

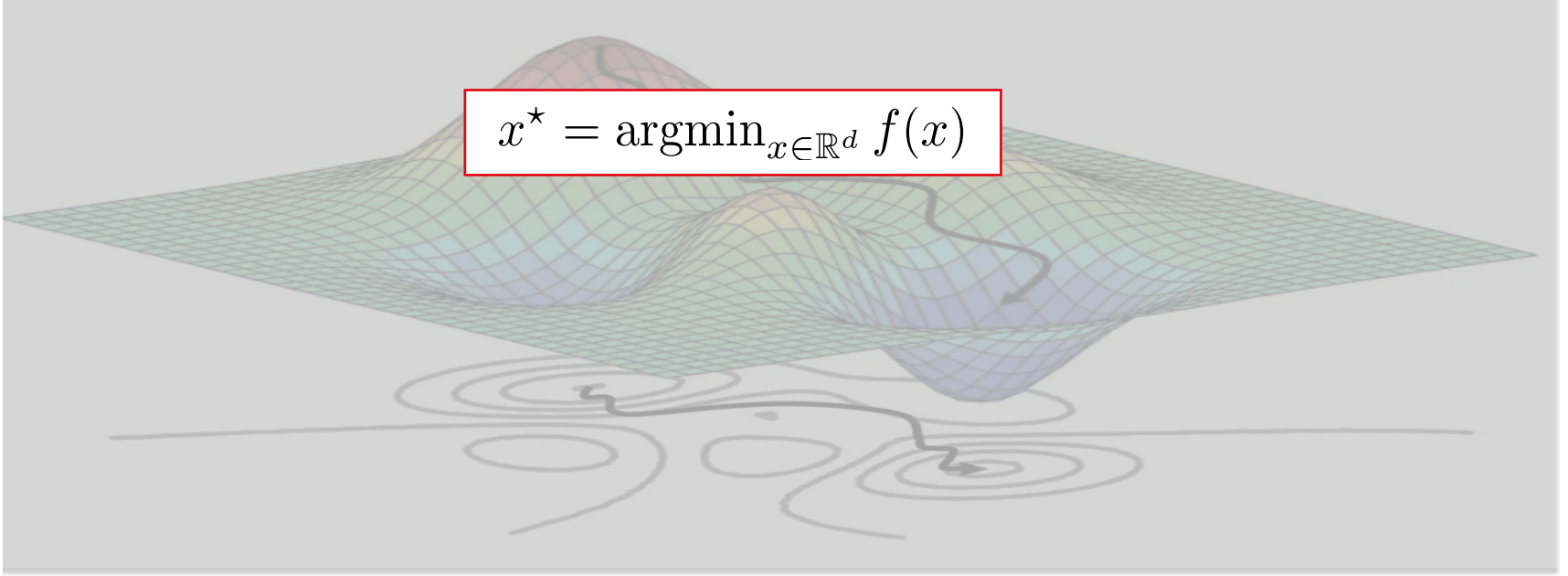
Learning to Optimize with Convergence Guarantees using nonlinear system theory^[1]

Luca Furieri

[1] Andrea Martin and Luca Furieri, “Learning to optimize with convergence guarantees using nonlinear system theory”, IEEE Control System Letters, 2024



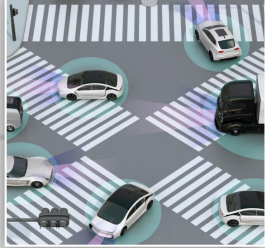
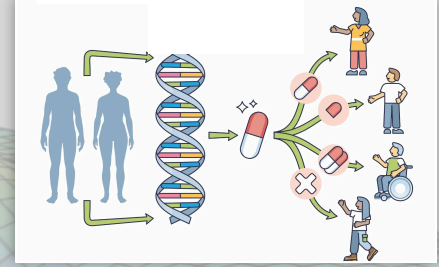
Optimization across disciplines



Optimization across disciplines



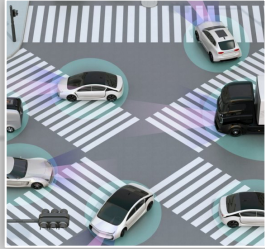
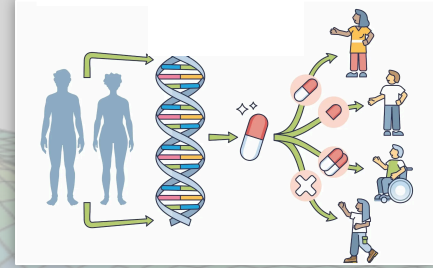
$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$$



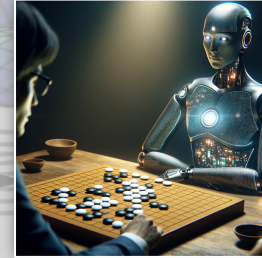
Optimization across disciplines



$$x^* = \operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$$



Nonconvex
Large scale
No analytical solution



Iterative optimization algorithms

$$x_{t+1} = x_t + \operatorname{update}_t(f, x_t) \quad \dots \text{until convergence}$$

(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial guesses x_0 ?
2. **How many iterations** for the algorithm to converge?
3. Does the algorithm find a «**good**» solution?



(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial guesses x_0 ?
2. **How many iterations** for the algorithm to converge?
3. Does the algorithm find a «**good**» solution?

Examples



(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial guesses x_0 ?
2. **How many iterations** for the algorithm to converge?
3. Does the algorithm find a «**good**» solution?

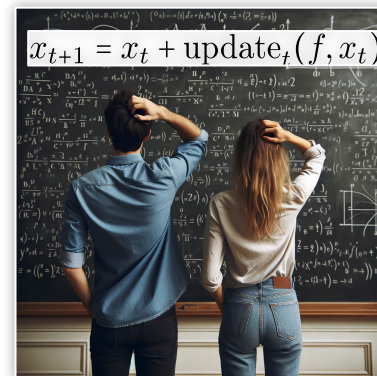


Examples

- Steepest descent algorithm: $x_{t+1} = x_t - \eta \nabla f(x_t)$
 1. **Any x_0**
 2. **Linear rate.**
 3. **Closest stationary point**

(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial guesses x_0 ?
2. **How many iterations** for the algorithm to converge?
3. Does the algorithm find a «**good**» solution?



Examples

- Steepest descent algorithm: $x_{t+1} = x_t - \eta \nabla f(x_t)$
 1. **Any x_0**
 2. **Linear rate.**
 3. **Closest stationary point**
- Newton-Raphson algorithm: $x_{t+1} = x_t - (\nabla^2 f(x_t))^{-1} \nabla f(x_t)$
 1. **(x_0 close to x^* , global variants)**
 2. **Quadratic rate**
 3. **Stationary point**

(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial conditions?
2. **How many iterations** to reach a certain accuracy?
3. **Does the algorithm** converge to a local minimum?



Hyper-parameter tuning

- Steepest descent
 - 1. **Any x_0**
 - 2. **Step-size, 1° and 2° momentum...**
 - 3. **Lack of a general theory beyond convex**
 - In ML, best practices and know-how



Design of new algorithms

- Newton-Raphson
 - 1. **(x_0 close to stationary point)**
- Adam algorithm (Kingma et al., 2014)
 - 1. **May diverge on convex**
 - 2. **Empirical**
 - 3. **Deep learning, empirical**



(Hard) questions for the algorithm designer

1. **Convergence** as $t \rightarrow \infty$? From which initial conditions?
2. **How many iterations** to reach a certain accuracy?
3. **Does the algorithm** converge to a local minimum?



Hyper-parameter tuning

- Steepest descent
 - Step-size, 1° and 2° momentum...
 - Lack of a general theory beyond convex
 - In ML, best practices and know-how

1. **Any x_0**

- Newton-Raphson

1. (x_0 close to



Design of new algorithms

- Adam algorithm (Kingma et al., 2014)

1. **May diverge on convex** 2. **Empirical** 3. Deep learning, empirical



tationary point

Design of new optimization algorithms

Numerical
Experiments



04

01



System Theory
for algorithm design

Theoretical
Results



03

02



Machine Learning
for algorithm design

Design of new optimization algorithms

Numerical
Experiments



04

01



System Theory
for algorithm design

Theoretical
Results



03

02

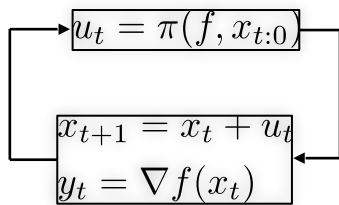


Machine Learning
for algorithm design



Success stories from system theory

Classical Algorithms \equiv Control policies^{[1],[2]}



$$u_t = \begin{cases} -\eta \nabla f(x_t) & \text{(gradient descent)} \\ -\eta \nabla f(x_t) + \beta(x_t - x_{t-1}) & \text{(heavy ball)}^{[1]} \\ -(\nabla^2 f(x_t))^{-1} \nabla f(x_t) & \text{(Newton-Raphson)} \\ \dots & \end{cases}$$

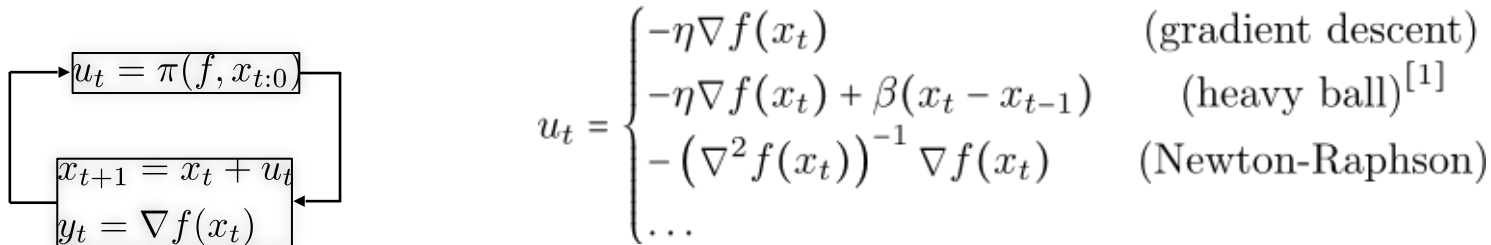
[1] B. Polyak. *Introduction to Optimization*. Optimization Software, Inc., New York, 1987

[2] J. Wang and N. Elia. *A control perspective for centralized and distributed convex optimization*. In IEEE Conference on Decision and Control and European Control Conference Orlando, FL, USA. IEEE, 2011



Success stories from system theory

Classical Algorithms \equiv Control policies^{[1],[2]}



Algorithm design \equiv **Robust** control problem

For every cost function f in a class \mathcal{F} , designed algo must achieve:

1. **Convergence:** regulate $y_t = \nabla f(x_t)$ to zero
2. **Speed:** $\mathcal{H}_2/\mathcal{H}_\infty$ guarantees, exponential stability...
3. **Solution quality:** e.g., performance of equilibria

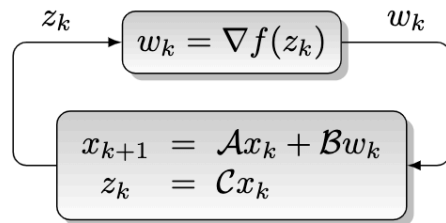
[1] B. Polyak. *Introduction to Optimization*. Optimization Software, Inc., New York, 1987

[2] J. Wang and N. Elia. *A control perspective for centralized and distributed convex optimization*. In IEEE Conference on Decision and Control and European Control Conference Orlando, FL, USA. IEEE, 2011



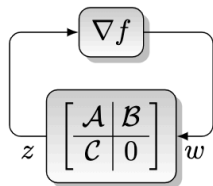
Success stories from system theory

- Classical algorithms equivalent to uncertain LTIs



$$\left(\begin{array}{c|c} \mathcal{A} & \mathcal{B} \\ \hline \mathcal{C} & 0 \end{array} \right) = \left(\begin{array}{cc|c} (1+\beta)I_d & -\beta I_d & -\alpha I_d \\ I_d & 0 & 0 \\ \hline (1+\gamma)I_d & -\gamma I_d & 0 \end{array} \right)$$

- Call algorithm by Polyak:



- ... is a Lure's system^[1,2]! We know a lot about them...

