
Geometric Signatures of Compositionality in Language Models

Thomas Jiralerspong*

Université de Montréal
Montreal, Canada

thomas.jiralerspong@mila.quebec

Jin Hwa Lee*

University College London
London, United Kingdom

jin.lee.22@ucl.ac.uk

Lei Yu

University of Toronto
Toronto, Canada

jadeleiyu@cs.toronto.edu

Emily Cheng

Universitat Pompeu Fabra
Barcelona, Spain

emilyshana.cheng@upf.edu

Abstract

1 Compositionality, the notion that the meaning of an expression is constructed from
2 the meaning of its parts and syntactic rules, permits the infinite productivity of
3 human language. For the first time, artificial language models (LMs) are able
4 to match human performance in a number of compositional generalization tasks.
5 However, much remains to be understood about the computational mechanisms
6 underlying these abilities. We take a high-level geometric approach to this problem,
7 relating the degree of compositionality in a dataset to the intrinsic dimensionality
8 of their representations under an LM, a measure of feature complexity. We find that
9 the degree of dataset compositionality is reflected in the intrinsic dimensionality of
10 data representations, where greater combinatorial complexity of the data results in
11 higher representational dimensionality. Finally, we compare linear and nonlinear
12 methods of computing dimensionality, showing that they capture different but
13 complementary aspects of compositional complexity.

14 1 Introduction

15 By virtue of compositionality, few syntactic rules and a finite lexicon can generate an unbounded
16 number of sentences [11]. That is, language, though seemingly high-dimensional, can be explained
17 using relatively few degrees of freedom. A great deal of effort has been made to test whether
18 neural language models (LMs) exhibit human-like compositionality [23, 4]. We take a geometric
19 perspective towards this question, asking how an LM’s representational structure reflects and supports
20 compositional understanding over the course of training.

21 If an LM is a good model of language, we expect its internal representations to exhibit the low-
22 dimensional structure of the latter. That is, representations should reflect the *manifold hypothesis*,
23 or the notion that real-life, high-dimensional data lie on a low-dimensional manifold [20]. The
24 dimension of this manifold, or *intrinsic dimension* (ID), is then the minimal number of degrees of
25 freedom required to describe it without information loss [20, 8].

*Equal contribution

26 The manifold hypothesis has been attested for linguistic representations: LMs indeed compress
27 inputs to an ID orders-of-magnitude lower than their extrinsic dimension [7, 9, 34]. However, despite
28 their conceptual similarity, no work has explicitly linked the degree of linguistic compositionality to
29 representational ID. To bridge this gap, we provide initial experimental insights into the relationship
30 between compositional complexity of inputs and the ID of their representations. In a series of
31 controlled experiments on the Pythia family of language models [6] and a carefully designed synthetic
32 dataset, we confirm that (1) LMs represent linguistic inputs on low-dimensional, nonlinear manifolds,
33 and (2) representational ID predictably reflects degree of input compositionality.

34 2 Background

35 **Compositionality** It has long been a topic of debate whether neural networks also exhibit human-
36 like compositionality when processing natural language [16, 33, 28]. This debate has fueled an
37 extensive line of empirical exploration aimed at assessing the compositionality of neural networks
38 in language modeling via synthetic data [5, 25, 3]. After the recent introduction of large language
39 models with human-level linguistic capability, researchers have shown via mechanistic interpretability
40 analyses that LMs often extract individual word meanings from early layer multi-layer perceptron
41 modules, and compose them via upper-layer self-attention heads to construct semantic representations
42 for multi-word expressions [21, 19]. Our work takes a different approach to understand language
43 model compositionality by connecting it with the geometric properties of a model’s embedding space.

44 **Manifold hypothesis** Deep learning problems are often considered high-dimensional, but research
45 suggests that they are governed by low-dimensional structures. In computer vision, studies have
46 demonstrated that common learning objectives and natural image data reside on low-dimensional
47 manifolds [27, 30, 34, 31]. Similarly, the learning dynamics of neural LMs have been shown to occur
48 within low-dimensional parameter subspaces [1, 35]. The nonlinear, low-dimensional structure that
49 emerges in the semantic space of these models likely follows from the training objective of predicting
50 sequential observations [32], which can simplify transfer learning to new tasks and datasets [9].

51 3 Setup

52 **Models** We evaluate pre-trained Transformer-based LMs of sizes $\in \{70\text{m}, 140\text{m}, 1.4\text{b}, 6.9\text{b}, 12\text{b}\}$
53 from the Pythia family [6]. Models were trained on the causal language modeling task on The Pile, a
54 natural language corpus comprising encyclopedic text, books, social media, code, and reviews [17].

55 **Dataset** As we investigate compositional generalization of the LM, we construct a dataset consisting
56 of nonce sentences from a toy grammar. To create the grammar, we set 12 semantic categories and
57 randomly sample a 50-word vocabulary for each category, where the categories’ vocabularies are
58 disjoint. The categories include 6 adjective types (quality, nationality, size, color, texture), 2 noun
59 types (job, animal) and 1 verb type. We use a simple, fixed syntax by ordering the word categories:

60 The [quality₁.ADJ][nationality₁.ADJ][job₁.N] [action₁.V] the [size₁.ADJ][texture.ADJ]
[color.ADJ][animal.N] then [action₂.V] the [size₂.ADJ][quality₂.ADJ][nationality₂.ADJ]
[job₂.N].

61 The vocabularies are found in Appendix D. The syntax is chosen so that sentences are grammatical
62 and that adjective order complies with the accepted order for English [12]. Although the syntactic
63 structure and vocabulary items are likely seen during training, words are sampled independently for
64 each category without considering the sentence’s global semantic coherence. Therefore, sentences
65 are unlikely seen during training. When encountering them for the first time, a frozen LM must
66 successfully construct their meanings from the meanings of their parts, or compositionally generalize.

67 **Controlling compositionality** We are interested in two types of compositionality: (1) combina-
68 torial dataset complexity, where a dataset is more compositional if it contains more unique word
69 combinations; (2) sentence-level compositional semantics, where sentence meaning is composed, via
70 syntax, from word meanings.

