

---

# On the Relationship between the Choice of Representation and In-context Learning

---

Anonymous Authors<sup>1</sup>

## Abstract

In-context learning (ICL) is the ability of a large language model (LLM) to learn a new task from a few demonstrations presented as part of the context. Prior work has attributed much of ICL's success to the representation of in-context demonstrations, particularly to the choice of labels in classification tasks. At the same time, evidence for ICL's learning capacity, i.e., whether additional demonstrations improve performance, has been mixed, and ICL is often thought to occur only under specific conditions. The interaction between representation and learning in ICL remains underexplored. We hypothesize that these two aspects influence ICL performance in distinct ways: the representation of demonstrations determines the baseline accuracy of ICL, while learning from additional demonstrations improves performance on top of this baseline. We test this hypothesis by developing an optimization algorithm that enumerates label sets with varying semantic relevance, and performing ICL with varying numbers of demonstrations for each label set. We observe that learning occurs regardless of label set quality, although its efficiency, measured by the slope of improvement over demonstrations, depends on both label set quality and the parameter count of the underlying language model. Despite the emergence of learning, the relative accuracy of different label sets is largely preserved throughout learning, confirming our hypothesis. Our results reveal a previously underexplored aspect of ICL: the distinct roles of representation and learning in determining ICL performance.

## 1. Introduction

LLMs are able to learn a new task from a few examples, an ability known as in-context learning (ICL) (Brown et al., 2020; Dong et al., 2024). A model is prompted with input-output pairs (demonstrations) illustrating the task and then asked to make a prediction for a novel input. The ICL paradigm is appealing as the models appear to learn something new without updating any weights, in contrast with the typical way in which a neural network learns via back-propagation. However, the performance of ICL depends heavily on properties of the given demonstrations (Perez et al., 2021), such as the the distribution of input text, the label space (Min et al., 2022), the number and order of examples (Lu et al., 2021; Liu et al., 2024; Chen et al., 2023; Bertsch et al., 2025), and the overall format of the sequence (Zhao et al., 2021). It remains unclear whether ICL truly constitutes learning, and if so, how learning interacts with elements of the prompt.

Prior work has studied learning and representation in ICL separately, not considering the interaction between the two, which may have led to incomplete conclusions. According to earlier studies, different kinds of in-context learning happen depending on the choice of how labels are represented. In particular, two types of labeling schemes have been studied extensively: gold (or semantically-meaningful) labeling and abstract (or semantically-void) labeling. Pan et al. (2023) found that with an abstract set of labels, smaller models perform similarly regardless of how many demonstrations were presented, while larger models showed increased performance with more demonstrations. This led them to conclude that the emergence of in-context learning depends on the model size. More recently, Kirsanov et al. (2025) observed that LLMs are sensitive to the representation of labels and perform better with gold labels than with abstract labels. In their study, the accuracy improved with an increasing number of demonstrations for both gold

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

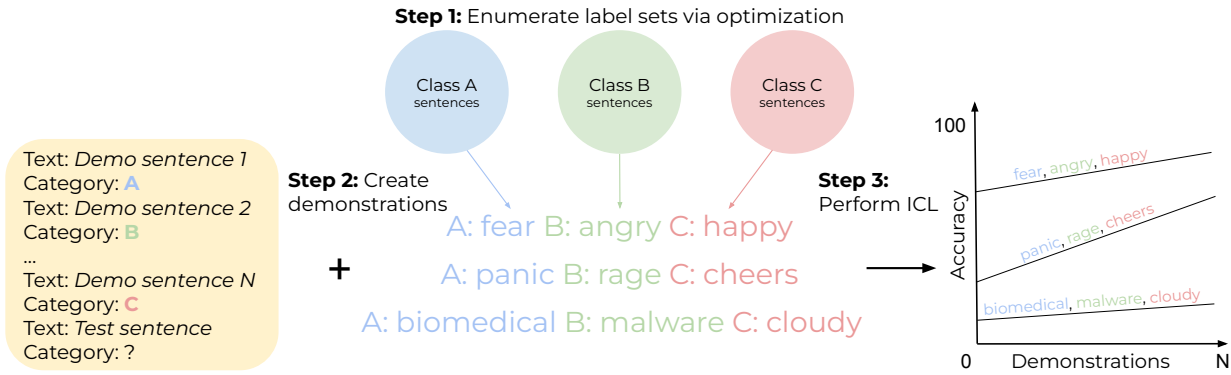


Figure 1. Method overview. Step 1: We develop an optimization algorithm to enumerate a list of possible label sets for a classification task. Step 2: We label demonstration sentences according to the label sets found. Step 3: We use these demonstrations in ICL tasks and evaluate the performance obtained with each label set on the same set of test sentences.

and abstract labels, even with a smaller model. Both Min et al. (2022) and Pan et al. (2023) observed that breaking the input-output correspondence while preserving the set of labels had a minimal effect for small models, suggesting that the representation is the sole driver of performance, rather than the demonstration pairings themselves. These findings highlight the need to investigate the interaction between learning and representation in ICL.

In this work, we propose that in classification tasks ICL performance is influenced by two separate components: representation—the choice of class names or labels, and learning—the number of examples presented in context. To quantify the role of representation, we evaluate label sets with varying degrees of semantic relevance to the task. We develop an optimization algorithm to enumerate such label sets. We then use these representations to label input sentences and to create demonstrations for ICL. We conduct experiments on sentiment and question classification tasks: 3-way and 5-way, across three model sizes from three model families. For each label set we analyze ICL performance while varying the number of demonstrations. We show an overview of our method in Figure 1.

We found that representation steers learning, although learning typically happens regardless of representation and model size. The ranking of representations in terms of accuracy is constant across different number of demonstrations, following the initial order (without any demonstrations). Moreover, the accuracy range attainable with a given representation is largely determined by the zero-shot accuracy. For most label sets, the  $N$ -shot accuracy generally increases with  $N$ , although we found that learning efficiency, that is the slope of improvement, depends on the model size. This characterization of the relationship between learning and representation in ICL suggests that it is possible to improve ICL performance by carefully choosing an appropriate label set representation for the task.

## 2. Related work

There have been a flurry of academic studies on ICL that have revealed its properties and characterized ICL as a new class of learning, since Brown et al. (2020) demonstrated the (surprising) effectiveness of ICL with a large-scale language model. In this section, we list up some of these studies that have shed light on ICL over the past few years.

**Content effects.** Recent studies suggest that LLMs are not fully-abstract reasoners, that is, they do not always learn a function which they can apply to an arbitrary input (Lampinen et al., 2024). Instead, these models show content effects similar to those of humans who reason more accurately about familiar or grounded situations, compared to unfamiliar or abstract ones. McCoy et al. (2024) found that LLM accuracy is influenced by the probability of the task to be performed, the probability of the target output, and the probability of the provided input. The bias towards outputs that have a high prior probability occurs in ICL as well. LLMs do not always identify a unique input-output mapping across the demonstrations, in order to apply it to the test input. They rely instead on the combination of their prior knowledge and presented demonstrations. There are several factors influencing ICL, such as the order (Lu et al., 2021) and number of demonstrations (Chen et al., 2023), input and output distributions, and the overall format of the prompt (Min et al., 2022). According to these studies, ICL may ignore the task defined by the demonstrations and instead resort to using the prior obtained from pretraining. This implies that ICL may not be considered learning under a strict definition, wherein learning must capture the input-output correspondence in a given training set.

**Learning mechanisms.** Theoretical work has explained ICL as implicit Bayesian inference by training language models from scratch on controlled synthetic data (Xie et al.,

2022; Wies et al., 2023; Panwar et al., 2024; Jiang, 2023). Arora et al. (2025) have shown that Bayesian scaling laws are a good fit for the ICL curve. Another line of studies has interpreted ICL as implicitly performing gradient descent (Von Oswald et al., 2023; Ahn et al., 2023) and/or other types of learning algorithms (Akyürek et al., 2023; Garg et al., 2022; Bai et al., 2023; Li et al., 2023). All these mathematical observations encourage us to view ICL as a real learning algorithm and to perform careful empirical investigations to study its properties in real-world settings.

**Pretraining data distribution.** ICL is known to emerge from pretraining when the pretraining data, or its distribution, exhibits a particular set of properties. Chan et al. (2022) found that ICL emerges when data exhibits burstiness (items appear in clusters rather than being uniformly distributed over time) and follows a skewed Zipfian distribution. Raventós et al. (2023) identified a task diversity threshold during pretraining beyond which language models can perform well on unseen ICL regression tasks. Hahn & Goyal (2023) found that ICL arises from generic next-token prediction when the pretraining distribution has a sufficient amounts of compositional structure.

**Prompt optimization.** By deepening our theoretical understanding of the interaction between representation and learning, we can further improve ICL. A common approach to improving LLMs’ performance without any extra weight update is via “prompt engineering,” that is, by crafting prompts manually. Recent studies introduce prompt optimizers that search over strings to identify high-performing prompts (Yuksekgonul et al., 2025; Zhou et al., 2023; Yang et al., 2024b; Guo et al., 2024; Agrawal et al., 2025). These approaches typically optimize one prompt at a time. For ICL classification tasks, we propose a method to optimize the class names on a separate “labeling” set of sentences, and directly use them as labels in new ICL prompts.

### 3. Method

#### 3.1. In-context learning formulation

We formulate the goal of an ICL task as solving

$$\arg \max_{y \in \mathcal{C}} p(\tau(y)|x, D_\tau), \quad (1)$$

where  $D_\tau = \{(x_n, \tau(y_n))\}_{n=1}^N$  refers to a (small) number of input-output pairs.  $\tau(y)$  defines a label set or how we represent each class  $y \in \{1, 2, \dots, \mathcal{C}\}$  as a token in a pre-defined vocabulary, i.e.,  $\tau : \{1, 2, \dots, \mathcal{C}\} \rightarrow V$ , where  $V$  is a vocabulary of unique tokens.  $D_\tau$  refers to presenting the dataset  $D$  using  $\tau$  to encode the classes. By properly formatting  $D$ ,  $x$  and  $\tau(y)$ , LLMs have been found to be

able to implicitly learn to predict the correct label associated with a new instance  $x$ .

Prior work has observed that ICL achieves better performance with gold labels than with abstract labels (Pan et al., 2023). For example, Kirsanov et al. (2025) analyzed a sentiment classification task. The model performed better on an ICL task with gold labels such as  $\{joy, anger, fear\}$  than with abstract labels such as  $\{A, B, C\}$ , even if the input-output correspondence was the same for both label sets.

While abstract labels lead to worse performance than gold labels, the accuracy increases with more examples for either of the label sets. Based on this observation, and taking into account the content effects revealed by Lampinen et al. (2024), we propose that the ICL predictive probability depends on two components:

$$p(\tau(y) | x, D_\tau) = f(q(\tau(y) | x, D_\tau), p(\tau(y) | x)). \quad (2)$$

The first term  $q(\tau(y)|x, D)$  corresponds to *learning*, and the second term  $p(\tau(y)|x)$  corresponds to *prior* knowledge learned by the language model during pretraining. We assume that the first component, *learning*, is largely invariant to how we represent the classes. In other words,

$$q(\tau(y)|x, D_\tau) \approx q(\tau'(y)|x, D_{\tau'}). \quad (3)$$

On the other hand, the prior knowledge must be sensitive to the choice of  $\tau$ , as it lacks the context which is presented in the form of in-context demonstrations. Unless  $\tau(y)$  is *meaningful* under the pretraining corpus, the language model cannot work with an arbitrary representation of a class  $a$  *priori*. That is, it is almost certain that

$$p(\tau(y)|x) \neq p(\tau'(y)|x), \quad (4)$$

for  $\tau \neq \tau'$ .

In this work, we investigate how the contributions of learning and and prior knowledge are disentangled in ICL. We design a readily actionable way to find a good label map  $\tau$  systematically, in order to facilitate this investigation.

#### 3.2. Class representation optimization

We describe a systematic method to choose a label set  $\tau$  that will maximize the performance of ICL across any set of inputs from the same task family. For example, for a sentiment classification task, we can find optimal labels for the classes, and then use these labels as the outputs in ICL demonstrations (input-output pairs) for any other set of inputs.

We assume access to a set of  $K$  examples, which we refer to as a labeling set, and knowledge of the class that each example belongs to (how the examples are clustered). The

goal is to find, for each class, a name, that is represented by a single token in the vocabulary, that is meaningful under the pretraining corpus. To name  $\mathcal{C}$  classes, we want to choose a set of  $\mathcal{C}$  tokens from  $|V|$  possible tokens in a given vocabulary,  $\tau = (l_1, l_2, \dots, l_{\mathcal{C}}) \in V^{\mathcal{C}}$ . A good representation map  $\tau$  should maximize the probability assigned to the correct class  $y^*$ , when represented as  $\tau(y^*)$ . We can write this directly as an objective function:

$$\max_{(l_1, l_2, \dots, l_{\mathcal{C}}) \in V^{\mathcal{C}}} \sum_{k=1}^K \left( f_{\theta}(x_k, l_{y_k}) - \log \sum_{c=1}^{\mathcal{C}} \exp(f_{\theta}(x_k, l_c)) \right), \quad (5)$$

where  $x_k$  are the input examples,  $y_k \in \{1, 2, \dots, \mathcal{C}\}$  are the classes they belong to,  $l_{y_k} = \tau(y_k)$  is the label assigned to class  $y_k$ , and  $f_{\theta}$  is the language model’s logit. Since the tokens in the label set represent class names and appear after the phrase “Category:”, we restrict the vocabulary to tokens that start with the character  $\bar{G}$  (which marks a space and the beginning of a new word).

We optimize this objective via hill climbing, shown in Algorithm 1: we start with an initial random token assignment for each class and iterate the following until no improvements can be made: (1) for each class, try all possible alternative tokens while keeping the rest of class names fixed, (2) evaluate the objective under the current assignment, (3) pick the best token if it improves the overall objective, (4) if there is an improvement, repeat. We run this algorithm ten times while varying random seeds and pick the assignment out of up to ten that maximizes the objective in Equation 5.

As  $K$ , the number of examples used to find a label assignment, increases it becomes harder to find an assignment for which the labels have high probability for many input sentences. To maximize the objective, that assignment should be generalizable: class names should be meaningful for other possible inputs. Thus, as  $K$  increases, we expect the semantics of the labels to be closer to those of gold labels. Equivalently, those labels’ zero-shot accuracy for new inputs would be higher with larger  $K$ . By exploiting the dependence of quality on  $K$ , we obtain a diverse set of label groups that vary in their semantic relevance to the given classification task.

