**References**

- [Rethinking the Role of Demonstrations: What Makes In-Context Learning Work? (https://arxiv.org/pdf/2202.12837.pdf)](https://arxiv.org/pdf/2202.12837.pdf)

# What makes In-context Learning work ?

- blog (http://ai.stanford.edu/blog/understanding-incontext/)
- paper (https://arxiv.org/pdf/2202.12837.pdf)
    - more empirical
    - various models
        - MetaICL: trained with InContextLearning objective
        - 2 methods: Direct vs Channel ???
    - gold-label vs random (uniform sampling) label: ground-truth not necessary
        - gold improves over zero shot
        - random: small decrease vs gold
            - **very** small for MetaICL
        - sampling for true label distribution: smaller decrease

# How does In Context Learning work ?

*In-context Learning* describes a means of using a fixed LLM to solve a task

- by supplying some number $k$ of *exemplars* (or *demonstrations*) of the new task
$$\langle \backslash\mathbf{x}^{(1)}, \backslash\mathbf{y}^{(1)} \rangle, \ldots, \langle \backslash\mathbf{x}^{(k)}, \backslash\mathbf{y}^{(k)} \rangle$$
- as a *context* C
    - that describes the new task's relationship between input $\backslash\mathbf{x}^{\backslash\mathrm{ip}}$ and output $\backslash\mathbf{y}^{\backslash\mathrm{ip}} *$
- and presenting a prompt $\backslash\mathbf{x}$ to the model
- expecting the model to produce a $\backslash\hat{\mathbf{y}}$
- that is the correct "response" to the new task on input $\backslash\mathbf{x}$

So the prompt (context plus example's feature $\mathbf{x}$) might look like

$$\text{Input:} \quad \mathbf{x}^{(1)}$$
$$\text{Output:} \quad \mathbf{y}^{(1)}$$

$$\vdots$$

$$\text{Input:} \quad \mathbf{x}^{(k)}$$
$$\text{Output:} \quad \mathbf{y}^{(k)}$$

$$\text{Input:} \quad \mathbf{x}$$
$$\text{Output:}$$

and expect the continuation to be the prediction $\hat{\mathbf{y}}$ corresponding to test input $\mathbf{x}$

For example:



**Demonstrations**

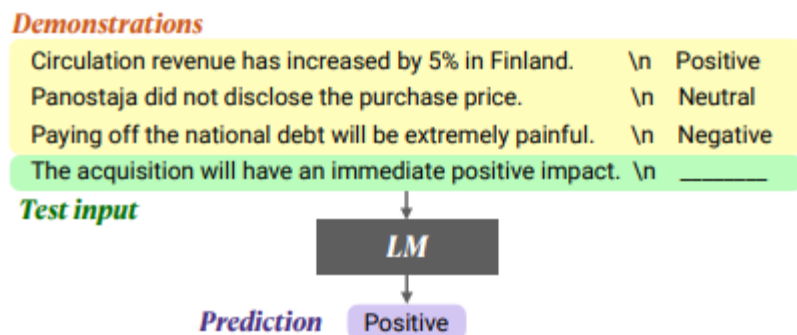| Circulation revenue has increased by 5% in Finland. | \n | Positive |
| Panostaja did not disclose the purchase price. | \n | Neutral |
| Paying off the national debt will be extremely painful. | \n | Negative |
| The acquisition will have an immediate positive impact. | \n | _____ |

*Test input*

LM

*Prediction* Positive

Figure 2: An overview of in-context learning. The demonstrations consist of $k$ input-label pairs from the training data ($k = 3$ in the figure).

Atttribution: https://arxiv.org/pdf/2202.12837.pdf#page=2 (https://arxiv.org/pdf/2202.12837.pdf#page=2)

In-Context Learning uses a pre-trained LLM and the trick of using the Universal Text API

- to turn the new task
- into a text-continuation ("predict the next") task

It appears to be a way

- of extending a LM
- *without* further training
    - as opposed to Fine-Tuning
- since
    - the exemplars are given at *test* time
    - no parameter updates to the LLM occur

But why should this work ?

More interestingly

- what is a good theory
- and how can we test it

We will present a [paper (https://arxiv.org/pdf/2202.12837.pdf)](https://arxiv.org/pdf/2202.12837.pdf) that attempts to present some insights into the process.

