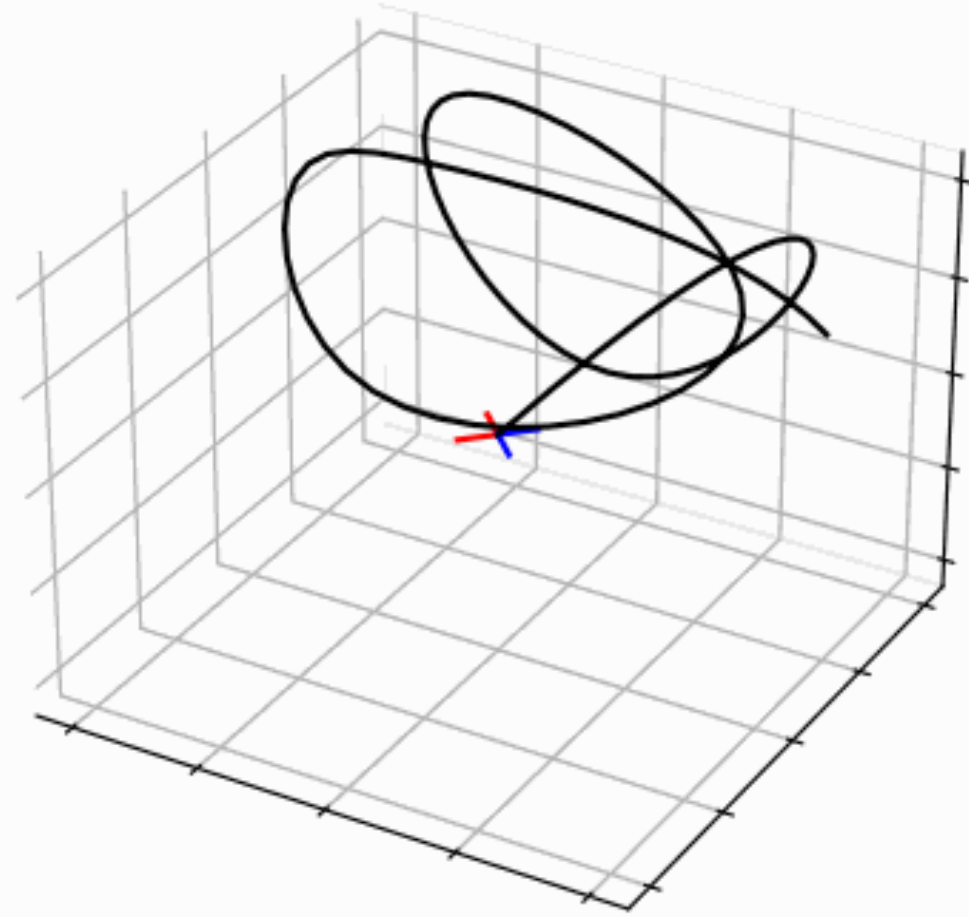


Learning to Warm-Start Fixed-Point Optimization Algorithms

Yale Robotics Seminar 2023
Rajiv Sambharya



Tracking a reference trajectory with a quadcopter



Model predictive control

Model predictive controller

minimize $\sum_{t=1}^T (x_t - x_t^{\text{ref}})^T Q (x_t - x_t^{\text{ref}})$

subject to $x_{t+1} = Ax_t + Bu_t$

$$x_t \in \mathcal{X}, \quad u_t \in \mathcal{U}$$

$$x_0 = x_{\text{init}}$$

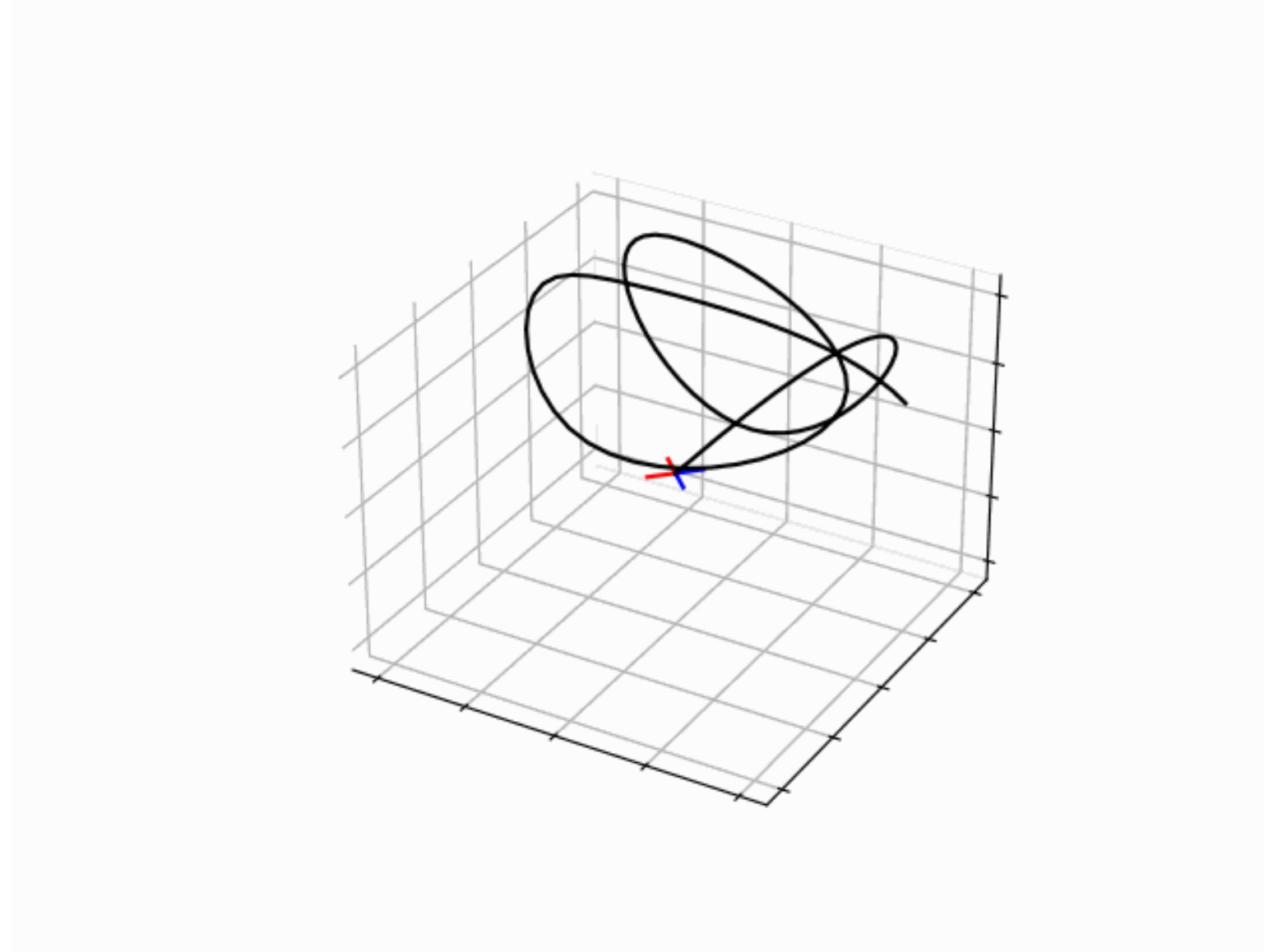
Current state,
reference trajectory



Optimal
controls

Challenge: we need faster methods to solve optimization problems

Robotics and control



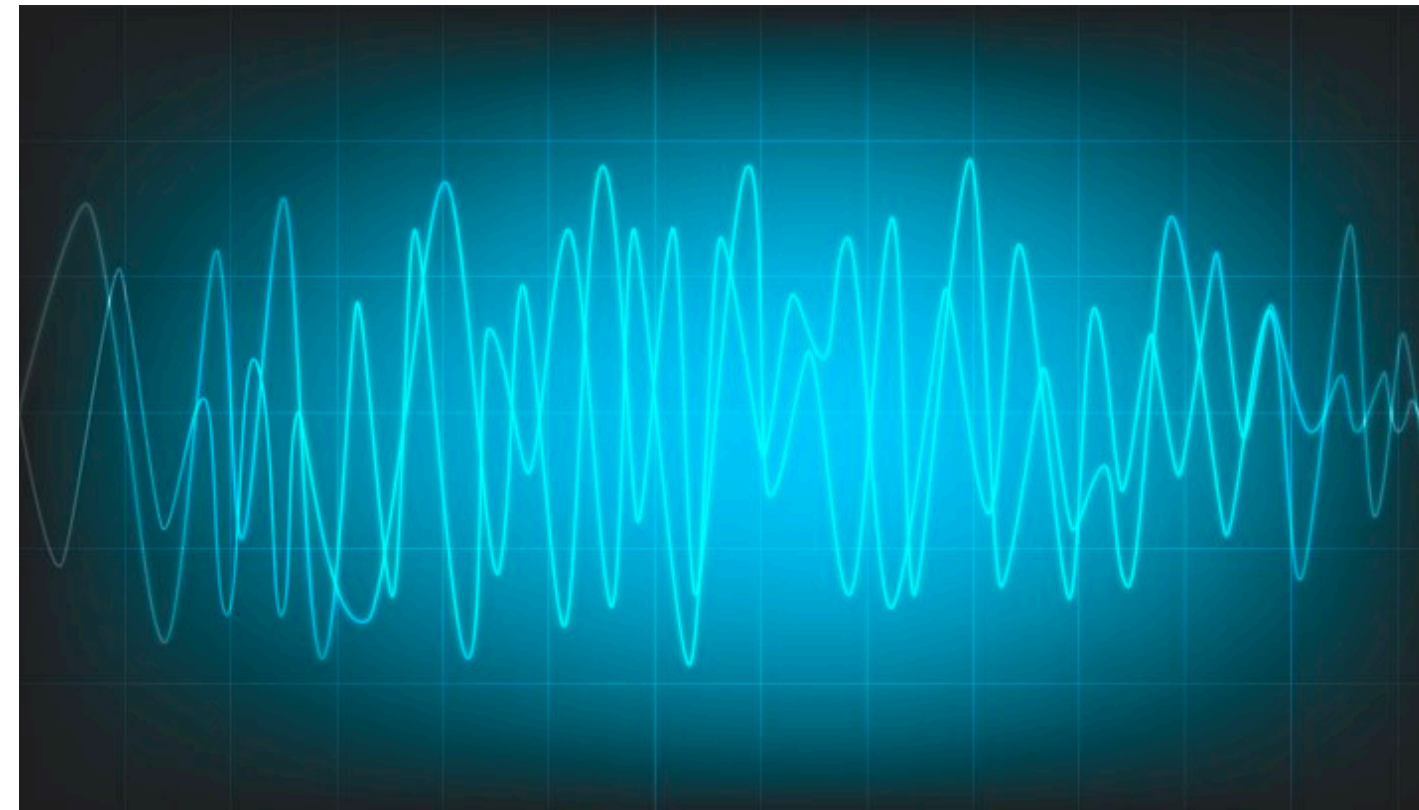
We sometimes need solutions in ~10 milliseconds or less

Claim: Real-world optimization is parametric

Robotics and control



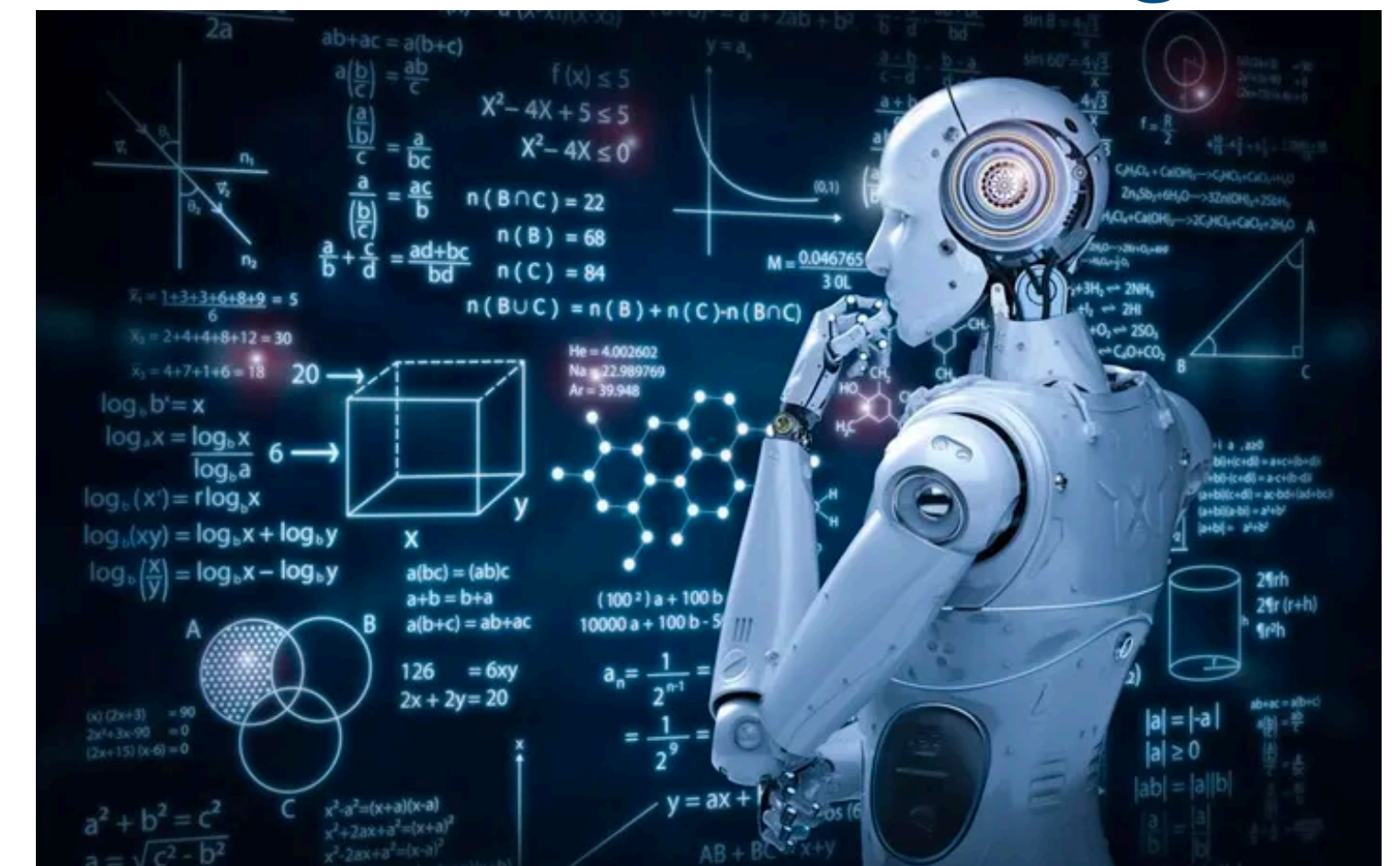
Signal processing



Energy



Machine learning



Can machine learning speed up parametric optimization?

