HIERARCHICAL ANALYSIS: MONOTONICITY OF LAYER PERFORMANCE IN LARGE LANGUAGE MOD ELS

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce a quantitative framework to evaluate how Large Language Models (LLMs) learn tasks across all layers, revealing a 'monotonicity phenomenon'. Specifically: i) performance at each layer consistently improves from one layer to the next on the pre-training set, and ii) this improvement is consistently observed across various downstream tasks. This monotonicity phenomenon indicates that LLMs effectively capture complex hierarchical features across diverse datasets. For example, our study on the abstraction of concepts using linear representations in word embeddings shows that the clarity of these abstractions progressively increases with each layer. Finally, by leveraging this monotonicity, we can significantly reduce inference time and memory requirements by selecting the most appropriate layer, thereby enhancing the efficiency of LLMs in real-world applications.

024 025

026

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

Large Language Models (LLMs) have achieved remarkable success across a wide range of domains (Brown, 2020; Bubeck et al., 2023; Chowdhery et al., 2023). Nevertheless, they face significant challenges related to computational and memory demands during inference. To address these challenges and enhance inference efficiency, various techniques have been proposed, such as quantization (Liu et al., 2023), pruning (Sun et al., 2023), and weight sparsification (Frantar & Alistarh, 2023). In addition to these efficiency concerns, LLMs are frequently criticized for their "black box" nature, which has led to numerous studies (Park et al., 2023; Tigges et al., 2023) aiming to investigate and shed light on the underlying mechanisms of how these models function.

Embedding learning plays a crucial role in understanding why LLMs are effective (Park et al., 2023;
Tigges et al., 2023; Hernandez et al., 2023; Yan et al., 2024). For instance, (Yan et al., 2024) demonstrated a sequential development pattern in which cognitive abilities are primarily established during
pretraining, while expressive abilities are mainly refined through supervised fine-tuning (SFT) and
reinforcement learning with human feedback (RLHF). Additionally, embedding learning has been
successfully applied across various neural networks with diverse architectures (Kim et al., 2018;
Yang & Hu, 2020; Vyas et al., 2024).

In the context of embedding learning in LLMs, the Linear Representation Hypothesis is a crucial concept, suggesting that high-level ideas are represented linearly within the model's representation space. This linear representation indicates that LLMs have the capability to comprehend textual meaning. For example, (Tigges et al., 2023) showed that sentiment is represented linearly in the output layers of LLMs, where a single direction in activation space effectively captures this feature across multiple tasks, with one end corresponding to positive sentiment and the other to negative.

The output from the last layer is typically used as the embedding for downstream tasks (Devlin, 2018; Dosovitskiy, 2020; Radford et al., 2021). However, for a given downstream task, it's worth questioning whether the last layer output of pre-trained models might be overly comprehensive, and whether the outputs from the middle layers could be sufficient as the embedding for that task.
Similar to CNNs, where earlier layers capture fundamental visual features like lines and curves, and later layers extract more abstract representations, it is worth exploring whether each layer of pre-trained models could serve as an effective embedding for various tasks.



Figure 1: Monotonicity and Its Applications: (a) Layerwise accuracy of GPT-2 at different stages of training, using the OpenWebText dataset (Gokaslan & Cohen, 2019); (b) Demonstration of the monotonicity phenomenon in three pre-trained LLMs on the Banking77 dataset; (c) Comparison of the performance of the three LLMs after selecting the most suitable layer embedding on Banking77 dataset.

1.1 CONTRIBUTIONS

 In this paper, we conduct a hierarchical analysis of LLM performance. We introduce a quantitative and precise characterization of how LLMs learn tasks across all layers by defining the concept of "layerwise performance". Our contribution is as follows:

