

# Exact Verification of First-order Methods via Mixed-Integer Linear Programming

Bartolomeo Stellato

Department of Operations Research and Financial Engineering

Department of Electrical and Computer Engineering

Department of Computer Science



PRINCETON  
UNIVERSITY









**Most applications require fast and effective  
decisions in real-time**



# Real-time optimization can help us

$$\begin{array}{l} \text{minimize} \\ \text{subject to} \end{array} \quad \begin{array}{l} \text{decisions} \\ \downarrow \\ f(z, x) \\ z \in C(x) \\ \uparrow \\ \text{parameters} \end{array}$$

**objective**  $f$ : energy consumption, costs  
**constraints**  $C$ : dynamics, physical limits

re-planning in real-time  
is the key to effective  
decision-making

How do we solve  
such problems?



# ...and they can solve many constrained convex problems!

## Linear Programs



**PDLP**

Applegate, Díaz, Hinder, Lu, Lubin,  
O'Donoghue, Schudy (2021)

## Quadratic Programs



**OSQP**

Stellato, Banjac, Goulart,  
Bemporad, Boyd (2020)

## Conic Programs



**SCS**  
SPLITTING CONIC SOLVER

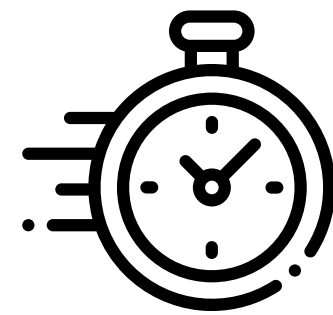
O'Donoghue, Chu, Parikh, Boyd (2016)

many more solvers available by the day: cuOPT, PDCS, ...

# But they can converge slowly

major issue in safety-critical applications with

real-time requirements



limited computing power



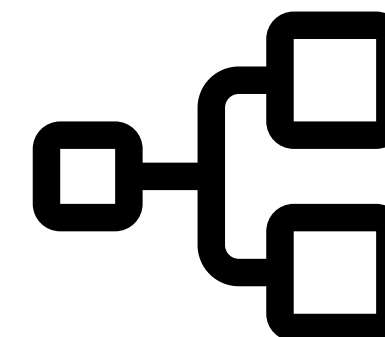
## 💡 main idea

in most applications we repeatedly solve the **same problem** with **varying parameters**

minimize  $f(z, x)$   
subject to  $z \in C(x)$



large amount of data  
(e.g., instances, solutions)



structured problems  
(e.g., parameters -> solutions)

# Performance verification of first-order methods



# Convergence of first-order methods

iterations

$$z^{k+1} = T(z^k, x) \quad \text{for } k = 0, 1, \dots$$

$T$  operator  
(e.g., contractive, averaged)

goal: find fixed-points

$$z^* = T(z^*, x)$$

example  
gradient descent

problem  $\longrightarrow$  optimality conditions

$$\text{minimize } f(z, x) \longrightarrow \nabla f(z^*, x) = 0$$

iterations

$$z^{k+1} = z^k - \theta \nabla f(z^k, x)$$

fixed-points

$$z^* = z^* - \theta \nabla f(z^*, x)$$

same as

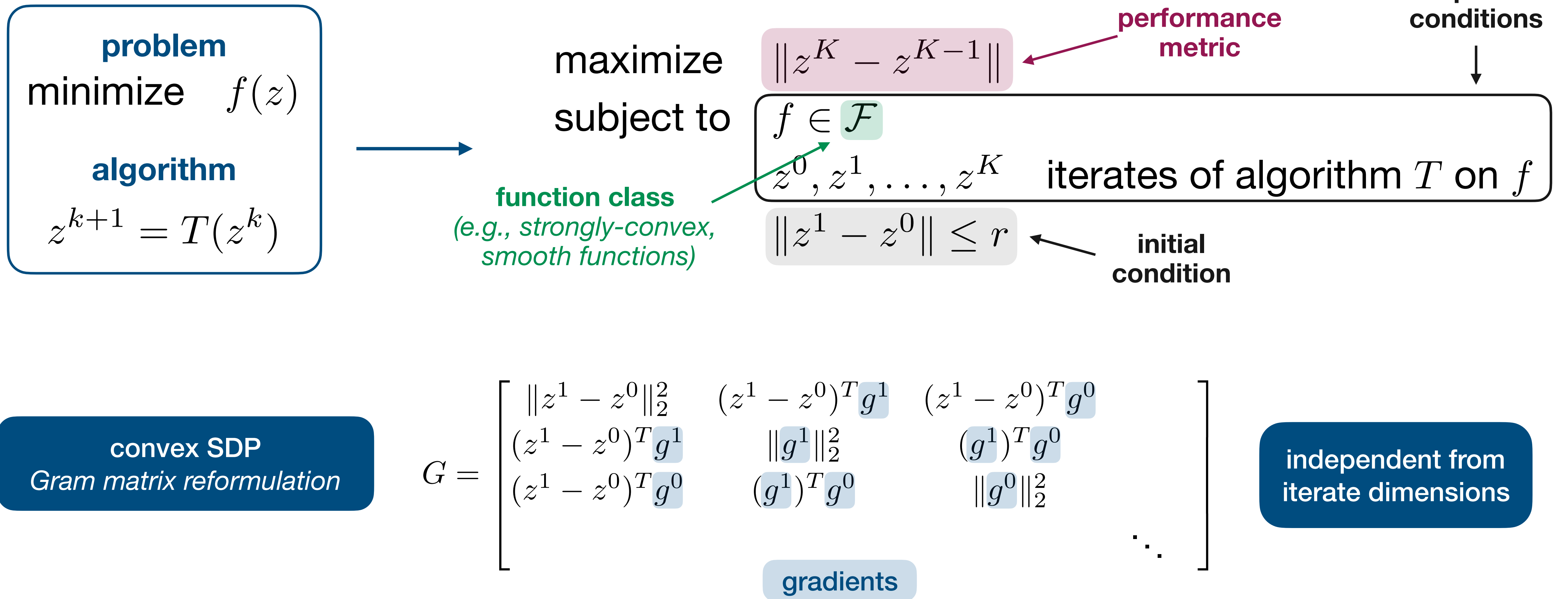
performance metric

$$r^k(x) = \|T(z^{k-1}) - z^{k-1}\| = \|z^k - z^{k-1}\|$$

fixed-point residual  
(converges to 0)



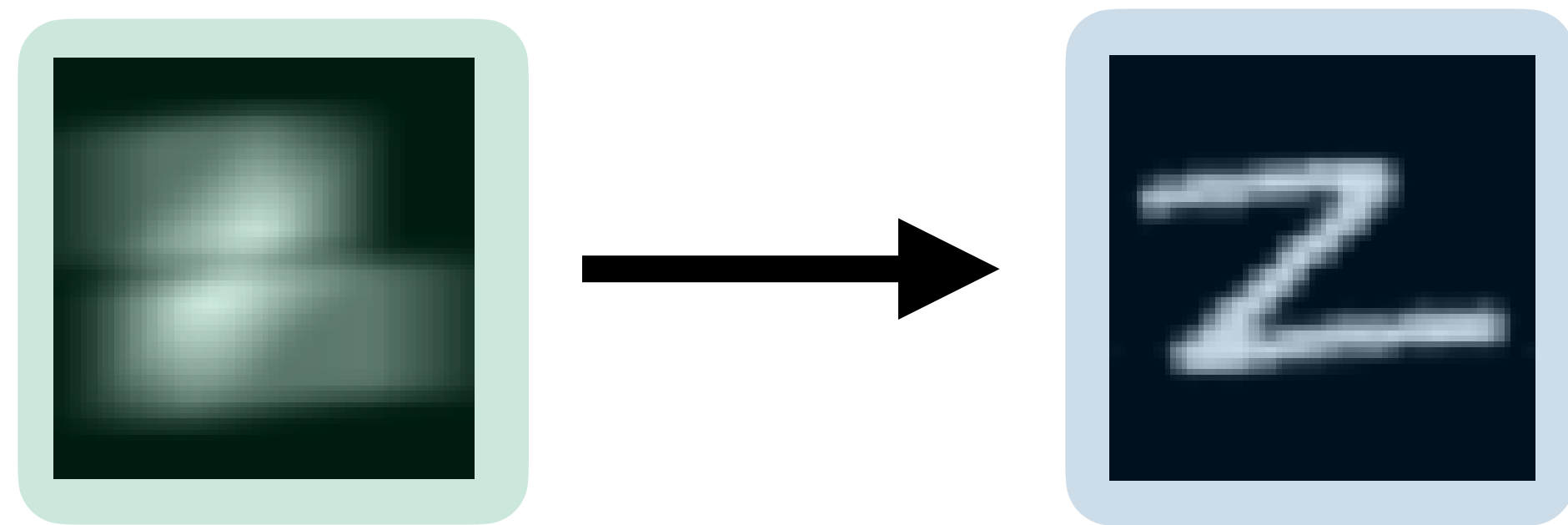
# Best known convergence bounds via Performance Estimation





# Classical worst-case convergence bounds can be very loose

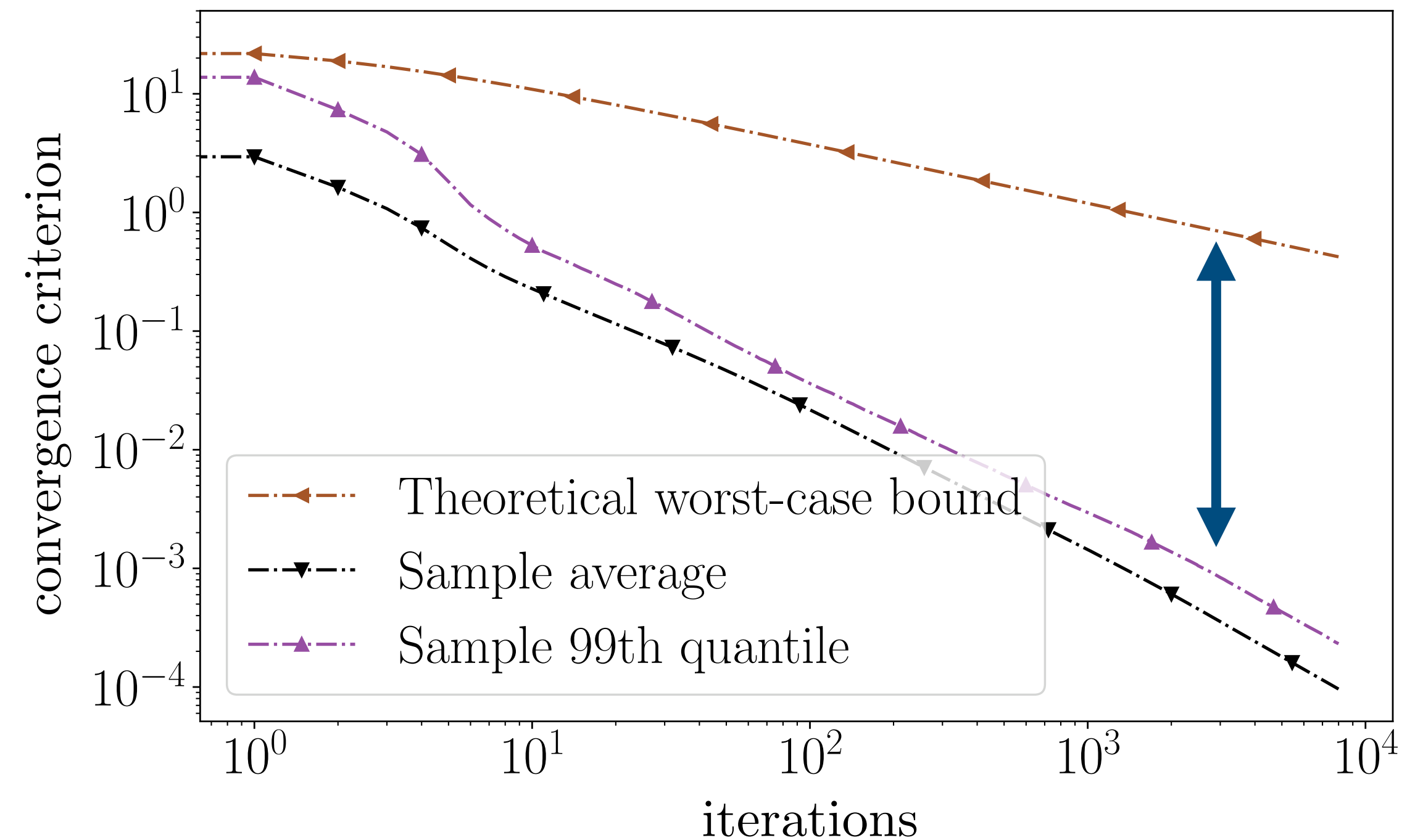
image deblurring problem  
*emnist dataset*



minimize  $\|Az - x\|_2^2 + \lambda \|z\|_1$   
subject to  $0 \leq z \leq 1$

deblurred image (pointing to  $z$ )

blurred image (pointing to  $x$ )

