

Thesis Defense

# Differentiable Optimization-Based Modeling for Machine Learning

Brandon Amos • Carnegie Mellon University

Thesis Committee:

J. Zico Kolter, Chair

Barnabás Póczos

Jeff Schneider

Vladlen Koltun (Intel Labs)

# My Ph.D. Contributions

■ Secondary Contribution

[CMU 2016] OpenFace

[ICML 2016] Collapsed Variational Inference for SPNs

[ICML 2017] Input Convex Neural Networks

[ICML 2017] OptNet

[NeurIPS 2017] Task-Based Model Learning

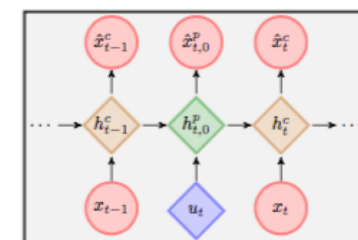
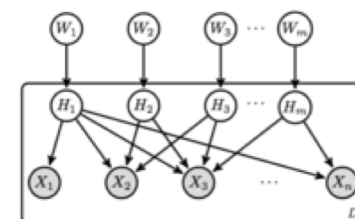
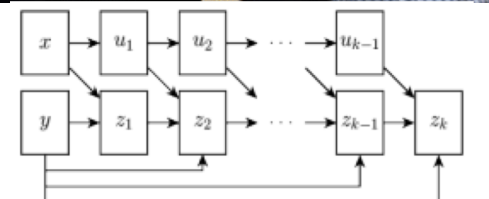
[ICLR 2018] Learning Awareness Models

[NeurIPS 2018] Imperfect-Information Game Solving

[NeurIPS 2018] Differentiable MPC

The Limited Multi-Label Projection Layer

Differentiable `cvxpy` Optimization Layers



# This Talk

■ Secondary Contribution

[CMU 2016] OpenFace

[ICML 2016] Collapsed Variational Inference for SPNs

[ICML 2017] Input Convex Neural Networks

[ICML 2017] OptNet

[NeurIPS 2017] Task-Based Model Learning

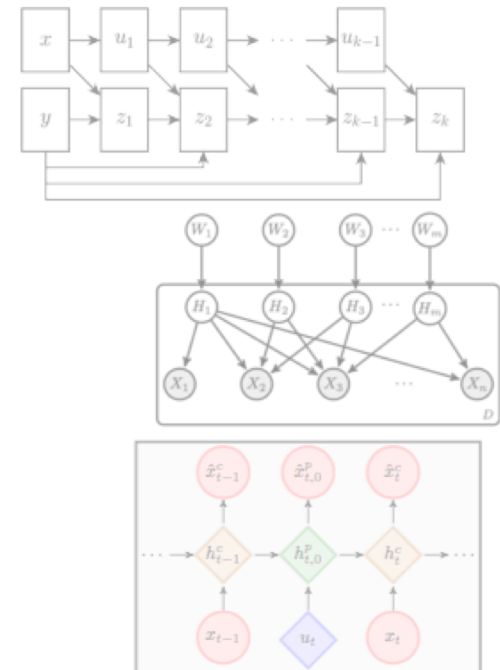
[ICLR 2018] Learning Awareness Models

[NeurIPS 2018] Imperfect-Information Game Solving

[NeurIPS 2018] Differentiable MPC

The Limited Multi-Label Projection Layer

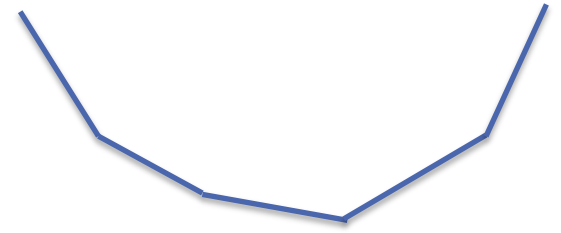
Differentiable `cvxpy` Optimization Layers



# Input Convex Neural Networks

A quick glimpse

# Input Convex Neural Networks (ICNNS)



**Definition** Scalar-valued network  $f(x, y; \theta)$  such that  $f$  is **convex** in  $y$  for all values of  $x$  (note that these networks are still **not convex** in  $\theta = \{W_i, b_i\}$ )

We can efficiently **optimize over some inputs** to the network **given other inputs**

Efficiently captures dependencies in the output space for prediction

It turns out, we don't need very many restrictions on the network to achieve this property

# How to achieve input convexity?

Most networks can be “trivially” modified to **guarantee input convexity**

Consider a simple **feedforward ReLU network**:

$$z_{i+1} = \max\{0, W_i z_i + b_i\}, \quad i = 1, \dots, k$$
$$f(y; \theta) = z_{k+1}, \quad z_1 = y$$

**Proposition.**  $f$  is convex in  $y$  provided that the  $W_i$  are non-negative for  $i > 1$

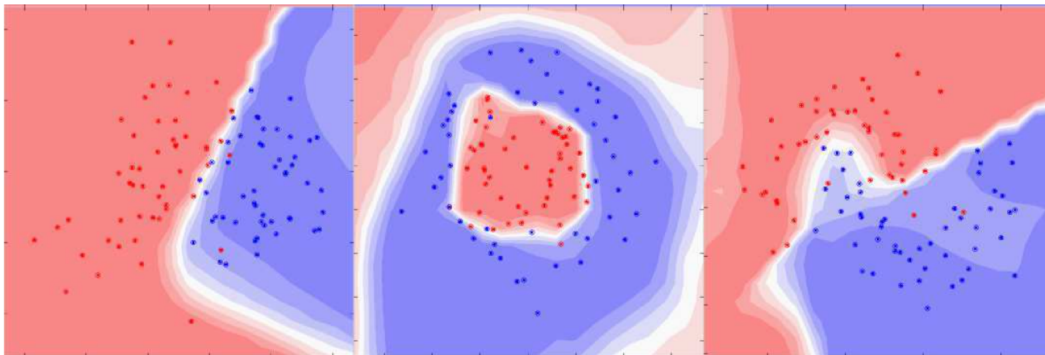
More generally, any activation function that is convex and non-decreasing also has this property.

# Is convexity restrictive?

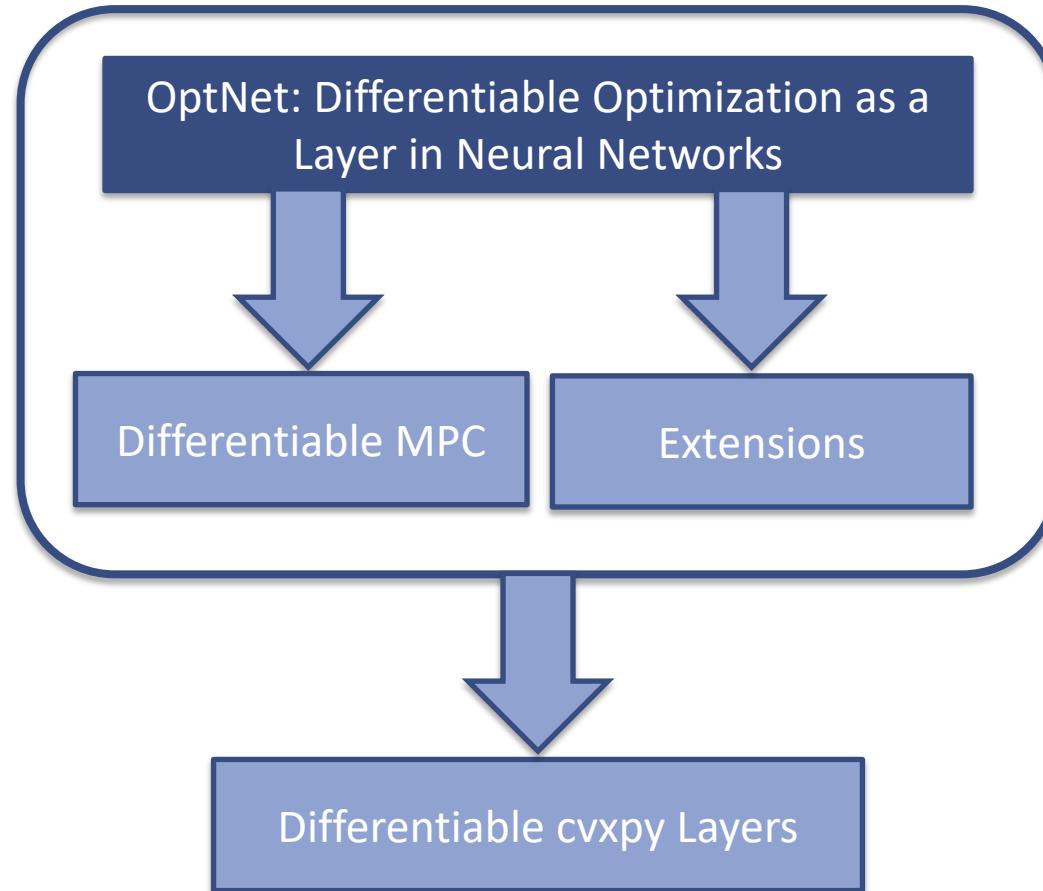
Yes (by definition, the functions are restricted to be convex), but not that bad in practice

Proposition. ICNNs trivially subsume any feedforward network  $\tilde{f}(x)$  with the network  $f(x, y) = (y - \tilde{f}(x))^2$

More complex convex portion adds additional structure over  $y$ , which can still be “easily” optimized over

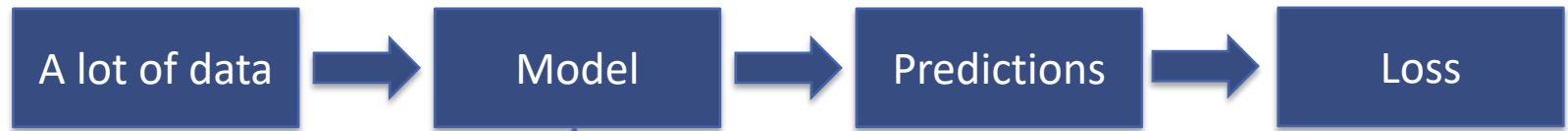


# Overview for the remainder of this talk





# Today's Machine Learning Systems



**Current Primitive Operations:** Linear maps, convolutions, activation functions, random sampling, simple projections (e.g. onto the simplex or Birkhoff polytope)

## How can the modeling part be improved?

Black-box neural networks don't work everywhere and when they fail, task-specific domain knowledge can provide useful modeling priors

My work mostly focuses on ways to **use optimization to inject domain knowledge into the modeling process**

# Optimization and Machine Learning

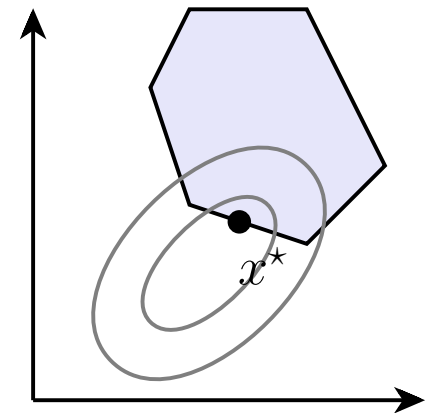
**Non-convex optimization is thriving** in machine learning for parameter optimization and architecture search. **This is not what this talk covers.**

**In this talk**, we argue that optimization is also a useful operation for inference and control.

We consider optimization as another potential **layer**, to be **composed with others**

Why? Optimization is an extremely powerful paradigm for decision-making.

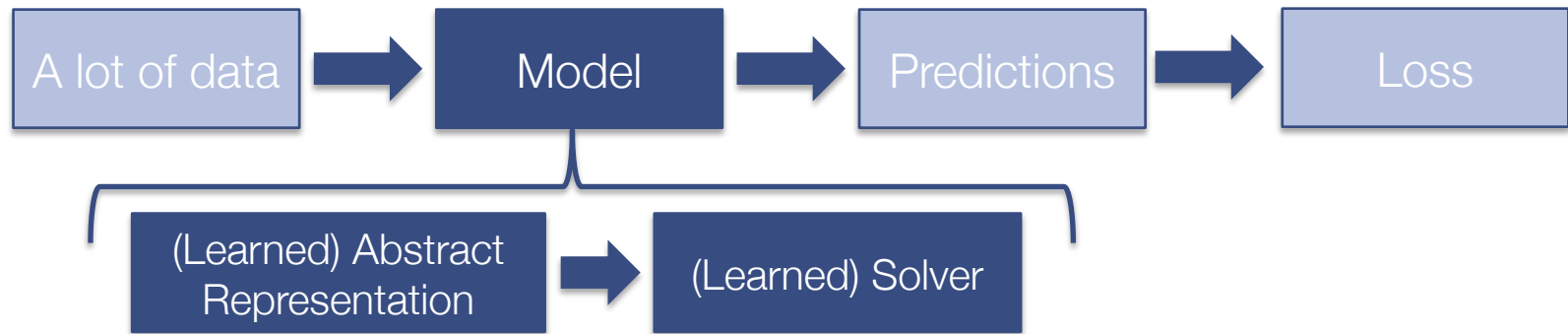
- Applications in **finance** (Markowitz portfolio optimization), **machine learning** (support vector machines), **control** (linear-quadratic model predictive control), **geometry** (projections onto polyhedra)



# Why is optimization a useful primitive operation in learning systems?

We have incomplete domain knowledge about what we want to model

- Fill in parts of the optimization problem that we know
- Use data to learn the parts that we don't



Also subsumes many standard layers (ReLU, sigmoid, softmax)