- Our results show that layerwise performance on pre-trained data progressively improves from one layer to the next, a phenomenon we refer to as "monotonicity". Notably, this phenomenon is absent at the initialization stage of LLMs but becomes increasingly evident as training progresses, as illustrated in Figure 1 (a);
- We then show that for pre-trained LLMs, monotonicity is not limited to the pre-training datasets but also extends to other datasets, indicating that layerwise performance improves progressively across layers even when evaluated on new data. This observation indicates that the ability of LLMs to capture and refine complex features is not limited to the data they were initially trained on but is a more generalizable characteristic. Our experiments provide strong evidence of monotonicity in several widely used LLMs, including Llama 3.1 Instruct (Dubey et al., 2024), Phi-3.5 (Abdin et al., 2024), and GPT-Neo (Black et al., 2022), when tested on diverse datasets such as Banking77 (Casanueva et al., 2020), AG News (Zhang et al., 2015), and the Twitter Sentiment Analysis dataset (Kharde et al., 2016), as illustrated in Figures 1 (b) and 2. To further understand this consistent pattern, we utilize the concept of linear representation, which helps explain why monotonicity reliably emerges across different datasets and model architectures. This analysis not only highlights the robustness and adaptability of monotonicity within LLMs but also offers deeper insights into how these models maintain their performance across various tasks and data domains;
- · The concept of monotonicity offers a practical approach to reduce inference time and mem-ory requirements by allowing us to select the most appropriate layer embedding for a given task. When applying this idea to a specific downstream task, we begin by assessing the layerwise performance of the LLMs. If the layerwise performance gains in the final layers start to plateau, it indicates that the model's capacity is likely sufficient for the task. In this scenario, we can identify the layer at the inflection point—where performance gains become marginal-and use this as the cutoff for inference. For instance, as illustrated in Figure 1 (b), the optimal layers are found to be the 5th for Llama 3.1 and the 9th for Phi-3.5. By limiting inference to only the necessary preceding layers, as shown in Figure 1 (c) and 4, we can substantially accelerate the inference process without sacrificing per-formance. This selective approach enables a more efficient use of LLMs, particularly for resource-constrained applications.

108 1.2 RELATED WORKS

110 **Hierarchical analysis** Exploring the intermediate layers of neural networks and LLMs is essential 111 for gaining a deeper understanding of these models and identifying potential issues (Alain, 2016; He 112 & Su, 2024). (Alain, 2016) investigated how features evolve across layers in neural networks by applying linear classifiers independently to each layer's features, revealing that the linear separability 113 of features increases monotonically with depth, highlighting the importance of deeper layers in en-114 hancing the model's ability to differentiate data points. In a similar vein, (He & Su, 2024) introduced 115 a quantitative law describing how contextualized token embeddings are learned through intermediate 116 layers in pre-trained LLMs for next-token prediction, a task relevant to our context. Their findings 117 show that prediction accuracy consistently improves from one layer to the next, confirming our main 118 result.

119 120

Linear representation in LLMs The concept of linear representation in LLMs pertains to the 121 idea that high-level features can be expressed in a linear manner, a topic that has been explored in 122 several previous studies (Tigges et al., 2023; Park et al., 2023; Jiang et al., 2024). (Park et al., 2023) 123 formalized this notion by examining how high-level concepts are linearly represented in LLMs, 124 demonstrating that these representations are closely linked to interpretability and control through a 125 non-Euclidean inner product that aligns with the underlying language structure. In a complemen-126 tary study, (Jiang et al., 2024) investigated the origins of linear representations in LLMs, revealing 127 that both the softmax objective and the implicit bias introduced by gradient descent promote the 128 formation of linear structures in concept representations, as demonstrated through a latent variable 129 model.

130

131 **Layer pruning** Several studies have explored techniques to streamline large language models by 132 addressing layer redundancy and optimizing computation. Research on layer pruning focuses on 133 identifying and removing redundant or unimportant layers, thereby reducing the number of parame-134 ters and accelerating inference without significant performance degradationSong et al. (2024); Kim 135 et al. (2024); Chen et al. (2024); Men et al. (2024). Additionally, LayerSkip Elhoushi et al. (2024) introduces a specialized training method that enables the use of only the first half of the model's 136 137 layers while maintaining high performance. These works collectively highlight the importance of addressing redundancy in LLMs to improve efficiency. 138

139 140

141 142

143

144

2 MONOTONICITY IN LLM

Given a sequence of tokens as input, LLMs function as nonlinear models that iteratively transform the token embeddings into new sequences at each layer using attention mechanisms and other operations.

145 146 147

2.1 LAYERWISE PERFORMANCE ON PRE-TRAINED DATA

148 149 Consider an LLM composed of L transformer layers, with the pre-training data consisting of S input 150 sequences. Each input sequence, denoted as X_s for $1 \le s \le S$, is composed of T_s tokens. For each 151 token position $1 \le t \le T_s$, let x_t^s represent the t-th token within the sequence X_s , expressed as a 152 one-hot encoded vector to capture the token's identity in a high-dimensional space.

During the pre-training phase, the LLM is tasked with learning to predict the next token $x_{T_s}^s$ based on the sequence of preceding tokens $x_1^s, x_2^s, ..., x_{T_s-1}^s$. This predictive process enables the model to gradually capture complex patterns and dependencies inherent in the data. As the input sequence progresses through each of the *L* transformer layers, richer and more abstract representations of the tokens are formed.

