

Lifelong Learning

CS 330

Course Reminders

Optional homework 4 due **today**.

Project milestone due **Wednesday**.

Guest lecture on Wednesday!

Hanie Sedghi

Please try to show up in person & on-time.

Plan for Today

The lifelong learning **problem statement**

Basic approaches to lifelong learning

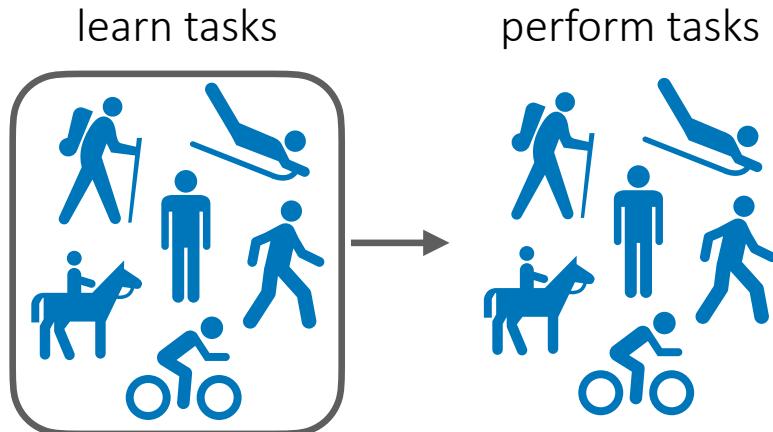
Can we do **better** than the basics?

Revisiting the problem statement
from **the meta-learning perspective**

A brief review of problem statements.

Multi-Task Learning

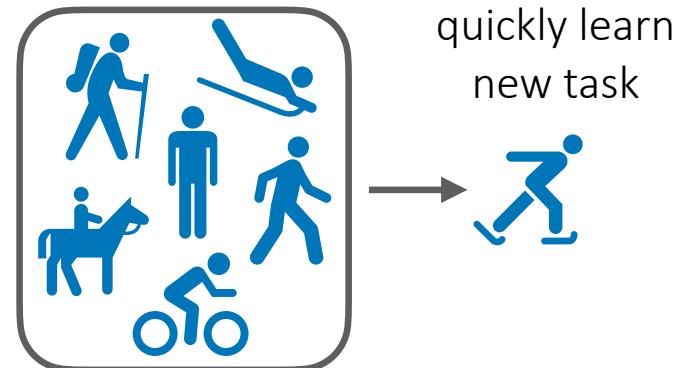
Learn to solve a set of tasks.



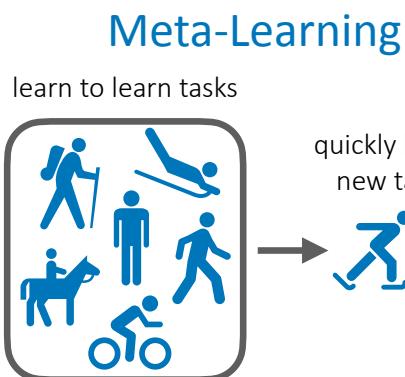
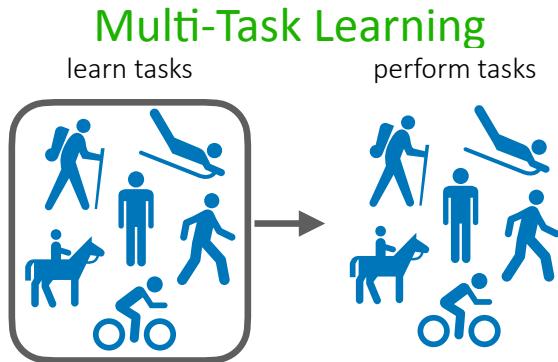
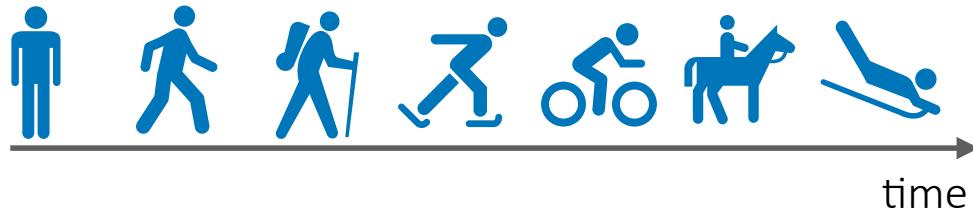
Meta-Learning

Given i.i.d. task distribution,
learn a new task efficiently

learn to learn tasks



In contrast, many real world settings look like:



Some examples:

- a student learning concepts in school
- a deployed **image classification system** learning from a stream of images from users
- a **robot** acquiring an increasingly large set of skills in different environments
- a **virtual assistant** learning to help different users with different tasks at different points in time
- a **doctor's assistant** aiding in medical decision-making

→ curriculum that increases in difficulty. ~ subsequent tasks build on previous tasks

Some Terminology

mostly
synonyms

Sequential learning settings

online learning, lifelong learning, continual learning, incremental learning, streaming data

≠

distinct from sequence data and sequential decision-making

↳ Data points have
a sequential structure

↳ system is making multiple decisions in sequence

≠

Learning as the
data comes in
(sequential / continual / lifelong learning)

What is the lifelong learning *problem statement*?

