# On PI controllers for updating Lagrange multipliers in constrained optimization

Jose Gallego-Posada

March 2024

gallego-posada.github.io

# Today's agenda

‣ Constrained optimization

‣ Dynamics of gradient descent-ascent

‣ The $\nu$PI controller

‣ Applications of $\nu$PI in constrained optimization

REPUBLIQUE FRANÇAISE

POSTES

8F

1736 LAGRANGE 1813

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

1

Collaborators

Motahareh Sohrabi*

Juan Ramirez*

Helen Zhang

Simon Lacoste-Julien

2

# Constrained optimization

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

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

**Feasible set**
$$\mathcal{X} = \left\{ \boldsymbol{x} \in \mathbb{R}^d \,\middle|\, \boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0} \text{ and } \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0} \right\}$$

**Optimality condition** (necessary)
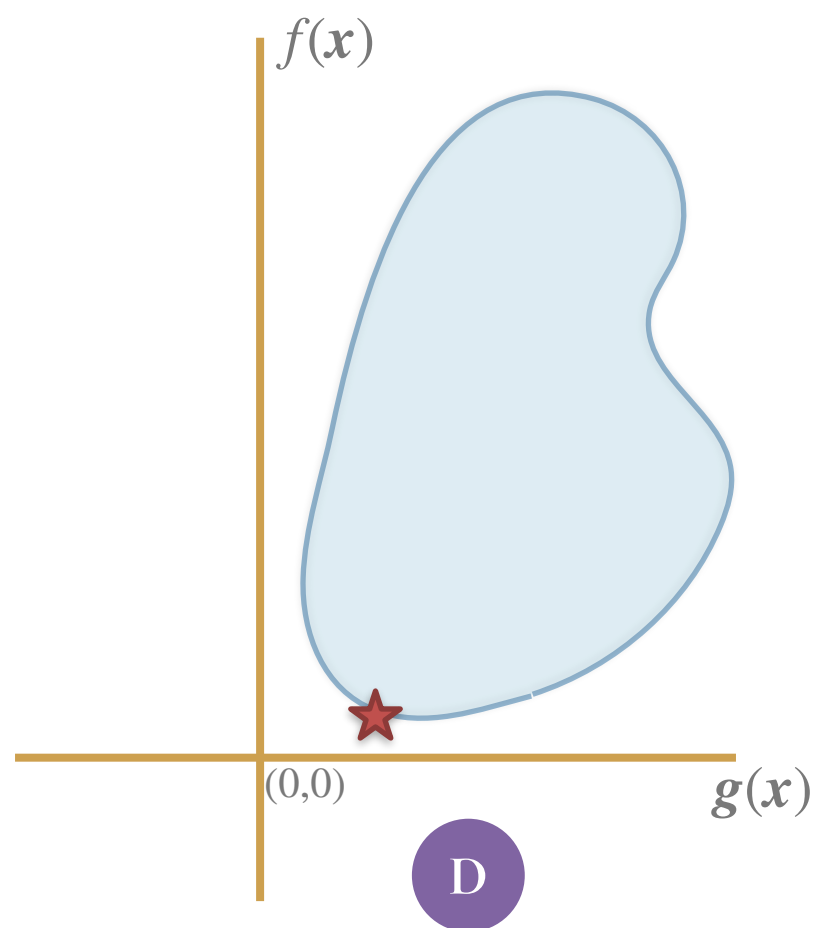If $\boldsymbol{x}*$ is a local minimum of $f$ over $\mathcal{X}$, then
$$\nabla f(\boldsymbol{x}*)^\top \boldsymbol{z} \geq 0 \;\; \forall \boldsymbol{z} \in \mathcal{F}(\boldsymbol{x}*)$$

*feasible directions at* $\boldsymbol{x}*$

$f(\boldsymbol{x})$

*realizable configurations*
$\left\{ (\boldsymbol{g}(\boldsymbol{x}), f(\boldsymbol{x})) \mid \boldsymbol{x} \in \mathbb{R}^d \right\}$

*constrained*
*optimum*

(0,0)

$\boldsymbol{g}(\boldsymbol{x})$

*unconstrained*
*optimum*

3

# Lagrangian problem

$$\min_{x} f(x)$$

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

$$\Longleftrightarrow$$

**Lagrangian**

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

*"Lagrange multipliers" or "dual variables"*

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

$$\nabla f(x^*) + \sum_{i=1}^{m} \lambda_i^* \nabla g_i(x^*) + \sum_{i=1}^{n} \mu_i^* \nabla h_i(x^*) = 0$$

**Algorithmic approach**

Saddle points of the Lagrangian correspond to constrained optima. Find them!

# Gradient Descent-Ascent (GDA)

**Lagrangian** $\quad \min\limits_{x} \max\limits_{\lambda \geq \mathbf{0}, \boldsymbol{\mu}} \mathfrak{L}(\boldsymbol{x}, \lambda, \boldsymbol{\mu}) \triangleq f(\boldsymbol{x}) + \lambda^\top \boldsymbol{g}(\boldsymbol{x}) + \boldsymbol{\mu}^\top \boldsymbol{h}(\boldsymbol{x})$

**Algorithm**

**Initialize** $\boldsymbol{x}_0$, $\lambda_0 = \mathbf{0}$ and $\boldsymbol{\mu}_0 = \mathbf{0}$

**Repeat**

$$\boldsymbol{\mu}_{k+1} \leftarrow \boldsymbol{\mu}_k + \alpha_d \nabla_{\boldsymbol{\mu}} \mathfrak{L}(\boldsymbol{x}_k, \boldsymbol{\mu}_k, \boldsymbol{\mu}_k) = \boldsymbol{\mu}_k + \alpha \boldsymbol{h}(\boldsymbol{x}_k)$$

$$\lambda_{k+1} \leftarrow \left[\lambda_k + \alpha_d \nabla_\lambda \mathfrak{L}(\boldsymbol{x}_k, \lambda_k, \boldsymbol{\mu}_k)\right]^+ = \left[\lambda_k + \alpha \boldsymbol{g}(\boldsymbol{x}_k)\right]^+$$

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

$$\boldsymbol{x}_{k+1} \leftarrow \boldsymbol{x}_k - \alpha_p \nabla_{\boldsymbol{x}} \mathfrak{L}(\boldsymbol{x}_k, \lambda_k, \boldsymbol{\mu}_k)$$