For any given layer $1 \le l \le L$, let h_l^s denote the embedding corresponding to the final token in the *l*-th layer for the sequence $x_1^s, x_2^s, ..., x_{T_s-1}^s$. This embedding serves as a refined representation of the input sequence after processing through the first *l* layers. Consequently, this iterative embedding process across all layers results in the formation of a series of datasets $\{\mathcal{D}_l\}_{l=1}^L$, where each dataset \mathcal{D}_l is defined as $\mathcal{D}_l = \{(h_l^s, x_{T_s}^s) \mid 1 \le s \le S\}$. 162 These datasets collectively represent L distinct classification tasks, one for each layer, effectively 163 capturing how the model's ability to encode and predict token information evolves as it progresses 164 through the layers. This layered structure offers insights into the intermediate representations 165 learned by the LLM during the pre-training process, highlighting how each layer contributes to 166 refining the model's understanding of sequential data.

167 Given the L distinct datasets $\{\mathcal{D}_l\}_{l=1}^L$ generated from the pre-trained LLM, we can evaluate the 168 performance of each layer of the LLM through the following procedure. First, we split each dataset 169 \mathcal{D}_l into a training set and a testing set to facilitate model evaluation. For each layer l, we then apply 170 a logistic regression model to the training set of dataset \mathcal{D}_l to learn the relationship between the 171 layer's embedding and the corresponding target tokens.

172 Once the logistic regression model is trained on this training set, we evaluate its predictive accuracy 173 using the testing set. The accuracy obtained from this evaluation serves as a quantitative measure 174 of the layerwise performance, reflecting how well the embeddings produced by the *l*-th layer of 175 the LLM capture the necessary information for predicting the next token. By repeating this pro-176 cess for all L layers, we gain a comprehensive understanding of how the predictive capability of 177 the LLM evolves across different layers, providing insights into the effectiveness of each layer's 178 representations in contributing to the overall learning process.

- 179 180
- 181
- 182 183

2.2 MONOTONICITY OF LAYERWISE PERFORMANCE IN GPT-2

184 185

191

186 To investigate the layerwise performance in LLMs, we conducted an analysis using the widely 187 adopted GPT-2 model (Radford et al., 2019) in conjunction with the OpenWebText dataset. Specifi-188 cally, we trained the small variant of GPT-2 using this dataset to observe how layerwise performance 189 evolves over time during the training process. At various training stages, specifically at epochs 0, 190 1,500, 5,000, and 165,000, we recorded the performance of each layer to gain insights into the model's learning progression, presented in Figure 1 (a). 192

Figure 1 (a) illustrates the layer-wise performance of a pre-trained GPT-2 model, demonstrating a 193 clear increase in testing accuracy across layers, particularly at epoch 165,000. This trend, where 194 the accuracy progressively improves from the lower layers to the upper layers, is referred to as 195 monotonicity. As the layers deepen, the model's ability to capture and represent complex patterns 196 in the data becomes more pronounced, leading to higher accuracy rates. This phenomenon suggests 197 that deeper layers in the fully trained GPT-2 model contribute more significantly to the model's 198 overall predictive capability, reinforcing the importance of multi-layered architectures in LLMs like 199 GPT-2.

200 Furthermore, Figure 1 (a) captures the evolution of this monotonicity throughout the training pro-201 cess. Notably, at the initialization stage (epoch 0), the model exhibits no clear monotonic trend, 202 with the testing accuracy remaining relatively low and flat across all layers. As training progresses, 203 however, monotonicity becomes increasingly evident, with significant gains in layerwise perfor-204 mance emerging by epochs 1,500, 5,000, and ultimately, 165,000. This progressive emergence of 205 monotonicity underscores how the training process enables the model to develop a more refined un-206 derstanding of language, leading to enhanced representation and predictive capabilities as training 207 deepens. This observation highlights that monotonicity is not an inherent characteristic of LLMs at initialization but rather a phenomenon that emerges and strengthens as the model learns from data 208 over time. 209

210 The concept of monotonicity has been explored in prior research (Alain, 2016; He & Su, 2024). For 211 instance, He et al. (2024) applied a least squares fit on the dataset \mathcal{D}_l and introduced the concept of 212 the prediction residual (PR) to measure the LLM's capability for next-token prediction, with their 213 findings closely aligning with our results. Moreover, the phenomenon of monotonicity is not limited to LLMs; it can also be observed in neural networks with different architectures, such as FCNs and 214

CNNs (Alain, 2016). This suggests that monotonicity is a fundamental characteristic of hierarchical 215 network models trained by gradient descent algorithm.