- Exercise:**
1. Pick an example setting.
 2. Discuss problem statement in small groups:
 - (a) how would you set-up an experiment to develop & test your algorithm?
 - (b) what are desirable/required properties of the algorithm?
 - (c) how do you evaluate such a system?

 - A. a student learning concepts in school
 - B. a deployed **image classification system** learning from a stream of images from users
 - C. a **robot** acquiring an increasingly large set of skills in different environments
 - D. a **virtual assistant** learning to help different users with different tasks at different points in time
 - E. a doctor's **assistant** aiding in medical decision-making

What is the lifelong learning *problem statement*?

Desireable Properties

- good performance as # examples grows
(\rightarrow asymptote to perfect) - ability to learn quickly.
- good perf on past tasks - adaptability to new objectives
- don't forget past tasks/envs
- capacity - depth + breadth of knowledge.

Evaluation

- # of images for a user (learning efficiency)
- performance on both past tasks & new tasks
- generalizability
- delta in performance.
- performance on last task

What is the lifelong learning *problem statement*?

Problem variations:

- task/data order: i.i.d. vs. predictable vs. curriculum vs. adversarial
 - e.g.) detect spam over time
 - spammers might make ways to get them through filters
 - Adversarial setting
- discrete task boundaries vs. continuous shifts (vs. both)
 - New users
 - e.g.) people's sentiments on certain topics, seasons, etc.
- known task boundaries/shifts vs. unknown

Some considerations:

- model performance
- data efficiency → ability to learn tasks efficiently over time
- computational resources
- memory → May not be able to simultaneously see all accumulated data (history)
 - too big
 - may not be able to store data over a fixed period of time.
 - Data silos
- others: privacy, interpretability, fairness, test time compute & memory

"Differential Privacy"
⇒ Set up the model s.t.
you can't get the data
back out from the model.

Substantial variety in problem statement!

What is the lifelong learning *problem statement*?

General [supervised] online learning problem:

for $t = 1, \dots, n$

observe x_t

predict \hat{y}_t → Estimated label

observe label y_t → For some problem there might be a delay in observation, or the obs. might be limited/constrained.

i.i.d. setting: $x_t \sim p(x)$, $y_t \sim p(y|x)$

p not a function of t

otherwise: $x_t \sim p_t(x)$, $y_t \sim p_t(y|x)$

changing over time

<— if observable task boundaries: observe x_t, z_t

→ Images, patient record, etc.

Q: For deep learning models, couldn't this take an unreasonably long time
dec. there are so many params?

A: Yes, so in practice you'll warm-start this process w/ an existing dataset
to initialize the predictor rather than start completely from scratch.

streaming setting: cannot store (x_t, y_t)

- lack of memory → too much data (too big)
- lack of computational resources → too much data coming in at the same time.
- privacy considerations → e.g. studying the brain
- want to study neural memory mechanisms

true in some cases, but not in many cases!

9 -

recall: replay buffers

If you're seeing these datapoints from a sequence of tasks, then instead of observing just x_t , you can also observe the task identifier that corresponds to that task.

What do you want from your lifelong learning algorithm?

Q: Calculate best learner (in hindsight)
Using the last (best) parameter θ ?

A: Note, we can't evaluate regret in practice b/c.
we don't have access to the best possible params
 (θ^*) . → Approximate θ^* @ the end of
training by training a model using all
observed data.

Q: Could the regret value < 0 ?

A: Typically in i.i.d. settings, there exists a
single model that can do well on all
of the data. But in non-i.i.d. setting,
you can get a < 0 regret.

minimal regret (that grows slowly with t)

regret: cumulative loss of learner — cumulative loss of best learner in hindsight

$$\text{Regret}_T := \sum_1^T \mathcal{L}_t(\theta_t) - \min_{\theta} \sum_1^T \mathcal{L}_t(\theta)$$

↑ no subscript!

→ Function of T (capital)

(cannot be evaluated in practice, useful for analysis)

Regret that grows linearly in t is trivial. Why?

↪ Constant θ_t through t will give this result → Model's not being updated.

What do you want from your lifelong learning algorithm?

minimal regret (that grows slowly with t)

regret: cumulative loss of learner — cumulative loss of best learner in hindsight



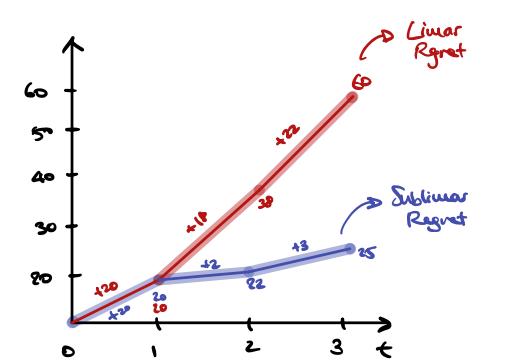
$$\text{Regret}_T := \sum_t^T \mathcal{L}_t(\theta_t) - \min_{\theta} \sum_1^T \mathcal{L}_t(\theta)$$

@ $t=T (=3)$,

$$20 + 18 + 22$$

$$20 + 2 + 3$$

t	A	B
1	20	20
2	38	22
3	60	25



We want Sublinear regret!