# Testing some theories

In order to test a theory

- various aspects of the exemplars are proposed as variables
- one variable at a time is perturbed
- the effect of the perturbations is measured across a range of benchmarks

- **and compare to measurements before the perturbation**

The results are summarized in the following diagram

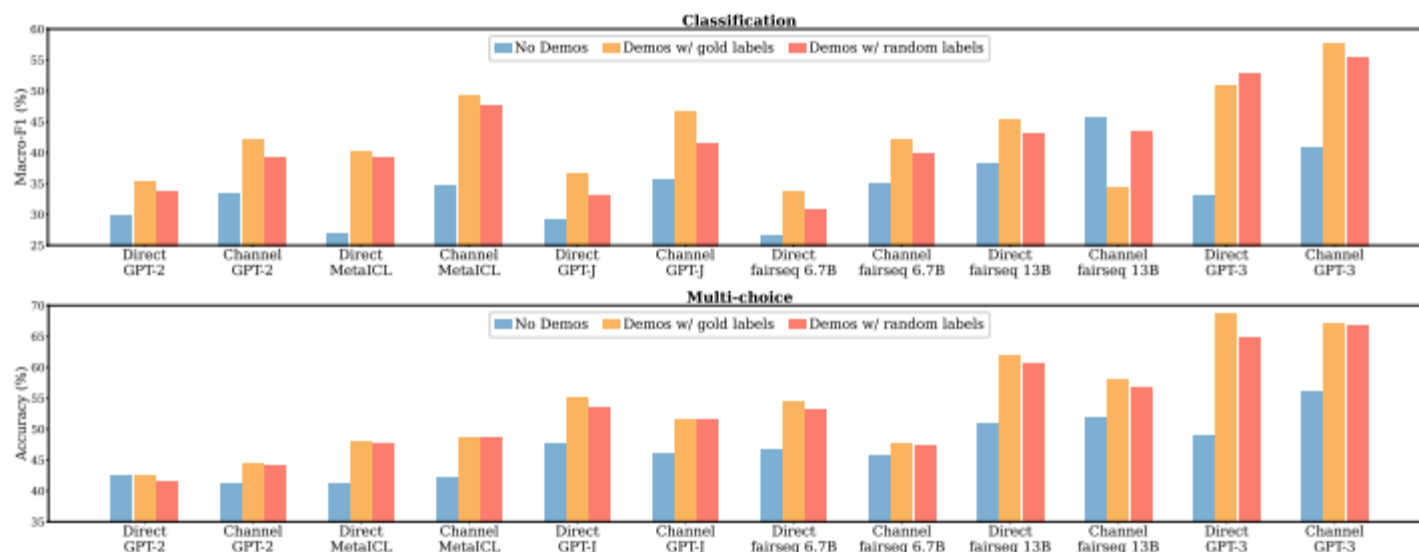- that we will subsequently refer to for each experiment



Figure 3: Results when using no-demonstrations, demonstrations with gold labels, and demonstrations with random labels in classification (top) and multi-choice tasks (bottom). The first eight models are evaluated on 16 classification and 10 multi-choice datasets, and the last four models are evaluated on 3 classification and 3 multi-choice datasets. See Figure 11 for numbers comparable across all models. **Model performance with random labels is very close to performance with gold labels** (more discussion in Section 4.1).

This chart shows the result of perturbations

- run across a variety of models
    - of sizes ranging from 774M to 175B parameters
- each experiment is averaged across multiple benchmarks

This summarizes a variety of experiments (to be described) involving the exemplars.

The 3 different colored bars represent variations on the exemplars

| | |
|---|---|
| blue | no exemplars |
| gold | exemplars $k = 16$ |
| red | exemplars with perturbed labels $k = 16$ |

# Zero shot verus $k \geq 1$ shot

The first experiment measures the effect of the presence/absence of exemplars

- blue vs. orange

**Conclusion**

$k \geq 1$ exemplars *improves performance* relative to zero-shot.

# Parts of the Context

The next set of experiments varies parts of the context (exemplars and prompt).