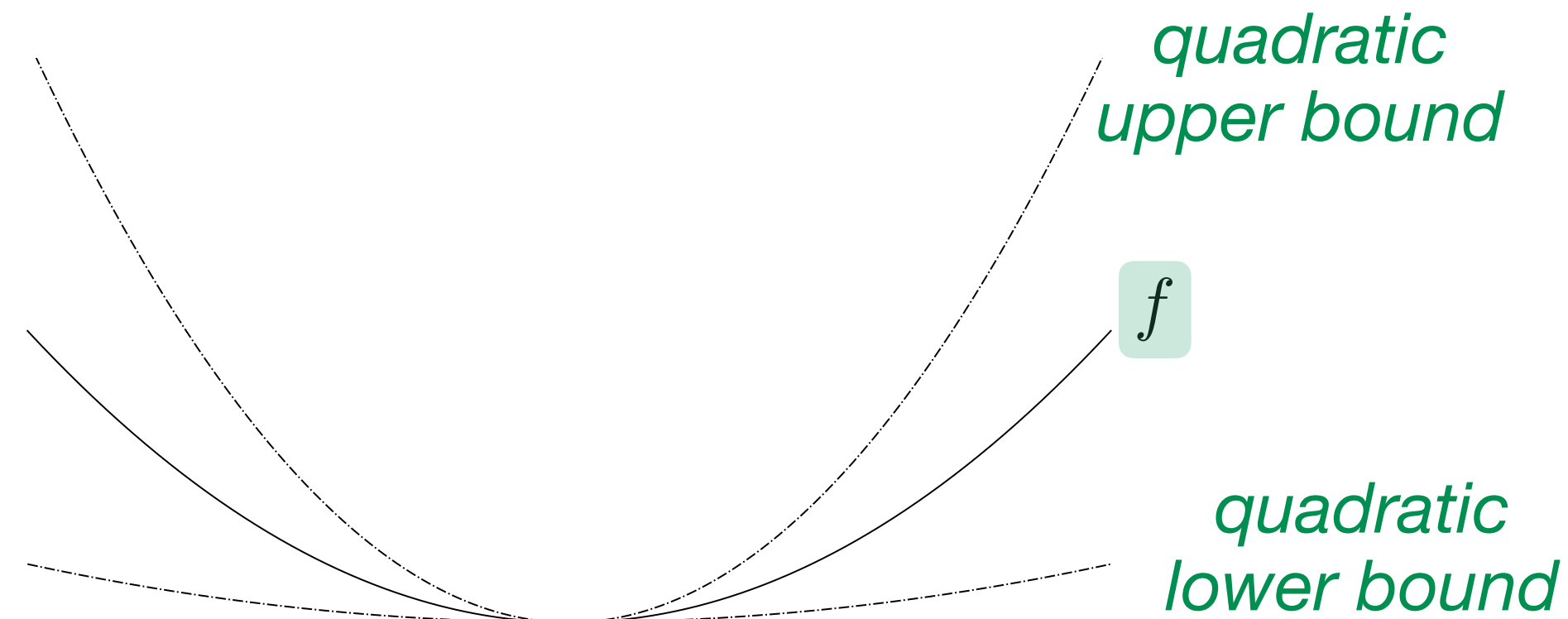


why are worst-case bounds pessimistic?



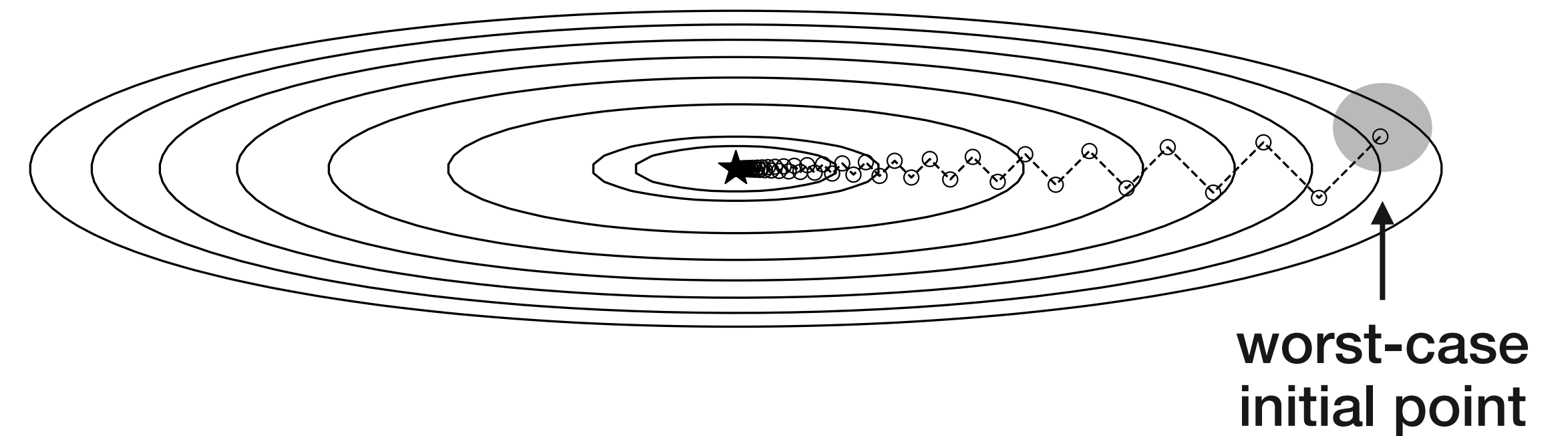
# Issues with classical convergence analysis

general function classes  
( $f$  is strongly convex and smooth...)



we may never encounter that function

pessimistic bounds



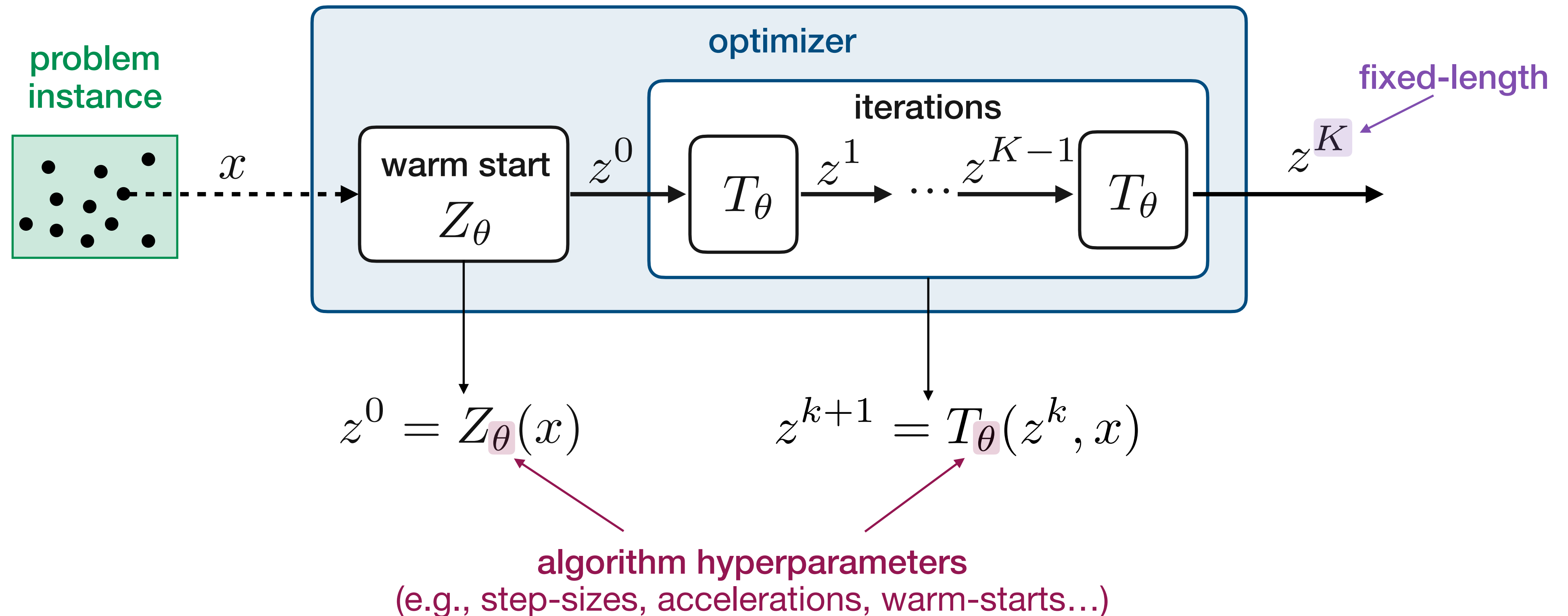
we may never start from that point

practical settings

$$\begin{aligned} &\text{minimize} && f(z, \mathbf{x}) \\ &\text{subject to} && z \in C(\mathbf{x}) \end{aligned}$$

same problem with varying parameters  
 $\longrightarrow x \sim \mathbf{P}$   
(unknown distribution)

# Algorithms as fixed-length computational graphs



**example**  
*projected gradient descent*

$$z^{k+1} = \Pi_{C(x)}(z^k - \theta \nabla_z f(z^k, x))$$



# Verifying the algorithm performance after $K$ iterations

goal

*estimate norm of fixed-point residual*

$$r^K(x) = \|z^K - z^{K-1}\|$$

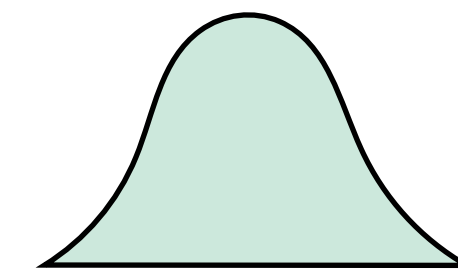


worst-case

$$\max_{x \in \mathcal{X}} r^K(x) \leq \epsilon$$

problem instances

convergence tolerance



probabilistic

$$\mathbf{P}(r^K(x) > \epsilon) \leq \eta$$

problem instances

convergence tolerance

probability bound

# Verification via mixed-integer linear programming



# Worst-case algorithm verification

parametric  
quadratic optimization

$$\text{minimize} \quad (1/2)z^T Pz + q(x)^T z$$

$$\text{subject to} \quad Az \leq b(x)$$

↑  
problem instances



algorithm

(ADMM, PDHG,...)

$$z^{k+1} = T_\theta(z^k, x)$$

verification problem

$$\max_{x \in \mathcal{X}} r^K(x) = \text{maximize} \quad \|z^K - z^{K-1}\|_\infty$$

↑  
performance metric

$$\text{subject to} \quad z^{k+1} = T_\theta(z^k, x), \quad k = 0, \dots, K-1$$
$$z^0 = Z_\theta(x), \quad x \in \mathcal{X}$$

↑  
problem instances

NP-hard problem!