[1] Lessard, L., Recht, B., & Packard, A. (2016). Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1), 57-95.

[2] Scherer, C., & Ebenbauer, C. (2021). Convex synthesis of accelerated gradient algorithms. *SIAM Journal on Control and Optimization*, 59(6), 4615-4645.



Success stories from system theory

- **Breakthrough**^[1,2]: **convex design** of algorithm $\mathcal{A}, \mathcal{B}, \mathcal{C}$ such that

$$\|x_k - x^*\| \leq K \rho^k \|x_0 - x^*\| \quad K > 0, \rho \in (0, 1)$$

for all cost functions that are **(strongly)-convex and smooth**

- Integral Quadratic Constraints (IQCs) for **optimal worst-case convergence rates**^[1,2]
- Extremum-seeking control, algorithms with delayed info...^[2]

[1] Lessard, L., Recht, B., & Packard, A. (2016). Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1), 57-95.

[2] Scherer, C., & Ebenbauer, C. (2021). Convex synthesis of accelerated gradient algorithms. *SIAM Journal on Control and Optimization*, 59(6), 4615-4645.



Success stories from system theory

- **Breakthrough**^[1,2]: **convex design** of algorithm $\mathcal{A}, \mathcal{B}, \mathcal{C}$ such that

$$\|x_k - x^*\| \leq K \rho^k \|x_0 - x^*\| \quad K > 0, \rho \in (0, 1)$$

for all cost functions that are **(strongly)-convex and smooth**

- Integral Quadratic Constraints (IQC) for **optimal worst-case convergence rates**^[1,2]
- Extremum-seeking control, algorithms with delayed info...^[2]

A main take-away

***Crucial role of nonlinear, robust system theory
for the next generation of optimization algorithms***

[1] Lessard, L., Recht, B., & Packard, A. (2016). Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1), 57-95.

[2] Scherer, C., & Ebenbauer, C. (2021). Convex synthesis of accelerated gradient algorithms. *SIAM Journal on Control and Optimization*, 59(6), 4615-4645.



Success stories from system theory

- **Breakthrough**^[1,2]: **convex design** of algorithm A, B, C such that

for all cost function

- Integral Quadratic
- Extremum-seeking

$\in (0, 1)$

convergence rates^[1,2]



Open challenges

- Algorithms for **non-convex** cost functions
- Scaling to very large problems
- Beyond worst-case rates? (e.g. solution quality)

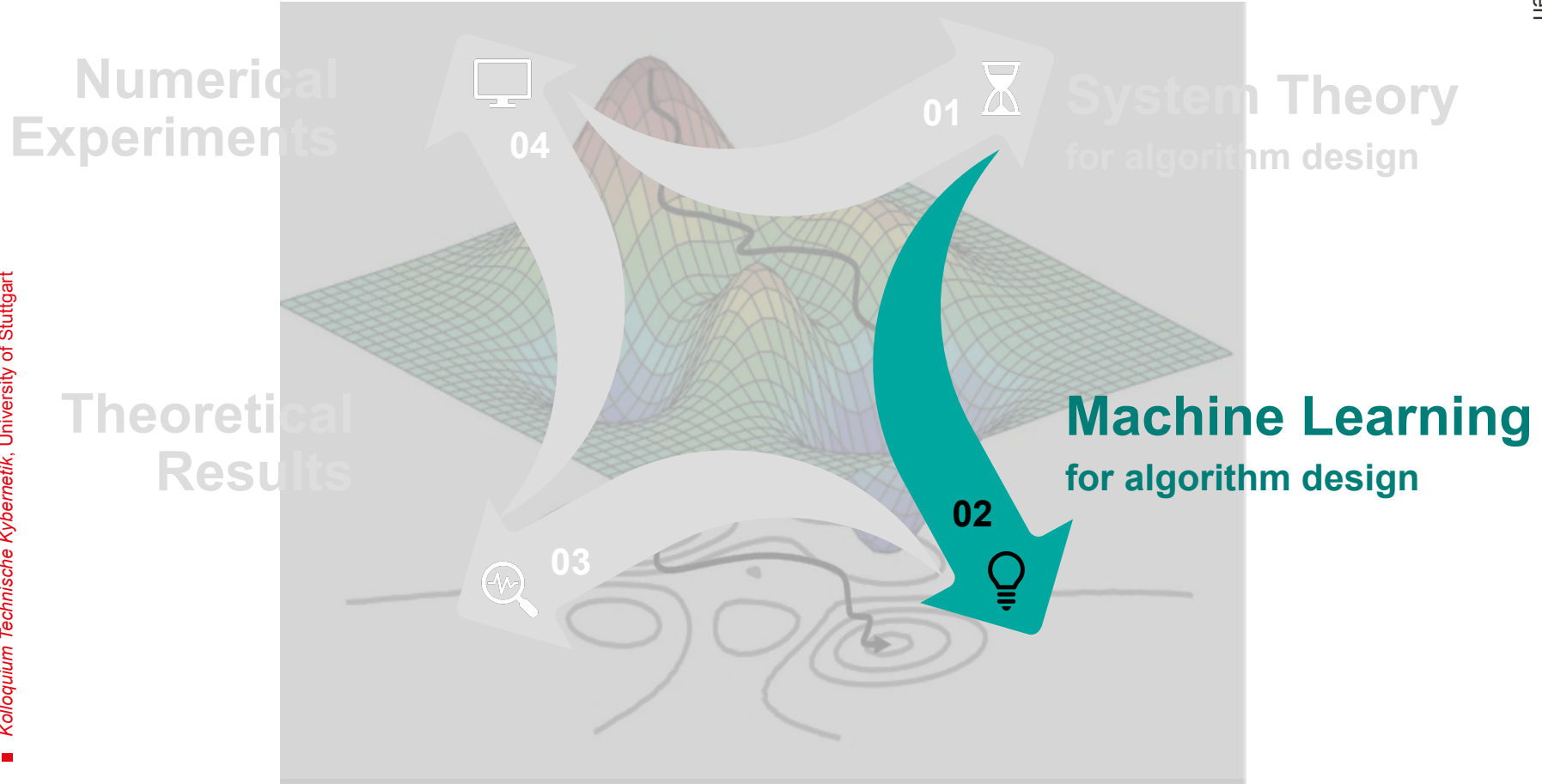
A main take-away

***Crucial role of nonlinear, robust system theory
for the next generation of optimization algorithms***

[1] Lessard, L., Recht, B., & Packard, A. (2016). Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM Journal on Optimization*, 26(1), 57-95.

[2] Scherer, C., & Ebenbauer, C. (2021). Convex synthesis of accelerated gradient algorithms. *SIAM Journal on Control and Optimization*, 59(6), 4615-4645.

Design of new optimization algorithms





Different philosophies for algorithm design

Offline process

Online process

**System theory
point-of-view**



**Machine learning
point-of-view**



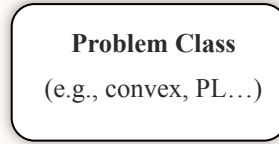


Different philosophies for algorithm design

System theory
point-of-view

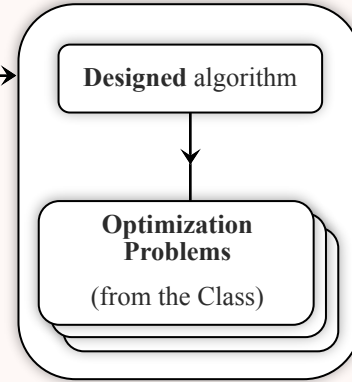


Offline process

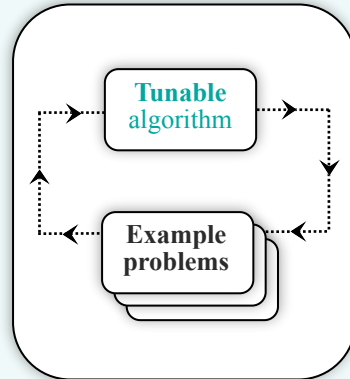


Analytic design
Formal guarantees

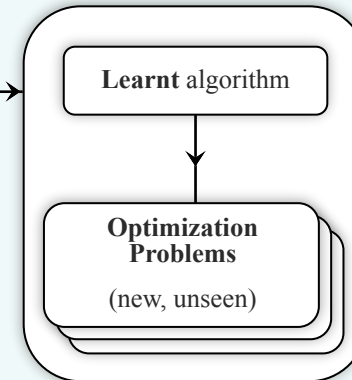
Online process



Machine learning
point-of-view

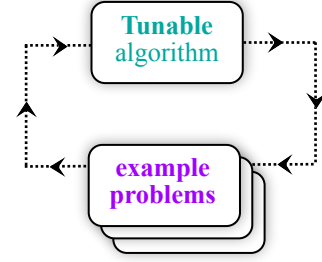


Training
Minimize empirical loss





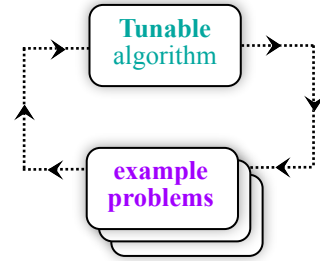
«Learning To Optimize» (L2O) for optimal control

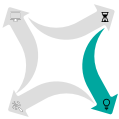




«Learning To Optimize» (L2O) for optimal control

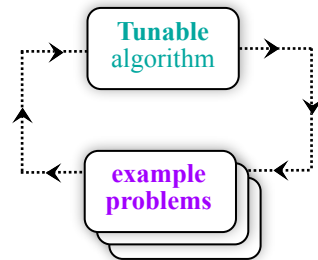
- **Crucial:** «distribution» of **example problems**
 - *generalizability VS performance*
- Rule of thumb
 - *effective when **example problems** are representative of future instances*





«Learning To Optimize» (L2O) for optimal control

- **Crucial:** «distribution» of **example problems**
 - *generalizability VS performance*
- Rule of thumb
 - *effective when **example problems** are representative of future instances*