- We will show this later

# Convex optimization viewpoint of standard layers

ReLU

$$y = \max\{0, x\}$$

$$y^* = \underset{y}{\operatorname{argmin}} \|y - x\|_2^2$$

subject to  $y \geq 0$

Sigmoid

$$y = \frac{1}{1 + e^{-x}}$$

$$y^* = \underset{y}{\operatorname{argmin}} -y^T x - H_b(y)$$

subject to  $0 \leq y \leq 1$

Softmax

$$y_j = \frac{e^{x_j}}{\sum e^{x_k}}$$

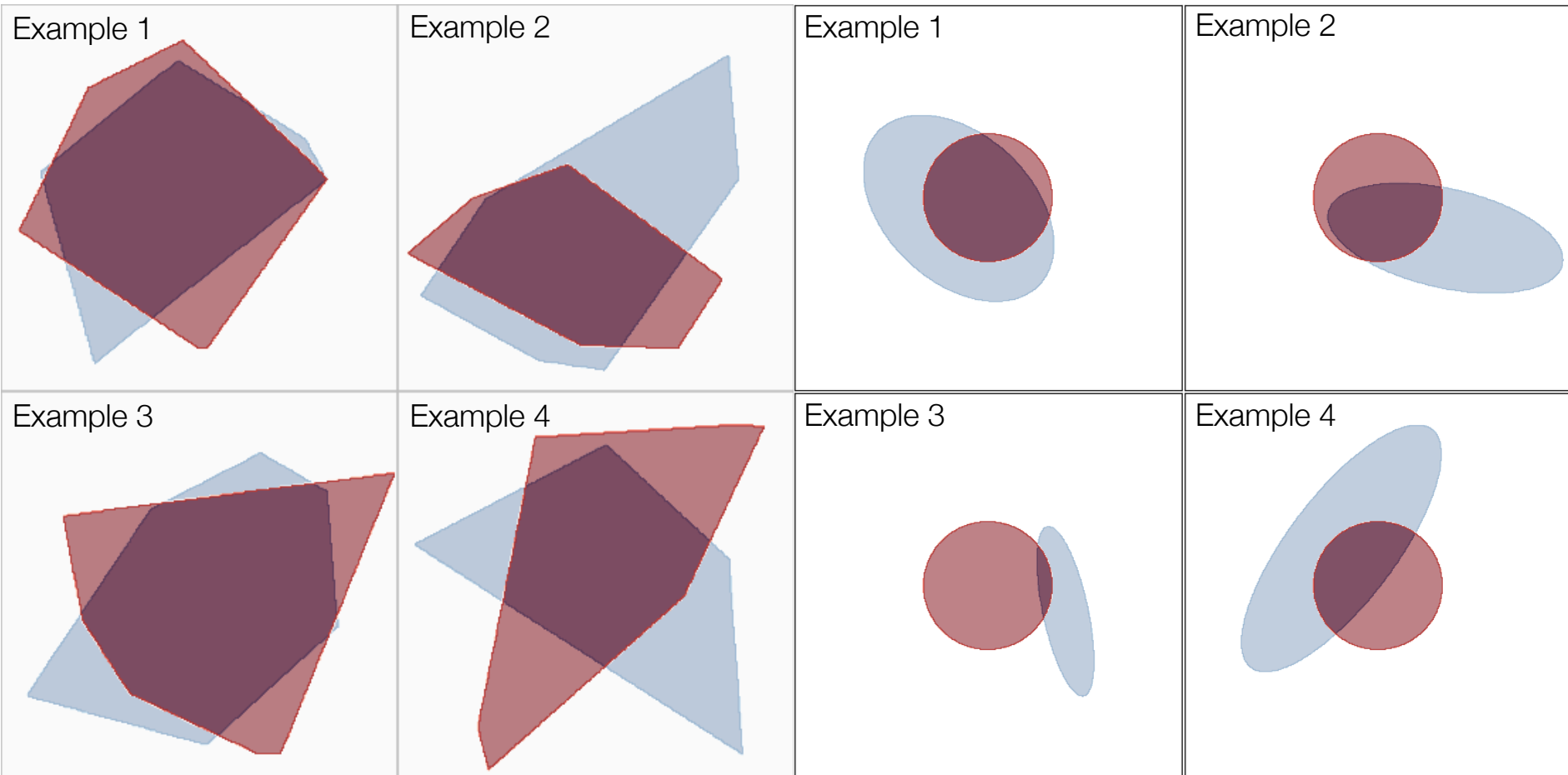
$$y^* = \underset{y}{\operatorname{argmin}} -y^T x - H(y)$$

subject to  $0 \leq y \leq 1$   
 $1^T y = 1$

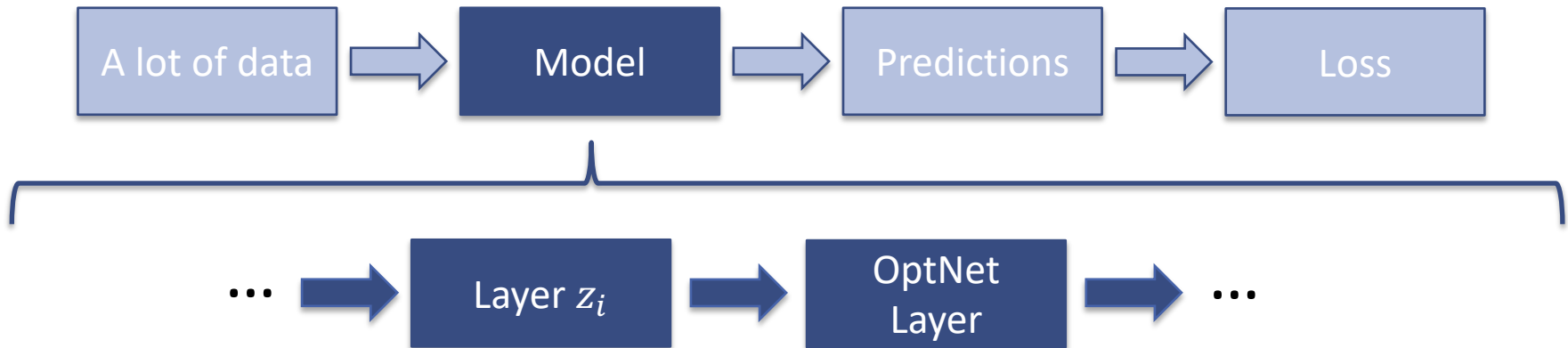
# OptNet Application: Modeling Constraints

■ True Constraint (Unknown to the model)

■ Constraint Predictions During Training



# The OptNet Layer



$$z_{i+1} = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$$

subject to  $A(z_i)z = b(z_i)$   
 $G(z_i)z \leq h(z_i)$

Parameters/Submodules :  $Q, q, A, b, G, h$

The matrix  $Q(z_i)$  depends on the previous layer  $z_i$

# Differentiating a quadratic argmin

Consider the optimization problem:

$$z^* = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^T Q z + q^T z$$

subject to  $Az = b, Gz \leq h$

From convex optimization theory, the **Karush-Kuhn-Tucker conditions** provide necessary and sufficient equations for optimality.

$$\begin{array}{ll} \text{stationarity} & Qz^* + q + A^T v^* + G^T \lambda^* = 0 \\ \text{primal feasibility} & Az^* - b = 0 \\ \text{complementary slackness} & D(\lambda^*)(Gz^* - h) = 0 \end{array}$$

To obtain  $\partial z^* / \partial \theta$  implicitly differentiate the KKT conditions.

This also works for any convex optimization problem (not just QPs)

# Implicitly differentiating the KKT conditions

Implicitly differentiate them (using differentials here):

$$dQz^* + qdz + dq + dA^T v^* + A^T dv + dG^T \lambda^* + G^T d\lambda = 0$$

$$dAz^* + Adz - db = 0$$

$$D(Gz^* - h)d\lambda + D(\lambda^*)(dGz^* + Gdz - dh) = 0$$

Fill in **desired differentials**, form a linear system, **solve for unknowns**

If done naively, takes **many** linear system solves

If done correctly, just requires a single solve to compute **all** gradients



# A Simple Application: Sudoku

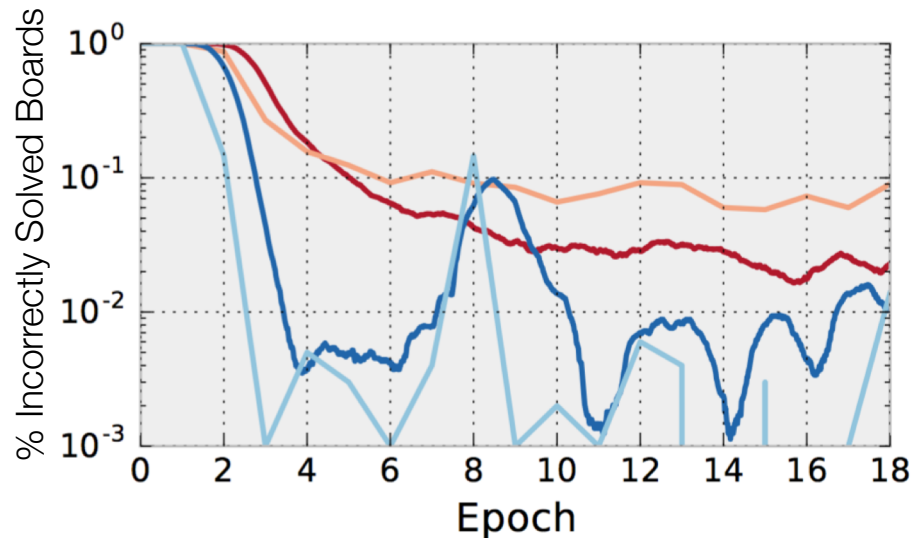
5	3			7				
6			1	9	5			
	9	8					6	
8				6				3
4			8		3			1
7				2				6
	6					2	8	
			4	1	9			5
				8			7	9

5	3	4	6	7	8	9	1	2
6	7	2	1	9	5	3	4	8
1	9	8	3	4	2	5	6	7
8	5	9	7	6	1	4	2	3
4	2	6	8	5	3	7	9	1
7	1	3	9	2	4	8	5	6
9	6	1	5	3	7	2	8	4
2	8	7	4	1	9	6	3	5
3	4	5	2	8	6	1	7	9

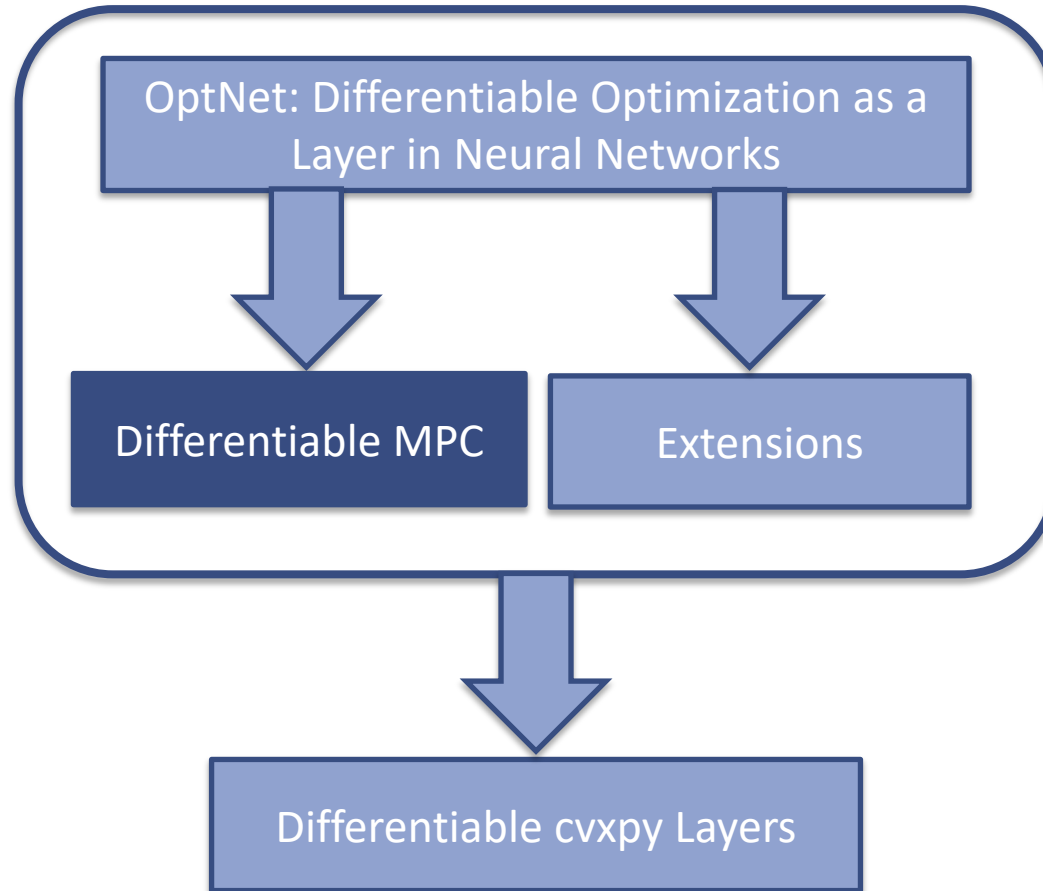
# OptNet Learns Sudoku

$$\begin{aligned} x^* = \operatorname{argmin}_x \operatorname{dist}(x, p) \\ \text{subject to } Ax = b \end{aligned}$$

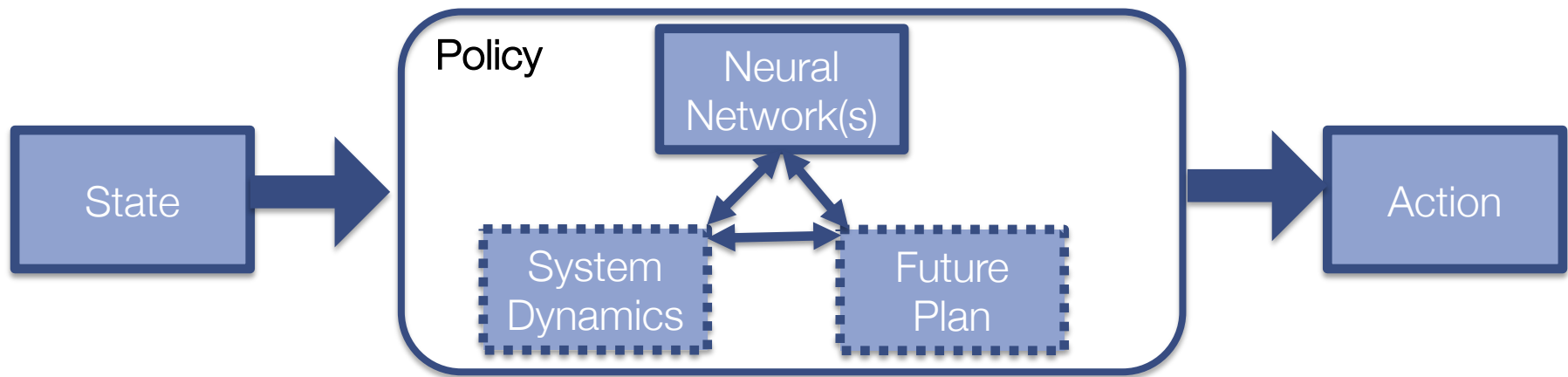
The OptNet layer exactly learns the mini-Sudoku constraints from data!  
**Baseline:** A deep convolutional feed-forward network



# Overview for the remainder of this talk



# Should RL policies have a system dynamics model or not?



## Model-free RL

More general, doesn't make as many assumptions about the world

Rife with poor data efficiency and learning stability issues

## Model-based RL (or control)

A useful prior on the world if it lies within your set of assumptions

# Combining model-based and model-free RL

Recently there has been a lot of interest in model-based priors for model-free reinforcement learning:

Among others: Dyna-Q (Sutton, 1990), GPS (Levine and Koltun, 2013), Imagination-Augmented Agents (Weber et al., 2017), Value Iteration Networks (Tamar et al., 2016), TreeQN (Farquhar et al., 2017)

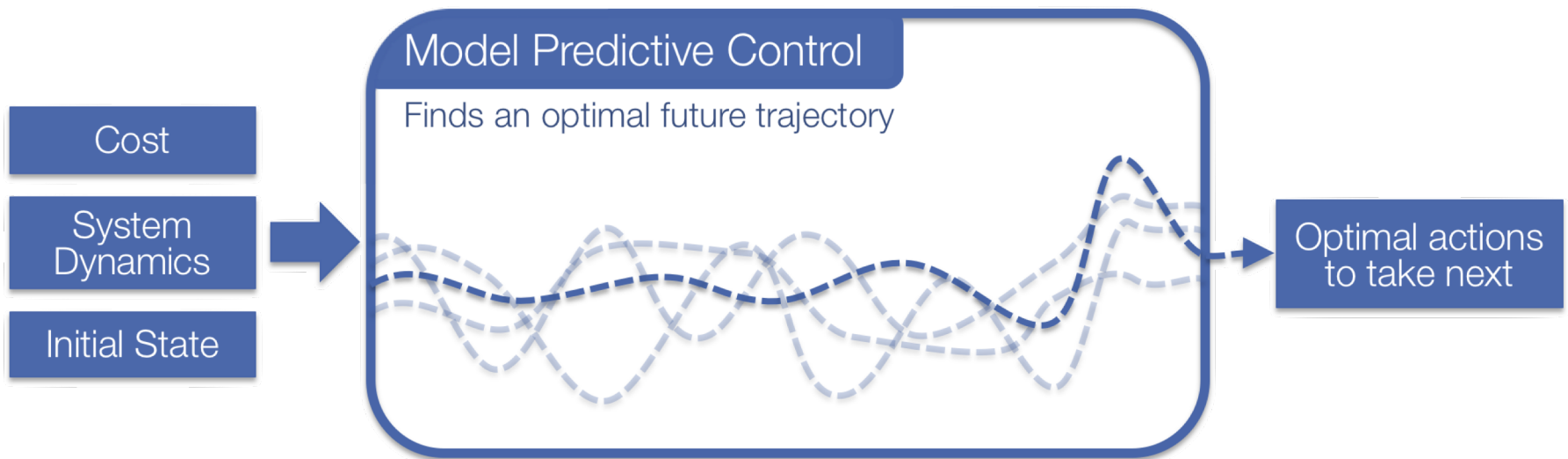
These typically involve:

1. **Using an RNN:** Efficient but not as expressive and general as MPC/iLQR
2. **Unrolling an LQR or gradient-based solver:** Expressive/general but inefficient

**Our approach:** Differentiable Model-Predictive Control

- **Explicitly** solves a control problem

# Our Approach: Model Predictive Control



# Our Approach: Model Predictive Control

Traditionally viewed as a pure **planning problem** given known (potentially non-convex) **cost** and **dynamics**:

$$\begin{aligned} \tau_{1:T}^* = \operatorname{argmin}_{\tau_{1:T}} \quad & \sum_t \mathcal{C}_\theta(\tau_t) \text{Cost} \\ \text{subject to } \quad & x_1 = x_{\text{init}} \\ & x_{t+1} = f_\theta(\tau_t) \text{Dynamics} \\ & \underline{u} \leq u \leq \bar{u} \end{aligned}$$

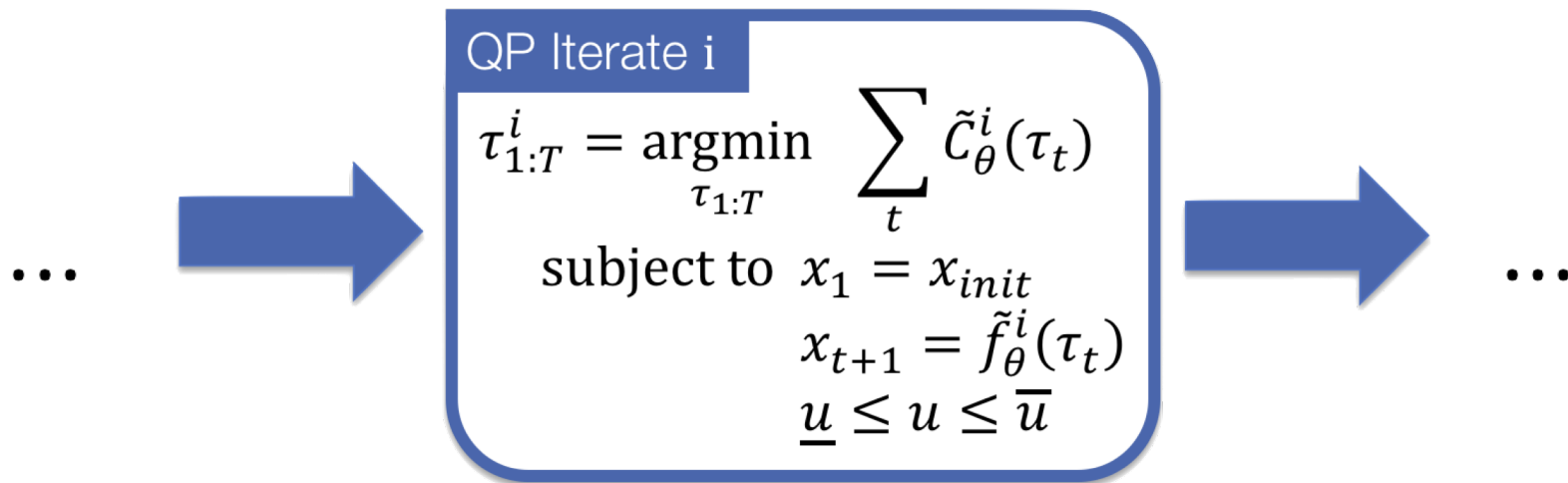
where  $\tau_t = \{x_t, u_t\}$

Execute  $u_1$  in the environment, observe the next observation, and repeat.

Cost and dynamics explicitly represented and learned.

# Model Predictive Control with SQP

- The standard way of solving MPC is to use **sequential quadratic programming (SQP)**, using LQR in most cases
- **Form approximations** to the cost and dynamics around the current iterate
- Repeat until a **fixed point** is reached and **differentiate through it**





# LQR, KKT Systems, and Differentiation

Solving LQR with dynamic Riccati recursion efficiently solves the KKT system

$$\overbrace{\begin{bmatrix} \dots & & & & \\ & C_t & F_t^\top & & \\ & F_t & [-I & 0] & \\ & & \begin{bmatrix} -I \\ 0 \end{bmatrix} & C_{t+1} & F_{t+1}^\top \\ & & & F_{t+1} & \\ & & & & \dots \end{bmatrix}}^K \begin{bmatrix} \vdots \\ \tau_t^* \\ \lambda_t^* \\ \tau_{t+1}^* \\ \lambda_{t+1}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ c_t \\ f_t \\ c_{t+1} \\ f_{t+1} \\ \vdots \end{bmatrix}$$

**Backwards Pass:** Use the OptNet approach from [Amos and Kolter, 2017] to implicitly **differentiate** the LQR KKT conditions:

$$\frac{\partial \ell}{\partial C_t} = \frac{1}{2} (d_{\tau_t}^* \otimes \tau_t^* + \tau_t^* \otimes d_{\tau_t}^*)$$

$$\frac{\partial \ell}{\partial F_t} = d_{\lambda_{t+1}}^* \otimes \tau_t^* + \lambda_{t+1}^* \otimes d_{\tau_t}^*$$

$$\frac{\partial \ell}{\partial c_t} = d_{\tau_t}^*$$

$$\frac{\partial \ell}{\partial f_t} = d_{\lambda_t}^*$$

$$\frac{\partial \ell}{\partial x_{\text{init}}} = d_{\lambda_0}^*$$

where  $K \begin{bmatrix} \vdots \\ d_{\tau_t}^* \\ d_{\lambda_t}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ \nabla_{\tau_t^*} \ell \\ 0 \\ \vdots \end{bmatrix}$



Just another LQR problem!

# LQR, KKT Systems, and Differentiation

Solving LQR with dynamic Riccati recursion efficiently solves the KKT system

$$\overbrace{\begin{bmatrix} \ddots & & & & \\ & C_t & F_t^\top & & \\ & F_t & [-I & 0] & \\ & & \begin{bmatrix} -I \\ 0 \end{bmatrix} & C_{t+1} & F_{t+1}^\top \\ & & & F_{t+1} & \\ & & & & \ddots \end{bmatrix}}^K \begin{bmatrix} \vdots \\ \tau_t^* \\ \lambda_t^* \\ \tau_{t+1}^* \\ \lambda_{t+1}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ c_t \\ f_t \\ c_{t+1} \\ f_{t+1} \\ \vdots \end{bmatrix}$$

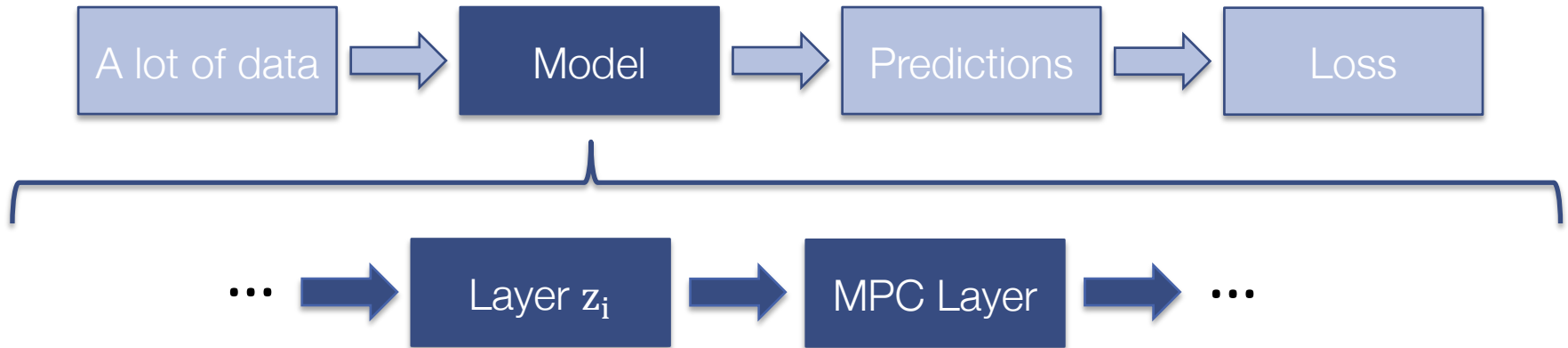
**Backwards Pass:** Use the OptNet approach from [Amos and Kolter, 2017] to implicitly differentiate the LQR KKT conditions:

$$\begin{aligned}
 \frac{\partial \ell}{\partial C_t} &= \frac{1}{2} (d_{\tau_t}^* \otimes \tau_t^* + \tau_t^* \otimes d_{\tau_t}^*) & \frac{\partial \ell}{\partial c_t} &= d_{\tau_t}^* & \frac{\partial \ell}{\partial x_{\text{init}}} &= d_{\lambda_0}^* & \text{where } K \begin{bmatrix} \vdots \\ d_{\tau_t}^* \\ d_{\lambda_t}^* \\ \vdots \end{bmatrix} &= - \begin{bmatrix} \vdots \\ \nabla_{\tau_t^*} \ell \\ 0 \\ \vdots \end{bmatrix} \\
 \frac{\partial \ell}{\partial F_t} &= d_{\lambda_{t+1}}^* \otimes \tau_t^* + \lambda_{t+1}^* \otimes d_{\tau_t}^* & \frac{\partial \ell}{\partial f_t} &= d_{\lambda_t}^* & & & & 
 \end{aligned}$$

**Just another LQR problem!**

# A Differentiable MPC Module

We can differentiate through (non-convex) MPC with a single (convex) LQR solve by differentiating the SQP fixed point



## What can we do with this now?

Replace **neural network policies** in model-free algorithms with MPC policies, and also **replace the unrolled controllers** in other settings (hindsight plan, universal planning networks)