Given exemplars

$$\langle \backslash\mathbf{x}^{(1)}, \backslash\mathbf{y}^{(1)} \rangle, \ldots, \langle \backslash\mathbf{x}^{(k)}, \backslash\mathbf{y}^{(k)} \rangle$$

the authors posit some salient characteristics

- the *input distribution $I$* from which the exemplar *features* are drawn $\backslash\mathbf{x}^{(1)}, \ldots, \backslash\mathbf{x}^{(k)}$
- the distribution $L$ of the exemplar *labels* $\backslash\mathbf{y}^{(1)}, \ldots, \backslash\mathbf{y}^{(k)}$
- the feature/label mapping relationship $M$
    - i.e., the pair of $\backslash\mathbf{x}^{\backslash\mathrm{ip}}$ and $\backslash\mathbf{y}^{\backslash\mathrm{ip}}$, for $1 \leq i \leq k$
- formatting
    - the encoding of the exemplars and prompt into $\backslash\mathbf{\dot{x}}$

# Feature/label mapping relationship

- gold vs. red

Let $\mathcal{C}$ denote the set from which exemplar labels are drawn.

In this experiment, replace

- correct label $\mathbf{y}^{\backslash ip}$ for exemplar $i$
- with label $\tilde{\mathbf{y}}^{\backslash ip}$ drawn at random (uniformly) from $\mathcal{C}$.

That is, we preserve $I$ and $L$, but break $M$.

**Conclusions**

- Correct ("gold") labels improve performance over random labels
    - but not as much as expected, perhaps
- But Random labels *improves performance* over *no exemplars* !
    - "Ground truth" matters surprisingly little !

The fact that an *incorrect* $M$ improves performance relative to no exemplars is surprising.

This suggests

- that the exemplars are used to infer the *task to be performed*

- once the task has been identified

    - the exemplar mis-labeling is ignored
- the model is able to perform the task as it was *trained* during training

See the [Signifier theory in the module](Prompt_Engineering_Suggestions.ipynb#Signifier:-direct-specification)
[(Prompt_Engineering_Suggestions.ipynb#Signifier:-direct-specification)](Prompt_Engineering_Suggestions.ipynb#Signifier:-direct-specification)

# Input distribution

This experiment measures a *shift in the distribution* of the exemplar features $\mathbf{x}^{\backslash ip}$

In this experiment

- each exemplar input $\mathbf{x}^{\backslash ip}$ is replaced by
- a random $\mathbf{x}^{\backslash ip}_{\mathrm{rand}}$ drawn from a text corpus *other than the one used for Training*

We note that this experiment *also* breaks the feature/label relationship $M$

- we preserve the original labels $\backslash\mathbf{y}^{\backslash\text{ip}}$ for exemplar $i$
- which is not necessarily related to $\backslash\mathbf{x}^{\backslash\text{ip}}_{\text{rand}}$

We can contrast the results of this experiment the effect of breaking $M$ alone.
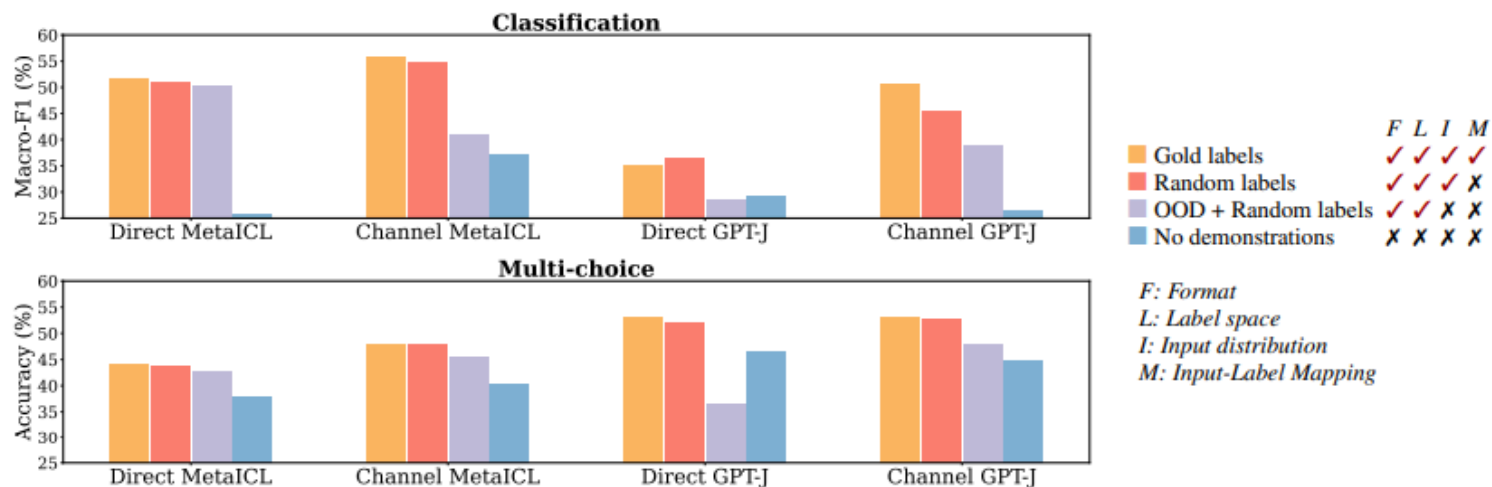