Prime example: Nonlinear Model Predictive Control (NMPC)

Repeatedly solve...

$$\begin{aligned}
 u^*(x(t)) = \arg \min_{\{u_k\}_{k=0}^{N-1}} & \sum_{k=0}^{N-1} l(x_k, u_k) + l_f(x_N) \\
 \text{subject to} & \quad x_0 = x(t), \quad x_{k+1} = f(x_k, u_k) \\
 & \quad x_k \in \mathcal{X}, \quad u_k \in \mathcal{U}, \quad x_N \in \mathcal{X}_f
 \end{aligned}$$

“examples”

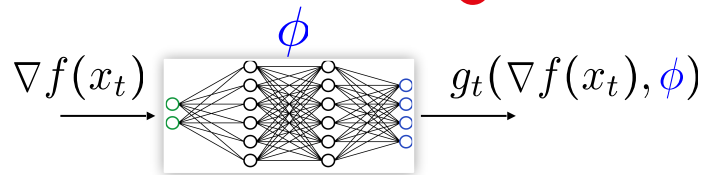
A learnt algorithm can discover shortcuts unknown to standard interior-point solvers



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

[2] Li, Ke, and Jitendra Malik. "Learning to optimize." *ICLR*, 2017



Success of L2O in machine learning^[1,2,3]

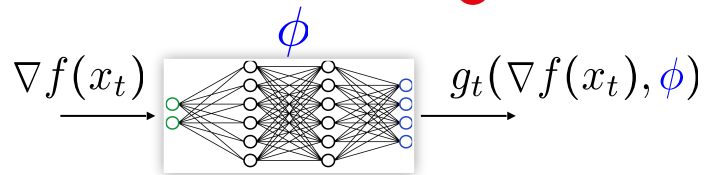
- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

- Measure performance of an algorithm over T steps:

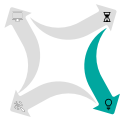
$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

where \mathcal{F} is the set of «*example problems*»



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

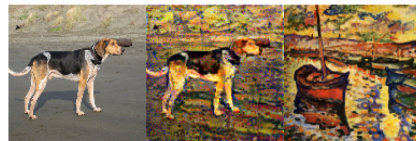
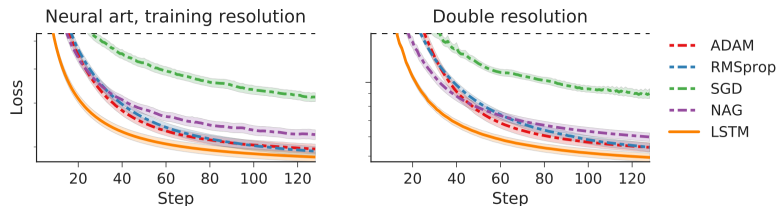
$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

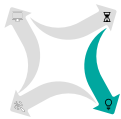
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

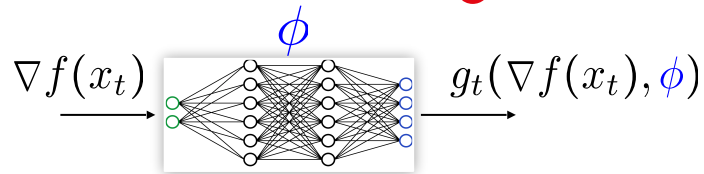
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

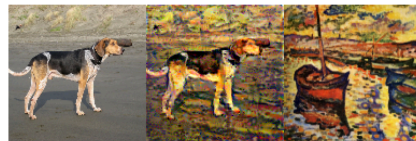
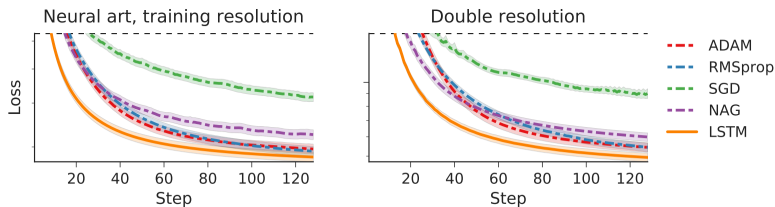


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

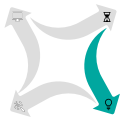
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

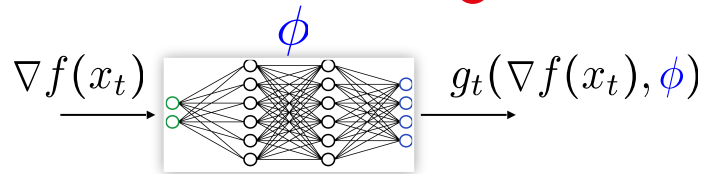
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

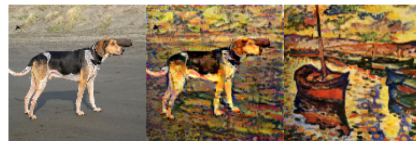
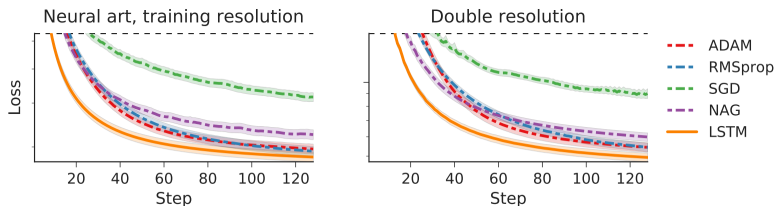


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

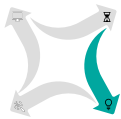
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

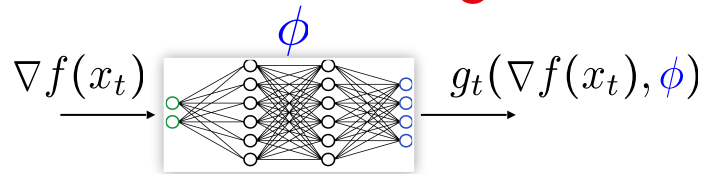
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

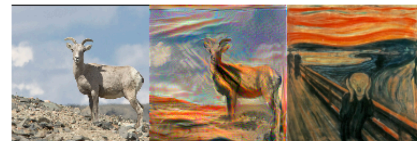
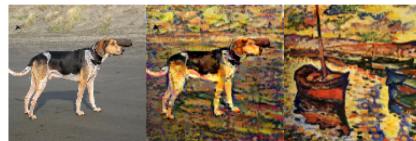
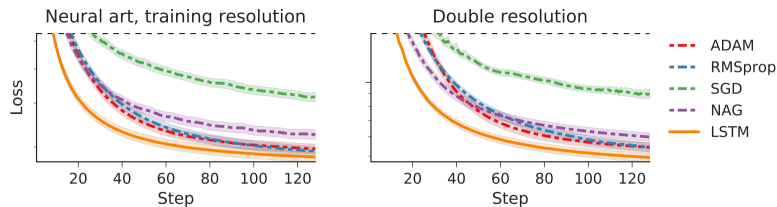


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

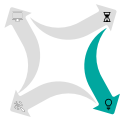
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

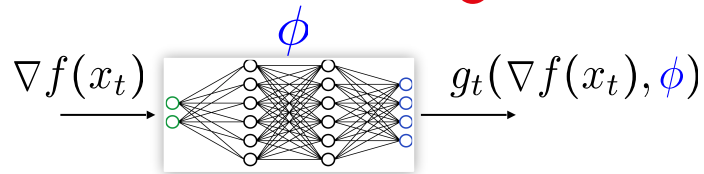
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

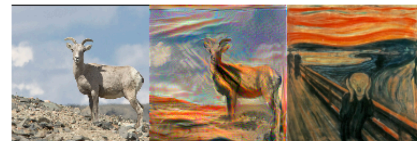
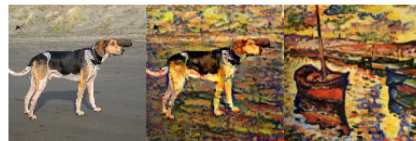
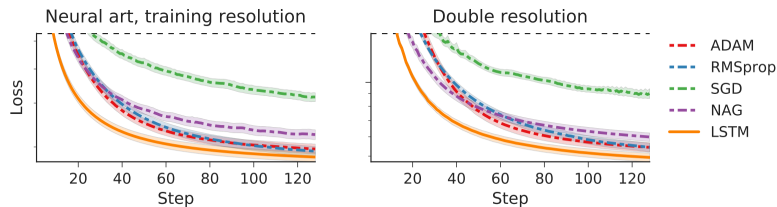


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

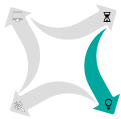
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

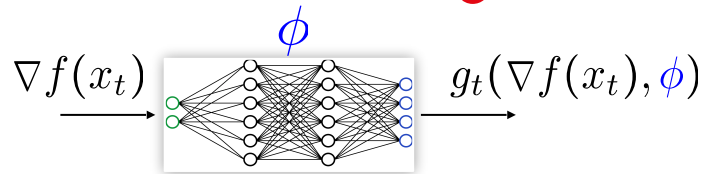
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

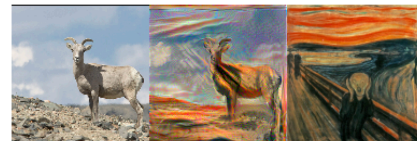
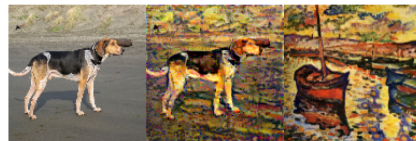
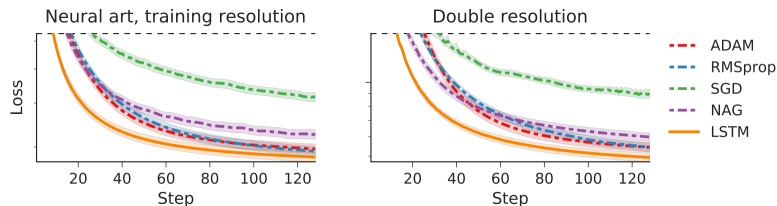


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

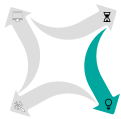
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

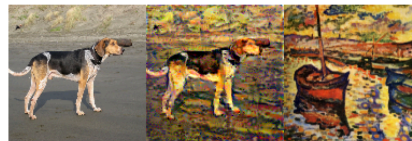
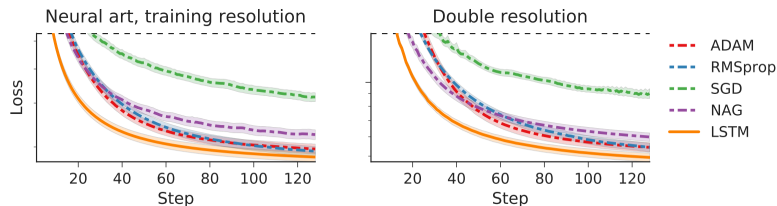
$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

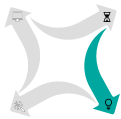
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