=> Means we're getting better at making predictions on your online learning prob.

↳ Avg. of all prev. values is a better predictor than constant (0) predictor

What do you want from your lifelong learning algorithm?

positive & negative transfer

positive **forward** transfer: previous tasks cause you to do better on future tasks compared to learning future tasks from scratch

positive **backward** transfer: current tasks cause you to do better on previous tasks compared to learning past tasks from scratch

→ e.g.) For each task you don't have much data (e.g. 100 data pts/task)
 → Once you accumulate enough data from all other tasks observed so far, you may do better on the 1st task.

positive \rightarrow negative : better \rightarrow worse

→ i.e. negative forward transfer : prev. tasks cause you to do worse on future tasks.
 etc.

→ aka "catastrophic forgetting"

Negative Backward Transfer is quite common if you have a small memory
 (basically start forgetting the past tasks)

Positive Forward Transfer is also common → kind of like pretraining on the previous tasks rather than starting from scratch.

Plan for Today

The lifelong learning **problem statement**

Basic approaches to lifelong learning

Can we do **better** than the basics?

Revisiting the problem statement
from **the meta-learning perspective**

Approaches

$$\left\{ \begin{array}{l} @ t=T, \\ D_T = \bigcup_{t=1}^T (x_t, y_t) \\ \min_{\theta} \mathcal{L}(\theta, D_T) \end{array} \right.$$

Store all the data you've seen so far, and train on it. → follow the leader algorithm

+ will achieve very strong performance

- computation intensive → Continuous fine-tuning can help.

- can be memory intensive [depends on the application]

Take a gradient step on the datapoint you observe.

→ stochastic gradient descent

$$@ t=T, \\ \theta = \theta - \nabla_{\theta} \mathcal{L}(\theta, \{x_T, y_T\}) \rightarrow \text{Update the model using just the dataset retrieved at } t=T.$$

+ computationally cheap

+ requires 0 memory

⌚ - subject to negative backward transfer

“forgetting”

- slow learning

Common when you have
small memory

sometimes referred to as
catastrophic forgetting

Can we do better?