## 4. Experimental setup

We conduct a series of experiments to test the hypothesis that learning and representations are largely disentangled in ICL. First, we want to test whether learning emerges regardless of the choice of label representation. For this to be true, for any label set, the  $N$ -shot accuracy should be increasing with  $N$ . Second, we want to see how representations influence the learning trajectory. For this, we look at how the  $N$ -shot accuracy relates to the zero-shot accuracy (for

---

### Algorithm 1 Hill Climbing for Token Assignment Optimization

---

**Input:** Initial token assignment for each class  
**Input:** Set of candidate tokens, training sentences with labels  
**Output:** Optimized token assignments  
 $assignments \leftarrow initial\_assignments$   
 $objective \leftarrow CalculateObjective(assignments)$   
**repeat**  
    $improved \leftarrow false$   
   **for** each  $class$  in classes **do**  
      $candidates \leftarrow$  all tokens except current token for  $class$   
     **for** each  $token$  in  $candidates$  **do**  
       Compute objective value assigning this token to  $class$  (Eq. 5)  
     **end for**  
      $best\_token \leftarrow$  token with highest objective  
     **if**  $best\_token$  improves objective **then**  
        $assignments[class] \leftarrow best\_token$   
       Update  $objective$   
        $improved \leftarrow true$   
     **break**  
   **end if**  
   **end for**  
**until** not  $improved$  or max iterations reached  
**return**  $assignments, objective$

---

the test input) across the different label representations. We conduct the main experiments with three different size open-weight models: Llama 3.2 1B, Llama 3.1 8B, Llama 3.1 70B Instruct (Grattafiori et al., 2024). We confirm that the findings hold with Mistral-7B-v0.3 (Jiang et al., 2023) and Qwen2.5-7B (Yang et al., 2024a) as well. We first apply the optimization Algorithm 1 to obtain a series of label sets with varying quality for a classification task. Then, we sample demonstrations and name the outputs according to the label set. We prompt a model with the relabeled and concatenated demonstrations to evaluate the ICL performance on these new inputs.

**Data and prompting.** We use a synthetic sentiment classification dataset from Kirsanov et al. (2025), which contains 1,000 sentences split equally among 5 classes for 5-way classification. We also use a subset of 600 sentences covering only 3 of the classes for 3-way classification. We split the dataset into a labeling set (25%), a demonstration set (25%), and a test set (50%). The labeling set is used to enumerate class name assignments, the demonstration set is used for the support examples for ICL, and the test set is used for the query inputs in ICL. For each  $N$ -shot classification task, the task is presented in a minimal format with no explicit instructions, only  $N$  demonstrations and a query sentence. We use a similar setup for the question classification task on TREC (Hovy et al., 2001; Li & Roth, 2002).

**Label sets.** We evaluate different label sets in ICL. These label sets do not break the original input-output correspondence and only replace the original label names, i.e.

the assignment of the classes remains the same. Each label set is obtained by optimizing Equation 5 using  $K \in \{10, 20, \dots, 100\}$  examples. We show the label sets found with each of the three models in Appendix A Table 1 for 3-way classification and Table 2 for 5-way classification. The examples used for finding a label set are the same for each fixed  $K$  across all model sizes. Some of the  $K$  values (adjacent ones) resulted in the same label set.

We illustrate a few of the label sets obtained for 3-way classification with the 70B model. Naturally, using a small  $K = 10$  leads to overfitting on labels that have a high zero-shot probability only for those labeling examples. This yielded random words as labels such as  $\{\textit{biomedical, malware, cloudy}\}$ . With a small  $K$ , we cannot find label sets that appear relevant for a sentiment classification task. For a medium value of  $K = 40$ , the labels obtained are more general  $\{\textit{panic, rage, Cheers}\}$ , a much better fit for the task. While these labels are clearly descriptive, they are slightly odd choices for class names. Finally, using a large  $K = 70$  leads to natural category names for a sentiment classification task such as  $\{\textit{fear, angry, happy}\}$ .

The label sets obtained with the same  $K$  value vary with different models. For instance, for  $K = 100$ , the 1B model found  $\{\textit{spectacle, dance, condolences, peril, pissed}\}$ , the 8B model found  $\{\textit{surprising, joyful, sorrow, fears, anger}\}$ , and the 70B model found  $\{\textit{surprise, happy, sad, anxious, ang}\}$ . In general, the label sets found by larger models appear to be more semantically meaningful.

**In-context learning.** We sample  $N \in \{0, 1, \dots, 40\}$  examples from the demonstration set and name them according to one of the label sets previously obtained. For the 70B model, we only ran experiments with  $N \in \{0, 10, 20, 30, 40\}$  due to compute limitations. For the 1B and 8B models, we ran experiments with  $N$  up to 100, as shown in Appendix B. We create demonstrations with a given label set by using that set to label the inputs in a context, and preserve the original input-output mapping from the dataset. These input-output pairs are concatenated, and, together with a query, are given as a prompt to a model. The model then predicts the class for a novel input selected from the test set, which has not been shown in any of the demonstrations and was not used to compute the label sets. We report the average accuracy for the test set, over 10 runs, in which the inputs of the demonstrations are resampled every time.

## 5. Results

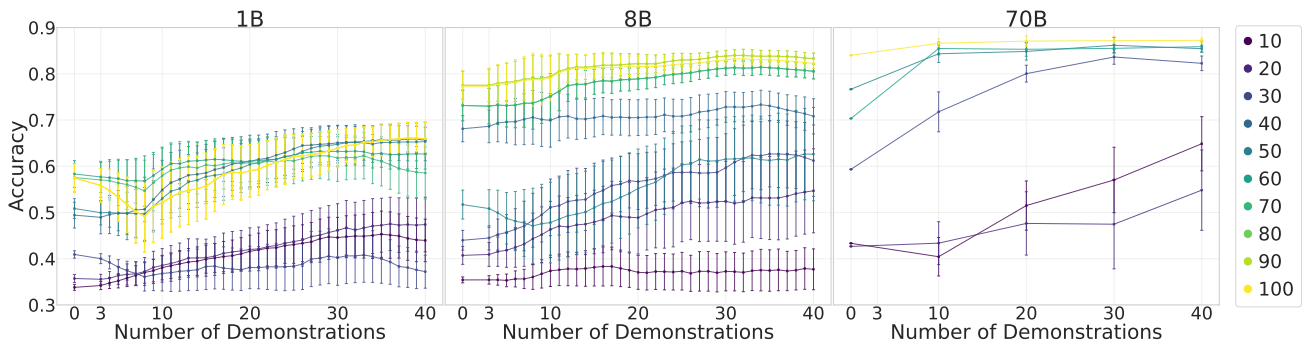
Figure 2 shows the accuracy vs. number of demonstrations in ICL tasks with different label sets for the Llama 1B, 8B, and 70B models, for 3-way (Figure 2a) and 5-way classification (Figure 2b) for the sentiment analysis task. In Appendix C we show results with Mistral-7B-v0.3 (Jiang

et al., 2023) and Qwen2.5-7B (Yang et al., 2024a; Team, 2024). In Appendix D we also show the results on TREC dataset. Across all experimental conditions, we observe that the accuracy is generally increasing with the number of demonstrations. There are exceptions, such as when the label set found has a very small zero-shot test accuracy, most curves stay flat, especially for the harder task of 5-way classification. The zero-shot accuracies span a wide range from chance to ceiling: 33% to 87% for 3-way classification and 20% to 76% for 5-way classification. The representations with a lower zero-shot accuracy typically resulted from optimization on a small  $K$  labeling examples, while those with a high zero-shot accuracy resulted from a larger  $K$ . The ordering of the label sets as determined by their zero-shot accuracy generally stays constant across  $N$ -shot tasks, suggesting a consistent ranking of label sets in terms of ICL performance, regardless of the number of demonstrations.

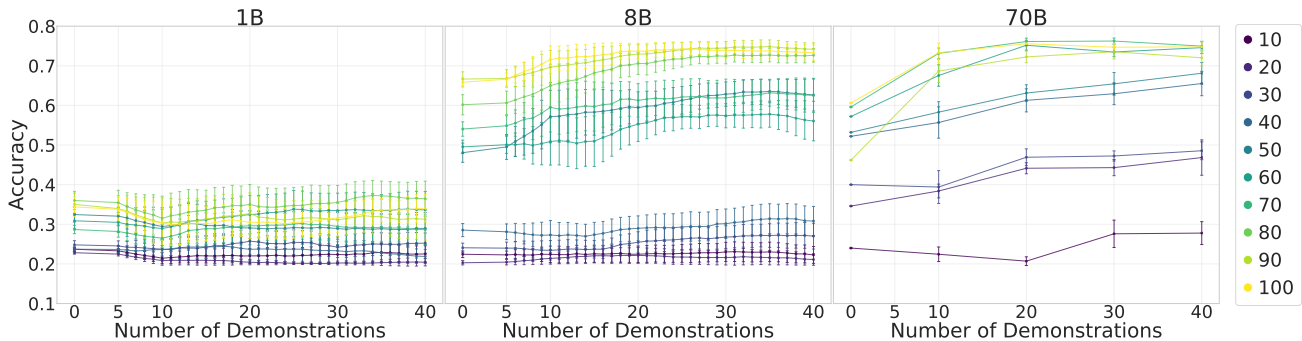
### 5.1. Role of representation in ICL

**Consistent label set ranking.** The  $N$ -shot accuracy of an ICL task using a label set depends on the zero-shot accuracy with that label set: the  $N$ -shot accuracy is typically higher for label sets with higher zero-shot accuracy and can only grow up to a limit. This is consistent across label sets. ICL performs better if the label set is meaningful under the pretraining corpus. We observe that for each  $N$ -shot classification task, the accuracies for ICL with different label sets are ordered according to their initial zero-shot accuracy. We compute the ranking correlation between the zero-shot accuracies and the  $N$ -shot accuracies (of all the label sets) with  $N \in \{\text{num classes}, \dots, 40\}$  for 1B and 8B models,  $N \in \{10, 20, 30, 40\}$  for 70B model. We find that the correlations are indeed high across all model sizes, for both 3-way and 5-way classification (see Figure 3), although there is a lot of variance for the 1B model.

**Representation limits the accuracy range.** If the zero-shot accuracy of a given label set is low, it is very difficult for ICL to reach a high accuracy regardless of how many demonstrations are used. Reaching a high accuracy with a low zero-shot accuracy label set might require a very large number of demonstrations. Most of the curves appear to increase more slowly around 40 demonstrations, indicating a possible upper bound. The chosen label set thus largely determines the range of accuracies attainable with that representation. However, there are exceptions where the accuracy has not yet plateaued with 40 demonstrations (see Figure 2a 70B model,  $K = 10$ ), suggesting that it is possible to overcome the limits of the representation with a large number of demonstrations and a larger model. Our findings indicate that the choice of representation is an essential factor when studying ICL and the role of demonstrations, and they shed



(a) 3-way classification. 1B:  $K \in \{80, 90, 100\}$  examples resulted in the same set; 8B:  $K \in \{60, 70\}$  and  $K \in \{80, 90\}$  same set; 70B:  $K \in \{40, 50\}$  and  $K \in \{70, 80, 90, 100\}$  same set



(b) 5-way classification. 70B:  $K \in \{70, 80\}$  same set

Figure 2. Accuracy vs. number of demonstrations across model sizes for (a) 3-class and (b) 5-class settings. The curves were smoothed with a window size of 10, with error bars showing 95% CI over 10 runs. The legend shows the number of labeling examples  $K$  used to fit the label set. Different  $K$  values may result in the same label sets. For these sets, the color shown is that of the higher  $K$ .

light on some earlier findings. For example, Pan et al. (2023) found that an abstract label set underperformed random allocation of the gold labels to the inputs of the demonstrations, and claimed that this meant that the models could not truly learn the task, but rather relied on their priors. We instead attribute their finding to the fact that the abstract label set has a much lower zero-shot accuracy than a gold label set, and the accuracy increase from learning from additional demonstrations was insufficient to overcome the baseline limitation, which is typically the case for smaller models.

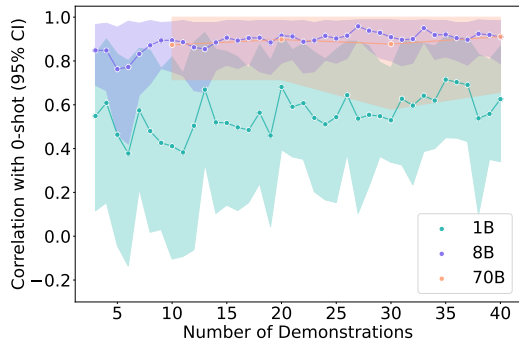
## 5.2. When does ICL learn?

**Learning almost always happens.** We observe that if the zero-shot accuracy is above some threshold, the curves are always increasing regardless of the model size. For the 3-way classification task (Figure 2a), the threshold zero-shot accuracy is very low (33%, chance level), and all curves increase monotonically. For the 5-way classification (Figure 2b), the threshold is higher (40%, double the chance accuracy). Smaller language models can only in-context learn under specific conditions (Schick & Schütze, 2021), and sometimes not at all. We also observe that this is the case, especially for the harder task of 5-way classification,

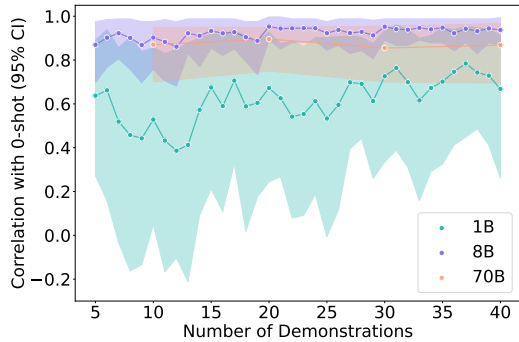
when labels that are not sufficiently semantically meaningful are used. It appears that with a sufficiently good representation, all models, regardless of size, are able to benefit (to different extents) from more demonstrations.

**Model size influences the learning rate.** From Figure 2 we observe that most learning curves are increasing. In Figure 4 we show that the slope depends on model size and zero-shot accuracy. The larger 70B model is more efficient; it makes more use of fewer examples and thus exhibits steeper curves (such as Figure 2a, 70B model,  $K = 30$ ). The  $N$ -shot accuracies for this curve are highly correlated with  $N$  (see Figure 4a orange curve, zero-shot accuracy 59%). With representations of a similar zero-shot accuracy (40%-60% range), the smaller models can also learn, but their curves increase more slowly (and thus have a lower correlation between  $N$  and  $N$ -shot accuracy), suggesting that it would take many more demonstrations to attain the same accuracy that the 70B model achieves with 20-30 demonstrations.

