IMPROVING NEURON-LEVEL INTERPRETABILITY WITH WHITE-BOX LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Neurons in auto-regressive language models like GPT-2 can be interpreted by analyzing their activation patterns. Recent studies have shown that techniques such as dictionary learning, a form of post-hoc sparse coding, enhance this neuron-level interpretability. In our research, we are driven by the goal to fundamentally improve neural network interpretability by embedding sparse coding directly within the model architecture, rather than applying it as an afterthought. In our study, we introduce a white-box transformer-like architecture named Coding RAte TransformEr (CRATE), explicitly engineered to capture sparse, low-dimensional structures within data distributions. Our comprehensive experiments showcase significant improvements (up to 106% relative improvement) in neuron-level interpretability across a variety of evaluation metrics. Detailed investigations confirm that this enhanced interpretability is steady across different layers irrespective of the model size, underlining CRATE's robust performance in enhancing neural network interpretability. Further analysis shows that CRATE's increased interpretability comes from its enhanced ability to consistently and distinctively activate on relevant tokens. These findings point towards a promising direction for creating white-box foundation models that excel in neuron-level interpretation.

026 027

004

010 011

012

013

014

015

016

017

018

019

021

023

025

028 029 1 INTRODUCTION

Representation learning aims to learn a continuous mapping, to transform a random vector in a high 031 dimensional space that is sampled from a dataset, to a feature vector in another (typically lowerdimensional) space (Bengio et al., 2013). Recently, deep learning has witnessed tremendous empir-033 ical success in modeling massive amounts of high-dimensional data, and the predominant practice 034 has been to learn first a task-agnostic representation by pre-training a large neural network, which is commonly known as the *foundation model* (Devlin et al., 2019; Radford et al., 2019). Among 035 language foundation models, the transformers architecture (Vaswani et al., 2017) with Generative Pre-Training (Radford et al., 2019) (GPT) has recently demonstrated a strong capability of model-037 ing sequential data and thus predicting subsequent tokens (Brown et al., 2020; Ouyang et al., 2022). Such strong capability has emerged significant success in downstream applications (Lewis et al., 2020; Yang et al., 2023), yet the large neural network is known to be *black-box*, where the represen-040 tations in the model are not independently interpretable, introducing difficulty in designing effective 041 paradigms for major known challenges of (visual) language models like hallucination (Ji et al., 042 2023; Tong et al., 2024), bias (Garrido-Muñoz et al., 2021; Nadeem et al., 2020), and catastrophic 043 forgetting (Kemker et al., 2018; Zhai et al., 2024). 044

To interpret the functions of individual modules in the language models, *mechanistic interpretation* was proposed to reverse-engineer such models, through identifying meaningful patterns in the data 046 representations and computational mechanisms of the model components (Olah, 2022; Meng et al., 047 2022a). Recent studies on auto-regressive models like GPT-2 have delved into neuron-level inter-048 pretation, where the focus is on understanding the activations within the model's MLP layers (Yun et al., 2021). This approach helps to reveal the specific roles of individual neurons, which is crucial for precise model editing and control (Meng et al., 2022a;b). Recent research has also proposed 051 that sparse auto-encoder (SAE)-based dictionary learning effectively promotes mono-semanticity of neurons, thus enhancing neuron-level interpretability (Bricken et al., 2023). However, as a post-hoc 052 method, sparse auto-encoders always introduces a non-negative reconstruction loss, which introduces noise for steering the model and imperfect fidelity when interpreting the neurons (Bricken



075

077

Figure 1: Instances are systematically identified where the interpretability of \mathfrak{B} CRATE (ours, row 1) outperforms GPT-2 (row 2). For each neuron (rounded box), we show two top activated text 074 excerpts (excerpt 1 and 2) and one randomly activated excerpt (excerpt 3). Results show that CRATE consistently activates on and only on semantically relevant text excerpts (first two excerpts), leading 076 to more precise explanations predicted by agents like Mistral.

078 et al., 2023). More recent studies also show that sparse auto-encoders are hard to scale up, as there 079 exists a significant amount of directions in a neuron in larger language models, which makes the decomposition difficult (Kissane, 2024; Templeton, 2024; Rajamanoharan et al., 2024).

081 *Can we instead build sparsity directly into the language model?* In this paper, we develop the CRATE language model, a GPT-2-size language model that builds sparse coding into the model with a 083 mathematically principled way. CRATE handles the problems introduced by sparse encoders at scale: 084 it (i) escapes the loss introduced in reconstructing the language model representations, enabling loss-085 free steering, and (ii) escapes the unsteady process of training a sparse auto-encoder. To avoid adding inconsistency into the language model, we develop on top of a mathematically principled white-box 087 model framework, named CRATE (Yu et al., 2023a).¹ After encoding the text tokens into numbers, we apply language-domain-specific modifications to the original CRATE architecture and obtain the token representations. The final representations are then used to predict the next token, while the intermediate representations gets interpreted.

To this end, the main contribution of this work is to propose a causal language model architecture 092 based on the CRATE model framework that builds sparsity inherently, that achieves significantly better neuron-level interpretability (106% relative increase) than language models with the GPT archi-094 tecture under a similar configuration. CRATE forms a family of models from single-layer model up to a 12-layer configuration. Comparative qualitative analysis of neuron activations between CRATE and GPT-2 is provided in Figure 1, alongside extensive quantitative evaluations demonstrating that 096 by explicitly integrating sparse coding into the language model, CRATE achieves markedly improved interpretability across layers compared to GPT-2, applicable across a wide range of model sizes 098 under a variety of evaluation metrics.

2 **RELATED WORK**

102 **Neuron-level Interpretation.** Recent studies have provided insights into how auto-regressive mod-103 els like GPT-2 work at the level of individual neurons, named *neuron-level interpretation*. These 104 studies focus on analyzing *activations*, which are the outputs from the activation functions within 105 the model's multi-layer perceptron (MLP) layers (Yun et al., 2021). Analyzing activations helps

106 107

101

¹In the remaining parts of the paper, we use "CRATE" or "CRATE language model" to refer to our language model architecture, while "original CRATE" denotes the architecture framework described in the literature.

108 uncover the roles of language model components, which is significant to applications like precise 109 modifications and control on the model, known as model editing (Meng et al., 2022a;b). Neuron-110 level interpretation is crucial to understanding the mechanisms in a model, including what concepts 111 are learned in the neurons of the network, whether specific neurons are learning particular concepts, 112 and how localized/distributed and redundantly the knowledge is preserved within neurons of the network. A higher neuron-level interpretability indicates that more neurons are interpretable or neu-113 rons are more interpretable (Sajjad et al., 2022). As interpretations of the neurons can help localize 114 the knowledge obtained in a neural network, neuron-level interpretation can be used for editing the 115 knowledge in models (Meng et al., 2022a;b), model pruning and distillation (Belinkov et al., 2020), 116 adapting the model to different domains and steering the output (Erhan et al., 2009; Rimsky et al., 117 2023), and debugging model errors (Hernandez et al., 2021). Improved neuron-level interpretability 118 increases reliability and performance in the applications above. 119

Sparse Auto-encoders. To enhance interpretability, post-hoc sparse coding methods like dictionary 120 learning (Kreutz-Delgado et al., 2003) are used, but these techniques result in imperfect reconstruc-121 tions and thus always introduces loss when steering the model (Conmy, 2023; Bricken et al., 2023). 122 Literature also indicates that sparse-autoencoders are hard to scale, i.e., a dramatic drop in inter-123 pretable features can be observed when models becomes deeper (Kissane, 2024). Additionally, tun-124 ing SAE models for larger L values involves extensive hyperparameter tuning and time-consuming 125 training, requiring multiple metrics (reconstruction rate, L1 loss, number of dead neurons) for re-126 liable judgment, which can't be easily optimized with automatic engineering tricks (Bricken et al., 127 2023). 128

Evaluation of neuron-level interpretability. Metrics now exist to evaluate neuron-level inter-129 pretability in language models, examing if neurons trigger on relevant tokens in given contexts (Bills 130 et al., 2023; Bricken et al., 2023). Recent works have demonstrated that a small number of circuits 131 in language models are interpretable (Wang et al., 2022; Chughtai et al., 2023), but comprehend-132 ing each neuron, out of millions, is vital for thorough model safety audits. Given the prohibitive 133 cost of human evaluation on such a scale, OpenAI introduced an automated metric using large lan-134 guage models for interpretability assessment (Bills et al., 2023), which Anthropic later refined for 135 sparse activations (Bricken et al., 2023). These methods align closely with human judgment and have gained broad acceptance within the research community (Conmy et al., 2024; Liu et al., 2023; 136 Burns et al., 2023; Lieberum et al., 2023). These metrics show that neuron-level interpretability 137 in auto-regressive models is limited (Sajjad et al., 2022), where the popular hypothesis is that neu-138 rons are superpositions of simpler semantics, which makes them *fire* (produce a high activation) at 139 multiple semantically distinct sets of tokens (Elhage et al., 2022). 140

141 White-box models and structured representation learning. In the domain of structured repre-142 sentation learning, white-box models stand out for their ability to generate explicit, structured data representations that adhere to specific, desirable configurations such as sparsity and piece-wise lin-143 earity, as discussed by Gregor and LeCun (2010) and Chan et al. (2022). Within this framework, Yu 144 et al. (2023a) introduced an innovative approach to constructing deep networks based on unrolled 145 optimization. Specifically, Yu et al. (2023a) proposed the CRATE model, utilizing an information-146 theoretic objective aimed at promoting the *compression and sparsity* of data towards a predefined 147 statistical structure. Recently, empirical experiments suggest that the white-box design of CRATE 148 inherently develops segmentation capabilities from the data representations at both holistic and 149 component levels with supervised training in the vision domain (Yu et al., 2023b), which directly 150 motivates us to further explore the data representations within such architecture for language mod-151 els. Recent work has also shown that the CRATE framework is scalable: it can be effectively scaled 152 up to comparable performance as Vision Transformer (ViT) with careful engineering (Yang et al., 2024). Furthermore, the fine-tuning performance of the pretrained CRATE model is also proven to 153 be comparable in both the language domain (Yu et al., 2023a) and vision domain (Yang et al., 2024). 154

155 156

3 Preliminaries

This section introduces the original CRATE architecture introduced in Yu et al. (2023a).

159 Notations. In this paper, we denote the one-hot input tokens by $X = [x_1, ..., x_N] \in \mathbb{R}^{V \times N}$, **160** where $x_i \in \mathbb{R}^{V \times 1}$ represents the *i*-th one-hot token, N is the total number of input tokens, and **161** V is the vocabulary size. We use $f \in \mathcal{F} : \mathbb{R}^{V \times N} \to \mathbb{R}^{d \times N}$ to denote the mapping induced by the model, which is a composition of L + 1 operators (layers) $f = f^L \circ \cdots f^\ell \circ \cdots f^1 \circ f^{\text{pre}}$,

189

195 196

197

198

199

200

206

207 208 209

162 163 where $f^{\ell} : \mathbb{R}^{d \times N} \to \mathbb{R}^{d \times N} (1 \le \ell \le L)$ represents the mapping of the ℓ -th operator, and f^{pre} : 163 $X \in \mathbb{R}^{V \times N} \to Z^1 \in \mathbb{R}^{d \times N}$ represents the pre-processing layer that transforms the one-hot token 164 representations $X = [x_1, \ldots, x_N]$ to semantic embeddings $Z^1 = [z_1^1, \ldots, z_N^1]$. We let Z^{ℓ} denote 165 the input token representations of the ℓ -th operator f^{ℓ} for $1 \le \ell \le L$, so that $z_i^{\ell} \in \mathbb{R}^d$ denotes the 166 representation of the *i*-th token x_i before the ℓ -th layer. We denote $Z = Z^{L+1}$ as the output token 167 representations of the last (*L*-th) layer.