**The cost** can also be learned! No longer have to hard-code in a known value.

# Imitation learning with a linear model

Linear dynamics:  $f(x_t, u_t) = Ax_t + Bu_t$

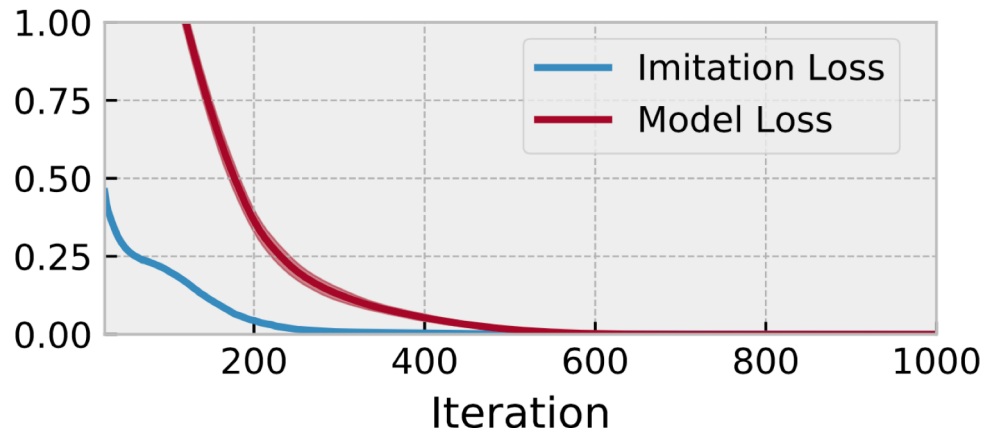
Parameters:  $\theta = \{A, B\}$

Trajectory:  $\tau_\theta(x_{\text{init}})$  obtained by MPC

Given known  $\theta$  and sample trajectories, learn  $\hat{\theta}$

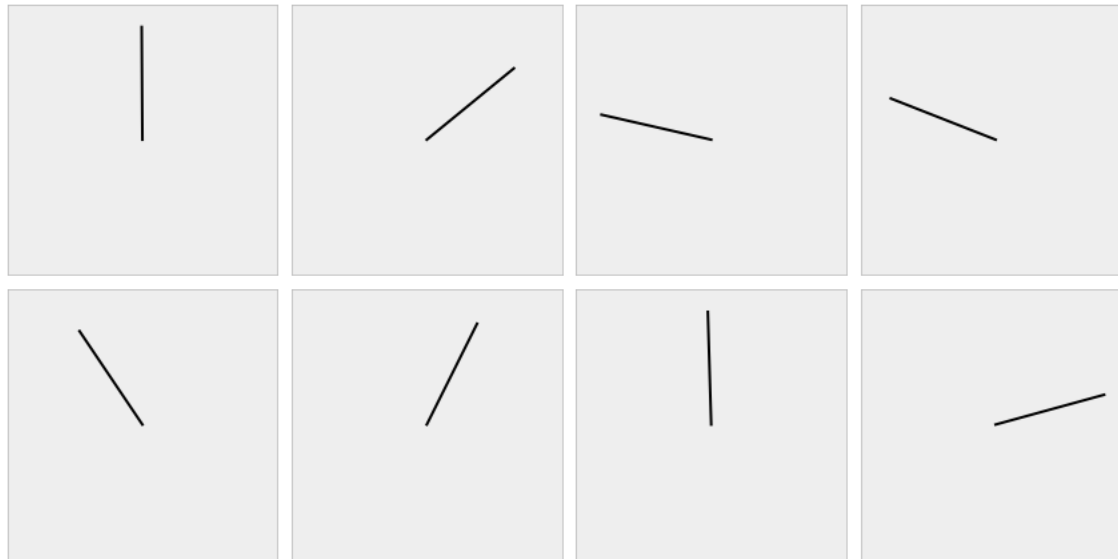
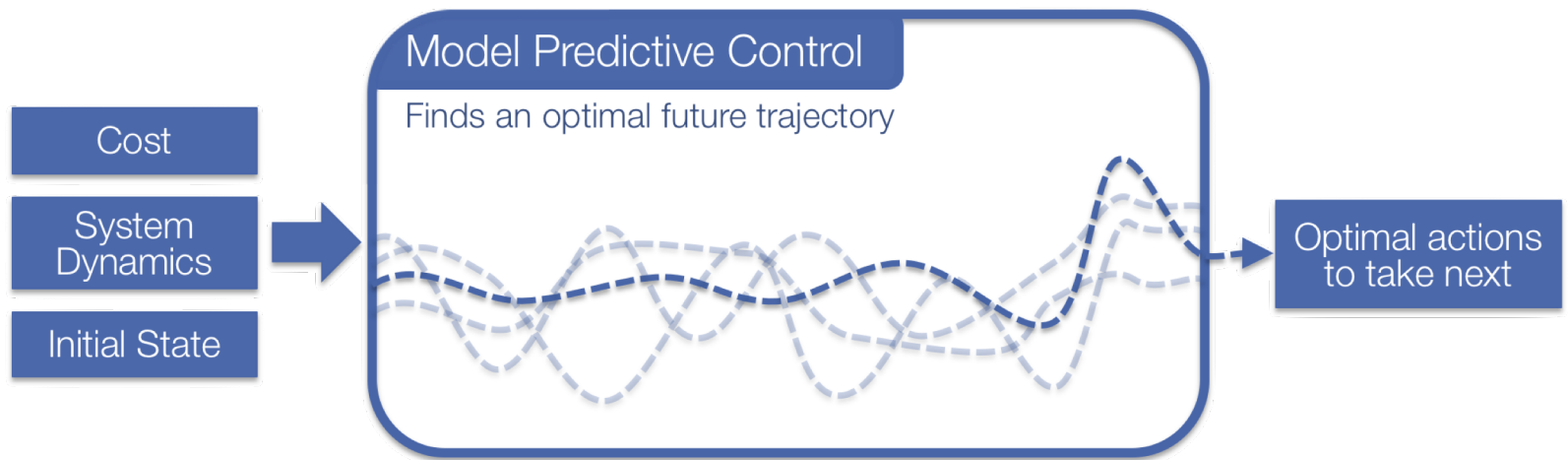
Trajectory (Training) Loss:  $\text{MSE}(\tau_\theta(x_{\text{init}}), \tau_{\hat{\theta}}(x_{\text{init}}))$

Model Loss:  $\text{MSE}(\theta, \hat{\theta})$



Not guaranteed to converge, but a good sanity check that it does in small cases.

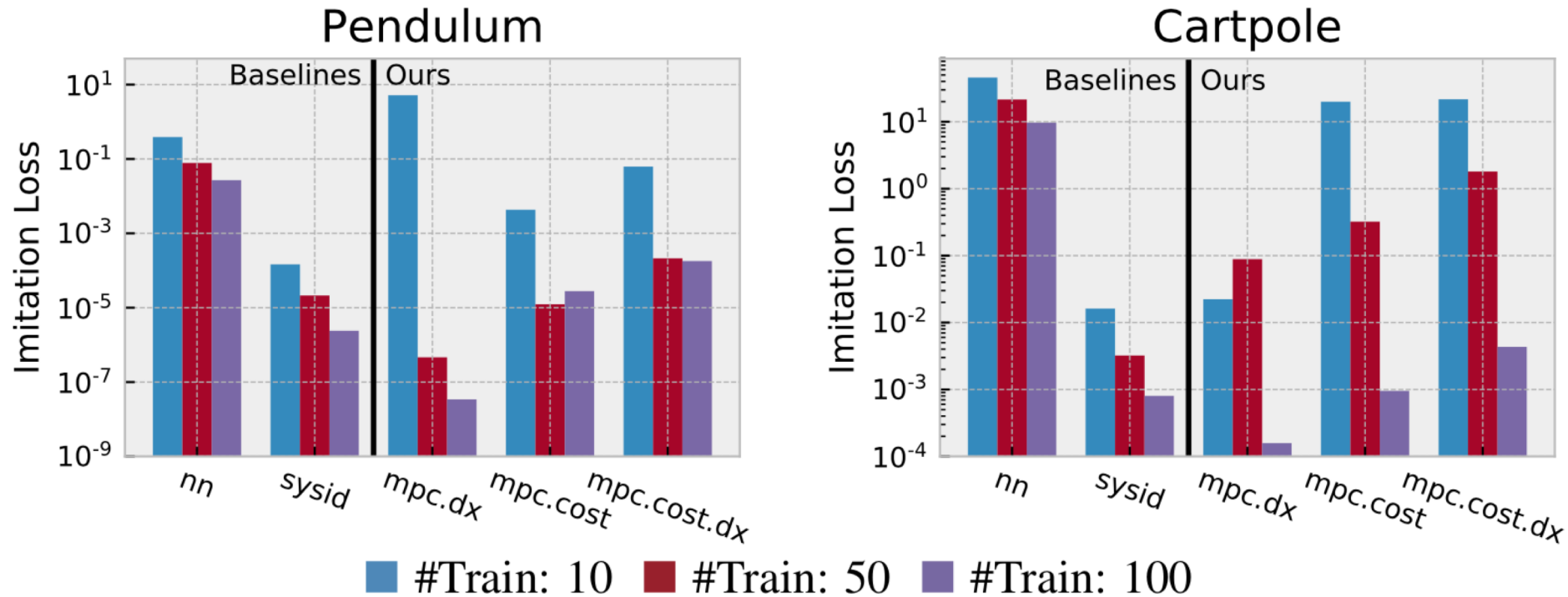
# Simple Pendulum Control



# Imitation learning with the pendulum/cartpole

Again optimizes the imitation loss with respect to the controller's parameters

Using only **action trajectories** we can recover the true parameters

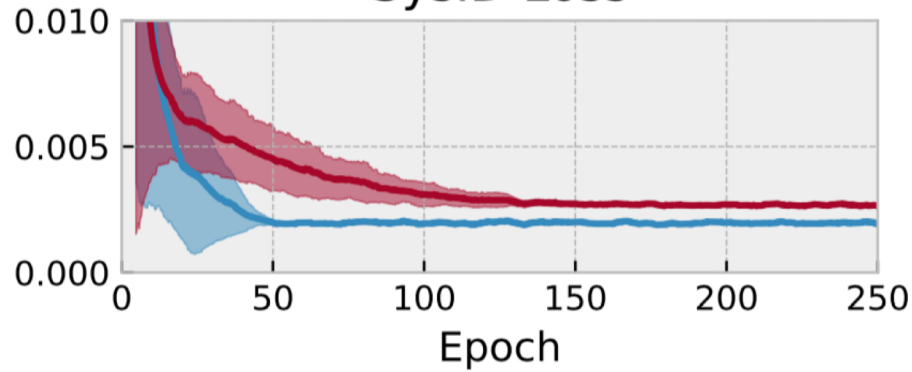


# Optimizing the task loss is often better than SysID in the unrealizable case

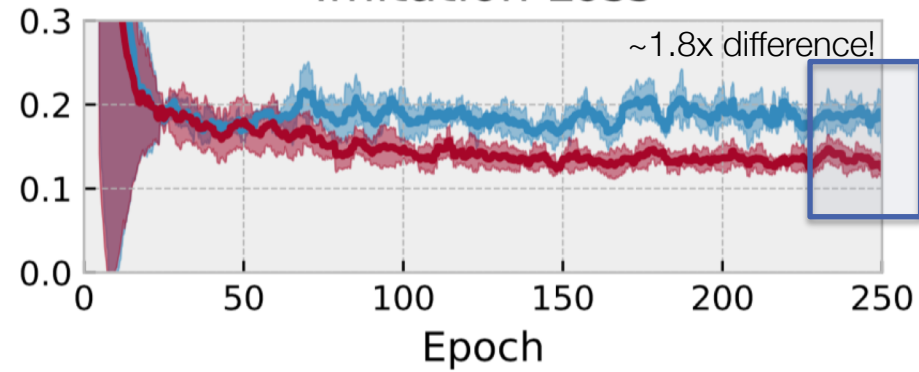
True System: Pendulum environment with noise (damping and a wind force)

Approximate Model: Pendulum without the noise terms

### SysID Loss



### Imitation Loss



■ Vanilla SysId Baseline ■ (Ours) Directly optimizing the Imitation Loss





# A PyTorch MPC Solver

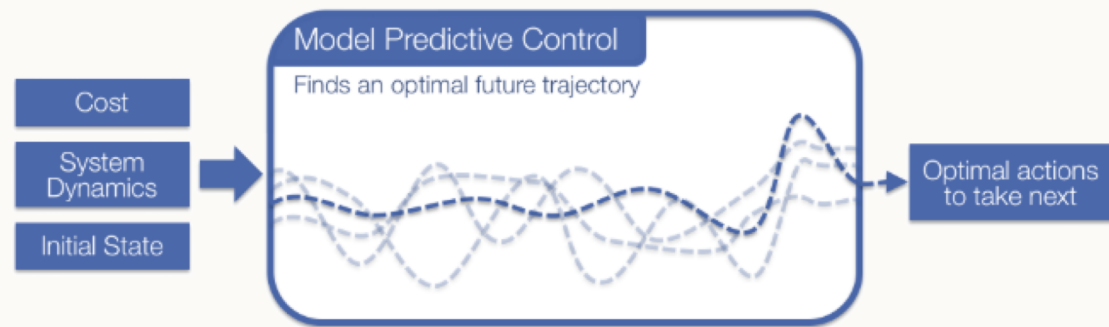
<https://locuslab.github.io/mpc.pytorch>

## mpc.pytorch

A fast and differentiable model predictive control (MPC) solver for PyTorch. Crafted by [Brandon Amos](#), [Ivan Jimenez](#), [Jacob Sacks](#), [Byron Boots](#), and [J. Zico Kolter](#). For more context and details, see our [ICML 2017 paper](#) on [OptNet](#) and our (forthcoming) [NIPS 2018 paper](#) on differentiable MPC.

 [View On GitHub](#)

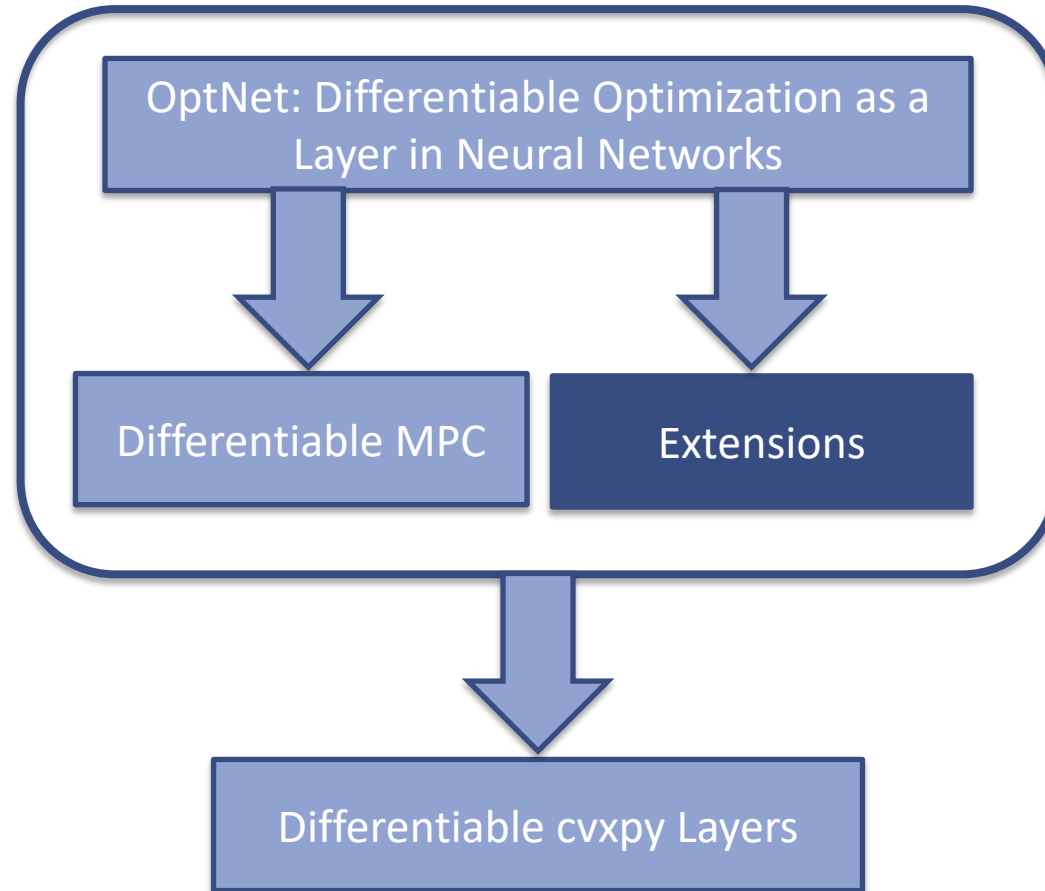
## Control is important!



Optimal control is a widespread field that involve finding an optimal sequence of future actions to take in a system or environment. This is the most useful in domains when you can analytically model your system and can easily define a cost to optimize over your system. This project focuses on solving [model predictive control](#) (MPC) with the [box-DDP](#) heuristic. MPC is a powerhouse in many real-world domains ranging from short-time horizon robot control tasks to long-time horizon control of chemical processing plants. More recently, the reinforcement learning community, [strife](#) with poor sample-complexity and instability issues in model-free learning, has been actively searching for useful model-based applications and priors.



