

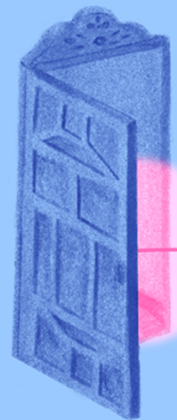
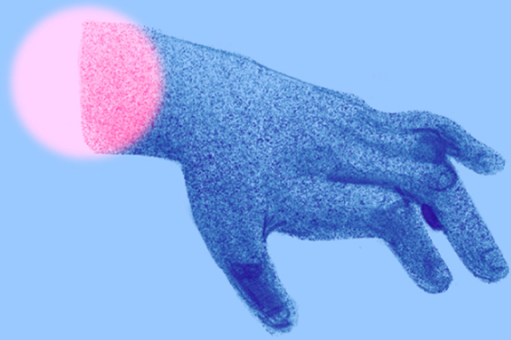


CONSTRAINED

OPTIMIZATION

*for*

MACHINE LEARNING



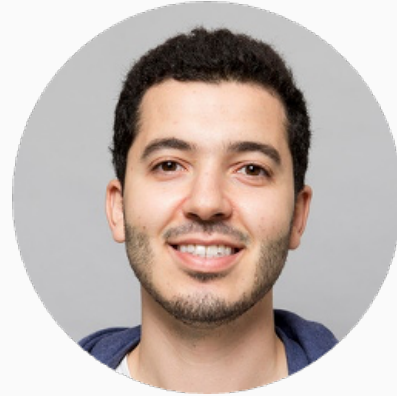
# Today's agenda

- ▶ Why constrained optimization?
- ▶ Foundations of constrained optimization
- ▶ Three core papers
- ▶ Conclusions and perspectives



*“If I had been rich, I probably would not have devoted myself to mathematics.”*

# Collaborators



Akram Erraqabi



Tianyue H. Zhang



Juan Ramirez



Meraj Hashemizadeh



Motahareh Sohrabi



Rohan Sukumaran



Simon Lacoste-Julien



Golnoosh Farnadi



Yoshua Bengio





## Controlled Sparsity via Constrained Optimization

**Gallego-Posada**, Ramirez, Erraqabi, Bengio, Lacoste-Julien

NeurIPS 2022

## L0onie: Compressing COINs with L0-constraints

Ramirez, **Gallego-Posada\***

Sparsity in Neural Networks Workshop 2022

## Balancing Act: Constraining Disparate Impact in Sparse Models

Hashemizadeh\*, Ramirez\*, Sukumaran, Farnadi, Lacoste-Julien, **Gallego-Posada**

ICLR 2024



## On PI controllers for updating Lagrange multipliers in constrained optimization

Sohrabi\*, Ramirez\*, Zhang, Lacoste-Julien, **Gallego-Posada**

ICML 2024



## Cooper: A Library for Constrained Optimization in Deep Learning

**Gallego-Posada\***, Ramirez\*, Hashemizadeh, Lacoste-Julien

JMLR MLOSS 2024 (under submission)





## **GAIT: A Geometric Approach to Information Theory**

**Gallego-Posada**, Vani, Schwarzer, Lacoste-Julien

*AISTATS 2020*

## **Equivariant Mesh Attention Networks**

Basu\*, **Gallego-Posada\***, Vigano\*, Rowbottom\*, Cohen

*TMLR 2022*

## **AI & Cities: Risks, Applications and Governance**

Koseki, Jameson et al.

*Tech Report - Mila and UN-Habitat 2022*

## **A Distributed Data-Parallel PyTorch Implementation of the Distributed Shampoo Optimizer for Training Neural Networks At-Scale**

Shi, Lee, Iwasaki, **Gallego-Posada**, Li, Rangadurai, Mudigere, Rabbat

*Tech Report 2024*

# Introduction





Widespread deployment of powerful machine learning models has resulted in mounting pressures to enhance the robustness, safety and fairness of such models—often arising from regulatory and ethical considerations



# ~~“Build now, fix later”~~

- ▶ **Inability to guarantee compliance** with industry standards and governmental regulations **limits implementation** of ML solutions in real-world applications
- ▶ Retro-fitting safety measures as afterthoughts!
- ▶ Continuous **incurrence of technical debt hinders long-term progress** of the field

# Secure by design

- ▶ We advocate for a paradigm shift in which **constraints are an integral part** of the model development process
- ▶ Constrained optimization offers a **rich conceptual framework** accompanied by **algorithmic tools** for reliably enforce complex properties on ML models

# *Recent works on CO for ML*

---

- ▶ **Fairness:** Zafar et al. (2017); Cotter et al. (2019); Hashemizadeh et al. (2024)
- ▶ **Safe reinforcement learning:** Stooke et al. (2020)
- ▶ **Sparse neural network training:** Gallego-Posada et al. (2022)
- ▶ **Active learning:** Elenter et al. (2022)
- ▶ **Model quantization:** Hounie et al. (2023)
- ▶ **Dynamics of constrained learning:** Sohrabi et al. (2024)
- ▶ **Safe RLHF:** Dai et al. (2024)

Problem

Formulation

Algorithm

Analysis

# Problem

- ▶ **Controllability**
- ▶ **Hyperparameter interpretability**
- ▶ **Better exploration of trade-offs**
- ▶ **Experimental accountability**

Formulation

Algorithm

Analysis

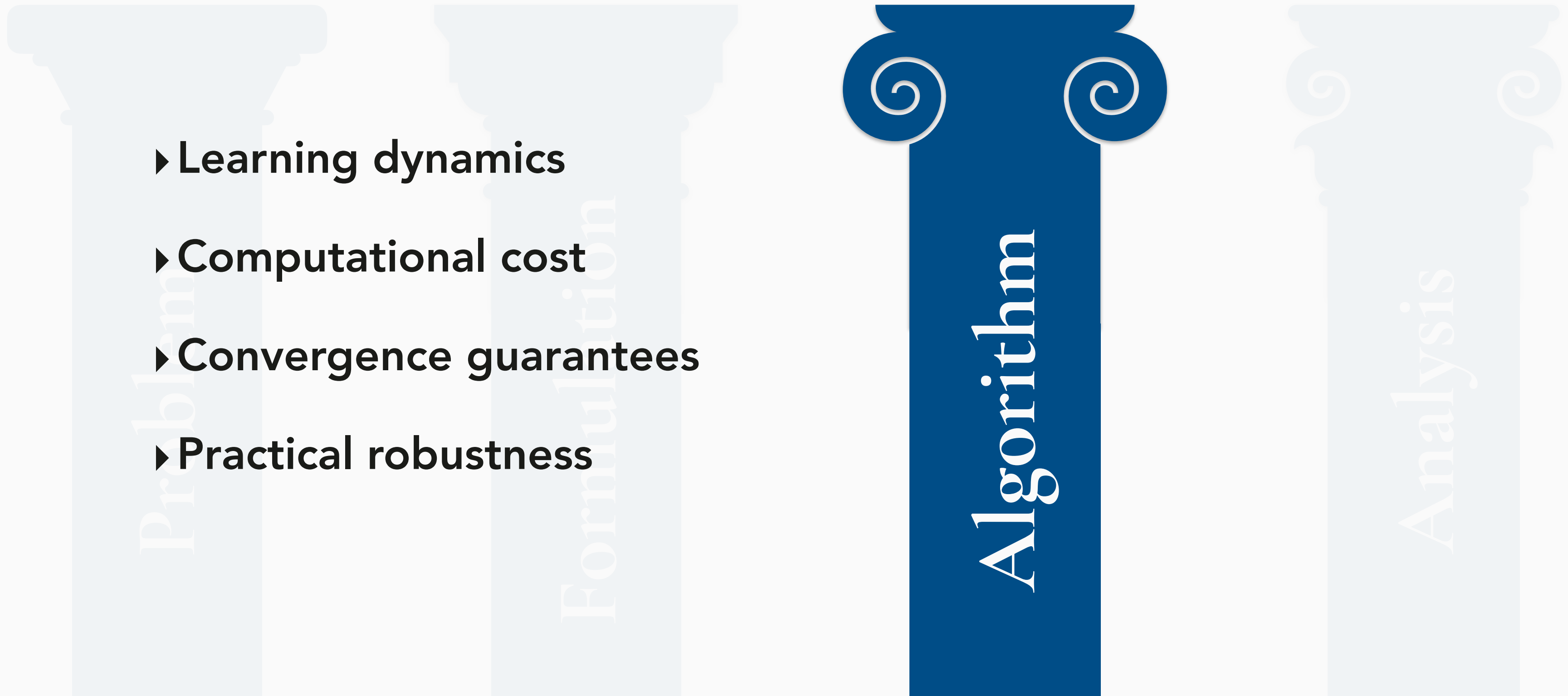
Problem

Formulation

- ▶ **Game structure**
- ▶ **Functional representation  
of the problem**

Algorithm

Analysis

- 
- ▶ Learning dynamics
  - ▶ Computational cost
  - ▶ Convergence guarantees
  - ▶ Practical robustness

Algorithm

Analysis

Problem

- ▶ Feasibility reigns
- ▶ Two axis for generalization
- ▶ How fast to become feasible?  
How fast to improve the loss?

Formulation

Algorithm

Analysis



# Constrained optimization

minimize  $f(\mathbf{x})$   
 $\mathbf{x}$

subject to  $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}_m$  and  $\mathbf{h}(\mathbf{x}) = \mathbf{0}_n$

## Feasible set

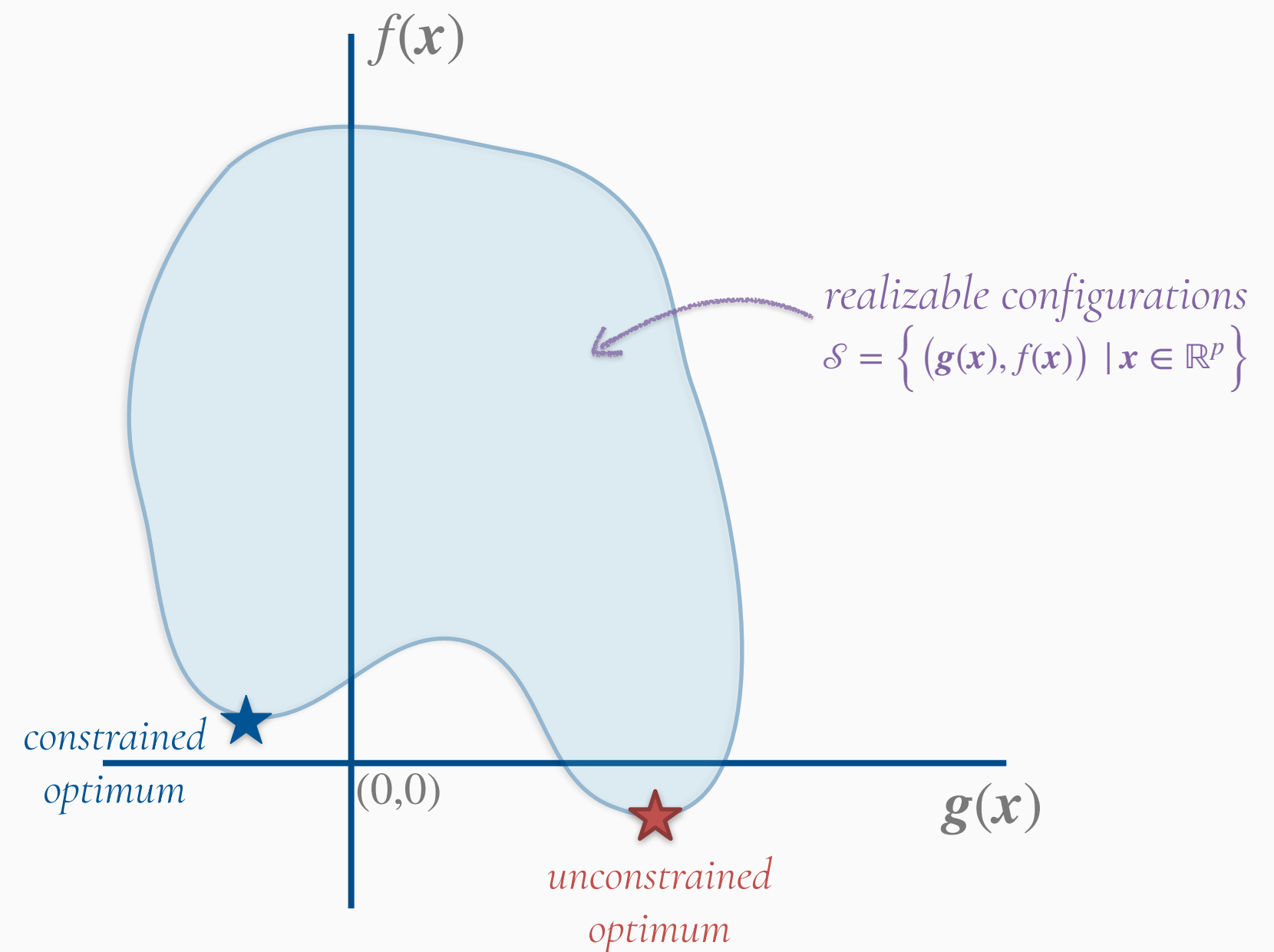
$$\Omega = \{ \mathbf{x} \in \mathbb{R}^p \mid \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \text{ and } \mathbf{h}(\mathbf{x}) = \mathbf{0} \}$$

## Optimality condition (necessary)

If  $\mathbf{x}^*$  is a local minimum of  $f$  over  $\Omega$ , then

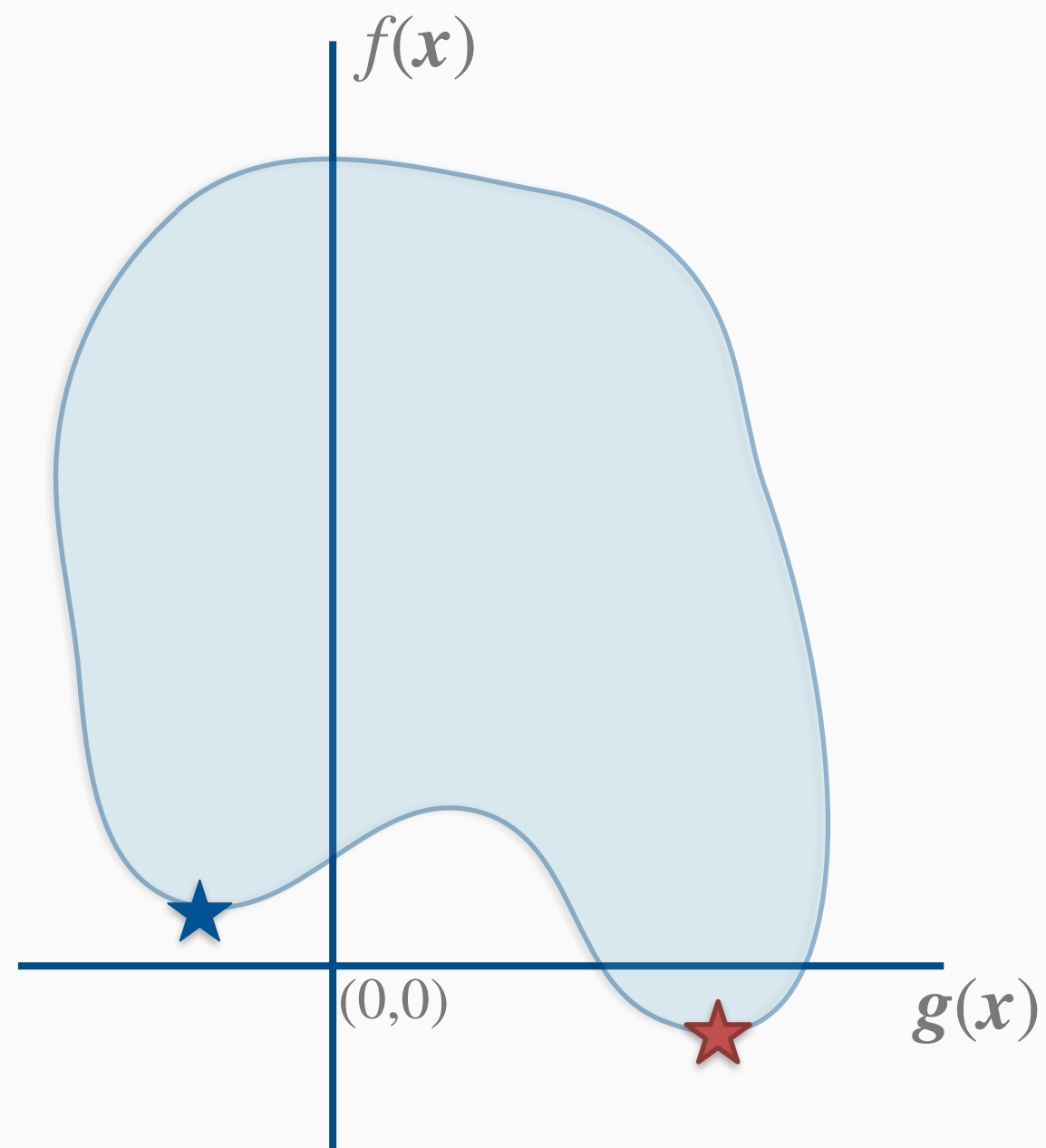
$$\nabla f(\mathbf{x}^*)^\top \mathbf{z} \geq 0 \quad \forall \mathbf{z} \in \text{FD}(\mathbf{x}^*)$$

