

Proximal Algorithms for Large-scale Convex Nonsmooth Optimization

Laurent Condat

King Abdullah University of Science and Technology (KAUST)
Thuwal, Saudi Arabia



Feb. 14, 2022

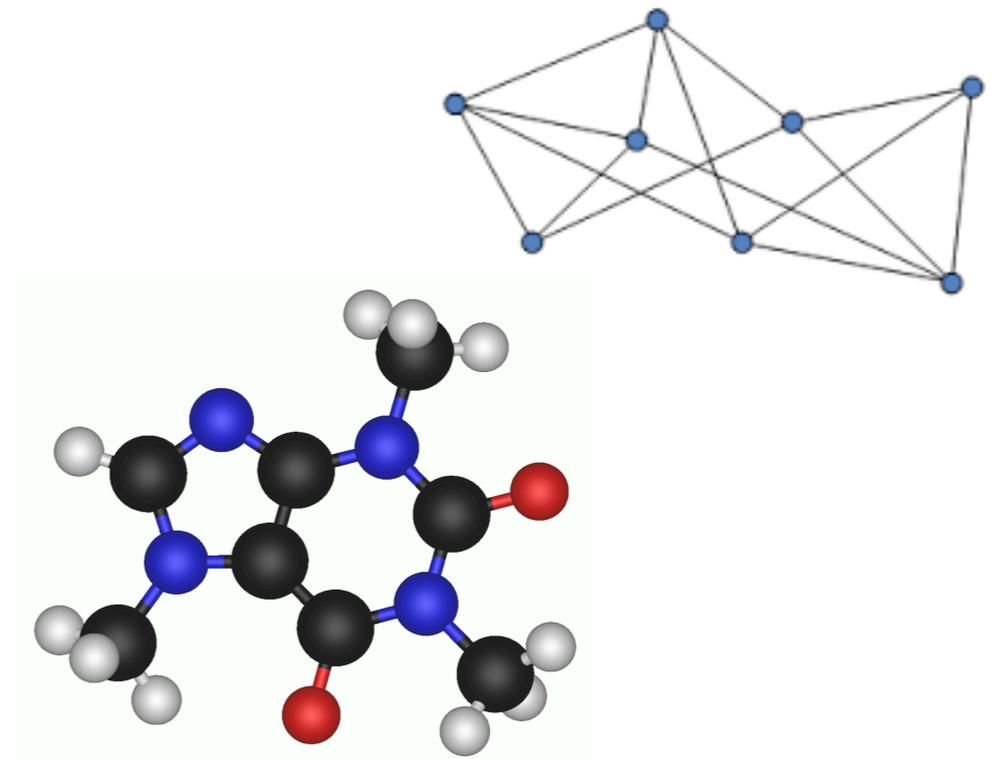
Optimization

Find $x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x)$

\mathcal{X} is a real Hilbert space:



$$\begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$



Optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

with

- linear operators $L_m : \mathcal{X} \rightarrow \mathcal{U}_m$
- real Hilbert spaces $\mathcal{X}, \mathcal{U}_m$
- functions g_m

Optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

$$\text{Example: } \psi(x) = \|Ax - y\|^2 + \|x\|_1$$

Motivation: image processing

Given data $y \approx Ax^\#$



Motivation: image processing

Given data $y \approx Ax^\#$

estimate the unknown image $x^\#$ by solving

$$\text{Find } x^\star \in \arg \min_{x \in \mathbb{R}^{N_1 \times N_2}} \left\{ D(Ax, y) + R(x) \right\}$$

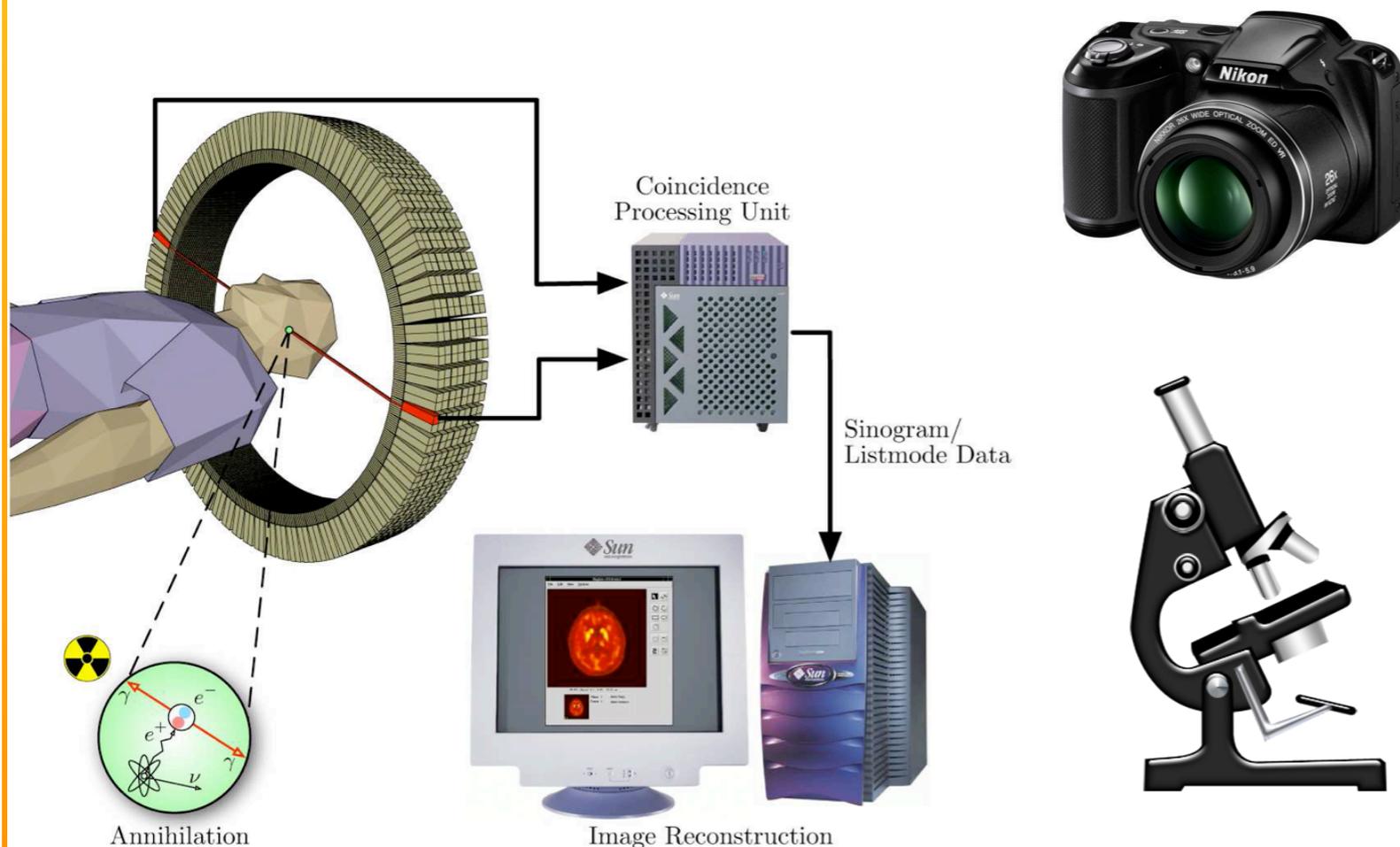


Motivation: image processing

Given data $y \approx Ax^\#$

estimate the unknown image $x^\#$ by solving

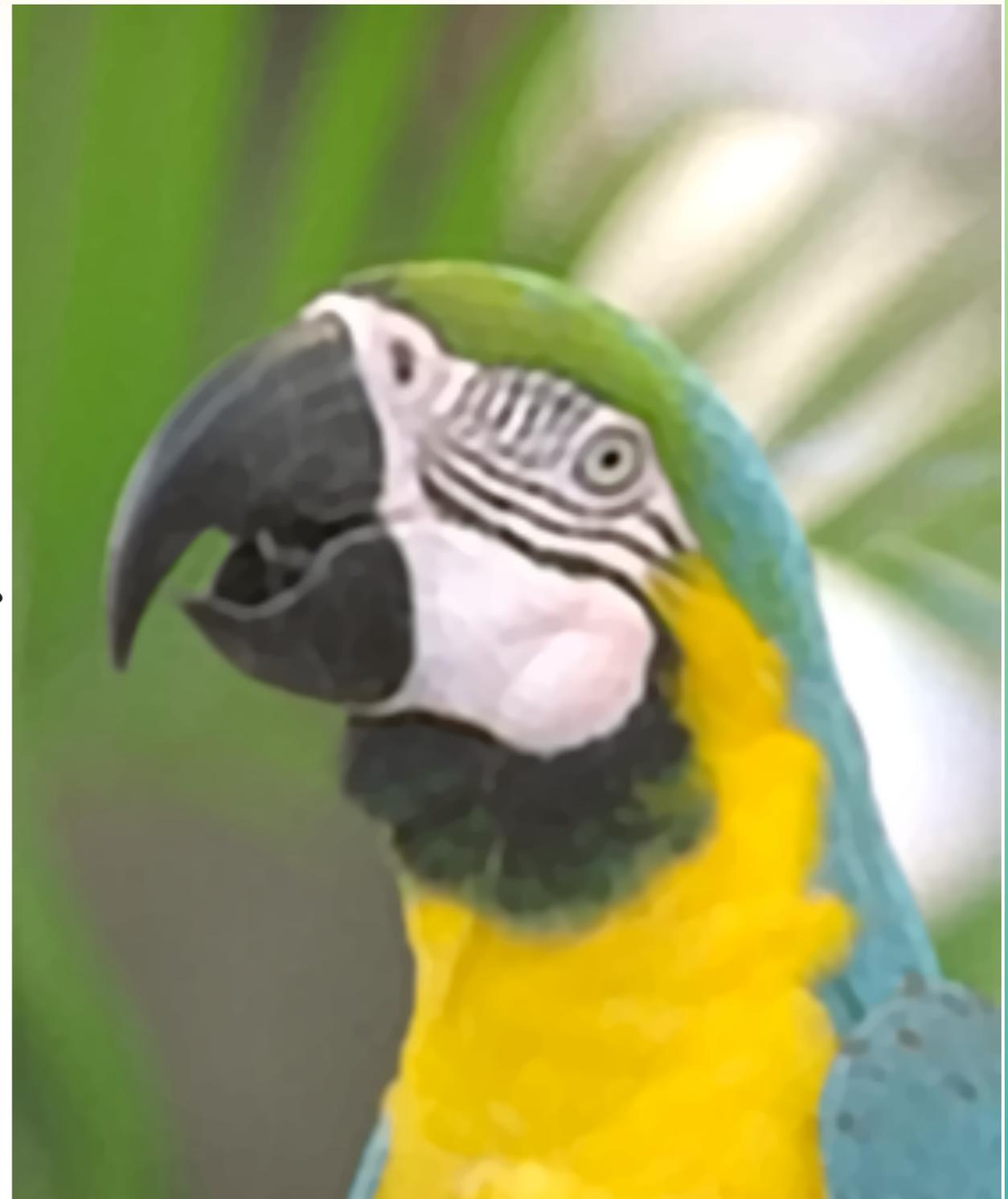
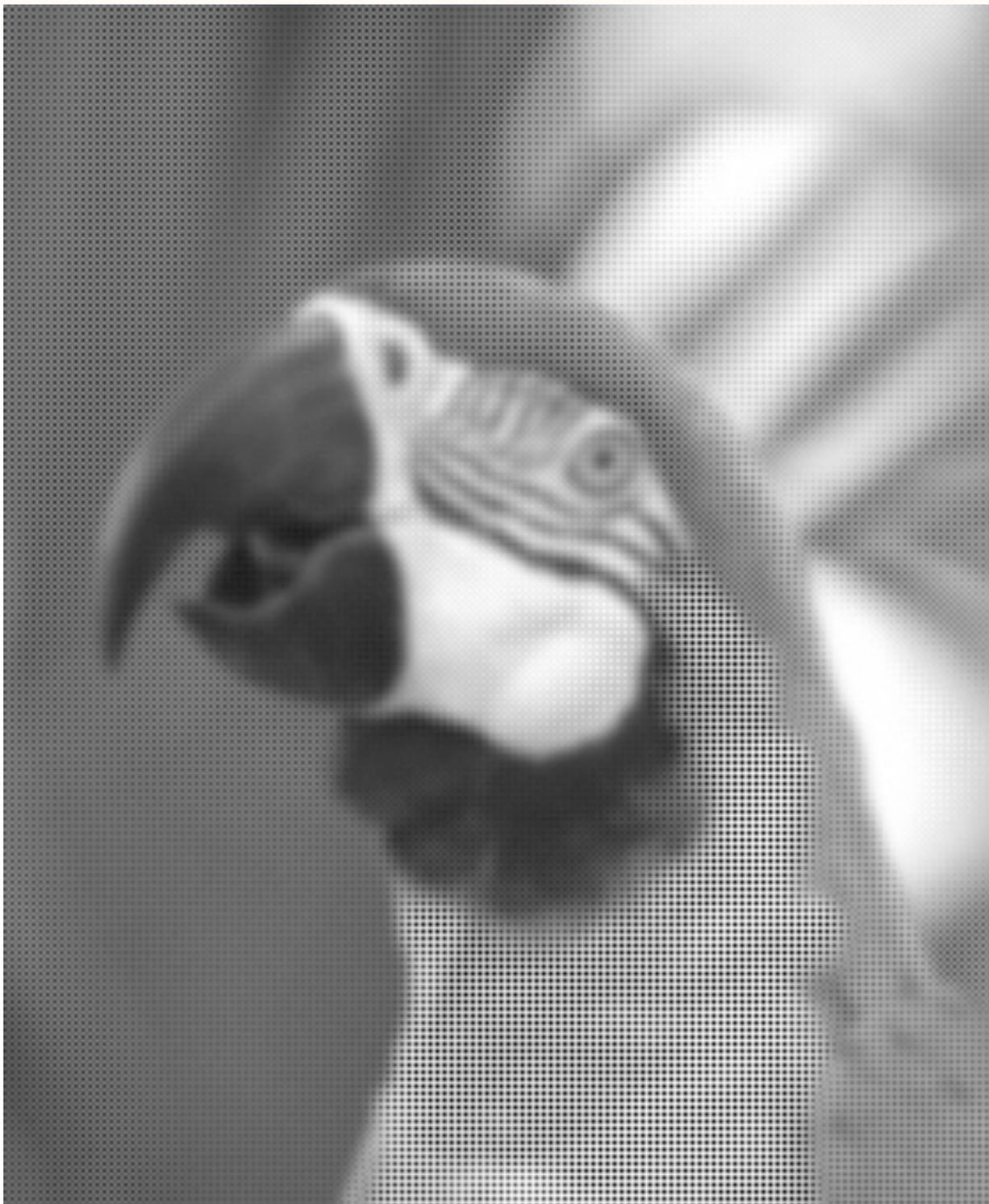
$$\text{Find } x^\star \in \arg \min_{x \in \mathbb{R}^{N_1 \times N_2}} \left\{ D(Ax, y) + R(x) \right\}$$