**Learning is conditioned by representation.** Most of the learning curves typically increase, but there is a lot of variance in how much ICL improves with more demonstrations.



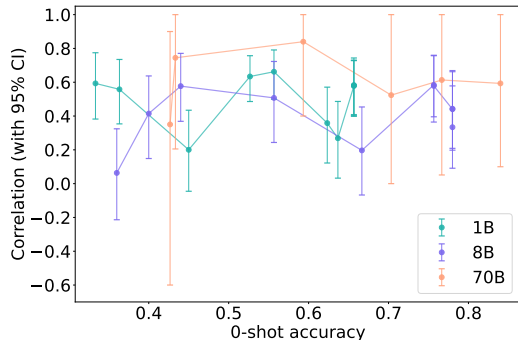
(a) 3-way classification



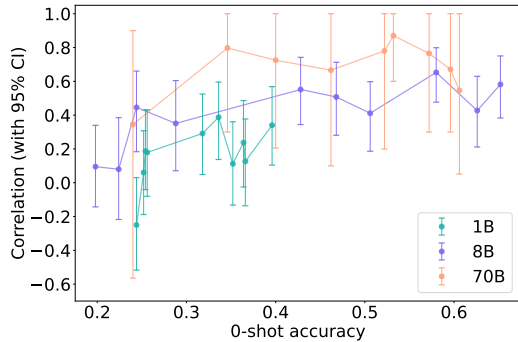
(b) 5-way classification

Figure 3. Ranking correlation coefficient between the zero-shot accuracy and the  $N$ -shot accuracy vs.  $N$  number of demonstrations.  $N \in \{\text{num classes}, \dots, 40\}$  for 1B and 8B models,  $N \in \{10, 20, 30, 40\}$  for 70B model. The CI are computed over 1000 bootstrapping samples from 10 runs per  $N$ -shot accuracy. **The order of label sets in terms of quality stays consistent across  $N$ -shot experiments.**

The increase between the minimum accuracy (zero-shot) and the maximum accuracy (40-shot) ranges from 0% up to 25%. We observe that the representations fall into three categories: small, medium, and high zero-shot accuracy. The small, zero-shot accuracy representations are usually found with a small  $K$  number of labeling examples and are not intuitive or appropriate names for the task. This type of label set makes the task challenging: the model may have to infer the true nature of the task (possibly by inferring more suitable class names) and then map the unintuitive labels onto them. It is not always apparent from the sentences that they illustrate a sentiment classification task. For example, a sentence like “In the upcoming season, I’ll be in the zone every time I step onto the court.” labeled with “cloudy,” might distract the model from the clustering of sentences into appropriate classes. Typically for representations like this, the models start with near-chance zero-shot accuracy, and the accuracy increases only very little regardless of how many demonstrations are presented (e.g. Figure 2a 8B model,  $K = 10$  and Figure 2b 8B model,  $K \in \{10, 20\}$ ). The representations with a medium (40%–60%) zero-shot



(a) 3-way classification



(b) 5-way classification

Figure 4. **Evaluation of learning curves for label sets obtained with different  $K$  labeling examples.** Ranking correlation coefficient between  $N$  and  $N$ -shot accuracy vs. zero-shot accuracy for each curve.  $N \in \{\text{num classes}, \dots, 40\}$  for 1B and 8B models,  $N \in \{10, 20, 30, 40\}$  for 70B model. Higher correlation indicates that the accuracy for that curve is often strictly increasing with  $N$  (steeper curve), while lower accuracy indicates that the accuracy can be plateauing or decreasing on some intervals (flatter curve). The CI are computed over 1000 bootstrapping samples from 10 runs per  $N$ -shot accuracy.

accuracy benefit the most from demonstrations. They can get 15%–25% improvement from the baseline by seeing demonstrations. These labels are sufficiently suggestive of the task  $\{\text{medically, offending, celebrations}\}$  that the model can eventually determine the mapping.

The last group of representations consists of the high zero-shot accuracy representations, those that match or are very close to gold labels. These label sets are already close to the ceiling accuracy possible for each model size. In Figure 2a, for all the models, the curve corresponding to the highest zero-shot accuracy only obtains a 3–5% increase from the baseline before it plateaus. In this group, we observed one exception. In Figure 2a, for 3-way classification with a 1B model, the curve corresponding to  $K \in \{80, 90, 100\}$  initially decreases before increasing. One of the labels in this set is the translation of the word *danger* in Nepali. The ICL task may be harder because it requires multilingual reasoning, which can involve translation as a first step before

figuring out the input-output mapping. In the zero-shot case, there is only one test sentence, and no labels appear in the prompt. The model simply predicts the high probability token, even if it is in a different language than the input, since it does not “know” that the other labels are in a different language. It appears that for  $N < 9$  examples, the model is confused and thus the accuracy decreases. This might happen since as soon as the demonstrations are shown ( $N \geq 3$ ), all three labels appear, so one of them being in a different language than the rest adds the complexity of translation to the original classification task. With enough demonstrations, the model recovers and achieves a high accuracy toward  $N = 40$  demonstrations, as expected for the corresponding zero-shot accuracy.

## 6. Conclusion

The success of ICL has previously been attributed to how the in-context demonstrations are represented, and prior work has questioned whether true learning is, in fact, happening (Perez et al., 2021; Min et al., 2022). Previous observations show that ICL performance improves with the number of demonstrations for both gold and abstract labels (Kirsanov et al., 2025), with gold labels consistently outperforming abstract ones. Based on this, we hypothesized that the choice of representation influences the learning trajectory in ICL. We developed an algorithm to enumerate a spectrum of label representations varying in semantic relevance and tested the performance of these label sets in ICL. We found that the representation of demonstrations determines the baseline accuracy of ICL, as measured by zero-shot performance. The relative quality of the label sets is consistent across demonstrations, and follows the order determined by the baseline accuracies. Furthermore, this baseline typically limits the range of attainable accuracies. It is possible to overcome the limits of the representation, but only with a large number of demonstrations and larger models. The efficiency of learning, measured as the slope of improvement over in-context demonstrations, is influenced both by the quality of representation and model size. Representations with a medium zero-shot accuracy typically benefit the most from seeing more demonstrations and have a higher slope, and larger models can learn faster. In summary, our work reveals the relationship between number of demonstrations and how they are represented on ICL performance, and highlights the importance of considering the representation when studying properties of in-context learning from demonstrations.

Our findings on the interaction between learning and representation in LLMs closely reflect what we know of more conventional neural network learning. The search for high-performing prompts for LLMs is in spirit similar to hyperparameter search (Bengio & LeCun, 2007; Liu et al., 2019)

for neural network classifiers that learn via backpropagation. Perez et al. (2021) found that good prompts are effective because they are chosen using large validation sets. The prompts influence the model behavior similarly to how a choice of initialization influences neural network training. In particular, the choice of label representation in ICL is analogous to the feature selection for the inputs of a neural network classifier. The different choices of representation determine the learning trajectory in both cases: for LLMs, a high quality representation leads to a high zero-shot accuracy and faster convergence; for neural network classifiers, a good set of features can lead to efficient learning (LeCun et al., 2012).

Our framework has a few limitations. First, we assume that the labels are limited to one token, which might not be the case in more complex tasks. We chose to only explore this option since the search space for label sets with our algorithm grows exponentially with the number of tokens. Future work should study multi-token labels. Second, we focus only on how labels are represented. Variations in input representations or prompt format may also influence learning. For instance, choosing inputs with high priors under the model could improve efficiency (Ceballos-Arroyo et al., 2024). Although we treated inputs as fixed, real-world datasets often allow for choosing among many possible inputs. Thus, future work should examine how the inputs and prompt structure affect learning. Third, the in-context learning setting we studied is restricted to one round of prompting. It would be interesting to study the how multi-turn interactions influence ICL in cases where the demonstrations are presented in an online manner, possibly together with extra information in each round (Lee et al., 2023).

Beyond in-context learning, LLMs have shown high performance on complex reasoning tasks, such as programming and mathematical problem solving (Guo et al., 2025; Ruis et al., 2025). Our study also has potential implications about the role of representation in such reasoning. The finding that the representation determines both a baseline accuracy and the efficiency of in-context learning suggests that LLMs already have useful priors, but in order to make the most use of them, we need to present the task in an appropriate manner. Extending these findings about ICL to more complex reasoning tasks could offer a more nuanced understanding about memorization vs. reasoning in LLMs (Bowen et al., 2024; Jin et al., 2025; Salido et al., 2025). Moreover, our findings could explain LLM reasoning failures when changing parameters of an original problem such as document length or the number of variables in a math problem (Malek et al., 2025). Such changes in the prompt, despite attempting to preserve the fundamental difficulty of a problem, result in a significant change in the representation, which lowers the baseline accuracy.

## References

- 440  
441  
442 Agrawal, L. A., Tan, S., Soylu, D., Ziems, N., Khare, R.,  
443 Opsahl-Ong, K., Singhvi, A., Shandilya, H., Ryan, M. J.,  
444 Jiang, M., Potts, C., Sen, K., Dimakis, A. G., Stoica, I.,  
445 Klein, D., Zaharia, M., and Khattab, O. Gepa: Reflective  
446 prompt evolution can outperform reinforcement learn-  
447 ing, 2025. URL [https://arxiv.org/abs/2507.](https://arxiv.org/abs/2507.19457)  
448 [19457](https://arxiv.org/abs/2507.19457).
- 449 Ahn, K., Cheng, X., Daneshmand, H., and Sra,  
450 S. Transformers learn to implement preconditioned  
451 gradient descent for in-context learning. In Oh, A.,  
452 Naumann, T., Globerson, A., Saenko, K., Hardt, M.,  
453 and Levine, S. (eds.), *Advances in Neural Information  
454 Processing Systems*, volume 36, pp. 45614–45650. Curran  
455 Associates, Inc., 2023. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2023/file/8ed3d610ea4b68e7afb30ea7d01422c6-Paper-Conference.pdf)  
456 [cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/8ed3d610ea4b68e7afb30ea7d01422c6-Paper-Conference.pdf)  
457 [8ed3d610ea4b68e7afb30ea7d01422c6-Paper-Conference.](https://proceedings.neurips.cc/paper_files/paper/2023/file/8ed3d610ea4b68e7afb30ea7d01422c6-Paper-Conference.pdf)  
458 [pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8ed3d610ea4b68e7afb30ea7d01422c6-Paper-Conference.pdf).
- 459  
460  
461 Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and  
462 Zhou, D. What learning algorithm is in-context learn-  
463 ing? investigations with linear models. In *The Eleventh  
464 International Conference on Learning Representations*,  
465 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=0g0X4H8yN4I)  
466 [id=0g0X4H8yN4I](https://openreview.net/forum?id=0g0X4H8yN4I).
- 467  
468 Arora, A., Jurafsky, D., Potts, C., and Goodman, N.  
469 Bayesian scaling laws for in-context learning. In *Second  
470 Conference on Language Modeling*, 2025. URL [https:](https://openreview.net/forum?id=U2ihVSREUb)  
471 [//openreview.net/forum?id=U2ihVSREUb](https://openreview.net/forum?id=U2ihVSREUb).
- 472  
473 Bai, Y., Chen, F., Wang, H., Xiong, C., and Mei, S.  
474 Transformers as statisticians: Provable in-context  
475 learning with in-context algorithm selection. In Oh,  
476 A., Naumann, T., Globerson, A., Saenko, K., Hardt,  
477 M., and Levine, S. (eds.), *Advances in Neural Infor-  
478 mation Processing Systems*, volume 36, pp. 57125–  
479 57211. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2023/file/b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.pdf)  
480 [cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.pdf)  
481 [b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.](https://proceedings.neurips.cc/paper_files/paper/2023/file/b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.pdf)  
482 [pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.pdf).
- 483  
484 Bengio, Y. and LeCun, Y. Scaling learning algorithms  
485 towards AI. In *Large Scale Kernel Machines*. MIT Press,  
486 2007.
- 487  
488 Bertsch, A., Ivgi, M., Xiao, E., Alon, U., Berant, J., Gorm-  
489 ley, M. R., and Neubig, G. In-context learning with long-  
490 context models: An in-depth exploration. In Chiruzzo,  
491 L., Ritter, A., and Wang, L. (eds.), *Proceedings of the  
492 2025 Conference of the Nations of the Americas Chapter  
493 of the Association for Computational Linguistics: Hu-  
494 man Language Technologies (Volume 1: Long Papers)*,  
pp. 12119–12149, Albuquerque, New Mexico, April  
2025. Association for Computational Linguistics. ISBN  
979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.  
605. URL [https://aclanthology.org/2025.](https://aclanthology.org/2025.naacl-long.605/)  
[naacl-long.605/](https://aclanthology.org/2025.naacl-long.605/).
- Bowen, C., Sætre, R., and Miyao, Y. A comprehensive  
evaluation of inductive reasoning capabilities and prob-  
lem solving in large language models. In Graham, Y.  
and Purver, M. (eds.), *Findings of the Association for  
Computational Linguistics: EACL 2024*, pp. 323–339, St.  
Julian’s, Malta, March 2024. Association for Computa-  
tional Linguistics. URL [https://aclanthology.](https://aclanthology.org/2024.findings-eacl.22/)  
[org/2024.findings-eacl.22/](https://aclanthology.org/2024.findings-eacl.22/).
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D.,  
 Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G.,  
 Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G.,  
 Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J.,  
 Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M.,  
 Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S.,  
 Radford, A., Sutskever, I., and Amodei, D. Language  
models are few-shot learners. In Larochelle, H.,  
 Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.),  
 *Advances in Neural Information Processing Systems*,  
 volume 33, pp. 1877–1901. Curran Associates, Inc.,  
 2020. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
[cc/paper\\_files/paper/2020/file/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
[1457c0d6bfc4967418bfb8ac142f64a-Paper.](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
[pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- Ceballos-Arroyo, A. M., Munnangi, M., Sun, J., Zhang,  
 K., McInerney, J., Wallace, B. C., and Amir, S. Open  
(clinical) LLMs are sensitive to instruction phrasings. In  
Demner-Fushman, D., Ananiadou, S., Miwa, M., Roberts,  
 K., and Tsujii, J. (eds.), *Proceedings of the 23rd Work-  
shop on Biomedical Natural Language Processing*, pp.  
50–71, Bangkok, Thailand, August 2024. Association  
for Computational Linguistics. doi: 10.18653/v1/2024.  
bionlp-1.5. URL [https://aclanthology.org/](https://aclanthology.org/2024.bionlp-1.5/)  
[2024.bionlp-1.5/](https://aclanthology.org/2024.bionlp-1.5/).
- Chan, S. C. Y., Santoro, A., Lampinen, A. K., Wang, J. X.,  
 Singh, A., Richmond, P. H., McClelland, J., and Hill, F.  
Data distributional properties drive emergent in-context  
learning in transformers, 2022. URL [https://arxiv.](https://arxiv.org/abs/2205.05055)  
[org/abs/2205.05055](https://arxiv.org/abs/2205.05055).
- Chen, J., Chen, L., Zhu, C., and Zhou, T. How  
many demonstrations do you need for in-context learn-  
ing? In Bouamor, H., Pino, J., and Bali, K.  
(eds.), *Findings of the Association for Computational  
Linguistics: EMNLP 2023*, pp. 11149–11159, Singa-  
pore, December 2023. Association for Computational  
Linguistics. doi: 10.18653/v1/2023.findings-emnlp.