71 First, to control for dataset compositionality, we couple the values of k word positions for $k = 1 \dots 4$.
72 When k positions are coupled, the sequence’s atomic units are sets of k contiguous words, constraining

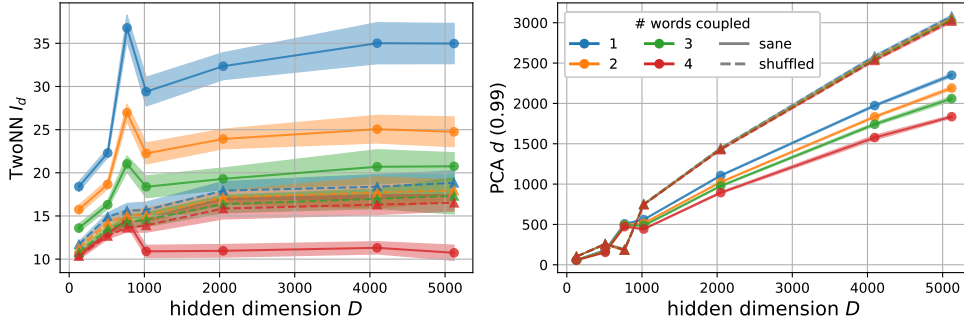


Figure 1: **Mean dimensionality over model size.** Mean nonlinear I_d (left) and linear d (right) over layers is shown for increasing LM hidden dimension D . While the nonlinear I_d does not depend on extrinsic dimension D (flat lines), the PCA d scales roughly linearly in D . Curves are averaged over 5 random seeds, shown with ± 1 SD.

73 the number of degrees of freedom to l/k where $l = 12$ is the number of variable words in the sequence.
 74 For instance, in the 1-coupled setting, words are sampled independently, thus 12 degrees of freedom;
 75 if 2-coupled, bigrams are sampled independently, hence 6 degrees of freedom. Increasing k maintains
 76 the dataset’s unigram distribution, but constrains its combinatorial complexity.

77 Second, to investigate *compositional semantics*, we randomly shuffle the words in each sequence.
 78 This destroys syntactic coherence, and in turn, the composed meaning of the sentence; it instead
 79 preserves distributional properties like sequence length and unigram frequencies. Then, LM behavior
 80 on grammatically sane vs. shuffled sequences proxies compositional vs. lexical-only semantics.

81 **Dimensionality estimation** We are interested in whether the geometry of representations reflects
 82 their underlying degree of compositionality. In particular, we consider representations in the *residual*
 83 *stream* of the Transformer [14]. Because sequence lengths vary, in line with prior work [9], we
 84 aggregate over the sequence by taking the last token representation, as it is the only to attend to the
 85 entire context. For each layer and dataset, we compute both a nonlinear and a linear measure of
 86 dimensionality, which have key conceptual differences. The nonlinear I_d is the number of degrees of
 87 freedom, or latent features, needed to describe the underlying representation manifold [8, 2, 15]. This
 88 differs from the *linear* effective dimensionality d , or the dimension of the minimal linear subspace
 89 needed to contain the set of representations. Throughout, we will use *dimensionality* to refer to both
 90 nonlinear and linear estimates. When appropriate, we will specify I_d as the nonlinear ID, d as the
 91 *linear* effective dimension, and D as the extrinsic dimension, or hidden dimension of the model.

92 We report the nonlinear I_d using the popular TwoNN estimator of [15], and we estimate the linear
 93 effective dimensionality d using Principal Component Analysis [24] with a variance cutoff of 99%.
 94 Though in the main paper we focus on TwoNN and PCA, we also tested the Maximum Likelihood
 95 Estimator of [26] and the Participation Ratio [32]. For mathematical details, see Appendix C.

96 4 Results

97 We find representational dimensionality to reflect compositionality in ways that are predictable across
 98 model scale. First, we demonstrate that linear and nonlinear dimensionality measures behave differ-
 99 ently across model scale. Then, we show that dimensionality reflects the degree of compositionality
 100 of its inputs, highlighting the difference between nonlinear and linear measures. For brevity, we focus
 101 on model sizes 410m, 1.4b, and 6.9b in the main text, with full results in the appendix.

102 **Nonlinear and linear ID scale differently with model size** Like in previous work [7, 34, 10,
 103 22, 13], we confirm that inputs are represented in a nonlinear manifold with orders-of-magnitude
 104 lower dimension than the ambient dimension. In particular, we find that $I_d \sim O(10)$ across models
 105 sizes (see Figure 1 left). We find, moreover, that larger models tend to have higher representational
 106 dimensionality, but that the scaling is not uniform. Figure 1 shows that while the linear d scales
 107 linearly with hidden dimension D , nonlinear I_d instead stabilizes to the mentioned range $\sim O(10)$

