

# When Parts are Greater Than Sums: Individual LLM Components Can Outperform Full Models

Anonymous ACL submission

## Abstract

This paper studies in-context learning (ICL) by decomposing the output of large language models into the individual contributions of attention heads and MLPs (*components*). We observe curious components: good-performing ones that individually do well on a classification task, even when the model performs poorly; bad-performing ones that do much worse than chance; and label-biased components that always predict the same label. We find that component accuracies are well-correlated across different demonstration sets and perturbations of prompt templates, even when the full-model accuracy varies greatly. Based on our findings, we propose component reweighting, which learns to linearly re-scale the component activations from a few labeled examples. Given 24 labeled examples, our method improves by an average of 6.0% accuracy points over 24-shot ICL across 8 tasks on Llama-2-7B. Overall, this paper both enriches our understanding of ICL and provides a practical method for improvement by examining model internals.

## 1 Introduction

Rapid changes in large language models (LLMs) have popularized prompting, which guides LLMs to perform tasks with instructions or examples. Notably, in-context learning (ICL; Brown et al., 2020) adapts LLMs to a new task using only a few labeled examples without parameter updates. However, how LLMs react to the in-context examples is sometimes unintuitive (Min et al., 2022). Recently, Sclar et al. (2024) find that even for LLMs with instruction tuning (Ouyang et al., 2022) and large size, adding a space or newline in prompts can greatly affect accuracy.

We look into the LLM internals to understand what causes the surprising behavior across various ICL settings. Our work stands in contrast to prior studies on ICL, which treat LLMs as black boxes and alter either the input (Liu et al., 2022; Bertsch

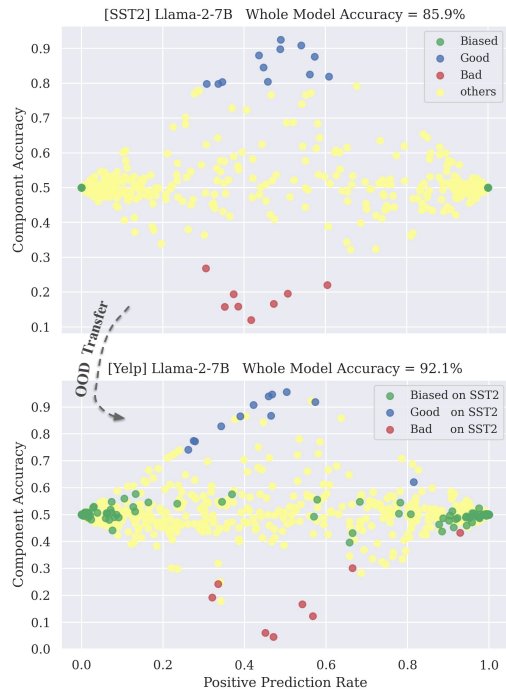


Figure 1: Each dot represents a component (attention head or MLP) under 4-shot ICL on Llama-2-7B. The x-axis shows how often a component predicts “positive” on the test set. **Up:** We discover good-performing (blue), bad-performing (red), and label-biased (green) components. **Down:** Most components identified on SST2 show similar characteristics on Yelp-polarity.

et al., 2024) or output (Zhao et al., 2021; Holtzman et al., 2021). We introduce a new view of ICL by decomposing the output of an LLM into the sum of individual contributions of MLPs and attention heads, denoted “components.” Figure 1 reveals three types of curious components: good-performing (blue) that individually perform well or even outperform the full model, bad-performing (red) that perform below chance, and label-biased (green) that predict the same label on the entire test set. We observe these three types on Llama-2-7B, Llama-2-13B (Touvron et al., 2023), and Mistral-Instruct-7B (Jiang et al., 2023) across 8 tasks.

We study the sensitivity of LLM components to multiple prompts, formed by different demonstrations and templates. We also construct contrast sets of templates—pairs of similar templates that yield large differences in ICL accuracy. Despite large variance in full-model accuracy, we find that component accuracies correlate well across different demonstrations ( $r = 0.79$  on average) and contrast set templates ( $r = 0.57$ ). The top-performing components in contrast set pairs overlap and achieve decent accuracy even when the full model performs near random (Figure 2). In contrast, the component accuracies of two sampled templates are less correlated ( $r = 0.34$ ). Further, good-performing components generalize well to out-of-distribution test sets. For instance, the top-1 component for MNLI outperforms the full Llama-2-13B model by 9.1% on MedNLI; Figure 1 also shows that components are transferrable from SST2 to Yelp. We conclude that components are relatively consistent in their behavior across prompts and datasets.

Inspired by our findings, we propose component reweighting. Compared to prior work that selects prompts to improve ICL accuracy (Rubin et al., 2022), component reweighting softly selects components by learning weights from the same few-shot examples to scale component activations. Training these weights only involves learning a linear layer—less than a minute on one CPU. Overall, component reweighting better utilizes the same labeled examples, improving over 24-shot ICL by 6.0%, 2.2%, 5.1% on Llama-2-7B, Llama-2-13B, and Mistral-Instruct-7B, respectively, while enjoying similar inference speed as 4-shot ICL.

Finally, we study the training dynamics of components using the Pythia pretraining checkpoints (Biderman et al., 2023). During pretraining, good-performing components emerge well before the full model performs well. These findings suggest that LLMs acquire the internal ability to perform ICL early in training, but this ability only surfaces in the full model’s behavior later on. Our work conducts extensive analysis of LLM internals, which motivates a practical method to improve ICL; we hope to inspire future work that further sheds light on LLM internals in order to improve performance.

## 2 Decomposing the Transformer in ICL

We introduce a new view of in-context learning by decomposing the Transformer architecture (Vaswani et al., 2017). Our decomposition is

exact—a mathematically equivalent formula for the model’s outputs—and enables us to analyze model internals without training additional parameters (unlike, e.g., probing). We first discuss what our new view offers over the standard view of ICL, and then walk through the mathematical details.

### 2.1 A New View of In-Context Learning

**Standard view.** An LLM performs in-context learning (ICL) on a task based on a few demonstrations without training, where each demonstration is a templated example  $(x, y)$  consisting of an input  $x$  and a label word  $y$ . We refer to a sequence of  $K$  demonstrations  $[x_1, y_1, \dots, x_K, y_K]$  as a *prompt*. The LLM makes predictions on a test input  $x_{\text{test}}$  conditioned on the prompt, denoted by  $\arg \max_{y \in \mathcal{Y}} P(y | \text{prompt}, x_{\text{test}})$ , where  $\mathcal{Y}$  is the set of possible label words in a classification task.

**Our view.** The residual stream of an LLM directly carries the information of the initial hidden state, every attention head, and every MLP, collectively named “components,” towards the output layer. We view this information as the direct contributions<sup>1</sup> of components to the output logits, and derive a formula for logits,  $\sum_j \mathbf{g}_j$ , where  $\mathbf{g}_j$  is the direct contribution of the component indexed by  $j$ . We can obtain the predictions of component  $j$  with  $\arg \max_{y \in \mathcal{Y}} \mathbf{g}_j$ , and then calculate its individual ICL accuracy. Specifically, we derive  $\mathbf{g}_j = U \cdot C_j$  in Eq. 8 below, where  $U$  is the output embedding matrix and  $C_j$  is the post-layernorm activations of component  $j$ . We name the operation  $(C_j \mapsto U \cdot C_j)$  as early decode, sharing the same spirit as *nostalgebraist* (2020); Geva et al. (2022), which interpret hidden representations by decoding through  $U$ . Compared to the standard view, we can directly study the behavior of individual components (Figure 2), characterizing them and scaling their contributions to the model output.

### 2.2 A Walkthrough of the Decomposition

A Transformer of  $L$  layers consists of a multi-headed attention (MHA) and MLP in every layer. Let  $a^{(l)} \in \mathbb{R}^d$  and  $m^{(l)} \in \mathbb{R}^d$  be the output of the MHA and MLP at layer  $l$ , respectively. Due to

<sup>1</sup>In comparison, a component has indirect contributions to the output by affecting other components in later layers.

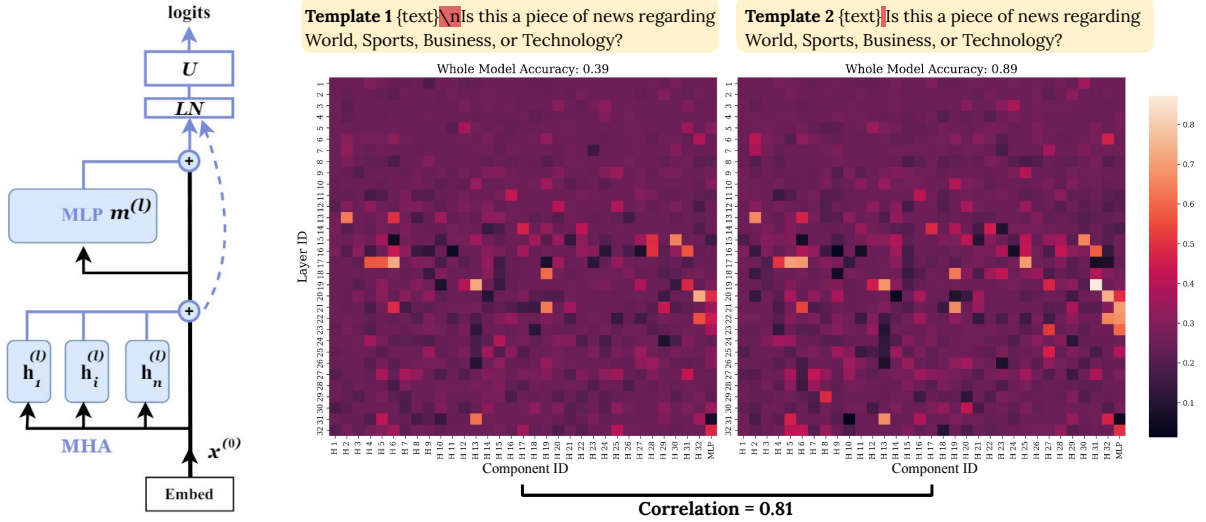


Figure 2: **Left:** Transformer decomposition. The components, MLPs and attention heads, are filled with blue and the blue lines show the flow of early decoding. **Right:** We can calculate the individual accuracy of every component after decomposition. Although a pair of templates that only differ slightly yield contrasting accuracies (0.39 vs. 0.89 on AGNews with Llama-2-7B), the accuracies of their internal components are highly correlated. The top components for Template 1 overlap with the ones for Template 2 and achieve  $> 0.7$  accuracy despite the poor full-model accuracy.

residual connections, the hidden state  $x^{(l)} \in \mathbb{R}^d$  is:

$$x^{(l)} = x^{(l-1)} + a^{(l)} + m^{(l)}, \quad (1)$$

$$x^{(L)} = x^{(0)} + \sum_{l=1}^L (a^{(l)} + m^{(l)}). \quad (2)$$

Note that GPT2-like LLMs apply layernorm before MHA and MLP (Radford et al., 2019); thus, layernorm is already taken into account as part of the formula for computing  $a^{(l)}$  and  $m^{(l)}$  (see A.3).