**If** convergence check satisfied; **stop**

# Dynamics of GDA

$$\lambda_{k+1} = \left[ \lambda_k + \alpha_d \nabla_\lambda \mathfrak{L}(x_k, \lambda_k, \mu_k) \right]^+ = \left[ \lambda_k + \alpha g(x_k) \right]^+$$
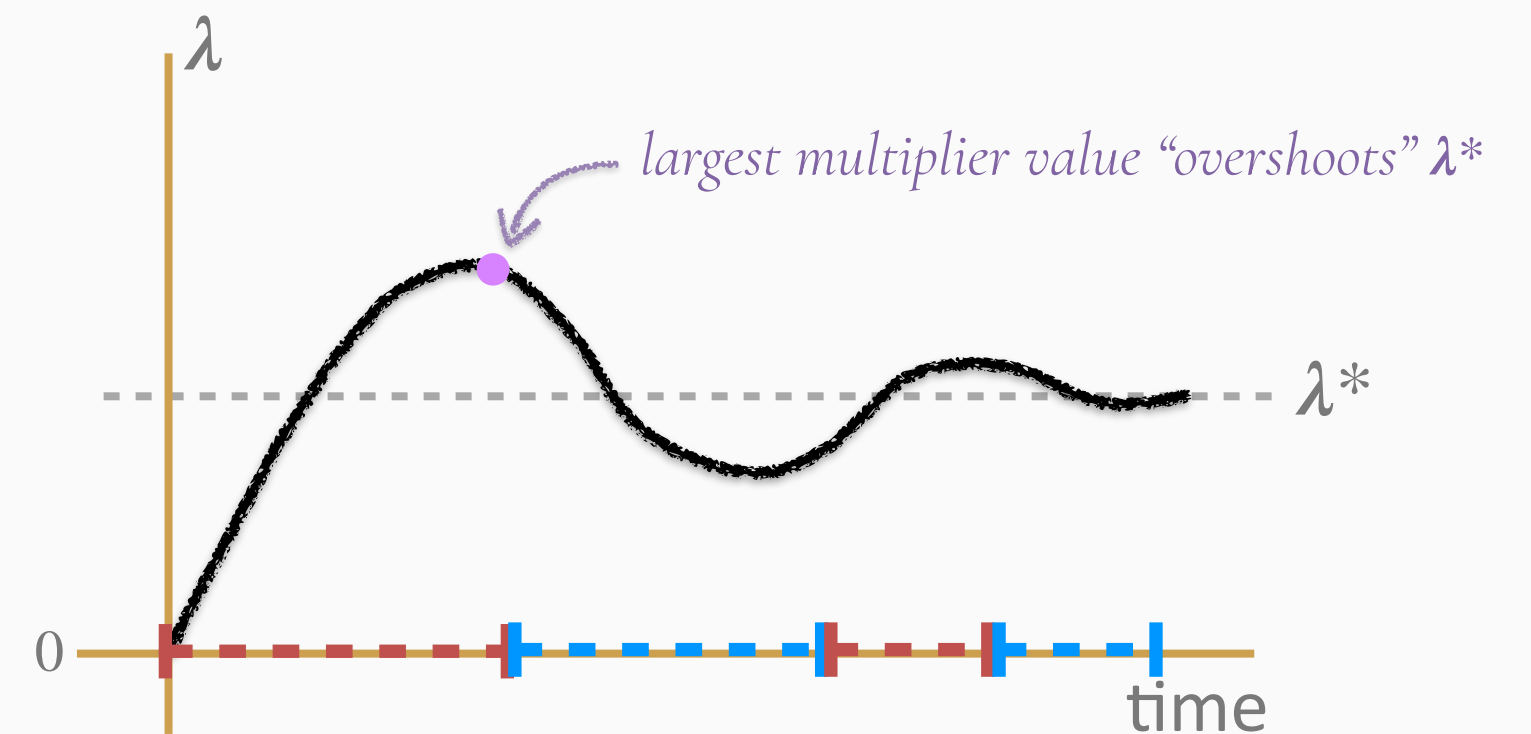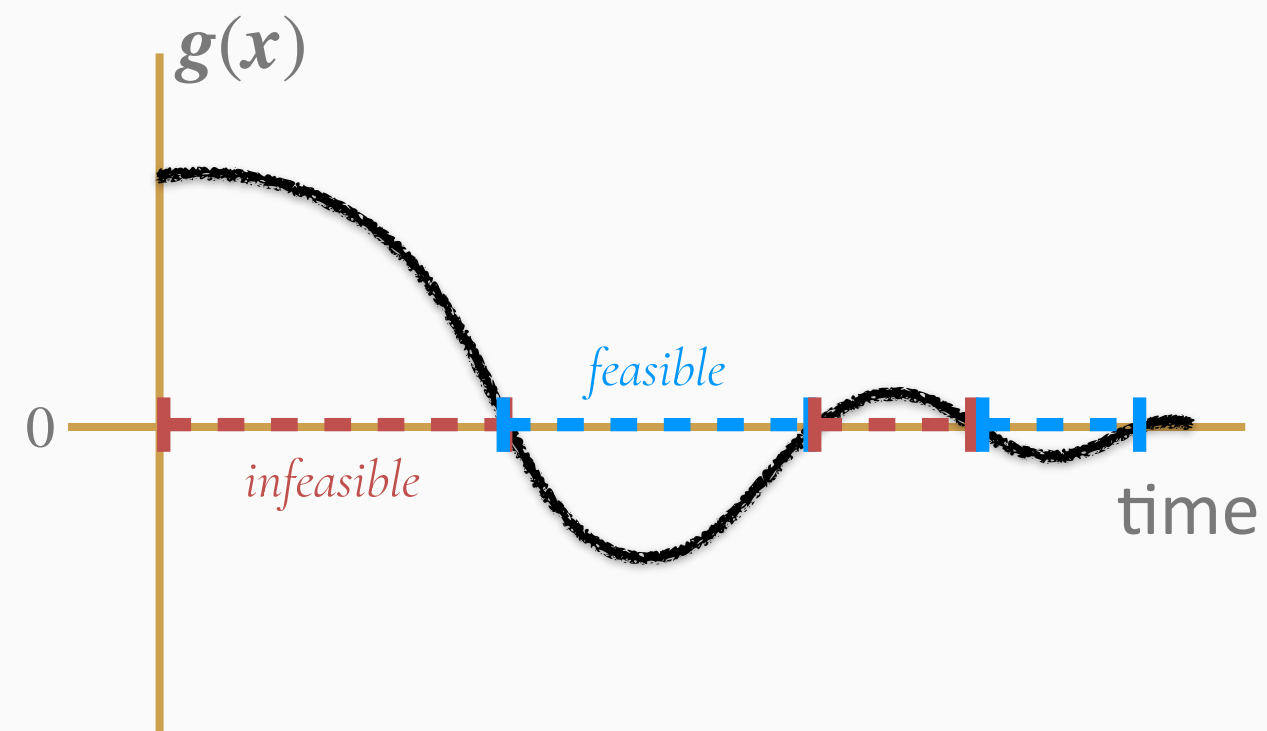


