
Successor Heads: Recurring, Interpretable Attention Heads In The Wild

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 In this work we present successor heads: attention heads that increment tokens
2 with a natural ordering, such as numbers, months, and days. For example, succes-
3 sor heads increment ‘Monday’ into ‘Tuesday’. We explain the successor head be-
4 havior with an approach rooted in mechanistic interpretability, the field that aims
5 to explain how models complete tasks in human-understandable terms. Existing
6 research in this area has found interpretable language model components in small
7 toy models. However, results in toy models have not yet led to insights that explain
8 the internals of frontier models and little is currently understood about the internal
9 operations of large language models. In this paper, we analyze the behavior of
10 successor heads in large language models (LLMs) and find that they implement
11 abstract representations that are common to different architectures. They form
12 in LLMs with as few as 31 million parameters, and at least as many as 12 bil-
13 lion parameters, such as GPT-2, Pythia, and Llama-2. We find a set of ‘mod 10’
14 features¹ that underlie how successor heads increment in LLMs across different
15 architectures and sizes. We perform vector arithmetic with these features to edit
16 head behavior and provide insights into numeric representations within LLMs.
17 Additionally, we study the behavior of successor heads on natural language data,
18 identifying interpretable polysemanticity in a Pythia successor head.

19 1 Introduction

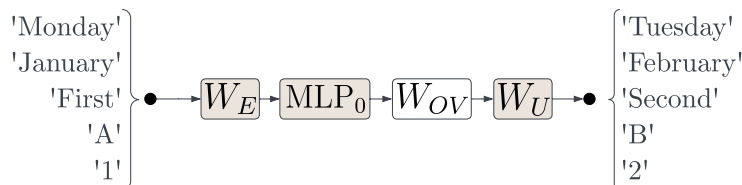


Figure 1: A *successor head* with OV matrix W_{OV} maps numbered tokens in embedding space (e.g. ‘Monday’) to their successor values in unembedding space (e.g. ‘Tuesday’). The circuit consists of the embedding matrix, the first MLP block, one attention head, and the unembedding matrix.

20 Mechanistic interpretability [2] is the process of reverse-engineering the algorithms that trained neu-
21 ral networks have learned. Recently, much attention has been paid to interpreting transformer-based
22 large language models (LLMs), as while these models have demonstrated impressive capabilities
23 [3], there is little understanding of how these models produce their outputs. Existing interpretability

¹Here, following Elhage et al. [1], a ‘feature’ refers to an interpretable (linear) direction in activation space.

24 research includes comprehensive reverse-engineering efforts into toy models [4] and small language
25 models [5, 6], though few insights have been gained about how frontier LLMs function.

26 In mechanistic interpretability, universality [7, 8] is a hypothesis that there are *common representa-*
27 *tions* in neural networks – specifically, that neural networks with different architectures and scales
28 form common internal representations. In this work, we use **abstraction** to refer to how common
29 representations are used for different *tasks* models perform. Strong evidence for (or against) the
30 universality hypothesis and task abstraction could significantly affect research priorities in inter-
31 pretability. If common representations form across different language models and tasks, then re-
32 search on small or toy language models [1, 9] and narrow tasks [6, 10, 11] may be the best way to
33 gain insights into LLM capabilities. Conversely, if the representations used by language models do
34 not generalize to different model scales and/or tasks, then developing methods that can be applied to
35 larger language models and don’t rely on lessons from small models generalizing (such as Wu et al.
36 [12], Bills et al. [13], Conmy et al. [14]) may be the most important direction for interpretability.

37 In this work, we find an interpretable set of attention heads we call **successor heads** in models of
38 many different scales and architectures. Successor heads are attention heads that perform incre-
39 mentation in language models. The input to a successor head is the representation of a token in an
40 ordinal sequence such as ‘Monday’, ‘first’, ‘January’, or ‘one’. The output of a successor head as-
41 signs a higher likelihood to the incremented token, such as ‘Tuesday’, ‘second’, ‘February’, or ‘two’.
42 We find successor heads in the smaller and larger Pythia language models [15] with between 30
43 million and 12 billion parameters. We can understand the role of successor heads through a simple,
44 end-to-end path through language models (Figure 1) and we identify transferrable features that these
45 attention heads use across different tasks (Section 3). In our work, we find evidence for a weak form
46 of universality ([16]; points 1. and 2.) as well as abstraction (point 3.) in language models, as

- 47 1. Successor heads use a *common representation*, such as a linear mod-10 features.
- 48 2. Successor heads form in LLMs with as few as 31M and as many as 12B parameters.
- 49 3. Successor heads perform incrementation across several different tasks.

50 Related work is discussed in Appendix P. Our contributions can be summarised as follows:

51 1. Introducing and interpreting successor heads (Section 2-3)

- 52 (a) We introduce and explain successor heads, which to the best of our knowledge are the most
53 closely studied components in LLMs that occur in both small and large models.
- 54 (b) Using findings from 3., we edit successor heads’ numeric inputs with vector arithmetic.

55 2. Showing that the succession mechanism is important in the wild (Appendix J)

- 56 (c) We show that successor heads play an important role in incrementation-based tasks in natural
57 language datasets – for instance, predicting the next number in a numbered list of items.

58 3. Finding abstract numeric representations in language models (Section 3)

- 59 (d) We isolate a common *numeric subspace* within embedding space, that for any given token
60 (e.g. ‘February’) encodes the index of that token within its ordinal sequence (e.g. months).
- 61 (e) Moreover, we find that this numeric subspace has *interpretable features*, as an unsupervised de-
62 composition of token representations yields a crucial set of features we call the **mod-10 features**
63 $\{f_0, \dots, f_9\}$. f_n is present in all tokens whose numerical index $\equiv n \pmod{10}$, e.g. f_2 is present
64 in the model’s representations of ‘2’, ‘32’, ‘172’, ‘February’, ‘second’ and ‘twelve’.

65 2 Successor Heads

66 LLMs are able to increment elements in an ordinal sequence. For instance, Pythia-1.4B will com-
67 plete the prompt “If this is 1, the next is” with “ 2”, and the prompt “If this is January, the next is”
68 with “ February”. Given this observation, we find evidence for attention heads within LLMs (which
69 we refer to as **successor heads**) responsible for performing this type of incrementation. To get evi-
70 dence for successor heads we require three definitions: i) the **succession dataset** of tasks involving
71 abstract numeric representations, ii) an **effective OV circuit** to measure how attention heads affect
72 model outputs, and finally iii) **successor score**.

73 The **succession dataset** consists of tokens from eight different ordinal sequences such as numbers,
74 days and months (see Appendix Q for more details). We also include different forms of the tokens,

75 as language model tokenizers often have several tokens corresponding to the same word, such as
 76 words with/without a space at the start being different tokens.

77 We determine whether attention heads perform succession by studying their *effective OV circuit*,
 78 which measures how the direct effect of input tokens when multiplied by an OV matrix W_{OV} . The
 79 (non-effective) OV circuit $W_U W_{OV} W_E$ (1) from Elhage et al. [9] is the inspiration for our *effective*
 80 *OV circuit* $W_U W_{OV} \text{MLP}_0(W_E)$ (2). Intuitively, (2)’s columns represent input tokens to the head
 81 and the rows represent the logits on each possible output token. We then operationalize *successor*
 82 *heads* by considering an input token T from our succession dataset (e.g. ‘Monday’). If the effective
 83 OV circuit column for input T has a larger output on the successor to T (‘Tuesday’) than on any other
 84 of the tokens in that task (‘Monday’ or ‘Wednesday’ or ‘Thursday’ ...) then we consider the head
 85 to have performed succession in this case. *Successor Heads* are then defined as the attention heads
 86 that pass this test for more than half of the tokens in the succession dataset. We call the proportion
 87 of succession dataset tokens on which an attention head performs succession the *succession score*.
 88 We plot successor scores across models of varying size in Figure 2.

89 To better understand the role successor heads fulfill in real datasets, we perform a case study of
 90 the successor head L12H0 in Pythia-1.4B (see Appendix J): we find the head displays interpretable
 91 polysemanticity, performing incrementation but also copying, acronym, and greater-than behaviour.

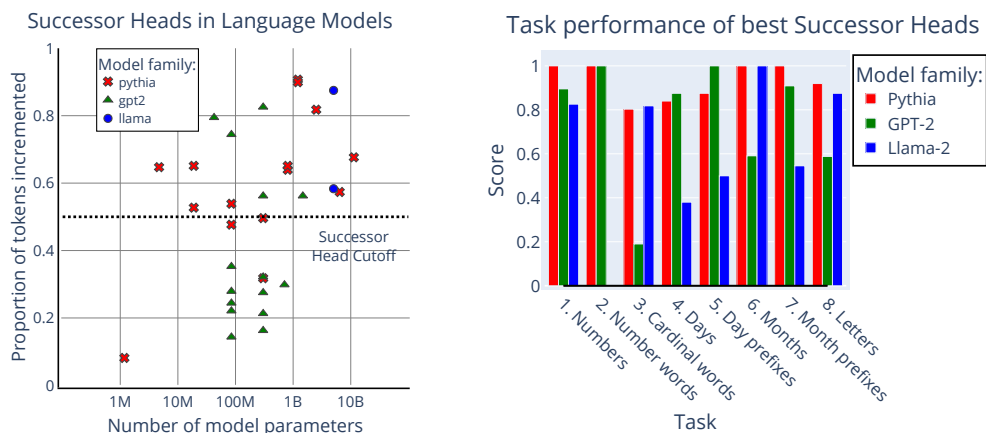


Figure 2: Plots of successor scores (proportion of tokens where succession occurs) for each model tested. A plot of the highest successor score observed across all attention heads for each model tested (left) and successor scores of the best successor heads in models (Pythia-1.4B, GPT-2 XL, Llama-2 7B) across different tasks (right).

