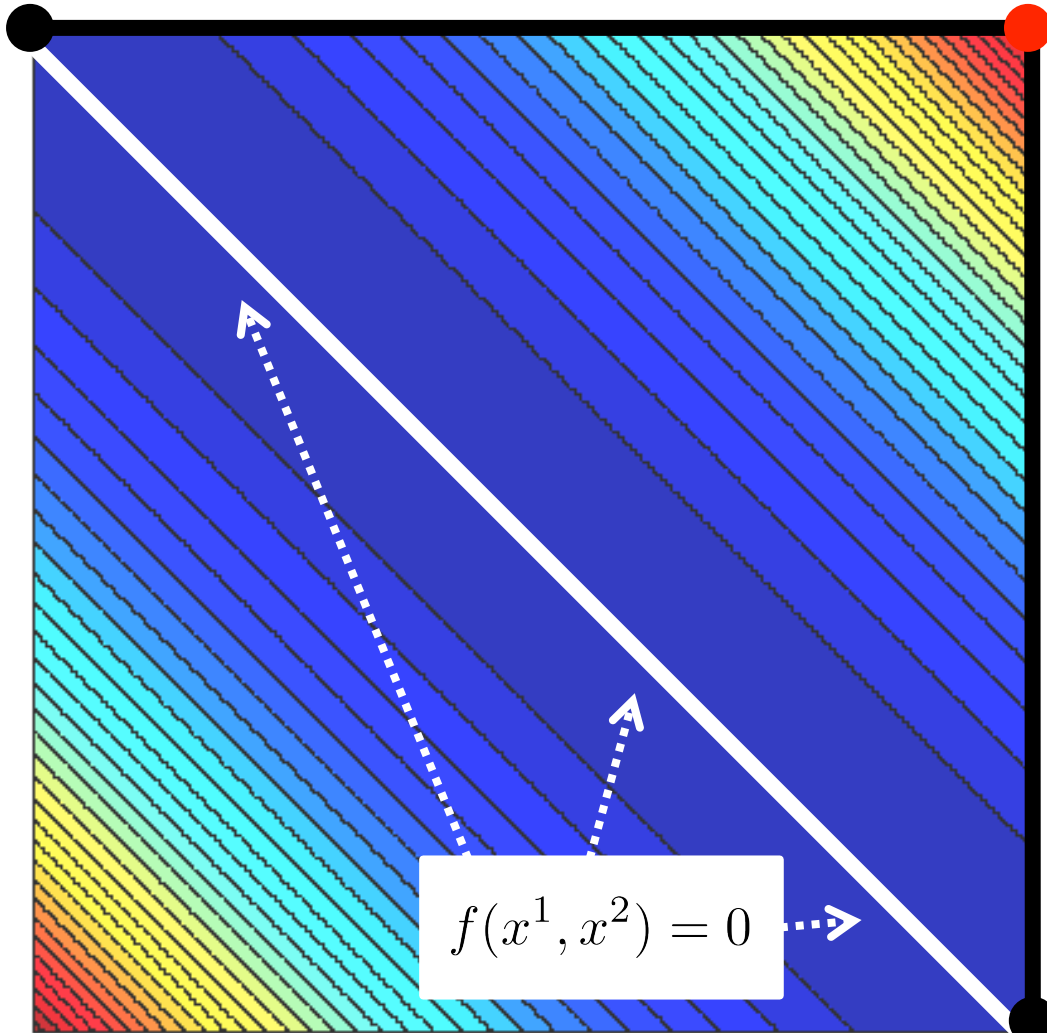
$$f(x^1, x^2) = (x^1 + x^2 - 1)^2$$



Failure of naive parallelization

$$f(x^1, x^2) = (x^1 + x^2 - 1)^2$$

2b



1

$$f(1, 1) = 1$$

$$f(0, 0) = 1$$

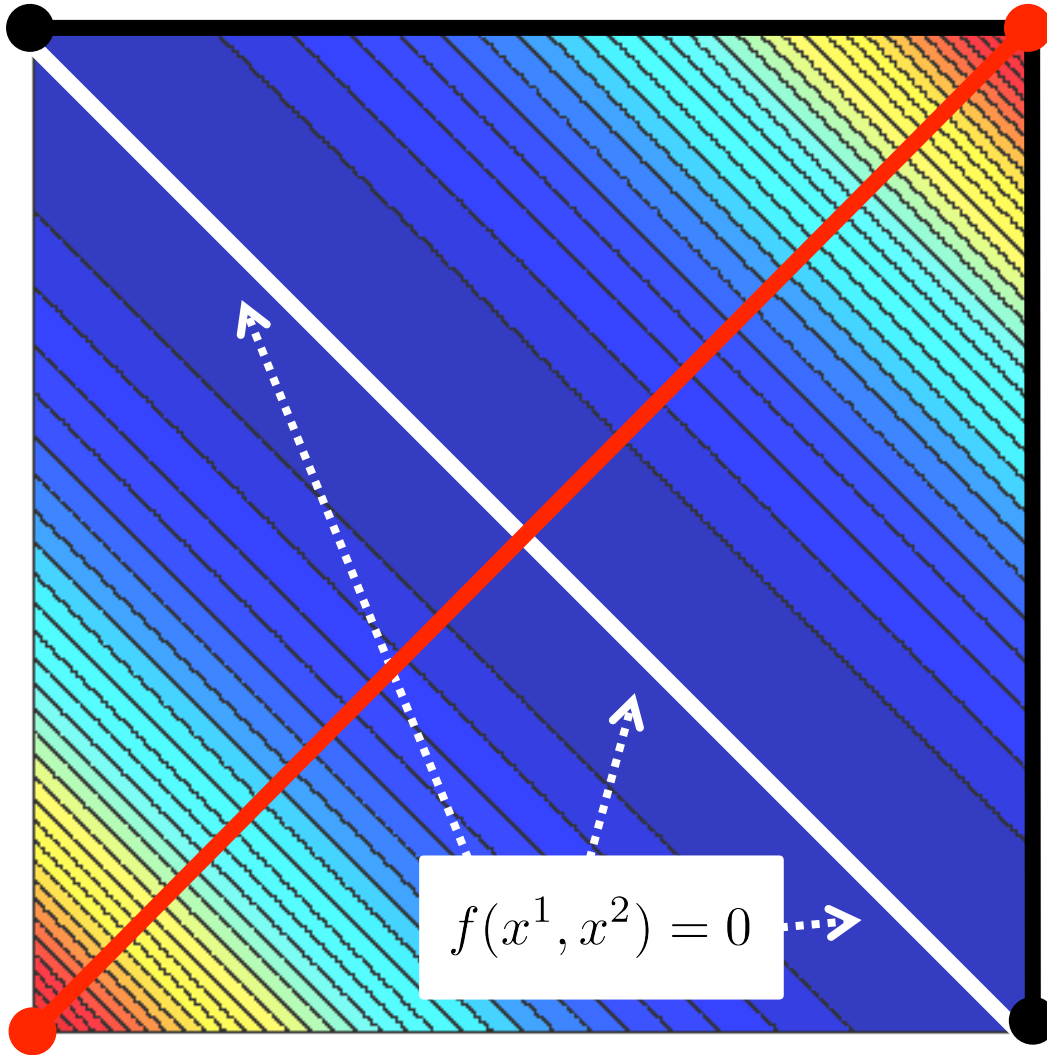
$$f(x^1, x^2) = 0$$

2a

Failure of naive parallelization

$$f(x^1, x^2) = (x^1 + x^2 - 1)^2$$

2b



1

$$f(1, 1) = 1$$

$$f(0, 0) = 1$$

2

$$f(x^1, x^2) = 0$$

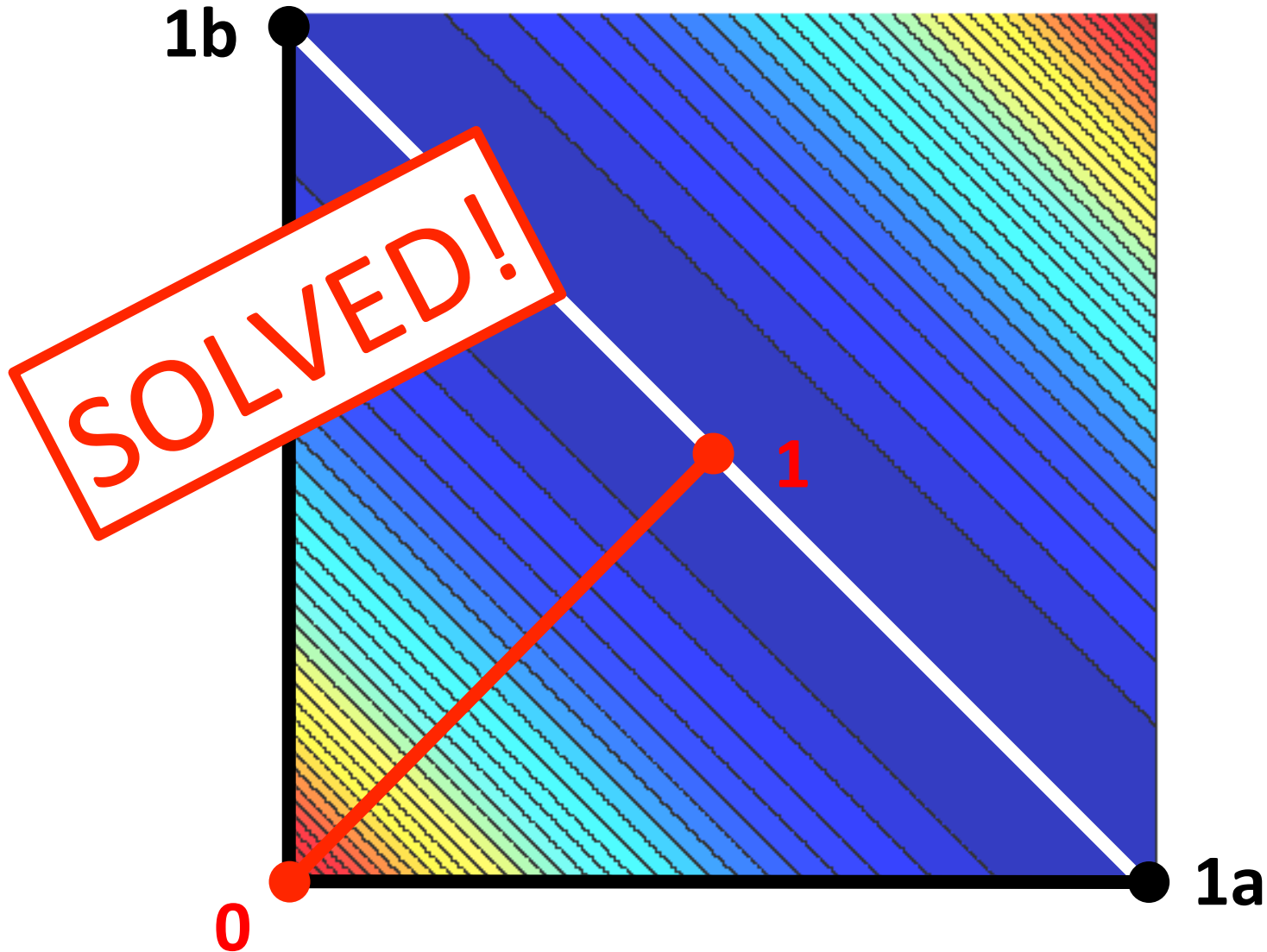
2a

Failure of naive parallelization

$$f(x^1, x^2) = (x^1 + x^2 - 1)^2$$

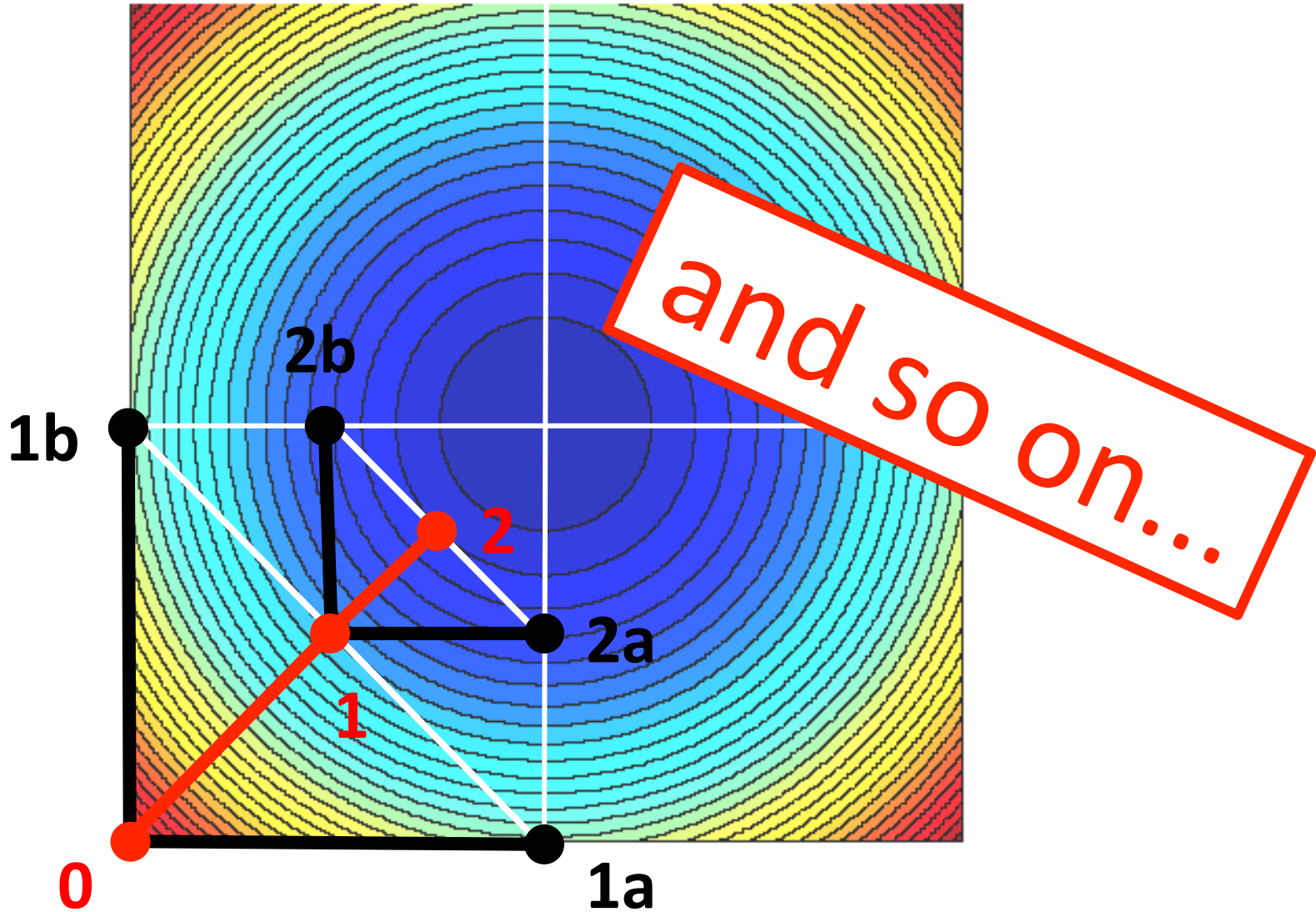


Idea: averaging updates may help



Averaging can be too conservative

$$f(x^1, x^2) = (x^1 - 1)^2 + (x^2 - 1)^2$$



Averaging may be too conservative 2

$$f(x) = (x^1 - 1)^2 + (x^2 - 1)^2 + \dots + (x^n - 1)^2$$

But we wanted:

$$O\left(\frac{n\beta}{\tau} \times \xi(\epsilon)\right)$$

$$k \geq \frac{n}{2} \log\left(\frac{n}{\epsilon}\right)$$

BAD!!!

$$f(x_k) = n \left(1 - \frac{1}{n}\right)^{2k} \leq \epsilon$$

WANT

What to do?

$$x_+ \leftarrow x + \frac{1}{\beta} \sum_{i=1}^n h^i e_i$$

Update to coordinate i

i -th unit coordinate vector

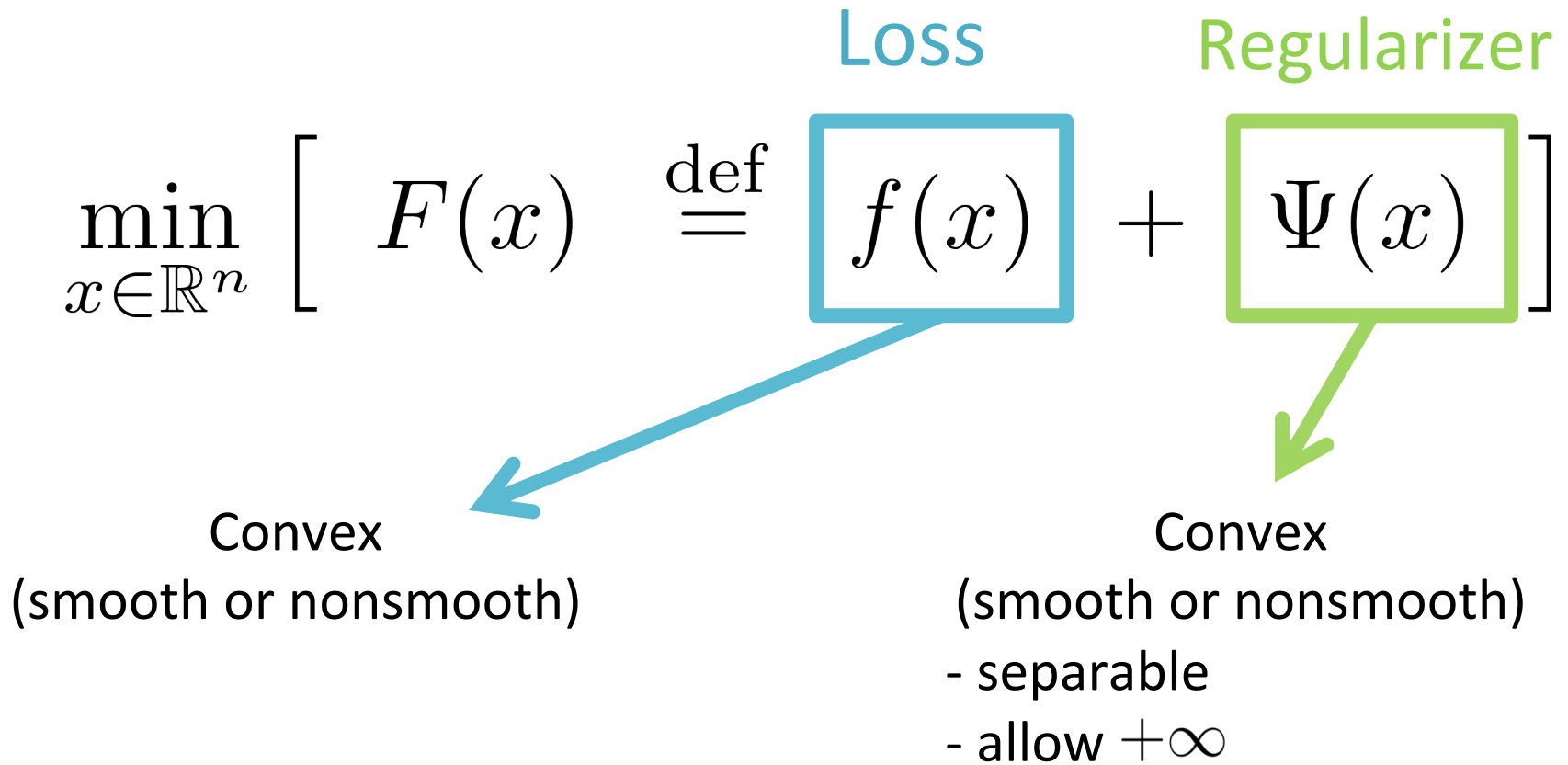