LC, "Discrete total variation: New definition and minimization", 2017

LC, "A generic proximal algorithm...", 2014

Example: joint demosaicking-deblurring



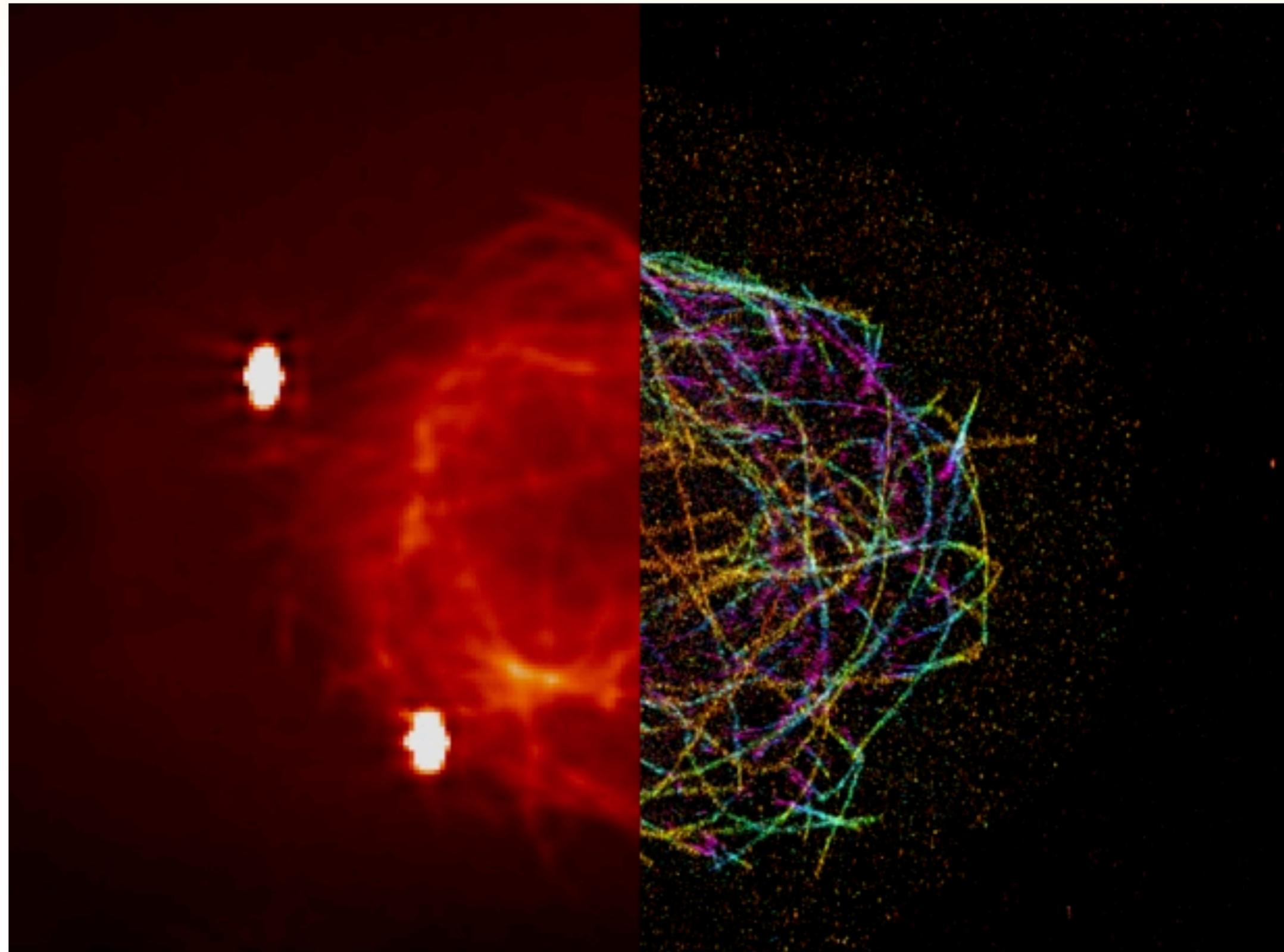
LC, "A simple, fast and efficient approach to denoising-deblurring...", 2010

Example: pansharpening/fusion



He, LC et al., "A new pansharpening method...", 2014

Super-resolution

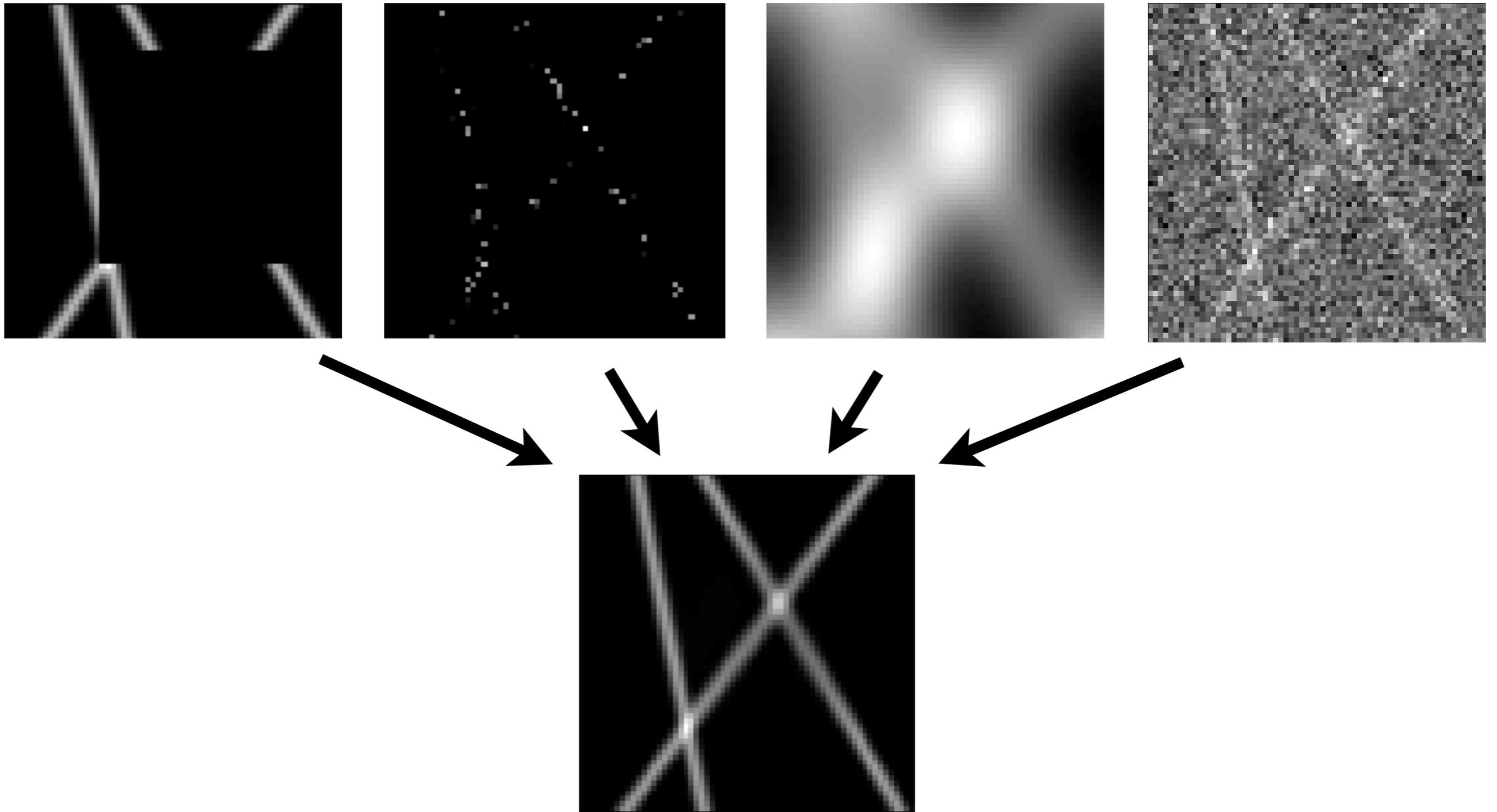


regular

PALM

[Galbraith et al.]

Example: super-resolution of lines



Polisano, LC et al., "A convex approach to superresolution...", 2019

Convex optimization

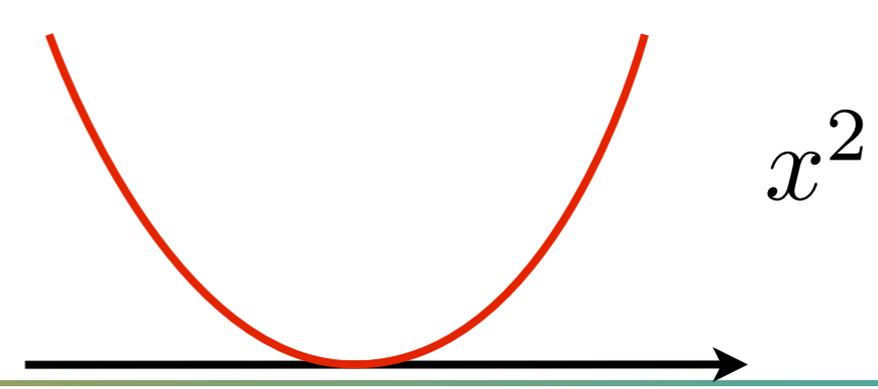
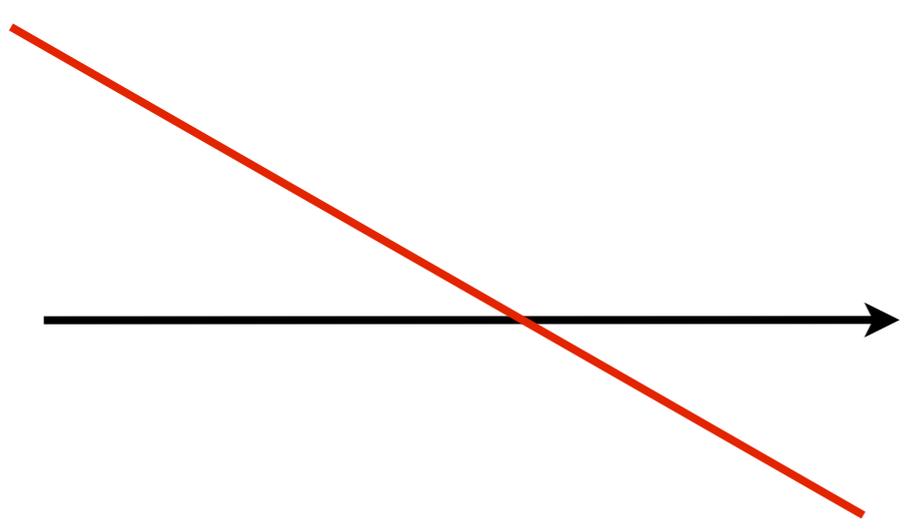
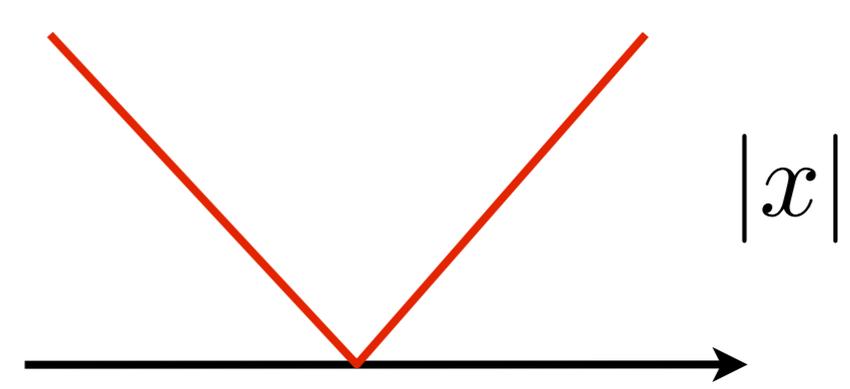
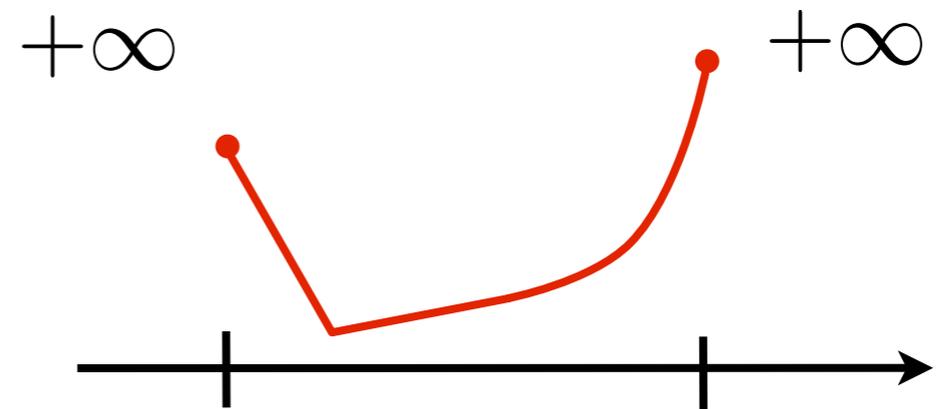
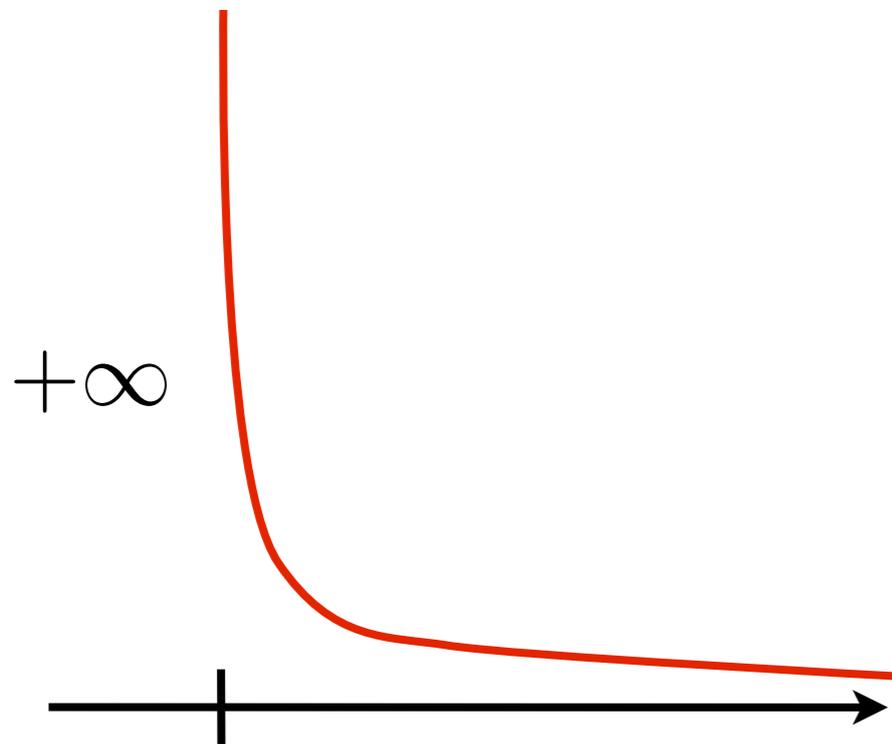
$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \psi(x) = \sum_{m=1}^M g_m(L_m x)$$