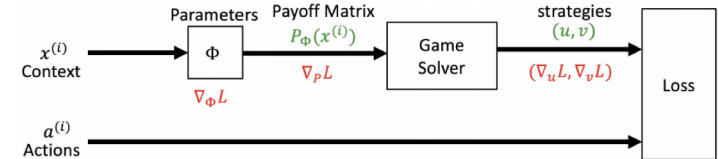
# Overview for the remainder of this talk



# Extensions

Section 2 and Section 8 of my thesis document contain a more complete set of references

**Game Theory [Ling, Fang, and Kolter; IJCAI 2017]: Distinguished Paper Award**



$$\begin{bmatrix} 0 \\ 0 \\ a \\ \sigma \\ \zeta \end{bmatrix} - \begin{bmatrix} M & -J_c & -J_c & -J_f & 0 \\ J_c & 0 & 0 & 0 & 0 \\ J_c & 0 & 0 & 0 & 0 \\ J_f & 0 & 0 & 0 & E \\ 0 & 0 & 0 & -E^T & 0 \end{bmatrix} \begin{bmatrix} v_{t+dt} \\ \lambda_c \\ \lambda_f \\ \gamma \end{bmatrix} = \begin{bmatrix} Mv_t + dtf_t \\ 0 \\ c \\ 0 \\ 0 \end{bmatrix}$$

$$\text{subject to } \begin{bmatrix} a \\ \sigma \\ \zeta \end{bmatrix} \geq 0, \begin{bmatrix} \lambda_c \\ \lambda_f \\ \gamma \end{bmatrix} \geq 0, \begin{bmatrix} a \\ \sigma \\ \zeta \end{bmatrix}^T \begin{bmatrix} \lambda_c \\ \lambda_f \\ \gamma \end{bmatrix} = 0,$$

**Stochastic optimization and end-to-end learning [Donti, Amos, and Kolter; NeurIPS 2017]**

**Reinforcement learning and control**

Safety [Dalal et al. 2018], physics-based modeling [Peres et al. NeurIPS 2018], inverse cost and reward learning, multi-agent systems, learnable embeddings

**Discrete, combinatorial, and submodular optimization**

[Djolonga and Krause 2017, Niculae and Blondel 2017, Mensch and Blondel 2018]  $y^* = \min_{x \in \mathcal{B}(G)} \frac{1}{2} \|y - y'\|$ , where  $y' = \arg \min_y f(y) + \frac{1}{2} \|y - z\|^2$ .

$$\Pi_{\Omega}(x) := \arg \max_{y \in \Delta^d} y^T x - \gamma \Omega(y) = \nabla \max_{\Omega}(x)$$

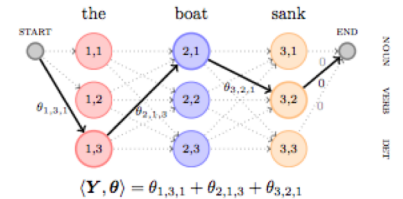
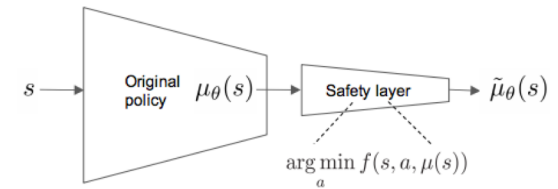
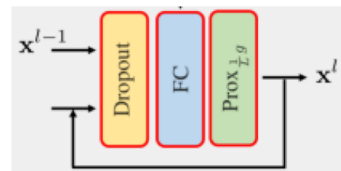


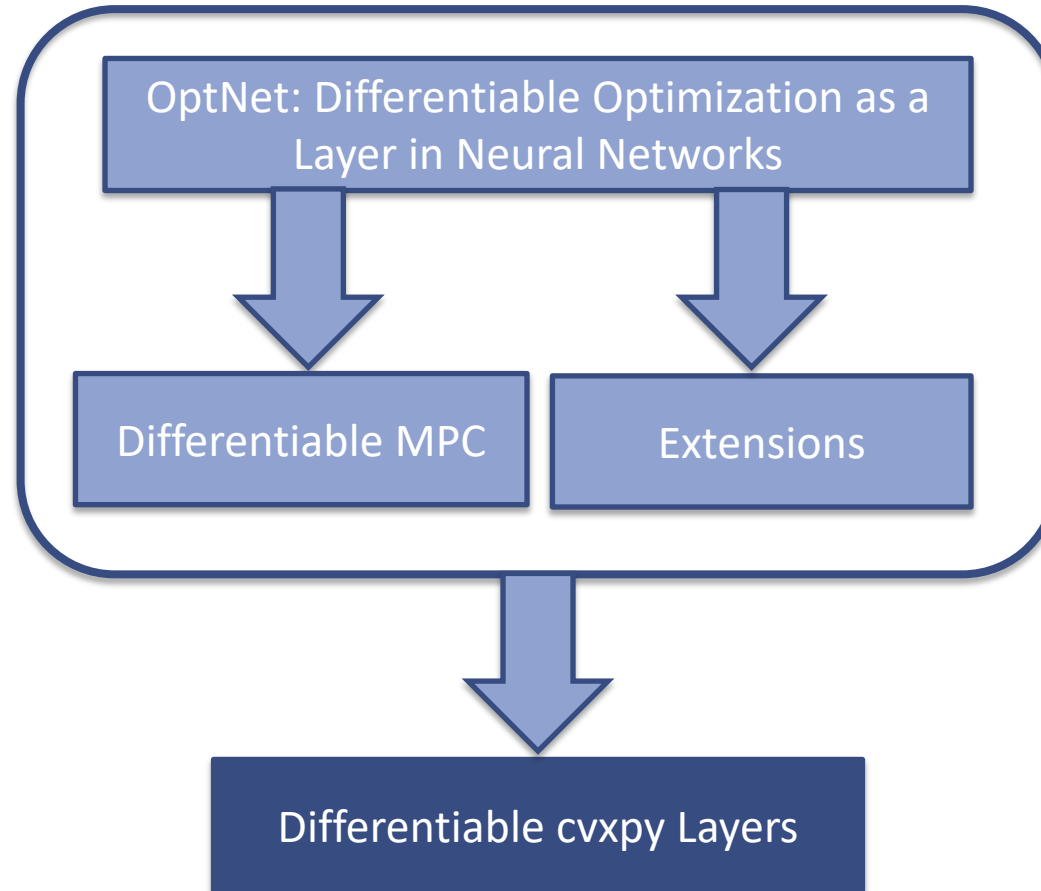
Figure 2. Computational graph of the Viterbi algorithm.

**Optimization viewpoints of standard components**

[Bibi et al. ICLR 2019]



# Overview for the remainder of this talk



# Background: cvxpy

<http://cvxpy.org>

(constrained LASSO)

[Diamond2018]

$$\begin{aligned} & \text{minimize} && \|Ax - b\|_2^2 + \gamma \|x\|_1 \\ & \text{subject to} && \mathbf{1}^T x = 0, \quad \|x\|_\infty \leq 1 \end{aligned}$$

with variable  $x \in \mathbf{R}^n$

---

```
from cvxpy import *
x = Variable(n)
cost = sum_squares(A*x-b) + gamma*norm(x,1)
obj = Minimize(cost)
constr = [sum_entries(x) == 0, norm(x,"inf") <= 1]
prob = Problem(obj, constr)
opt_val = prob.solve()
solution = x.value
```

# Hand-Implementing Optimization Layers is Non-Trivial

$$\begin{aligned} dQz^* + Qdz + dq + dA^T\nu^* + \\ A^T d\nu + dG^T\lambda^* + G^T d\lambda = 0 \\ dAz^* + Adz - db = 0 \end{aligned}$$

$$D(Gz^* - h)d\lambda + D(\lambda^*)(dGz^* + Gdz - dh) = 0$$

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = \begin{bmatrix} -dQz^* - dq - dG^T\lambda^* - dA^T\nu^* \\ -D(\lambda^*)dGz^* + D(\lambda^*)dh \\ -dAz^* + db \end{bmatrix}$$

$$\begin{aligned} \nabla_Q \ell &= \frac{1}{2}(d_z z^T + z d_z^T) & \nabla_q \ell &= d_z \\ \nabla_A \ell &= d_\nu z^T + \nu d_\nu^T & \nabla_b \ell &= -d_\nu \\ \nabla_G \ell &= D(\lambda^*)(d_\lambda z^T + \lambda d_\lambda^T) & \nabla_h \ell &= -D(\lambda^*)d_\lambda \end{aligned}$$

$$\begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} Q & G^T D(\lambda^*) & A^T \\ G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} \nabla_{z^*} \ell \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} Q & A^T & \tilde{G}^T \\ A & 0 & 0 \\ \tilde{G} & 0 & 0 \end{bmatrix} \begin{bmatrix} d_x^* \\ d_\lambda^* \\ d_\nu^* \end{bmatrix} = - \begin{bmatrix} \nabla_{x^*} \ell \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{aligned} \nabla_Q \ell &= \frac{1}{2}(d_x^* \otimes x^* + x^* \otimes d_x^*) & \nabla_p \ell &= d_x^* \\ \nabla_A \ell &= d_\lambda^* \otimes x^* + \lambda^* \otimes d_x^* & \nabla_b \ell &= -d_\lambda^* \end{aligned}$$

$$\begin{array}{c} \overbrace{\hspace{10em}}^K \\ \begin{bmatrix} \ddots & & & & & \\ & C_t & F_t^\top & & & \\ & F_t & & [-I & 0] & \\ & & & [-I & ] & \\ & & & & C_{t+1} & F_{t+1}^\top \\ & & & & F_{t+1} & \\ & & & & & \ddots \end{bmatrix} \begin{bmatrix} \vdots \\ \tau_t^* \\ \lambda_t^* \\ \tau_{t+1}^* \\ \lambda_{t+1}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ c_t \\ f_t \\ c_{t+1} \\ f_{t+1} \\ \vdots \end{bmatrix} \end{array}$$

$$K \begin{bmatrix} \vdots \\ d_{\tau_t}^* \\ d_{\lambda_t}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ \nabla_{\tau_t^*} \ell \\ 0 \\ \vdots \end{bmatrix}$$

$$\begin{aligned} \frac{\partial \ell}{\partial C_t} &= \frac{1}{2} (d_{\tau_t}^* \otimes \tau_t^* + \tau_t^* \otimes d_{\tau_t}^*) \\ \frac{\partial \ell}{\partial F_t} &= d_{\lambda_{t+1}^*} \otimes \tau_t^* + \lambda_{t+1}^* \otimes d_{\tau_t}^* \end{aligned}$$

$$\begin{aligned} \frac{\partial \ell}{\partial c_t} &= d_{\tau_t}^* \\ \frac{\partial \ell}{\partial f_t} &= d_{\lambda_t}^* \end{aligned}$$

$$\frac{\partial \ell}{\partial x_{\text{init}}} = d_{\lambda_0}^*$$

# Why should practitioners care?

$$\begin{aligned}
 dQz^* + Qdz + dq + dA^T\nu^* + \\
 A^T d\nu + dG^T\lambda^* + G^T d\lambda = 0 \\
 dAz^* + Adz - db = 0 \\
 D(Gz^* - h) + D(\lambda^*)(dGz^* + Gdz - dh) = 0
 \end{aligned}$$

$$\begin{aligned}
 \nabla_{Q\ell} &= \frac{1}{2}(d_z z^T + z d_z^T) & \nabla_{q\ell} &= d_z \\
 \nabla_{A\ell} &= d_\nu z^T + \nu d_z^T & \nabla_{b\ell} &= -d_\lambda \\
 \nabla_{G\ell} &= D(\lambda^*)(d_\lambda z^T + \lambda d_z^T) & &= -D(\lambda^*)d_\lambda
 \end{aligned}$$

$$\begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{bmatrix} dz \\ d\lambda \\ d\nu \end{bmatrix} = - \begin{bmatrix} -dQz^* - dq - dG^T\lambda^* - dA^T\nu^* \\ -D(\lambda^*)dGz^* + D(\lambda^*)dh \\ -dAz^* + db \end{bmatrix}$$

$$\begin{bmatrix} d_z \\ d_\lambda \\ d_\nu \end{bmatrix} = - \begin{bmatrix} Q & G^T & A^T \\ D(\lambda^*)G & D(Gz^* - h) & 0 \\ A & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} \nabla_{z^*}\ell \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} \dots & \tau_t & \lambda_t & \tau_{t+1} & \lambda_{t+1} \\ \dots & C_t & F_t^T & & \\ C_t & F_t & & [-I & 0] \\ & [-I] & & C_{t+1} & F_{t+1}^T \\ & 0 & & F_{t+1} & \dots \end{bmatrix} \begin{bmatrix} \lambda_t^* \\ \tau_{t+1}^* \\ \lambda_{t+1}^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} f_t \\ c_{t+1} \\ f_{t+1} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} Q & A^T & \tilde{G}^T \\ A & 0 & 0 \\ \tilde{G} & 0 & 0 \end{bmatrix} \begin{bmatrix} d_x^* \\ d_\lambda^* \\ d_\nu^* \end{bmatrix} = - \begin{bmatrix} \nabla_{x^*}\ell \\ 0 \\ 0 \end{bmatrix}$$