Framework, objective, and optimization. The transformation of input data into *parsimonious* (piecewise linearized and compact) representations is accomplished by adopting a local signal model for the marginal distribution of the tokens z_i . This statement suggests that the tokens can be approximately considered to occupy a union of several (identified as K) low-dimensional spaces, each with a dimension $p \ll d$. These spaces are characterized by orthonormal bases, represented as $U_{[K]} = (U_k)_{k=1}^K, U_k \in \mathbb{R}^{d \times p}$. Within the framework of this local signal model, CRATE aims to optimize the *sparse rate raduction* objective:

$$\max_{f \in \mathcal{F}} \mathbb{E}_{\boldsymbol{Z}} \left[\Delta R(\boldsymbol{Z} \mid \boldsymbol{U}_{[K]}) - \lambda \| \boldsymbol{Z} \|_{0} \right] = \max_{f \in \mathcal{F}} \mathbb{E}_{\boldsymbol{Z}} \left[R(\boldsymbol{Z}) - \lambda \| \boldsymbol{Z} \|_{0} - R^{c}(\boldsymbol{Z}; \boldsymbol{U}_{[K]}) \right].$$
(1)

where λ is the sparsification regularizer and $\mathbf{Z} = f(\mathbf{X})$. The coding rate $R(\mathbf{Z})$ serves as a close 177 estimate (following Ma et al. (2007)) for the average amount of bits necessary for encoding the 178 tokens z_i to a precision level ε using a Gaussian codebook. Additionally, $R^c(Z \mid U_{[K]})$ represents 179 the theoretical maximum average amount of bits needed to encode the projection of the tokens onto 180 each low dimensional subspace defined in the local signal model, specifically $U_k^* z_i$, to the same 181 precision level ε utilizing a Gaussian codebook, as outlined by Yu et al. (2023a). If the subspaces 182 are adequately incoherent from each other, the solutions that minimize the object function, viz. 183 Equation (1), in terms of Z, are associated with subspace configurations that are both incoherent and aligned with the axes, as pointed out by Yu et al. (2020). 185

A network aimed at optimizing the sparse coding rate reduction objective through unrolled optimization gradually shifts the distribution of X towards the intended canonical forms, where each iteration of the unrolled optimization process acts as a layer.

$$f: \mathbf{X} \xrightarrow{f^{\text{pre}}} \mathbf{Z}^1 \to \dots \to \mathbf{Z}^{\ell} \xrightarrow{f^{\ell}} \mathbf{Z}^{\ell+1} \to \dots \to \mathbf{Z}^{L+1} = \mathbf{Z} \xrightarrow{f^{\text{head}}} \mathbf{Y},$$
(2)

The iterative optimization framework incorporates multiple design choices, among which is a twostep alternating minimization approach grounded in robust theoretical principles (Yu et al., 2023a). This approach delineates two distinct blocks: the MSSA and the ISTA block, collectively defining a single CRATE layer:

$$\boldsymbol{Z}^{\ell+1/2} \doteq \boldsymbol{Z}^{\ell} + \mathsf{MSSA}(\boldsymbol{Z}^{\ell} \mid \boldsymbol{U}_{[K]}^{\ell}), \qquad f^{\ell}(\boldsymbol{Z}^{\ell}) = \boldsymbol{Z}^{\ell+1} \doteq \mathsf{ISTA}(\boldsymbol{Z}^{\ell+1/2} \mid \boldsymbol{D}^{\ell}). \tag{3}$$

Compression operator: Multi-Head Subspace Self-Attention (MSSA). Given local models $U_{[K]}^{\ell}$, to form the incremental transformation f^{ℓ} optimizing Equation (1) at layer ℓ , CRATE first compresses the token set Z^{ℓ} against the subspaces by minimizing the coding rate $R^{c}(\cdot \mid U_{[K]}^{\ell})$. As Yu et al. (2023a) show, doing this minimization locally by performing a step of gradient descent on $R^{c}(\cdot \mid U_{[K]}^{\ell})$ leads to the so-called multi-head subspace self-attention (MSSA) operation, defined as

$$\mathsf{MSSA}(\boldsymbol{Z} \mid \boldsymbol{U}_{[K]}) \doteq \frac{p}{(N+1)\varepsilon^2} \left[\boldsymbol{U}_1, \dots, \boldsymbol{U}_K \right] \begin{bmatrix} (\boldsymbol{U}_1^*\boldsymbol{Z}) \operatorname{softmax} \left((\boldsymbol{U}_1^*\boldsymbol{Z})^* (\boldsymbol{U}_1^*\boldsymbol{Z}) \right) \\ \vdots \\ (\boldsymbol{U}_K^*\boldsymbol{Z}) \operatorname{softmax} \left((\boldsymbol{U}_K^*\boldsymbol{Z})^* (\boldsymbol{U}_K^*\boldsymbol{Z}) \right) \end{bmatrix}, \quad (4)$$

In practice, the calculation of the intermediate representations $Z^{\ell+1/2}$ with the output from the MSSA block is calculated in a weighted form:

$$\boldsymbol{Z}^{\ell+1/2} \approx \left(1 - \kappa \cdot \frac{p}{(N+1)\varepsilon^2}\right) \boldsymbol{Z}^{\ell} + \kappa \cdot \frac{p}{(N+1)\varepsilon^2} \cdot \mathsf{MSSA}(\boldsymbol{Z}^{\ell} \mid \boldsymbol{U}_{[K]}), \tag{5}$$

where $\kappa > 0$ is a learning rate hyperparameter. This block resembles to GPT's multi-head selfattention block, but the query, key, and value projection matrices within a single head are all identical in the MSSA block.

214 Sparsification operator: Iterative Shrinkage-Thresholding Algorithm (ISTA). The remaining 215 term to optimize in Equation (1) is the difference of the global coding rate R(Z) and the ℓ^0 norm of the tokens, which together encourage the representations to be both sparse and non-collapsed. Yu et al. (2023a) show that local minimization of this objective in a neighborhood of the intermediate representations $Z^{\ell+1/2}$ is approximately achieved by a LASSO problem with respect to a sparsifying orthogonal dictionary $D^{\ell} \in \mathbb{R}^{d \times h}$. Taking an iterative step towards solving this LASSO problem gives the iterative shrinkage-thresholding algorithm (ISTA) block (Wright and Ma, 2022; Yu et al., 2023a). The ReLU nonlinearity appearing in this block arises from an additional nonnegativity constraint on the representations in the LASSO program, motivated by the goal of better separating distinct modes of variability in the token distribution:

226

234

240

247

253

254 255 256

257 258

$$\boldsymbol{Z}^{\ell+1} = f^{\ell}(\boldsymbol{Z}^{\ell}) = \operatorname{ReLU}(\boldsymbol{Z}^{\ell+1/2} + \eta \boldsymbol{D}^{\ell*}(\boldsymbol{Z}^{\ell+1/2} - \boldsymbol{D}^{\ell}\boldsymbol{Z}^{\ell+1/2}) - \eta\lambda \boldsymbol{1}) \doteq \operatorname{ISTA}(\boldsymbol{Z}^{\ell+1/2} \mid \boldsymbol{D}^{\ell}).$$
(6)

4 THE CRATE LANGUAGE MODEL

This section introduces the difference between our work and the original CRATE paper (Yu et al., 2023a), thus introducing what changes we made to the CRATE architecture. We first note that the task in this work is different: we apply the CRATE architecture to the *next-token prediction task* in the language domain. while the original CRATE paper applies the architecture to the image classification task in the vision domain. This difference leads to differences in the architecture design: (i) we apply a causal mask to the original CRATE model to avoid the model seeing tokens after the current token, and (ii) we change the embedding layer and heads of the original CRATE model.

Second, we're interested in interpreting the neurons within CRATE on the next-token prediction task and making direct comparisons to the GPT architecture. As neuron-level interpretation is commonly evaluated on the hidden states in the FFN block of the GPT model (Bills et al., 2023; Bricken et al., 2023), (iii) we increase the hidden dimension of the ISTA block of the original CRATE model to align with the hidden dimension of the FFN block of the GPT model. We thus call the new ISTA block the ISTA-overcomplete block.

Below we show specific definitions of the modifications we made. We illustrate the architecture in Figure 6, show implementation details in Appendix A, and discuss about the learning process in Appendix A.5.

Embedding and Head. In order to apply the CRATE architecture to the language domain, we define the pre-processing layer f^{pre} that transforms tokens into position-aware semantic embeddings, and define post-processing head $f^{\text{head}}(Z)$ that maps the representations to output token distributions:

$$F^{\text{pre}}(\boldsymbol{X}) = \boldsymbol{E}_{\text{sem}}(\boldsymbol{X}) + \boldsymbol{E}_{\text{pos}}, \qquad f^{\text{head}}(\boldsymbol{Z}) = \boldsymbol{W}^{\text{head}}\boldsymbol{Z},$$
 (7)

where E_{sem} is a semantic embedding matrix that maps input tokens x_i to embedding vectors in \mathbb{R}^d , $E_{\text{pos}} \in \mathbb{R}^{d \times N}$ is a positional embedding matrix, and $W^{\text{head}} \in \mathbb{R}^{V \times d}$ maps the (contextualized) token representations Z^{L+1} to the distribution of the next token. All parameters mentioned are learnable.

MSSA Block. To align with the next word prediction task used in GPT (Radford et al., 2019), we replace the attention matrix in MSSA (Equation (4)) with a *causally masked* self-attention, defined as

$$\operatorname{softmax}\left((\boldsymbol{U}_{k}^{*}\boldsymbol{Z})^{*}(\boldsymbol{U}_{k}^{*}\boldsymbol{Z})\right) \to \operatorname{softmax}\left(\operatorname{CausalMask}\left((\boldsymbol{U}_{k}^{*}\boldsymbol{Z})^{*}(\boldsymbol{U}_{k}^{*}\boldsymbol{Z})\right)\right),$$
where
$$\operatorname{CausalMask}(\boldsymbol{M})_{ij} = \begin{cases} \boldsymbol{M}_{ij}, & i \leq j \\ -\infty, & i > j. \end{cases}$$
(8)

ISTA Block. To investigate the neuron interpretability of the activation matrix $A \in \mathbb{R}^{h \times N}$, we design an *overcomplete* version of the ISTA block (Equation (6)) with $D^{\ell} \in \mathbb{R}^{d \times h}$ where h = nd, and n = 4 to keep a fair comparison to GPT (also same as standard transformer architecture proposed in Vaswani et al. (2017)):

264

265 266

$$\boldsymbol{A}_{t} \doteq \text{ISTA}(\boldsymbol{Z}^{\ell+1/2} \mid \boldsymbol{D}^{\ell}),$$

$$\boldsymbol{A}_{k} = \text{ReLU}(\boldsymbol{A}_{k-1} - \eta(\boldsymbol{D}^{\ell})^{*}(\boldsymbol{D}^{\ell}\boldsymbol{A}_{k-1} - \boldsymbol{Z}) - \eta \cdot \lambda \cdot \boldsymbol{1}), \quad \boldsymbol{A}_{0} = \boldsymbol{0}, \ k \in [t], \qquad (9)$$

$$\boldsymbol{Z}^{\ell+1} = \boldsymbol{D}^{\ell}\boldsymbol{A}_{t}$$

267

Here, $\eta > 0$ is the step size, $\lambda > 0$ is the sparsification regularizer, and t is the number of ISTA iterations. In practice, we set t = 2 to keep computation efficient. The ISTA block resembles the MLP block in the GPT model, but with a relocated skip connection.