Averaging: $\beta = n$

Summation: $\beta = 1$

Figure out when one can safely use: $\beta \approx 1$

Optimization Problems

Problem



Regularizer: examples

$$\Psi(x) = \sum_{i=1}^n \Psi_i(x^i), \quad x = (x^1, x^2, \dots, x^n)^T$$

No regularizer

$$\Psi_i(x^i) \equiv 0$$

Box constraints

$$\Psi_i(x^i) = \begin{cases} 0, & x^i \in X_i, \\ +\infty, & \text{otherwise.} \end{cases}$$

e.g., SVM dual

e.g., LASSO

Weighted L1 norm

$$\Psi_i(x^i) = \lambda_i |x^i| \quad (\lambda_i > 0)$$

Weighted L2 norm

$$\Psi_i(x^i) = \lambda_i (x^i)^2 \quad (\lambda_i > 0)$$

Loss: examples

$$f(x)$$

Quadratic loss

$$\frac{1}{2} \|Ax - y\|_2^2 = \frac{1}{2} \sum_{j=1}^m (A_{j:}x - y_j)^2$$

Logistic loss

$$\sum_{j=1}^m \log(1 + \exp(-y_j A_{j:}x))$$

Square hinge loss

$$\frac{1}{2} (\max\{0, 1 - y_j A_{j:}x\})^2$$

L-infinity

$$\|Ax - y\|_\infty = \max_{1 \leq j \leq m} |A_{j:}x - y_j|$$

L1 regression

$$\|Ax - b\|_1 = \sum_{j=1}^m |A_{j:}x - y_j|$$

Exponential loss

$$\log \left(\frac{1}{m} \sum_{j=1}^m \exp(y_j A_{j:}x) \right)$$

BKBG'11

RT'11b

TBRS'13

RT'13a

FR'13

3 models for f with small β

1

Smooth partially separable f [RT'11b ]