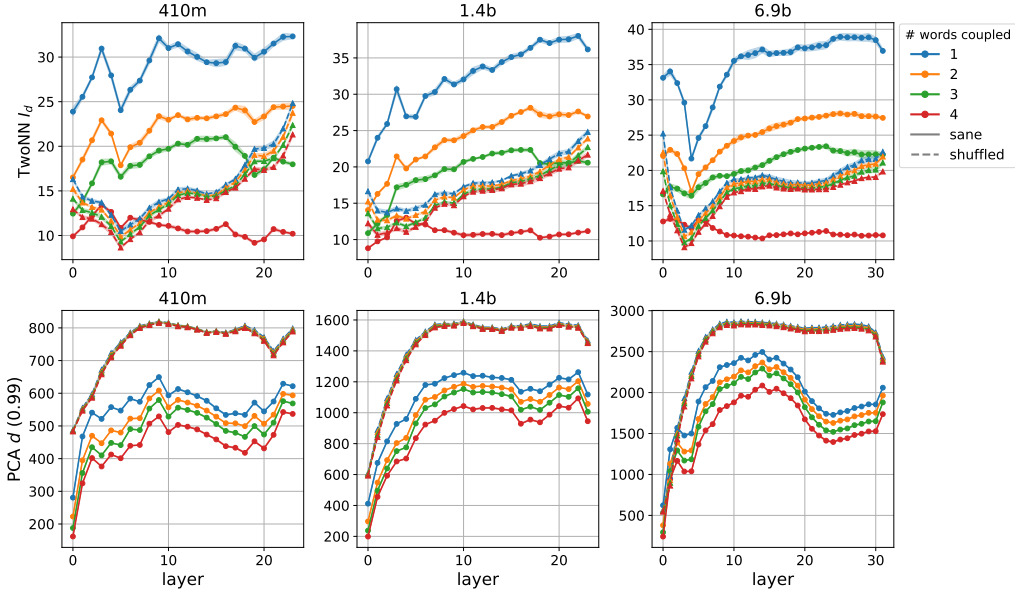


Figure 2: **Dimensionality over layers.** Nonlinear I_d (top) and linear d (bottom) over layers are shown for three sizes: 410m, 1.4b, and 6.9b (left to right). Each color corresponds to a coupling length $k \in 1 \dots 4$. Solid curves denote sane sequences, and dotted curves denote shuffled sequences. For all models, lower k results in higher I_d and d for both normal and shuffled settings. For all models, shuffling results in lower I_d but higher d . Curves are averaged over 5 random seeds, shown with ± 1 SD.

108 regardless of extrinsic dimension. This result highlights key differences in how linear and nonlinear
 109 dimensions are recruited: LMs *globally* distribute representations to occupy $d \propto D$ dimensions of the
 110 space, but *locally* constrains their shape to a low-dimensional (I_d) manifold.

111 **Representational ID reflects input compositionality** Representational dimensionality preserves
 112 relative data combinatorial complexity. Figure 2 shows I_d and d over LM layers for $k = 1 \dots 4$
 113 coupling lengths (different colors). For both sane and shuffled settings, both I_d and d increase
 114 predictably with input complexity: the highest curves correspond to the 1-coupled dataset, or 12
 115 degrees of freedom, while the lowest correspond to the 4-coupled dataset, or 3 degrees of freedom.

116 Now, we consider sequence-level compositional semantics. See Figure 2 again for the dimensionality
 117 over layers in sane (solid curves) and shuffled (dotted curves) settings. Intriguingly, nonlinear and
 118 linear dimensionalities of shuffled examples show opposing patterns: compared to the sane text,
 119 shuffled text I_d generally decreases and is compressed to a small range, while d *increases*. These
 120 diverging patterns do not necessarily contradict each other, however. We interpret the discrepancy in
 121 line with Recanatesi et al. [32]. Predictive coding requires an LM to encode the vast space of inputs
 122 and outputs, as well as extract latent semantic features to support the former. Recanatesi et al. [32]
 123 argue that encoding all possible sequences makes use of the *global* representation space \mathbb{R}^D ; instead,
 124 encoding semantic relationships between sequences, i.e., latent features, occurs via *local* correlations
 125 that give rise to a I_d -dimensional manifold. In our setting, randomly permuting words in a length- l
 126 sequence increases the implied input space by a factor of $\sim l!$, which puts an upward pressure on d .
 127 But, permuting words destroys the semantics of the sequence, exerting a downward pressure on I_d .

128 5 Discussion

129 We have studied the computational mechanism of LM compositionality from a geometric perspective.
 130 Using a carefully designed synthetic dataset, we found strong relationships between the compositionality
 131 of linguistic expressions and the geometric complexity of their representations. In particular,
 132 dataset combinatorial compositionality is positively correlated to both nonlinear and linear dimensionality.
 133 On the other hand, sequences with high semantic compositionality exhibit high nonlinear I_d but

134 a low linear d . Crucially, nonlinear complexity measures have been underexplored in the literature
135 compared to linear ones; we demonstrate their empirical differences, highlighting a need to further
136 investigate nonlinear measures to proxy feature learning in deep neural models. We hypothesize that
137 linear d proxies a dataset’s implied size, and nonlinear I_d its meaningful semantic variability.

138 **Limitations** Our analysis is limited to the Pythia family of models. Though it has been suggested
139 that causal LMs have similar representational geometry [29, 10], experiments on a wider range of
140 LMs and grammars, as well as theoretical work, will be necessary to draw general conclusions.