Often, we solve **parametric** optimization problems from the same family

Goal: Do mapping quickly and accurately

Parameter

$\theta \longrightarrow$

minimize $f_\theta(z)$
subject to $g_\theta(z) \leq 0$

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

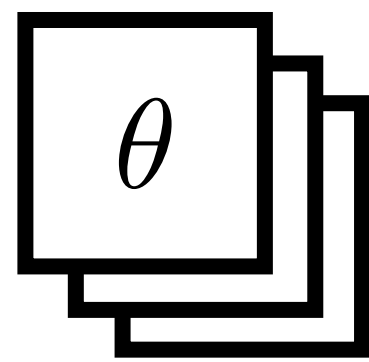
$$\begin{aligned} &\text{minimize} && f_{\theta}(z) \\ &\text{subject to} && g_{\theta}(z) \leq 0 \end{aligned}$$

Offline

Data collection

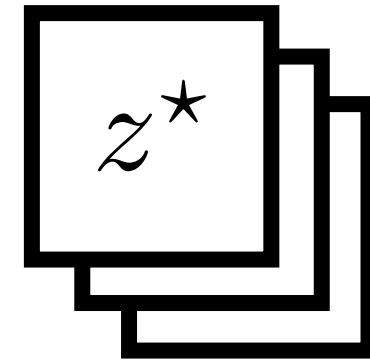


Parameters



Solve

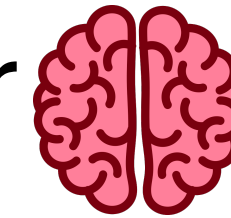
Optimal solutions



Training

Training parameter θ

Learnable Optimizer



Learn

Candidate solution

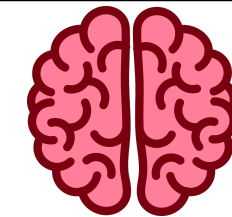
$$\hat{z}(\theta)$$

Loss

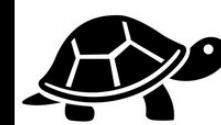
Online evaluation

Unseen parameter θ

Learned Optimizer



Classical Optimizer



High-quality solution

Learning to optimize is a growing research area

Inverse problems

(Gregor and LaCun 2010)

(Liu et. al 2018)

(Wu et. al 2020)

Convex optimization

(Venkataraman and Amos 2021)

(Ichnowski et. al 2021)

(Heaton et. al 2020)

(Jung et. al 2022)

Integer programming

Excellent tutorials

Learning to Optimize: A Primer and a Benchmark (Chen et. al 2021)

Tutorial on Amortized Optimization (Amos 2022)

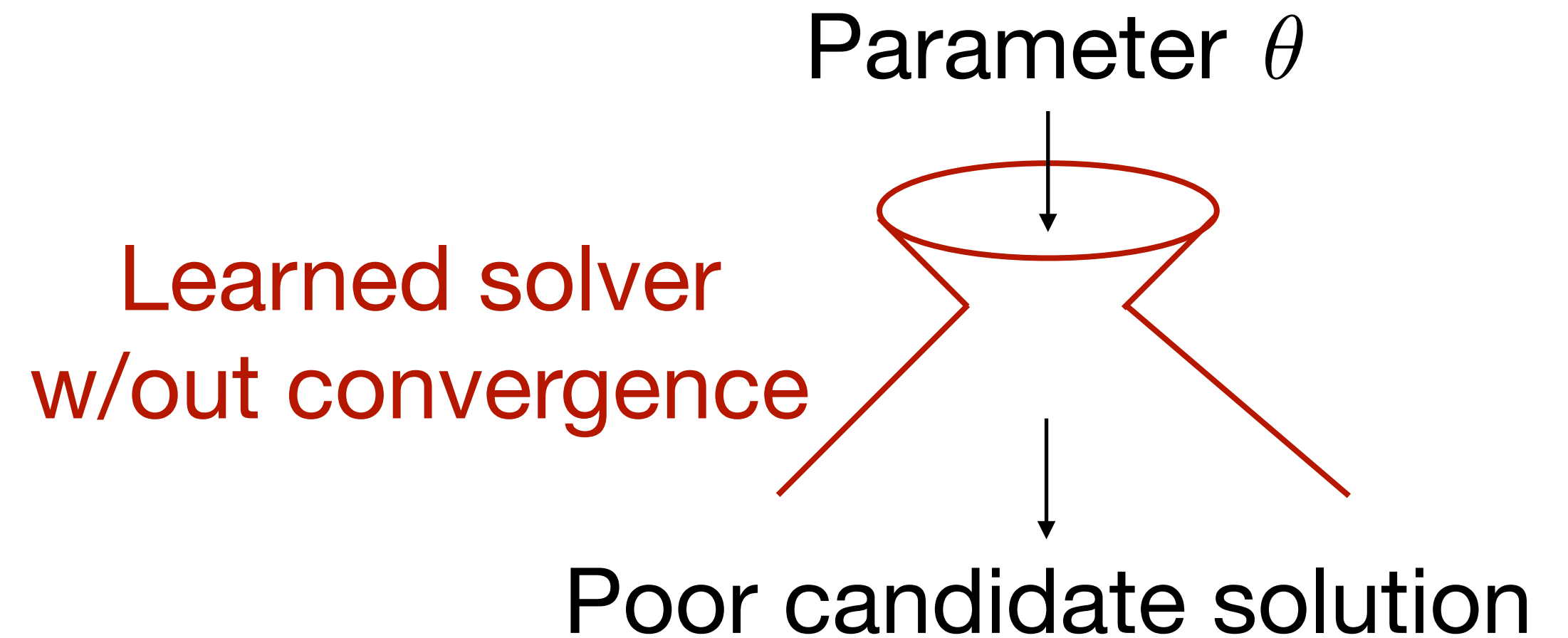
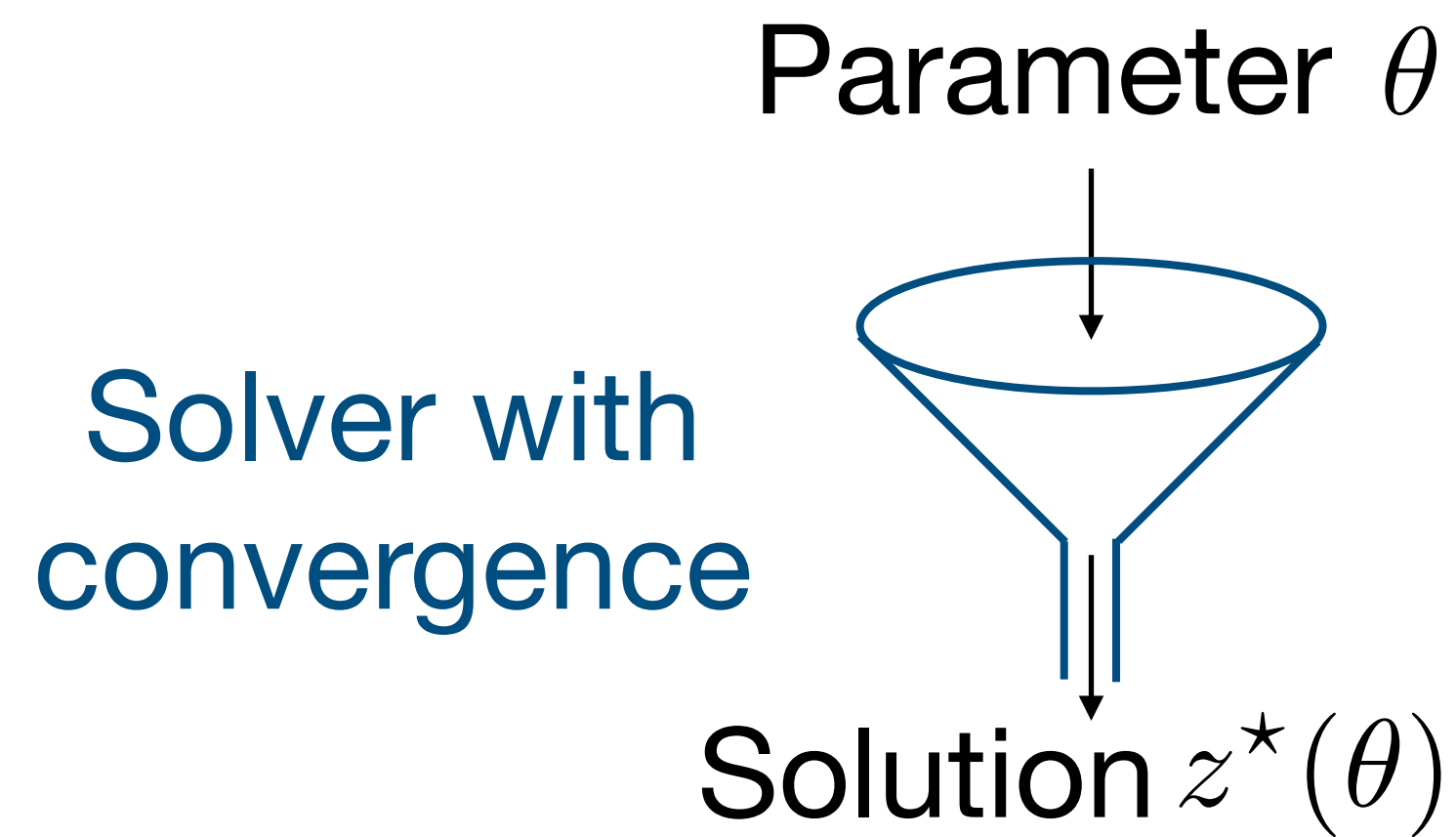
Issues in learning to optimize (L2O) methods

- **I: Lack convergence guarantees**
- **II: Lack generalization guarantees**
- **III: Incompatibility with state-of-the-art solvers**

