

Black-Box Meta-Learning

CS 330

Logistics

Project group form due **Monday, October 10**

Homework 1 due **Wednesday October 12**

Plan for Today

Meta-Learning

- Problem formulation
- General recipe of meta-learning algorithms
- Black-box adaptation approaches
- Case study of GPT-3 (time-permitting)

}

Topic of Homework 1!

Goals for by the end of lecture:

- Training set-up for few-shot meta-learning algorithms
- How to implement black-box meta-learning techniques

Plan for Today

Meta-Learning

- **Problem formulation**

- General recipe of meta-learning algorithms
- Black-box adaptation approaches
- Case study of GPT-3 (time-permitting)

Meta-Learning Problem

Transfer Learning with Many Source Tasks

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, solve new task $\mathcal{T}_{\text{test}}$ more quickly / proficiently / stably

Key assumption: meta-training tasks and meta-test task drawn i.i.d. from same task distribution

$$\mathcal{T}_1, \dots, \mathcal{T}_n \sim p(\mathcal{T}), \mathcal{T}_{\text{test}} \sim p(\mathcal{T})$$

Like before, tasks must share structure. → &

What do the tasks correspond to?

- recognizing handwritten digits from different languages (see homework 1!)
 - giving feedback to students on different exams
 - classifying species in different regions of the world
 - a robot performing different tasks

Nice thing about meta learning is that you can get away w/ much less data per task than if you were to train completely from scratch ::



as long as
they share
structure!

How many tasks do you need?

The more the better.
5

(analogous to more data in ML)

Two ways to view meta-learning algorithms

Mechanistic view

- Deep network that can read in an entire dataset and make predictions for new datapoints
- Training this network uses a meta-dataset, which itself consists of many datasets, each for a different task

Probabilistic view

- Extract prior knowledge from a set of tasks that allows efficient learning of new tasks
- Learning a new task uses this prior and (small) training set to infer most likely posterior parameters

How does meta-learning work? An example.

Given 1 example of 5 classes:



training data $\mathcal{D}_{\text{train}}$

Classify new examples



test set \mathbf{x}_{test}

How does meta-learning work? An example.



Given 1 example of 5 classes:

meta-testing $\mathcal{T}_{\text{test}}$

Target
Tasks

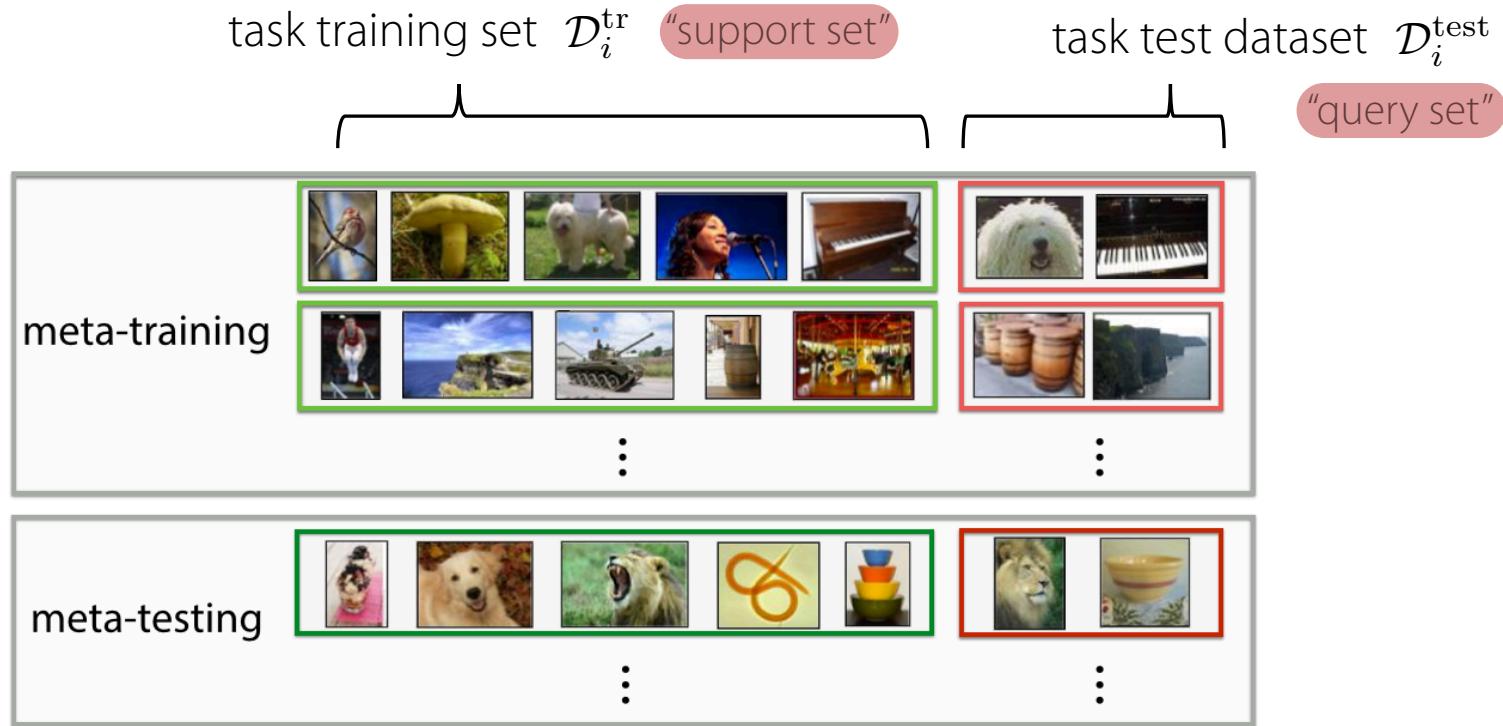


Classify new examples

any ML
problem

Can replace image classification with: regression, language generation, skill learning,

