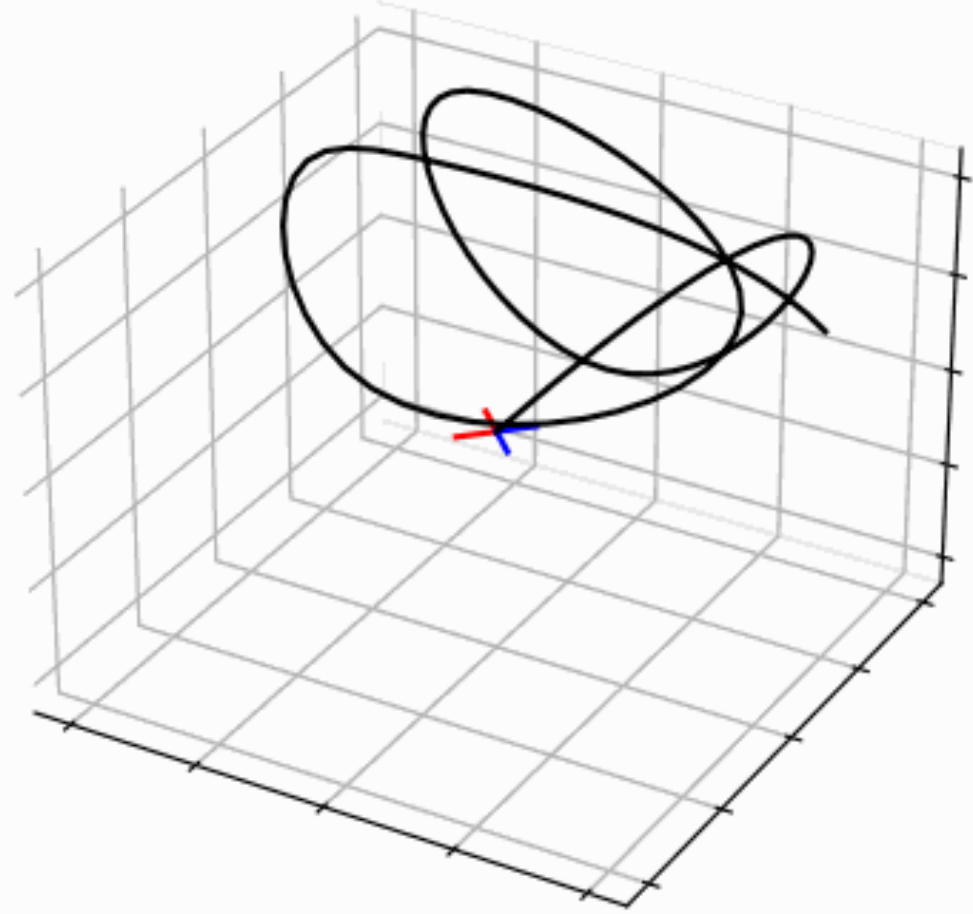# Learning to Accelerate Optimizers with Guarantees

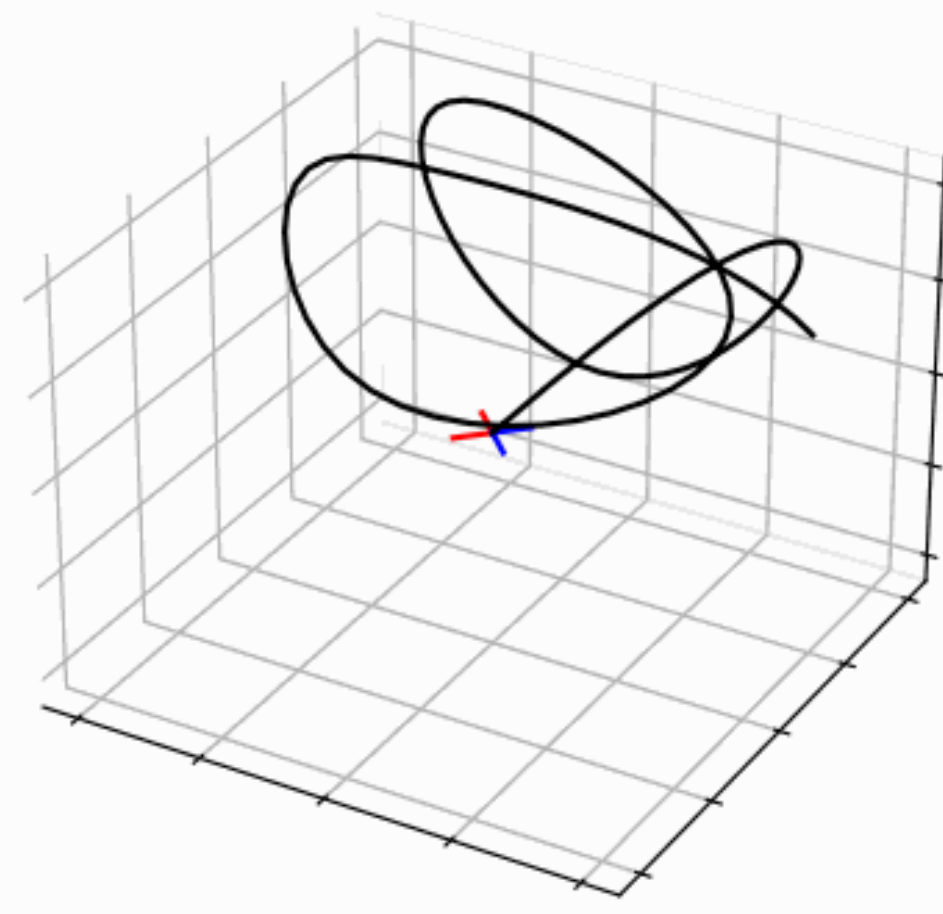**MIT REALM Talk 2024**

**Rajiv Sambharya**

# 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

Current state,
reference trajectory $\longrightarrow$

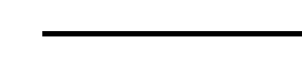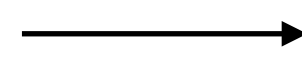**Model predictive controller**

minimize     $\sum_{t=1}^{T} \|x_t - x_t^{\mathrm{ref}}\|_2^2$

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

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

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

$\longrightarrow$ Control
inputs

# Challenge: we need faster methods for optimization

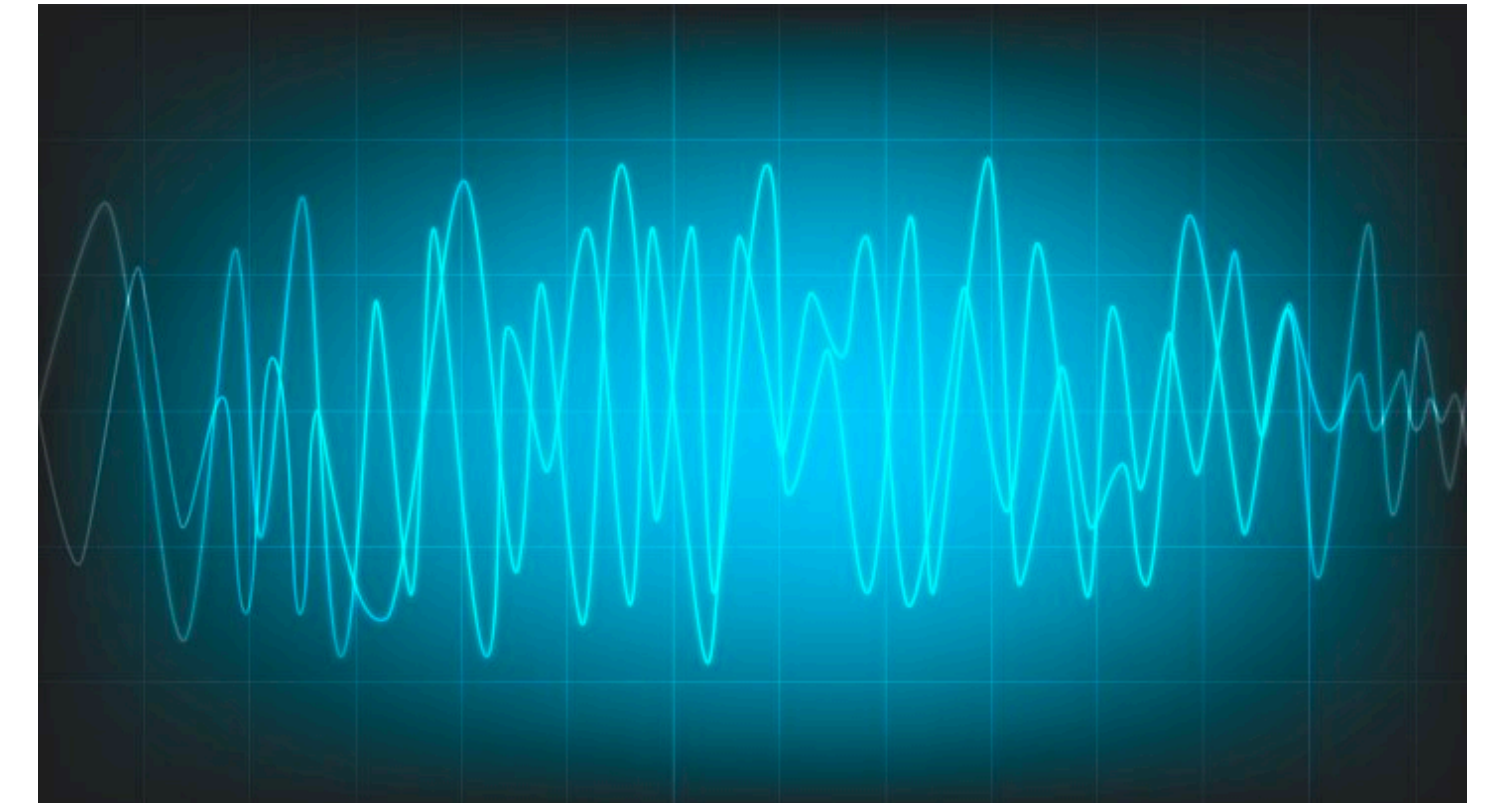# Claim: real-world optimization is parametric