92 3 Decomposing Numeric Representations

93 In the rest of this work, we perform a case study on the attention head (L12H0) with the maximal
 94 successor score in Pythia-1.4B. (Note that we also observe similar results across other models too –
 95 see Appendix C.3.) In Section 3.1, we find evidence for the existence of a *shared numeric subspace*
 96 within MLP0 representation-space. In Section 3.2, we find ‘mod-10 features’ in a decomposition of
 97 these representations, and we use these features to steer model behavior across different tasks.

98 3.1 Ordinal sequences are represented compositionally

99 We find evidence that ordinal sequences share a numeric subspace encoding the index of a token
 100 within its ordinal sequence. These findings are described in Appendix A.

101 3.2 Exploring mod-10 features

102 **Finding mod-10 features.** Now, given this evidence for a shared numeric subspace, a natural ques-
 103 tion to ask is whether these numeric embeddings themselves have any structure. To answer this
 104 question, we train a Sparse Autoencoder (SAE) (Ng [17], Elhage et al. [1], and Cunningham et al.

[18]) to recover the significant features across all tasks using reconstruction loss on the MLP0 outputs (see Appendix C for more details). Given a number token T and trained SAE, we define T 's **most important feature** as the SAE feature that, when ablated from the MLP0 output reconstruction, causes the biggest decrease in the probability of the successor of T (by calculating probabilities from the logits obtained by multiplying by $W_U W_{OV}$).

We find that the most important feature for a number is usually a common feature across other numbers in the same mod-10 class (on average, the most common feature in a mod class is shared by 5.85 numbers out of the 10). For example, most numbers amongst 3, 13, ..., 93 share the same important feature. This property gives rise to a mod-10 pattern in feature activations across the numbers, when we visualize the components of each most important feature across all the input prompts (Figure 3). We also verify that these mod-10 features have causal importance for the successor head's incrementation: multiplying these features by $W_U W_{OV}$ (in Figure 4) shows strong modulo 10 bands that are shifted by one to show that the successor head increments these features.

To define the **mod-10 features** we take the modal most important SAE feature for tokens within a particular mod class, and average this feature over 100 training runs, denoting the resulting features f_0, \dots, f_9 . We can then multiply these features by $W_U W_{OV}$ to study their effect on logits, as shown in Figure 5. We find unsurprisingly that f_i increases the logits on $f_{i+1 \pmod{10}}$. Further, the increase on logits for single-digit numbers is larger than the increase in logits for double-digit numbers.

Transferability of mod-10 features. Are our mod-10 features simply an artifact of the SAE technique? We provide evidence that the mod-10 features are natural, causally important features by using two independent methods to recover them; (1) linear probing and (2) identifying MLP0 neurons. We also show that the features transfer to other tasks in the succession dataset (Section 2). These findings are described in Appendix F.

Arithmetic with mod-10 features. The generalization of the linear probe to unseen numeric tasks provides us with evidence that there is a shared mod 10 structure across tasks. In this section, we stress test our understanding of this shared structure by trying to alter the index of ordinal sequence

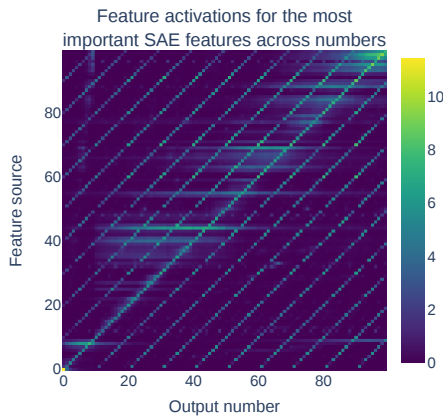


Figure 3: Feature activations, where for each number token (y-axis), we observe how its most important feature activates across all number tokens (x-axis). Values averaged over 100 SAE training runs.

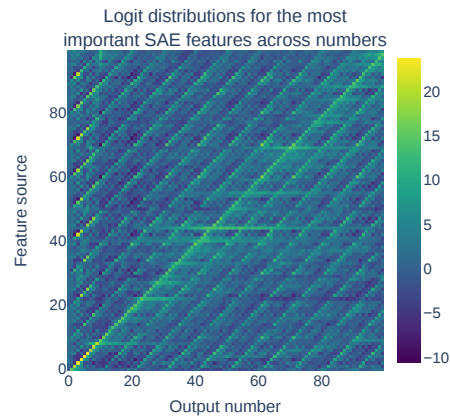


Figure 4: Logit distribution of important features, where for each number token (y-axis), we multiply its important feature by $W_U W_{OV}$ and observe how logits are attributed across all number tokens (x-axis).

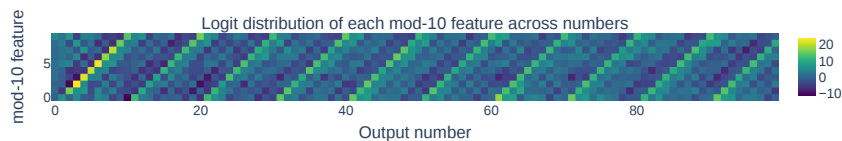


Figure 5: Logit distribution $W_U W_{OV} f_i$ for each mod-10 feature f_i .

131 tokens via vector arithmetic with the mod 10 features. For example, we expect $\text{MLP}_0(W_E(\text{'fifth'})) -$
 132 $f_5 + f_7$ (3) to be causally used by the model in a similar way to how the model would use a
 133 representation of the token 'seventh'. We use the successor head to test this hypothesis. If (3)
 134 behaves like $\text{MLP}_0(W_E(\text{'seventh'}))$, we would expect that multiplying (3) by $W_U W_{OV}$ attributes
 135 more logits on 'eighth' than any of the other tokens from task 2 in the succession dataset (Section 2).
 136 Indeed, this is correct as indicated by the circled checkmark (✓) in Figure 6. We can perform a similar
 137 experiment with all tokens in the succession dataset and with features other than f_7 added. The cases
 138 where the max logits are on the successor token are check-marked in Figure 6. We describe the
 139 experiment in more detail in Appendix G. We find that when the mod 10 addition feature is larger
 140 than the source value (modulo 10), vector arithmetic works on 53% (for months) and 89% (for digits
 141 20-29) of cases. Below the diagonal, we see mostly failures, and we find that this is due to successor
 142 heads possessing a *greater-than bias*, described in Appendix H.

	0	1	2	3	4	5	6	7	8	9		1	2	3	4	5	6	7	8	9	0	
20	✓	✓		✓	✓	✓	✓	✓	✓	✓	January	✓	✓	✓								
21		✓	✓	✓	✓	✓			✓	✓	February		✓	✓	✓							
22			✓	✓	✓	✓	✓	✓	✓	✓	March			✓	✓	✓						
23				✓	✓	✓	✓	✓		✓	April			✓	✓	✓						
24					✓	✓	✓	✓	✓	✓	May				✓	✓	✓					
25						✓	✓	✓	✓	✓	June			✓		✓	✓	✓	✓	✓		
26							✓	✓		✓	July							✓	✓	✓		
27								✓			August									✓	✓	
28										✓	September			✓							✓	✓
29										✓	first	✓	✓	✓	✓							
ten	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	second		✓	✓	✓	✓						
eleven		✓	✓	✓	✓	✓	✓	✓	✓	✓	third			✓	✓	✓	✓	✓				
twelve			✓	✓	✓	✓	✓		✓		fourth				✓	✓	✓	✓				
thirteen				✓	✓	✓	✓	✓	✓	✓	fifth					✓	✓	✓	✓	✓		
fourteen					✓	✓	✓	✓	✓		sixth				✓	✓	✓	✓	✓			
fifteen					✓	✓	✓		✓	✓	seventh					✓	✓	✓	✓			
sixteen						✓	✓	✓	✓	✓	eighth						✓	✓	✓	✓		
seventeen							✓	✓	✓	✓	ninth										✓	
eighteen								✓	✓	✓												
nineteen									✓	✓												
twenty									✓	✓												

Figure 6: Table displaying in which cases where vector arithmetic such as (3) are successful. Rows: source tokens. Columns: target residues modulo 10.

143 4 Conclusion

144 In this work, we discovered and interpreted a class of attention heads we call *successor heads*. We
 145 showed these heads increment tokens like numbers, months, and days partly by mapping them to
 146 an abstract mod-10 numeric representation. We provided evidence that successor heads exhibit a
 147 weak form of universality, arising in models across different architectures and scales, and using
 148 similar underlying mod-10 features in all cases. Additionally, we validated our understanding by
 149 demonstrating that a successor head reduced the loss on training data by predicting successor tokens.

150 Additional findings that stemmed from our work include:

- 151 1. Finding a greater-than bias, where a language model was much more likely to predict numer-
 152 ic answers larger than the values in the prompt, compared to smaller values than tokens
 153 present in the prompt, that was observable by a weights-level analysis.
- 154 2. Surprisingly interpretable individual MLP0 neurons on this narrow task.
- 155 3. A novel example of attention head polysemanticity (successor heads predicting acronyms).