270 5 EMPIRICAL EXPERIMENTS 271

We examine the next token prediction performance and neuron interpretability of CRATE in this section. We detail the architecture, size, pre-training recipe (Section 5.1), performance (Section 5.2), and neuron-level interpretability (Section 5.3) of CRATE compared to the standard transformer architecture. In this section, we denote K as the number of attention heads, d as the dimension of the residual stream in the model, and h as the hidden (inner) dimension of the ISTA/MLP module.

5.1 Setup

278

279

287

289

291 292 293

294

Model architecture and size. The CRATE model is designed with various sizes $L \in \{1, 2, 3, 6, 12\}$, where each size matches the GPT configurations for direct comparisons, as shown in Section 5.1. Configurations for $L \in \{1, 2, 3\}$ adhere to GPT models as per Bricken et al. (2023), while $L \in \{6, 12\}$ follow configurations from Sanh et al. (2019) and Radford et al. (2019), respectively. Notably, CRATE maintains approximately 2/3 the size of GPT at scale. Both models utilize the Bytelevel BPE tokenizer with a 50,257 vocabulary size, following Radford et al. (2019). We explain the difference between CRATE and GPT in parameter size to GPT in Appendix A.1.

Model Config	d	K	L	h	CRATE	GPT
1L	128	4	1	512	6.54M	6.64M
2L	128	4	2	512	6.64M	6.83M
3L	128	4	3	512	6.74M	7.03M
S(mall)	768	12	6	3,072	55.9M	81.1M
B (ase)	768	12	12	3,072	81.2M	123.6M

Table 1: Model configuration of CRATE and model size comparison to GPT.

Datasets and optimization. We pre-train both models using the next token prediction task on the uncopyrighted *Pile* dataset (Gao et al., 2020) using the Adam optimizer (Kingma and Ba, 2015). Following Bricken et al. (2023), we pre-train both CRATE and GPT of smaller sizes ($L \in \{1, 2, 3\}$) on 100 billion tokens with a context window of 1,024 tokens. Following the pre-training setup in Karpathy (2022) and scaling law in Touvron et al. (2023), we pre-train using 100 billion tokens for the Small models, and 160 billion tokens for the Base models.² It takes 4 days to pre-train CRATE-Base on 160 billion tokens with 32 A5000 GPUs.

301 302 303

304

5.2 Performance

This section demonstrates that CRATE, despite not outperforming GPT-2, still generates reasonable predictions, as evidenced through quantitative and qualitative comparisons.

We observe that both training and validation loss curve of CRATE-Base on the Pile dataset *converges* well, as presented in Figure 2 (*left*). Although the convergence is slower than GPT, the loss curve of CRATE keeps decreasing after training on 160 billion tokens, while GPT already tends to converge.

We also demonstrate the zero-shot validation loss curve of CRATE evaluated on OpenWebText as 311 well as other datasets (Radford et al., 2019) in Figure 2 (right). Results show that CRATE effectively 312 learns transferable representations across a number of datasets, and achieves comparable perfor-313 mance to GPT after full training on the 160 billion tokens. We also demonstrate the scalability of 314 the CRATE architecture by comparing the validation loss of CRATE and GPT with respect to the 315 model size in Figure 3 (*left*). Results show that the performance of CRATE is close to GPT across 316 all model sizes. However, we do recognize that forcing sparsification in a model potentially leads 317 to a higher compute cost on the next-token-prediction objective, which aligns with observations 318 in Bricken et al. (2023) that enabling monosemanticity might hurt model performance.

Qualitative examples from CRATE and GPT are demonstrated in Figure 3 (*right*). We conclude that CRATE can make reasonable predictions, encouraging us to further look into its neuron-level interpretability.

³²² 323

²Practicaly, we train with a batch size of 768 for 125,000 iterations for $L \in \{1, 2, 3, 6\}$, and a batch size of 256 for 600,000 iterations for L = 12.



Figure 2: Left: loss curve when pre-training CRATE-Base and GPT-Base on the Pile dataset. Right: zero-shot validation loss of CRATE evaluated on a variety of datasets (Pile, LAMBADA, OpenWebText and WikiText).



Figure 3: Left: Validation loss of CRATE compared to GPT on the Pile dataset, with respect to the model size. *Right:* Qualitative examples of predictions made by our models and the official models. The tokens in blue are considered good. We compare I CRATE-Base to GPT2-Base on the next word prediction task.

5.3 INTERPRETABILITY

In order to quantitatively evaluate the interpretability of the neuron activations, we adopt the large language model-based approach introduced in Bills et al. (2023) and Bricken et al. (2023). We demonstrate the algorithm for scoring interpretations in Algorithm 1. We retrieve the sparse code A_t (activations after the ReLU unlinearly in the ISTA block) of CRATE for interpretation, and compare with activations from the MLP block of GPT.

Algorithm 1 Interpretability Evaluation Algorithm (Bills et al., 2023) 1: Inputs: Input token set S (in text form) and its activation matrix $A \in \mathbb{R}^{h \times T \times B}$ at ℓ -th layer, where T is the length of a single text excerpt, and B is the number of text excerpts in the corpus.

Models: Explanation model \mathcal{F}_1 , simulation model \mathcal{F}_2 . 2:

3: for $i \in [d]$ do

 $S' \sim S, A' \in \mathbb{R}^{h \times T \times b} \sim A$: Retrieve b text excerpts of T tokens, together with the 4: corresponding activation matrices.

 $k_i = \mathcal{F}_1(S', A'_{i,*})$: Explain common patterns retrieved activations of *i*-th neuron. 5:

 $\tilde{A}'_{i*} = \mathcal{F}_2(k_i, S')$: Use the explanation to *simulate* scores given only the tokens, not 6: including true activations.

 $\rho_i = \rho(A'_{i,*}, A'_{i,*})$: Calculate *correlation* between the accurate and simulated activations. 7: 8: end for

9: **Output:** Averaged interpretation score over all neurons $s = \mathbb{E}_{i \in [d]}(\rho_i)$.

374 In practice, we adopt three evaluation metrics: two from OpenAI (top-and-random and random-375 only) (Bills et al., 2023) and one from Anthropic (Bricken et al., 2023). We adopt the official 376 implementation from Wu et al. (2023), where details on the implementation are elaborated in Appendix B. Note that the Anthropic metric has much shorter text excerpts than the OpenAI metrics, 377 so it is biased to sparse activations. For all evaluations, we discard the last layer of CRATE and

348

349

350 351 352

353

354

355

356

357

358 359

360

361

362

364

366

367

368

369

370

372 373

336

337

378		Mean (↑, darker green means more interpretable)							Variance $(\downarrow$, darker red means less steady)					
379		Top-and-	Random	Randor	n-only	Anthr	Anthropic T		Top-and-Random		Random-only		ropic	
200		CRATE	GPT	CRATE	GPT	CRATE	GPT	CRATE	GPT	CRATE	GPT	CRATE	GPT	
300	1L	3.9	8.8	4.8	8.9	10.1	14.2	0.0	0.0	0.0	0.0	0.0	0.0	
381	2L	8.05	4.2	6.95	1.95	11.35	10.2	0.06	0.01	1.1	0.12	0.0	0.25	
382	3L	9.1	3.57	8.43	1.37	11.23	9.2	0.26	7.51	1.2	1.93	1.14	19.21	
202	6L	7.96	5.4	6.36	3.14	10.4	8.52	2.29	20.85	1.87	18.39	2.01	32.56	
303	12L	6.8	6.34	5.12	2.67	8.88	8.65	7.09	11.35	2.83	7.48	18.3	24.65	



Figure 4: Interpretation scores evaluated using the OpenAI Random-only metric, Top-and-Random metric, 412 and Anthropic metric, respectively. Top: interpretation scores of CRATE and GPT for L = 12. Middle: 413 interpretation scores of CRATE and GPT for L = 6. Bottom: interpretation scores of CRATE, GPT, and GPT 414 with sparse auto-encoder for $L \in \{1, 2, 3\}$. Variance bars are normalized to 1/10 of its original size.

379 380 381

383 384 385

GPT, according to the empirical observation that the last layer of CRATE is biased to the pre-training 417 task (Yu et al., 2023b). 418

419 CRATE achieves markedly improved and more steady neuron-level interpretability across 420 layers compared to GPT-2, applicable across a wide range of model sizes. We show evalua-421 tion results of the interpretability of CRATE and GPT averaged across layers in Table 2 (*left*). We 422 observe that the interpretability of CRATE comprehensively outperforms GPT on all metrics for 423 $L \in \{2, 3, 6, 12\}$. When averaging the mean interpretability across all metrics, CRATE outperforms GPT up to strikingly 45.1% when L = 6, and up to 16.3% when L = 12. We also present the layer-424 wise interpretation scores in Figure 4, which shows that CRATE has higher interpretability than GPT 425 on almost all layers using the OpenAI metrics, and is slightly better than GPT using the Anthropic 426 metric. For detailed distributions of the layer-wise scores of CRATE-Base compared to GPT-Base on 427 different metrics, refer to Appendix D. 428

429 The variances of the average interpretation scores of CRATE and GPT across layers are shown in Table 2 (right). From the results we draw a solid conclusion that the interpretability of CRATE is much 430 more steady than GPT across all model sizes. Figure 4 further demonstrates a clear pattern that, for 431 all model sizes, CRATE maintains a higher interpretability than GPT among almost all layers.

432 The built-in sparse coding approach introduces consistent and specific neuron-level behaviors. 433 The strong interpretability of CRATE on the OpenAI top-and-random metric and the Anthropic met-434 ric, as shown in Figure 4, indicates its consistent behavior on relevant tokens. These two methods 435 contain a large portion of top-activated text excerpts, so they are valid for measuring whether the 436 activations are consistent with the summarized explanation (Bills et al., 2023; Bricken et al., 2023). Additionally, the larger interpretability gap of CRATE and GPT on the OpenAI random-only metric 437 versus the top-and-random metric highlights the specificity of CRATE in avoiding firing on irrele-438 vant tokens. The random-only metric usually includes highly irrelevant text excerpts, so it effec-439 tively measures the capability of the language model to avoid activating on semantically irrelevant 440 tokens (Bills et al., 2023). 441

442 Qualitatively, we refer back to the qualitative examples shown in Figure 1. We list three neurons from CRATE (row 1) and GPT (row 2), respectively. For each neuron, we show two top-activated 443 text excerpts and one random excerpt. Results show that CRATE is able to consistently activate on 444 sementically similar tokens within the most relevant text excerpts, and does not activate on random 445 tokens that are semantically distinguished from the top tokens. This promotes a more precise ex-446 planation given by the explanation model (Mistral in the figure). On the other hand, GPT is much 447 worse at distinguishing tokens from different contexts, because it also has high activations on ran-448 dom text excerpts where the semantic meanings deviate from the top activations a lot. As a side 449 note, we also analyze the activation sparsity of CRATE and GPT in Appendix C. 450

Comparing CRATE to GPT with post-hoc sparse auto-encoders. We follow Bricken et al. (2023) and train SAEs for models with layers $L \in \{1, 2, 3\}$, using output activations from GPT on the Pile dataset's training split, leading to the GPT-SAE model. Details on the SAEs' architecture and training are in Appendix E.