Verification of First-Order Methods for Parametric Quadratic Optimization

V. Ranjan and B. Stellato

arXiv e-prints:2403.03331 (2025)

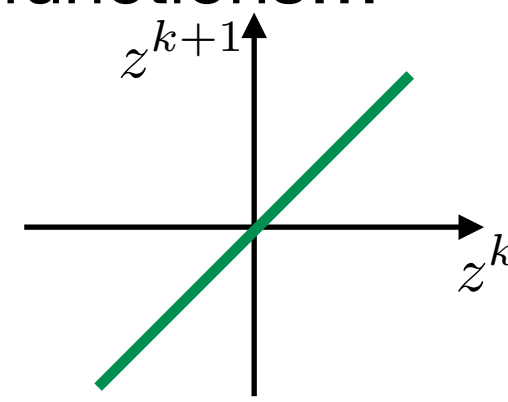
[github.com/stellatogrp/sdp\\_algo\\_verify](https://github.com/stellatogrp/sdp_algo_verify)

# Algorithm steps as mixed-integer linear constraints

## Linear steps → Linear constraints

e.g., gradient, momentum, restarts, anchors, prox of quadratic functions...

$$Mz^{k+1} = Az^k + Bx$$

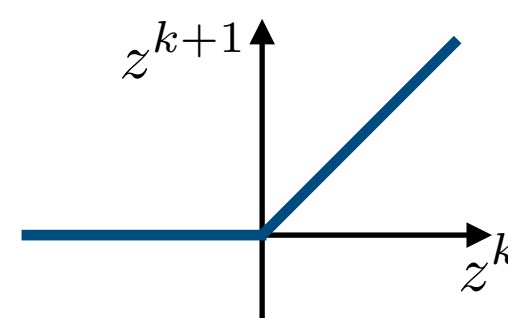


## Piecewise affine steps → Mixed-integer constraints

### Elementwise maximum (ReLU)

e.g., one-sided projections

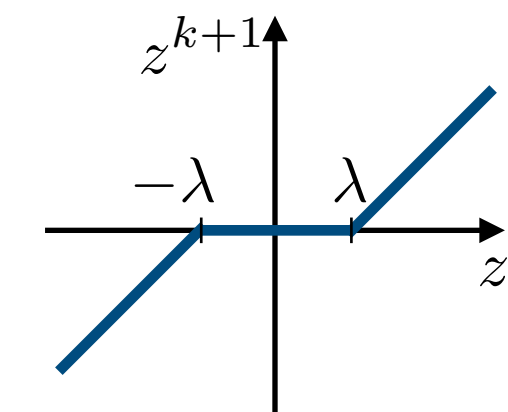
$$z^{k+1} = (z^k)_+ = \max\{z^k, 0\}$$



### Soft-thresholding

e.g., prox of 1-norm function

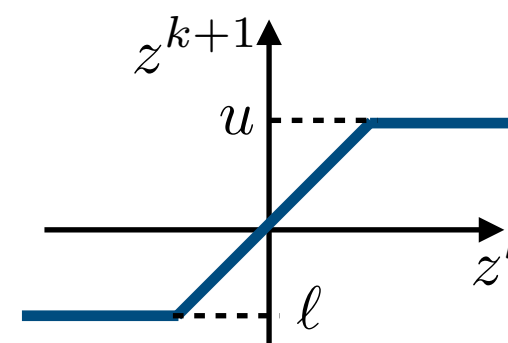
$$z^{k+1} = \phi_\lambda(z^k) = \max\{z^k, \lambda\} - \max\{-z^k, \lambda\}$$



### Saturated linear unit (SatLin)

e.g., box projections

$$z^{k+1} = \mathcal{S}_{\ell, u}(z^k) = \min\{\max\{z^k, \ell\}, u\}$$



gradient  
step

**Example:**

**Nonnegative least squares**

minimize  $(1/2)\|Dz - x\|_2^2$

subject to  $z \geq 0$

**Projected Gradient Descent**

$$w^{k+1} = (I - \theta D^T D)z^k + \theta D^T x$$

$$z^{k+1} = \max\{w^{k+1}, 0\}$$

projection  
step

similar MIP constraints in  
neural network verification

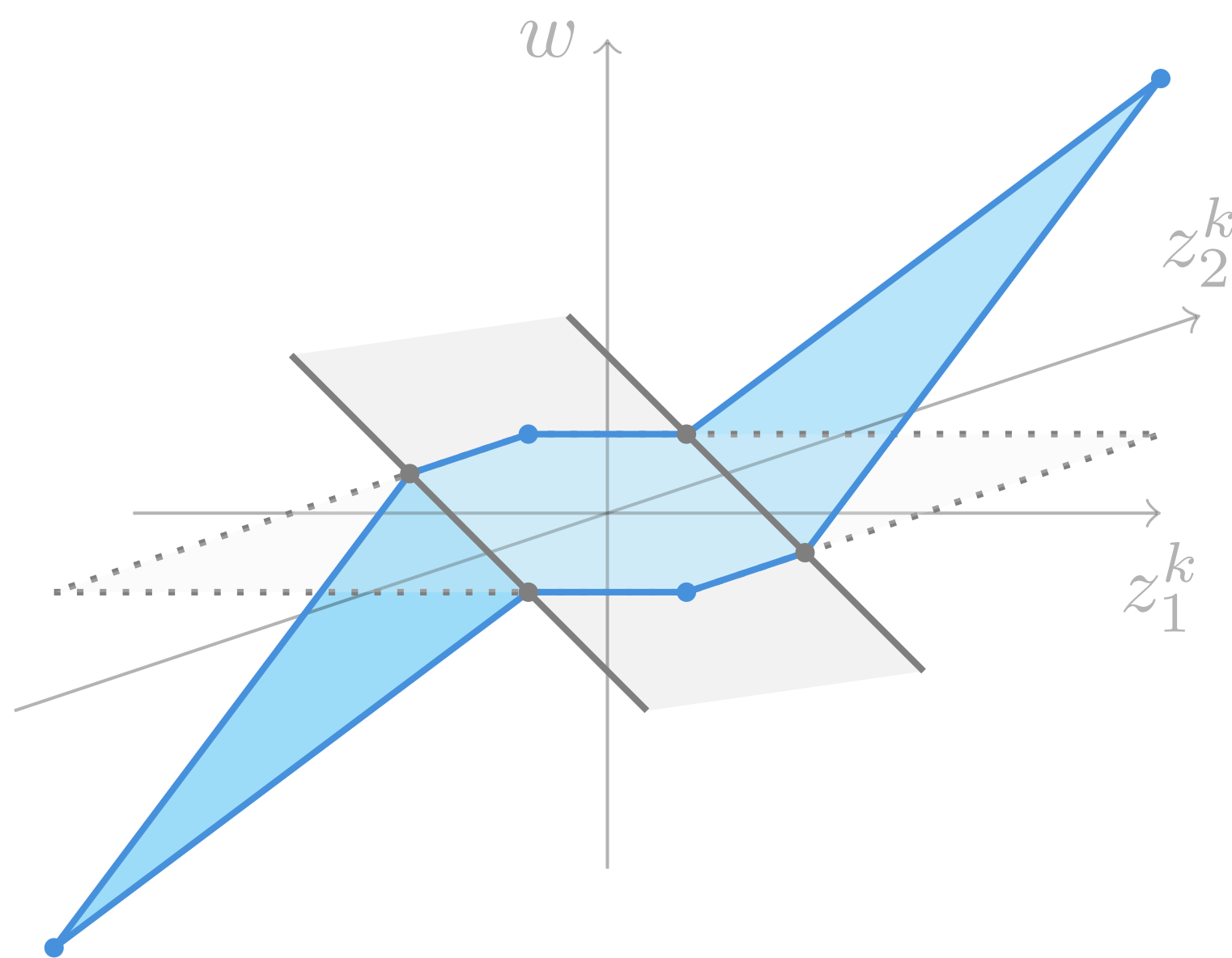
Liu et al. (2021), Albarghouthi (2021),  
Ceccon et al. (2022), Fischetti and Jo (2018),  
Tjeng et al. (2019)



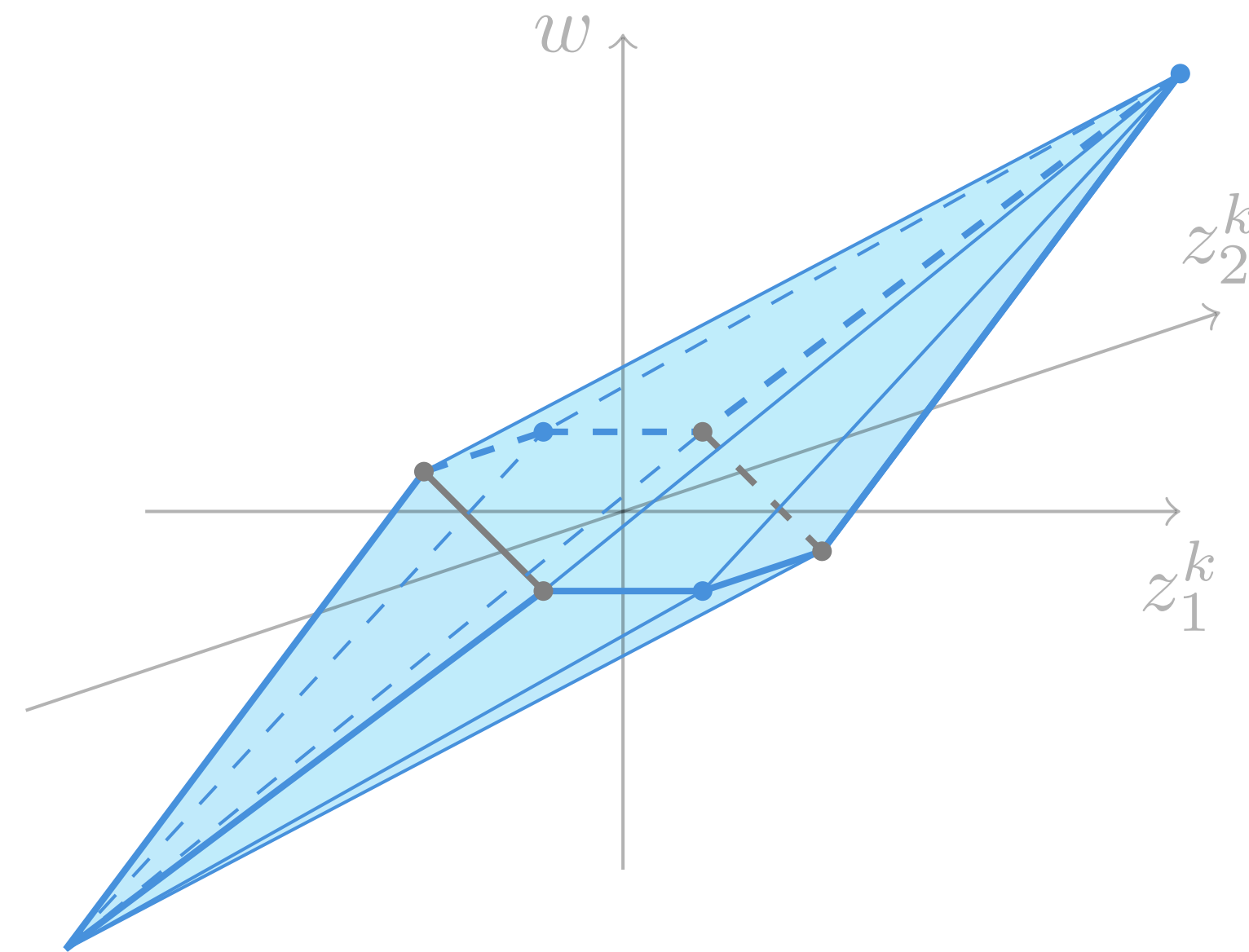
# Constructing strong MIP formulations

piecewise affine steps  
*soft-thresholding operator*

$$w = \phi_\gamma(z_1^k + z_2^k)$$



convex hull



Inspired by: Anderson et al. (2020), Tjandraatmadja et al. (2020), Tsay et al. (2021), Hojny et al. (2024), Huchette et al. (2025)



exponential  
number of  
inequalities!