The multiplier accumulates/*integrates* the sequence of observed constraint violations

# What we are looking for

## Shortcomings of GDA

‣ GDA may result in overshoot and oscillations (Gidel et al. 2019; Stooke at 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!**

Gidel, G., Askari, R., Pezeshki, M., LePriol, R., Huang, G., Lacoste-Julien, S., and Mitliagkas, I. *Negative Momentum for Improved Game Dynamics.* In AISTATS, 2019.
Stooke, A., Achiam, J., and Abbeel, P. *Responsive Safety in Reinforcement Learning by PID Lagrangian Methods.* In ICML, 2020.
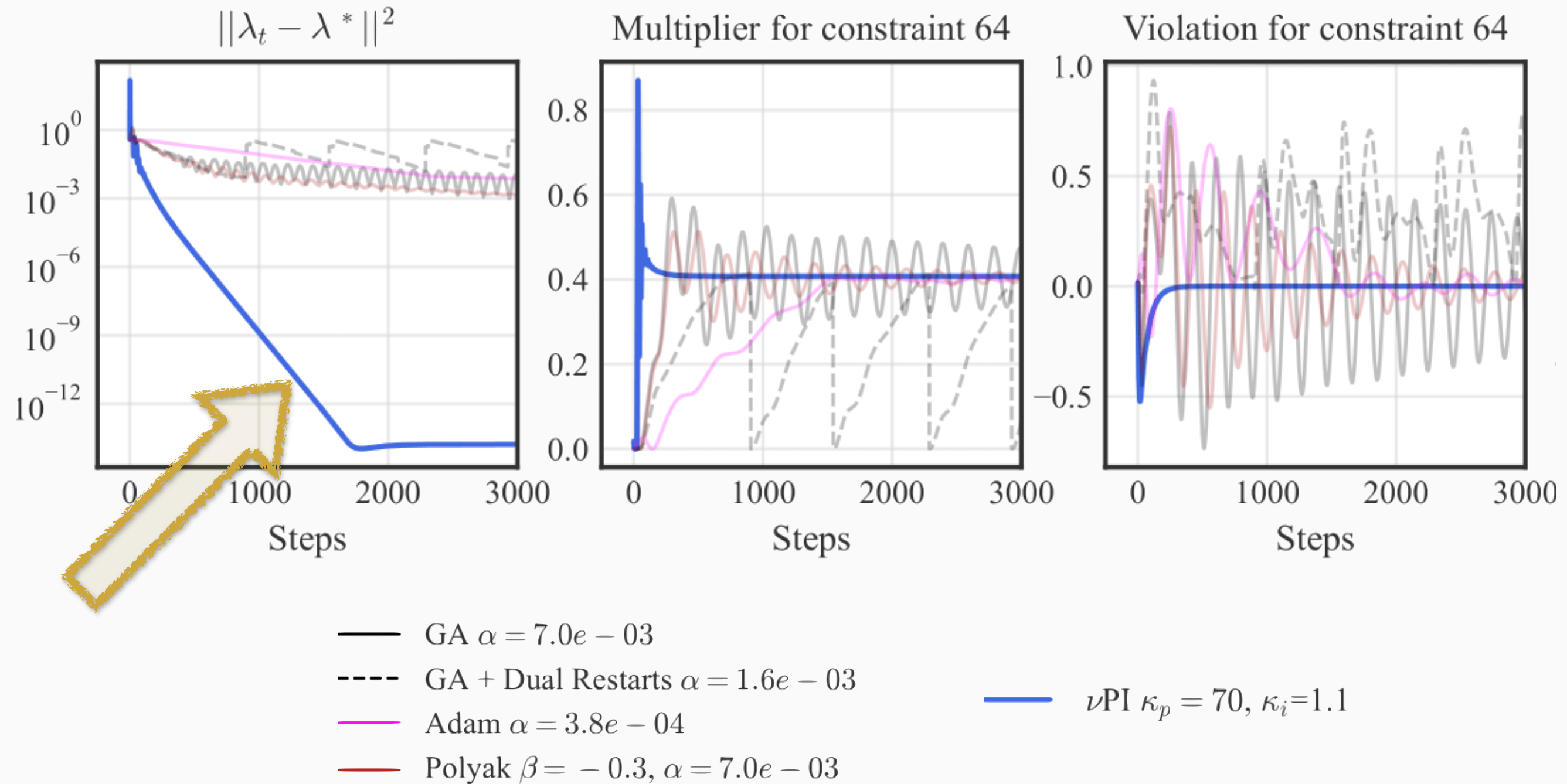
# TLDR of our paper

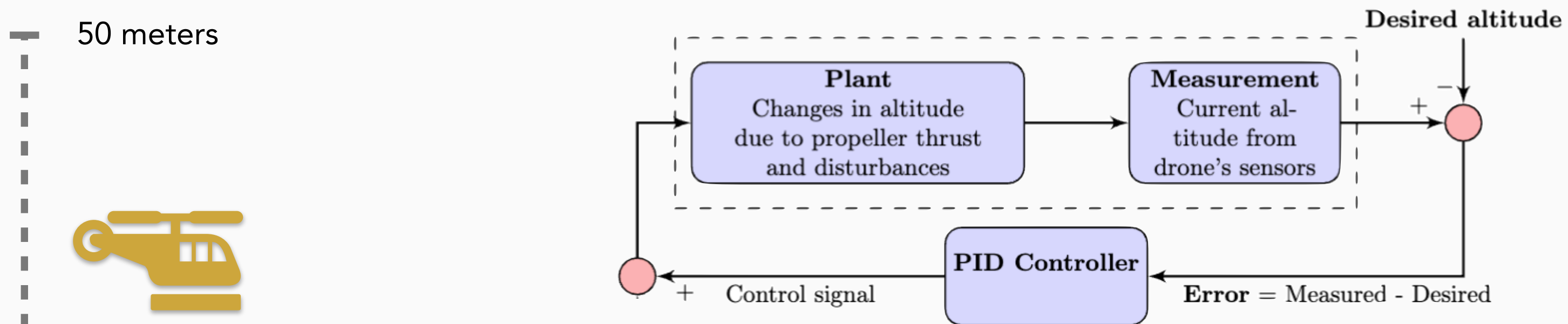‣ Stooke et al. (2020) propose **updating the Lagrange multipliers based on PID control**, improving stability on RL tasks with safety constraints.

