# COUNTING HYPOTHESIS: POTENTIAL MECHANISM OF IN-CONTEXT LEARNING

**Jung H. Lee**
Pacific Northwest National Laboratory
Seattle, WA
jung.lee@pnnl.gov

**Sujith Vijayan**
School of Neuroscience
Virginia Tech
Blacksburg, VA
neuron99@vt.edu

February 3, 2026

## ABSTRACT

In-Context Learning (ICL) indicates that large language models (LLMs) pretrained on a massive amount of data can learn specific tasks from input prompts' examples. ICL is notable for two reasons. First, it does not need modification of LLMs' internal structure. Second, it enables LLMs to perform a wide range of tasks/functions with a few examples demonstrating a desirable task. ICL opens up new ways to utilize LLMs in more domains, but its underlying mechanisms still remain poorly understood, making error correction and diagnosis extremely challenging. Thus, it is imperative that we better understand the limitations of ICL and how exactly LLMs support ICL. Inspired by ICL properties and LLMs' functional modules, we propose 'the counting hypothesis' of ICL, which suggests that LLMs' encoding strategy may underlie ICL, and provide supporting evidence.

## 1 Introduction

Large language models (LLMs) have demonstrated that they can learn diverse functions/tasks in a wide range of domains, but fine-tuning and adapting them involve extensive resources, which creates challenges when deploying them. Notably, a line of studies suggested that LLMs pretrained on vast amounts of data could perform novel tasks without being trained on them. Brown et al. [1] showed that LLMs infer necessary tasks based on examples included in input prompts and find correct answers to prompted queries, which is often referred to as In-Context Learning (ICL); see [2, 3] as well. For instance, pretrained LLMs function as a thesaurus to an input prompt query, if the input prompt provides a pair of synonyms.

As ICL allows users to ask LLMs for specific tasks without training, it facilitates adopting LLMs in a wide range of domains. Then, how do LLMs realize ICL? Earlier studies proposed LLMs' mechanisms that support ICL. First, Xie et al. [4] suggested that transformers (i.e., backbone of language models) learn concept distributions during pretraining and use Bayesian inference to infer the best matching concept from given ICL prompts; see [5] as well. Second, an earlier study [6] proposed that dual gradient descent can underlie ICL, in which inputs, instead of model parameters, are updated to reduce errors. They further showed that transformers can implement dual gradient descent when they are trained for autoregression tasks. The empirical evidence of dual-gradients underlying ICL in pretrained LLMs was presented [7], but it was challenged by a later study [8].

Despite the theories of ICL discussed above, the exact mechanisms underlying ICL remain elusive. As transformers (the backbone of LLMs) [9], use simple fixed vector arithmetic operations (after pretraining), these operations should enable LLMs to realize ICL. That is, the vector operations occurring in transformer blocks should be able to extract task information embedded in the examples and apply them to the prompted queries. Interestingly, earlier studies [10, 11, 12] suggested that the geometry of word embeddings may be linked to semantic meanings. For example, in the embedding space, the word 'king' can be converted to the word 'queen' when it is subtracted by 'man' and added by 'woman' [10]. A more recent [13] study showed that adding steering vectors to layer activations can change LLMs' behaviors. These studies indicate that arithmetic operations occurring in transformers can directly manipulate LLMs

behaviors, leading us to hypothesize that ICL evokes internal representations working like steering vectors. If true, how do fixed operations in pretrained LLMs create steering like inputs, when they face diverse tasks and examples from ICL prompts?

In this study, we probe LLMs' functional modules and the observations on ICL to gain insights into how LLMs can support ICL, which leads to our proposed counting hypothesis.

## 2 Counting Hypothesis

We note that 1) ICL prompts have unique structure and unnatural sequence [4], 2) LLMs use a sequence of transformer blocks, and their outputs are added to so-called 'residual' streams [14] and 3) examples included in ICL prompts are not restricted and can be drawn from a large pool of candidates. Below, we provide more details and the implications of LLMs' mechanisms underlying ICL.

### 2.1 ICL prompts

ICL prompts contain multiple examples, each of which contains a query and an answer separated by delimiters/separators. Here, we use the dataset generator released by an earlier study [11] to create ICL prompts for 6 different tasks, antonyms, synonyms, countries-their capital cities (country-capital), English-French words, products-their producers (product-company), players-sports (person-sport). For all tasks, ICL prompts consist of 5 example pairs of queries and answers and a single test query (without a corresponding answer), and LLMs are required to find its counterpart using example pairs. Each query and answer are marked by the separators, 'Q:' and 'A:' respectively.

### 2.2 Transformations through transformer blocks

LLMs first convert the words in input prompts into tokens (i.e., sub-words) and use a sequence of transformer layers to process tokens, each of which consists of self-attention (SA) layer and Feed-Forward Network (FFN) [9]. A transformer layer $l$ receives inputs $h^{l-1}$ from its previous layer and generates $h^l$ using its own SA and FFN. Notably, the transfer layer uses residual connections, and consequently $h^l$ can be recursively expressed as $h^l = h^{l-1} + a^l + m^l(a^l + h^{l-1})$ [15], where $a^l$ and $m^l$ denote the output of SA layer and FFN in transformer layer $l$, respectively.

FFNs in LLMs consist of two synaptic layers, ($\boldsymbol{w}^{1st}$ and $\boldsymbol{w}^{2nd}$) and one hidden layer of cells $m_i$ (Eq. 1).

$$m_i = \Sigma_j \boldsymbol{w}_{ij}^{1st} x_j,$$
$$O_k = \Sigma_i \boldsymbol{w}_{ki}^{2nd} g(m_i), \tag{1}$$

where $g$ is the activation function. SA layers evaluate attention score $A_{ij}^l$, which determines the strength of the influence of token $i$ onto token $j$ (Eq, 2).

$$A_{ij}^l = (\boldsymbol{K}\vec{h}_i^{l-1})^T \cdot \boldsymbol{Q}\vec{h}_j^{l-1} \tag{2}$$

Then, SA layer's outputs are modulated by value matrix $\boldsymbol{V}$.

$$a_{ij}^l = \frac{1}{\sqrt{d}} softmax(A_{ij}^l)\boldsymbol{V}h_i^{l-1} \tag{3}$$

, where $softmax$ is estimated over $i$ to normalize the influence from all tokens $i$[1].

Then, how do LLMs support ICL? We made two notes on FFNs' potential roles. First, two consecutive tokens in the language prompts are semantically related. That is, for a given token, there are a finite number of possible choices for the next token. Second, earlier studies [16, 15, 17] suggested that FFNs work as associative memory, storing pairs of keys and values. Specifically, Geva et al. [16] found that a unit in the hidden layer, referred to as memory cell $m_i$, could detect a distinct linguistic feature $f_i$.

Interestingly, if linguistic features $\vec{f_i}$s and corresponding values $\vec{v_i}$s are known/determined, a FFN can be explicitly built as an associative memory by setting a row of $\boldsymbol{w}^{1st}$ to be $\vec{f_i}$ and a column of $\boldsymbol{w}^{2nd}$ to be $\vec{v_i}$ (Supplemental Fig. S1A). When an input $\vec{x}$ is similar to $\vec{f_i}$ ($\vec{x} \approx \vec{f_i}$), memory cell $m_i$ in the hidden layer will be activated strongly. If all other memory cells $m_{j\neq i}$s are quiescent, the outputs will be determined by $m_i$. Specifically. $\vec{O}_j = \boldsymbol{w}_{ji}^{2nd} g(m_i)$, where $g(x)$ is an activation function of a FFN's hidden layer. In this way, FFNs can detect features $f_i$ and return desirable key vectors $\vec{v_i}$ (Supplemental Fig. S1B). More importantly, with this FFN, it becomes clear that FFNs can process multiple

---

[1]For brevity, we ignore the output matrix $\boldsymbol{O}$ here

keys in parallel. The residual stream of token $i$ in layer $l-1$ ($\vec{h}_i^{l-1}$) can be the superposition of multiple features (Eq. 4),

$$\vec{h}_i^{l-1} = \Sigma_m \alpha_m \vec{f}_m. \tag{4}$$

This will activate multiple memory cells of FFN in layer $l$ simultaneously and then return multiple values at once (Supplemental Fig. S1C). This line of thoughts leads us to hypothesize that FFNs may encode possible answers to given inputs (i.e., $h^{l-1}$): Eq. 5.

$$
\begin{aligned}
h^l &= h^{l-1} + a^l + m^l(a^l + h^{l-1}) \\
&\equiv h^{l-1} + a^l + \Sigma_j \beta_j An_j^l
\end{aligned}
\tag{5}
$$

, where $An_j^l$ is a possible answer to $a^l + h^{l-1}$, and $\beta_j$ is a coefficient constant. Further, due to the recursive nature of residual connections, $h^l$ can be further expanded as in Eq. 6.

$$
\begin{aligned}
h^l &= h^0 + a^1 + \Sigma_{k=2,l-1} a^k + \Sigma_{k=1,l-1} m^k \\
&= h^0 + \Sigma_j \alpha_j h_k^0 + \Sigma_{k=2,l-1} a^k + \Sigma_{k=1,n} m^k
\end{aligned}
\tag{6}
$$

As $a^l$ is a linear combination of $h^{l-1} \equiv m^{l-1} + a^{l-1} + h^{l-2}$, only $h^0$ and $\Sigma_j \alpha_j h_k^0$ are completely independent from FFNs. That is, FFNs' outputs $m^l$, which can be decomposed into potential answers, would dominate $h^l$ in deeper layers $l \gg 1$. Thus, we hypothesize that the residual streams are superpositions of possible answers ($An_k^l$) stored in FFNs (i.e., values in the terminology of associative memory), and LLMs' final answers can be determined via the competition between them ($An_k^l$). Interestingly, an earlier study [18] found that two consecutive transformer layers are interchangeable with minimal impact on LLMs' answers, suggesting that multiple layers could create the same functional space together. Additionally, as self-attention makes the competition across tokens possible, LLMs may track the competition of a large number of possible answers in parallel using the neighboring layers.