216 3 **MONOTONICITY ACROSS DATASETS** 217

218 The previous section demonstrated that monotonicity is present in pre-trained data and becomes pro-219 gressively more pronounced as training advances. Given that pre-trained LLMs are widely favored 220 for their strong generalization abilities on downstream tasks and datasets, it raises an intriguing question: Does monotonicity persist when the pre-trained model is applied to new tasks? Ad-222 dressing this question is crucial, as it can provide valuable insights into the behavior and adaptability of LLMs, informing their effective application across a broader range of real-world tasks.

224 225

226

221

3.1 LAYERWISE PERFORMANCE ACROSS DATASETS

227 Consider a pre-trained LLM with L transformer layers, and a new dataset consisting of S input sequences $\{X_s\}_{s=1}^S$ along with their corresponding labels $\{Y_s\}_{s=1}^S$. For any layer $1 \le l \le L$, let h_l^s 228 denote the embedding of the final token in the *l*-th layer for the input sequence X_s . This embedding 229 serves as a progressively refined representation of the input sequence as it is processed through the 230 first l layers of the LLM. Through this process, we construct a series of datasets $\{\mathcal{D}_l\}_{l=1}^L$, where 231 each dataset \mathcal{D}_l is defined as $\mathcal{D}_l = \{(h_l^s, Y_s) \mid 1 \le s \le S\}.$ 232

233 Similarly, to evaluate the performance of each layer, we further divide each dataset \mathcal{D}_l into separate 234 training and testing sets. For each layer l, a linear model—either logistic regression for discrete labels or linear regression for continuous labels—is then fitted to the training set of dataset \mathcal{D}_l and 235 evaluates the linear model on the testing set, which is considered as the **layerwise performance** 236 **across the dataset**. This approach allows us to learn and assess the relationship between the layer's 237 embeddings and the corresponding labels, providing insights into how well the representations at 238 each layer capture the relevant features for the given task. 239

- 240
- 3.2 MONOTONICITY ACROSS DATASETS 241
- 242

To investigate the layerwise performance across datasets, we conducted an analysis using some 243 famous LLMs including Llama 3.1 Instruct, Phi-3.5 and GPT-Neo, and in conjunction with some 244 NLP datasets including AG News, Banking77 and Twitter Sentiment Analysis dataset. The results

245 is presented in Figure 2. 246

Figure 2 demonstrates that all three pre-trained LLMs-Llama 3.1 Instruct, Phi-3.5, and GPT-247 Neo-exhibit the phenomenon of monotonicity across the three NLP datasets: AG News, Bank-248 ing77, and the Twitter Sentiment Analysis dataset, indicating their ability to generalize effectively 249 to new tasks. monotonicity, which we first observed in pre-training data, is characterized by the 250 progressive improvement in accuracy as data passes through successive layers of the model. The 251 fact that this trend continues when the LLMs are applied to entirely new datasets suggests that 252 the representational power of these models, acquired during pre-training, transfers well to differ-253 ent downstream tasks. This generalization ability is crucial because it highlights the adaptability 254 of LLMs to learn complex features even outside the original training domain, confirming that their 255 deeper layers can refine and enhance understanding across various contexts and tasks, beyond what was learned during pre-training. 256

257 When examining the degree of monotonicity for a given pre-trained LLM, it becomes evident that 258 this phenomenon varies in intensity. For instance, in some cases, the performance improvement 259 across layers is relatively flat in the last few layers, while in others, it is steep. A flatter monotonicity 260 curve may imply that the model's current size is sufficient for effectively capturing the information 261 needed for the task, suggesting that adding more layers might yield diminishing returns in terms of performance gains. Conversely, a steeper monotonicity curve suggests that the model is still gaining 262 significant representational power as layers are added, indicating that the current model size might 263 not be large enough to fully capture the complexity of the task at hand. Therefore, the shape of the 264 monotonicity curve provides insights into whether a model has reached its optimal capacity or if 265 there is room for further improvement with additional layers. 266

When comparing monotonicity across different models on the same dataset, it is apparent that each 267 LLM exhibits unique learning capabilities. For instance, the Llama 3.1 Instruct model shows a dif-268 ferent pattern of layerwise performance compared to Phi-3.5 and GPT-Neo, reflecting their varying abilities to extract and understand the semantic meaning of sentences within each dataset. These dif-



Figure 2: Monotonicity Across Datasets. The x-axis corresponds to the layer index, while the y-axis represents the accuracy for both training and testing sets. Monotonicity is consistently present across the various new tasks. However, this phenomenon is not uniform; it differs based on the interaction between each dataset and the specific pre-trained LLM being used.