The above methods
actually work pretty
well in a lot of cases

Applying a simple continual learning algorithm to robotics



7 robots collected 580k grasps

86%
↳ Grasp success rate

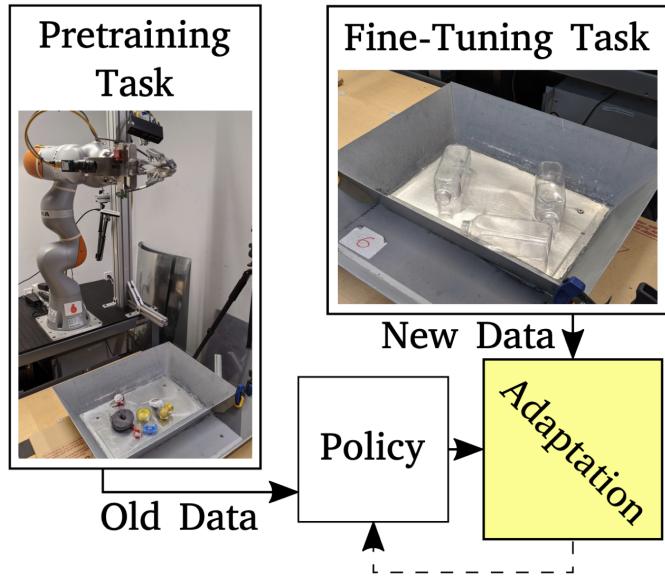


49%

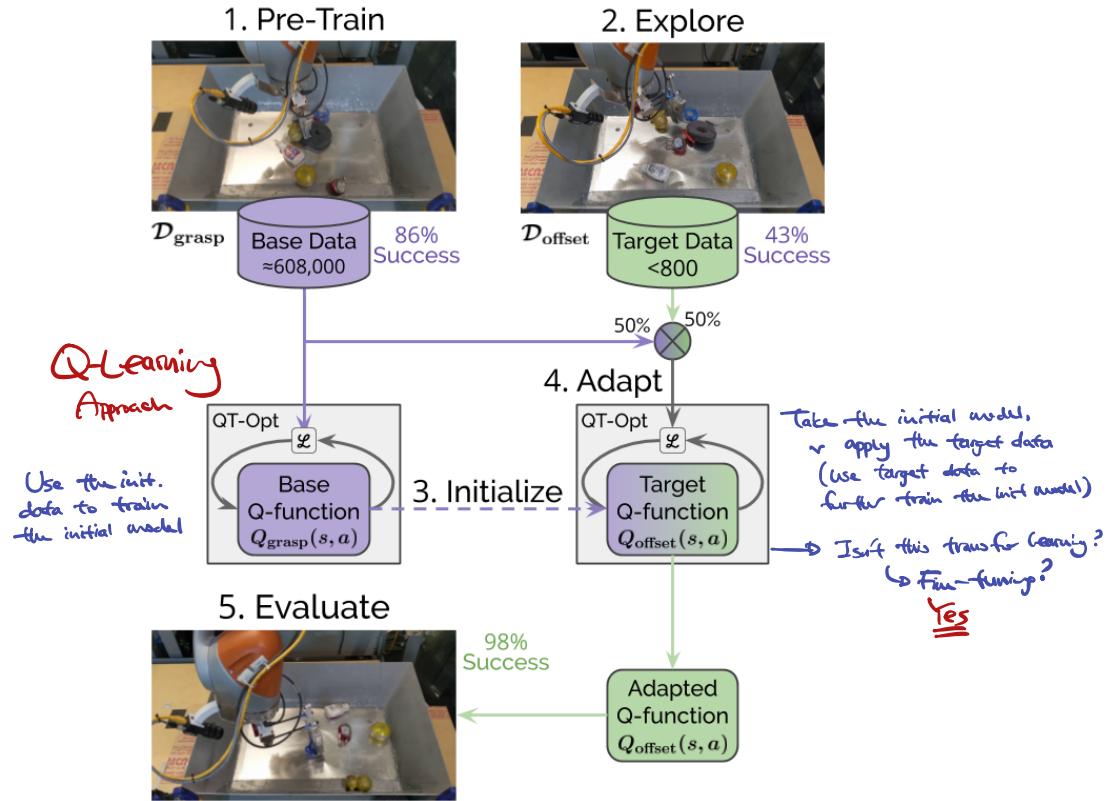
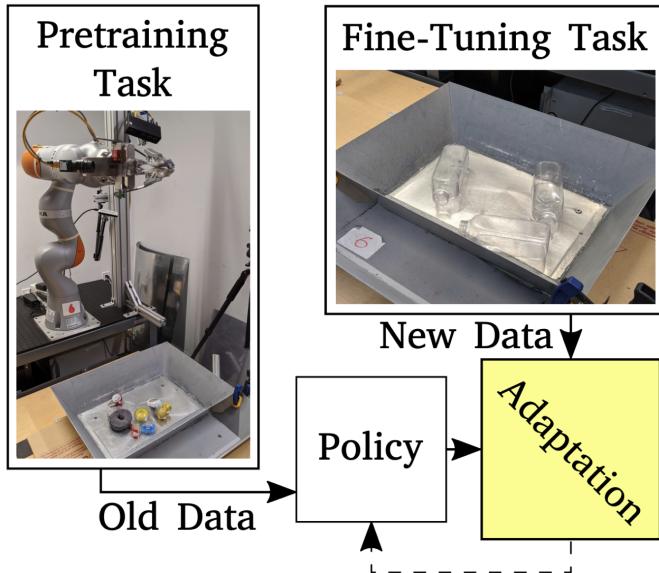


↳ Lowered success rate
for different-looking (new)
objects

Applying a simple continual learning algorithm to robotics



Applying a simple continual learning algorithm to robotics

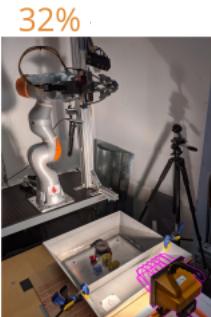


Applying a simple continual learning algorithm to robotics

Pre-Train



Object
Grasping



Harsh
Lighting



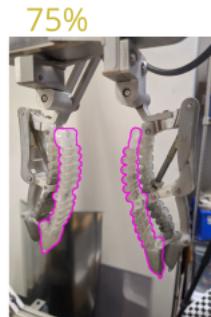
Transparent
Bottles

Fine-Tune

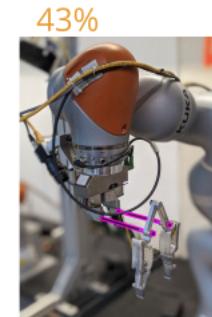


Checkerboard
Backing

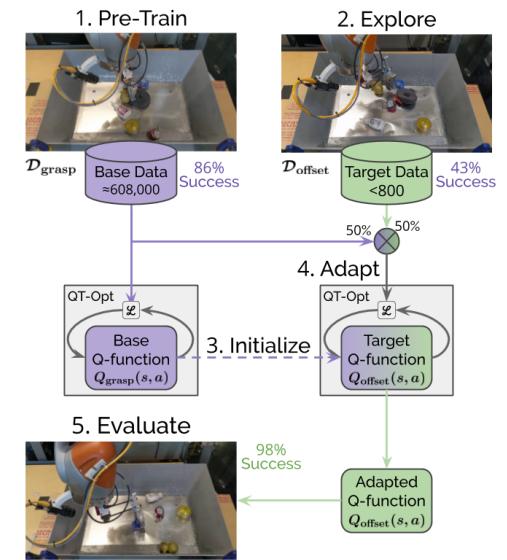
Yeah, it's just
fine-tuning.
⇒ Works pretty well.