with

- linear operators $L_m : \mathcal{X} \rightarrow \mathcal{U}_m$
- real Hilbert spaces $\mathcal{X}, \mathcal{U}_m$
- **convex** functions $g_m : \mathcal{U}_m \rightarrow \mathbb{R} \cup \{+\infty\}$

Convex functions

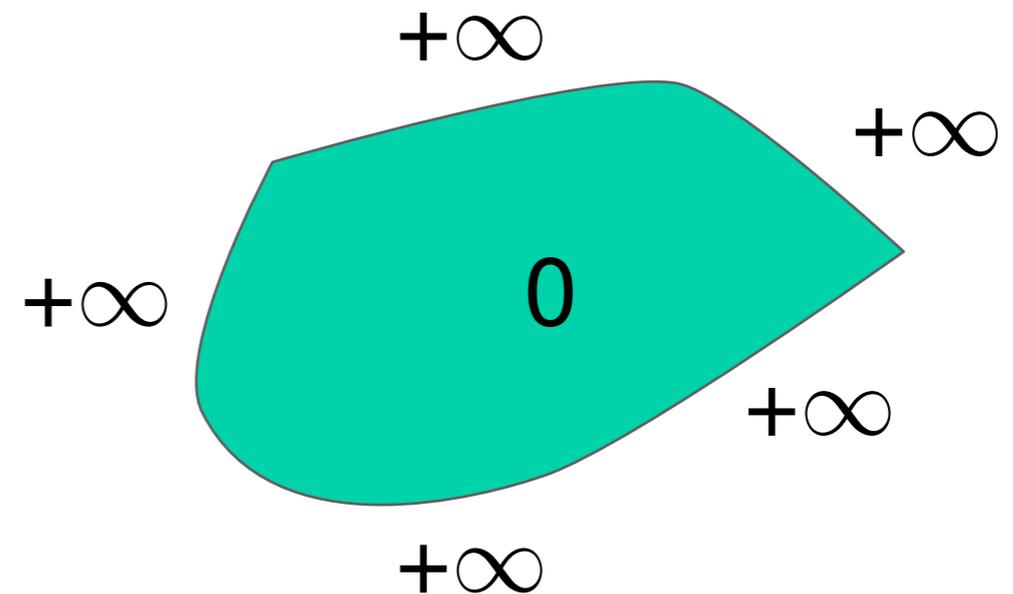
Some convex functions: $\mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$



Indicator functions

For a closed convex set $\Omega \subset \mathcal{X}$, its **indicator function** is

$$I_{\Omega}(x) = \begin{cases} 0 & \text{if } x \in \Omega, \\ +\infty & \text{else.} \end{cases}$$



I_{Ω} is convex.

Constrained optimization

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$

encompasses the presence of constraints:

$$\begin{aligned} \underset{x \in \Omega}{\text{minimize}} \ f(x) &\equiv \underset{x \in \mathcal{X}}{\text{minimize}} \ f(x) \ \text{s.t.} \ x \in \Omega \\ &\equiv \underset{x \in \mathcal{X}}{\text{minimize}} \ f(x) + I_{\Omega}(x) \end{aligned}$$

Optimization algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \Psi(x) = \sum_{m=1}^M g_m(L_m x)$$