Some terminology



k-shot learning: learning with **k** examples per class
(or **k** examples total for regression)

N-way classification: choosing between **N** classes

Question: What are k and N for the above example?

$$\begin{array}{l} k=1 \\ N=5 \end{array}$$

Plan for Today

Transfer Learning

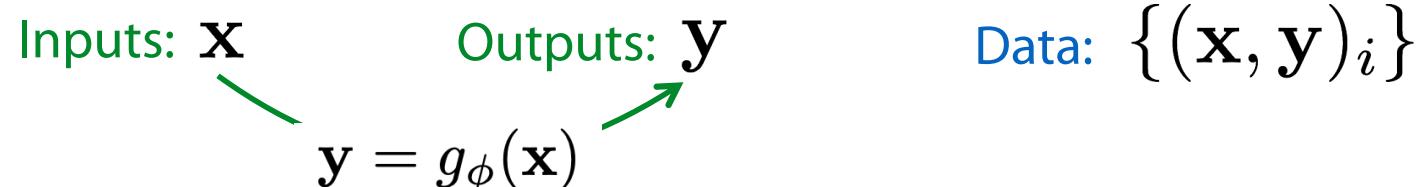
- Problem formulation
- Fine-tuning

Meta-Learning

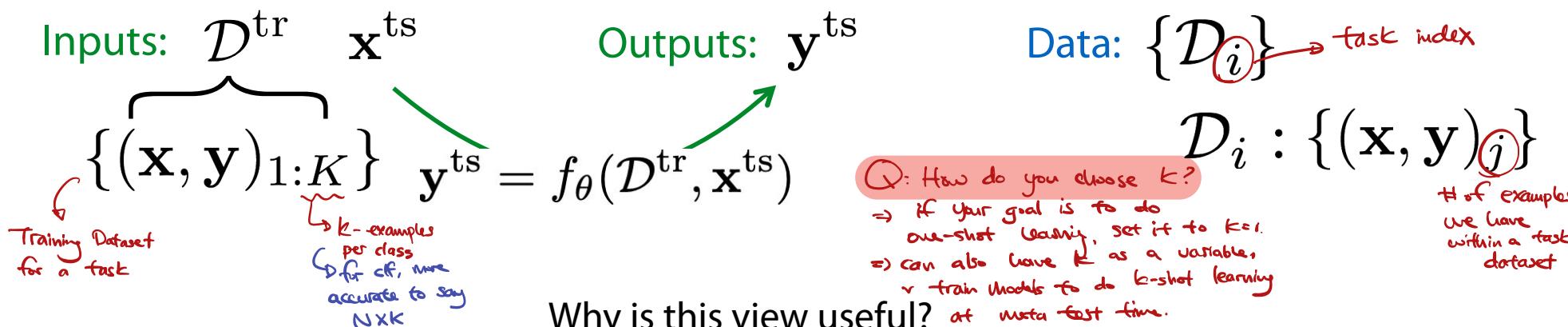
- Problem formulation
- **General recipe of meta-learning algorithms**
- Black-box adaptation approaches
- Case study of GPT-3 (time-permitting)

One View on the Meta-Learning Problem

Supervised Learning:



Meta Supervised Learning:



Why is this view useful?

Reduces the meta-learning problem to the design & optimization of f .

General recipe

How to design a meta-learning algorithm

1. Choose a form of $f_{\theta}(\mathcal{D}^{\text{tr}}, \mathbf{x}^{\text{ts}})$
2. Choose how to optimize θ w.r.t. max-likelihood objective using meta-training data

meta-parameters

Plan for Today

Meta-Learning

- Problem formulation
- General recipe of meta-learning algorithms
- **Black-box adaptation approaches**
- Case study of GPT-3 (time-permitting)

Running example

Omniglot dataset Lake et al. Science 2015

1623 characters from 50 different alphabets



20 instances of each character

Typically, supervised learning would have a lot more data points → This type might be a little more reflective of the real world

whiteboard

More few-shot image recognition datasets: tieredImageNet, CIFAR, CUB, CelebA, ORBIT, others

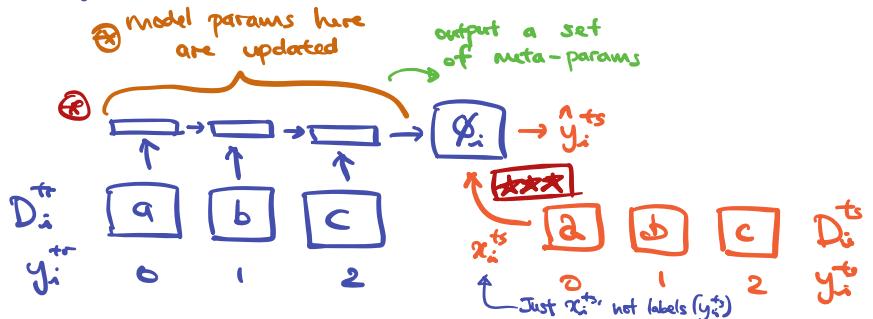
More benchmarks: molecular property prediction (Nguyen et al. '20), object pose prediction (Yin et al. ICLR '20), channel coding (Li et al. '21)

3-Way, 1-Shot

→ 3 :: it's a 3-way clfr. prob.