455 The interpretability scores of GPT-SAE compared to CRATE and GPT, as depicted in Figure 4, reveal that under the long-context OpenAI metrics, GPT-SAE matches GPT but falls short of CRATE. This 456 is attributed to its neuron activations becoming nearly 99% sparse after sparse auto-encoding, dimin-457 ishing interpretability in long contexts. Conversely, under the Anthropic metric, GPT-SAE surpasses 458 both GPT and CRATE in interpretability, corroborating findings in Bricken et al. (2023) that post-hoc 459 approaches enhance short-context interpretability, often a sign of mono-semanticity. However, the 460 interpretability of GPT-SAE on the Anthropic metric decreases significantly when ℓ increases, while 461 CRATE remains steady. Further qualitative comparisons are can be found in Appendix F. 462

Besides its good performance on the Anthropic metric, the post-hoc dictionary learning approach requires considerable *manual effort*. To get a taste, training a sparse auto-encoder for a single GPT layer takes 4 hours when h = 512 and a day when h = 3072 on an A100 GPU.

	OpenAI TaR	Anthropic
CRATE-SAE - CRATE	-10.2	+34.8
GPT-SAE - GPT	+6.5	+38.1

Table 3: Interpretability improvement of CRATE and GPT after applying SAE. Results are obtained by subtracting the interpretation scores of the language model and the SAE model trained on that language model. Results consistently show that the interpretability improvement of CRATE-SAE over CRATE is smaller than GPT-SAE over GPT, indicating more optimal representations of CRATE over GPT.

Does CRATE have more optimal representation than GPT in terms of interpretability? Alternatively, we train SAE models upon the CRATE model, and compare the interpretability improvement of CRATE-SAE over CRATE to the interpretability improvement of GPT-SAE over GPT. As shown in Table 3, the improvement of interpretability of CRATE-SAE over CRATE is smaller than GPT-SAE over GPT under both OpenAI and Anthropic metrics. This suggests that CRATE has more optimal representations than GPT in terms of interpretability. Experimental details can be found in Appendix E.

481

482 5.4 DISENTANGLING INTERPRETABILITY FROM OTHER FACTORS483

Is the improved interpretability due to performance gap? We first investigate whether the interpretability improvement is due to worse performance by comparing the interpretability of two different checkpoints of CRATE (intermediate checkpoint and full checkpoint). Results in Table 4

491

500

501

502

Checkpoint	Loss	0	1	2	3	4	5	6	7	8	9	10	11	Avg
79B Tokens	2.38	13.9	8.3	6.5	6.0	4.3	6.7	5.1	4.6	3.7	5.2	8.0	2.8	6.3
158B Tokens	2.29	13.6	8.2	6.7	5.4	5.0	7.2	5.1	5.1	3.6	5.5	9.4	5.3	6.7





Figure 5: Qualitative examples on logit effects of manually activating feature (i) in CRATE. Text shown on the right side are the most positive changes in token prediction probability. The logit effects align with feature interpretations.

show that the interpretability of the intermediate checkpoint is lower than the full checkpoint, which
 suggests that "sacrificing" performance does *not* necessarily introduce better interpretation scores.

Another piece of evidence is that CRATE-2L has higher interpretability scores than CRATE-1L. As shown in Figure 3 (*left*), CRATE-2L has much better performance than CRATE-1L on the next-token prediction task. On the other hand, as shown in Table 2, the interpretability of CRATE-2L is also much higher than CRATE-1L. Thus, a lower performance does not necessarily introduce higher interpretability.

Is the interpretability gap due to number of parameters? We observe that CRATE-Base (81.2M) has a similar number of parameters as GPT-Small (81.1M). However, results in Table 2 indicate that the interpretability of CRATE-Base is higher than GPT-Small on all metrics, and their layer-wise interpretation scores are also different in Figure 4. This evidence suggests that two models with similar number of parameters does not necessarily have similar interpretability.

Another piece of evidence is that the interpretability of CRATE/GPT-1L all the way up to
 CRATE/GPT-12L does not have a consistent trend of increasing/decreasing interpretability, but their
 number of parameters both monotonously increases. This indicates that a model with larger number
 of parameters does not necessarily has better/worse interpretability.

Steering the CRATE model. Following Bricken et al. (2023), we manually activate some neurons and observe the *logit effects* (changes of the token probability of the language model head). Some qualitative examples are shown in Figure 5. Compared to the lossy steering of the SAE models, CRATE are steered without loss. Discussions on the lossy steering process can be found in Appendix E.6.

⁵²⁵ 6 CONCLUSION, LIMITATION, AND FUTURE WORK

In this paper, we demonstrated that replacing the standard transformer architecture with the white box model CRATE as a foundational architecture significantly improves the interpretability. Our
 empirical findings on the capability of CRATE to be consistent and distinctive on the neuron-level
 activations underscore the importance of the white-box design in developing better language foun dation models, fostering optimism that the introduction of built-in sparse coding approaches will
 catalyze further advancements in neuron-level interpretations.

Despite these findings, we acknowledge that the performance of CRATE is not as good as GPT on
the next-token prediction task, which is potentially due to the introduction of the ISTA operator
that introduces sparsity. This aligns with previous work suggesting that the performance might drop
when explicitly introducing sparsity (Bricken et al., 2023). Future work should investigate towards a
better trade-off between performance and interpretability of language models with built-in sparsity.
It would also be meaningful to research on more qualitative mechanisms in the white-box language
model and how to use these mechanisms for downstream edits.

540 REPRODUCIBILITY STATEMENT

541 542 543

544

546 547

548

556 557 To facilitate reproducibility of our work, we will open-source the model checkpoints and training infrastructure. We have included the model architecture in Section 4, pre-training recipe (including dataset and hyper-parameters) in Section 5.1, interpretability evaluation in Section 5.3, and SAE training setup in Appendix E.2.

ETHICS STATEMENT

⁵⁴⁹ By improving the interpretability of language models, our work promotes a deeper understanding of their mechanisms, aiding in the identification and mitigation of potential risks, thereby supporting transparency and responsible AI development. On the language model side, this research pretrains a GPT-2-sized language model on publicly available data, with no intentional inclusion of harmful content. The model's moderate size and data scope reduce the likelihood of generating harmful or out-of-distribution outputs. However, risks associated with intentionally training models on harmful datasets, which can lead to biased or unsafe generations, must be considered.

- References
- Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Hassan Sajjad, and James Glass. On the linguistic representational power of neural machine translation models. *Computational Linguistics*, 46(1): 1–52, 2020.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. Language models can explain neurons in language models. URL https://openaipublic. blob. core. windows. net/neuron-explainer/paper/index. html.(Date accessed: 14.05. 2023), 2023.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, page 2, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-573 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-574 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 575 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, 576 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCan-577 dlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-shot 578 Learners. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, 579 and Hsuan-Tien Lin, editors, Advances in Neural Information Processing Systems 33: An-580 nual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html. 582
- ⁵⁸³ Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner, Bowen Baker, Leo Gao, Leopold Aschenbrenner, Yining Chen, Adrien Ecoffet, Manas Joglekar, Jan Leike, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023.
- Kwan Ho Ryan Chan, Yaodong Yu, Chong You, Haozhi Qi, John Wright, and Yi Ma. ReduNet: A White-box Deep Network from the Principle of Maximizing Rate Reduction. *Journal* of Machine Learning Research, 23(114):1–103, 2022. URL http://jmlr.org/papers/v23/ 21-0631.html.
- Bilal Chughtai, Lawrence Chan, and Neel Nanda. A toy model of universality: Reverse engineering
 how networks learn group operations. In *International Conference on Machine Learning*, pages
 6243–6267. PMLR, 2023.

594	Arthur	Conmy	Mv	hest	oness	at	the	important	tricks	for	train-
595	ing	11 saes	IVIY	httn	s·//www	less	wrong	com/posts	/fifPCos	6ddsmT	YahD/
596	mv-be	est-guess-at-t	he-impo	rtant-	tricks-	for-t	raini	ng-11-saes	. 2023.	oudoino	ranb,
597				· · · · · · · · · · · · · · · · · · ·					, 20201		
598	Arthur C	Conmy. sae. http	s://git	hub.co	m/Arthu	rConm	y/sae	, 2024.			
599											
600	Arthur C	Conmy, Augustine	e Mavor-	Parker,	Aengus	Lynch,	Stefa	n Heimershe	im, and A	Adrià Ga	arriga-
601	Alons	o. Towards autom	nated circ	uit disc	overy for	mecha	anistic	interpretabil	ity. Advai	nces in l	Veural
602	Inform	nation Processing	Systems	, 36, 20	024.						
603	Jacob D	evlin Ming Wei (⁷ hang K	Centon I	ee and	Kristin	a Tou	tanova BEP	T. Dra tra	ining of	Deen
604	Bidire	ectional Transform	ners for	Langua	are Unde	retand	ing I	n Iill Burste	in Christ	v Dora	n and
605	Tham	ar Solorio editor	s Proce	edinos	of the 20	13000	nferer	ice of the No	orth Amer	rican C	hanter
600	of the	Association for (S, 1 roce Computa	tional I	inguistic.	$s \cdot Hun$	nan L	nguage Tech	nologies	NAACI	L-HLT
000	2019.	Minneapolis. MN	J. USA	lune 2-1	7. 2019. V	Volume	$\frac{1}{2}$ 1 (L	ong and Shor	t Papers)	. pages	4171_
607	4186.	Association for	Computa	ational	Linguisti	cs. 20	19. d	oi: 10.1865.	3/V1/N19	-1423.	URL
608	https	s://doi.org/10.	18653/\	/1/n19·	-1423.	,					
609		U U									
610	Nelson H	Elhage, Tristan Hu	ime, Cat	herine (Olsson, N	ichola	s Schie	efer, Tom He	nighan, Sl	hauna K	ravec,
611	Zac H	latfield-Dodds, Ro	obert Las	senby, I	Dawn Dra	ain, Ca	rol Cł	nen, et al. To	y models	of supe	erposi-
612	tion. a	arXiv preprint arX	Kiv:2209	.10652,	2022.						
613	D	Eshan Vashaa I	Demeio	A	C	a an al	D 1	Vices	7:1::	. 1.: . 1	. 1
614	footur	Ernan, rosnua i	orly Uni	Aaron	of Montry	a_{al} 12	$A_1(2)$	1 2000	isualizing	g mgnei	r-layer
615	Teatur	es of a deep netwo	OIK. Om	versity	oj monire	eai, 15	41(5).	1, 2009.			
616	Leo Gao	o. Stella Biderma	n. Sid I	Black, I	Laurence	Goldi	ng. Ti	ravis Hoppe.	Charles	Foster.	Jason
617	Phang	, Horace He, Anis	sh Thite,	Noa Na	abeshima	, et al.	The p	ile: An 800g	b dataset	of diver	se text
618	for la	nguage modeling.	arXiv p	reprint	arXiv:21	01.000	27, 20	20.			
619		0 0 0									
620	Ismael C	Garrido-Muñoz, A	rturo Mo	ontejo-H	Ráez, Fer	nando	Martí	nez-Santiago	, and L A	lfonso I	Jreña-
621	López	z. A survey on bia	is in deep	o nlp. A	pplied Sc	ciences	, 11(7):3184, 2021	•		
622	Varal C		Com	T	Track A		:		Calina	L. L.L	
623	Karor G	regor and Thorston	LeCun.	Learnii na adit	ig rast A	Approx	imatic	his of sparse	e Counig.		lannes
624	Fullik	ing Lagrange (ICA	I JOaching III JOaching III JOaching III JOaching III JOACHING	118, eutonometricon 100 m m m m m m m m m m m m m m m m m m	018, F100 24 2010	eeung Haife	s Of II	$\frac{10}{2}$	406 Om	,onjerei ninross	2010
625	IIRI	https://icml.co	(Confe	rences	/2010/0	, mars	ι, Israe / <u>/ </u>	ndf	-400. Om	mpress,	2010.
626	UKL		c/ com c	renees	7 20107 p	aper 5	/ ++5.	pur.			
627	Evan He	ernandez, Sarah S	chwettm	ann, Da	avid Bau,	Teona	a Baga	shvili, Antoi	nio Torral	ba, and	Jacob
628	Andre	as. Natural langu	age des	cription	s of deep	visual	l featu	res. In Intern	national (Conferen	nce on
620	Learn	ing Representatio	ns, 2021		1					U	
620						_					_
621	Ziwei Ji	, Nayeon Lee, R	ita Fries	ke, Tie	zheng Yu	ı, Dan	Su, Y	ran Xu, Etsu	iko Ishii,	Ye Jin	Bang,
031	Andre	a Madotto, and P	ascale Fi	ing. Su	rvey of ha	allucin	atıon i	n natural lang	guage ger	eration.	ACM
632	Comp	uting Surveys, 55	(12):1-3	8, 2023	•						
633	Andrei I	Carnathy nanoGE	T http	s · //ai	thub co	m/karı	nathy	/nanoCPT 20	122		
634	/ marcj i	arpany. nanoor	1. 1100	5.7751	1100.00		Jucity	/ 101001 1, 20	<i>JLL</i> .		
635	Ronald I	Kemker, Marc Mc	Clure, A	ngelina	Abitino,	Tyler	Hayes	, and Christo	pher Kana	an. Mea	suring
636	catast	rophic forgetting	in neura	l netwo	rks. In F	Proceed	dings of	of the AAAI	conferenc	e on ar	tificial
637	intelli	gence, 2018.						•	•		
638											
639	Diederik	P. Kingma and	Jimmy I	Ba. Ad	am: A N	/lethod	for S	tochastic Op	otimizatio	n. In Y	oshua
640	Bengi	o and Yann LeC	un, edito	ors, $3rd$	Internat	ional (Confei	rence on Lea	rning Re	presente	ations,
641	ICLR	2015, San Diego	, CA , US	A, Ma	y 7-9, 20	13, Ca	onferer	ice Irack Pro	oceedings	, 2015.	URL
642	http:	//arxiv.org/ab	s/1412	6980.							
643	Connor	Kissane Attent	ion saes	scale	to ont-?	small	htt	DS. / / WIMIM 14	esswrong		oste/
644	FSTRe	dtiuHa4Gfdbr/a	attentio	n-sae	s-scale-	-to-o	ot-2-	small 2024		. com/ p	0303/
645	1 3110				c coure	61	~~ ~				
646	Kenneth	Kreutz-Delgado,	Joseph	F Murra	ay, Bhask	ar D I	Rao, K	jersti Engan	, Te-Won	Lee, an	d Ter-
647	rence 15(2):	J Sejnowski. Dict 349–396, 2003.	ionary le	earning	algorithm	ns for s	parse	representatio	n. <i>Neurai</i>	сотри	tation,