Figure 8: Impact of the distribution of the inputs. Evaluated in classification (top) and multi-choice (bottom). The impact of the distribution of the input text can be measured by comparing ■ and ■. The gap is substantial, with an exception in Direct MetaICL (discussion in Section 5.1).

In the above diagram, compare

- the lavender (third bar from left): perturbed $I$ and $M$
- the red bar (second bar from left): perturbed $M$ alone

**Conclusions**

The $M$ relationship is broken in both cases. But

- preserving the original distribution $I$ of exemplar features
- improves performance relative to changing the distribution

Why might this be ?

The suggestion is that the model was trained with the LLM objective ("predict the next")

- from a training distribution
- and $\mathbf{x}_{\mathrm{rand}}$ is from a *different* distribution
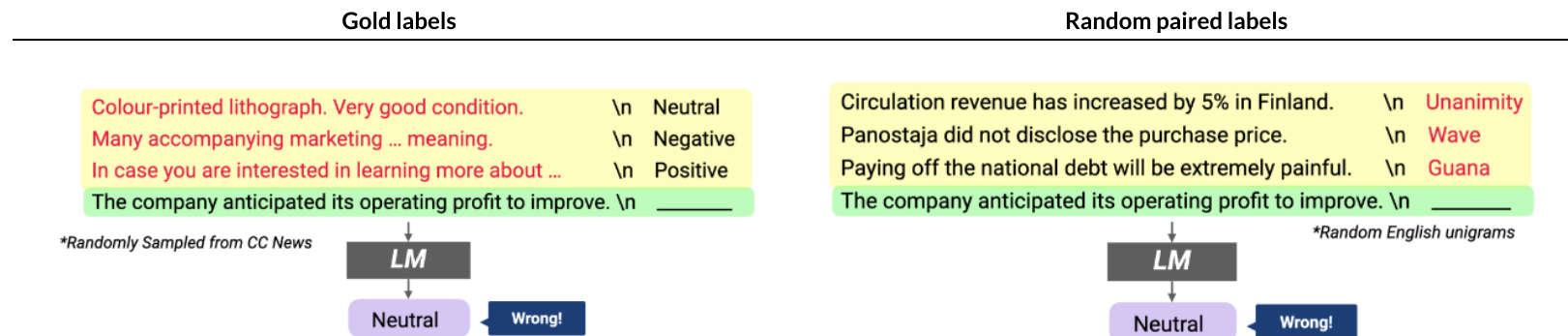- so the model struggles on non-training input