1. Sample task :: (3 char)
2. Sample 2 images per char → 1 tr, 1 ts
3. Put it through the process below ~ backpropagate

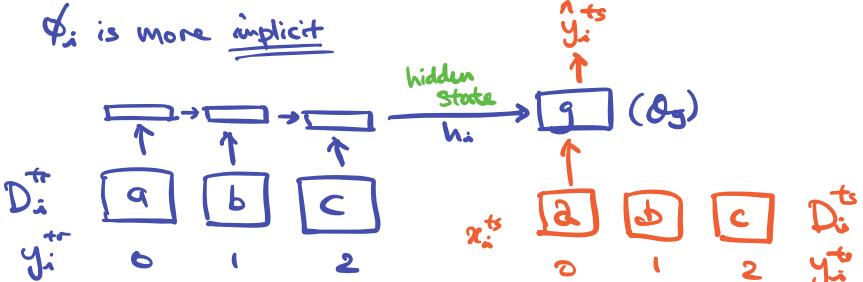
[Ver. 1]



Sequence models, like RNN, LSTM or transformers

might be a better choice than RNN or LSTM bcs. of permutation invariance

[Ver. 2]



Similar to how h_i (hidden state, activation) is not updated in Ver. 2.

ϕ_i is not updated! Model params here are.

ϕ_i is more of the activation of NN, as opposed to the model's weights

⇒ output of the neural network

Not a hypernetwork → Any NN that outputs the weights of another NN (can be used for meta learn.)

ϕ_i can be v. high dimensional

Here, we are not just inputting x_i^ts ~ comparing \hat{y}_i^ts vs y_i^ts ; we are training the model using D_i^ts , & testing w/ \hat{y}_i^ts vs y_i^ts

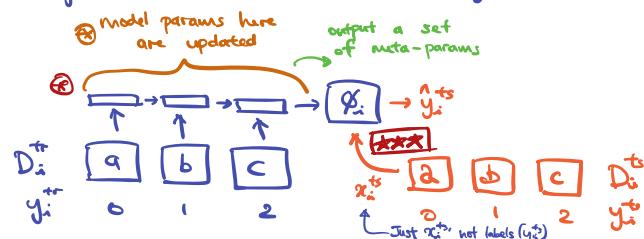
⇒ In meta learning, real tests are new tests!!!

- Ver. 2 is used more in practice bcs. it doesn't have to output large no. of params
- Task-specific parameters (ϕ_i) are the output of these NN models (both ver.)

3-Way, 1-Shot

1. Sample task \sim (3 char)
2. Sample 2 images per char \rightarrow 1 tr, 1 ts
3. Put it through the process below \sim backpropagate

[Ver. 1]



* Sequence models, like RNN, LSTM or transformer

Similar to how h_i (hidden state, activation) is not updated in Ver. 2.

ϕ_i is not updated! Model params here are.

② ϕ_i is more of the activation of NN, as opposed to the model's weights
 ⇒ output of the neural network

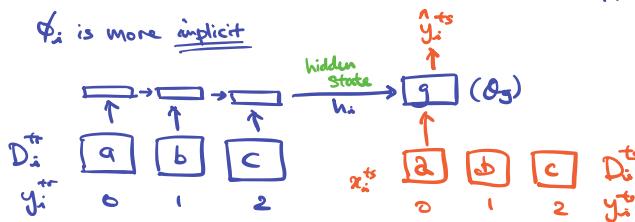
Not a hypernetwork \rightarrow Any NN that outputs the weights of another NN (can be used for meta learn.)

ϕ_i can be v. high dimensional

* Here, we are not just inputting x_i^ts & comparing y_i^ts vs y_i^tr ; we are training the model using D_i^tr , & testing w/ y_i^ts .
 ⇒ In meta learning, real tests are new tests!!!

Might be a better choice than RNN or LSTM b/c. of permutations invariance

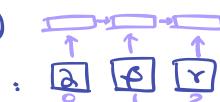
[Ver. 2]



$$\phi_i = \{h_i, \theta_g\}$$

- Ver. 2 is used more in practice b/c. it doesn't have to output large no. of params
- Task-specific parameters (ϕ_i) are the output of these NN models (both ver.)

Meta-Test

1. Given D_i^ts (test task)
2. Given D_i^tr Notation!: 
3. Given x_i^ts : 

Q: Is there any parameter updating in this process?

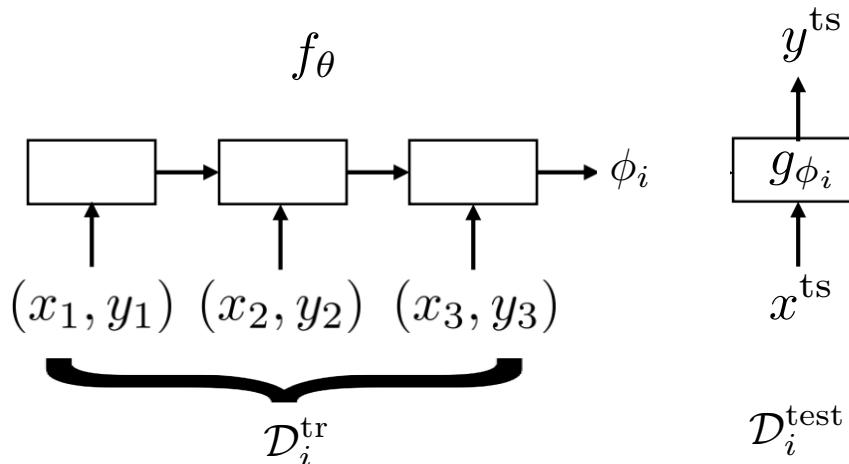
A: No. Just doing a forward pass through RNN