- 495 745. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-emnlp.745/)  
 496 [findings-emnlp.745/](https://aclanthology.org/2023.findings-emnlp.745/).  
 497
- 498 Dong, Q., Li, L., Dai, D., Zheng, C., Ma, J., Li, R., Xia, H.,  
 499 Xu, J., Wu, Z., Chang, B., Sun, X., Li, L., and Sui, Z. A  
 500 survey on in-context learning. In Al-Onaizan, Y., Bansal,  
 501 M., and Chen, Y.-N. (eds.), *Proceedings of the 2024 Con-*  
 502 *ference on Empirical Methods in Natural Language Pro-*  
 503 *cessing*, pp. 1107–1128, Miami, Florida, USA, November  
 504 2024. Association for Computational Linguistics. doi:  
 505 10.18653/v1/2024.emnlp-main.64. URL [https://](https://aclanthology.org/2024.emnlp-main.64/)  
 506 [aclanthology.org/2024.emnlp-main.64/](https://aclanthology.org/2024.emnlp-main.64/).
- 507 Garg, S., Tsipras, D., Liang, P. S., and Valiant, G. What  
 508 can transformers learn in-context? a case study of  
 509 simple function classes. In Koyejo, S., Mohamed, S.,  
 510 Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.),  
 511 *Advances in Neural Information Processing Systems*,  
 512 volume 35, pp. 30583–30598. Curran Associates, Inc.,  
 513 2022. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2022/file/c529dba08a146ea8d6cf715ae8930cbe-Paper-Conference.pdf)  
 514 [cc/paper\\_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/c529dba08a146ea8d6cf715ae8930cbe-Paper-Conference.pdf)  
 515 [c529dba08a146ea8d6cf715ae8930cbe-Paper-Conference-](https://proceedings.neurips.cc/paper_files/paper/2022/file/c529dba08a146ea8d6cf715ae8930cbe-Paper-Conference.pdf)  
 516 [pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c529dba08a146ea8d6cf715ae8930cbe-Paper-Conference.pdf).
- 517 Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian,  
 518 A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A.,  
 519 Vaughan, A., Yang, A., Fan, A., Goyal, A., Hartshorn,  
 520 A., Yang, A., Mitra, A., Sravankumar, A., Korenev,  
 521 A., Hinsvark, A., Rao, A., Zhang, A., Rodriguez, A.,  
 522 Gregerson, A., Spataru, A., Roziere, B., Biron, B., Tang,  
 523 B., Chern, B., Caucheteux, C., Nayak, C., Bi, C., Marra,  
 524 C., McConnell, C., Keller, C., Touret, C., Wu, C., Wong,  
 525 C., Ferrer, C. C., Nikolaidis, C., Allonsius, D., Song, D.,  
 526 Pintz, D., Livshits, D., Wyatt, D., Esiobu, D., Choudhary,  
 527 D., Mahajan, D., Garcia-Olano, D., Perino, D., Hupkes,  
 528 D., Lakomkin, E., AlBadawy, E., Lobanova, E., Dinan,  
 529 E., Smith, E. M., Radenovic, F., Guzmán, F., Zhang, F.,  
 530 Synnaeve, G., Lee, G., Anderson, G. L., Thattai, G., Nail,  
 531 G., Mialon, G., Pang, G., Cucurell, G., Nguyen, H., Ko-  
 532 revaar, H., Xu, H., Touvron, H., Zarov, I., Ibarra, I. A.,  
 533 Kloumann, I., Misra, I., Evtimov, I., Zhang, J., Copet, J.,  
 534 Lee, J., Geffert, J., Vranes, J., Park, J., Mahadeokar, J.,  
 535 Shah, J., van der Linde, J., Billock, J., Hong, J., Lee, J.,  
 536 Fu, J., Chi, J., Huang, J., Liu, J., Wang, J., Yu, J., Bitton,  
 537 J., Spisak, J., Park, J., Rocca, J., Johnstun, J., Saxe, J., Jia,  
 538 J., Alwala, K. V., Prasad, K., Upasani, K., Plawiak, K., Li,  
 539 K., Heafield, K., Stone, K., El-Arini, K., Iyer, K., Malik,  
 540 K., Chiu, K., Bhalla, K., Lakhota, K., Rantala-Yeary,  
 541 L., van der Maaten, L., Chen, L., Tan, L., Jenkins, L.,  
 542 Martin, L., Madaan, L., Malo, L., Blecher, L., Landzaat,  
 543 L., de Oliveira, L., Muzzi, M., Pasupuleti, M., Singh,  
 544 M., Paluri, M., Kardas, M., Tsimpoukelli, M., Oldham,  
 545 M., Rita, M., Pavlova, M., Kambadur, M., Lewis, M.,  
 546 Si, M., Singh, M. K., Hassan, M., Goyal, N., Torabi, N.,  
 547 Bashlykov, N., Bogoychev, N., Chatterji, N., Zhang, N.,  
 548 Duchenne, O., Çelebi, O., Alrassy, P., Zhang, P., Li, P.,  
 549 Vasic, P., Weng, P., Bhargava, P., Dubal, P., Krishnan,  
 P., Koura, P. S., Xu, P., He, Q., Dong, Q., Srinivasan,  
 R., Ganapathy, R., Calderer, R., Cabral, R. S., Stojnic,  
 R., Raileanu, R., Maheswari, R., Girdhar, R., Patel, R.,  
 Sauvestre, R., Polidoro, R., Sumbaly, R., Taylor, R., Silva,  
 R., Hou, R., Wang, R., Hosseini, S., Chennabasappa, S.,  
 Singh, S., Bell, S., Kim, S. S., Edunov, S., Nie, S., Narang,  
 S., Raparthy, S., Shen, S., Wan, S., Bhosale, S., Zhang,  
 S., Vandenhende, S., Batra, S., Whitman, S., Sootla, S.,  
 Collot, S., Gururangan, S., Borodinsky, S., Herman, T.,  
 Fowler, T., Sheasha, T., Georgiou, T., Scialom, T., Speck-  
 bacher, T., Mihaylov, T., Xiao, T., Karn, U., Goswami, V.,  
 Gupta, V., Ramanathan, V., Kerkez, V., Gonguet, V., Do,  
 V., Vogeti, V., Albiero, V., Petrovic, V., Chu, W., Xiong,  
 W., Fu, W., Meers, W., Martinet, X., Wang, X., Wang,  
 X., Tan, X. E., Xia, X., Xie, X., Jia, X., Wang, X., Gold-  
 schlag, Y., Gaur, Y., Babaei, Y., Wen, Y., Song, Y., Zhang,  
 Y., Li, Y., Mao, Y., Coudert, Z. D., Yan, Z., Chen, Z.,  
 Papakipos, Z., Singh, A., Srivastava, A., Jain, A., Kelsey,  
 A., Shajnfeld, A., Gangidi, A., Victoria, A., Goldstand,  
 A., Menon, A., Sharma, A., Boesenberg, A., Baevski, A.,  
 Feinstein, A., Kallet, A., Sangani, A., Teo, A., Yunus, A.,  
 Lupu, A., Alvarado, A., Caples, A., Gu, A., Ho, A., Poul-  
 ton, A., Ryan, A., Ramchandani, A., Dong, A., Franco,  
 A., Goyal, A., Saraf, A., Chowdhury, A., Gabriel, A.,  
 Bharambe, A., Eisenman, A., Yazdan, A., James, B.,  
 Maurer, B., Leonhardi, B., Huang, B., Loyd, B., Paola,  
 B. D., Paranjape, B., Liu, B., Wu, B., Ni, B., Hancock,  
 B., Wasti, B., Spence, B., Stojkovic, B., Gamido, B.,  
 Montalvo, B., Parker, C., Burton, C., Mejia, C., Liu, C.,  
 Wang, C., Kim, C., Zhou, C., Hu, C., Chu, C.-H., Cai, C.,  
 Tindal, C., Feichtenhofer, C., Gao, C., Civin, D., Beaty,  
 D., Kreymer, D., Li, D., Adkins, D., Xu, D., Testuggine,  
 D., David, D., Parikh, D., Liskovich, D., Foss, D., Wang,  
 D., Le, D., Holland, D., Dowling, E., Jamil, E., Mont-  
 gomery, E., Presani, E., Hahn, E., Wood, E., Le, E.-T.,  
 Brinkman, E., Arcaute, E., Dunbar, E., Smothers, E., Sun,  
 F., Kreuk, F., Tian, F., Kokkinos, F., Ozgenel, F., Cag-  
 gioni, F., Kanayet, F., Seide, F., Florez, G. M., Schwarz,  
 G., Badeer, G., Swee, G., Halpern, G., Herman, G., Sizov,  
 G., Guangyi, Zhang, Lakshminarayanan, G., Inan, H.,  
 Shojanazeri, H., Zou, H., Wang, H., Zha, H., Habeeb, H.,  
 Rudolph, H., Suk, H., Aspegren, H., Goldman, H., Zhan,  
 H., Damaj, I., Molybog, I., Tufanov, I., Leontiadis, I.,  
 Veliche, I.-E., Gat, I., Weissman, J., Geboski, J., Kohli,  
 J., Lam, J., Asher, J., Gaya, J.-B., Marcus, J., Tang, J.,  
 Chan, J., Zhen, J., Reizenstein, J., Teboul, J., Zhong, J.,  
 Jin, J., Yang, J., Cummings, J., Carvill, J., Shepard, J.,  
 McPhee, J., Torres, J., Ginsburg, J., Wang, J., Wu, K., U,  
 K. H., Saxena, K., Khandelwal, K., Zand, K., Matosich,  
 K., Veeraraghavan, K., Michelena, K., Li, K., Jagadeesh,  
 K., Huang, K., Chawla, K., Huang, K., Chen, L., Garg,  
 L., A. L., Silva, L., Bell, L., Zhang, L., Guo, L., Yu, L.,

- 550 Moshkovich, L., Wehrstedt, L., Khabsa, M., Avalani, M.,  
 551 Bhatt, M., Mankus, M., Hasson, M., Lennie, M., Reso,  
 552 M., Groshev, M., Naumov, M., Lathi, M., Keneally, M.,  
 553 Liu, M., Seltzer, M. L., Valko, M., Restrepo, M., Patel,  
 554 M., Vyatskov, M., Samvelyan, M., Clark, M., Macey,  
 555 M., Wang, M., Hermoso, M. J., Metanat, M., Rastegari,  
 556 M., Bansal, M., Santhanam, N., Parks, N., White, N.,  
 557 Bawa, N., Singhal, N., Egebo, N., Usunier, N., Mehta,  
 558 N., Laptev, N. P., Dong, N., Cheng, N., Chernoguz, O.,  
 559 Hart, O., Salpekar, O., Kalinli, O., Kent, P., Parekh, P.,  
 560 Saab, P., Balaji, P., Rittner, P., Bontrager, P., Roux, P.,  
 561 Dollar, P., Zvyagina, P., Ratanchandani, P., Yuvraj, P.,  
 562 Liang, Q., Alao, R., Rodriguez, R., Ayub, R., Murthy, R.,  
 563 Nayani, R., Mitra, R., Parthasarathy, R., Li, R., Hogan,  
 564 R., Battey, R., Wang, R., Howes, R., Rinott, R., Mehta,  
 565 S., Siby, S., Bondu, S. J., Datta, S., Chugh, S., Hunt, S.,  
 566 Dhillon, S., Sidorov, S., Pan, S., Mahajan, S., Verma,  
 567 S., Yamamoto, S., Ramaswamy, S., Lindsay, S., Lindsay,  
 568 S., Feng, S., Lin, S., Zha, S. C., Patil, S., Shankar, S.,  
 569 Zhang, S., Zhang, S., Wang, S., Agarwal, S., Sajuyigbe,  
 570 S., Chintala, S., Max, S., Chen, S., Kehoe, S., Satter-  
 571 field, S., Govindaprasad, S., Gupta, S., Deng, S., Cho,  
 572 S., Virk, S., Subramanian, S., Choudhury, S., Goldman,  
 573 S., Remez, T., Glaser, T., Best, T., Koehler, T., Robinson,  
 574 T., Li, T., Zhang, T., Matthews, T., Chou, T., Shaked,  
 575 T., Vontimitta, V., Ajayi, V., Montanez, V., Mohan, V.,  
 576 Kumar, V. S., Mangla, V., Ionescu, V., Poenaru, V., Mi-  
 577 hailescu, V. T., Ivanov, V., Li, W., Wang, W., Jiang, W.,  
 578 Bouaziz, W., Constable, W., Tang, X., Wu, X., Wang, X.,  
 579 Wu, X., Gao, X., Kleinman, Y., Chen, Y., Hu, Y., Jia, Y.,  
 580 Qi, Y., Li, Y., Zhang, Y., Zhang, Y., Adi, Y., Nam, Y., Yu,  
 581 Wang, Zhao, Y., Hao, Y., Qian, Y., Li, Y., He, Y., Rait,  
 582 Z., DeVito, Z., Rosnbrick, Z., Wen, Z., Yang, Z., Zhao,  
 583 Z., and Ma, Z. The llama 3 herd of models, 2024. URL  
 584 <https://arxiv.org/abs/2407.21783>.  
 585
- 586 Guo, D., Yang, D., Zhang, H., Song, J., Wang, P., Zhu, Q.,  
 587 Xu, R., Zhang, R., Ma, S., Bi, X., Zhang, X., Yu, X., Wu,  
 588 Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A.,  
 589 Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C.,  
 590 Deng, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin,  
 591 F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H.,  
 592 Xu, H., Ding, H., Gao, H., Qu, H., Li, H., Guo, J., Li,  
 593 J., Chen, J., Yuan, J., Tu, J., Qiu, J., Li, J., Cai, J. L., Ni,  
 594 J., Liang, J., Chen, J., Dong, K., Hu, K., You, K., Gao,  
 595 K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L.,  
 596 Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang,  
 597 M., Zhang, M., Tang, M., Zhou, M., Li, M., Wang, M.,  
 598 Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen,  
 599 Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen,  
 600 R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye,  
 601 S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S., Zhou, S.,  
 602 Wu, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Liu,  
 603 W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L.,  
 604 An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X.,  
 Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X.,  
 Li, X. Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X.,  
 Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K., Wang,  
 Y. Q., Wei, Y. X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun,  
 Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y.,  
 Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou,  
 Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo,  
 Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Huang, Y., Li,  
 Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y.,  
 Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang,  
 Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu,  
 Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu,  
 Z., Zhang, Z., and Zhang, Z. Deepseek-r1 incentivizes  
 reasoning in llms through reinforcement learning. *Nature*,  
 645(8081):633–638, Sep 2025. ISSN 1476-4687. doi:  
 10.1038/s41586-025-09422-z. URL <https://doi.org/10.1038/s41586-025-09422-z>.
- Guo, Q., Wang, R., Guo, J., Li, B., Song, K., Tan, X., Liu,  
 G., Bian, J., and Yang, Y. Connecting large language  
 models with evolutionary algorithms yields powerful  
 prompt optimizers. In *The Twelfth International Confer-  
 ence on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ZG3RaNIso8>.
- Hahn, M. and Goyal, N. A theory of emergent in-context  
 learning as implicit structure induction, 2023. URL  
<https://arxiv.org/abs/2303.07971>.
- Hovy, E., Gerber, L., Hermjakob, U., Lin, C.-Y., and  
 Ravichandran, D. Toward semantics-based answer  
 pinpointing. In *Proceedings of the First Interna-  
 tional Conference on Human Language Technology Re-  
 search*, 2001. URL <https://aclanthology.org/H01-1069/>.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C.,  
 Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel,  
 G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-  
 A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix,  
 T., and Sayed, W. E. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Jiang, H. A latent space theory for emergent abilities in  
 large language models, 2023. URL <https://arxiv.org/abs/2304.09960>.
- Jin, M., Luo, W., Cheng, S., Wang, X., Hua, W., Tang, R.,  
 Wang, W. Y., and Zhang, Y. Disentangling memory and  
 reasoning ability in large language models. In Che, W.,  
 Nabende, J., Shutova, E., and Pilehvar, M. T. (eds.), *Pro-  
 ceedings of the 63rd Annual Meeting of the Association  
 for Computational Linguistics (Volume 1: Long Papers)*,  
 pp. 1681–1701, Vienna, Austria, July 2025. Association