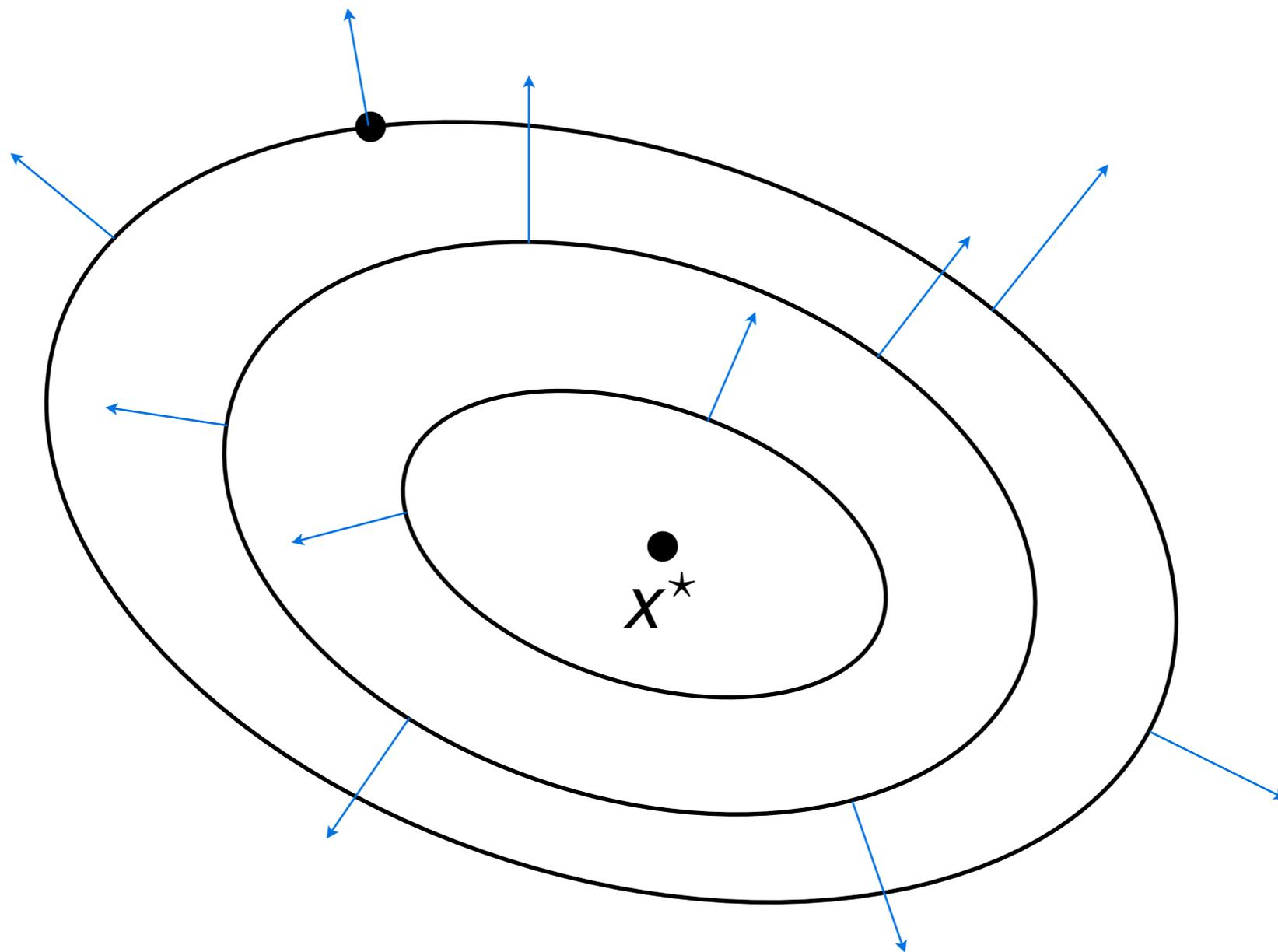
iterative algorithm computing

$$x^{k+1} = T(x^k)$$

with x^k converging to $x^* = T(x^*)$

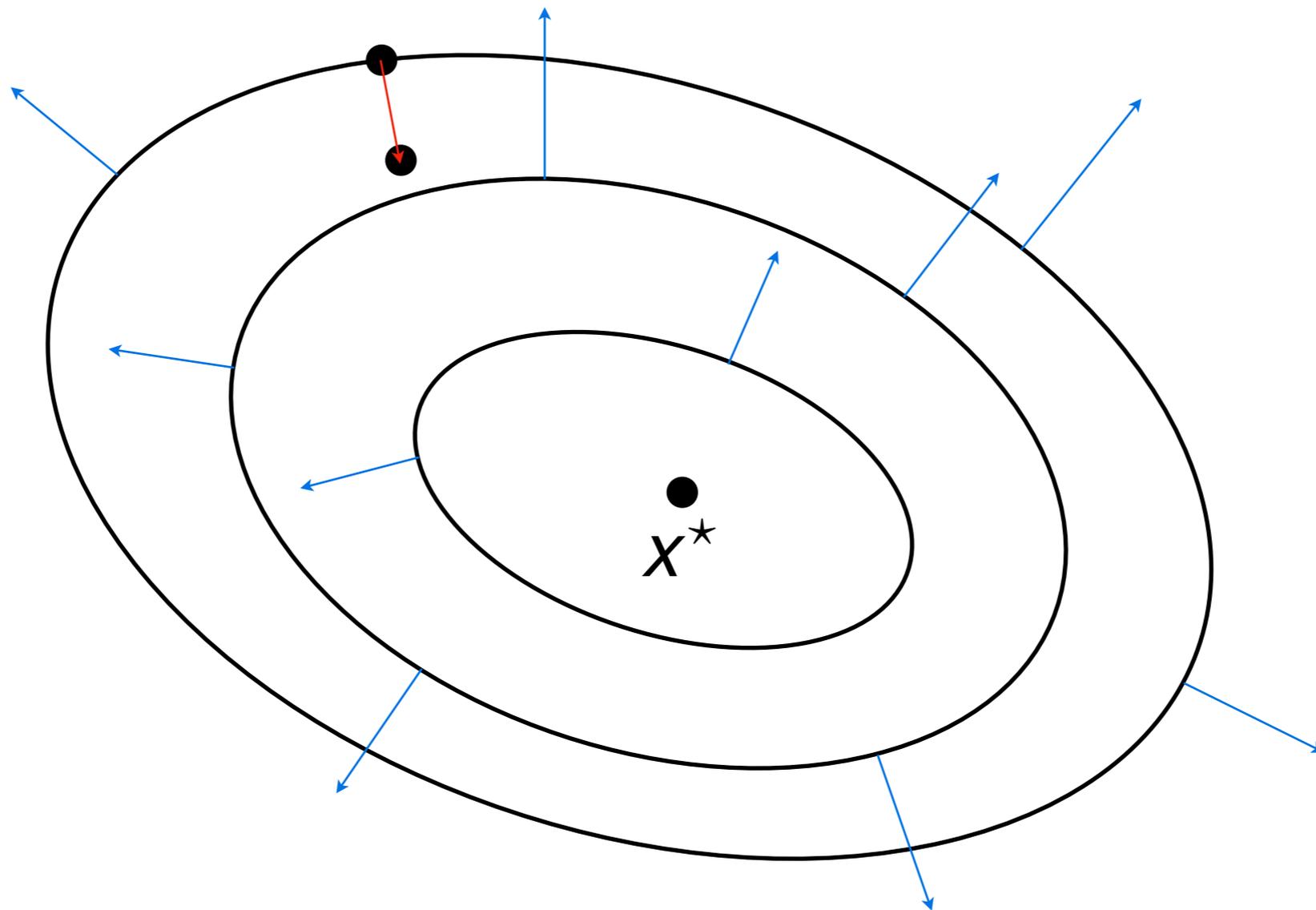
Smooth minimization: gradient descent

$$x^{k+1} = x^k - \gamma \nabla f(x^k)$$



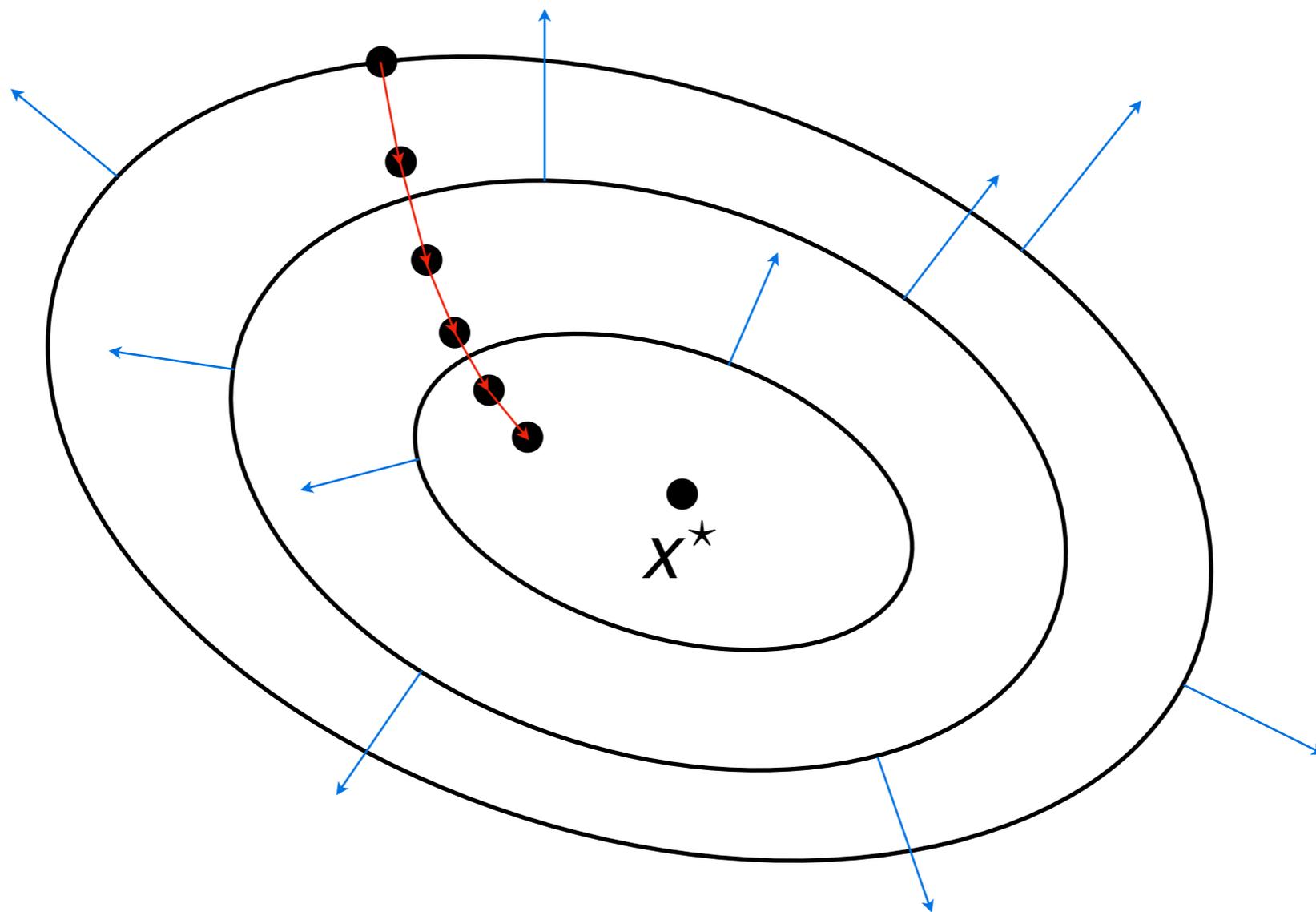
Smooth minimization: gradient descent

$$x^{k+1} = x^k - \gamma \nabla f(x^k)$$

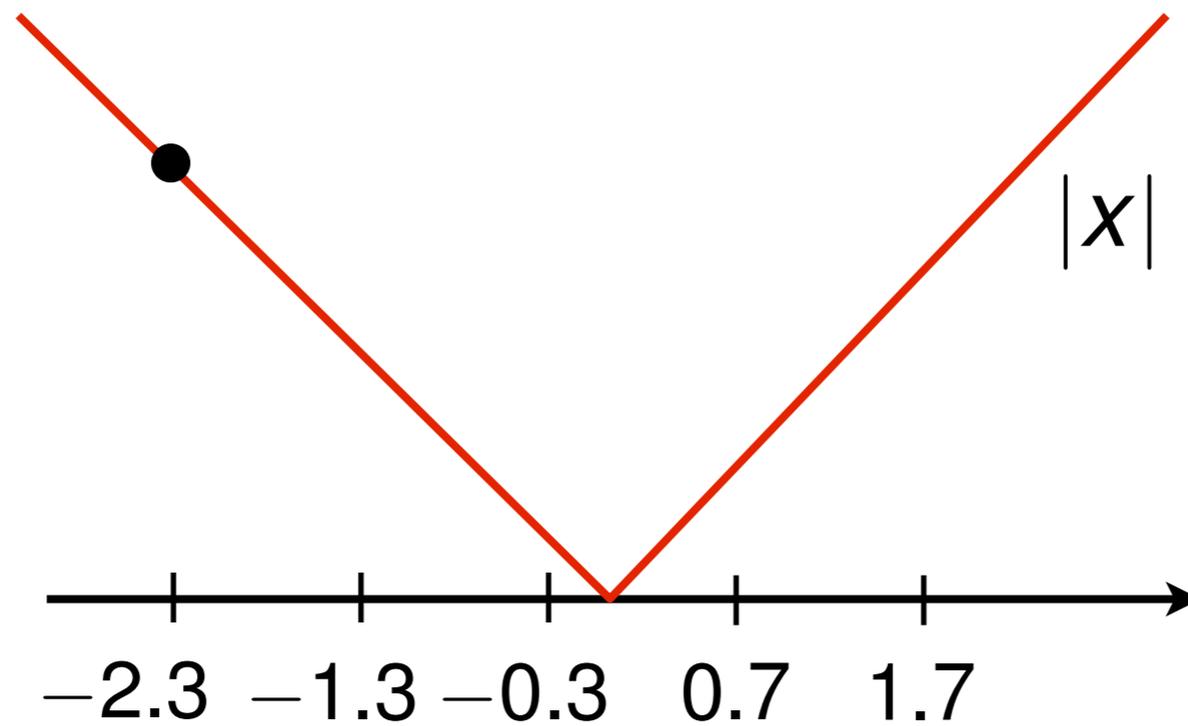


Smooth minimization: gradient descent

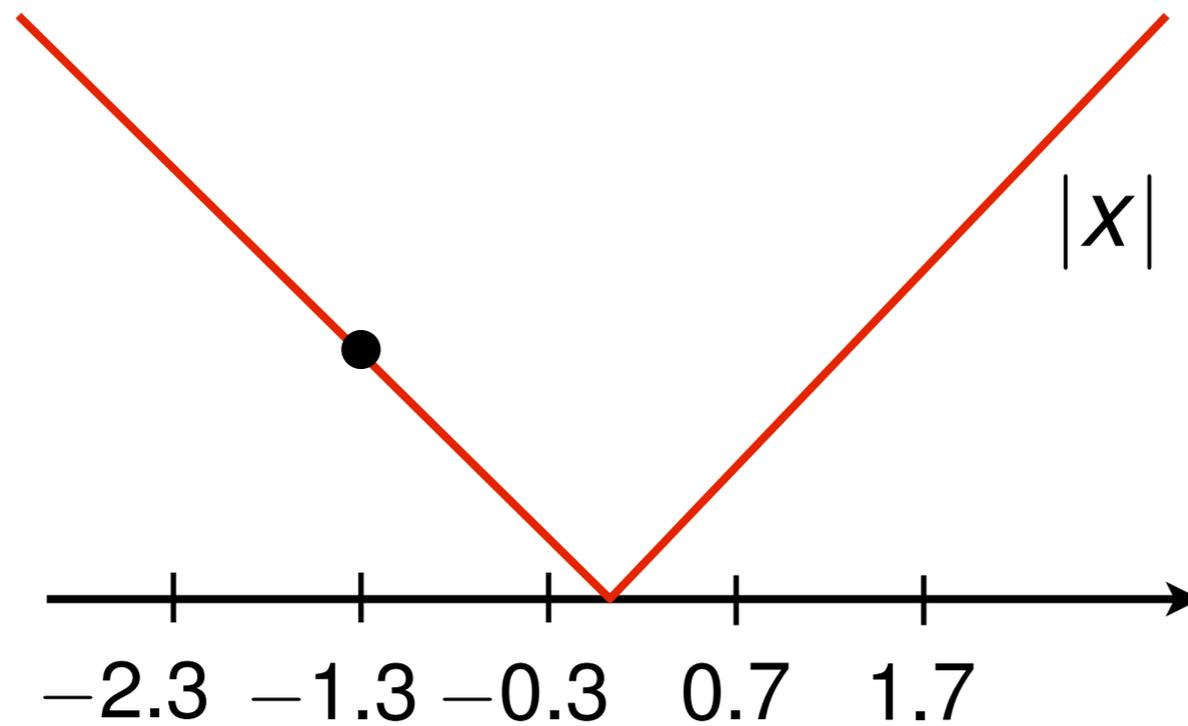
$$x^{k+1} = x^k - \gamma \nabla f(x^k)$$



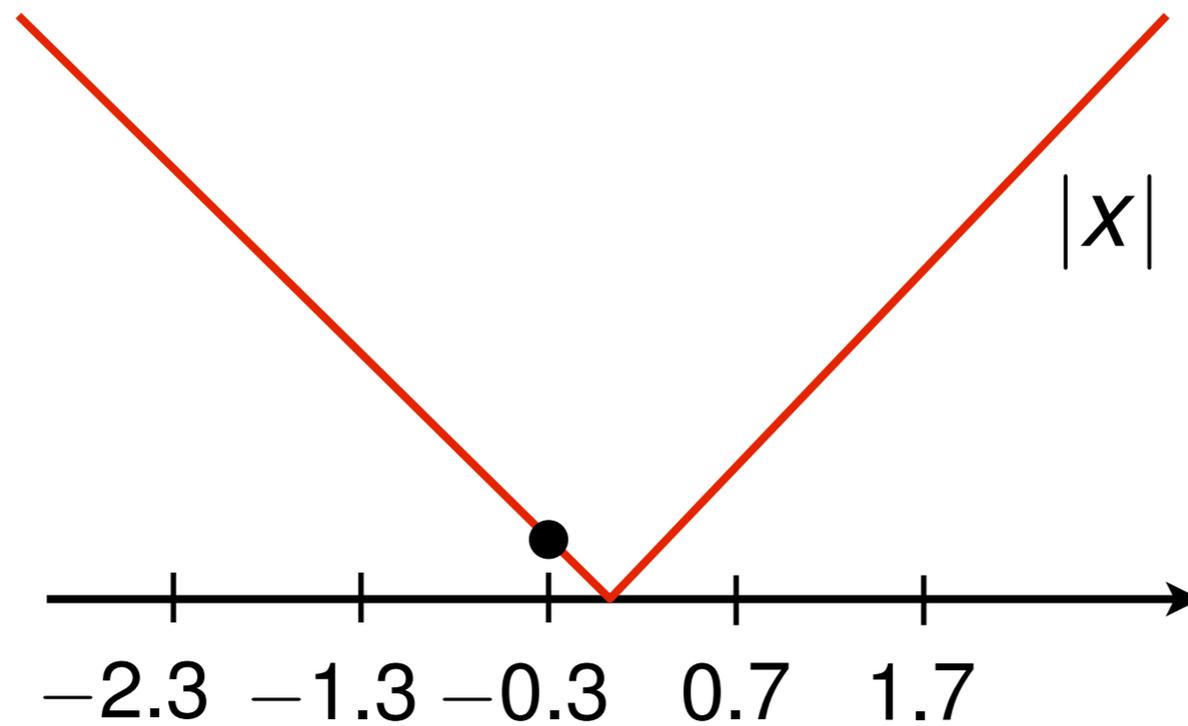
Nonsmooth minimization?