141 References

- 142 [1] Armen Aghajanyan, Sonal Gupta, and Luke Zettlemoyer. Intrinsic dimensionality explains
143 the effectiveness of language model fine-tuning. In *Proceedings of the 59th Annual Meeting
144 of the Association for Computational Linguistics and the 11th International Joint Conference
145 on Natural Language Processing (Volume 1: Long Papers)*, pages 7319–7328, Online, August
146 2021. Association for Computational Linguistics.
- 147 [2] Alessio Ansuini, Alessandro Laio, Jakob H Macke, and Davide Zoccolan. Intrinsic dimension
148 of data representations in deep neural networks. In *Advances in Neural Information Processing
149 Systems*, volume 32. Curran Associates, Inc., 2019.
- 150 [3] Dzmitry Bahdanau, Shikhar Murty, Michael Noukhovitch, Thien Huu Nguyen, Harm de Vries,
151 and Aaron Courville. Systematic generalization: what is required and can it be learned? *arXiv
152 preprint arXiv:1811.12889*, 2018.
- 153 [4] Marco Baroni. Linguistic generalization and compositionality in modern artificial neural
154 networks. *Philosophical Transactions of the Royal Society B*, 375, 2019. URL [https://api.
155 semanticscholar.org/CorpusID:90260325](https://api.semanticscholar.org/CorpusID:90260325).
- 156 [5] Luisa Bentivogli, Arianna Bisazza, Mauro Cettolo, and Marcello Federico. Neural versus
157 phrase-based machine translation quality: a case study. In *Proceedings of the Conference on
158 Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational
159 Linguistics (ACL), 2016.
- 160 [6] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien,
161 Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward
162 Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In
163 *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- 164 [7] Xingyu Cai, Jiayi Huang, Yuchen Bian, and Kenneth Church. Isotropy in the contextual embed-
165 ding space: Clusters and manifolds. In *International Conference on Learning Representations*,
166 2021.
- 167 [8] P. Campadelli, E. Casiraghi, C. Ceruti, and A. Rozza. Intrinsic dimension estimation: Relevant
168 techniques and a benchmark framework. *Mathematical Problems in Engineering*, 2015:e759567,
169 Oct 2015. ISSN 1024-123X.
- 170 [9] Emily Cheng, Corentin Kervadec, and Marco Baroni. Bridging information-theoretic and
171 geometric compression in language models. In *Proceedings of EMNLP*, pages 12397–12420,
172 Singapore, 2023.
- 173 [10] Emily Cheng, Diego Doimo, Corentin Kervadec, Iuri Macocco, Jade Yu, Alessandro Laio, and
174 Marco Baroni. Emergence of a high-dimensional abstraction phase in language transformers,
175 2024. URL <https://arxiv.org/abs/2405.15471>.
- 176 [11] Noam Chomsky. *Syntactic Structures*. Mouton and Co., The Hague, 1957.
- 177 [12] Robert Mw Dixon. Iwhere have all the adjectives gone. *Studies in Language*, 1:19–80, 1976.
- 178 [13] Diego Doimo, Alessandro Serra, Alessio Ansuini, and Alberto Cazzaniga. The representation
179 landscape of few-shot learning and fine-tuning in large language models, 2024. URL [https://
180 arxiv.org/abs/2409.03662](https://arxiv.org/abs/2409.03662).

- 181 [14] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
182 Amanda Aspell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep
183 Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt,
184 Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and
185 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,
186 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 187 [15] Elena Facco, Maria d’Errico, Alex Rodriguez, and Alessandro Laio. Estimating the intrinsic
188 dimension of datasets by a minimal neighborhood information. *Scientific Reports*, 7(1):12140,
189 Sep 2017. ISSN 2045-2322. doi: 10.1038/s41598-017-11873-y.
- 190 [16] Jerry A Fodor and Zenon W Pylyshyn. Connectionism and cognitive architecture: A critical
191 analysis. *Cognition*, 28(1-2):3–71, 1988.
- 192 [17] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason
193 Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The Pile: An 800GB dataset of diverse
194 text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- 195 [18] Peiran Gao, Eric M. Trautmann, Byron M. Yu, Gopal Santhanam, Stephen I. Ryu, Krishna V.
196 Shenoy, and Surya Ganguli. A theory of multineuronal dimensionality, dynamics and measure-
197 ment. *bioRxiv*, 2017. URL <https://api.semanticscholar.org/CorpusID:19938440>.
- 198 [19] Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual
199 associations in auto-regressive language models. In *Proceedings of the 2023 Conference on*
200 *Empirical Methods in Natural Language Processing*, pages 12216–12235, 2023.
- 201 [20] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.
202 <http://www.deeplearningbook.org>.
- 203 [21] Adi Haviv, Ido Cohen, Jacob Gidron, Roei Schuster, Yoav Goldberg, and Mor Geva. Understand-
204 ing transformer memorization recall through idioms. In *Proceedings of the 17th Conference of*
205 *the European Chapter of the Association for Computational Linguistics*, pages 248–264, 2023.
- 206 [22] Evan Hernandez and Jacob Andreas. The low-dimensional linear geometry of contextualized
207 word representations. In *Proceedings of the 25th Conference on Computational Natural*
208 *Language Learning*, pages 82–93, Online, November 2021. Association for Computational
209 Linguistics. doi: 10.18653/v1/2021.conll-1.7. URL <https://aclanthology.org/2021.conll-1.7>.
- 211 [23] Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed:
212 How do neural networks generalise? *J. Artif. Intell. Res.*, 67:757–795, 2019. URL <https://api.semanticscholar.org/CorpusID:211259383>.
- 214 [24] Ian Jolliffe. *Principal Component Analysis*. Springer, 1986.
- 215 [25] Brenden Lake and Marco Baroni. Generalization without systematicity: On the compositional
216 skills of sequence-to-sequence recurrent networks. In *International conference on machine*
217 *learning*, pages 2873–2882. PMLR, 2018.
- 218 [26] Elizaveta Levina and Peter Bickel. Maximum likelihood estimation of intrinsic
219 dimension. In *Advances in Neural Information Processing Systems*, volume 17.
220 MIT Press, 2004. URL https://papers.nips.cc/paper_files/paper/2004/hash/74934548253bcab8490ebd74afed7031-Abstract.html.
- 222 [27] Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic
223 dimension of objective landscapes. In *International Conference on Learning Representations*,
224 2018.
- 225 [28] Gary F Marcus. *The algebraic mind: Integrating connectionism and cognitive science*. MIT
226 press, 2003.
- 227 [29] Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello,
228 and Emanuele Rodolà. Relative representations enable zero-shot latent space communication.
229 In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=SrC-nwieGJ>.
- 230