## 2.3 Observation of ICL

We note that ICL has two properties. First, the order of ICL examples does not significantly impact LLMs' answers. Second, when more examples are provided, LLMs provide more accurate answers. Below, we probe SA layers to gain insights into how LLMs support ICL.

### 2.3.1 Invariance of the order of examples

By definition, attention score $A_{ij}^l$ (Eq. 2) is location specific and sensitive to the order of examples. As the final separator predicts the answer to the test query, we consider the attention score $A_{j=last}^l$ between token $i$ and the last separator. If all examples are equally important and can be shuffled with minimal impact on LLMs' answers, $A_{j=last}^l$ can be approximated to be constant. We note that ICL prompts have three different types of tokens, queries ($Q$), separators ($S$) and answers ($A$) and assume that they would have different attention scores. With this assumption, we group them accordingly and approximate attention scores using constants. Then, $A_{j=last}^l$ can be written as in Eq. 7,

$$
\begin{aligned}
A_{j=last}^l &= A_j^{l,Q} + A_j^{l,S} + A_j^{l,A} \\
&\approx \Sigma_{m=1,N_Q} \alpha + \Sigma_{n=1,N_S} \beta + \Sigma_{o=1,N_A} \gamma
\end{aligned}
\tag{7}
$$

, where $N_Q$, $N_S$ and $N_A$ denote the number of queries, separators and answer tokens, indices $m$, $n$ and $o$ run for queries, separators and answer tokens and $\alpha$, $\beta$ and $\gamma$ denote the attention score of individual tokens, which we hypothesize to be constant due to the invariance of ICL examples.

With these constant attention scores, we can approximate SA layers' outputs $a_{j=last}^l$ as linear systems (Eq. 8).

$$
\begin{aligned}
a_{j=last}^l &= \frac{e^\alpha}{M} \Sigma_{m=1,N_Q} \boldsymbol{V} h_m^{l-1} \\
&+ \frac{e^\beta}{M} \Sigma_{n=1,N_S} \boldsymbol{V} h_n^{l-1} + \frac{e^\gamma}{M} \Sigma_{o=1,N_A} \boldsymbol{V} h_o^{l-1}, \\
M &= \sqrt{d}(N_Q e^\alpha + e^\beta + N_S e^{\gamma N_A}).
\end{aligned}
\tag{8}
$$

, where $h_i^l$ denotes the residual stream of token $i$ in layer $l$. That is, self-attention outputs can also be grouped according to token categories, queries, separators and answers.

### 2.3.2 More examples ensure more accurate answers

LLMs' predictions on ICL prompts become more accurate, as the number of examples increases. To investigate this property, we consider a simple and ideal case, in which an ICL prompt contains two pairs of examples and a single test query, and all words are tokenized into single tokens. That is, the input prompt is a sequence of 8 tokens $\{q_1, s_1, a_1, q_2, s_2, a_2, q_3, s_3\}$, where $q_i$, $s_i$ and $a_i$ denote a query, a separator and an answer token.

If $A_{ij}$ is a constant depending on the type of token (Eq. 7), the attention output $a^l_{j=last}$ at the last token ($s_3$) can be written as in Eq. 9.

$$
\begin{aligned}
a^l_{j=last} =& \alpha' \boldsymbol{V}(h_0^{l-1} + h_3^{l-1} + h_6^{l-1}) + \beta' \boldsymbol{V}(h_1^{l-1} + h_4^{l-1}) \\
& + \gamma' \boldsymbol{V}(h_2^{l-1} + h_5^{l-1}) \\
& (\alpha', \beta', \gamma') := \frac{1}{M}(\alpha, \beta, \gamma) \\
& M = \sqrt{d}(3e^\alpha + 2e^\beta + 2e^\gamma)
\end{aligned}
\tag{9}
$$

That is, $a^l_{j=last}$ can be considered the weighted sum of 7 vectors ($h_0^{l-1}, \ldots, h_6^{l-1}$). Our assumption above suggests that $h_i^l$ can be decomposed into potential answers $An_k$. If all $h_i^{l-1}$ share a set of components (i.e., $\tilde{A}n_k$s), $a^l_{j=last}$ can be reorganized as follows.

$$
\begin{aligned}
a^l_{j=last} =& \Sigma_k \tilde{A}n_k \boldsymbol{V}(\alpha' + \beta' + \gamma') \\
& + \alpha' \Sigma_k An_k^Q \boldsymbol{V} + \beta' \Sigma_k An_k^S \boldsymbol{V} + \gamma' \Sigma_k An_k^A \boldsymbol{V} \\
\equiv& \frac{1}{\sqrt{d}} \Sigma_k \tilde{A}n_k \boldsymbol{V} \\
& + \alpha' \Sigma_k An_k^Q \boldsymbol{V} + \beta' \Sigma_k An_k^S \boldsymbol{V} + \gamma' \Sigma_k An_k^A \boldsymbol{V}
\end{aligned}
\tag{10}
$$

, where $An_k^Q$, $An_k^S$ and $An_k^A$ denote the components (i.e., potential answers stored in a FFN) of queries, separators and answer tokens that are not common across them. Also, we use $\alpha' + \beta' + \gamma' = 1/\sqrt{d}$. Eq. 10 suggests that the common component $\tilde{A}n_k$ can constructively accumulate the contributions from all token types and thus, $a^l_{j=last}$ is likely to be dominated by the common components. Interestingly, this would more likely hold when more examples are presented, which is consistent with the ICL property that more examples can make LLMs' answers more accurate. Conversely, if LLMs use common components to encode ICL examples, this may explain how LLMs improve answers with more ICL examples in given prompts.

One may ask, how do LLMs create common components for all ICL examples? We note that all examples in ICL prompts are expected to mediate the same context. For instance, examples of antonyms mediate semantic definitions of words, and Country-Capital examples mediate the geographic meanings of words. With this in mind, we hypothesize that LLMs use context-dependent embedding space to encode contextual information, which allows LLMs to support ICL. These context-dependent embeddings, which rely on a set of components, allow LLMs to use a subspace of full embedding space whose dimension equals the model dimension $d$. Notably, this is consistent with earlier studies suggesting that LLMs use subspace to perform various tasks [19, 20] and vectors to encode concepts [21, 22].

### 2.4 Three assumptions of the counting hypothesis

Based on the observations regarding ICL, we made three assumptions.

- **Assumption 1**: FFN's output can be decomposed into possible answers for a given input. Thus, in LLMs, these possible answers compete with one another, and LLMs track the likelihood of these answers in parallel due to residual connections.
- **Assumption 2**: As all examples are expected to be equivalent in terms of importance, attention scores can be approximated as constants.
- **Assumption 3**: LLMs can encode multiple answers using distinct subspace, each of which reflects different contextual meanings of inputs.

One may ask, is it possible to explain ICL using these assumptions? Let's consider an input prompt of two pairs of examples and one query: "Paris: France, London: UK, Washington DC:". During pretraining, LLMs learn that Paris can be associated with multiple entities including 'France' and 'Fashion' depending on the context. According to

4

Assumption 1, LLMs recall France and Fashion among other things. According to Assumption 3, these answers are encoded in different subspaces. For instance, France is encoded by using 'geographic' subspace, and Fashion is encoded by using 'art' subspace. Similarly, LLMs recall 'UK' associated with London using geographic subspace . According to Assumption 2, France and London have the same level of influence on the last colon. As the geographic subspace appears twice, the last colon favors the geographic subspace, and as a result, its answer leans towards 'US', which is encoded geographical subspace of Washington DC.