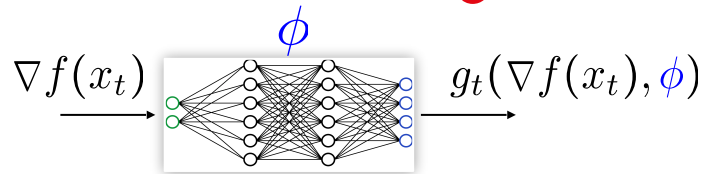
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

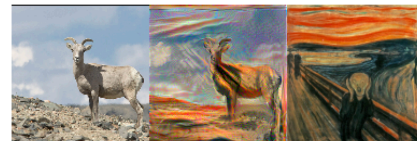
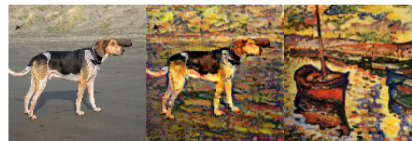
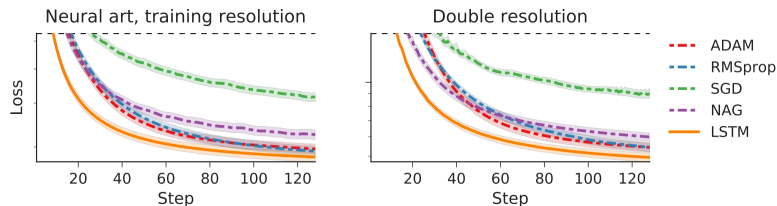


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

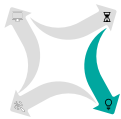
where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

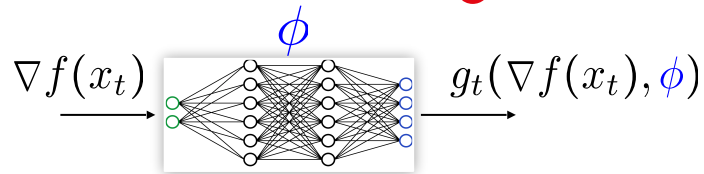
[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017



Success of L20 in machine learning^[1,2,3]

- Parametrize algorithms as follows:

$$x_{t+1} = x_t + g_t(\nabla f(x_t), \phi)$$

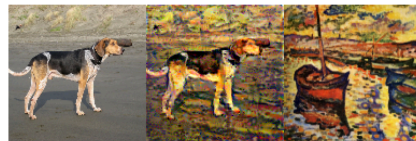
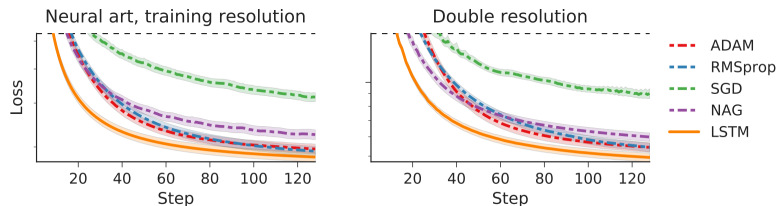


- Measure performance of an algorithm over T steps:

$$L(\phi) = \mathbb{E}_{f \sim \mathcal{F}, x_0 \sim X_0} \sum_{t=0}^T f(x_t)$$

where \mathcal{F} is the set of «*example problems*»

- Generalization capabilities!** Train over MNIST problems, test on NEURAL ART problems



[1] Andrychowicz, Marcin, et al. "Learning to learn by gradient descent by gradient descent." NeurIPS, 2016

[2] Li, Ke, and Jitendra Malik. "Learning to optimize." ICLR, 2017

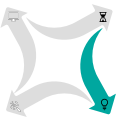


Open challenges for L20



Generalization: no guarantees!

- Performance on unseen problem?
- As hard as a *transfer-learning* problem



Open challenges for L20



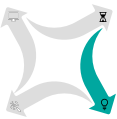
Generalization: no guarantees!

- Performance on unseen problem?
- As hard as a *transfer-learning* problem



Scalability: expensive training

- Repeatedly roll-out algorithm on example problems... **before** improving parameters



Open challenges for L2O



Generalization: no guarantees!

- Performance on unseen problem?
- As hard as a *transfer-learning* problem



Scalability: expensive training

- Repeatedly roll-out algorithm on example problems... **before** improving parameters

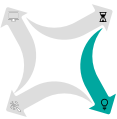


Convergence: no guarantees!^[1]

- Requires conservative **fallback mechanisms**... ..even for the convex case^[2]

[1] Harrison, James, Luke Metz, and Jascha Sohl-Dickstein. "A closer look at learned optimization: Stability, robustness, and inductive biases." NeurIPS, 2022

[2] Heaton, Howard, et al. "Safeguarded learned convex optimization." AAAI Conference, 2023.



Open challenges for L2O



Generalization: no guarantees!

- Performance on unseen problem?
- As hard as a *transfer-learning* problem



Scalability: expensive training

- Repeatedly roll-out algorithm on example problems... **before** improving parameters

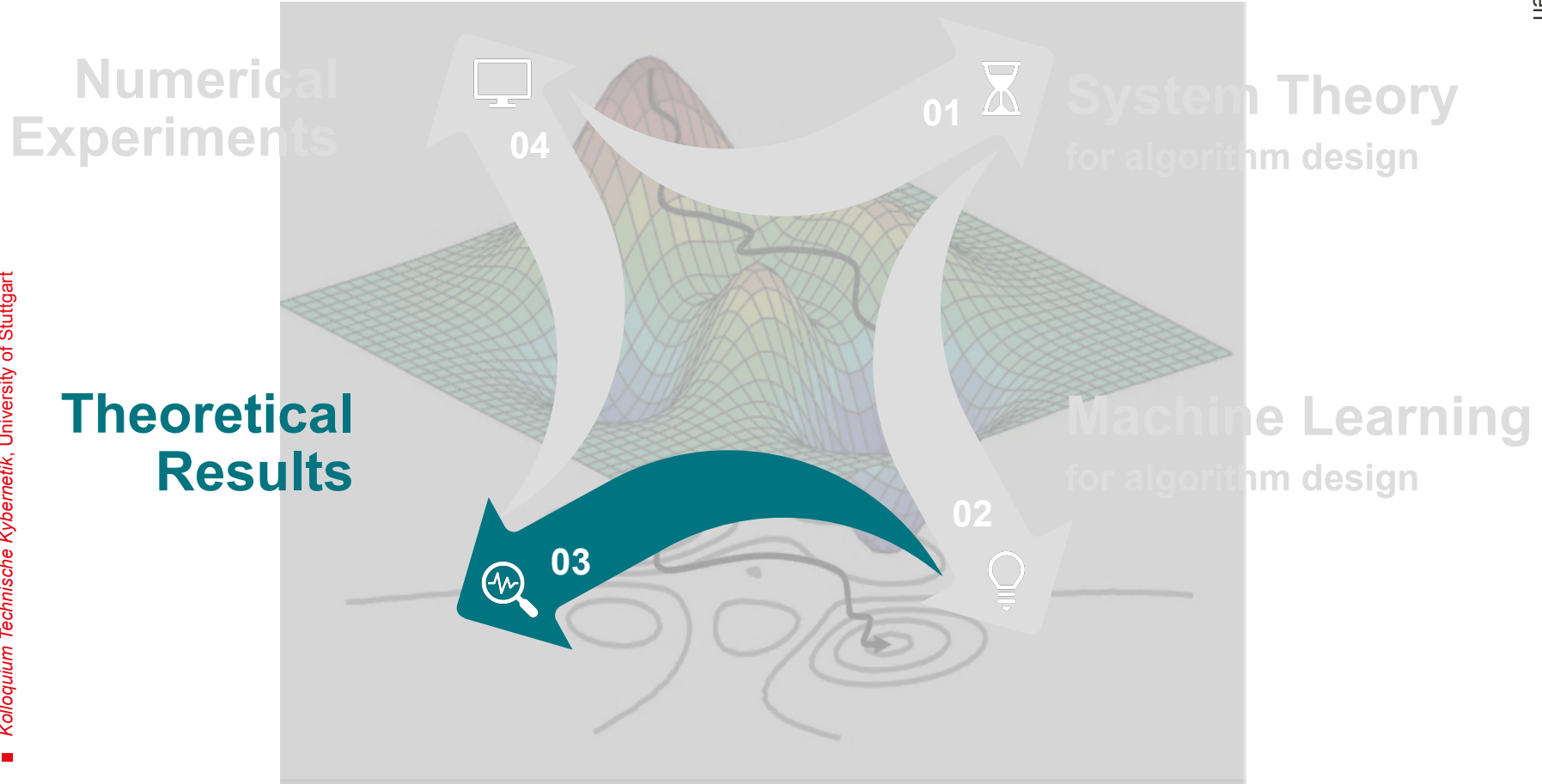
Focus of this work

- **Learn convergent algorithms** for smooth non-convex functions
 - ...beyond conservative heuristics^[1,2]

[1] Harrison, James, Luke Metz, and Jascha Sohl-Dickstein. "A closer look at learned optimization: Stability, robustness, and inductive biases." NeurIPS, 2022

[2] Heaton, Howard, et al. "Safeguarded learned convex optimization." AAAI Conference, 2023.

Design of new optimization algorithms





Problem Formulation (1)

- Denote as \mathcal{S}_β the set of **non-convex** cost functions with β -Lipschitz gradients:

$$f \in \mathcal{S}_\beta \iff |\nabla f(x) - \nabla f(y)| \leq \beta|x - y|$$



Problem Formulation (1)

- Denote as \mathcal{S}_β the set of **non-convex** cost functions with β -Lipschitz gradients:

$$f \in \mathcal{S}_\beta \iff |\nabla f(x) - \nabla f(y)| \leq \beta|x - y|$$

- Signal space notation:

$$\begin{cases} x_{t+1} = x_t + \pi_t(f, x_t) \\ x_0 = x_0 \end{cases} \iff z\mathbf{X} = \mathbf{x} + \boldsymbol{\pi}(f, \mathbf{x}) + z\boldsymbol{\delta}_0$$

$z \equiv$ one-time-step shift

$\mathbf{x} = (x_0, x_1, x_2, \dots)$

$\boldsymbol{\delta}_0 = (x_0, 0, 0, \dots)$

$\boldsymbol{\pi}(f, \mathbf{x}) \equiv$ algorithm



Problem Formulation (1)

- Denote as \mathcal{S}_β the set of **non-convex** cost functions with β -Lipschitz gradients:

$$f \in \mathcal{S}_\beta \iff |\nabla f(x) - \nabla f(y)| \leq \beta|x - y|$$

- Signal space notation:

$$\begin{cases} x_{t+1} = x_t + \pi_t(f, x_t) \\ x_0 = x_0 \end{cases} \iff z\mathbf{x} = \mathbf{x} + \boldsymbol{\pi}(f, \mathbf{x}) + z\boldsymbol{\delta}_0$$