- 231 [30] Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic
 232 dimension of images and its impact on learning. In *International Conference on Learning*
 233 *Representations*, 2021.
- 234 [31] Michael Psenka, Druv Pai, Vishal Raman, Shankar Sastry, and Yi Ma. Representation learning
 235 via manifold flattening and reconstruction. *Journal of Machine Learning Research*, 25(132):
 236 1–47, 2024.
- 237 [32] Stefano Recanatesi, Matthew Farrell, Guillaume Lajoie, Sophie Deneve, Mattia Rigotti, and Eric
 238 Shea-Brown. Predictive learning as a network mechanism for extracting low-dimensional latent
 239 space representations. *Nature Communications*, 12(1):1417, March 2021. ISSN 2041-1723.
 240 doi: 10.1038/s41467-021-21696-1.
- 241 [33] Paul Smolensky. Tensor product variable binding and the representation of symbolic structures
 242 in connectionist systems. *Artificial intelligence*, 46(1-2):159–216, 1990.
- 243 [34] Lucrezia Valeriani, Diego Doimo, Francesca Cuturello, Alessandro Laio, Alessio Ansuini, and
 244 Alberto Cazzaniga. The geometry of hidden representations of large transformer models.
 245 (arXiv:2302.00294), Feb 2023. doi: 10.48550/arXiv.2302.00294. URL [http://arxiv.org/](http://arxiv.org/abs/2302.00294)
 246 [abs/2302.00294](http://arxiv.org/abs/2302.00294). arXiv:2302.00294 [cs, stat].
- 247 [35] Zhong Zhang, Bang Liu, and Junming Shao. Fine-tuning happens in tiny subspaces: Exploring
 248 intrinsic task-specific subspaces of pre-trained language models. In *Proceedings of the 61st*
 249 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
 250 pages 1701–1713, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:
 251 10.18653/v1/2023.acl-long.95. URL <https://aclanthology.org/2023.acl-long.95>.

252 A Computing resources

253 All experiments were run on a cluster with 12 nodes with 5 NVIDIA A30 GPUs and 48 CPUs each.
 254 Extracting LM representations took a few wall-clock hours per model-dataset computation. ID
 255 computation took approximately 0.5 hours per model-dataset computation. Taking parallelization
 256 into account, we estimate the overall wall-clock time taken by all experiments, including failed runs,
 257 preliminary experiments, etc., to be of about 10 days.

258 B Assets

259 **Pythia** <https://huggingface.co/EleutherAI/pythia-6.9b-deduped>; license: apache-2.0
 260 **scikit-dimension** <https://scikit-dimension.readthedocs.io/en/latest/>; license: bsd-
 261 3-clause
 262 **PyTorch** <https://scikit-learn.org/>; license: bsd

263 C ID Estimation

264 **TwoNN Estimator** A number of methods have been proposed to estimate the nonlinear ID of high-
 265 dimensional point clouds [8]. State-of-the-art ID estimators work by exploiting known relationships
 266 between points in d -dimensions, then fitting d using maximum likelihood estimation from data. We
 267 considered the commonly used TwoNN estimator of 15, which has been found to highly correlate to
 268 other state-of-the-art estimators [9, 8].

269 The TwoNN method works as follows. In brief, points on the underlying manifold are assumed to
 270 follow a locally homogeneous Poisson point process. Local, in this case, refers to neighborhoods
 271 about each point x which encompass x 's first and second nearest neighbors. Let $r_k^{(i)}$ be the Euclidean
 272 distance between point x_i and its k th nearest neighbor. Then, under the mentioned assumptions, the
 273 distance ratios $\mu_i := r_1^{(i)}/r_2^{(i)}$ follow the cumulative distribution function $F(\mu) = 1 - \mu^{-I_d}$. Finally,
 274 I_d is numerically estimated from data.

275 **Maximum Likelihood Estimator** In addition to TwoNN, we considered Levina and Bickel [26]’s
 276 Maximum Likelihood Estimator (MLE), a similar, nonlinear measure of I_d . MLE has been used
 277 in prior works on representational geometry such as [7, 9, 30], and similarly models the number of
 278 points in a neighborhood around a reference point x to follow a Poisson point process. For details
 279 we refer to the original paper [26]. Like past work [15, 9], we found MLE and TwoNN to be highly
 280 correlated, producing results that were nearly identical: compare Figure 1 left to Figure E.3 left, and
 281 Figure E.1 top to Figure E.2 top).

282 **Participation Ratio** For our primary linear measure of dimensionality d , we computed PCA and
 283 took the number of components that explain 99% of the variance. In addition to PCA, we computed
 284 the Participation Ratio (PR), defined as $(\sum_i \lambda_i)^2 / (\sum_i \lambda_i^2)$ [18]. We found PR to give results that
 285 were incongruous with intuitions about linear dimensionality. In particular, it produced a lower
 286 dimensionality estimate than the nonlinear estimators we tested; see, e.g., Figure E.3, where the PR- d
 287 for sane text is less than that of TwoNN. This contradicts the mathematical relationship that $I_d \leq d \leq$
 288 D . This may be because, empirically, PR- d corresponded to explained variances of 60 – 80%, which
 289 are inadequate to describe the bounding linear subspace for the representation manifold. Therefore,
 290 while we report the mean PR- d over model size in Figure E.3 and the dimensionality over layers in
 291 Figure E.2 for completeness, we do not attempt to interpret them.

292 D Toy Grammar

293 The grammar is composed of sentences of the form

294 The [quality₁.ADJ][nationality₁.ADJ][job₁.N] [action₁.V] the [size₁.ADJ][texture.ADJ]
 [color.ADJ][animal.N] then [action₂.V] the [size₂.ADJ][quality₂.ADJ][nationality₂.ADJ]
 [job₂.N].

295 Each category, colored and enclosed in brackets, is sampled from a vocabulary of 50 possible words,
 296 listed in the table below:

Category	Words
job ₁	teacher, doctor, engineer, chef, lawyer, plumber, electrician, accountant, nurse, mechanic, architect, dentist, programmer, photographer, painter, firefighter, police, pilot, farmer, waiter, scientist, actor, musician, writer, athlete, designer, carpenter, librarian, journalist, psychologist, gardener, baker, butcher, tailor, cashier, barber, janitor, receptionist, salesperson, manager, tutor, coach, translator, veterinarian, pharmacist, therapist, driver, bartender, security, clerk
job ₂	banker, realtor, consultant, therapist, optometrist, astronomer, biologist, geologist, archaeologist, anthropologist, economist, sociologist, historian, philosopher, linguist, meteorologist, zoologist, botanist, chemist, physicist, mathematician, statistician, surveyor, pilot, steward, dispatcher, ichthyologist, oceanographer, ecologist, geneticist, microbiologist, neurologist, cardiologist, pediatrician, surgeon, anesthesiologist, radiologist, dermatologist, gynecologist, urologist, psychiatrist, physiotherapist, chiropractor, nutritionist, personal trainer, yoga instructor, masseur, acupuncturist, paramedic, midwife
animal	dog, cat, elephant, lion, tiger, giraffe, zebra, monkey, gorilla, chimpanzee, bear, wolf, fox, deer, moose, rabbit, squirrel, raccoon, beaver, otter, penguin, eagle, hawk, owl, parrot, flamingo, ostrich, peacock, swan, duck, frog, toad, snake, lizard, turtle, crocodile, alligator, shark, whale, dolphin, octopus, jellyfish, starfish, crab, lobster, butterfly, bee, ant, spider, scorpion

