

Learning to Accelerate Optimizers with Guarantees

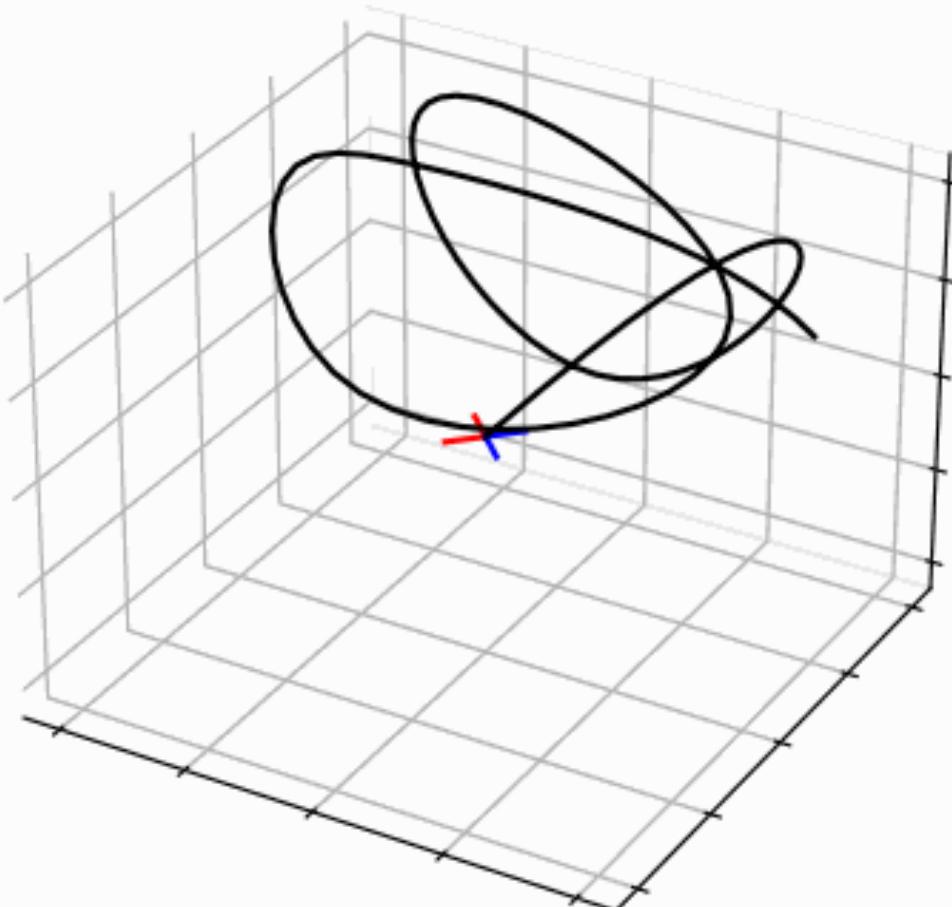
Penn Talk 2024
Rajiv Sambharya



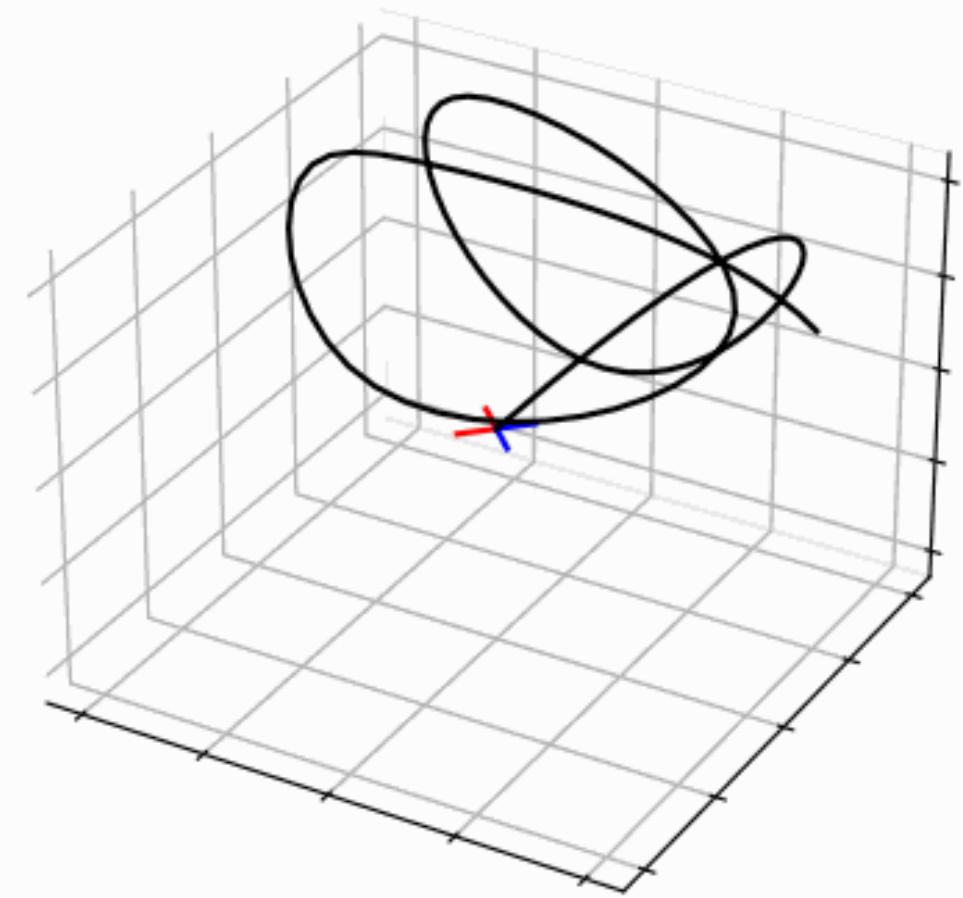
PRINCETON
UNIVERSITY



Tracking a reference trajectory with a quadcopter



Success!
(If given enough time)



Failure: not enough time to solve

Model predictive control

optimize over a smaller horizon (T steps),
implement first control,
repeat

Model predictive controller

$$\begin{aligned} & \text{minimize} && \sum_{t=1}^T \|x_t - \underline{x}_t^{\text{ref}}\|_2^2 \\ & \text{subject to} && x_{t+1} = Ax_t + Bu_t \\ & && x_t \in \mathcal{X}, \quad u_t \in \mathcal{U} \\ & && x_0 = \underline{x}_{\text{init}} \end{aligned}$$

Current state,
reference trajectory

Control
inputs

Challenge: we need faster methods for optimization

Claim: real-world optimization is parametric

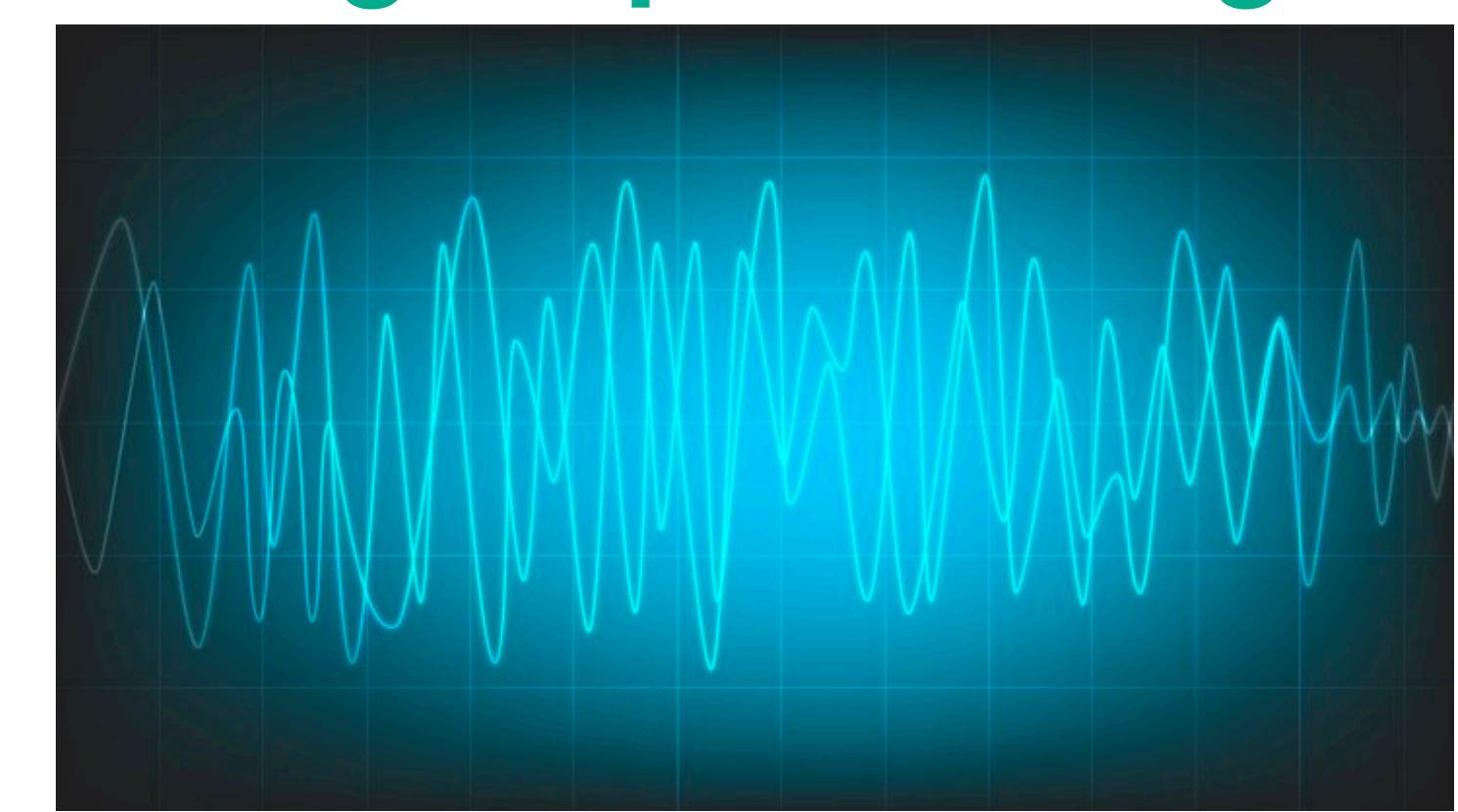
Robotics and control



Energy



Signal processing



Can machine learning speed up parametric optimization?

Goal: Do mapping quickly and accurately

Parameter

$$\theta \longrightarrow$$

$$\begin{aligned} & \text{minimize} && f_{\theta}(z) \\ & \text{subject to} && g_{\theta}(z) \leq 0 \end{aligned}$$

Optimal solution

$$\longrightarrow z^*(\theta)$$

$$\theta \longrightarrow$$



Only Optimization

$$\longrightarrow \hat{z}^{\text{Opt}}(\theta)$$



$$\theta \longrightarrow$$



Only Machine Learning

$$\longrightarrow \hat{z}^{\text{ML}}(\theta)$$



$$\theta \longrightarrow$$



Optimization Machine Learning

$$\longrightarrow \hat{z}^{\text{Opt/ML}}(\theta)$$

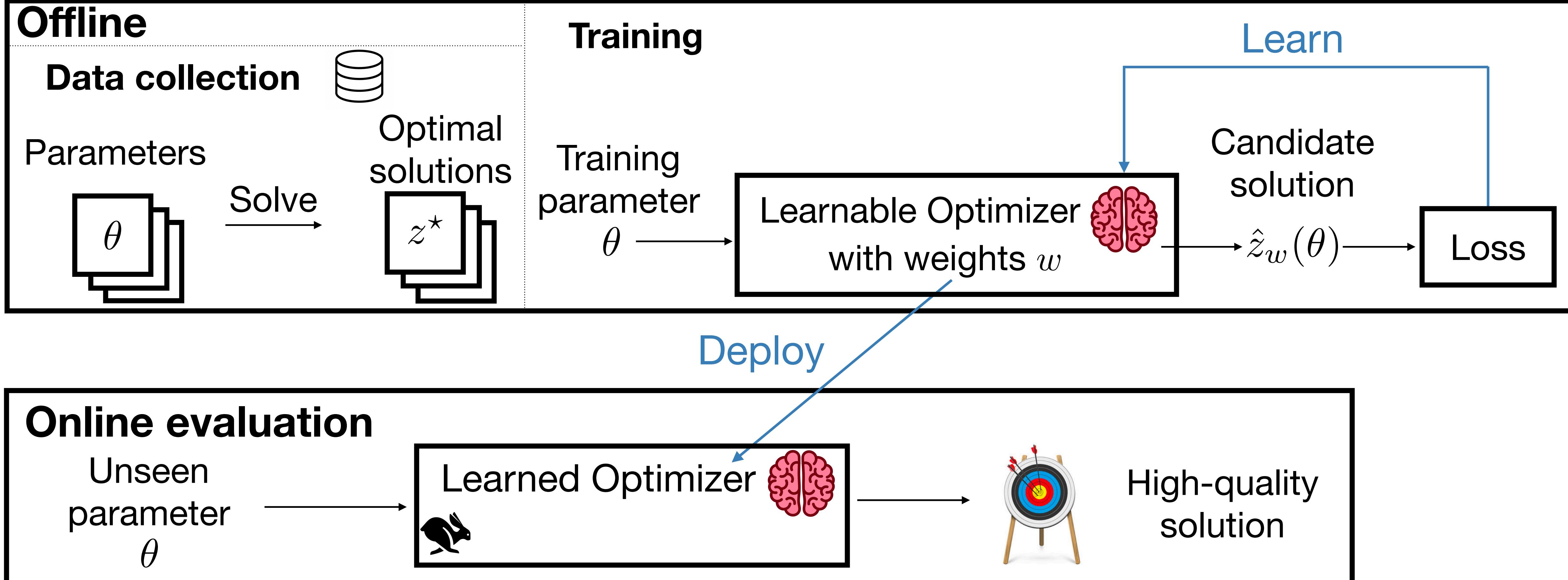


Learning to Optimize

The learning to optimize paradigm

Goal: solve the parametric optimization problem fast

$$\begin{aligned} & \text{minimize} && f_{\theta}(z) \\ & \text{subject to} && g_{\theta}(z) \leq 0 \end{aligned}$$



Challenges in learning to optimize methods

- I. Lack convergence guarantees
- II. Lack generalization guarantees
- III. Lack scalability



We need more **reliable** methods

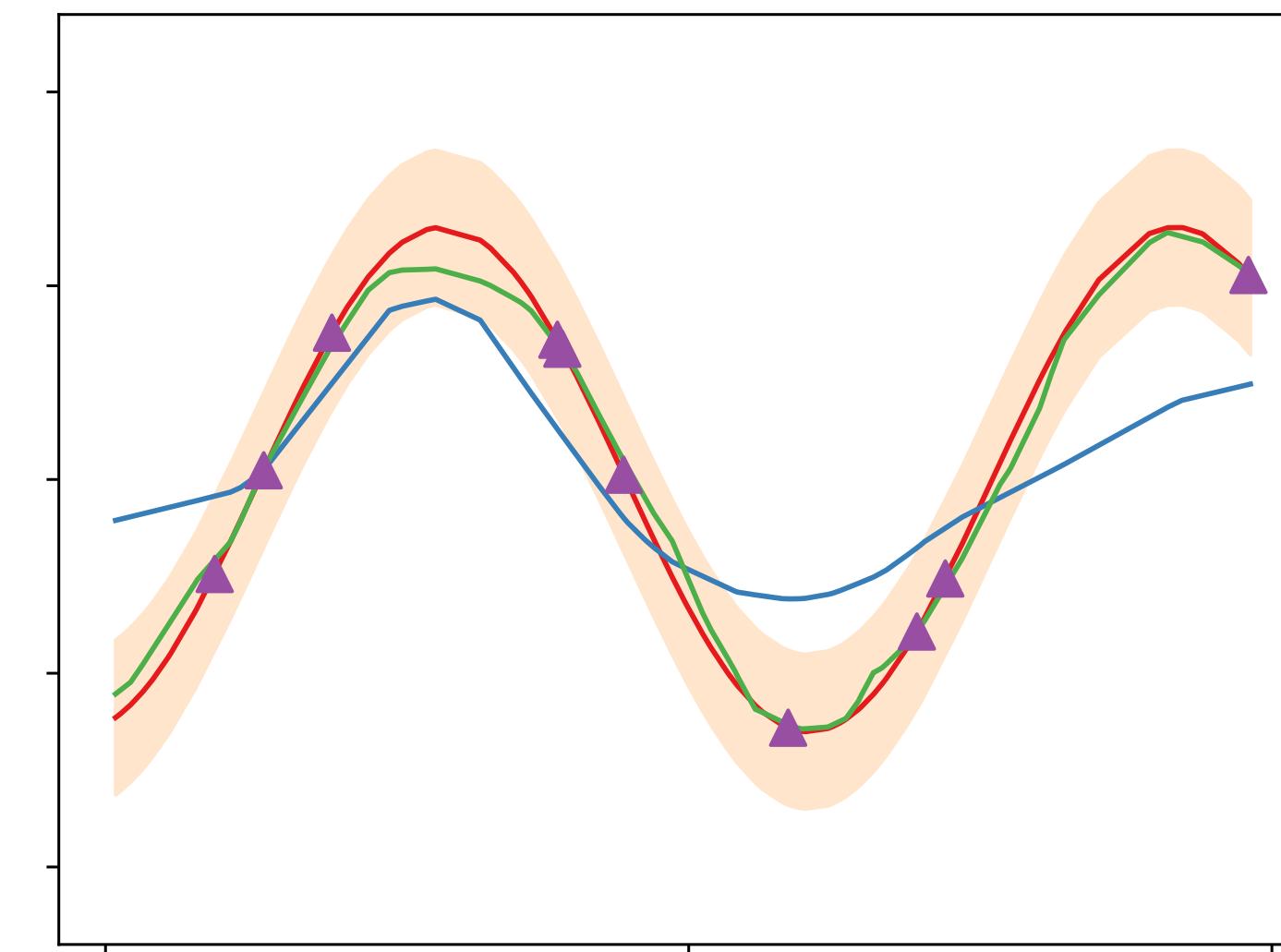
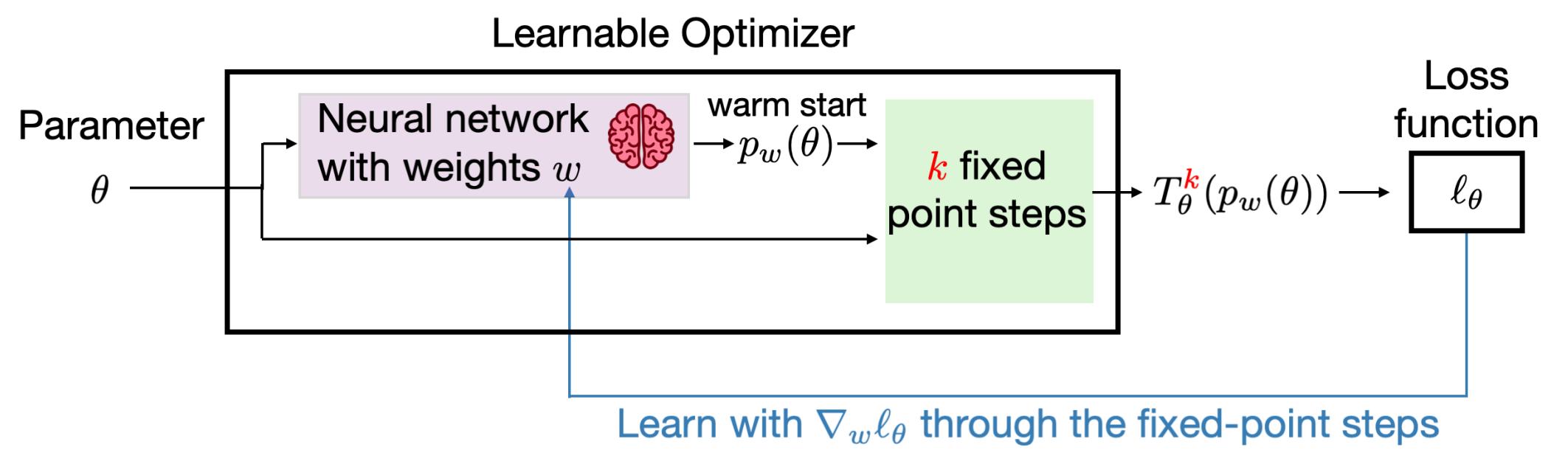


Learning to Optimize: A Primer and A Benchmark [Chen. et al 2021]

“So, to conclude this article, let us quote Sir Winston Churchill: ‘Now this is not the end. It is not even the beginning of the end. But it is, perhaps, the end of the beginning.’”