We need **reliable** L2O methods

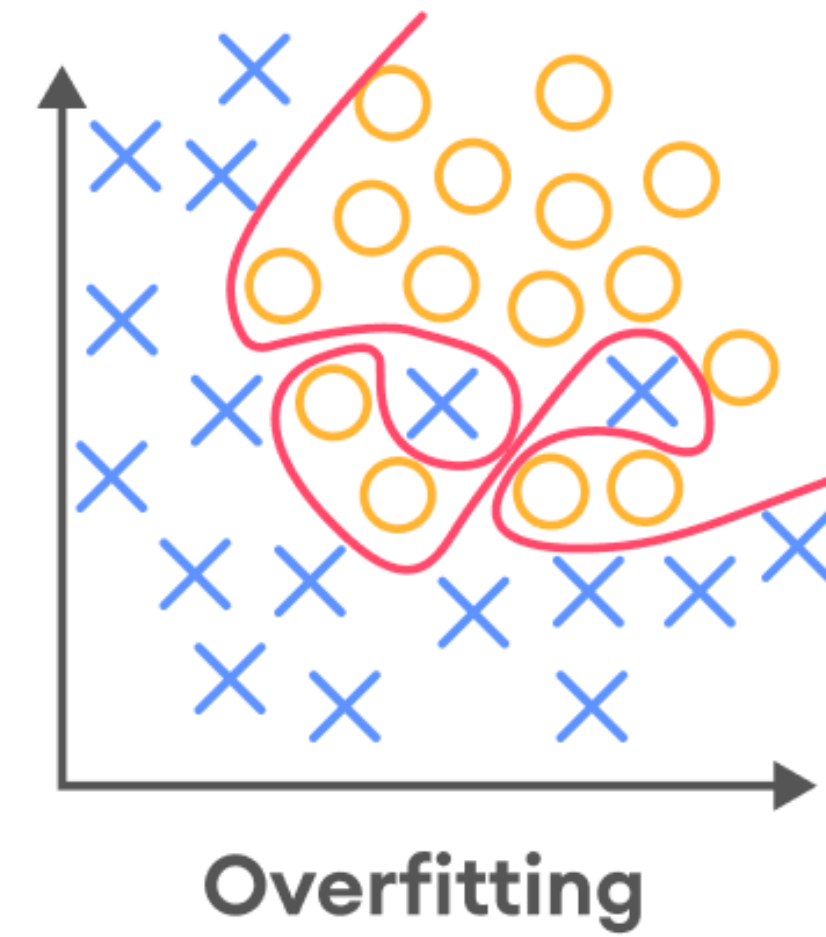
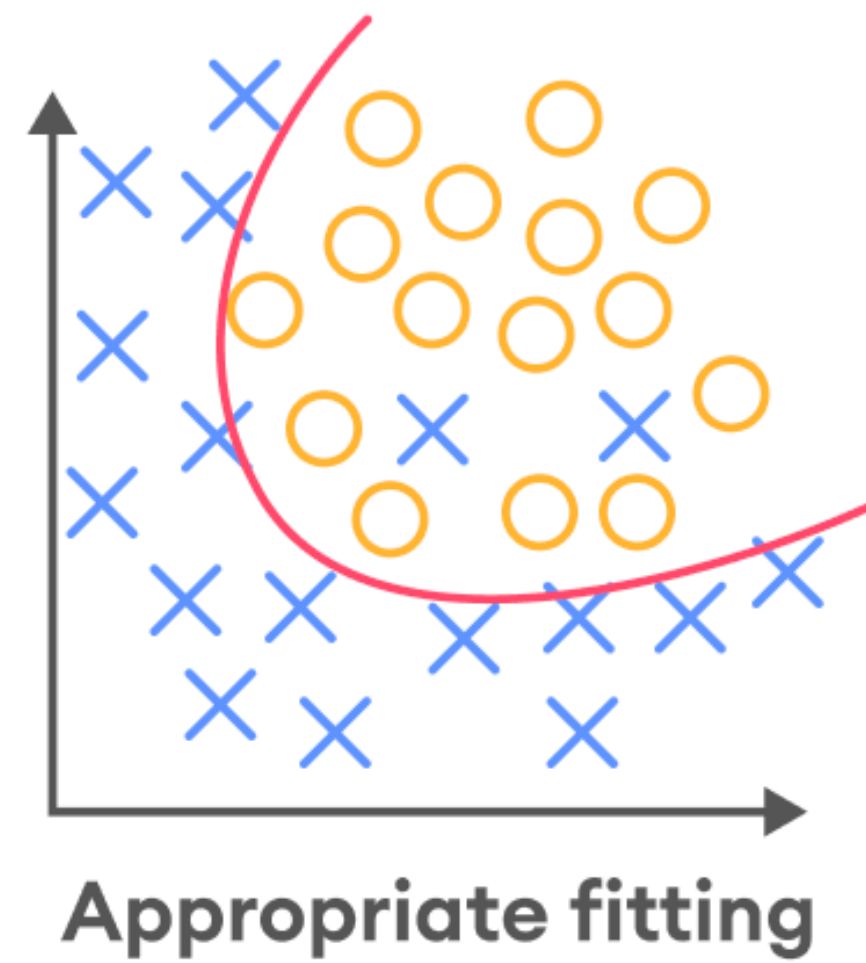


L2O Challenge I: Convergence guarantees



L2O Challenge II: Generalization guarantees

To unseen problems



L2O Challenge III: Incompatibility with state-of-the-art solvers

Existing state-of-the-art solvers are highly optimized

Written in low-level languages



SCS
SPLITTING CONIC SOLVER



Clarabel.jl



When learning the algorithm steps, we cannot use these solvers

Learning to Warm-Start Fixed-Point Optimization Algorithms

Collaborators



Georgina
Hall



Brandon
Amos



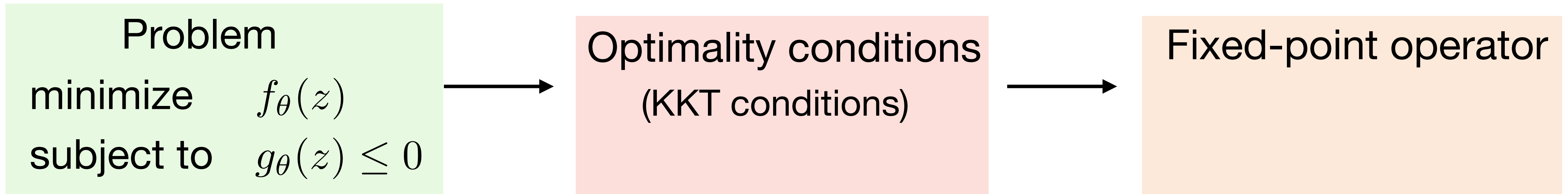
Bartolomeo
Stellato



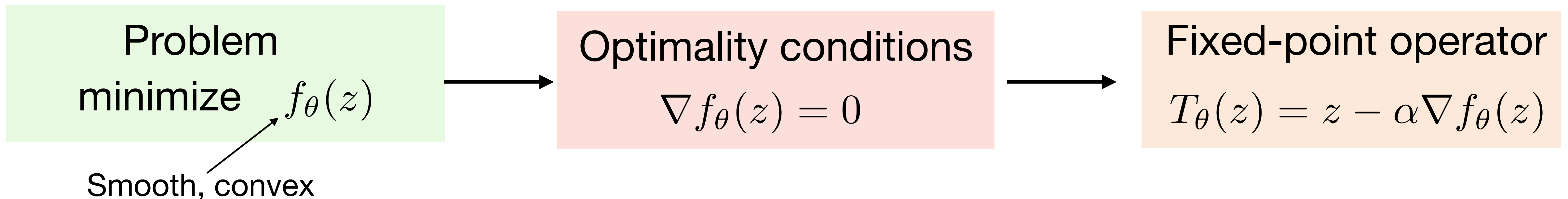
Fixed-point optimization problems are ubiquitous

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

Convex optimization

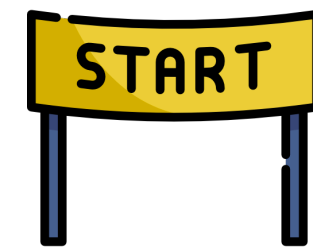


Unconstrained, smooth convex optimization



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 $f_\theta(z^i) = \|T_\theta(z^i) - z^i\|_2$ is small

Fixed point residual

Example: Proximal gradient descent

minimize $g_\theta(z) + h_\theta(z)$

Convex
Smooth

Convex
Non-smooth

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

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

Operator $T_\theta(z) = \text{prox}_{\alpha h_\theta}(z - \alpha \nabla g_\theta(z))$



Problem: limited iteration budget



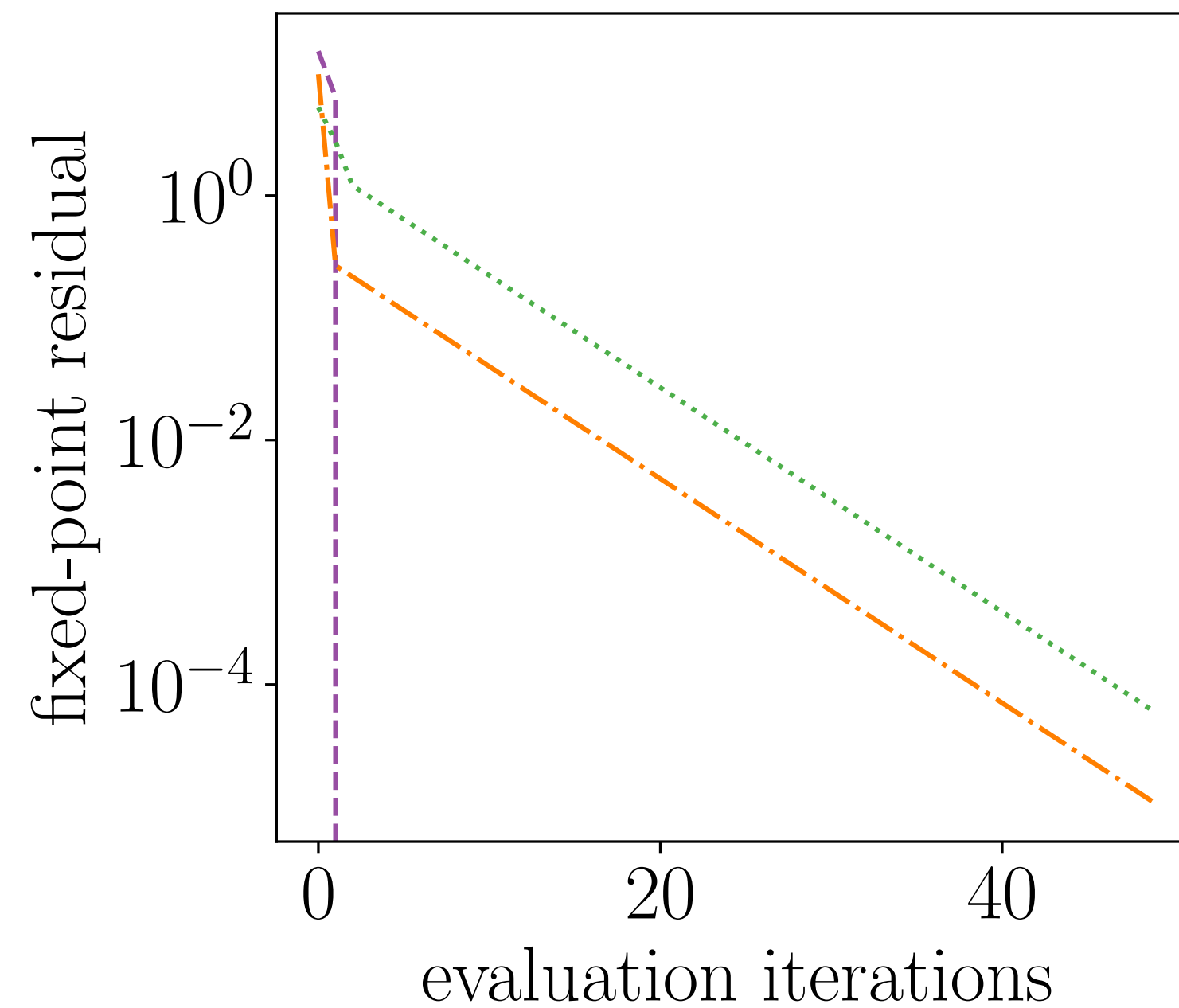
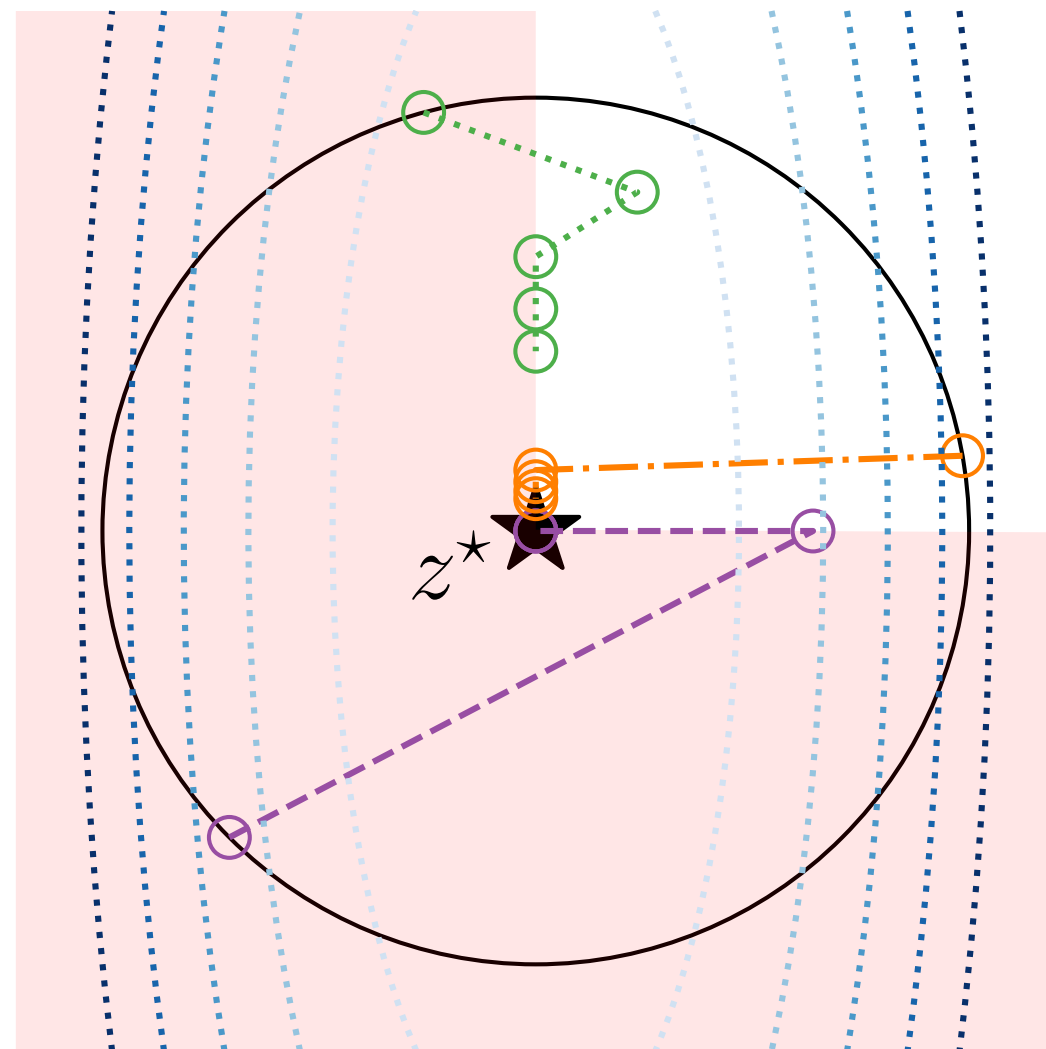
Solution: learn the warm-start to improve the solution within budget