separation  
problem solvable  
in linear time

# Using operator theory to tighten MIP formulations

$$\begin{aligned} \max_{x \in \mathcal{X}} r^K(x) = & \text{maximize} & & \|z^K - z^{K-1}\|_\infty & \longleftarrow & \text{performance metric} \\ & \text{subject to} & & z^{k+1} = T_\theta(z^k, x), & & k = 0, \dots, K-1 \\ & & & z^0 = Z_\theta(x), & & x \in \mathcal{X} \end{aligned}$$

operator theory bound

$$\|z^K - z^{K-1}\|_\infty \leq \alpha_K$$

e.g., linear convergence

$$\alpha_K = C \tau^K \longleftarrow \text{rate}$$

previous iterate bounds

(bound tightening, interval propagation, etc)

$$\underline{z}^{K-1} \leq z^{K-1} \leq \bar{z}^{K-1}$$

combine

bounds on latest iterate

$$\underline{z}^{K-1} - \alpha_K \leq z^K \leq \bar{z}^{K-1} + \alpha_K$$

main idea

solve verification problem for increasing  $K$  and exploit bounds

Examples



# Sparse coding for signal reconstruction

minimize  $(1/2) \| Dz - x \|_2^2 + \lambda \| z \|_1$

known dictionary  $D$ , noisy signal  $x$ , reconstructed signal  $z$

## Iterative Soft-Thresholding Algorithm (ISTA)

$$z^{k+1} = \phi_{\lambda\theta} \left( (I - \theta D^T D) z^k + \theta D^T x \right)$$

soft-thresholding operator

## Fast Iterative Soft-Thresholding Algorithm (FISTA)

$$w^{k+1} = \phi_{\lambda\theta} \left( (I - \theta D^T D) z^k + \theta D^T x \right)$$

$$z^{k+1} = w^{k+1} + \frac{(\beta_k - 1)}{\beta_{k+1}} (w^{k+1} - w^k)$$

soft-thresholding operator, momentum

# Verification results for sparse coding example

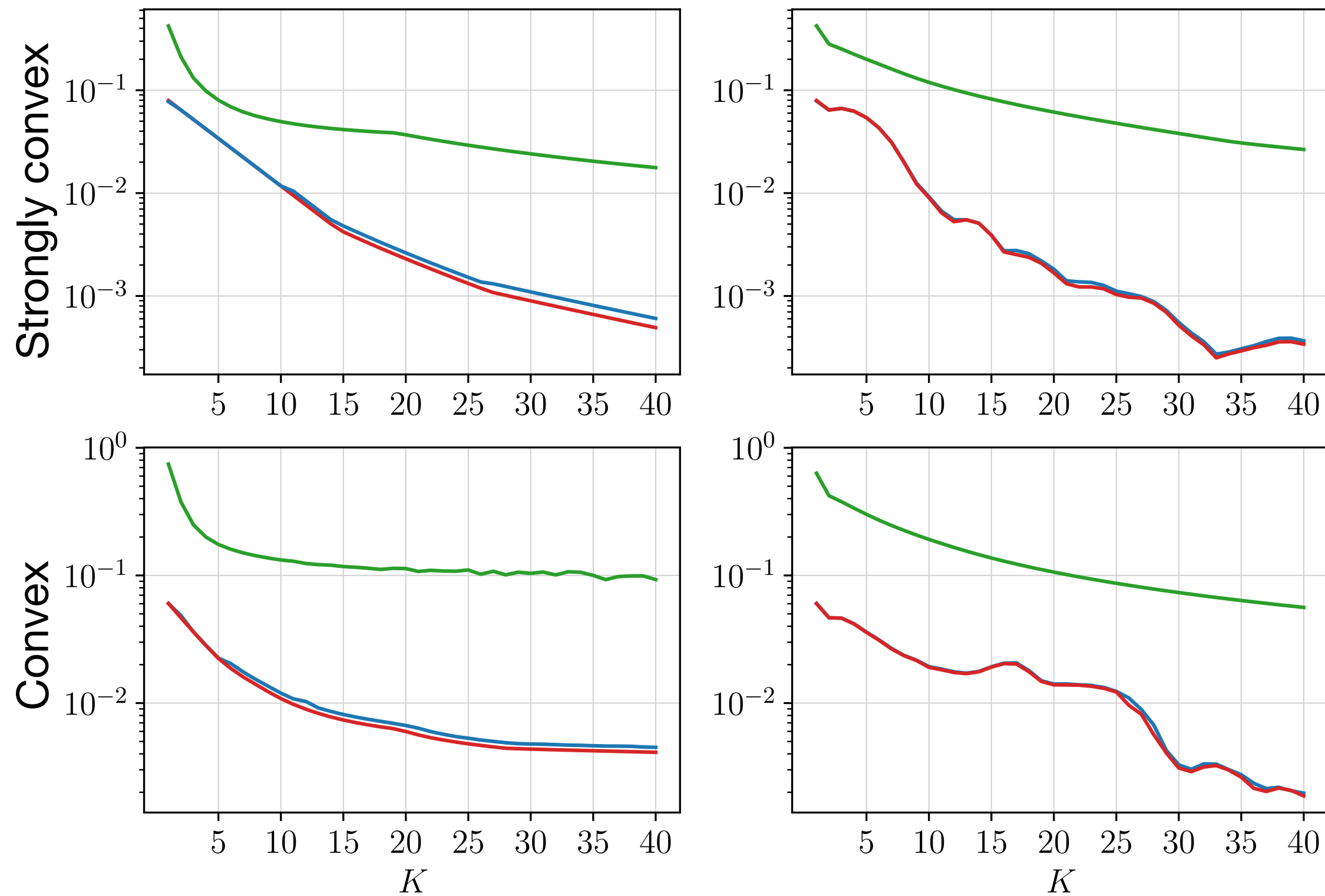
## ISTA

## FISTA

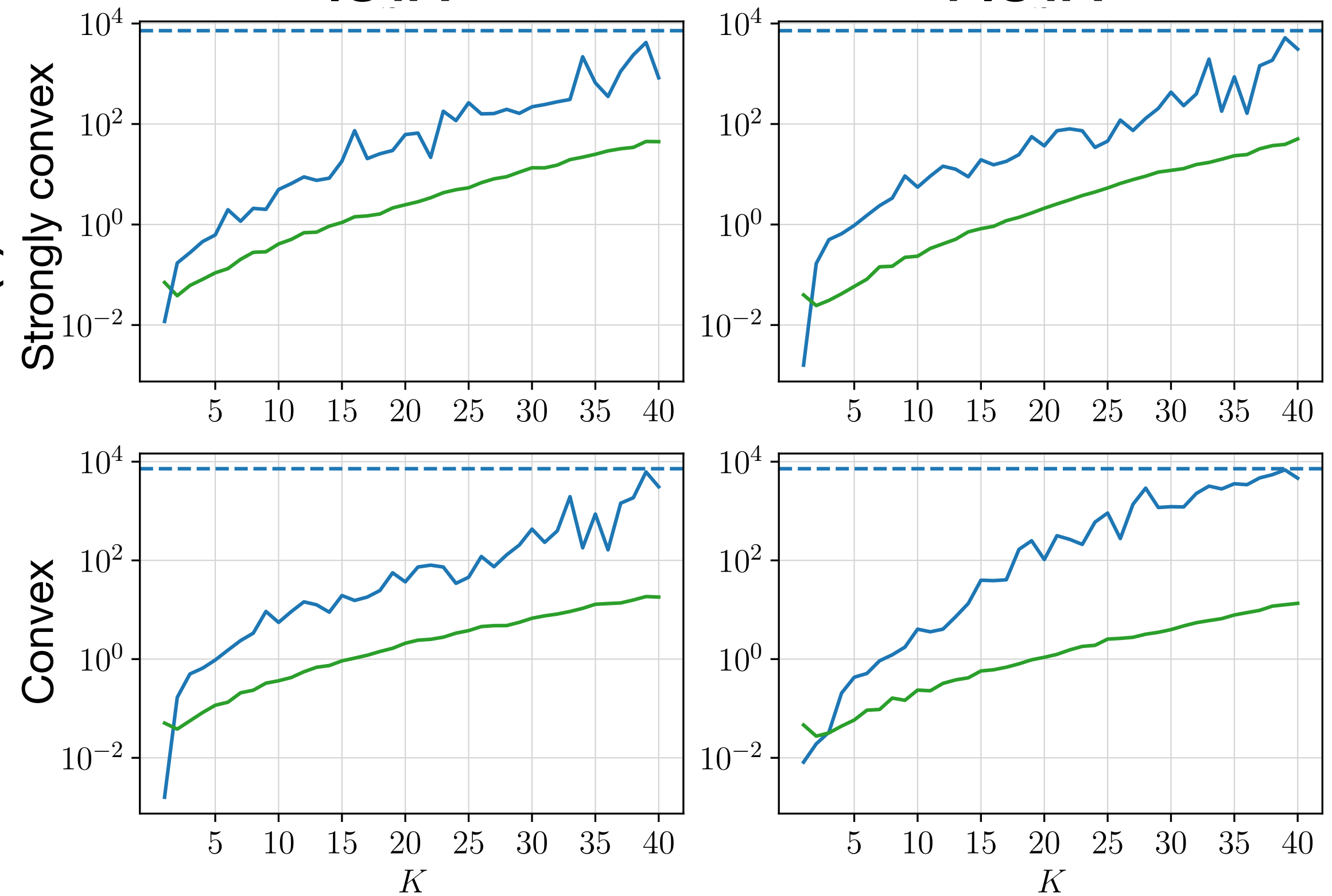
## ISTA

## FISTA

Worst-case fixed-point residual



Solve time (s)