- 605 for Computational Linguistics. ISBN 979-8-89176-251-  
606 0. doi: 10.18653/v1/2025.acl-long.84. URL <https://aclanthology.org/2025.acl-long.84/>.
- 607  
608  
609 Kirsanov, A., Chou, C.-N., Cho, K., and Chung, S. The  
610 geometry of prompting: Unveiling distinct mechanisms  
611 of task adaptation in language models. In Chiruzzo, L.,  
612 Ritter, A., and Wang, L. (eds.), *Findings of the Association for Computational Linguistics: NAACL 2025*,  
613 pp. 1855–1888, Albuquerque, New Mexico, April 2025.  
614 Association for Computational Linguistics. ISBN 979-  
615 8-89176-195-7. doi: 10.18653/v1/2025.findings-naacl.  
616 100. URL [https://aclanthology.org/2025.  
617 findings-naacl.100/](https://aclanthology.org/2025.findings-naacl.100/).
- 618  
619 Lampinen, A. K., Dasgupta, I., Chan, S. C. Y., Sheahan,  
620 H. R., Creswell, A., Kumaran, D., McClelland, J. L., and  
621 Hill, F. Language models, like humans, show content  
622 effects on reasoning tasks. *PNAS Nexus*, 3(7):pgae233, 07  
623 2024. ISSN 2752-6542. doi: 10.1093/pnasnexus/pgae233.  
624 URL [https://doi.org/10.1093/pnasnexus/  
625 pgae233](https://doi.org/10.1093/pnasnexus/pgae233).
- 626  
627 LeCun, Y. A., Bottou, L., Orr, G. B., and Müller, K.-R.  
628 *Efficient BackProp*, pp. 9–48. Springer Berlin Heidelberg,  
629 Berlin, Heidelberg, 2012. ISBN 978-3-642-35289-8. doi:  
630 10.1007/978-3-642-35289-8\_3. URL [https://doi.  
631 org/10.1007/978-3-642-35289-8\\_3](https://doi.org/10.1007/978-3-642-35289-8_3).
- 632  
633 Lee, D., Kim, S., Lee, M., Lee, H., Park, J., Lee, S.-W., and  
634 Jung, K. Asking clarification questions to handle ambi-  
635 guity in open-domain QA. In Bouamor, H., Pino, J., and  
636 Bali, K. (eds.), *Findings of the Association for Computa-  
637 tional Linguistics: EMNLP 2023*, pp. 11526–11544, Sin-  
638 gapore, December 2023. Association for Computational  
639 Linguistics. doi: 10.18653/v1/2023.findings-emnlp.  
640 772. URL [https://aclanthology.org/2023.  
641 findings-emnlp.772/](https://aclanthology.org/2023.findings-emnlp.772/).
- 642  
643 Li, X. and Roth, D. Learning question classifiers. In  
644 *COLING 2002: The 19th International Conference on  
645 Computational Linguistics*, 2002. URL [https://  
646 aclanthology.org/C02-1150/](https://aclanthology.org/C02-1150/).
- 647  
648 Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. Trans-  
649 formers as algorithms: Generalization and stability in  
650 in-context learning. In Krause, A., Brunskill, E., Cho,  
651 K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.),  
652 *Proceedings of the 40th International Conference on Ma-  
653 chine Learning*, volume 202 of *Proceedings of Machine  
654 Learning Research*, pp. 19565–19594. PMLR, 23–29 Jul  
655 2023. URL [https://proceedings.mlr.press/  
656 v202/li231.html](https://proceedings.mlr.press/v202/li231.html).
- 657  
658 Liu, H., Simonyan, K., and Yang, Y. DARTS: Differ-  
659 entiable architecture search. In *International Confer-  
ence on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S1eYHoC5FX>.
- Liu, Y., Liu, J., Shi, X., Cheng, Q., Huang, Y., and Lu, W. Let’s learn step by step: Enhancing in-context learning ability with curriculum learning, 2024. URL <https://arxiv.org/abs/2402.10738>.
- Lu, Y., Bartolo, M., Moore, A., Riedel, S., and Stenetorp, P. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Annual Meeting of the Association for Computational Linguistics*, 2021. URL <https://api.semanticscholar.org/CorpusID:233296494>.
- Malek, A., Ge, J., Lazic, N., Jin, C., György, A., and Szepesvári, C. Frontier llms still struggle with simple reasoning tasks, 2025. URL <https://arxiv.org/abs/2507.07313>.
- McCoy, R. T., Yao, S., Friedman, D., Hardy, M. D., and Griffiths, T. L. Embers of autoregression show how large language models are shaped by the problem they are trained to solve. *Proceedings of the National Academy of Sciences*, 121(41):e2322420121, 2024. doi: 10.1073/pnas.2322420121. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2322420121>.
- Min, S., Lyu, X., Holtzman, A., Artetxe, M., Lewis, M., Hajishirzi, H., and Zettlemoyer, L. Rethinking the role of demonstrations: What makes in-context learning work? In Goldberg, Y., Kozareva, Z., and Zhang, Y. (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.759. URL [https://aclanthology.org/2022.  
emnlp-main.759/](https://aclanthology.org/2022.emnlp-main.759/).
- Pan, J., Gao, T., Chen, H., and Chen, D. What in-context learning “learns” in-context: Disentangling task recognition and task learning. In *Annual Meeting of the Association for Computational Linguistics*, 2023. URL <https://api.semanticscholar.org/CorpusID:258740972>.
- Panwar, M., Ahuja, K., and Goyal, N. In-context learning through the bayesian prism. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=HX5ujdsSon>.
- Perez, E., Kiela, D., and Cho, K. True few-shot learning with language models. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*,

- 660 volume 34, pp. 11054–11070. Curran Associates, Inc.,  
 661 2021. URL [https://proceedings.neurips.  
 662 cc/paper\\_files/paper/2021/file/  
 663 5c04925674920eb58467fb52ce4ef728-Paper.  
 664 pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/5c04925674920eb58467fb52ce4ef728-Paper.pdf).
- 665 Raventós, A., Paul, M., Chen, F., and Ganguli, S.  
 666 Pretraining task diversity and the emergence of  
 667 non-bayesian in-context learning for regression.  
 668 In Oh, A., Naumann, T., Globerson, A., Saenko,  
 669 K., Hardt, M., and Levine, S. (eds.), *Advances  
 670 in Neural Information Processing Systems*, vol-  
 671 ume 36, pp. 14228–14246. Curran Associates, Inc.,  
 672 2023. URL [https://proceedings.neurips.  
 673 cc/paper\\_files/paper/2023/file/  
 674 2e10b2c2e1aa4f8083c37dfe269873f8-Paper-Conference.  
 675 pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/2e10b2c2e1aa4f8083c37dfe269873f8-Paper-Conference.pdf).
- 676 Ruis, L., Mozes, M., Bae, J., Kamalakara, S. R., Gnanesh-  
 677 war, D., Locatelli, A., Kirk, R., Rocktäschel, T., Grefen-  
 678 stette, E., and Bartolo, M. Procedural knowledge in  
 679 pretraining drives reasoning in large language models.  
 680 In *The Thirteenth International Conference on Learning  
 681 Representations*, 2025. URL [https://openreview.  
 682 net/forum?id=1hQKHU5Mx](https://openreview.net/forum?id=1hQKHU5Mx).
- 683 Salido, E. S., Gonzalo, J., and Marco, G. None of  
 684 the others: a general technique to distinguish reason-  
 685 ing from memorization in multiple-choice llm evalua-  
 686 tion benchmarks, 2025. URL [https://arxiv.org/  
 687 abs/2502.12896](https://arxiv.org/abs/2502.12896).
- 688 Schick, T. and Schütze, H. It’s not just size that mat-  
 689 ters: Small language models are also few-shot learn-  
 690 ers. In Toutanova, K., Rumshisky, A., Zettlemoyer,  
 691 L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cot-  
 692 terell, R., Chakraborty, T., and Zhou, Y. (eds.), *Pro-  
 693 ceedings of the 2021 Conference of the North Ameri-  
 694 can Chapter of the Association for Computational Lin-  
 695 guistics: Human Language Technologies*, pp. 2339–  
 696 2352, Online, June 2021. Association for Computa-  
 697 tional Linguistics. doi: 10.18653/v1/2021.naacl-main.  
 698 185. URL [https://aclanthology.org/2021.  
 699 naacl-main.185/](https://aclanthology.org/2021.naacl-main.185/).
- 700 Team, Q. Qwen2.5: A party of foundation models, Septem-  
 701 ber 2024. URL [https://qwenlm.github.io/  
 702 blog/qwen2.5/](https://qwenlm.github.io/blog/qwen2.5/).
- 703 Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento,  
 704 J., Mordvintsev, A., Zhmoginov, A., and Vladymy-  
 705 ro, M. Transformers learn in-context by gradient  
 706 descent. In Krause, A., Brunskill, E., Cho, K., En-  
 707 gelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Pro-  
 708 ceedings of the 40th International Conference on Ma-  
 709 chine Learning*, volume 202 of *Proceedings of Machine  
 710 Learning Research*, pp. 35151–35174. PMLR, 23–29 Jul  
 711 2023. URL [https://proceedings.mlr.press/  
 712 v202/von-oswald23a.html](https://proceedings.mlr.press/v202/von-oswald23a.html).
- 713 Wies, N., Levine, Y., and Shashua, A. The learnability of  
 714 in-context learning. In Oh, A., Naumann, T., Globerson,  
 A., Saenko, K., Hardt, M., and Levine, S. (eds.),  
*Advances in Neural Information Processing Systems*,  
 volume 36, pp. 36637–36651. Curran Associates, Inc.,  
 2023. URL [https://proceedings.neurips.  
 cc/paper\\_files/paper/2023/file/  
 73950f0eb4ac0925dc71ba2406893320-Paper-Conference  
 pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/73950f0eb4ac0925dc71ba2406893320-Paper-Conference.pdf).
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An ex-  
 planation of in-context learning as implicit bayesian infer-  
 ence. In *International Conference on Learning Represen-  
 tations*, 2022. URL [https://openreview.net/  
 forum?id=RdJVFCHjUMI](https://openreview.net/forum?id=RdJVFCHjUMI).
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C.,  
 Li, C., Li, C., Liu, D., Huang, F., Dong, G., Wei, H., Lin,  
 H., Tang, J., Wang, J., Yang, J., Tu, J., Zhang, J., Ma, J.,  
 Xu, J., Zhou, J., Bai, J., He, J., Lin, J., Dang, K., Lu, K.,  
 Chen, K., Yang, K., Li, M., Xue, M., Ni, N., Zhang, P.,  
 Wang, P., Peng, R., Men, R., Gao, R., Lin, R., Wang, S.,  
 Bai, S., Tan, S., Zhu, T., Li, T., Liu, T., Ge, W., Deng,  
 X., Zhou, X., Ren, X., Zhang, X., Wei, X., Ren, X., Fan,  
 Y., Yao, Y., Zhang, Y., Wan, Y., Chu, Y., Liu, Y., Cui, Z.,  
 Zhang, Z., and Fan, Z. Qwen2 technical report. *arXiv  
 preprint arXiv:2407.10671*, 2024a.
- Yang, C., Wang, X., Lu, Y., Liu, H., Le, Q. V., Zhou, D., and  
 Chen, X. Large language models as optimizers. In *The  
 Twelfth International Conference on Learning Representa-  
 tions*, 2024b. URL [https://openreview.net/  
 forum?id=Bb4VGOWELI](https://openreview.net/forum?id=Bb4VGOWELI).
- Yuksekgonul, M., Bianchi, F., Boen, J., Liu, S., Lu, P.,  
 Huang, Z., Guestrin, C., and Zou, J. Optimizing gener-  
 ative ai by backpropagating language model feedback.  
*Nature*, 639:609–616, 2025.
- Zhao, T., Wallace, E., Feng, S., Klein, D., and Singh,  
 S. Calibrate before use: Improving few-shot per-  
 formance of language models. In *International  
 Conference on Machine Learning*, 2021. URL [https://api.semanticscholar.org/CorpusID:  
 231979430](https://api.semanticscholar.org/CorpusID:231979430).
- Zhou, Y., Muresanu, A. I., Han, Z., Paster, K., Pitis,  
 S., Chan, H., and Ba, J. Large language models are  
 human-level prompt engineers. In *The Eleventh In-  
 ternational Conference on Learning Representations*,  
 2023. URL [https://openreview.net/forum?  
 id=92gvk82DE-](https://openreview.net/forum?id=92gvk82DE-).