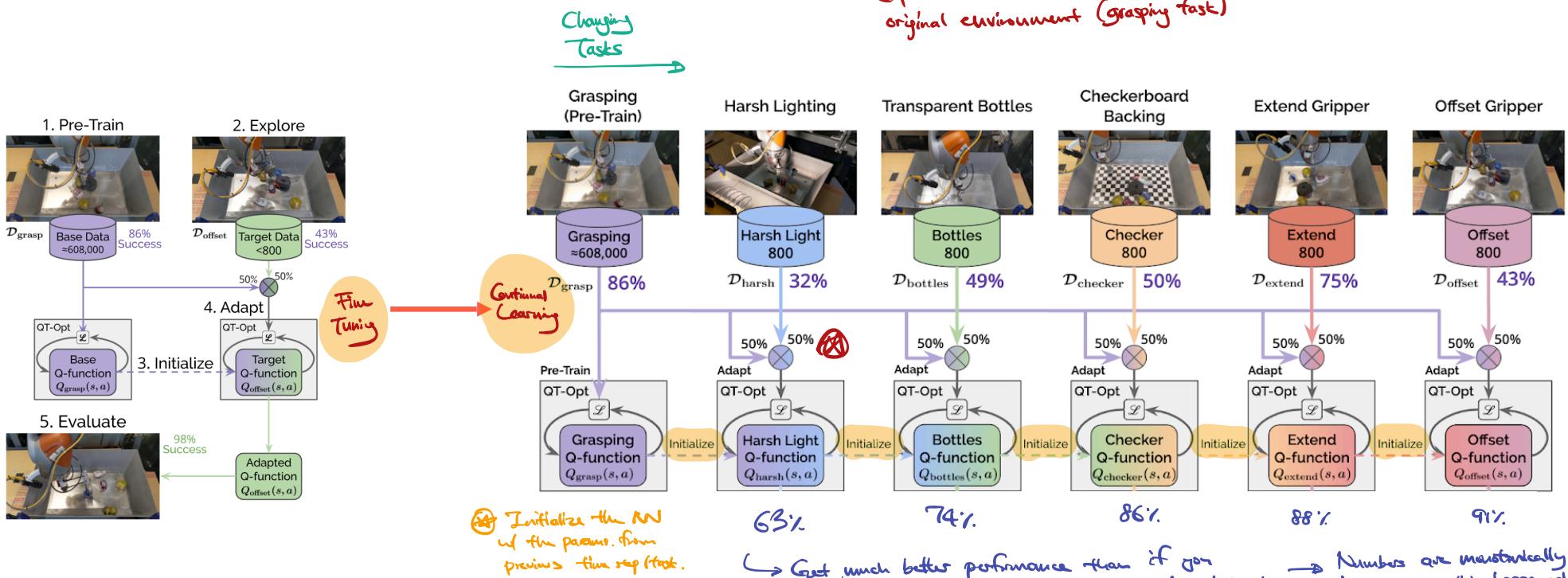
Extend Gripper
1cm



Offset Gripper
10cm



Applying a simple continual learning algorithm to robotics



What about backward transfer?

↳ aka "catastrophic forget"

Can we do better?

↳ Avoid ←ve backward transfer while using a small amount of memory

Plan for Today

The lifelong learning **problem statement**

Basic approaches to lifelong learning

Can we do **better** than the basics?

Revisiting the problem statement
from **the meta-learning perspective**

Case Study: Can we modify vanilla SGD to avoid negative backward transfer?
(from scratch)



Idea:

(1) store small amount of data per task in memory

(2) when making updates for new tasks, ensure that they don't unlearn previous tasks

How do we accomplish (2)?

learning predictor $y_t = f_\theta(x_t, z_t)$

memory: \mathcal{M}_k for task z_k

To as few as 5 examples

For $t = 0, \dots, T$

minimize $\mathcal{L}(f_\theta(\cdot, z_t), (x_t, y_t))$

subject to $\mathcal{L}(f_\theta, \mathcal{M}_k) \leq \mathcal{L}(f_\theta^{t-1}, \mathcal{M}_k)$ for all $k < t$

Instead of placing a constraint on the loss function directly, place it on the gradients.

Assume local linearity:

If the gradients are pointing in similar directions, you'll get (true) backwards transfer.

Ensure that the gradient of the func. points in the same direction ~ orthogonal to gradients for the previous task.

inner product

$$\langle g_t, g_k \rangle := \left\langle \frac{\partial \mathcal{L}(f_\theta, (x_t, y_t))}{\partial \theta}, \frac{\partial \mathcal{L}(f_\theta, \mathcal{M}_k)}{\partial \theta} \right\rangle$$

Quadratic Programming

Can formulate & solve as a QP.