*feasible directions at  $\mathbf{x}^*$*



# Feasibility and accountability

---

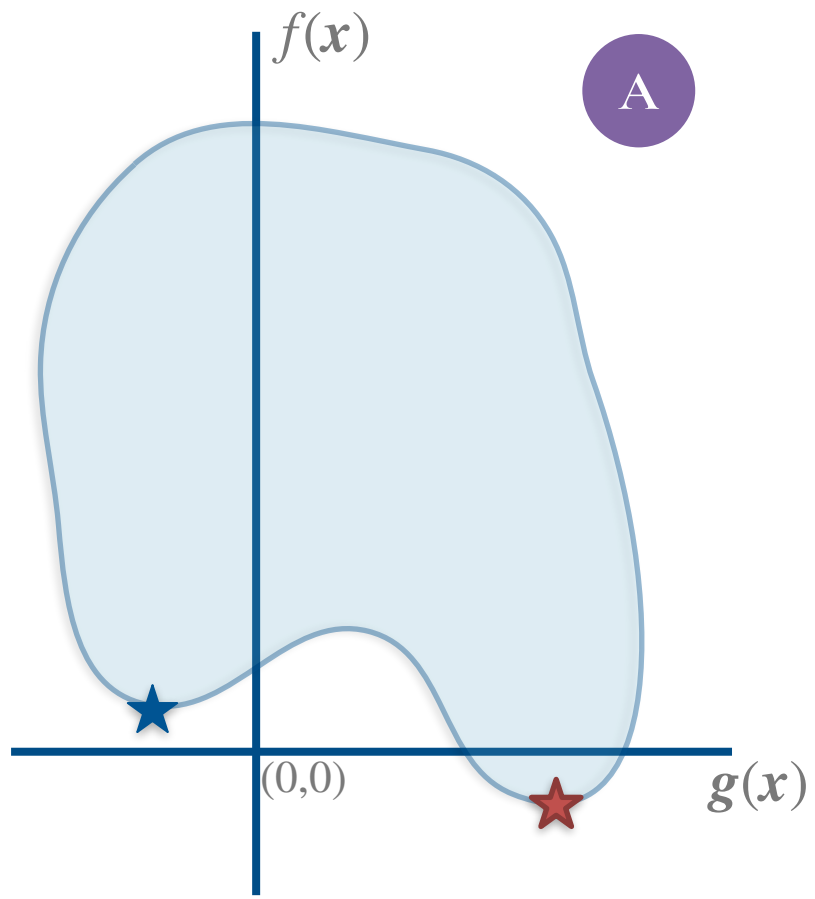


Although in the unconstrained setting ★ would be preferred, not valid solution for constrained problem since it is infeasible.

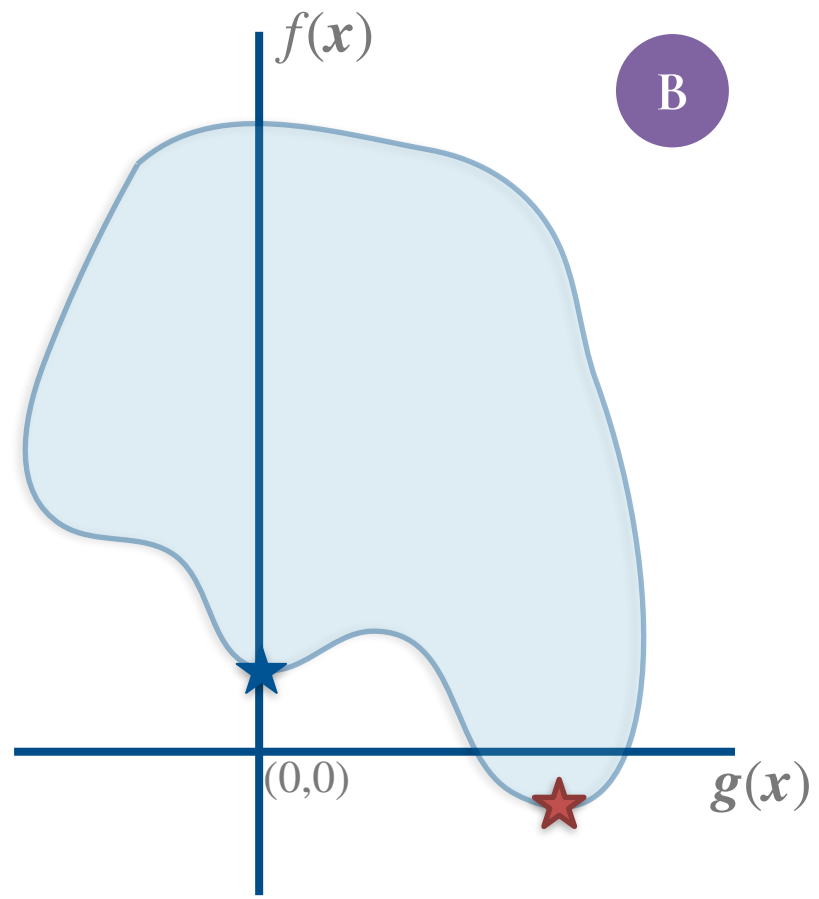
*usually informed by problem-dependent requirements*

Choosing the constraint level **beforehand** ensures experimental accountability.