‣ We provide an **optimization-oriented analysis of $\nu$PI**, our proposed PI controller

   ‣ $\nu$PI yields stable dynamics and allows for monotonic control on the degree of overshoot

   ‣ Conceptual insights explaining *why* using $\nu$PI helps

   ‣ Experimental evidence in SVMs and sparsity-constrained ResNets

‣ We prove that $\nu$**PI generalizes standard optimization techniques** (including momentum)

   ‣ We provide insights as to why momentum methods may aggravate the issues of GDA

Stooke, A., Achiam, J., and Abbeel, P. *Responsive Safety in Reinforcement Learning by PID Lagrangian Methods*. In ICML, 2020.

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



$||\lambda_t - \lambda^*||^2$     Multiplier for constraint 64     Violation for constraint 64

—— GA $\alpha = 7.0e - 03$

---- GA + Dual Restarts $\alpha = 1.6e - 03$

—— Adam $\alpha = 3.8e - 04$

—— Polyak $\beta = -0.3$, $\alpha = 7.0e - 03$

—— $\nu$PI $\kappa_p = 70$, $\kappa_i = 1.1$

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

# PID control in one slide

50 meters



**Desired altitude**

| Plant | Measurement |
|---|---|
| Changes in altitude due to propeller thrust and disturbances | Current altitude from drone's sensors |

**PID Controller**

Control signal

**Error** = Measured - Desired

## Continuous-time (Analog)

$$u_t = \kappa_p e_t + \kappa_i \int_0^t e_\tau \, \mathrm{d}\tau + \kappa_d \frac{\mathrm{d}e_t}{\mathrm{d}t}$$

## Discrete-time (Digital)

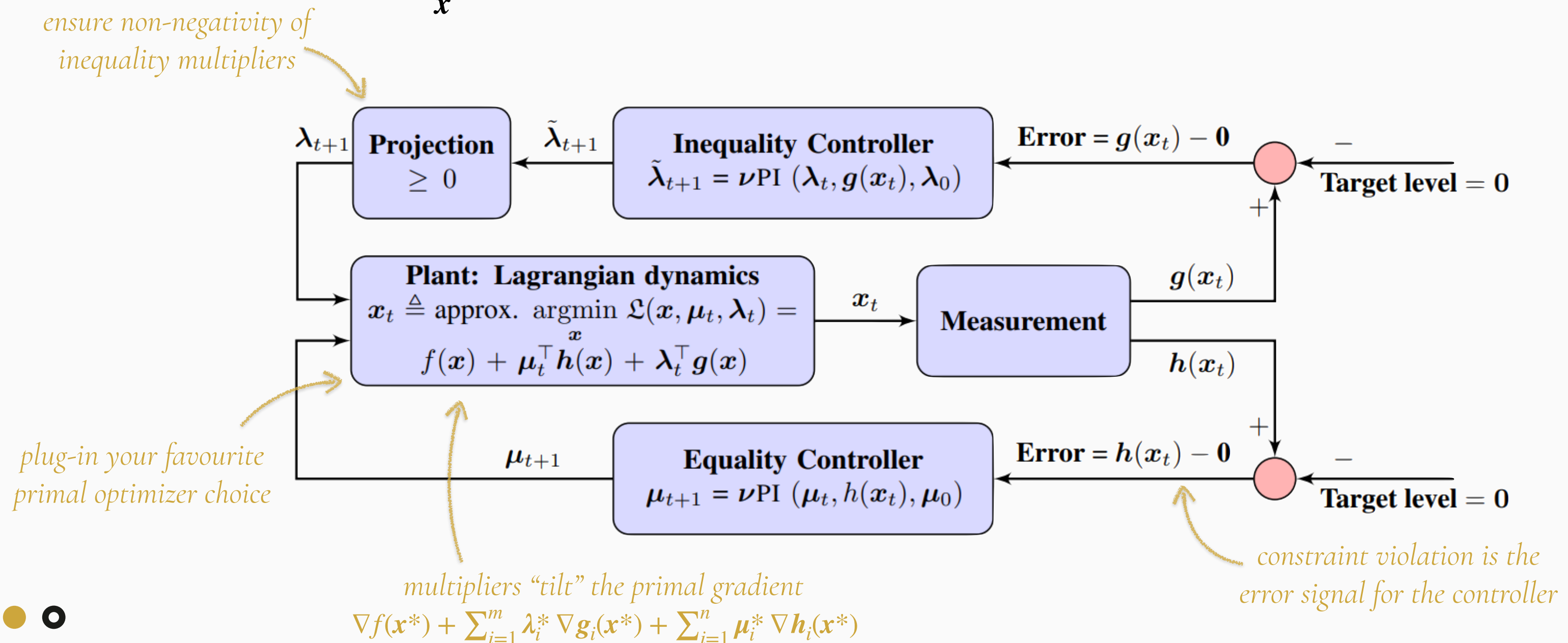$$u_t = \kappa_p e_t + \kappa_i \sum_{\tau=0}^{t} e_\tau + \kappa_d \left( e_t - e_{t-1} \right)$$

# A control theory view of constrained optimization

# Dynamical system's view of CO

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



*ensure non-negativity of inequality multipliers*

*plug-in your favourite primal optimizer choice*

*multipliers "tilt" the primal gradient*
$$\nabla f(\boldsymbol{x}^*) + \sum_{i=1}^{m} \lambda_i^* \nabla \boldsymbol{g}_i(\boldsymbol{x}^*) + \sum_{i=1}^{n} \mu_i^* \nabla \boldsymbol{h}_i(\boldsymbol{x}^*)$$

*constraint violation is the error signal for the controller*

**Projection** $\geq 0$

**Inequality Controller** $\tilde{\boldsymbol{\lambda}}_{t+1} = \nu \text{PI}(\boldsymbol{\lambda}_t, \boldsymbol{g}(\boldsymbol{x}_t), \boldsymbol{\lambda}_0)$