An MHA  $a^{(l)}$  is composed of  $n$  attention heads:

$$a^{(l)} = W_o^{(l)} \cdot \text{Concat}([h_1^{(l)}, \dots, h_n^{(l)}]) \quad (3)$$

for  $h_i^{(l)} \in \mathbb{R}^{d_{\text{head}}}$  a head and  $W_o^{(l)} \in \mathbb{R}^{d \times n d_{\text{head}}}$  the output projection in MHA aggregating all heads. Elhage et al. (2021) rewrite Eq. 3 by segmenting  $W_o^{(l)}$  into  $n$  matrices  $W_{o_i}^{(l)} \in \mathbb{R}^{d \times d_{\text{head}}}$ :

$$a^{(l)} = \sum_{i=1}^n (W_{o_i}^{(l)} \cdot h_i^{(l)}) = \sum_{i=1}^n \tilde{h}_i^{(l)}, \quad (4)$$

$$\text{where } [W_{o_1}^{(l)}, \dots, W_{o_n}^{(l)}] = W_o^{(l)} \quad (5)$$

Thus, we can treat each head as a single component adding  $\tilde{h}_i^{(l)} = W_{o_i}^{(l)} \cdot h_i^{(l)}$  to the residual stream.

Finally, through the output embedding matrix

$U \in \mathbb{R}^{|\text{Vocab}| \times d}$ , the output logits are:

$$\text{logits} = U \cdot \text{LN}(x^{(L)})$$

$$= U \cdot \text{LN} \left( x^{(0)} + \sum_{l=1}^L \sum_{i=1}^n \tilde{h}_i^{(l)} + \sum_{l=1}^L m^{(l)} \right)$$

$$= U \cdot \text{LN} \left( \sum_{j=1}^{1+L \times n + L} z_j \right), \quad (6)$$

where  $z = [x^{(0)}, \mathbf{h}_1^{(1)}, \dots, \mathbf{h}_n^{(L)}, m^{(1)}, \dots, m^{(L)}]$  in Eq. 6 and we index every term in the summation with  $j$ .  $\text{LN}(\cdot)$  denotes the final layernorm, specifically, RMSNorm (Zhang and Sennrich, 2019) for LLMs in our paper (see A.3). In Eq. 6,  $\text{LN}(\sum_j z_j) = \frac{\sum_j z_j}{\text{RMS}(\sum_j z_j)} \odot \gamma$ , where RMS denotes root mean square,  $\odot$  denotes element-wise multiplication, and  $\gamma \in \mathbb{R}^d$  is the affine parameters. By pre-computing  $\hat{\gamma} = \frac{\gamma}{\text{RMS}(\sum_j z_j)}$ , we have:

$$\text{logits} = U \cdot \left( \sum_j z_j \odot \hat{\gamma} \right) \quad (7)$$

$$= \sum_j U \cdot C_j, \text{ where } C_j = z_j \odot \hat{\gamma} \quad (8)$$

We refer to all  $C_j \in \mathbb{R}^d$  as the component activations, which include the activations of attention heads and MLPs after the final layernorm.<sup>2</sup> Now

<sup>2</sup>Empirically, we find that  $x^{(0)}$  has near-random ICL accuracy on all the tasks, so we omit it in the rest of the paper.

		SST2	BoolQ	QQP	WiC	RTE	MNLI	AGNews	ARC-Easy	Avg.
Llama2 7B	FULL	75.8 <sub>18.1</sub>	69.2 <sub>12.0</sub>	61.3 <sub>9.9</sub>	52.4 <sub>3.0</sub>	<b>68.9</b> <sub>3.2</sub>	34.4 <sub>1.7</sub>	70.0 <sub>19.9</sub>	<b>57.5</b> <sub>14.4</sub>	61.2
	ORACLE-T1	<b>91.7</b> <sub>0.9</sub>	<b>69.7</b> <sub>7.7</sub>	<b>67.8</b> <sub>4.3</sub>	<b>57.8</b> <sub>1.1</sub>	64.6 <sub>2.7</sub>	<b>46.3</b> <sub>3.3</sub>	<b>80.8</b> <sub>5.2</sub>	54.5 <sub>10.1</sub>	<b>66.6</b>
	ORACLE-B1	12.1 <sub>2.7</sub>	34.1 <sub>7.3</sub>	32.5 <sub>3.9</sub>	42.9 <sub>1.2</sub>	34.7 <sub>2.8</sub>	24.1 <sub>2.4</sub>	3.0 <sub>1.1</sub>	12.7 <sub>4.2</sub>	24.5
Llama2 13B	FULL	89.0 <sub>5.3</sub>	<b>77.6</b> <sub>6.8</sub>	71.0 <sub>6.8</sub>	55.0 <sub>3.8</sub>	75.1 <sub>2.3</sub>	45.7 <sub>7.9</sub>	70.8 <sub>20.6</sub>	<b>73.2</b> <sub>13.7</sub>	69.7
	ORACLE-T1	<b>92.5</b> <sub>0.6</sub>	77.5 <sub>6.0</sub>	<b>73.5</b> <sub>2.9</sub>	<b>60.4</b> <sub>1.2</sub>	<b>75.7</b> <sub>2.3</sub>	<b>56.4</b> <sub>4.7</sub>	<b>84.6</b> <sub>3.6</sub>	73.1 <sub>7.9</sub>	<b>74.2</b>
	ORACLE-B1	8.2 <sub>1.0</sub>	27.1 <sub>9.7</sub>	31.8 <sub>3.4</sub>	39.5 <sub>1.6</sub>	27.9 <sub>2.8</sub>	18.6 <sub>2.6</sub>	1.8 <sub>0.9</sub>	5.4 <sub>3.5</sub>	20.0
Mistral 7B	FULL	90.1 <sub>2.9</sub>	<b>81.3</b> <sub>2.1</sub>	70.9 <sub>7.2</sub>	58.5 <sub>4.2</sub>	80.5 <sub>1.7</sub>	56.1 <sub>5.0</sub>	83.0 <sub>5.7</sub>	<b>79.8</b> <sub>1.4</sub>	75.0
	ORACLE-T1	<b>91.9</b> <sub>0.7</sub>	80.8 <sub>2.0</sub>	<b>75.6</b> <sub>2.6</sub>	<b>60.6</b> <sub>2.2</sub>	<b>81.3</b> <sub>0.8</sub>	<b>61.5</b> <sub>3.3</sub>	<b>83.7</b> <sub>4.3</sub>	78.5 <sub>2.2</sub>	<b>76.7</b>
	ORACLE-B1	8.1 <sub>0.9</sub>	19.5 <sub>2.5</sub>	25.8 <sub>4.1</sub>	39.3 <sub>2.8</sub>	20.0 <sub>1.7</sub>	14.6 <sub>2.9</sub>	1.8 <sub>0.7</sub>	4.6 <sub>1.3</sub>	16.7
RANDOM		50.0	50.0	50.0	50.0	50.0	33.3	25.0	25.0	41.7

Table 1: {3, 4}-shot ICL accuracy of 8 tasks and the average accuracy (**Avg.**). We run 15 prompts for each task (see §3) and report the mean accuracy and standard deviation. We show the existence of good components (ORACLE-T1) inside LLMs that individually perform on par with the full model (FULL) on diverse tasks. Similarly, there exist bad components (ORACLE-B1) that perform substantially below chance (RANDOM).

that we have broken down the Transformer output into simple additions in Eq. 8, we can easily analyze the direct contribution of each component to the logits through the residual stream,  $\mathbf{g}_j = U \cdot C_j$ .

In ICL, we only need to do the decomposition when LLMs start to generate, i.e, when processing the last token of the input. The computations on the other tokens are the same as the standard ICL. In all our experiments, we use single-token label words. We use multiple templates from Bach et al. (2022) that cover diverse label words for each task.

### 3 Characterizing Components for ICL

We conduct in-context learning across 8 classification tasks on 3 LLMs: Llama-2-7B, Llama-2-13B, and Mistral-Instruct-7B. ICL is sensitive to prompts, so we randomly sample 5 disjoint sets of demonstrations formatted with 3 templates and report the standard deviation across the 15 runs. To avoid majority and recency biases (Zhao et al., 2021), each prompt consists of the same number of demonstrations from every class in shuffled order. We use  $K = 3$  demonstrations for 3-way classification tasks and  $K = 4$  for the other tasks. Except for §5.1, we refer to  $K = \{3, 4\}$  without further notice. We sample 2000 examples with balanced labels as the test set for every task. Please see A.1 for details about the tasks and templates.

