

Internal Representation Dynamics in Transformers

Anonymous EACL submission

Abstract

In this study, we present an investigation into the anisotropy dynamics and intrinsic dimension of embeddings in transformer architectures, focusing on the dichotomy between encoders and decoders. Our findings reveal that the anisotropy profile in transformer decoders exhibits a distinct bell-shaped curve, with the highest anisotropy concentrations in the middle layers. This pattern diverges from the more uniformly distributed anisotropy observed in encoders. In addition, we found that the intrinsic dimension of embeddings increases in the initial phases of training, indicating an expansion into higher-dimensional space. This fact is then followed by a compression phase towards the end of training with dimensionality decrease, suggesting a refinement into more compact representations. Our results provide fresh insights to the understanding of encoders and decoders embedding properties.

1 Introduction

Introduced by Vaswani et al. (2017), transformers have underpinned many breakthroughs, ranging from language modeling to text-to-image generation. As the adoption of transformers has grown, so has the pursuit to understand the intricacies of their internal mechanisms, particularly in the realm of embeddings.

Embeddings in transformers are intricate structures, encoding vast amounts of linguistic nuances and patterns. Historically, researchers have mainly examined embeddings for their linguistic capabilities (Ettinger et al., 2016; Belinkov et al., 2017; Pimentel et al., 2022). Yet, more nuanced properties lie beyond these traditional scopes, like anisotropy and intrinsic dimensionality, which can offer critical insights into the very nature and behavior of these embeddings.

Anisotropy, essentially representing the non-uniformity of a distribution in space, provides

a lens through which we can study the orientation and concentration of embeddings (Ethayarajh, 2019; Biś et al., 2021). A higher degree of anisotropy suggests that vectors are more clustered or directed in specific orientations. In contrast, the intrinsic dimension offers a measure of the effective data dimensionality, highlighting the essence of information that embeddings capture. Together, these metrics can serve as pivotal tools to probe into the black-box nature of transformers.

Our investigation uncovers striking contrasts in anisotropy dynamics between transformer encoders and decoders. By analyzing the training phases of various transformer models, we shed light on the consistent yet previously unrecognized patterns of anisotropy growth. Even more, our analysis revealed a unique dynamic of the averaged intrinsic dimension across layers in decoders: an initial growth during the early stages of training is followed by a decline towards the end. This suggests a two-phase learning strategy where the model initially tries to unfold information in higher dimensional spaces and subsequently compresses it into more compact concepts, possibly leading to more refined representations.

Main Contributions:

- Uncovered a distinct bell-shaped curve for the anisotropy profile¹ in transformer decoders, contrasting with the uniformly distributed anisotropy in encoders.
- Confirmed that anisotropy increases progressively in decoders as the training proceeds.
- Identified a two-phase dynamic in the intrinsic dimension of decoder embeddings: an initial expansion into higher-dimensional space, followed by a compression phase indicating a shift towards compact representations.