**A. List of label sets**

We show the label sets obtained from Algorithm 1 on the sentiment analysis task.

<b>K</b>	<b>1B</b>	<b>8B</b>	<b>70B</b>
10	Nutrition, Giz Legends	Gluten, Laptop clouds	biomedical, malware cloudy
20	diabetes, Hacker Presbyterian	Diabetes, Revenge spirit	fitness, computer joyful
30	overweight, annoy scholarships	FDA, console celebration	Obesity, rage celebration
40	medically, offending celebrating	fearful, malicious celebration	panic, rage Cheers
50	medically, offending celebrations	digestive, insulting accomplishments	panic, rage Cheers
60	panicked, offending celebrations	fears, insults joyful	worry, complain celebration
70	hazardous, offending celebrations	fears, insults joyful	fear, angry happy
80	འཕྲིད་ཆུང་ལྔ་ (danger, Nepali), offending celebrations	fears, complaints joyful	fear, angry happy
90	འཕྲིད་ཆུང་ལྔ་ (danger, Nepali), offending celebrations	fears, complaints joyful	fear, angry happy
100	འཕྲིད་ཆུང་ལྔ་ (danger, Nepali), offending celebrations	fears, complaint joyful	fear, angry happy

Table 1. Label sets obtained from running Algorithm 1 on  $K$  labeling examples for 3-way classification. The gold labels are “fear, anger, joy.”

<b>K</b>	<b>1B</b>	<b>8B</b>	<b>70B</b>
10	movie, Musik Causes, Roller NRL	theater, COLOR HEALTH, ride Offensive	Marvel, MUSIC HEALTH, roller veh
20	witches, audition bere, adip Messi	cinema, Broadcasting Deng, nut Rugby	Magical, positive loss, Dietary Baseball
30	trick, Dresses bere, hysteria Messi	surprising MÃ©d ( <i>Med</i> , French) ÑġD¾D½ ( <i>dream</i> , Russian) snack, soccer	surprise, celebration tragedy, amused ìĤĤĤ-ìĤĤ ( <i>sports</i> , Korean)
40	puzzle, Ventures àĤĤ (A, Marathi), carniv Penalty	surprising ÑĤDµDĤĤ ( <i>tel</i> , Russian) resignation, aliment soccer	surprising, positives heartbreaking DĤD,ÑĤ ( <i>shout</i> , Bulgarian) frustrated
50	spectacle, talent mourn, endanger offense	surprised, baĀŁarĀ± ( <i>success</i> , Turkish) sadness, xen, Rage	surprising, positives heartbreaking, Brussels frustrated
60	spectacle, production mourn, peril offense	amazed, D½D°ÑĤD° ( <i>science</i> , Ukrainian) mourn, scare, brawl	surprising, positives heartbreaking, nerv agg
70	spectacle, productions mourn, peril racket	surprising, ĤĤ (?) sorrow, terror hostile	surprise, pleasant sorrow, fears rage
80	spectacle, dance mourn, peril criticizing	surprising, celebrates condolences, terror rage	surprise, pleasant sorrow, fears rage
90	magician, dancer mourning, risking wrath	surprising, joyful sorrow, fears rage	surprise, Lift broken, fears rage
100	spectacle, dance condolences, peril pissed	surprising, joyful sorrow, fears anger	surprise, happy sad, anxious ang

Table 2. Label sets obtained from from running Algorithm 1 on  $K$  labeling examples for 5-way classification. The gold labels are “surprise, joy, sadness, fear, anger.”

B. Learning curves

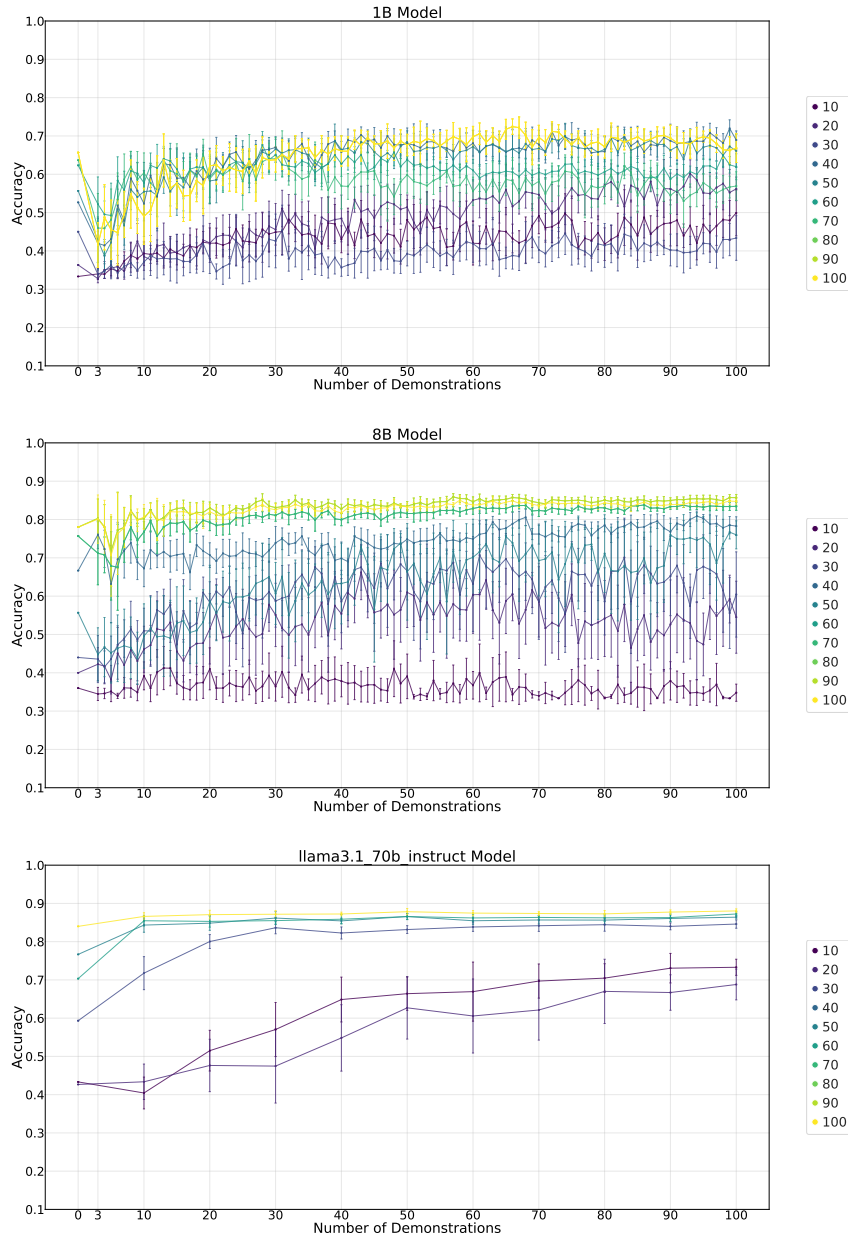


Figure 5. Full raw (unsmoothed) learning curves for up to 100 demonstrations for Llama models for 3-way classification for the sentiment analysis task.

## Relationship between Representation and In-context learning

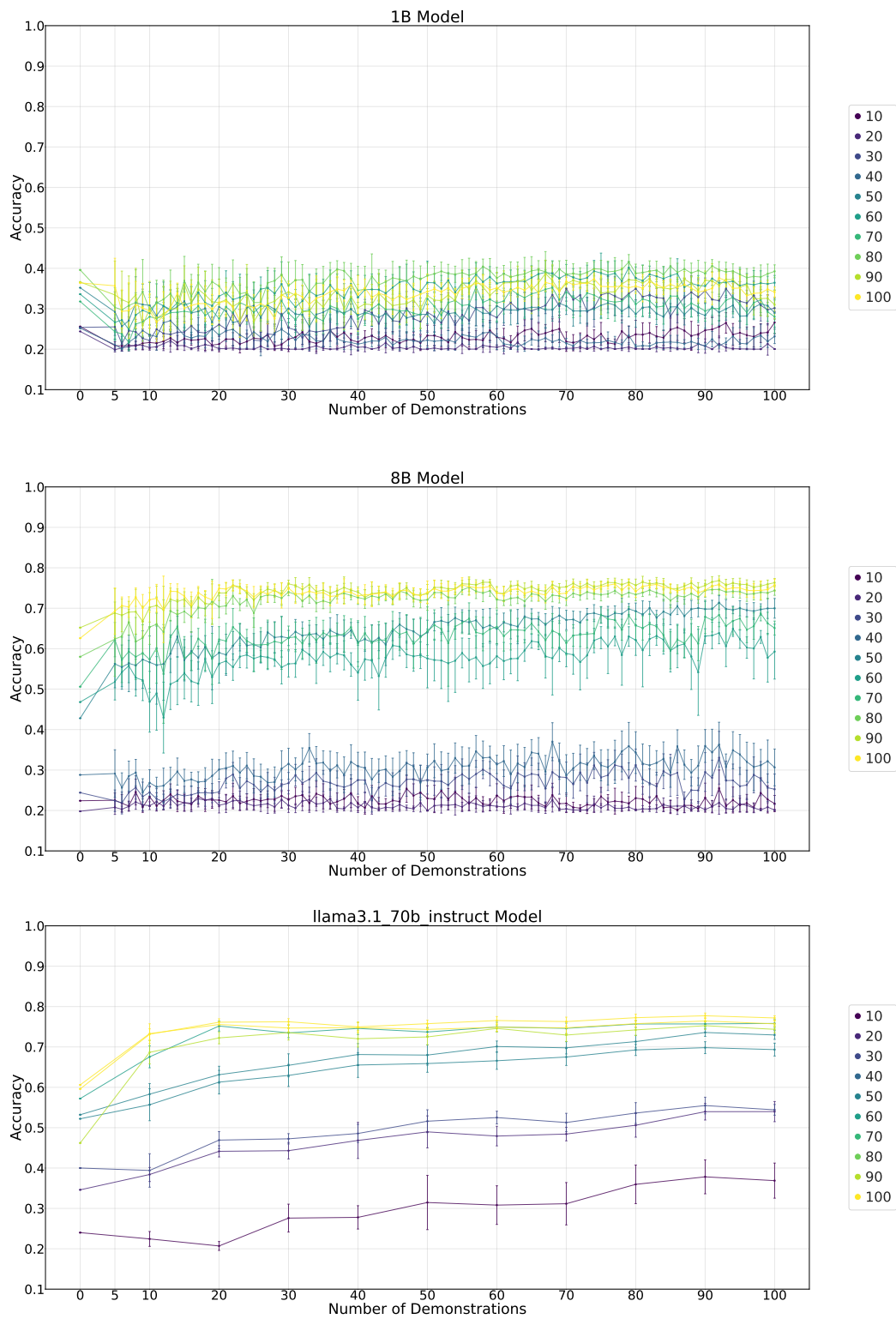


Figure 6. Full raw (unsmoothed) learning curves for up to 100 demonstrations for Llama models for 5-way classification for the sentiment analysis task.

### C. Additional models: Mistral-7B-v0.3 and Qwen2.5-7B

We show the label sets and learning curves for the sentiment analysis task using models from different families Mistral-7B-v0.3 and Qwen2.5-7B for 3-way and 5-way classification.

<b>K</b>	<b>Mistral</b>	<b>Qwen</b>
10	diet, hack, wholes	Sick, offensive, Ath
20	nutrition, hack, excitement	JsonResult, parliamentary, applause
30	protein, complaint, excited	weighing, warn, congrat
40	fear, angry, insp	ArgumentError, bitch, Wonderful
50	fear, angry, joy	fears, hatred, celebration
60	fear, angry, joy	noir, bitch, Wonderful
70	panic, rage, smiles	fears, abusive, congrat
80	fear, anger, joy	nightmares, offender, luz
90	fear, angry, joy	terror, offending, brag
100	fear, angry, joy	terror, offending, brag

Table 3. Label sets obtained from running Algorithm 1 on  $K$  examples for 3-way classification

<b>K</b>	<b>Mistral</b>	<b>Qwen</b>
10	photography, festival, AU, roll, Bull	Dialogue, SUM, clinical, Roller, negatives
20	aston, rehe, diseases, stomach, managers	Eye, academy, diagnosed, kaufen, threats
30	amazing, dress, defeat, monster, hockey	surprising, moda, diagnoses, fart, unlawful
40	surprise, competition, loss, monster, soccer	surprising, RoundedRectangle, privacy, immature, mud
50	surprise, dress, grief, monster, angry	surprising, RoundedRectangle, failures, Moo, FUCK
60	amazing, invent, depress, fright, piss	exploding, ExecutionContext, Ø§ÙØ, onOptionsItemSelected, misogyn
70	amazing, next, depress, terror, angry,	astonishing, Validates, distressed, nightmares, mud
80	aston, inspire, despair, fears, rage	astonishing, remains, distressed, feared, abuses
90	aston, pose, unhappy, afraid, complaint	surprising, remainder, failures, feared, bitter
100	aston, inspire, despair, fears, rage	surprising, remainder, failures, feared, bitter

Table 4. Label sets obtained from running Algorithm 1 on  $K$  examples for 5-way classification