#### 3.1 Good and Bad-Performing Components

Across all the tasks and LLMs, we observe good-performing components that perform well or even outperform the full model and bad-performing components, which individually perform much worse than chance (blue and red dots in Figure 1,

respectively). Table 1 compares the full model (FULL) with the top-1 (ORACLE-T1) and bottom-1 (ORACLE-B1) components selected on the test set. On average, ORACLE-T1 outperforms FULL by 5.4%, 4.5%, 1.7% on Llama-2-7B, Llama-2-13B, and Mistral-Instruct-7B, respectively; ORACLE-B1 underperforms random guessing (RANDOM) by 17.2%, 20.7%, 25.0%, respectively.

#### 3.2 Label-Biased Components

Besides good and bad-performing components, we also observe label-biased components, which predict a certain label on the entire test set (the green dots in Figure 1). These components exist in all the tasks and LLMs we study, accounting for 29.1%, 26.4%, 22.8% of components on average in Llama-2-7B, Llama-2-13B, and Mistral-Instruct-7B, respectively (see Table 6). When we prompt the model with all demonstrations of positive labels, a biased component still insists on predicting “negative” on the entire test set, and vice versa (A.2).

### 4 Transferability of Components

We observe moderate to high component transferability across demonstrations, minimally contrastive templates, and data distributions, whereas there is little transferability across randomly sampled templates. Our decomposition uncovers hidden abilities of individual components when the full model performs poorly.

#### 4.1 Transfer across Prompt Variants

We first measure the agreement in component accuracies between (1) two disjoint sets of demonstrations with a fixed template, (2) two randomly sampled templates with fixed demonstrations, and

		SST2	BoolQ	QQP	AGNews	ARC
Corr	(1) Demo	0.81	0.84	0.60	0.89	0.88
	(2) Temp	0.40	0.16	0.03	0.68	0.44
	(3) Cst T	0.72	0.63	0.23	0.82	0.46
IoU	(1) Demo	0.36	0.74	0.27	0.63	0.70
	(2) Temp	0.12	0.01	0.01	0.20	0.20
	(3) Cst T	0.40	0.23	0.02	0.36	0.45

Table 2: The average correlation and IoU between (1) two random sets of demonstrations, (2) two random templates, and (3) two minimally contrastive templates.

(3) two minimally-contrastive templates with fixed demonstrations. Recall that we have 5 sets of demonstrations and 3 templates in total (§3); here, we calculate the average agreement between every pair. For (3), we construct contrast sets (Gardner et al., 2020) by minimally editing the worst-performing template out of the 3 templates into a good template, which yields at least 10% improvement in average accuracy. Our edits include adding a space, removing a newline, or changing label words (see Table 8). We use two metrics to measure the agreement between each pair: Pearson correlation of the accuracies of all components and the intersection over union (IoU) on the sets of top-5 components, which measures whether the top-performing components of the pair overlap.

Table 2 summarizes the results on Llama-2-7B. A.6 shows similar findings on other models. (1) The accuracies of the internal components are highly consistent across different choices of demonstrations, having strong correlations and an average of 0.54 IoU. (2) The components have much weaker agreement across randomly sampled templates, having a near 0 IoU on BoolQ and QQP. (3) Nevertheless, there is agreement between minimally contrastive templates (Cst T), with an average correlation of 0.572 across tasks, despite contrasting full-model accuracy. For example, Figure 2 demonstrates that full-model accuracy changes dramatically (39% vs 89%) in a minimal pair of templates, but internal components have a high correlation of 0.81 and the pair shares top-performing components. Combining (2) and (3) suggests components behave similarly on similar templates, but this similarity decreases as the templates diverge.

## 4.2 Transfer to Out-of-Distribution Test Sets

We further study whether the best component selected on the test set can still perform well on an out-of-distribution (OOD) test set. We name

		Yelp-polarity	MedNLI	BoolQ Cst
Llama2 7B	FULL	84.7 <sub>15.4</sub>	34.3 <sub>1.7</sub>	<b>64.9</b> <sub>9.8</sub>
	TRANSFER-1	<b>94.9</b> <sub>3.1</sub>	<b>42.6</b> <sub>4.7</sub>	64.3 <sub>7.9</sub>
	ORACLE-1	96.9 <sub>0.7</sub>	48.8 <sub>2.3</sub>	66.2 <sub>5.7</sub>
Llama2 13B	FULL	95.9 <sub>1.4</sub>	46.8 <sub>9.6</sub>	72.0 <sub>7.6</sub>
	TRANSFER-1	<b>96.0</b> <sub>1.8</sub>	<b>55.9</b> <sub>4.0</sub>	<b>72.3</b> <sub>6.5</sub>
	ORACLE-1	97.1 <sub>0.4</sub>	57.0 <sub>3.7</sub>	73.0 <sub>6.1</sub>
Mistral 7B	FULL	<b>97.0</b> <sub>0.5</sub>	57.3 <sub>5.7</sub>	<b>74.6</b> <sub>3.5</sub>
	TRANSFER-1	95.6 <sub>1.6</sub>	<b>61.9</b> <sub>4.8</sub>	73.7 <sub>3.7</sub>
	ORACLE-1	97.1 <sub>0.4</sub>	62.7 <sub>4.1</sub>	74.5 <sub>3.6</sub>

Table 3: The average ICL accuracy and standard deviation on OOD test sets. The components selected on the in-distribution test sets (TRANSFER-1) can transfer to OOD sets, performing similarly to the oracle components (ORACLE-1) directly selected on the OOD sets.

this method, which uses a single component to make predictions, as TRANSFER-1. Specifically, we study component transferability from SST2 to Yelp-polarity, MNLi to MedNLI, and BoolQ to BoolQ Contrast Set. We compare TRANSFER-1 with using the full model (FULL) on the OOD test sets. To understand the best possible TRANSFER-1 accuracy, we also report the best component accuracy directly selected on the OOD set, ORACLE-1.

Table 3 shows that TRANSFER-1 performs very close to ORACLE-1 overall, suggesting that the top-performing components are consistent over data distribution. Moreover, TRANSFER-1 sometimes substantially outperforms FULL, especially on MedNLI (random: 33.3%), suggesting that when LLMs perform poorly in ICL, it is likely that some internal components can do much better.

## 4.3 Transfer between Two Opposite Tasks

We conduct a case study of component transferability across instructions using Task069 and Task070 of Super-NaturalInstructions (Wang et al., 2022). The instruction for Task 069 asks for correct answers, while Task070 asks for incorrect ones (“pick the one that makes less sense;” Figure 6). The examples for the two tasks are not parallel.

We find that Mistral-Instruct-7B achieves good accuracy across 15 runs on Task069 ( $76.8 \pm 2.4$ ), but below chance on Task070 ( $40.6 \pm 5.4$ ). We observe a strong negative correlation ( $-0.60$  on average) between accuracies of every pair of runs from the two tasks. The worst-performing components in Task069 become the top-performing in Task070 and vice versa. The correlation suggests that the model has the ability to solve Task070, but misunderstands negation. Thus, we apply the

TRANSFER-1 method (§4.2) but select the *worst*-performing component from Task069 and then calculate its individual accuracy on Task070, achieving  $58.7 \pm 4.8$  accuracy across the 15 runs, an **improvement of 18.1%** over the full model. These results suggest that components behave consistently even across tasks with opposite instructions.

## 5 Component Reweighting

### 5.1 Proposed Method

Our findings in §4 show the promising direction of selecting internal components to improve ICL. Therefore, we propose a method that reweights components by learning a weight  $w_j \in \mathbb{R}$  on every component activation  $C_j$ . Reweighting is a soft version of selection, which can be learned by gradient descent on very few examples.

Given  $K$  labeled examples, instead of using all of them as ICL demonstrations, we divide them into a demonstration set  $\mathcal{D}_{\text{demo}}$  and a training set  $\mathcal{D}_{\text{train}}$ . We first randomly sample  $K' = \{3, 4\}$  examples with balanced labels as demonstrations and use the remaining examples as  $\mathcal{D}_{\text{train}}$  to train the component weights. Specifically, we can rewrite Eq. 8 as logits =  $\sum_j w_j (U \cdot C_j)$ , where  $w_j = 1$  for all  $j$ . Because of the existence of good and bad-performing components, weighing all components equally may not be optimal. Therefore, we tune the weights  $w \in \mathbb{R}^N$  of  $N$  components on  $\mathcal{D}_{\text{train}}$  with cross-entropy loss and  $L_1$  regularization, while keeping the LLM frozen:

$$\mathcal{L} = \sum_{(x,y) \in \mathcal{D}_{\text{train}}} -\log P_{rw}(y|x) + \lambda \|w\|_1, \quad (9)$$

$$P_{rw}(y|x) = \text{softmax} \left( \sum_{j=1}^N w_j (U_{\mathcal{Y}} \cdot C_j) \right)_y,$$

where  $U_{\mathcal{Y}} \in \mathbb{R}^{|\mathcal{Y}| \times d}$  is a submatrix of  $U$  that comprises the output embeddings of label words,  $P_{rw}$  is the probability distribution of the LLM after reweighting, and  $\lambda$  is the hyperparameter of the  $L_1$  loss to encourage sparsity on the component weights. We obtain the activations  $\{C_j\}_{j=1}^N$  of all components in one  $K'$ -shot forward pass, computed on the prompt derived from  $\mathcal{D}_{\text{demo}}$ , followed by  $x$ . Our method scales each component’s direct contributions to the logits ( $U_{\mathcal{Y}} \cdot C_j \in \mathbb{R}^{|\mathcal{Y}|}$ ) by  $w_j$ . In practice, we cache these contributions on the entire training set as input features to the linear layer  $w$ , which allows us to discard the entire LLM while