$\boldsymbol{\lambda}_{t+1}$  $\tilde{\boldsymbol{\lambda}}_{t+1}$

**Error** $= \boldsymbol{g}(\boldsymbol{x}_t) - \mathbf{0}$  —  **Target level** $= 0$  +

**Plant: Lagrangian dynamics**
$$\boldsymbol{x}_t \triangleq \text{approx.} \underset{\boldsymbol{x}}{\text{argmin}} \, \mathfrak{L}(\boldsymbol{x}, \boldsymbol{\mu}_t, \boldsymbol{\lambda}_t) =$$
$$f(\boldsymbol{x}) + \boldsymbol{\mu}_t^\top \boldsymbol{h}(\boldsymbol{x}) + \boldsymbol{\lambda}_t^\top \boldsymbol{g}(\boldsymbol{x})$$

$\boldsymbol{x}_t$

**Measurement**

$\boldsymbol{g}(\boldsymbol{x}_t)$

$\boldsymbol{h}(\boldsymbol{x}_t)$

**Equality Controller** $\boldsymbol{\mu}_{t+1} = \nu \text{PI}(\boldsymbol{\mu}_t, \boldsymbol{h}(\boldsymbol{x}_t), \boldsymbol{\mu}_0)$

$\boldsymbol{\mu}_{t+1}$

**Error** $= \boldsymbol{h}(\boldsymbol{x}_t) - \mathbf{0}$  +  —  **Target level** $= 0$

# $\nu$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 $\boldsymbol{\xi}_0$ and $\boldsymbol{\theta}_0$

1. Measure the current system error $\boldsymbol{e}_t$

2. $\boldsymbol{\xi}_t \leftarrow \nu\boldsymbol{\xi}_{t-1} + (1-\nu)\boldsymbol{e}_t$   (for $t \geq 1$)

3. $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_0 + \kappa_p\boldsymbol{\xi}_t + \kappa_i \sum_{\tau=0}^{t} \boldsymbol{e}_\tau$

Recursively, $\boldsymbol{\theta}_1 \leftarrow \boldsymbol{\theta}_0 + \kappa_p\boldsymbol{\xi}_0 + \kappa_i\boldsymbol{e}_0$

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \kappa_i\boldsymbol{e}_t + \kappa_p\left(\boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1}\right)$$

General case

*like $\nabla$-ascent*

*new term looks at* **change** *in constraint satisfaction!*

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \kappa_i\boldsymbol{e}_t + \kappa_p\left(\boldsymbol{e}_t - \boldsymbol{e}_{t-1}\right)$$

Case $\nu = 0$

# Two low-hanging fruits

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \kappa_i \boldsymbol{e}_t + \kappa_p \left( \boldsymbol{\xi}_t - \boldsymbol{\xi}_{t-1} \right) = \boldsymbol{\theta}_t + \kappa_i \boldsymbol{e}_t + \kappa_p (1 - \nu) \left( \boldsymbol{e}_t - \boldsymbol{\xi}_{t-1} \right)$$

Suppose that the error signal is the negative gradient of a loss function $f$: $\boldsymbol{e}_t = - \nabla_{\boldsymbol{\theta}} f$

**Gradient descent:** $\kappa_p = 0$

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_0 + \kappa_i \boldsymbol{e}_t$$

**Optimistic gradient method** (Popov, 1980): $\kappa_p = \kappa_i; \nu = 0$

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \kappa_i \left[ \boldsymbol{e}_t + \left( \boldsymbol{e}_t - \boldsymbol{e}_{t-1} \right) \right]$$

Popov, L. D. *A modification of the Arrow-Hurwicz method for search of saddle points.* Mathematical notes of the Academy of Sciences of the USSR, 1980.

# The updates of $\nu$PI

The entries of $\boldsymbol{\theta}$ can be updated in parallel, but evolve collectively!

$$\theta_{t+1}^{\nu\text{PI}} \leftarrow \theta_t + \kappa_i e_t + \kappa_p(1-\nu)\big(e_t - \xi_{t-1}\big)$$

$$\theta_{t+1}^{\text{GA}} \leftarrow \theta_t + \kappa_i e_t$$

$$\frac{\Delta\nu\text{PI}}{\Delta\text{GA}} = \frac{\theta_{t+1}^{\nu\text{PI}} - \theta_t}{\theta_{t+1}^{\text{GA}} - \theta_t} = \frac{1}{1-\kappa}\left[1 - \frac{\kappa\xi_{t-1}}{e_t}\right]$$

*constant that depends on $\kappa_i$ and $\kappa_p$*

*"how large is the $\nu$PI update compared to GA?"*



Faster    Slower    Opposite direction

15

# $\nu$PI generalizes momentum methods

**Theorem 1**

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

Under the same initialization $\boldsymbol{\theta}_0$, UnifiedMomentum($\alpha$, $\beta \neq 1$, $\gamma$) is a special case of the $\nu$PI algorithm with the hyperparameter choices:

$$\nu \leftarrow \beta \qquad\qquad \boldsymbol{\xi}_0 \leftarrow (1 - \beta)\boldsymbol{e}_0$$

$$\kappa_i \leftarrow \frac{\alpha}{1 - \beta} \qquad\qquad \kappa_p \leftarrow -\frac{\alpha\beta}{(1 - \beta)^2}[1 - \gamma(1 - \beta)]$$

Shen, L., Chen, C., Zou, F., Jie, Z., Sun, J., and Liu, W. A. *Unified Analysis of AdaGrad with Weighted Aggregation and Momentum Acceleration.* In IEEE TNNLS, 2018.

# $\nu$PI generalizes momentum methods

| Algorithm | $\boldsymbol{\xi}_0$ | $\kappa_p$ | $\kappa_i$ | $\nu$ |
|---|---|---|---|---|
| UnifiedMomentum$(\alpha, \beta, \gamma)$ | $(1-\beta)\boldsymbol{e}_0$ | $-\dfrac{\alpha\beta}{(1-\beta)^2}\left[1 - \gamma(1-\beta)\right]$ | $\dfrac{\alpha}{1-\beta}$ | $\beta$ |
| Polyak$(\alpha, \beta)$ | $(1-\beta)\boldsymbol{e}_0$ | $-\dfrac{\alpha\beta}{(1-\beta)^2}$ | $\dfrac{\alpha}{1-\beta}$ | $\beta$ |
| Nesterov$(\alpha, \beta)$ | $(1-\beta)\boldsymbol{e}_0$ | $-\dfrac{\alpha\beta^2}{(1-\beta)^2}$ | $\dfrac{\alpha}{1-\beta}$ | $\beta$ |
| PI | $\boldsymbol{e}_0$ | $\kappa_p$ | $\kappa_i$ | $0$ |
| OptimisticGradientAscent$(\alpha)$ | $\boldsymbol{e}_0$ | $\alpha$ | $\alpha$ | $0$ |
| $\boldsymbol{\nu}$PI $(\kappa_i, \kappa_p, \nu)$ in practice | $\boldsymbol{0}$ | $\kappa_i$ | $\kappa_p$ | $\nu$ |
| GradientAscent$(\alpha)$ | $-$ | $0$ | $\alpha$ | $0$ |