156 Findings 1-3 could prompt further future work into how language models represent numeric con-
 157 cepts, particularly as 2 and 3 were surprising given existing evidence from existing work. Our work
 158 in finding a language model component that arises in models of many different scales (and uses
 159 abstract underlying numeric representations) may be a valuable contribution toward understanding
 160 the inner workings of frontier LLMs.

References

- [1] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. “Toy Models of Superposition”. In: *arXiv preprint arXiv:2209.10652* (2022).
- [2] Chris Olah. *Mechanistic Interpretability, Variables, and the Importance of Interpretable Bases*. <https://www.transformer-circuits.pub/2022/mech-interp-essay>. 2022.
- [3] OpenAI. *GPT-4 Technical Report*. 2023. arXiv: 2303.08774 [cs.CL].
- [4] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. “Progress measures for grokking via mechanistic interpretability”. In: *The Eleventh International Conference on Learning Representations*. 2023. URL: <https://openreview.net/forum?id=9XFSbDPmdW>.
- [5] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. *In-context learning and induction heads*. 2022. URL: <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- [6] Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. “Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 Small”. In: *The Eleventh International Conference on Learning Representations*. 2023. URL: <https://openreview.net/forum?id=NpsVSN6o4uL>.
- [7] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. “Zoom In: An Introduction to Circuits”. In: *Distill* (2020). DOI: 10.23915/distill.00024.001.
- [8] Yixuan Li, Jason Yosinski, Jeff Clune, Hod Lipson, and John Hopcroft. *Convergent Learning: Do different neural networks learn the same representations?* 2016. arXiv: 1511.07543 [cs.LG].
- [9] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. “A Mathematical Framework for Transformer Circuits”. In: *Transformer Circuits Thread* (2021). URL: <https://transformer-circuits.pub/2021/framework/index.html>.
- [10] Stefan Heimersheim and Jett Janiak. *A circuit for Python docstrings in a 4-layer attention-only transformer*. 2023. URL: <https://www.alignmentforum.org/posts/u6KXXmKFbXfWzoAXn/a-circuit-for-python-docstrings-in-a-4-layer-attention-only>.
- [11] Michael Hanna, Ollie Liu, and Alexandre Variengien. *How does GPT-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model*. 2023. arXiv: 2305.00586 [cs.CL].
- [12] Zhengxuan Wu, Atticus Geiger, Christopher Potts, and Noah D. Goodman. *Interpretability at Scale: Identifying Causal Mechanisms in Alpaca*. 2023. arXiv: 2305.08809 [cs.CL].
- [13] 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*. <https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html>. 2023.
- [14] Arthur Conmy, Augustine N. Mavor-Parker, Aengus Lynch, Stefan Heimersheim, and Adrià Garriga-Alonso. *Towards Automated Circuit Discovery for Mechanistic Interpretability*. 2023. arXiv: 2304.14997 [cs.LG].
- [15] Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. *Pythia: A Suite for Analyzing Large Language Models Across Training and Scaling*. 2023. arXiv: 2304.01373 [cs.CL].
- [16] Bilal Chughtai, Lawrence Chan, and Neel Nanda. *A Toy Model of Universality: Reverse Engineering How Networks Learn Group Operations*. 2023. URL: <https://arxiv.org/abs/2302.03025>.

- 217 [17] Andrew Ng. *Sparse Autoencoder*. Online Course Notes, Stanford University. 2011. URL:
218 <https://web.stanford.edu/class/cs294a/sparseAutoencoder.pdf>.
- 219 [18] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. *Sparse*
220 *Autoencoders Find Highly Interpretable Features in Language Models*. 2023. arXiv: 2309.
221 08600 [cs.LG].
- 222 [19] Lee Sharkey, Dan Braun, and Beren Millidge. *[Interim research report] Taking features out*
223 *of superposition with sparse autoencoders*. 2022. URL: [https://www.lesswrong.com/](https://www.lesswrong.com/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)
224 [posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-](https://www.lesswrong.com/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)
225 [of-superposition](https://www.lesswrong.com/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition).
- 226 [20] Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney,
227 Stella Biderman, and Jacob Steinhardt. *Eliciting Latent Predictions from Transformers with*
228 *the Tuned Lens*. 2023. arXiv: 2303.08112 [cs.LG].
- 229 [21] Stanford CRFM. *Mistral: A framework for transparent and accessible large-scale language*
230 *model training*. <https://github.com/stanford-crfm/mistral>. 2021.
- 231 [22] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. *Network Dissec-*
232 *tion: Quantifying Interpretability of Deep Visual Representations*. 2017. arXiv: 1704.05796
233 [cs.CV].
- 234 [23] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. “Feature Visualization”. In: *Dis-*
235 *till* (2017). <https://distill.pub/2017/feature-visualization>. DOI: 10.23915/distill.00007.
- 236 [24] Nicholas Goldowsky-Dill, Chris MacLeod, Lucas Sato, and Aryaman Arora. *Localizing*
237 *Model Behavior with Path Patching*. 2023. arXiv: 2304.05969 [cs.LG].
- 238 [25] Kevin Du, Lucas Torroba Hennigen, Niklas Stoehr, Alexander Warstadt, and Ryan Cotterell.
239 *Generalizing Backpropagation for Gradient-Based Interpretability*. 2023. arXiv: 2307.
240 03056 [cs.LG].
- 241 [26] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. *Distributed Rep-*
242 *resentations of Words and Phrases and their Compositionality*. 2013. arXiv: 1310.4546
243 [cs.CL].
- 244 [27] Jack Merullo, Carsten Eickhoff, and Ellie Pavlick. *Language Models Implement Simple*
245 *Word2Vec-style Vector Arithmetic*. 2023. arXiv: 2305.16130 [cs.CL].
- 246 [28] Nishant Subramani, Nivedita Suresh, and Matthew E. Peters. *Extracting Latent Steering Vec-*
247 *tors from Pretrained Language Models*. 2022. arXiv: 2205.05124 [cs.CL].
- 248 [29] Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte
249 MacDiarmid. *Activation Addition: Steering Language Models Without Optimization*. 2023.
250 arXiv: 2308.10248 [cs.CL].
- 251 [30] Neel Nanda. *A Comprehensive Mechanistic Interpretability Explainer & Glossary*. <https://www.neelnanda.io/glossary>. Dec. 2022.
252

253 A Factoring ordinal sequences

254 Let i_s denote the i th token in ordinal sequence s (such that e.g. 2_{Month} corresponds to the token
255 ‘February’), and let $\llbracket i_s \rrbracket = \text{MLP}_0(W_E(i_s))$ denote the model’s *internal representation* of token i_s
256 (the output of MLP_0 in Figure 1). Given successor heads $S = W_{OV}$ can increment tokens s_i from a
257 range of ordinal sequences s (e.g. numbers, months, days of the week), one might hypothesise that
258 numeric representations have *compositional structure* – i.e. that information about a token’s *position*
259 i in its ordinal sequence is encoded independently from information about *which ordinal sequence*
260 s it comes from. More precisely, we claim that we can decompose representations $\llbracket i_s \rrbracket$ into features
261 \mathbf{v}_i living in some ‘index space’ and \mathbf{v}_s living in some ‘domain space’, such that $\llbracket i_s \rrbracket = \mathbf{v}_i + \mathbf{v}_s$.

262 **Method.** To test this compositionality hypothesis, we wish to learn two linear maps – an *index-space*
263 *projection* $\pi_N : \mathbb{R}^{d_{\text{model}}} \rightarrow \mathbb{R}^{d_{\text{model}}}$ and a *domain-space projection* $\pi_D : \mathbb{R}^{d_{\text{model}}} \rightarrow \mathbb{R}^{d_{\text{model}}}$ – such that,
264 for all pairs of tokens i_s and j_t (with i_t a valid token), $\llbracket \hat{i}_t \rrbracket := \pi_N(\llbracket i_s \rrbracket) + \pi_D(\llbracket j_t \rrbracket) \approx \llbracket i_t \rrbracket$. To do
265 so, we enforce that $\pi_N + \pi_D = I$, and ensure predicted representations $\llbracket \hat{i}_t \rrbracket$ ‘behave like’ ground
266 truth representations $\llbracket i_t \rrbracket$ for randomly sampled pairs of tokens i_s and j_t – in other words, that there
267 is low L2-distance between $\llbracket \hat{i}_t \rrbracket$ and $\llbracket i_t \rrbracket$, that $\llbracket \hat{i}_t \rrbracket$ decodes to i_t (**output-space decoding**), and that
268 $S(\llbracket i_t \rrbracket)$ decodes to $(i + 1)_t$ (**successor decoding**). For full experimental details see Appendix B.