Some warm-starts are better than others

$$\begin{aligned} &\text{minimize} && 10z_1^2 + z_2^2 \\ &\text{subject to} && z \geq 0 \end{aligned}$$

Optimal solution at the origin

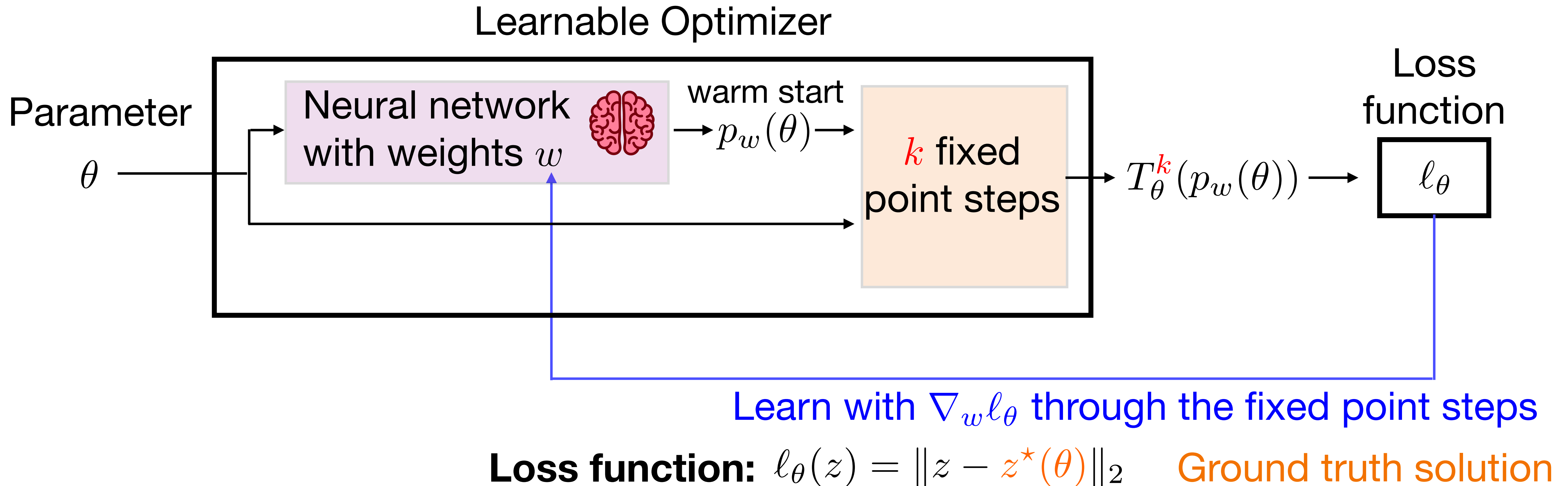
Run proximal gradient descent to solve



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

Learning Framework

End-to-end learning architecture



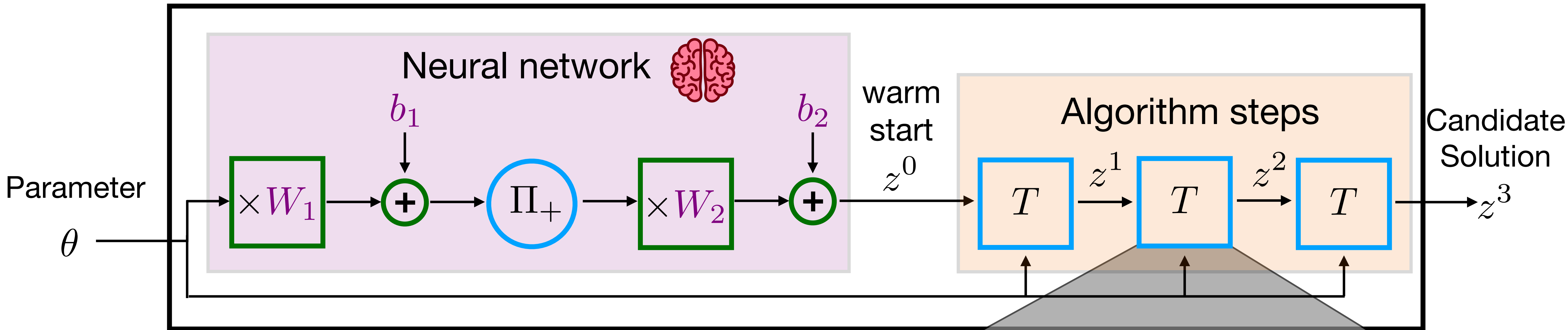
Learned warm start tailored for downstream algorithm

An example architecture

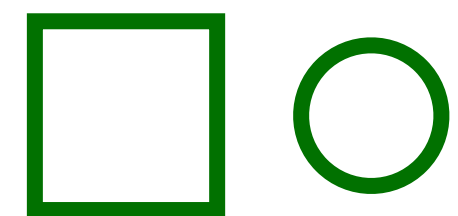
minimize $1/2 z^T P z + \theta^T z$
 subject to $z \geq 0$

Fixed-point operator: $T_\theta(z) = \Pi_+ \left((I - \alpha P)z - \alpha \theta \right)$

Projection onto \mathbb{R}_+ Step size



Computational Graph

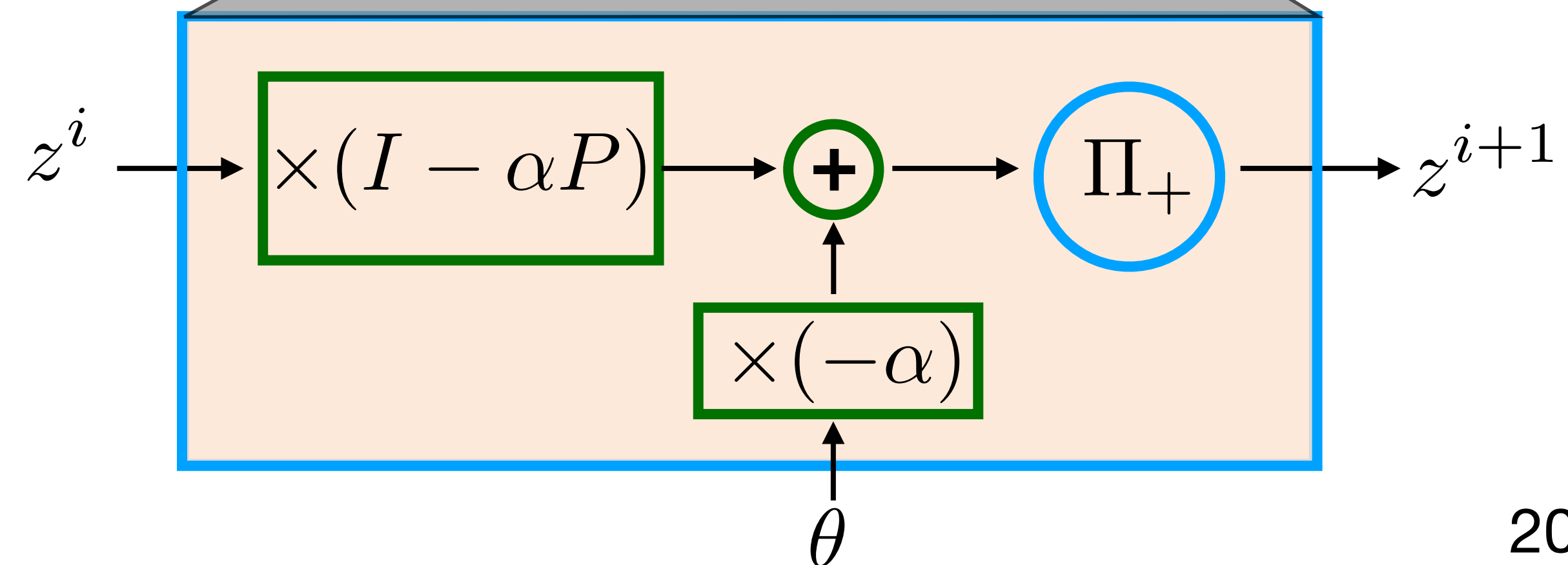


Linear components



Non-linear components

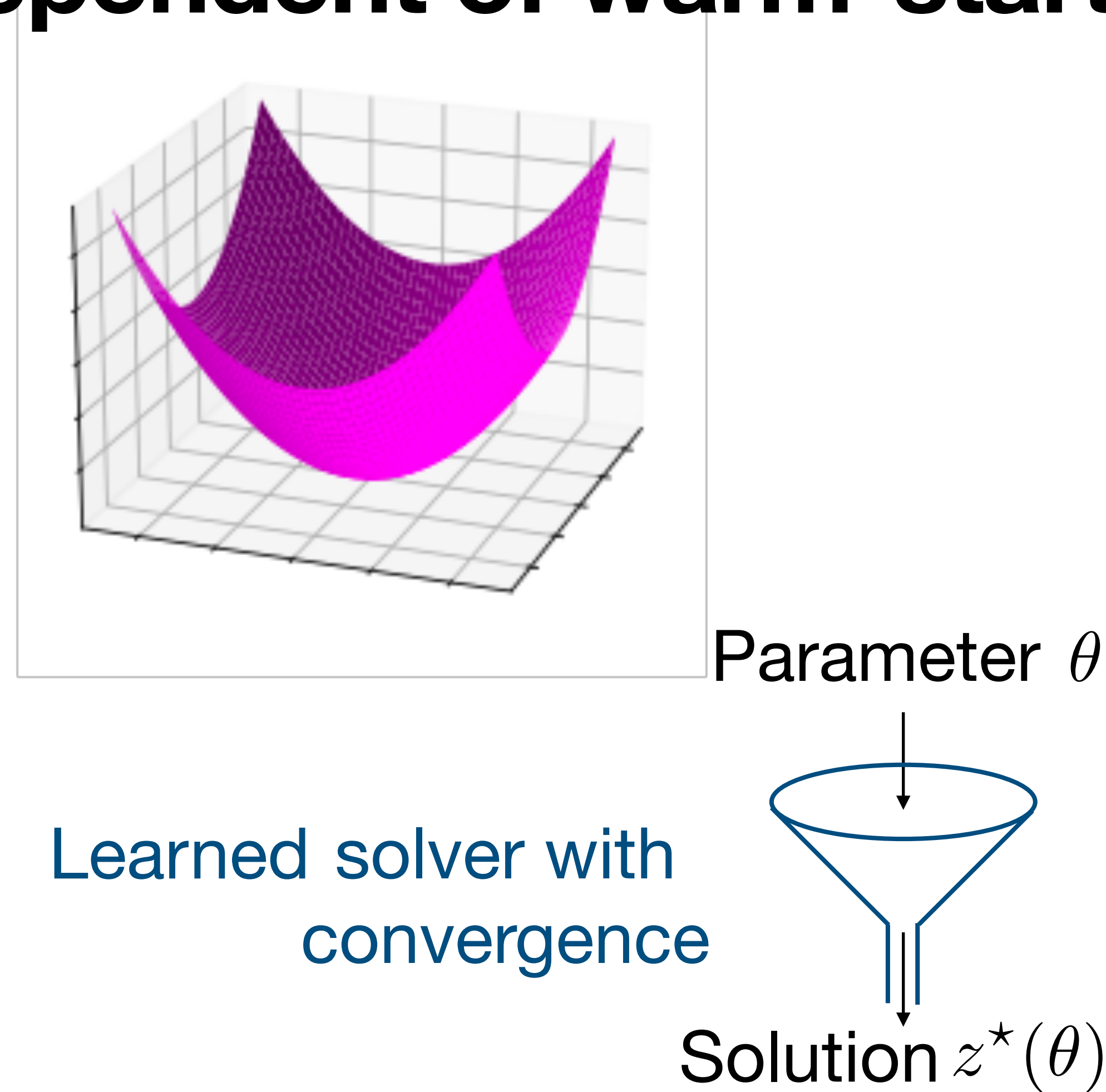
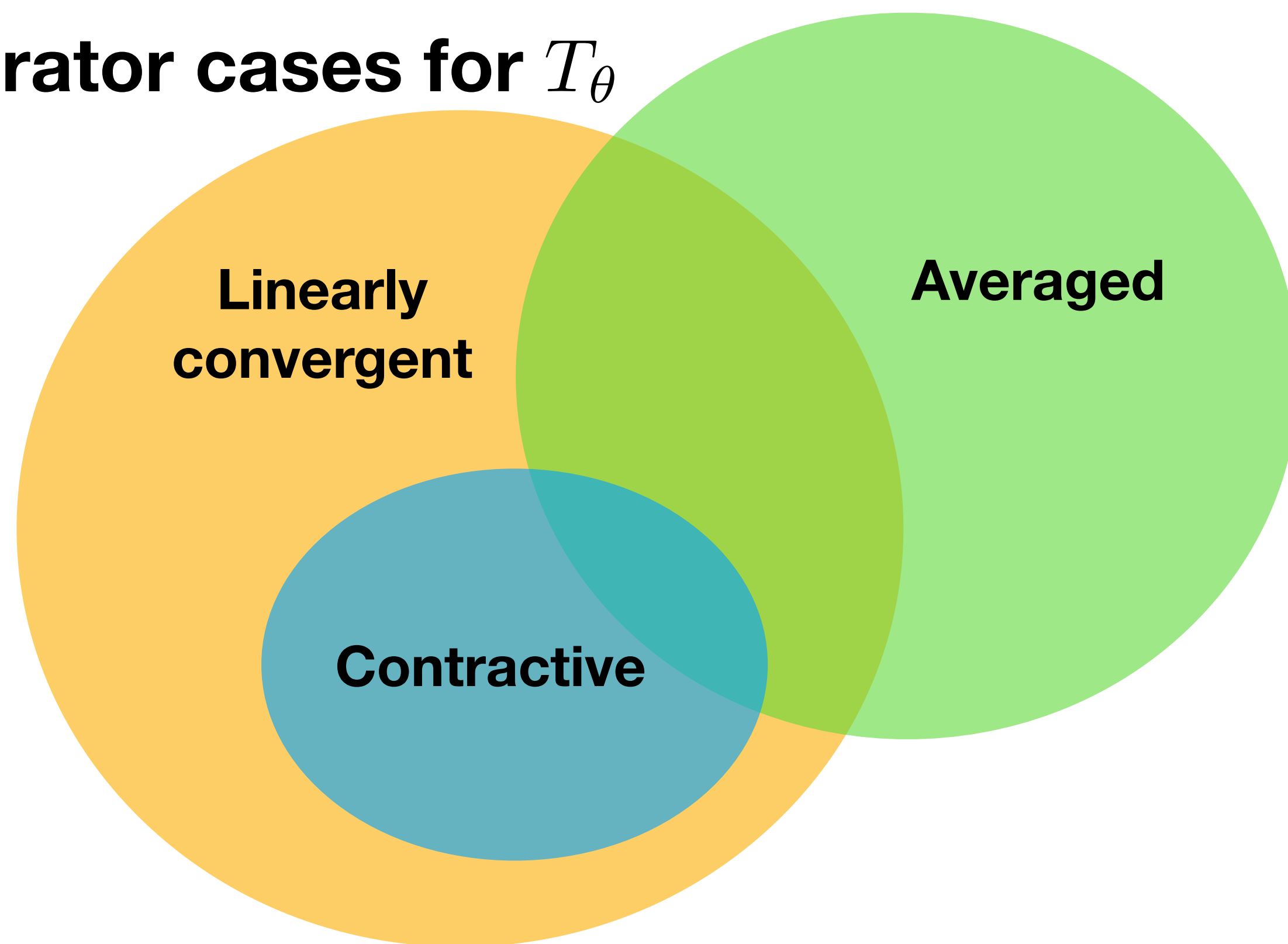
W_1, b_1, W_2, b_2 Learnable components