# $\nu$PI generalizes momentum methods

$$\kappa_i \leftarrow \frac{\alpha}{1-\beta}$$

Note the sign of the $\kappa_p$ coefficient for Polyak and Nesterov:

$$\kappa_p^{\text{Polyak}} \leftarrow -\frac{\alpha\beta}{(1-\beta)^2}$$

*non-positive for both positive and negative momentum*

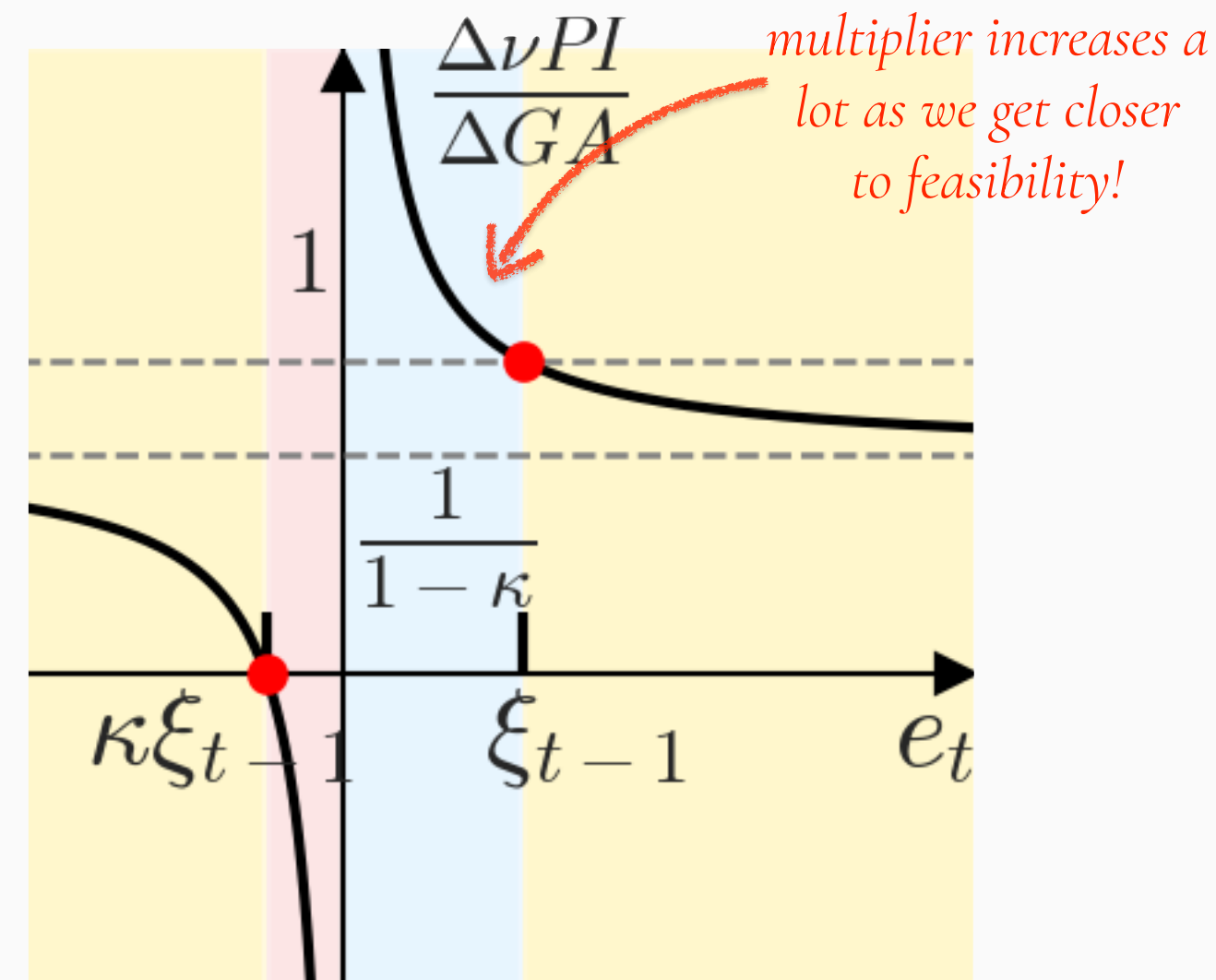$$\kappa_p^{\text{Nesterov}} \leftarrow -\frac{\alpha\beta^2}{(1-\beta)^2} \leq 0$$



*Spanned $\kappa_i$ and $\kappa_p$ coefficients for fixed $\alpha$ and changing $\beta$*