		Source token (index)											
No succ		I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII
Source token (sequence)	1	9	2	3	4	22	6	7	8	9	24	11	12
	one	nine	two	three	four	twenty	six	seven	eight	nine	twenty	eleven	twelve
	first	ninth	second	third	fourth	second	sixth	seventh	eighth	ninth	fourth	-	-
	Monday	Sunday	Tuesday	Wednesd	Thursday	Tuesday	Saturday	Sunday	-	-	-	-	-
	Mon	Sep	Tue	Wed	Thu	Tue	Sat	Sun	-	-	-	-	-
	January	Septemb	February	March	April	February	June	July	August	Septemb	Decembe	Novembe	Decembe
	Jan	Sep	Feb	Mar	Apr	Feb	Jun	Jul	Aug	Sep	Dec	Nov	Dec
	A	I	B	C	D	V	F	G	H	I	X	W	L
		Source token (index)											
With succ		I	II	III	IV	V	VI	VII	VIII	IX	X	XI	XII
Source token (sequence)	1	10	Third	Fourth	5	23	VII	8	9	10	25	12	13
	one	10	Third	fourth	fifth	23	seventh	eighth	ninth	10	25	12	13
	first	tenth	Third	fourth	fifth	VI	seventh	ninth	ninth	tenth	-	-	-
	Monday	RS	Third	MC	Friday	MV	VII	-	-	-	-	-	-
	Mon	RS	HM	HM	HM	MV	MTP	-	-	-	-	-	-
	January	RS	cs	cs	May	23	cs	cs	Septembe	October	25	cs	-
	Jan	CS	cs	cs	May	23	cs	cs	Sept	rs	35	cs	-
	A	III	III	MC	Fifth	VV	VII	VIII	9	cx	Y	XII	iii

Table 1: A table presenting top-1 predicted tokens (under output-space projection and successor decoding) from representations $\pi_D(1_s) + \pi_N(i_{Rom})$. Green cells denote predictions which match their target exactly; yellow cells denote predictions which match their target modulo formatting (e.g. ‘six’ rather than ‘6’); red cells denote incorrect predictions.

269 **Results.** On our held-out dataset of token pairs, we obtained a top-1 output-space decoding accuracy
270 of 1.00, and a top-1 successor decoding accuracy of 0.90.² To explore out-of-distribution perfor-
271 mance, we also test whether π_N can project out the numeric component of Roman numerals (which
272 weren’t in the successor dataset), by taking Roman numerals $i_{Rom} \in \{‘I’, \dots, ‘XII’\}$ and tokens 1_s
273 from sequences s in the successor dataset, and testing whether $\pi_D(1_s) + \pi_N(i_{Rom})$ decodes to i_s .
274 We present the top-1 predicted tokens (under both output-space and successor decoding) in Table 1.
275 Observe that, while output-space decoding yields perfect top-1 accuracy (apart from $i \in \{1, 5, 10\}$),
276 which we can attribute to the Roman numerals I, V and X being single-letter and impossible to dis-
277 ambiguate from 9_{Letter} , 22_{Letter} and 24_{Letter}), successor decoding achieves an accuracy of 0.125 (or
278 0.29 if we allow for incorrectly formatted predictions).

279 These results – in particular, our ability to project the numeric component out of tokens from unseen
280 sequences and transfer indices across domains – suggest that there is a shared numeric subspace
281 storing the index of a token within its ordinal sequence. Indeed, informal testing suggests that this
282 numeric subspace may be interpretable even for tokens not part of an ordinal sequence: for instance,
283 $d(\pi_N(\llbracket ‘triangle’ \rrbracket)) + \pi_D(1_{Num})$ yields 3_{Num} , and $d(\pi_N(\llbracket ‘week’ \rrbracket)) + \pi_D(1_{Num})$ yields 7_{Num} .

284 Despite this, the drop in performance when applying the successor head to our constructed repre-
285 sentations (and in particular, the leakage of Roman-numeral information into $\pi_D(1_s) + \pi_N(i_{Rom})$ –
286 notice the VII and VIII in Table 1) suggests our numeric projection π_N might be capturing slightly
287 more than just the numeric subspace. Specifically, there may be some components of domain-space
288 which are ignored by output-space decoding, but which our successor head lifts into output-space.

289 B Training details for compositionality experiments

290 **Remark: obtaining a decoding function.** Recall that we wish to learn two linear maps – an *index-*
291 *space projection* $\pi_N : \mathbb{R}^{d_{model}} \rightarrow \mathbb{R}^{d_{model}}$ and a *domain-space projection* $\pi_D : \mathbb{R}^{d_{model}} \rightarrow \mathbb{R}^{d_{model}}$ – such
292 that, for all pairs of tokens i_s and j_t (with i_t a valid token), $\llbracket i_t \rrbracket := \pi_N(\llbracket i_s \rrbracket) + \pi_D(\llbracket j_t \rrbracket) \approx \llbracket i_t \rrbracket$.
293 To evaluate the above identity, we want a decoding function $d : \mathbb{R}^{d_{model}} \rightarrow \text{Logits}$, such that
294 $d(\llbracket i_s \rrbracket) = i_s$. Given the informal observation that directly unembedding $\llbracket i_s \rrbracket$ yields next-token
295 predictions for i_s , whereas unembedding $S(i_s)$ yields $(i+1)_s$, we hypothesise that the unembedding
296 matrix W_U reads from some ‘output space’ \mathcal{O} and the embedding transform $\llbracket \cdot \rrbracket$ writes to some ‘input
297 space’ \mathcal{I} – and that the successor head reads from \mathcal{I} and writes to \mathcal{O} . Indeed, when training an
298 *output-space projection* $\pi_O : \mathcal{I} \rightarrow \mathcal{O}$ over tokens in the vocabulary such that $W_U(\pi_O(\llbracket i_s \rrbracket)) = i_s$,

²Note that, as our successor dataset contains 1041 tokens, a random classifier (even when restricted to tokens in the successor dataset) would achieve an accuracy of 0.001.

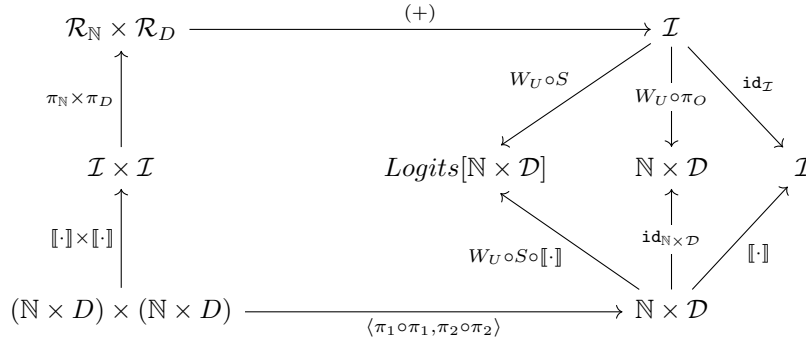
299 we obtain 97.4% top-1 accuracy on a set of 1000 held-out tokens – which both confirms the output-
 300 space hypothesis, and gives us a decoding function $d(\mathbf{x}) = W_U(\pi_O(\mathbf{x}))$.

301 **Method.** With our decoding function in hand, we can train π_N and π_D to satisfy our identity.
 302 Specifically, we define π_N and π_D to be matrices such that $\pi_N + \pi_D = I$. For valid token pairs i_s
 303 and j_t , we obtain predicted representations $\llbracket \hat{i}_t \rrbracket = \pi_N(\llbracket i_s \rrbracket) + \pi_D(\llbracket j_t \rrbracket)$, and minimise a combination
 304 of ‘closeness metrics’:

$$\|\llbracket \hat{i}_t \rrbracket - \llbracket i_t \rrbracket\|^2 + \mathcal{L}(W_U(\pi_O(\llbracket \hat{i}_t \rrbracket)), i_t) + \mathcal{L}(W_U(S(\llbracket \hat{i}_t \rrbracket)), W_U(S(\llbracket i_t \rrbracket)))$$

305 for \mathcal{L} the cross-entropy loss. Specifically, we ensure that predicted and ground truth representations
 306 ‘behave in the same way’ – in other words, that they are close together, that predicted representations
 307 $\llbracket \hat{i}_t \rrbracket$ decode to tokens i_t (**output-space decoding**), and that the logit distribution when decoding in-
 308 cremented predicted representations $S(\llbracket \hat{i}_t \rrbracket)$ matches that when decoding incremented ground truth
 309 representations $S(\llbracket i_t \rrbracket)$ (**successor decoding**).

310 More succinctly, we can frame the training procedure as learning π_N, π_D such that the following
 311 diagram commutes:



312 We trained for 10 epochs over valid token pairs sampled from the succession dataset, and evaluated
 313 on a held-out dataset of 500 randomly-sampled token pairs.

314 C Sparse auto-encoders

315 C.1 Definition

316 We refer to a single-layer autoencoder with a sparsity regularization term in its loss as a **sparse**
 317 **auto-encoder**.

318 For a dataset generated from a set of underlying vectors (each dataset example is a sparse linear
 319 combination of such vectors), it has been empirically observed [19, 18] that sparse auto-encoders are
 320 capable of retrieving the underlying set of vectors. We hope to obtain a set of sparse, interpretable
 321 features from the SAEs that decompose some of the structure of MLP0 space that we can use to
 322 analyze the way numeric operations are performed.

323 C.2 Training process for mod 10 features

324 Training a sparse auto-encoder with D features and regularization coefficient λ on a dataset of tokens
 325 in MLP0 space results in a map $F : \text{str} \rightarrow (\mathbb{R}^d \times \mathbb{R}_+)^D$, with $F(x) = \{(f_1, a_1), \dots, (f_D, a_D)\}$,
 326 mapping a token to a set of feature and feature-activation pairs, with reconstruction $R(x) =$
 327 $\sum_{i=1}^D a_i f_i$.