Convergence and Generalization Bounds

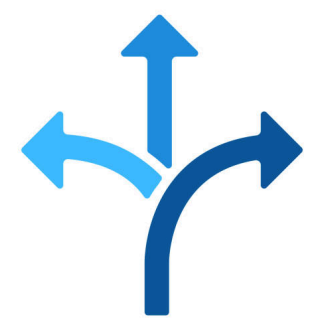
Guaranteed convergence independent of warm-start

Operator cases for T_θ

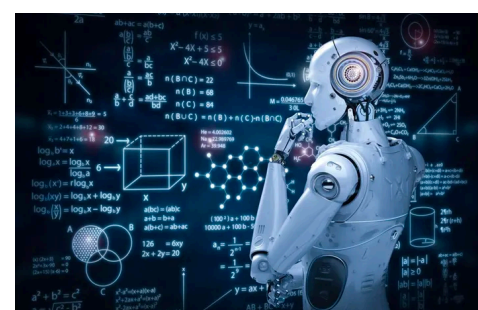
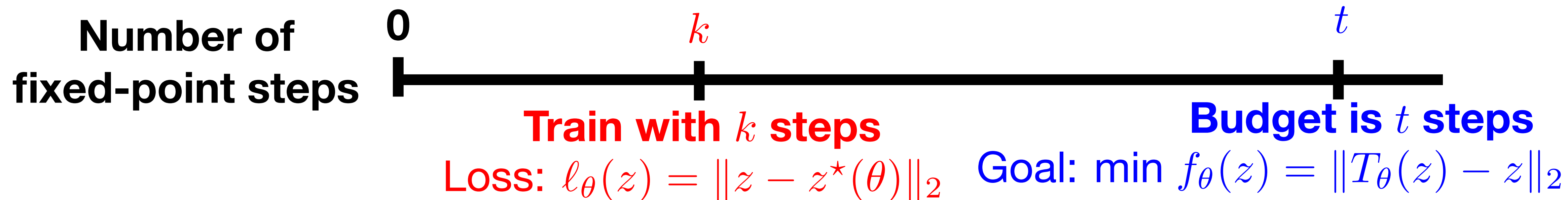


Convergence always guaranteed for learned warm starts

Generalization bounds: train for k , evaluate for t



Flexibility: # of evaluation steps can differ from # of train steps



Guarantees from k training steps to t evaluation steps

$$\beta\text{-contractive case } f_{\theta}(T_{\theta}^t(z)) \leq 2\beta^{t-k} l_{\theta}(T_{\theta}^k(z))$$

Generalization bounds to unseen data

β -contractive case

Theorem 1. *With high probability over a training set of size N , for any γ ,*

$$\mathbf{E} f_{\theta}(T_{\theta}^t(p_w(\theta))) \leq \underbrace{\frac{1}{N} \sum_{i=1}^N f_{\theta_i}(T_{\theta_i}^t(p_w(\theta_i)))}_{\text{Empirical risk}} + \underbrace{2\beta^t \gamma + \mathcal{O}\left(c_1(t) \sqrt{\frac{c_2(w) + \log(\frac{LN}{\delta})}{\gamma^2 N}}\right)}_{\text{Penalty term}}$$

$c_1(t)$: worst-case fixed-point residual after t steps

As $N \rightarrow \infty$, the **penalty term** decreases

As $t \rightarrow \infty$, the **penalty term** goes to zero

Derived from the PAC-Bayes framework

Non-contractive case: we provide similar bounds

Learned warm-start can easily interface with solvers

Written in 



Quadratic programs

$$\begin{aligned} \text{minimize} \quad & (1/2)x^T P x + c^T x \\ \text{subject to} \quad & \ell \leq A x \leq u \end{aligned}$$

```
sol = osqp_solver.warm_start(x=x0, y=y0)
```



Conic programs

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

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

We code exact replicas of OSQP and SCS

Allows us to make timing comparisons for QPs and conic programs

Numerical Experiments

We evaluate the gain over a cold-start


Baseline initializations

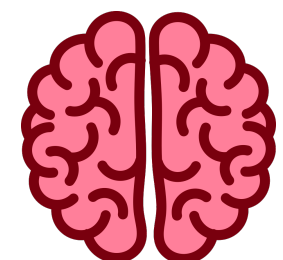
1. Cold-start: initialize at zero 
2. Nearest neighbor: initialize with solution of nearest training problem 

Metrics plotted

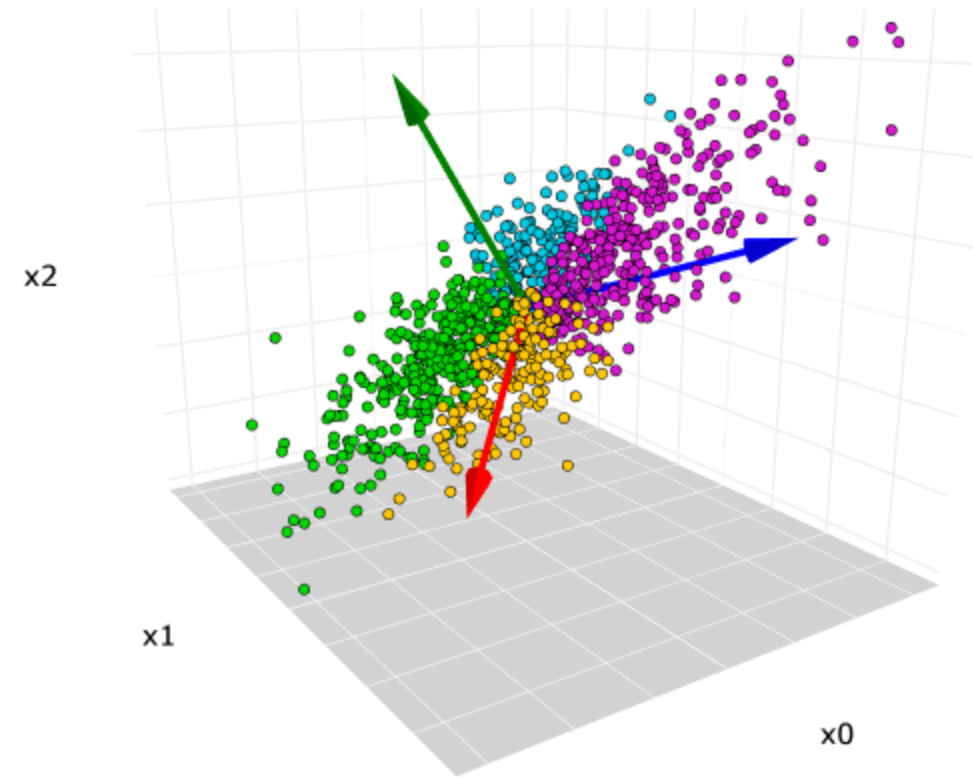
1. Fixed-point residual
2. Gain over the cold-start

$$\text{gain} = \frac{f_{\theta}(T_{\theta}^t(0))}{f_{\theta}(T_{\theta}^t(p_w(\theta)))}$$