# The (undesirable?) effect of momentum



Polyak $\nu = -0.3$

Polyak $\nu = +0.3$

*multiplier increases a lot as we get closer to feasibility!*

Faster   Slower   Opposite direction

# $\nu$PI in context

# Positive momentum

- Is a special case of $\nu$PI

- Induces a negative $\kappa_p$

- Has been shown to be counterproductive for (bi-linear) games
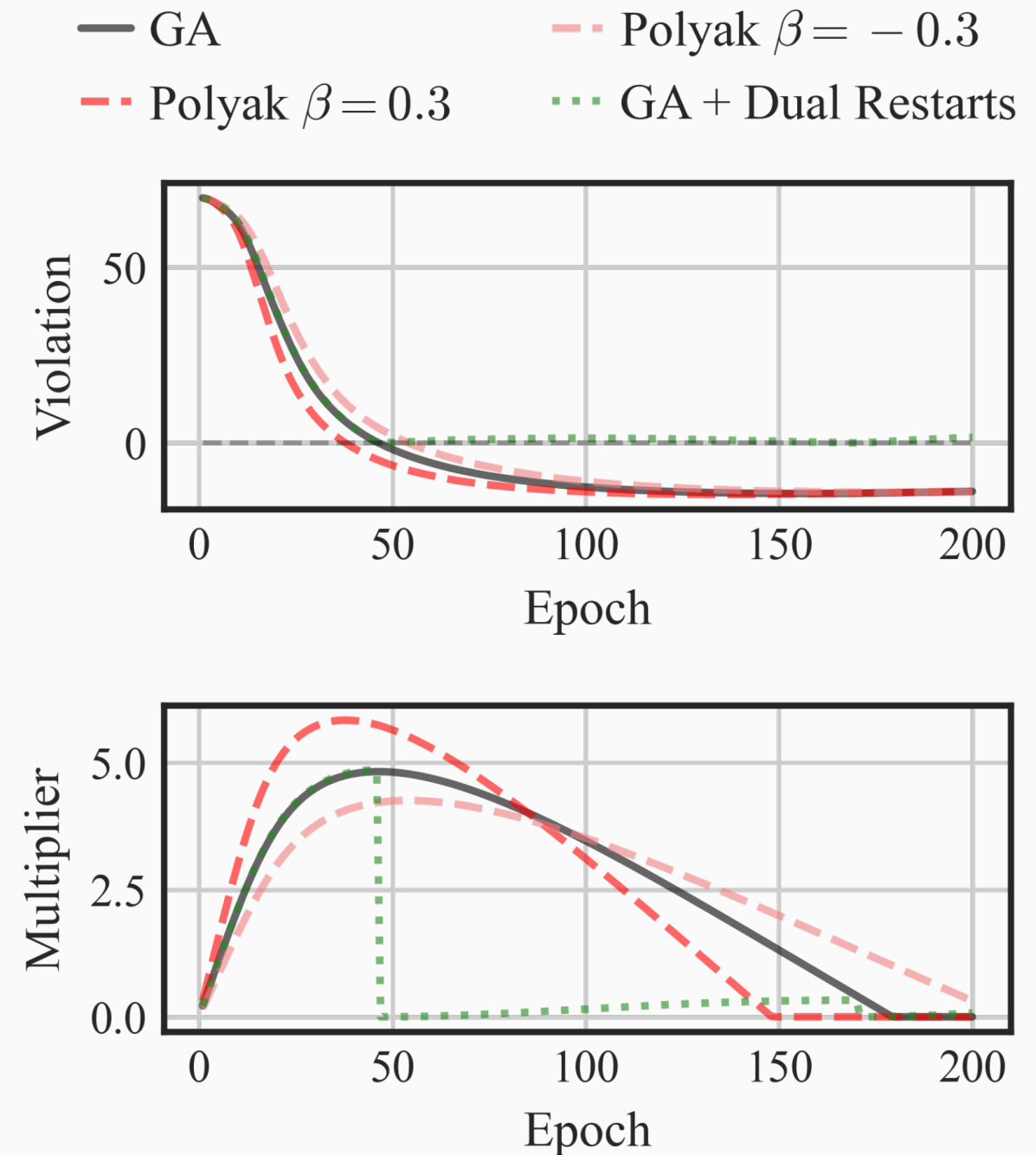
- **Makes overshoot problem worse**

# Negative momentum

- Is a special case of $\nu$PI

- Induces a positive $\kappa_p$

- Suboptimal for strongly convex games
  (Zhang et al., 2021)

- Alleviates multiplier overshooting, but not
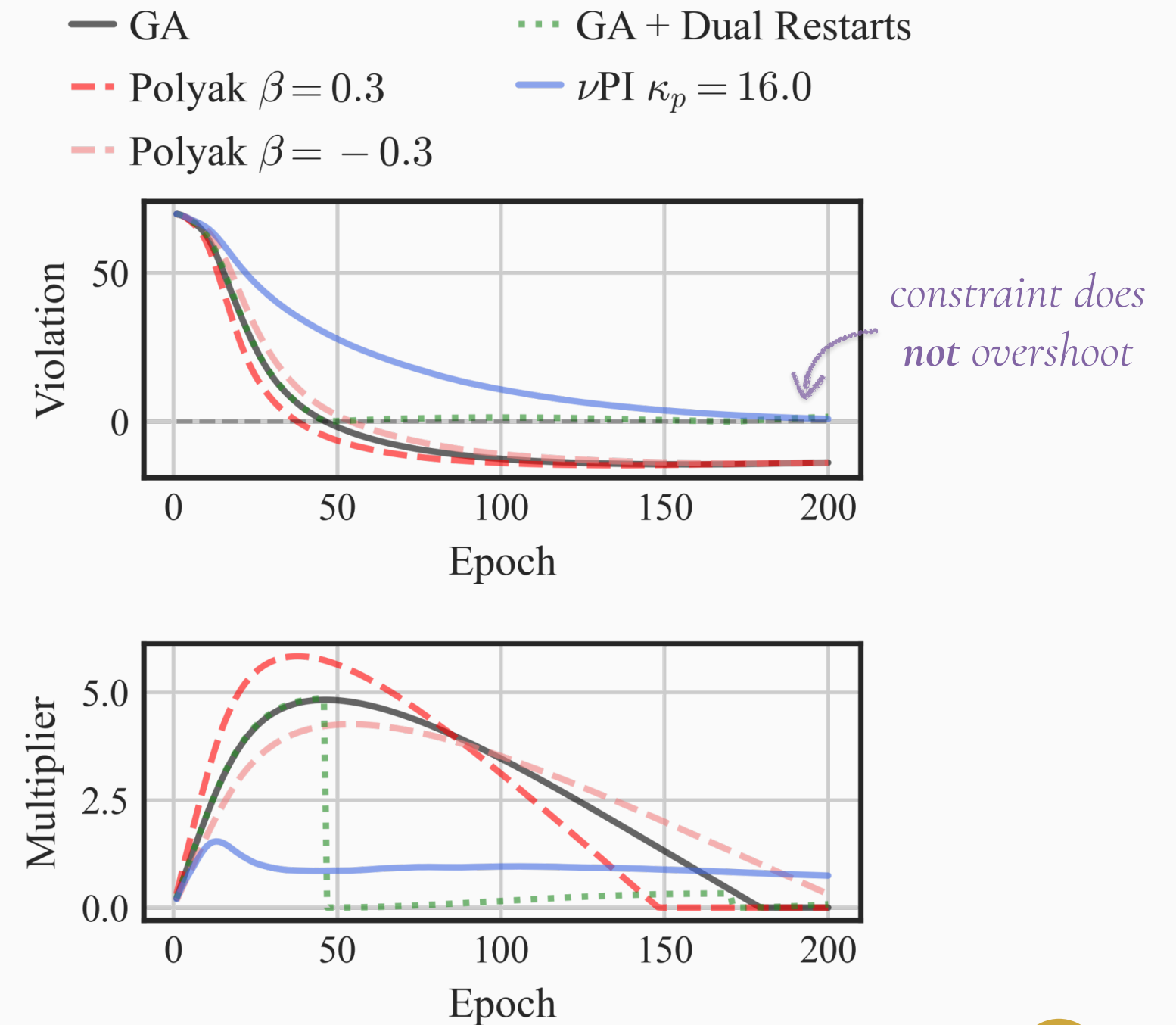  "over-enforcement" of the constraint



Zhang, G. and Wang, Y. On the Suboptimality of Negative Momentum for Minimax Optimization. In AISTATS, 2021.