That is, LLMs can effectively count the strengths of contexts (or subspace) in the prompt to solve ICL. Thus, we refer to the set of three assumptions as the 'counting hypothesis'. We note that Assumption 2 may not accurately hold in pretrained LLMs [23, 24]. In this case, counting of subspace may become less accurate, but we expect ICL can still emerge (with increased error rates) as long as LLMs use context-dependent subspace .

## 3 Evidence supporting the counting hypothesis

Our counting hypothesis postulates that LLMs use context-dependent subspace to encode potential answers (i.e., next tokens). If true, subspaces of two consecutive tokens can overlap. That is, the residual streams of two consecutive tokens can be decomposed into the same components. Thus, to address our hypothesis, we use component analysis to analyze separator and answer tokens and compare them to determine whether they share the same component(s). Below, we discuss the analysis obtained by dictionary learning and Independent Component Analysis (ICA) [25][2].

The dictionary learning is chosen because it does not impose any specific conditions on the components, which are referred to as 'atoms' in dictionary learning. ICA is chosen because the $2^{nd}$, $3^{rd}$ and $4^{th}$ terms in Eq. 10, representing the sums of $An_i^Q$, $An_j^S$ and $An_j^A$, can be easily canceled out, if $An_i^Q$, $An_j^S$ and $An_j^A$ are independent of each other. That is, if $An_i$s are independent of one another, the common components $\tilde{An}_i$ would be promoted more strongly. Inspired by this, we assume that residual streams can be decomposed to independent components, and ICA is an ideal component analysis.

In this analysis we probe 6 pretrained LLMs publicly available: GPT-j-6B [27], Meta-Llama-3.1-8B[28], OLMo-2-0325-32B[29], Phythia-12B/Phythia-6.9B[30] and GPT-NEOX-20B[31]. All models are run with publicly available machine learning libraries, Pytorch [32] and Transformers [33].

### 3.1 Dictionary learning and ICA

Both dictionary learning and ICA decompose input $\boldsymbol{X} = [\vec{x_1}, \vec{x_2}, ..., \vec{x_m}]$ into some basis vectors $\vec{v_i}$s. They are referred to as atoms for dictionary learning and independent components for ICA. The dictionary learning minimizes reconstruction errors shown in Eq. 11

$$\epsilon = 0.5 \left\| \boldsymbol{X} - \boldsymbol{V}\boldsymbol{R} \right\|_2^2 + \lambda \left\| \boldsymbol{V} \right\|_1^1 \tag{11}$$

, where $\boldsymbol{V} \equiv [\vec{v_1}, \vec{v_2}, ..., \vec{v_n}]$ denotes atoms, and $\boldsymbol{R}$ denotes representations (i.e., coding coefficients). ICA aims to decompose $\boldsymbol{X}$ into non-Gaussian independent signal $\vec{s_i}$ (Eq. 12). In our analysis, we obtain signals $\vec{s_i}$ using FastICA implemented in sckit-learn [26].

$$\vec{x_i} = \boldsymbol{A}\vec{s_i} \tag{12}$$

As the magnitude of the residual streams grows with a deeper layer, we normalize them to have a unit norm and use dictionary learning and ICA to obtain atoms and independent components, which will be referred to as $v_i$ below. In this study, we collect the residual streams $h_i^l$ from all layers using 200 ICL prompts and use them to obtain $v_i$. With the estimated $v_i$, we evaluate the coding coefficients for each component ($v_i$) by minimizing the reconstruction error in Eq. 11 for both atoms and independent components.

### 3.2 Residual streams of separators and answer tokens

In ICL prompts used in this study, queries and answers are marked by 'Q:' and 'A:', and all input prompts end with 'A:'. Thus, the colon right after 'A:' indicates the expectation to provide answers corresponding to queries of prompts. For brevity, we refer to the colon (after A) as a separator below. In this analysis, we exclude the last separator token and consequently record residual streams of 5 separators and answer tokens after separators. We note that answer examples can be either single or multiple tokens, and thus, the total number ($m$) of answer tokens can be greater than that ($n$) of separator tokens. From the recorded residual streams, we create two 2D arrays. $A_S \in R^{n \times d}$ is obtained from separator

---

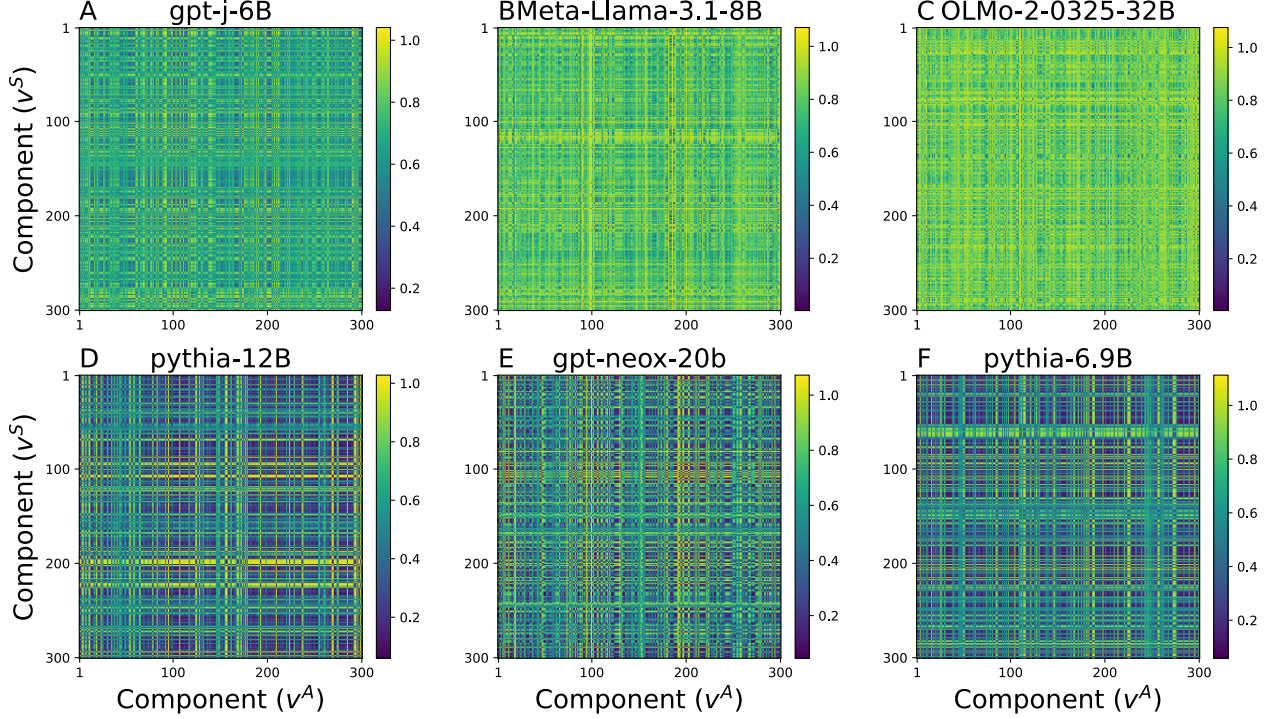[2]In this study, we use scikit-learn [26] to implement both dictionary learning and ICA.

Figure 1: Cosine distance ($D$) between 300 atoms obtained from dictionary learning. $x$-axis and $y$-axis denote $v_i^S$ and $v_j^A$.

tokens, and $A_A \in R^{m \times d}$ is obtained from answer tokens (between 'A:' and the next query tokens). $A_S$ and $A_A$ are decomposed using dictionary learning and independent component analysis (ICA) [25][3].

Using dictionary learning and ICA, we obtain components ($v_i^S \in R^d$ from $A_S$ and $v_j^A \in R^d$ from $A_A$) and compare cosine distance $D$ between $v_i^S$ and $v_j^A$ (Eq. 13).

$$D(v_i^S, v_j^A) = 1 - \frac{v_i^S \cdot v_j^A}{\left\| v_i^S \right\| \left\| v_j^A \right\|} \tag{13}$$

Figs. 1 and 2 show $D$ between $v_i^S$ and $v_j^A$. As shown in the figures, when we use dictionary learning, most distances are below 1, whereas the distances between some independent components (i.e., $v^S$, $v^A$ obtained by ICA) are close to 2. This is because ICA can only consider independency between the components and is not sensitive to the direction of the components. For instance, $-\vec{a}$ and $\vec{a}$ can be both described using the same set of independent components with opposite signs of the coding coefficient. Thus, when we use ICA, we use $D'(v_i^S, v_j^A) = min(D, 2 - D)$ to measure the similarity below.

We test all 6 models performing 6 ICL tasks and find the minimum distance $D$ using dictionary learning and $D'$ using ICA. As shown in Fig. 3, we observe the pairs of similar $v^S$ and $v^A$ in the majority of cases, supporting that the residual streams of separator and answer tokens share the common components. Also, $D$ and $D'$ appear to depend on the number of components (hyperparameters for both dictionary learning and ICA). Based on these results, we use 20 independent components and 300 atoms (components of dictionary learning) for further analysis, since $D$ and $D'$ are the smallest with them (Fig. 3).