— Sample Maximum    — Theory Bound (PEP SDP)    — Verification Problem

10x-100x reduction in worst-case fixed-point residual (exploiting parametric structure)

exactly captures the ripples of the FISTA acceleration

# Network flow optimization

minimize  $c^T z$  ← network flow

subject to  $A_s z \leq b_s$  ← supplies

arc-node matrices  
(supply and demand)  $A_d z = x$  ← demands

$0 \leq z \leq u$

## Primal-dual hybrid gradient (PDHG)

$$z^{k+1} = \mathcal{S}_{[0,u]}(z^k - \eta(c + A_s^T v^k - A_d^T w^k))$$

$$v^{k+1} = (v^k + \eta(-b_s + A_s(2z^{k+1} - z^k)))_+$$

$$w^{k+1} = w^k + \eta(x - A_d(2z^{k+1} - z^k))_+$$

↑ saturated linear unit
↑ one-sided projection

## Primal-dual hybrid gradient with momentum (mPDHG)

$$z^{k+1} = \mathcal{S}_{[0,u]}(z^k - \eta(c + A_s^T v^k - A_d^T w^k))$$

$$\tilde{z}^{k+1} = z^k + \frac{k}{k+3}(z^{k+1} - z^k)$$

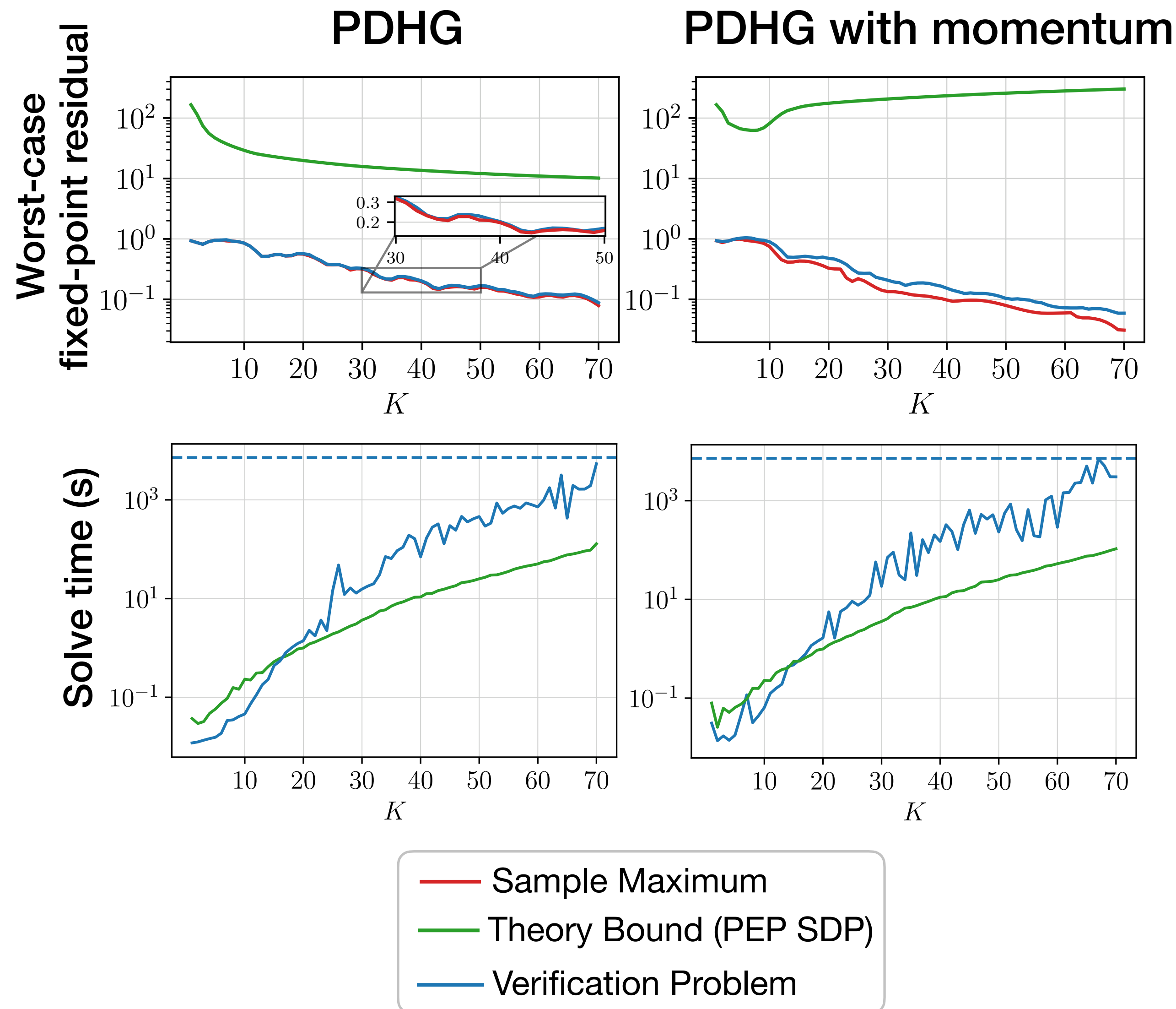
$$v^{k+1} = (v^k + \eta(-b_s + A_s(2\tilde{z}^{k+1} - z^k)))_+$$

$$w^{k+1} = w^k + \eta(x - A_d(2\tilde{z}^{k+1} - z^k))_+$$

↑ saturated linear unit
momentum
↑ one-sided projection



# Verification results for network flow optimization example



we verify convergence even when PEP doesn't!

known behavior in momentum/heavy-ball methods  
[L. Lessard, B. Recht, and A. Packard, (2016)]  
[Goujaud, Taylor, Dieuleveut (2023)]

can be faster than SDP for few iterations

# Optimal control

optimal sequence of states and controls

minimize 
$$\sum_{t=0}^T s_t^T Q s_t + u_t^T R u_t$$

subject to 
$$s_{t+1} = A^{\text{dyn}} s_t + B^{\text{dyn}} u_t, \quad t = 1, \dots, T - 1$$

linear dynamics 
$$s_{\min} \leq s_t \leq s_{\max}, \quad t = 1, \dots, T$$

$$u_{\min} \leq u_t \leq u_{\max}, \quad t = 1, \dots, T$$

$s_0 = x$  initial state

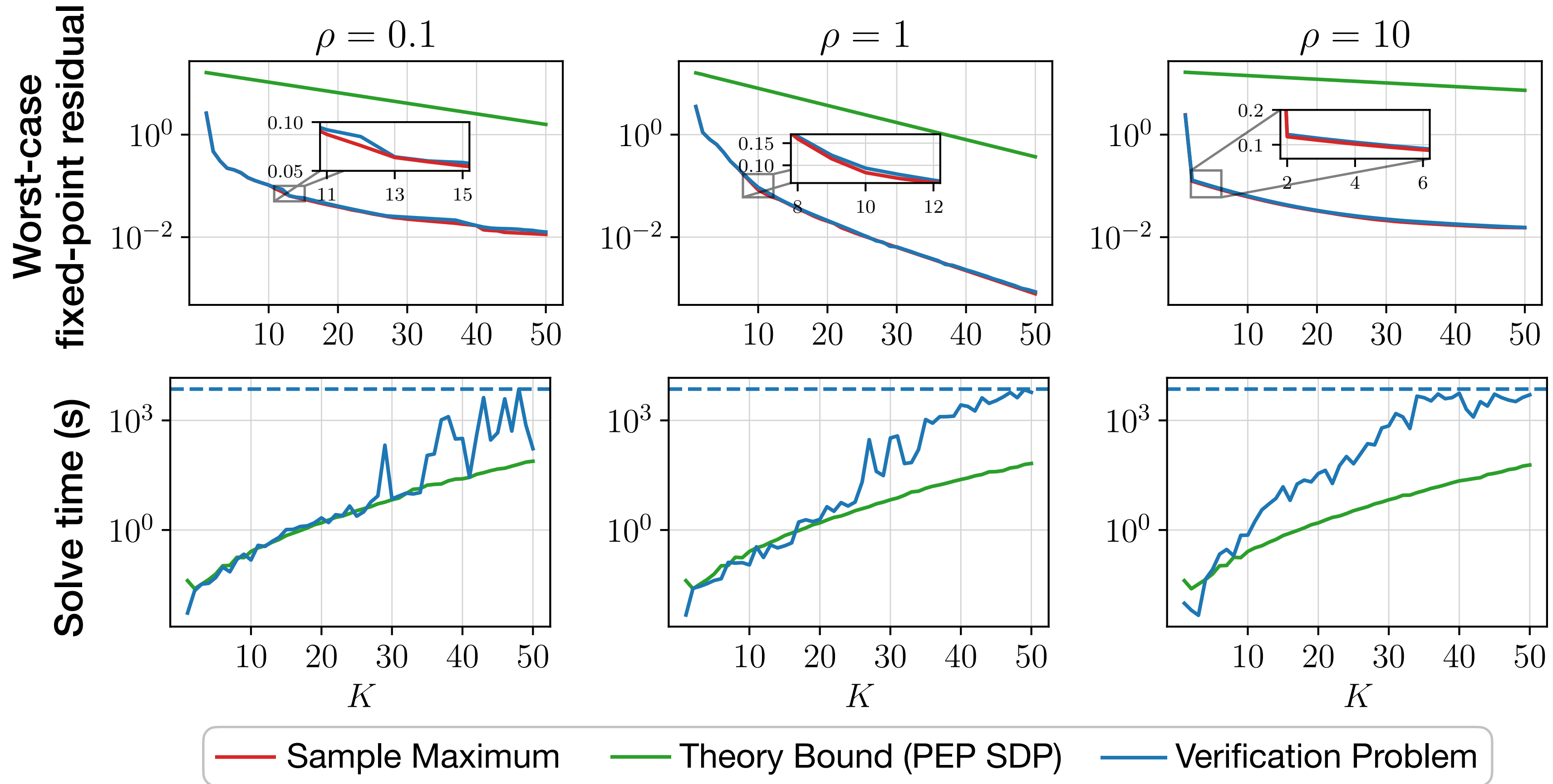
## OSQP ADMM splitting

saturated linear unit 
$$w^{k+1} = \mathcal{S}_{[l(x), u(x)]}(v^k)$$

Solve 
$$(P + \sigma I + \rho M^T M) z^{k+1} = \sigma z^k - q(x) + \rho M^T (2w^{k+1} - v^k)$$

$$v^{k+1} = v^k + M z^{k+1} - w^{k+1}$$

# Verification results for optimal control problem



can be used to design (tune) custom algorithms!

exactly quantifies the iterations required

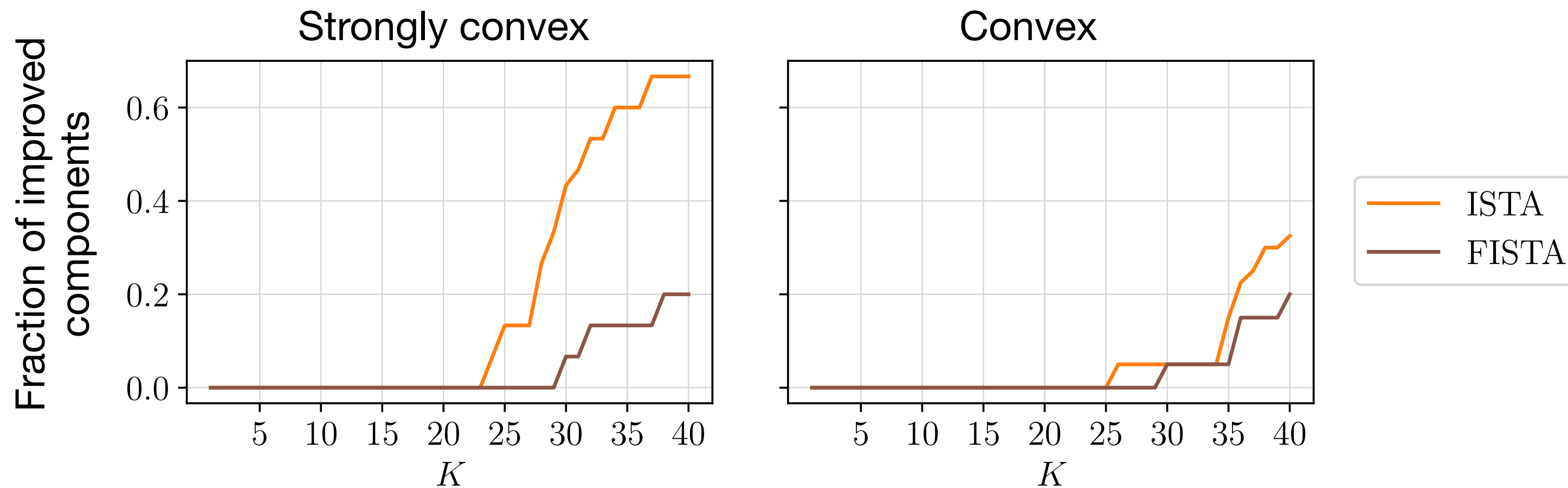


crucial in real-time settings!