328 We train the SAE using number tokens from 0 to 500, both with and without a space (‘123’ and
 329 ‘ 123’), alongside other tasks, such as number words, cardinal words, days, months, etc. 90% of
 330 these tokens go into the train set, and the remaining 10% to the validation set. Even with the other
 331 tasks, the dataset is dominated by numbers, but creating a more balanced dataset would give us less
 332 data to work with, and without enough data, the SAE fails to generalize to the validation set. Hence,
 333 we only concern ourselves with the features that the SAE learns for number tokens, and we then

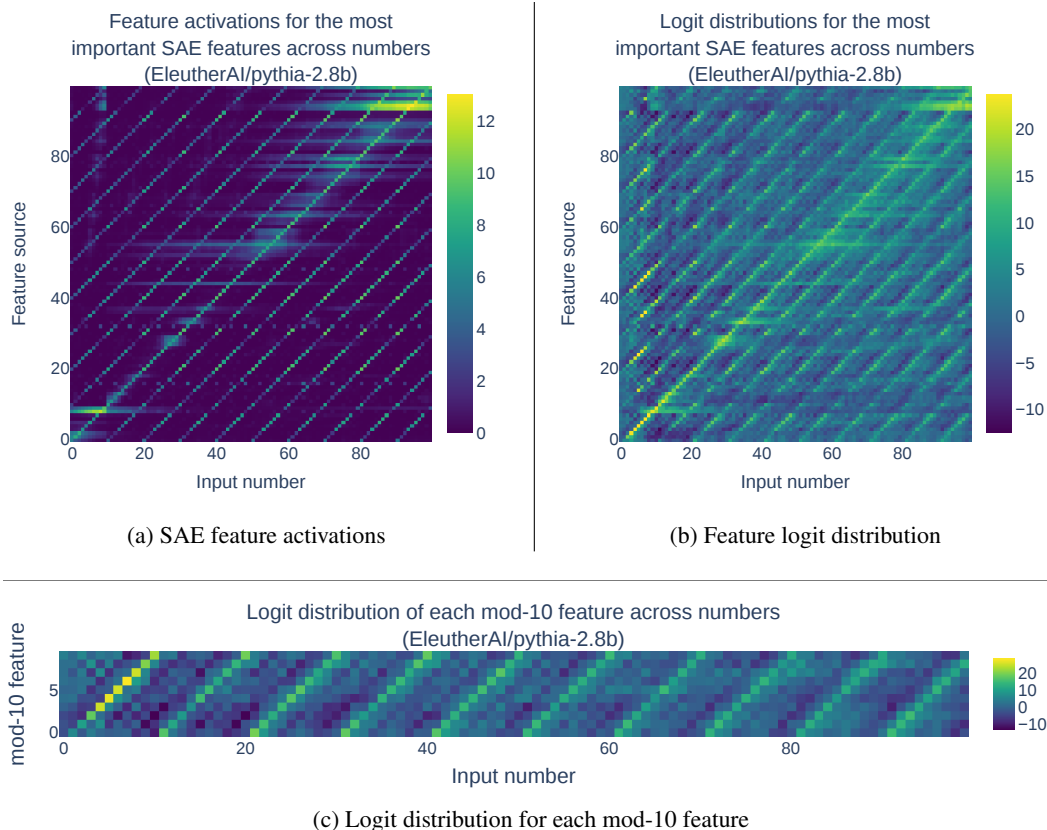


Figure 7: SAE plots for Pythia-2.8b analogous to Figure 3, Figure 4, and Figure 5

334 separately check whether these features generalize to the other tasks on the basis of logits, rather
 335 than SAE activations.

336 We used the hyperparameters $D = 512$ and $\lambda = 0.3$, with a batch size of 64, and trained for 100
 337 epochs. To find these hyperparameters, we used the metric of mean max cosine similarity between
 338 two trained SAEs, as described in Sharkey, Braun, and Millidge [19] and Cunningham et al. [18].

339 C.3 Universality of mod-10 results

340 We also observe the mod 10 structure via SAEs across models other than Pythia-1.4B, without any
 341 finetuning of SAE parameters to these models. We reproduce the SAE figures seen in Section 3.2 for
 342 other models, with Appendix C.3 for Pythia-2.8B, and Appendix C.3 for celebrimbor-gpt2-medium-
 343 x81.

344 D Linear probing

345 We train a linear probe to predict the mod 10 value of tokens. Specifically, we train on number
 346 tokens from '0' to '500', both with and without a space, assigning 90% of tokens to a train set, and
 347 the remaining 10% to a valid set. We use a learning rate of 0.001, and a batch size of 32, for 100
 348 epochs.

349 We then evaluate on a dataset of unseen tasks, including number words (from 'one' to 'nineteen'),
 350 placements, numerals, months, days, and any valid spaced and capitalized variants. Out of the total
 351 102 such examples, 94/102 are correct, and the 8 failures are: ['January', 'December', 'Friday',
 352 'Saturday', 'Sunday', 'V', 'X', 'XV'].

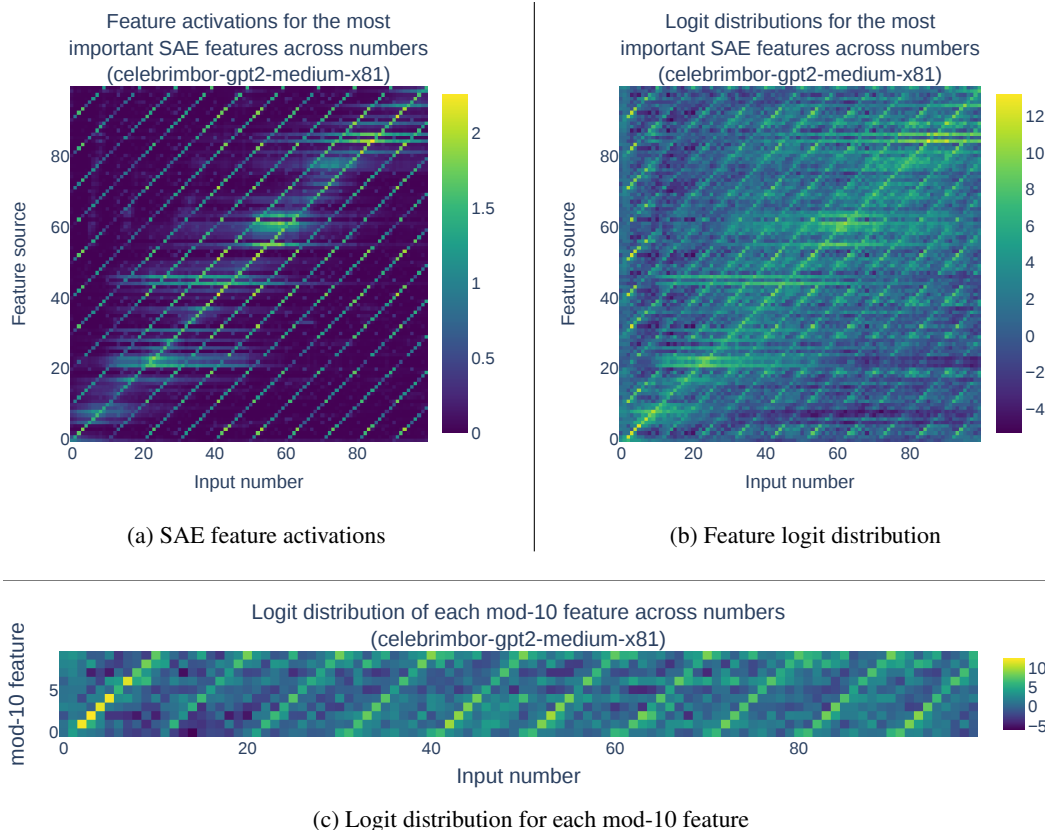


Figure 8: SAE plots for celebrimbor-gpt2-medium-x81 analogous to Figure 3, Figure 4, and Figure 5

353 The failures of 3 out of 7 days are consistent with our inability to interpret the day task well with
 354 our mod 10 features. Additionally, we see ‘January’ and ‘December’ as failure cases, which is also
 355 consistent with our finding that there does not seem to be a mod 10 feature that corresponds to any
 356 of them: f_1 behaves as ‘November’ rather than ‘January’, and f_2 as ‘February’.

357 E MLP0 neurons

358 In our MLP0 neuron experiments, we do the following: for each $T \in \{‘0’, ‘1’, \dots, ‘99’\}$, we ablate
 359 each neuron from the final activation in MLP0 (the final activation is just before the final linear layer
 360 of MLP0), and store the probability attributed to the successor of T after passing the modified (due
 361 to ablation) MLP0 output through the successor.

362 Averaging the correct probabilities across all 100 prompts then gives an averaged correct probability
 363 for each neuron after ablation. We then look at the intensities and logits (across number tokens) for
 364 neurons with the lowest correct probability after ablation, meaning they have the most impact on
 365 successorship when ablated. This gives us the plots seen in Figure 9.