- You can think of the RNN as kind of implementing the learning process (give it a big enough RNN, & it will learn in that way)
- You can sort of think of the hidden state as being "updated" as the D_i^tr passes through.
 ↳ Just the D_i^tr

Black-Box Adaptation

↳ BB b/c we don't see the contents of NN

Key idea: Train a neural network to represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$ "learner"
 Predict test points with $\mathbf{y}^{\text{ts}} = g_{\phi_i}(\mathbf{x}^{\text{ts}})$



Train with standard supervised learning!

$$\min_{\theta} \sum_{\mathcal{T}_i} \sum_{(x,y) \sim \mathcal{D}_i^{\text{test}}} -\log g_{\phi_i}(y | x)$$

↗ NLL or
 cross-entropy
 $\Rightarrow \sum_{\text{all data}} \text{NLL} = \text{CE}$

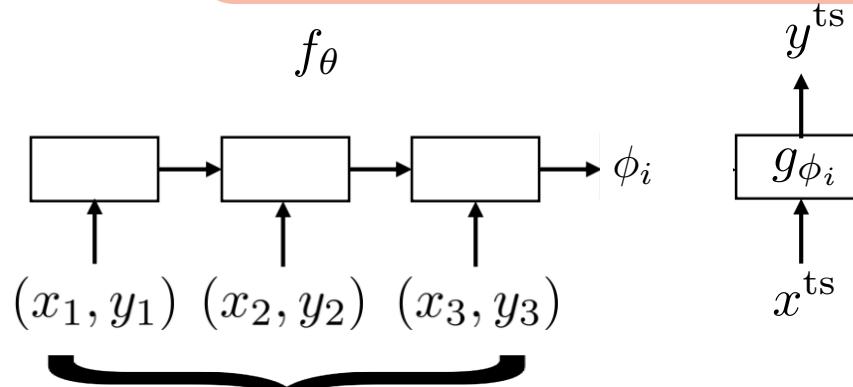
Loss function
 for Task i

$$\mathcal{L}(\phi_i, \mathcal{D}_i^{\text{test}})$$

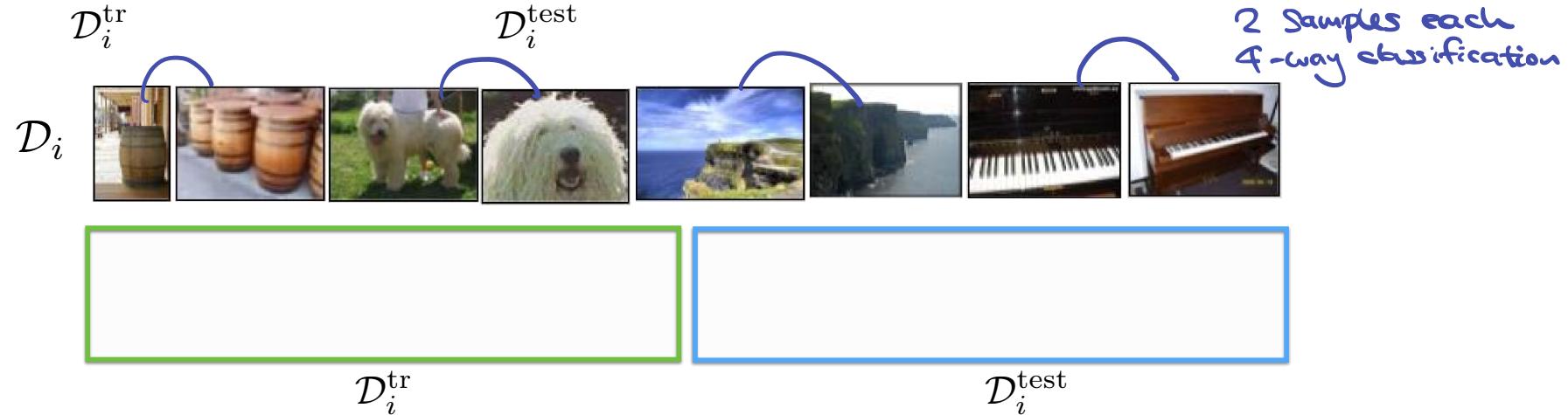
$$\min_{\theta} \sum_{\mathcal{T}_i} \mathcal{L}(f_\theta(\mathcal{D}_i^{\text{tr}}), \mathcal{D}_i^{\text{ts}})$$

Black-Box Adaptation

Key idea: Train a neural network to represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$.