**Robotics and control**

**Energy**

**Signal processing**

# Can machine learning speed up parametric optimization?
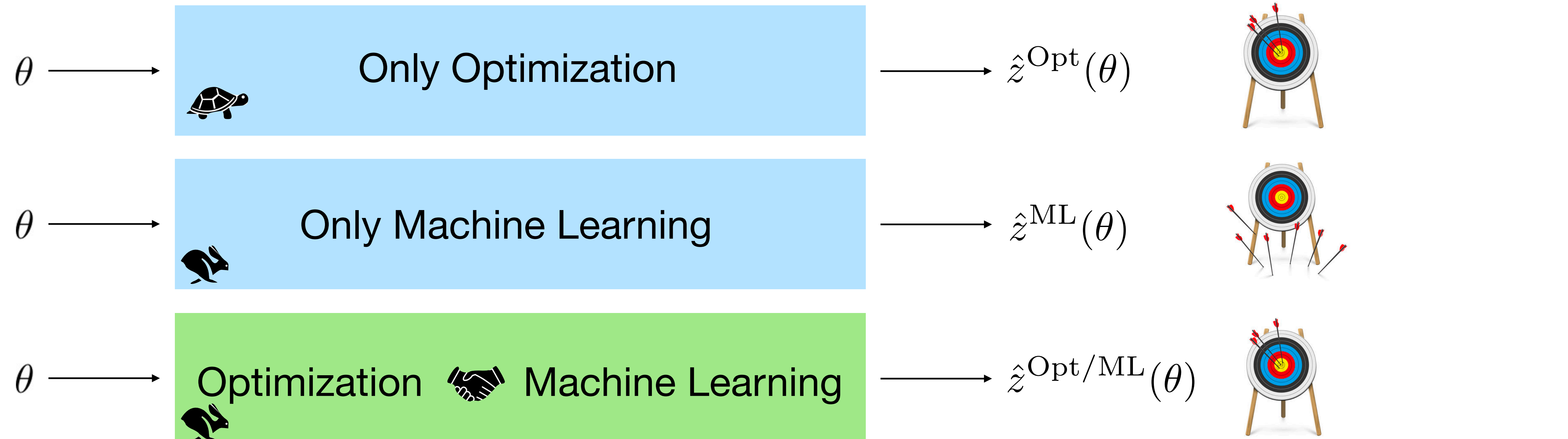
**Goal: Do mapping quickly and accurately**

Parameter

$$\theta \longrightarrow$$

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

Optimal solution

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

$\theta \longrightarrow$  Only Optimization  $\longrightarrow \hat{z}^{\mathrm{Opt}}(\theta)$

$\theta \longrightarrow$  Only Machine Learning  $\longrightarrow \hat{z}^{\mathrm{ML}}(\theta)$

$\theta \longrightarrow$  Optimization 🤝 Machine Learning  $\longrightarrow \hat{z}^{\mathrm{Opt/ML}}(\theta)$

# Learning to Optimize

# The learning to optimize paradigm
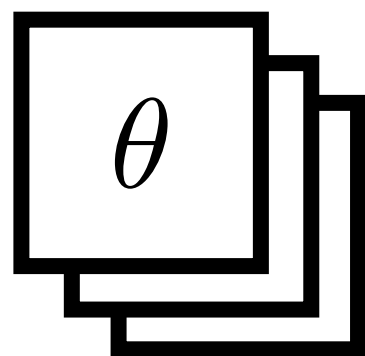
**Goal**: solve the parametric optimization problem fast

$$\text{minimize} \quad f_\theta(z)$$
$$\text{subject to} \quad g_\theta(z) \leq 0$$

**Offline**

**Training**

**Data collection**

Parameters

$\theta$

$\xrightarrow{\text{Solve}}$

Optimal solutions

$z^\star$

Training parameter $\theta$ $\longrightarrow$

Learnable Optimizer with weights $w$

$\longrightarrow \hat{z}_w(\theta) \longrightarrow$ Loss

Candidate solution

Learn

Deploy

**Online evaluation**

Unseen parameter $\theta$ $\longrightarrow$

Learned Optimizer

$\longrightarrow$ High-quality solution

# Challenges in learning to optimize methods

- **I: Lack convergence guarantees**

- **II: Lack generalization guarantees**

- **III: Hard to integrate with state-of-the-art solvers**

<p align="center"><strong>We need <span style="color:green">reliable</span> L2O methods</strong></p>