Cold-start 

Learned warm-start 

Sparse PCA



Non-convex problem

maximize $x^T A x$

subject to $\|x\|_2 \leq 1$

Card $(x) \leq c$

Covariance matrix



Sparse PCA



Sparse principal component

$\theta = \text{vec}(A)$



Semidefinite program

maximize $\text{Tr}(A X)$

subject to $\text{Tr}(X) = 1$

$\mathbf{1}^T |X| \mathbf{1} \leq c$

$X \succeq 0$

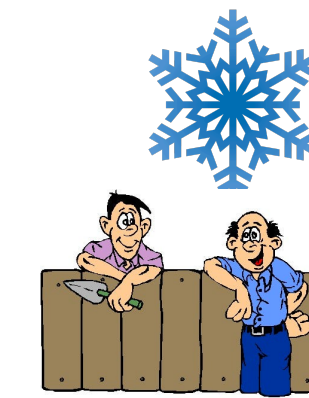


X^*

Sparse PCA results

Different initializations

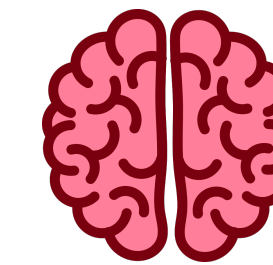
Baselines



▼ Cold-start

◄ Nearest neighbor

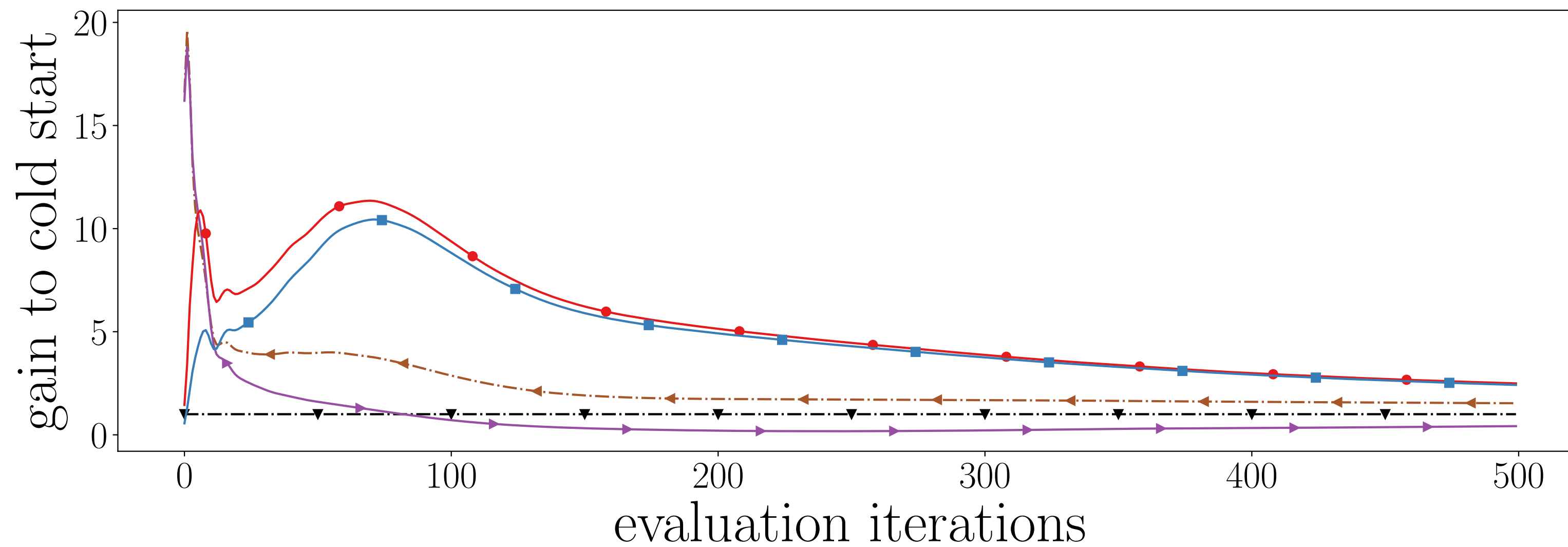
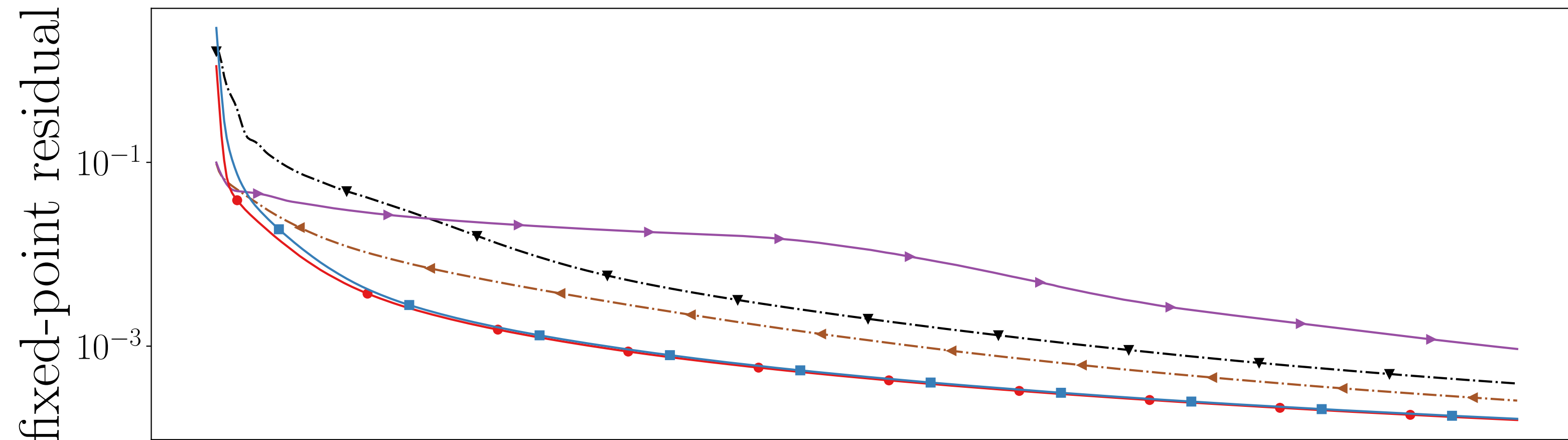
Learned



► $k = 0$

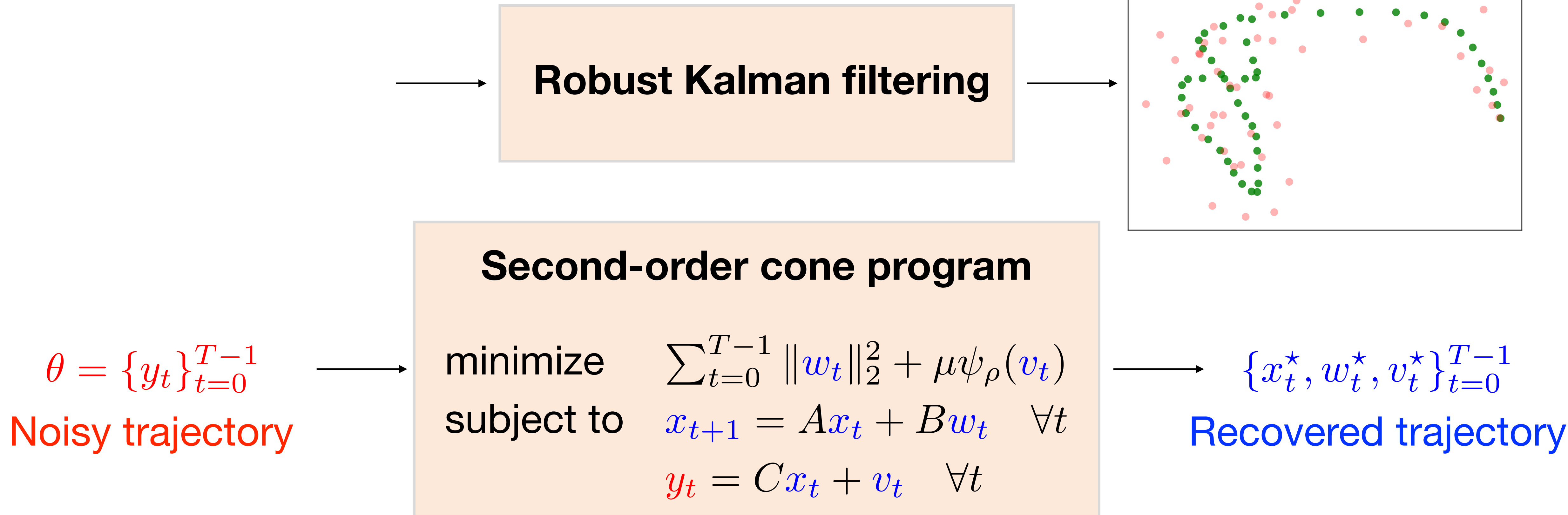
● $k = 5$

■ $k = 15$



Picking $k > 0$ is essential to improve convergence

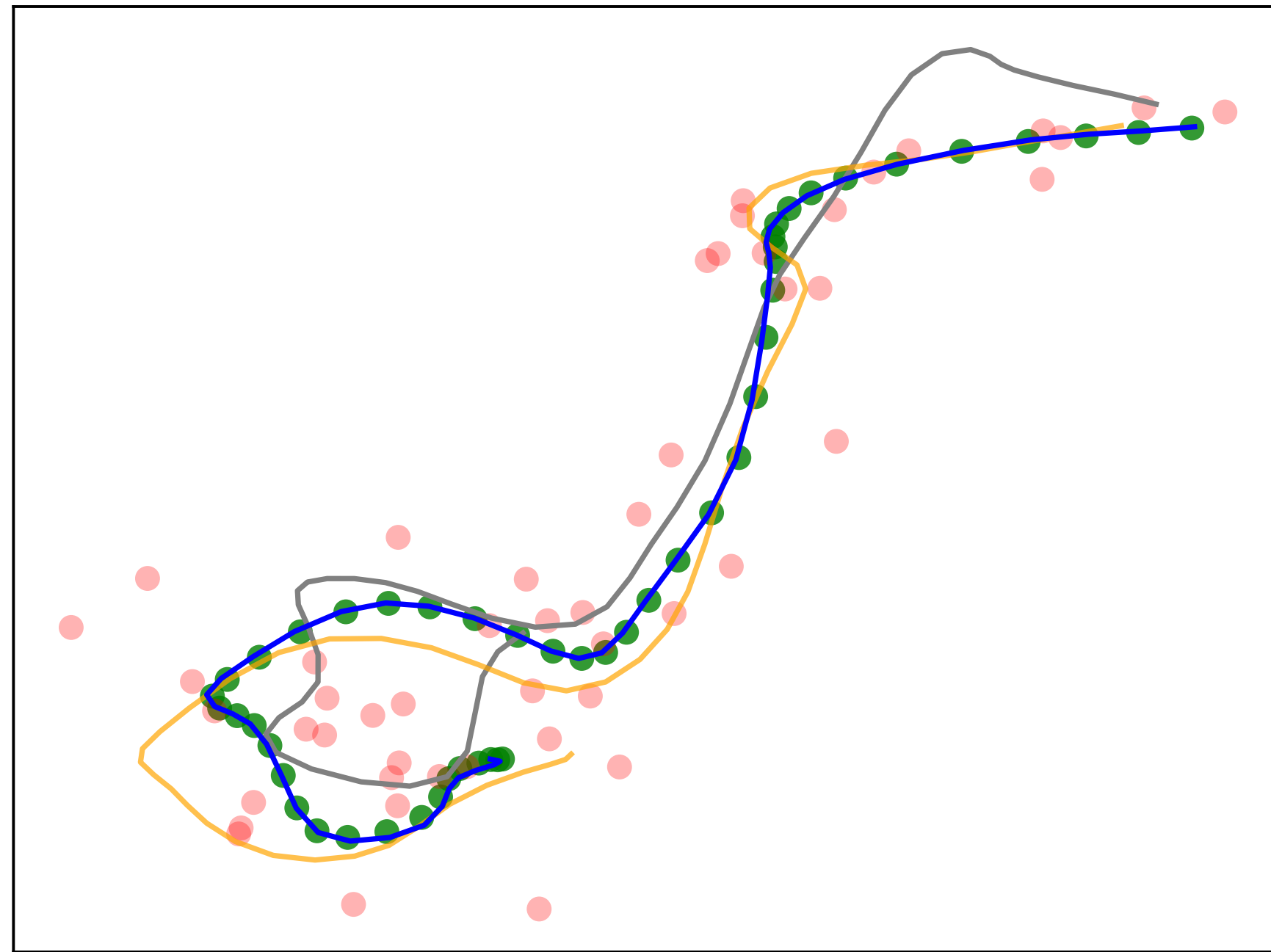
Robust Kalman filtering



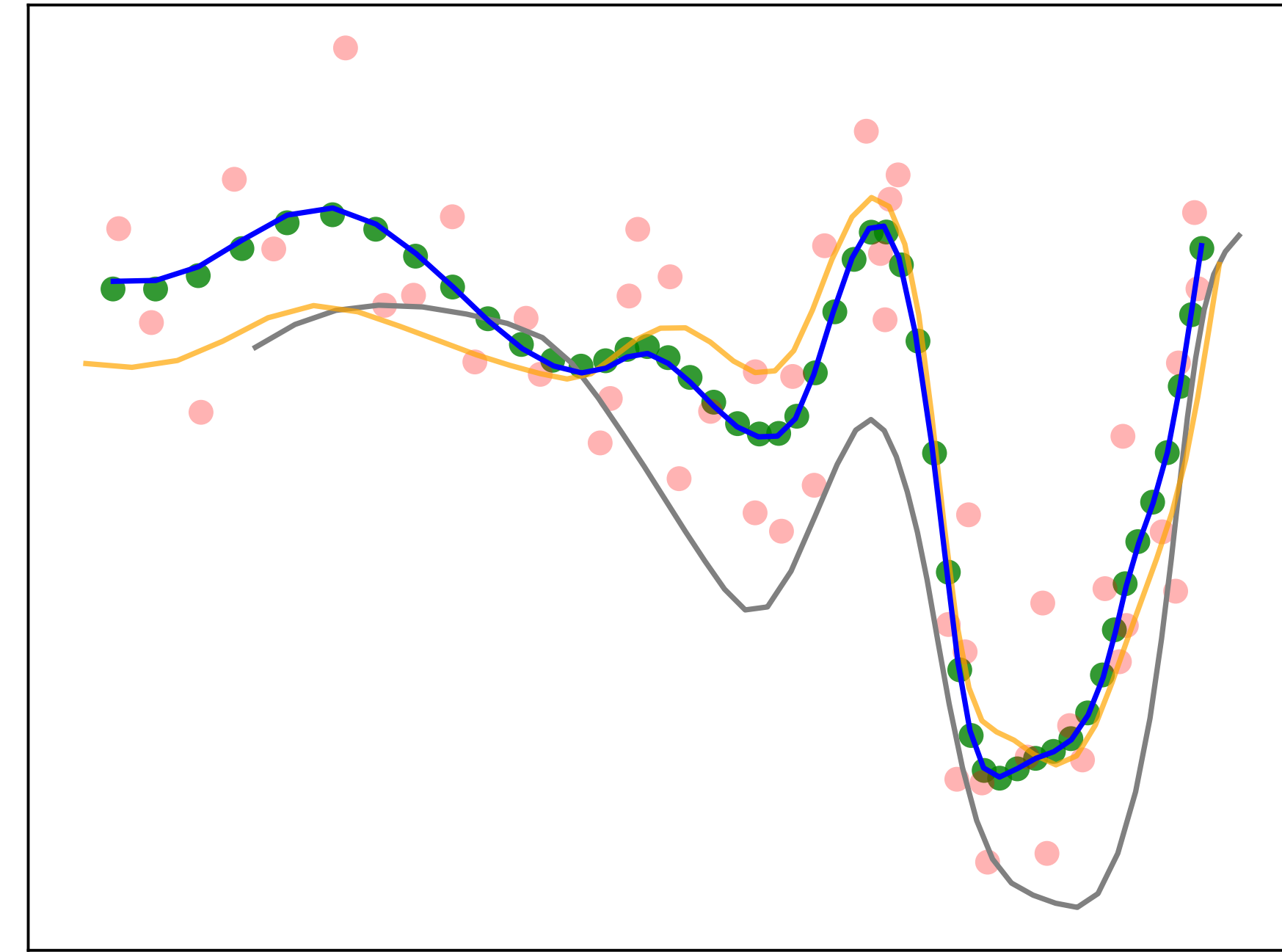
Dynamics matrices: A, B

Observation matrix: C






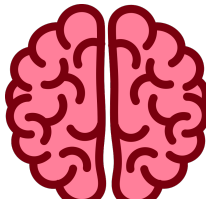
Robust Kalman filtering visuals



-  Noisy trajectory
-  Optimal solution



Solution after 5 fixed-point steps with different initializations

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

With learning, we can estimate the state well

Model predictive control (MPC) of a quadcopter

Current state,
previous control
reference trajectory



Optimal
controls

$$\theta = (x_{\text{init}}, u_{\text{prev}}, \{x_t^{\text{ref}}\}_{t=1}^T)$$



Quadratic program

minimize $\sum_{t=1}^T (x_t - x_t^{\text{ref}})^T Q (x_t - x_t^{\text{ref}}) + \sum_{t=0}^{T-1} u_t^T R u_t$