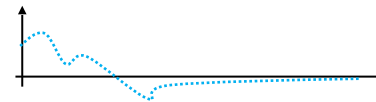
$z \equiv$ one-time-step shift

$\mathbf{x} = (x_0, x_1, x_2, \dots)$

$\boldsymbol{\delta}_0 = (x_0, 0, 0, \dots)$

$\boldsymbol{\pi}(f, \mathbf{x}) \equiv$ algorithm

- ℓ_2 -stable *signals* defined as $\mathbf{x} \in \ell_2 \iff \sum_{t=0}^{\infty} |x_t|_2^2 < \infty$
- \mathcal{L}_2 -stable operators: $\mathbf{A}(\mathbf{x}) \in \ell_2, \forall \mathbf{x} \in \ell_2$





Problem Formulation (1)

- Denote as \mathcal{S}_β the set of **non-convex** cost functions with β -Lipschitz gradients:

$$f \in \mathcal{S}_\beta \iff |\nabla f(x) - \nabla f(y)| \leq \beta|x - y|$$

- Signal space notation:

$$\begin{cases} x_{t+1} = x_t + \pi_t(f, x_t) \\ x_0 = x_0 \end{cases} \iff z\mathbf{x} = \mathbf{x} + \boldsymbol{\pi}(f, \mathbf{x}) + z\boldsymbol{\delta}_0$$

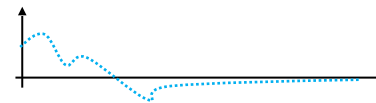
$z \equiv$ one-time-step shift

$\mathbf{x} = (x_0, x_1, x_2, \dots)$

$\boldsymbol{\delta}_0 = (x_0, 0, 0, \dots)$

$\boldsymbol{\pi}(f, \mathbf{x}) \equiv$ algorithm

- ℓ_2 -stable *signals* defined as $\mathbf{x} \in \ell_2 \iff \sum_{t=0}^{\infty} |x_t|_2^2 < \infty$
- \mathcal{L}_2 -stable operators: $\mathbf{A}(\mathbf{x}) \in \ell_2, \forall \mathbf{x} \in \ell_2$



- Algorithms that convergence for a function

$$f \in \mathcal{S}_\beta$$

$$\begin{aligned} \text{“}\boldsymbol{\pi} \text{ is sum-square convergent for } f\text{”} &\iff \|\nabla f(\mathbf{x})\|^2 < \infty \quad \|\boldsymbol{\pi}(f, \mathbf{x})\|^2 < \infty, \\ \text{“}\boldsymbol{\pi} \in \Sigma(f)\text{”} & \quad \left(= \sqrt{\sum_{t=0}^{\infty} |\nabla f(x_t)|^2} \right) \end{aligned}$$



Problem Formulation (2)

- **Problem:** design an *optimal robustly convergent* algorithm

$$\begin{aligned}
 \min_{\pi} \quad & \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\
 \text{s. t.} \quad & x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\
 & \pi(f, \mathbf{x}) \in \Sigma(f), \quad \forall f \in \mathcal{S}_\beta,
 \end{aligned}$$



Problem Formulation (2)

- **Problem:** design an *optimal robustly convergent* algorithm

$$\begin{aligned} \min_{\pi} \quad & \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t.} \quad & x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\ & \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta, \end{aligned}$$

- **Constraints:** limit the search to algorithms that converge for all smooth functions
 - E.g.: $\pi(f, x_t) = -\eta \nabla f(x_t)$ belongs to $\Sigma(f)$ for every $f \in \mathcal{S}'_\beta$ if $\eta < \beta^{-1}$ ^[1]
 - Beyond gradient descent, **case-by-case mathematical analysis**

[1] D. P. Bertsekas and J. N. Tsitsiklis, "Gradient convergence in gradient methods with errors," SIAM Journal on Optimization, vol. 10, no. 3, pp. 627–642, 2000.



- **Problem:** design an *optimal robustly convergent* algorithm

$$\begin{aligned} \min_{\pi} \quad & \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t.} \quad & x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\ & \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta, \end{aligned}$$

- **Constraints:** limit the search to algorithms that converge for all smooth functions

- E.g.: $\pi(f, x_t) = -\eta \nabla f(x_t)$ belongs to $\Sigma(f)$ for every $f \in \mathcal{S}'_\beta$ if $\eta < \beta^{-1}$ ^[1]
 - Beyond gradient descent, **case-by-case mathematical analysis**

- **Performance metric:**

$$\text{AlgoPerf}(f, \mathbf{x}) = \sum_{t=0}^T \alpha_t |\nabla f(x_t)|^2 + \gamma_t f(x_t),$$

to balance fast convergence and solution quality

[1] D. P. Bertsekas and J. N. Tsitsiklis, "Gradient convergence in gradient methods with errors," SIAM Journal on Optimization, vol. 10, no. 3, pp. 627–642, 2000.



Main result 1: a separation principle

- Consider the following algorithm:

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v}$$

- **Key idea:** separate algorithms into two addends



Main result 1: a separation principle

- Consider the following algorithm:

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v}$$

- Key idea:** separate algorithms into two addends
 - Gradient descent, aiming to ensure convergence on \mathcal{S}_β
 - \mathbf{v} is a performance-boosting term to be designed, *without ruining convergence*



Main result 1: a separation principle

- Consider the following algorithm:

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v}$$

- Key idea:** separate algorithms into two addends
 - Gradient descent, aiming to ensure convergence on \mathcal{S}_β
 - \mathbf{v} is a performance-boosting term to be designed, *without ruining convergence*

Result: The algorithm above is sum-square convergent for all $f \in \mathcal{S}_\beta$ if:

- $\eta < \beta^{-1}$
- \mathbf{v} is a bounded-energy ℓ_2 signal, that is, $\sum_{t=0}^{\infty} \|v_t\|^2 \leq \infty$



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - **Nonlinear system theory**: check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounded-Input-Bounded-Output...

Proof sketch^[2]

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - **Nonlinear system theory**: check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounded-Input-Bounded-Output...

Proof sketch^[2]

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - **Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounded-Input-Bounded-Output...

- β - smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} |x_{t+1} - x_t|^2$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - **Nonlinear system theory**: check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounded-Input-Bounded-Output...

- β - smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} |x_{t+1} - x_t|^2$$
- Highlight the role of \mathbf{v} :

$$\frac{\overbrace{2\eta\epsilon(1 - \beta\eta) - 1}^{>0}}{2\epsilon} |\nabla_t|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) |v_t|^2.$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounder-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounder-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furiere, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounder-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounder-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounded-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 1: sketch of the proof

- Nonlinear dynamics: $x_{t+1} = x_t - \eta \nabla f(x_t) + v_t$
- Exponential convergence if $v_t = 0 \dots$ What if $\mathbf{v} \in \ell_2$?
 - Nonlinear system theory:** check **Input-to-state (ISS) stability** property! [1]
 - Non-trivial: multiple equilibria, stronger than Bounder-Input-Bounded-Output...

- β -smoothness of the gradients: **Proof sketch**[2]

$$f(x_{t+1}) \leq f(x_t) + \nabla_t^\top (x_{t+1} - x_t) + \frac{\beta}{2} \|x_{t+1} - x_t\|^2$$
- Highlight the role of \mathbf{v} :

$$\underbrace{2\eta\epsilon(1 - \beta\eta) - 1}_{>0} \|\nabla_t\|^2 \leq f(x_t) - f(x_{t+1}) + \left(\frac{\epsilon}{2} + \beta\right) \|v_t\|^2.$$
- Summing over t :

$$\sum_{t=0}^{\infty} \|\nabla_t\|^2 \leq M (f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)) + N \sum_{t=0}^{\infty} \|v_t\|^2.$$

$$\implies \|\nabla f(\mathbf{x})\|_2^2 \leq K_1 + K_2 \|\mathbf{v}\|_2^2 \leq \infty$$

[1] Khalil, Hassan K. *Control of nonlinear systems*. Prentice Hall, New York, NY, 2002.

[2] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024



Main result 2: completeness

(Crucial) question: can *any* algorithm that robustly converges on \mathcal{S}_β be expressed as

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v} \quad ?$$



Main result 2: completeness

(Crucial) question: can *any* algorithm that robustly converges on \mathcal{S}_β be expressed as

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v} \quad ?$$

Result: Let π be any robustly convergent algorithm:

$$\pi(f, \mathbf{x}) \in \Sigma(f), \quad \forall f \in \mathcal{S}_\beta.$$

Then, there exists an \mathcal{L}_2 operator \mathbf{V} such that the algorithm

$$\pi'(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{V}(x_0)$$

yields the same trajectories as π , for every x_0 and every $f \in \mathcal{S}_\beta$.



Main result 2: completeness

(Crucial) question: can *any* algorithm that robustly converges on \mathcal{S}_β be expressed as

$$\pi(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{v} \quad ?$$

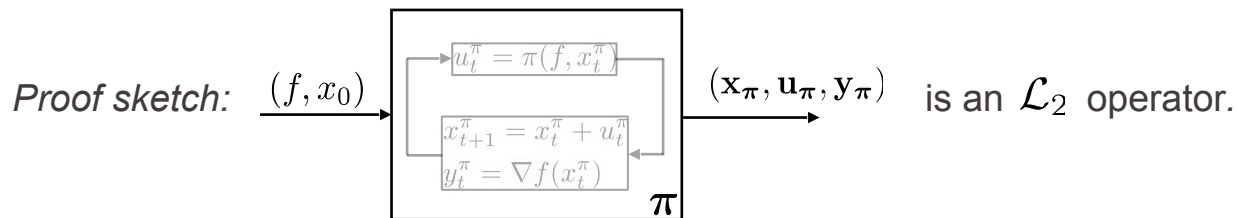
Result: Let π be any robustly convergent algorithm:

$$\pi(f, \mathbf{x}) \in \Sigma(f), \quad \forall f \in \mathcal{S}_\beta.$$

Then, there exists an \mathcal{L}_2 operator \mathbf{V} such that the algorithm

$$\pi'(f, \mathbf{x}) = -\eta \nabla f(\mathbf{x}) + \mathbf{V}(x_0)$$

yields the same trajectories as π , for every x_0 and every $f \in \mathcal{S}_\beta$.