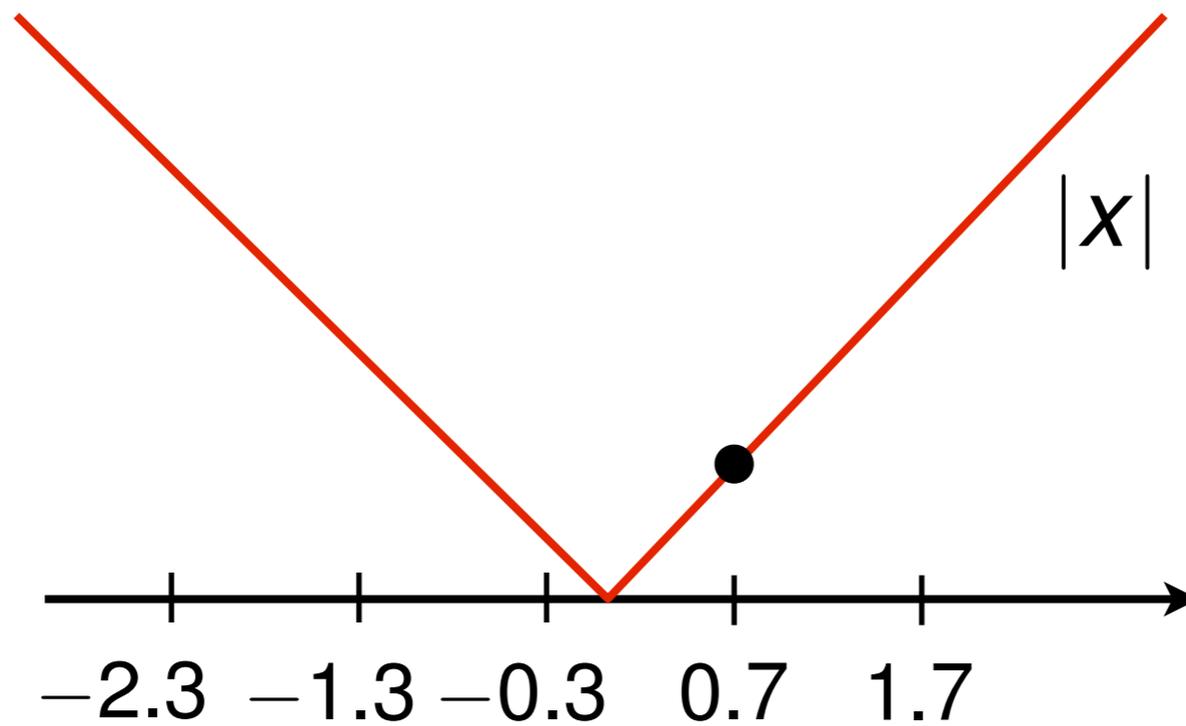
Nonsmooth minimization?



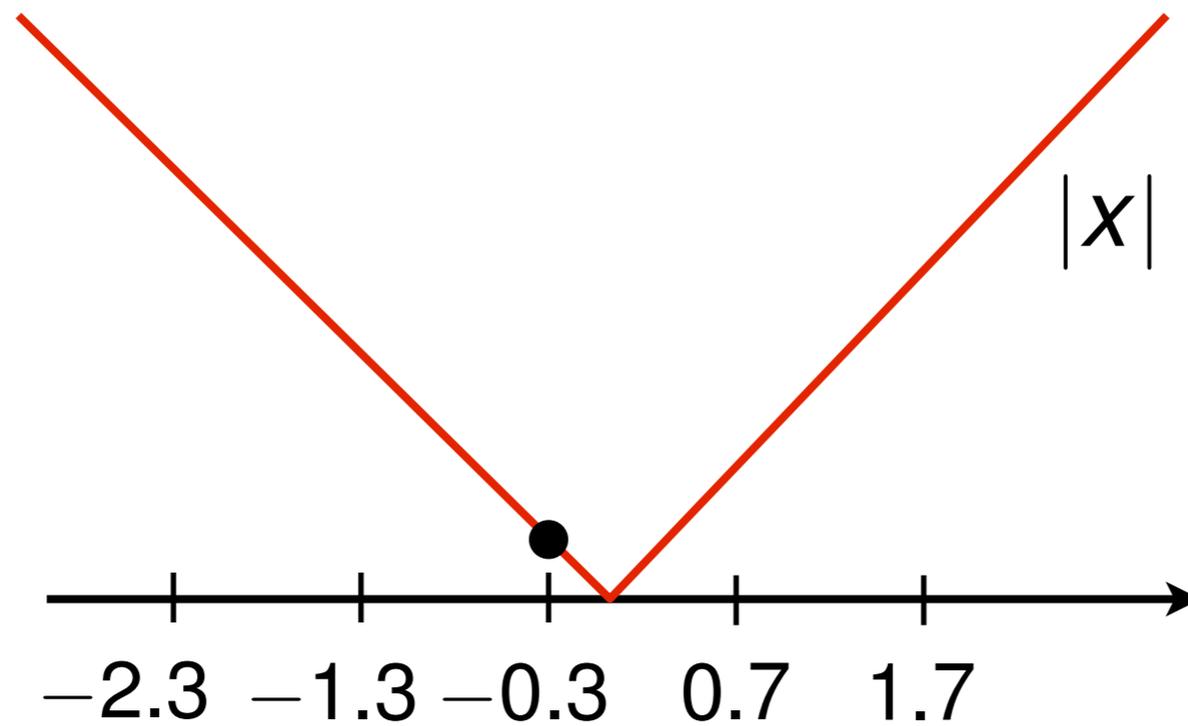
Nonsmooth minimization?



Nonsmooth minimization?

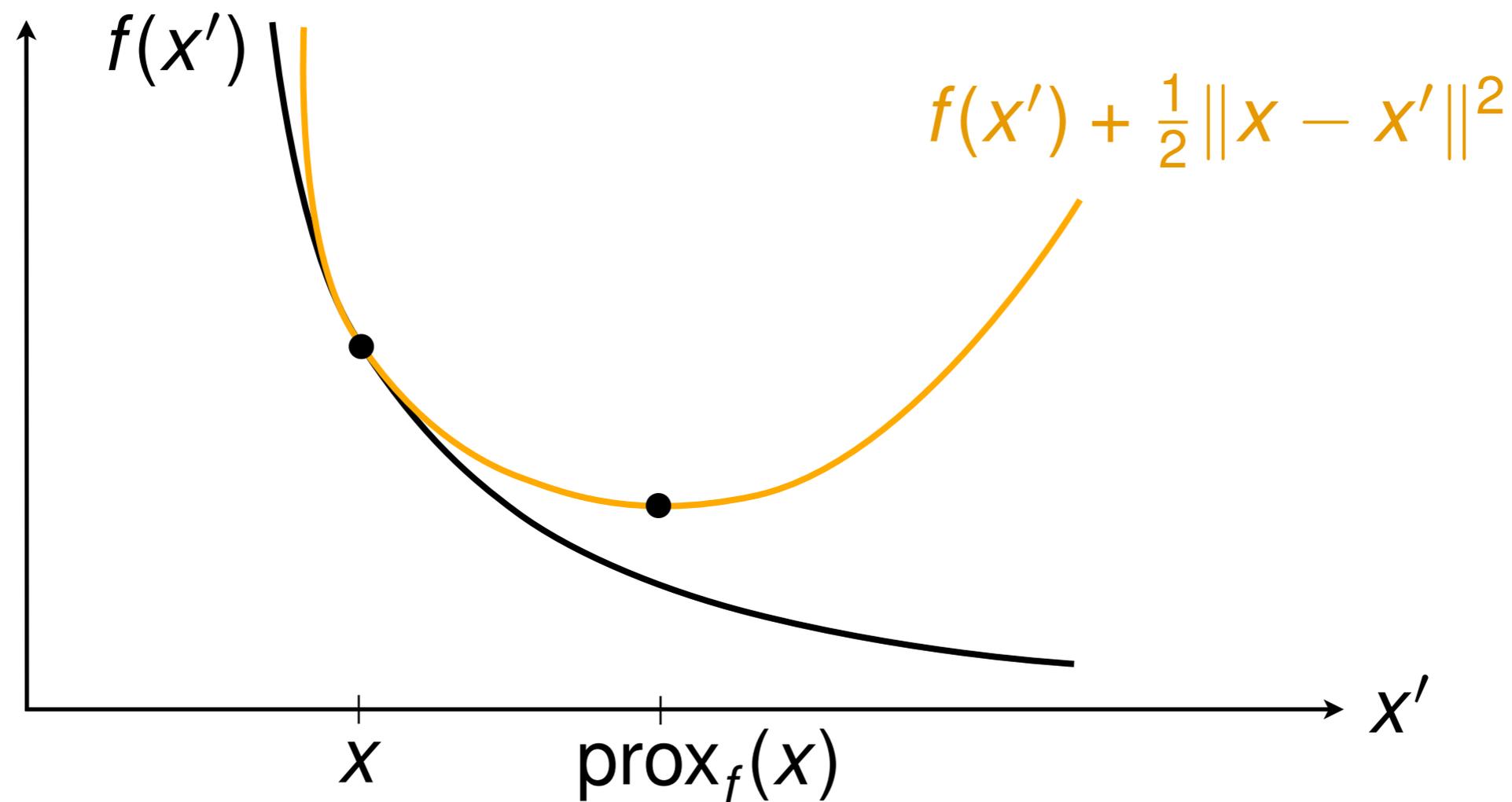


Nonsmooth minimization?



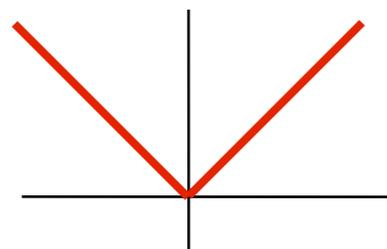
The proximity operator

$$\text{prox}_f : \mathcal{H} \rightarrow \mathcal{H} : x \mapsto \arg \min_{x' \in \mathcal{H}} f(x') + \frac{1}{2} \|x - x'\|^2$$

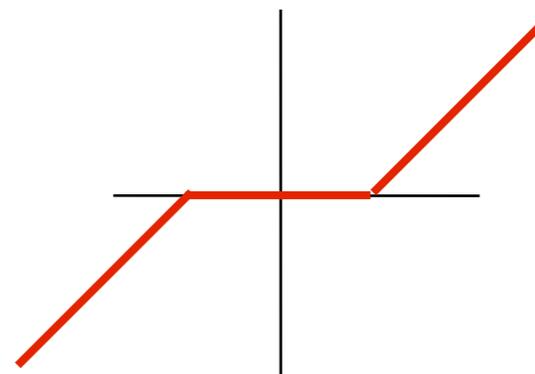


The proximity operator

$$\text{prox}_f : \mathcal{H} \rightarrow \mathcal{H} : x \mapsto \arg \min_{x' \in \mathcal{H}} f(x') + \frac{1}{2} \|x - x'\|^2$$



$$f(x) = |x|$$

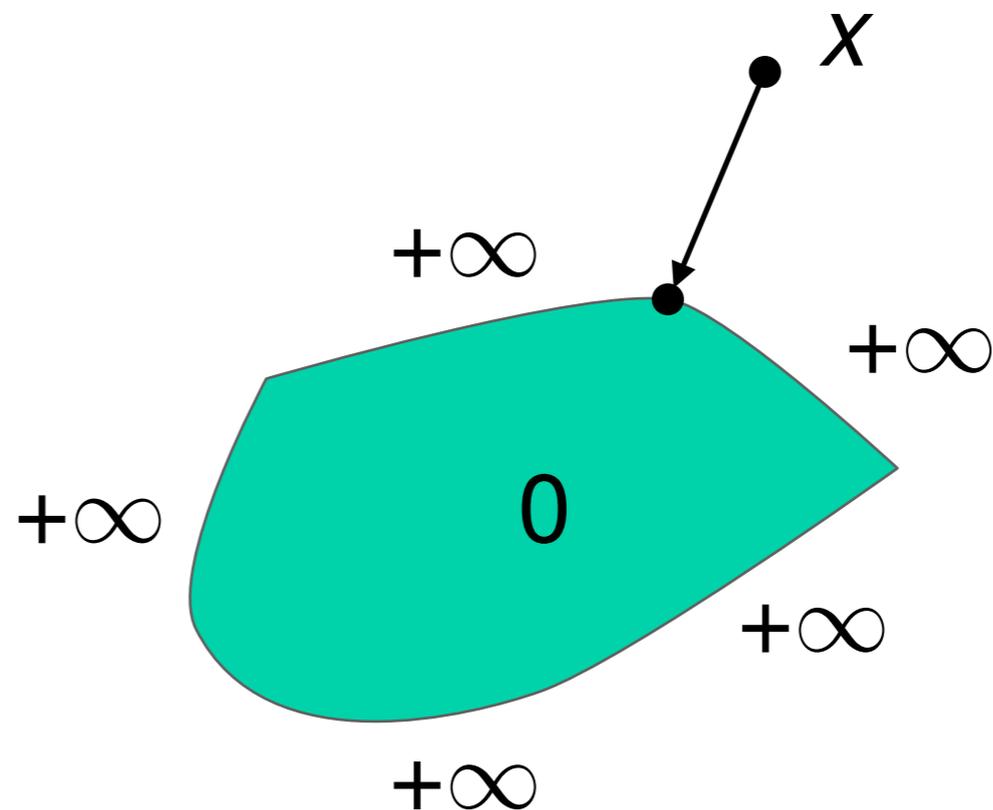


$$\text{prox}_f(x) = \text{sgn}(x) \max(|x| - 1, 0)$$

The proximity operator

$$\text{prox}_f : \mathcal{H} \rightarrow \mathcal{H} : x \mapsto \arg \min_{x' \in \mathcal{H}} f(x') + \frac{1}{2} \|x - x'\|^2$$

$$\text{prox}_{I_\Omega} = \text{proj}_\Omega$$



The proximity operator

Exact, finite time, algorithms are available to compute the proximity operator of:

- $\|X\|_* \rightarrow \text{SVD}, O(N^3)$
- 1-D TV \rightarrow taut-string alg., $O(N)$
- 2-D anisotropic TV \rightarrow graph cuts
- proj. onto the simplex $\rightarrow O(N)$

...