color	red, blue, green, yellow, purple, orange, pink, brown, gray, black, white, cyan, magenta, turquoise, indigo, violet, maroon, navy, olive, teal, lime, aqua, coral, crimson, fuchsia, gold, silver, bronze, beige, tan, khaki, lavender, plum, periwinkle, mauve, chartreuse, azure, mint, sage, ivory, salmon, peach, apricot, mustard, rust, burgundy, mahogany, chestnut, sienna, ochre
size ₁	big, small, large, tiny, huge, giant, massive, microscopic, enormous, colossal, miniature, petite, compact, spacious, vast, wide, narrow, slim, thick, thin, broad, expansive, extensive, substantial, boundless, considerable, immense, mammoth, towering, titanic, gargantuan, diminutive, minuscule, minute, hulking, bulky, hefty, voluminous, capacious, roomy, cramped, confined, restricted, limited, oversized, undersized, full, empty, half, partial
size ₂	lengthy, short, tall, long, deep, shallow, high, low, medium, average, moderate, middling, intermediate, standard, regular, normal, ordinary, sizable, generous, abundant, plentiful, copious, meager, scanty, skimpy, inadequate, sufficient, ample, excessive, extravagant, exorbitant, modest, humble, grand, majestic, imposing, commanding, dwarfed, diminished, reduced, enlarged, magnified, amplified, expanded, contracted, shrunken, swollen, bloated, inflated, deflated
nationality ₁	American, British, Canadian, Australian, German, French, Italian, Spanish, Japanese, Chinese, Indian, Russian, Brazilian, Mexican, Argentinian, Turkish, Egyptian, Nigerian, Kenyan, African, Swedish, Norwegian, Danish, Finnish, Icelandic, Dutch, Belgian, Swiss, Austrian, Greek, Polish, Hungarian, Czech, Slovak, Romanian, Bulgarian, Serbian, Croatian, Slovenian, Ukrainian, Belarusian, Estonian, Latvian, Lithuanian, Irish, Scottish, Welsh, Portuguese, Moroccan, Algerian
nationality ₂	Vietnamese, Thai, Malaysian, Indonesian, Filipino, Singaporean, Nepalese, Bangladeshi, Maldivian, Pakistani, Afghan, Iranian, Iraqi, Syrian, Lebanese, Israeli, Saudi, Emirati, Qatari, Kuwaiti, Omani, Yemeni, Jordanian, Palestinian, Bahraini, Tunisian, Libyan, Sudanese, Ethiopian, Somali, Ghanaian, Ivorian, Senegalese, Malian, Cameroonian, Congolese, Ugandan, Rwandan, Tanzanian, Mozambican, Zambian, Zimbabwean, Namibian, Botswanan, New Zealander, Fijian, Samoan, Tongan, Papuan, Marshallese
action ₁	feeds, walks, grooms, pets, trains, rides, tames, leashes, bathes, brushes, adopts, rescues, shelters, houses, cages, releases, frees, observes, studies, examines, photographs, films, sketches, paints, draws, catches, hunts, traps, chases, pursues, tracks, follows, herds, corrals, milks, shears, breeds, mates, clones, dissects, stuffs, mounts, taxidermies, domesticates, harnesses, saddles, muzzles, tags, chips, vaccinates
action ₂	hugs, kisses, loves, hates, admires, respects, befriends, distrusts, helps, hurts, teaches, learns from, mentors, guides, counsels, advises, supports, undermines, praises, criticizes, compliments, insults, congratulates, consoles, comforts, irritates, annoys, amuses, entertains, bores, inspires, motivates, discourages, intimidates, impresses, disappoints, surprises, shocks, delights, disgusts, forgives, resents, envies, pities, understands, misunderstands, trusts, mistrusts, betrays, protects

quality ₁	good, bad, excellent, poor, superior, inferior, outstanding, mediocre, exceptional, sublime, superb, terrible, wonderful, awful, great, horrible, fantastic, dreadful, marvelous, atrocious, splendid, appalling, brilliant, dismal, fabulous, lousy, terrific, abysmal, incredible, substandard, amazing, disappointing, extraordinary, stellar, remarkable, unremarkable, impressive, unimpressive, admirable, despicable, praiseworthy, blameworthy, commendable, reprehensible, exemplary, subpar, ideal, flawed, perfect, imperfect
quality ₂	acceptable, unacceptable, satisfactory, unsatisfactory, sophisticated, insufficient, adequate, exquisite, suitable, unsuitable, appropriate, inappropriate, fitting, unfitting, proper, improper, correct, incorrect, right, wrong, accurate, inaccurate, precise, imprecise, exact, inexact, flawless, faulty, sound, unsound, reliable, unreliable, dependable, undependable, trustworthy, untrustworthy, authentic, fake, genuine, counterfeit, legitimate, illegitimate, valid, invalid, legal, illegal, ethical, unethical, moral, immoral
texture	smooth, rough, soft, hard, silky, coarse, fluffy, fuzzy, furry, hairy, bumpy, lumpy, grainy, gritty, sandy, slimy, slippery, sticky, tacky, greasy, oily, waxy, velvety, leathery, rubbery, spongy, springy, elastic, pliable, flexible, rigid, stiff, brittle, crumbly, flaky, crispy, crunchy, chewy, stringy, fibrous, porous, dense, heavy, light, airy, feathery, downy, woolly, nubby, textured