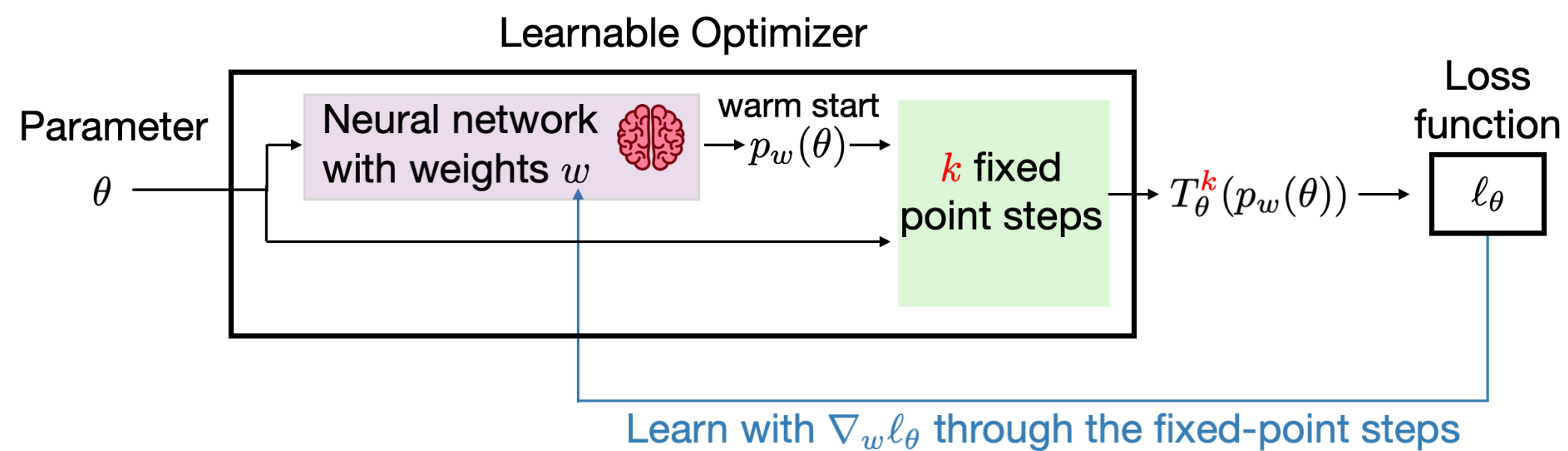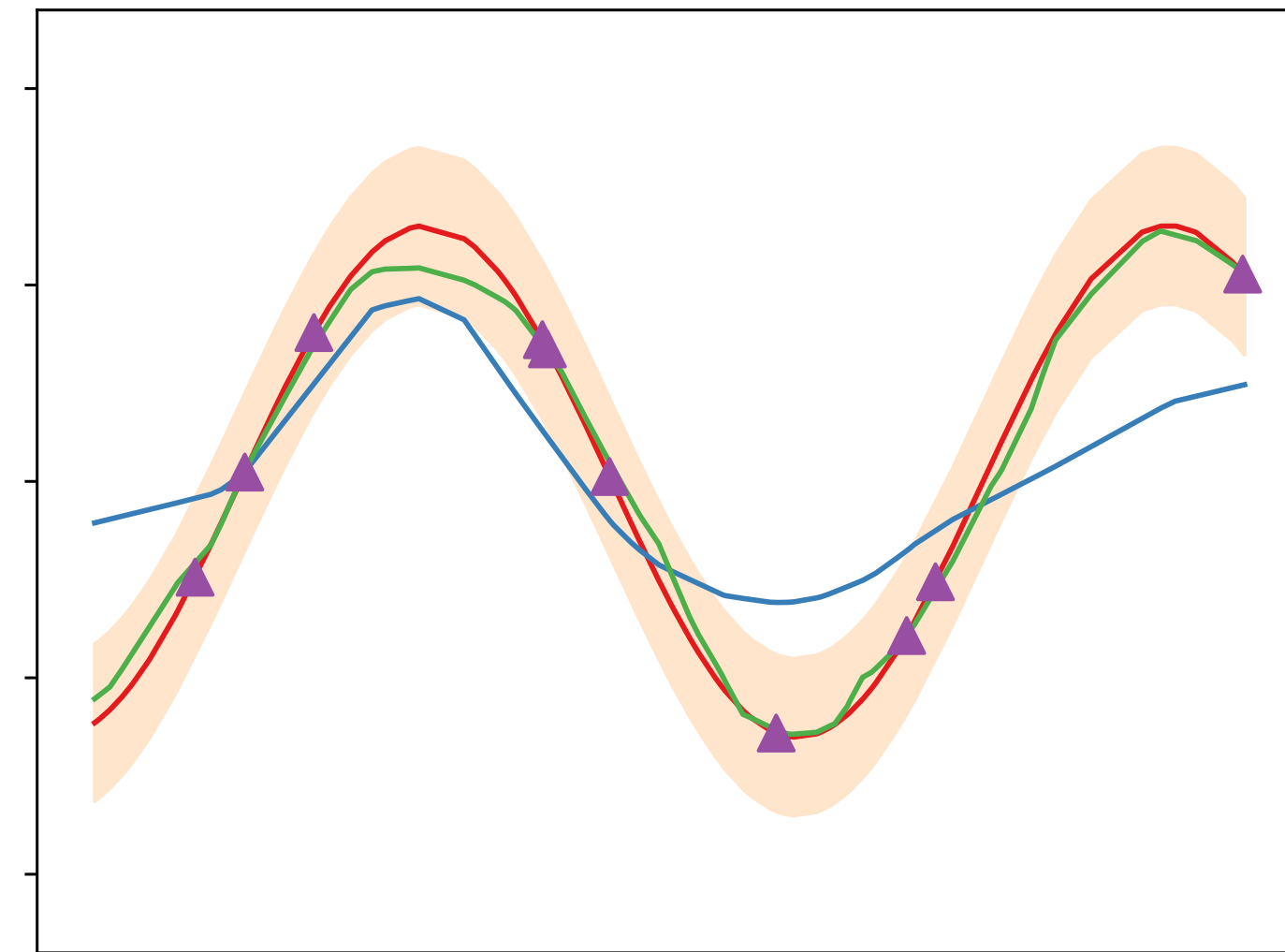<p align="center"><strong>Learning to Optimize: A Primer and A Benchmark</strong> [Chen. et al 2021]</p>

*"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: Practical Performance Guarantees for Classical and Learned Optimizers**
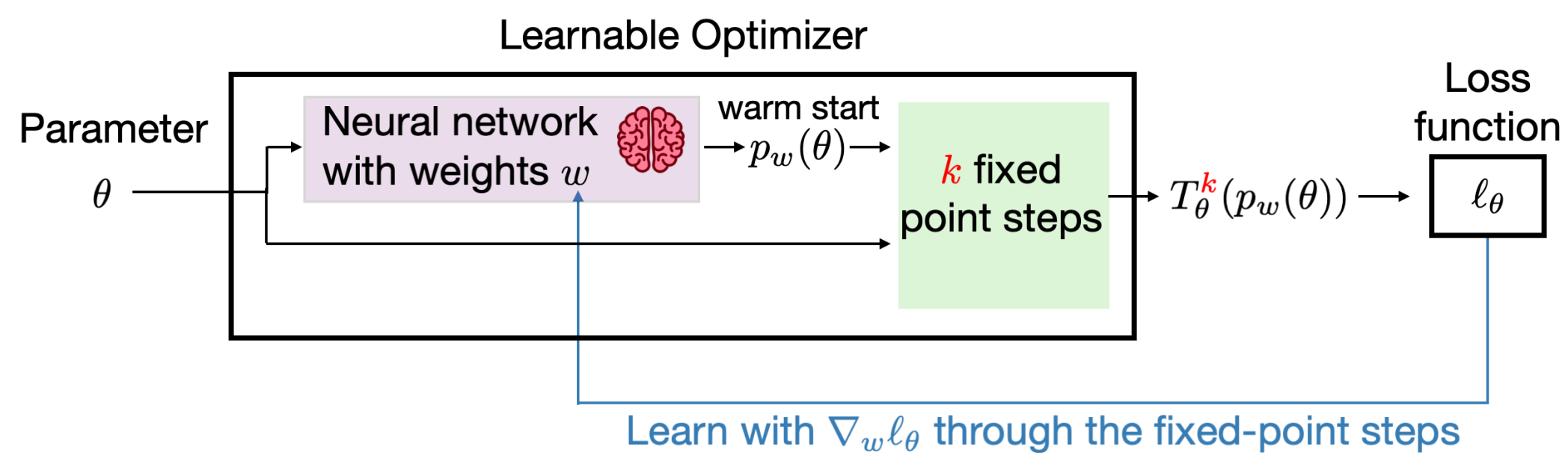
# Collaborators