### 3.3 Alignment between the last separators and answer tokens

If ICL relies on context-dependent subspace existing across all examples in an ICL prompt, the last separator, which predicts the answer to the test token, can use the same subspace to make a prediction. To test this possibility, we decompose the residual stream of the last separator using $v_i^A$ and evaluate the coding coefficients (Eq. 11). Figs. 4, 5

---

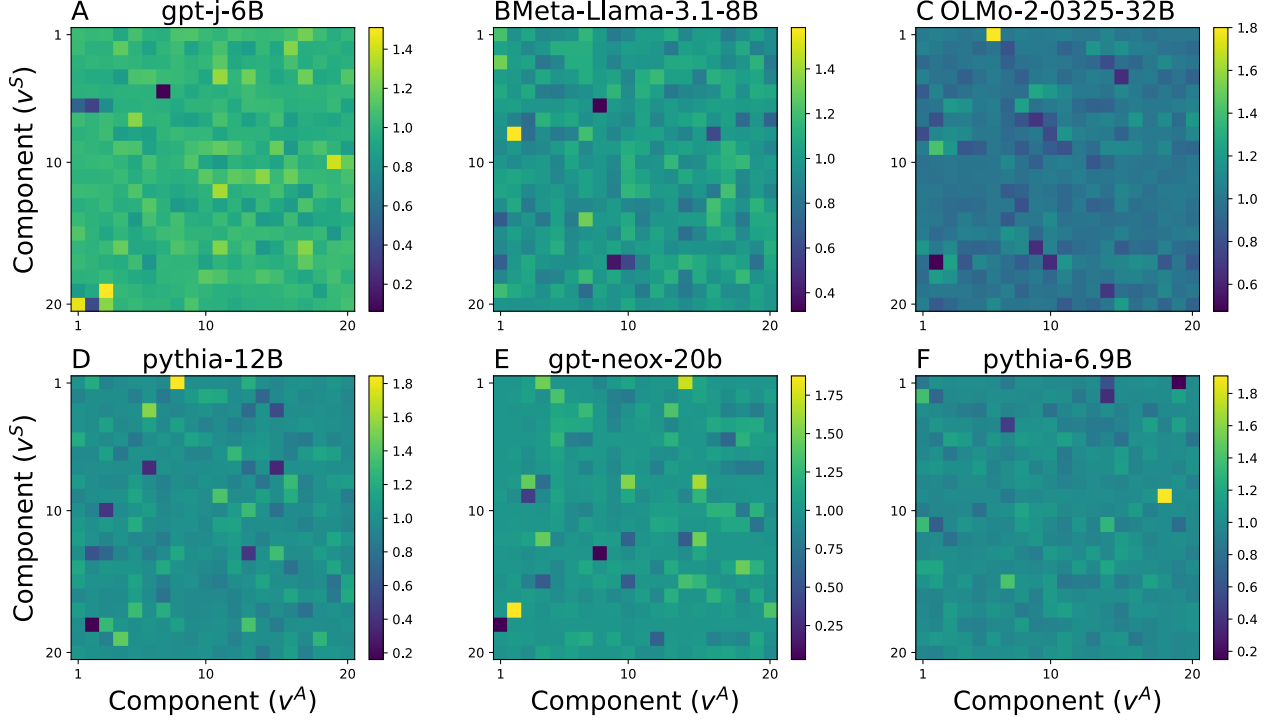[3]When ICA is used, $A_1^T$ and $A_2^T$ are used due to the convention of sckit-learn

Figure 2: Cosine distance $D$ between 20 independent components obtained from ICA. $x$-axis and $y$-axis denote $v_i^S$ and $v_j^A$, respectively.

show the coding coefficients of the last separator across all layers, which are obtained by dictionary learning and ICA, respectively.

As shown in the figures, we observe that the last separators' residual streams align with $v_i^A$. Further, we note that the degree of alignments measured via coding coefficients between the last separators and some of $v^A$ is increasing, as the layers go deeper, supporting that the last separator and answer tokens can share the same components, and thus, the last separator can utilize the same subspace to make a prediction.

Our hypothesis postulates that residual streams of answer tokens and separators share the same context and that the component encoding the common context can explain both answer tokens and separators. Conversely, if $v_i^A$ is close to one of $v_j^S$, this can be a good candidate for the descriptor of subspace used for ICL and the last separator should be well aligned.

We use this possibility to further test our hypothesis. Specifically we correlate the minimum distance $D_{min}$ and $D'_{min}$ (Eq. 14) and the coding coefficients for the residual streams of the last separator obtained from dictionary learning and ICA.

$$D_{min}(v_i^A) = min_j(D(v_j^S, v_i^A)), \text{for dict. learning}$$
$$D'_{min}(v_i^A) = min_j(D'(v_j^S, v_i^A)), \text{for ICA}$$

(14)

$D_{min}$ and $D'_{min}$ are small if any pair of $v_i^S$ and $v_j^A$ are close to each other. That is, if $D_{min}$ (or $D'_{min}$) for a given $v_i^A$ is small, $v_i^A$ can be a good descriptor of subspace, and thus, the coding coefficient is likely to be high. Fig. 6, shows correlations between the coding coefficient and $D_{min}$ and $D'_{min}$. As shown in the figure, when the coding coefficient of $v_i^A$ is large, $D'_{min}$ and $D_{min}$ are relatively small.

## 3.4 Implications of LLMs operations diagnosis

If LLMs support ICL using the context-dependent embedding, as our counting hypothesis suggests, we speculate that the embedding can be used to predict the accuracy of LLMs' operations. More specifically, we assume that the embeddings can be different depending on the accuracy of LLMs' answers. To address this assumption, we approximate the embeddings using the coding coefficients of $v_i^A$s and compare them when LLMs make correct or incorrect predictions.
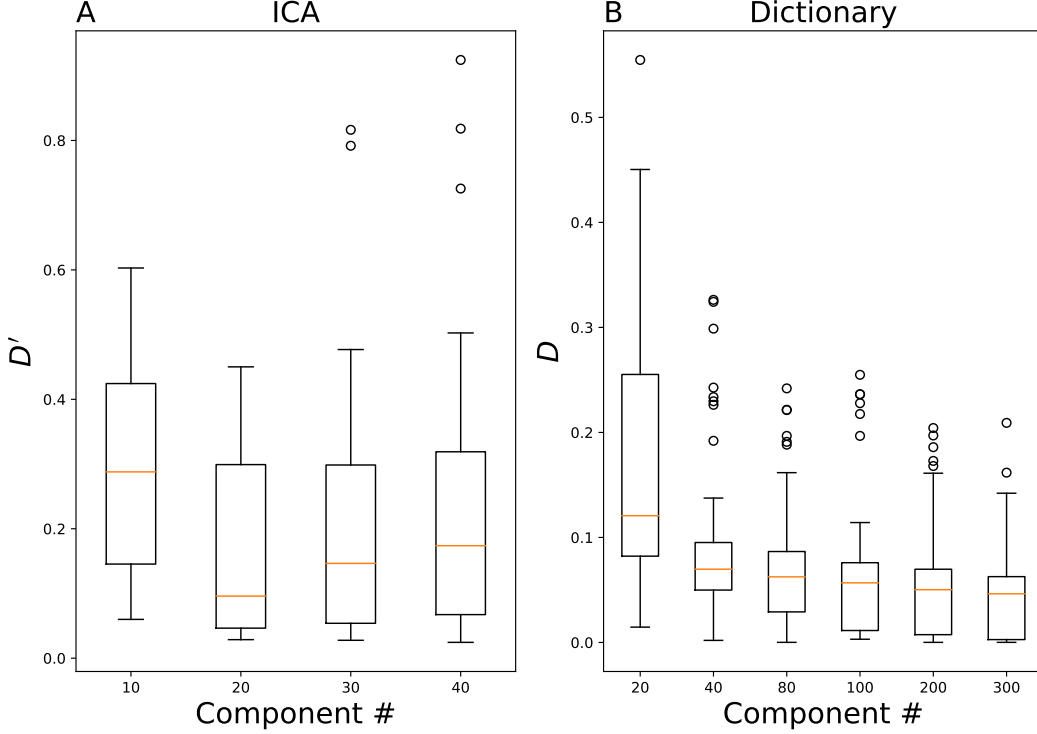
7

Figure 3: The minimum distance between $v_i^S$ and $v_j^A$. (A), $D'$ between independent components from separators and answer tokens. (B), $D$ estimated using dictionary learning.