297 **E Additional Results**

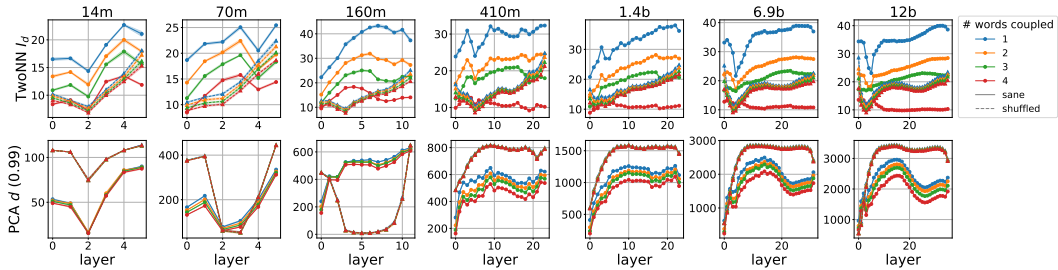


Figure E.1: **Dimensionality over layers.** TwoNN nonlinear I_d (top) and PCA linear d (bottom) over layers are shown for all sizes (left to right). Each color corresponds to a coupling length $k \in 1 \dots 4$. Solid curves denote sane sequences, and dotted curves denote shuffled sequences. For all models, lower k results in higher I_d and d for both normal and shuffled settings. For all models, shuffling results in lower I_d but higher d . Curves are averaged over 5 random seeds, shown with ± 1 SD.

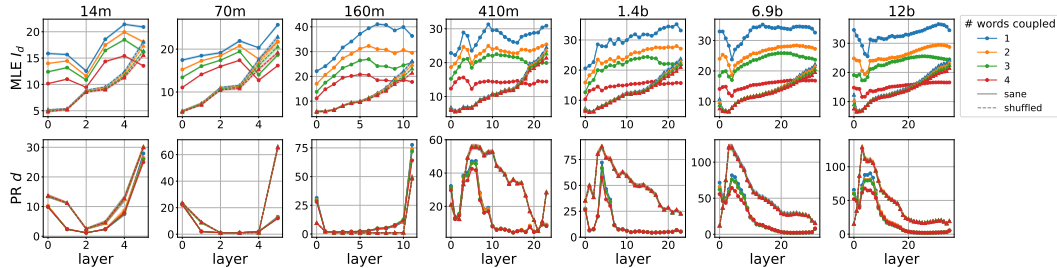


Figure E.2: **Other dimensionality metrics over layers.** MLE nonlinear I_d (top) and PR linear d (bottom) over layers are shown for all model sizes (left to right). Each color corresponds to a coupling length $k \in 1 \dots 4$. Solid curves denote sane sequences, and dotted curves denote shuffled sequences. For all models, lower k results in higher I_d for both normal and shuffled settings. For all models, shuffling results in lower I_d . The PR- d produced nonsensical results, with linear dimensionality higher than nonlinear dimensionality. Curves are averaged over 5 random seeds, shown with ± 1 SD.

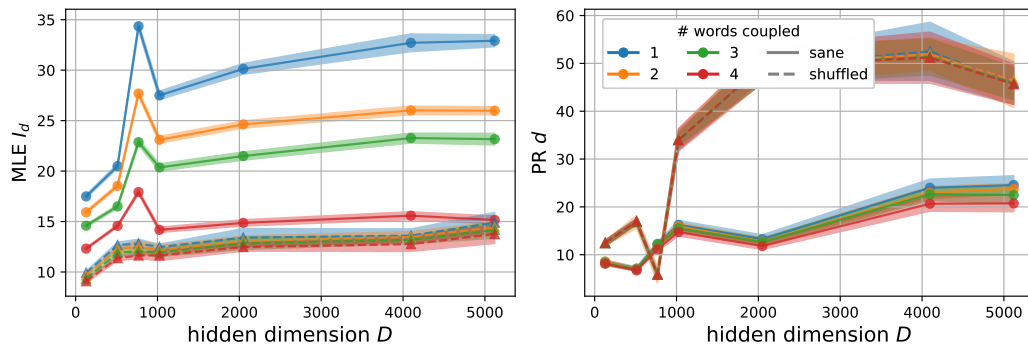


Figure E.3: **Mean dimensionality over model size (other metrics).** Mean nonlinear I_d computed with MLE (left) and linear d computed with PR (right) over layers is shown for increasing LM hidden dimension D . MLE I_d does not depend on extrinsic dimension D (flat lines). PR d produces nonsensical values, higher than the nonlinear I_d . Curves are averaged over 5 random seeds, shown with ± 1 SD.

298 NeurIPS Paper Checklist

299 1. Claims

300 Question: Do the main claims made in the abstract and introduction accurately reflect the
301 paper's contributions and scope?

302 Answer: [Yes]

303 Justification: Sections 4.1-3 support claims in the Abstract and Introduction.

304 Guidelines:

- 305 • The answer NA means that the abstract and introduction do not include the claims
306 made in the paper.
- 307 • The abstract and/or introduction should clearly state the claims made, including the
308 contributions made in the paper and important assumptions and limitations. A No or
309 NA answer to this question will not be perceived well by the reviewers.
- 310 • The claims made should match theoretical and experimental results, and reflect how
311 much the results can be expected to generalize to other settings.
- 312 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
313 are not attained by the paper.

314 2. Limitations

315 Question: Does the paper discuss the limitations of the work performed by the authors?

316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368

Answer: [Yes]

Justification: The "Limitations" subsection in Section 4 discusses the limitations of this work.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This is an empirical work.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Data preprocessing steps are provided in the Appendix. Code will be deanonymized upon acceptance.

369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [\[Yes\]](#)

Justification: We will release the code on github after the notification decision.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- 423 • At submission time, to preserve anonymity, the authors should release anonymized
424 versions (if applicable).
425 • Providing as much information as possible in supplemental material (appended to the
426 paper) is recommended, but including URLs to data and code is permitted.