648 649 650 651	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33: 9459–9474, 2020
652 653 654 655	Tom Lieberum, Matthew Rahtz, János Kramár, Geoffrey Irving, Rohin Shah, and Vladimir Miku- lik. Does circuit analysis interpretability scale? evidence from multiple choice capabilities in chinchilla. <i>arXiv preprint arXiv:2307.09458</i> , 2023.
656 657 658	Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. Trustworthy llms: a survey and guideline for evaluating large language models' alignment. <i>arXiv preprint arXiv:2308.05374</i> , 2023.
659 660 661 662	Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of Multivariate Mixed Data via Lossy Data Coding and Compression. <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> , 29(9):1546–1562, 2007. doi: 10.1109/TPAMI.2007.1085. URL https://doi.org/10.1109/TPAMI.2007.1085.
663 664	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual asso- ciations in gpt. <i>Advances in Neural Information Processing Systems</i> , 35:17359–17372, 2022a.
665 666	Kevin Meng, Arnab Sen Sharma, Alex Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. <i>arXiv preprint arXiv:2210.07229</i> , 2022b.
668 669	Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained language models. <i>arXiv preprint arXiv:2004.09456</i> , 2020.
670 671 672	Chris Olah. Mechanistic interpretability, variables, and the importance of interpretable bases. <i>Trans- former Circuits Thread</i> , page 2, 2022.
673 674 675 676	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35: 27730–27744, 2022.
677 678	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9, 2019.
679 680 681 682	Senthooran Rajamanoharan, Tom Lieberum, Nicolas Sonnerat, Arthur Conmy, Vikrant Varma, János Kramár, and Neel Nanda. Jumping ahead: Improving reconstruction fidelity with jumprelu sparse autoencoders. <i>arXiv preprint arXiv:2407.14435</i> , 2024.
683 684	Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. Steering llama 2 via contrastive activation addition. <i>arXiv preprint arXiv:2312.06681</i> , 2023.
685 686 687	Hassan Sajjad, Nadir Durrani, and Fahim Dalvi. Neuron-level interpretation of deep nlp models: A survey. <i>Transactions of the Association for Computational Linguistics</i> , 10:1285–1303, 2022.
688 689	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> , 2019.
690 691 692	Adly Templeton. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. Anthropic, 2024.
693 694 695	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. <i>arXiv preprint arXiv:2401.06209</i> , 2024.
696 697 698	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
700 701	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.

- 702 Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter-703 pretability in the wild: a circuit for indirect object identification in gpt-2 small. arXiv preprint 704 arXiv:2211.00593, 2022. 705 John Wright and Yi Ma. High-Dimensional Data Analysis with Low-Dimensional Models: Princi-706 ples, Computation, and Applications. Cambridge University Press, 2022. 707 708 Jeff Wu, Hank Tillman, and Steven Bills. Automated interpretability. https://github.com/ 709 openai/automated-interpretability, 2023. 710 711 Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, 712 and Xia Hu. Harnessing the power of llms in practice: A survey on chatgpt and beyond. arXiv 713 preprint arXiv:2304.13712, 2023. 714 715 Jinrui Yang, Xianhang Li, Druv Pai, Yuyin Zhou, Yi Ma, Yaodong Yu, and Cihang Xie. Scaling 716 white-box transformers for vision. arXiv preprint arXiv:2405.20299, 2024. 717 Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning Di-718
- Yaodong Yu, Kwan Ho Ryan Chan, Chong You, Chaobing Song, and Yi Ma. Learning Diverse and Discriminative Representations via the Principle of Maximal Coding Rate Reduction. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/6ad4174eba19ecb5fed17411a34ff5e6-Abstract.html.
- Yaodong Yu, Sam Buchanan, Druv Pai, Tianzhe Chu, Ziyang Wu, Shengbang Tong, Hao Bai, Yuexiang Zhai, Benjamin D Haeffele, and Yi Ma. White-box transformers via sparse rate reduction: Compression is all there is? *arXiv preprint arXiv:2311.13110*, 2023a.
- Yaodong Yu, Tianzhe Chu, Shengbang Tong, Ziyang Wu, Druv Pai, Sam Buchanan, and Yi Ma. Emergence of segmentation with minimalistic white-box transformers. *arXiv preprint arXiv:2308.16271*, 2023b.
- Zeyu Yun, Yubei Chen, Bruno A Olshausen, and Yann LeCun. Transformer visualization via dictionary learning: contextualized embedding as a linear superposition of transformer factors. *arXiv preprint arXiv:2103.15949*, 2021.
- Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. Investigating
 the catastrophic forgetting in multimodal large language model fine-tuning. In *Conference on Parsimony and Learning*, pages 202–227. PMLR, 2024.
- 739 740

A DETAILS ON THE CRATE ARCHITECTURE

743 A.1 PARAMETER SIZE OF CRATE AND GPT

CRATE is smaller than GPT because of the architecture difference. The vanilla GPT architecture has two main parameterized blocks: Attention block and MLP block.

Parameter size of the MSSA Block. In CRATE, the MSSA block resembles the Attention block, but
 instead of K, Q, V matrices, we only have one matrix. Therefore, compared to standard transformers,
 CRATE uses 1/3 of the parameters for the multi-head attention part.

Parameter size of the ISTA block. The MLP block in vanilla GPT has one parametric matrix that transforms the input representations to the inner space (usually 4x larger), and another parametric matrix that transforms the inner representations back to the output space (as large as the input space).
 In CRATE, the MLP block is replaced by the ISTA-overcomplete block, which transforms the input representation to the overcomplete basis (4x larger) and transforms back with the same parametric matrix. Therefore, compared to standard transformers, CRATE uses 1/2 of the parameters for the MLP part.



810 A.4 OVER-COMPLETE ISTA BLOCK

To give a better idea of how Equation (9) works, we expand the two-iteration process (t = 2). Given $D^{\ell} \in \mathbb{R}^{d \times h}$, we expand the first ISTA step to

$$A_{0} = \mathbf{0},$$

$$A_{1} = S_{\lambda} \left(A_{0} - \eta \cdot (\mathbf{D}^{\ell})^{*} (\mathbf{D}^{\ell} A_{0} - \operatorname{LN}(\mathbf{Z}^{\ell+1/2})) \right)$$

$$= \operatorname{ReLU} \left(\eta \cdot (\mathbf{D}^{\ell})^{*} \operatorname{LN}(\mathbf{Z}^{\ell+1/2}) - \eta \lambda \right).$$
(10)

The second ISTA step continues the process from the initialized sparse code A_1 :

$$A_{2} = S_{\lambda} \left(A_{1} - \eta \cdot (D^{\ell})^{*} (D^{\ell} A_{1} - \operatorname{LN}(Z^{\ell+1/2})) \right)$$

= ReLU $\left(A_{1} - \eta \cdot (D^{\ell})^{*} (D^{\ell} A_{1} - \operatorname{LN}(Z^{\ell+1/2})) - \eta \lambda \right),$ (11)

which can be decomposed to:

832
833
834
835
836

$$G_1 = (D^{\ell})^* D^{\ell} A_1$$

 $G_2 = (D^{\ell})^* \cdot LN(Z^{\ell+1/2})$
 $G = \eta \cdot (G_2 - G_1) - \eta \cdot \lambda$
 $A_2 = ReLU(A_1 + G)$
(12)

where A_2 is the output sparse code. At last, we convert the output sparse code from the coding rate space back to the original representation space:

$$\boldsymbol{Z}^{\ell+1} = \boldsymbol{D}^{\ell} \boldsymbol{A}_2 \tag{13}$$

This process can be implemented by PyTorch-like code shown in Algorithm 3.