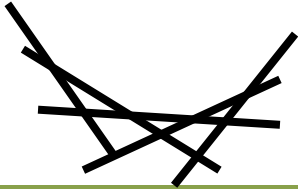
$$f(x + te_i) \leq f(x) + (\nabla f(x))^T te_i + \frac{L_i}{2} t^2$$

$$f(x) = \sum_{J \in \mathcal{J}} f_J(x), \quad f_J \text{ depends on } x^i \text{ for } i \in J \text{ only}$$

$$\omega \stackrel{\text{def}}{=} \max_{J \in \mathcal{J}} |J|$$

2

Nonsmooth max-type f [FR'13]

$$f(x) = \max_{z \in Q} \{z^T Ax - g(z)\}$$


$$\omega \stackrel{\text{def}}{=} \max_{1 \leq j \leq m} |\{i : A_{ji} \neq 0\}|$$

3

f with 'bounded Hessian' [BKBG'11, RT'13a ]

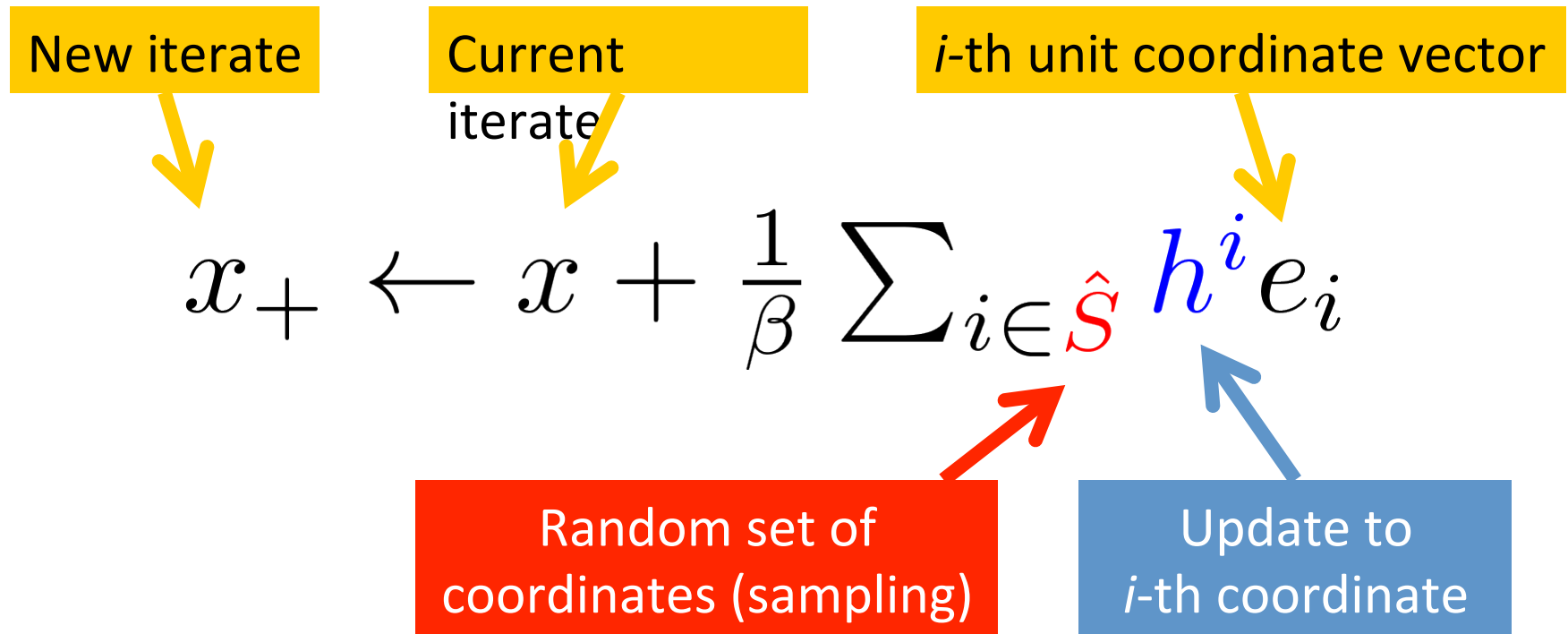


$$f(x + h) \leq f(x) + (\nabla f(x))^T h + \frac{1}{2} h^T A^T A h$$

$$\sigma \stackrel{\text{def}}{=} \lambda_{\max}(L^{-1/2} A^T A L^{-1/2}) \quad L = \text{Diag}(A^T A)$$

General Theory

Randomized Parallel Coordinate Descent Method



The update h^i depends on F , x and on the law describing \hat{S}

ESO: Expected Separable Overapproximation

$$x_+ \leftarrow x + \frac{1}{\beta} \sum_{i \in \hat{S}} \frac{1}{w_i} \nabla_i f(x) e_i$$

Minimize in h

$$\|h\|_w^2 \stackrel{\text{def}}{=} \sum_{i=1}^n w_i (h^i)^2$$

$$\mathbf{E} \left[f(x + h_{[\hat{S}]}) \right] \leq f(x) + \frac{\mathbf{E}[|\hat{S}|]}{n} \left((\nabla f(x))^T h + \frac{\beta}{2} \|h\|_w^2 \right)$$

$$h_{[\hat{S}]} = \sum_{i \in \hat{S}} h^i e_i$$

1. **Separable** in h
2. Can minimize in **parallel**
3. Can compute updates for $i \in \hat{S}$ **only**

Convergence rate: convex f

Theorem [RT'11b] If $(f, \hat{S}) \sim ESO(\beta, w)$, then

The diagram illustrates the convergence rate formula with several annotations:

- stepsize parameter**: points to β in the denominator βn .
- # coordinates**: points to n in the denominator βn .
- # iterations**: points to k .
- average # updated coordinates per iteration**: points to $\mathbf{E}[|\hat{S}|]$.
- error tolerance**: points to ϵ in the denominator of the log argument and $\epsilon \rho$ in the numerator of the log argument.

$$k \geq \left(\frac{\beta n}{\mathbf{E}[|\hat{S}|]} \right) \left(\frac{2R_w^2(x_0, x_*)}{\epsilon} \right) \log \left(\frac{F(x_0) - F(x_*)}{\epsilon \rho} \right)$$

implies

$$\mathbf{P}(F(x_k) - F(x_*) \leq \epsilon) \geq 1 - \rho$$

Convergence rate: strongly convex f

Theorem [RT'11b] If $(f, \hat{S}) \sim ESO(\beta, w)$, then

If $\mu_\psi(w)$ is large, then the slowdown effect of β is eliminated

Strong convexity constant of the regularizer Ψ

$$k \geq \left(\frac{n}{\mathbf{E}[|\hat{S}|]} \right) \left(\frac{\beta + \mu_\Psi(w)}{\mu_f(w) + \mu_\Psi(w)} \right) \log \left(\frac{F(x_0) - F(x_*)}{\epsilon \rho} \right)$$