¹Layer-wise anisotropy

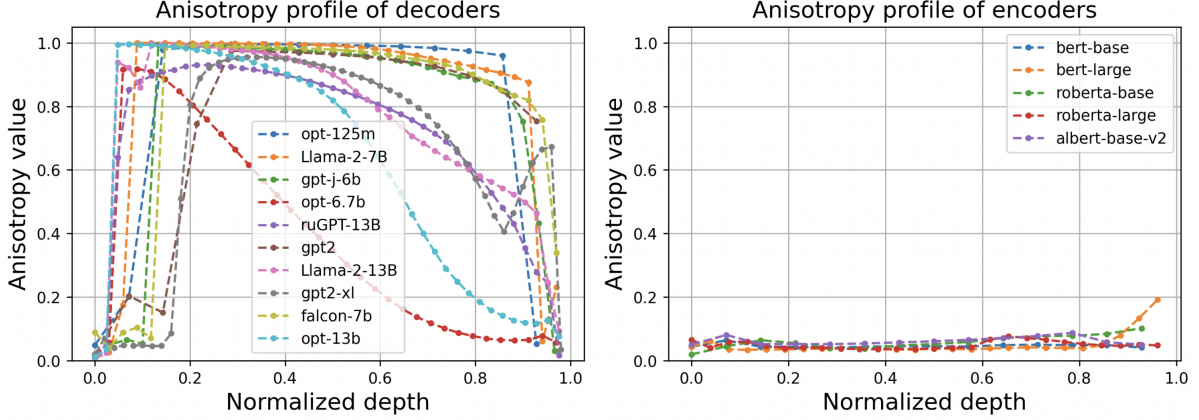


Figure 1: Different anisotropy profiles for transformer-based encoders and decoders.

2 Methodology

2.1 Datasets

As our source for embedding we chose enwik8 dataset (English Wikipedia²) that contains the initial 100 million bytes of Wikipedia, making it a rich source of diverse textual content. It is publicly available through the Hutter Prize website³. The preprocessing includes the removal of all code, media, and HTML tags, resulting in a clean and structured dataset with a vocabulary of 205 distinct characters.

2.2 Embeddings

The vectors are grouped into batches, each with a minimum of 4096 elements. We apply the selected method to determine anisotropy or intrinsic dimension. Prior to assessing intrinsic dimension, embeddings are shuffled (before batching) to mitigate potential correlations. Results from individual batches are then averaged to gauge the metric for that layer, also capturing the standard deviation.

2.3 Anisotropy

To compute anisotropy, we employed singular value decomposition (SVD).

Let $X \in \mathbb{R}^{n_{\text{samples}} \times \text{emb}_{\text{dim}}}$ represent the centered matrix of embeddings, where $\sigma_1, \dots, \sigma_k$ are its singular values. The anisotropy score of X is given by:

$$\text{anisotropy}(X) = \frac{\sigma_1^2}{\sum_{i=1}^k \sigma_i^2}.$$

²<https://www.wikipedia.org/>

³<http://prize.hutter1.net>

Equivalently, this can be deduced using the eigenvalues $\sigma_1^2, \dots, \sigma_k^2$ of the covariance matrix:

$$C = \frac{X^T X}{n_{\text{samples}} - 1}.$$

2.4 Intrinsic Dimension

To determine the intrinsic dimension of a set of embeddings, we utilized the approach proposed by Facco et al. (2018). This method explores how the volume of an n -dimensional sphere (representing the count of embeddings) scales with dimension d .

For each data point within our embeddings, we determine the distances r_1 and r_2 to its two closest neighboring points. This process generates a set of pairs $\{(r_1, r_2)\}$. Using this set, the intrinsic dimension d can be estimated. Hence, we define:

$$\mu_i = \frac{r_2}{r_1},$$

for each point i .

The cumulative distribution function (CDF) of $\{\mu_i\}$ is provided by:

$$F(\mu) = (1 - \mu^{-d})\mathbf{1}_{[1, +\infty)}(\mu).$$

This expression for F is based on the derivations and proofs presented by the authors of the referenced paper. From the CDF, we deduce:

$$\frac{\log(1 - F(\mu))}{\log(\mu)} = d.$$

To estimate d , linear regression $y = kx$ is applied on the plane (x, y) , with:

$$x_i = \log(\mu_i) \quad \text{and} \quad y_i = 1 - F_{\text{emp}}(\mu_i),$$

where F_{emp} signifies the empirical CDF for $\{\mu_i\}$.

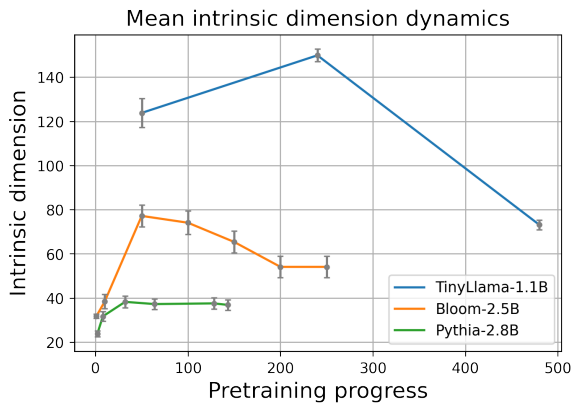


Figure 2: Intrinsic dimension averaged across layers for different pretraining stages.