LC, "A direct algorithm for 1-D total variation...", 2013

LC, "Fast projection onto the simplex...", 2016

Pustelnik, LC, "Proximity operator of a sum...", 2017

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{m=1}^M g_m(L_m x)$$



No easy form of $\text{prox}_{g_1+g_2}$ or $\text{prox}_{g \circ L}$

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{m=1}^M g_m(L_m x)$$



We want **full splitting**, with individual activation of L_m , L_m^* , the gradient or proximity operator of g_m .

only fast operations in $\tilde{\mathcal{O}}(N)$, with $N =$ dimension



typically, $N \sim 10^6 - 10^9$

Proximal splitting algorithms

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + \sum_{m=1}^M g_m(L_m x) + h(x) \right)$$

with:

- h smooth with β -Lipschitz continuous gradient
→ calls to ∇h
- simple functions f and g_m → calls to $\text{prox}_{\gamma f}$ and $\text{prox}_{\sigma_m g_m}$
- calls to L_m, L_m^*

Product space trick

Find $x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$

$$g(u) = \sum_{m=1}^M g_m(u_m)$$



$$g(Lx) = \sum_{m=1}^M g_m(L_m x)$$

$$Lx = (L_1 x, \dots, L_m x)$$

Minimization of 3 functions

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \left(f(x) + g(Lx) + h(x) \right)$$

with:

- h smooth with β -Lipschitz continuous gradient
→ calls to ∇h
- simple functions f and g → calls to $\text{prox}_{\gamma f}$ and $\text{prox}_{\sigma g}$
- calls to L, L^*

Minimization of 3 functions

minimize $f + g \circ L + h$

	$f + h$		forward-backward alg.
1979	$f + g$		Douglas-Rachford alg. = ADMM
2011	$f + g \circ L$		Chambolle-Pock alg.
2013	$f + g \circ L + h$		LC, "A primal-dual splitting method for convex optimization...", 2013 Vu, "A splitting algorithm for dual monotone inclusions...", 2013

Minimization of 3 functions

minimize $f + g \circ L + h$

1979

$$f + h$$



forward-backward alg.

$$f + g$$



Douglas-Rachford alg.
= ADMM

2011

$$f + g \circ L$$



Chambolle-Pock alg.

2011

$$g \circ L + h$$



Loris-Verhoeven alg.

Combettes, LC, Pesquet, Vu, "A forward-backward view of some primal-dual optimization...", 2014

Minimization of 3 functions

minimize $f + g \circ L + h$

	$f + h$		forward-backward alg.	
1979	$f + g$		Douglas-Rachford alg. = ADMM	
2011	$f + g \circ L$		Chambolle-Pock alg.	
2011	$g \circ L + h$		Loris-Verhoeven alg.	
2013	$f + g \circ L + h$		Condat-Vu alg.	
2018			PD3O alg. (Yan)	
2020			PDDY alg.	Salim, LC et al., "Dualize, split, randomize...", 2020

4 Primal-dual algorithms

Condat–Vu algorithm form I

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma L(2x^{k+1} - x^k)) \end{cases}$$

minimize
 $f + g \circ L + h$

Condat–Vu algorithm form II

$$\begin{cases} u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

PD3O algorithm

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma L(2x^{k+1} - x^k - \gamma \nabla h(x^{k+1}) + \gamma \nabla h(x^k))) \end{cases}$$

PDDY algorithm

$$\begin{cases} u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k - \gamma L^*(u^{k+1} - u^k)) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

4 Primal-dual algorithms

Condat–Vu algorithm form I

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma L(2x^{k+1} - x^k)) \end{cases}$$

Salim, LC et al., "Dualize, split, randomize: Fast nonsmooth optimization algorithms", 2020

Condat–Vu algorithm form II

$$\begin{cases} u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

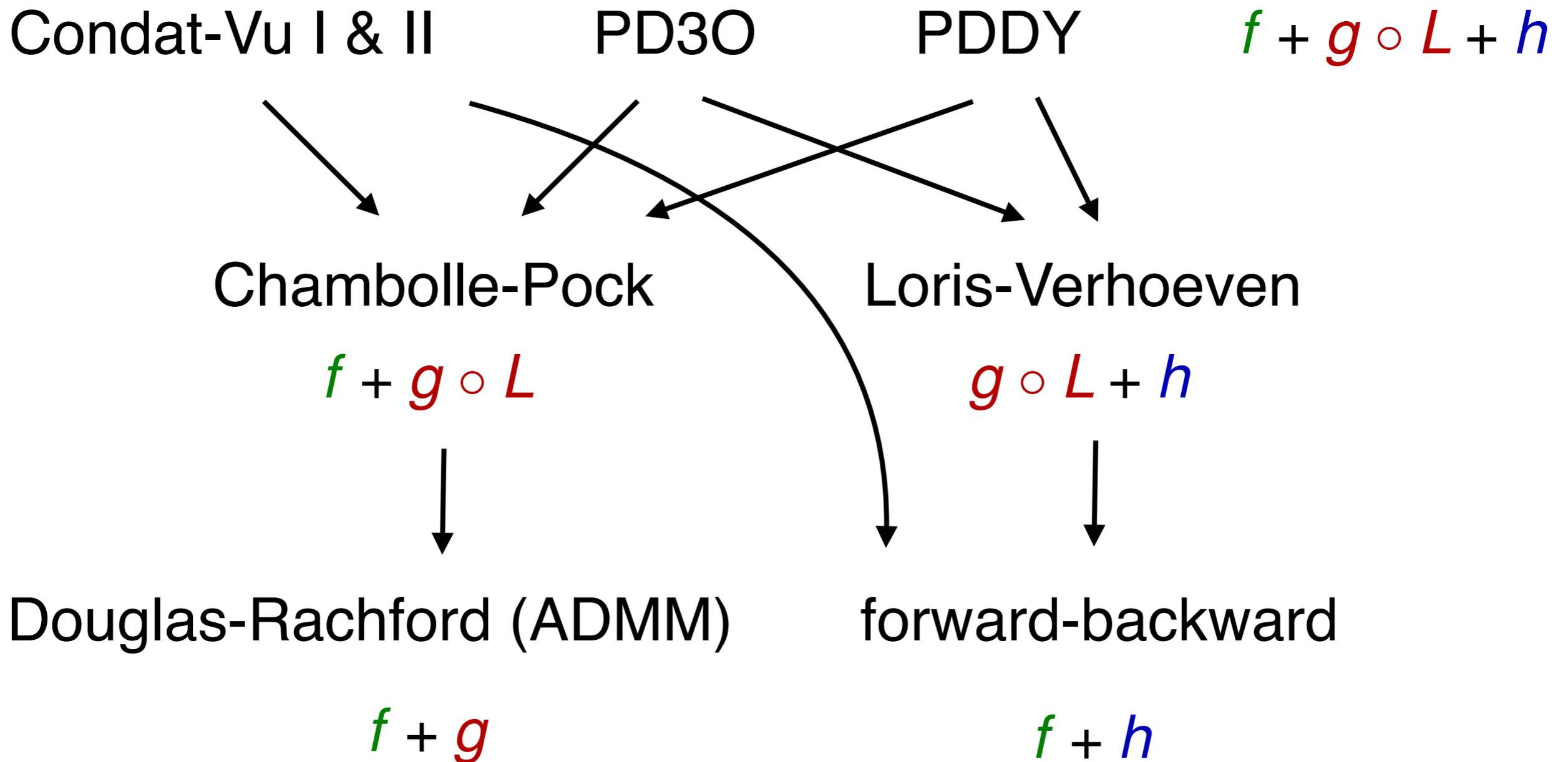
PD3O algorithm

$$\begin{cases} x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k) - \gamma L^* u^k) \\ u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma L(2x^{k+1} - x^k - \gamma \nabla h(x^{k+1}) + \gamma \nabla h(x^k))) \end{cases}$$

PDDY algorithm

$$\begin{cases} u^{k+1} = \text{prox}_{\sigma g^*} (u^k + \sigma Lx^k) \\ x^{k+1} = \text{prox}_{\gamma f} (x^k - \gamma \nabla h(x^k - \gamma L^*(u^{k+1} - u^k)) - \gamma L^*(2u^{k+1} - u^k)) \end{cases}$$

Primal-dual algorithms



Primal-dual algorithms

Condat-Vu I & II

PD3O

PDDY

$$f + g \circ L + h$$

Chambolle-Pock

Loris-Verhoeven

$$f + g \circ L$$

$$g \circ L + h$$

Douglas-Rachford (ADMM)

$$f + g$$

LC et al. "Proximal Splitting Algorithms for Convex Optimization: A Tour of Recent Advances, with New Twists", 2019

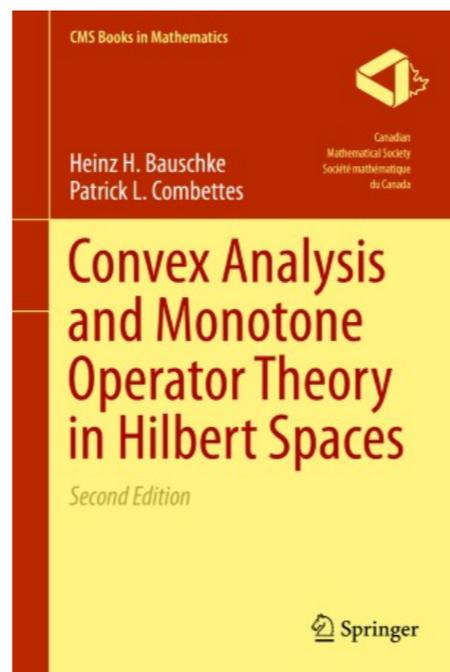
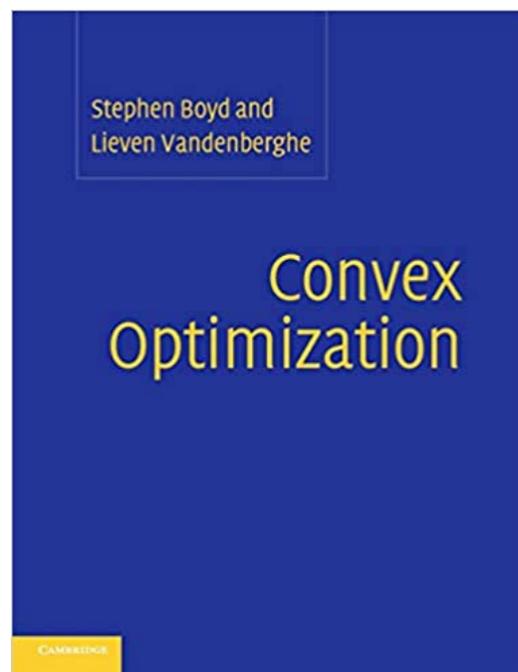
Principle

convex analysis
monotone and fixed-point
operator theory

allow us
to design

can be
analyzed
with

algorithms for large-scale
nonsmooth optimization



LC et al. "Proximal Splitting Algorithms for Convex Optimization: A Tour of Recent Advances, with New Twists", 2019

Summary



nonsmooth functions



large scale



proximal splitting algorithms

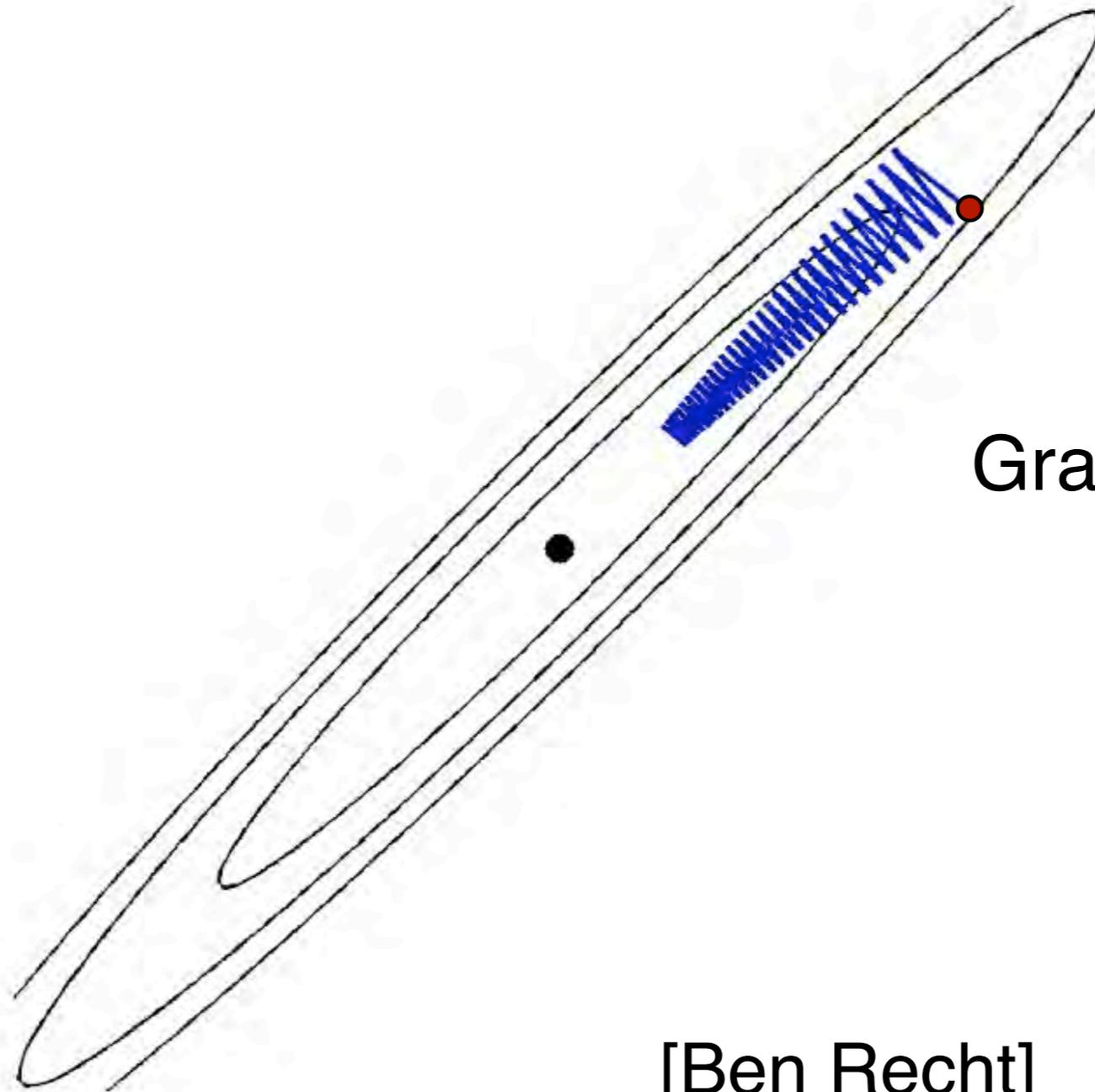


Condat-Vu, PD3O, PDDY: versatile primal-dual algorithms

LC et al. "Proximal Splitting Algorithms for Convex Optimization: A Tour of Recent Advances, with New Twists", 2019



Speed?