Assuming local linearity, \rightarrow means that the loss for prev. task hasn't increased.

(i.e. s.t. loss on previous tasks doesn't get worse)

\hookrightarrow If you're making a small update to the model, and if it doesn't make the predictions on some of the datasets worse, it probably wouldn't make other things a lot worse either.

≥ 0 for all $z_k < z_t$

② It's possible that this constraint leads to an infeasible constrained optimization. (i.e. current prev. task's gradients might be pointing to opposite ways). \rightarrow may be increasingly infeasible as the no. of past tasks T . (\because increase the no. of constraints on the optimization).

Experiments

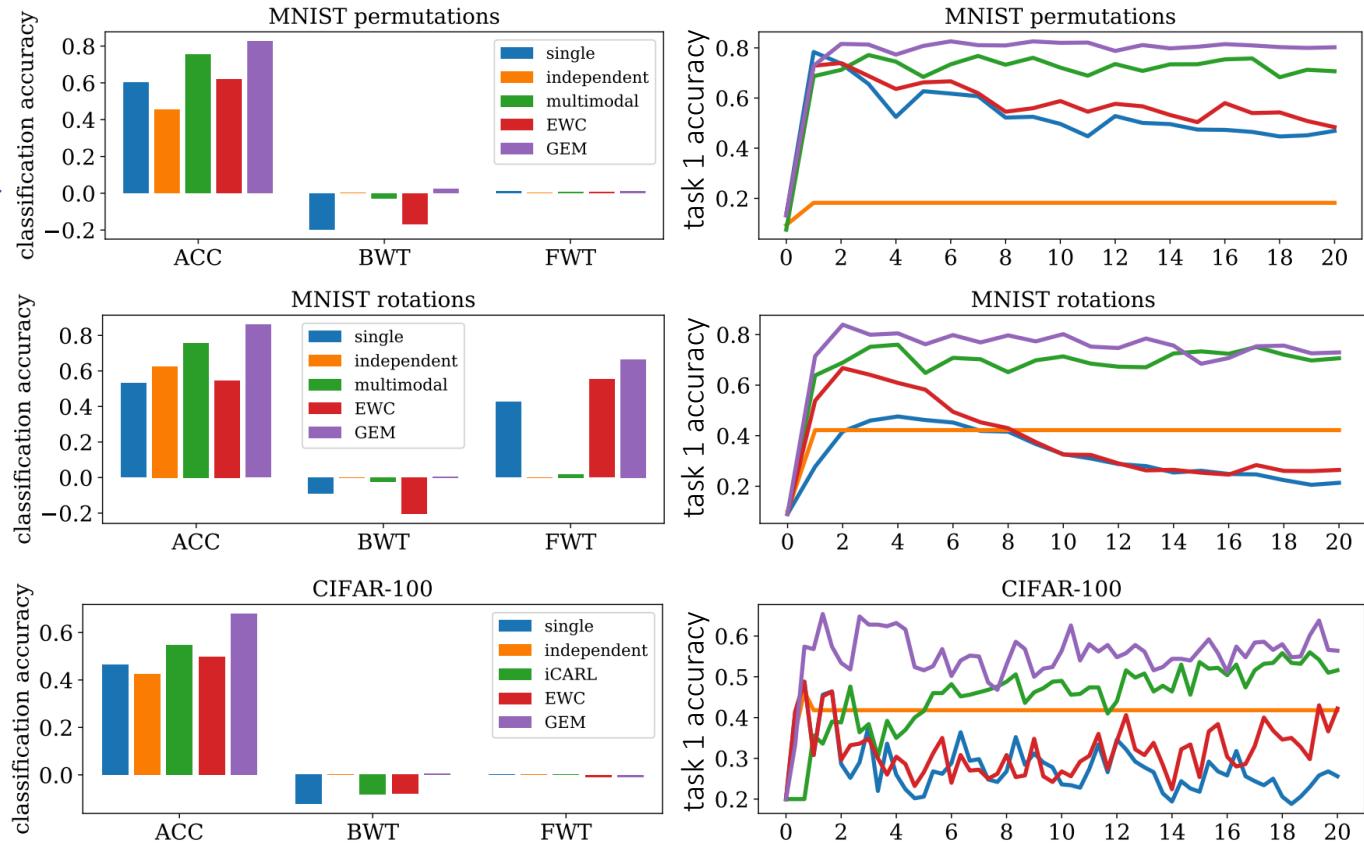
Problems:

- MNIST permutations
- MNIST rotations
- CIFAR-100 (5 new classes/task)

Each task has a diff. permutation of the pixels in the image
Randomly rearranging the pixels

BWT: backward transfer,
FWT: forward transfer

Total memory size:
5012 examples



If we take a step back...

do these experimental domains make sense?