Talk Outline

- Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms
- Part 2: Data-Driven Performance Guarantees for Classical and Learned Optimizers



Collaborators



Georgina
Hall



Brandon
Amos

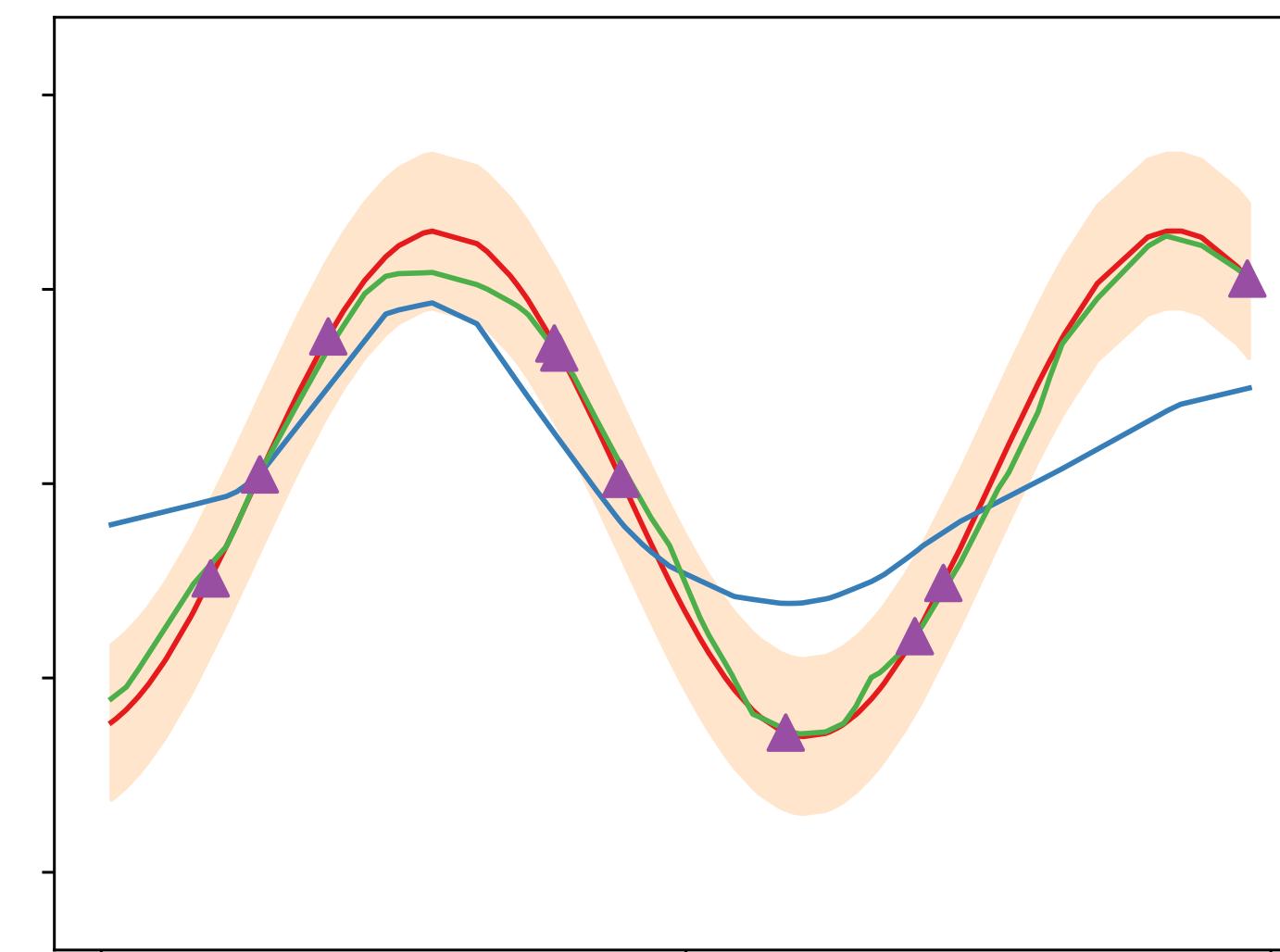
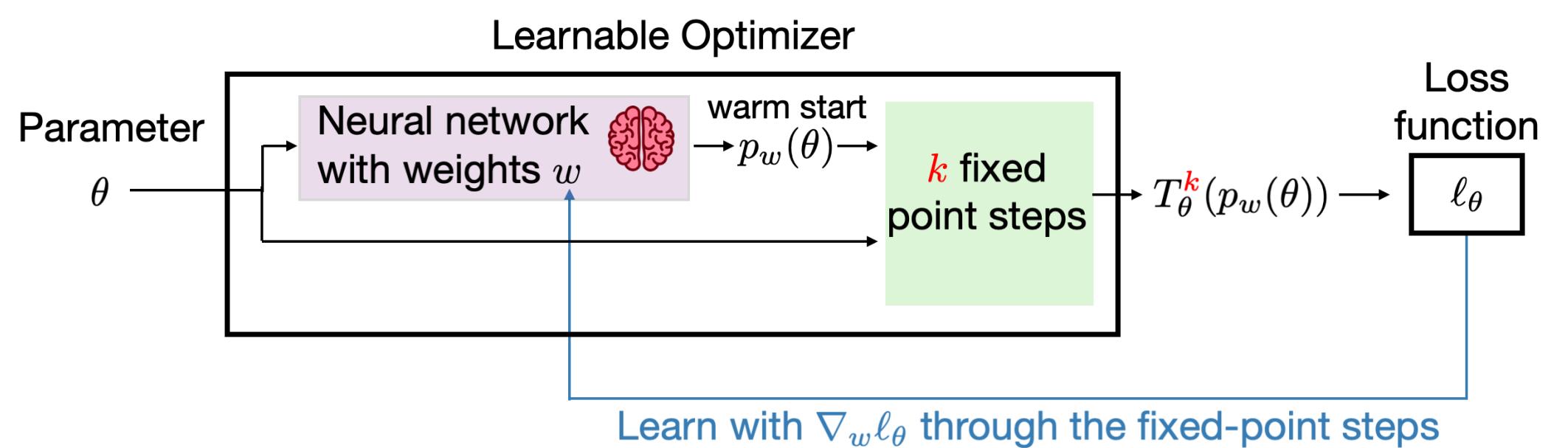


Bartolomeo
Stellato



Talk Outline

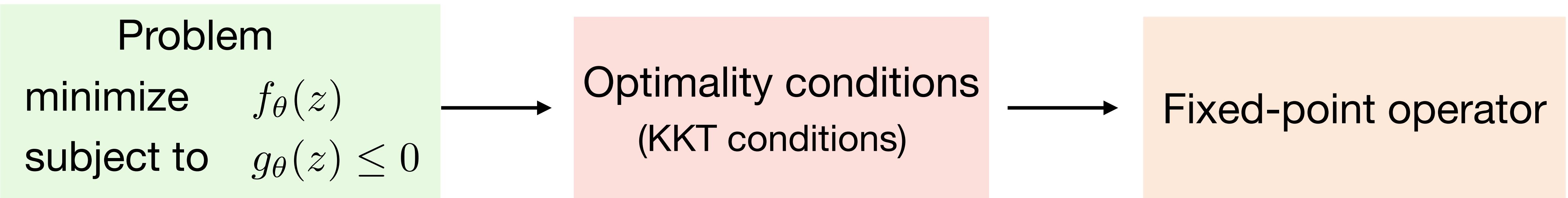
- Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms
- Part 2: Data-Driven Performance Guarantees for Classical and Learned Optimizers



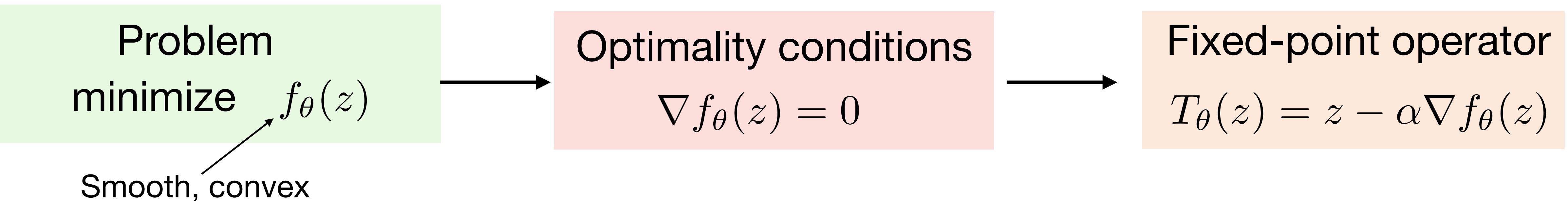
Fixed-point optimization problems are ubiquitous

Parametric fixed-point problem: find z such that $z = T_\theta(z)$

Convex optimization

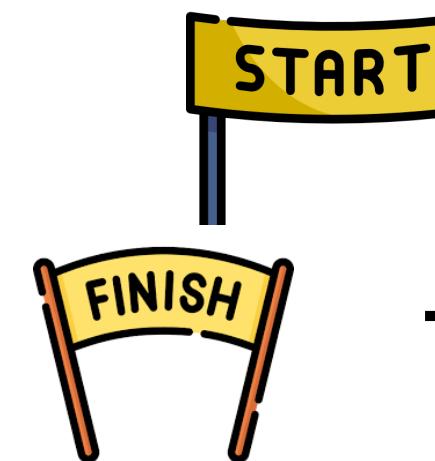


Unconstrained, smooth convex optimization



Many optimization algorithms are fixed-point iterations

Fixed-point iterations: $z^{i+1} = T_\theta(z^i)$



Initialize with z^0 (a warm-start)

Terminate when $\|T_\theta(z^i) - z^i\|_2$ is small

Fixed-point residual

Example: Proximal gradient descent

$$\begin{array}{ll} \text{minimize} & g_\theta(z) + h_\theta(z) \\ & \begin{array}{ll} \text{Convex} & \text{Convex} \\ \text{Smooth} & \text{Non-smooth} \end{array} \end{array}$$

Iterates $z^{i+1} = \text{prox}_{\alpha h_\theta}(z^i - \alpha \nabla g_\theta(z^i))$

$$\text{prox}_s(v) = \arg \min_x \left(s(x) + \frac{1}{2} \|x - v\|_2^2 \right)$$



Problem: limited iteration budget

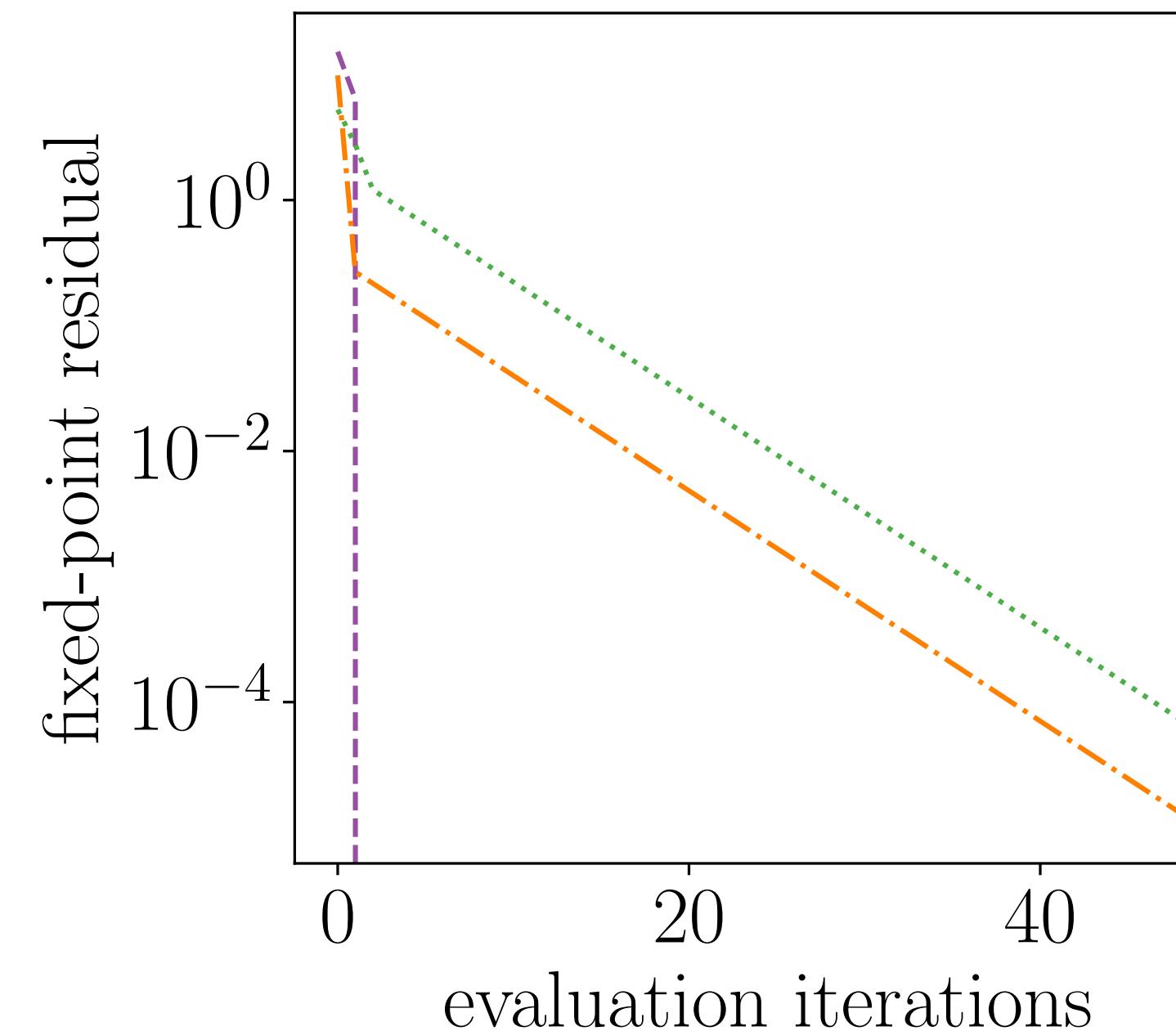
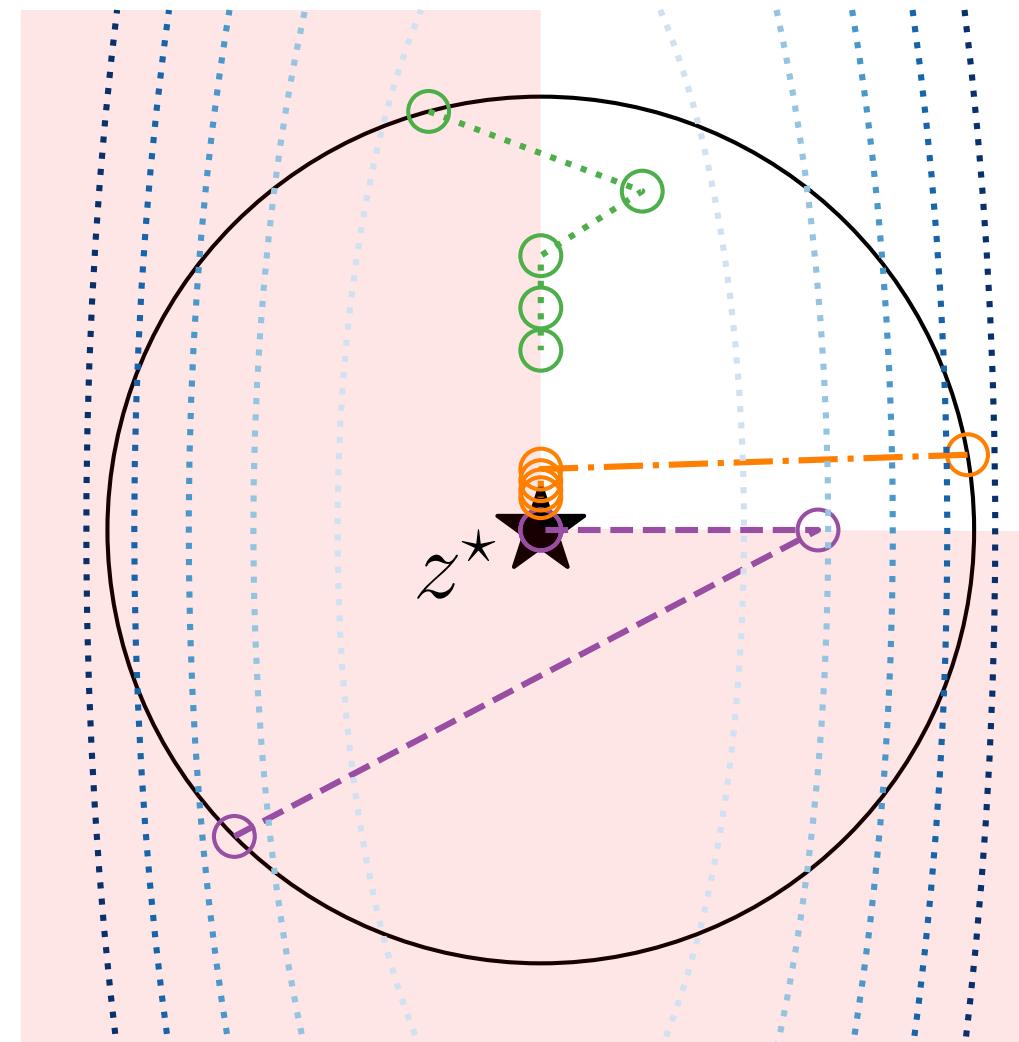