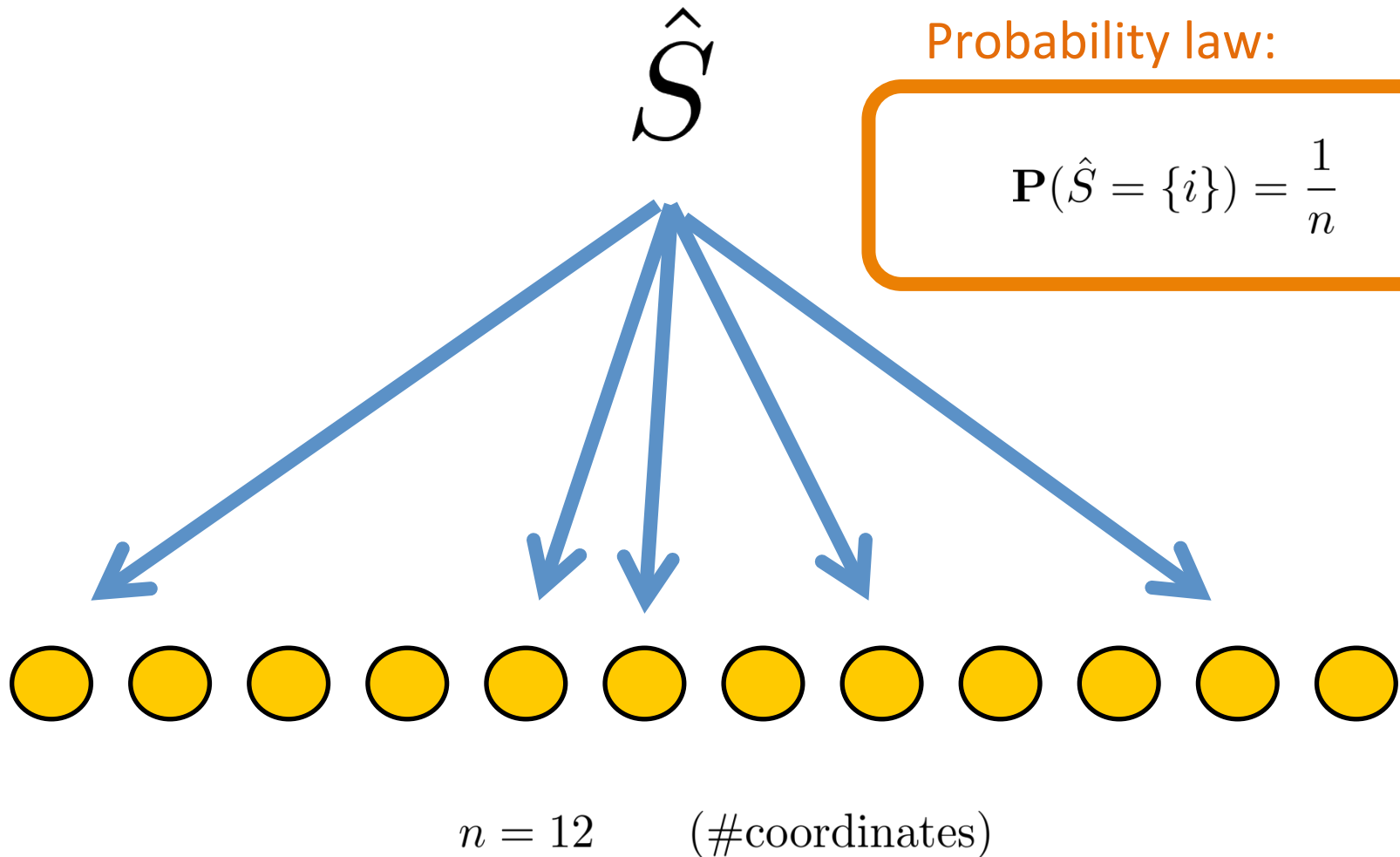
Strong convexity constant of the loss f

implies

$$\mathbf{P}(F(x_k) - F(x_*) \leq \epsilon) \geq 1 - \rho$$

Partial Separability
and
Doubly Uniform
Samplings

Serial uniform sampling



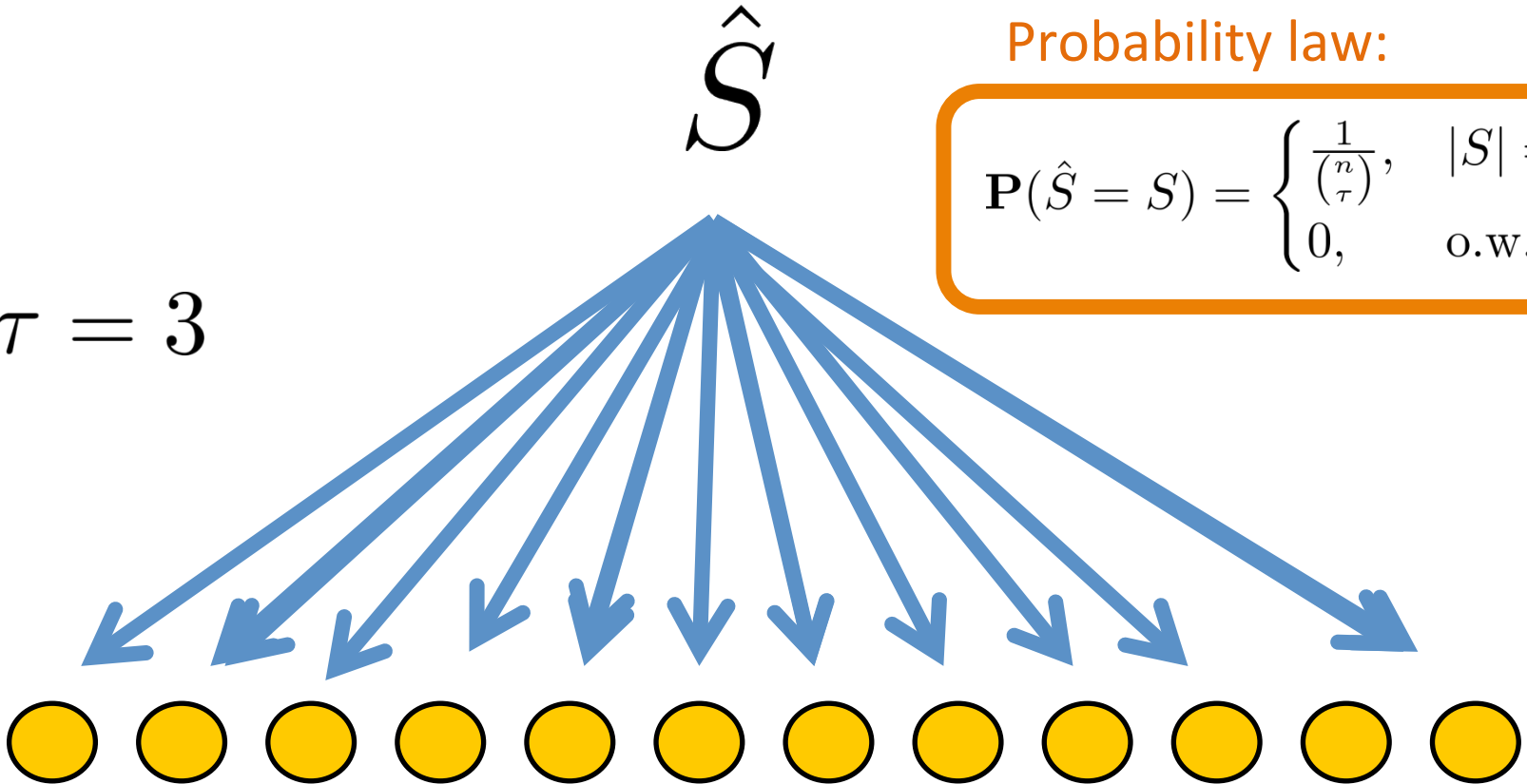
\mathcal{T} -nice sampling

Good for shared memory systems

Probability law:

$$\mathbf{P}(\hat{S} = S) = \begin{cases} \frac{1}{\binom{n}{\tau}}, & |S| = \tau \\ 0, & \text{o.w.} \end{cases}$$

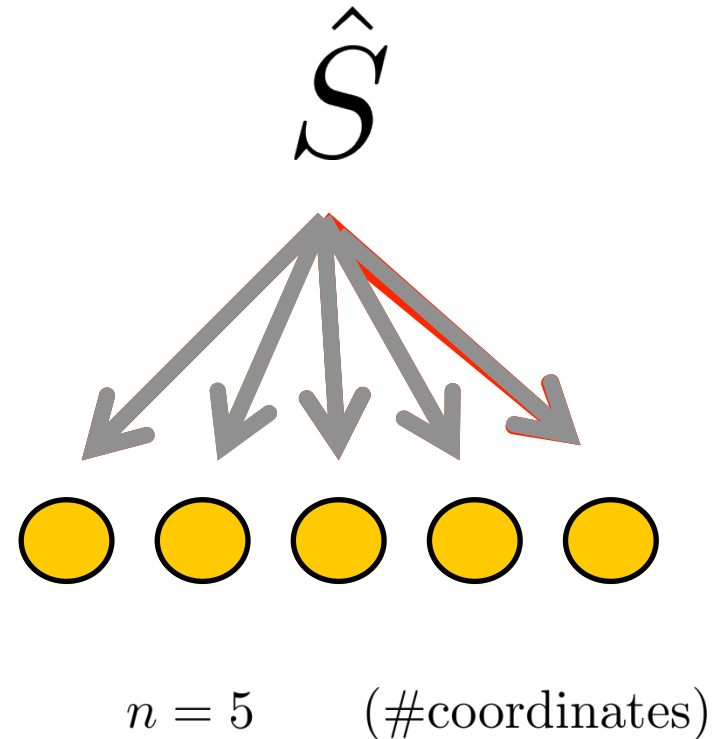
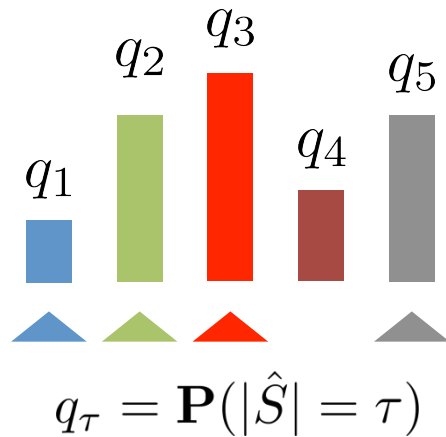
$\tau = 3$



$n = 12$ (#coordinates)

Doubly uniform sampling

Can model unreliable processors / machines



Probability law:

$$\mathbf{P}(\hat{S} = S) = \frac{q_{|S|}}{\binom{n}{|S|}}$$

ESO for partially separable functions and doubly uniform samplings

1

Smooth partially separable f [RT'11b ]

$$f(x + te_i) \leq f(x) + (\nabla f(x))^T te_i + \frac{L_i}{2} t^2$$

$$f(x) = \sum_{J \in \mathcal{J}} f_J(x), \quad f_J \text{ depends on } x^i \text{ for } i \in J \text{ only}$$

$$\omega \stackrel{\text{def}}{=} \max_{J \in \mathcal{J}} |J|$$

Theorem [RT'11b]

If \hat{S} is doubly uniform

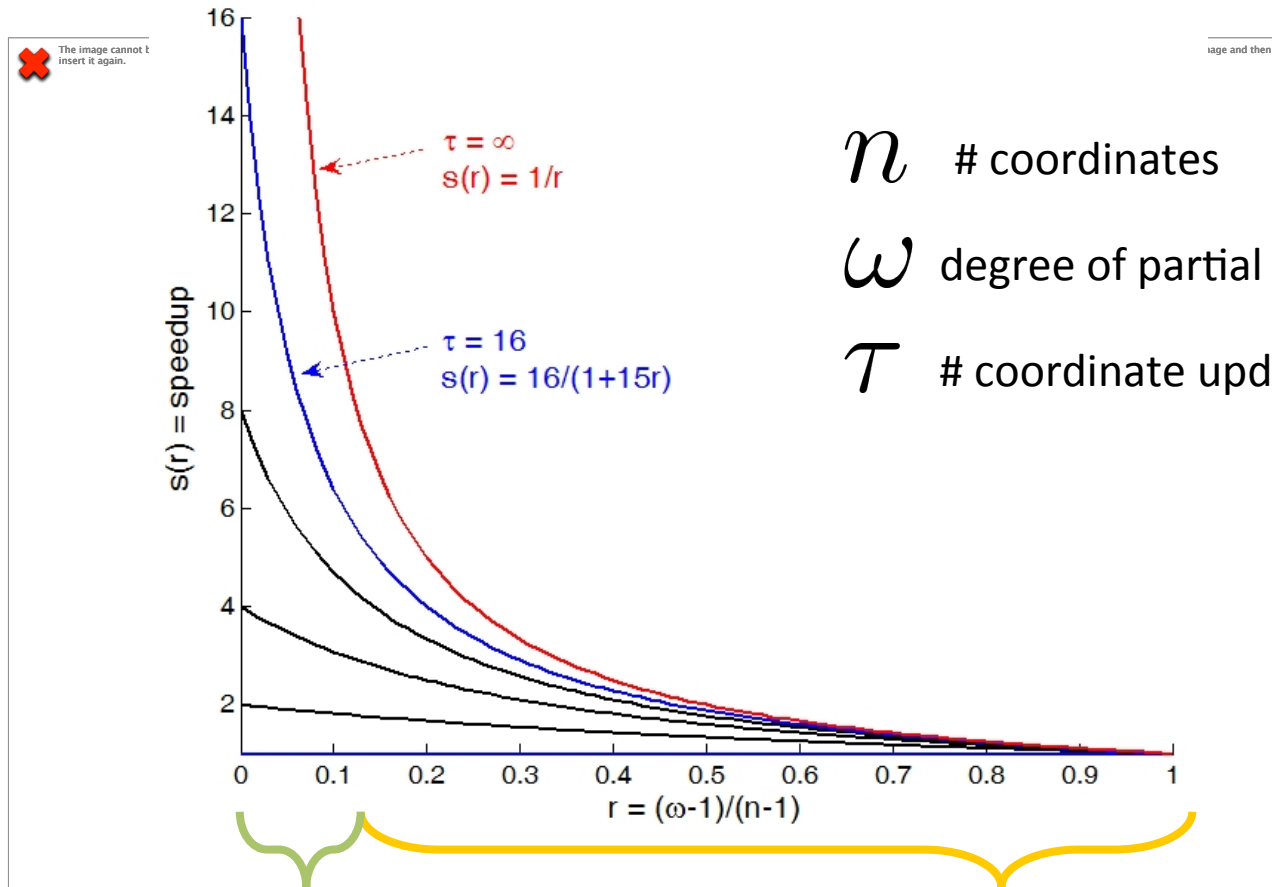
$$\mathbf{E} \left[f(x + h_{[\hat{S}]}) \right] \leq f(x) + \frac{\beta}{2} \|h\|_{w^2}$$