3 Related Work

3.1 Isotropy of Hidden Representations

Gao et al. (2019) introduce the *representation degeneration problem*. This is the phenomena of degenerating in the representation of the learned embeddings in the generative models, particularly when they are tied. The authors conclude that unlike fixed word embeddings (e.g., word2vec (Mikolov et al., 2013)), the vanilla transformer embeddings are clustered within the narrow cone.

Recent research reveals that global anisotropy is a common trait among all transformer-based architectures (Ait-Saada and Nadif, 2023; Godey et al., 2023; Tyshchuk et al., 2023). However, within local subspaces, isotropy prevails, enhancing model expressiveness and contributing to high performance in downstream tasks.

Ding et al. (2022) conducted an extensive empirical evaluation of modern anisotropy calibration methods, showing no statistically significant improvements in downstream tasks. They conclude that the local isotropy of the hidden space of transformers may lead to the high level of model’s expressiveness (Cai et al., 2021). While most isotropy findings are observed in encoder-only or encoder-decoder architectures, Cai et al. (2021) brought to light an interesting variation. Their study demonstrated that the cosine similarity among embeddings varies across different transformer architectures. The authors conducted experiments on various architectures, evaluating the reduced effective embedding dimension using PCA, and observed high cosine values across layers, especially in models such as GPT-2 (decoder).

The work (Ait-Saada and Nadif, 2023) supports

previous research through extensive experimental evaluation. This study was motivated by the presence of local isotropy in hidden representations, suggesting that anisotropy does not necessarily compromise the expressiveness of these representations.

Godey et al. (2023) investigates the potential causes of anisotropy, particularly its connection to rare words in the transformer’s vocabulary. They explore character-level models to eliminate the influence of rare tokens, but these models do not show significant improvements in experiments. The authors also uncovered that adding common bias term to the inputs can lead to increased attention score variance, promoting the emergence of categorical patterns in self-attention softmax distributions. Increasing input embeddings norm shows signs of anisotropy based on query and key values.

3.2 Intrinsic Dimensionality

Following the idea of local isotropy of the hidden representations, the investigation of the intrinsic task-specific subspaces offers new insights into fine-tuning and the potential to improve model efficiency. Li et al. (2018) suggested that the training trajectory of the Transformer architectures occurs in a low-dimensional subspace. Zhang et al. (2023) demonstrated that fine-tuning engages only a small portion of the model’s parameter, and it is possible to identify the principal directions of these intrinsic task-specific subspaces. Using their method of identifying the training direction they achieved performance similar to fine-tuning in the full parameter space.

3.3 Encoder and Decoder Architecture

The original transformer architecture consists of both encoder and decoder blocks, and each of these blocks can operate independently. The self-attention mechanism is a shared key feature, with decoders utilizing causal self-attention. Decoders are typically trained for language modeling tasks, focusing on generating coherent sequences of text. In contrast, encoders are aimed to produce contextual representations (i.e., embeddings), from the input text.

With limited previous research on the distinctions between the inner representations of encoders and decoders, our study analyzes multiple encoder-based models (such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), and ALBERT (Lan et al., 2020)) and decoder-based models (includ-

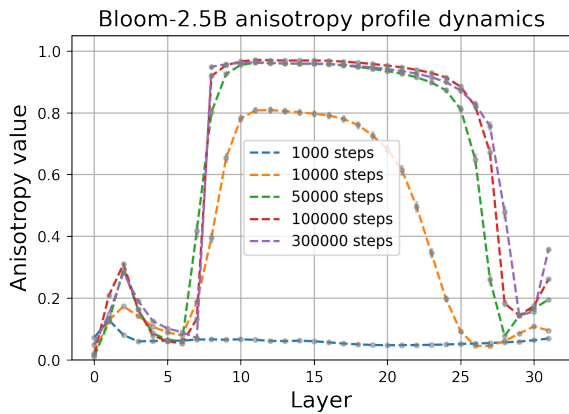


Figure 3: Anisotropy profile for Bloom-2.5 at different number of pretraining steps.