In this experiment, we let LLMs to generate new tokens and record the embeddings of the first new tokens in the final layer via open-source deep learning library 'NNsight' [34]. The embeddings of the first new tokens are projected to 20 ICA components obtained above. Then, we compare the top-4 components of the residual stream of the last colon in the last layer by using the ratio $R$ of the coding coefficient $(C_1, C_2, C_3, C_4)$ for the top 4 components (Eq. 15).

$$R = \frac{C_2 + C_3 + C_4}{3C_1} \tag{15}$$

We estimate $R$s from 400 examples, independently generated from 200 examples for component analysis, and compare them according to the accuracy of LLMs' answers. If the 4-top components are strongly associated with LLMs' decision-makings, $R$ would be significantly different from correct to incorrect answers. Fig. 7 shows top-4 cases, in which $R$ is significantly different (t-test, $p < 0.05$) depending on the accuracy of answers (correct and incorrect).

We compare $R$ for all 6 tasks used above and 3 models. $p$-values observed in all these cases are presented in Table 1. As shown in the table, significantly different $R$ is observed throughout the datasets and the models, which indicates that the components of the embeddings are associated with LLMs' decision-making. This result raises the possibility that the decomposition of embeddings can be used to monitor LLMs' operations and make them more reliable.

# 4 Discussion

## 4.1 Related works

In this section, we briefly review the earlier works on ICL and the potential links between our work and them.

First, Xie et al. [4] proposed that ICL can be considered the implicit Bayesian inference. Specifically, they noted that LLMs may learn the distribution of latent concepts $\Theta$, and for a $\theta \in \Theta$, LLMs may generate output tokens using latent concept space: $p(O_1, O_2, \ldots, O_T) = p(O_1, O_2, \ldots, O_T | \theta)p(\theta)$. In principle, Bayesian inference rule allows us to $p(\theta | O_1, O_2, \ldots, O_T)$ from $p(O_1, O_2, \ldots, O_T | \theta)p(\theta)$. They provided the theoretical arguments to derive a formal
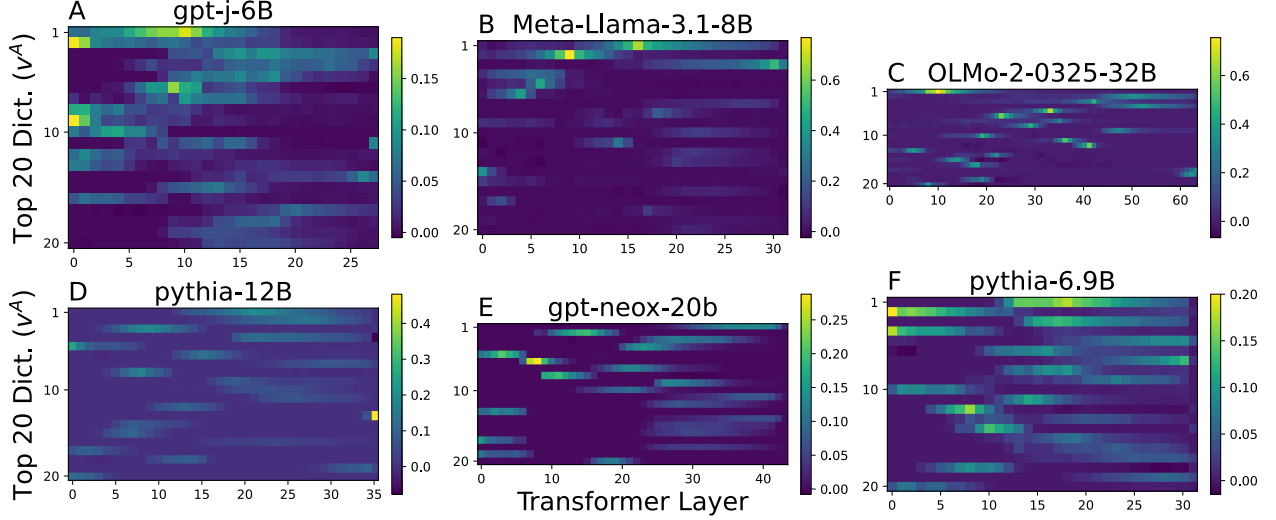
Figure 4: Alignment between the last separator and answer tokens, which is analyzed by dictionary learning. The residual streams are recorded from LLMs performing the 'country-capital' task. We choose the top 20 atoms out of 300 for 6 models and show their coding coefficients across layers shown in $x$-axis. $y$-axis denotes the top 20 atoms. The names of the models are displayed above the plots.
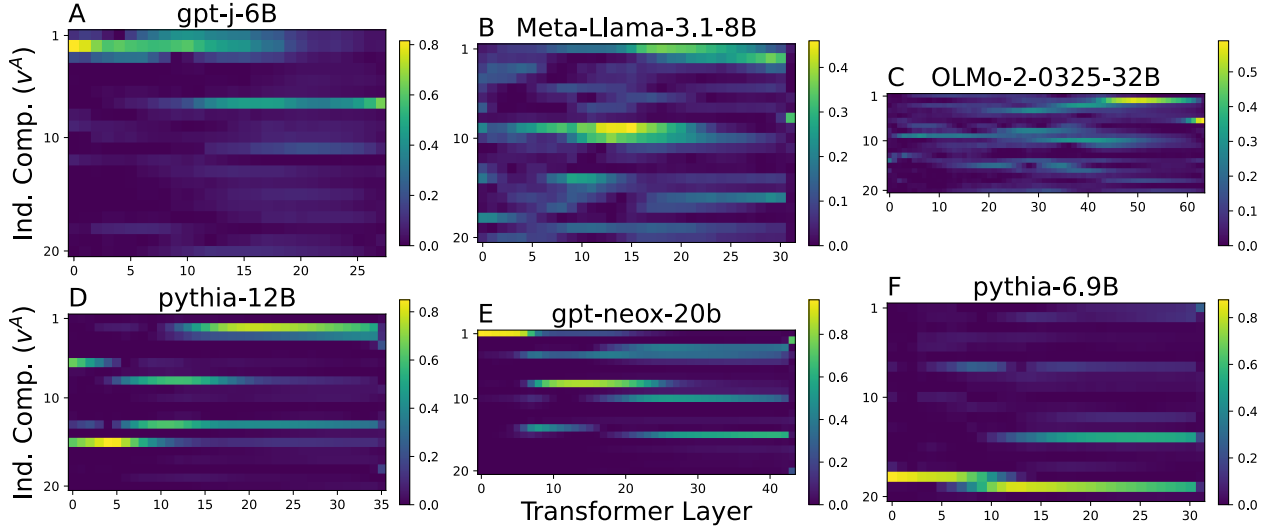


Figure 5: Alignment between the last separator and answer tokens, which is analyzed by ICA. The residual streams are recorded from LLMs performing the 'country-capital' task. Each plot shows the coding coefficients of 20 independent components for 6 models. $x$-axis and $y$-axis denote independent components and layers, respectively. As ICA is not sensitive to the component and consequently coding coefficient, we display the absolute value of coding coefficients. We present the alignment observed during different tasks in Supplemental Figs. S2, S3, S4, S5 and S6.

framework for ICL from Bayesian rule. The study, however, did not explain how LLMs can recover a latent concept $\theta$ from ICL examples. Our hypothesis suggests that LLMs use FFNs and SA layers to extract contextual information. As residual streams store them in parallel, simple vector operations can easily recover the contextual information from the residual streams, which can be associated with latent concept space $\Theta$ proposed in Xie et al. [4].