Georgina
Hall

Brandon
Amos

Bartolomeo
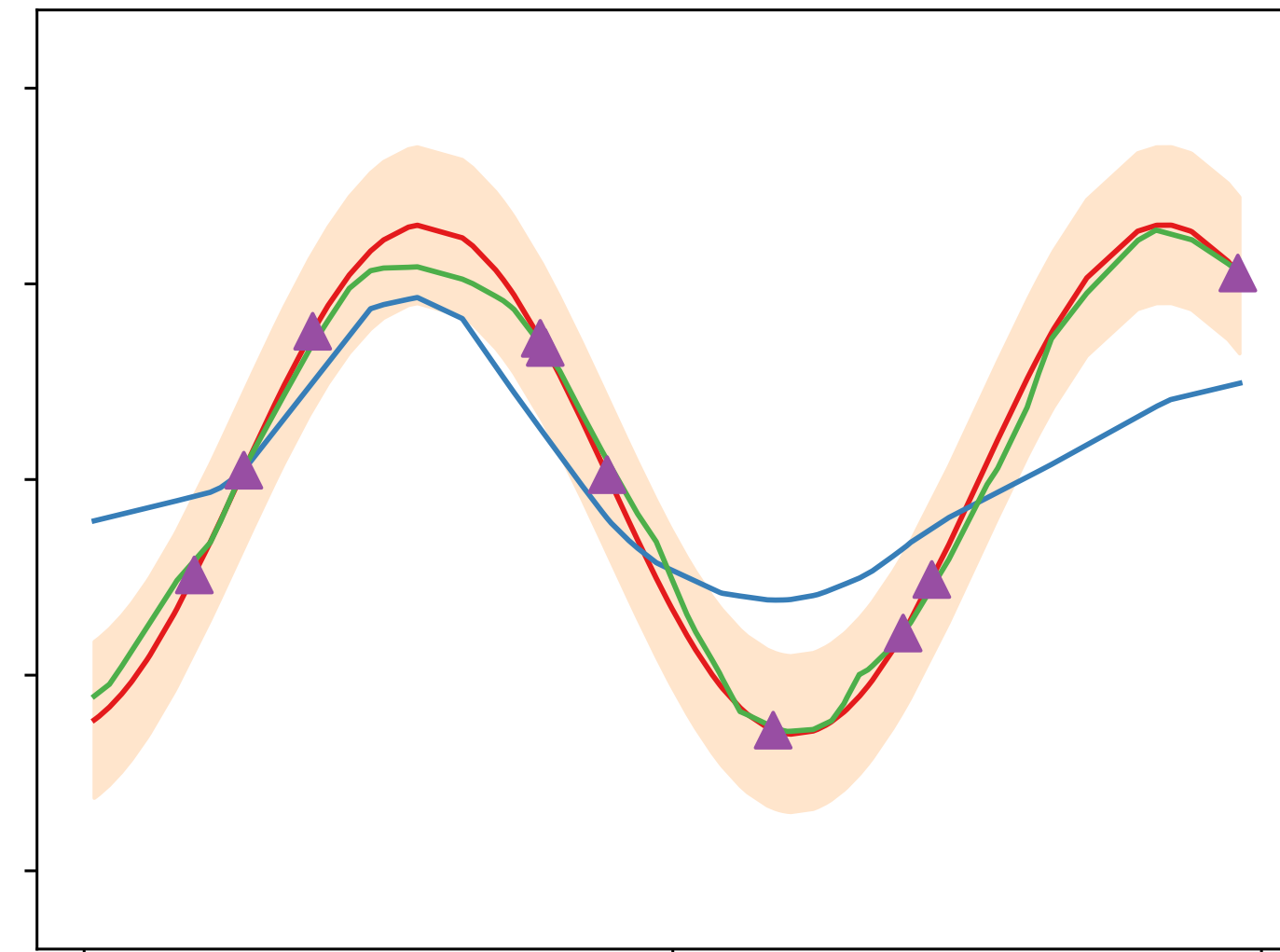Stellato

INSEAD

∞ Meta

PRINCETON
UNIVERSITY

# Talk Outline

- **Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms**

- **Part 2: Practical Performance Guarantees for Classical and Learned Optimizers**
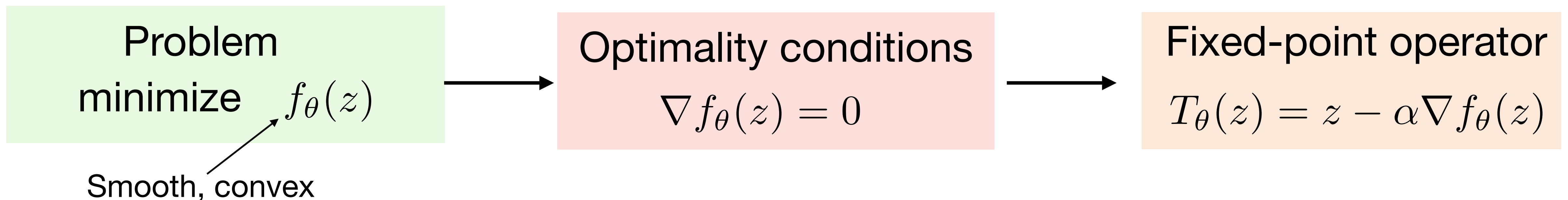


Learnable Optimizer

Parameter $\theta$ → Neural network with weights $w$ → warm start $p_w(\theta)$ → $k$ fixed point steps → $T_\theta^k(p_w(\theta))$ → Loss function $\ell_\theta$

Learn with $\nabla_w \ell_\theta$ through the fixed-point steps

# Fixed-point optimization problems are ubiquitous

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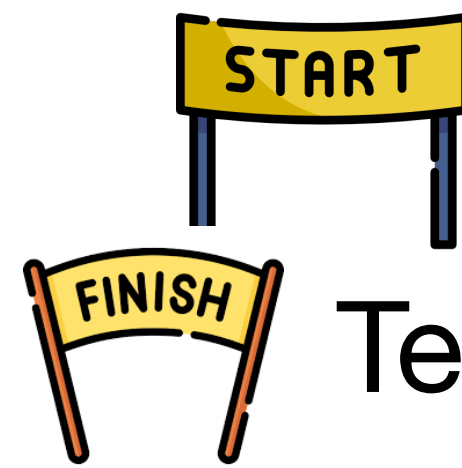
**Convex optimization**

Problem

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

$\longrightarrow$

Optimality conditions
(KKT conditions)

$\longrightarrow$

Fixed-point operator

**Unconstrained, smooth convex optimization**

Problem
minimize $f_\theta(z)$

Smooth, convex

$\longrightarrow$

Optimality conditions
$\nabla f_\theta(z) = 0$

$\longrightarrow$

Fixed-point operator
$T_\theta(z) = z - \alpha \nabla f_\theta(z)$

# 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 $\|T_\theta(z^j) - z^j\|_2$ is small

**Fixed-point residual**

---

**Example: Proximal gradient descent**

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

$\qquad\qquad$ Convex $\qquad$ Convex
$\qquad\qquad$ Smooth $\qquad$ Non-smooth

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

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

**Problem: limited iteration budget**

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

12

# Some warm starts are better than others

minimize     $10z_1^2 + z_2^2$

subject to   $z \geq 0$

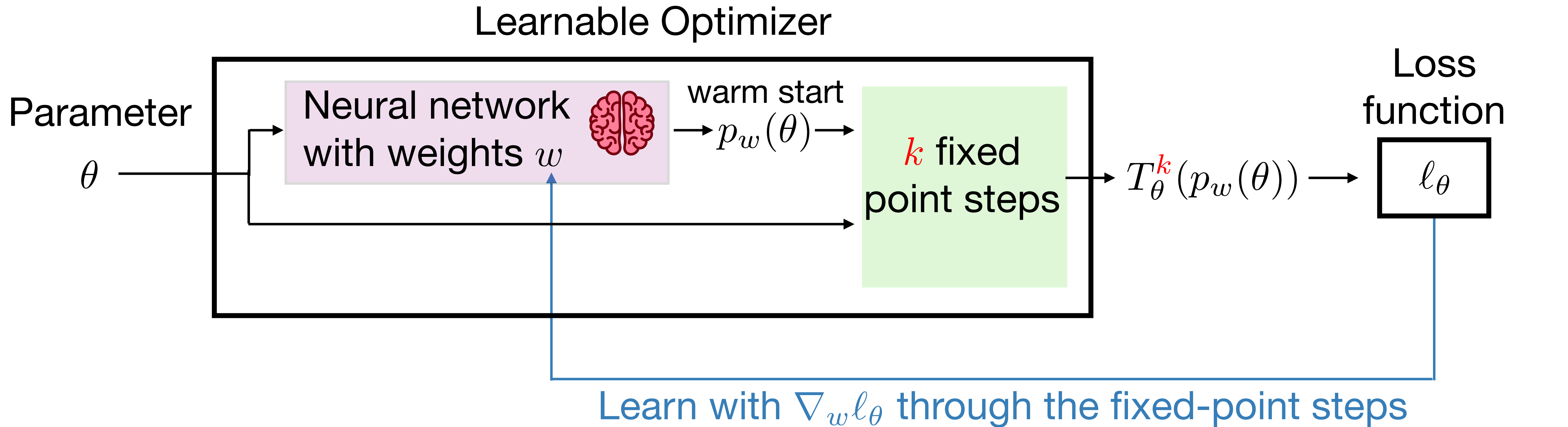★  Optimal solution at the origin

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**
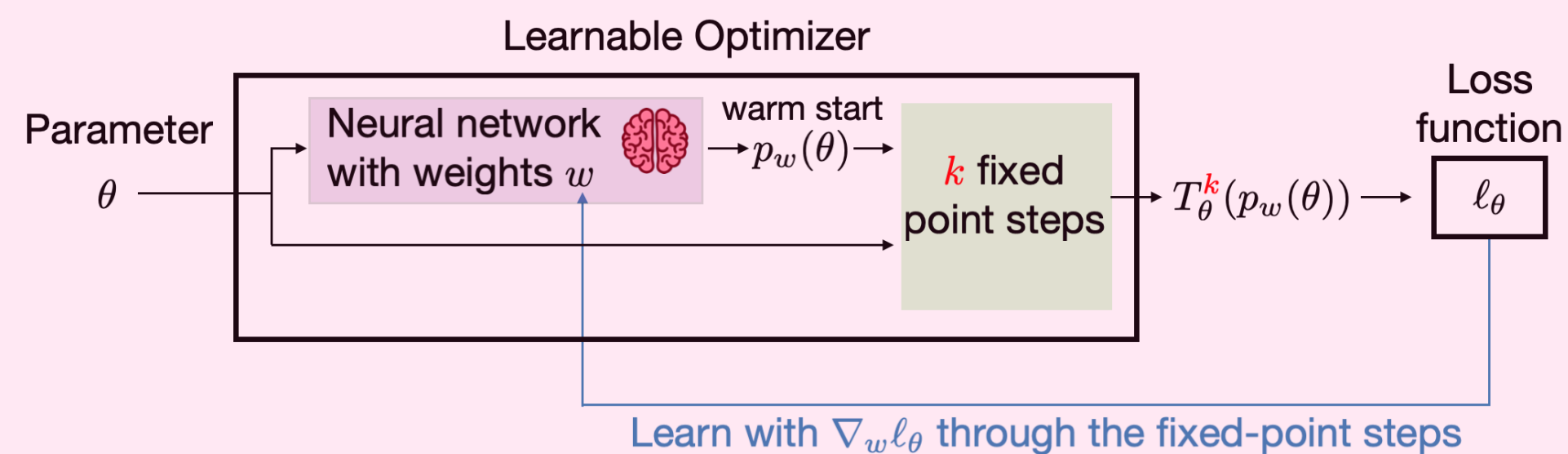
# End-to-end learning architecture