Gradient descent



[Ben Recht]

Convergence rates

Theorem – PD3O, with g continuous around Lx^* :

$$\psi(x^k) - \psi(x^*) = o(1/\sqrt{k})$$

Theorem – accelerated PD3O and PDDY
when h or f strongly convex, with varying stepsizes:

$$\|x^k - x^*\|^2 = O(1/k^2)$$

Theorem – linear convergence of PD3O and PDDY
when h or f strongly convex and g smooth:

$$\|x^k - x^*\|^2 \leq (1 - \rho)^k c_0$$

LC et al. "Distributed Proximal Splitting Algorithms with Rates and Acceleration", 2022

Minimization with an affine constraint

$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} h(x) \quad \text{s.t. } Lx = y$$



h strongly convex
linear convergence

Salim, LC et al., "Dualize, split, randomize: Fast nonsmooth optimization algorithms", 2020

optimal acceleration:

Salim, LC et al., "An Optimal Algorithm for Strongly Convex Minimization under Affine Constraints", AISTATS 2022

Decentralized setting

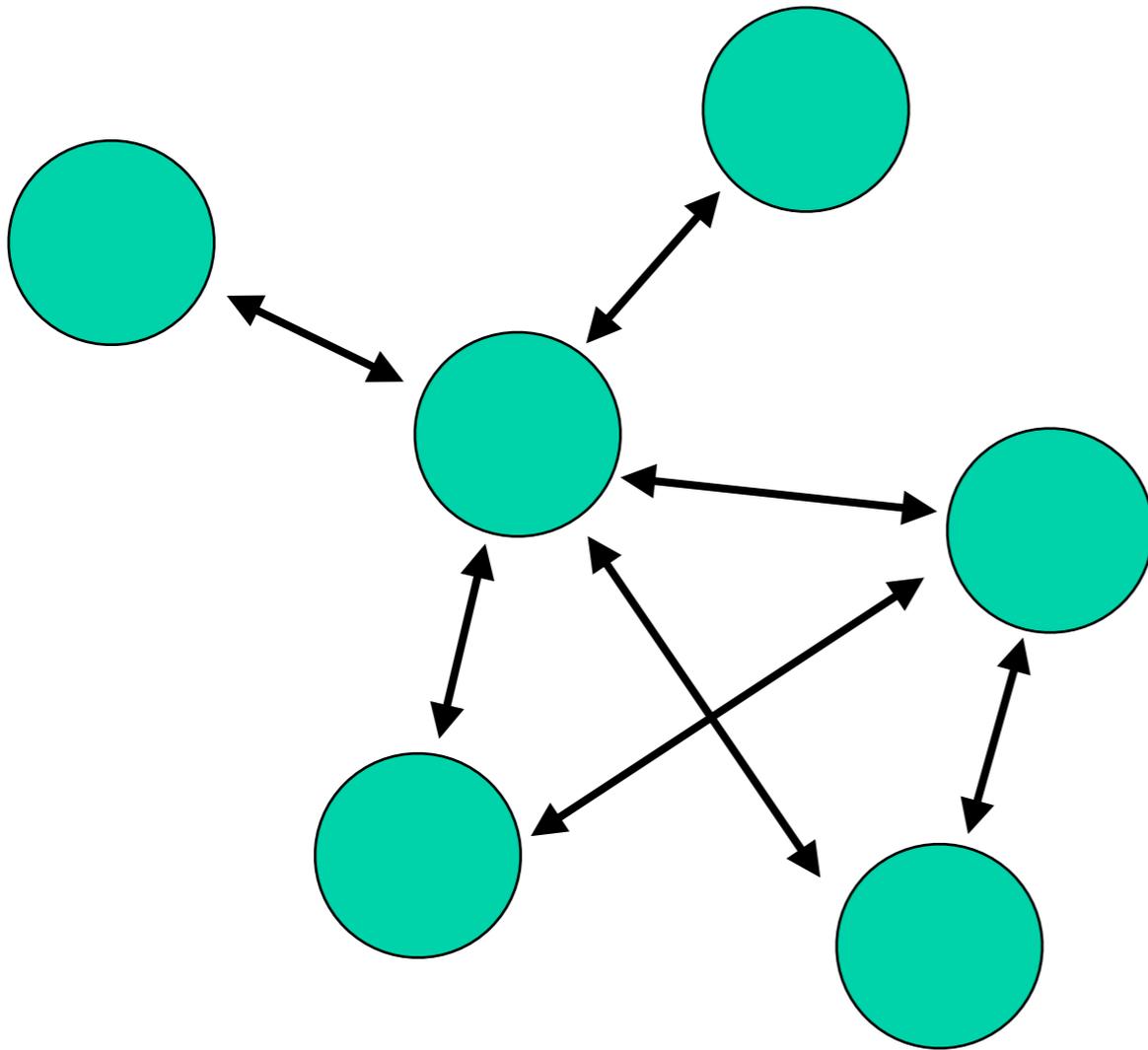
$$\text{Find } x^* \in \arg \min_{x \in \mathcal{X}} \sum_{m=1}^M h_m(x_m) \quad \text{s.t. } x_1 = \dots = x_M$$

DESTROY algorithm:

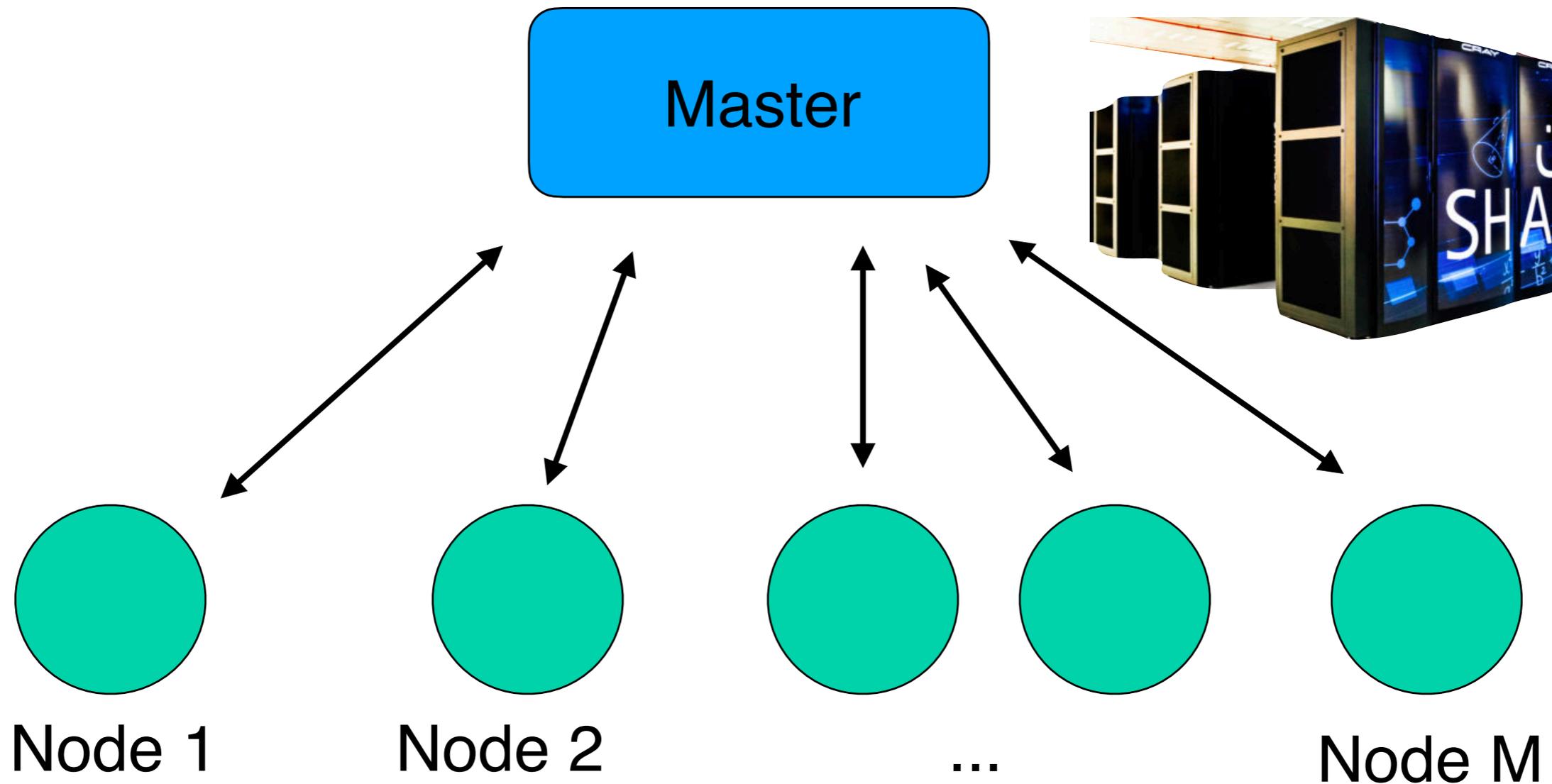
Salim, LC et al., "Dualize, split, randomize: Fast nonsmooth optimization algorithms", 2020

optimal acceleration:

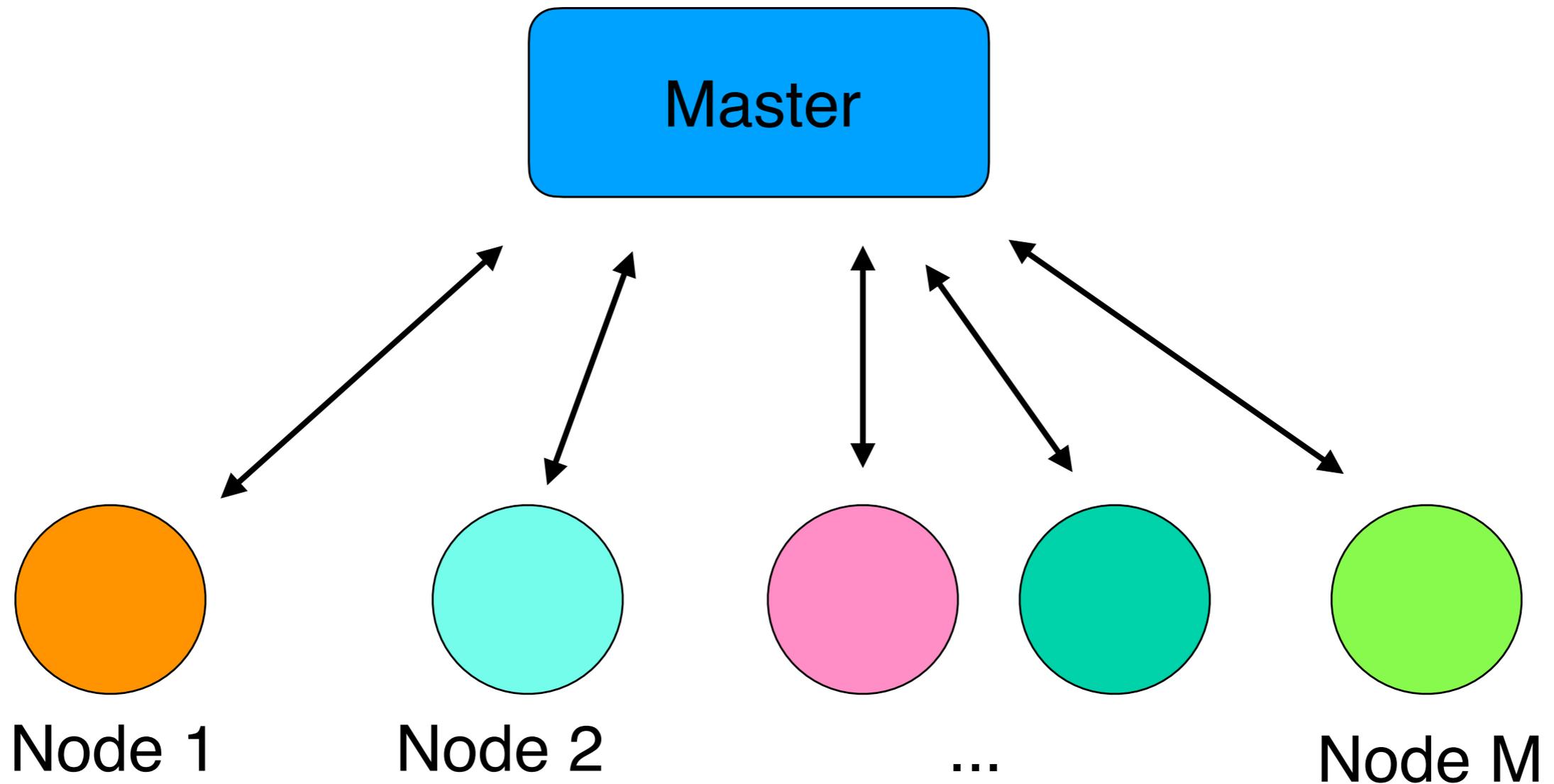
Salim, LC et al., "An Optimal Algorithm for Strongly Convex Minimization under Affine Constraints", AISTATS 2022