Second, another theoretical view is based on dual gradient updates [6, 7]. Dai et al. [7] presented the evidence supporting strong 'ICL-GD (gradient descent) correspondence' in pretrained LLMs by probing the outputs and updates of attention heads. Specifically, they found that the updates and outputs of attention heads are well aligned between ICL and fine-tuning. However, this evidence for ICL-GD correspondence was questioned by a more recent study [8]. Although it remains unclear if the dual gradient descent underlies ICL, it may explain how LLM refines subspace, which encodes contexts of inputs, to support ICL.
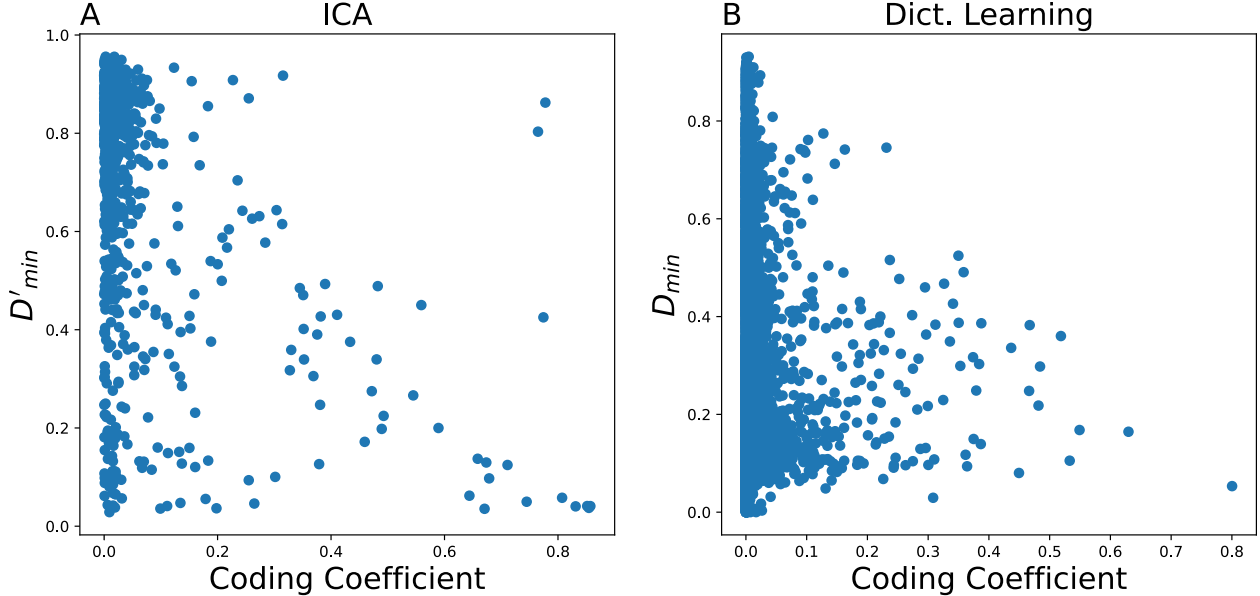
Figure 6: Correlations between $D_{min}/D'_{min}$ and the coding coefficient (Score). For all 6 models, we aggregate the coding coefficients of top 300 atoms and 20 independent components. For each atom and independent component, the minimum distance between a chosen component ($v_i^A$) of answer tokens and all of the components ($v_j^S$) of the separators. We use the absolute magnitude of the coding coefficient for independent components.

Table 1: p-values observed from 3 models and 6 ICL tasks. We compare $R$ when LLMs make correct predictions and when they make incorrect predictions. GPT, LLaMa and OLMO denote GPT-j-6B, Meta-Llama-3.1-8B and OLMo-2-0325-32B, respectively.

|  | GPT | LLaMa | OLMO |
|---|---|---|---|
| antonym | 0.013 | 0.113 | 0.025 |
| synonym | 0.001 | 0.863 | 0.958 |
| country-capital | 0.002 | 0.000 | 0.001 |
| english-french | 0.645 | 0.091 | 0.772 |
| product-company | 0.000 | 0.000 | 0.021 |
| person-sport | 0.784 | 0.000 | 0.000 |

Third, the two earlier studies [11, 12] noted that the hidden activity patterns can encode the task information of ICL tasks, suggesting that pretraining LLMs enables them to readily support ICL. Interestingly, our hypothesis suggests that LLMs use context-dependent subspace to encode multiple contextual information of input prompts in parallel. Consequently, if the hidden activity patterns are aggregated over multiple prompts, most non-relevant information would be cancelled out, and only the information associated with a task would remain, which may explain how task/function vectors can be obtained.

## 4.2 Limitations

We list two main limitations of our study. First, in section 2, we use a simple and ideal case of ICL prompts. For instance, we assume the invariance of the order of examples, which may not be always the case in practice, but our intention is not to formally prove that LLMs support the invariance of the order of examples or that more examples ensure more accurate answers in LLMs. Instead, using ideal arguments, we seek potential conditions, under which LLMs support the properties of and observations on ICL. In our analysis, this ideal case leads us to our counting hypothesis.

Second, We use common component analyses to address our hypothesis. Although ICA and dictionary learning provide supporting evidence, they may not be the most optimal methods to extract subspace that LLMs use to encode the inputs' contextual information. We note that residual streams are superpositions of many different components, and thus, the identified components (atoms for dictionary learning and independent components for ICA) can explain only a subset
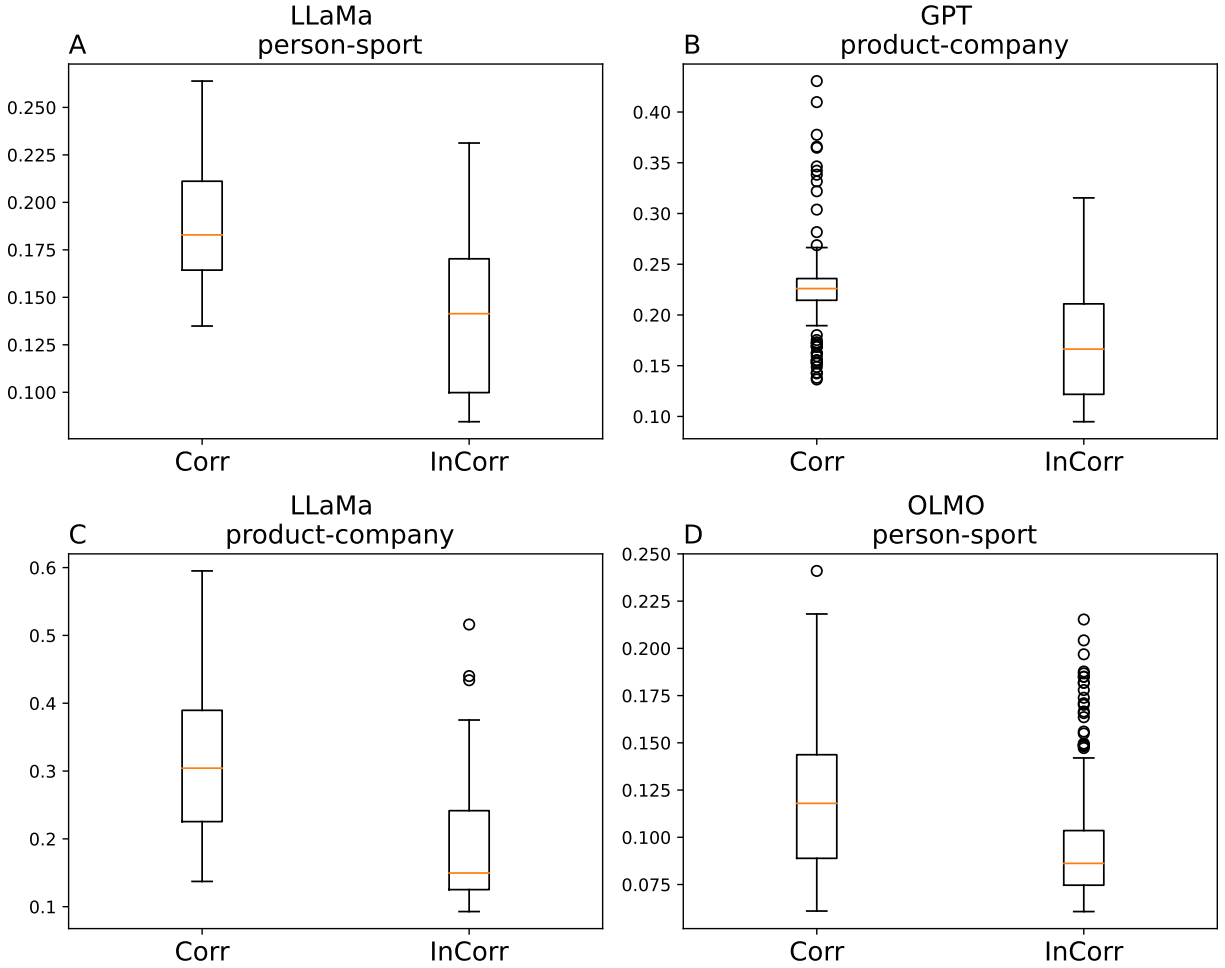
Figure 7: Comparing $R$. We choose 4 cases, in which $R$ is the most strikingly different between correct and incorrect predictions. The model and the task of 4 cases are displayed above the plots. GPT, LLaMa and OLMO denote GPT-j-6B, Meta-Llama-3.1-8B and OLMo-2-0325-32B, respectively.

of all components. In the future, we plan to explore effective algorithms to investigate LLMs' subspace associated with LLMs' decision-making.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

[1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.

[2] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022. Survey Certification.

[3] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. A survey on in-context learning. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 1107–1128, Miami, Florida, USA, November 2024. Association for Computational Linguistics.

[4] Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*, 2022.

[5] Fabian Falck, Ziyu Wang, and Christopher C. Holmes. Is in-context learning in large language models bayesian? a martingale perspective. In *Forty-first International Conference on Machine Learning*, 2024.

[6] Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org, 2023.

[7] Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4005–4019, Toronto, Canada, July 2023. Association for Computational Linguistics.