Learnable Optimizer

Parameter

$\theta$

Neural network with weights $w$ 🧠

warm start
$\rightarrow p_w(\theta) \rightarrow$

$k$ fixed point steps

$T_\theta^k(p_w(\theta)) \rightarrow$

Loss function

$\ell_\theta$

Learn with $\nabla_w \ell_\theta$ through the fixed-point steps

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

Parameter $\theta$

⚡ Learned solver with convergence

Solution $z^\star(\theta)$

**Generalization guarantees**



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

**Easy integration with popular solvers**


SCS
SPLITTING CONIC SOLVER

minimize $\quad (1/2)x^T P x + c^T x$

subject to $\quad Ax + s = b$

$\qquad\qquad s \in \mathcal{K}$

Conic programs

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

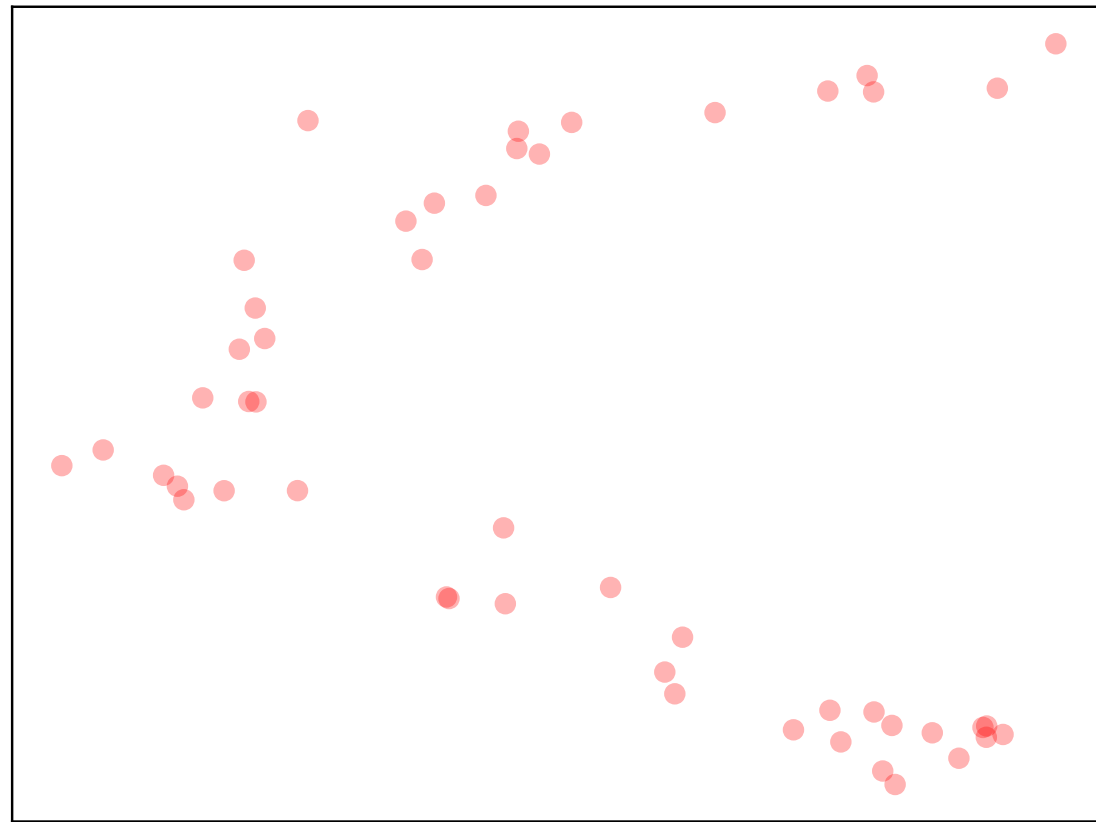Allows us to quantify solve time in seconds

# Numerical Experiments

**Comparing our learned warm starts** 🧠 **against**

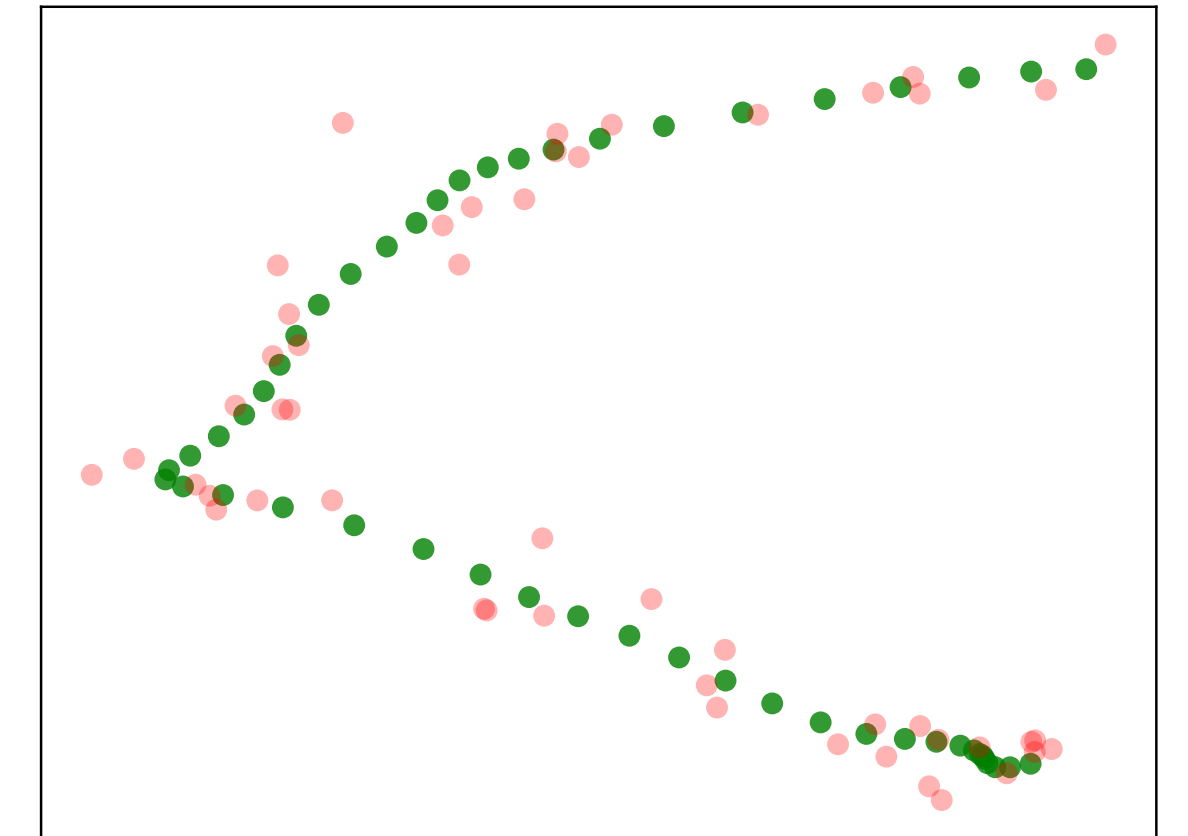**Baseline initializations**

  **1. Cold-start: initialize at zero** ❄️

  **2. Nearest neighbor: initialize with solution of nearest training problem**

# Robust Kalman filtering



Robust Kalman filtering

## Second-order cone program

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

Noisy trajectory

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

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

$\qquad\qquad\; y_t = Cx_t + v_t \quad \forall t$

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

Recovered trajectory

Dynamics matrices: $A, B$

Observation matrix: $C$

Huber loss: $\psi_\rho$

# Robust Kalman filtering visuals



Solution after 5 fixed-point steps
with different initializations

**Noisy trajectory**

**Optimal solution**

**With learning, we can estimate the state well**

Nearest neighbor

Previous solution

Learned: $k = 5$

# Model predictive control (MPC) of a quadcopter

Current state, previous control reference trajectory $\longrightarrow$

**Controller**

$\longrightarrow$ Control inputs

**Quadratic program**

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

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 $x_{t+1} = A(\theta)x_t + B(\theta)u_t$