In practice we can't store the data in RAM/VRAM, so the above results may not apply. → 7116731

Can we meta-learn how to avoid negative backward transfer?

→ Update rules s.t. it doesn't get ^{→ve} backwards transfer

Javed & White. *Meta-Learning Representations for Continual Learning*. NeurIPS '19

Beaulieu et al. *Learning to Continually Learn*. '20

Plan for Today

The lifelong learning **problem statement**

Basic approaches to lifelong learning

Can we do **better** than the basics?

Revisiting the problem statement
from **the meta-learning perspective**

Formulation of online learning when faced with sequence of tasks

Online Learning

(Hannan '57, Zinkevich '03)

Perform sequence of tasks
while minimizing static regret.

perform perform perform perform perform perform perform



zero-shot performance

↳ often difficult to perform well zero-shot
on a completely new task.

More realistically:

Given
or
small
amount
of data

learn learn learn learn learn learn



slow learning

rapid learning

Faster learning over time

Formulation of online learning when faced with sequence of tasks

Online Learning

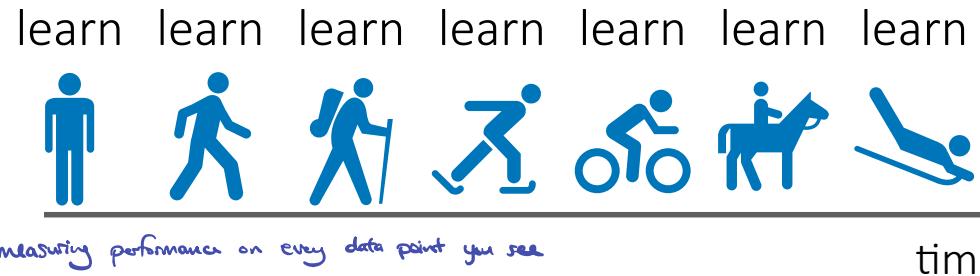
(Hannan '57, Zinkevich '03)

Perform sequence of tasks
while minimizing static regret.



Online Meta-Learning

Efficiently learn a sequence of tasks
from a non-stationary distribution.



→ instead of measuring performance on every data point you see

evaluate performance after seeing a small amount of data

Primarily a difference in *evaluation*, rather than the *data stream*.

The Online Meta-Learning Setting

for task $t = 1, \dots, n$

observe $\mathcal{D}_t^{\text{tr}}$

use update procedure $\Phi(\theta_t, \mathcal{D}_t^{\text{tr}})$ to produce parameters ϕ_t

observe x_t

predict $\hat{y}_t = f_{\phi_t}(x_t)$

observe label y_t

Standard online learning setting \rightarrow Identical

Initially given a period
to try to learn the task
w/ a small amount of data

Goal: Learning algorithm with sub-linear

$$\text{Regret}_T := \sum_{t=1}^T \ell_t(\Phi_t(\theta_t)) - \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\Phi_t(\theta))$$

Loss of algorithm

Loss of best algorithm
in hindsight

Can we apply meta-learning in lifelong learning settings?

Recall the **follow the leader** (FTL) algorithm:

Store all the data you've seen so far, and train on it.

Deploy model on current task.

Follow the *meta*-leader (FTML) algorithm:

Store all the data you've seen so far, and **meta-train** on it.

Run update procedure on the current task.

What meta-learning algorithms are well-suited for FTML?

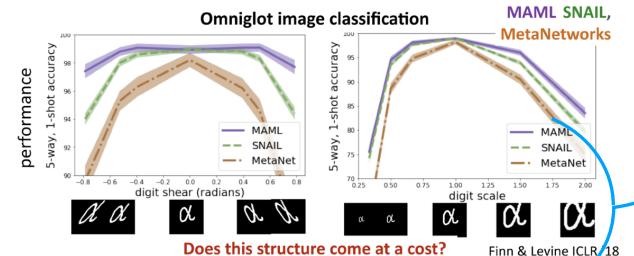
What if $p_t(\mathcal{T})$ is non-stationary?

Use optimization-based meta learners
 \therefore it's expected to do well on out-of-distrib. tasks

Warm-start the metaparams w/ the metaparams from the previous time step.

Optimization vs. Black-Box Adaptation

How well can learning procedures generalize to similar, but extrapolated tasks?