[8] Gilad Deutch, Nadav Magar, Tomer Natan, and Guy Dar. In-context learning and gradient descent revisited. In Kevin Duh, Helena Gomez, and Steven Bethard, editors, *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1017–1028, Mexico City, Mexico, June 2024. Association for Computational Linguistics.

[9] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.

[10] Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings*, 2013.

[11] Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

[12] Roee Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9318–9333, Singapore, December 2023. Association for Computational Linguistics.

[13] Jacob Dunefsky and Arman Cohan. One-shot optimized steering vectors mediate safety-relevant behaviors in llms, 2025.

[14] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html.

[15] Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2023.

[16] Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, 2021.

[17] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. *Advances in Neural Information Processing Systems*, 35, 2022.

[18] Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of LLMs: Stages of inference? In *ICML 2024 Workshop on Mechanistic Interpretability*, 2024.

[19] Xingyu Cai, Jiaji Huang, Yuchen Bian, and Kenneth Church. Isotropy in the contextual embedding space: Clusters and manifolds. In *International Conference on Learning Representations*, 2021.

[20] Yibo Jiang, Goutham Rajendran, Pradeep Kumar Ravikumar, Bryon Aragam, and Victor Veitch. On the origins of linear representations in large language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 21879–21911. PMLR, 21–27 Jul 2024.

[21] Haiyan Zhao, Heng Zhao, Bo Shen, Ali Payani, Fan Yang, and Mengnan Du. Beyond single concept vector: Modeling concept subspace in LLMs with gaussian distribution. In *The Thirteenth International Conference on Learning Representations*, 2025.

[22] Nedjma Ousidhoum, Xinran Zhao, Tianqing Fang, Yangqiu Song, and Dit-Yan Yeung. Probing toxic content in large pre-trained language models. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *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 4262–4274, Online, August 2021. Association for Computational Linguistics.

[23] Qi Guo, Leiyu Wang, Yidong Wang, Wei Ye, and Shikun Zhang. What makes a good order of examples in in-context learning. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 14892–14904, Bangkok, Thailand, August 2024. Association for Computational Linguistics.

[24] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland, May 2022. Association for Computational Linguistics.

[25] A. Hyvärinen and E. Oja. Independent component analysis: algorithms and applications. *Neural Netw.*, 13(4–5):411–430, May 2000.

[26] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

[27] Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. `https://github.com/kingoflolz/mesh-transformer-jax`, May 2021.

[28] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun

Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The llama 3 herd of models, 2024.

[29] Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2 OLMo 2

Furious, 2024.

[30] Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Moham-mad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.

[31] Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. Gpt-neox-20b: An open-source autoregressive language model, 2022.

[32] Adam Paszke, Sam Gross, Soumith Chintala, Edward Chanan, Gregory Yang, Zachary DeVito, Alban Lin, Zeming Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.

[33] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface's transformers: State-of-the-art natural language processing, 2020.

[34] Jaden Fiotto-Kaufman, Alexander R Loftus, Eric Todd, Jannik Brinkmann, Caden Juang, Koyena Pal, Can Rager, Aaron Mueller, Samuel Marks, Arnab Sen Sharma, Francesca Lucchetti, Michael Ripa, Adam Belfki, Nikhil Prakash, Sumeet Multani, Carla Brodley, Arjun Guha, Jonathan Bell, Byron Wallace, and David Bau. Nnsight and ndif: Democratizing access to foundation model internals. 2024.

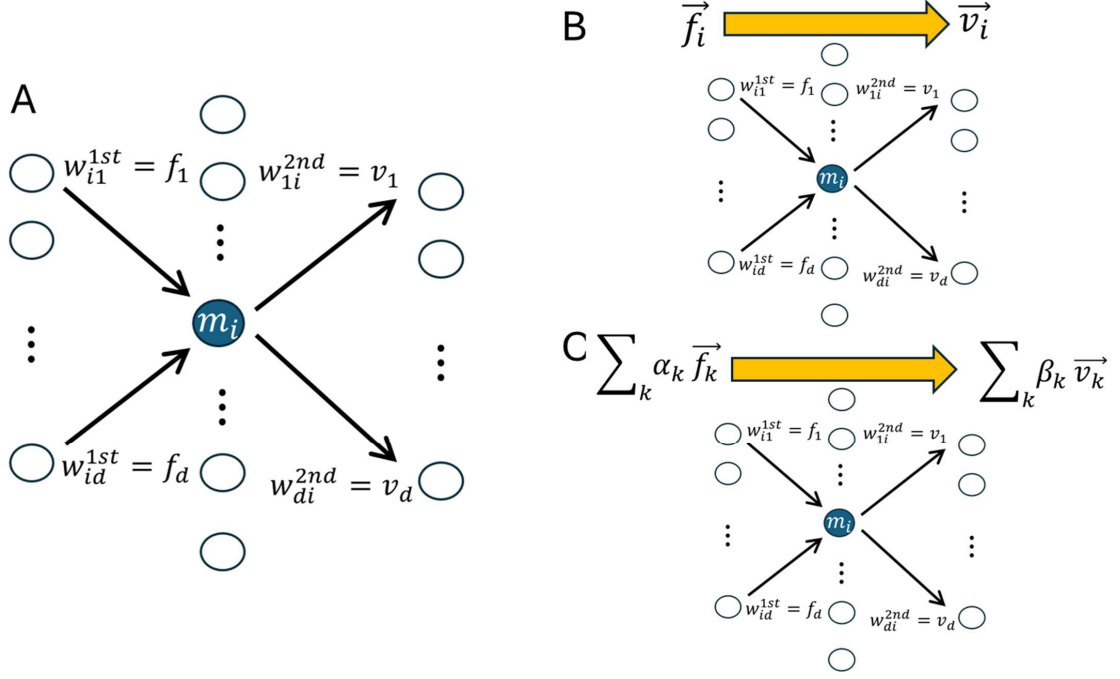## Supplemental Figures

Below, we list supplemental figures.



Fig. S1: FFN as an associative memory (A), schematics. (B), when an input is one of the features. (C), when an input is a linear combination of features.
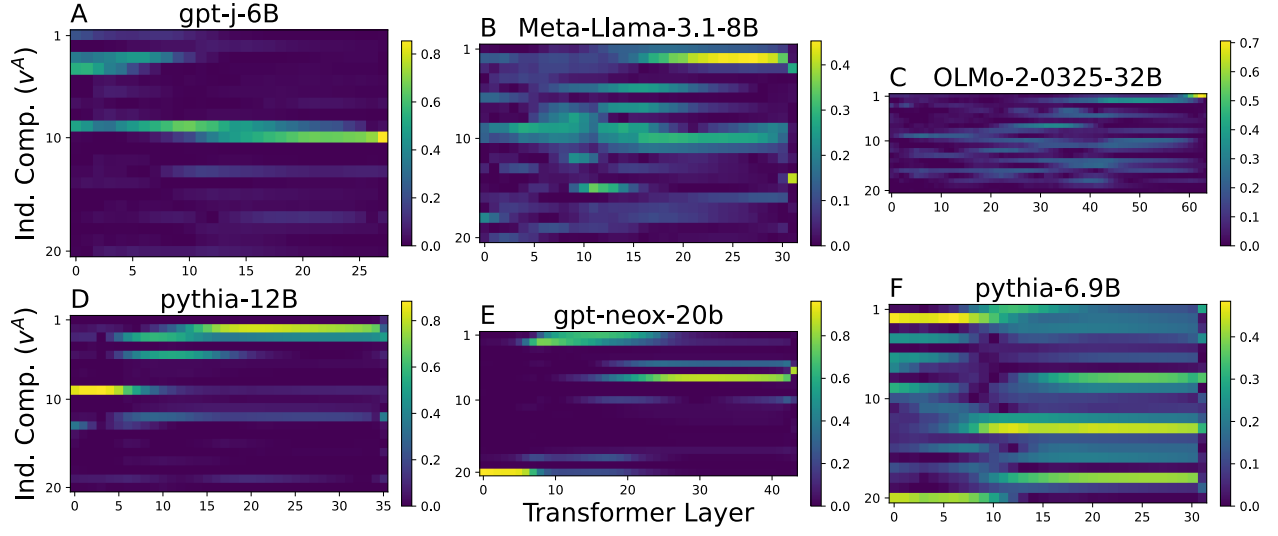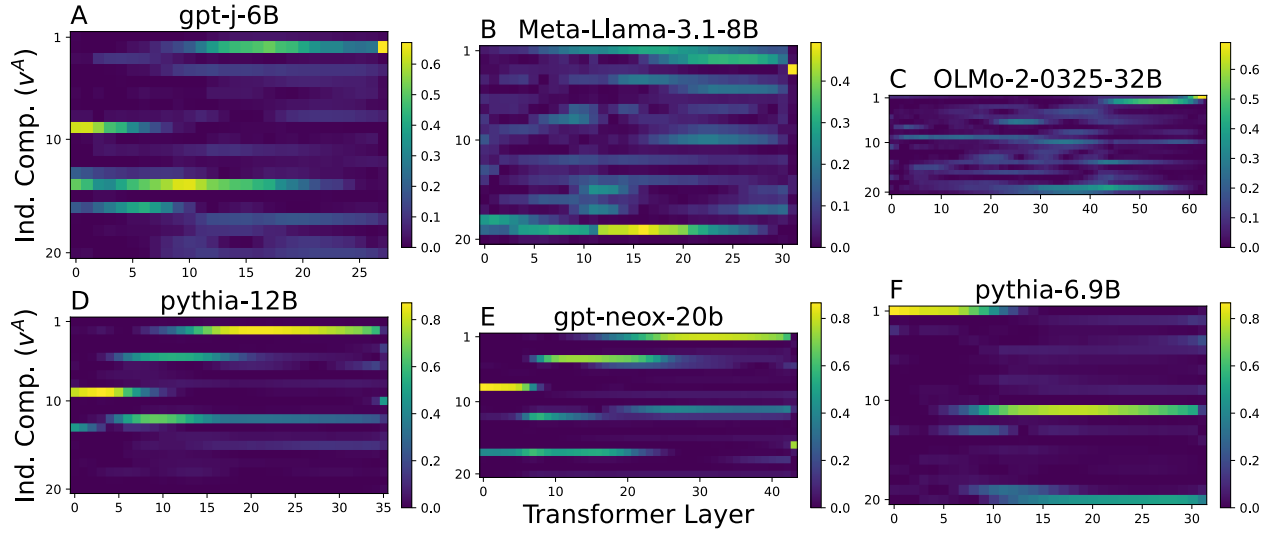
Fig. S2: Same as Fig. 5, but task is antonym.