324 ferences in monotonicity indicate that some models are more adept at learning the relevant features 325 for a given task, while others may require more layers or different architectures or even different 326 training processes to achieve similar levels of performance. Thus, Figure 2 not only highlights the 327 existence of monotonicity but also underscores how this phenomenon varies depending on the model 328 and inherent ability to generalize across different NLP tasks.

329 330

331

349 350 351

352

358

359

360 361

3.3 EXPLANATION OF MONOTONICITY ACROSS DATASETS

332 **Linear representation** The concept of linear representation in LLMs refers to the idea that high-333 level features and relationships between words can be captured and expressed in a linear fashion, 334 an area that has been investigated in several studies (Tigges et al., 2023; Park et al., 2023; Jiang 335 et al., 2024). A classic example of this phenomenon is found in word embeddings, where it has 336 been empirically observed that pairs like embedding("woman") - embedding("man") and embed-337 ding("queen") - embedding("king") are nearly parallel in the vector space, indicating that the model has effectively captured the underlying semantic relationship between these words (Mikolov et al., 338 2013). This parallelism suggests that the embedding space represents consistent linear relationships 339 that encode concepts such as word inflection, gender or other analogies. 340

341 In our study, we applied the idea of linear representation to LLMs by examining whether similar 342 semantic pairs maintain such linear relationships across different layers of the model. For each pair 343 of related words, such as word A and word B, we extracted their embeddings from each layer of 344 the LLM and computed the vector difference between them. We then calculated the inner product of these differences across all similar word pairs. If the resulting inner product is close to 1, this 345 indicates that the embeddings effectively capture the intended concept, contributing to the model's 346 ability to understand the meanings of sentences. We evaluated this phenomenon in two models, 347 GPT-Neo and Llama 3.1 Instruct, across all layers, with the results presented in Figure 3. 348



364 Figure 3: Linear Representation in GPT-Neo and Llama 3.1: The 0th layer represents the input embedding, which transforms the one-hot encoded input into continuous, lower-dimensional vectors. 366 The first row of the figures corresponds to GPT-Neo, while the second row pertains to Llama 3.1. 367 The impact of linear representation becomes increasingly evident in both GPT-Neo and Llama 3.1. 368 This effect is significantly more pronounced in Llama 3.1, where the values approach 1, suggesting that its layerwise embeddings capture and comprehend the text more effectively than those in 369 GPT-Neo. This indicates that Llama 3.1's embeddings are more adept at representing the underlying 370 relationships and concepts within the text as they advance through the layers. 371

372

373 Figure 3 demonstrates that the impact of linear representation becomes increasingly evident as we 374 progress through the layers of LLMs, indicating a greater ability to capture and comprehend the 375 meaning of the text in the deeper layers. This trend suggests that as data traverses through successive layers, the model continuously refines its understanding and representation of the underlying 376 semantics. This refinement process helps to explain why the phenomenon of monotonicity is ob-377 served across different new datasets.

378 Additionally, this effect is notably more pronounced in Llama 3.1, where the values consistently 379 approach 1 across layers, signifying a stronger alignment in capturing semantic relationships. This 380 indicates that Llama 3.1's layerwise embeddings are more adept at understanding and encoding the 381 nuances of the text compared to those of GPT-Neo. Consequently, Llama 3.1 exhibits a superior 382 capability to accurately represent the intended meaning and context when applied to new datasets.

384 385

387

391

397

399

400

401

409 410

411

412

413

414

415

416

417

418

419

420 421 422

423

424

425

4 APPLICATION OF MONOTONICITY

386 The previous section demonstrated the phenomenon of Monotonicity across various datasets, revealing that it varies depending on the interaction between each dataset and the specific pre-trained 388 LLM. This means that the degree of performance improvement across layers is not uniform. For 389 certain datasets and LLMs, the increase in accuracy flattens in the last few layers. According to 390 the definition of layer performance, this suggests that the embedding from the final layer might be overly sufficient for the datasets, even though it is common practice to use the output from the last 392 layer as the primary embedding for downstream tasks (Devlin, 2018; Dosovitskiy, 2020; Radford 393 et al., 2021). To expedite the inference process with the given dataset and LLM, we can proceed with the following steps: 394

- Step 1: Assess the layer-wise performance of the LLM on this dataset.
- Step 2: Identify the layer that produces embeddings comparable in effectiveness to those generated by the final layer.
- Step 3: Utilize the embeddings from the identified layer for downstream tasks associated with this dataset.