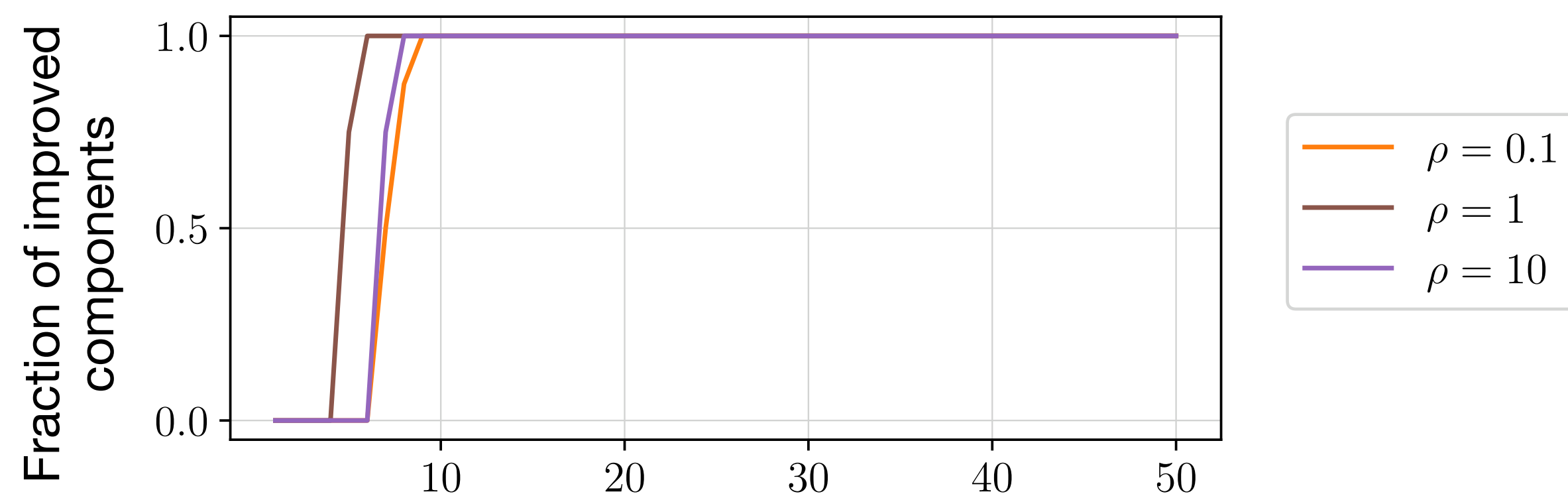
# Operator theory can tighten bounds on the iterates

## Sparse coding



strongly convex problems have the largest benefits  
*(because of linear convergence)*

## Optimal control



operator theory bounds did not help PDHG  
*(usually requires thousands of iterations)*

# Acknowledgements



**Vinit Ranjan**  
Princeton



**Jisun Park**  
Princeton

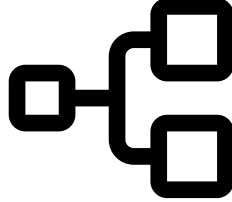

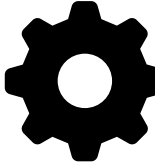


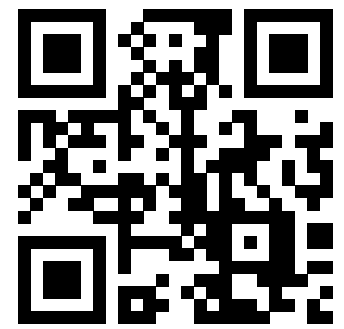
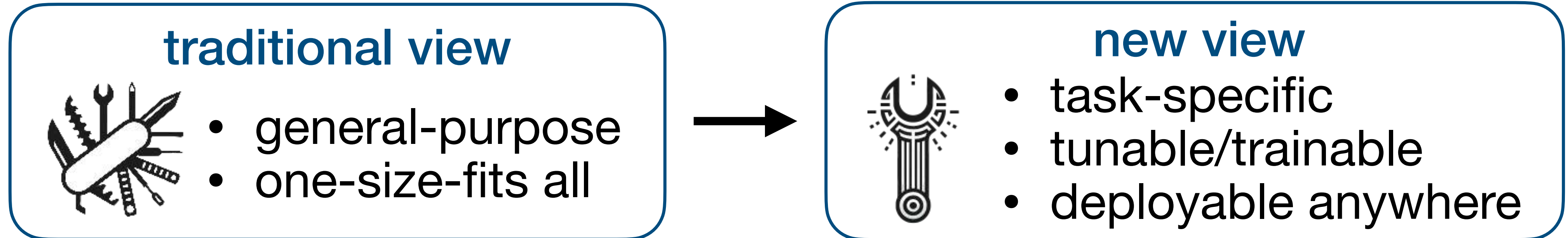
**Andrea Lodi**  
CornellTech



**Stefano Gualandi**  
Univ. of Pavia

# Verification of First-order Methods via Mixed-Integer Linear Programming

1. **parametric** structure matters 
2. **algorithm verification**  operator theory
3. useful to **design new algorithms** 



**Exact Verification of First-Order Methods via Mixed-Integer Linear Programming**

V. Ranjan, J. Park, S. Gualandi, A. Lodi, and B. Stellato

*arXiv e-prints:2412.11330 (2025)*

 [github.com/stellatogrp/mip\\_algo\\_verify](https://github.com/stellatogrp/mip_algo_verify)



**Verification of First-Order Methods for Parametric Quadratic Optimization**

V. Ranjan and B. Stellato

*arXiv e-prints:2403.03331 (2025)*

 [github.com/stellatogrp/sdp\\_algo\\_verify](https://github.com/stellatogrp/sdp_algo_verify)



**Backup**

# Unconstrained QP

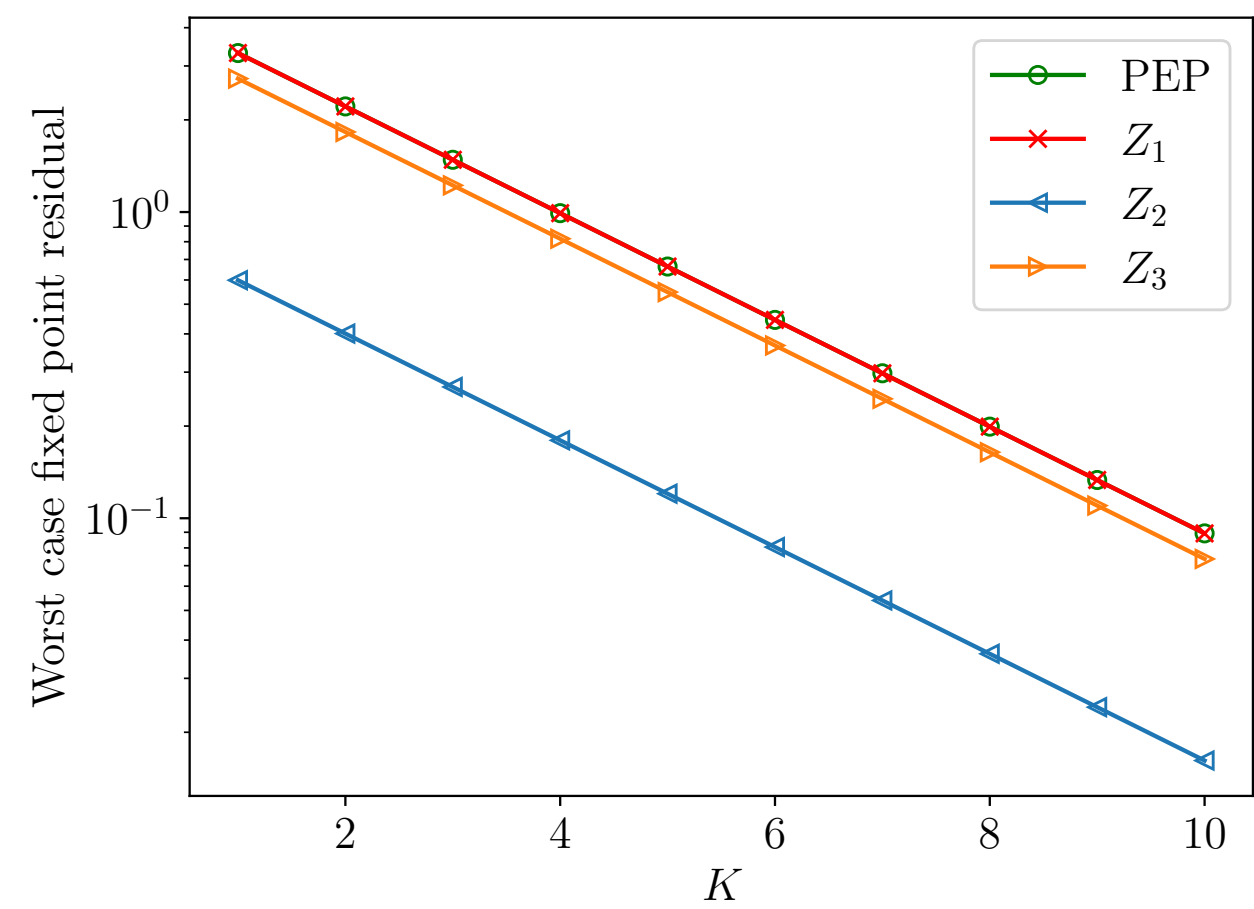
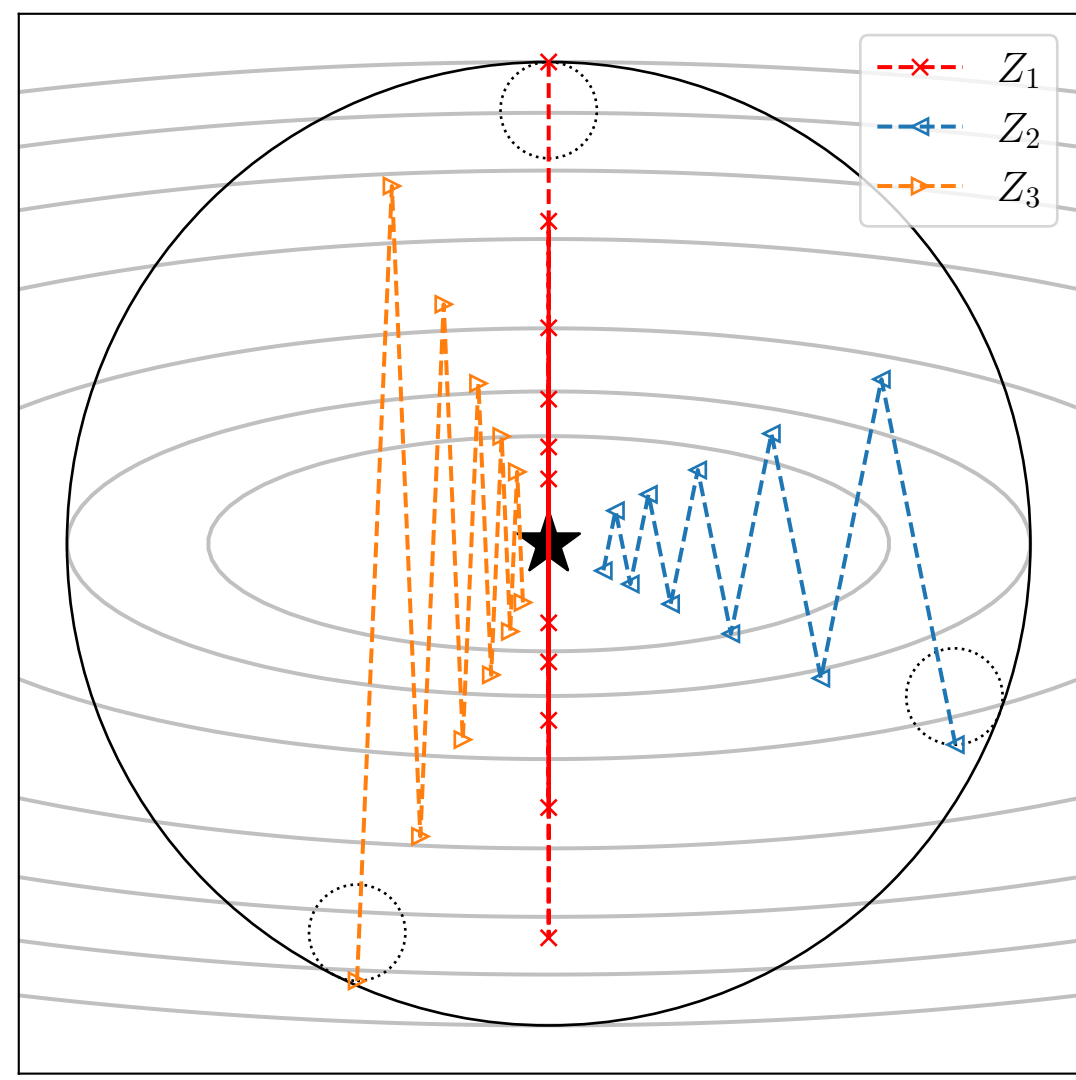
minimize  $(1/2)z^T Pz + x^T z$  ← parameters