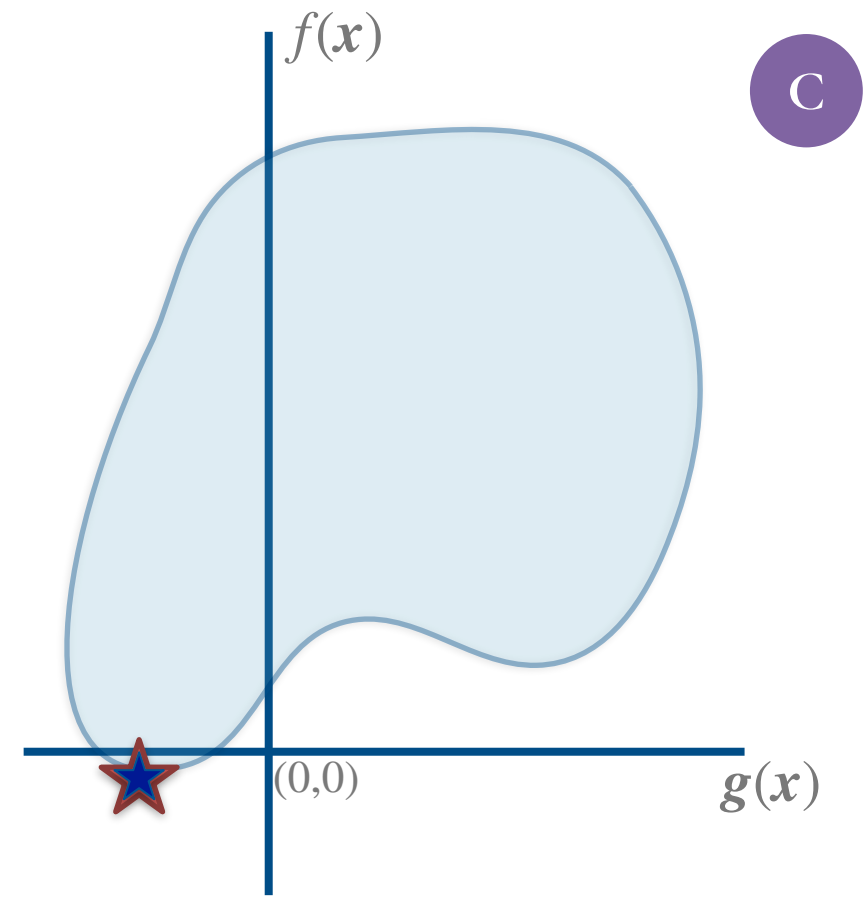
*“not allowed to cheat”*



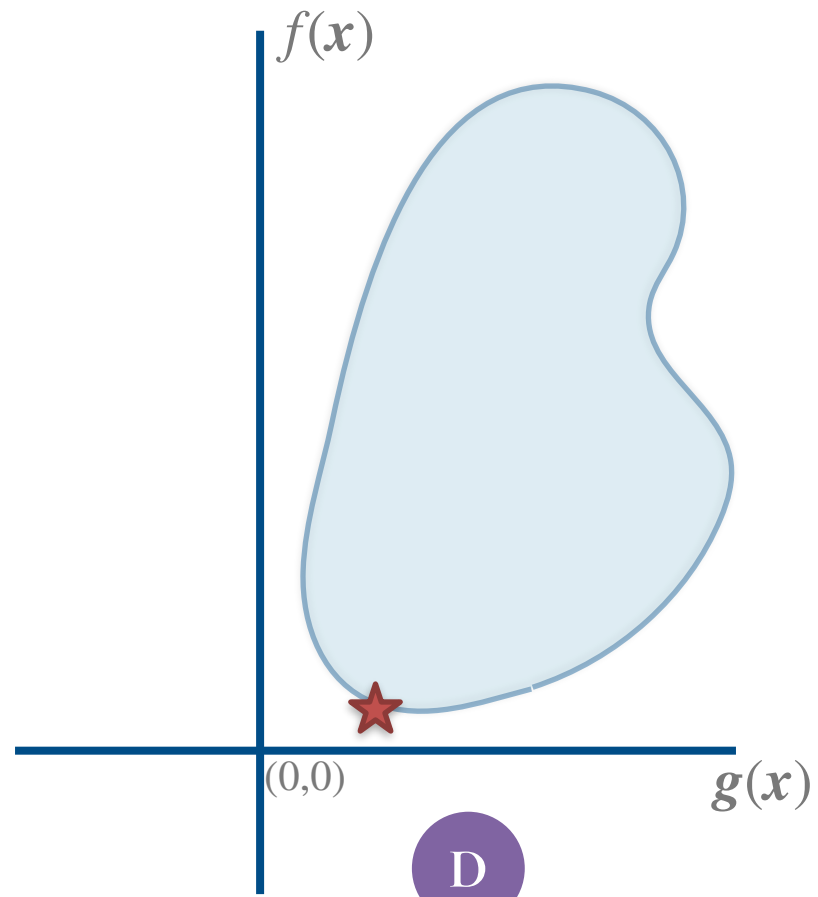
A



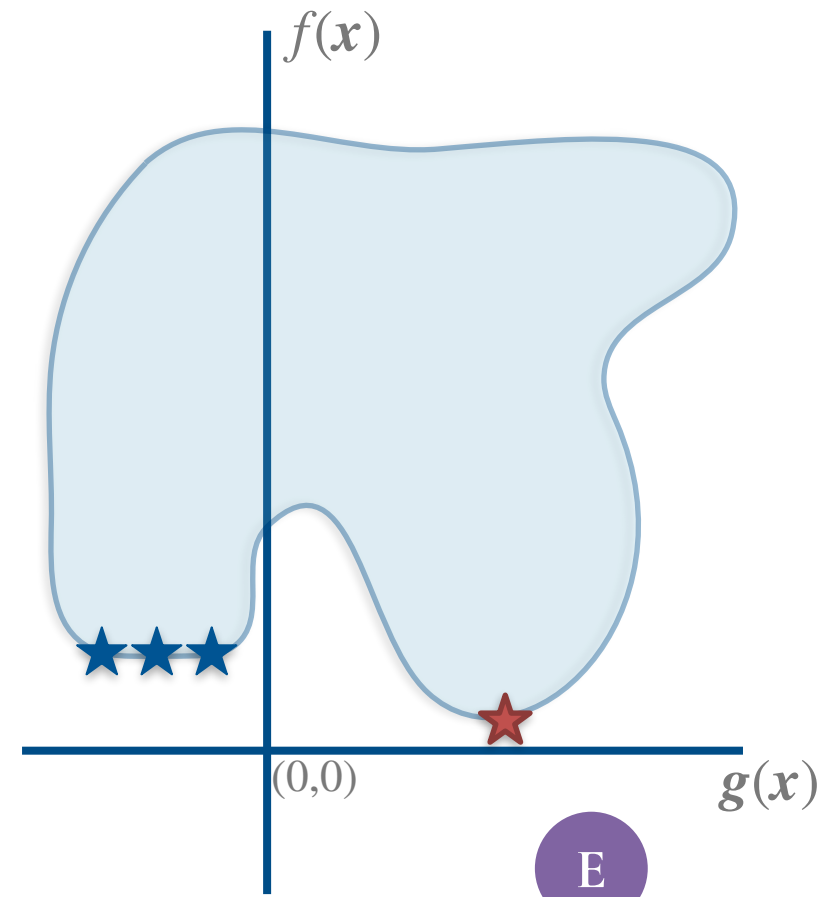
B



C

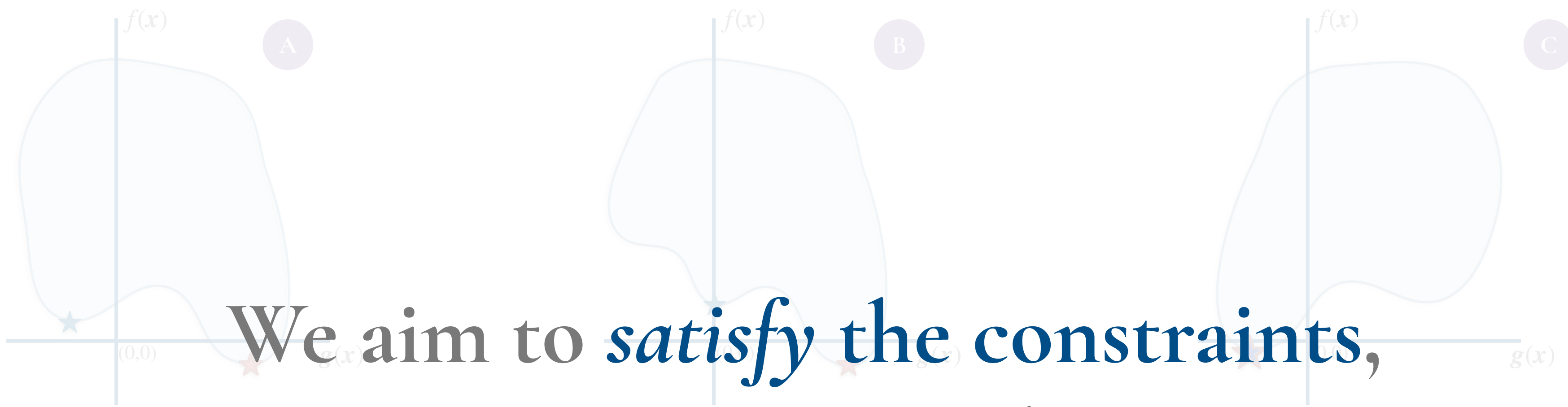


D



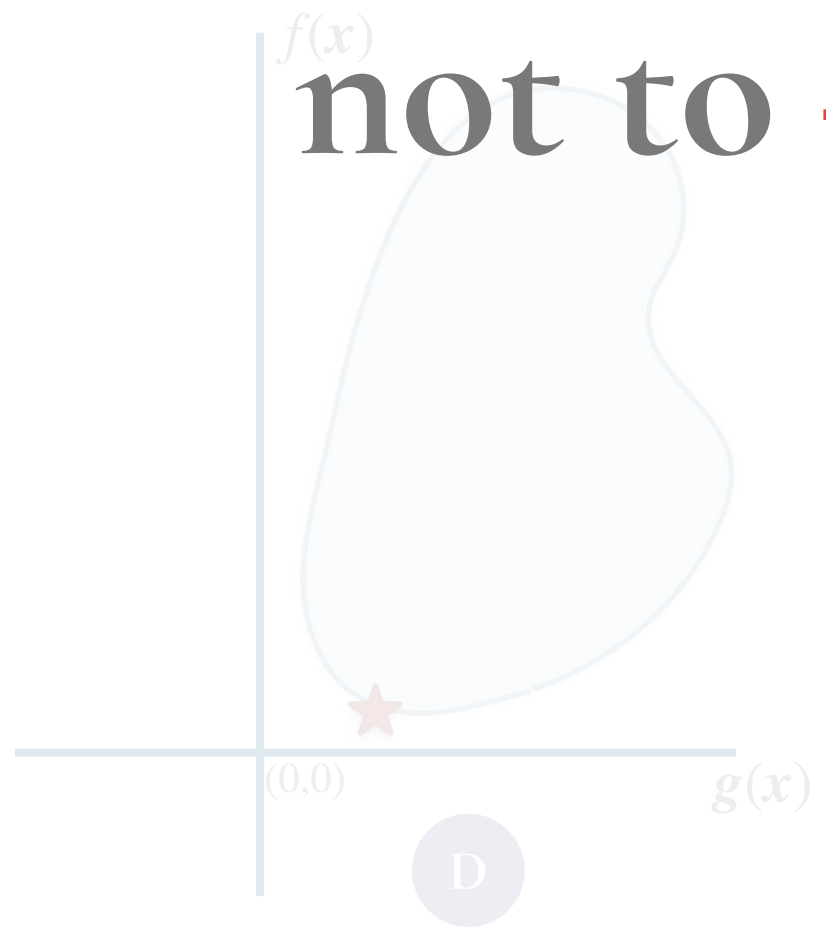
E



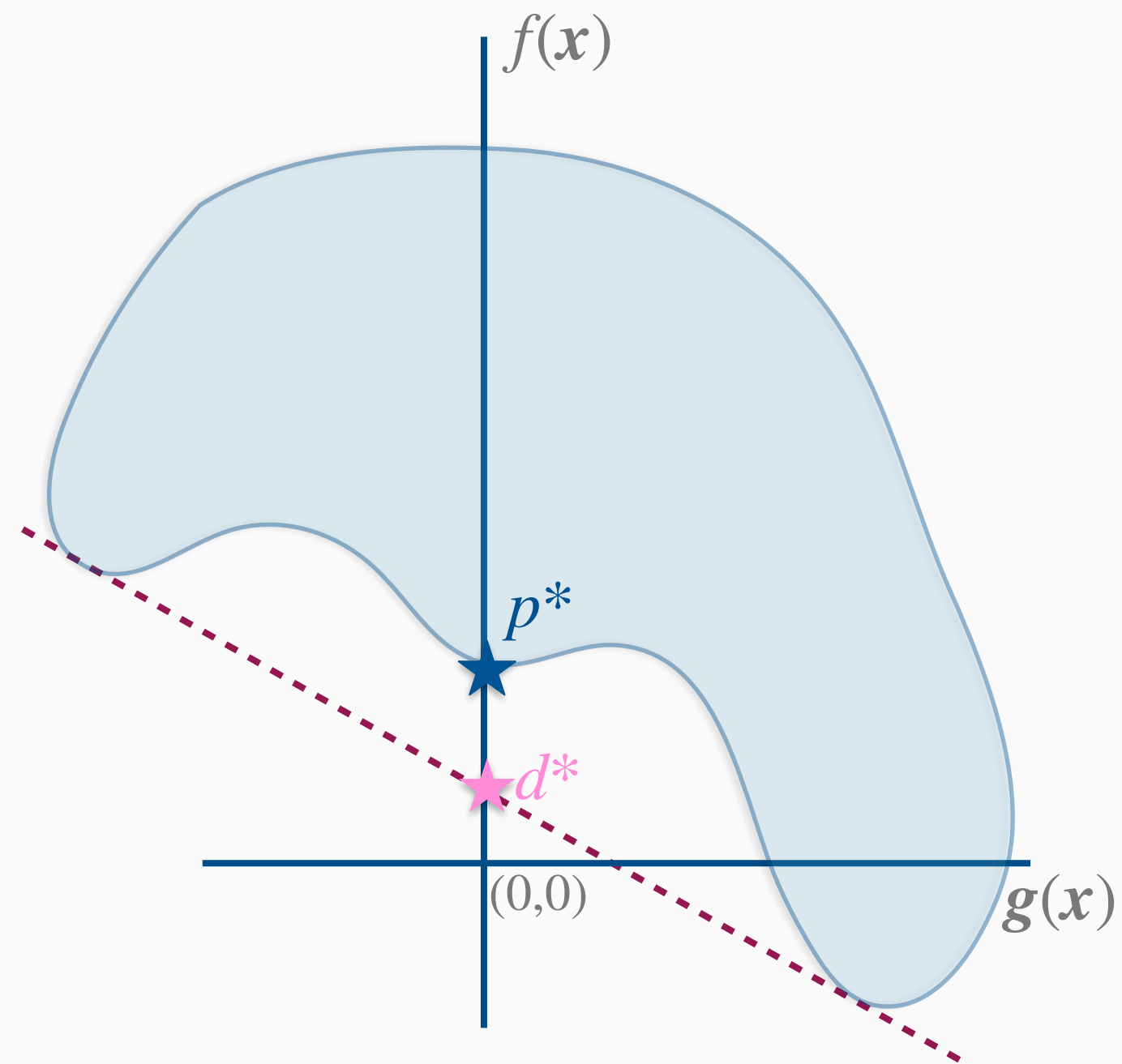


We aim to *satisfy* the constraints,

not to ~~optimize~~ them!



# Why not just penalize?



*tunable hyperparameter* ↖

$$\underset{x}{\text{minimize}} \quad f(x) + \lambda_{\text{pen}} g(x)$$

Tuning  $\lambda_{\text{pen}}$  typically requires a trial-and-error search!

In non-convex problems, there may be trade-offs between the objective and constraints that are **not reachable** using a penalized formulation.

● ● ●  $Duality\ gap \triangleq p^* - d^* \geq 0$

# Lagrangian problem

---

## Lagrangian

$$\begin{array}{l} \min_x f(\mathbf{x}) \\ \text{subject to } \mathbf{g}(\mathbf{x}) \leq \mathbf{0}_m \text{ and } \mathbf{h}(\mathbf{x}) = \mathbf{0}_n \end{array} \iff \min_x \max_{\lambda \geq \mathbf{0}, \mu} \mathcal{L}(\mathbf{x}, \lambda, \mu) \triangleq f(\mathbf{x}) + \lambda^\top \mathbf{g}(\mathbf{x}) + \mu^\top \mathbf{h}(\mathbf{x})$$

*“Lagrange multipliers” or “dual variables”*

**Role of the multipliers** (cf. Karush-Kuhn-Tucker necessary conditions)

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla \mathbf{g}_i(\mathbf{x}^*) + \sum_{i=1}^n \mu_i^* \nabla \mathbf{h}_i(\mathbf{x}^*) = \mathbf{0}$$

## Algorithmic approach

Saddle points of the Lagrangian correspond to constrained optima, but may not exist.

Find a min-max point!

# Gradient Descent-Ascent (GDA)

---

Lagrangian  $\min_x \max_{\lambda \geq 0, \mu} \mathcal{L}(x, \lambda, \mu) \triangleq f(x) + \lambda^\top g(x) + \mu^\top h(x)$

## Algorithm

Initialize  $x_0, \lambda_0 = \mathbf{0}$  and  $\mu_0 = \mathbf{0}$

Repeat

$$\mu_{k+1} \leftarrow \mu_k + \eta_{\text{dual}} \nabla_{\mu} \mathcal{L}(x_k, \lambda_k, \mu_k) = \mu_k + \eta_{\text{dual}} \mathbf{h}(x_k)$$

$$\lambda_{k+1} \leftarrow [\lambda_k + \eta_{\text{dual}} \nabla_{\lambda} \mathcal{L}(x_k, \lambda_k, \mu_k)]^+ = [\lambda_k + \eta_{\text{dual}} \mathbf{g}(x_k)]^+$$

$$x_{k+1} \leftarrow x_k - \eta_{\text{primal}} \nabla_x \mathcal{L}(x_k, \lambda_k, \mu_k)$$

If convergence check satisfied; **stop**

*projected gradient ascent  
maintains non-negativity  
of inequality multipliers*

# Gradient Descent-Ascent (GDA)

---

$$\boldsymbol{\mu}_{k+1} \leftarrow \boldsymbol{\mu}_k + \eta_{\text{dual}} \nabla_{\boldsymbol{\mu}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$$

$$\boldsymbol{\lambda}_{k+1} \leftarrow [\boldsymbol{\lambda}_k + \eta_{\text{dual}} \nabla_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)]^+$$

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \eta_{\text{primal}} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \boldsymbol{\lambda}_k, \boldsymbol{\mu}_k)$$

## Extensibility

Simplest possible first-order strategy. Can be combined with **more sophisticated updates**.

*pick your favourite  
primal optimizer!*

## Negligible computational overhead

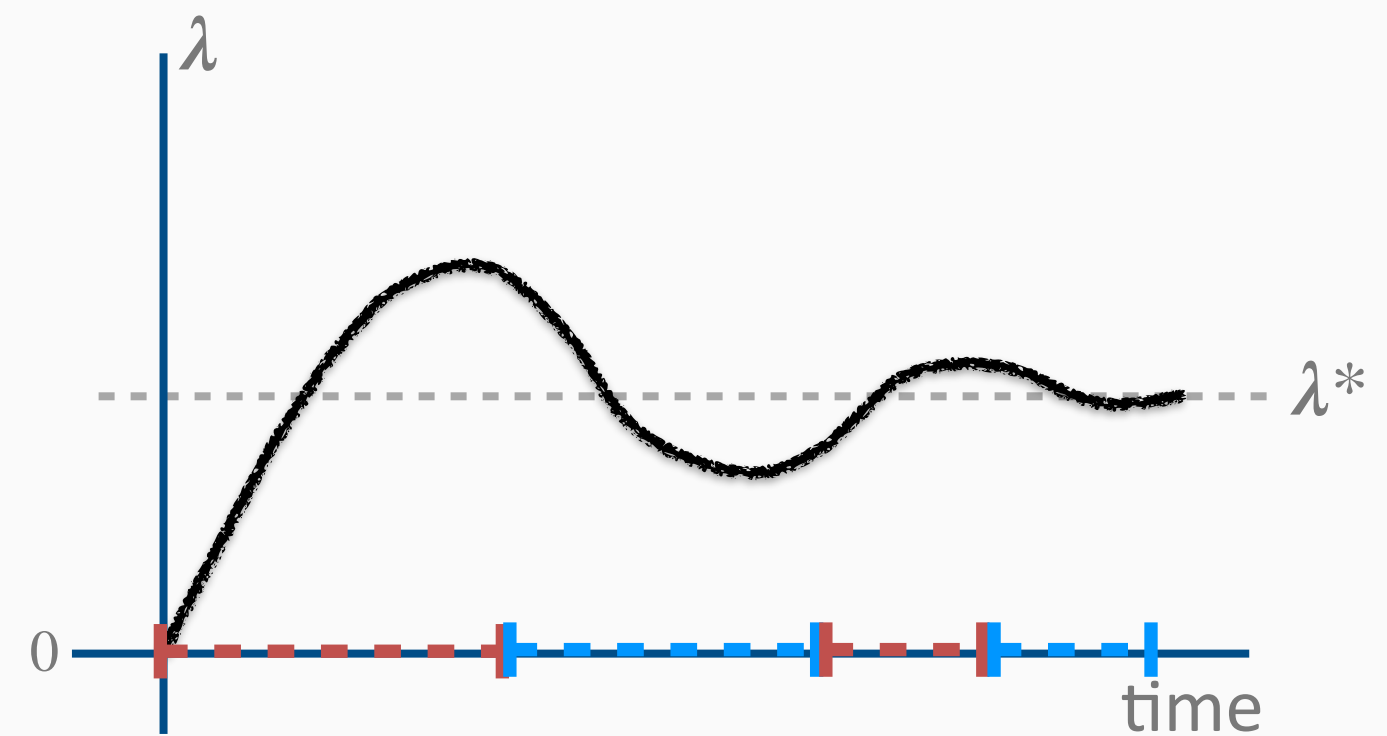
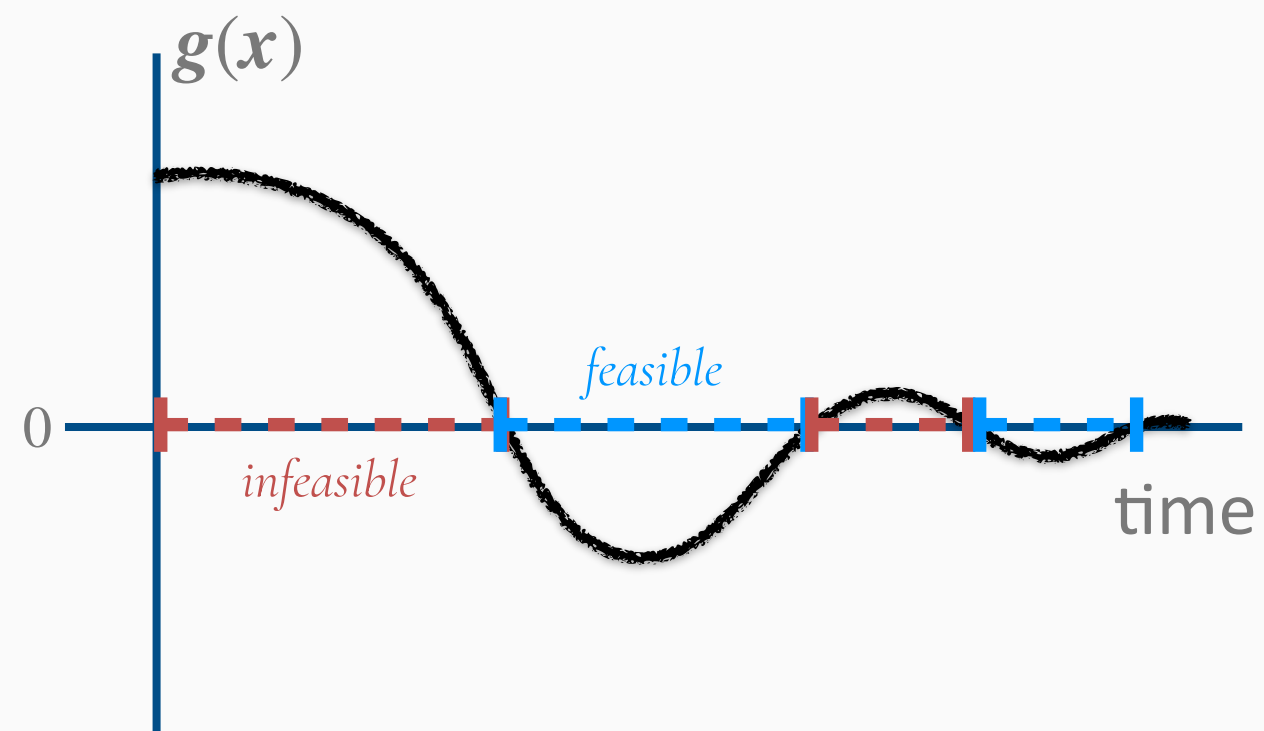
Compared to the penalized approach: only need to update value of the multipliers.