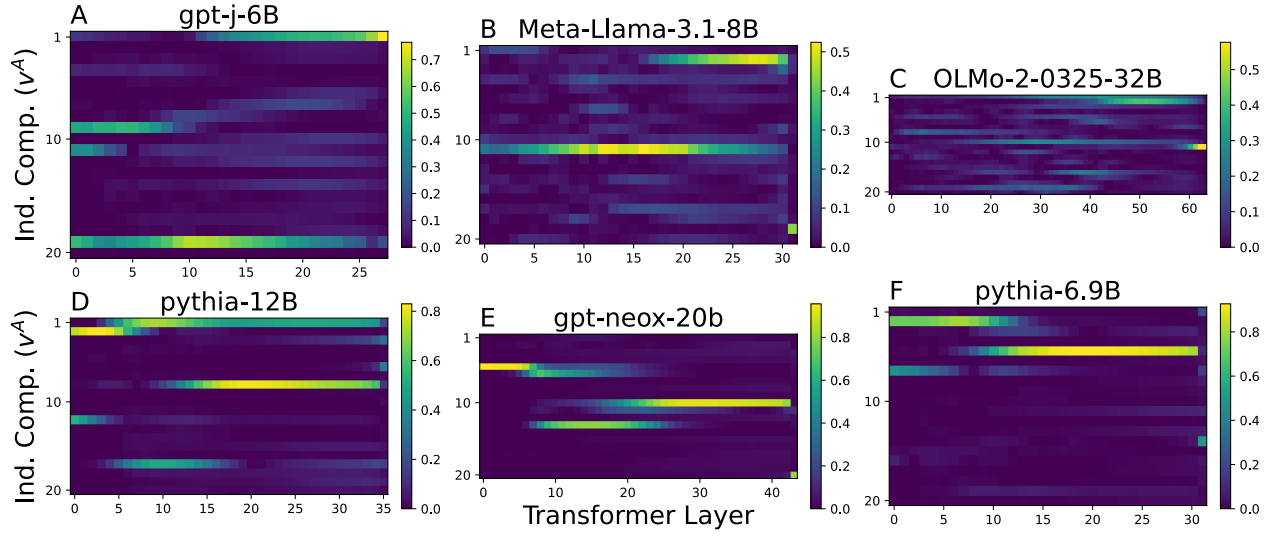
Fig. S3: Same as Fig. 5, but task is English-French.
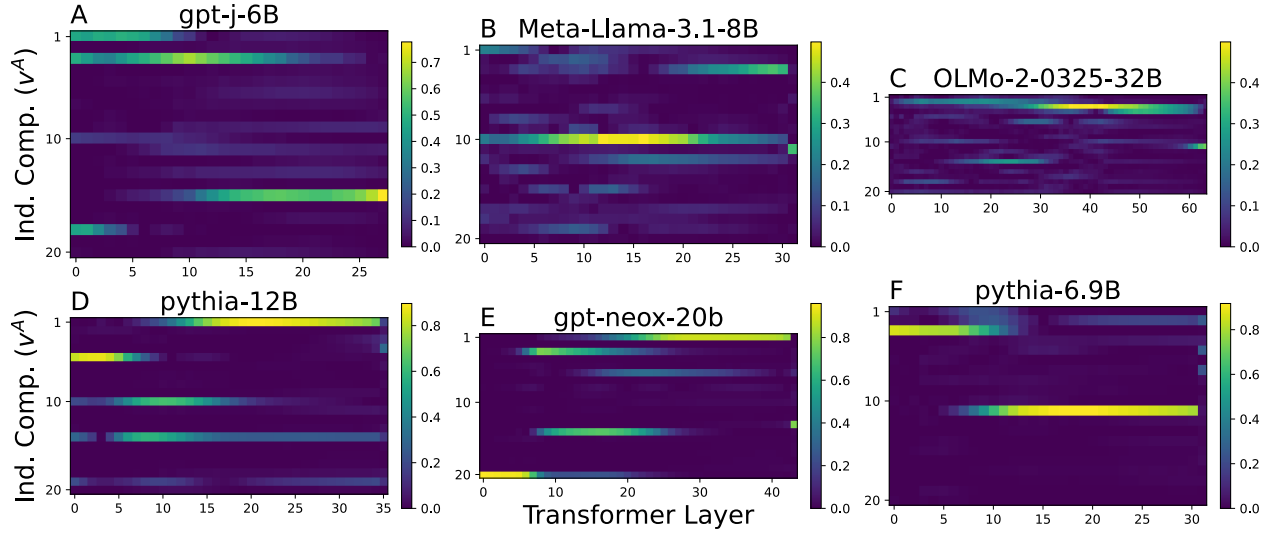
Fig. S4: Same as Fig. 5, but task is person-sport.
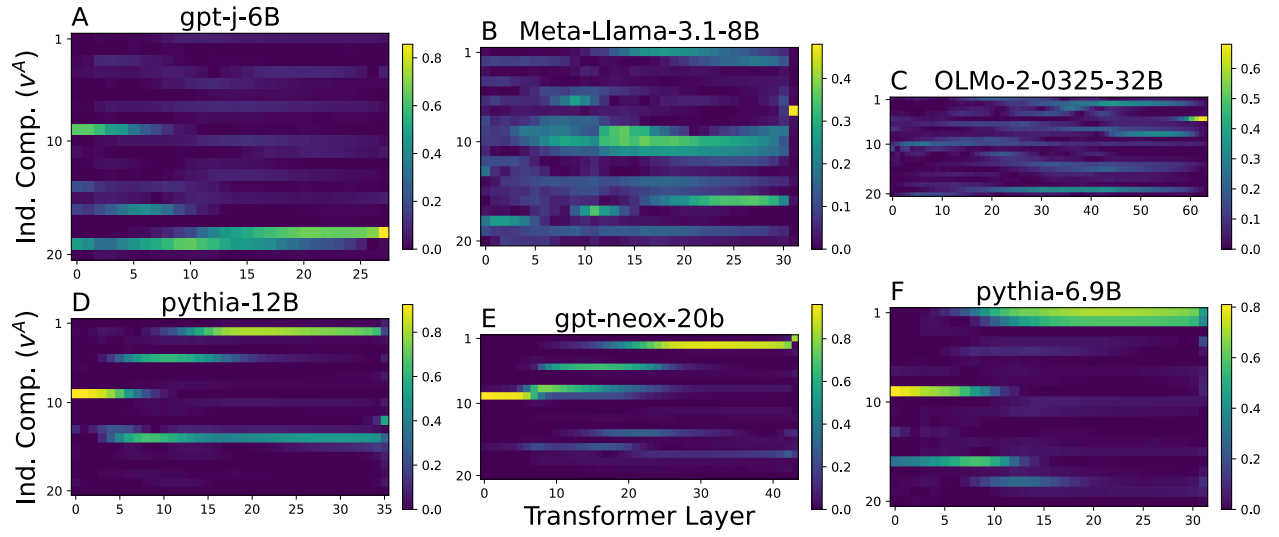


Fig. S5: Same as Fig. 5, but task is product-company.

Fig. S6: Same as Fig. 5, but task is synonym.