---

### Algorithm 1 Component Reweighting

---

```

1: Input:  $K$  labeled examples, a test set  $\mathcal{D}_{\text{test}}$ , a set of label
   words  $\mathcal{Y}$ , an LLM  $\mathcal{M}$ , the number of components  $N$ 
2: Output:  $\mathcal{Z}$ , the predictions of  $\mathcal{M}$  on  $\mathcal{D}_{\text{test}}$ 
3: Split  $K$  examples into a prompt consists of  $K'$  demon-
   strations and a training set  $\mathcal{D}_{\text{train}}$  of  $K - K'$  examples
4:  $U_{\mathcal{Y}} \leftarrow$  concatenate the output embeddings of  $\mathcal{Y}$  in  $\mathcal{M}$ 
5: Initialize  $\mathcal{G}^{\text{dev}} \leftarrow \emptyset$ 
6: for  $(x, y) \in \mathcal{D}_{\text{train}}$  do
7:    $\{C_j\}_{j=1}^N \leftarrow \mathcal{M}(\text{prompt}, x)$  ▷  $K'$ -shot ICL
8:   for  $j \leftarrow 1$  to  $N$  do
9:      $\mathcal{G}^{\text{dev}} \leftarrow \mathcal{G}^{\text{dev}} \cup (U_{\mathcal{Y}} \cdot C_j)$  ▷ early decode
10:  end for
11: end for
12: Initialize  $w \leftarrow [1, \dots, 1] \in \mathbb{R}^N$ 
13: Train the weights  $w$  on  $\mathcal{G}^{\text{dev}}$  with Eq. 9
14: Initialize  $\mathcal{Z} \leftarrow \emptyset$  ▷ Start Inference
15: for  $(x, y) \in \mathcal{D}_{\text{test}}$  do
16:    $\{C_j\}_{j=1}^N \leftarrow \mathcal{M}(\text{prompt}, x)$  ▷  $K'$ -shot ICL
17:   Initialize  $\mathbf{g} \leftarrow [0, \dots, 0] \in \mathbb{R}^{|\mathcal{Y}|}$ 
18:   for  $j \leftarrow 1$  to  $N$  do ▷ Test-Time Overhead
19:      $\mathbf{g} \leftarrow \mathbf{g} + w_j (U_{\mathcal{Y}} \cdot C_j)$  ▷ early decode
20:   end for
21:    $\hat{y} \leftarrow \arg \max_{y \in \mathcal{Y}} \mathbf{g}$ 
22:    $\mathcal{Z} \leftarrow \mathcal{Z} \cup \hat{y}$ 
23: end for
24: return  $\mathcal{Z}$ 

```

---

training  $w$  (line 9 and 13 in Algorithm 1), saving tremendous training time and GPU memory. The cache only requires  $O(|\mathcal{Y}| \times N \times |\mathcal{D}_{\text{train}}|)$  space. At inference, the overhead of our method over  $K'$ -shot ICL is to early decode  $N$  components and apply the learned weights, i.e.,  $\sum_{j=1}^N w_j (U_{\mathcal{Y}} \cdot C_j)$ . As both  $|\mathcal{Y}|$  and  $N$  are small<sup>3</sup>, the overhead is negligible compared to the computation of the LLM itself.

### 5.2 Baselines

**Standard ICL.** The simplest baseline is to use all the  $K$  labeled examples as demonstrations. Since the other methods use  $K'$  examples as demonstrations, we report the accuracy of standard  $K'$ -shot ICL using the same  $\mathcal{D}_{\text{demo}}$  for reference.

**Prompt Selection.** Liu et al. (2022) improve ICL accuracy by selecting demonstrations from a pool of labeled data for each test example. Here, we select from the given  $K$  labeled examples. Following Rubin et al. (2022), we use SBERT (Reimers and Gurevych, 2019) to encode examples into sentence embeddings and select the  $K' = \{3, 4\}$  nearest neighbors under cosine similarity as the demonstrations for each test example.

**Calibration.** As LLMs tend to predict a certain class over others, Zhao et al. (2021) reweight the output class probabilities. They use context-free

<sup>3</sup> $|\mathcal{Y}| \leq 4, N < 2000$  for all tasks and LLMs in this paper.

		SST2	BoolQ	QQP	WiC	RTE	MNLI	AGNews	ARC-Easy	Avg.
Llama-2-7B	STANDARD 3, 4	75.8 <sub>18.1</sub>	69.2 <sub>12.0</sub>	61.3 <sub>9.9</sub>	52.4 <sub>3.0</sub>	68.9 <sub>3.2</sub>	34.4 <sub>1.7</sub>	70.0 <sub>19.9</sub>	57.5 <sub>14.4</sub>	61.2
	STANDARD 12	77.8 <sub>19.6</sub>	<b>71.6</b> <sub>8.0</sub>	63.6 <sub>7.8</sub>	52.5 <sub>2.4</sub>	<b>71.1</b> <sub>2.1</sub>	37.0 <sub>2.8</sub>	69.0 <sub>20.8</sub>	<b>59.6</b> <sub>13.9</sub>	62.8
	PROMPTS 12	73.8 <sub>19.2</sub>	69.4 <sub>10.5</sub>	62.2 <sub>6.1</sub>	53.1 <sub>2.7</sub>	65.5 <sub>1.8</sub>	35.5 <sub>1.6</sub>	59.1 <sub>28.7</sub>	58.7 <sub>11.9</sub>	59.7
	CALIB+ 12	85.1 <sub>6.0</sub>	69.2 <sub>13.6</sub>	<b>73.6</b> <sub>6.1</sub>	55.1 <sub>5.1</sub>	70.3 <sub>2.7</sub>	45.5 <sub>7.8</sub>	77.8 <sub>12.2</sub>	58.6 <sub>14.6</sub>	66.9
	COMPRW 12	<b>88.5</b> <sub>2.8</sub>	70.4 <sub>11.2</sub>	71.4 <sub>5.4</sub>	<b>56.3</b> <sub>3.4</sub>	70.0 <sub>2.8</sub>	<b>48.3</b> <sub>4.8</sub>	<b>87.4</b> <sub>2.3</sub>	58.3 <sub>13.6</sub>	<b>68.8</b>
	STANDARD 24	77.8 <sub>19.5</sub>	71.6 <sub>7.3</sub>	66.4 <sub>5.0</sub>	53.2 <sub>3.3</sub>	<b>71.9</b> <sub>1.5</sub>	39.9 <sub>3.6</sub>	71.1 <sub>20.0</sub>	58.3 <sub>16.2</sub>	63.8
	PROMPTS 24	74.2 <sub>20.4</sub>	68.9 <sub>10.2</sub>	62.1 <sub>4.9</sub>	53.6 <sub>1.9</sub>	64.8 <sub>0.9</sub>	36.4 <sub>1.5</sub>	57.5 <sub>30.2</sub>	58.0 <sub>12.5</sub>	59.4
	COMPRW 24	87.6 <sub>5.0</sub>	70.3 <sub>11.9</sub>	<b>73.4</b> <sub>5.5</sub>	55.8 <sub>4.9</sub>	70.4 <sub>2.7</sub>	46.4 <sub>6.7</sub>	78.4 <sub>11.8</sub>	<b>59.2</b> <sub>14.4</sub>	67.7
Llama-2-13B	STANDARD 3, 4	89.0 <sub>5.3</sub>	77.6 <sub>6.8</sub>	71.0 <sub>6.8</sub>	55.0 <sub>3.8</sub>	75.1 <sub>2.3</sub>	45.7 <sub>7.9</sub>	70.8 <sub>20.6</sub>	73.2 <sub>13.7</sub>	69.7
	STANDARD 12	<b>91.3</b> <sub>1.9</sub>	78.1 <sub>7.4</sub>	70.5 <sub>7.3</sub>	<b>59.6</b> <sub>2.4</sub>	74.4 <sub>3.5</sub>	55.1 <sub>6.2</sub>	84.7 <sub>7.8</sub>	71.2 <sub>16.4</sub>	73.1
	PROMPTS 12	83.8 <sub>10.2</sub>	74.9 <sub>6.6</sub>	64.6 <sub>5.7</sub>	57.0 <sub>2.1</sub>	69.5 <sub>3.5</sub>	48.1 <sub>5.4</sub>	64.4 <sub>29.6</sub>	74.2 <sub>9.3</sub>	67.1
	CALIB+ 12	89.4 <sub>3.2</sub>	<b>78.4</b> <sub>6.1</sub>	72.1 <sub>4.1</sub>	58.1 <sub>5.1</sub>	75.3 <sub>1.9</sub>	57.3 <sub>4.5</sub>	81.5 <sub>8.7</sub>	74.7 <sub>9.3</sub>	73.3
	COMPRW 12	89.1 <sub>3.2</sub>	77.7 <sub>6.7</sub>	<b>72.7</b> <sub>3.3</sub>	58.7 <sub>4.0</sub>	<b>76.2</b> <sub>2.0</sub>	<b>60.2</b> <sub>3.7</sub>	<b>88.1</b> <sub>1.7</sub>	<b>76.2</b> <sub>6.8</sub>	<b>74.9</b>
	STANDARD 24	<b>91.9</b> <sub>0.6</sub>	77.7 <sub>8.2</sub>	69.5 <sub>8.5</sub>	60.6 <sub>1.6</sub>	74.7 <sub>3.3</sub>	58.2 <sub>7.0</sub>	85.8 <sub>4.4</sub>	69.1 <sub>17.7</sub>	73.5
	PROMPTS 24	81.9 <sub>13.2</sub>	75.1 <sub>5.7</sub>	64.9 <sub>4.8</sub>	57.3 <sub>1.8</sub>	69.5 <sub>1.7</sub>	49.8 <sub>5.1</sub>	65.2 <sub>28.9</sub>	74.2 <sub>9.4</sub>	67.2
	COMPRW 24	90.7 <sub>2.1</sub>	<b>78.6</b> <sub>6.2</sub>	73.1 <sub>4.3</sub>	<b>59.5</b> <sub>3.2</sub>	75.9 <sub>1.9</sub>	58.4 <sub>2.8</sub>	82.0 <sub>8.4</sub>	75.2 <sub>9.1</sub>	74.2
Mistral-Instruct-7B	STANDARD 3, 4	90.1 <sub>2.9</sub>	81.3 <sub>2.1</sub>	70.9 <sub>7.2</sub>	58.5 <sub>4.2</sub>	80.5 <sub>1.7</sub>	56.1 <sub>5.0</sub>	83.0 <sub>5.7</sub>	79.8 <sub>1.4</sub>	75.0
	STANDARD 12	91.4 <sub>0.9</sub>	81.2 <sub>2.2</sub>	67.9 <sub>8.7</sub>	57.7 <sub>2.8</sub>	79.1 <sub>1.6</sub>	57.2 <sub>3.6</sub>	85.4 <sub>3.6</sub>	77.7 <sub>5.6</sub>	74.7
	PROMPTS 12	90.3 <sub>2.5</sub>	81.1 <sub>1.9</sub>	68.7 <sub>5.8</sub>	57.1 <sub>2.7</sub>	79.1 <sub>1.6</sub>	56.7 <sub>3.2</sub>	84.9 <sub>3.0</sub>	79.0 <sub>3.0</sub>	74.6
	CALIB+ 12	<b>91.5</b> <sub>1.6</sub>	<b>81.3</b> <sub>1.8</sub>	<b>75.8</b> <sub>2.6</sub>	58.3 <sub>6.6</sub>	81.0 <sub>1.3</sub>	61.9 <sub>4.7</sub>	85.4 <sub>4.0</sub>	<b>79.6</b> <sub>1.6</sub>	76.9
	COMPRW 12	89.9 <sub>2.7</sub>	80.7 <sub>2.7</sub>	75.1 <sub>2.9</sub>	<b>60.0</b> <sub>4.9</sub>	<b>81.1</b> <sub>1.3</sub>	<b>64.7</b> <sub>4.6</sub>	<b>87.6</b> <sub>2.1</sub>	79.2 <sub>1.2</sub>	<b>77.3</b>
	STANDARD 24	91.2 <sub>1.0</sub>	80.8 <sub>2.3</sub>	65.3 <sub>8.4</sub>	57.4 <sub>4.0</sub>	75.6 <sub>1.7</sub>	56.6 <sub>6.5</sub>	85.8 <sub>4.3</sub>	68.8 <sub>16.9</sub>	72.7
	PROMPTS 24	90.5 <sub>2.6</sub>	<b>81.3</b> <sub>2.0</sub>	68.9 <sub>5.6</sub>	57.1 <sub>2.1</sub>	79.1 <sub>1.7</sub>	57.4 <sub>3.1</sub>	86.0 <sub>2.1</sub>	78.7 <sub>3.3</sub>	74.9
	COMPRW 24	<b>91.6</b> <sub>1.5</sub>	80.9 <sub>2.0</sub>	76.1 <sub>2.4</sub>	59.5 <sub>5.4</sub>	81.2 <sub>0.9</sub>	62.7 <sub>4.3</sub>	85.9 <sub>3.7</sub>	<b>80.1</b> <sub>1.2</sub>	77.2
	90.8 <sub>1.8</sub>	80.6 <sub>2.1</sub>	<b>76.4</b> <sub>1.7</sub>	<b>60.7</b> <sub>4.4</sub>	<b>81.6</b> <sub>1.0</sub>	<b>65.3</b> <sub>3.4</sub>	<b>88.0</b> <sub>1.8</sub>	79.0 <sub>1.6</sub>	<b>77.8</b>	