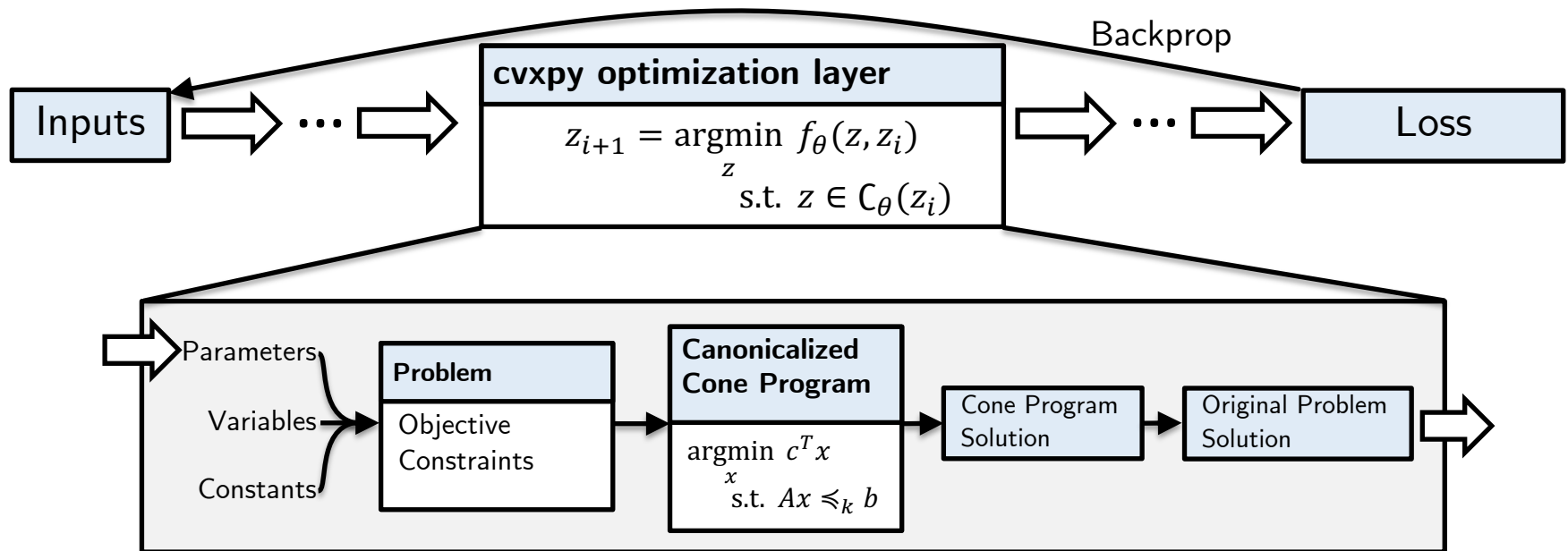
$$\begin{aligned}
 \nabla_{Q\ell} &= \frac{1}{2}(d_x^* \otimes x^* + x^* \otimes d_x^*) & \nabla_{p\ell} &= d_x^* \\
 \nabla_{A\ell} &= d_\nu^* \otimes x^* + \lambda^* \otimes d_x^* & \nabla_{b\ell} &= -d_\lambda^*
 \end{aligned}$$

$$\begin{bmatrix} \vdots \\ d_x^* \\ \vdots \\ \lambda_t^* \\ \vdots \end{bmatrix} = - \begin{bmatrix} \vdots \\ \nabla_{\tau_t^*}\ell \\ 0 \\ \vdots \end{bmatrix}$$

$$\begin{aligned}
 \frac{\partial \ell}{\partial C_t} &= \frac{1}{2} (d_{\tau_t}^* \otimes \tau_t^* + \tau_t^* \otimes d_{\tau_t}^*) \\
 \frac{\partial \ell}{\partial F_t} &= d_{\lambda_{t+1}}^* \otimes \tau_t^* + \lambda_{t+1}^* \otimes d_{\tau_t}^*
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \ell}{\partial c_t} &= d_{\tau_t}^* & \frac{\partial \ell}{\partial \lambda_0} &= d_{\lambda_0}^* \\
 \frac{\partial \ell}{\partial f_t} &= d_{\lambda_t}^*
 \end{aligned}$$

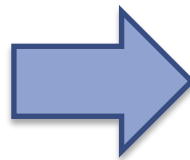
# A new way of rapidly prototyping optimization layers



# Full source code example: OptNet QP

**Before:** 1k lines of code

Hand-implemented and optimized PyTorch GPU-capable batched primal-dual interior point method



**Now:** 10 lines of code

Same speed

$$z_{i+1} = \underset{z}{\operatorname{argmin}} \frac{1}{2} z^T Q(z_i) z + q(z_i)^T z$$

subject to  $A(z_i)z = b(z_i)$   
 $G(z_i)z \leq h(z_i)$

Parameters/Submodules :  $Q, q, A, b, G, h$

```
1 Q = cp.Parameter((n, n), PSD=True)
2 p = cp.Parameter(n)
3 A = cp.Parameter((m, n))
4 b = cp.Parameter(m)
5 G = cp.Parameter((p, n))
6 h = cp.Parameter(p)
7 x = cp.Variable(n)
8 obj = cp.Minimize(0.5*cp.quad_form(x, Q) + p.T * x)
9 cons = [A*x == b, G*x <= h]
10 prob = cp.Problem(obj, cons)
11 layer = CvxpyLayer(prob, params=[Q, p, A, b, G, h], out=[x])
```

`import cvxpy as cp`  
`from cvxpyth import CvxpyLayer`



# Full source code example: The sigmoid

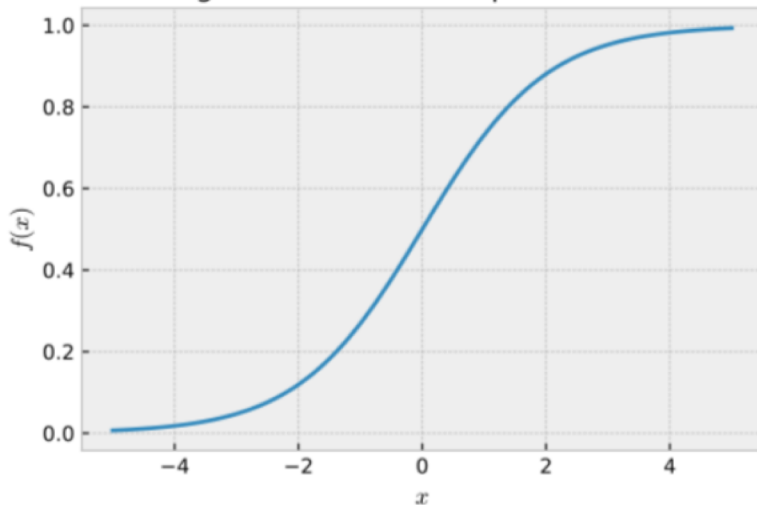
$$y = \frac{1}{1 + e^{-x}}$$

$$y^* = \underset{y}{\operatorname{argmin}} -y^T x - H_b(y)$$

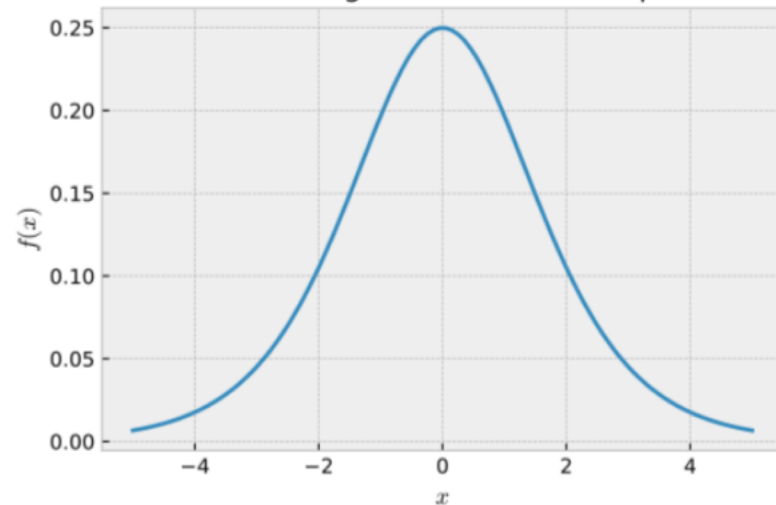
subject to  $0 \leq y \leq 1$

```
1 x = cp.Parameter(n)
2 y = cp.Variable(n)
3 obj = cp.Minimize(-x.T*y - cp.sum(cp.entr(y) + cp.entr(1.-y)))
4 prob = cp.Problem(obj)
5 layer = CvxpyLayer(prob, params=[x], out_vars=[y])
```

The Sigmoid Function in Optimization Form



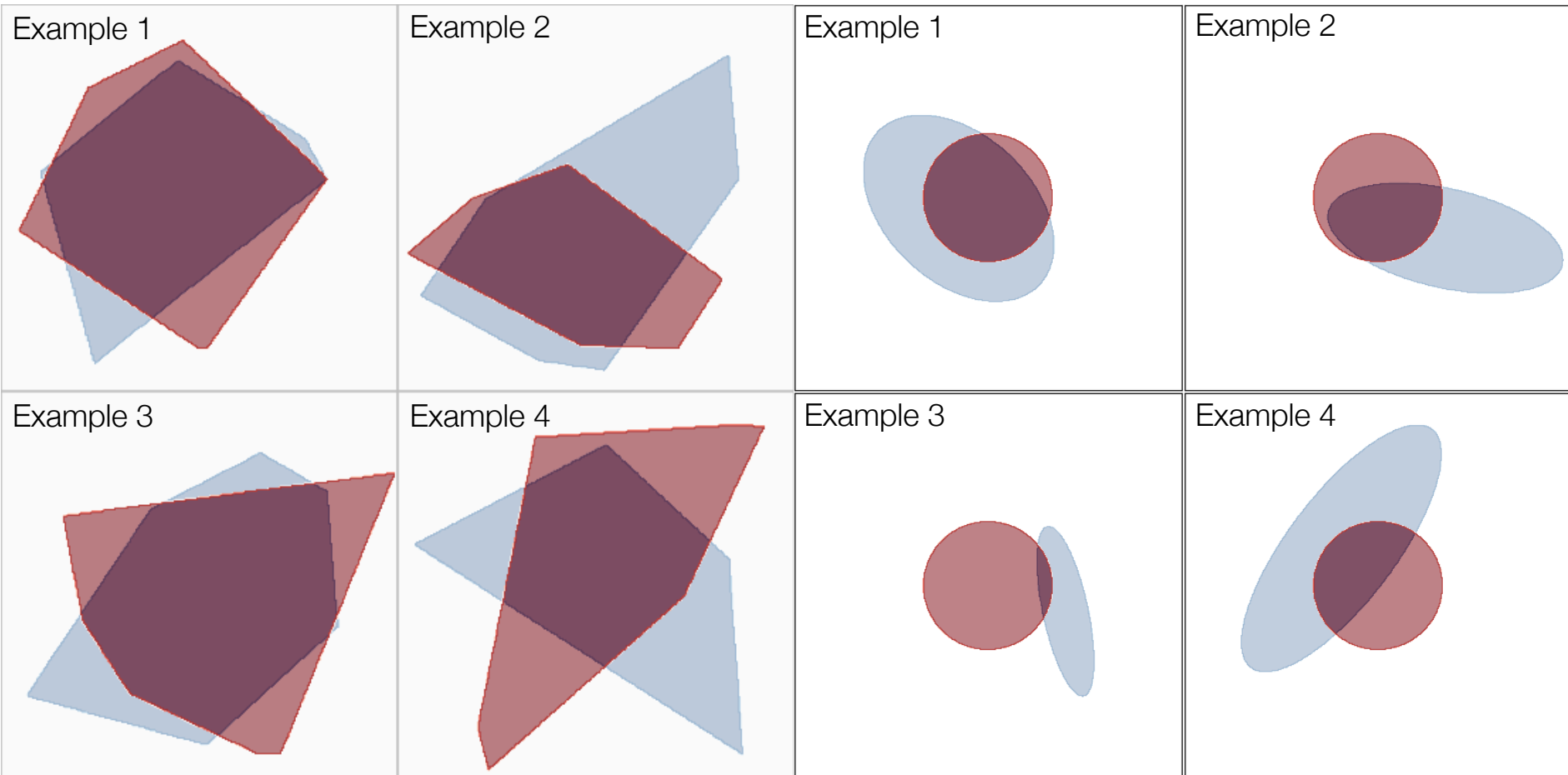
The Derivative of the Sigmoid Function in Optimization Form



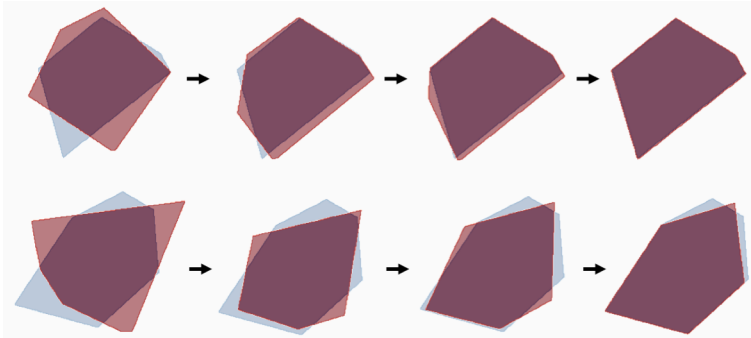
# OptNet Application: Modeling Constraints

■ True Constraint (Unknown to the model)

■ Constraint Predictions During Training

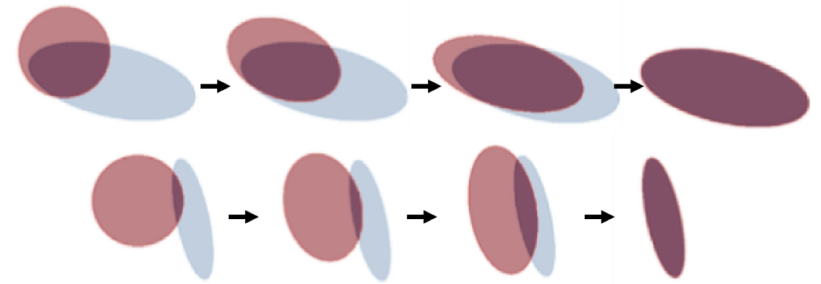


# Full source code example: Constraint modeling



$$\hat{y} = \operatorname{argmin}_y \frac{1}{2} \|p - y\|_2^2$$

$$\text{s.t. } Gy \leq h$$



$$\hat{y} = \operatorname{argmin}_y \frac{1}{2} \|p - y\|_2^2$$

$$\text{s.t. } \frac{1}{2} (y - z)^\top A (y - z) \leq 1$$

```

1 G = cp.Parameter((m, n))
2 h = cp.Parameter(m)
3 p = cp.Parameter(n)
4 y = cp.Variable(n)
5 obj = cp.Minimize(0.5*cp.sum_squares(y-p))
6 cons = [G*y <= h]
7 prob = cp.Problem(obj, cons)
8 layer = CvxpyLayer(prob, params=[p, G, h], out=[y])

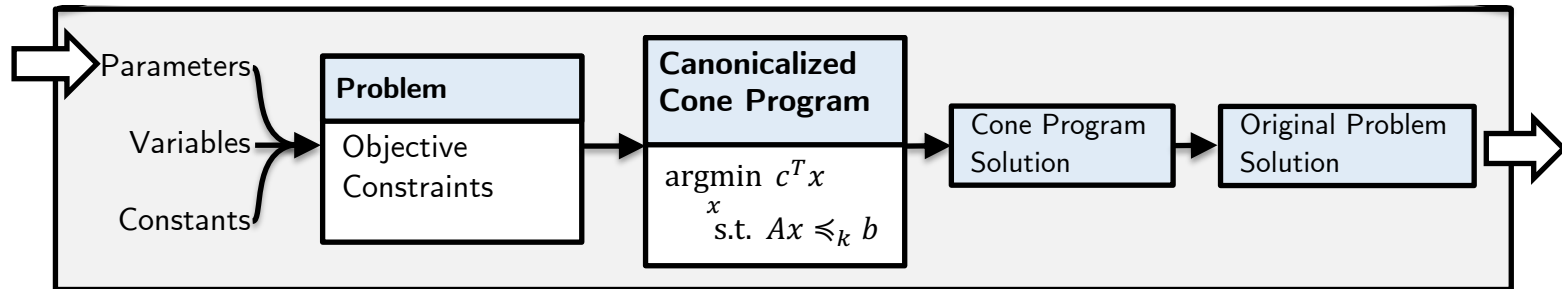
```

```

1 A = cp.Parameter((n, n), PSD=True)
2 z = cp.Parameter(n)
3 p = cp.Parameter(n)
4 y = cp.Variable(n)
5 obj = cp.Minimize(0.5*cp.sum_squares(y-p))
6 cons = [0.5*cp.quad_form(y-z, A) <= 1]
7 prob = cp.Problem(obj, cons)
8 layer = CvxpyLayer(prob, params=[p, A, z], out=[y])

```

# What's going on behind the scenes?



## Cone Program Differentiation

Much more general than the QPs we considered in OptNet

Question from my thesis proposal: How to differentiate non-polyhedral cones?