# Output distribution

This experiment measures a *shift in the distribution* of the exemplar labels $\mathbf{y}^{\backslash ip}$
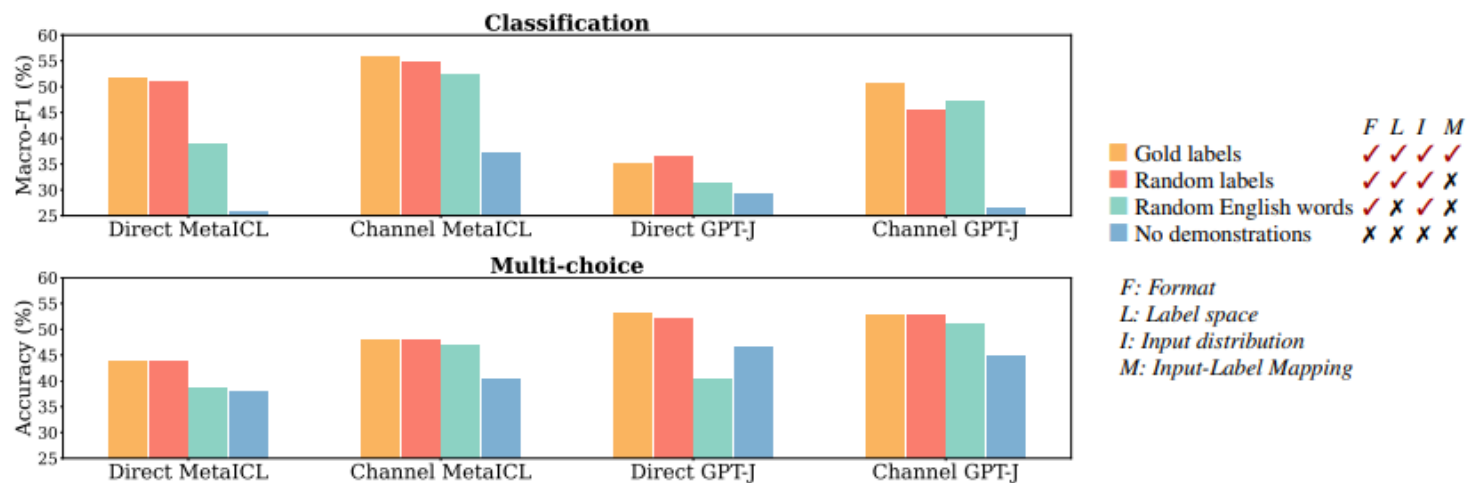
In this experiment

- each exemplar label $\mathbf{y}^{\backslash ip}$ (from $\mathcal{C}$) is replaced by
- $\tilde{\mathbf{y}}^{\backslash ip}$ a *random* English word from $\mathcal{C}_{\text{rand}}$

Note that we *also break $M$*

- e.g., $\mathbf{y}^{\backslash ip} = \mathbf{y}^{(i')}$ does not imply $\tilde{\mathbf{y}}^{\backslash ip} = \tilde{\mathbf{y}}^{(i)}$

| Gold labels | | Random paired labels | |
|---|---|---|---|



Attribution: http://ai.stanford.edu/blog/understanding-incontext/
(http://ai.stanford.edu/blog/understanding-incontext/)

**Classification**

**Multi-choice**

Gold labels — F ✓ L ✓ I ✓ M ✓
Random labels — F ✓ L ✓ I ✓ M ✗
Random English words — F ✓ L ✗ I ✓ M ✗
No demonstrations — F ✗ L ✗ I ✗ M ✗

F: Format
L: Label space
I: Input distribution
M: Input-Label Mapping

In the above diagram, compare

- the red bar: random labels
- the turquoise bar (third from left): random words

The difference: Although we break $M$ in both cases

- in the "random labels" case: the labels are chosen from the correct output distribution $\mathcal{C}$
- in the "random words" case: the labels come from a distribution other than $\mathcal{C}$

**Conclusions**

The $M$ relationship is broken in both cases. But

- preserving the original distribution $L$ of exemplar labels
- improves performance relative to changing the distribution of labels

# Formatting

In this experiment

- *format* is defined as the *pairing* of a feature and label within an exemplar
- not necessarily a *correct pairing*: mapping $M$ not necessarily correct

One experiment is run

- with *only* exemplar features (and no exemplar labels): $\backslash \mathbf{x}^{(1)}, \ldots, \backslash \mathbf{x}^{(k)}$
- natural comparison is with experiment of correct format
    - $\backslash \mathbf{x}^{\backslash \mathrm{ip}}$
    - paired with random English words (from $\mathcal{C}_{\mathrm{rand}}$) as labels

A second experiment is run

- with *only* exemplar labels (and no exemplar features): $\backslash\mathbf{y}^{(1)}, \ldots, \backslash\mathbf{y}^{(k)}$
- natural comparison is with experiment of correct format
  - a random $\backslash\mathbf{x}_{\text{rand}}^{\backslash\text{ip}}$ drawn from a text corpus
  - paired with $\backslash\mathbf{y}^{\backslash\text{ip}}$