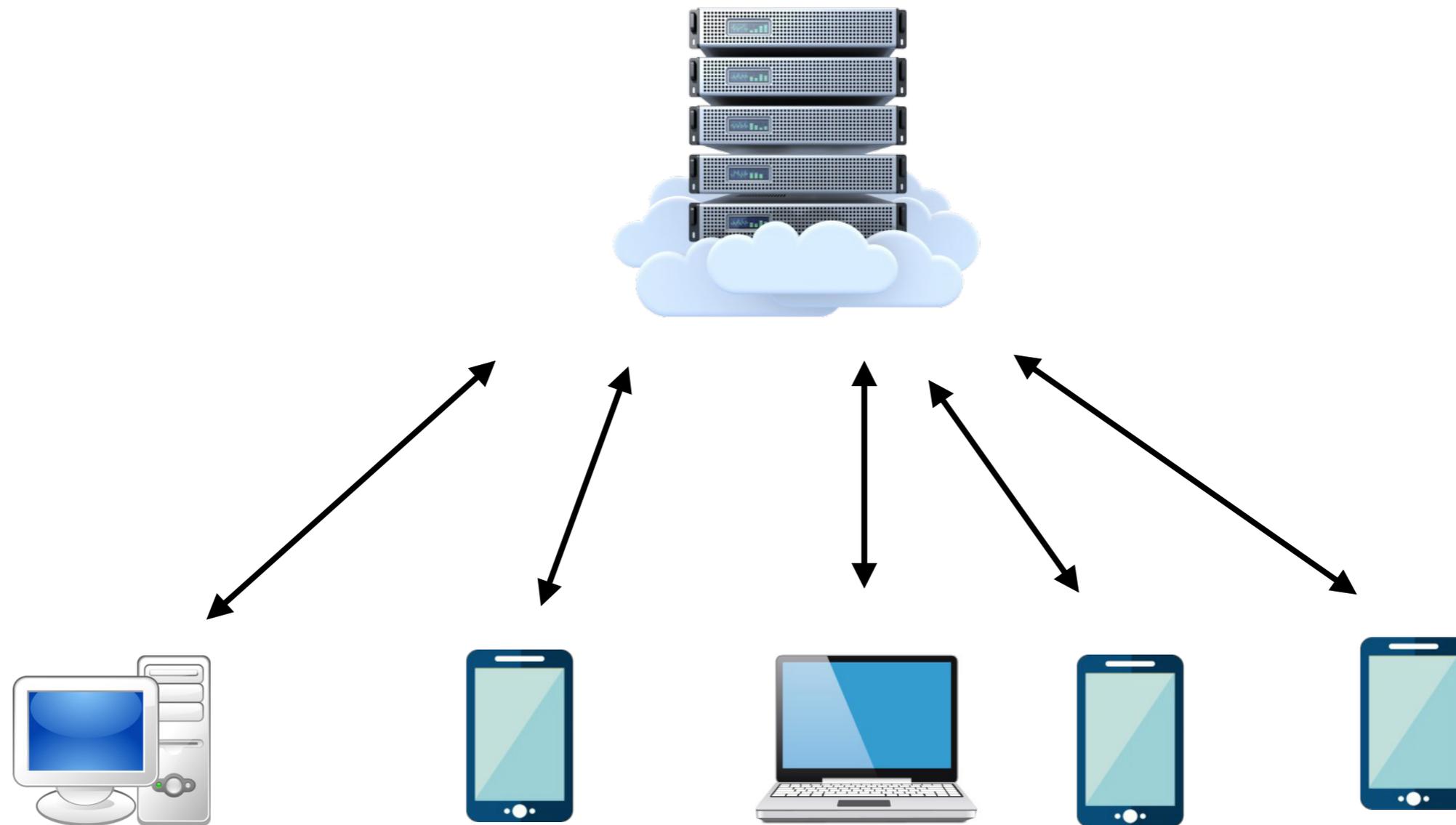
Distributed optimization



Distributed optimization



Federated learning



Distributed algorithms

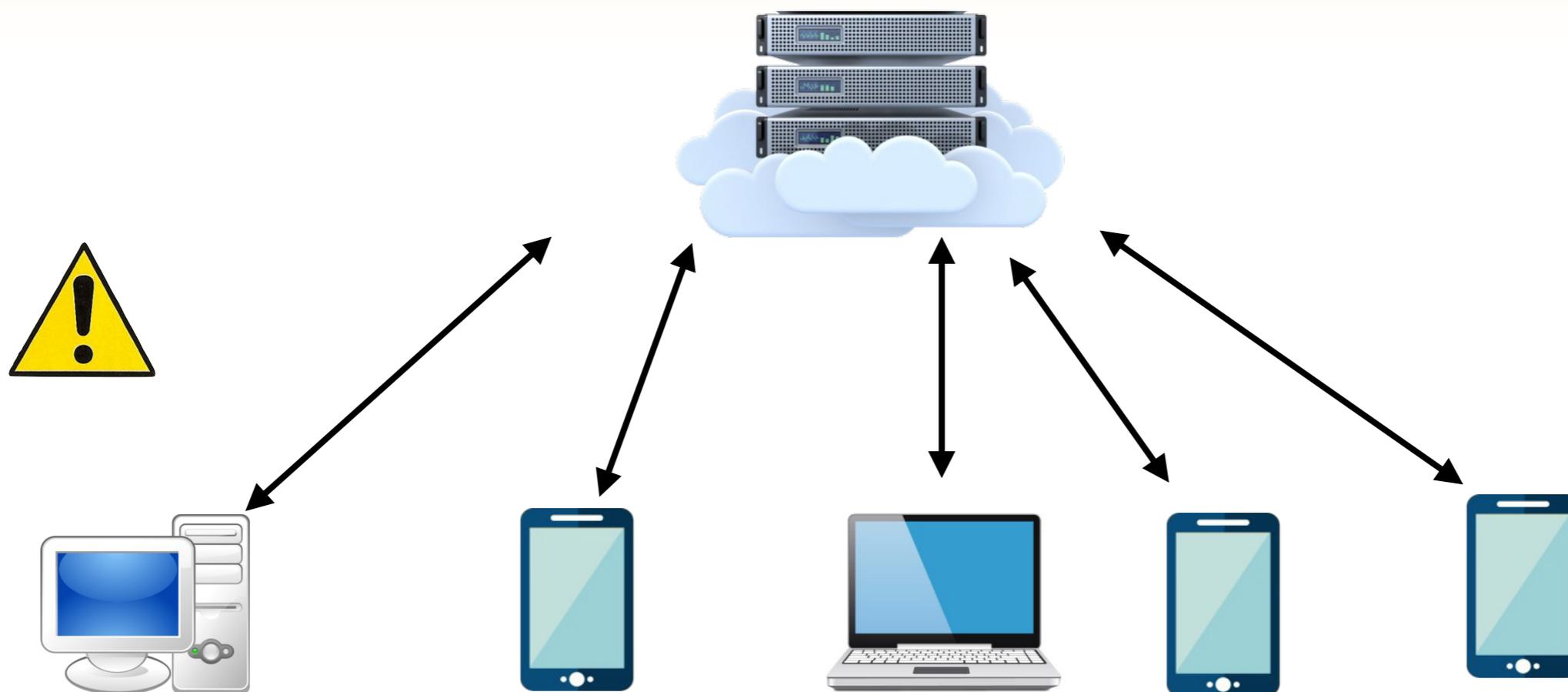
$$\underset{x \in \mathcal{X}}{\text{minimize}} \quad f(x) + \frac{1}{M} \sum_{m=1}^M g_m(L_m x) + \frac{1}{M} \sum_{m=1}^M h_m(x),$$

Distributed PDDY Algorithm

input: $(\gamma_k)_{k \in \mathbb{N}}, \eta \geq \|\widehat{L}\|^2, (\omega_m)_{m=1}^M,$
 $x_f^0 \in \mathcal{X}, (u_m^0)_{m=1}^M \in \widehat{\mathcal{U}}$
initialize: $p_m^0 := L_m^* u_m^0, m = 1, \dots, M$
for $k = 0, 1, \dots$ **do**
 at all nodes, for $m = 1, \dots, M,$ **do**
 $u_m^{k+1} := \text{prox}_{M\omega_m g_m^*/(\gamma_k \eta)}(u_m^k$
 $+ \frac{M\omega_m}{\gamma_k \eta} L_m x_f^k)$
 $p_m^{k+1} := L_m^* u_m^{k+1}$
 $x_m^{k+1} := x_f^k - \frac{\gamma_k}{M\omega_m} (p_m^{k+1} - p_m^k)$
 $a_m^k := M\omega_m x_m^{k+1} - \gamma_{k+1} \nabla h_m(x_m^{k+1})$
 $- \gamma_{k+1} p_m^{k+1}$
 transmit a_m^k to master
 at master, **do**
 $x_f^{k+1} := \text{prox}_{\gamma_{k+1} f}(\frac{1}{M} \sum_{m=1}^M a_m^k)$
 broadcast x_f^{k+1} to all nodes
end for

LC et al. "Distributed Proximal Splitting Algorithms with Rates and Acceleration", 2022

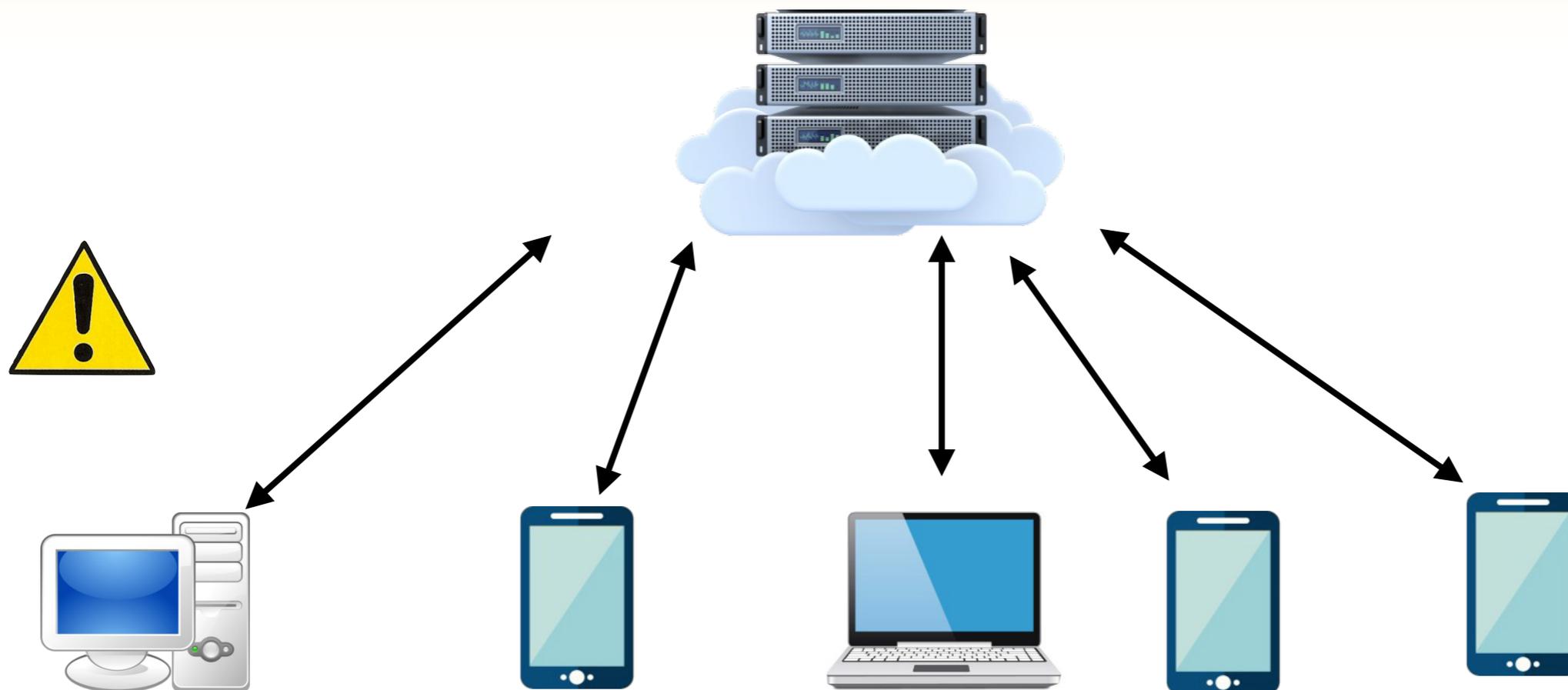
Communication bottleneck



less communication with **local steps**

Malinovsky et al., "From local SGD to local fixed point methods for federated learning", ICML 2020

Communication bottleneck



compression

Albasyoni et al. "Optimal Gradient Compression for Distributed and Federated Learning", 2020

LC and Richtárik, "MURANA: A Generic Framework for Stochastic Variance-Reduced Optimization", 2021

LC et al., "EF-BV: A unified theory of error feedback...", 2022

Summary



nonsmooth large-scale optimization



proximal splitting algorithms



speed



- ▶ (math.) acceleration
- ▶ cheaper iterations (computations, communication)



Perspectives

Beyond SGD: randomized fixed-point algorithms:
random/inexact/cheap estimates of operators to call
or variables/coordinates to update

Acceleration of fixed-point algorithms

Applications to communication networks, data processing
on graphs, federated learning