**Non-trivial** because we can't easily differentiate the KKT conditions of cone programs because of non-trivial cone constraints

# Cone Program Differentiation

Take the homogenous self-dual embedding of the cone program

$$Qu = v \quad \text{where} \quad Q = \begin{bmatrix} 0 & A^\top & c \\ -A & 0 & b \\ -c^\top & -b^\top & 0 \end{bmatrix} \quad \begin{array}{l} u \in \mathcal{K}, \quad v \in \mathcal{K}^*, \quad u_{m+n+1} + v_{m+n+1} > 0, \\ \mathcal{K} = \mathbb{R}^n \times \mathcal{K}^* \times \mathbb{R}_+, \quad \mathcal{K}^* = \{0\}^n \times \mathcal{K} \times \mathbb{R}_+, \end{array}$$

**Definition:** Minty's projection onto the embedding space

$$M: \mathbb{R}^{m+n+1} \rightarrow \mathcal{C} \quad M(z) = (\Pi z, -\Pi^* z) \quad \text{where } \mathcal{C} = \{(u, v) \in \mathcal{K} \times \mathcal{K}^* \mid u^\top v = 0\}$$

Take the **residual map** of Minty's parameterization:

$$\mathcal{R}(z) = Q\Pi z + \Pi^* z$$

Implicitly differentiate  $\mathcal{R}$ :

$$D_\theta(z) = -(D_z \mathcal{R}(z^*))^{-1} D_\theta R(z^*)$$

Captures KKT differentiation as a special case

# Closing Thoughts And Future Directions

Optimization is a powerful primitive to use within larger systems

- This thesis has uncovered **theoretical** and **engineering** foundations
- Can be **propagated through and learned**, just like any layer
- Provides a **perspective to analyze** existing models and layers
- Can be used to **project onto sets in a differentiable way**

Even if a closed form solution doesn't exist

Applications in:

- **Model-based RL and control**
  - In the **policy** or for **exploration**
  - Inverse control, **cost learning**
  - **Learning embedded state spaces** for planning
  - **Multi-agent systems**  
Interpret other agents as solving optimization problems
- **Meta-Learning**
- **Energy-based learning and structured prediction**

# Closing Thoughts And Future Directions

Optimization is a powerful primitive to use within larger systems

- This thesis has uncovered **theoretical** and **engineering** foundations
- Can be **propagated through and learned**, just like any layer
- Provides a **perspective to analyze** existing models and layers
- Can be used to **project onto sets in a differentiable way**  
Even if a closed form solution doesn't exist



Applications in:

- **Model-based RL and control**
  - In the **policy** or for **exploration**
  - Inverse control, **cost learning**
  - **Learning embedded state spaces** for planning
  - **Multi-agent systems**  
Interpret other agents as solving optimization problems
- **Meta-Learning**
- **Energy-based learning and structured prediction**

Thesis Defense

# Differentiable Optimization-Based Modeling for Machine Learning

Brandon Amos • Carnegie Mellon University

 brandondamos  
 bamos.github.io



The source code behind all of my work is free and publicly available:

**<http://github.com/bamos/thesis>**

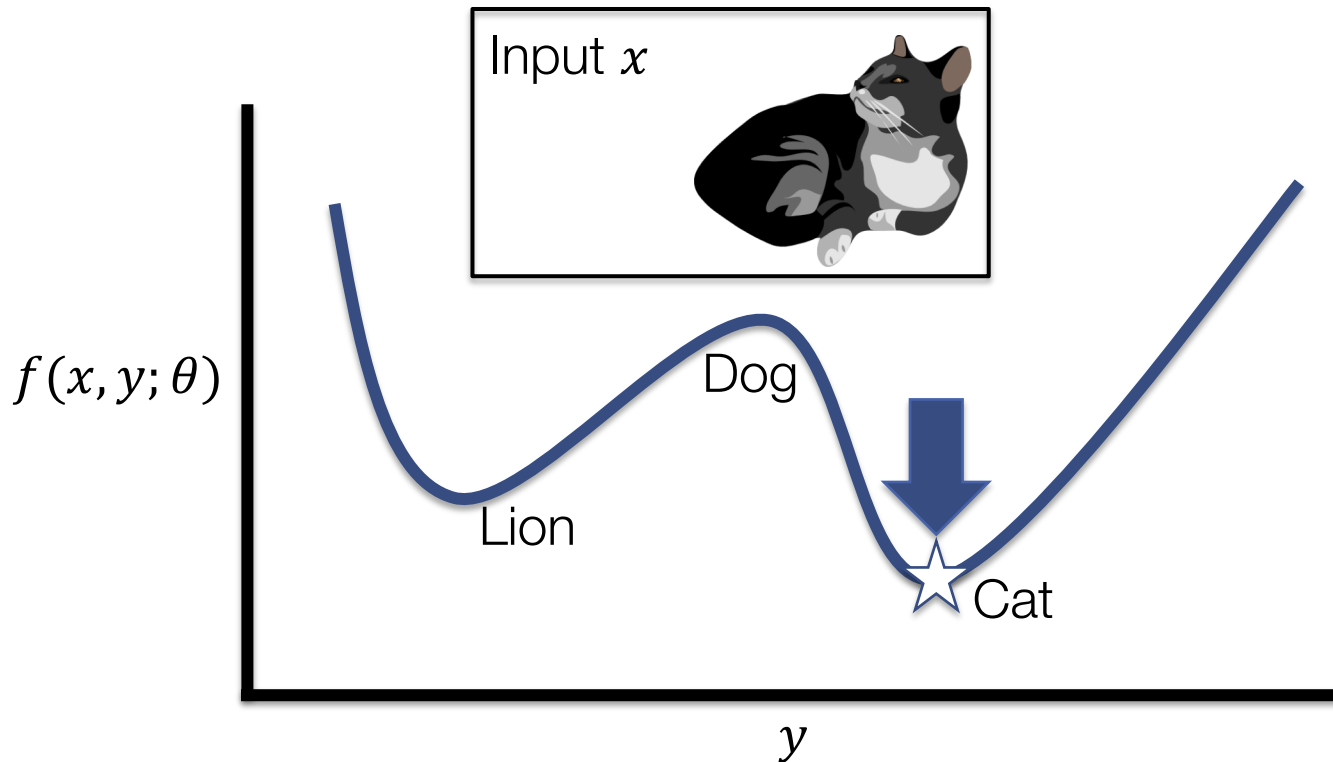


# Extra Slides

# Optimization-Based Inference

Structured prediction: define a network over  $\mathcal{X} \times \mathcal{Y}$  and predict via  
$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$

\*This is also called energy-based modeling



# Structured prediction models nicely capture dependencies in the output space

Especially useful for high-dimensional, correlated output spaces

- Multi-label classification
- Semantic segmentation
- Scene-graph generation

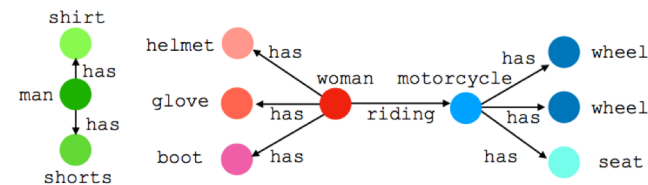
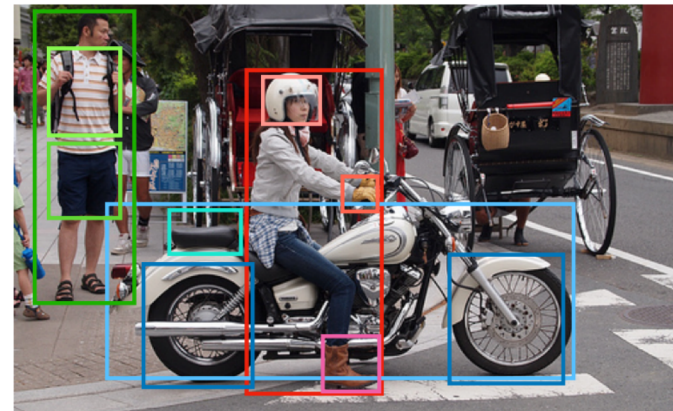
Difficult to capture with most feed-forward models

Intractable in many graphical models if a special structure is not imposed

- Like in MRFs/CRFs

Easy with energy-based models

- Just add them to the energy  $f_{\theta}(x, y)$

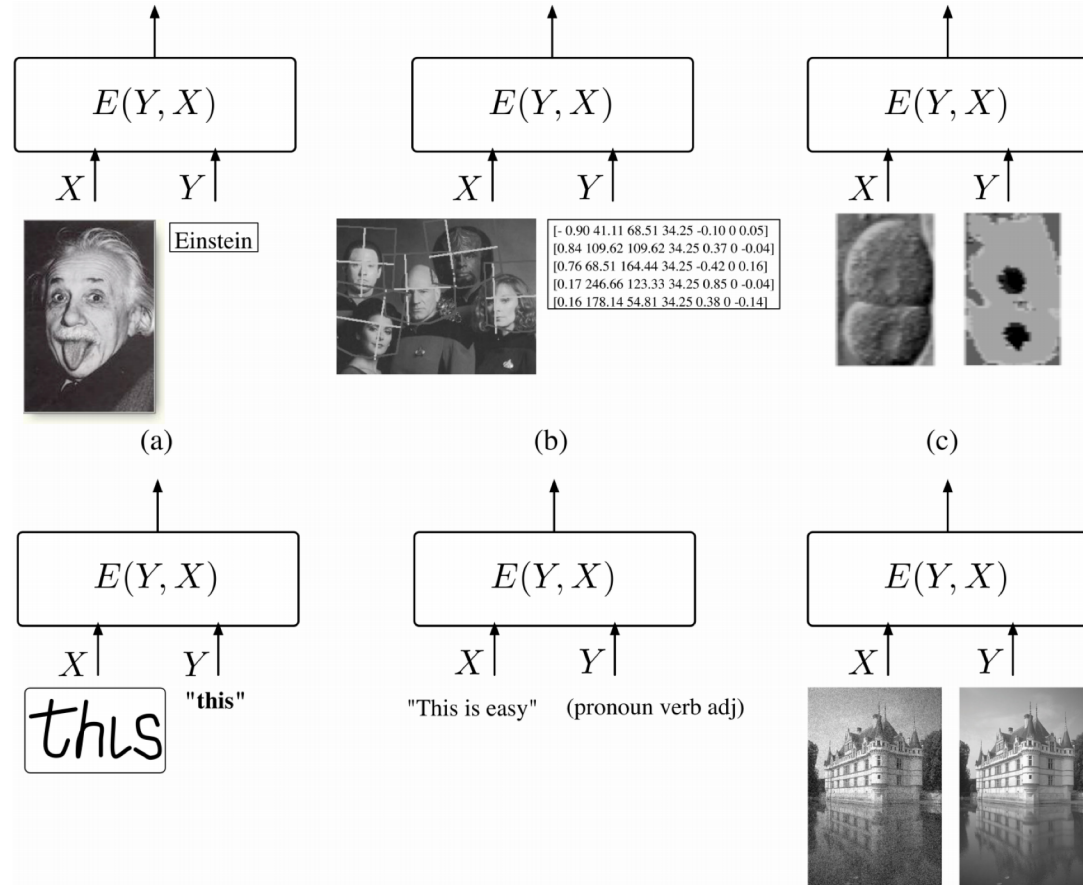


[Zellers2018]

# Energy-based models have historically been used for many tasks

Historically these have relied on **shallow energy functions** and **hand-engineered features**

We show how to use a **deep convex energy-based model** with **learned features**



[LeCun2006]

# Optimization-Based Inference

**Structured prediction:** define a network over  $\mathcal{X} \times \mathcal{Y}$  and predict via

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$

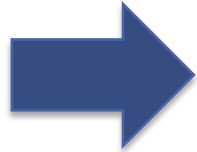
**Data imputation:** build a network over only over  $\mathcal{Y}$ , given  $y_j$  populate the remaining entries via

$$\hat{y}_{\bar{j}} = \operatorname{argmin}_{y_{\bar{j}}} f(y_{\bar{j}}, y_j; \theta)$$

**Continuous action reinforcement learning:** Represent  $Q$  function as  $Q^*(s, a) = -f(s, a; \theta)$ , policy becomes

$$\pi^*(s) = \operatorname{argmin}_a f(s, a; \theta)$$

# ICNN Portion Overview



Our Contribution: Input Convex Neural Networks

Challenges: Inference and Learning

Experiments

- Synthetic

- Multi-label Classification

- Image Completion

- Continuous-Action Q-Learning

# Input Convex Neural Networks (ICNNS)

**Definition** Scalar-valued network  $f(x, y; \theta)$  such that  $f$  is **convex** in  $y$  for all values of  $x$  (note that these networks are still **not convex** in  $\theta = \{W_i, b_i\}$ )

We can efficiently **optimize over some inputs** to the network **given other inputs**

Efficiently captures dependencies in the output space for prediction

It turns out, we don't need very many restrictions on the network to achieve this property

# Applications of Optimization for Inference

**With ICNNs:** All of these problems are convex, “easy” to solve globally

**Structured prediction:** define a network over  $\mathcal{X} \times \mathcal{Y}$  and predict via

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$

**Data imputation:** build a network over only over  $\mathcal{Y}$ , given  $y_j$  populate the remaining entries via

$$\hat{y}_{\bar{j}} = \operatorname{argmin}_{y_{\bar{j}}} f(y_{\bar{j}}, y_j; \theta)$$