# Dynamics of GDA

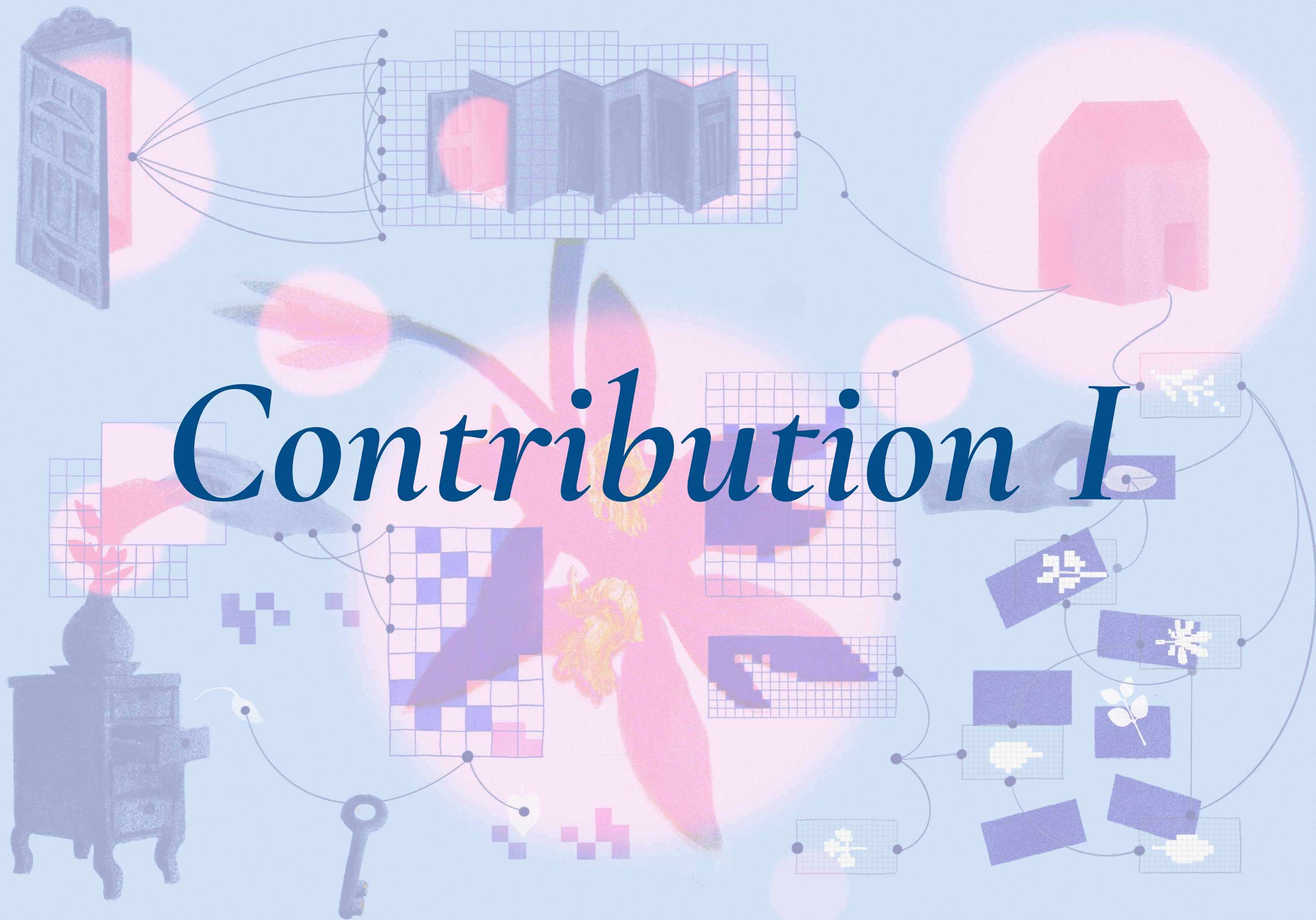
$$\lambda_{k+1} = [\lambda_k + \eta_{\text{dual}} \nabla_{\lambda} \mathcal{L}(\mathbf{x}_k, \lambda_k, \mu_k)]^+ = [\lambda_k + \eta_{\text{dual}} \mathbf{g}(\mathbf{x}_k)]^+$$

*constraint violation*



The multiplier accumulates the sequence of observed constraint violations.

# *Contribution I*





# Controlled Sparsity via Constrained Optimization

*How I Learned to Stop Tuning Penalties & Love Constraints*



Jose Gallego-Posada



Juan Ramirez



Akram Erraqabi



Yoshua Bengio



Simon Lacoste-Julien

NeurIPS 2022



*(We originally wanted to write a paper on  
constrained optimization.*

*The result was a “case study” on the use of  
constrained optimization for training  
sparse neural networks.)*

# Sparsity via $L_0$ regularization

---

$$\min_{\tilde{\theta}, \phi} \mathbb{E}_{z|\phi} \left[ L_{\mathcal{D}} \left( \tilde{\theta} \odot z \right) \right] + \lambda_{\text{pen}} \mathbb{E}_{z|\phi} \left[ \|z\|_0 \right]$$

*$L_0$ -“norm” penalty induces sparsity*

## $L_0$ reparametrization

Louizos et al. (2018) introduced a stochastic, differentiable reparametrization  $\theta = \tilde{\theta} \odot z$  for training sparse neural networks

*stochastic binary gates*

## Challenges with $\lambda_{\text{pen}}$

- ▶ Strength of the regularization is mediated by coefficient  $\lambda_{\text{pen}}$ .
- ▶ Tuning  $\lambda_{\text{pen}}$  to achieve a pre-determined sparsity level is expensive.



Instead of **penalizing**, formulate sparsity goals as  $L_0$ -norm constraints and **solve the Lagrangian min-max problem**

$$\min_{\tilde{\theta}, \phi} \mathbb{E}_{z|\phi} \left[ L_{\mathcal{D}} \left( \tilde{\theta} \odot z \right) \right] \quad \text{subject to} \quad \frac{\mathbb{E}_{z|\phi} [\|z\|_0]}{\#(\theta)} \leq \epsilon$$

- ✓ Interpretable hyperparameter semantics: target sparsity level
- ✓ Reliable control over the model sparsity

# Contributions

---

- ▶ Proposed a constrained approach for learning models with controllable levels of sparsity, highlighting its benefits with respect to the popular penalized approach
- ▶ Introduced a heuristic called "*dual restarts*" to avoid excessive sparsity caused by accumulation of constraint violations in the multipliers
- ▶ Through simple experimental adjustments, we **managed to successfully train sparse (Wide)ResNets** — prior experimental studies had failed at this!
- ▶ Demonstrated that **we can reliably achieve controllable sparsity levels** across many different architectures and datasets — without compromising performance

# Dual Restarts

---

When using GDA, the multipliers can be excessively large, even at a feasible primal iterate.

## Motivation of dual restarts as a “conditional” best response

The game-theoretic best response of the dual player to a primal action  $(\tilde{\theta}, \phi)$  is:

$$\lambda_{\text{CO}}^{\text{BR}}(\tilde{\theta}, \phi) = \underset{\lambda_{\text{CO}} \geq 0}{\operatorname{argmax}} \mathcal{L}(\tilde{\theta}, \phi, \lambda_{\text{CO}}) = \underset{\lambda_{\text{CO}} \geq 0}{\operatorname{argmax}} f(\tilde{\theta}, \phi) + \lambda_{\text{CO}}^{\text{T}} (g(\phi) - \epsilon)$$

This is a linear program whose solution depends purely on the feasibility of  $(\tilde{\theta}, \phi)$ :

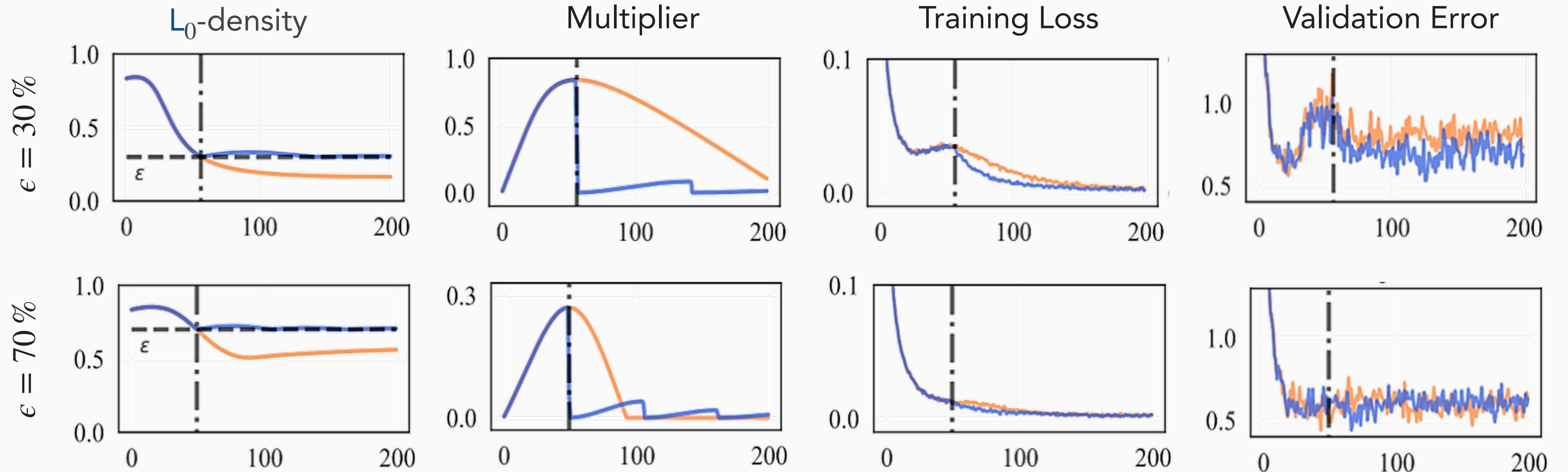
$$\lambda_{\text{CO}}^{\text{BR}}(\tilde{\theta}, \phi) = \begin{cases} \infty & \text{if } g(\phi) > \epsilon \\ \mathbb{R}^+ & \text{if } g(\phi) = \epsilon \\ 0 & \text{if } g(\phi) < \epsilon \end{cases}$$



# Training Dynamics

Dataset: MNIST

Model: LeNet5



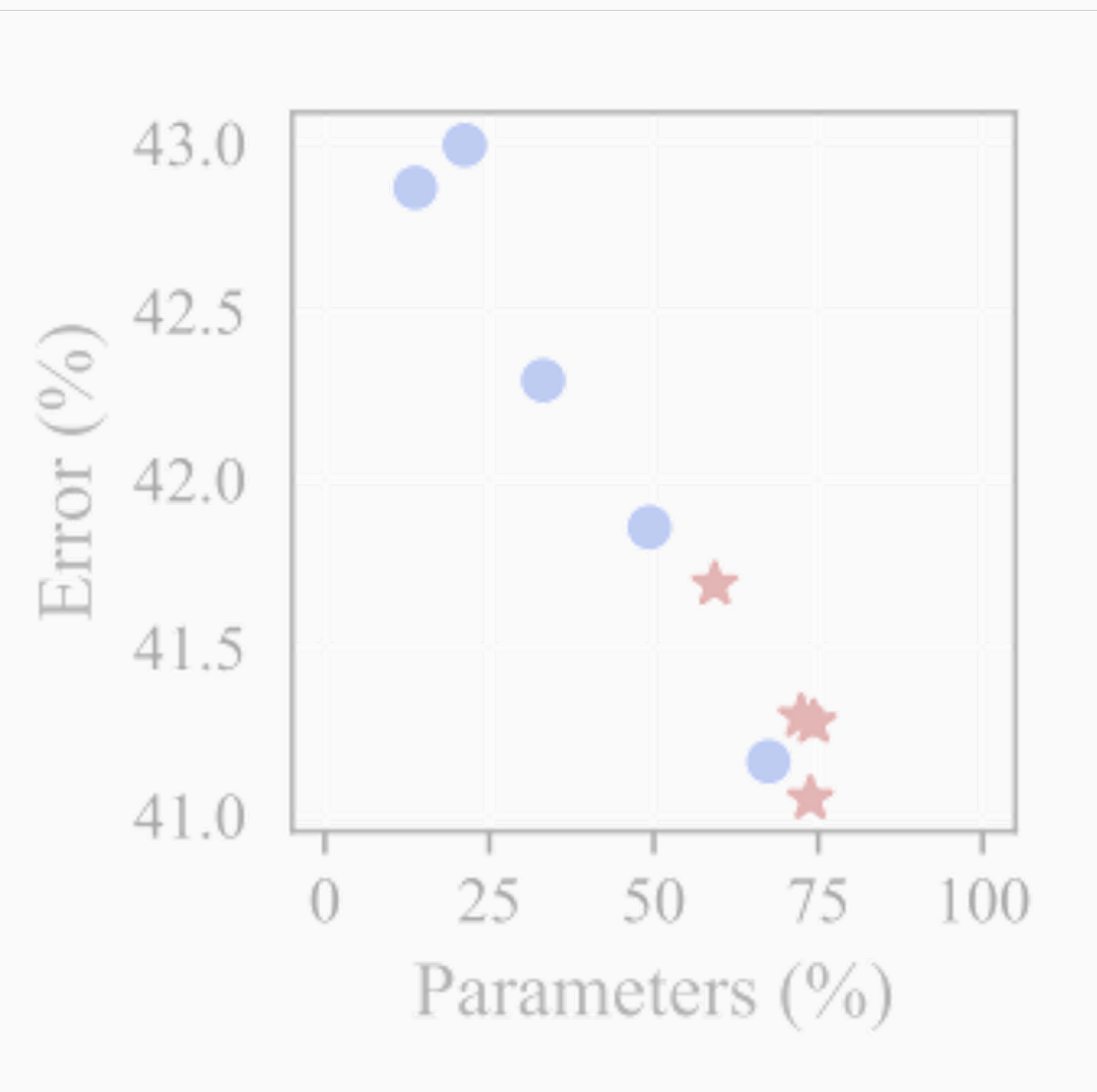
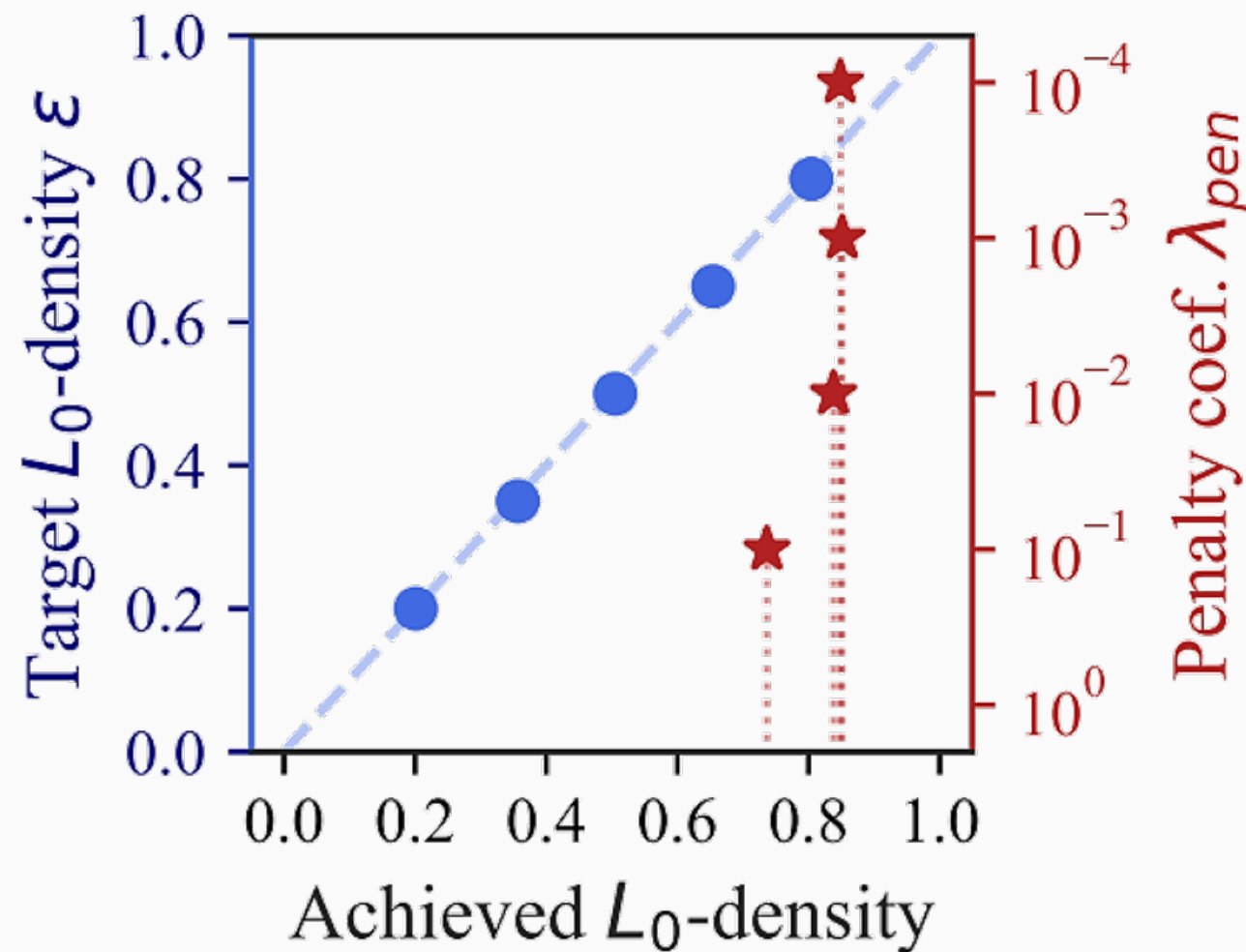
Dual Restarts