990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044

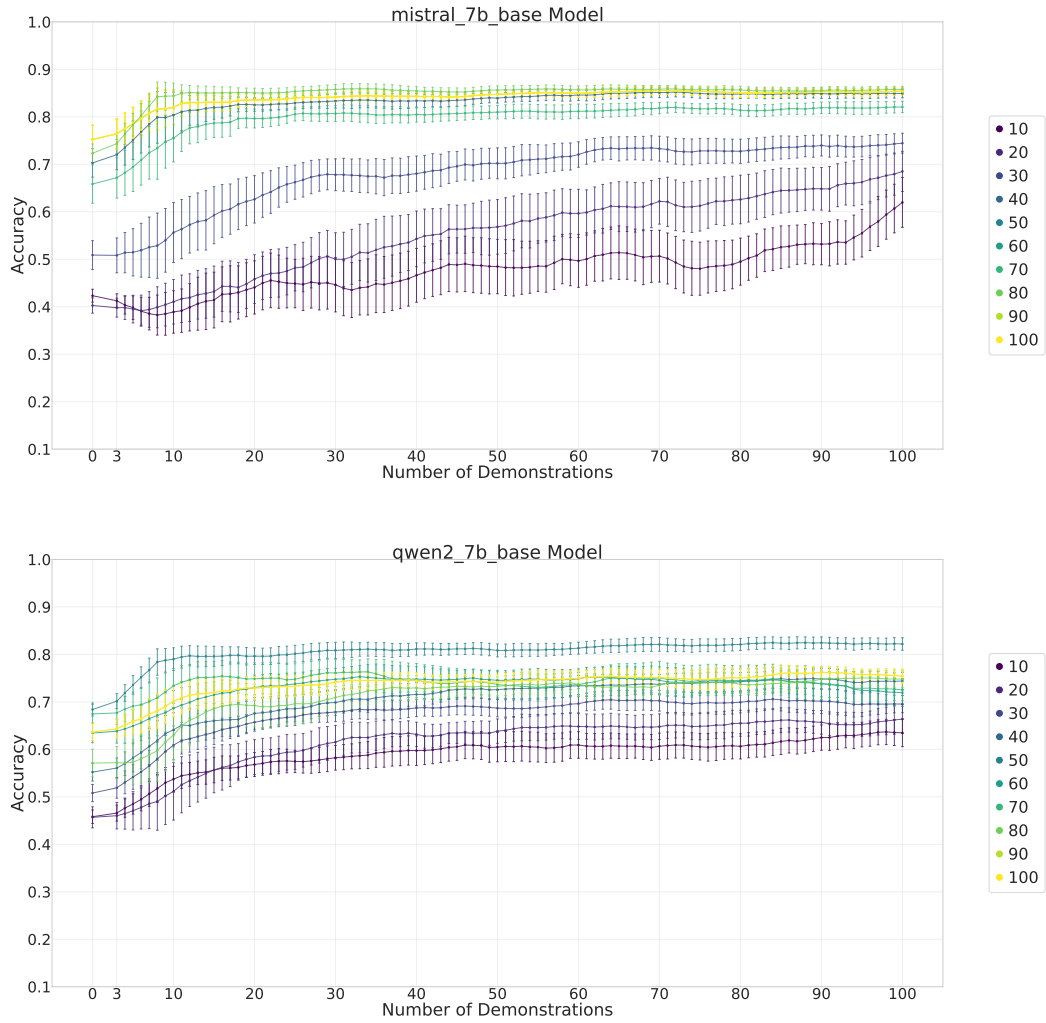


Figure 7. Mistral and Qwen 3-way classification for the sentiment analysis task

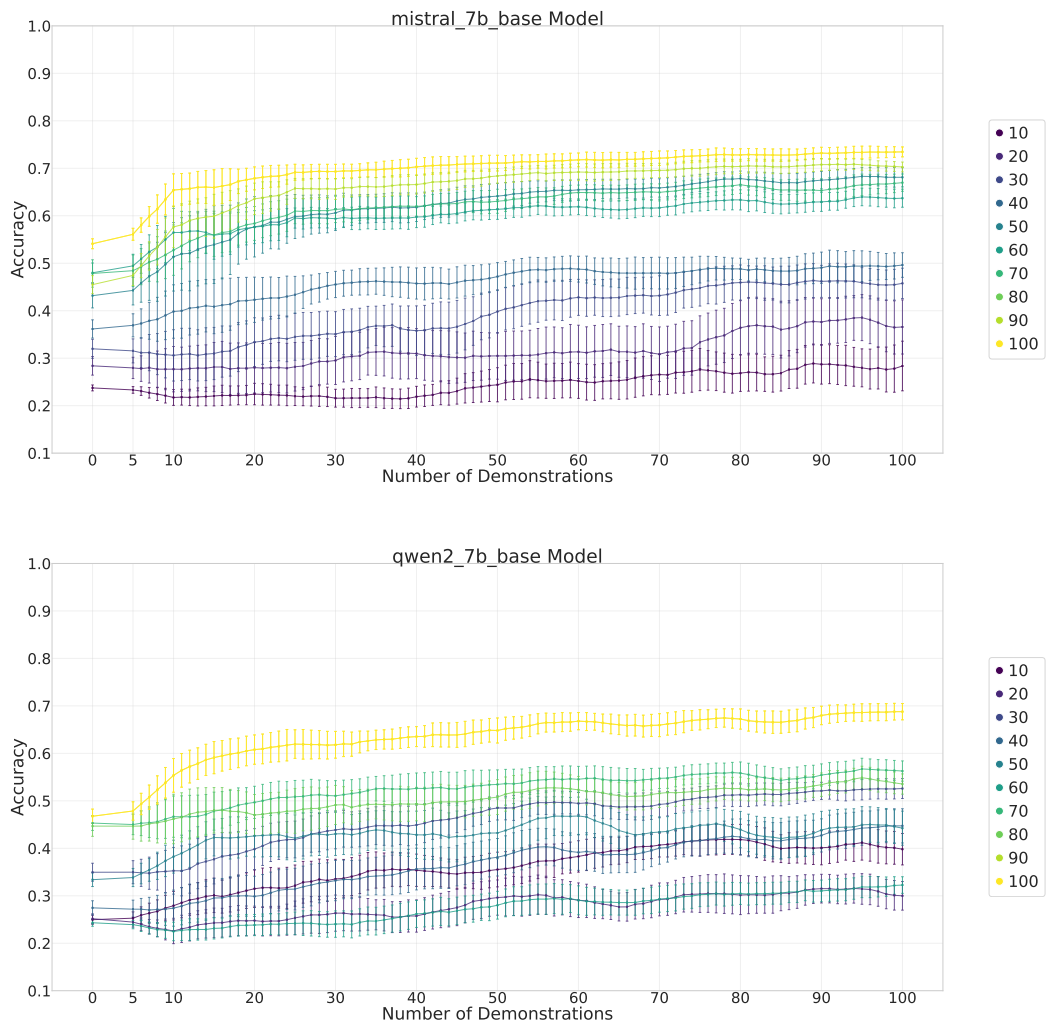


Figure 8. Mistral and Qwen 5-way classification for the sentiment analysis task

**D. Additional dataset: TREC**

We conducted additional experiments using the TREC dataset (Li & Roth, 2002; Hovy et al., 2001) for question classification. The gold labels are “Entity, Description, Human, Location, Numeric.” The input sentences in this dataset are questions. This makes applying our framework to this task challenging since high probability next tokens following a question would be answers rather than class names. Despite this fact, our findings from sentiment analysis hold true on the TREC dataset.

**D.1. Llama 3.1 8B**

K	Llama 8B
10	gods, Derm, Philippe, Helsinki, Wade
20	Judaism, Skin, Christians, Ukraine, Watts
30	easiest, Carb, Isaiah, Antarctica, dollars
40	GENERIC, SCC, quienes, Antarctica, dollars
50	which, Explain, who, Nations, Rate
60	Taste, explain, Malcolm, maps, Timing
70	Taste, explain, qui, maps, Timing
80	taste, development, personalities, whereabouts, D°D¾D¿¿D,ÑiDµÑgÑHD²D¾
90	taste, development, quienes, destinations, timed
100	Taste, explanations, qui, locations, amount

Table 5. Label sets obtained from running Algorithm 1 on K examples for 5-way classification

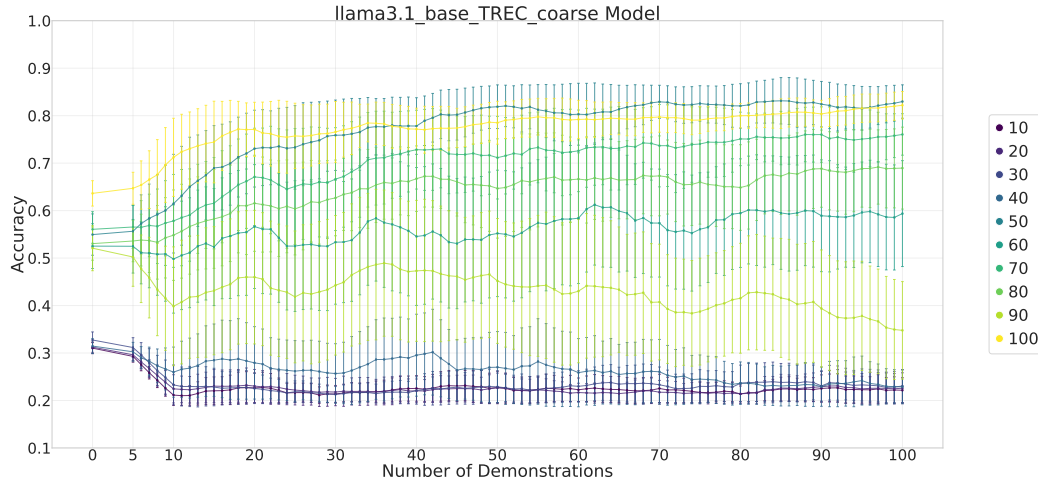


Figure 9. Llama 8B 3-way classification on question classification task

D.2. Mistral 7B

K	Mistral 7B
10	universe, Fer, artists, Tokyo, gover
20	Jung, advice, Leonard, Finland, college
30	matching, explaining, kings, continent, accounting
40	filling, explaining, athletes, UEFA, amount
50	night, explan, whom, UEFA, amount
60	conj, explan, whom, locations, amount
70	eating, explan, whom, locations, amount
80	drinking, explan, whom, locations, amount
90	eating, explan, whom, locations, amount
100	cul, explan, whom, locations, amount

Table 6. Label sets obtained from running Algorithm 1 on  $K$  examples for 5-way classification

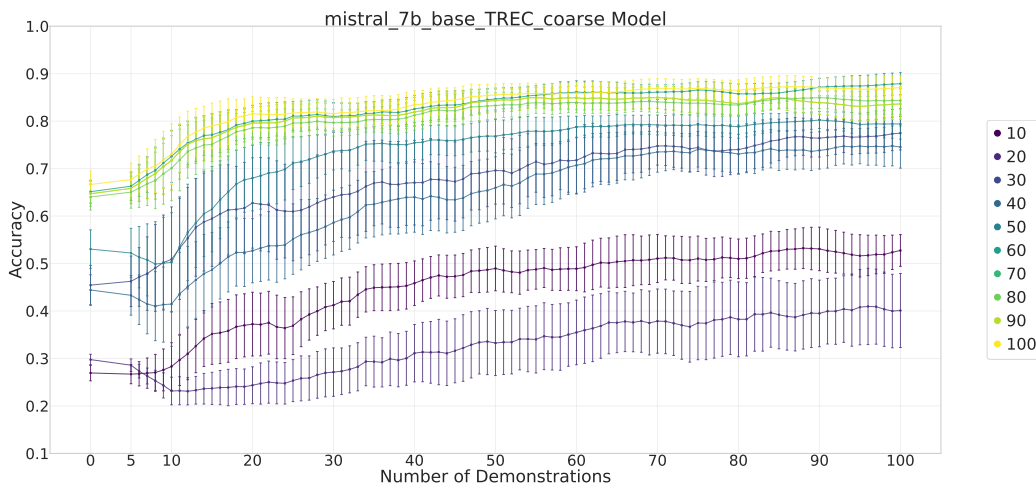


Figure 10. Mistral 7B 3-way classification on question classification task

D.3. Qwen 7B

<b>K</b>	<b>Qwen 7B</b>
10	—, Bronx, Kaz, PhoneNumber, PRES
20	comparator, Chung, Boris, wherever, æķ°
30	ConfigurationManager, Geh, quoted, wherever, amounts
40	meinen, NDEBUG, Santa, wherever, amounts
50	me, Technologies, Vi, geographical, Num
60	animal, Anatomy, Vi, geographical, Num
70	animal, Lecture, Vi, geographical, Num
80	mt, Lecture, Vi, WHERE, Num
90	serif, ×ç, entreprise, gĀ©, numb
100	serif, ×ç, entreprise, gĀ©, numb

Table 7. Label sets obtained from running Algorithm 1 on  $K$  examples for 5-way classification

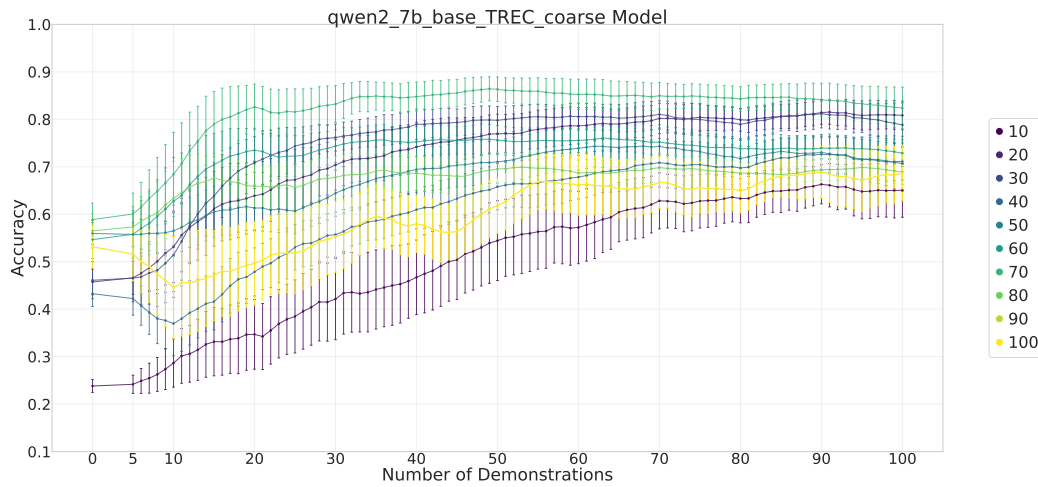


Figure 11. Qwen 7B 5-way classification on question classification task

### E. Transfer of labels across models

In order to strictly isolate the effect of model size from the choice of representation, we evaluated different models on the exact same label sets. We fixed the optimized label sets discovered by the Llama 1B model and the Llama 70B model, and then evaluated both of these sets across all three model sizes (1B, 8B, and 70B) on the 3-way and 5-way sentiment classification tasks. We show the results with the labels found by Llama 70B in Figure 12 for 3-way and in Figure 13 for 5-way classification. Figure 14 shows the results with the labels from Llama 1B on 3-way classification. By holding the representation constant, we can now clearly observe the isolated effect of model scale. We found that (1) when prompted with the exact same label set, the larger 70B model exhibits a steeper learning curve than the 1B and 8B models and (2) the effect holds regardless of which model generated the labels.