**Continuous action reinforcement learning:** Represent  $Q$  function as

$Q^*(s, a) = -f(s, a; \theta)$ , policy becomes

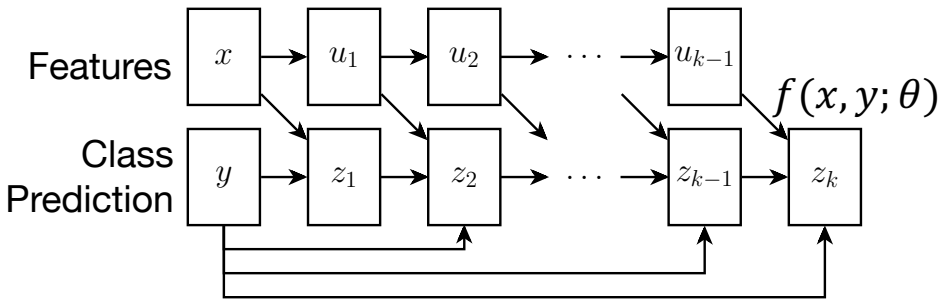
$$\pi^*(s) = \operatorname{argmin}_a f(s, a; \theta)$$



# Example Networks

ICNN for structured prediction:

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$



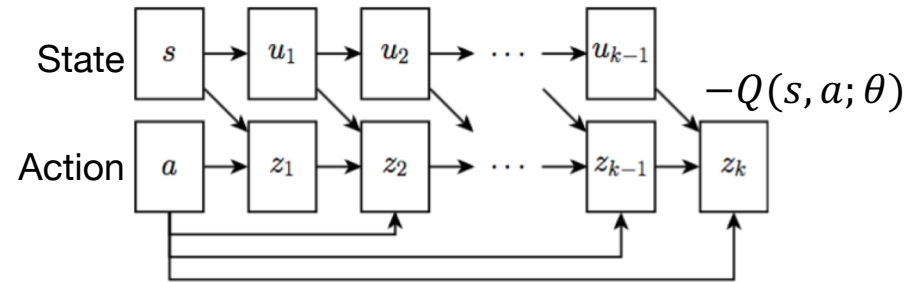
$$u_{i+1} = \tilde{g}_i(\tilde{W}_i u_i + \tilde{b}_i)$$

$$z_{i+1} = g_i \left( W_i^{(uz)}(u_i \circ z_i) + W_i^{(u)} u_i + W_i^{(z)} z_i + W_i^{(y)} y_i + b_i \right)$$

$$f(x, z; \theta) = z_k$$

ICNN for Q learning:

$$\pi^*(s) = \operatorname{argmin}_a -Q(s, a; \theta)$$



$$u_{i+1} = \tilde{g}_i(\tilde{W}_i u_i + \tilde{b}_i)$$

$$z_{i+1} = g_i \left( W_i^{(z)}(z_i \circ [W_i^{(zu)} u_i + b_i^{(z)}]_+) + \right.$$

$$\left. W_i^{(a)}(a \circ [W_i^{(au)} u_i + b_i^{(a)}]) + W_i^{(u)} u_i + b_i \right)$$

$$-Q(s, a; \theta) = f(s, a; \theta) = z_k, u_0 = s, z_0 = a$$

# How to achieve input convexity?

Most networks can be “trivially” modified to **guarantee input convexity**

Consider a simple **feedforward ReLU network**:

$$z_{i+1} = \max\{0, W_i z_i + b_i\}, \quad i = 1, \dots, k$$
$$f(y; \theta) = z_{k+1}, \quad z_1 = y$$

**Proposition.**  $f$  is convex in  $y$  provided that the  $W_i$  are non-negative for  $i > 1$

More generally, any activation function that is convex and non-decreasing also has this property.

# Is convexity restrictive?

Yes (by definition, the functions are restricted to be convex), but **not that bad** in practice

**Proposition.** ICNNs trivially subsume any feedforward network

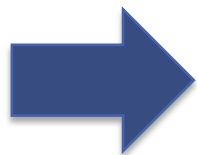
$$\tilde{f}(x) \text{ with the network } f(x, y) = (y - \tilde{f}(x))^2$$

More complex convex portion **adds additional structure** over  $y$ , which can still be “easily” optimized over

We’ll see more evidence for this later

# ICNN Portion Overview

Our Contribution: Input Convex Neural Networks



Challenges: Inference and Learning

Experiments

- Synthetic

- Multi-label Classification

- Image Completion

- Continuous-Action Q-Learning

# Challenges for ICNNs

**Inference:** how do we efficiently perform the optimization?

$$y^*(x; \theta) = \operatorname{argmin}_y f(x, y; \theta)$$

**Learning:** How do we train the network (find  $\theta$ ) such that it gives good predictions?

$$\operatorname{minimize}_{\theta} \sum_{i=1}^n \ell(y_i, y^*(x_i; \theta))$$

# Inference in ICNNs

In theory, inference in ICNNs is just a linear program

$$\begin{aligned} \min_y f(y; \theta) &= \min_{y,z} z_{k+1} \\ \text{s.t. } z_{i+1} &\geq W_i z_i + b_i \\ z_i &\geq 0 \text{ for } i > 1 \\ z_1 &= y \end{aligned}$$

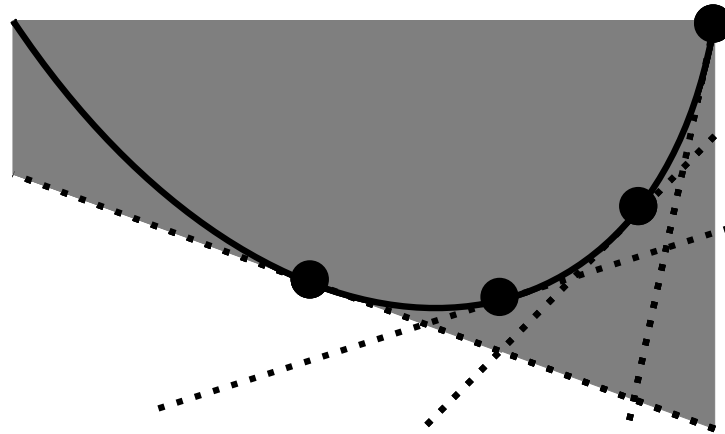
This program has **as many variables as hidden units** in the network, exact solution methods require that we invert the  $W_i^T W_i$  matrices

Instead, exploit the fact that we can easily compute the gradient of  $f(x, y; \theta)$  with respect to  $y$  (this is just **backprop**), and **optimize using gradient-based methods**

We found that the **bundle method** (defined on the next slide) performs better than gradient descent in some cases

# Inference with the Bundle Method

Repeatedly minimize a lower bound on the function



Uses convexity to minimize more quickly than gradient descent

Boundary constraints are difficult, so we actually use an entropy penalty  
$$\tilde{f}(x, y; \theta) + y \log y + (1 - y) \log(1 - y)$$

# ICNN Learning

Two possibilities for training networks

1. **Max-margin structured prediction:** enforce constraint that

$$f(x_i, y_i; \theta) \leq \operatorname{argmin}_y (f(x_i, y; \theta) + \Delta(y, y_i))$$

Common structured prediction approach

Margin-scaling term  $\Delta(y, y_i)$  can be finicky

2. **Argmin differentiation,** directly compute

$$\nabla_{\theta} \ell(y_i, y^*(x_i; \theta))$$

Can be approximated by unrolling an optimization procedure

Plays nicely with bundle method and approximate optimization

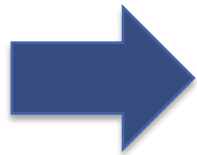
May require some differential calculus (nothing too nasty)



# ICNN Portion Overview

Our Contribution: Input Convex Neural Networks

Challenges: Inference and Learning



Experiments

Synthetic

Multi-label Classification

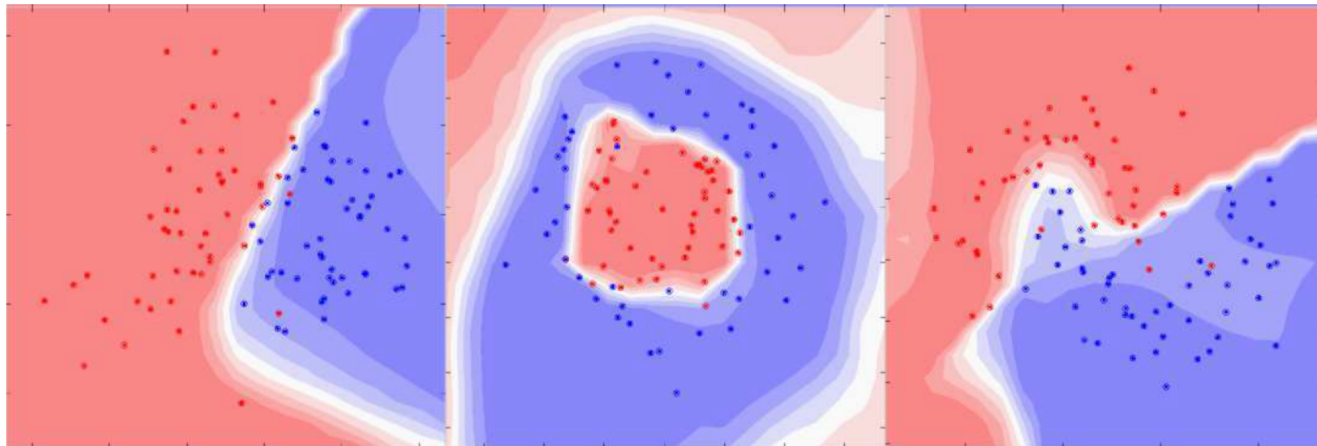
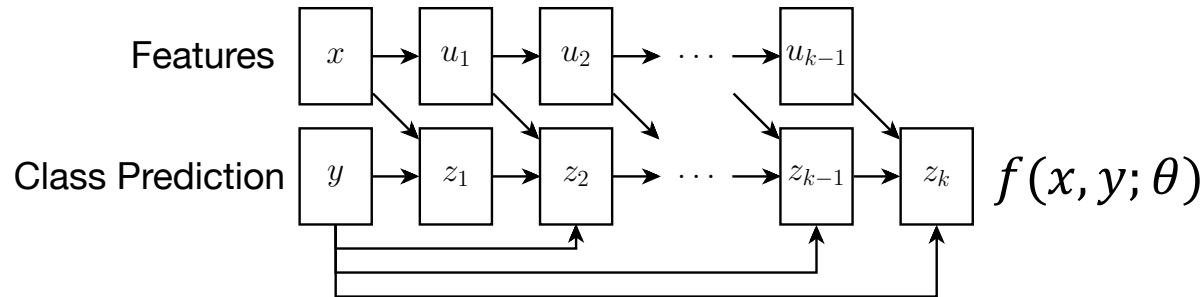
Image Completion

Continuous-Action Q-Learning

# Results: toy example

Partially input convex neural network trained to classify points in 2D space

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$



Only point to remember from this: **convex energy function does *not* imply a convex decision boundary**; **argmin operator is a powerful one**

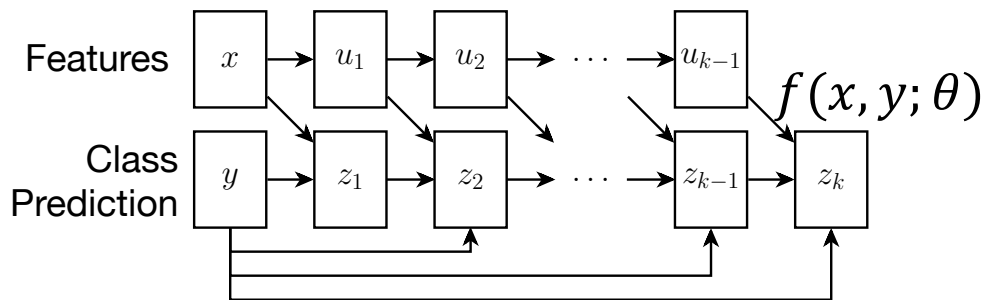
# Results: multi-label classification

**Task:** Predict tags for bibtex entries from bag of words features

Used in Belanger and McCallum, 2016: Structured Prediction Energy Networks

ICNNs almost recover the same performance as SPENs despite the convexity restrictions

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$



(Higher = Better)

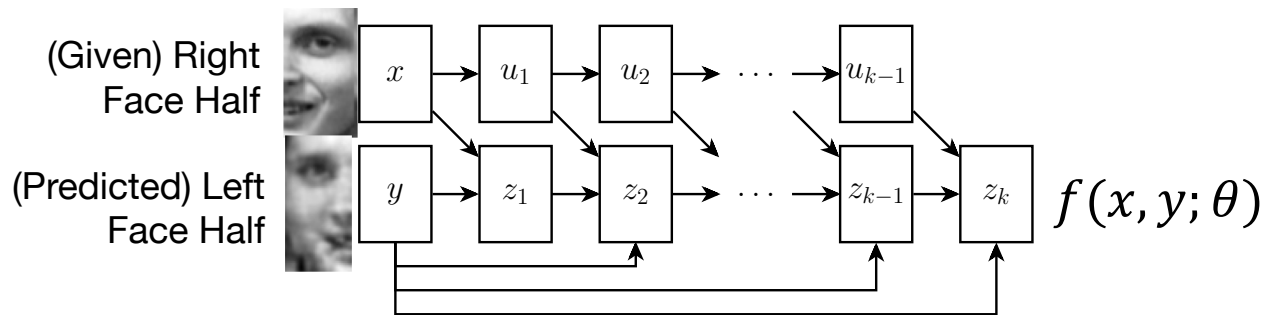
Method	Test Macro-F1
NN (Baseline)	0.396
<b>SPEN</b>	<b>0.422</b>
ICNN	0.415

# Results: image completion

**Task:** Predict the left side of the image given the right side. Used in Poon and Domingos 2011; Sum-Product Networks

**ICNN:** DQN-like network over both input and output

$$\hat{y}(x) = \operatorname{argmin}_y f(x, y; \theta)$$



ICNN Test Set Completions



Method	MSE
Sum-Product Network Baseline [PD11]	942.0
Dilated CNN Baseline [YK15]	800.0
FCN Baseline [LSD15]	<b>795.4</b>
ICNN - Bundle Entropy	833.0
ICNN - Gradient Decent	872.0
ICNN - Nonconvex	850.9

# Input Convex Neural Networks

Brandon Amos   Lei Xu   J. Zico Kolter  
Carnegie Mellon University  
School of Computer Science

---

Our Contribution: **Input Convex Neural Networks**

Challenges: Inference and Learning

Experiments

1. Synthetic
2. Multi-label Classification
3. Image Completion
4. Continuous–Action Q-Learning



The full TensorFlow source code to reproduce all of our experiments is available online at <https://github.com/locuslab/icnn>