366 F Transferability of mod-10 features

367 **(1) Linear probing:** we train a linear probe to predict the mod-10 value across numbers from
 368 their MLP0 representations. We find that the i th row of the linear probe matrix has a high cosine
 369 similarity (on average 0.70764) to the corresponding mod-10 feature f_i obtained with the SAE. This
 370 suggests that the features f_i are likely to be directly finding the mod-10 value of input tokens rather
 371 than picking up on a property that is correlated with mod classes. Further, the probe generalizes,

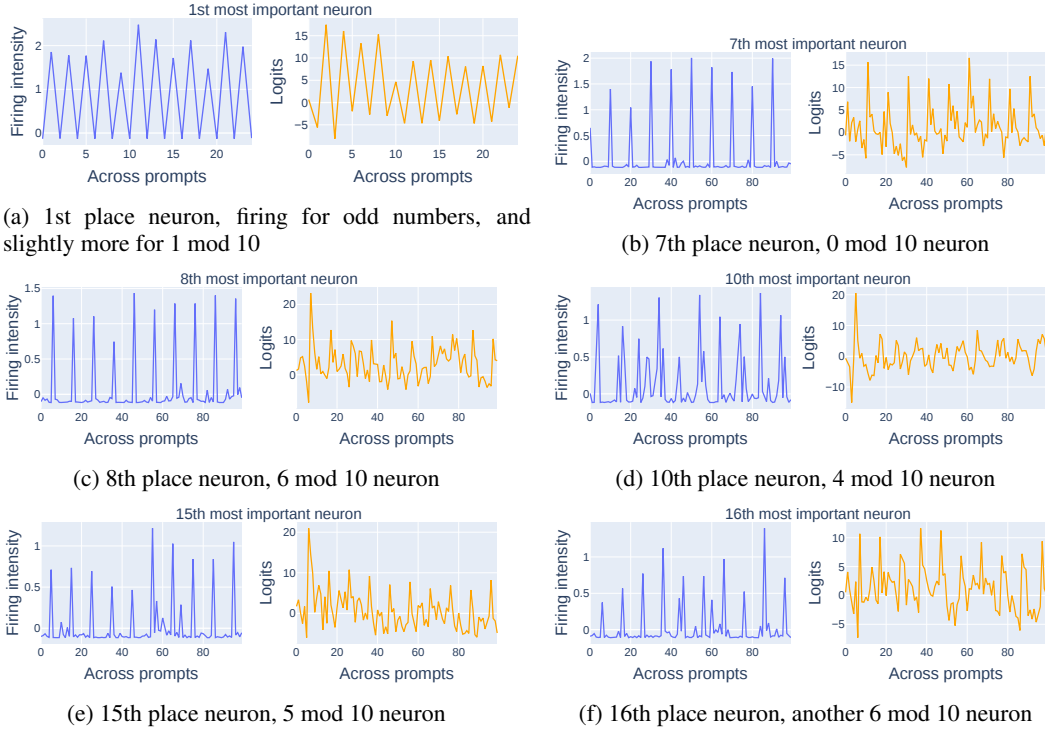


Figure 9: Some examples of neurons firing strongly in modulo 10 patterns out of the top 16 most important MLP0 neurons for successorship.

372 correctly predicting the mod-10 value for 94/102 examples from the succession dataset tasks 2-8
 373 (Section 2). Appendix D describes our full experimental setup.

374 **(2) MLP0 neurons:** we perform ablative experiments on individual MLP0 neurons to find the most
 375 important neurons for successorship across numbers, as measured by probability decrease as in the
 376 SAE experiments. This reveals neurons that strongly fire in a mod-10 pattern, which is the most
 377 common frequency amongst the top 16 most important neurons. Some examples of such neurons in
 378 the top 16 most important neurons can be seen in Figure 9. We also verify that the neurons indeed
 379 increase probability on successor tokens by multiplying their output values by $W_U W_{OV}$ in the same
 380 figures. Further technical details can be found in Appendix E.

381 Our results on the interpretability of individual neurons are surprising given recent research sug-
 382 gesting that the individual neurons of language models may be inappropriate as the units of un-
 383 derstanding [1]. However, our results do not contradict previous finding that understanding MLPs
 384 requires understanding distributed behaviors, since for example the 6 mod 10 feature appears to be
 385 in superposition across at least two neurons (Figure 9c, 9f).

386 **G Arithmetic experiments**

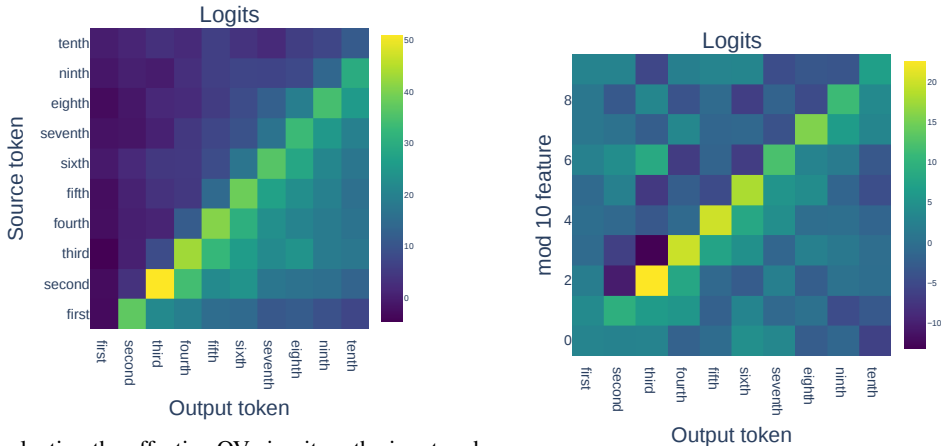
387 For a token $T \in V$ (row of arithmetic table) in numeric class V and feature f_i (column of arithmetic
 388 table), we consider how $x := \text{MLP}_0(W_E(T)) + k(-f_{\text{ord}(T)} + f_i)$ attributes logits to tokens in V ,
 389 with $k \in \mathbb{R}_+$ a scaling, and $\text{ord}(T)$ the numeric order of T with respect to V . We denote whether
 390 x correctly attributes maximal logits to the token $U \in V$ with $\text{ord}(U) = i + 1$ by a checkmark,
 391 giving Figure 6. Since our mod-10 features $\{f_i\}_i$ obtained from the SAE are normalized to unit
 392 norm, hence some scaling is necessary in order to have an effect on the order of the tokens. We pick
 393 $k = \text{MLP}_0(W_E(T)) \cdot f_{\text{ord}(T)}$ everywhere other than months, where we observe that we must use a
 394 scaling of 2 times this to affect the order while keeping the task identity intact.

395 **H Greater-than bias**

396 The vector arithmetic experiments (Figure 6) mostly fail below the diagonal, when the mod-10
 397 addition is smaller than the source tokens’s ordinal sequence position mod 10. This is because
 398 successor heads are biased towards values greater than the successor, compared to values less than
 399 the successor.

400 This effect can be seen in Figure 10a on the tokens ‘first’, ‘second’, ..., ‘tenth’. However, our mod-
 401 10 features do not exhibit a greater-than bias, as seen in Figure 10b. We survey these effects across
 402 all tasks in Appendix M. As a result, using the mod-10 features to shift logits towards tokens of a
 403 lower order than the input token fails, as the changes in logits are not significant compared to the
 404 large logits on higher-order tokens. In the case of numbers, this leads to the effect that, for example,
 405 $W_U W_{OV}(\text{MLP}_0(W_E('35'))) - f_5 + f_3$ has high logits on ‘44’, rather than ‘34’ (this ‘+10’ effect
 406 occurs for 2/3 of entries below the diagonal in the 20-29 numbers table of Figure 6).

407 **Limitations:** The absence of a strong greater-than bias in our mod 10 features suggests this feature-
 408 level description is missing some details – specifically that successor heads must use other numeric
 409 information to produce the greater-than bias we observe. Additionally, while we see a good general-
 410 ization of the mod 10 features across various tasks in the table in Figure 6, the mod 10 features are
 411 not able to steer the Days and Letters tasks from Section 2. We describe this in Appendix I.



(a) Evaluating the effective OV circuit on the input and output tokens ‘first’, ‘second’, ..., ‘tenth’. (b) Multiplying all mod 10 features f_i by $W_U W_{OV}$.

Figure 10: The Successor Head OV circuit displays a clear bias against decrementation (Figure 10a), i.e. the logits on or above the main diagonal are less than the logits below the main diagonal. This bias isn’t captured in the mod 10 feature (Figure 10b).

412 **I Failure cases of mod 10 features**

413 For the day and alphabet task, analogously to Figure 10, we look at the logits across the task and the
 414 mod 10 features. These are displayed in Figure 11, and demonstrate that our mod 10 features are not
 415 very interpretable in the context of days and the alphabet in terms of logits, with no clear diagonal
 416 of high logits.

417 **J Successor Heads in the Wild**

418 In this section, we analyze the behaviour of successor heads within natural-language datasets, and
 419 observe that they aren’t simply responsible for incrementation: indeed, we identify four distinct,
 420 interpretable categories of successor head behavior, highlighting successor heads as an example of
 421 an *interpretablely polysemantic* attention head ‘in the wild’.