β is small if ω is small!
(i.e., more separable problems are easier)

$$\beta = 1 + \frac{(\omega - 1) \left(\frac{\mathbf{E}[|\hat{S}|^2]}{\mathbf{E}[|\hat{S}|]} - 1 \right)}{n - 1}$$

$$w_i = L_i, \quad i = 1, 2, \dots, n$$

Theoretical speedup



n # coordinates

ω degree of partial separability

\mathcal{T} # coordinate updates / iter

LINEAR OR GOOD SPEEDUP:

Nearly separable (sparse) problems

Much of Big Data is here!

WEAK OR NO SPEEDUP:

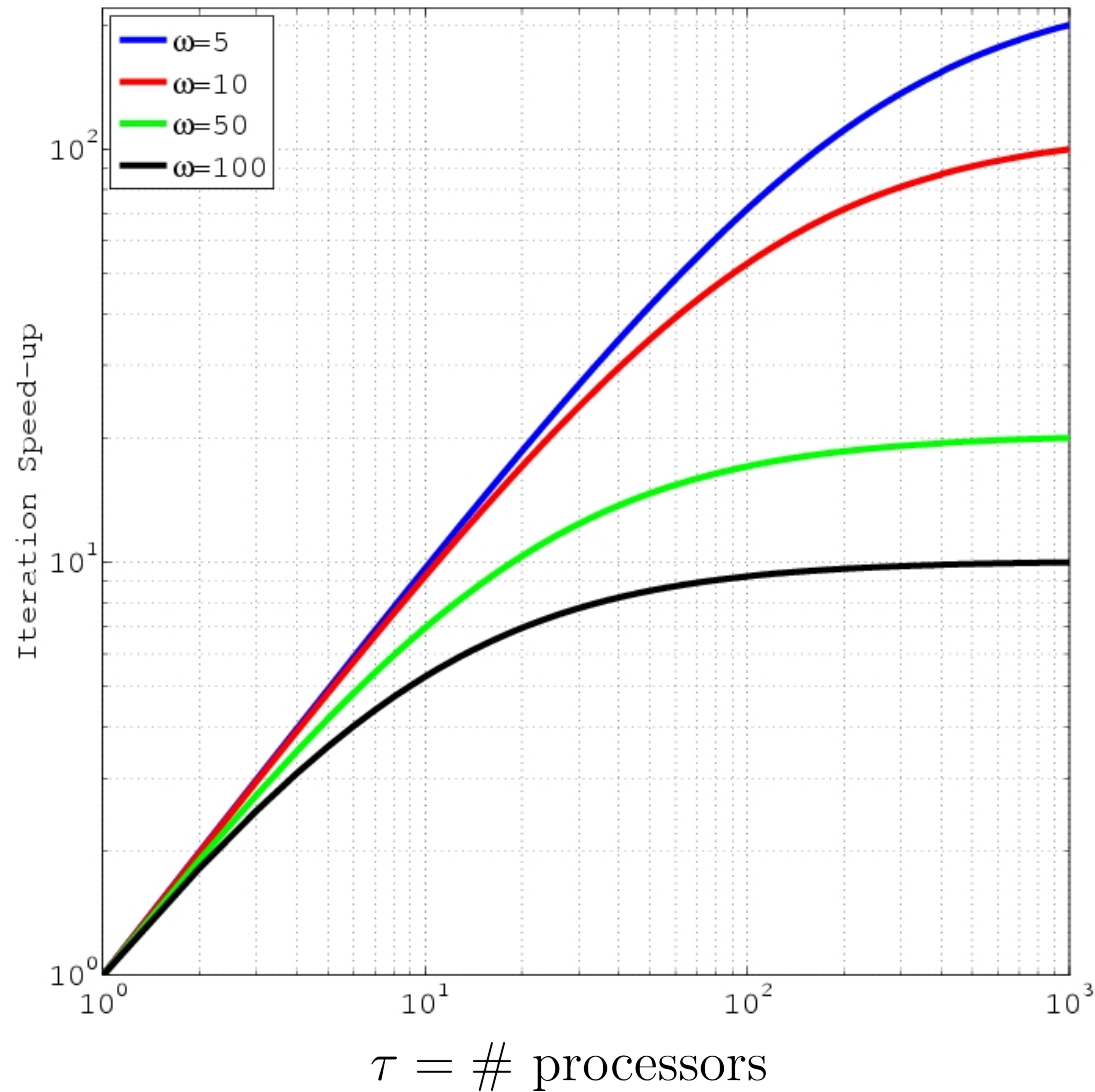
Non-separable (dense) problems

Theory is when you know everything but nothing works.

Practice is when everything works but no one knows why.

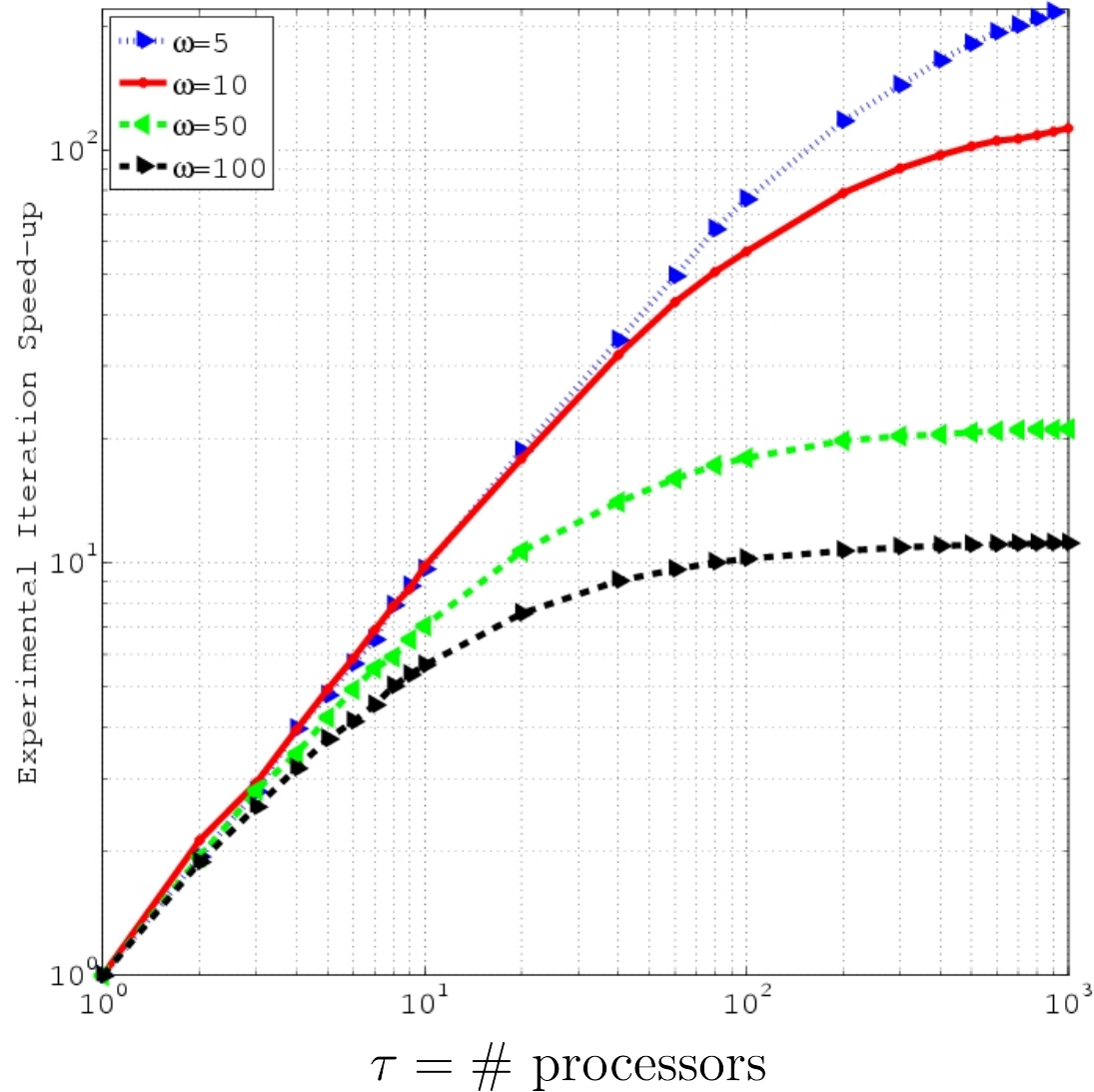
In our lab, theory and practice are combined: nothing works and no one knows why.