subject to $x_{t+1} = A x_t + B u_t$

$u_{\min} \leq u_t \leq u_{\max}$

$x_{\min} \leq x_t \leq x_{\max}$

$|u_{t+1} - u_t| \leq \Delta u$

$x_0 = x_{\text{init}}$

$u_{-1} = u_{\text{prev}}$

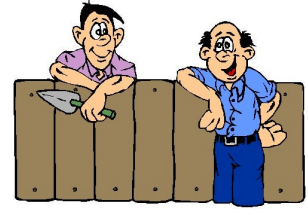


$$\{x_t^*, u_t^*\}_{t=0}^T$$

Linearized dynamics

MPC of a quadcopter in a closed loop

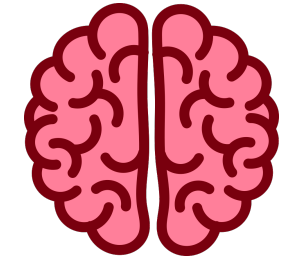
Budget of 15 fixed-point steps



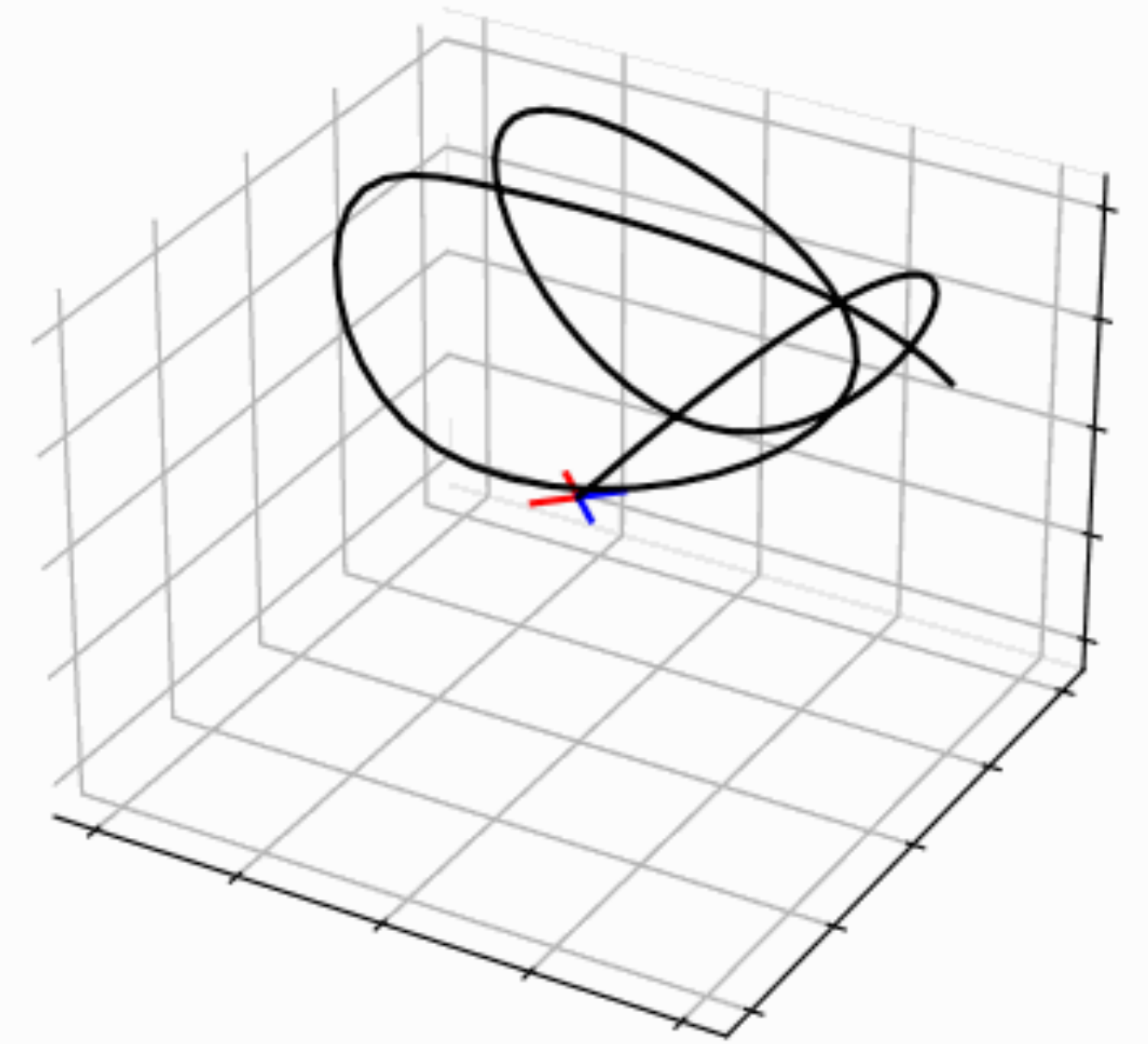
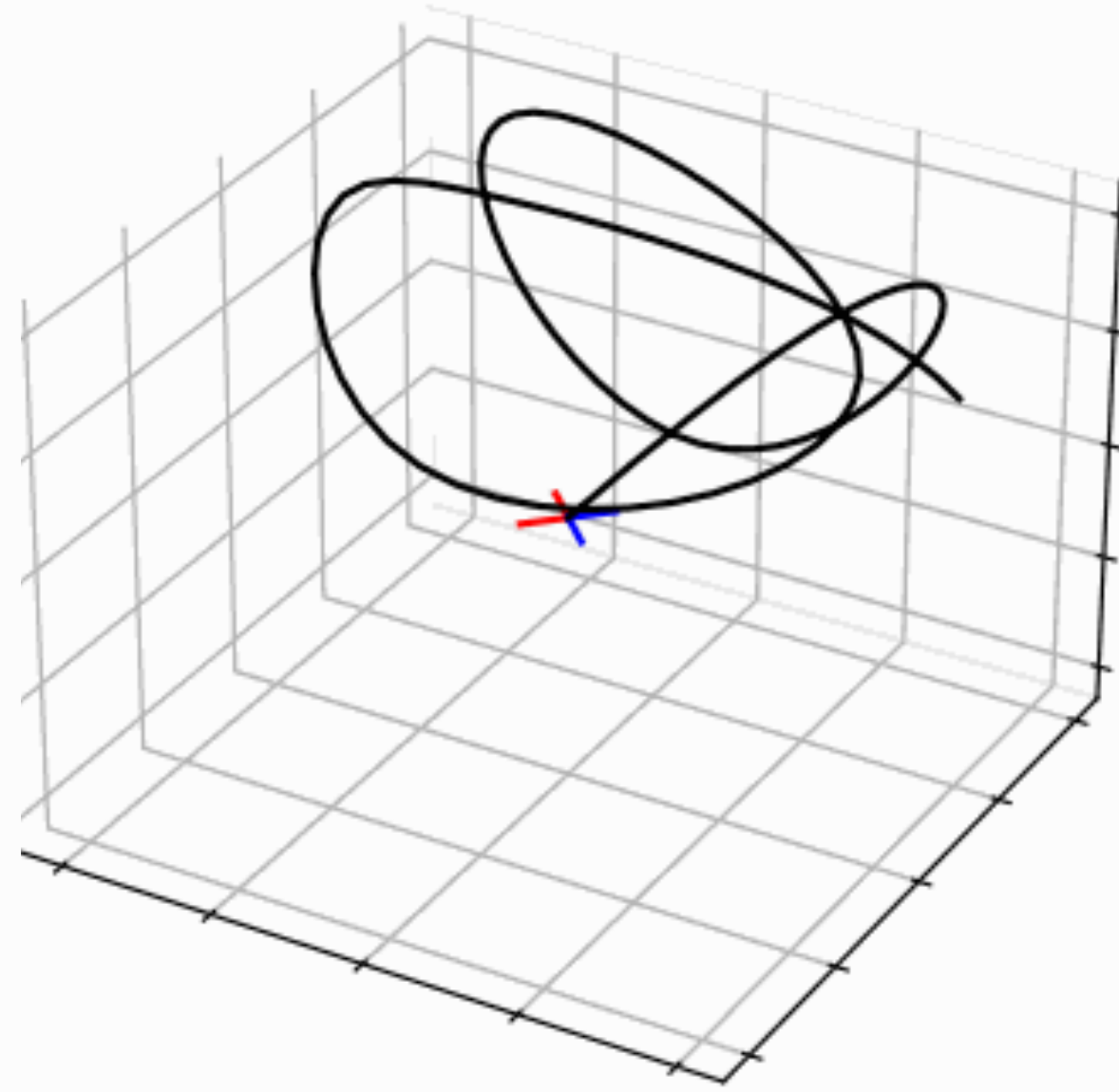
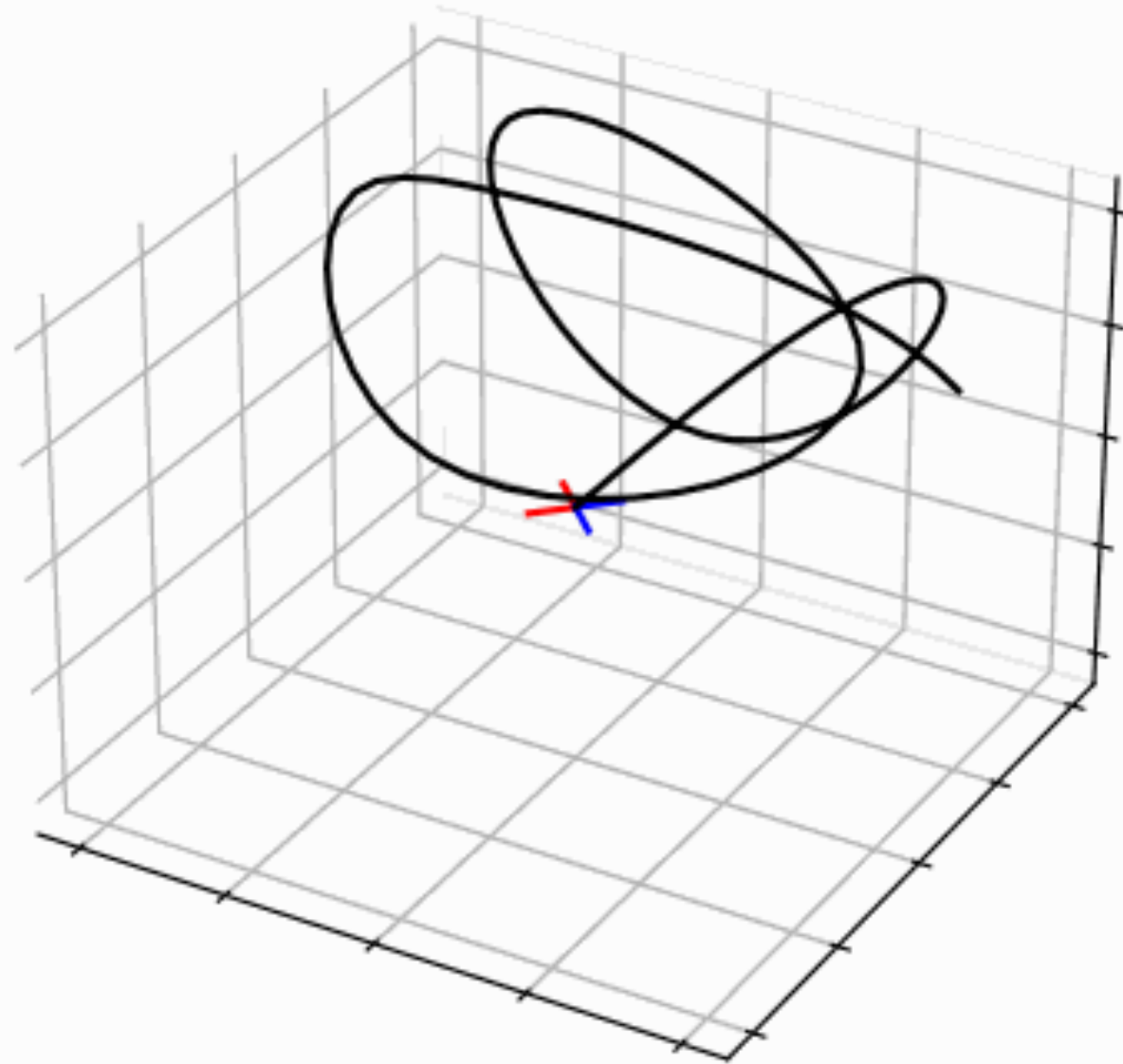
Nearest neighbor