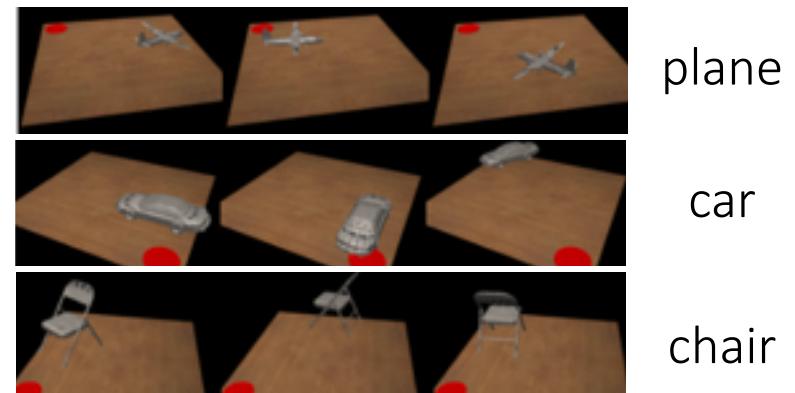


Online meta-learning experiments

Experiment with **sequences of tasks**:

- Colored, rotated, scaled **MNIST**
- **3D object pose prediction**
- **CIFAR-100** classification

Example pose prediction tasks



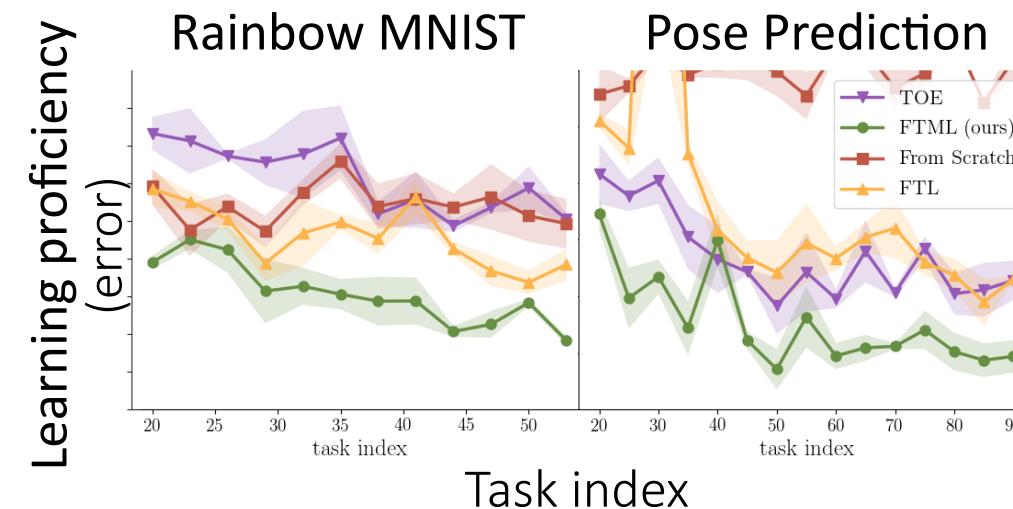
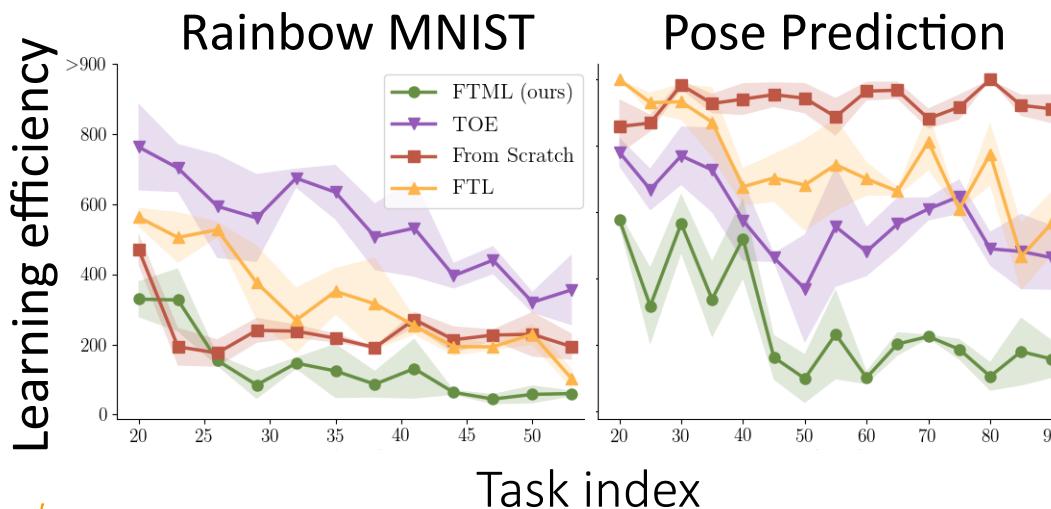
plane

car

chair

Online meta-learning experiments

- Comparisons:
- **TOE (train on everything)**: train on all data so far
 - **FTL (follow the leader)**: train on all data so far, fine-tune on current task
 - **From Scratch**: train from scratch on each task



↳ Meta-learning approaches applicable to online learning

Follow The Meta-Leader

learns each new task faster & with greater proficiency,
approaches **few-shot learning** regime

Takeaways

Many flavors of lifelong learning, all under the same name.

Defining the problem statement is often the hardest part

Meta-learning can be viewed as a slice of the lifelong learning problem.

A very open area of research.

Course Reminders

Optional homework 4 due **today**.

Project milestone due **Wednesday**.

Guest lecture on Wednesday!

Hanie Sedghi

Please try to show up in person & on-time.