Solution: learn the warm-start to improve the solution within budget

Some warm starts are better than others

minimize $10z_1^2 + z_2^2$
subject to $z \geq 0$

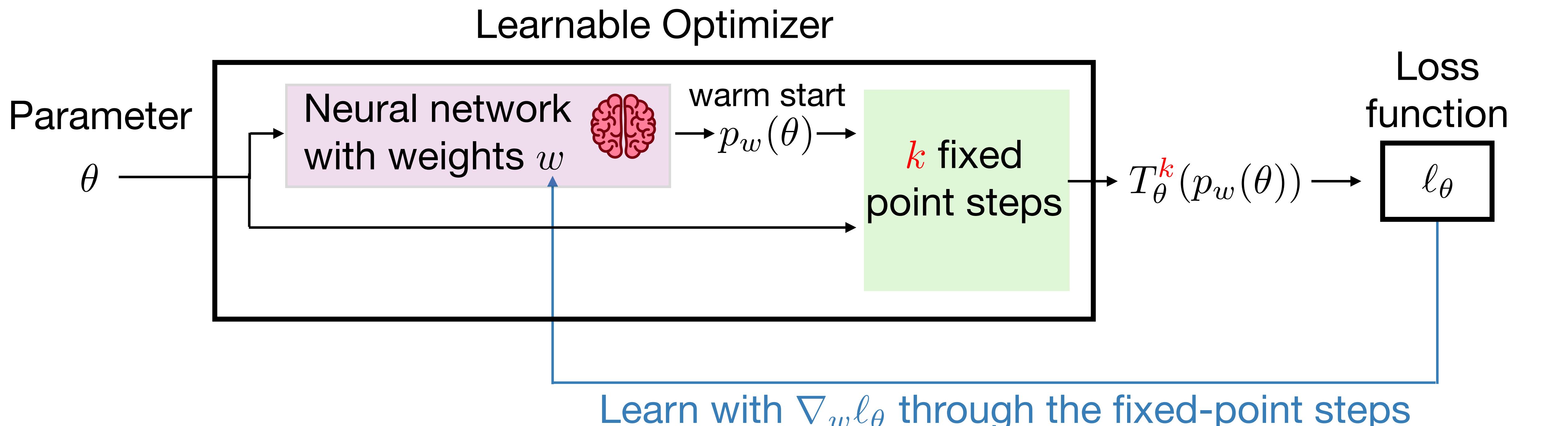
★ Optimal solution at the origin
Run proximal gradient descent to solve



All three warm starts appear to be
equally suboptimal but converge
at very different rates

The quality of the warm start depends on the algorithm

End-to-end learning architecture

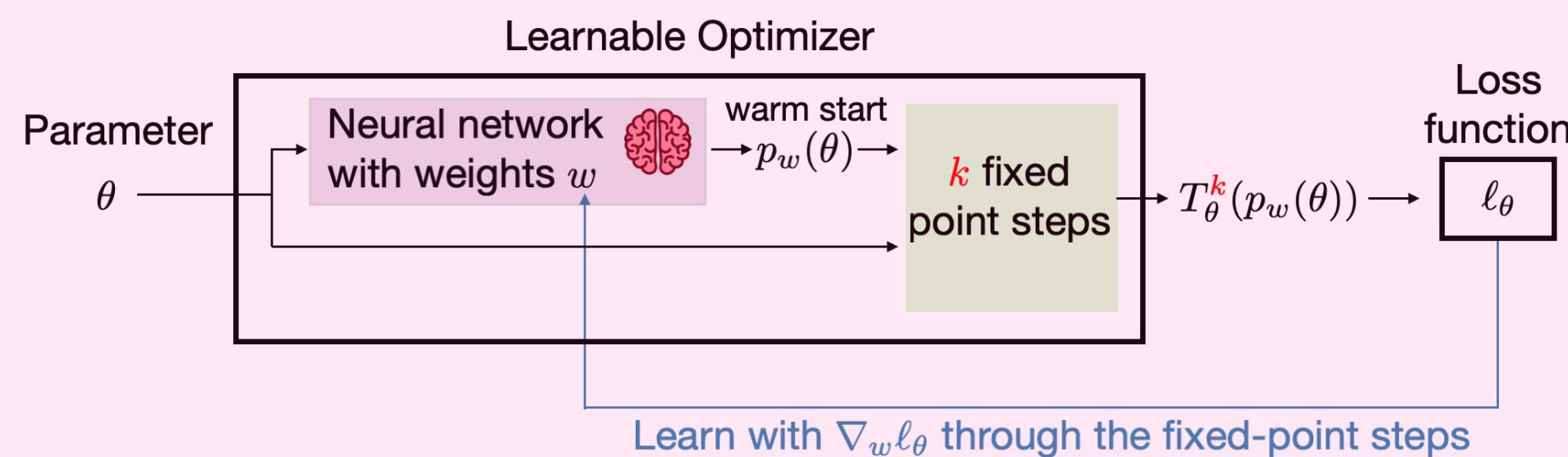


Loss function: $\ell_\theta(z) = \|z - z^*(\theta)\|_2$ **Ground truth solution**

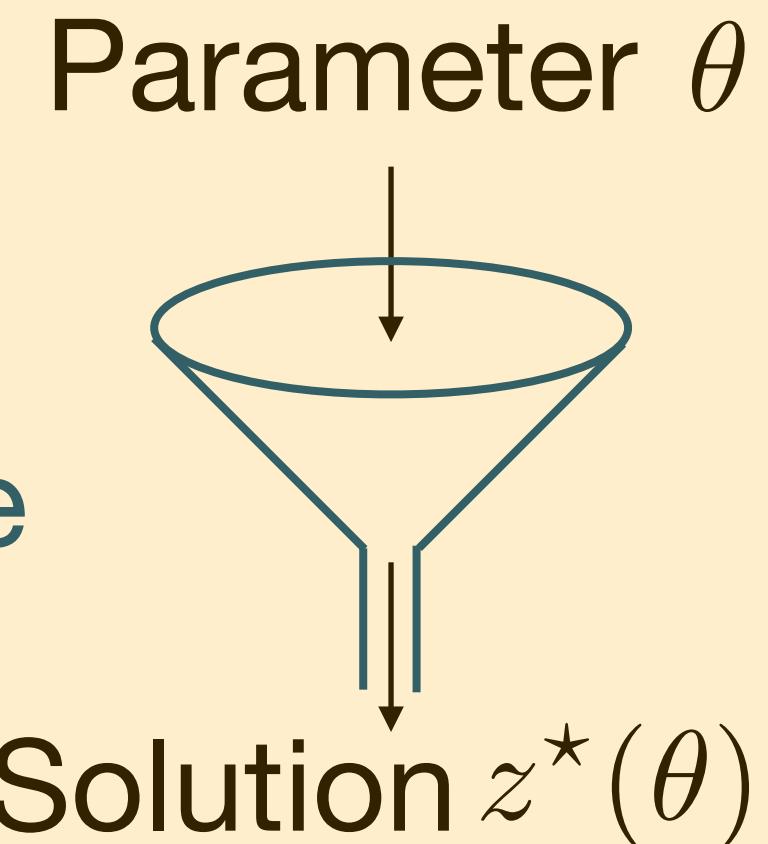
Learned warm start tailored for downstream algorithm

Benefits of our learning framework

End-to-end learning: warm-start predictions tailored to downstream algorithm



Guaranteed convergence



Generalization guarantees



- I. Guarantees from k training steps to t evaluation steps
- II. Guarantees to unseen data

Easy integration with popular solvers



Conic programs

```
sol = scs_solver.solve(warm_start=True,  
                      x=x0, y=y0, s=s0)
```

$$\begin{aligned} & \text{minimize} && (1/2)x^T P x + c^T x \\ & \text{subject to} && Ax + s = b \\ & && s \in \mathcal{K} \end{aligned}$$

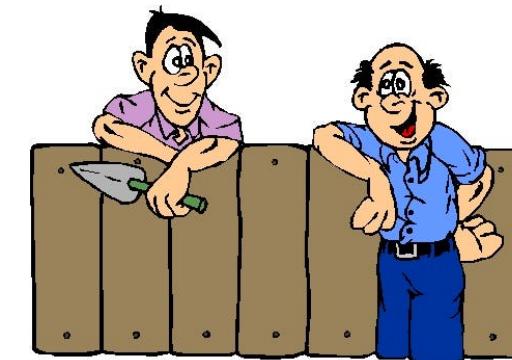
Allows us to quantify solve time in seconds

Numerical Experiments

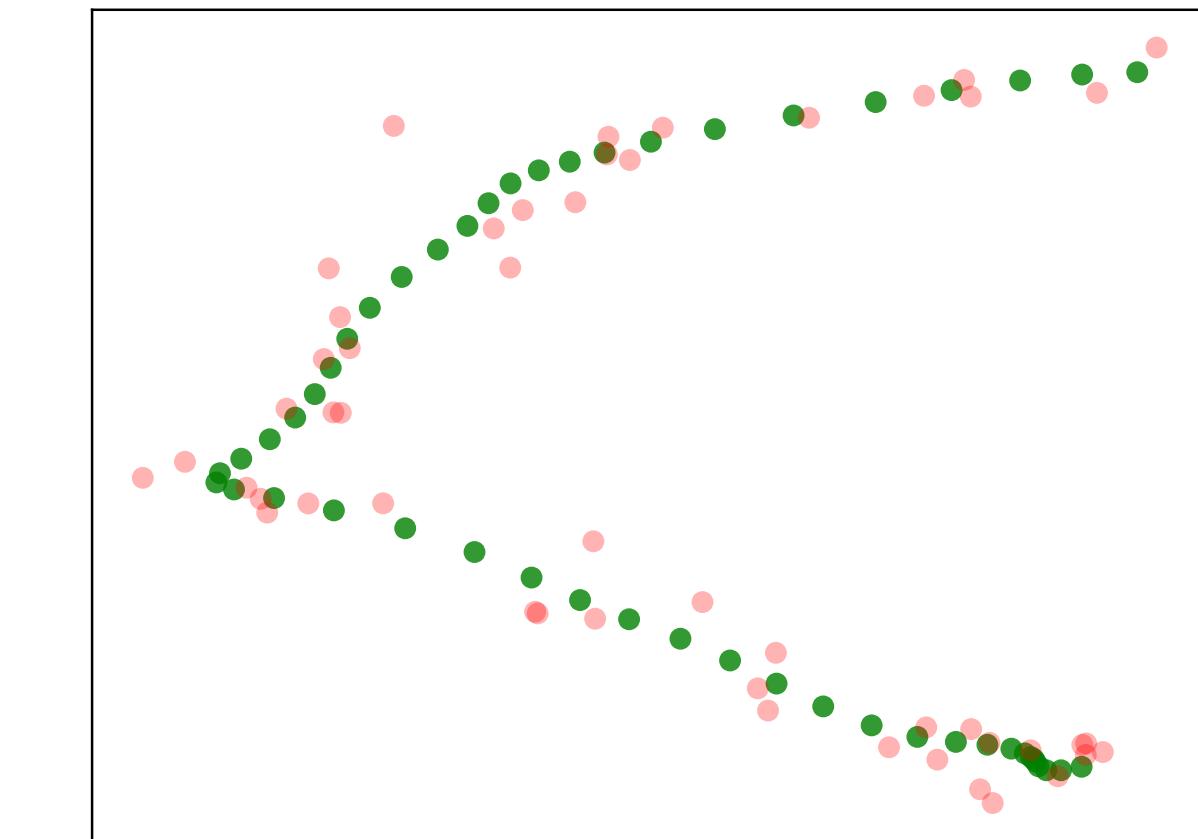
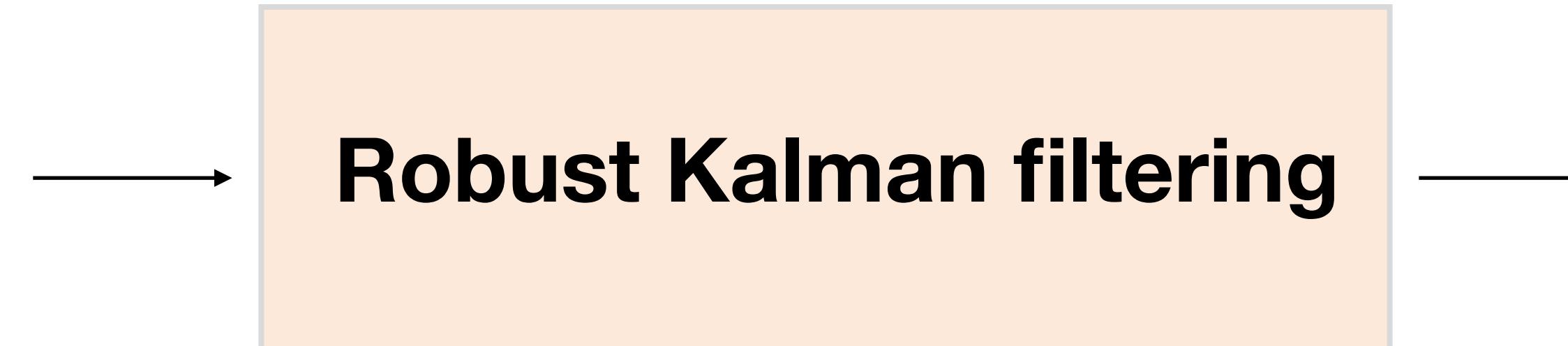
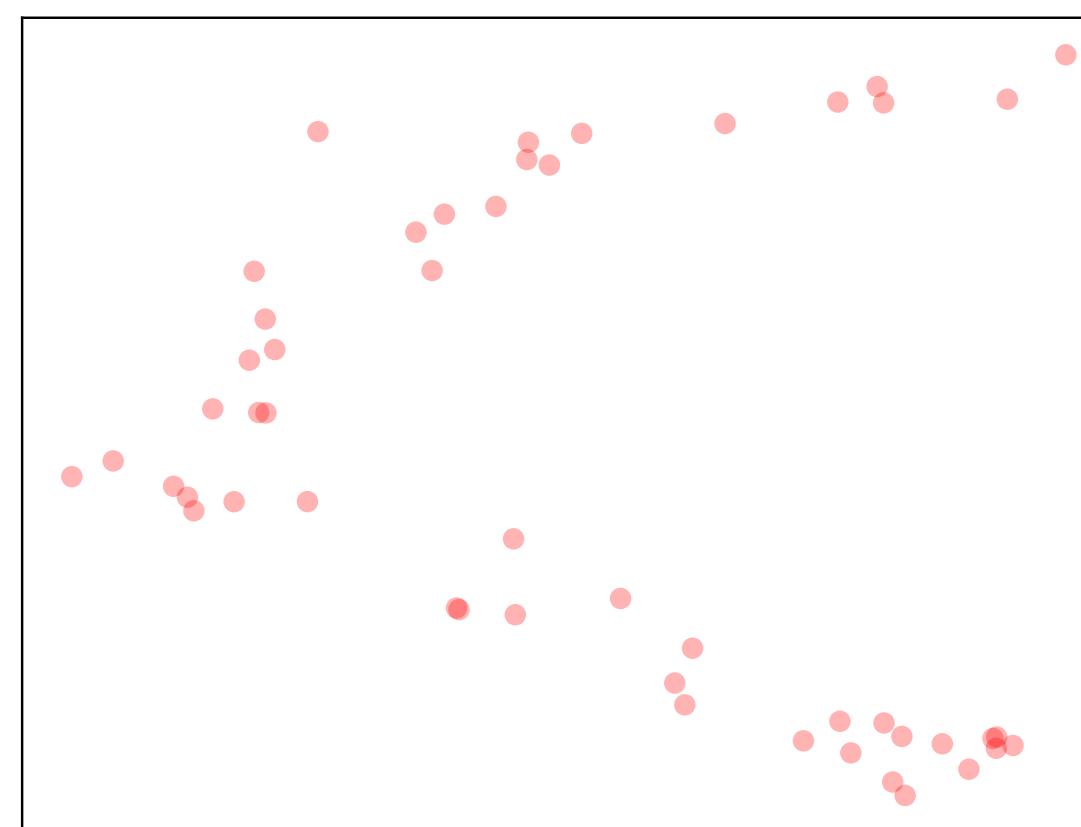
Comparing our learned warm starts  against