Previous solution



Learned: $k = 5$



With learning, we can track the trajectory well

Image deblurring

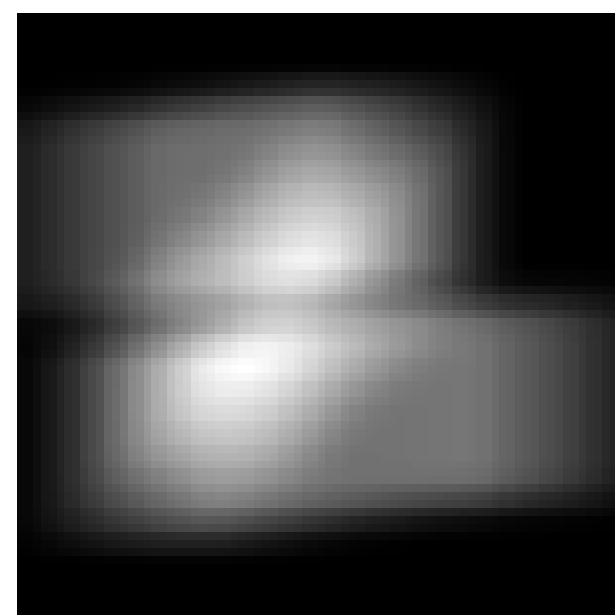
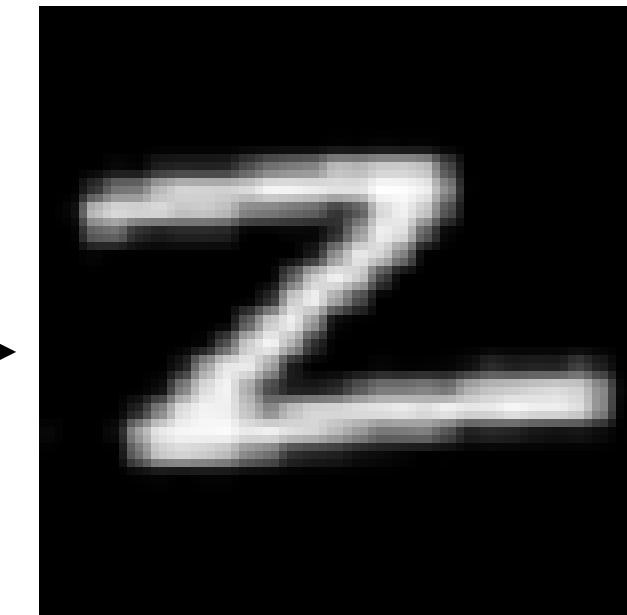


Image deblurring



$\theta = b$
Blurred image

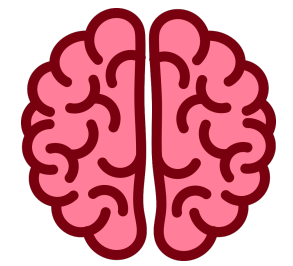
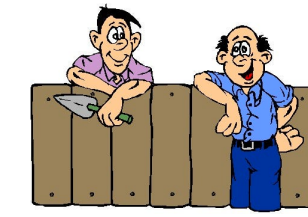


Quadratic program
minimize $\|Ax - b\|_2^2 + \lambda \|x\|_1$
subject to $0 \leq x \leq 1$



x^*
Deblurred image

Image deblurring



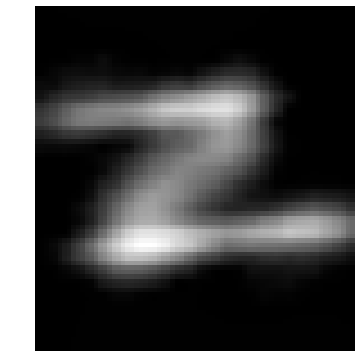
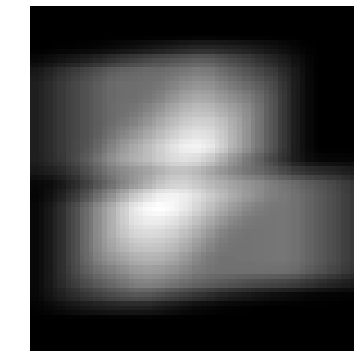
percentile optimal blurred cold-start nearest neighbor learned

50 fixed-point steps

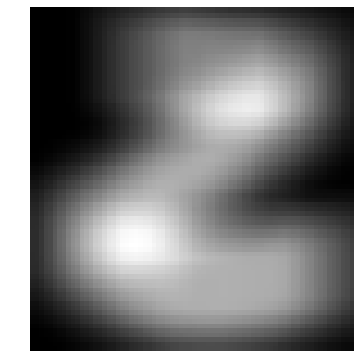
Distance to nearest neighbor increases



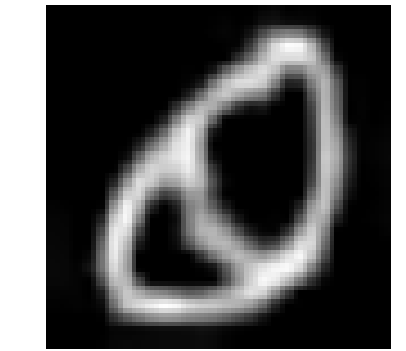
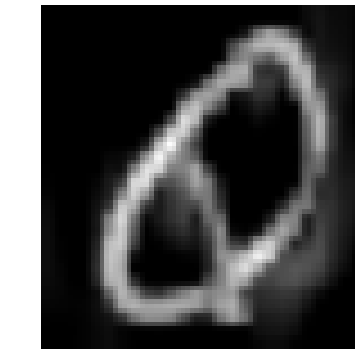
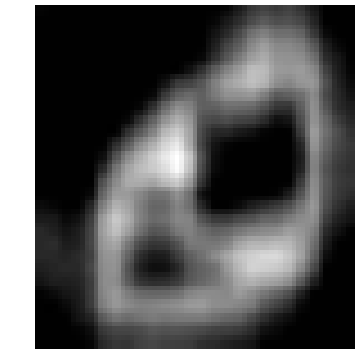
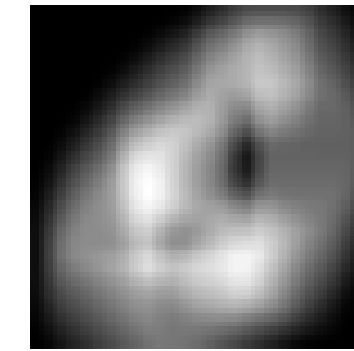
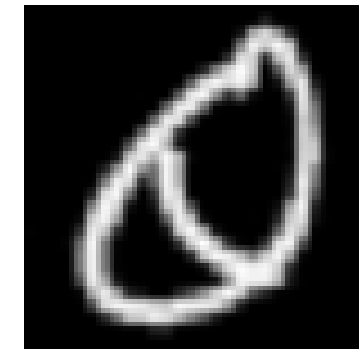
10th



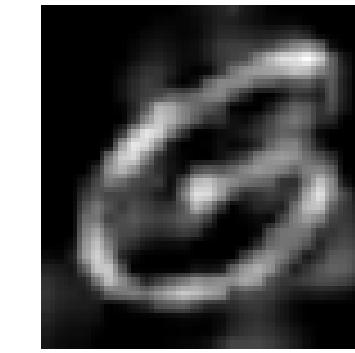
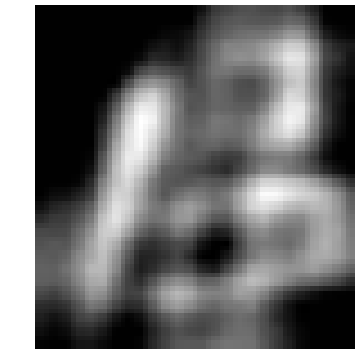
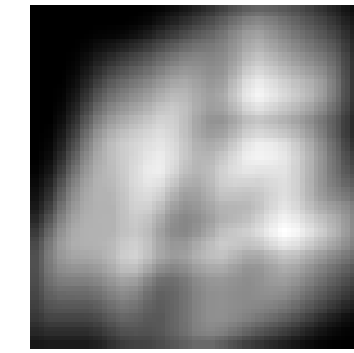
50th



90th



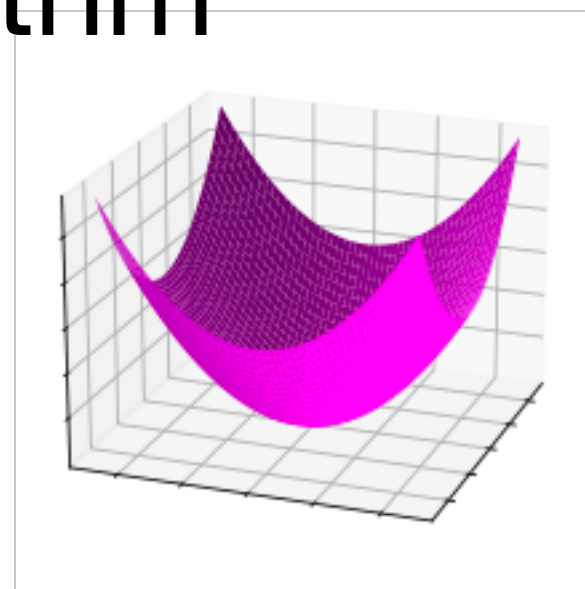
99th



With learning, we can deblur all of the images quickly

Benefits of our learning framework

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



Guaranteed convergence

Can interface with state-of-the-art solvers

Generalization to

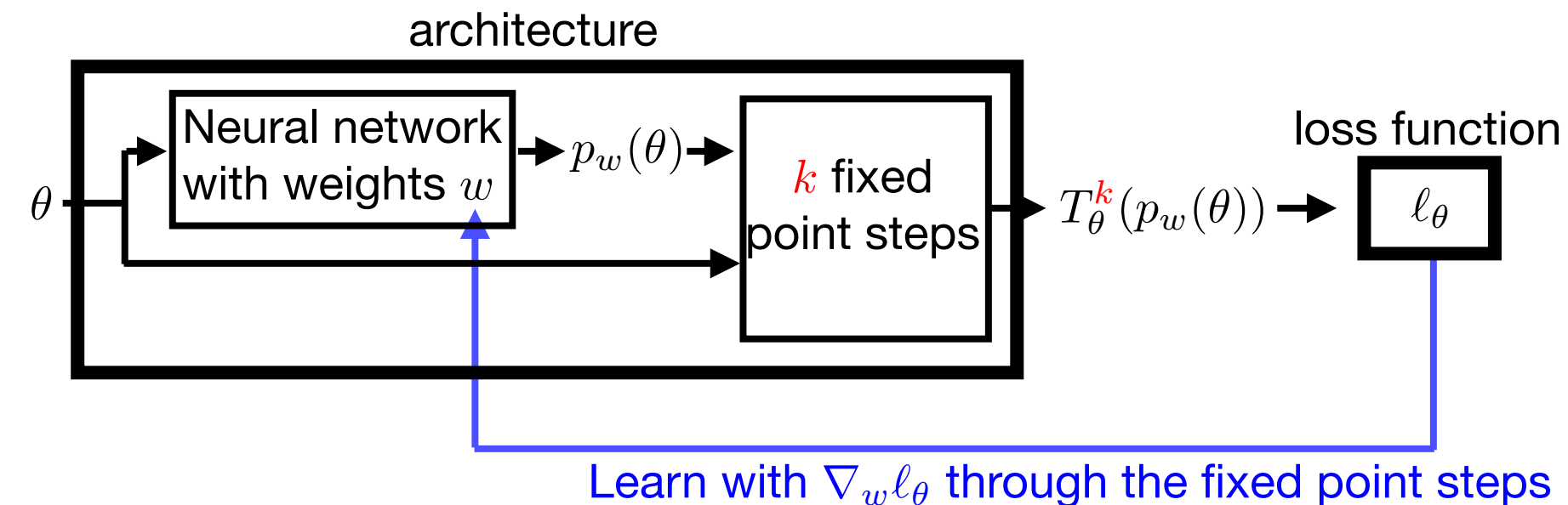
Future iterations
Unseen data



rajivs@princeton.edu



[rajivsambharya.github.io](https://github.com/rajivsambharya)



Quadratic programs



Conic programs



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