No Dual Restarts

First Feasibility

# Achieving controlled sparsity

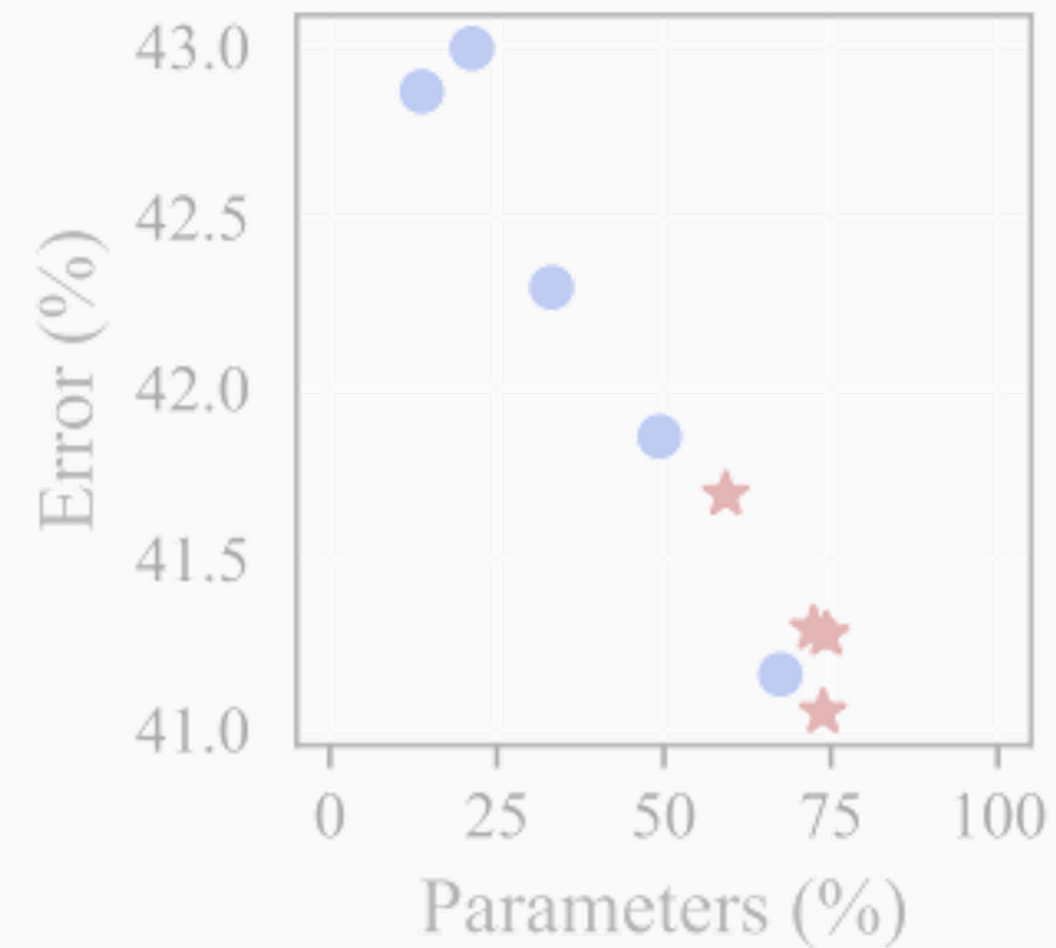
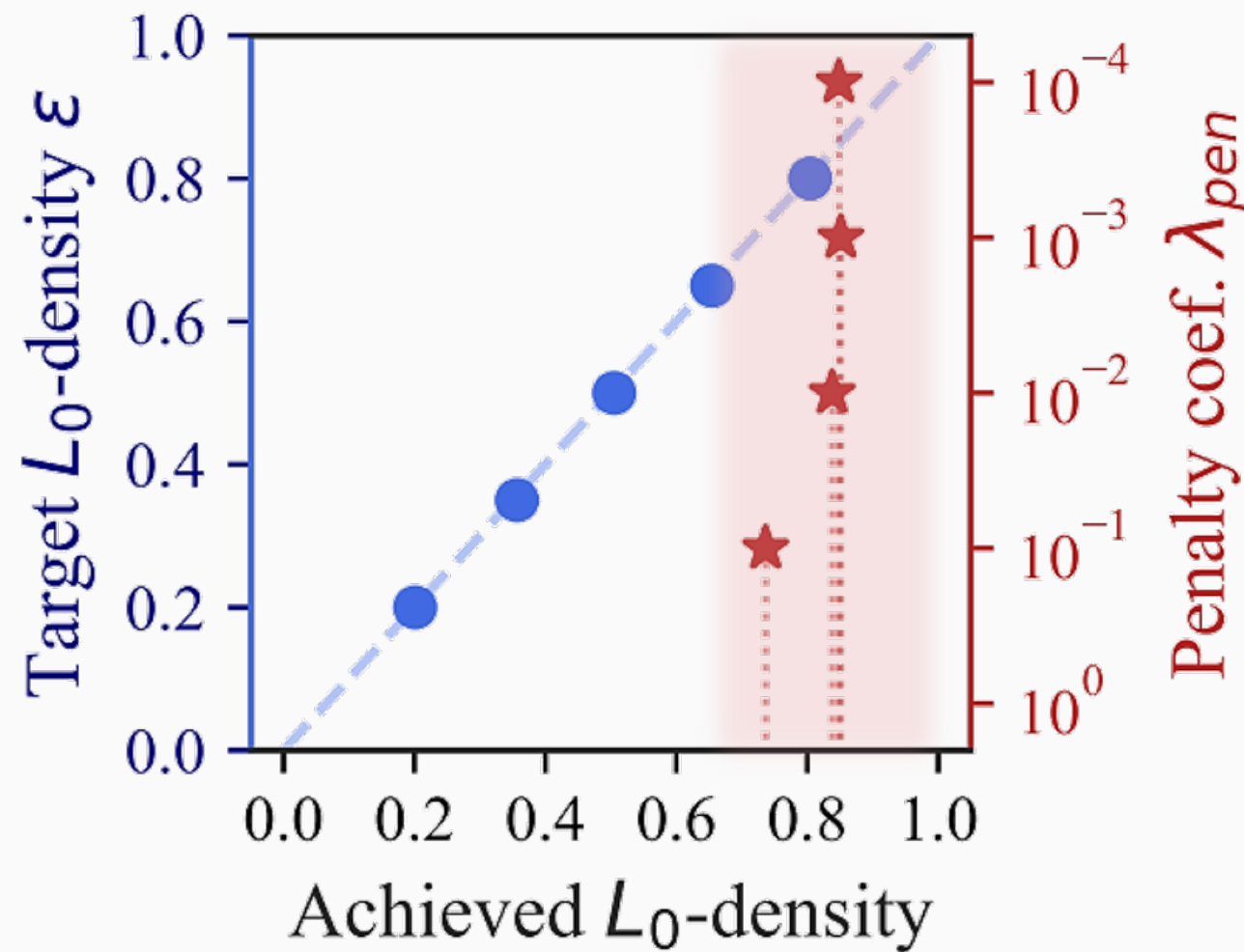
Dataset: Tiny-ImageNet Model: ResNet18



● Constrained ★ Penalized

# Achieving controlled sparsity

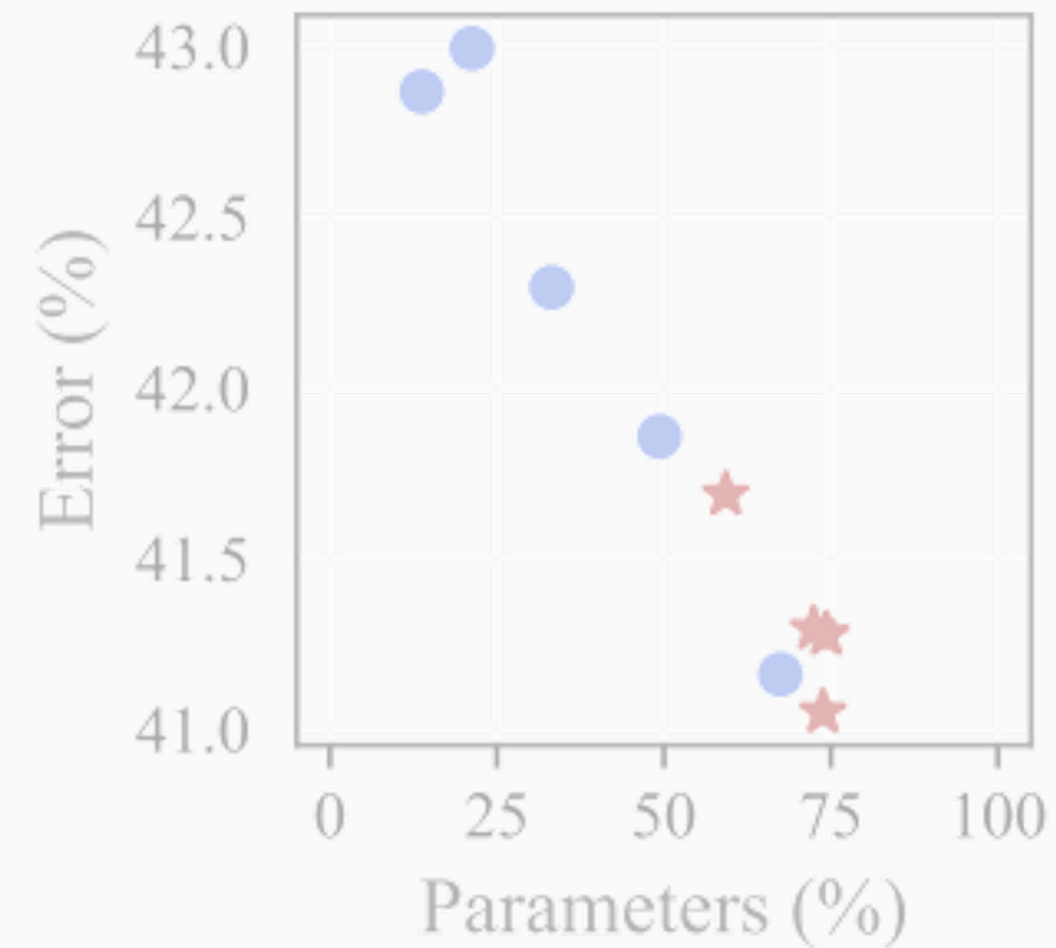
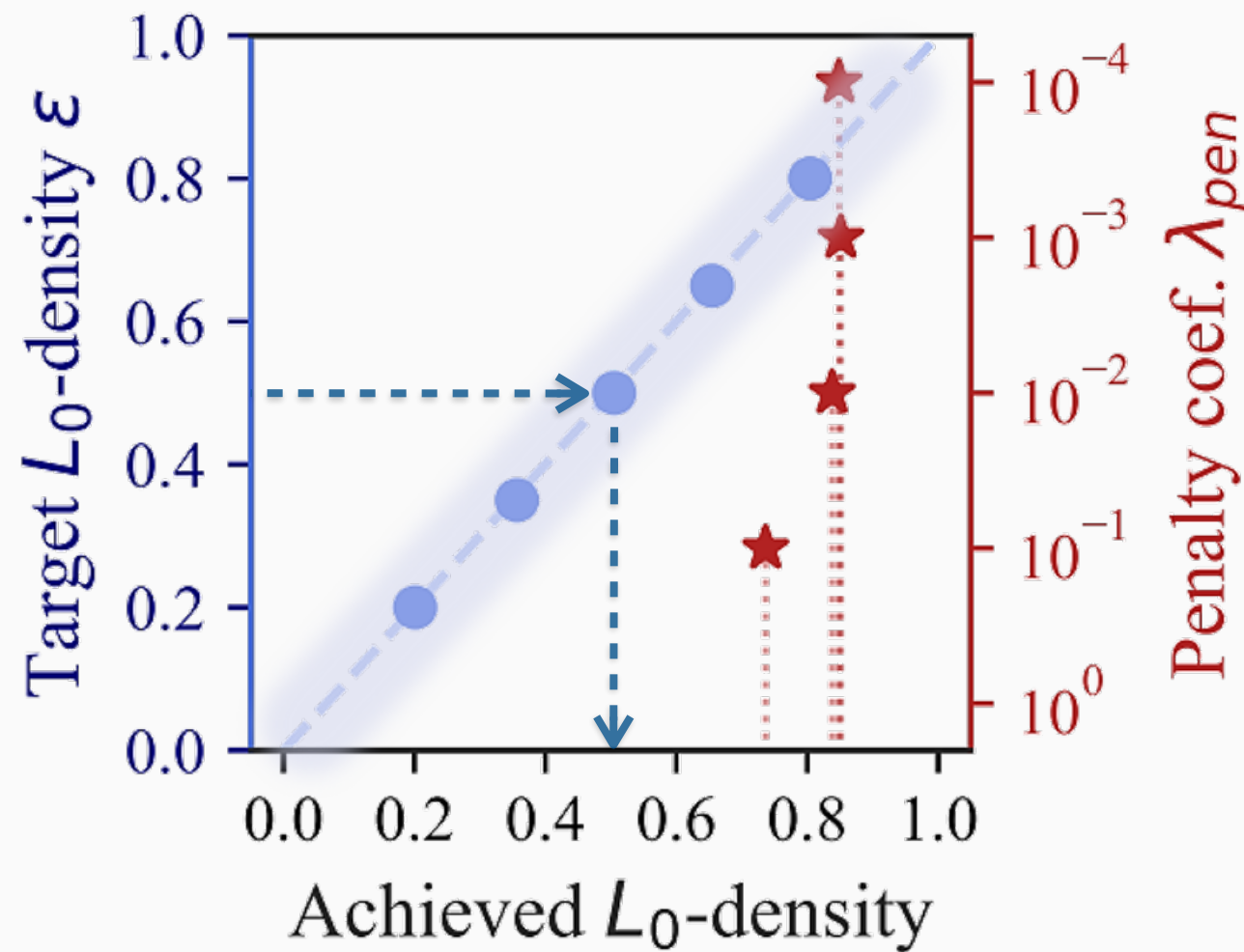
Dataset: Tiny-ImageNet Model: ResNet18