Both comparison experiments

- preserve the format: feature/exemplar pairs
- without preserving $M$
- or the distribution $I$ in the first case, and $L$ in the second case
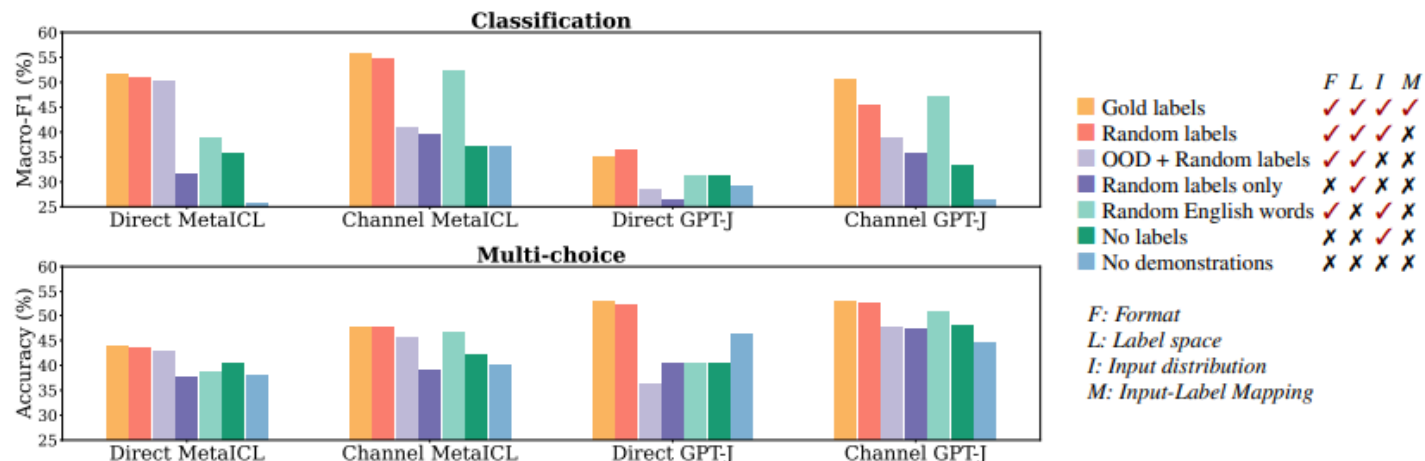
Figure 10: Impact of the format, i.e., the use of the input-label pairs. Evaluated in classification (top) and multi-choice (bottom). Variants of demonstrations without keeping the format (■ and ■) are overall not better than no demonstrations (■). Keeping the format is especially significant when it is possible to achieve substantial gains with the label space but without the inputs (■ vs. ■ in Direct MetaICL), or with the input distribution but without the labels (■ vs. ■ in Channel MetaICL and Channel GPT-J). More discussion in Section 5.3.

**Conclusions**

- Not keeping the format has performance *on par* with **no demonstrations** at all
- Keeping the format retains most of the benefit achievable with either
  - correct $I$ (but incorrect $M$)
  - or correct $L$ (but incorrect $M$)

The suggestion is that correct format an important feature

- to enable the LLM to recognize the task from exemplars

# Exemplars that differ from the LLM

The on-line articles makes another interesting observation.

Observe the encoding of the exemplars

$$\langle \backslash \mathbf{x}^{(1)}, \backslash \mathbf{y}^{(1)} \rangle, \ldots, \langle \backslash \mathbf{x}^{(k)}, \backslash \mathbf{y}^{(k)} \rangle$$

and example feature $\backslash \mathbf{x}$

into the prompt

$$C, \backslash \mathbf{x}$$

The distribution of prompt $C, \backslash\mathbf{x}$

- is probably *much different* than
- the distribution (Internet text documents) on which the LLM was trained

in several ways

- syntax
    - structured list of $\backslash\mathbf{x}^{\backslash\mathrm{ip}}$, $\backslash\mathbf{y}^{\backslash\mathrm{ip}}$ pairs
    - versus natural sentences
- coherence
    - consecutive exemplars may be from different *topics*
        - that illustrate the *concept* of the new task
    - consecutive natural sentences may share the same topic

The article posits that

- these encoding anomalies
- are low-frequency noise
- that the LLM is able to ignore
- providing there is more "signal" in the exemplars