422 In order to characterize the behavior of Pythia-1.4B’s successor head on natural-language data, we
 423 sample 128 length-512 contexts from *The Pile*, and for each prefix of each context, we assess

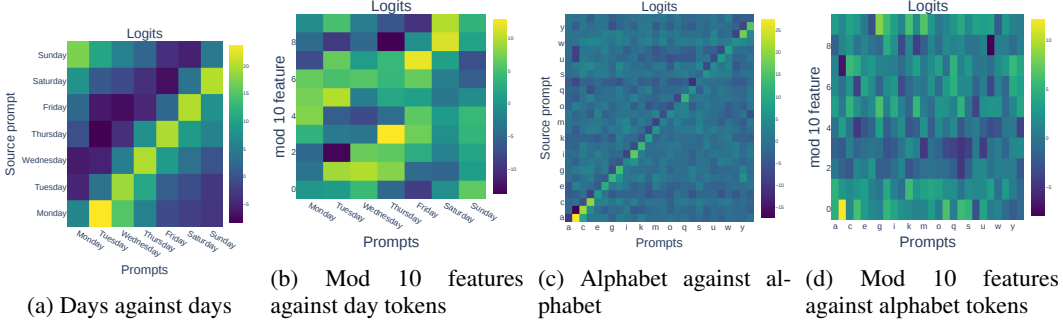


Figure 11: Logit plots across day and alphabet tasks, where attempting to steer the model with mod 10 features fails.

424 whether the successor head is important for the model’s ability to predict the correct next token.
 425 We evaluate prefixes using two different metrics for per-prompt successor head importance:

426 **Winning cases.** We identify prefixes where the head that most decreases the logit for the correct
 427 next token under direct effect mean ablation is the successor head, denoting them as *winning cases*.

428 **Loss-reducing cases.** We identify prefixes p where direct effect mean ablation of the successor head
 429 increases next-token prediction loss (by $\Delta\mathcal{L}(p)$), denoting them as *loss-reducing cases*.

430 J.1 Interpretable Polysemanticity in Successor Heads

431 On analyzing prefixes where the successor head is particularly important for next-token prediction
 432 – i.e. loss-reducing and winning cases – we observe four main categories of behavior, which we
 433 operationalize as follows (denoting a *top- n -attended token* as a token at one of the top n positions
 434 to which the successor head attends most strongly):

435 **Successorship behavior:** *the successor head pushes for the successor of a token in the context.* We
 436 say this behavior occurs when one of the top-5-attended tokens is in the successorship dataset, and
 437 the correct next token is the successor of t .

438 **Acronym behavior:** *the successor head pushes for an acronym of words in the context.* We say this
 439 behavior occurs when the correct next token is an acronym whose last letter corresponds to the first
 440 letter of the top-1-attended token. (For example, if the successor head attends most strongly to the
 441 token ‘Defense’, and the correct next token is ‘OSD’.)

442 **Copying behavior:** *the successor head pushes for a previous token in the context.* We say this
 443 behavior occurs when the correct next token t has already occurred in the prompt, and token t is one
 444 of the top-5-attended tokens.

445 **Greater-than behavior:** *the successor head pushes for a token greater than a previous token in the*
 446 *context.* We say this behavior occurs when we do not observe successorship behavior, but when the
 447 correct next token is still part of an ordinal sequence and has greater order than some top-5-attended
 448 token (e.g. if the successor head attends to the token ‘first’ and the model predicts the token ‘third’.)

449 We plot the proportions of each behavior observed across winning cases in Figure 12, and the frac-
 450 tion of total reduced loss over all contexts ($\Delta\mathcal{L}$) attributable to contexts of each behavior in Fig-
 451 ure 14. We also illustrate a random sample of 5 winning cases in Figure 13, and of 5 loss-reducing
 452 cases in Figure 15. We observe that, while successorship is the predominant behavior across both
 453 winning and loss-reducing cases, acronym and greater-than behaviors also form a non-negligible
 454 fraction of successor head behavior. In other words, the successor head is an example of an attention
 455 head with *interpretable polysemanticity*³. While polysemanticity has been observed in both vision
 456 models [7] and toy models trained to perform simple tasks [1], to the best of our knowledge the
 457 presence of both successorship and acronym behavior in head L12H0 is the cleanest example of
 458 polysemantic behavior identified so far in an LLM.

³A component of a network is said to be *interpretable polysemantic* if it performs multiple distinct, inter-
 pretable functions.

Proportions of behaviours in winning cases

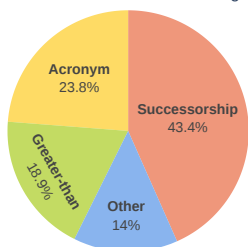


Figure 12: Proportions of three dominant behaviours across winning cases (copying behavior is negligible here).

Prompt	Completion
(...) [B2] Hence, bonding to ceramic requires strict attention to detail for optimal clinical outcomes.[@B	3
(...) designated as boxazomycin A and (... called Generalized Single Step Single Solve (B GS
(...) More than two-	thirds
(...) where one or more access points (AP

Figure 13: Random sample of 5 winning cases.

Proportions of behaviours in most loss reducing cases

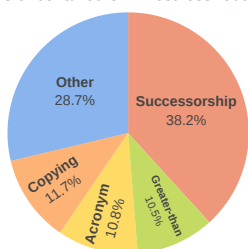


Figure 14: Proportions of reduced loss ($\Delta\mathcal{L}$) attributable to prompts of each behaviour.

Prompt	Completion
(...) low-dose administration of 14- and	15
(...) in the second round, let's get it to the	third
(...) for 1) investigators to ask the missionaries,	2
(...) You are using Microsoft Test Manager (MT
(...) known as the Medical Research Council Technology (MR

Figure 15: The top 5 most loss reducing examples.

459 Note that in this section, while we identified that successor heads are often used in tasks involving
 460 incrementation, we did not explicitly demonstrate that successor heads are *necessary* for incremen-
 461 tation. In Appendix O we describe an experiment that reveals that successor heads are necessary for
 462 a specific incrementation task (*numbered listing*).

463 K Residual Connections Are Not Important For Succession

464 To show that there is no relevant information in the residual stream, i.e. the path $W_U \text{MLP}_0(W_E)$
 465 is not sufficient to predict successors, we perform an experiment using the Tuned Lens [20], which
 466 approximates the optimal predictions after a given layer inside a transformer.

467 For all tasks in the succession dataset (Section 2), we used prompt formats (where | denotes a gap
 468 between tokens)

- 469 1. |Here| is| a| list|:| alpha| beta| gamma| and| here| is| another|:|<token1>|
- 470 2. |The|Monday|Tuesday|Wednesday| and| The|<token1>|

471 in order to measure how well models were able to predict the successor <token2> (e.g ‘February’)
 472 given the final token of the prompt was <token1> (e.g ‘January’), as LLMs, predict successors
 473 given these prompts.

474 We then took GPT-2 Small and Pythia-1.4B’s output after MLP0 and used the Tuned Lens to get
 475 logits on output tokens.⁴ The resulting successor score was <1% and commonly predicted bigrams,
 476 such as <token1>=“first” giving “time” as a completion and <token1>=‘Sunday’ giving ‘morn-
 477 ing’ as a prediction. This suggests that MLP0 information is insufficient for incrementation and the
 478 successor head is critical for succession.

⁴Note we did run with attention layer 0 to maximise the model’s chances are being able to perform suc-
 sion.

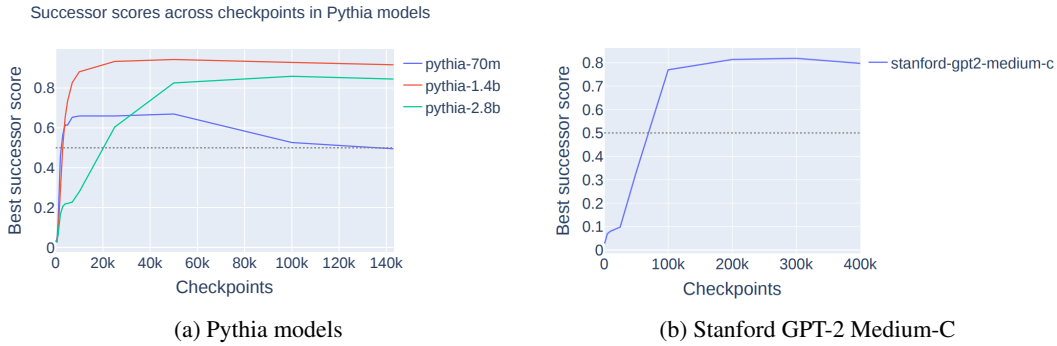


Figure 16: Best successor scores across successor heads throughout training checkpoints for Pythia and stanford-gpt2 models.

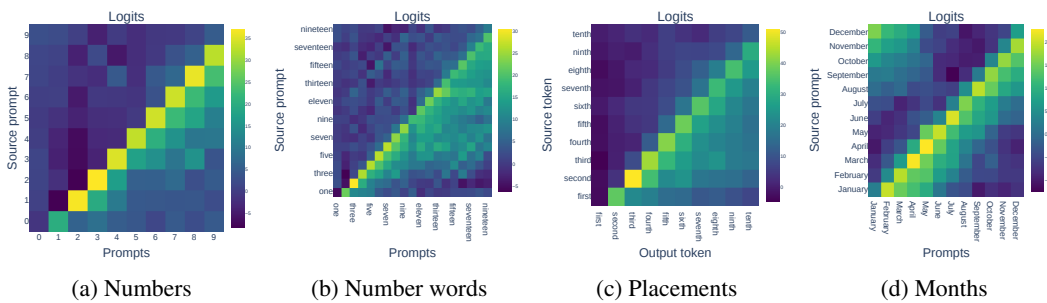


Figure 17: Plots of logits across various numeric classes, analogous to Figure 10a

479 L Testing successor score over training steps

480 Another line of evidence that Successor Heads are an important model component for low training
 481 loss can be found by studying successor scores across training points. We study a Pythia model [15]
 482 as well as a Stanford GPT model [21], as these models have training checkpoints. The emergence
 483 of Successor Heads throughout training is displayed in Figure 16.