Baseline initializations

1. Cold-start: initialize at zero 
2. Nearest neighbor: initialize with solution of nearest training problem



Robust Kalman filtering



Second-order cone program

$$\theta = \{y_t\}_{t=0}^{T-1}$$

Noisy trajectory

minimize $\sum_{t=0}^{T-1} \|w_t\|_2^2 + \mu\psi_\rho(v_t)$

subject to $x_{t+1} = Ax_t + Bw_t \quad \forall t$
 $y_t = Cx_t + v_t \quad \forall t$

$$\{x_t^*, w_t^*, v_t^*\}_{t=0}^{T-1}$$

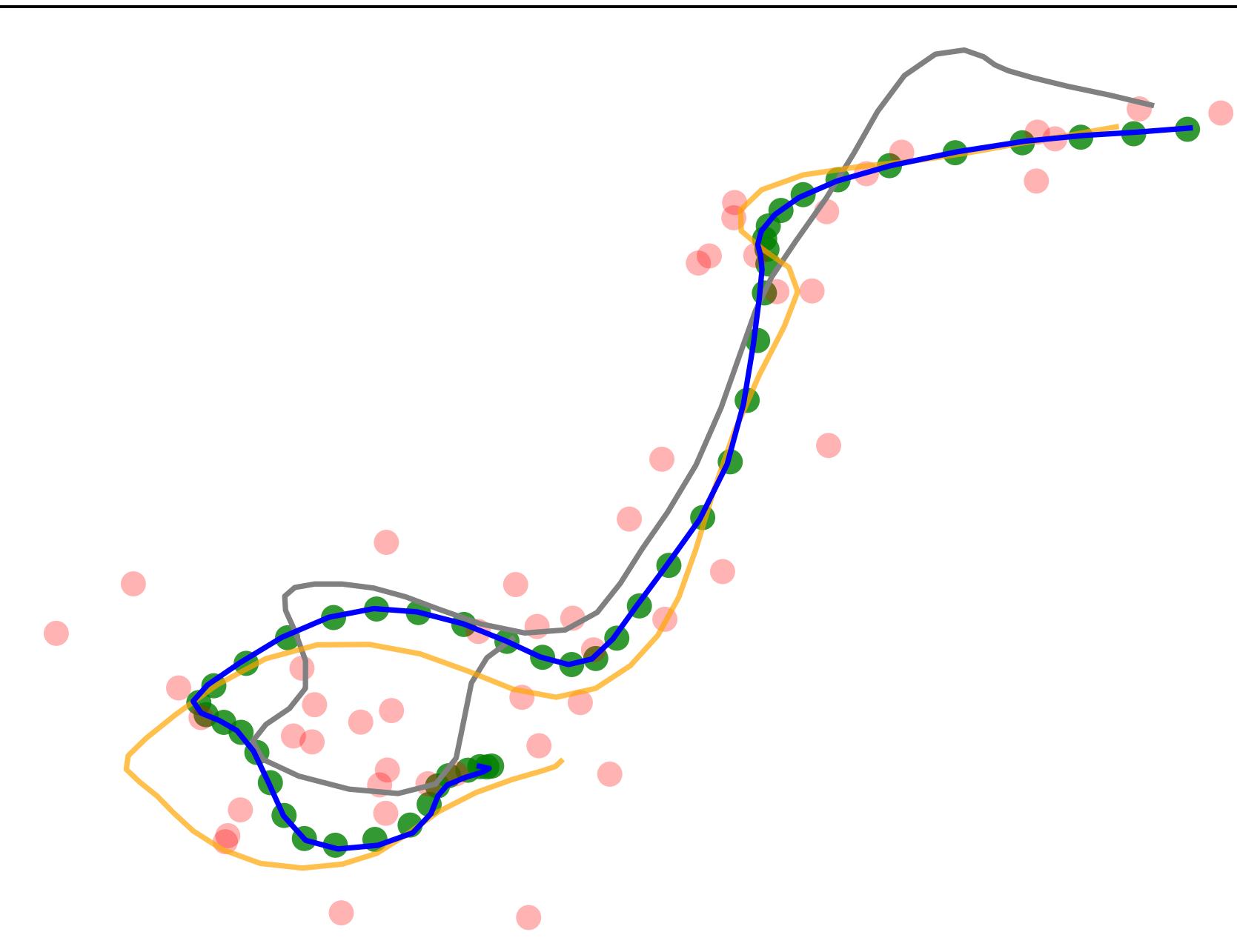
Recovered trajectory

Dynamics matrices: A, B

Observation matrix: C

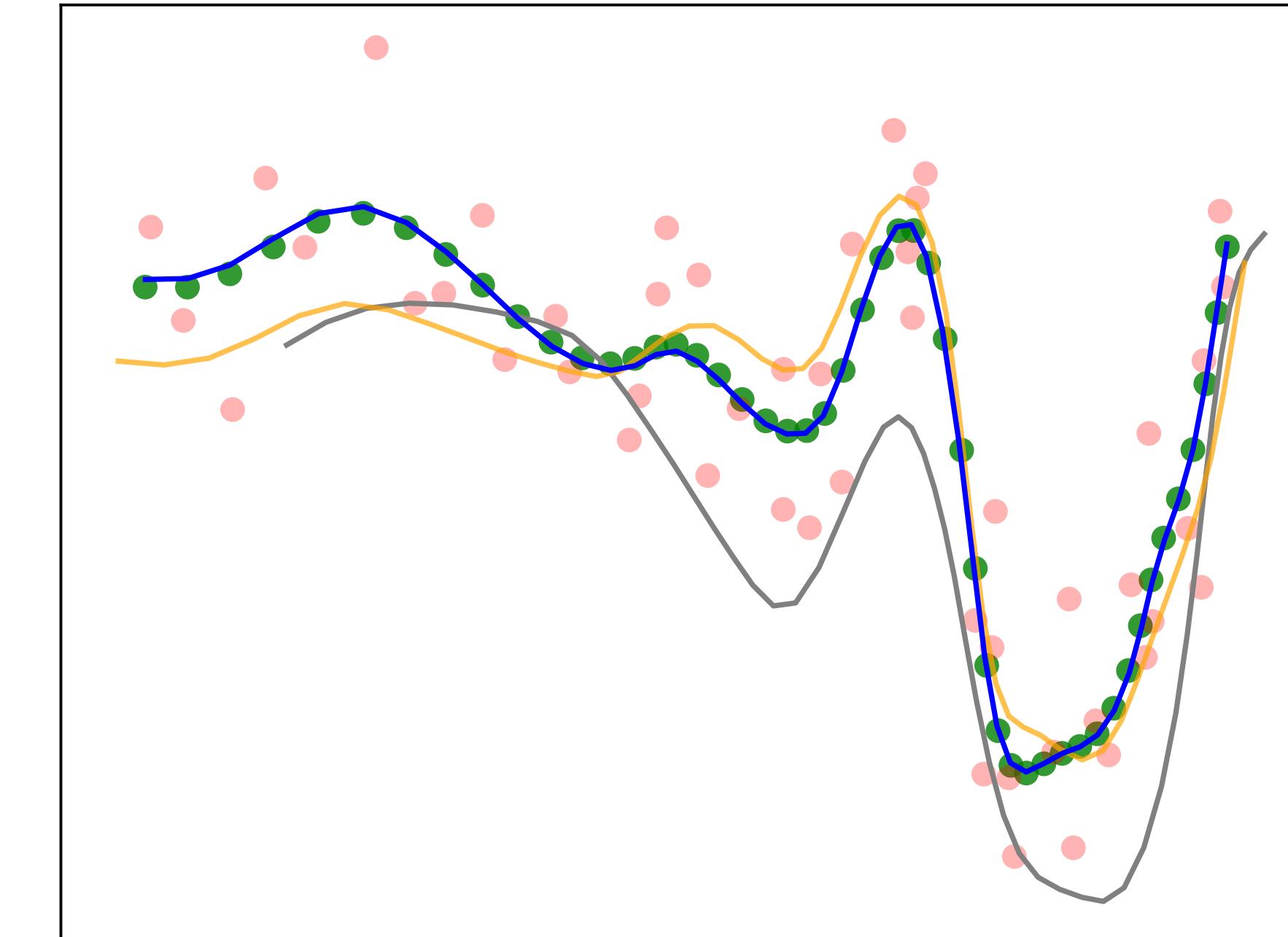
Huber loss: ψ_ρ

Robust Kalman filtering visuals



- Noisy trajectory
- Optimal solution

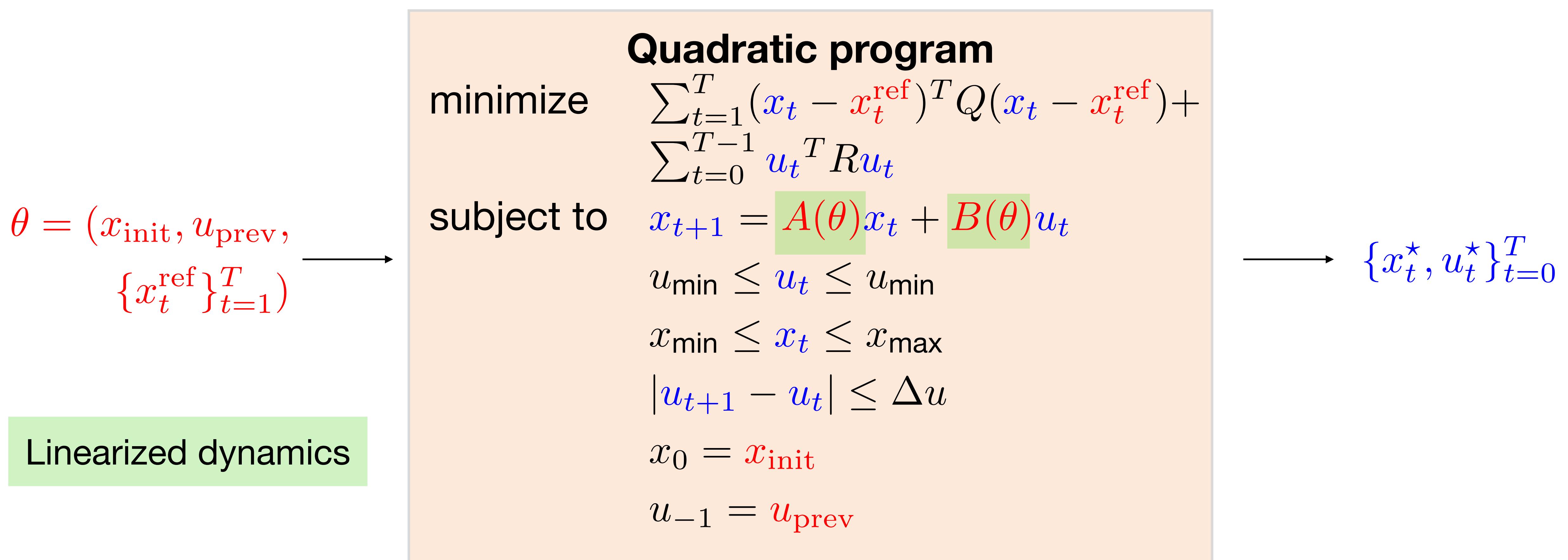
With learning, we can estimate the state well



Solution after 5 fixed-point steps
with different initializations

- Nearest neighbor 
- Previous solution 
- Learned: $k = 5$ 

Model predictive control (MPC) of a quadcopter

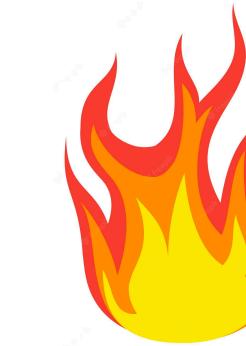
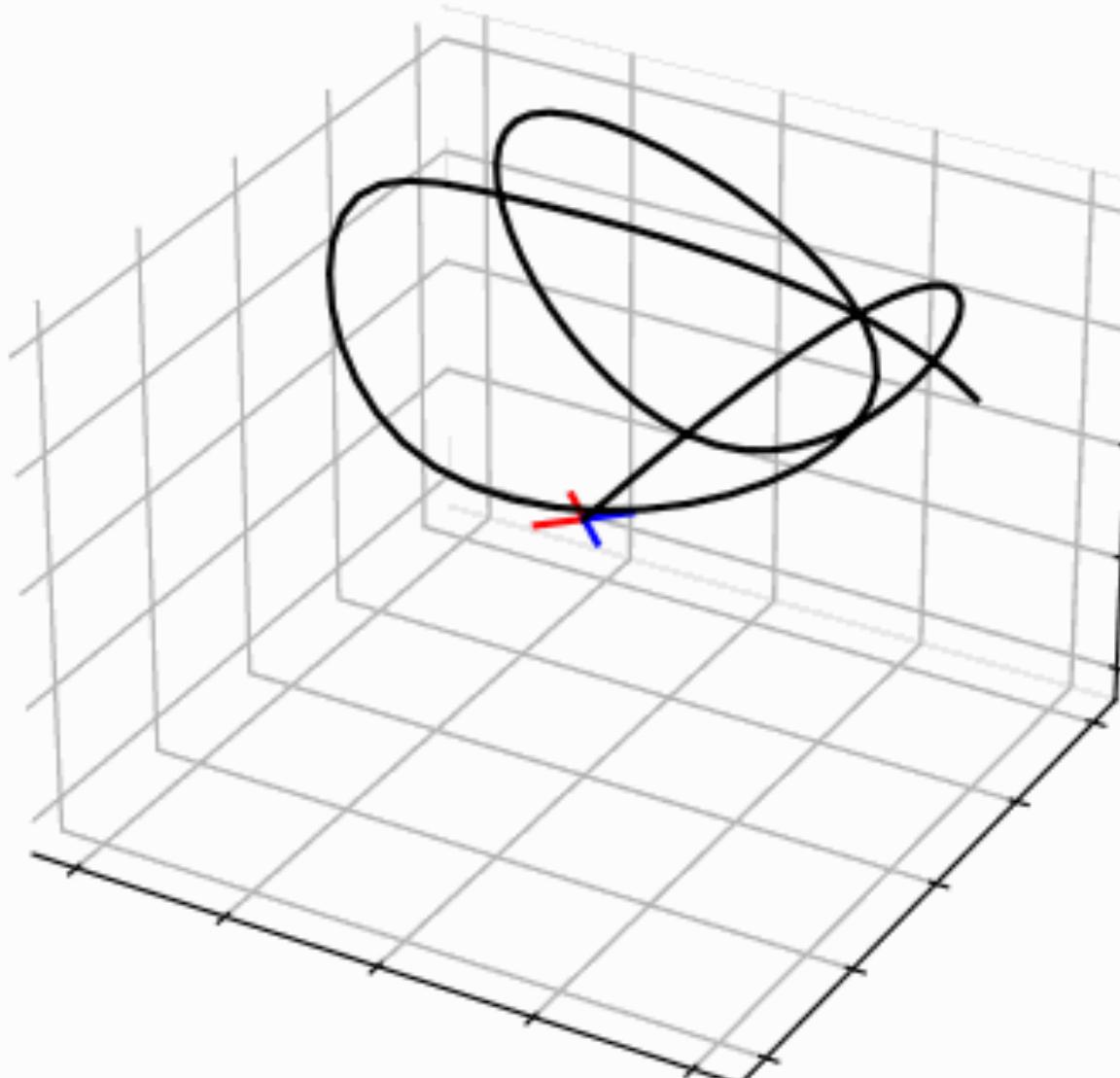


MPC of a quadcopter in a closed loop

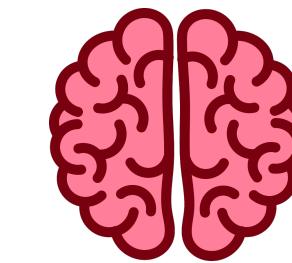
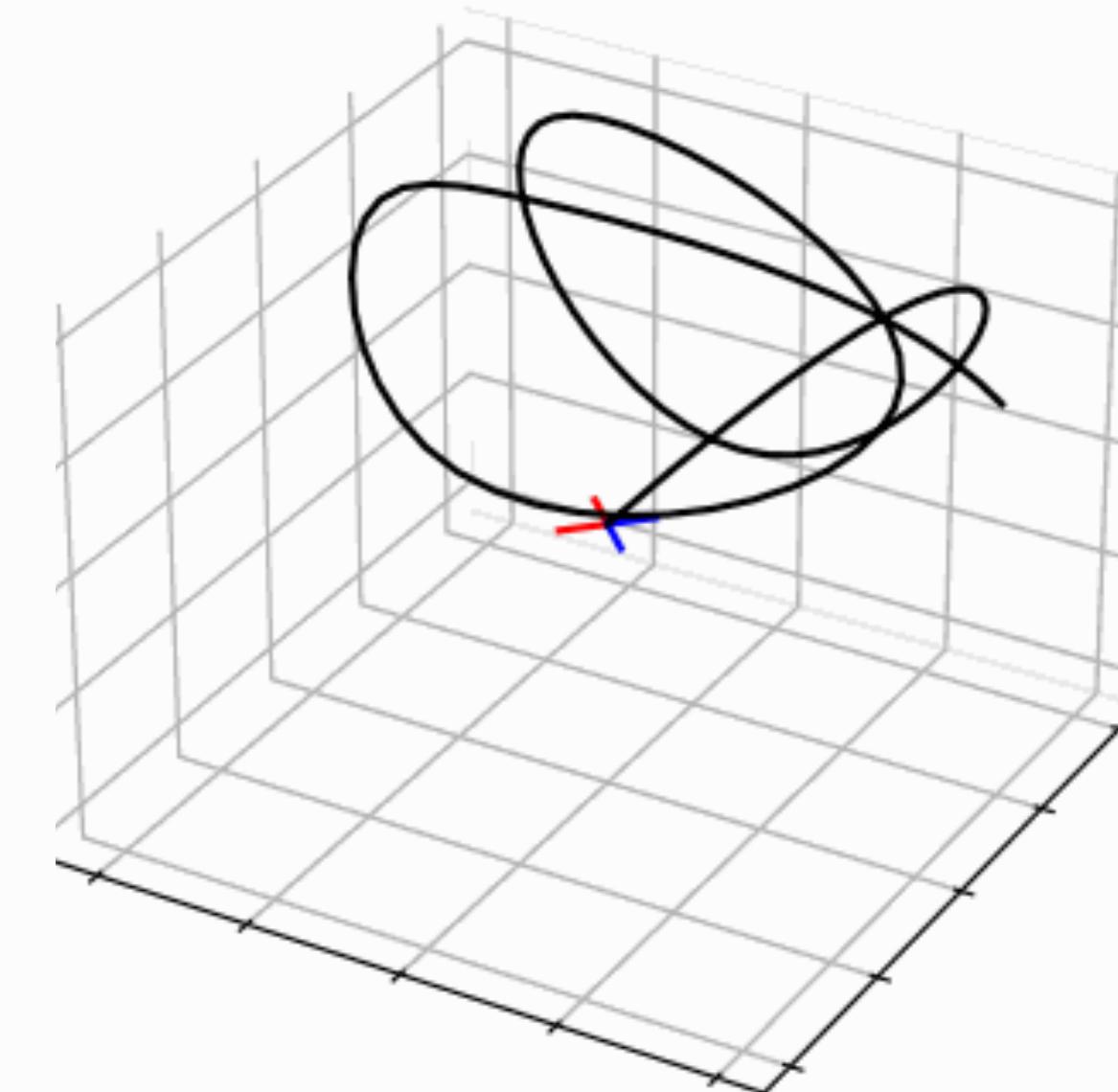
Budget of 15 fixed-point steps