verification problem

maximize  $\|z^K - z^{K-1}\|$   
 subject to  $z^{k+1} = z^k - \theta(Pz^k + x), \quad k = 0, \dots, K-1$  ← gradient descent  
 $z^0 = Z_\theta(x), \quad x \in \mathcal{X}$

warm-starts

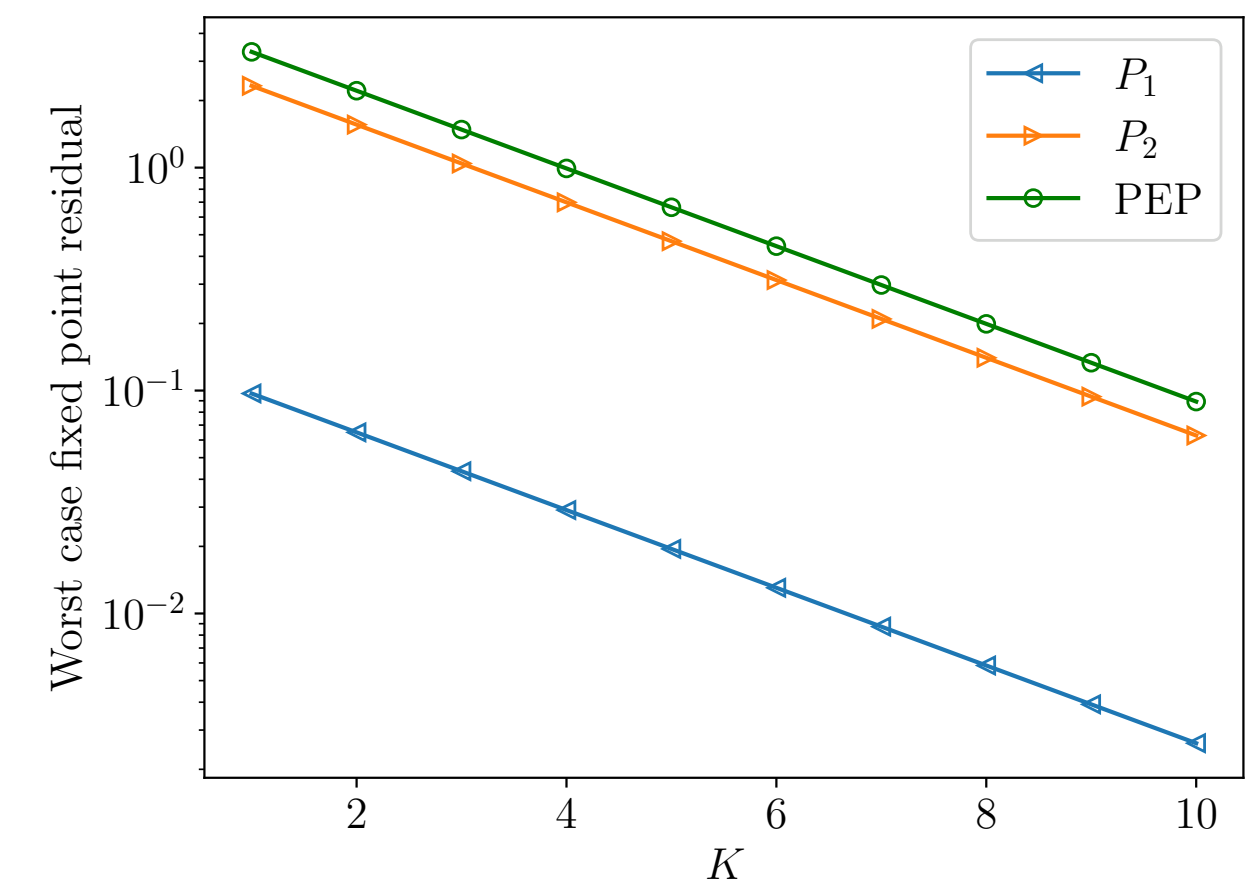
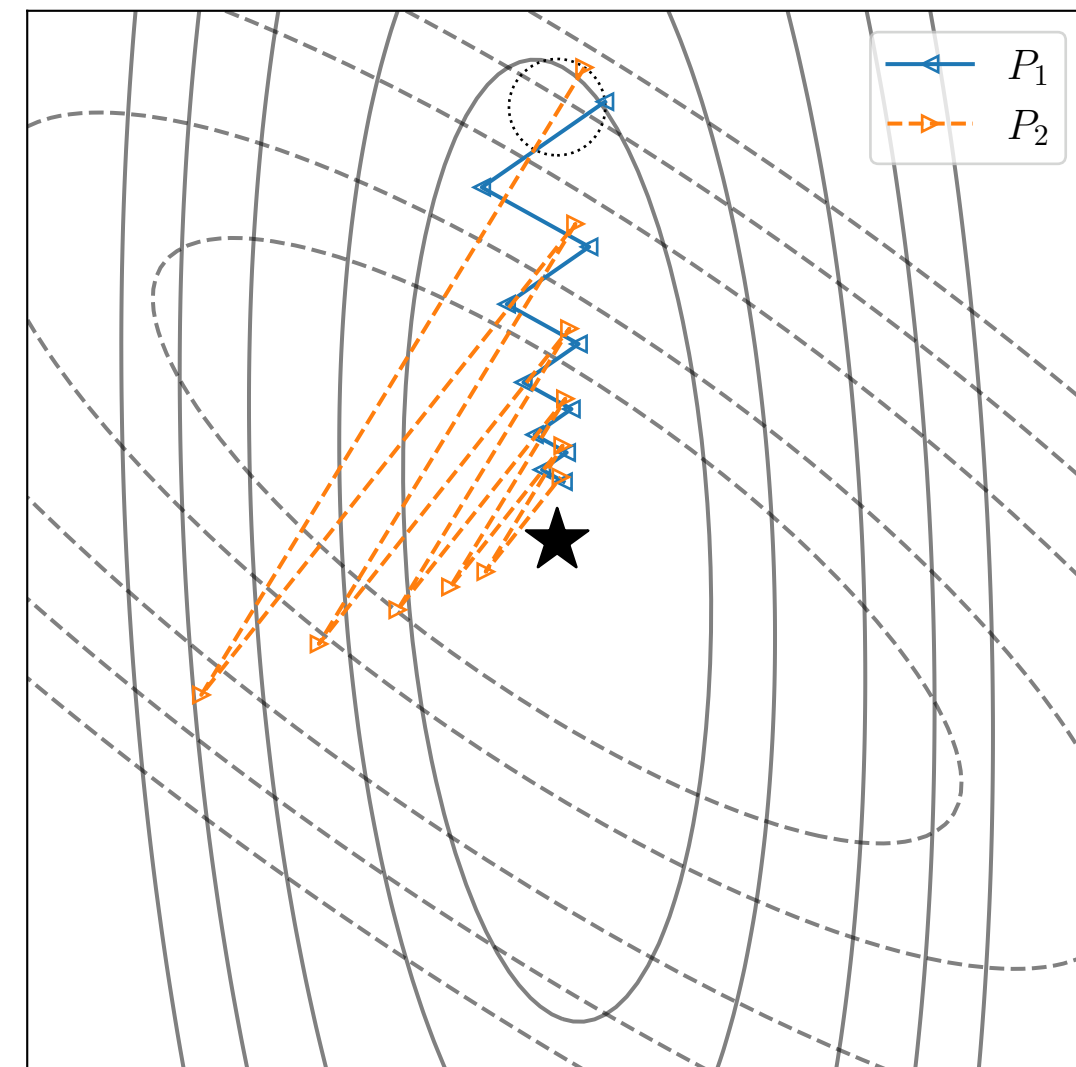
**case I**  
 $Z_\theta(x) = Z_1, Z_2, \text{ or } Z_3$   
 $x \in \mathcal{X} = \{0\}$



PEP-SDP cannot distinguish warm-starts

rotated functions

**case II**  
 $Z_\theta(x) = \{z \mid \|z - 0.9 \cdot \mathbf{1}\| \leq 0.1\}$   
 $x \in \mathcal{X} = \{0\}$   
 $P_1, P_2$  rotations of  $P$



PEP-SDP cannot distinguish quadratic functions

# Objective of verification problem as MIP

$$\|s^K - s^{K-1}\|_\infty = \|t\|_\infty = \delta_K$$

- lower bounds  $\underline{s}^{K-1}, \underline{s}^K$
- upper bounds,  $\bar{s}^{K-1}$  and  $\bar{s}^K$



- lower bound  $\underline{t} = \underline{s}^K - \bar{s}^{K-1}$
- upper bound  $\bar{t} = \bar{s}^K - \underline{s}^{K-1}$