# A theory of In-Context Learning

A more [theoretical paper (https://arxiv.org/pdf/2111.02080.pdf)](https://arxiv.org/pdf/2111.02080.pdf) and accompanying [online article (http://ai.stanford.edu/blog/understanding-incontext/)](http://ai.stanford.edu/blog/understanding-incontext/)
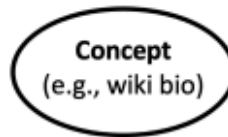
- combine these experimental insights
- into a theory
- and mathematical model of the theory
- that is consistent with the experimental results

The authors posit

- during training, the LLM learns "concepts", for example
    - abstract ideas
        - question answering
        - sentiment
    - plans
        - how to solve a multi-step task: travel directions

# Concepts

## 1. Pretraining documents are conditioned on a latent concept (e.g., biographical text)

Concept (e.g., wiki bio) → Albert Einstein was a German theoretical physicist, widely acknowledged to be one of the greatest physicists of all time. Einstein is best known for developing the theory of relativity, but he also ....

## 2. Create independent examples from a shared concept. If we focus on full names, wiki bios tend to relate them to nationalities.

| Input (x) | Output (y) | Delimiter |
|---|---|---|
| Albert Einstein was | German | \n |
| Mahatma Gandhi was | Indian | \n |
| Marie Curie was | ? | ...brilliant? ...Polish? |

Concept (e.g., wiki bio)

## 3. Concatenate examples into a prompt and predict next word(s). Language model (LM) implicitly infers the shared concept across examples despite the unnatural concatenation

Albert Einstein was German \n Mahatma Gandhi was Indian \n Marie Curie was → LM → Polish

Figure 1: In-context learning can emerge from modeling long-range coherence in the pretraining data. During pretraining, the language model (LM) implicitly learns to infer a latent concept (e.g., wiki bios, which typically transition between name (Albert Einstein) → nationality (German) → occupation (physicist) → ...) shared across sentences in a document. Although prompts are unnatural sequences that concatenate independent examples, in-context learning occurs if the LM can still infer the shared concept across examples to do the task (name → nationality, which is part of wiki bios).

The model's LLM "predict the next" training objective did not specify to goal of learning concepts.

But

- summarizing a large number of similar training documents (e.g., collection of biographies)
- in a parameters-constrained model
- logically suggests that concepts emerge as a way of reducing parameter usage

The authors suggest that the LLM's probability of outputting $\backslash\mathbf{y}$ given prompt $\backslash\mathbf{x}$ is formed by

$$\backslash\mathrm{pr}\backslash\mathbf{y}|\backslash\mathbf{x} = \int_{c \in \mathrm{Concepts}} \backslash\mathrm{pr}\backslash\mathbf{y}|\backslash\mathbf{x}, c\backslash\mathrm{pr}c \; d(c)$$

That is, the output

- is the sum over all concepts
- of the probability of outputting $\backslash\mathbf{y}$ given prompt $\backslash\mathbf{x}$ and concept $c$

Furthermore: the context (i.e., exemplars) of in-context learning

- helps the LLM identify the concept $c$
- to which the prompt $\backslash \mathbf{x}$ implicitly refers

The experimental results seem to suggest that the exemplars

- don't need to be fully accurate
    - the model tolerates inaccurate mappings $M$ between feature input space $I$ and label space $L$
- that correctly identifying $I$ and $L$ through the exemplars is
    - advantageous
    - but not completely necessary
- that the *format* of the exemplar
    - paired features and labels
    - is important

Under this theory

- the exemplars **are not teaching** new concepts
    - hence $M$ can be inaccurate
- but serving to help the LLM **identify** a concept learned in training

That is, the encoded exemplars in $C$

- are related to $\backslash \mathrm{pr}c$

Once the concept $c$ is identified, the output $\hat{y}$ depends on

- the distributions $I_{\text{train}}$ and $L_{\text{train}}$
  - of features/labels
- the mapping $M_{\text{train}}$
  - of features to labels

that were learned during **training**.

**NOT**

- on the $I, L, M$ of the exemplars

That is

$$\hat{y} = \pr{y | x, c} \quad \text{concept } c$$

not

$$\hat{y} = \pr{y | C, x} \quad \text{exemplar context } C$$

```python
In [2]: print("Done")
```

Done