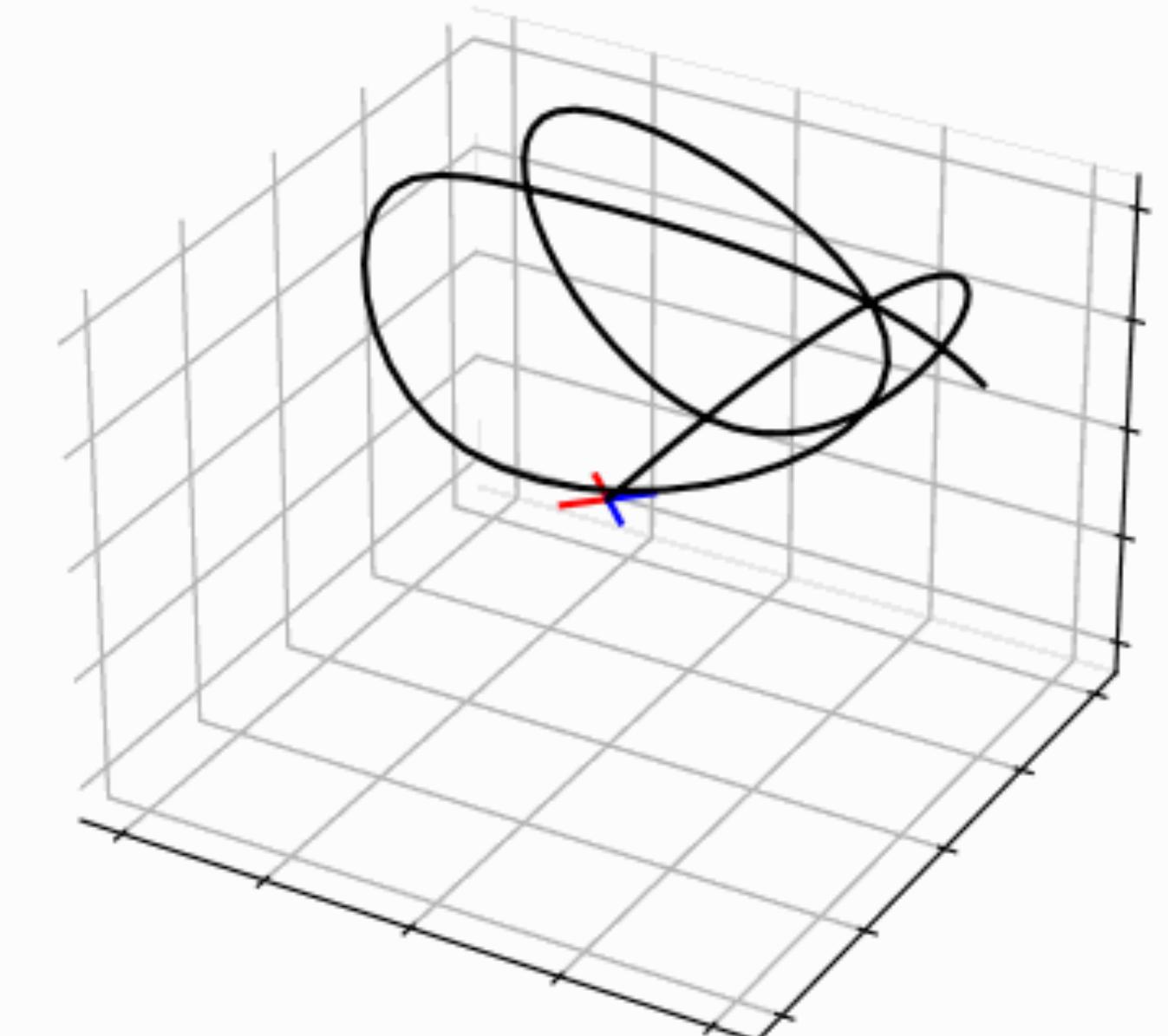
Nearest neighbor



Previous solution



Learned: $k = 5$



With learning, we can track the trajectory well

Image deblurring

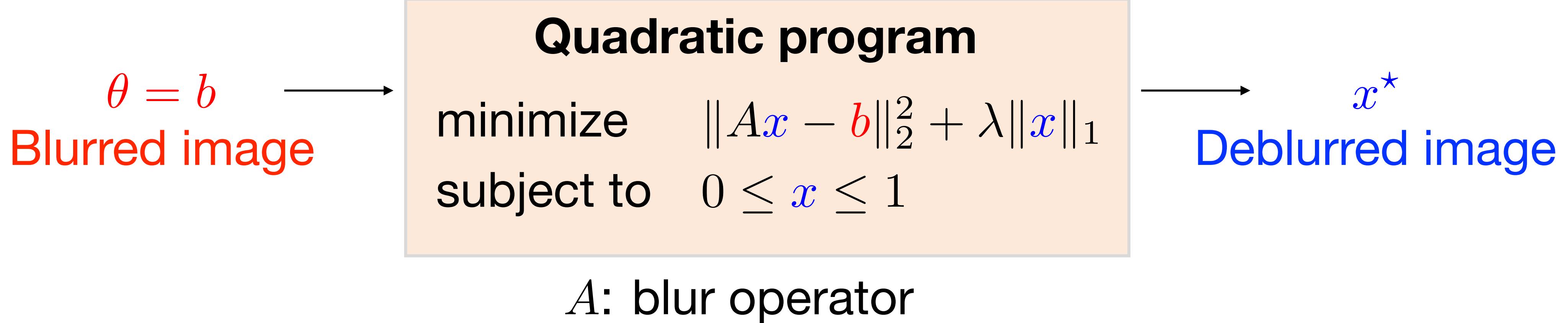
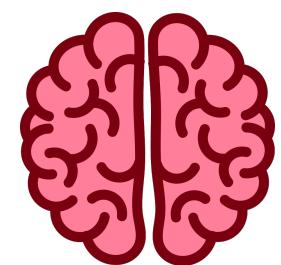


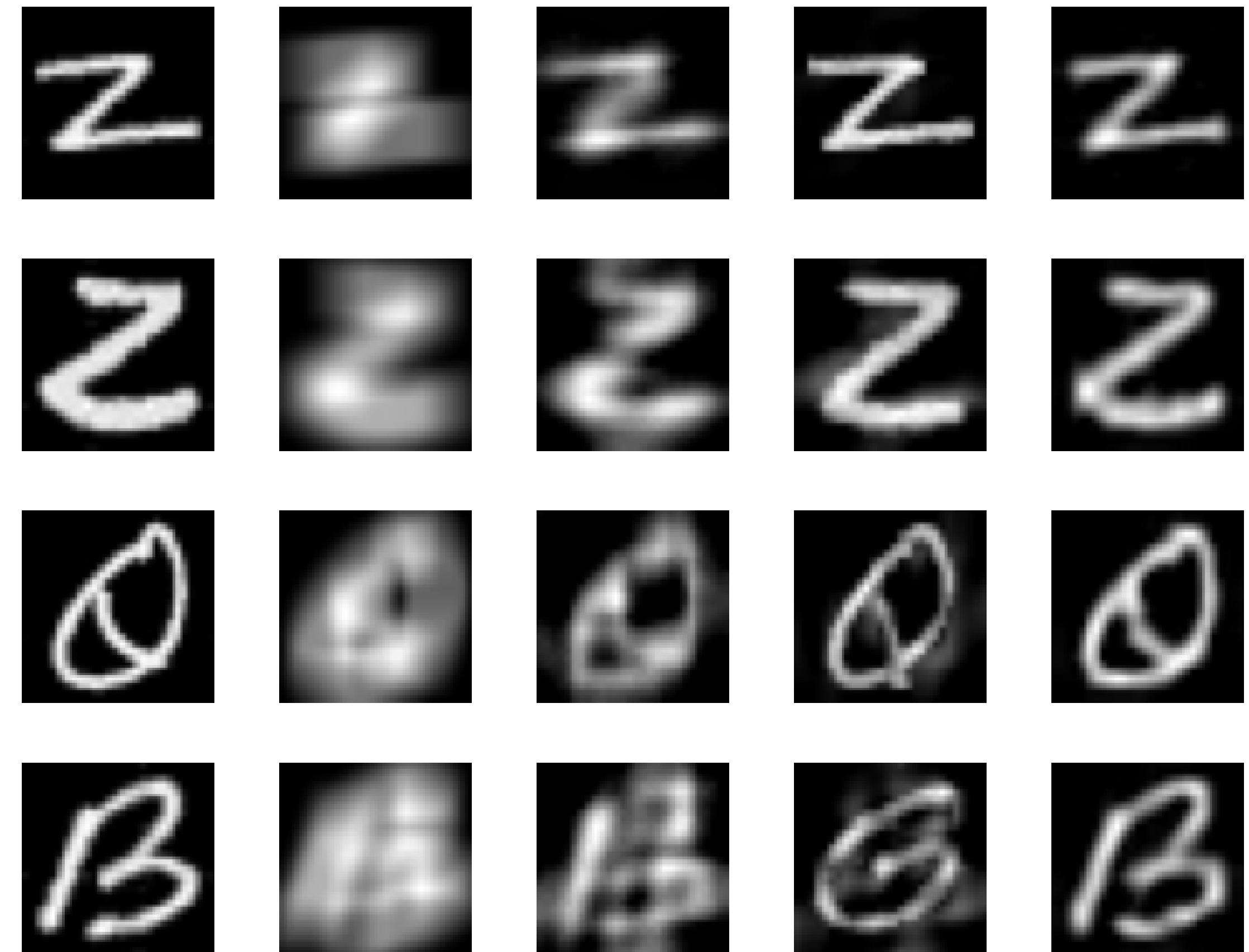
Image deblurring



50 fixed-point steps

Distance to nearest neighbor increases

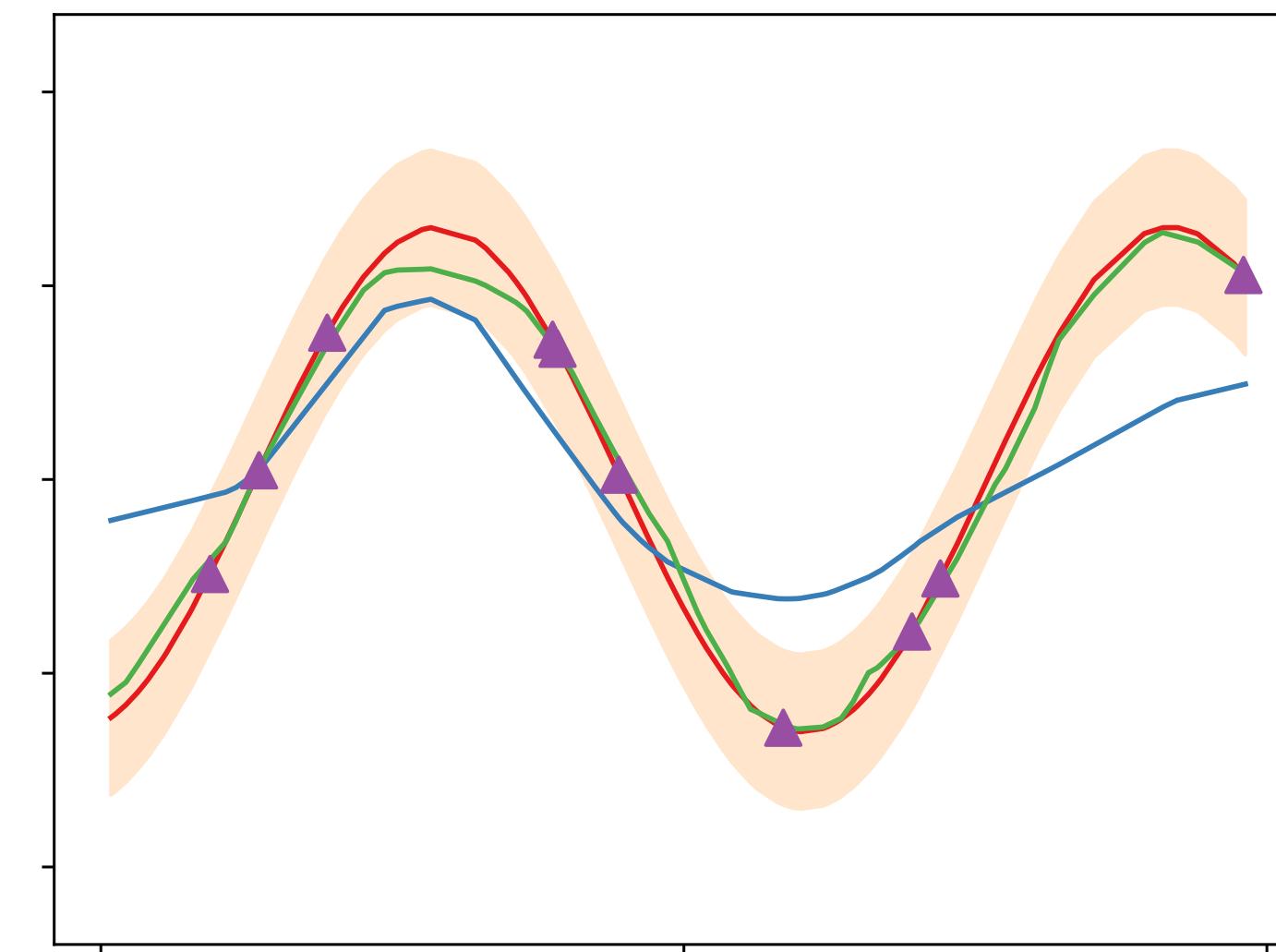
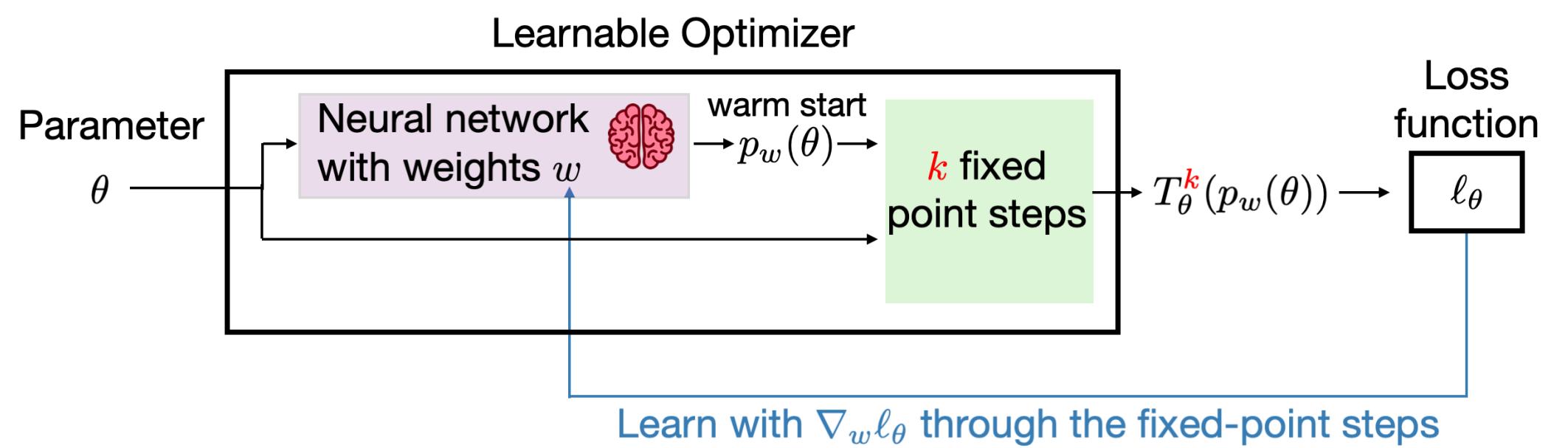
percentile optimal blurred cold-start nearest learned
neighbor



With learning, we can deblur all of the images quickly

Talk Outline

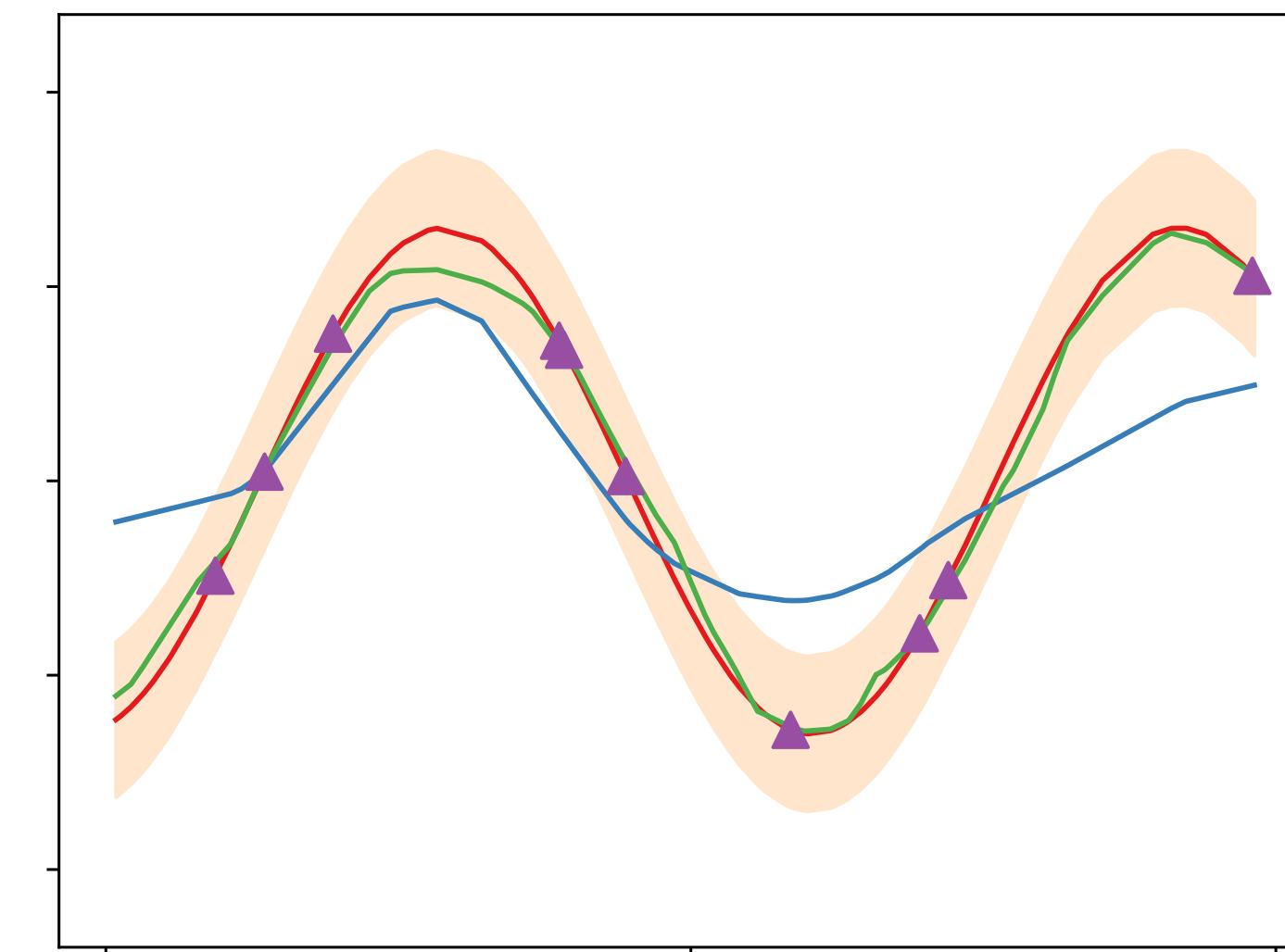
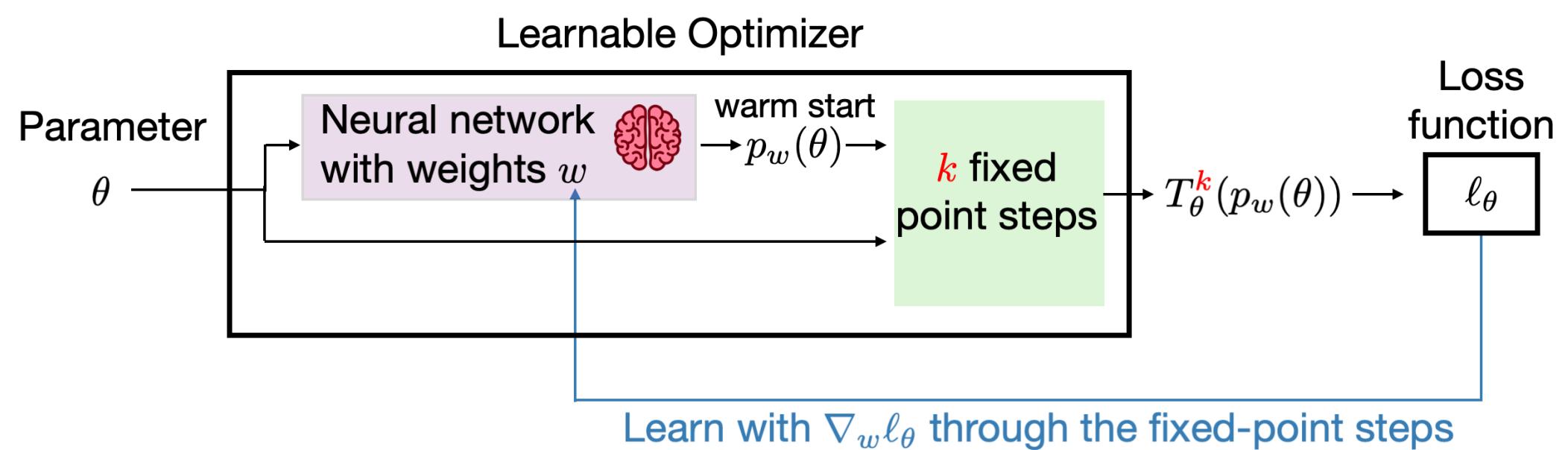
- Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms
- Part 2: Data-Driven Performance Guarantees for Classical and Learned Optimizers