● Constrained ★ Penalized

# Achieving controlled sparsity

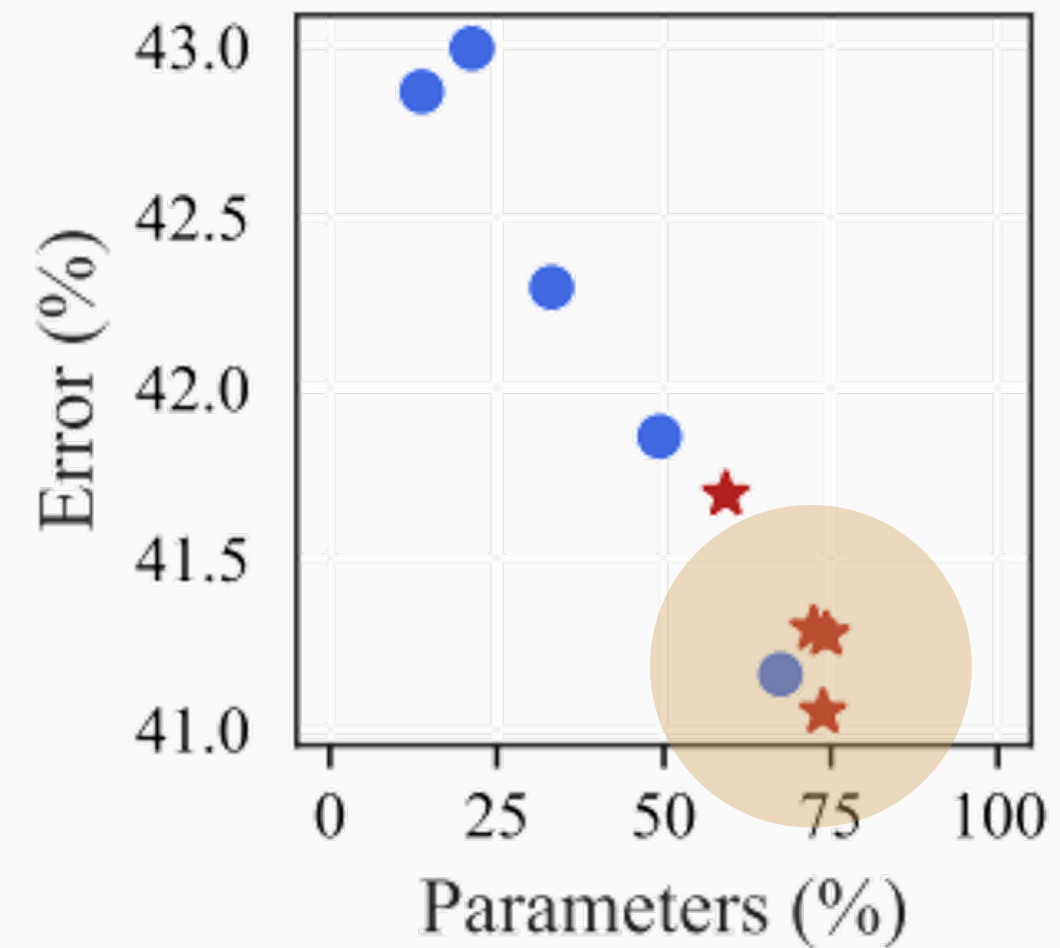
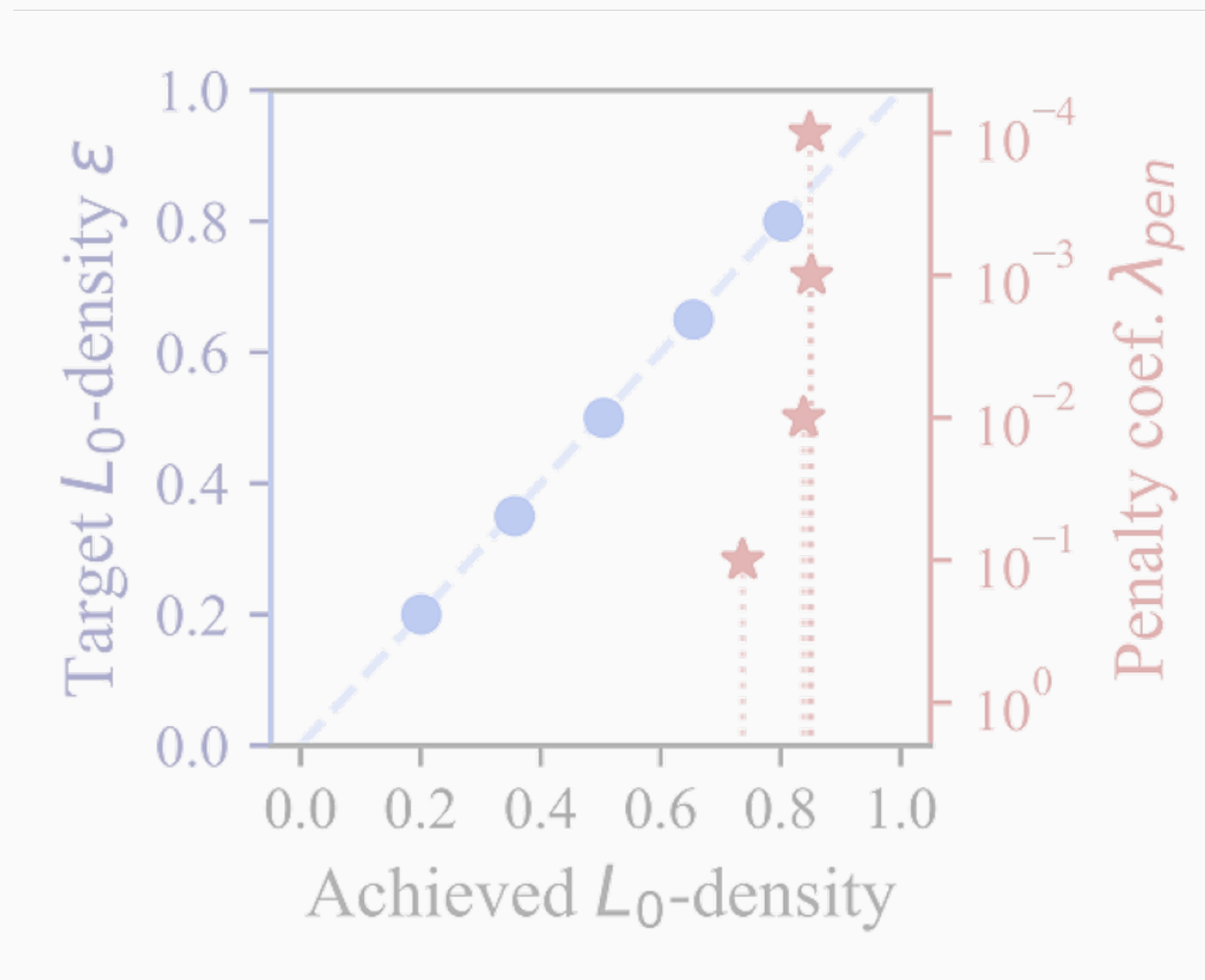
Dataset: Tiny-ImageNet Model: ResNet18



● Constrained ★ Penalized

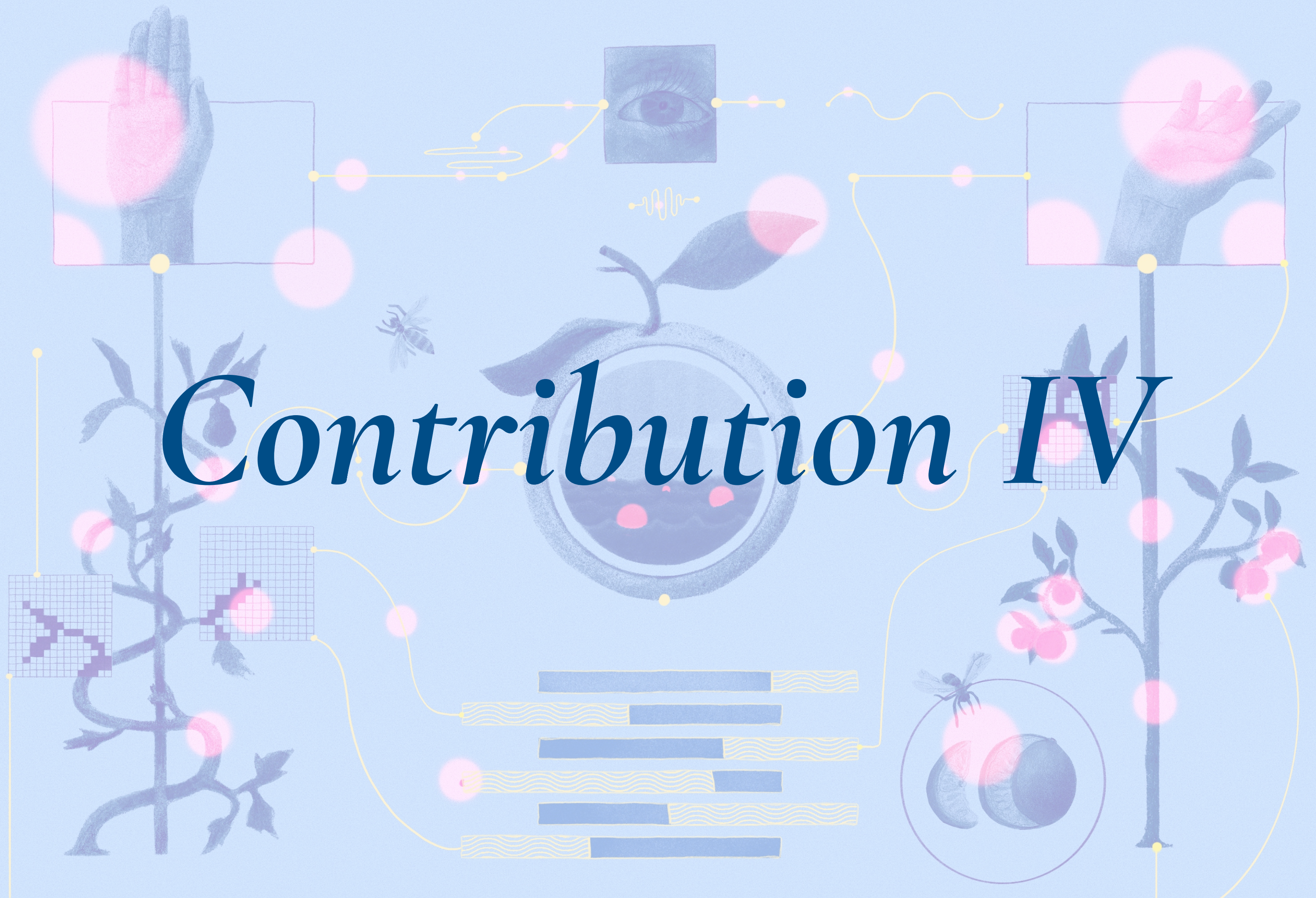
# ... while retaining performance!

Dataset: Tiny-ImageNet Model: ResNet18



● Constrained ★ Penalized

# Contribution IV





# On PI controllers for updating Lagrange multipliers in constrained optimization



Motahareh Sohrabi



Juan Ramirez



Tianyue H. Zhang



Simon Lacoste-Julien

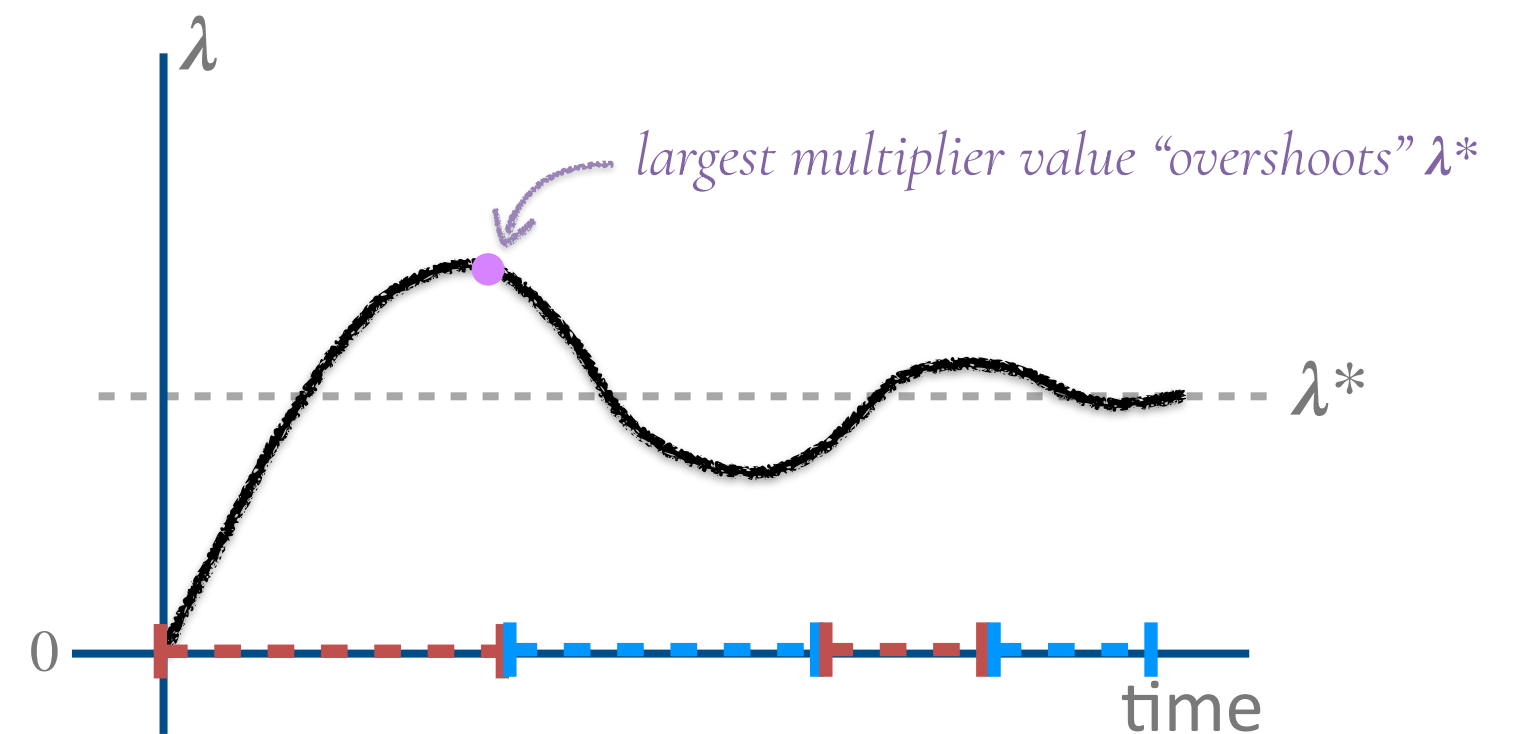
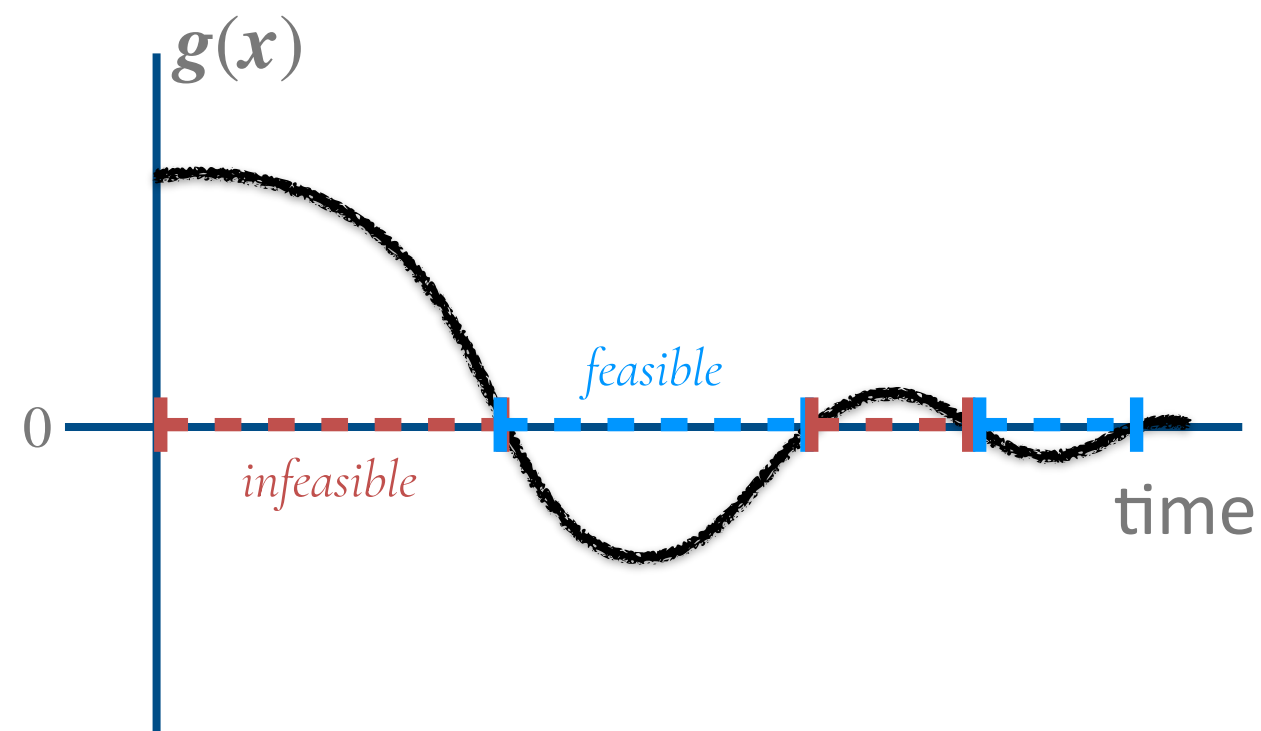