402 A more concrete example of this method can be observed with the Banking77 dataset when using 403 Llama 3.1, as shown in Figure 1 (b) and Figure 2. In this case, the performance of the embeddings 404 from the 5th layer is nearly as effective as those from the last layer, indicating that the deeper layers do not always provide significant additional benefits for certain tasks. Therefore, it may be 405 advantageous to use only the first five layers of Llama 3.1 during the inference process, resulting in 406 substantial savings in both time and memory without sacrificing model performance. More detailed 407 examples supporting this observation can be found in Figure 4. 408



Figure 4: Reducing Inference Time While Maintaining Performance Through Strategic Layer Embedding Selection: By carefully selecting the most appropriate layer embeddings, it can significantly reduce inference time without compromising the model's performance.

426 Figure 4 demonstrates that this approach significantly optimizes computational efficiency. For in-427 stance, in the case of the Banking77 dataset, selecting the appropriate layer embedding from Llama 428 3.1 can lead to an inference time reduction of nearly 85%, with virtually no loss in performance. Similarly, for the AG News dataset, choosing the optimal layer embedding from Phi-3.5 results in 429 a reduction of inference time by approximately 70%, while still maintaining high accuracy. Fur-430 thermore, by leveraging the concept of monotonicity, we can not only identify the most suitable 431 pre-trained LLM for a given downstream task but also select the optimal layer embeddings based on the specific performance and inference time requirements, as illustrated in Figures 1 (b) and (c).
 This makes the approach highly adaptable and efficient for various real-world applications.

435 436

437

5 EXPERIMENTS

Settings in Figure 1 (a) We pre-trained the GPT-2 small model, which consists of 12 transformer
 blocks and approximately 124 million parameters, using the OpenWebText dataset. The training was
 conducted on an 8X A100 80GB node. The process continued until the model's training loss reached
 approximately 3.15, taking around 4 days to complete. Throughout the training, we recorded the
 model parameters at specific checkpoints—epochs 0, 1,500, 5,000, and 165,000. For each of these
 checkpoints, we evaluated the layerwise accuracy on the OpenWebText dataset to assess the model's
 performance at various stages of training, as presented in Figure 1 (a).

444

Settings in Figure 1 (b) and Figure 2 The model sizes used in our experiments are as follows:
1.3 billion parameters for GPT-Neo, 3.8 billion parameters for Phi-3.5-mini-Instruct, and 8 billion parameters for Llama 3.1 Instruct. We conducted the experiments according to the procedure outlined in Section 3. Specifically, the experiments for GPT-Neo and Phi-3.5-mini-Instruct were run on a single RTX 4090 GPU with 24GB of memory, while the experiment for Llama 3.1 Instruct was conducted on an A100 GPU with 40GB of memory.

The corresponding datasets are detailed as follows: For the AG News dataset, we randomly selected 10,000 samples from the training set and used all 7,600 samples from the testing set. This dataset contains four label categories, $Y = \{0, 1, 2, 3\}$. The Banking77 dataset comprises 10,003 training samples and 3,080 testing samples, with a total of 77 label categories, $Y = \{0, 1, ..., 76\}$. Lastly, the Twitter Sentiment Analysis dataset includes 10,223 training samples and 4,382 testing samples, with two label categories, $Y = \{0, 1\}$.

Settings in Figure 3 We utilized the word analogy dataset introduced by (Gladkova et al., 2016)
and specifically selected the "noun - plural" analogy category to conduct the linear representation
experiment. The experiments were carried out following the procedure outlined in Section 3. For
the GPT-Neo model, the experiment was run on a single RTX 4090 GPU with 24GB of memory,
while the Llama 3.1 Instruct experiment was conducted on an A100 GPU with 40GB of memory.

463
 464
 464
 464
 465
 466
 466
 467
 468
 468
 469
 469
 460
 460
 460
 460
 461
 462
 462
 463
 464
 464
 465
 465
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466
 466

- 6 DISCUSSION
- 467 468 469

In this paper, we presented a detailed hierarchical analysis of LLMs, uncovering the phenomenon 470 of monotonicity in layerwise performance. Our study reveals that this progressive improvement in 471 accuracy, observed across successive layers, is not only evident in pre-training data but also extends 472 to a variety of downstream tasks. This demonstrates the ability of LLMs to adapt and generalize 473 to new datasets, making them highly versatile across different applications. We introduced the con-474 cept of linear representation to provide further insights into why monotonicity emerges, highlighting 475 the structured way in which LLMs capture and refine complex features. By identifying that layer-476 wise performance can plateau in the later stages, we also showed how this understanding can be harnessed to optimize inference time and memory usage. Selecting the optimal layer embeddings 477 for specific tasks enables significant efficiency gains without compromising model performance, 478 offering a practical solution for deploying LLMs in resource-constrained environments. 479