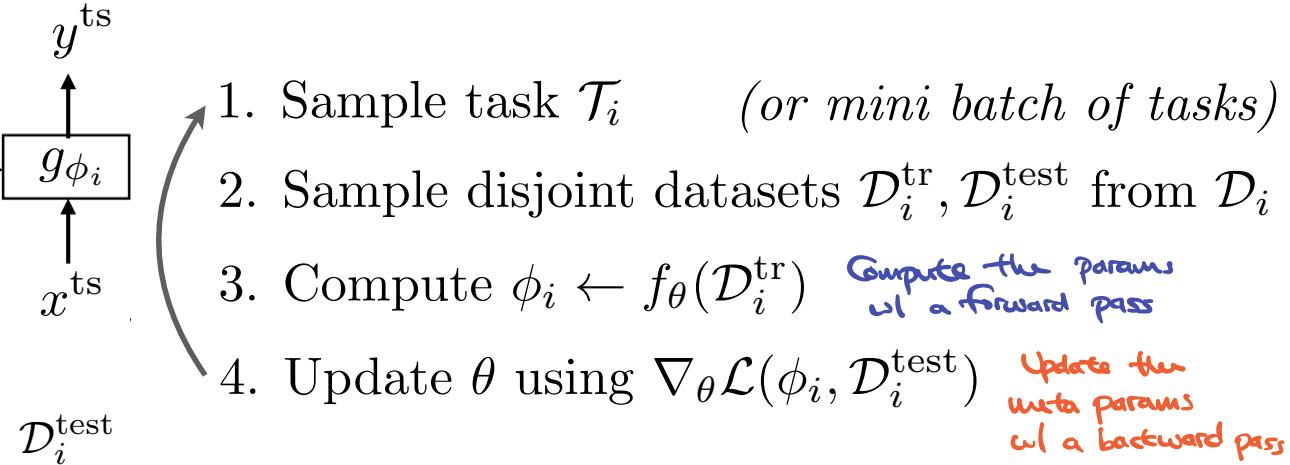
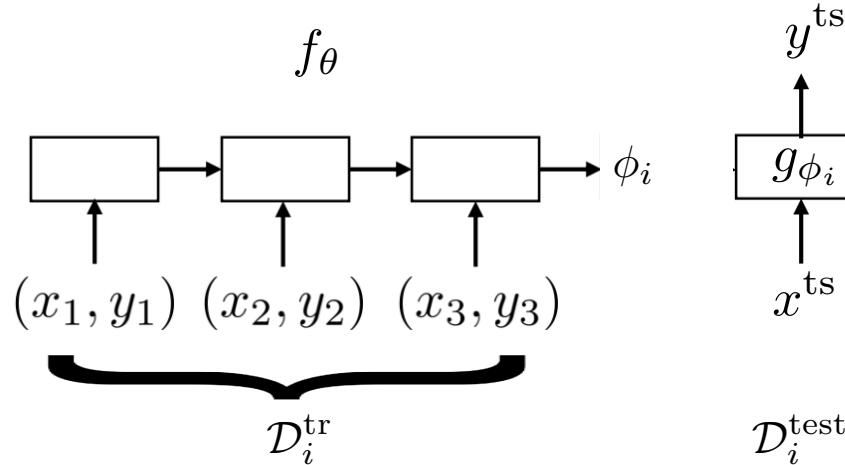


1. Sample task \mathcal{T}_i (or mini batch of tasks)
2. Sample disjoint datasets $\mathcal{D}_i^{\text{tr}}, \mathcal{D}_i^{\text{test}}$ from \mathcal{D}_i

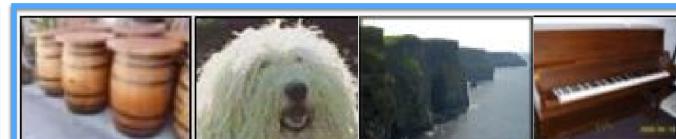


Black-Box Adaptation

Key idea: Train a neural network to represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$.



$\mathcal{D}_i^{\text{tr}}$



$\mathcal{D}_i^{\text{test}}$

Black-Box Adaptation

Key idea: Train a neural network to represent

Q: How do we learn ϕ_g for ϕ_i ?
 A: Think of ϕ_g as a part of the meta-params.
 Optimize ϕ_g alongside all of the rest of the params of ϕ .
 ↳ How we make it possible to have h_i be only a low dim. vector.
 ϕ_g may also have some params that are shared w/ the other
 part of the RNN → ϕ_g same image encoder for training & testing samples

Challenge

Outputting all neural net parameters does not seem scalable?

Idea: Do not need to output **all** parameters of neural net, only sufficient statistics

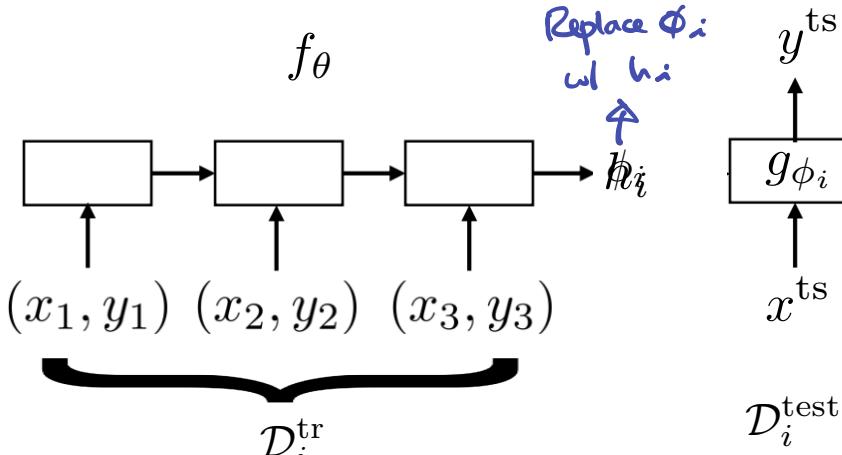
(Santoro et al. MANN, Mishra et al. SNAIL)

low-dimensional vector h_i

represents contextual task information

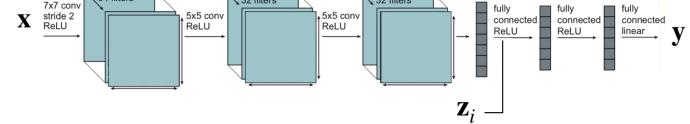
$$\phi_i = \{h_i, \theta_g\}$$

↳ Can think of this as a task descriptor in multi-task learning



Q: How do you know if you have enough tasks in D_i^{test} ?
 A: Look at the variance of accuracy across tasks.
 Too high? → Not enough test tasks.