Jose Gallego-Posada

ICML 2024

# Dynamics of GDA

$$\lambda_{k+1} = [\lambda_k + \eta_{\text{dual}} \nabla_{\lambda} \mathcal{L}(\mathbf{x}_k, \lambda_k, \boldsymbol{\mu}_k)]^+ = [\lambda_k + \eta_{\text{dual}} \mathbf{g}(\mathbf{x}_k)]^+$$



The multiplier accumulates/*integrates* the sequence of observed constraint violations



# What we are looking for

---

## Shortcomings of GDA

- GDA may result in overshoot and oscillations (Gidel et al. 2019; Stooke et al. 2020)
- Especially problematic in safety-related applications

## Goal and scope

- **Reliable and robust** approach for solving Lagrangian optimization problems
- That **does not modify** training “recipe” for primal variables

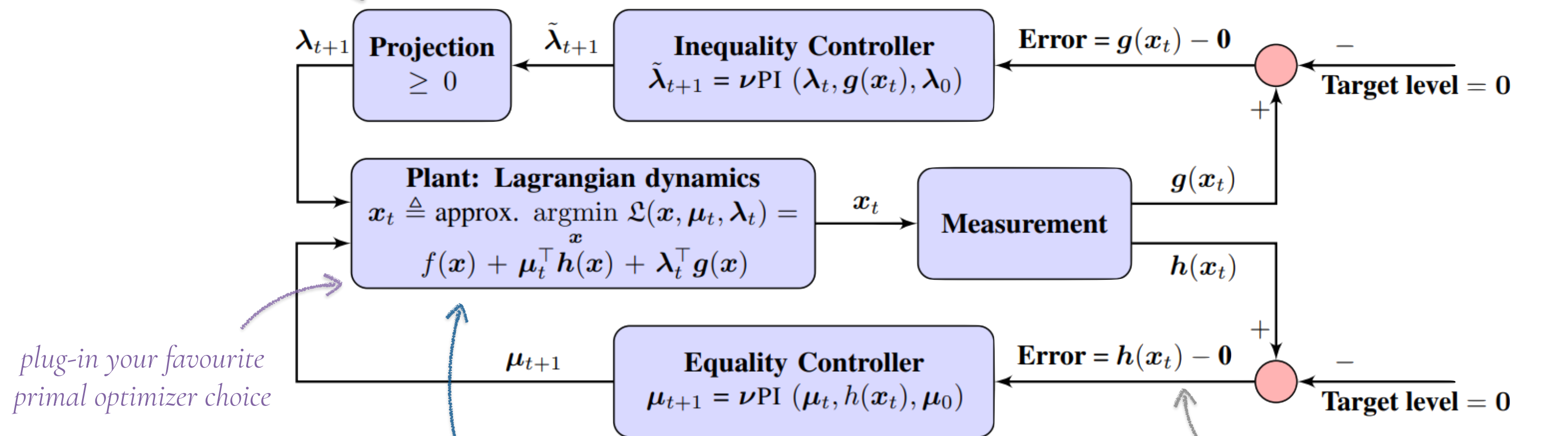
**Achieving this goal enables wider adoption of Lagrangian optimization in deep learning!**



# Dynamical system's view of CO

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to } \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \text{ and } \mathbf{h}(\mathbf{x}) = \mathbf{0}$$

ensure non-negativity of inequality multipliers



plug-in your favourite primal optimizer choice

multipliers "tilt" the primal gradient

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla g_i(\mathbf{x}^*) + \sum_{i=1}^n \mu_i^* \nabla h_i(\mathbf{x}^*)$$

constraint violation is the error signal for the controller

# $\nu$ PI control for constrained optimization

## Algorithm: $\nu$ PI update on parameter $\theta$

**Args:** EMA coefficient  $\nu$ , proportional ( $\kappa_p$ ) and integral ( $\kappa_i$ ) gains; initial conditions  $\xi_0$  and  $\theta_0$

1. Measure the current system error  $e_t$

2.  $\xi_t \leftarrow \nu \xi_{t-1} + (1 - \nu)e_t$  (for  $t \geq 1$ )

3.  $\theta_{t+1} \leftarrow \theta_0 + \kappa_p \xi_t + \kappa_i \sum_{\tau=0}^t e_\tau$

Recursively,  $\theta_1 \leftarrow \theta_0 + \kappa_p \xi_0 + \kappa_i e_0$

$$\theta_{t+1} \leftarrow \theta_t + \kappa_i e_t + \kappa_p (\xi_t - \xi_{t-1})$$

General case

like  $\nabla$ -ascent

new term looks at  
change in constraint  
satisfaction!

$$\theta_{t+1} \leftarrow \theta_t + \kappa_i e_t + \kappa_p (e_t - e_{t-1})$$

Case  $\nu = 0$

# $\nu$ PI generalizes momentum methods

## Theorem

*Polyak  $\gamma = 0$ ; Nesterov  $\gamma = 1$*

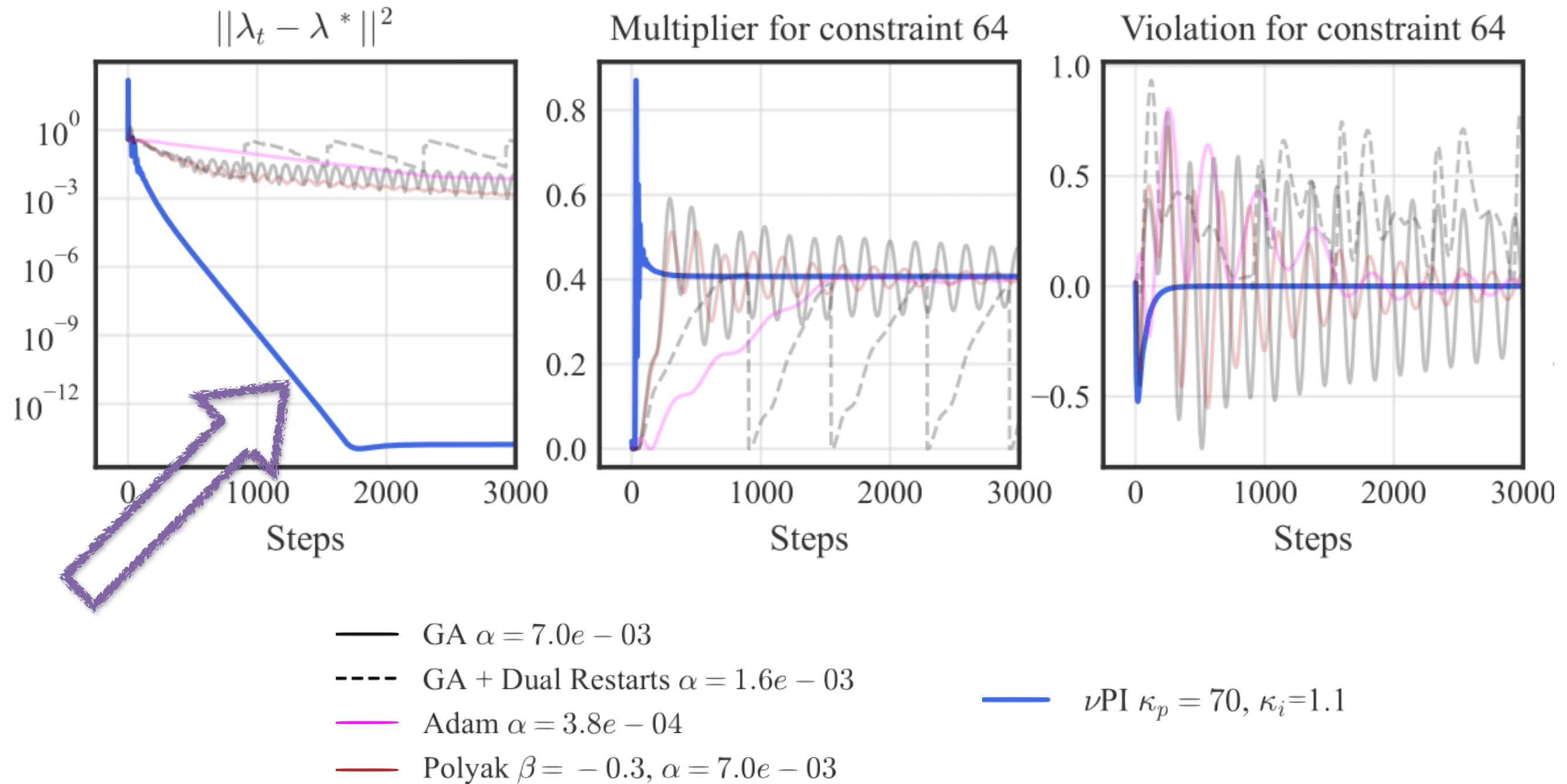
Under the same initialization  $\theta_0$ , UnifiedMomentum( $\alpha, \beta \neq 1, \gamma$ ) is a special case of the  $\nu$ PI algorithm with the hyperparameter choices:

$$\begin{aligned} \nu &\leftarrow \beta & \xi_0 &\leftarrow (1 - \beta)e_0 \\ \kappa_i &\leftarrow \frac{\alpha}{1 - \beta} & \kappa_p &\leftarrow -\frac{\alpha\beta}{(1 - \beta)^2}[1 - \gamma(1 - \beta)] \end{aligned}$$

# $\nu$ PI generalizes momentum methods

Algorithm	$\xi_0$	$\kappa_p$	$\kappa_i$	$\nu$
UNIFIEDMOMENTUM( $\alpha, \beta, \gamma$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta}{(1 - \beta)^2} [1 - \gamma(1 - \beta)]$	$\frac{\alpha}{1 - \beta}$	$\beta$
POLYAK( $\alpha, \beta$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta}{(1 - \beta)^2}$	$\frac{\alpha}{1 - \beta}$	$\beta$
NESTEROV( $\alpha, \beta$ )	$(1 - \beta)\mathbf{e}_0$	$-\frac{\alpha\beta^2}{(1 - \beta)^2}$	$\frac{\alpha}{1 - \beta}$	$\beta$
PI	$\mathbf{e}_0$	$\kappa_p$	$\kappa_i$	0
OPTIMISTICGRADIENTASCENT( $\alpha$ )	$\mathbf{e}_0$	$\alpha$	$\alpha$	0
$\nu$ PI ( $\kappa_i, \kappa_p, \nu$ ) in practice	<b>0</b>	$\kappa_i$	$\kappa_p$	$\nu$
GRADIENTASCENT( $\alpha$ )	–	0	$\alpha$	0

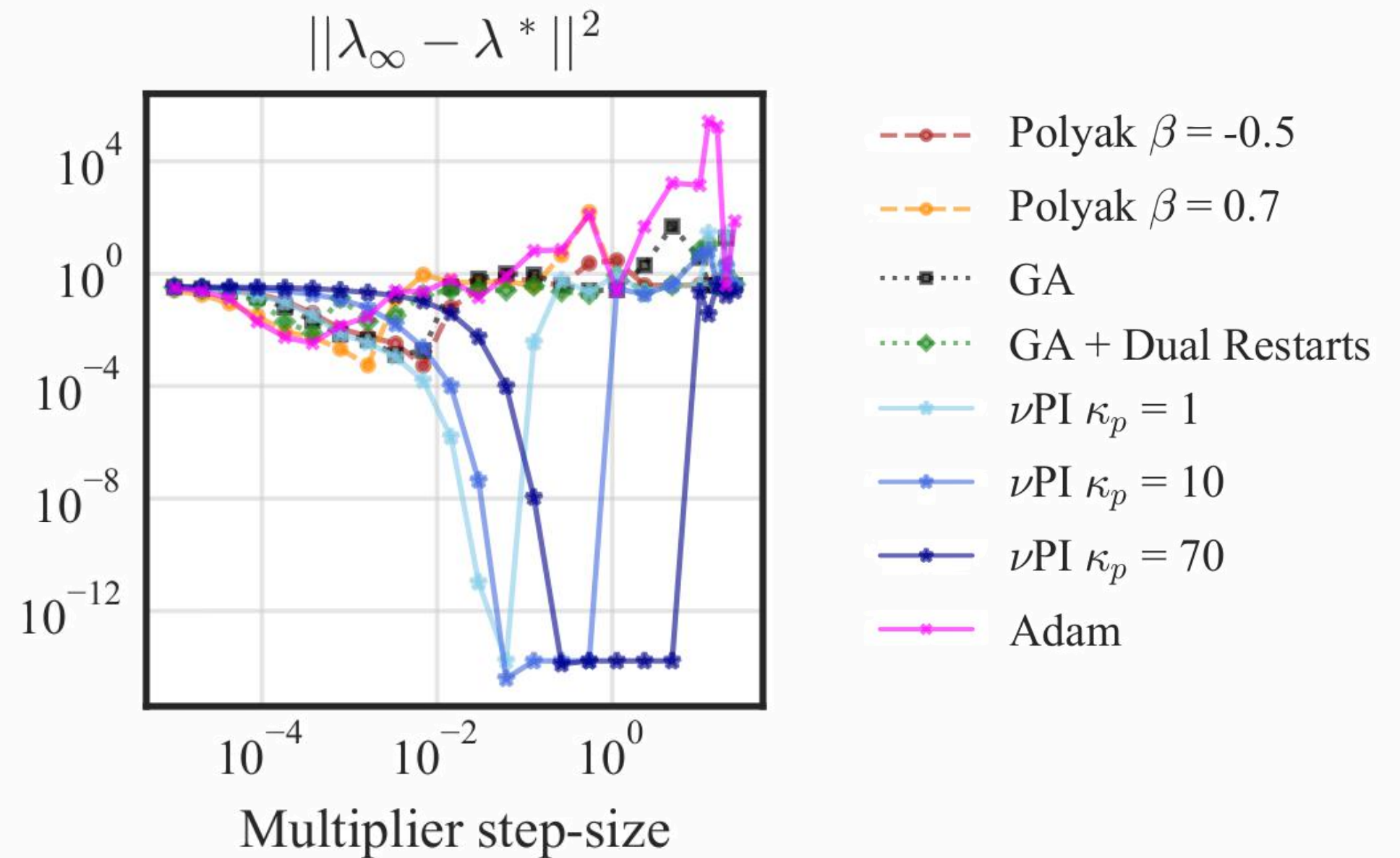
Of all attempted optimizers\*, **only  $\nu$ PI converged successfully to the true solution!**

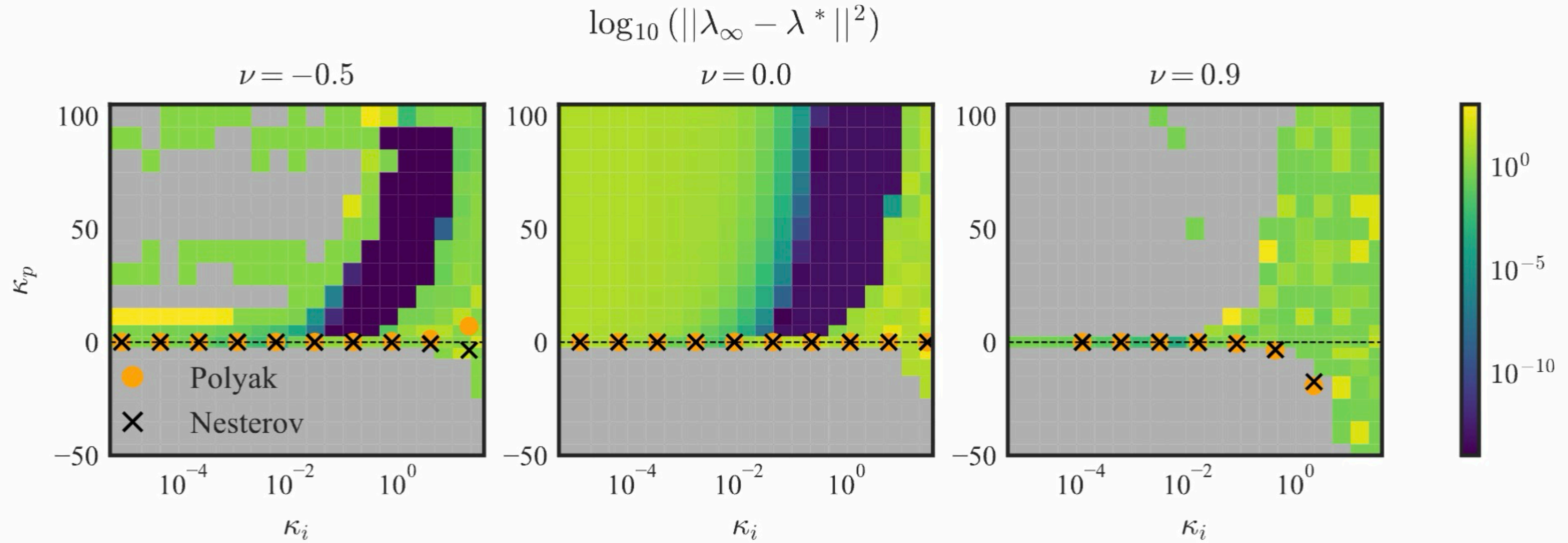


\*Showing best hyperparameters for each optimizer after grid-search aiming to minimize the distance to  $\lambda^*$  after 5.000 iterations

# Robustness

Higher values of  $\kappa_p$  allow for choosing **larger values of  $\kappa_i$**  (multiplier step-size) and **over a wider range**, while still achieving convergence.