Talk Outline

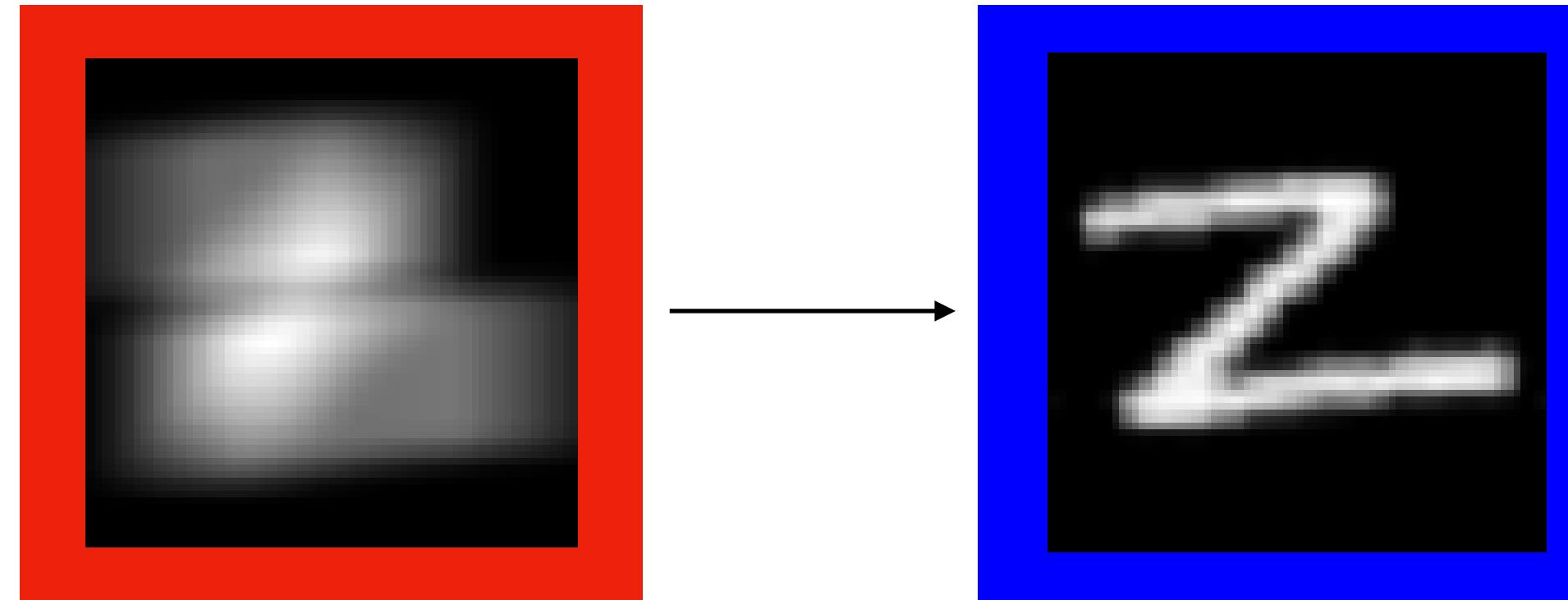
- Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms
- Part 2: Data-Driven Performance Guarantees for Classical and Learned Optimizers

Classical = no learning



Worst-case bounds can be very loose

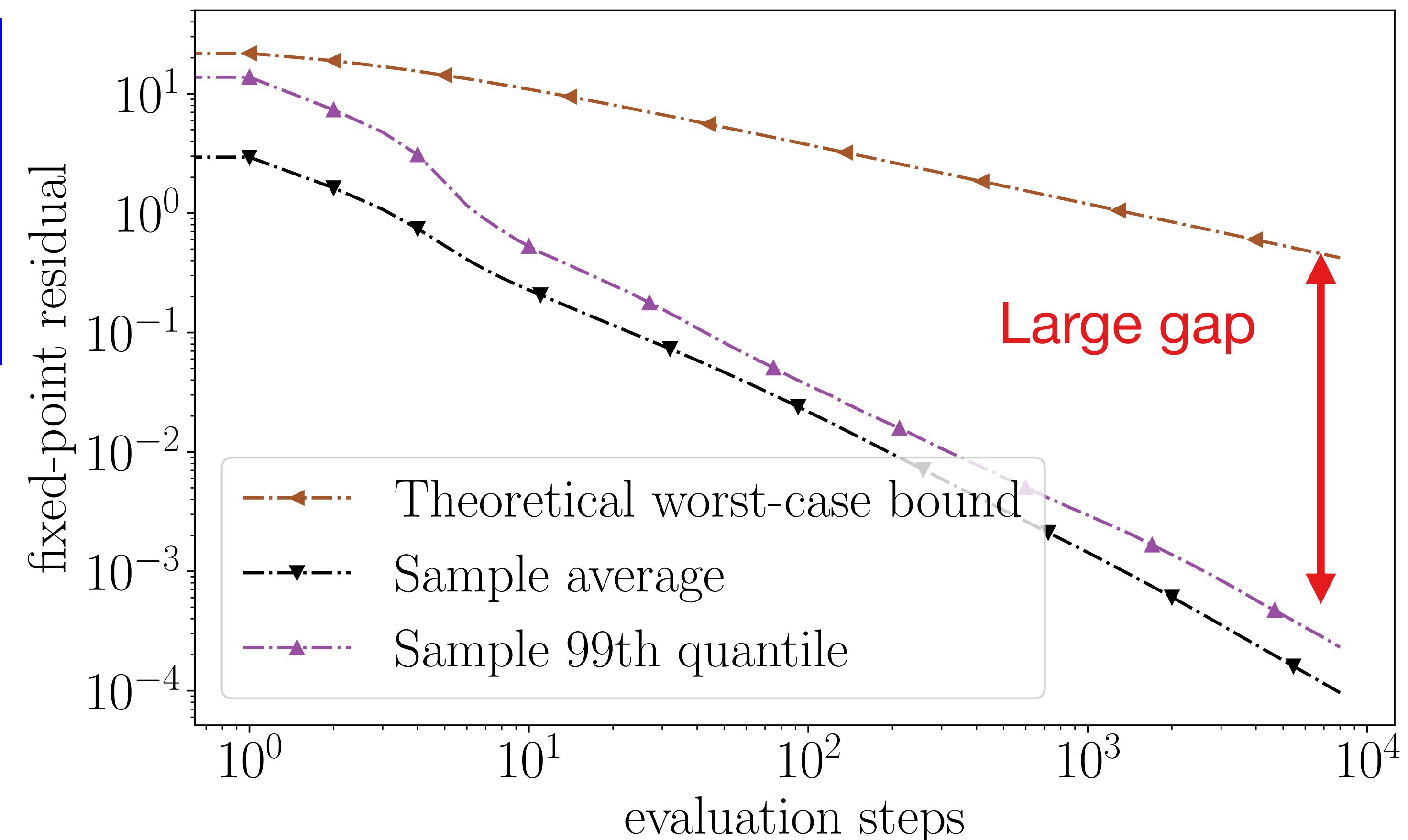
Example: image deblurring



Quadratic program

minimize $\|Ax - b\|_2^2 + \lambda\|\mathbf{x}\|_1$
subject to $0 \leq \mathbf{x} \leq 1$

1000 problems solved with



Worst-case bounds are pessimistic and do not consider the parametric structure

Approach: probabilistically bound over the parametric distribution

Our recipe for guarantees for classical optimizers

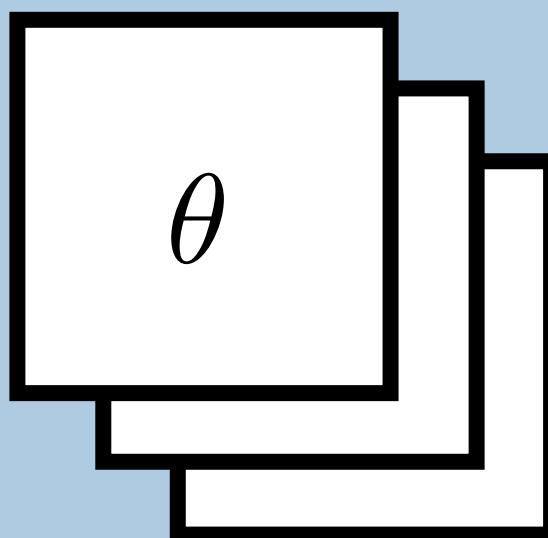
algorithm steps
 $e(\theta) = \mathbf{1}(\ell^k(\theta) > \epsilon)$
tolerance
Any metric
(e.g., fixed-point residual)

Step 1
Run k steps
for N parametric problems

Step 2
Evaluate the empirical risk

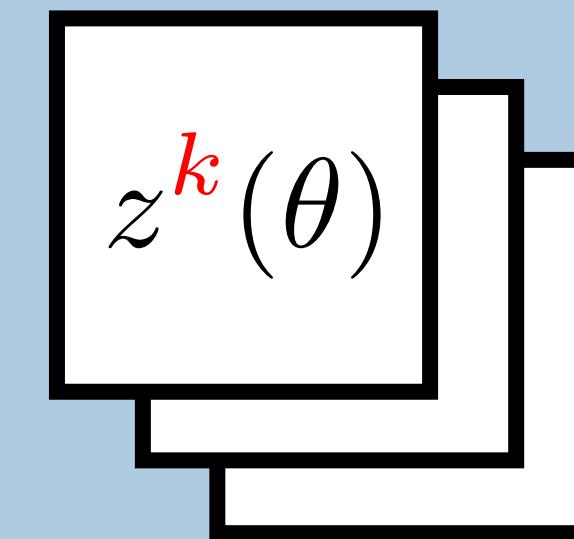
Step 3
Bound the risk
(Next slide)

Parameters



Candidate solutions

Run k steps



$$\frac{1}{N} \sum_{i=1}^N e(\theta_i)$$

$$\text{risk} = \mathbf{E}_{\theta \sim \mathcal{X}} e(\theta) \leq \text{bound}$$

Statistical learning theory can provide probabilistic guarantees

$$e(\theta) = \mathbf{1}(\ell^k(\theta) > \epsilon)$$

algorithm steps →
tolerance ↓

Sample convergence bound: with probability $1 - \delta$ [Langford et. al 2001]

$$\mathbf{E}_{\theta \sim \mathcal{X}} e(\theta) \leq \text{KL}^{-1} \left(\frac{1}{N} \sum_{i=1}^N e(\theta_i) \middle| \frac{\log(2/\delta)}{N} \right)$$

$\mathbf{P}(\ell^k(\theta) > \epsilon) = \text{risk} \leq \text{KL}^{-1} (\text{empirical risk} \mid \text{regularizer})$

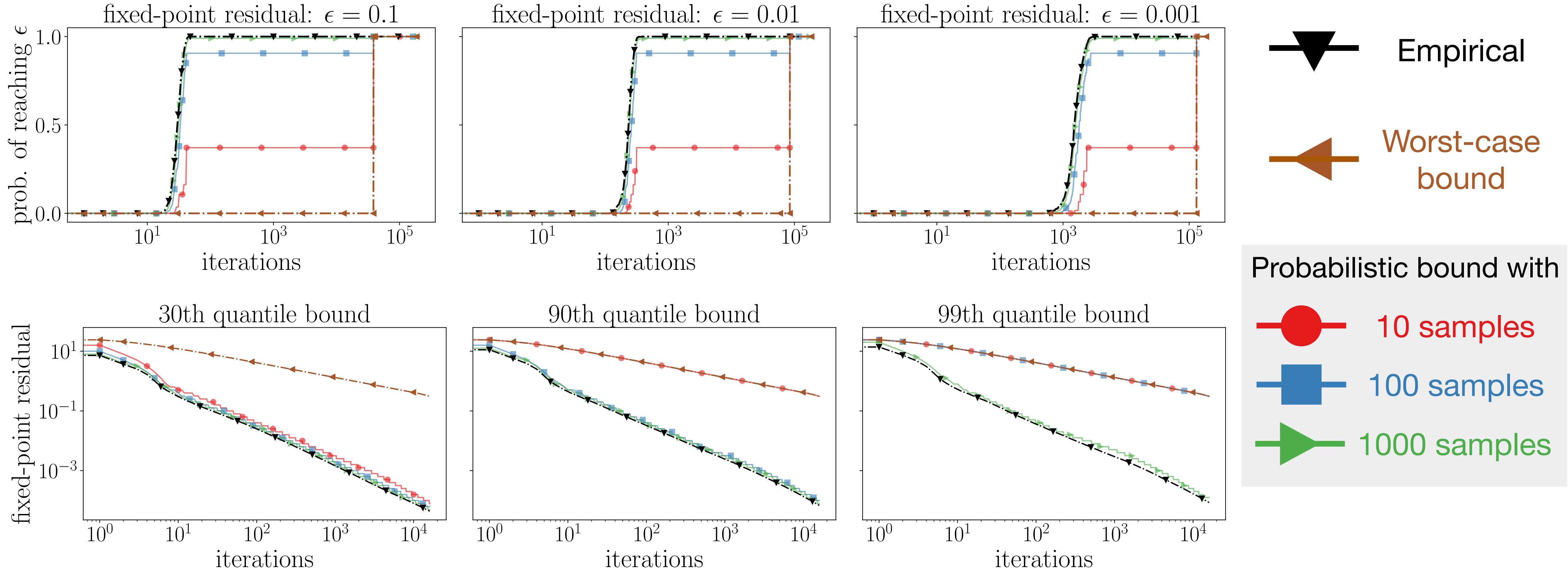
Number of problems

1D convex
optimization problem

```
graph TD; A["E_theta ~ X e(theta) ≤ KL⁻¹ (1/N ∑ᵢ¹⁹⁸⁵ e(θᵢ) | log(2/δ)/N)"] -- "Number of problems" --> B["P(ℓᵏ(θ) > ε) = risk ≤ KL⁻¹ (empirical risk | regularizer)"]; B -- "1D convex optimization problem" --> C["1D convex optimization problem"]
```

"With probability $1 - \delta$, 90% of the time the fixed-point residual is below $\epsilon = 0.01$ after $k = 20$ steps"

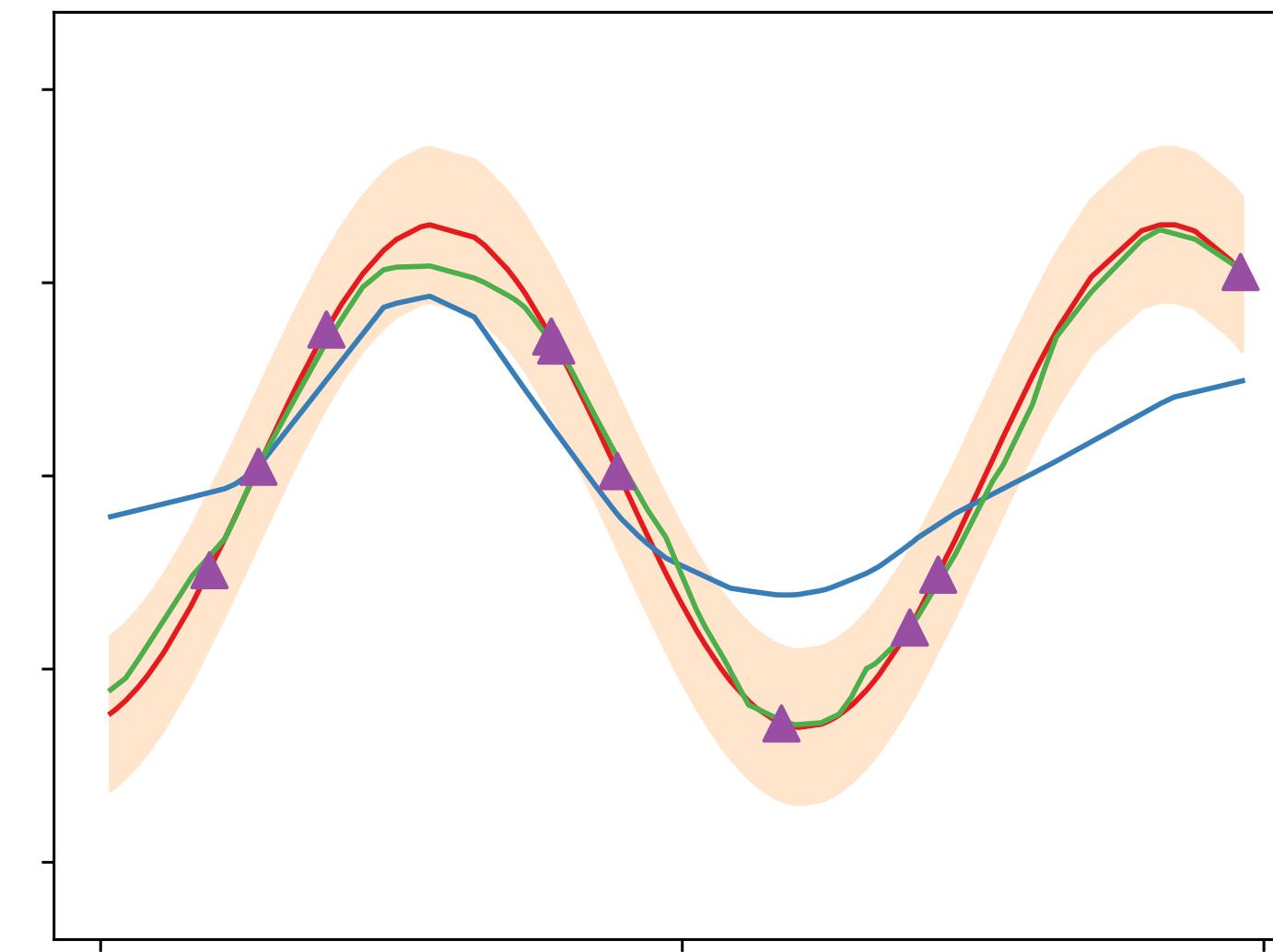
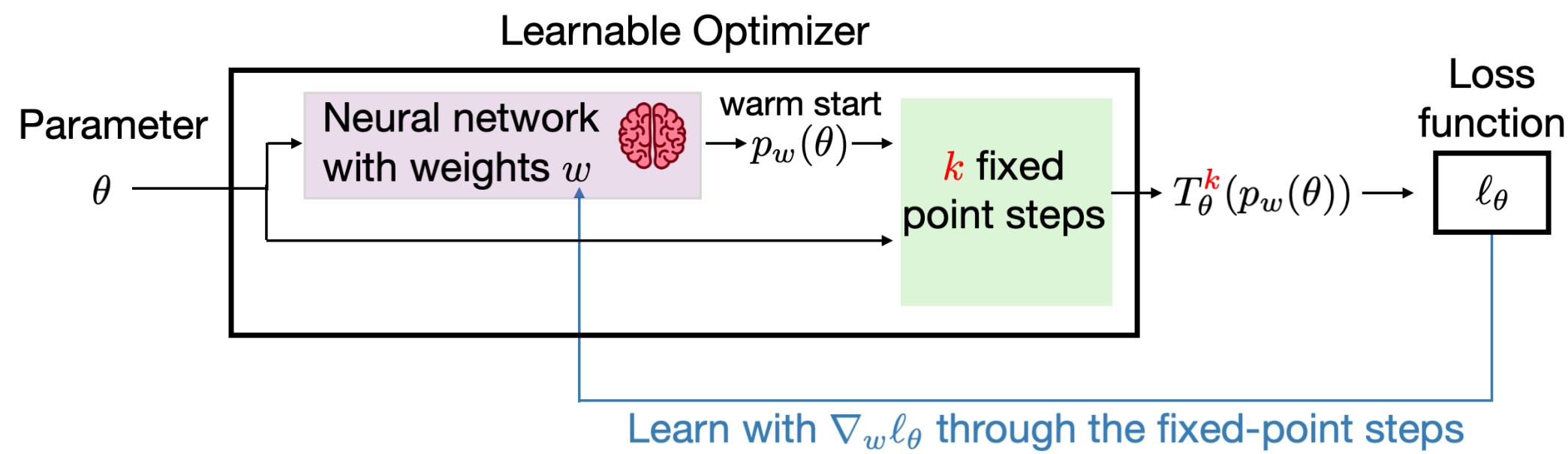
Image deblurring guarantees



With 1000 samples, we provide strong probabilistic guarantees on the 99th quantile

Talk Outline

- Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms
- Part 2: Data-Driven Performance Guarantees for Classical and Learned Optimizers



Tutorial on Amortized Optimization [Amos 2023]

“Despite having the capacity of surpassing the convergence rates of other algorithms, oftentimes in practice amortized optimization methods can deeply struggle to generalize and converge to reasonable solutions.”