427 6. Experimental Setting/Details

428 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
429 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
430 results?

431 Answer: [Yes]

432 Justification: The details necessary to understand the results can be found in appendices A
433 and B

434 Guidelines:

- 435 • The answer NA means that the paper does not include experiments.
436 • The experimental setting should be presented in the core of the paper to a level of detail
437 that is necessary to appreciate the results and make sense of them.
438 • The full details can be provided either with the code, in appendix, or as supplemental
439 material.

440 7. Experiment Statistical Significance

441 Question: Does the paper report error bars suitably and correctly defined or other appropriate
442 information about the statistical significance of the experiments?

443 Answer: [Yes]

444 Justification: Standard errors are plotted for all results (although they are often too small to
445 see).

446 Guidelines:

- 447 • The answer NA means that the paper does not include experiments.
448 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
449 dence intervals, or statistical significance tests, at least for the experiments that support
450 the main claims of the paper.
451 • The factors of variability that the error bars are capturing should be clearly stated (for
452 example, train/test split, initialization, random drawing of some parameter, or overall
453 run with given experimental conditions).
454 • The method for calculating the error bars should be explained (closed form formula,
455 call to a library function, bootstrap, etc.)
456 • The assumptions made should be given (e.g., Normally distributed errors).
457 • It should be clear whether the error bar is the standard deviation or the standard error
458 of the mean.
459 • It is OK to report 1-sigma error bars, but one should state it. The authors should
460 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
461 of Normality of errors is not verified.
462 • For asymmetric distributions, the authors should be careful not to show in tables or
463 figures symmetric error bars that would yield results that are out of range (e.g. negative
464 error rates).
465 • If error bars are reported in tables or plots, The authors should explain in the text how
466 they were calculated and reference the corresponding figures or tables in the text.

467 8. Experiments Compute Resources

468 Question: For each experiment, does the paper provide sufficient information on the com-
469 puter resources (type of compute workers, memory, time of execution) needed to reproduce
470 the experiments?

471 Answer: [Yes]

472 Justification: In Appendix A.

473 Guidelines:

- 474 • The answer NA means that the paper does not include experiments.
- 475 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 476 or cloud provider, including relevant memory and storage.
- 477 • The paper should provide the amount of compute required for each of the individual
- 478 experimental runs as well as estimate the total compute.
- 479 • The paper should disclose whether the full research project required more compute
- 480 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 481 didn't make it into the paper).

482 9. Code Of Ethics

483 Question: Does the research conducted in the paper conform, in every respect, with the
484 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

485 Answer: [Yes]

486 Justification: The authors have reviewed the NeurIPS Code of Ethics and made sure to
487 follow it

488 Guidelines:

- 489 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 490 • If the authors answer No, they should explain the special circumstances that require a
- 491 deviation from the Code of Ethics.
- 492 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 493 eration due to laws or regulations in their jurisdiction).

494 10. Broader Impacts

495 Question: Does the paper discuss both potential positive societal impacts and negative
496 societal impacts of the work performed?

497 Answer: [Yes]

498 Justification: The "Broader Impacts" subsection in section 4 discusses broader impacts

499 Guidelines:

- 500 • The answer NA means that there is no societal impact of the work performed.
- 501 • If the authors answer NA or No, they should explain why their work has no societal
- 502 impact or why the paper does not address societal impact.
- 503 • Examples of negative societal impacts include potential malicious or unintended uses
- 504 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 505 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 506 groups), privacy considerations, and security considerations.
- 507 • The conference expects that many papers will be foundational research and not tied
- 508 to particular applications, let alone deployments. However, if there is a direct path to
- 509 any negative applications, the authors should point it out. For example, it is legitimate
- 510 to point out that an improvement in the quality of generative models could be used to
- 511 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 512 that a generic algorithm for optimizing neural networks could enable people to train
- 513 models that generate Deepfakes faster.
- 514 • The authors should consider possible harms that could arise when the technology is
- 515 being used as intended and functioning correctly, harms that could arise when the
- 516 technology is being used as intended but gives incorrect results, and harms following
- 517 from (intentional or unintentional) misuse of the technology.
- 518 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 519 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 520 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 521 feedback over time, improving the efficiency and accessibility of ML).

522 11. Safeguards

523 Question: Does the paper describe safeguards that have been put in place for responsible
524 release of data or models that have a high risk for misuse (e.g., pretrained language models,
525 image generators, or scraped datasets)?

526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576

Answer: [NA]

Justification: The paper poses no such risks

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Assets used are the estimator implementations and the pretrained Pythia models Appendix B.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets are introduced in the paper

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

577 Question: For crowdsourcing experiments and research with human subjects, does the paper
578 include the full text of instructions given to participants and screenshots, if applicable, as
579 well as details about compensation (if any)?

580 Answer: [NA]

581 Justification: No human crowdsourcing.

582 Guidelines:

- 583 • The answer NA means that the paper does not involve crowdsourcing nor research with
584 human subjects.
- 585 • Including this information in the supplemental material is fine, but if the main contribu-
586 tion of the paper involves human subjects, then as much detail as possible should be
587 included in the main paper.
- 588 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
589 or other labor should be paid at least the minimum wage in the country of the data
590 collector.

591 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
592 **Subjects**

593 Question: Does the paper describe potential risks incurred by study participants, whether
594 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
595 approvals (or an equivalent approval/review based on the requirements of your country or
596 institution) were obtained?

597 Answer: [NA]

598 Justification: No human experiments.

599 Guidelines:

- 600 • The answer NA means that the paper does not involve crowdsourcing nor research with
601 human subjects.
- 602 • Depending on the country in which research is conducted, IRB approval (or equivalent)
603 may be required for any human subjects research. If you obtained IRB approval, you
604 should clearly state this in the paper.
- 605 • We recognize that the procedures for this may vary significantly between institutions
606 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
607 guidelines for their institution.
- 608 • For initial submissions, do not include any information that would break anonymity (if
609 applicable), such as the institution conducting the review.