Algorithm 3 PyTorch-Like Code for Over-complete ISTA Forward Pass

```
1 class ISTA(nn.Module):
848
            def __init__(self, config):
849
                super().__init__()
850
                self.weight = nn.Parameter(torch.Tensor(4 * config.n_embd, config.n_embd)) # h*d
      4
851
                with torch.no_grad():
      5
852
                    init.kaiming_uniform_(self.weight)
      6
                self.step_size = 0.1
853
                self.lambd = 0.1
      8
854
      9
855
            def forward(self, x, enhanced_feature_id=None):
     10
856
     11
                z_init = F.relu(self.step_size * F.linear(x, self.weight, bias=None) - self.
857
            step_size * self.lambd) # A1
                x1 = F.linear(z_init, self.weight.t(), bias=None)
     12
858
                grad_1 = F.linear(x1, self.weight, bias=None)
     13
859
                grad_2 = F.linear(x, self.weight, bias=None)
     14
860
                grad_update = self.step_size * (grad_2 - grad_1) - self.step_size * self.lambd
     15
861
     16
                output_sparse_code = F.relu(z_init + grad_update) # A2
                output = F.linear(output_sparse_code, self.weight.t(), bias=None)
862
     17
                return output
863
     18
```



Figure 7: CRATE iteratively compresses (MSSA block) and sparsifies (ISTA block) the token representations (*colored points*) across its layers from 1 to *L*, transforming them into parsimonious representations aligned on axes (*colored lines*) with distinct semantic meanings.

A.5 DETAILS OF THE LEARNING PROCESS

873

874

875

876

884

885

887

889

890 891

892 893 894

895 896

897

899

905 906

The desired optimization process is illustrated in Figure 7. The process starts with random token representations (Z^1). Through successive layers, the representations (Z^ℓ) are *compressed* to align with the axis via the MSSA block, forming $Z^{\ell+1/2}$ that are semantically more consistent among relevant tokens. This is then refined by *sparse coding* (the ISTA block) to produce the representations $Z^{\ell+1}$ that align on incoherent axes, leading to semantically more specified token representations. Repeated across layers, this culminates in distinct token representations Z^{L+1} aligned on unique semantic axes. More detailed explanation of this optimization process can be found in Appendix A.5.

We elaborate the learning process of CRATE in this section, with a close reference to Figure 7.

In Figure 7, what is the space the points are drawn in? The space is the representation space (of layer ℓ). Because the model is pretrained with next-token prediction in the language domain, the space is specifically a semantic space. Thus, each point (token) has a semantic representation in this high-dimensional space.



Figure 8: Illustration of the concepts of the activation matrix and poly-semanticity.

How do the layouts of the points suggest mono or poly-semanticity? First, each axis (red/yellow) in the figure represents a neuron/feature visualized in the semantic space. We visualized the activation matrix A_t in Figure 8. For example, when $L \in \{1, 2, 3\}$, the model dimension is 128, which means that the overcomplete basis of ISTA will have a dimension of 512, introducing 512 features. Now if we input a sequence of 256 tokens, the activation matrix will have a shape of [512, 256].

Poly-semanticity means that token representations in the semantic space are clustered as a broader set of semantic meanings - that is, each neuron has a broader set of semantic meanings. For example, in the gray box on the top left of Figure 7, both yellow- and red-backgrounded tokens represent either a number or a capitalized token. This corresponds to multiple high activations in the feature, where the tokens that activated this feature can either be a number or a capitalized token, which is shown in Figure 8 (where pink squares represent high activations).

In the *compression* phase, the token representations are pushed towards the semantic axes, so that
 the tokens will activate on fewer features but will gain higher activations on these features - which
 is essentially an activation condensing process.

In the *sparsification* phase, the neurons (axes) are made further from each other, meaning that the features have less semantic overlap with each other. In this case, the results will become the gray box on the top right side of Figure 7, indicating that each neuron has distinct semantic meanings, like "numbers" or "capitalized tokens".

- 926 Note that this is a minimal example. In practice, tokens appear in context.
- 927 928

929

B DETAILS ON INTERPRETABILITY EVALUATIONS

⁹³⁰ This section details the implementation details of the interpretability evaluations.

In practice, we adopt three evaluation metrics: two from OpenAI (Bills et al., 2023) and one from Anthropic (Bricken et al., 2023). As the Anthropic metric is a closed-source follow-up of OpenAI, we start from the official implementation provided by OpenAI (Wu et al., 2023) for both metrics.

For each layer, we use randomly sampled 8,000 text excerpts of 1,024 tokens each, which sums up to 8M tokens in total, from the test split of the uncopyrighted Pile dataset, to evaluate the interpretability scores.

- 938 939
- **B.1** PARAMETERS OF EVALUATION METRICS

Comprehensive parameter settings are shown in Table 5. For the OpenAI metrics, each input text excerpt contains 64 tokens. For the *OpenAI top-and-random* metric, we use 5 top activated excerpts for explanation, and a mixture of 5 top activated and 5 randomly activated excerpts for simulation. For the *OpenAI random-only* metric, we only use 5 randomly activated excerpts for simulation.

For the *Anthropic* metric, each text excerpt contains only 8 tokens. For the explanation model, we input 15 top activated excerpts, 5 randomly activated excerpts, and 22 excerpts from different activation quantiles. To elaborate, we evenly divide the activation range into 11 quantiles, where we pick 2 excerpts from each of them. For the simulation model, we input 10 top activated excerpts, 5 randomly activated excerpts, 22 quantiled excerpts, and 10 top activated out-of-context (OOC) excerpts. Our implementation of the OOC excerpts is to cut the input text excerpt into length of only 3 tokens.

Table 5: Evaluation parameter settings of the OpenAI and Anthropic approach.

			Explanation				Simu	lation	
		#Token	#Top	#Rand	#Qua	#Top	#Rand	#Qua	#OOC
OpenAI	TaR	64	5			5	5		
OpenAl	Rand	64	5				5		
Anthro	opic	8	15	5	$2 \cdot 11$	10	5	$2 \cdot 11$	10

961

952

B.2 DISCUSSION ON FOCUS OF DIFFERENT MEASURES

962 The OpenAI random-only metric is the easiest to interpret. As noted by Bills et al. (2023), the 963 random-only metric considers an explanation's ability to capture the neuron's representation of features in the pre-training distribution, because the simulated tokens are uniformly randomly sampled 964 from the validation set of the pre-train dataset. However, the random-only scoring with small sam-965 ple size risks failing to capture behavior, due to lacking both tokens with high simulated activations 966 and tokens with high real activations. Top-and-random scoring addresses the latter, but causes us to 967 penalize falsely low simulations more than falsely high simulations, and thus tends to accept overly 968 broad explanations. 969

The Anthropic metric, on the other hand, puts more focus on the mono-semanticity of the activations,
as noted by Bricken et al. (2023). For sparse features, which don't fire on most random samples,
evaluating across a wide range of activations effectively tests the model's ability to distinguish a

feature's large activations from zero, and the short text excerpts make it easier for the simulation model to identify the sparse activations.

B.3 MORE ACCESSIBLE EVALUATION

To reduce compute cost, we use Mistral-7B-instruct as the explanation model, and LLaMA-2-7B as the simulation model. We empirically prove that these replacements does not affect the conclusions of apple-to-apple comparison between CRATE and GPT below.

Explanation model. In the official implementation (Wu et al., 2023), the explanation model is gpt-4. According to ablations described in Bills et al. (2023), it also makes sense to use the sligtly cheaper model gpt-3.5-instruct. Due to the high compute cost, we use the open-source model mistral-7b-instruct instead. We demonstrate the performance of gpt-3.5-turbo and mistral-7b-instruct using the OpenAI random-only and top-and-random metrics in Table 6. Results show that the change of model doesn't significantly change the scores, and doesn't affect conclusions at all.

Table 6: Interpretability measure of GPT, GPT-SAE and CRATE-GPT on the Pile dataset based on the OpenAI metrics. Explanation model: Mistral-7B-instruct/GPT-3.5-turbo. Simulation model: LLaMA-2-7B.