PAC-Bayes guarantees for learned optimizers

$$e_w(\theta) = \mathbf{1}(\ell_w^k(\theta) > \epsilon)$$

algorithm steps
tolerance
learnable weights

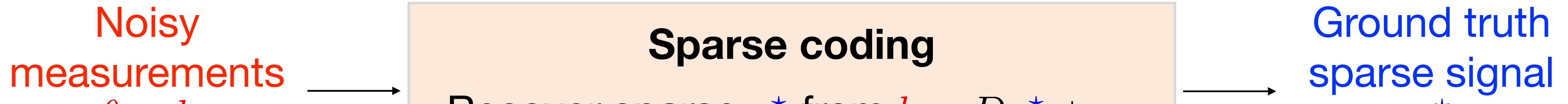
McAllester bound: given posterior and prior distributions [McAllester et. al 2003]
 P and P_0 , with probability $1 - \delta$

$$\mathbf{E}_{\theta \sim \mathcal{X}} \mathbf{E}_{w \sim P} e_w(\theta) \leq \text{KL}^{-1} \left(\frac{1}{N} \sum_{i=1}^N \mathbf{E}_{w \sim P} e_w(\theta_i) \middle| \frac{1}{N} (\text{KL}(P \parallel P_0) + \log(N/\delta)) \right)$$

risk $\leq \text{KL}^{-1} (\text{empirical risk} \mid \text{regularizer})$

Optimize the bounds directly

Learned algorithms for sparse coding



D : dictionary, σ : noise

Standard technique

$$\text{minimize } \|Dz - b\|_2^2 + \lambda \|z\|_1$$

ISTA (iterative shrinkage thresholding algorithm)
(Classical optimizer)

$$z^{j+1} = \text{soft threshold}_{\frac{\lambda}{L}} \left(z^j - \frac{1}{L} (Dz^j - b) \right)$$

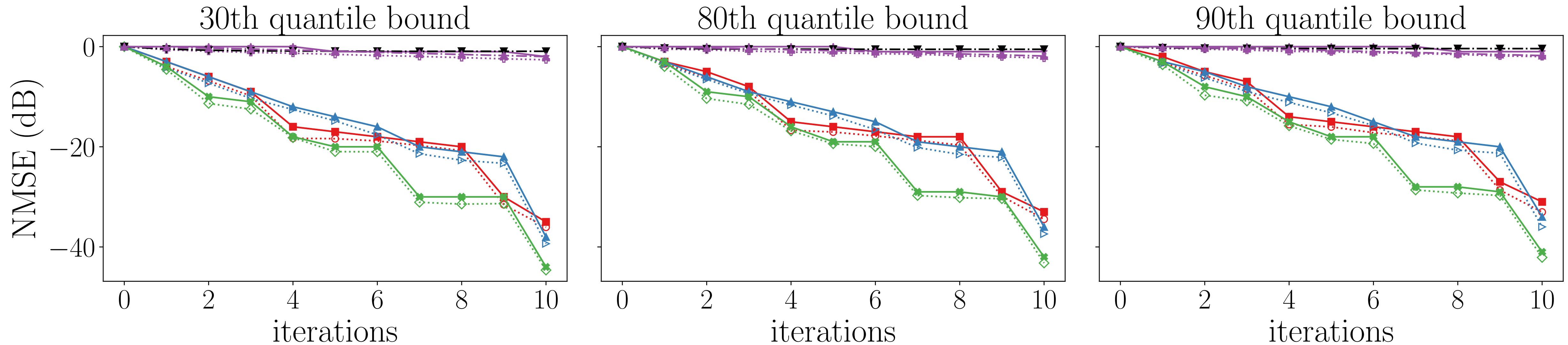
$$\text{soft threshold}_\psi(z) = \mathbf{sign}(z) \max(0, |z| - \psi)$$

Learned ISTA
(Learned optimizer)

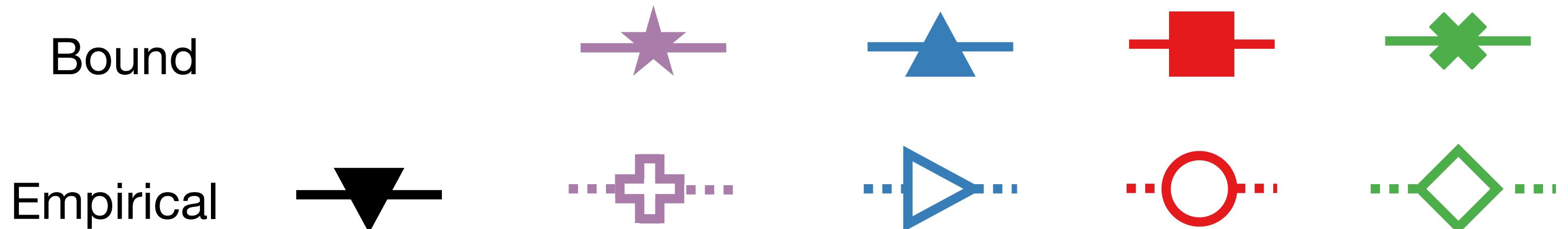
$$z^{j+1} = \text{soft threshold}_{\psi^j} \left(W_1^j z^j + W_2^j b \right)$$

+ variants [Gregor and LeCun 2010, Liu et. al 2019]

Learned ISTA results for sparse coding



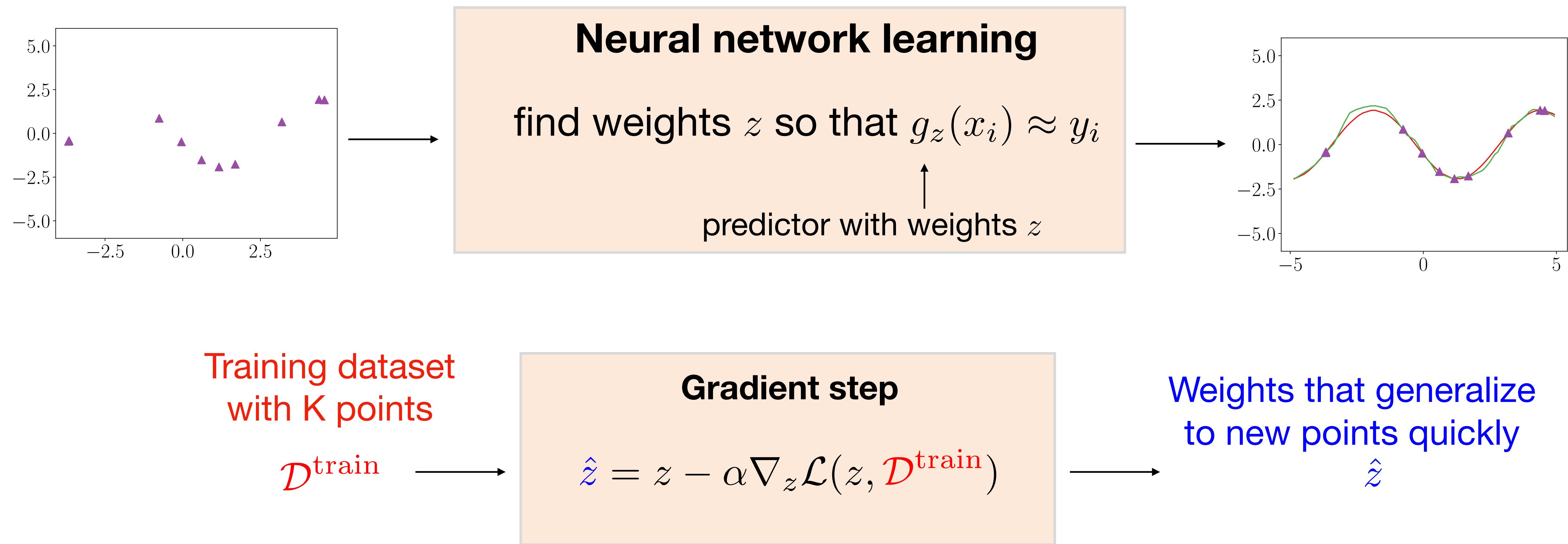
Not learned		Learned		
ISTA	LISTA	ALISTA	TiLISTA	GLISTA