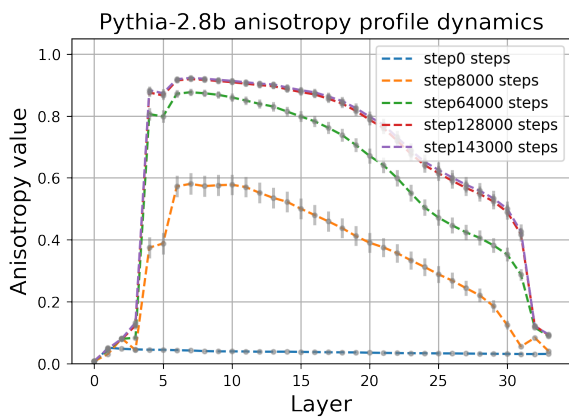


Figure 4: Anisotropy profile for Pythia-2.8B at different number of pretraining steps.

ing OPT 125M-13B (Zhang et al., 2022), Llama-2 7B-13B, Llama-2 7B Chat (Touvron et al., 2023), GPT2 (Radford et al., 2019), GPT-J (Wang and Komatsuzaki, 2021), Falcon-7B, and Falcon-7B-Instruct (Almazrouei et al., 2023)) to offer a comprehensive comparison of their behavior.

4 Results

In this section, we present our empirical findings concerning the anisotropy dynamics and intrinsic dimensionality of transformer embeddings at different layers. Our results span various pretrained transformer models, showcasing clear patterns in the behavior of encoders versus decoders and illuminating the transformation of their properties during training.

4.1 Anisotropy Across Pretrained Transformers

We began by comparing the anisotropy levels across various pretrained transformers, analyzing

both encoder and decoder models. Their anisotropy profiles can be found in the Figure 1.

Encoders: Anisotropy levels remain relatively consistent across models, with minor variations based on model size and training data.

Decoders: In contrast to the encoders, decoders showcase a unique bell-shaped structure, indicating that the middle layers tend to have a higher anisotropy concentration among all examined models.

4.2 Anisotropy Dynamics During Training

To further probe the evolution of anisotropy, we examined its progression through the training phases of various models.

Figure 3 and Figure 4 captures this trajectory by plotting anisotropy values for decoders at different training checkpoints at all internal layers. The consistent growth pattern, followed by stabilization, is observed across various models, suggesting an inherent characteristic of the language modeling training dynamics of decoders.

4.3 Intrinsic Dimensionality During Training

Our exploration into the intrinsic dimensionality revealed intriguing patterns: Figure 2 displays the averaged intrinsic dimension of models throughout the training process. The initial stages exhibit a sharp rise, indicating the model’s attempt to map the information to higher dimensional spaces. However, as training progresses, there is a notable decline, suggesting a subsequent phase where the model compresses this information, refining more compact concepts.

5 Conclusion

Our exploration into the anisotropy dynamics and intrinsic dimensionality of transformer embeddings has brought to light significant distinctions between encoder and decoder transformers. Notably, the intrinsic dimensionality showcased a two-phased training behaviour, where models initially expand information into higher-dimensional spaces and then refine it into compact concepts towards the end of training. These insights not only deepen our understanding of transformer architectures but also suggest new avenues for tailoring training approaches in future NLP research.

282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297

298

299
300
301
302
303
304

305

306
307
308
309
310
311

312
313
314
315
316
317
318

319
320
321
322
323
324
325

326
327
328
329
330
331

Limitations

While our study offers valuable insights into the behavior of transformer embeddings, there are a few limitations to consider.

Model Diversity: Our findings predominantly revolve around specific transformer models, and generalization to all transformer architectures is not guaranteed.