Table 4: ICL accuracy of 8 classification tasks and the average accuracy (Avg.). The number after a method denotes the number of labeled data used. We run 15 prompts for each task (5 disjoint sets of  $K$  labeled data and 3 templates) and report the mean accuracy and standard deviation. COMPRW achieves the best average accuracy in all setups.

inputs, such as “N/A”, to calibrate the probability distribution. However, Fei et al. (2023); Zhou et al. (2023) find context-free inputs sometimes ineffective, because in-domain context is important for calibration. Thus, we introduce CALIB+, which calibrates the original probabilities  $\mathbf{p} \in \mathbb{R}^{|\mathcal{V}|}$  with a training set of in-distribution labeled examples,  $\mathcal{D}_{\text{train}}$ . We train the calibration weights  $\mathbf{v} \in \mathbb{R}^{|\mathcal{V}|}$  on  $\mathcal{D}_{\text{train}}$  with cross-entropy loss and obtain the calibrated probabilities  $\hat{\mathbf{p}} = \text{softmax}(\mathbf{v} \cdot \mathbf{p})$ . For direct comparisons, we split the  $K$  examples into the same  $\mathcal{D}_{\text{demo}}$  and  $\mathcal{D}_{\text{train}}$  sets as component reweighting for CALIB+, where  $|\mathcal{D}_{\text{demo}}| = K'$ . We include the training details of both methods in A.4.

### 5.3 Results

We set  $K = \{12, 24\}$  following common practices. Table 4 compares our component reweighting (COMPRW) with standard ICL (STANDARD), prompt selection (PROMPTS), and calibration (CALIB+). First, we find that simply increasing the number of demonstrations from 4 to 24

has limited improvements in ICL accuracy, while the longer prompt greatly increases the inference time. For example, on Llama-2-7B, STANDARD 24 only improves the average accuracy by 2.6% over STANDARD 3, 4 and the accuracy even decreases on Mistral-Instruct. Second, PROMPTS performs the worst in most setups, likely because it is hard to find similar examples from a small pool of  $K$  examples, and a bad selection induces majority label biases. Third, both calibration (CALIB+) and component reweighting (COMPRW) achieve substantially better accuracy than STANDARD 3, 4 with little test-time overhead. Overall, COMPRW achieves the best average accuracy in all setups, outperforming STANDARD 12 by 6.0%, 1.8%, 2.6% on Llama-2-7B, Llama-2-13B, and Mistral-Instruct-7B, respectively, and outperforming STANDARD 24 by 6.0%, 2.2%, 5.1%. We run one-tailed paired t-tests comparing COMPRW with CALIB+ and find that p-values  $< 0.05$  in all 6 setups (see Table 5), showing that COMPRW performs significantly better CALIB+.

## 6 When Do Good Components Emerge?

We study the dynamics of components during pre-training by monitoring their accuracies on 32 checkpoints of Pythia-6.9B, uniformly sampled from the first to the last checkpoint. For each checkpoint, we run 4-shot ICL on AGNews with 3 templates  $\times$  3 sets of demonstrations. Figure 3 shows the average accuracy shaded by the standard deviation.

While the full model (green) fluctuates and has a large variance across prompts, the top-1 components (solid blue) achieve good accuracy at an early step and plateaus quickly. We also backtrack the top-1 components of different prompts at the last checkpoint (dashed blue), monitoring how they perform on average during pretrain. We observe that they are not the top components at the early stage (there are gaps between the two blue lines before the 75k steps), but start to perform steadily well from the middle stage. Our findings also hold on SST2 and Pythia-1.4B (Figure 7 in the appendix), suggesting that the model’s ability to do a task emerges before it is apparent from the full model.

## 7 Related Work

**Improving ICL.** Prior work shows that ICL performance varies greatly across different choices of demonstrations and templates (Zhao et al., 2021; Lu et al., 2022; Sclar et al., 2024; Voronov et al., 2024). While several approaches, such as prompt selection (Liu et al., 2022; Chang and Jia, 2023), template ensemble (Voronov et al., 2024), and many-shot ICL (Agarwal et al., 2024), address this issue, they treat LLMs like black boxes without studying the internals. Besides, they usually greatly increase inference time or require a large set of labeled data, which deviates from true few-shot learning (Perez et al., 2021). In comparison, our paper studies this problem by looking inside the LLMs. Rather than selecting prompts, we select components in a soft, learnable way. Our method only requires  $\{12, 24\}$  examples and has little computation overhead over 4-shot ICL at inference.

**Components Interpretation.** Components interpretation aims to study the function of different components in a trained model (Csordás et al., 2021; Elhage et al., 2021; Shah et al., 2024), where components considered in prior work include neurons (Radford et al., 2017; Gurnee et al., 2023), attention head (Voita et al., 2019; Olsson et al., 2022), and MLPs (Geva et al., 2021). To analyze

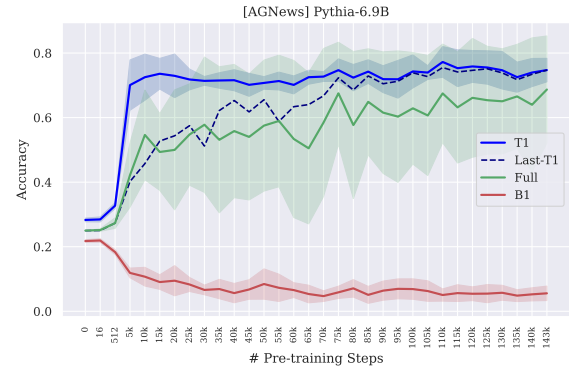


Figure 3: The ICL accuracy of the full model (green) fluctuates greatly during pretraining. However, good-performing components (T1) emerge in the early steps.

the components, pruning (Karnin, 1990), probing (Alain and Bengio, 2017), and early decoding (nostalgebraist, 2020; Geva et al., 2022) are widely used techniques. For example, Michel et al. (2019) prune away a large percentage of attention heads in Transformer models and show that only a few heads are critical to performance at test time. Li et al. (2024) train a linear probe on every attention head to discover the “truthful heads” inside LLMs. Early decoding interprets the investigated components in the textual space by projecting them through the output embedding matrix. Our model decomposition method is based on early decoding and we share some similarities with prior work (Yu et al., 2023; Merullo et al., 2024; Liu et al., 2023; Wang et al., 2023), especially in discovering individual components that perform well. Our contributions lie in providing a new view of ICL via decomposition, revealing the transferability of components across diverse ICL settings.

## 8 Conclusion

We introduce a new perspective of ICL via decomposing the model output into the sum of individual contributions of components. We then identify three types of component characteristics across 3 LLMs and 8 classification tasks. Our extensive analyses reveal consistency in component accuracy across prompts and suggest the promising direction of improving ICL by selecting components. To this end, we propose component reweighting, which learns to scale components differently on few-shot examples. Our method achieves the best average accuracy compared to prior methods. We hope this work can deepen our grasp of LLMs while motivating more methods for practical use.