Our bounds are close to empirical performance

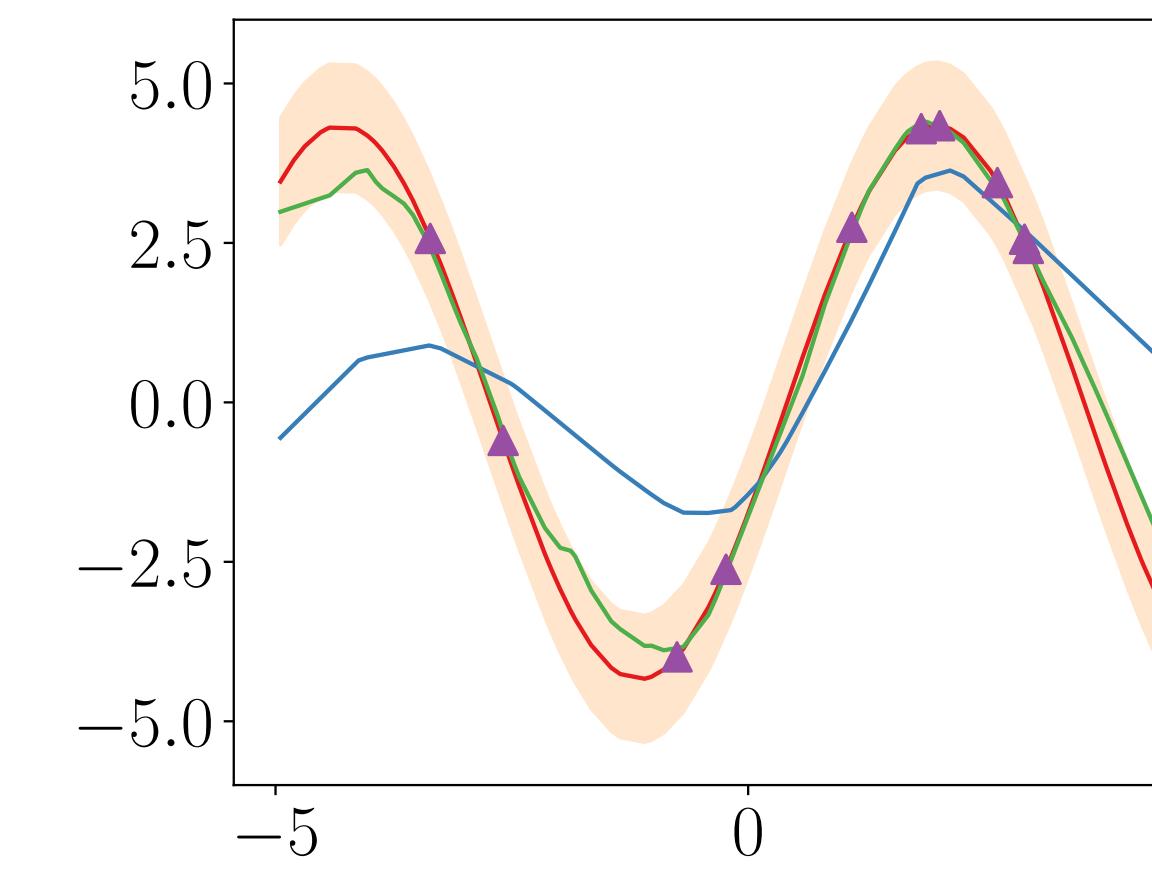
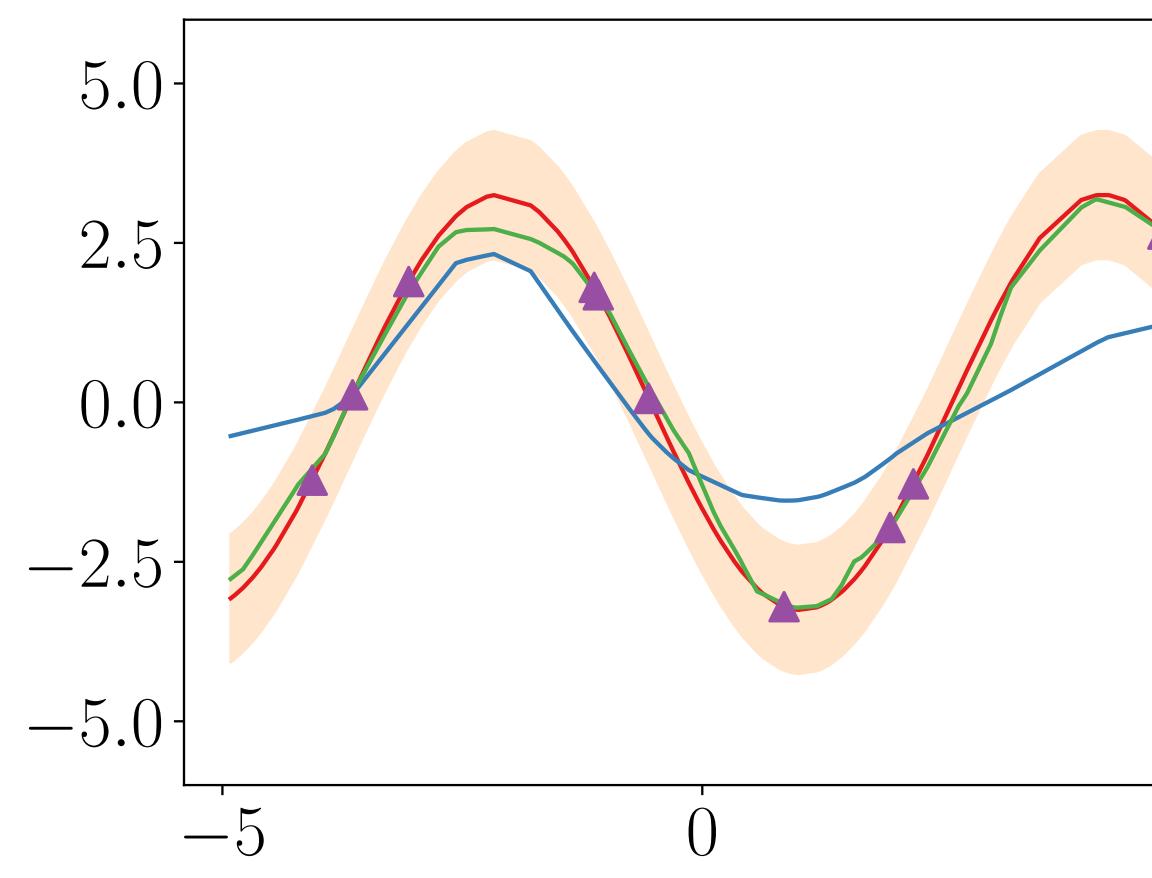
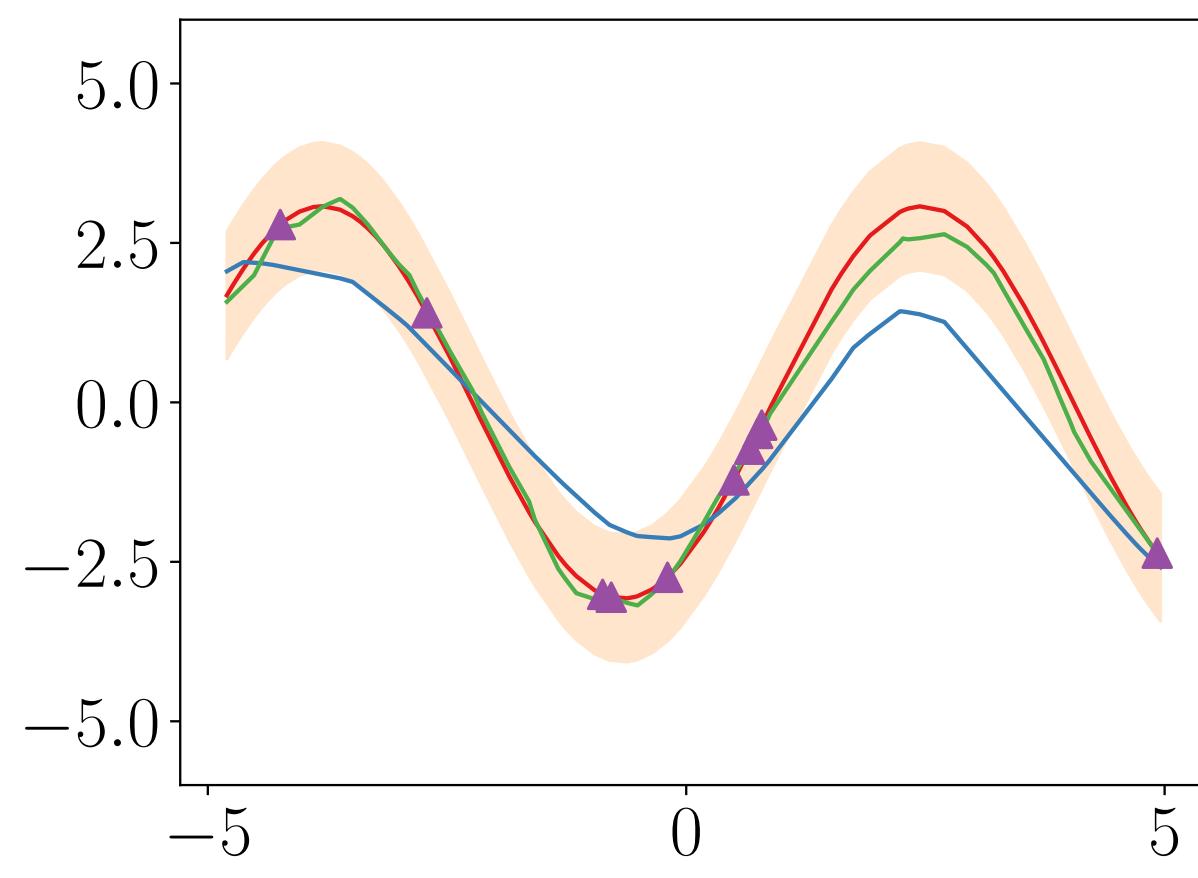
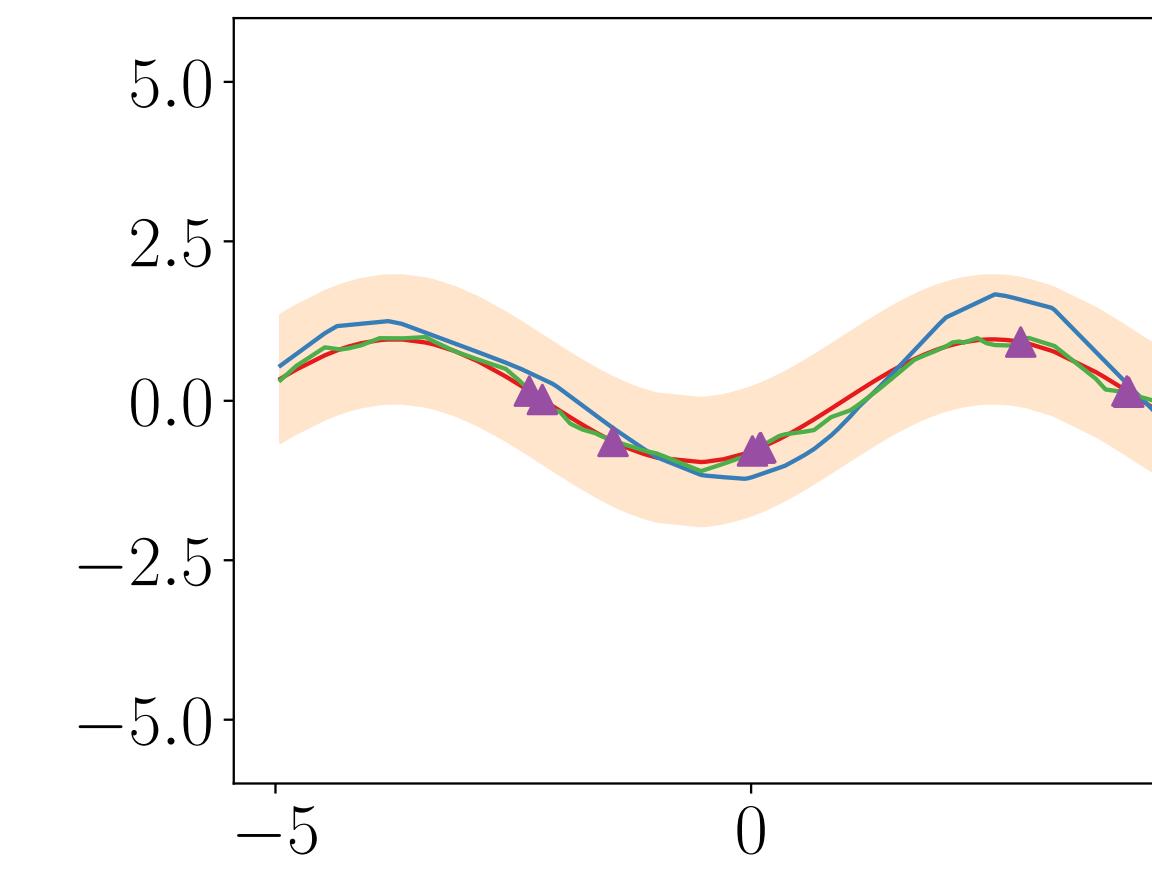
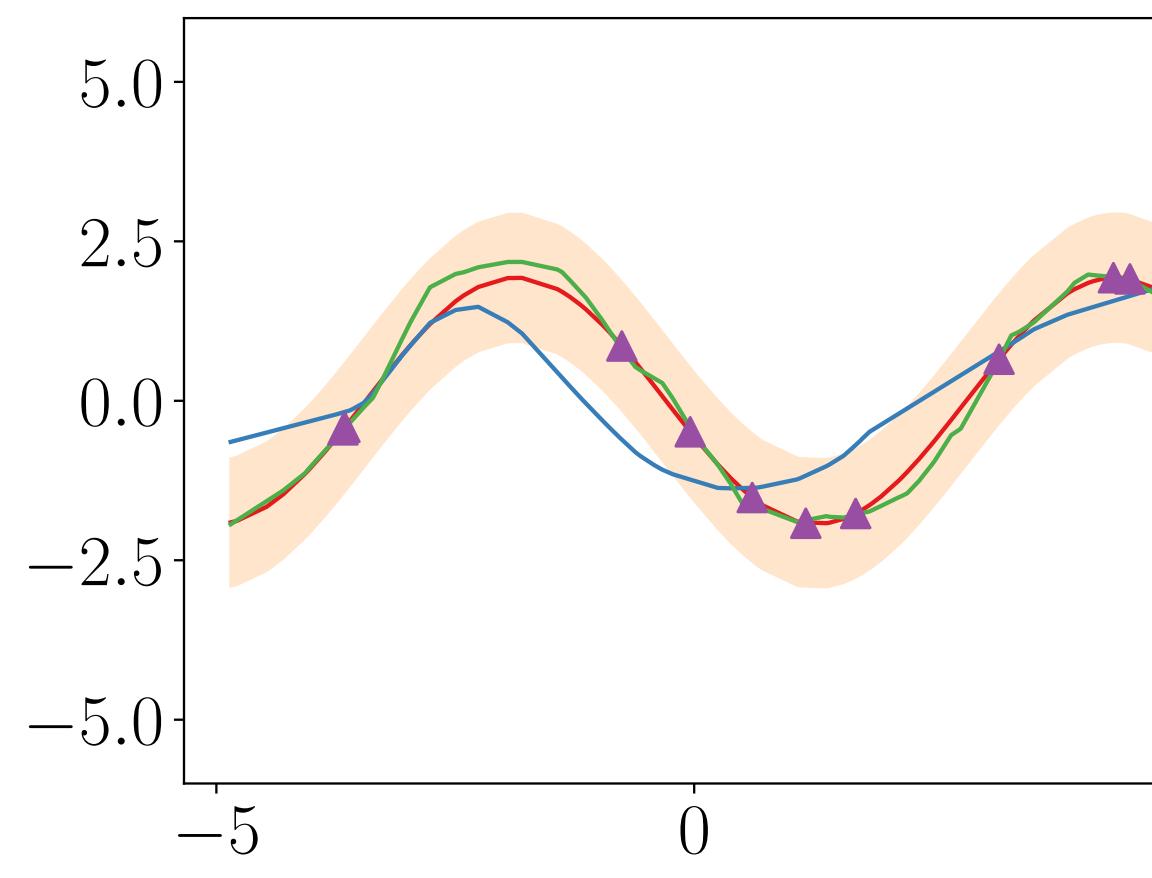
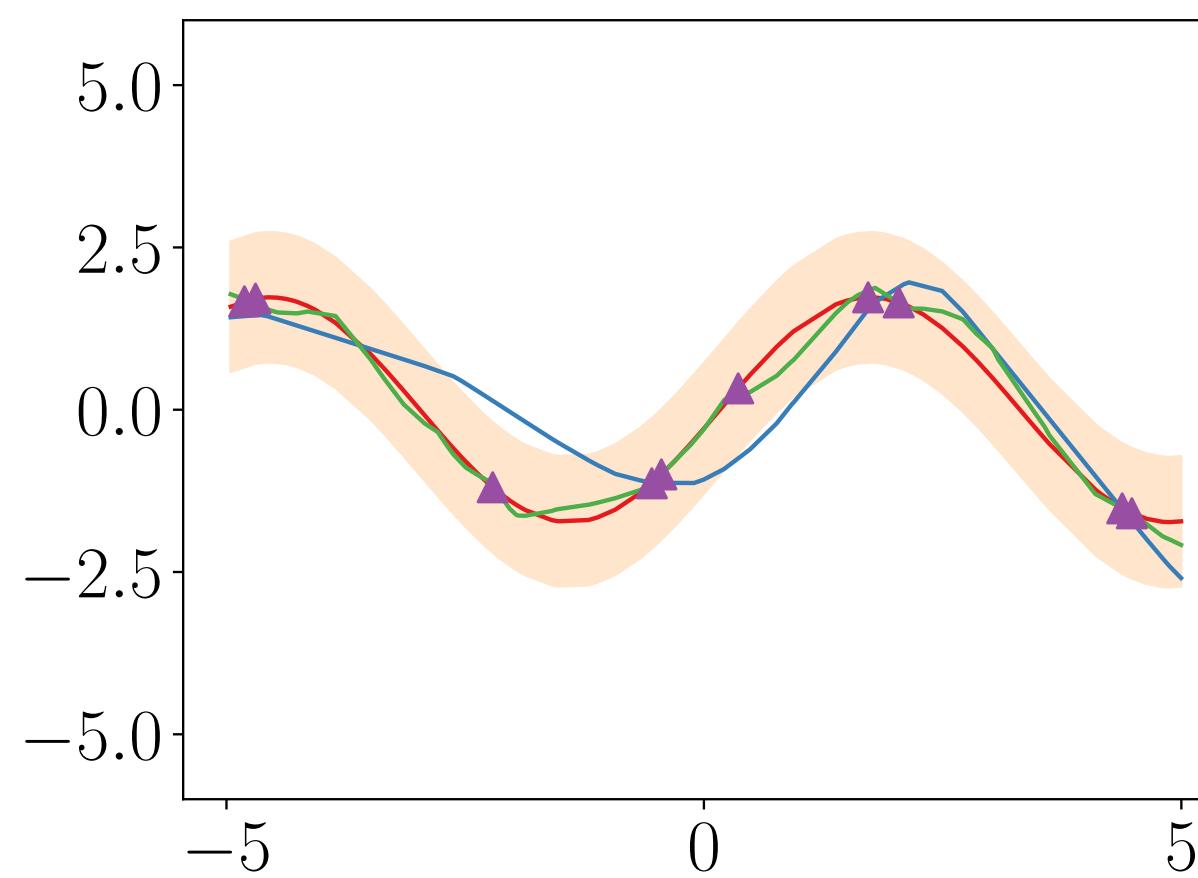
Learned optimizers provably perform well in just 10 steps

K-shot Meta-Learning for Sine Curves



Model-Agnostic Meta-Learning (MAML) [Finn et. al 2017]
MAML learns a shared initialization z so that \hat{z} performs well on test data

Visualizing Guarantees: K-shot Meta-Learning for Sine Curves

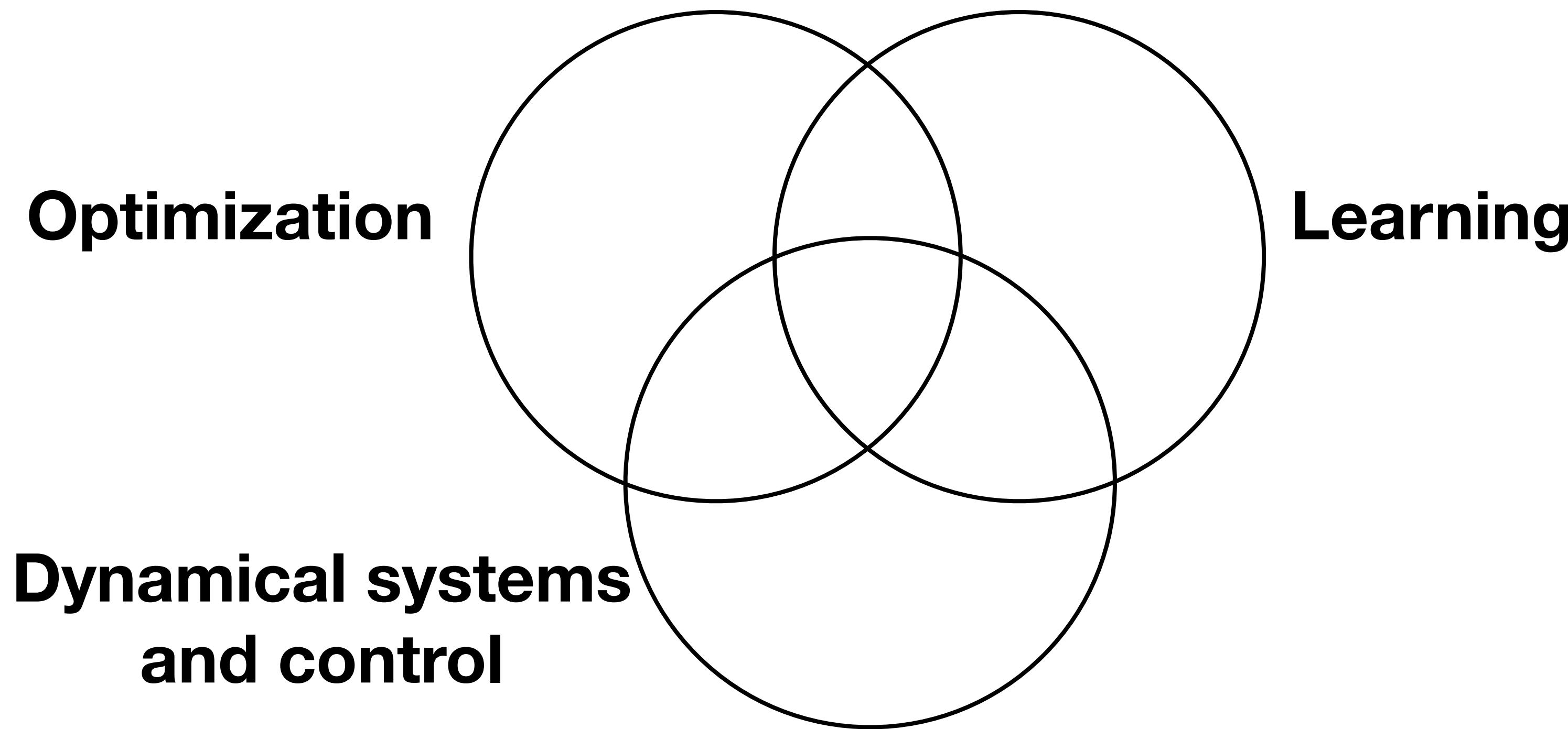


After 10 grad steps
Empirical MAML
Pretrained
Region with MAML guarantee

With high probability, 87% of the time
MAML stays within the band

The pretrained baseline only stays
within the band 33% of the time

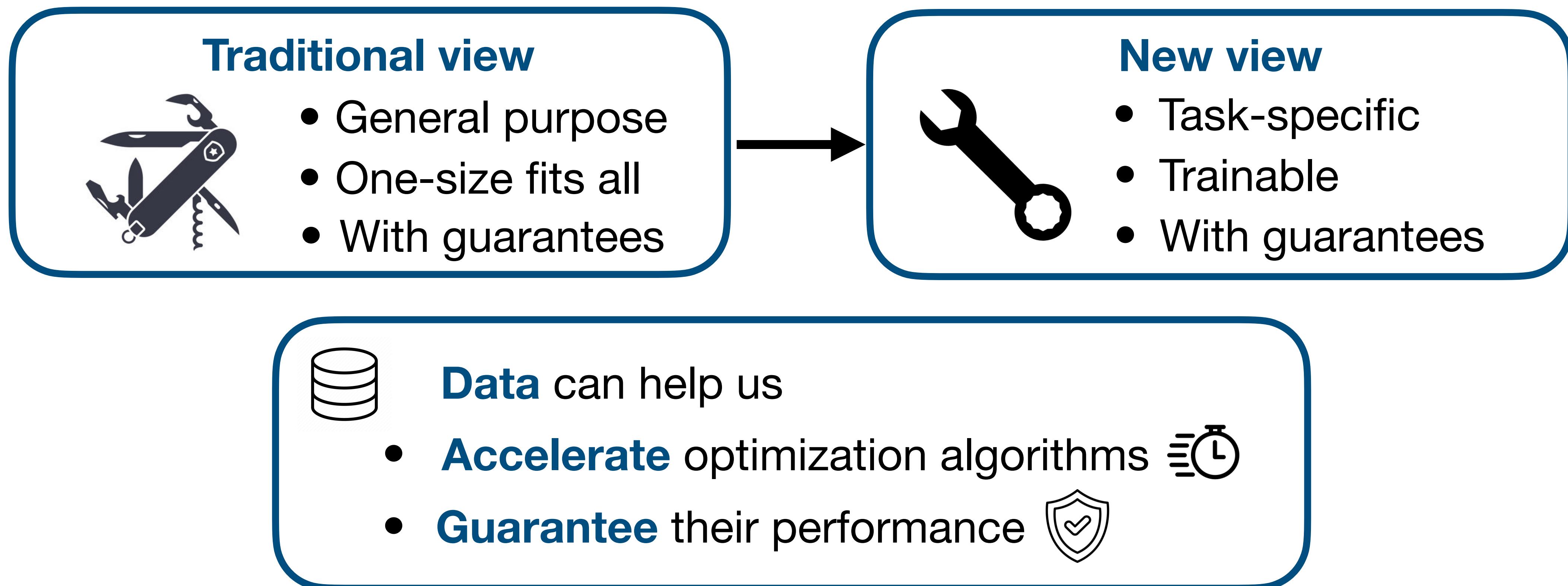
Future directions



Areas of interest

- Learning dynamical systems
- Learning safety and stability certificates
- Guarantees via statistical learning

Conclusions: Learning Optimizers with Guarantees



Data-Driven Performance Guarantees for Classical and Learned Optimizers

R. Sambharya, B. Stellato

Submitted: Journal of Machine Learning Research

<https://arxiv.org/pdf/2404.13831.pdf>



Learning to Warm-Start Fixed-Point Optimization Algorithms

R. Sambharya, G. Hall, B. Amos, B. Stellato

Journal of Machine Learning Research (to appear)

<https://arxiv.org/pdf/2309.07835.pdf>



rajivs@princeton.edu



rajivsambharya.github.io