recall:

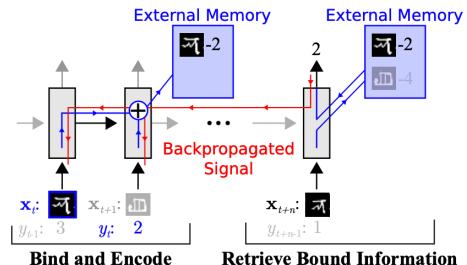


general form: $y^{ts} = f_\theta(D_i^{\text{tr}}, x^{ts})$

Black-Box Adaptation Architectures

LSTMs or Neural turing machine (NTM)

→ relic of the past;
don't use.



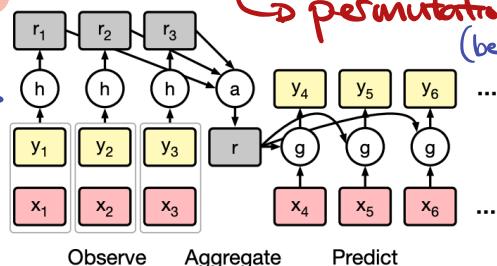
Meta-Learning with Memory-Augmented Neural Networks

Santoro, Bartunov, Botvinick, Wierstra, Lillicrap. ICML '16

Q) There are architectures called "deep sets"
⇒ can represent any permutation-invariant functions

Feedforward + average

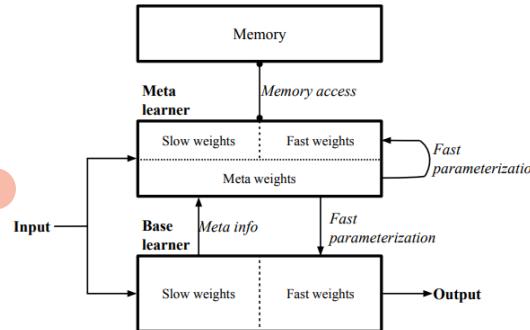
→ permutation-invariant.
(better than LSTMs)



Conditional Neural Processes. Garnelo, Rosenbaum, Maddison, Ramalho, Saxton, Shanahan, Teh, Rezende, Eslami. ICML '18

Question: Why might feedforward+average be better than a recurrent model?

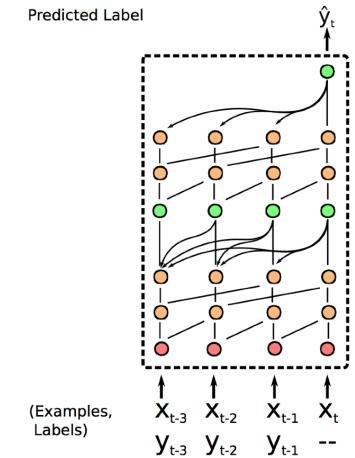
Other external
memory mechanisms



Meta Networks

Munkhdalai, Yu. ICML '17

Convolutions & attention



A Simple Neural Attentive Meta-Learner

Mishra, Rohaninejad, Chen, Abbeel. ICLR '18

Method
SNAIL, Ours

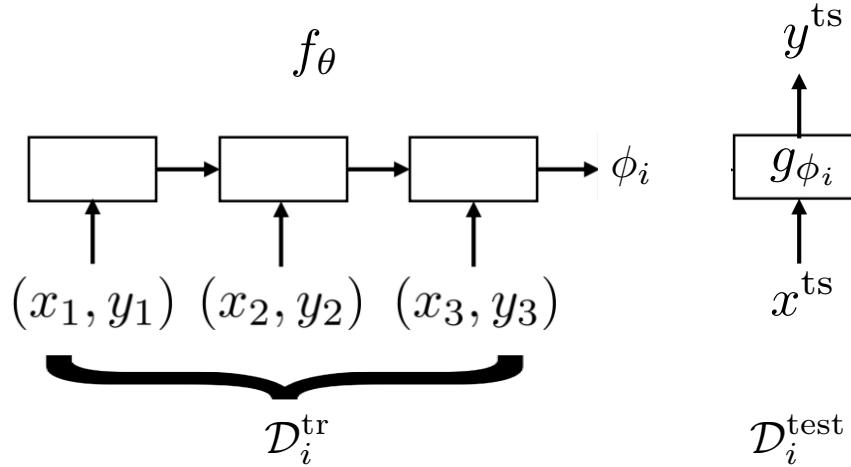
HW 1:

- implement data processing
- implement simple black-box meta-learner
- train few-shot Omniglot classifier

→ splitting up tasks, etc.

Black-Box Adaptation

Key idea: Train a neural network to represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$.



- + expressive *Can represent any function if they're large enough*
- + easy to combine with **variety of learning problems** (e.g. SL, RL) *Supervised learning, reinforcement learning*
- complex model w/ complex task: *v. large models challenging optimization problem*
- often data-inefficient *\because it's a hard optimization problem.*

How else can we represent $\phi_i = f_\theta(\mathcal{D}_i^{\text{tr}})$?