## exact reformulation

$$t = t^+ - t^-, \quad t^+ \leq \bar{t} \odot w, \quad t^- \leq -\underline{t} \odot (\mathbf{1} - w)$$

$$t^+ + t^- \leq \delta_K \leq t^+ + t^- + \max\{\bar{t}, -\underline{t}\} \odot (\mathbf{1} - \gamma)$$

$$\mathbf{1}^T \gamma = 1, \quad t^+ \geq 0, \quad t^- \geq 0$$

$w \in \{0, 1\}^n$  (absolute values of the components of  $t$ )  
 $\gamma \in \{0, 1\}^d$  (maximum inside the  $\ell_\infty$ -norm)



# Soft-thresholding operator

$$w = \phi_\lambda(a^T z) = \begin{cases} a^T z - \lambda & a^T z > \lambda \\ 0 & |a^T z| \leq \lambda \\ a^T z + \lambda & a^T z < -\lambda \end{cases}$$

region

$$\Phi = \{(z, w) \in [\underline{z}, \bar{z}] \times \mathbf{R} \mid w = \phi_\lambda(a^T z)\}$$

lower and upper bounds  
(needed for convex hull)

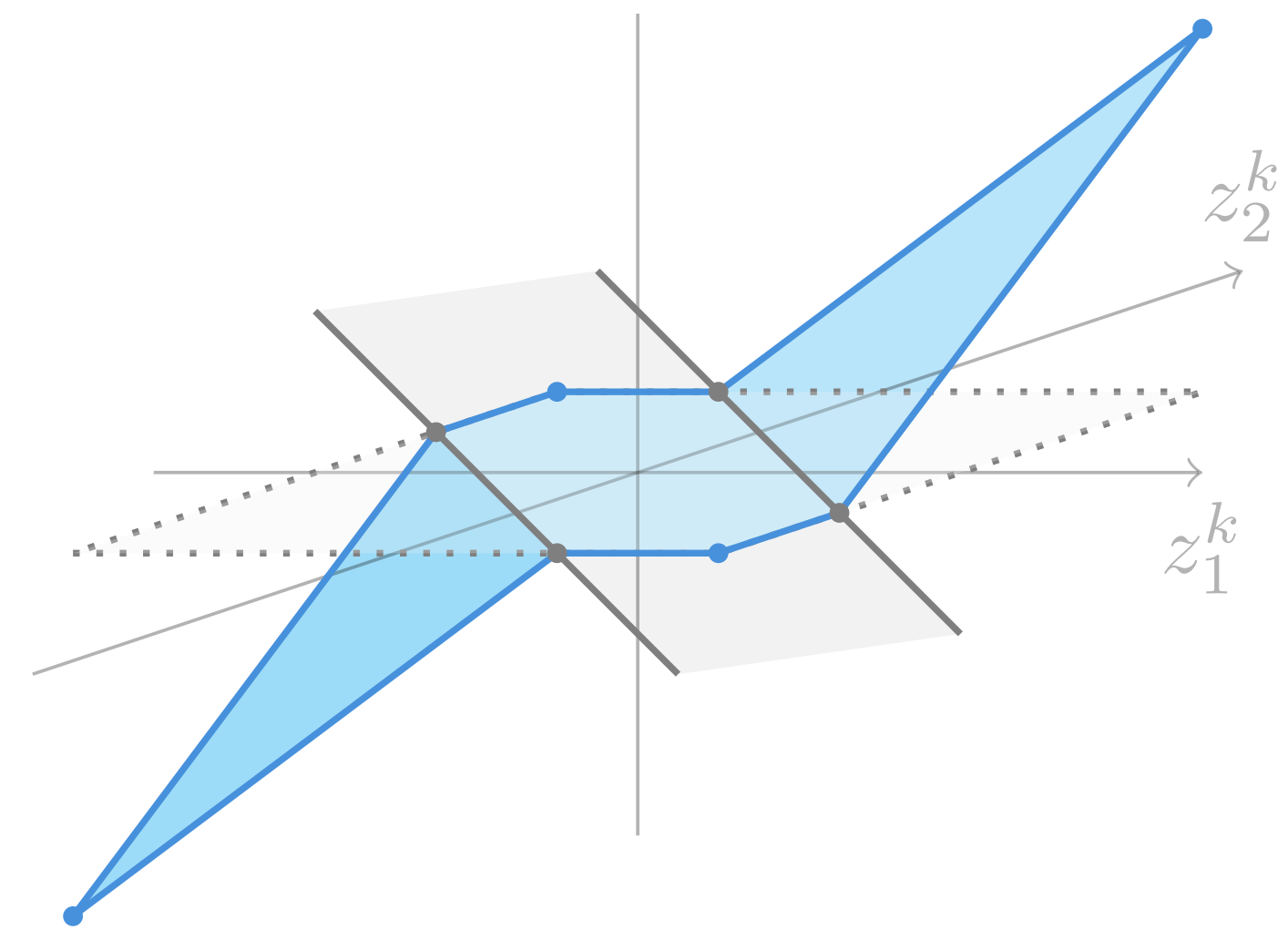
$$\ell_I = \sum_{i \in I} a_i \ell_i^0 + \sum_{i \notin I} a_i u_i^0$$

$$u_I = \sum_{i \in I} a_i u_i^0 + \sum_{i \notin I} a_i \ell_i^0$$

$$u_i^0 = \begin{cases} \bar{z}_i & a_i \geq 0 \\ \underline{z}_i & \text{otherwise} \end{cases}$$

$$\ell_i^0 = \begin{cases} \underline{z}_i & a_i \geq 0 \\ \bar{z}_i & \text{otherwise} \end{cases}$$

example:  $\phi_\lambda(z_1^k + z_2^k)$

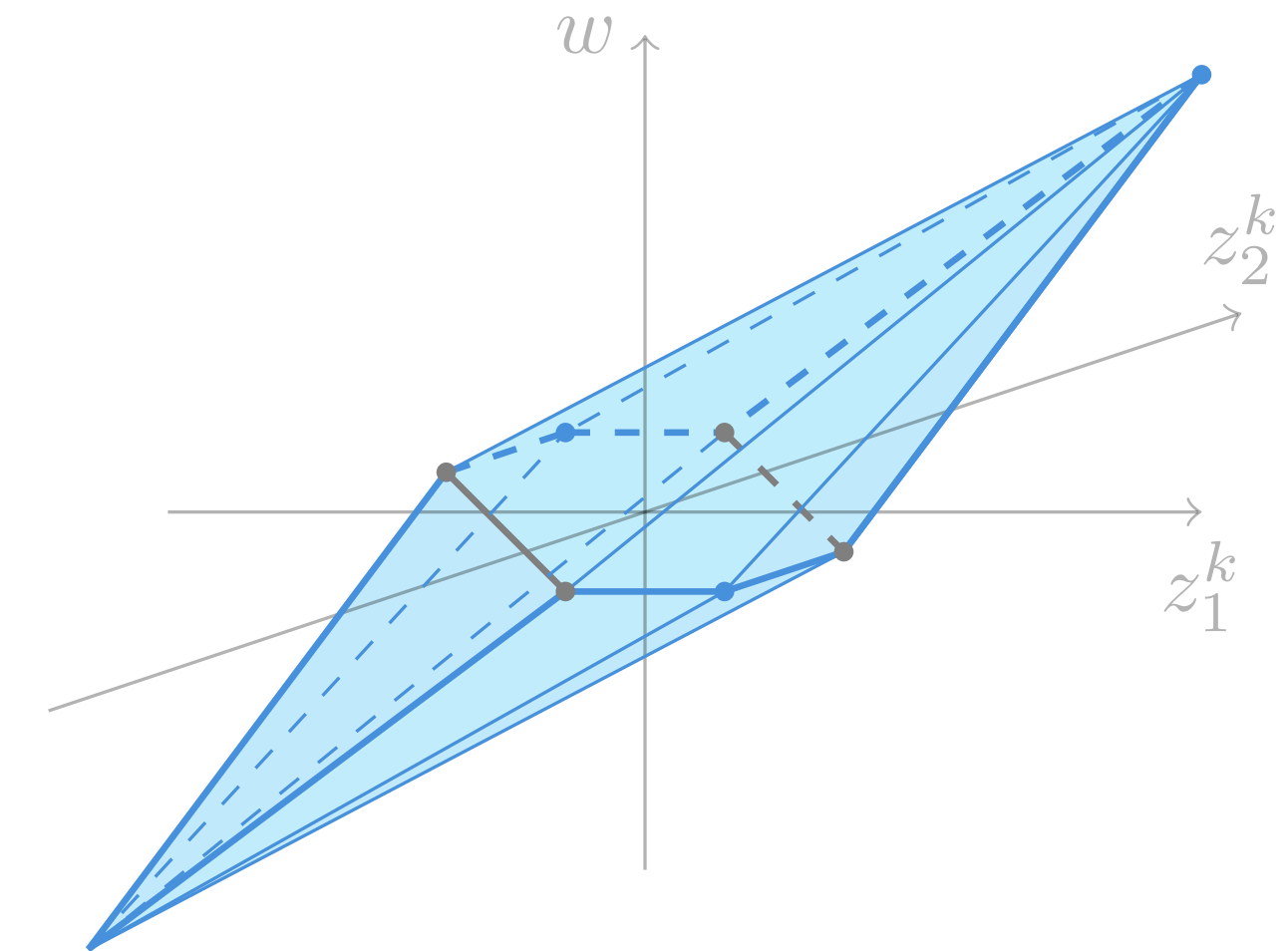


# Convex hull of soft-thresholding operator

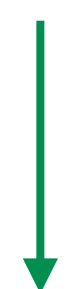
convex hull

$$\text{conv}(\Phi) = \left\{ (z, w) \in [\underline{z}, \bar{z}] \times \mathbf{R} \left| \begin{cases} w = a^T z - \lambda & \ell_{\{1, \dots, n\}} > \lambda \\ w = a^T z + \lambda & u_{\{1, \dots, n\}} < -\lambda \\ w = 0 & -\lambda \leq \ell_{\{1, \dots, n\}} \leq u_{\{1, \dots, n\}} \leq \lambda \\ (z, w) \in Q & \text{otherwise} \end{cases} \right. \right\}$$

example:  $\phi_\lambda(z_1^k + z_2^k)$



$$Q = \left\{ \begin{aligned} & a^T z - \lambda \leq w \leq a^T z + \lambda \\ & \frac{\ell_J + \lambda}{\ell_J - \lambda} (a^T z - \lambda) \leq w \leq \frac{u_J - \lambda}{u_J + \lambda} (a^T z + \lambda) \\ & w \leq \sum_{i \in I} a_i (z_i - \ell_i^0) + \frac{\ell_I - \lambda}{u_o^0 - \ell_o^0} (z_o - \ell_o^0), \quad \forall (I, o) \in \mathcal{I}^+ \\ & w \geq \sum_{i \in I} a_i (z_i - u_i^0) + \frac{u_I + \lambda}{\ell_o^0 - u_o^0} (z_o - u_o^0), \quad \forall (I, o) \in \mathcal{I}^- \end{aligned} \right\}$$



separation problem can be solved in linear time (by sorting)



exponential number of inequalities

$$\begin{aligned} \mathcal{I}^- &= \{(I, o) \in 2^{\{1, \dots, n\}} \times \{1, \dots, n\} \mid u_I \leq -\lambda < u_{I \cup \{o\}}, w_I \neq 0\} \\ \mathcal{I}^+ &= \{(I, o) \in 2^{\{1, \dots, n\}} \times \{1, \dots, n\} \mid \ell_{I \cup \{o\}} < \lambda \leq \ell_I, w_I \neq 0\} \end{aligned}$$