Theory



$n = 1000$
(# coordinates)

Practice



$n = 1000$
(# coordinates)

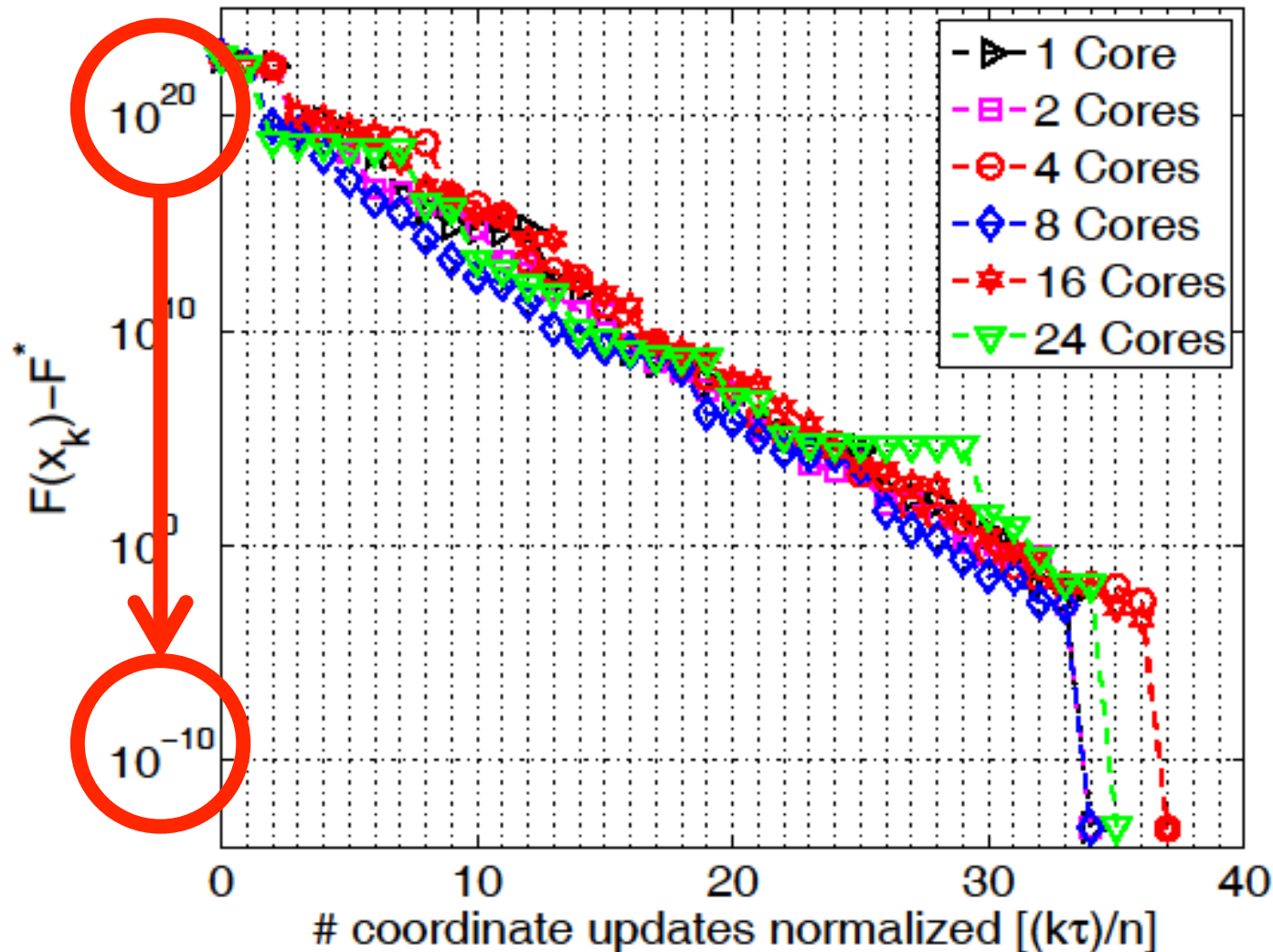
Experiment with a
1 billion-by-2 billion
LASSO problem

Optimization with Big Data = Extreme* Mountain Climbing

* in a billion dimensional space on a foggy day

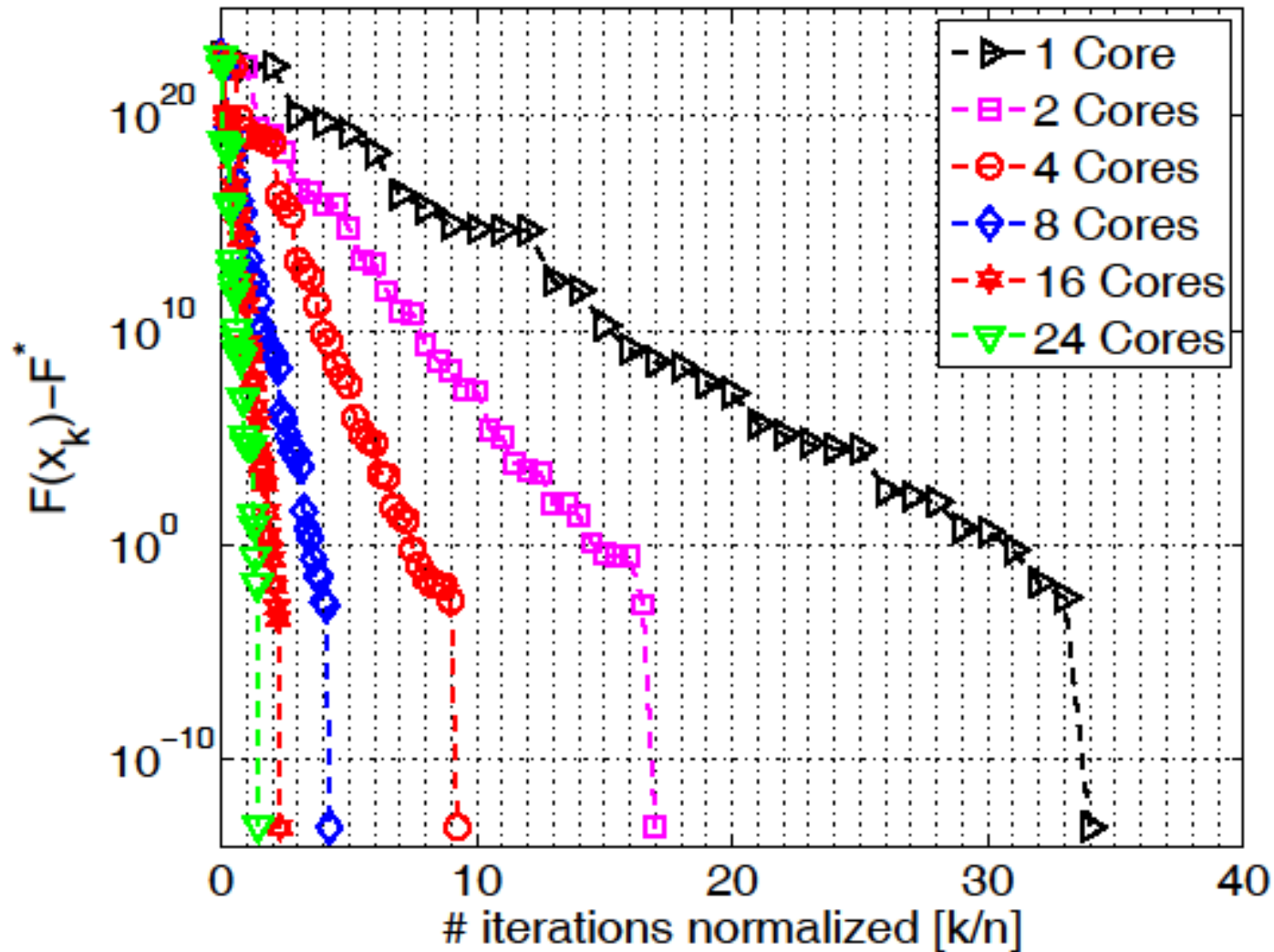


Coordinate Updates



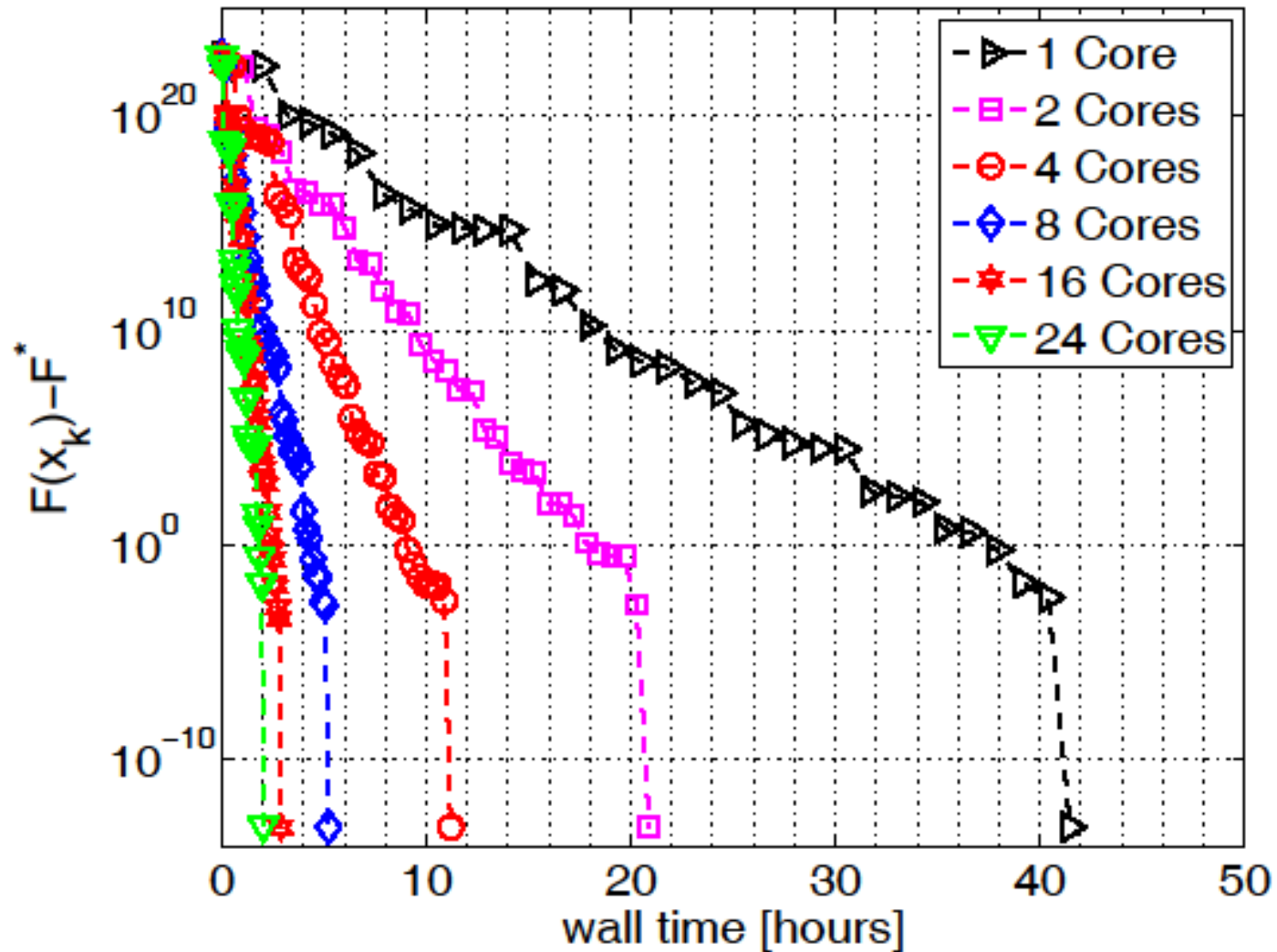
LASSO problem with $A \in \mathbb{R}^{m \times n}$, where $n = 10^9$ and $m = 2 \times 10^9$