Training Dynamics: The observed two-phased behavior in intrinsic dimensionality might be influenced by the datasets or specific training configurations.

Anisotropy Interpretation: While we identified distinct anisotropy patterns in encoders and decoders, the direct implications of these patterns on downstream tasks remain to be fully explored.

Ethics Statement

Our research focuses on analyzing transformer embeddings and does not involve human subjects or sensitive data. All findings are derived from publicly available models and datasets. We strive for transparency and reproducibility in our methods and analyses.

References

Mira Ait-Saada and Mohamed Nadif. 2023. [Is anisotropy truly harmful? a case study on text clustering](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1194–1203, Toronto, Canada. Association for Computational Linguistics.

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [Falcon-40B: an open large language model with state-of-the-art performance](#).

Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. 2017. [What do neural machine translation models learn about morphology?](#) In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 861–872, Vancouver, Canada. Association for Computational Linguistics.

Daniel Biś, Maksim Podkorytov, and Xiuwen Liu. 2021. [Too much in common: Shifting of embeddings in transformer language models and its implications](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*,

pages 5117–5130, Online. Association for Computational Linguistics. 332
333

Xingyu Cai, Jiaji Huang, Yuchen Bian, and Kenneth Church. 2021. [Isotropy in the contextual embedding space: Clusters and manifolds](#). In *International Conference on Learning Representations*. 334
335
336
337

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics. 338
339
340
341
342
343
344
345
346

Yue Ding, Karolis Martinkus, Damian Pascual, Simon Clematide, and Roger Wattenhofer. 2022. [On isotropy calibration of transformer models](#). In *Proceedings of the Third Workshop on Insights from Negative Results in NLP*, pages 1–9, Dublin, Ireland. Association for Computational Linguistics. 347
348
349
350
351
352

Kawin Ethayarajh. 2019. [How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65, Hong Kong, China. Association for Computational Linguistics. 353
354
355
356
357
358
359
360
361

Allyson Ettinger, Ahmed Elgohary, and Philip Resnik. 2016. [Probing for semantic evidence of composition by means of simple classification tasks](#). In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, pages 134–139, Berlin, Germany. Association for Computational Linguistics. 362
363
364
365
366
367

Elena Facco, Maria d’Errico, Alex Rodriguez, and Alessandro Laio. 2018. [Estimating the intrinsic dimension of datasets by a minimal neighborhood information](#). *CoRR*, abs/1803.06992. 368
369
370
371

Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tiejun Liu. 2019. [Representation degeneration problem in training natural language generation models](#). 372
373
374

Nathan Godey, Éric de la Clergerie, and Benoît Sagot. 2023. [Is anisotropy inherent to transformers?](#) 375
376

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [Albert: A lite bert for self-supervised learning of language representations](#). 377
378
379
380

Chunyuan Li, Heerad Farkhor, Rosanne Liu, and Jason Yosinski. 2018. [Measuring the intrinsic dimension of objective landscapes](#). *CoRR*, abs/1804.08838. 381
382
383

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. 384
385
386

387	Roberta: A robustly optimized bert pretraining approach.	443
388		444
389	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space.	445
390		446
391		447
392	Tiago Pimentel, Josef Valvoda, Niklas Stoehr, and Ryan Cotterell. 2022. Attentional probe: Estimating a module’s functional potential. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11459–11472, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	448
393		449
394		
395		
396		
397		
398		
399	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	
400		
401		
402	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.	
403		
404		
405		
406		
407		
408		
409		
410		
411		
412		
413		
414		
415		
416		
417		
418		
419		
420		
421		
422		
423		
424		
425	Kirill Tyshchuk, Polina Karpikova, Andrew Spiridonov, Anastasiia Prutianova, Anton Razzhigaev, and Alexander Panchenko. 2023. On isotropy of multimodal embeddings. <i>Inf.</i> , 14(7):392.	
426		
427		
428		
429	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	
430		
431		
432		
433		
434	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model.	
435		
436	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pre-trained transformer language models.	
437		
438		
439		
440		
441		
442		