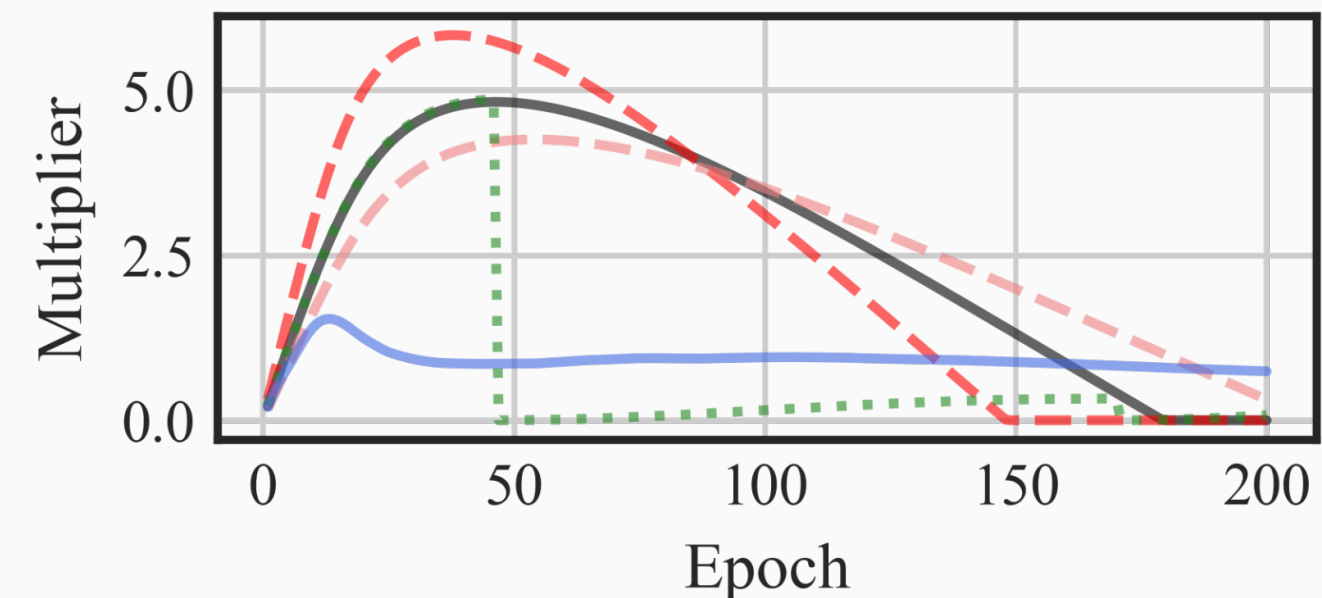
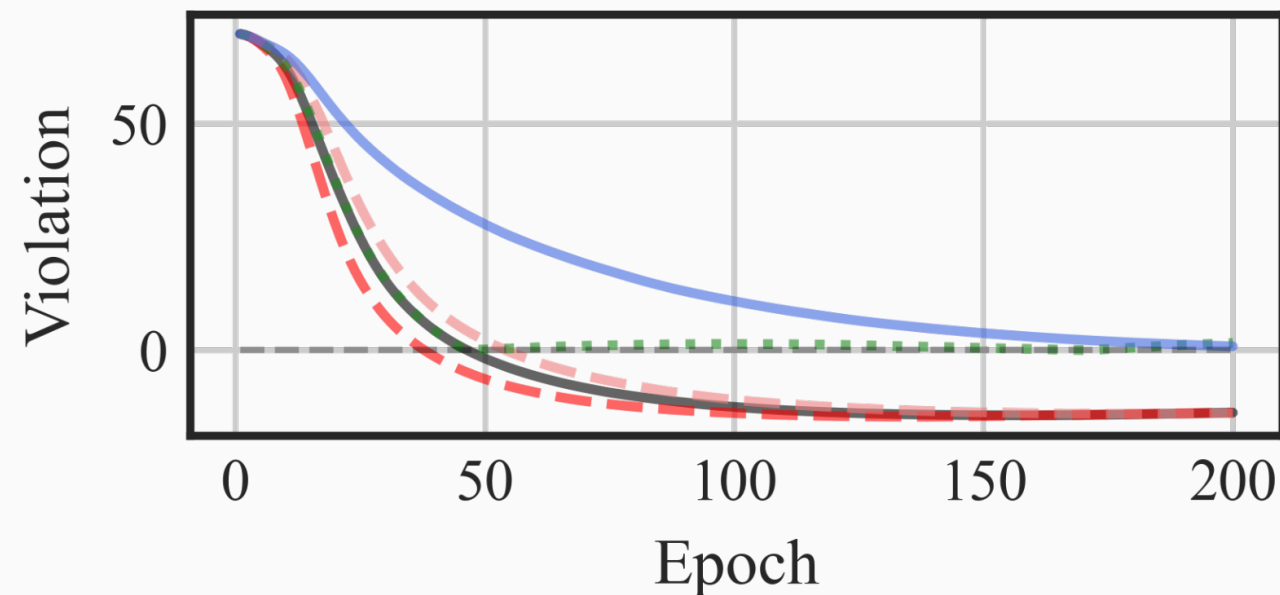
$\nu$ PI provides additional flexibility compared to Polyak and Nesterov which is crucial for achieving convergence in this task.



# Revisiting $L_0$ -constrained problem

$$\min_{\tilde{\theta}, \phi} \mathbb{E}_{z|\phi} \left[ L_{\mathcal{D}}(\tilde{\theta} \odot z) \right] \quad \text{subject to} \quad \frac{\mathbb{E}_{z|\phi} [\|z\|_0]}{\#(\theta)} \leq \epsilon$$

$\nu$ PI delivers **stable multiplier dynamics without constraint overshoot**

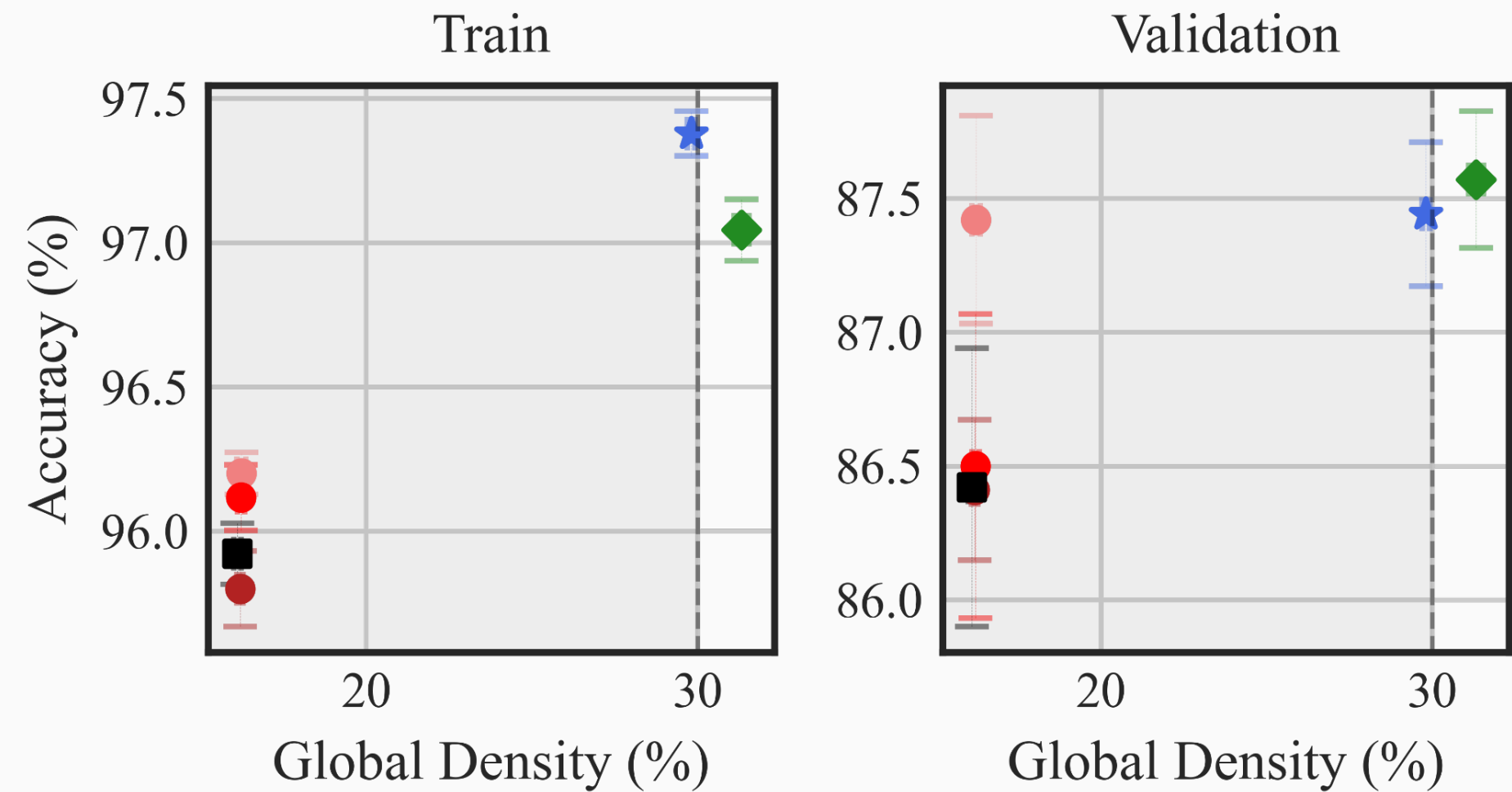


— GA    - - Polyak  $\beta = 0.3$     - - Polyak  $\beta = -0.3$     ··· GA + Dual Restarts    —  $\nu$ PI  $\kappa_p = 16.0$

# Impact on performance

$\nu$ PI achieves high accuracy and tightly respects the constraints, without overshooting

- Polyak  $\beta = -0.5$
- Polyak  $\beta = -0.3$
- Polyak  $\beta = 0.3$
- GA
- ◆ GA + Dual Restarts
- ★  $\nu$ PI  $\kappa_p = 14.4$



# *Contribution V*





# Cooper: A Library for Constrained Optimization in Deep Learning



Jose Gallego-Posada



Juan Ramirez



Meraj Hashemizadeh



Simon Lacoste-Julien

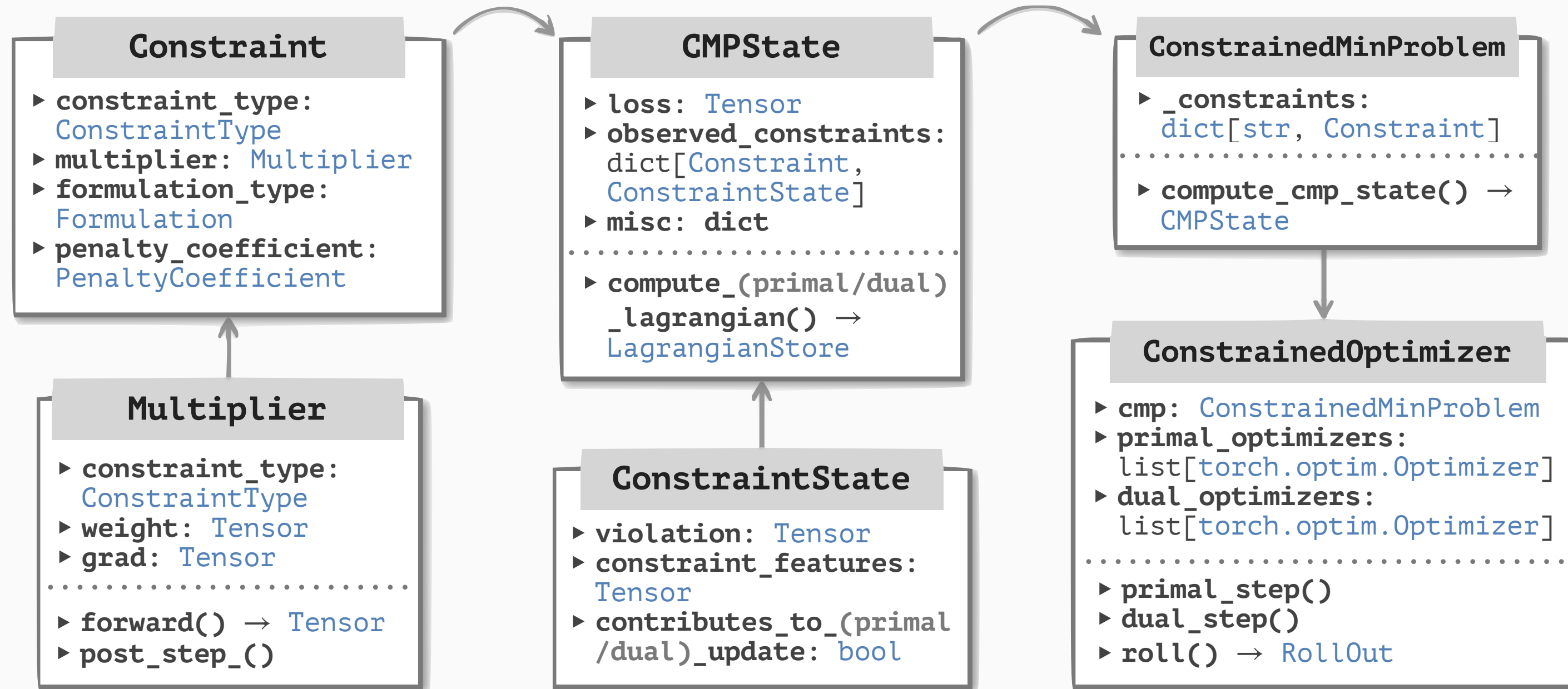
JMLR MLOSS 2024 (*under submission*)



# Cooper

*a general-purpose, deep learning-first  
library for constrained optimization,  
built on PyTorch.*

# Cooper's class overview



# Conclusions & Perspectives



(input)

(output)

(false)

(false)

(false)

<leafhead>

(false)

(false)

(false)

(yes)

(it)

(false)

(false)

(false)

(false)



# Constrained optimization is an up-and-coming research direction

- ▶ As ML becomes a “technology”, ensuring compliance with government regulations and industry standards is crucial next-step
- ▶ Constrained optimization is a rich field, ripe for integration by ML community
- ▶ Socially impactful research; inter-disciplinary relevance





# Main challenges when solving constrained problems in machine learning

- ▶ Optimization dynamics
- ▶ Non-differentiable constraints (“*proxy-constraints*” from Cotter et al. (2019))
- ▶ Generalization properties for loss **and** constraints
- ▶ Feasibility: always? if not, how fast?



# (Some) Open questions

- ▶ How to deal with constraints that are difficult to quantify?
- ▶ Why do GDA-like schemes work in practical Lagrangian problems?
- ▶ What is the role of overparametrization in constrained optimization?
- ▶ How can we improve the Lagrange multipliers further?
- ▶ How can we make constrained techniques usable “during inference”?
- ▶ What is next for Cooper?

A group of 12 people, including men and women of various ethnicities, are posing for a group photo on a wooden deck. They are dressed in winter clothing such as jackets, scarves, and hats. The background features a dense forest of evergreen trees and a waterfall cascading over rocks. The overall atmosphere is bright and cheerful.

*Thank you!*