484 M Decrementation bias across different tasks

485 We show the strength of the decrementation bias in figures Figure 17 and 18.

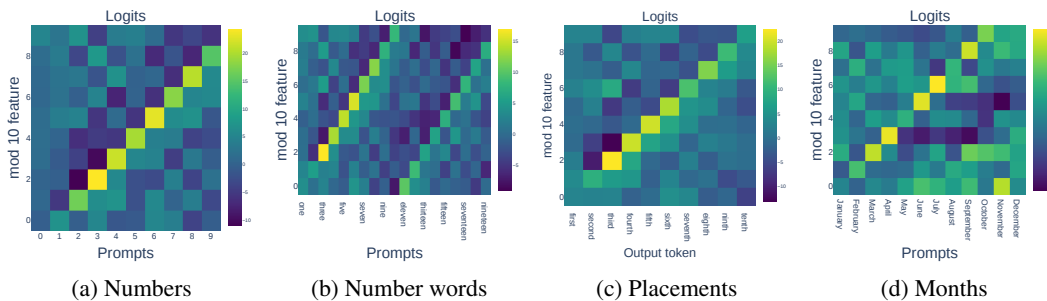


Figure 18: Plots of mod 10 feature logits across various numeric classes, analogous to Figure 10b

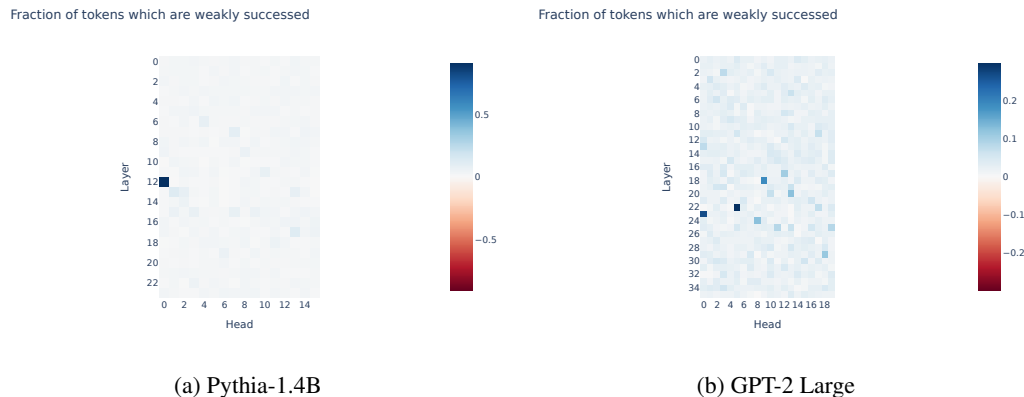


Figure 19: Successor scores for Pythia-1.4B and GPT-2 Large.

Prompt	Answer
(...) (A) Colony formation and (<	B
(...) (i) $f_g^*(y)$ equals the factual density $f(y)$ for all $g \in \mathbb{G}$; (<	ii
(...) [^2]: Conceived and designed the experiments (...) [^	3
(...) 6. Kirovsky Zavod Station – St. Petersburg, Russia (...) you can see a statue of Lenin here.	7
(...) [9] Minutes, Criminal Law Revision Commission, January 28, 1972, 16.[10

Figure 20: Some examples of numbered listing prompts from the Pile dataset.

486 N All successor scores in a model

487 In Figure 19 we find that for both Pythia-1.4B (the mainline model in the paper) and GPT-2 Large (a
 488 randomly selected model without a successor head from Figure 2 on the left), the heads with highest
 489 successor score are sparse: in Pythia-1.4B L12H0 has eight times as great a successor score to the
 490 next higher successor score and in GPT-2 Large only three heads have a successor score that’s above
 491 1/10.

492 O Case study: numbered listing

493 In Appendix J we demonstrate that when the successor head is contributing usefully, the prompts
 494 often required some kind of incrementation. However, we want to investigate whether the converse
 495 holds: are prompts requiring incrementation mostly solved by successor heads?

496 Numbered listing is widespread across real datasets and requires incrementation. Additionally, even
 497 small LLMs are capable of this task in the case of incrementing citations.⁵ Examples of prompts
 498 involving numbered listing can be seen in Figure 20.

499 We collect 64 such prompts and check for whether the successor head in Pythia-1.4b is the most
 500 important head for this prompt under direct effect mean ablation, and find that the successor head is
 501 indeed the most important head across **all** 64 prompts. Hence this provides some evidence prompts
 502 requiring incrementation in real datasets are indeed mostly solved by successor heads.

503 P Related Work

504 **Mechanistic Interpretability** research aims to reverse engineer trained neural networks analogously
 505 to how software binaries can be reverse-engineered [2]. This research was largely developed
 506 in vision models [22, 23] though most recent research has studied language models [9, 5] and trans-

⁵<https://www.lesswrong.com/posts/LkBmAGJgZX2tbwGKg>

507 formers [4]. Olah et al. [7] introduces the universality hypothesis and we use Chughtai, Chan, and
 508 Nanda [16]’s ‘weak universality’ notion in this work (Section 1).

509 **Transformer circuits.** More specifically, our work builds from the insights of Elhage et al. [9]’s
 510 framework for understanding circuits in transformers, including how autoregressive transformers
 511 have a residual stream. Due to the residual stream, different paths from input to output can bypass
 512 as many attention heads and MLPs as necessary. This has further been explored in specific case
 513 studies [6, 24] and generalizes to backwards passes [25]. One related case study to our work is
 514 Hanna, Liu, and Variengien [11] which studies a Greater-Than circuit in GPT-2 Small, similar to
 515 how we indirectly found the Greater-Than operation in Section 3. Hanna, Liu, and Variengien [11]
 516 focus mainly on numbers, not other tasks.

517 **LLMs and arithmetic.** Mikolov et al. [26]’s seminal work on word embedding arithmetic showed
 518 that latent language model representations had compositionality, e.g. ‘King’ – ‘Man’ + ‘Woman’
 519 approximated the embedding of ‘Queen’. Recently Merullo, Eickhoff, and Pavlick [27] showed
 520 some extension of these arithmetic results to LLMs. Additionally Subramani, Suresh, and Peters
 521 [28] and Turner et al. [29] use residual stream additions to steer models. Our work differs in that it
 522 considers shallow targeted paths through networks, rather than deep hidden states in networks.

523 Q Succession dataset

524 We present the full succession dataset in Table 2. Note that the days and months tasks are special
 525 as the final tokens in these classes (‘Sunday’ and ‘December’) have cyclical successors (‘Monday’
 526 and ‘January’); we don’t consider the end tokens of the other tasks to have cyclical successors. Full
 527 details of our dataset can be found in our open-sourced experiments.⁶

Task	Tokens
Numbers	‘1’, ‘2’, ..., ‘199’, ‘200’
Number words	‘one’, ‘two’, ..., ‘nineteen’, ‘twenty’
Cardinal words	‘first’, ‘second’, ..., ‘tenth’
Days	‘Monday’, ‘Tuesday’, ..., ‘Sunday’
Day prefixes	‘Mon’, ‘Tue’, ..., ‘Sun’
Months	‘January’, ‘February’, ..., ‘December’
Month prefixes	‘Jan’, ‘Feb’, ..., ‘Dec’
Letters	‘A’, ‘B’, ..., ‘Z’

Table 2: Tokens in the succession dataset

528 Glossary

529 W_{OV} The $\mathbb{R}^{d_{\text{model}}} \times \mathbb{R}^{d_{\text{model}}}$ matrix for an attention head that is the product of its W_O and W_V
 530 matrices. See Elhage et al. [9] for motivation.

531 **Direct effect** involves identifying the output attributed to the head irrespective of the behavior of
 532 other heads. This is different from the indirect effect, where the effects of ablating a head
 533 are hidden by a backup head..

534 **Feature** is a term we use to refer to an interpretable (linear) direction in activation space, inspired by
 535 definition 2 from Elhage et al. [1]. Since SAEs can be viewed as learning a set of directions
 536 which their W_{in} matrix reads in (i.e. their neurons pre-ReLU activations are dot products
 537 between these directions and the SAE input), we similarly refer to these directions as SAE
 538 features. SAEs have been found to be interpretable in prior and our paper (Section 3).

539 **Linear probe** is a learnable $\mathbb{R}^{10} \times \mathbb{R}^{d_{\text{model}}}$ matrix that acts on MLP0’s output space to predict the
 540 mod-10 value of the underlying numeric token.

541 **Mean ablation** involves patching the output of a head with its mean output over a chosen distribu-
 542 tion.

⁶We will release our code upon successful publication.

- 543 **Numbered listing** is a special case of incrementation where the token being incremented designates
544 items in a list, such as ‘A’ and ‘B’ in ‘A) ... B)’ or ‘i’ and ‘ii’ in ‘[i] ... [ii]’. See Figure 20
545 for in-the-wild examples.
- 546 **Ordinal sequence** refers to a series of words arranged in a specific, meaningful order, where the
547 position of each item is important.
- 548 **Patching** replaces part of a model’s forward pass with activations from a different input [6, 30] .
- 549 **SAE** stands for Sparse Autoencoder [17], see Appendix C.