## 9 Limitations

Our component reweighting method requires a small set of labeled data  $\mathcal{D}_{\text{train}}$  to train the component weights  $w$ . However, we believe it is not unreasonable to have at least  $K = 12$  labeled examples in total and we compare with baselines using the same  $K$  examples. On the other hand, we do not compare with fine-tuning-based baselines, such as LM-BFF(Gao et al., 2021) and LoRA (Hu et al., 2021), because they usually require a larger GPU memory and tend to overfit on few-shot examples. Another limitation is that we only experiment with classification tasks for easy evaluation. We leave it for future work to generalize our method to generation tasks by doing decomposition and reweighting at every token during generation.

## References

Rishabh Agarwal, Avi Singh, Lei M Zhang, Bernd Bohnet, Stephanie Chan, Ankesh Anand, Zaheer Abbas, Azade Nova, John D Co-Reyes, Eric Chu, et al. 2024. Many-shot in-context learning. *arXiv preprint arXiv:2404.11018*.

Guillaume Alain and Yoshua Bengio. 2017. [Understanding intermediate layers using linear classifier probes](#).

Stephen H. Bach, Victor Sanh, Zheng-Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, Maged S. Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Xiangru Tang, Mike Tian-Jian Jiang, and Alexander M. Rush. 2022. [Promptsources: An integrated development environment and repository for natural language prompts](#). *Preprint*, arXiv:2202.01279.

Amanda Bertsch, Maor Ivgi, Uri Alon, Jonathan Berant, Matthew R Gormley, and Graham Neubig. 2024. In-context learning with long-context models: An in-depth exploration. *arXiv preprint arXiv:2405.00200*.

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot

learners. *Advances in neural information processing systems*, 33:1877–1901. 576

Ting-Yun Chang and Robin Jia. 2023. [Data curation alone can stabilize in-context learning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8123–8144, Toronto, Canada. Association for Computational Linguistics. 578

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. [BoolQ: Exploring the surprising difficulty of natural yes/no questions](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics. 580

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*. 582

Róbert Csordás, Sjoerd van Steenkiste, and Jürgen Schmidhuber. 2021. [Are neural nets modular? inspecting functional modularity through differentiable weight masks](#). In *International Conference on Learning Representations*. 583

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. 2021. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1:1. 584

Yu Fei, Yifan Hou, Zeming Chen, and Antoine Bosselut. 2023. [Mitigating label biases for in-context learning](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14014–14031, Toronto, Canada. Association for Computational Linguistics. 585

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3816–3830, Online. Association for Computational Linguistics. 586

Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating models’ local decision boundaries via contrast sets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 587

633	1307–1323, Online. Association for Computational Linguistics.	688
634		689
635	Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. 2022. <a href="#">Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 30–45, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	690
636		691
637		692
638		693
639		694
640		695
641		696
642	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. <a href="#">Transformer feed-forward layers are key-value memories</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	697
643		698
644		699
645		700
646		701
647		702
648		703
649	Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. <a href="#">Finding neurons in a haystack: Case studies with sparse probing</a> . <i>Transactions on Machine Learning Research</i> .	704
650		705
651		706
652		707
653		708
654	Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. <a href="#">Surface form competition: Why the highest probability answer isn’t always right</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7038–7051, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	709
655		710
656		711
657		712
658		713
659		714
660		715
661		716
662	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2021. <a href="#">Lora: Low-rank adaptation of large language models</a> . In <i>International Conference on Learning Representations</i> .	717
663		718
664		719
665		720
666		721
667		722
668	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. <a href="#">Mistral 7b</a> . <i>arXiv preprint arXiv:2310.06825</i> .	723
669		724
670		725
671		726
672	E.D. Karnin. 1990. <a href="#">A simple procedure for pruning back-propagation trained neural networks</a> . <i>IEEE Transactions on Neural Networks</i> , 1(2):239–242.	727
673		728
674		729
675	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. <a href="#">Inference-time intervention: Eliciting truthful answers from a language model</a> . <i>Advances in Neural Information Processing Systems</i> , 36.	730
676		731
677		732
678		733
679		734
680	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What makes good in-context examples for GPT-3?</a> In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.	735
681		736
682		737
683		738
684		739
685		740
686		741
687		742
		743
	Sheng Liu, Lei Xing, and James Zou. 2023. <a href="#">In-context vectors: Making in context learning more effective and controllable through latent space steering</a> . <i>arXiv preprint arXiv:2311.06668</i> .	744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

744	<i>Volume 1 (Long and Short Papers)</i> , pages 1267–1273,	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	799
745	Minneapolis, Minnesota. Association for Computa-	Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz	800
746	tional Linguistics.	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	801
747	Alec Radford, Rafal Jozefowicz, and Ilya Sutskever.	<a href="#">you need</a> . In <i>Advances in Neural Information Pro-</i>	802
748	2017. <a href="#">Learning to generate reviews and discovering</a>	<i>cessing Systems</i> , volume 30.	803
749	<a href="#">sentiment</a> . <i>ArXiv preprint</i> , abs/1704.01444.	Elena Voita, David Talbot, Fedor Moiseev, Rico Sen-	804
750	Alec Radford, Jeffrey Wu, Rewon Child, David Luan,	rich, and Ivan Titov. 2019. <a href="#">Analyzing multi-head</a>	805
751	Dario Amodei, Ilya Sutskever, et al. 2019. Language	<a href="#">self-attention: Specialized heads do the heavy lift-</a>	806
752	models are unsupervised multitask learners. <i>OpenAI</i>	<a href="#">ing, the rest can be pruned</a> . In <i>Proceedings of the</i>	807
753	<i>blog</i> , 1(8):9.	<i>57th Annual Meeting of the Association for Computa-</i>	808
754	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-</a>	<i>tional Linguistics</i> , pages 5797–5808, Florence, Italy.	809
755	<a href="#">BERT: Sentence embeddings using Siamese BERT-</a>	Association for Computational Linguistics.	810
756	<a href="#">networks</a> . In <i>Proceedings of the 2019 Conference on</i>	Anton Voronov, Lena Wolf, and Max Ryabinin. 2024.	811
757	<i>Empirical Methods in Natural Language Processing</i>	Mind your format: Towards consistent evaluation of	812
758	<i>and the 9th International Joint Conference on Natu-</i>	in-context learning improvements. <i>arXiv preprint</i>	813
759	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	<i>arXiv:2401.06766</i> .	814
760	3982–3992, Hong Kong, China. Association for Com-	Alex Wang, Amanpreet Singh, Julian Michael, Felix	815
761	putational Linguistics.	Hill, Omer Levy, and Samuel Bowman. 2018. <a href="#">GLUE:</a>	816
762	Alexey Romanov and Chaitanya Shivade. 2018.	<a href="#">A multi-task benchmark and analysis platform for nat-</a>	817
763	<a href="#">Lessons from natural language inference in the clini-</a>	<a href="#">ural language understanding</a> . In <i>Proceedings of the</i>	818
764	<a href="#">cal domain</a> . In <i>Proceedings of the 2018 Conference</i>	<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>	819
765	<i>on Empirical Methods in Natural Language Process-</i>	<i>and Interpreting Neural Networks for NLP</i> , pages	820
766	<i>ing</i> , pages 1586–1596, Brussels, Belgium. Associa-	353–355, Brussels, Belgium. Association for Com-	821
767	tion for Computational Linguistics.	putational Linguistics.	822
768	Ohad Rubin, Jonathan Herzig, and Jonathan Berant.	Tony Wang, Miles Kai, Kaivalya Hariharan, and Nir	823
769	2022. <a href="#">Learning to retrieve prompts for in-context</a>	Shavit. 2023. <a href="#">Forbidden facts: An investigation of</a>	824
770	<a href="#">learning</a> . In <i>Proceedings of the 2022 Conference of</i>	<a href="#">competing objectives in llama 2</a> . In <i>Socially Respon-</i>	825
771	<i>the North American Chapter of the Association for</i>	<i>sible Language Modelling Research</i> .	826
772	<i>Computational Linguistics: Human Language Tech-</i>	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-	827
773	<i>nologies</i> , pages 2655–2671, Seattle, United States.	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva	828
774	Association for Computational Linguistics.	Naik, Arjun Ashok, Arut Selvan Dhanasekaran,	829
775	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane	Anjana Arunkumar, David Stap, Eshaan Pathak,	830
776	Suhr. 2024. <a href="#">Quantifying language models’ sensitiv-</a>	Giannis Karamanolakis, Haizhi Lai, Ishan Puro-	831
777	<a href="#">ity to spurious features in prompt design or: How i</a>	hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,	832
778	<a href="#">learned to start worrying about prompt formatting</a> .	Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,	833
779	In <i>The Twelfth International Conference on Learning</i>	Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,	834
780	<i>Representations</i> .	Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,	835
781	Harshay Shah, Andrew Ilyas, and Aleksander Madry.	Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,	836
782	2024. Decomposing and editing predictions by	Shailaja Keyur Sampat, Siddhartha Mishra, Sujan	837
783	modeling model computation. <i>arXiv preprint</i>	Reddy A, Sumanta Patro, Tanay Dixit, and Xudong	838
784	<i>arXiv:2404.11534</i> .	Shen. 2022. <a href="#">Super-NaturalInstructions: Generaliza-</a>	839
785	Richard Socher, Alex Perelygin, Jean Wu, Jason	<a href="#">tion via declarative instructions on 1600+ NLP tasks</a> .	840
786	Chuang, Christopher D. Manning, Andrew Ng, and	In <i>Proceedings of the 2022 Conference on Empiri-</i>	841
787	Christopher Potts. 2013. <a href="#">Recursive deep models for</a>	<i>cal Methods in Natural Language Processing</i> , pages	842
788	<a href="#">semantic compositionality over a sentiment treebank</a> .	5085–5109, Abu Dhabi, United Arab Emirates. As-	843
789	In <i>Proceedings of the 2013 Conference on Empiri-</i>	association for Computational Linguistics.	844
790	<i>cal Methods in Natural Language Processing</i> , pages	Adina Williams, Nikita Nangia, and Samuel Bowman.	845
791	1631–1642, Seattle, Washington, USA. Association	2018. <a href="#">A broad-coverage challenge corpus for sen-</a>	846
792	for Computational Linguistics.	<a href="#">tence understanding through inference</a> . In <i>Proceed-</i>	847
793	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	<i>ings of the 2018 Conference of the North American</i>	848
794	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>Chapter of the Association for Computational Lin-</i>	849
795	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	<i>guistics: Human Language Technologies, Volume</i>	850
796	Bhosale, et al. 2023. Llama 2: Open founda-	<i>1 (Long Papers)</i> , pages 1112–1122, New Orleans,	851
797	tion and fine-tuned chat models. <i>arXiv preprint</i>	Louisiana. Association for Computational Linguis-	852
798	<i>arXiv:2307.09288</i> .	tics.	853
		Qinan Yu, Jack Merullo, and Ellie Pavlick. 2023. <a href="#">Char-</a>	854
		<a href="#">acterizing mechanisms for factual recall in language</a>	855
		<a href="#">models</a> . In <i>Proceedings of the 2023 Conference on</i>	856

857 *Empirical Methods in Natural Language Processing*,  
858 pages 9924–9959, Singapore. Association for Com-  
859 putational Linguistics.

