DeepMind

## Data Distributional Properties Drive Emergent In-Context Learning in Transformers

Stephanie Chan, Adam Santoro, Andrew Lampinen, Jane Wang, Aaditya Singh, Pierre Richemond, Jay McClelland, Felix Hill

## Two distinct modes of learning in neural networks

### In-weights learning

- gradient-based
- slow: needs many examples
- standard supervised learning





### In-context learning

- no gradient updates
- rapid: from a few examples
- few-shot / one-shot learning



Where is another?

ಅ	с <mark>а</mark>	ബ	പ	$\overline{C}$
ભ	ಖ	Ч	ಬ	ಗ್
ন্দ্ৰ	ō	ខា	ല്	ជ
ಸ	സ	ಲ	മ	ಛ

#### Normally, only happens if we explicitly meta-train for it:



E.g.: Santoro et al, 2016; Vinyals et al, 2016; Wang et al, 2016



## But in-context (few-shot) learning can also *emerge* ...when trained on a very different objective!







### How do large transformer models achieve emergent in-context learning?

Hypothesis: Maybe it's because the distributions of naturalistic data have special properties



Natural data is long-tailed

Natural data is bursty



document in corpus

Maybe training on naturalistic data is like an interpolation between supervised and few-shot meta-training...



### **Standard Supervised**

- items recur and are uniform
- label mappings are fixed



## Naturalistic data (e.g. language)

- words do recur
- word meanings are somewhat fixed

but also:

- rare words do not recur often
- some rare words are bursty
- many-to-many relationships

Few-shot meta-training

explicitly train for few-shot learning

- items differ on every episode
- label mappings are only fixed within episodes



few-shot learning emerges, even without explicit training



## **Our Project**

*Hypothesis*: Certain non-uniformities in data distributions can lead to emergent few-shot (in-context) learning, and this is a general phenomenon.

*Experiments*: Modify a standard few-shot learning image dataset (Omniglot), to control these distributional properties and measure their effects on few-shot learning.

### Implications:

• understanding how we might design or collect datasets to achieve in-context learning in domains outside of language

## **General structure of the experiments**





## **Training data**

• labels are fixed across all of training

## Example "bursty" sequence

### Two ways to solve:

- 1. In-weights memorization
- 2. In-context learning



query





## **Evaluation data**

#### Example evaluation sequence for <u>in-context learning</u>

• Two holdout classes, randomly assigned to labels [0, 1]

Compositional binding of images and labels

#### Example evaluation sequence for <u>in-weights memorization</u>

• The query class was seen in training, and does not appear in the context. ? • 1136 h 45 d 1008 b 821 i 436 e 121 g 579 C 907 0 query DeepMind

# What kinds of training data promote in-context learning?



## **Importance of <b>burstiness** in the data



<u>Eval</u>



 More burstiness leads to better in-context learning In-weights learning on trained classes.



 In-context learning trades off against weights-based learning



## Transformers succeed at the Omniglot challenge: Importance of <u>number of classes</u> in the data







 More training classes leads to better in-context learning

#### (b) In-weights learning on trained classes.



 Again, in-context learning trades off against weights-based learning



## Other natural data-inspired distributional properties

Dynamic meaning:

### Multiplicity of item-label mappings



Figure 4: Dynamic meanings improve in-context learning. Increasing the number of labels per class ('label multiplicity') increases in-context learning.

### Within-class variation







# Can in-weights memorization and in-context learning co-exist in the same model?



### We can achieve both kinds of learning when we train on *skewed* distributions.



There is a sweet spot at Zipf exponent = 1, where we attain both few-shot learning and in-weights memorization

→ Intriguingly, Zipf exponent 1 corresponds approximately to the skew in natural languages





# But architecture does matter too...



## **Transformers vs RNNs**

In-context learning on holdout classes.

(a) Transformer.



(c) LSTM.

- each line is a hyperparameter setting
- matched on: # params, # layers, hidden size, training data
  - → Recurrent models never achieve few-shot learning, with the same training data
  - → But even though architecture matters, it's not enough we need the right data, too

(b) Vanilla RNN.



## Conclusions



## **Implications for compositionality**

We study in-context few-shot learning, which can be construed as a a narrow instantiation of compositionality.

Our findings on the drivers of this emergent behavior:

- Large-scale data and models are not necessary
- Certain distributions of training data promote it

   these distributional features are present in
   natural data like language
- Architecture matters too
  - Transformers > RNNs





## **Implications for compositionality**

We study in-context few-shot learning, which can be construed as a a narrow instantiation of compositionality.

Our findings on the drivers of this emergent behavior:

- Large-scale data and models are not necessary
- Certain distributions of training data promote it
   these distributional features are present in
  - natural data like language
- Architecture matters too
  - Transformers > RNNs

RNNs with memory augmentation \*could\* perform well on SCAN (Lake et al 2019)

> memory-augmented NNs help meta-learning because they have both long-term and short-term storage (that is stable + element-wise addressable) (Santoro et al 2013)

Transformers can perform well on certain kinds of composition

- transformers have the desired properties as well
- transformers may perform
   compositional operations
   (Elhage et al 2021; Olsson et al 2022)

## **Implications for compositionality**

We study in-context few-shot learning, which can be construed as a a narrow instantiation of compositionality.

Our findings on the drivers of this emergent behavior:

- Large-scale data and models are not necessary
- Certain distributions of training data promote it

   these distributional features are present in
   natural data like language
- Architecture matters too
  - Transformers > RNNs

We need both! Non-uniformity is

important