Then $\mathbf{V}(x_0) = \eta \nabla f(\mathbf{x}_\pi(x_0)) + \mathbf{u}_\pi(x_0)$ achieves desired behavior, and also lies in \mathcal{L}_2



Importance of main results 1 and 2

One-to-one parametrization

$$-\eta \nabla f(\mathbf{x}) + \mathbf{v} \in \mathcal{L}_2 \quad \equiv \quad \pi(f, \mathbf{x}) \in \Sigma(f), \quad \forall f \in \mathcal{S}_\beta.$$



Importance of main results 1 and 2

One-to-one parametrization

$$-\eta \nabla f(\mathbf{x}) + \mathbf{V} \in \mathcal{L}_2 \quad \equiv \quad \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta.$$

- Unconstrained learning over all converging algorithms

$$\begin{aligned} \min_{\mathbf{V} \in \mathcal{L}_2} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t. } x_{t+1} = x_t - \eta \nabla f(x_t) + V_t(x_0), \end{aligned} \quad \equiv \quad \begin{aligned} \min_{\pi} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t. } x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\ \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta, \end{aligned}$$



Importance of main results 1 and 2

One-to-one parametrization

$$-\eta \nabla f(\mathbf{x}) + \mathbf{V} \in \mathcal{L}_2 \quad \equiv \quad \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta.$$

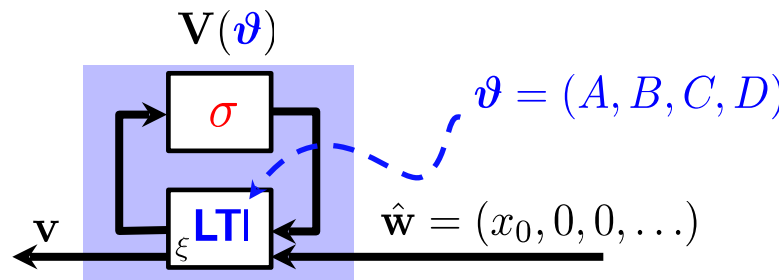
- Unconstrained learning over all converging algorithms

$$\begin{aligned} \min_{\mathbf{V} \in \mathcal{L}_2} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s.t. } x_{t+1} = x_t - \eta \nabla f(x_t) + V_t(x_0), \end{aligned} \quad \equiv \quad \begin{aligned} \min_{\pi} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s.t. } x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\ \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta, \end{aligned}$$

How do we search over these operators?



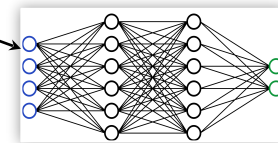
Unconstrained parametrizations of $V \in \mathcal{L}_2$



- Expressive models including

$$\xi_t = \hat{A}\xi_{t-1} + \hat{B} \text{NN}(\xi_{t-1}, \hat{\mathbf{w}}_t)$$

$$u_t = \hat{C}\xi_t + \hat{D} \text{NN}(\xi_t, \hat{\mathbf{w}}_t)$$



- $V(\vartheta) \in \mathcal{L}_2$ if there is a storage function $S(\xi) = \xi^\top P \xi$ verifying

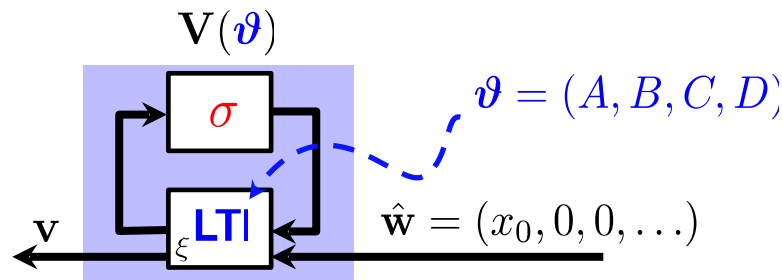
$$S(\xi_{t+1}) - S(\xi_t) \leq \gamma^2 |x_0| - |u_t|$$

[1] Kim, K. K., E. Ríos Patrón, and R. D. Braatz. *Standard representation and unified stability analysis for dynamic artificial neural network models*, Transactions on Neural Networks, 2018

[2] Revay, M, R. Wang, and I.R. Manchester. *Recurrent equilibrium networks: Flexible dynamic models with guaranteed stability and robustness*. IEEE TAC 2023



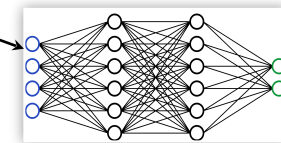
Unconstrained parametrizations of $V \in \mathcal{L}_2$



- Expressive models including

$$\xi_t = \hat{A}\xi_{t-1} + \hat{B} \text{NN}(\xi_{t-1}, \hat{\mathbf{w}}_t)$$

$$u_t = \hat{C}\xi_t + \hat{D} \text{NN}(\xi_t, \hat{\mathbf{w}}_t)$$



- $V(\vartheta) \in \mathcal{L}_2$ if there is a storage function $S(\xi) = \xi^\top P \xi$ verifying

$$S(\xi_{t+1}) - S(\xi_t) \leq \gamma^2 |x_0| - |u_t|$$

Free parametrization^[2]: explicit map $\Theta \rightarrow (\vartheta, P)$ such that $V(\vartheta) \in \mathcal{L}_2$, for any $\Theta \in \mathbb{R}^d$

[1] Kim, K. K., E. Ríos Patrón, and R. D. Braatz. *Standard representation and unified stability analysis for dynamic artificial neural network models*, Transactions on Neural Networks, 2018

[2] Revay, M, R. Wang, and I.R. Manchester. *Recurrent equilibrium networks: Flexible dynamic models with guaranteed stability and robustness*. IEEE TAC 2023



Unconstrained parametrizations of $V \in \mathcal{L}_2$

Free parametrization^[2]: explicit map $\Theta \rightarrow (\vartheta, P)$ such that $V(\vartheta) \in \mathcal{L}_2$, for any $\Theta \in \mathbb{R}^d$



$$\begin{aligned} \min_{\Theta \in \mathbb{R}^d} \quad & \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t.} \quad & x_{t+1} = x_t - \eta \nabla f(x_t) + V_t(x_0, \Theta) \end{aligned}$$

[1] Kim, K. K., E. Ríos Patrón, and R. D. Braatz. *Standard representation and unified stability analysis for dynamic artificial neural network models*, Transactions on Neural Networks, 2018

[2] Revay, M, R. Wang, and I.R. Manchester. *Recurrent equilibrium networks: Flexible dynamic models with guaranteed stability and robustness*. IEEE TAC 2023



Unconstrained parametrizations of $V \in \mathcal{L}_2$

Free parametrization^[2]: explicit map $\Theta \rightarrow (\vartheta, P)$ such that $V(\vartheta) \in \mathcal{L}_2$, for any $\Theta \in \mathbb{R}^d$



$$\begin{aligned} \min_{\Theta \in \mathbb{R}^d} \quad & \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t.} \quad & x_{t+1} = x_t - \eta \nabla f(x_t) + V_t(x_0, \Theta) \end{aligned}$$



Challenges for effective L2O

1. **Insufficient features**
 - Learning to map $x_0 \in \mathbb{R}^d \rightarrow \mathbf{v}^* \in \mathbb{R}^\infty$ is hard
2. **Full gradient measurements** for convergence

[1] Kim, K. K., E. Ríos Patrón, and R. D. Braatz. *S*
network models, Transactions on Neural Networks

[2] Revay, M, R. Wang, and I.R. Manchester. *Recu*
robustness. IEEE TAC 2023



1) Enriching the input features with $[I, f, \nabla f](x_t)$

- Parametrize additional signal with rich input features (**not in ℓ_2**):

$$\omega(\varphi) = \Omega(\mathbf{x}, f(\mathbf{x}), \nabla f(\mathbf{x}), \varphi)$$

- Construct a combined signal $\mathbf{z}(\varphi, \Theta)$ which is ℓ_2 by construction:

$$z_t(f, \nabla f, x_{t:0}, \varphi, \Theta) = \frac{\omega_t(f, \nabla f, x_{t:0}, \varphi)}{|\omega_t(f, \nabla f, x_{t:0}, \varphi)|} |v_t(x_0, \Theta)|$$



1) Enriching the input features with $[I, f, \nabla f](x_t)$

- Parametrize additional signal with rich input features (**not in ℓ_2**):

$$\omega(\varphi) = \Omega(\mathbf{x}, f(\mathbf{x}), \nabla f(\mathbf{x}), \varphi)$$

- Construct a combined signal $\mathbf{z}(\varphi, \Theta)$ which is ℓ_2 by construction:

$$z_t(f, \nabla f, x_{t:0}, \varphi, \Theta) = \frac{\omega_t(f, \nabla f, x_{t:0}, \varphi)}{|\omega_t(f, \nabla f, x_{t:0}, \varphi)|} |v_t(x_0, \Theta)|$$

Preserves one-to-one parametrization!

$$-\eta \nabla f(\mathbf{x}) + \frac{\Omega(f, \nabla f, \mathbf{x})}{|\Omega(f, \nabla f, \mathbf{x})|} |\mathbf{V}(x_0)|$$

\equiv

$$\pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta.$$

Proof of sufficiency: by construction

Proof of necessity: reduces to previous case by choosing $\Omega(f, \nabla f, \mathbf{x}) = \mathbf{V}(x_0)$



2) Convergence with incomplete gradients

- Typical scenario in machine learning

$$f(x) = \sum_{i \in \text{batches}} f_i(x) \quad \Rightarrow \quad \underline{\text{access to } \nabla f_i(x) \text{ only}}$$



2) Convergence with incomplete gradients

- Typical scenario in machine learning

$$f(x) = \sum_{i \in \text{batches}} f_i(x) \quad \Rightarrow \quad \underline{\text{access to } \nabla f_i(x) \text{ only}}$$

Result: Convergence (*asymptotic*) is preserved using the update rule

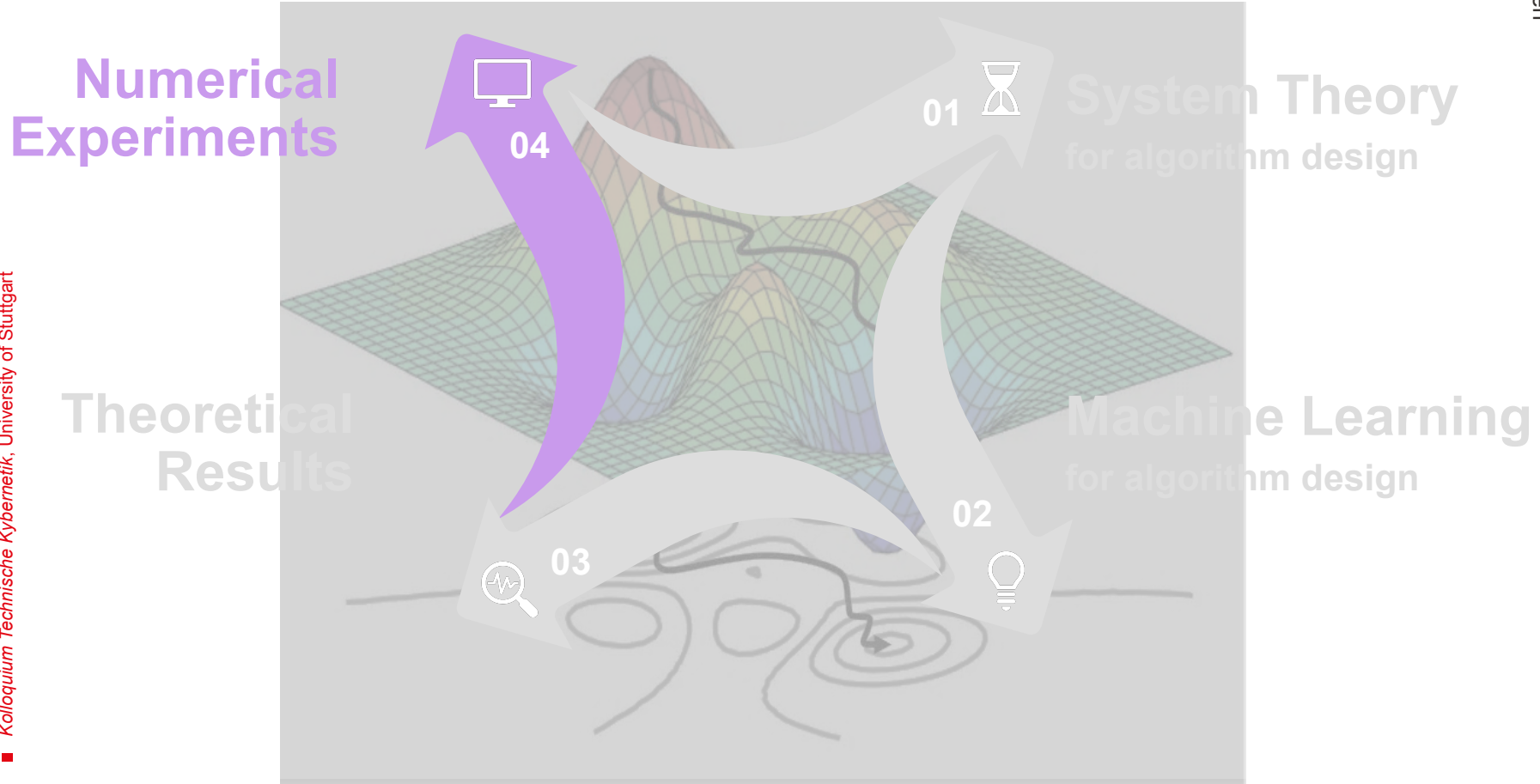
$$\pi_t(x_t) = -\eta_t \nabla f_t(x_t) + v_t$$