860 Biao Zhang and Rico Sennrich. 2019. Root mean square  
861 layer normalization. *Advances in Neural Information*  
862 *Processing Systems*, 32.

863 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.  
864 [Character-level convolutional networks for text clas-](#)  
865 [sification](#). In *Advances in Neural Information Pro-*  
866 *cessing Systems*, volume 28. Curran Associates, Inc.

867 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and  
868 Sameer Singh. 2021. Calibrate before use: Improv-  
869 ing few-shot performance of language models. In  
870 *International conference on machine learning*, pages  
871 12697–12706. PMLR.

872 Han Zhou, Xingchen Wan, Lev Proleev, Diana Mincu,  
873 Jilin Chen, Katherine Heller, and Subhrajit Roy.  
874 2023. Batch calibration: Rethinking calibration for  
875 in-context learning and prompt engineering. *arXiv*  
876 *preprint arXiv:2309.17249*.

## A Appendix

### A.1 Tasks and Templates

Table 7 summarizes all the 13 datasets we use in the paper, where we construct balanced test sets by randomly sampling 2000 examples in each task. We form the prompts by concatenating demonstrations in a randomly shuffled order. To avoid the recency bias (Zhao et al., 2021), we keep shuffling the demonstrations until the last two have different labels. For minimally conservative templates (§4.1), Table 8 compares the contrast sets we construct on Llama-2-7B. For our case study on Task069 and Task070, we sample 3 templates from Sclar et al. (2024). In §4.3, we calculate the correlation between component accuracies of every pair of runs from the two tasks under a fixed template. The correlation score ( $r = -0.6$ ) is averaged across runs and templates. Figure 6 compare the prompts of Task069 and Task070, which consist of an instruction followed by  $K$  templated demonstrations. Originally, the two tasks have  $\sim 4\%$  of parallel examples. To make our task transfer challenging, we discard these overlapped examples.

### A.2 Label-Biased Components

We say a component is label-biased when it always predicts a certain label on the entire test set §3.2. We focus on the most biased components in binary classification tasks, i.e., the two components that have the largest value of  $(\text{logit}_0 - \text{logit}_1)$  and  $(\text{logit}_1 - \text{logit}_0)$ , respectively, where  $\text{logit}_0 \in \mathbb{R}$  and  $\text{logit}_1 \in \mathbb{R}$  are the output logits on the two classes. We name these two components as Biased Component-0 and Biased Component-1, respectively. To understand how biased these two components are, we alter the demonstrations and observe their behavior. Specifically, we consider three settings: demonstrations balanced in labels (green in Figure 5), demonstrations of all positive labels ( $[1, 1, 1, 1]$ ; blue), and demonstrations of all negative labels ( $[0, 0, 0, 0]$ ; red). We fix the template and sample 5 disjoint sets of demonstrations for each setting. Each dot in Figure 5 shows the components' prediction on an example. We show that both Biased Component-0 and Biased Component-1 still insist on predicting a certain label on all examples, regardless of the labels in the prompts. These extremely biased components may be related to the majority label bias of ICL observed by Zhao et al. (2021) on the full model.

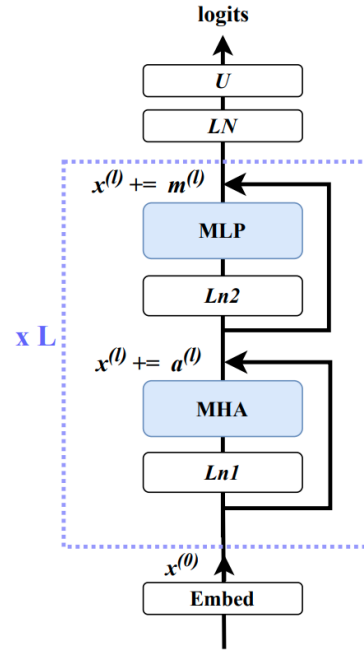


Figure 4: Transformer architecture in GPT2.

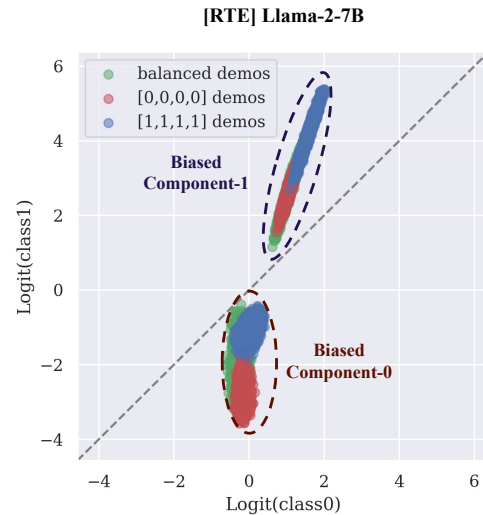


Figure 5: Each dot represents an example in the test set. The two most biased components still insist on predicting the same label on the entire test set regardless of the labels of the demonstrations.

### A.3 LayerNorms

Figure 4 shows the transformer architecture in GPT2-like models. Because the layernorms inside each block are before MHA and MLP, known as Pre-LN, Eq. 1 has already taken Ln2 and Ln2 into account, and Eq. 8 only has the term for the final layernorm,  $\text{LN}(\cdot)$ .

Both Llama-2 and Mistral families use RMSNorm (Zhang and Sennrich, 2019), a layer normalization variant without centering and adding

bias terms. Formally, let  $x \in \mathbb{R}^d$  be the input, the root mean square norm  $\text{LN}(x)$  is:

$$\text{LN}(x) = \frac{x}{\text{RMS}(x)} \odot \gamma, \quad (10)$$

$$\text{RMS}(x) = \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2}, \quad (11)$$

where  $\gamma \in \mathbb{R}^d$  is the affine transform parameters and  $\odot$  denotes element-wise multiplication.

#### A.4 Training Details and Hyperparameters

For both COMPRW and CALIB+ methods, we train a linear layer on  $\mathcal{D}_{\text{train}}$  with stochastic gradient descent. Because we do not have an additional dev set to tune the hyperparameters, we use the same hyperparameters on all the tasks and models and do early stopping based on the loss and accuracy on  $\mathcal{D}_{\text{train}}$ . Specifically, we set learning rate = 0.05 for both methods and  $\lambda = 0.1$  for the L1 regularization term in COMPRW. We run all our ICL experiments on a single RTX A6000 GPU (48G). Both the component reweighting and calibration training processes can be run on a single i7 CPU within a minute.

#### A.5 Tests of Significance

We run one-tailed paired t-tests to test whether COMPRW outperforms CALIB+ significantly. In Table 4, we have the results of 15 prompts for each task and 8 tasks in total. For each model, we aggregate the 120 accuracy scores of COMPRW and CALIB+, respectively, and then calculate the p-values. Table 5 shows that p-values  $< 0.05$  in 6/6 setups, suggesting that COMPRW performs significantly better than CALIB+.

	Llama-2-7B	Llama-2-13B	Mistral-Instruct-7B
$K = 12$	0.0010	0.0002	0.0470
$K = 24$	0.0003	0.0001	0.0027

Table 5: The p-values  $< 0.05$  in all 6 setups (3 LLMs, with  $K = \{12, 24\}$  labeled examples), showing that COMPRW performs significantly better than CALIB+.

#### A.6 More Results on Transferability

In §4, we study the transferability of components across different choices of demonstrations and templates. Here, Table 9 shows the full results on all LLMs and tasks. We observe the same findings as Table 2: component accuracies agree well across

randomly sampled demonstrations, but have much weaker agreements across randomly sampled templates. Because constructing minimally-contrastive templates requires non-trivial manual efforts, we only build contrast sets for 5 tasks on Llama-2-7B (shown in Table 2), where these tasks have the largest variances across templates.

	SST2	BoolQ	QQP	WiC	RTE	MNLI	AGNews	ARC-Easy
<b>Llama-2-7B</b>	37.8	18.9	43.4	44.2	35.4	28.2	13.4	11.5
<b>Llama-2-13B</b>	32.6	18.7	37.2	39.3	32.1	26.0	12.4	13.1
<b>Mistral-Instruct-7B</b>	31.9	14.0	32.3	32.4	27.6	20.9	14.5	8.4

Table 6: We report the average percentage of label-biased components across 15 prompts for each task. A label-biased component always predicts the same label on the entire test set.