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

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

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

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

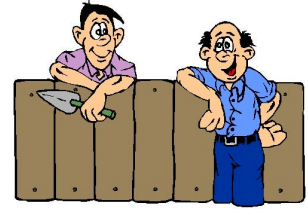$u_{-1} = u_{\text{prev}}$

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

Linearized dynamics

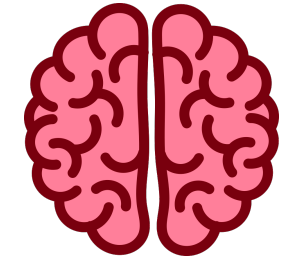# MPC of a quadcopter in a closed loop

**Budget of 15 fixed-point steps**
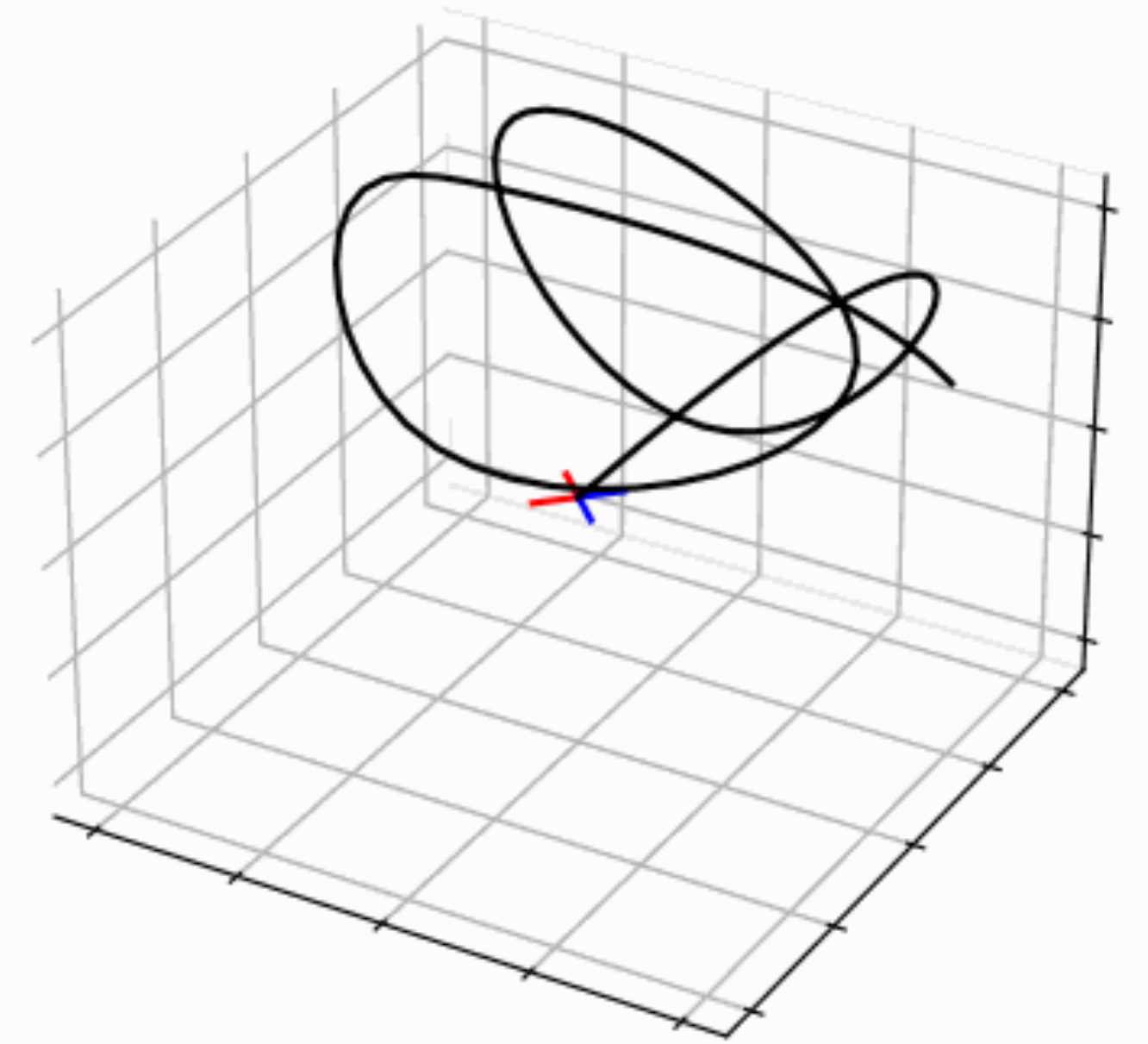
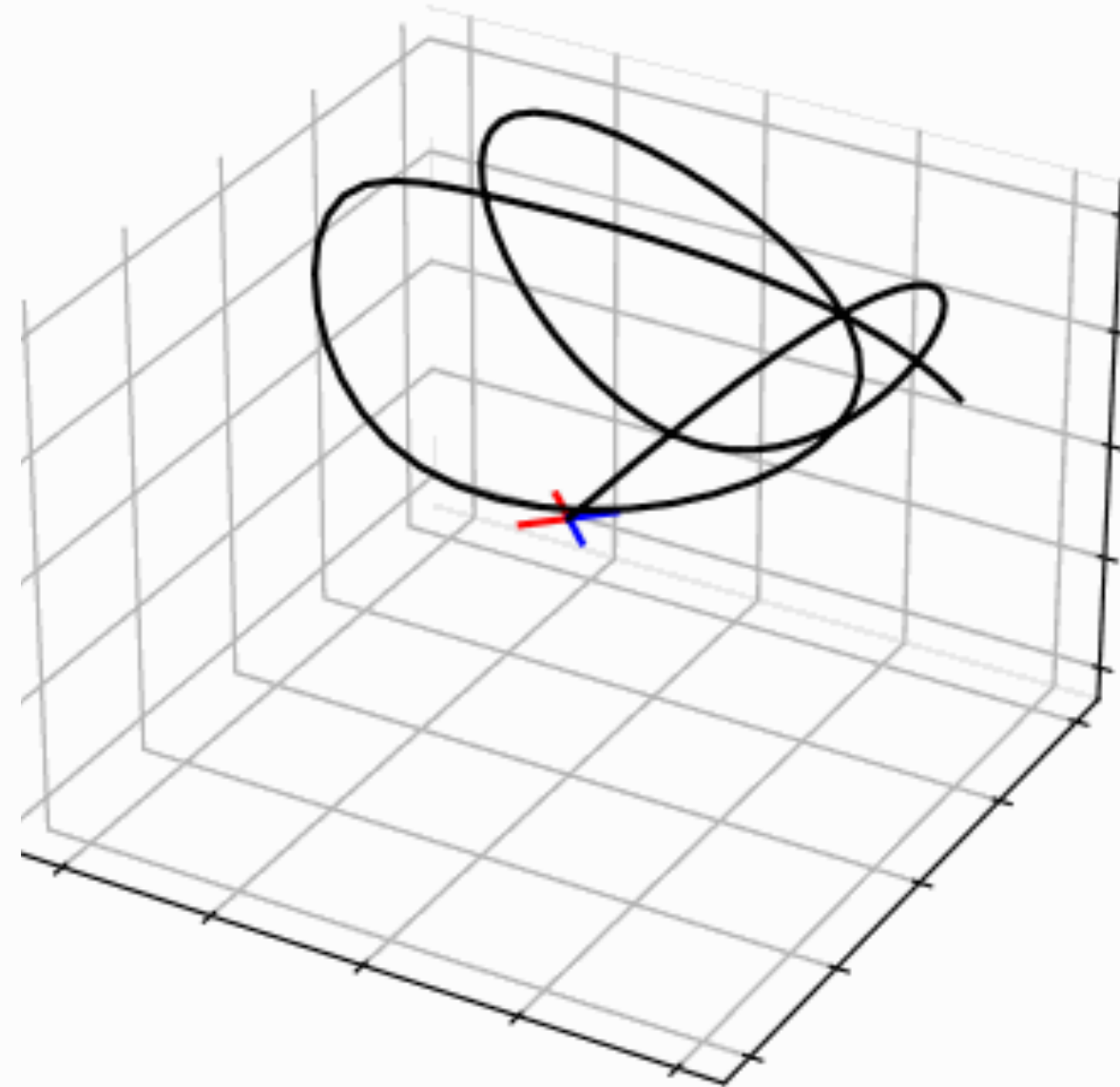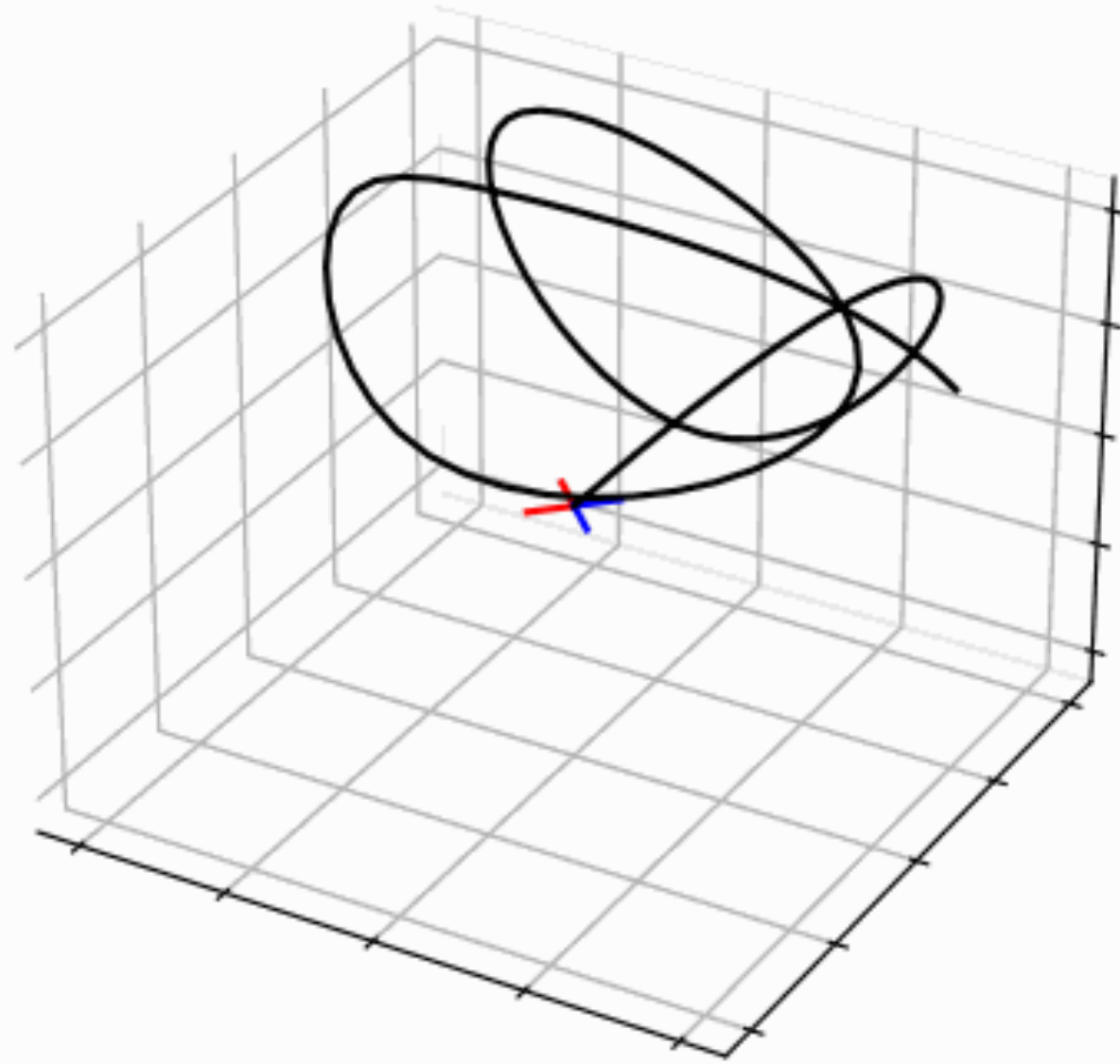Nearest neighbor                    Previous solution                    Learned: $k = 5$