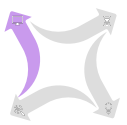
as long as $\eta \in \ell_2$ and $|v_t| \leq C\eta_t |\nabla f(x_t)|$.

- Convergence** guarantees of L2O algorithms **compatible with ML tasks!**
- Only sufficient, only asymptotic convergence**
 - necessity* relies on full gradient information
 - asymptotic convergence*: consistent with stochastic gradient descent results

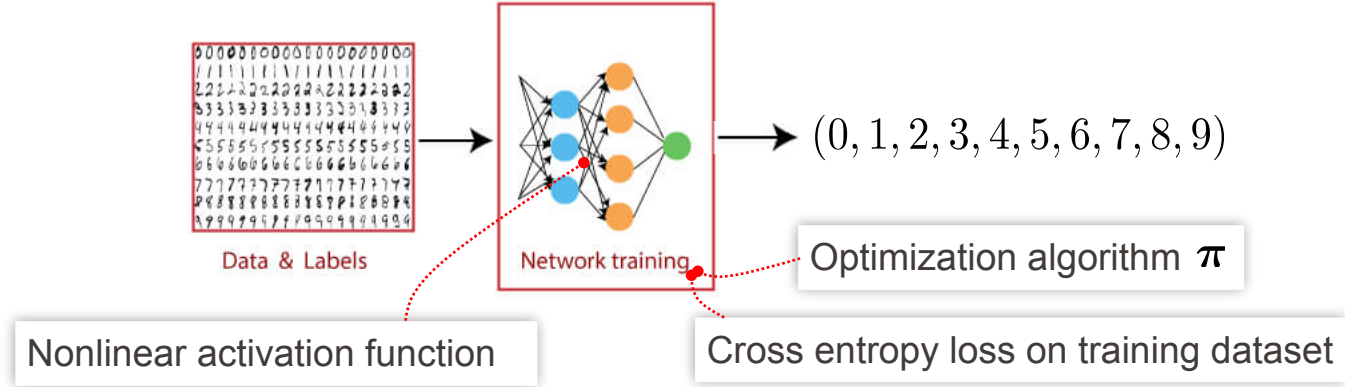
[1] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024

Design of new optimization algorithms

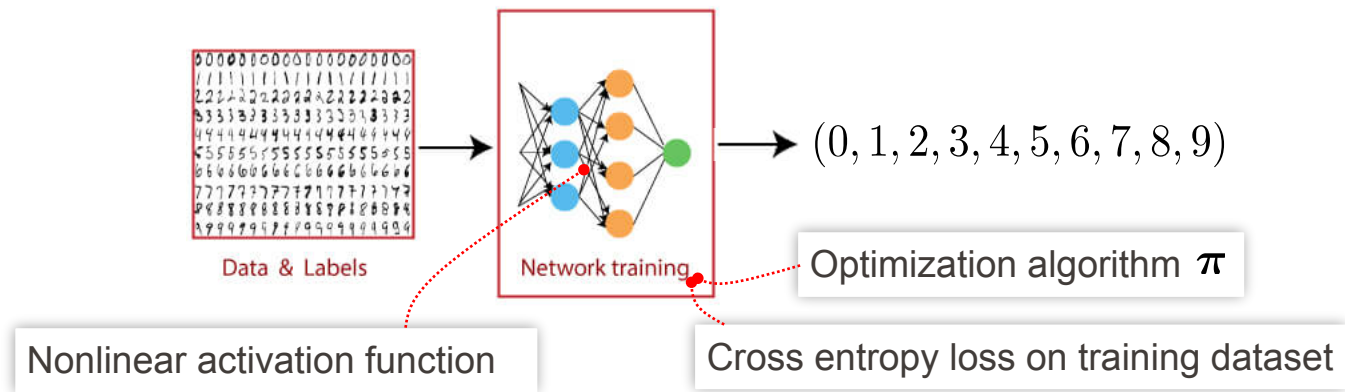





Numerical experiments: image classification^[1]



[1] LeCun, Y., Cortes, C. and Burges, C.J.C. (1998) The MNIST Database of Handwritten Digits. New York, USA. <http://yann.lecun.com/exdb/mnist/>

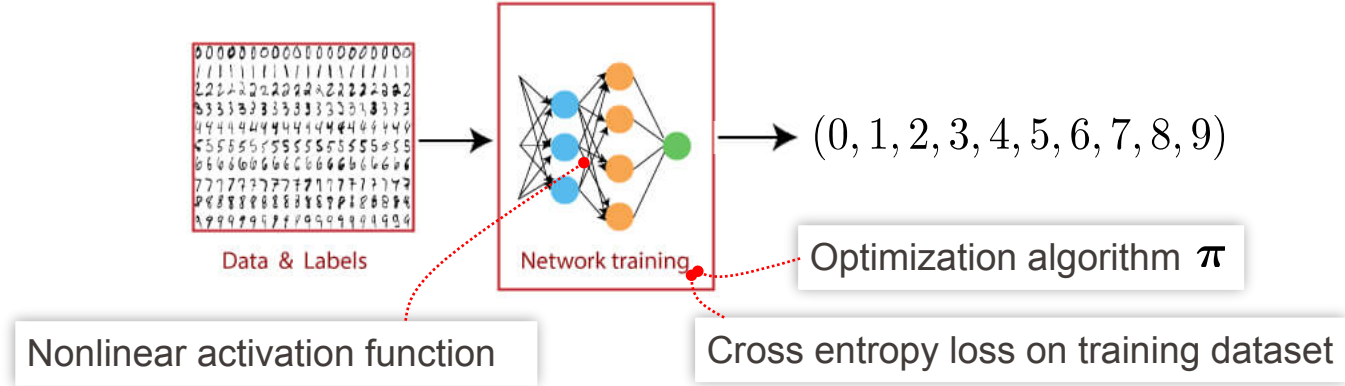


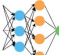
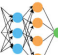
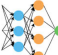
- **Algorithm training:** repeatedly optimize  from uniformly distributed initial weights
 - AlgoPerf: average algorithm performance (speed and solution quality)
 - We fix \tanh activation function on 

[1] LeCun, Y., Cortes, C. and Burges, C.J.C. (1998) The MNIST Database of Handwritten Digits. New York, USA. <http://yann.lecun.com/exdb/mnist/>

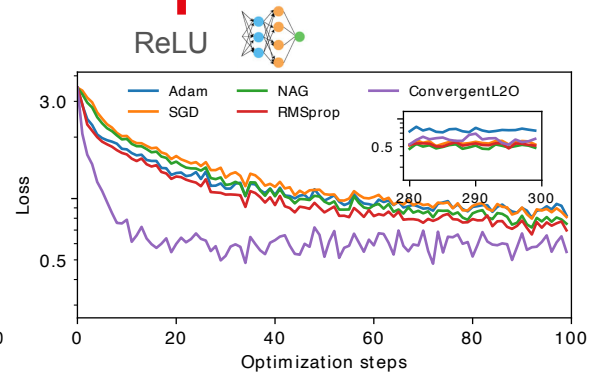
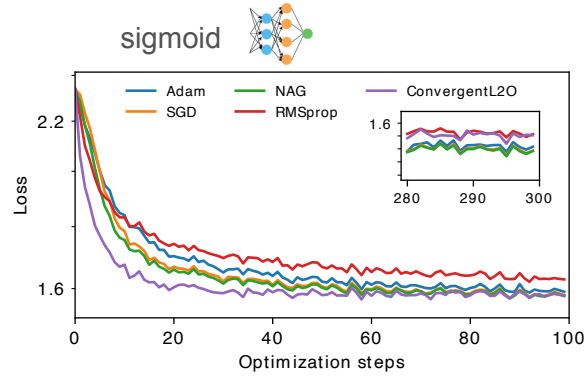
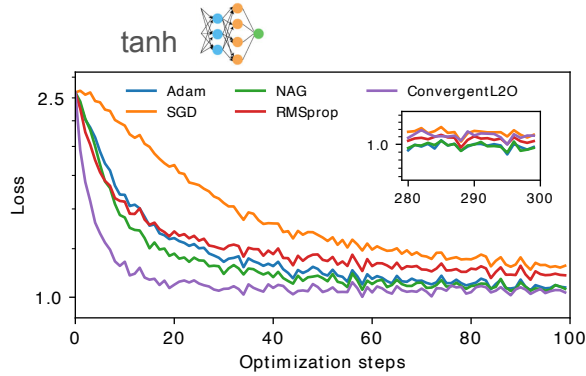