Dataset	Task	# Classes
<b>SST-2</b> (Socher et al., 2013)	Sentiment Analysis	2
<b>Yelp-polarity</b> (Zhang et al., 2015)	Sentiment Analysis	2
<b>BoolQ</b> (Clark et al., 2019)	Yes/No QA	2
<b>BoolQ Contrast Set</b> (Gardner et al., 2020)	Yes/No QA	2
<b>QQP</b> (Wang et al., 2018)	Paraphrase Identification	2
<b>WiC</b> (Pilehvar and Camacho-Collados, 2019)	Word Sense Disambiguation	2
<b>RTE</b> (Wang et al., 2018)	Natural Language Inference	2
<b>MNLI</b> (Williams et al., 2018)	Natural Language Inference	3
<b>MedNLI</b> (Romanov and Shivade, 2018)	NLI in Medical Domain	3
<b>AGNews</b> (Zhang et al., 2015)	Topic Classification	4
<b>ARC-Easy</b> (Clark et al., 2018)	Multiple-Choice QA	4
<b>Task069</b> (Mishra et al., 2022; Wang et al., 2022)	Abductive NLI	2
<b>Task070</b> (Mishra et al., 2022; Wang et al., 2022)	Abductive NLI	2

Table 7: Summary of all the datasets.

Task	Templates	Labels	Accuracy
SST-2	<b>T1</b> Review: {text}\n Do you think the review is positive or negative? {label}	negative/positive	50.6 ± 0.7
	<b>T2</b> Review: {text} <b>{space}</b> \n Do you think the review is positive or negative? {label}	negative/positive	72.7 ± 6.1
BoolQ	<b>T1</b> Based on the following passage, {question}? {passage}\n Answer: {label}	No/Yes	52.5 ± 2.0
	<b>T2</b> Based on the following passage, {question}? {passage} Answer: {label}	No/Yes	66.7 ± 2.1
QQP	<b>T1</b> Are the questions "{sent1}" and "{sent2}" asking the same thing? {label}	<b>no / yes</b>	54.3 ± 1.1
	<b>T2</b> Are the questions "{sent1}" and "{sent2}" asking the same thing? {label}	<b>No / Yes</b>	68.7 ± 4.1
AGNews	<b>T1</b> {text} \n Is this a piece of news regarding World, Sports, Business, or Technology? {label}	World/ Sports/ Business/ Technology	43.9 ± 8.7
	<b>T2</b> {text} Is this a piece of news regarding World, Sports, Business, or Technology? {label}	World/ Sports/ Business/ Technology	88.5 ± 0.8

Table 8: We construct minimally contrastive templates that only differ slightly (colored in red) but yield large differences in 4-shot ICL accuracy on Llama-2-7B. We report the average accuracy and standard deviation across 5 ICL runs under the same template.

**Task069**

In this task, you will be shown a short story with a beginning, two potential middles, and an ending. Your job is to choose the middle statement that makes the story **coherent / plausible by writing "1" or "2"** in the output. If both sentences are plausible, pick the one that **makes most sense**.

Beginning: The clown was blowing several bubbles to the kids. Middle 1: Isaiah kept on popping the bubbles. Middle 2: Isaiah kept eating the bubbles. Ending: He said that Isaiah is currently sick from ingesting too much soap. **Answer: 2**

*K* demonstrations :

**Task070**

In this task, you will be shown a short story with a beginning, two potential middles, and an ending. Your job is to choose the middle statement that makes the story **incoherent / implausible by indicating "1" or "2"** in the output. If both sentences are plausible, pick the one that **makes less sense**.

Beginning: Killy was 9 months pregnant and almost ready to pop. Middle 1: Luckily Killy's water broke when she was in hospital. Middle 2: Killy's water broke when she was on a walk. Ending: Five minutes later, she delivered her baby with the help of passersby. **Answer: 1**

*K* demonstrations :

Figure 6: Comparing the prompts of Task069 and Task070. We follow the formats of Sclar et al. (2024) and prepend the task instructions before *K* demonstrations. We ensure that the two tasks do not have parallel examples to make the transfer experiment (§4.3) challenging.

	SST2	BoolQ	QQP	WiC	RTE	MNLI	AGNews	ARC
<b>Correlation</b>	<i>Llama-2-7B</i>							
(1) Demo	0.81	0.84	0.60	0.65	0.75	0.65	0.89	0.88
(2) Temp	0.40	0.16	0.03	0.15	0.19	0.09	0.68	0.44
<b>IoU</b>								
(1) Demo	0.36	0.74	0.27	0.21	0.53	0.24	0.63	0.70
(2) Temp	0.12	0.01	0.01	0.03	0.05	0.01	0.20	0.20
<b>Correlation</b>	<i>Llama-2-13B</i>							
(1) Demo	0.83	0.84	0.63	0.67	0.78	0.73	0.91	0.91
(2) Temp	0.57	0.30	0.09	0.19	0.28	0.16	0.76	0.55
<b>IoU</b>								
(1) Demo	0.26	0.71	0.31	0.18	0.46	0.39	0.55	0.65
(2) Temp	0.21	0.11	0.07	0.01	0.21	0.07	0.25	0.30
<b>Correlation</b>	<i>Mistral-Instruct-7B</i>							
(1) Demo	0.88	0.91	0.72	0.75	0.87	0.82	0.92	0.97
(2) Temp	0.58	0.44	0.19	0.26	0.40	0.30	0.77	0.60
<b>IoU</b>								
(1) Demo	0.39	0.59	0.27	0.29	0.50	0.45	0.68	0.80
(2) Temp	0.10	0.17	0.06	0.05	0.17	0.09	0.29	0.22

Table 9: Full results of the average correlation and IoU between (1) two random sets of demonstrations and (2) two random templates. **ARC** is short for the ARC-Easy task.



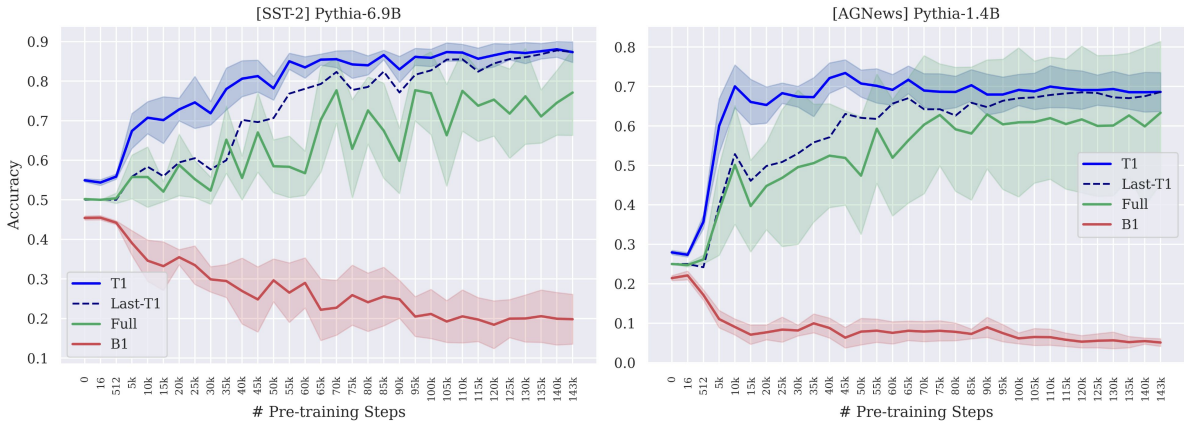


Figure 7: 4-shot ICL accuracy on different pretraining checkpoints. We compare the full model (green) with the top-1 (solid blue) and bottom-1 (red) components. The dashed blue line tracks how the top-1 components of the last checkpoint (Last-T1) perform across time.

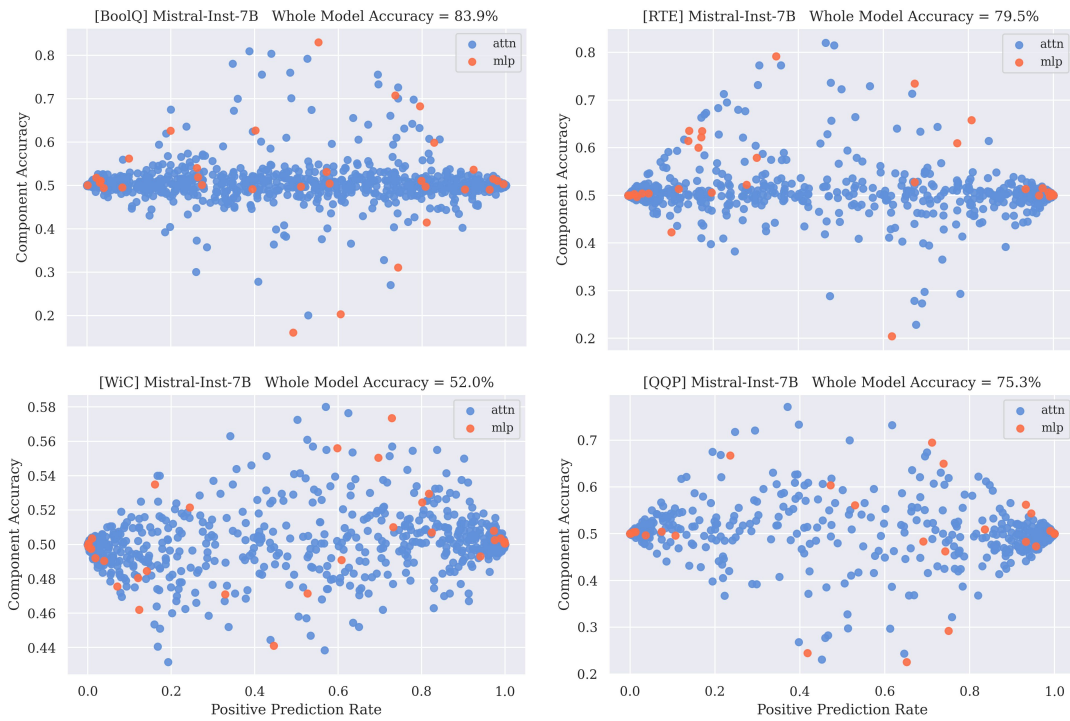


Figure 8: Each dot represents a component (attention head: blue; MLP: orange) under 4-shot ICL on Mistral-Instruct-7B. The x-axis shows how often a component predicts “positive” on the test set.