# Dual restarts



- Once a constraint is strictly satisfied, set its corresponding multiplier to zero (Gallego-Posada et al., 2022)

- Only applicable to (strictly feasible) inequality constraints.

- Relies on exact assessment of constraint satisfaction

  - Stochasticity; numerical precision; "temporary satisfaction"

Gallego-Posada, J., Ramirez, J., Erraqabi, A., Bengio, Y., Lacoste-Julien, S. *Controlled Sparsity via Constrained Optimization or: How I Learned to Stop Tuning Penalties and Love Constraints* . In NeurIPS, 2022.

# $\nu$PI controller

- Natural generalization of the optimistic gradient method, which is (near) optimal for games (Mokhtari et al., 2020)

- Monotonic effect of $\kappa_p$ on the degree of overshoot

- One fewer degree of freedom than full PID

Mokhtari, A., Ozdaglar, A. E., and Pattathil, S. *Convergence Rate of O(1/k) for Optimistic Gradient and Extragradient Methods in Smooth Convex-Concave Saddle Point Problems.* SIAM Journal on Optimization, 2020.

# Experiments

# Hard-margin SVMs

$$\min_{\boldsymbol{w}} \frac{1}{2} ||\boldsymbol{w}||^2 \text{ subject to } y_i \left( \boldsymbol{w}^\top \boldsymbol{x}_i + b \right) \geq 1 \text{ for } i = 1, \ldots, N$$

## Motivation

‣ Simple, well behaved convex problem with unique* KKT tuple

‣ Specialized solvers exist for QCQPs, we use this task for illustration

‣ Cheap experiment allows fine grid-search to test influence of hyperparameters of different algorithms

## Experimental setup

‣ Linearly-separable subset of the Iris dataset

‣ 70 training samples $\Rightarrow$ 70 inequality constraints

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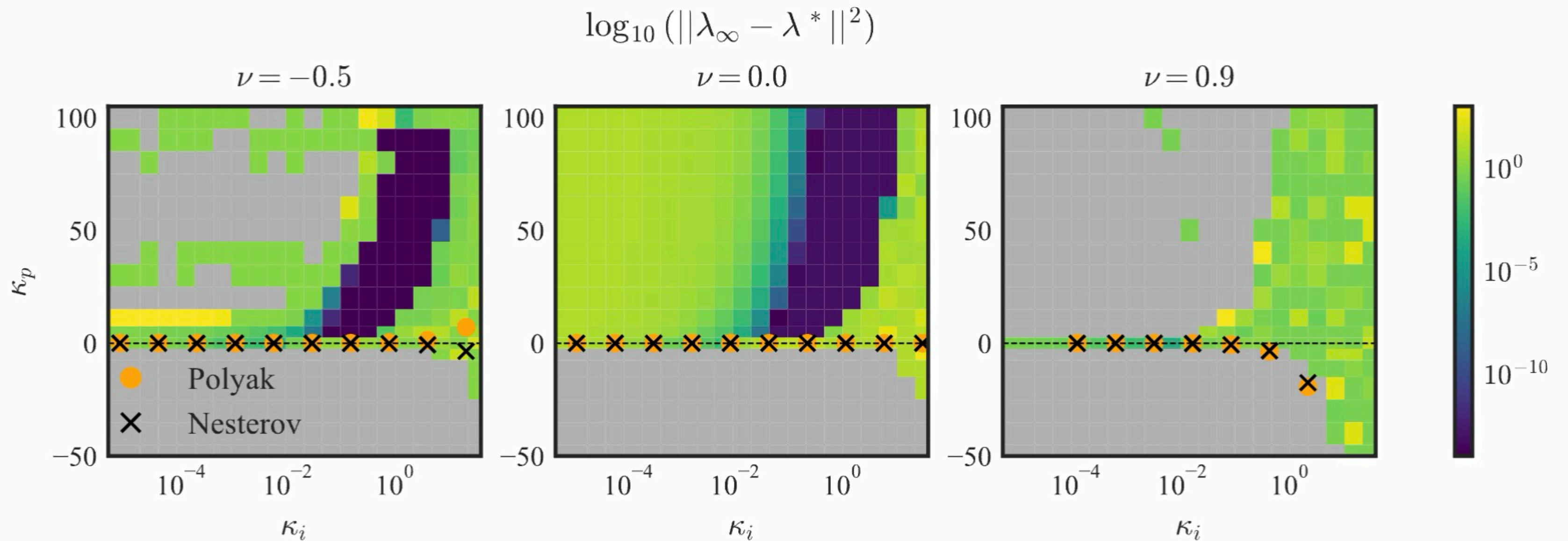


$||\lambda_t - \lambda^*||^2$      Multiplier for constraint 64      Violation for constraint 64

Legend:
- GA $\alpha = 7.0e - 03$
- GA + Dual Restarts $\alpha = 1.6e - 03$
- Adam $\alpha = 3.8e - 04$
- Polyak $\beta = -0.3, \alpha = 7.0e - 03$
- $\nu$PI $\kappa_p = 70, \kappa_i = 1.1$

*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.



$$||\lambda_\infty - \lambda^*||^2$$

Legend:
- Polyak $\beta = -0.5$
- Polyak $\beta = 0.7$
- GA
- GA + Dual Restarts
- $\nu$PI $\kappa_p = 1$
- $\nu$PI $\kappa_p = 10$
- $\nu$PI $\kappa_p = 70$
- Adam

Multiplier step-size

$$\log_{10}\left(||\lambda_\infty - \lambda *||^2\right)$$

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

# Training sparse ResNets

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

## Motivation

‣ More realistic deep application with non-convex constraints

‣ In our prior work we document the issue of overshoot and propose "dual restarts" heuristic

## Experimental setup

‣ Training a ResNet-18 model on CIFAR10

‣ Structured sparsity with layer-wise or model-wise constraints

# Addressing constraint overshooting

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

# Monotonicity on $\kappa_p$