Next time (Monday): What if we treat it as an **optimization** procedure?

Embed optimization into f , rather than having large RNNs.

Plan for Today

Meta-Learning

- Problem formulation
- General recipe of meta-learning algorithms
- Black-box adaptation approaches
- **Case study of GPT-3 (time-permitting)**

eg) { GPT-3
Chinchilla
Gopher
PaLM } → canonical example.

Case Study: GPT-3

Language Models are Few-Shot Learners

Tom B. Brown*

Benjamin Mann*

Nick Ryder*

Melanie Subbiah*

Jared Kaplan[†]

Prafulla Dhariwal

Arvind Neelakantan

Pranav Shyam

Girish Sastry

Amanda Askell

Sandhini Agarwal

Ariel Herbert-Voss

Gretchen Krueger

Tom Henighan

Rewon Child

Aditya Ramesh

Daniel M. Ziegler

Jeffrey Wu

Clemens Winter

Christopher Hesse

Mark Chen

Eric Sigler

Mateusz Litwin

Scott Gray

Benjamin Chess

Jack Clark

Christopher Berner

Sam McCandlish

Alec Radford

Ilya Sutskever

Dario Amodei

OpenAI

May 2020

“emergent” few-shot learning

Wasn't explicitly set up
to do few-shot learning

What is GPT-3?

a language model

black-box meta-learner trained on language generation tasks

$\mathcal{D}_i^{\text{tr}}$: sequence of characters
trained to be conditioned on

$\mathcal{D}_i^{\text{ts}}$: the following sequence of characters
Trained to generate

[meta-training] dataset: crawled data from the internet, English-language Wikipedia, two books corpora

architecture: giant "Transformer" network 175 billion parameters, 96 layers, 3.2M batch size



What do different tasks correspond to?

↳ where "emergent" few-shot learning comes into play.

spelling correction

simple math problems

translating between languages

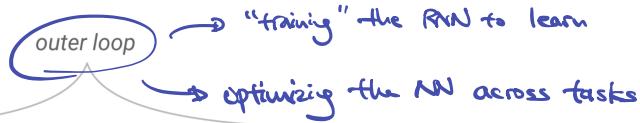
a variety of other tasks

↳ seen like v. different tasks

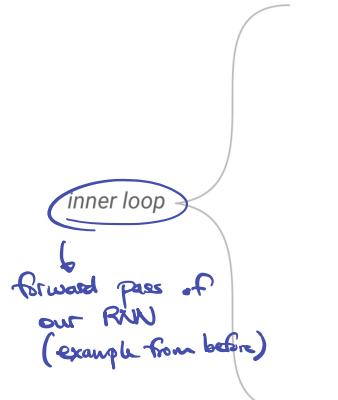
How can those tasks all be solved by a single architecture?

How can those tasks all be solved by a single architecture? Put them all in the form of text!

Why is that a good idea? Very easy to get a lot of meta-training data.



Learning via SGD during unsupervised pre-training



1	$5 + 8 = 13$
2	$7 + 2 = 9$
3	$1 + 0 = 1$
4	$3 + 4 = 7$
5	$5 + 9 = 14$
6	$9 + 8 = 17$

↑
sequence #1

simple math problems

1	gaot => goat
2	sakne => snake
3	brid => bird
4	fsih => fish
5	dcuk => duck
6	cmihp => chimp

↑
sequence #2

spelling correction

1	thanks => merci
2	hello => bonjour
3	mint => menthe
4	wall => mur
5	otter => loutre
6	bread => pain

↑
sequence #3

translating between languages

Some Results

One-shot learning from dictionary definitions:

Few-shot language editing:

3 examples of
Good-Bad English

To "screeg" something is to swing a sword at it. An example of a sentence that uses the word screeg is:

We screeghed at each other for several minutes and then we went outside and ate ice cream.

Poor English input: I eated the purple berries.

Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it.

Good English output: Thank you for choosing me as your designer. I appreciate it.

Poor English input: The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

Good English output: I'd be more than happy to work with you on another project.

Poor English input: Please provide me with a short brief of the design you're looking for and that'd be nice if you could share some examples or project you did before.

Good English output: Please provide me with a brief description of the design you're looking for and that would be nice if you could share some examples or projects you have done before.

Title: United Methodists Agree to Historic Split

Subtitle: Those who oppose gay marriage will form their own denomination

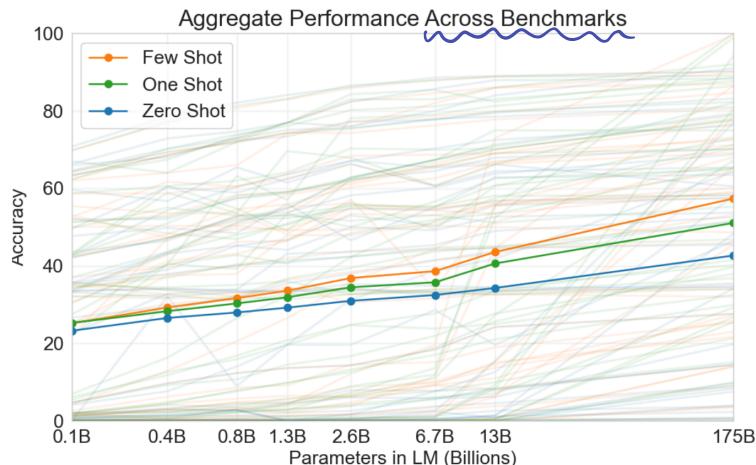
Article: After two days of intense debate, the United Methodist Church has agreed to a historic split - one that is expected to end in the creation of a new denomination, one that will be "theologically and socially conservative," according to The Washington Post. The majority of delegates attending the church's annual General Conference in May voted to strengthen a ban on the ordination of LGBTQ clergy and to write new rules that will "discipline" clergy who officiate at same-sex weddings. But those who opposed these measures have a new plan: They say they will form a separate denomination by 2020 calling their church the Christian Methodist

Non-few-shot learning tasks: → Write an article given the title + subtitle .