However, a limitation of our study is that it does not fully address how the interaction between the
data and the pre-trained LLMs influences monotonicity. For instance, if a dataset contains random
labels without any inherent structure, the phenomenon of monotonicity is unlikely to emerge, as
there is no meaningful relationship for the model to learn. Similarly, monotonicity may not manifest
when the dataset is unrelated to the LLM's pre-training domain, such as using non-NLP datasets
with LLMs that were pre-trained on language tasks. These observations suggest that monotonicity
is not a universal property but rather depends on the alignment between the dataset characteristics

486 and the LLM's learned representations. Future work should explore these interactions more deeply 487 to better understand the conditions under which monotonicity emerges, thereby providing a more 488 comprehensive framework for applying LLMs to a broader range of tasks. 489

490 References 491

499

500

501

502

517

518

519

521

522

523

524

525

527

528

529

538

- 492 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany 493 Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. arXiv preprint arXiv:2404.14219, 494 2024. 495
- 496 Guillaume Alain. Understanding intermediate layers using linear classifier probes. arXiv preprint 497 arXiv:1610.01644, 2016. 498
 - Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. Gpt-neox-20b: An open-source autoregressive language model. arXiv preprint arXiv:2204.06745, 2022.
 - Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- 504 Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Ka-505 mar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general 506 intelligence: Early experiments with gpt-4. arXiv preprint arXiv:2303.12712, 2023. 507
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. Efficient 509 intent detection with dual sentence encoders. arXiv preprint arXiv:2003.04807, 2020.
- 510 Xiaodong Chen, Yuxuan Hu, and Jing Zhang. Compressing large language models by streamlining 511 the unimportant layer. arXiv preprint arXiv:2403.19135, 2024. 512
- 513 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam 514 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: 515 Scaling language modeling with pathways. Journal of Machine Learning Research, 24(240): 516 1-113, 2023.
 - Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805, 2018.
- 520 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. arXiv preprint arXiv:2407.21783, 2024.
 - Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, et al. Layer skip: Enabling early exit inference and self-speculative decoding. arXiv preprint arXiv:2404.16710, 2024.
- 530 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In International Conference on Machine Learning, pp. 10323-10337. PMLR, 2023. 531
- 532 Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. Analogy-based detection of morphologi-533 cal and semantic relations with word embeddings: what works and what doesn't. In Proceedings 534 of the NAACL Student Research Workshop, pp. 8–15, 2016. 535
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/ OpenWebTextCorpus, 2019.
- Hangfeng He and Weijie J Su. A law of next-token prediction in large language models. arXiv preprint arXiv:2408.13442, 2024.

540 541 542	Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. Linearity of relation decoding in transformer language models. <u>arXiv preprint arXiv:2308.09124</u> , 2023.
543 544 545 546	Yibo Jiang, Goutham Rajendran, Pradeep Ravikumar, Bryon Aragam, and Victor Veitch. On the origins of linear representations in large language models. <u>arXiv preprint arXiv:2403.03867</u> , 2024.
547 548	Vishal Kharde, Prof Sonawane, et al. Sentiment analysis of twitter data: a survey of techniques. arXiv preprint arXiv:1601.06971, 2016.
549 550 551 552	Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In <u>International conference on machine learning</u> , pp. 2668–2677. PMLR, 2018.
553 554 555	Bo-Kyeong Kim, Geonmin Kim, Tae-Ho Kim, Thibault Castells, Shinkook Choi, Junho Shin, and Hyoung-Kyu Song. Shortened llama: A simple depth pruning for large language models. <u>arXiv</u> preprint arXiv:2402.02834, 11, 2024.
556 557 558	Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. <u>arXiv preprint arXiv:2305.17888</u> , 2023.
559 560 561 562	Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect. <u>arXiv preprint arXiv:2403.03853</u> , 2024.
563 564 565	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representa- tions of words and phrases and their compositionality. <u>Advances in neural information processing</u> <u>systems</u> , 26, 2013.
565 567 568	Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. <u>arXiv preprint arXiv:2311.03658</u> , 2023.
569 570	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <u>OpenAI blog</u> , 1(8):9, 2019.
571 572 573 574	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <u>International conference on machine learning</u> , pp. 8748–8763. PMLR, 2021.
575 576 577 578	Jiwon Song, Kyungseok Oh, Taesu Kim, Hyungjun Kim, Yulhwa Kim, and Jae-Joon Kim. Sleb: Streamlining llms through redundancy verification and elimination of transformer blocks. <u>arXiv</u> preprint arXiv:2402.09025, 2024.
579 580	Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. <u>arXiv preprint arXiv:2306.11695</u> , 2023.
581 582 583	Curt Tigges, Oskar John Hollinsworth, Atticus Geiger, and Neel Nanda. Linear representations of sentiment in large language models. <u>arXiv preprint arXiv:2310.15154</u> , 2023.
584 585 586	Nikhil Vyas, Alexander Atanasov, Blake Bordelon, Depen Morwani, Sabarish Sainathan, and Cen- giz Pehlevan. Feature-learning networks are consistent across widths at realistic scales. <u>Advances</u> in Neural Information Processing Systems, 36, 2024.
587 588 589	Yuzi Yan, Jialian Li, Yipin Zhang, and Dong Yan. Exploring the llm journey from cognition to expression with linear representations. <u>arXiv preprint arXiv:2405.16964</u> , 2024.
590 591	Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. <u>arXiv preprint</u> <u>arXiv:2011.14522</u> , 2020.
592	Xiang Zhang, Junho Zhao, and Yann LeCun. Character-level convolutional networks for text clas-

Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. <u>Advances in neural information processing systems</u>, 28, 2015.

594 A APPENDIX

In this section, we conduct more experiments to explore more metrics and datasets to check the monotonicity.

Perplexity in GPT-2 We evaluated the performance of GPT-2 using Perplexity on the OpenWeb-Text dataset in different epochs. The Perplexity is calculated as follows:

Perplexity = exp{
$$-\frac{1}{N} \sum_{i=1}^{N} \log p_{\theta}(w_i | w_1, ..., w_{i-1})$$
} (1)

where $p_{\theta}(w_i|w_1, ..., w_{i-1})$ is the probability assigned by the model to the *i*-th word, given the preceding words $w_1, ..., w_{i-1}$ and N is the total number of words in the sequence. The results are shown as follows: The results demonstrate that the monotonicity principle remains effective across



Figure 5: Monotonicity of Perplexity in GPT-2 during training process

different metrics, even after the training process.

Log-Loss for Classification To assess classification performance, we calculated the log-loss on the Banking77 dataset across three models (GPT-neo, Phi-3.5 and Llama 3.1). The log-loss is defined as follows:

$$\text{Log-loss} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{c=1}^{C} y_{i,c} \log p_{i,c}$$
(2)

where N is the total number of words in the sequence and C is the total number of classes. $y_{i,c}$ is a binary indicator (1 or 0) that specifies whether class c is the correct class for sample i. $p_{i,c}$ is the predicted probability of sample i belonging to class c. The results are shown as follows: The results indicate that the monotonicity property is consistently upheld across these models on this metric as well.

MMLU Question Answering - Humanities We conducted experiments on the MMLU Humani-ties question-answering benchmark to evaluate two models. The dataset consists of questions with four answer options, where only one option is correct. For instance, a sample input might be: "ques-tion: What is the embryological origin of the hyoid bone?. answer: A. The first pharyngeal arch, B. The first and second pharyngeal arches", C. The second pharyngeal arch, D. The second and third pharyngeal arches" and the label is "D". This problem is framed as a four-class classification task. We evaluated the monotonicity of two models, Phi-3.5 and Llama 3.1, with the results presented as follows: The findings confirm that the monotonicity principle remains valid for this task. However, due to the complexity of this dataset, there are no accuracy flattens in the last few layers, meaning the proposed method cannot be directly used to accelerate the inference process in this case.



Regression Data For regression tasks, we evaluated the Yelp dataset for rating prediction ($Y = \{1, 2, 3, 4, 5\}$) using Mean Squared Error (MSE) as the evaluation metric. We randomly select 10000 samples with balanced labels (the sample size of each label is 2000). We use 70% of the sample for training and 30% for testing. We test the monotonicity of three models (GPT-Neo, Phi-3.5 and Llama 3.1) and the results are shown as follows:



Figure 8: Monotonicity across three models for regression data

The results confirm that the monotonicity property holds consistently for all three models tested on this regression dataset.

Monotonicity when increasing training samples With a smaller dataset, the monotonicity property may not be immediately apparent due to the limited number of samples, which can introduce variability and noise into the results. However, as the sample size increases, the underlying patterns become more pronounced, and the monotonicity becomes more evident, which is presented as follows: The experiment described above utilizes the Banking77 dataset with varying training sample



Figure 9: Monotonicity across different training sample sizes

sizes to assess performance. The model used for these evaluations is Llama 3.1 8B.