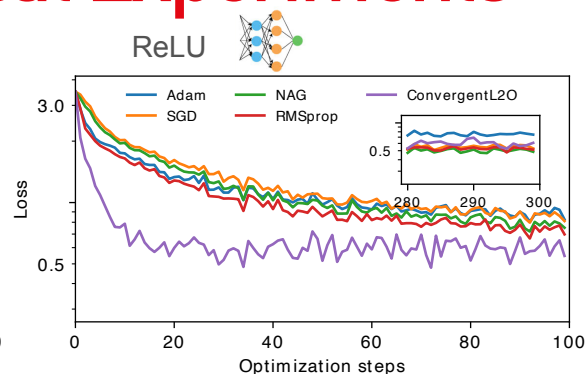
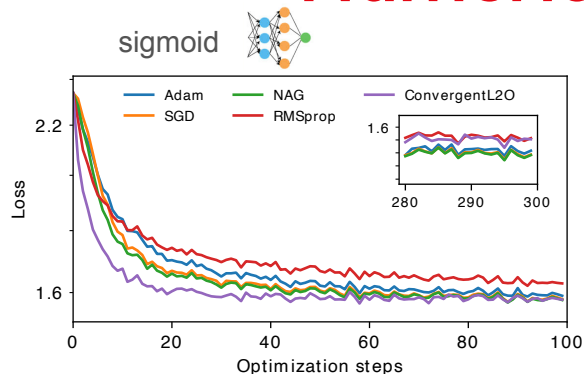
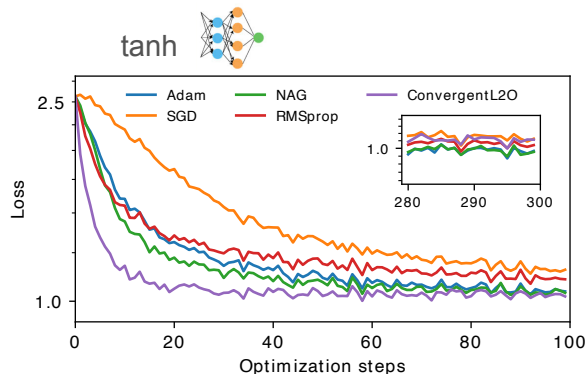
Numerical experiments: image classification^[1]



- **Algorithm training:** repeatedly optimize  from uniformly distributed initial weights
 - AlgoPerf: average algorithm performance (speed and solution quality)
 - We fix \tanh activation function on 
- **Algorithm testing:** compare with classical fine-tuned optimizers
 - Generalization to different activation functions on 
 - Convergence guaranteed by design

[1] LeCun, Y., Cortes, C. and Burges, C.J.C. (1998) The MNIST Database of Handwritten Digits. New York, USA. <http://yann.lecun.com/exdb/mnist/>

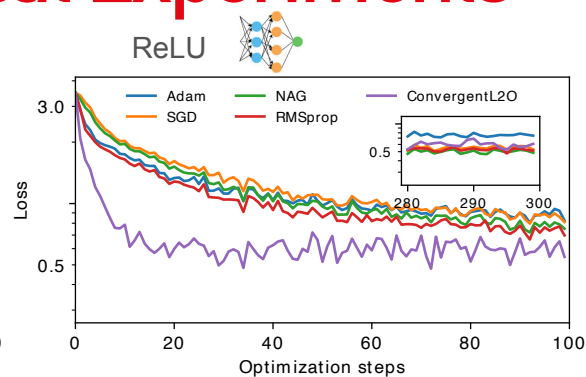
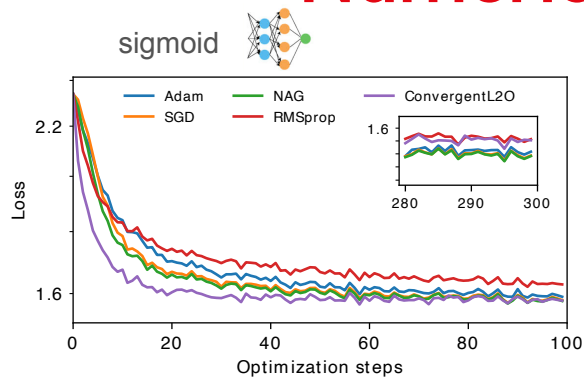
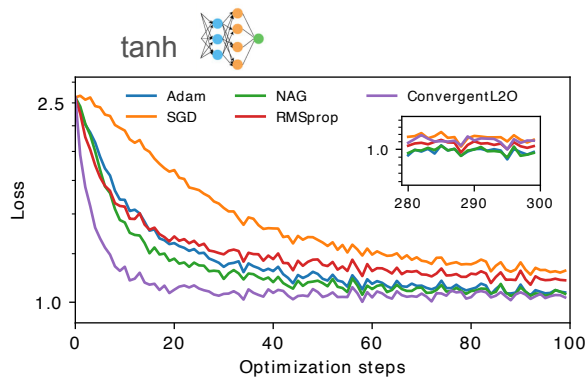




Classification accuracy on test data

Step $t = 20$	tanh	sigmoid	ReLU
Adam	$71.7 \pm 5.1\%$	$76.1 \pm 3.1\%$	$52.7 \pm 11.1\%$
SGD	$44.9 \pm 4.2\%$	$79.7 \pm 1.9\%$	$49.8 \pm 9.3\%$
NAG	$79.7 \pm 1.4\%$	$81.1 \pm 1.5\%$	$52.7 \pm 10.2\%$
RMSprop	$69.4 \pm 2.9\%$	$72.8 \pm 2.3\%$	$61.1 \pm 8.9\%$
ConvergentL2O	$87.0 \pm 0.5\%$	$86.8 \pm 0.6\%$	$86.3 \pm 0.6\%$
LSTM	$82.2 \pm 0.1\%$	$83.3 \pm 0.1\%$	$88.3 \pm 0.0\%$

Step $t = 300$	tanh	sigmoid	ReLU
Adam	$89.5 \pm 0.5\%$	$89.6 \pm 0.3\%$	$70.3 \pm 12.2\%$
SGD	$87.4 \pm 0.4\%$	$89.3 \pm 0.3\%$	$80.6 \pm 8.1\%$
NAG	$89.4 \pm 0.2\%$	$89.4 \pm 0.2\%$	$82.2 \pm 7.6\%$
RMSprop	$87.6 \pm 2.1\%$	$88.5 \pm 0.4\%$	$81.5 \pm 7.5\%$
ConvergentL2O	$88.5 \pm 0.2\%$	$88.4 \pm 0.3\%$	$87.7 \pm 0.2\%$
LSTM	$81.4 \pm 0.0\%$	$81.4 \pm 0.0\%$	$88.3 \pm 0.0\%$

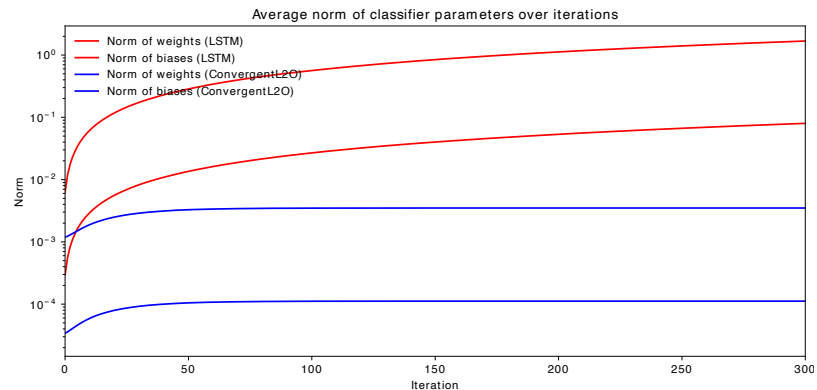


Classification accuracy on test data

Step $t = 20$	tanh	sigmoid	ReLU
Adam	71.7 ± 5.1%	76.1 ± 3.1%	52.7 ± 11.1%
SGD	44.9 ± 4.2%	79.7 ± 1.9%	49.8 ± 9.3%
NAG	79.7 ± 1.4%	81.1 ± 1.5%	52.7 ± 10.2%
RMSprop	69.4 ± 2.9%	72.8 ± 2.3%	61.1 ± 8.9%
ConvergentL2O	87.0 ± 0.5%	86.8 ± 0.6%	86.3 ± 0.6%
LSTM	82.2 ± 0.1%	83.3 ± 0.1%	88.3 ± 0.0%

Step $t = 300$	tanh	sigmoid	ReLU
Adam	89.5 ± 0.5%	89.6 ± 0.3%	70.3 ± 12.2%
SGD	87.4 ± 0.4%	89.3 ± 0.3%	80.6 ± 8.1%
NAG	89.4 ± 0.2%	89.4 ± 0.2%	82.2 ± 7.6%
RMSprop	87.6 ± 2.1%	88.5 ± 0.4%	81.5 ± 7.5%
ConvergentL2O	88.5 ± 0.2%	88.4 ± 0.3%	87.7 ± 0.2%
LSTM	81.4 ± 0.0%	81.4 ± 0.0%	88.3 ± 0.0%

LSTM from [1] diverges!



- **Automated synthesis of high-performance optimization algorithms**
 - With theoretical guarantees by design
 - For non-convex functions
- **What's next?**
 - *Analyzing generalization capabilities (robust performance)*
 - Online and distributed optimization
 - Stronger convergence guarantees for classes of non-convex functions
 - Applications in optimal control, e.g., NMPC

Nonlinear system
theory

L2O

- **Automated synthesis of high-performance optimization algorithms**

- With theoretical guarantees by design
- For non-convex functions

Nonlinear system
theory

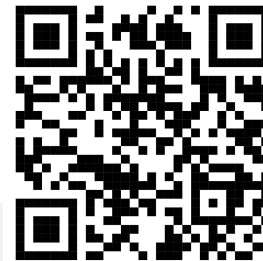


L2O

- **What's next?**

- *Analyzing generalization capabilities (robust performance)*
- Online and distributed optimization
- Stronger convergence guarantees for classes of non-convex functions
- Applications in optimal control, e.g., NMPC

Check out the paper:

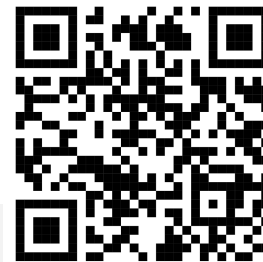


Thank you! Questions?

$$\begin{aligned} & \min_{\Omega, \mathbf{V} \in \mathcal{L}_2} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t. } & x_{t+1} = x_t - \eta \nabla f(x_t) + \frac{\omega_t(f, \nabla f, x_t)}{|\omega_t(f, \nabla, x_t)|} |v_t(x_0)|, \end{aligned}$$

Thank you! Questions?

Check out the paper:



$$\begin{aligned} & \min_{\Omega, \mathbf{V} \in \mathcal{L}_2} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t. } & x_{t+1} = x_t - \eta \nabla f(x_t) + \frac{\omega_t(f, \nabla f, x_t)}{|\omega_t(f, \nabla, x_t)|} |v_t(x_0)|, \end{aligned}$$

$$\begin{aligned} & \min_{\pi} \mathbb{E}_{f \sim \text{Examples}, x_0 \sim \mathcal{X}_0} [\text{AlgoPerf}(f, \mathbf{x})] \\ \text{s. t. } & x_{t+1} = x_t + \pi_t(f, x_{t:0}), \\ & \pi(f, \mathbf{x}) \in \Sigma(f), \forall f \in \mathcal{S}_\beta, \end{aligned}$$

[1] Andrea Martin and Luca Furieri, *Learning to optimize with convergence guarantees using nonlinear system theory*, IEEE Control System Letters, 2024