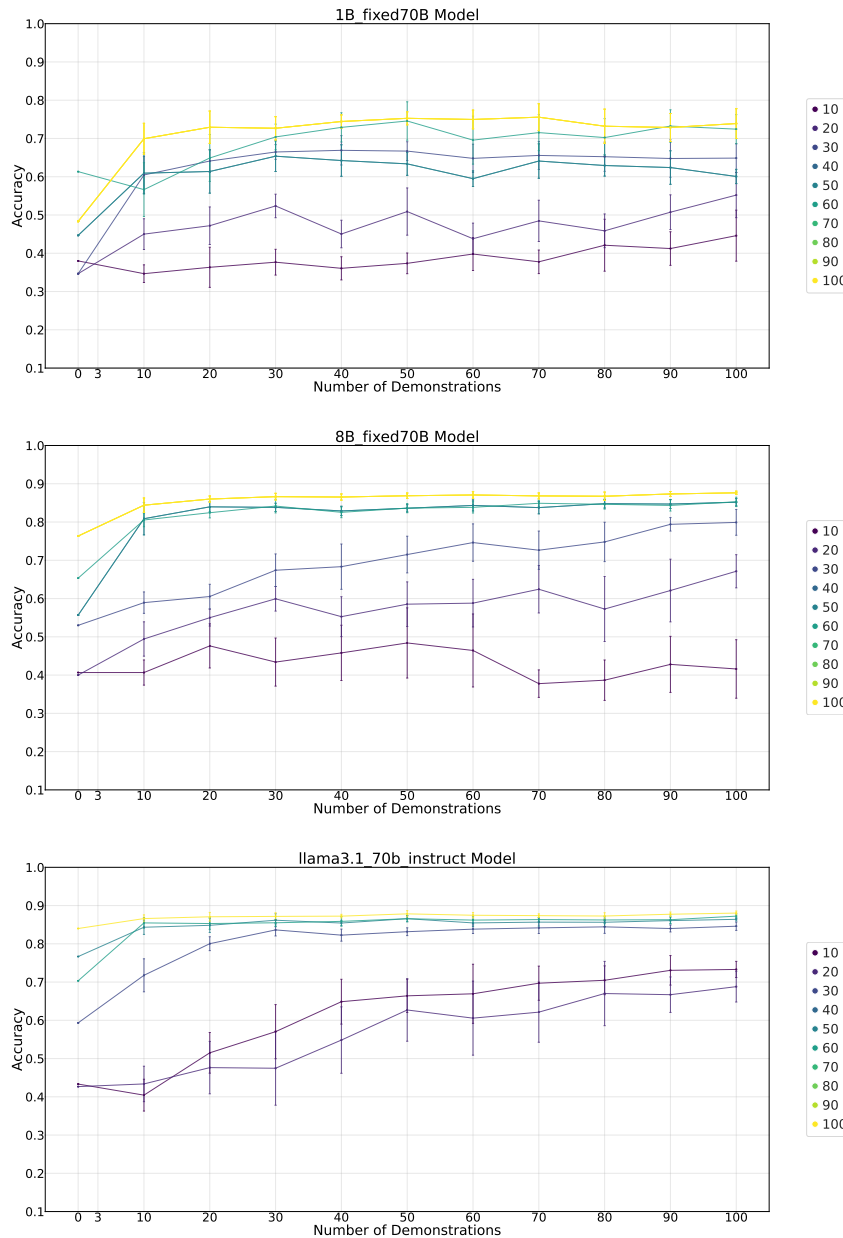


Figure 12. Learning curves for up to 100 demonstrations for Llama models for 3-way classification for the sentiment analysis task using label sets found with Llama 70B.

## Relationship between Representation and In-context learning

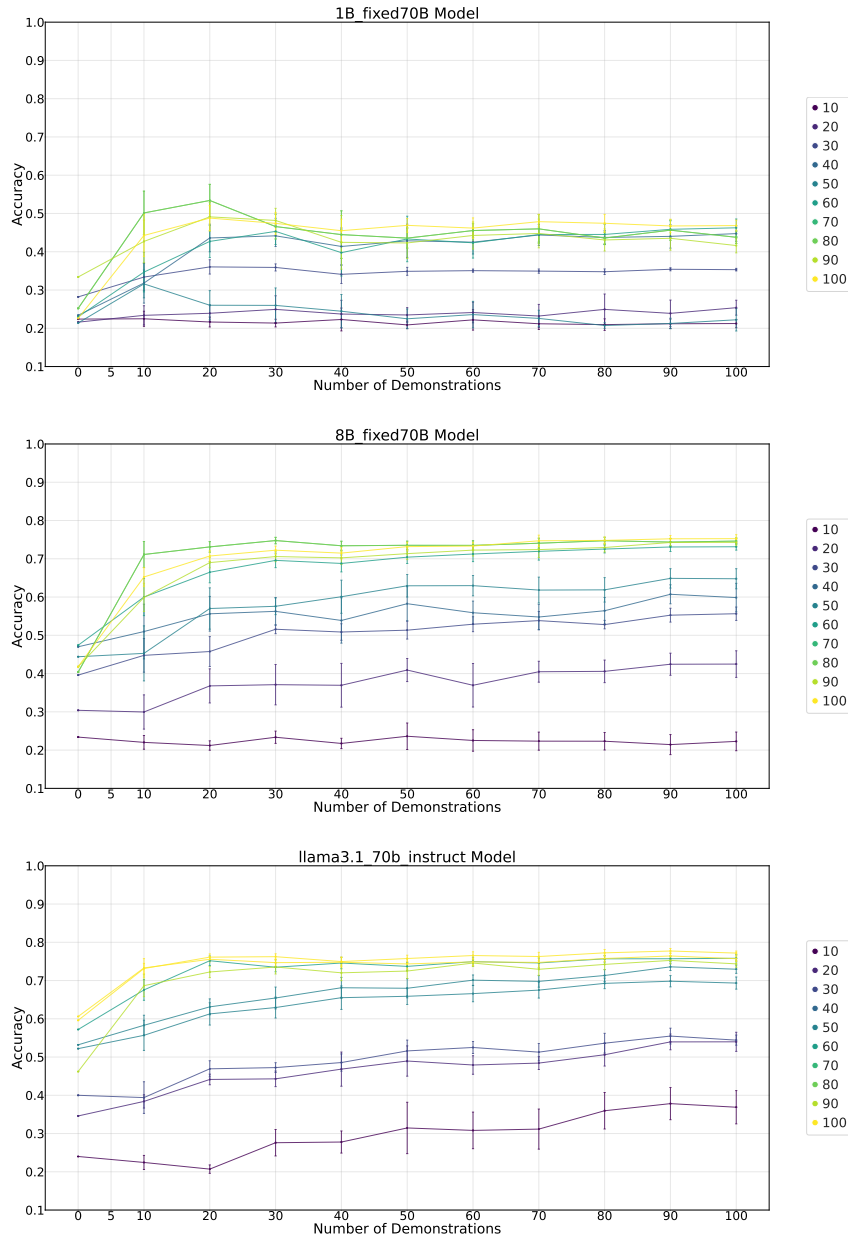


Figure 13. Learning curves for up to 100 demonstrations for Llama models for 5-way classification for the sentiment analysis task using label sets found with Llama 70B.

## Relationship between Representation and In-context learning

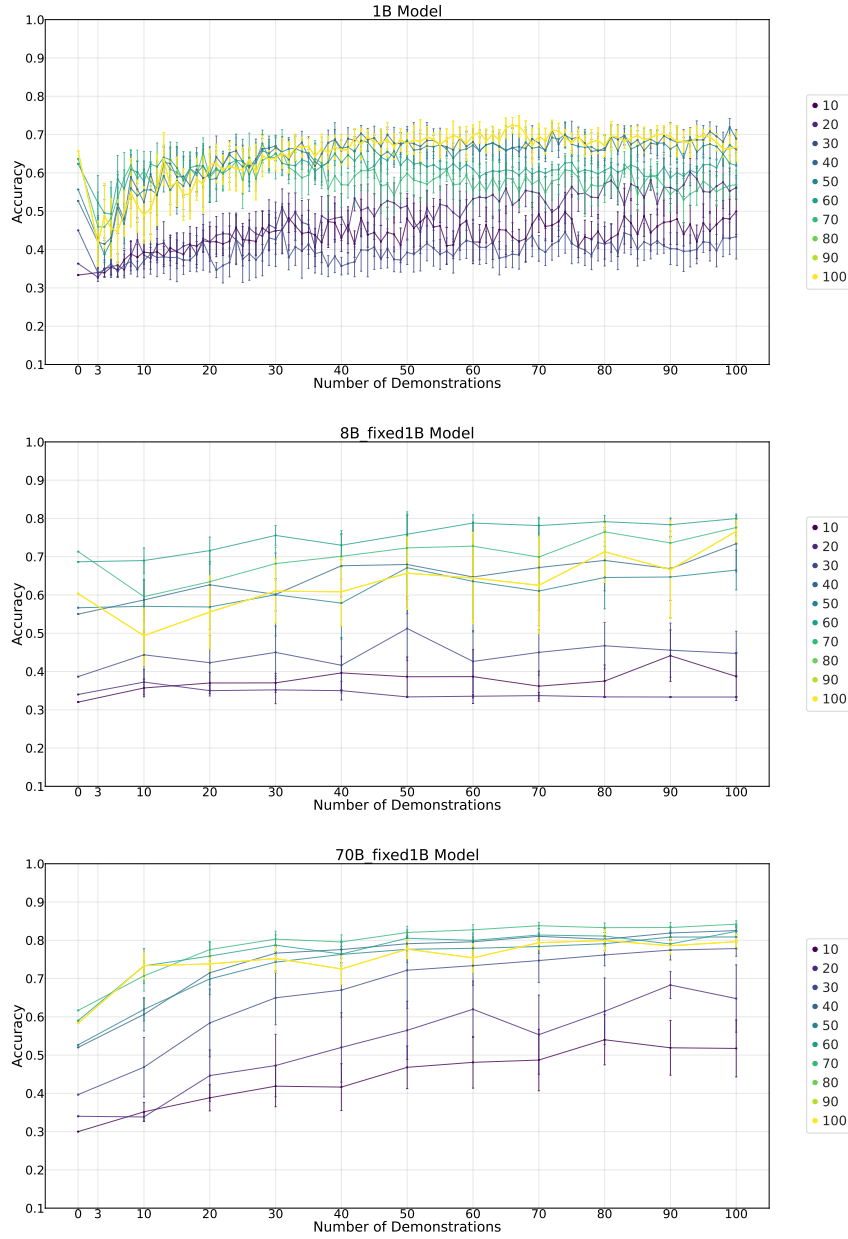


Figure 14. Learning curves for up to 100 demonstrations for Llama models for 3-way classification for the sentiment analysis task using label sets found with Llama 1B.

**F. Effect of prompt format**

In addition to our original prompt format which is shown in Figure 1 (Text/Category), we tried two other prompts: Sentence/Label and Arrow (Input → Output) with Qwen 7B. The results are comparable, as illustrated in Figure 15. showing rank preservation and confirming robustness to minor prompt changes.

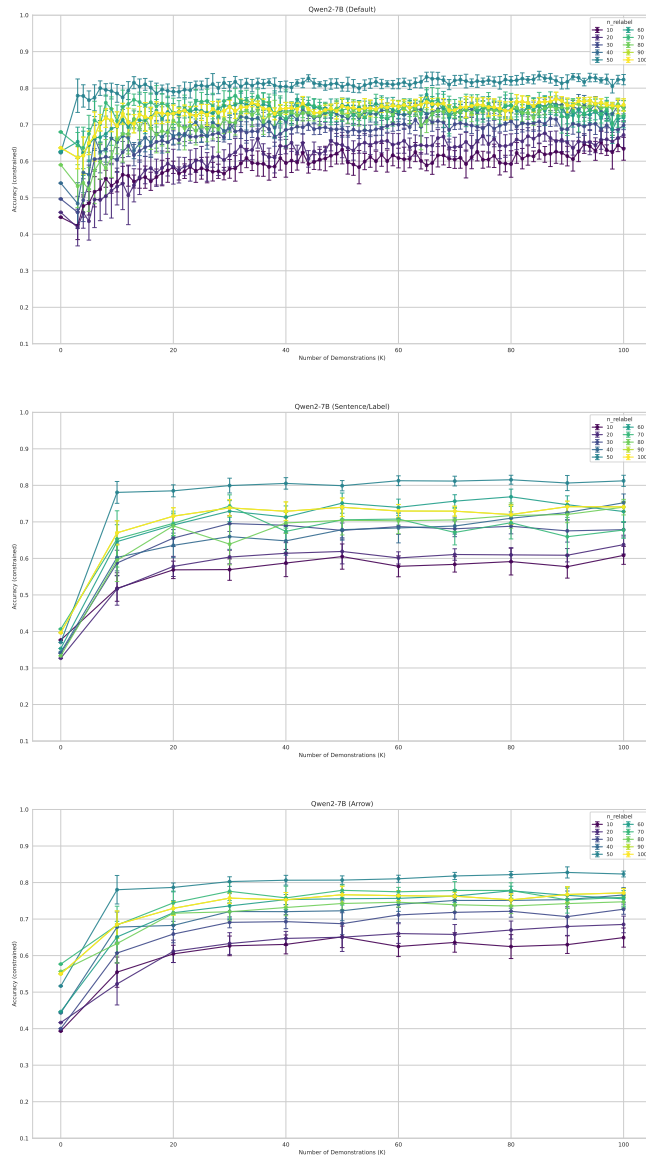


Figure 15. Learning curves for up to 100 demonstrations with Qwen 7B for 3-way classification for the sentiment analysis with three different prompt formats.

## G. Implementation details

### G.1. Algorithm 1

We first precompute the next token logits for all inputs with only one forward pass through the relabeling sentences in the dataset. During optimization, for each  $\tau$  we perform a simple lookup for the logits of the  $K$  sentences corresponding to the tokens in  $\tau$  and normalize. The algorithm finds “good enough” solutions in the sense that it allows us to enumerate labels sets covering a wide spectrum zero-shot accuracies. Moreover, the label sets found with large  $K$  are close (or identical) in meaning with the gold labels. These label sets already have a zero-shot accuracy close (within 3-5%) of the ceiling accuracy obtained after seeing demonstrations, further confirming their quality. While our algorithm is not guaranteed to find the global optimum, it is extremely simple and efficient, making it suitable for studying ICL from demonstrations or for applications which require high-quality labels.

### G.2. In-context learning

For ICL, we sample 10 different sets of  $N$  demonstrations and query the entire test set. We use prefix KV caching for demonstrations, so that during testing we only need to compute the attention between the query sentence and the demonstrations.

**Runtime.** We report approximate wall clock times for any fixed label set with  $N$  demonstrations for 300-500 test sentences, for 10 runs on an 80GB A100 or H100 GPU: 5-7 minutes for a 7B or 8B model and 20-24h for the 70B model.