Iterations



LASSO problem with $A \in \mathbb{R}^{m \times n}$, where $n = 10^9$ and $m = 2 \times 10^9$

Wall Time



LASSO problem with $A \in \mathbb{R}^{m \times n}$, where $n = 10^9$ and $m = 2 \times 10^9$

Distributed-Memory Coordinate Descent



Distributed τ -nice sampling

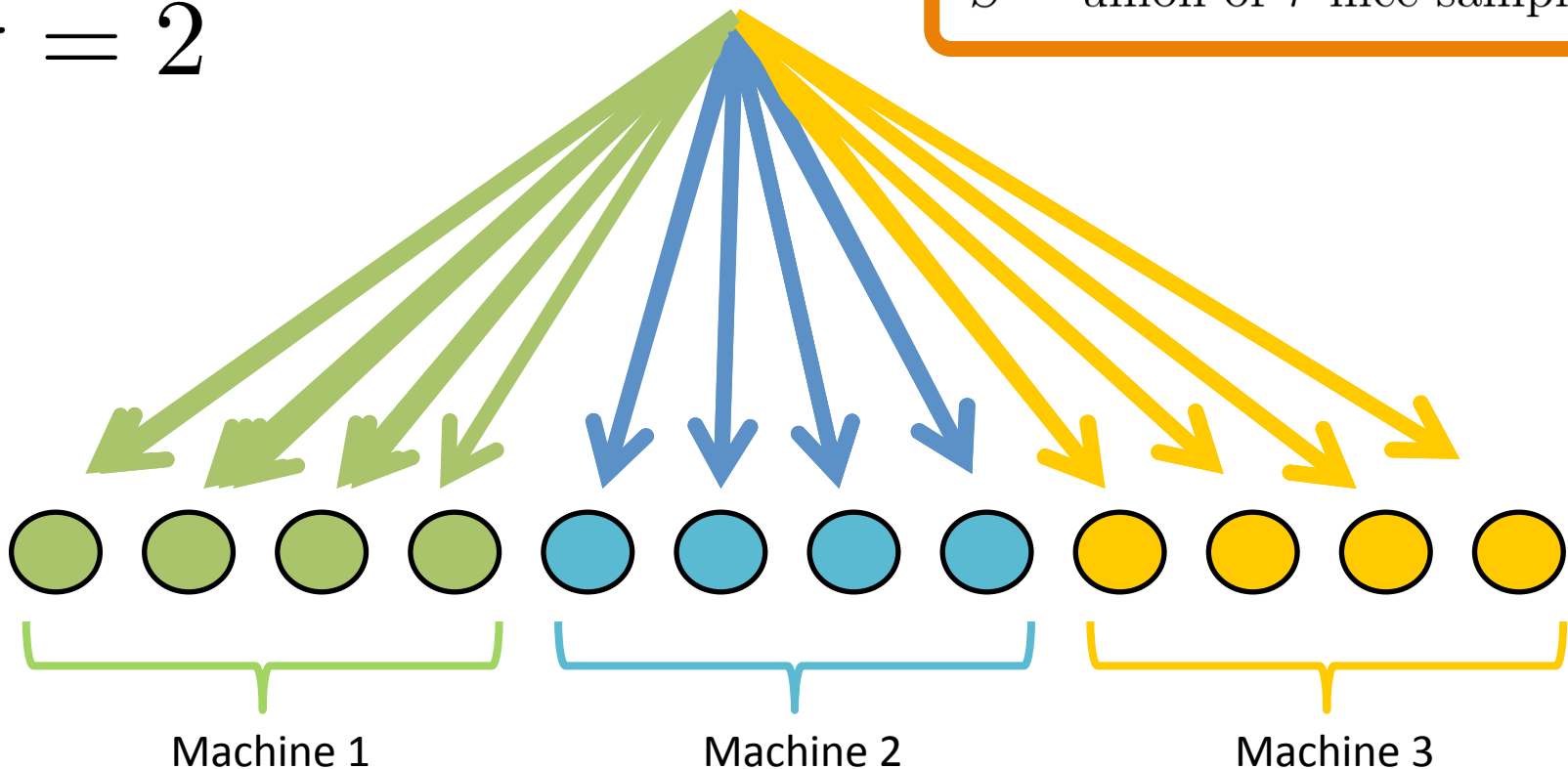
Good for a distributed version of coordinate descent

\hat{S}

Probability law:

$\hat{S} = \text{union of } \tau\text{-nice samplings}$

$\tau = 2$



ESO: Distributed setting

3

f with 'bounded Hessian' [BKBG'11, RT'13a



$$f(x+h) \leq f(x) + (\nabla f(x))^T h + \frac{1}{2} h^T A^T A h$$

$$L = \text{Diag}(A^T A)$$

$$\sigma \stackrel{\text{def}}{=} \lambda_{\max}(L^{-1/2} A^T A L^{-1/2})$$

Theorem [RT'13b]

spectral norm of the **data**

If \hat{S} is **distributed τ -nice sampling** (with c nodes, each owning $s \stackrel{\text{def}}{=} n/c$ coordinates) and Assumption 1 holds, then

$$\mathbf{E} \left[f(x + h_{[\hat{S}]}) \right] \leq f(x) + \frac{\mathbf{E}[|\hat{S}|]}{n} \left((\nabla f(x))^T h + \frac{\beta}{2} \|h\|_{w}^2 \right)$$

$$\beta = 1 + \frac{(\sigma - 1)(\tau - 1)}{n - 1} + \left(\frac{\tau}{s} - \frac{\tau - 1}{s - 1} \right) \frac{\sigma' - 1}{\sigma'}$$

$$w_i = L_{ii}, \quad i = 1, 2, \dots, n$$

Bad partitioning at most doubles # of iterations

spectral norm of the partitioning

$$\beta = 1 + \underbrace{\frac{(\sigma - 1)(\tau - 1)}{n - 1}}_{\beta_1} + \underbrace{\left(\frac{\tau}{s} - \frac{\tau - 1}{s - 1} \right) \frac{\sigma' - 1}{\sigma'}}_{\beta_2} \sigma$$

Theorem [RT'13b] If $\tau \geq 2$, then $\beta_2 \leq \beta_1$

implies

$$\# \text{ iterations} = O\left(\frac{2\beta_1 n}{\tau c}\right)$$

updates/node

nodes

1. PROBLEM FORMULATION

$$\min_{x \in \mathbb{R}^n} [F(x) \equiv f(x) + \Psi(x)]$$

- f convex, **partially separable of degree ω** and $\forall x \in \mathbb{R}^n, t \in \mathbb{R}$ and $i \in \{1, 2, \dots, n\}$ satisfying $|\nabla_i f(x) - \nabla_i f(x + te_i)| \leq L_i |t|$, where L_i are coordinate Lipschitz constants
- Ψ convex and separable ($\Psi(x) = \sum_i \Psi_i(x^i)$)
- Description of f so large that it does not fit onto a single computer! \Rightarrow a cluster of C nodes

2. THE ALGORITHM

Pre-processing: Partition coordinates $\{1, 2, \dots, n\}$ to C sets S_1, S_2, \dots, S_C

In one iteration computers $c = 1, 2, \dots, C$ in parallel do

- Choose random $\hat{S}_c \subset S_c$
- For each $i \in \hat{S}_c$ in parallel compute $t_i^* \leftarrow \arg \min_{t \in \mathbb{R}} \nabla_i f(x_k) t + \beta \frac{\omega}{2} t^2 + \Psi_i(x_k^i + t)$
- $x_{k+1} \leftarrow x_k + \sum_{i \in \hat{S}_c} t_i^* e_i$

2. DISTRIBUTED SAMPLING

We can analyze the above algorithm under the following assumptions:

- $|S_c| = \frac{n}{C}$ for all $c = 1, 2, \dots, C$
- \hat{S}_c is chosen uniformly as one of the subsets of S_c of cardinality τ

Distributed sampling: $\hat{S} = \cup_{c=1}^C \hat{S}_c$

$$\beta = 1 + \frac{-C(\tau - 1) + \omega[\tau C(1 + \frac{C-1}{m}) - 1]}{\max\{n - C, 1\}} \quad (\text{see [1]})$$

Special cases:

- $C = 1 \Rightarrow \beta = 1 + \frac{(\tau-1)(\omega-1)}{\max\{n-1, 1\}}$ (see [2])
- $C = \tau = 1 \Rightarrow \beta = 1$ (see [3])

However, we need new analysis for the $C > 1$ case.

4. COMPLEXITY THEOREM

$$k \geq \frac{\beta n}{\tau C} \frac{2R^2}{\epsilon} \log \left(\frac{F(x_0) - F^*}{\epsilon \rho} \right)$$

\Downarrow

$$\text{Prob}(F(x_k) - F(x_*) \leq \epsilon) \geq 1 - \rho$$

$$(R^2 \approx \sum_i L_i (x_0^i - z_i^*)^2)$$

5. AC/DC SOLVER

We developed a solver (<http://code.google.com/p/ac-dc/>) for

$$f(x) = \sum_{i=1}^m \text{Loss}(x; A_i, b_i), \quad \Psi(x) = \lambda \|x\|_1$$

3 supported losses	$\text{Loss}(x; A_j, b_j)$
square loss	$\frac{1}{2}(b_j - A_j x)^2$
logistic loss	$\log(1 + e^{-b_j A_j x})$
hinge square loss	$\frac{1}{2} \max\{0, 1 - b_j A_j x\}^2$

Note that $A_j \in \mathbb{R}^m$ is a row vector and later will represent the j -th row of matrix A .