General Notes & Takeaways

The results are extremely impressive.

The model is far from perfect.



The model fails in unintuitive ways.

Q: How many eyes does a giraffe have?
A: A giraffe has two eyes.

Q: How many eyes does my foot have?
A: Your foot has two eyes.

Q: How many eyes does a spider have?
A: A spider has eight eyes.

Q: How many eyes does the sun have?
A: The sun has one eye.

Source: <https://lacker.io/ai/2020/07/06/giving-gpt-3-a-turing-test.html>

The choice of \mathcal{D}_i^{tr} at test time is important. ("prompting")

Source: <https://github.com/shreyashankar/gpt3-sandbox/blob/master/docs/priming.md>

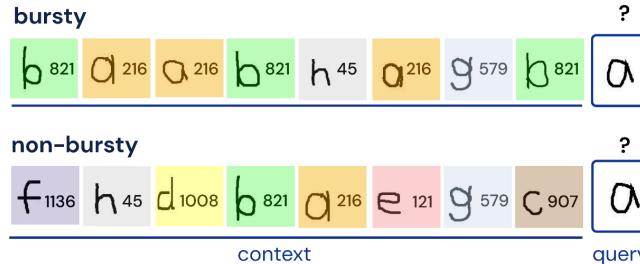
“Hallucinations”

What is needed for few-shot learning to emerge?

An active research topic!

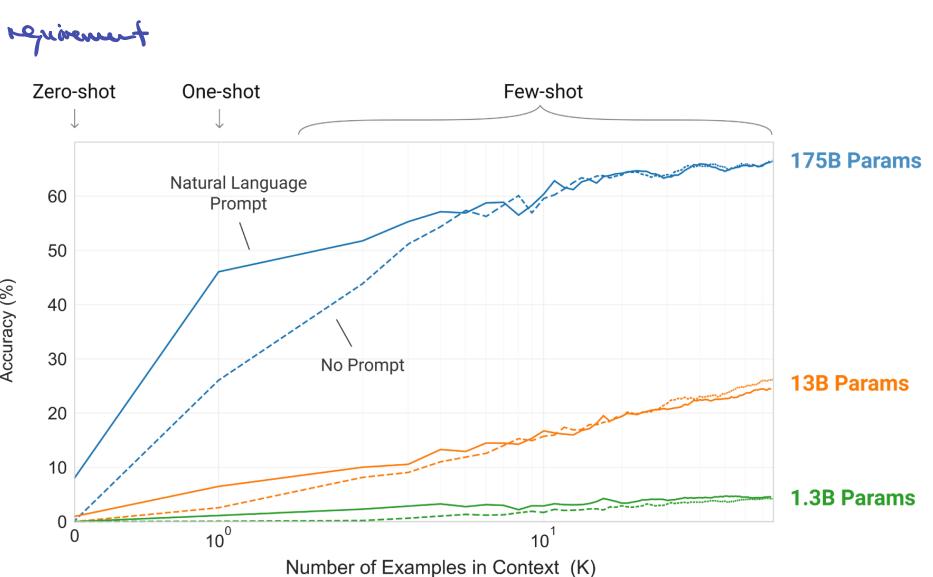
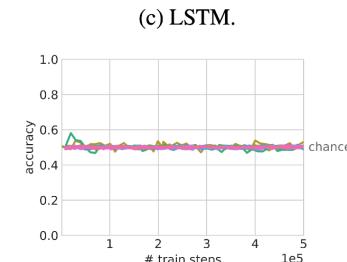
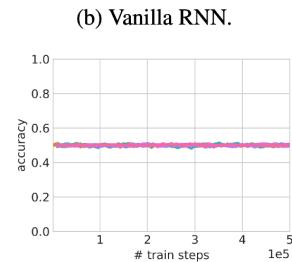
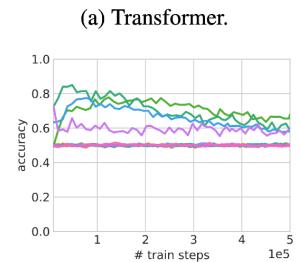
Data:

- temporal correlation
- dynamic meaning of words
(e.g.) "wicked", depending on context, can mean two different things



Model:

- large capacity models
- transformers > RNNs
- large models > small models



Plan for Today

Meta-Learning

- Problem formulation
- General recipe of meta-learning algorithms
- Black-box adaptation approaches
- Case study of GPT-3 (time-permitting)

}

Topic of Homework 1!

Goals for by the end of lecture:

- Training set-up for few-shot meta-learning algorithms
- How to implement black-box meta-learning techniques

Reminders

Project group form due **Monday, October 10**

Homework 1 due **Wednesday October 12**

Next time: Optimization-based meta-learning