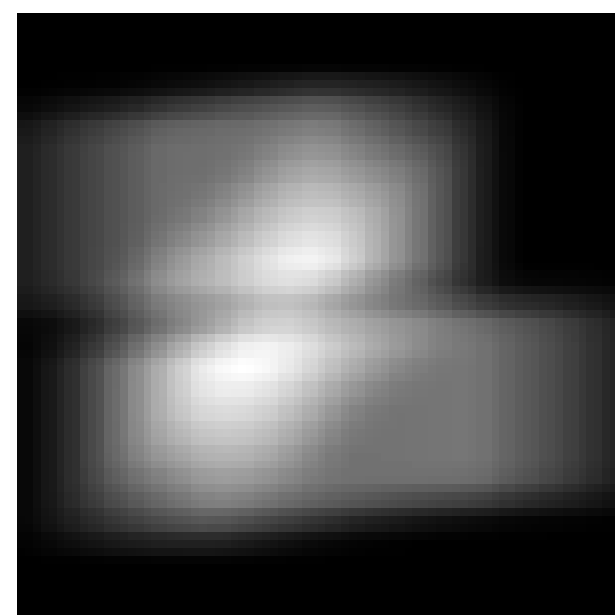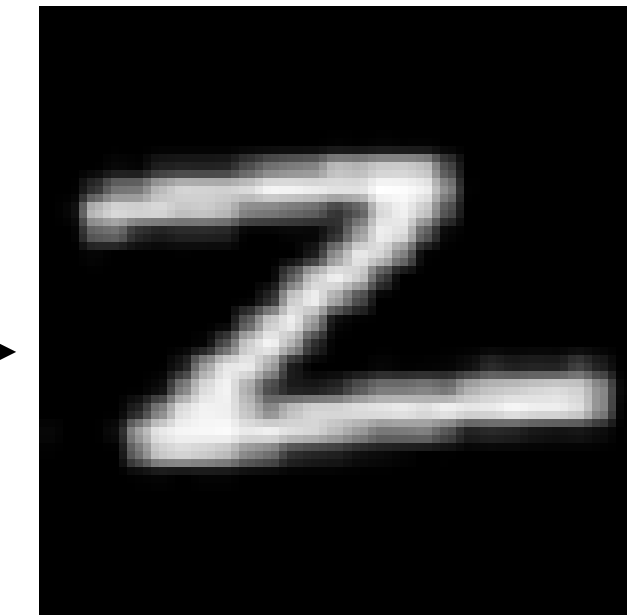**With learning, we can track the trajectory well**
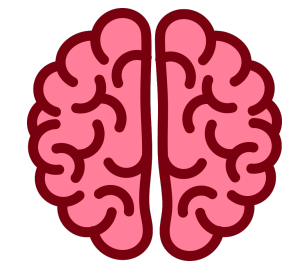
# Image deblurring



**Image deblurring**

**Quadratic program**

$\theta = b$
Blurred image

minimize $\quad \|Ax - b\|_2^2 + \lambda\|x\|_1$

subject to $\quad 0 \leq x \leq 1$

$x^\star$
Deblurred image

$A$: blur operator

# Image deblurring

| percentile | optimal | blurred | cold-start ❄️ | nearest neighbor 🚧 | learned 🧠 |
|---|---|---|---|---|---|

**50 fixed-point steps**

| | | | | | |
|---|---|---|---|---|---|
| 10th | | | | | |
| 50th | | | | | |
| 90th | | | | | |
| 99th | | | | | |

**Distance to nearest neighbor increases**

**With learning, we can deblur all of the images quickly**

22

# Talk Outline

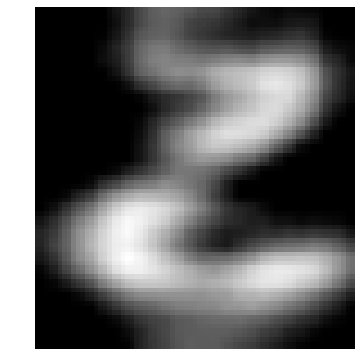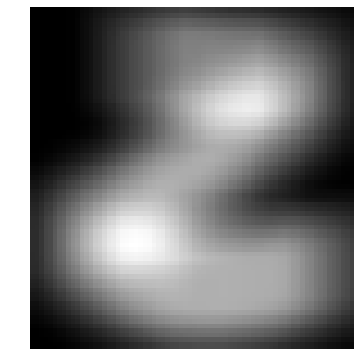- **Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms**

- **Part 2: Practical Performance Guarantees for Classical and Learned Optimizers**

# Talk Outline

- **Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms**

- **Part 2: Practical Performance Guarantees for Classical and Learned Optimizers**

**Classical = no learning**

# Worst-case bounds can be very loose



Example: robust Kalman filtering

**Second-order cone program**

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

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

$\quad\quad\quad\quad\quad y_t = Cx_t + v_t \quad \forall t$

▼ —— SCS empirical average performance over 1000 parametric problems

◄ —— Worst-case bound

In practice: **linear** convergence over the parametric family

Worst-case analysis: **sublinear** convergence

Worst-case bounds do not consider the **parametric** structure

Approach: solve N problems and then bound

# We will bound 0-1 error metrics

## We will provide guarantees for any measured quantity

algorithm steps

tolerance

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

## Standard metrics

*e.g.,* fixed-point residual

algorithm steps    cold start    tolerance

$$e(\theta) = \mathbf{1}(\ell^{\mathrm{fp}}_\theta(T^k_\theta(0)) > \epsilon)$$

## Task-specific metrics:

*e.g.,* quality of extracted states in robust Kalman filtering

recovered state    optimal state

$$e(\theta) = \mathbf{1}\left(\max_{t=1,\ldots,T} \|x_t - x^\star_t\|_2 > \epsilon\right)$$