7. IMPLEMENTATION DETAILS (SQUARE LOSS EXAMPLE)

If we can maintain $g_k = Ax_k - b$ on all computers, then since $\nabla_i f(x_k) = \langle a_i, g_k \rangle$ (a_i is the i -th column of matrix A), computer c can compute $\nabla_i f(x_k)$ for $i \in S_c$, and hence the algorithm can be run.

- Note that $g_{k+1} = Ax_{k+1} - b = A(x_k + \sum_{c=1}^C \sum_{i \in \hat{S}_c} t_i^* e_i) - b = g_k + \sum_{c=1}^C \sum_{i \in \hat{S}_c} a_i t_i^*$
- That is, computer c additively contributes $g_k[c] := \sum_{i \in \hat{S}_c} a_i t_i^*$ to the update of g_k
- So, we need to add up the distributed updates $g_k[c]$

Reduce All (RA)



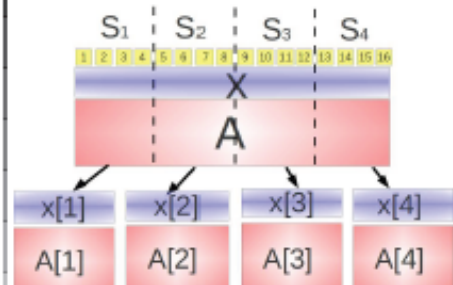
Asynchronous StreamLine (ASL)



- $G_k[c] = G_{k-1}[\text{Prev}(c)] + g_k[c] - g_{k-1}[c]$
- $g_k^c = g_k^c + g_k[c] + G_k[\text{Prev}(c)] - g_{k-1}[c]$
- ASL: **MUCH LESS** communication than RA!
- ASL: **asynchronous** (non-blocking) communication
- ASL: **communication only between two closest computers**

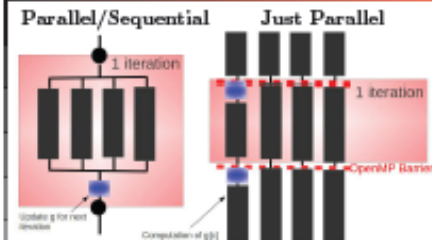
6. DATA DISTRIBUTION

Assume that we have $C = 4$ compute nodes and $n = 16$ coordinates. The coordinates can be partitioned into 4 balanced groups $\{S_1, S_2, S_3, S_4\}$.



On computer 1, only the first 4 coordinates of vector x are stored and also the corresponding 4 columns of matrix A . **Data distribution is crucial for problems whose size exceeds available memory of a single computer!**

8. HYBRID IMPLEMENTATIONS

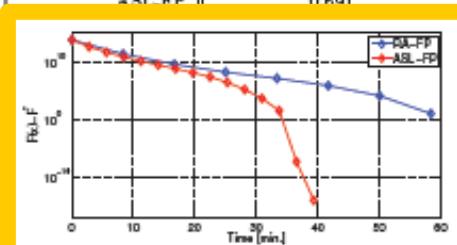


Left image shows **Parallel and Serial (PS)** approach, where each MPI process runs few OpenMP threads for computing t_i^* and $g_k[c]$ (black boxes) and afterwards, **MPI communication takes place** (blue boxes). Right image shows **Fully Parallel (FP)** approach in which one of the threads deals with communication and when waiting for a new communication, it helps the other threads to do some computation.

9. NUMERICAL EXPERIMENTS

All experiments were done on HECTOR - Cray XE6 using **2,048 cores**. Problem size $A \in \mathbb{R}^{10^6 \times 10^6}$ had **1.2 TBytes** and we used $\tau = 10^3$.

method	avg. time / iter.
RA-PS	2.252
RA-FP	2.052
ASL-FP	0.601



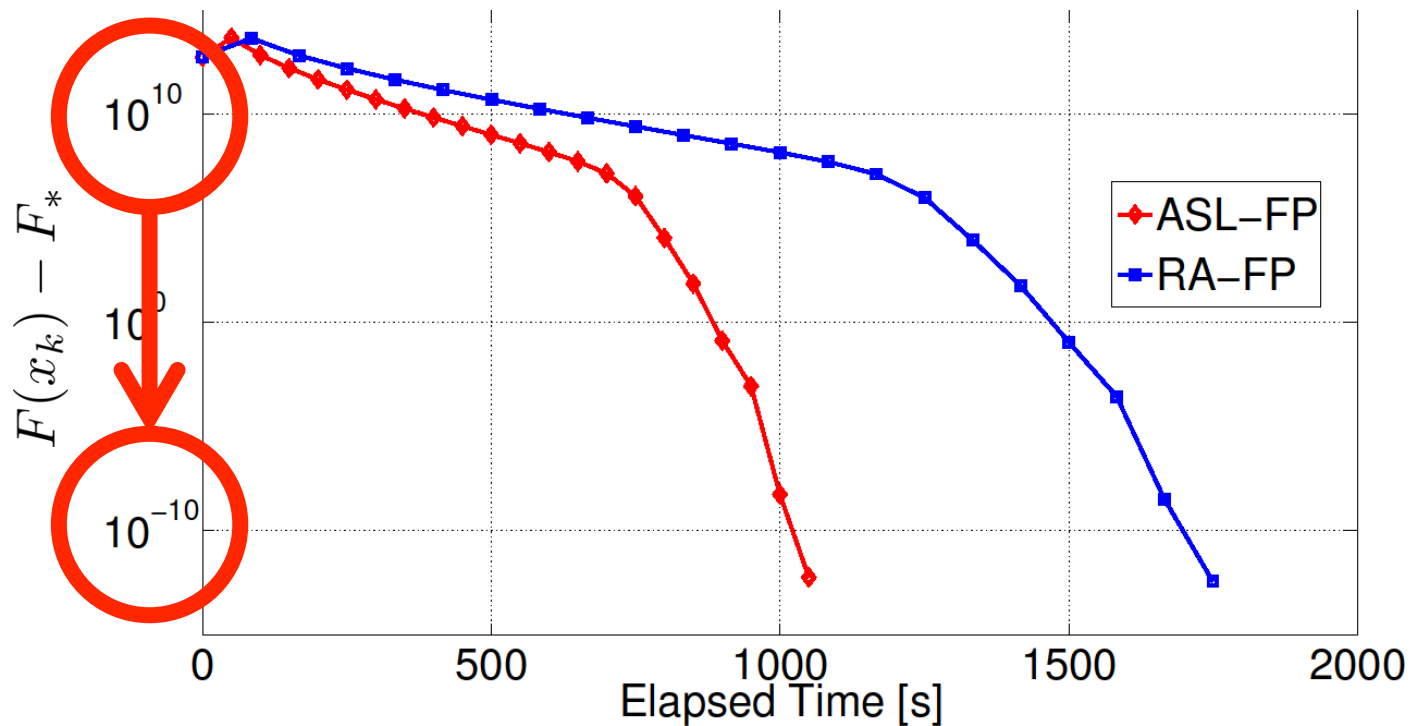
10. REFERENCES

- Takáč, M., Mareček, J. and Richtárik, P.: *Distributed coordinate descent methods for big data optimization*, 2013
- Richtárik, P., Takáč, M.: *Parallel coordinate descent methods for big data optimization*, 2012
- Richtárik, P., Takáč, M.: *Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function*, Mathematical Programming, 2012

LASSO with a 3TB data matrix

$$A = 10^9\text{-by-}0.5 \times 10^9$$

n = # coordinates





128 Cray XE6 nodes with 4 MPI processes ($c = 512$)

Each node: 2 x 16-cores with 32GB RAM

Conclusions

- Coordinate descent methods scale very well to big data problems of special structure
 - Partial separability (sparsity)
 - Small spectral norm of the data
 - Nesterov separability, ...
- Care is needed when combining updates (add them up? average?)

References: serial coordinate descent

- Shai Shalev-Shwartz and Ambuj Tewari, **Stochastic methods for L1-regularized loss minimization**. *JMLR 2011*.
- Yurii Nesterov, **Efficiency of coordinate descent methods on huge-scale optimization problems**. *SIAM Journal on Optimization*, 22(2):341-362, 2012.
-  [RT'11b] P.R. and Martin Takáč, **Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function**. *Mathematical Prog.*, 2012.
-  Rachael Tappenden, P.R. and Jacek Gondzio, **Inexact coordinate descent: complexity and preconditioning**, *arXiv: 1304.5530*, 2013.
- Ion Necoara, Yurii Nesterov, and Francois Glineur. **Efficiency of randomized coordinate descent methods on optimization problems with linearly coupled constraints**. *Technical report, Politehnica University of Bucharest*, 2012.
- Zhaosong Lu and Lin Xiao. **On the complexity analysis of randomized block-coordinate descent methods**. *Technical report, Microsoft Research*, 2013.

References: parallel coordinate descent



- [BKBG'11] Joseph Bradley, Aapo Kyrola, Danny Bickson and Carlos Guestrin, **Parallel Coordinate Descent for L1-Regularized Loss Minimization**. *ICML 2011*



- [RT'12] P.R. and Martin Takáč, **Parallel coordinate descent methods for big data optimization**. *arXiv:1212.0873*, 2012



- Martin Takáč, Avleen Bijral, P.R., and Nathan Srebro. **Mini-batch primal and dual methods for SVMs**. *ICML 2013*

- [FR'12] Olivier Fercoq and P.R., **Smooth minimization of nonsmooth functions with parallel coordinate descent methods**. *arXiv:1309.5885*, 2013



- [RT'13a] P.R. and Martin Takáč, **Distributed coordinate descent method for big data learning**. *arXiv:1310.2059*, 2013

- [RT'13b] P.R. and Martin Takáč, **On optimal probabilities in stochastic coordinate descent methods**. *arXiv:1310.3438*, 2013

Good entry point to the topic (4p paper)

References: parallel coordinate descent



- P.R. and Martin Takáč, **Efficient serial and parallel coordinate descent methods for huge-scale truss topology design.** *Operations Research Proceedings 2012.*
- Rachael Tappenden, P.R. and Burak Buke, **Separable approximations and decomposition methods for the augmented Lagrangian.** *arXiv:1308.6774, 2013.*
- Indranil Palit and Chandan K. Reddy. **Scalable and parallel boosting with MapReduce.** *IEEE Transactions on Knowledge and Data Engineering, 24(10): 1904-1916, 2012.*
- Shai Shalev-Shwartz and Tong Zhang, **Accelerated mini-batch stochastic dual coordinate ascent.** *NIPS 2013.*

TALK TOMORROW