995									
006	mistral-	-7b-instr	ruct	ρ (Rar	ndom-only)) (%,↑)	ρ (Top-a	nd-Rando	m) (%, ↑)
990	Model	Size	Loss	Layer 1	Layer 2	Layer 3	Layer 1	Layer 2	Layer 3
997	CRATE-1L	6.54M	4.06	4.8	-	-	3.9	-	-
998	CRATE-2L	6.64M	3.55	8.0	5.8	-	7.8	8.3	-
999	CRATE-3L	6.74M	3.46	9.0	9.4	6.9	9.6	9.3	8.4
1000	GPT-1L	6.64M	3.83	8.9	-	-	8.8	-	-
1001	GPT-2L	6.83M	3.23	2.3	1.6	-	4.3	4.1	-
1002	GPT-3L	7.03M	3.11	3.1	-0.3	1.3	7.3	0.8	2.6
1003	GPT-11	L (16x SA	E)	2.9	-	-	5.4	-	-
1004	GPT-2I	(16x SA	E)	3.5	1.8	-	7.4	4.2	-
1005	GPT-31	L (16x SA	E)	3.2	2.3	1.1	9.6	5.0	4.5
1006	GPT-3	.5-turbo)	ρ (Ran	dom-only)	$(\%,\uparrow)$	ρ (Top-a	nd-Randoi	n) (%, ↑)
1006 1007	GPT-3 Model	5-turbo Size	Loss	ρ (Ran Layer 1	dom-only) Layer 2	(%, ↑) Layer 3	ρ (Top-a) Layer 1	nd-Randor Layer 2	n) (%, †) Layer 3
1006 1007 1008	GPT-3 Model CRATE-1L	5-turbo Size 6.54M	Loss 4.06	ρ (Ran Layer 1 4.8	dom-only) Layer 2	(%, ↑) Layer 3	ρ (Top-a Layer 1 3.9	nd-Randor Layer 2	n) (%, ↑) Layer 3
1006 1007 1008 1009	GPT-3 Model CRATE-1L CRATE-2L	8.5-turbo Size 6.54M 6.64M	Loss 4.06 3.55	ρ (Ran Layer 1 4.8 8.2	dom-only) Layer 2 - 6.0	(%, ↑) Layer 3	ρ (Top-a Layer 1 3.9 7.5	nd-Randon Layer 2 - 8.0	n) (%, ↑) Layer 3
1006 1007 1008 1009 1010	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L	5.5-turbc Size 6.54M 6.64M 6.74M	Loss 4.06 3.55 3.46	ρ (Ran Layer 1 4.8 8.2 9.1	dom-only) Layer 2 - 6.0 9.2	(%, ↑) Layer 3 - 6.9	ρ (Top-a Layer 1 3.9 7.5 9.5	nd-Randor Layer 2 - 8.0 9.1	n) (%, ↑) Layer 3 - - 8.3
1006 1007 1008 1009 1010 1011	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L	8.5-turbo Size 6.54M 6.64M 6.74M 6.64M	Loss 4.06 3.55 3.46 3.83	ρ (Ran Layer 1 4.8 8.2 9.1 9.0	dom-only) Layer 2 - 6.0 9.2 -	(%, ↑) Layer 3 - - 6.9	ρ (Top-a Layer 1 3.9 7.5 9.5 9.0	nd-Randor Layer 2 - 8.0 9.1	n) (%, ↑) Layer 3 - 8.3 -
1006 1007 1008 1009 1010 1011 1012	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L GPT-2L	6.5-turbo Size 6.54M 6.64M 6.74M 6.64M 6.83M	Loss 4.06 3.55 3.46 3.83 3.23	ρ (Ran Layer 1 4.8 8.2 9.1 9.0 2.2	dom-only) Layer 2 - 6.0 9.2 - 1.6	(%, ↑) Layer 3 - - 6.9 -	 ρ (Top-at Layer 1 3.9 7.5 9.5 9.0 4.3 	nd-Randor Layer 2 - 8.0 9.1 - 4.4	n) (%, ↑) Layer 3 - - 8.3 - -
1006 1007 1008 1009 1010 1011 1012	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L GPT-2L GPT-3L	8.5-turbo Size 6.54M 6.64M 6.74M 6.64M 6.83M 7.03M	Loss 4.06 3.55 3.46 3.83 3.23 3.11	ρ (Ran Layer 1 4.8 8.2 9.1 9.0 2.2 3.0	dom-only) Layer 2 - 6.0 9.2 - 1.6 -0.3	(%, ↑) Layer 3 - 6.9 - 1.2	$\begin{array}{c} \rho \ ({\rm Top-ar}) \\ {\rm Layer} \ 1 \\ 3.9 \\ 7.5 \\ 9.5 \\ 9.0 \\ 4.3 \\ 7.0 \end{array}$	nd-Randor Layer 2 - 8.0 9.1 - 4.4 3.1	n) (%, ↑) Layer 3 - - 8.3 - - 3.0
1006 1007 1008 1009 1010 1011 1012 1013	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L GPT-2L GPT-3L GPT-1L	5.5-turbo Size 6.54M 6.64M 6.74M 6.64M 6.83M 7.03M 2.(16x SAI	Loss 4.06 3.55 3.46 3.83 3.23 3.11 E)	$\begin{array}{c} \rho \ ({\bf Ran} \\ {\bf Layer 1} \\ 4.8 \\ 8.2 \\ 9.1 \\ 9.0 \\ 2.2 \\ 3.0 \\ 2.6 \end{array}$	dom-only) Layer 2 - 6.0 9.2 - 1.6 -0.3 -	(%, ↑) Layer 3 - - 6.9 - 1.2	$ \begin{array}{c} \rho \ ({\rm Top-ar}\\ {\rm Layer \ 1} \\ 3.9 \\ 7.5 \\ 9.5 \\ 9.0 \\ 4.3 \\ 7.0 \\ 4.7 \end{array} $	nd-Randor Layer 2 - 8.0 9.1 - 4.4 3.1	n) (%, ↑) Layer 3 - - 8.3 - 3.0 -
1006 1007 1008 1009 1010 1011 1012 1013 1014	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L GPT-2L GPT-3L GPT-1L GPT-2L	8.5-turbo Size 6.54M 6.64M 6.74M 6.64M 6.83M 7.03M 2.(16x SAI 2.(16x SAI	Loss 4.06 3.55 3.46 3.83 3.23 3.11 E) E)	$\begin{array}{c} \rho \ (\textbf{Ram}\\ Layer \ 1 \\ 4.8 \\ 8.2 \\ 9.1 \\ 9.0 \\ 2.2 \\ 3.0 \\ 2.6 \\ 3.4 \end{array}$	dom-only) Layer 2 - 6.0 9.2 - 1.6 -0.3 - 1.6	(%, ↑) Layer 3 - - 6.9 - 1.2 -	$ \begin{array}{c} \rho \ ({\rm Top-ar}) \\ {\rm Layer \ 1} \\ 3.9 \\ 7.5 \\ 9.5 \\ 9.0 \\ 4.3 \\ 7.0 \\ 4.7 \\ 5.0 \end{array} $	nd-Randor Layer 2 - 8.0 9.1 - 4.4 3.1 - 2.9	n) (%, ↑) Layer 3 - - 8.3 - 3.0 - -
1006 1007 1008 1009 1010 1011 1012 1013 1014 1015	GPT-3 Model CRATE-1L CRATE-2L CRATE-3L GPT-1L GPT-2L GPT-3L GPT-1L GPT-2L GPT-3L	8.5-turbo Size 6.54M 6.64M 6.74M 6.64M 6.83M 7.03M 2.03M 2.03X 2.0	Loss 4.06 3.55 3.46 3.83 3.23 3.11 E) E) E) E)	$\begin{array}{c} \rho \ (\textbf{Ram}\\ Layer \ 1 \\ 4.8 \\ 8.2 \\ 9.1 \\ 9.0 \\ 2.2 \\ 3.0 \\ 2.6 \\ 3.4 \\ 2.8 \end{array}$	dom-only) Layer 2 - 6.0 9.2 - 1.6 -0.3 - 1.6 1.8	(%, ↑) Layer 3 - - - - - - - - - - - - - - - - - - -	$ \begin{array}{c} \rho \ ({\rm Top-ar}) \\ {\rm Layer \ 1} \\ 3.9 \\ 7.5 \\ 9.5 \\ 9.0 \\ 4.3 \\ 7.0 \\ 4.7 \\ 5.0 \\ 7.4 \end{array} $	nd-Randor Layer 2 - 8.0 9.1 - 4.4 3.1 - 2.9 3.8	n) (%, ↑) Layer 3 - - 8.3 - 3.0 - 3.2

Simulation model. The official implementation of the simulation model utilizes text-davinci-003 (now named gpt-3.5-turbo-instruct), which no longer supports retrieving the logprobs through the API, so we use LLaMA-2-70B as an equally capable replacement (Touvron et al., 2023). For more accessible evaluations, we use LLaMA-2-7B instead. We show the difference in interpretability caused by different simulation model size on LLaMA-2-7B and LLaMA-2-70B in Table 7. Empirical results show that although LLaMA-2-7B has overall lower scores and higher variance than LLaMA-2-70B, it doesn't affect essential conclusions about the apple-to-apple comparison between CRATE and GPT.

		Interpre	tability (7	B) (%, ↑)	Interpret	tability (70	B) (%, ↑)
		Layer 1	Layer 2	Layer 3	Layer 1	Layer 2	Layer 3
C	rate-1L	3.9	-	-	6.4	-	-
С	RATE-2L	7.5	8.0	-	7.4	7.1	-
С	rate-3L	9.5	9.1	8.3	10.4	7.4	6.5
(GPT-1L	9.0	-	-	13.4	-	-
(GPT-2L	4.3	4.4	-	6.4	7.8	-
	GPT-3L	7.0	3.1	3.0	10.1	3.2	6.3
		Layer	1 Layer	2 Layer	3 Layer	4 Layer	5 Laye
7D	CRATE-6L	. 10.5	5 7.4	5.8	8.0	8.1	5.7
/ D	GPT-6L	13.7	5.2	1.9	0.6	5.6	6.9
70R	CRATE-6L	. 10.1	6.5	7.0	8.3	9.4	0.7
700	GPT-6L	14.5	6.2	2.5	0.7	3.9	4.1

1026Table 7: Interpretability measure of GPT and CRATE-GPT on the Pile dataset based on the OpenAI Top-and-1027random metric. Explanation model: GPT-3.5-turbo. Simulation model: LLaMA-2-7B/LLaMA-2-70B.

C ANALYSIS ON ACTIVATION SPARSITY

We demonstrate the activation sparsity of CRATE compared to GPT in Figure 9. We observe that the activations of CRATE are higher than GPT. One might have the confusion about why CRATE is designed to be sparse but the activations evaluated is denser than GPT. Note that the sparsity evaluated in standard transformer model is output from the hidden layer of the MLP layer, which is the activation matrix A before applying to the residual stream, as shown in Figure 10. The actual representations in standard transformers, which are after applying the residual stream, are not sparse at all. In contrast, the sparsity evaluated in CRATE is the actual representations A_t (including the residual stream).



Figure 9: Layer-wise activation sparsity of CRATE and GPT. Left: 6L models. Right: 12L models.



Figure 10: Extracting sparse code A_t from CRATE and hidden layer output A from GPT.

1080 We also present the activation dynamics of CRATE and GPT with the progression of the pre-training 1081 process in Figure 11. We observe a strong trend that the sparsity of CRATE monotonically decreases 1082 in the early stage (trained on 1.6B tokens), which aligns with the design purpose. In the late stage 1083 (16B, 160B tokens), the sparsities in the early sites (L < 12) significantly decreases, which also 1084 aligns with the design purpose. On the other hand, GPT never appears to have a decreasing trend of activation sparsity over layers across the whole pre-training stage, indicating a systematic difference between the sparsity dynamics between CRATE and GPT. One counter-intuitive observation is that 1086 the decreasing trend fades as the stage moves on. Our hypothesis is that CRATE overfits on the next 1087 token prediction task due to the large amount of tokens trained. 1088



Figure 11: Layer-wise activation sparsity w.r.t. tokens trained. Left: CRATE language model. Right: GPT.

DETAILS ON INTERPRETATION SCORE DISTRIBUTIONS D

We visualize the distributions of layer-wise interpretation scores of CRATE and GPT with L = 121107 in Figure 12. We exclude cases where activations sampled from GPT-Base (random-only metric) 1108 are all zeroes, as in these cases the correlation ρ will be undefined. This results in a smaller number 1109 counted in the GPT activations in the first two rows. 1110

1111 1112

1113

1115

1091

1094

1101

1102 1103 1104

1105 1106

DETAILS ON SPARSE AUTO-ENCODER Ε

1114 E.1 SPARSE AUTOENCODER AND DICTIONARY LEARNING

The dictionary learning model is an MLP with a single hidden layer. It is trained as an auto-encoder 1116 using input weights as the encoder that maps the input activations to a higher dimension, and output 1117 weights as the decoder. Formally, given activation $a \in \mathbb{R}^h$ sampled from $A \in \mathbb{R}^{h \times N}$, the encoder 1118 W_1, b_1 with dimension multiplicator μ maps the activations to a hidden representation $h \in \mathbb{R}^{\mu h}$, 1119 whereas the decoder W_2, b_2 maps the representation back to the original dimension $\hat{a} \in \mathbb{R}^h$. The 1120 dictionary learning objective can thus be expressed as 1121

1122 1123

$$\bar{a} = a - b_2 \tag{14}$$

$$\boldsymbol{h} = \operatorname{ReLU}(\boldsymbol{W}_1 \bar{\boldsymbol{x}} + \boldsymbol{b}_1) \tag{15}$$

$$\hat{\boldsymbol{a}} = \boldsymbol{W}_2 \boldsymbol{h} + \boldsymbol{b}_2 \tag{16}$$

$$\mathcal{L} = \frac{1}{|A|} \sum_{a \in A} \|a - \hat{a}\|_{2}^{2} + \lambda \|h\|_{1}$$
(17)

1128 1129

1131

1130 E.2 DETAILED SETUP

We train the sparse auto-encoders on the train split of the uncopyrighted Pile dataset until conver-1132 gence. Following Bricken et al. (2023) and Conmy (2023), we adopt the resampling strategy to 1133 re-train the dead features, and the learning rate scheduling strategy to improve recovery rate. For



Figure 12: Distribution of the interpretation scores over CRATE-12L and GPT-12L. *x*-axis: interpretation score. *y*-axis: count of neurons falling in the corresponding interval of interpretation score.

implementation, we mainly follow Conmy (2024), with $\lambda_{\ell_1} = 1.6 \times 10^{-4}$, $\alpha = 1.2 \times 10^{-3}$ for all sizes of models. We evaluate using the average loss of randomly sampled batches on the validation split of the uncopyrighted Pile dataset.

1180 1181

1173

1174 1175 1176

1182

1183 E.3 LOSS CURVE

1184

The loss curves of training the sparse auto-encoders are shown in Figure 13. Generally, resampling boosts the performance of the reconvery score, which aligns with the conclusions shown in Bricken et al. (2023) and Conmy (2023). We also observe an increasing trend of performance with the increases of the SAE multiplication factor μ and model size L.