# Background: Kullback-Liebler Divergence

**KL divergence**: measures distance between distributions

$$\mathrm{KL}(q \parallel p) = \sum_{i=1}^{m} q_i \log \left( \frac{q_i}{p_i} \right)$$

Our bounds on the risk will take the form

$$\mathrm{KL}(\text{empirical risk} \parallel \text{risk}) \leq \text{regularizer}$$

**Invert** these bounds by solving

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

1D convex optimization problem

$$\mathrm{KL}^{-1}(q \mid c) = \text{maximize} \quad p$$

$$\text{subject to} \quad q \log \frac{q}{p} + (1-q) \log \frac{1-q}{1-p} \leq c$$

$$0 \leq p \leq 1$$

# Statistical learning theory can provide probabilistic guarantees

algorithm steps

tolerance

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

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

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

Number of problems

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

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

# Robust Kalman filtering guarantees



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

# Visualizing Robust Kalman filtering guarantees



**Task-specific error metric**

$$e(\theta) = \mathbf{1}\left(\max_{t=1,\ldots,T} \|x_t - x_t^\star\|_2 > \epsilon\right)$$

- Noisy trajectory
- Optimal solution
- Solution after 15 steps
- Region with guarantee

"With high probability, 90% of the time, all of the recovered states after 15 steps of problems drawn from the distribution will be within the correct ball with radius 0.1"

# Talk Outline

- **Part 1: Learning to Warm-Start Fixed-Point Optimization Algorithms**

- **Part 2: Practical 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

algorithm steps

tolerance

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

learnable weights

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

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

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

Optimize the bounds directly

32

# PAC-Bayes training architecture to optimize the guarantees

**PAC-Bayes Training**

Learn

Training parameter $\theta$ ⟶ 

Learnable Optimizer
$P, P_0$
$w \sim P$

Stochastic candidate solution
$\hat{z}_w(\theta)$ ⟶ L2O loss

posterior   prior

Regularizer ⟶ KL inverse ⟶ Guarantee

Use differentiable optimization

We show that the derivative always exists

We implement the learnable optimizer and train with this architecture

33

# Learned algorithms for sparse coding

<span style="color:red">Noisy</span>
<span style="color:red">measurements</span>
$\theta = b$

$\longrightarrow$

**Sparse coding**

Recover sparse $z^\star$ from $b = Dz^\star + \sigma$

$\longrightarrow$

<span style="color:blue">Ground truth</span>
<span style="color:blue">sparse signal</span>
$z^\star$

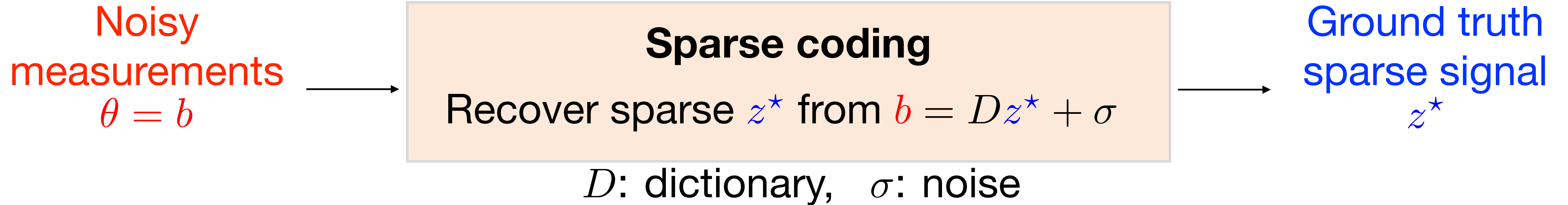$D$: dictionary,   $\sigma$: noise

Standard technique

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

ISTA (iterative shrinkage thresholding algorithm)

(Classical optimizer)

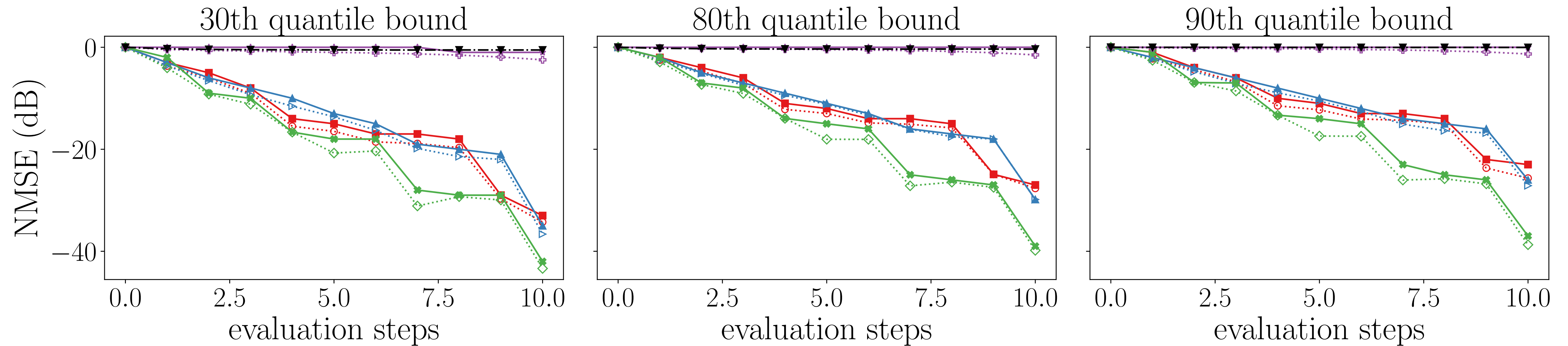$z^{j+1} = \text{soft threshold}_{\frac{\lambda}{L}}\left(z^j - \frac{1}{L}(Dz^j - b)\right)$

Learned ISTA
(Learned optimizer)

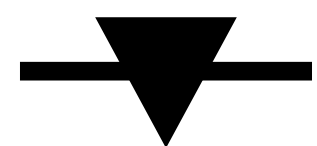$z^{j+1} = \text{soft threshold}_{\psi^j}\left(W_1^j z^j + W_2^j b\right)$

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

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

# Learned ISTA results for sparse coding

# K-shot Meta-Learning for Sine Curves



**Neural network learning**

find weights $z$ so that $g_z(x_i) \approx y_i$

predictor with weights $z$

Training dataset
with K points

$\mathcal{D}^{\mathrm{train}}$

**Gradient step**

$$\hat{z} = z - \alpha \nabla_z \mathcal{L}(z, \mathcal{D}^{\mathrm{train}})$$

Weights that generalize
to new points quickly

$\hat{z}$

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



Legend:
- **Used for gradients** (purple triangle)
- **Ground truth** (red)

After 10 grad steps
- **Stochastic MAML** (green)
- **Pretrained** (blue)
- **Region with MAML guarantee** (orange band)

With high probability, 90% of the time stochastic MAML after 10 steps will stay within the band

The pretrained baseline only stays within the band 30% of the time

# Future directions

**Optimization**

**Learning**

**Dynamical systems and control**

**Connections with REALM**

Learning dynamical systems, certificates for stability and safety

Focus on guarantees

**Safe Control with Learned Certificates: A Survey of Neural Lyapunov, Barrier, and Contraction Methods for Robotics and Control** [Dawson et. al 2023]

*"Closely related is the issue of generalization error, which relates a learned certificate's performance on a finite training set with its performance on the full state space… Some works have [obtained] probabilistic upper-bounds on the generalization error, but these bounds tend to be conservative."*

# Conclusions

We do not need to sacrifice **guarantees for learning-based systems**

Learning to Warm-Start
Fixed-Point Optimization Algorithms

<span style="color:#8B0000">**Journal of Machine Learning Research
(accepted conditioned
on minor revision)**</span>
**https://arxiv.org/pdf/2309.07835.pdf**

End-to-End Learning to Warm-Start for
Real-Time Quadratic Optimization

<span style="color:#8B0000">**Learning for Dynamics and
Control Conference**</span>
**https://arxiv.org/pdf/2212.08260.pdf**

Practical Performance Guarantees
for Classical and Learned Optimizers

<span style="color:#8B0000">**To be on Arxiv soon!**</span>

rajivs@princeton.edu

rajivsambharya.github.io