Figure 13: Left: The recovery scores of GPT-1L ($\ell = 0$) with SAE multiplication factors $\mu = \in \{1, 4, 16\}$. Right: The reconstruction loss of SAE with $\mu = 16$ on different sizes of GPT models $L \in \{1, 2, 3\}$, averaged across all layers.

1205 E.4 PERFORMANCE

1201

1202

1203 1204

1210

1211

1207 The performance of sparse auto-encoders of CRATE-GPT and GPT under a variety of settings (model 1208 L, ℓ and sparse autoencoder width multiplication factor μ) are shown in Table 8. The percentages of 1209 dead neurons for all layers of $L \in \{1, 2, 3\}$ are less than 1%.

Table 8: Reconstruction loss and recovery score of the sparse autoencoders on CRATE and GPT.

212									
010	$ \mu$	i = 16		Recon	struction I	Loss (\downarrow)	Recov	ery Score	$(\%,\uparrow)$
213		Size	Loss	Layer 0	Layer 1	Layer 2	Layer 0	Layer 1	Layer 2
214	GPT-1L	6.64M	3.83	4.35	-	-	95.0	-	-
215	GPT-2L	6.83M	3.23	3.50	3.45	-	95.2	92.2	-
216	GPT-3L	7.03M	3.11	3.38	3.39	3.29	94.6	94.8	92.4
217	CRATE-1L	6.54M	4.06	4.33	-	-	93.6	-	-
218	CRATE-2L	6.64M	3.55	4.12	3.80	-	95.7	95.7	-
219	CRATE-3L	6.74M	3.46	4.05	3.99	3.77	93.0	93.9	95.0
220		u = 4		Recon	struction I	Loss (\downarrow)	Recov	ery Score	(%, †)
221		Size	Loss	Layer 0	Layer 1	Layer 2	Layer 0	Layer 1	Layer 2
222	GPT-1L	6.64M	3.83	4.34	-	-	93.7	-	-
223	GPT-2L	6.83M	3.23	3.59	3.56	-	92.7	88.6	-
224	GPT-3L	7.03M	3.11	3.45	3.50	3.34	92.2	94.9	89.9
225	CRATE-1L	6.54M	4.06	4.39	-	-	92.1	-	-
226	CRATE-2L	6.64M	3.55	4.37	3.93	-	93.7	93.8	-
227	CRATE-3L	6.74M	3.46	4.03	4.11	3.81	92.6	92.6	92.6
228		$\iota = 1$		Recons	struction L	$loss(\downarrow)$	Recov	ery Score	$\overline{(\%,\uparrow)}$
1229	,	Size	Loss	Layer 0	Layer 1	Layer 2	Layer 0	Layer 1	Layer 2
230	GPT-1L	6.64M	3.83	4.93	-	-	95.0	-	-
1231	GPT-2L	6.83M	3.23	3.89	3.75	-	89.0	82.2	-
1232	GPT-3L	7.03M	3.11	3.63	3.61	3.58	86.9	91.0	84.9
1233	CRATE-1L	6.54M	4.06	4.69	-	-	86.2	-	-
1234	CRATE-2L	6.64M	3.55	4.68	4.29	-	90.4	88.9	-
1225	CRATE-3L	6.74M	3.46	4.39	4.38	4.16	97.0	89.2	88.1
LUJ									

¹²³⁶

1237 E.5 INTERPRETABILITY 1238

1239 Does CRATE have more optimal representation than GPT in terms of interpretability? As it's
1240 hard to decide how much interpretability gain it is from CRATE to CRATE-SAE directly (as explained in Section 5.3), we compare the interpretability *improvement* of CRATE-SAE over CRATE to the interpretability improvement of GPT-SAE over GPT.

The interpretability of GPT-SAE is already included in Figure 4. The interpretability of CRATE-SAE under the OpenAI TaR and Anthropic metrics are shown in Table 9.

Table 9: Interpretability of CRATE-SAE under the OpenAI TaR and Anthropic metics.

	Op	enAI TaR	(†)	Anthropic (†)				
	Layer 1	Layer 2	Layer 3	Layer 1	Layer 2	Layer 3		
1L	6.0	-	-	17.9	-	-		
2L	7.7	5.2	-	21.7	12.4	-		
3L	7.2	6.4	4.6	19.3	18.4	11.6		

1254 E.6 STEERING THE LM OR SAE

1255 In comparison to post-hoc trained SAEs, built-in sparsification processes, such as the one we pro-1256 posed in this paper, have the potential to be steered with perfect fidelity. As visualized in Figure 14, 1257 post-hoc approaches like SAE require steering the model with the decomposed hidden states h, 1258 whose encoding and decoding processes are both *lossy*. An imperfect reconstruction systematically 1259 leads to *distortions* of the steering signal upon the hidden states, and thus affects downstream appli-1260 cations of the GPT-SAE model. In contrast, CRATE doesn't include any approximation that distorts 1261 the steering signal, so the signal can be propagated without loss of fidelity. This conclusion does not 1262 change whether the performance of GPT-SAE outperforms CRATE or not.



Figure 14: Illustration of steering a language model directly or using SAE.

F FURTHER QUALITATIVE RESULTS ON INTERPRETABLE NEURONS

This section lists some further qualitative examples of tokens and their activations when L = 3, including examples of the GPT-SAE activations. Specifically, we demonstrate two neurons in each model. Tokens with a deeper blue background have a higher activation. Explanations k_i and scores s_i are obtained by Algorithm 1.

1286 1287

1278

1279 1280

1281 1282

F.1 CRATE-3L, LAYER 0, NEURON 288

1289 **OpenAI Evaluation** Score: 0.44478400135318646

Explanation: information related to the regulation of mRNA expression and its role in carbohydrate
 metabolism, with a focus on CRC cells and gene signaling in the context of cancer development.
 Top Activations

```
1293 animal models of cartilage degradation ([@b41-mmr-16-04-3841]). Among
1294 these cytokines, IL-1\beta is highly overexpressed in the cartilage and in
1295 the synovial tissue, while the expression of IL-1 receptor antagonist ([
@b42-mmr-
```

296	raft's Moreover ACOX1 overe yn ression atten wated the aug mentation of
297	migration and invasion of CRC cells by $miR = 15h = 5n$ overe vn ression. In
298	conclusion our study demonstrated a functional role of the SIRT1/miR-
299	15 b - 5 p / AC 0X 1 axis
300	Random Activations
301	ed with Peter Braun, the Moravian \n mission ary in Antigua; and to that
302	correspondence he owed in part his \n interest in missionary work. But that
202	was not the end of the Brethren's \n influence. At all meetings addressed
204	by the founders of the proposed \setminus n Society, the speaker repeatedly
205	ESA) and the American Association for the Advance ment of Science (AA AS
206) have well - developed and successful science policy fellows hips . These
207	programs acknowledge that scientists can play important roles in directing
208	new laws and policies in their field, and that their expertise is needed
200	for effective decision - making \ [[@B81 - in
10	
U U	
10	Anthropic Evaluation Score: 0.2813215497394413
12	Explanation: Phrases related to molecular biology and gene expression, specifically in the context
3	of mRNA transcription and its activation or inhibition. Additionally, there are mentions of certain
4	proteins (IPOTENT, cytokine, IRT), cellular processes (proliferation,
2	Top Activations
6	b - 5 p <mark>overe</mark> xpression.
7	cases showed EG FR overe xp ression .
3	, TLR2 overe xp ression in
	IL - 1 β mRNA expression in the
	Random Activations
	G. albidus * T MW
	had suffered from heat contact ur <mark>tic</mark> aria
	"alive": true,\n
	= _mm_packus_ep
	F - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 -
	E 2 GPT 31 LAVED 2 NELIDON 280
	1.2 OI 1-5L, LATER 2, NEURON 209
	Onen Al Eveluction Scenes 0.2226227046012070
	Explanation: The provided text contains multiple sections, but the activations given for Neuron 4.
	Explanation: The provided text contains indupple sections, but the activations given for Neuron 4 seems to be related to genetic and statistical data (e.g., population, CL percent, risk association, and
	recessive models). Given this, the main thing this neuron does is identify
	Ton Activations
	in Asian population Similarly in Caucasian population the rs 499 776
	polymorph ism attributes risk association in homozygote OR 0, 70 (95% CI
	[0, 50 - 0, 98]), dominant OR 3, 57 (95% CI [2, 34 - 5, 27]), and recess
	ive models OR 0
	* SE * = 0, 04 40, *t * = -1, 07 75, *p * \> 0 05 95% CT (-0 13 38
	0.0390) for the Slovakian villagers story 1, therefore indicating full
	mediation by exonerations and out-group focused emotions
	Random Activations
	longer the vortex, it's the smooth current of rotating air which is next
	to n the vortex, and we use the updraft of this air." Taking advantage of
	the free 📊 lift in this updraft of air is called 🔤 wake-energy retrieval
	." on long– <mark>\n</mark> haul flights, fuel savings of between
	ushing . $n \in S$ igh called her supervisor . Sergeant Sweeney and
	Deputy Ray responded $n n$ and moved Don ery so that S igh could search his
	cell. Don ery had been in his new $n \ln cell$ for less than five minutes
1	when the toilet overflowed and water began flowing out $\n \$

1350 Anthropic Evaluation Score: 0.06532542667915378 1351 Explanation: strings containing specific numbers and alphanumeric characters, such as "CI-50-.", 1352 "e-44-", "87-", and "f-". Additionally, it activates slightly for certain words like "cost", "weeks", 1353 "disability", and 1354 **Top Activations** 95% CI [0.50-1355 1356 95% CI, 0.13-1357 f\"], [0.2222 1358 95 % CI [0.50-1359 **Random Activations** { 8 }{ 45 }\ pi \\ \n 1360 1361 instead of that silly website. How 1362 free software; you can redistribute is 1363 1364 1365 F.3 GPT-3L-SAE-16x, LAYER 2, NEURON 57 1366 **OpenAI Evaluation** Score: 0.1427145260798203 1367 Explanation: months or the word "Bank" followed by a year. 1368 **Top Activations** 1369 No. 18-20609 February 21, 2020 \n 1370 1371 with the compact - open top ology, is a locally compact group.' n author 1372 <code>\n - 'Nicolas Radu[^1]'\n date: 'July 15, 2016'\n title: |\n</code> 1373 top ological characterization of the Moufang $\setminus n$ property for compact 1374 polyg ons 1375 **Random Activations** 1376 of DN \star db / db \star mice. $\langle n(\star \star \star \star)$ Ur inary album in to creat in ine ratio 1377 (** B **) Ser um u rea nitrogen . (** C **) Left kidney weight to body 1378 weight ratio. (**D**) HE staining. Bar = 1379 this email: ot isd ark o 60 @ yahoo . com $\ n \ E \ FIX$ THE FO LLOW ING PR OB LE MS 1380 TO ALL $\ln \ln CROSS$ THE GLOBE ON : $\ln \ln 1$. Getting your lover or husband back $\ln 2$. Spiritual bullet proof $\ln 3$. Training $\ln 4$. Money 1381 1382 Anthropic Evaluation Score: 0.2540425061001668 1384 **Explanation:** dates and specifically, the month and day for a given year. The neuron is not activated 1385 by the year alone, and it requires both the month and day for a complete activation. 1386 **Top Activations** February 21, 2020 1387 'July 15, 1388 \n date : 2016 1389 field, Missouri (December 15, 2014 1390 February 5, 1998 1391 **Random Activations** 1392 Can ola oil \n 1393 \n n . c ke1394 uana when a draw would have clinched 1395 . n t1396 1399

- 1400
- 1401
- 1402 1403