TRANSFORMERS USE CAUSAL WORLD MODELS IN
 MAZE-SOLVING TASKS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recent studies in interpretability have explored the inner workings of transformer models trained on tasks across various domains, often discovering that these networks naturally develop surprisingly structured representations. When such representations comprehensively reflect the task domain's structure, they are commonly referred to as "World Models" (WMs). In this work, we discover such WMs in transformers trained on maze tasks. In particular, by employing Sparse Autoencoders (SAEs) and analysing attention patterns, we examine the construction of WMs and demonstrate consistency between the circuit analysis and the SAE feature-based analysis. We intervene upon the isolated features to confirm their causal role and, in doing so, find asymmetries between certain types of interventions. Surprisingly, we find that models are able to reason with respect to a greater number of active features than they see during training, even if attempting to specify these in the input token sequence would lead the model to fail. Futhermore, we observe that varying positional encodings can alter how WMs are encoded in a model's residual stream. By analyzing the causal role of these WMs in a toy domain we hope to make progress toward an understanding of emergent structure in the representations acquired by Transformers, leading to the development of more interpretable and controllable AI systems.

#### 027 028 029

030

004

010 011

012

013

014

015

016

017

018

019

020

021

022

024

025

026

#### 1 INTRODUCTION

The study of world models (WMs) in AI systems has gained significant traction of late, yet much interpretability research focuses on large language models trained on diverse, complex datasets (Belrose et al., 2023; Lieberum et al., 2023; Olsson et al., 2022). In an attempt to seek a more comprehensive understanding, our work examines WMs acquired by Transformers (Vaswani, 2017) in a controlled, synthetic environment. In particular, we use Maze-solving tasks (Subsection 2.1) as an ideal testbed for understanding learned WMs due to their human-understandable structure, controllable complexity, and relevance to spatial reasoning. Using this constrained domain, we can rigorously analyze how transformers trained (Subsection 2.2) to solve mazes construct and utilize internal representations of their environment.

Our methodology leverages Sparse Autoencoders (SAEs) (Bricken et al., 2023) to overcome the limitations of linear probes in detecting WM features. While linear probes can and have been used to identify latent directions associated with features defined in an imposed ontology, SAEs are found to discover features actually used by our models to make decisions (Section 3). By manipulating specific features identified by the SAEs and observing the impact on our models' maze-solving behavior, we provide strong evidence that these features are causally involved in the model's decisionmaking process (Section 4). This stands in contrast to prior work analyzing the formation of WMs in maze settings where no causal features were able to be isolated (Ivanitskiy et al., 2024).

Our findings provide important considerations for AI interpretability and alignment. By investigat ing how transformers form causal WMs even in relatively simple tasks, we hope to provide new
 avenues for understanding representations and potentially steering behavior in more complex AI
 systems. This work lays the groundwork for future research into how we might intervene on WMs
 to better align transformer-based AI systems to desired constraints.



Figure 1: Overview of our methodology for discovering and validating world models in transformerbased maze solvers. (A) We analyze attention patterns in early layers, finding heads that consolidate maze connectivity information at semicolon tokens. (B) We train sparse autoencoders on the residual stream immediately following the first block, identifying interpretable features that encode maze connectivity. (C) We demonstrate the causal role of the world models in our transformers comparing the features extracted through both methods and validating them through causal interventions.

080

081

082

084

085

087

090

091

092

094

095 096

We outline our methodology in Figure 1. In short, we begin in Subsection 3.1, identifying attention heads that appear to construct world model features by examining their attention patterns across maze coordinate tokens (A). We validate these findings in Subsection 3.3 by training sparse autoencoders on the residual stream and demonstrating that the extracted features match those identified through attention analysis (B). Lastly, in Subsection 3.3 we establish the causal nature of these representations through targeted interventions, showing that perturbing specific features produces predictable changes in the model's maze-solving behavior (C).

## Contributions

- Empirical Findings: We show that transformers form WMs when solving mazes and that these WMs are causal: they can be intervened upon in the latent space of SAEs. Surprisingly, we find that interventions that activate features are more effective than those that remove them, suggesting an asymmetry in how transformers utilize WM features. In performing these interventions we also uncovered our models' abilities to reason in the presence of a out-of-distribution number of activated features than would naturally arise for a given token sequence length.
- **Methodological Insights:** By effectively utilizing decision trees to isolate WMs features in SAEs, we demonstrate that transformers utilizing different encodings schemes may use varyingly compositional codes to represent their WMs. More generally, our analyses suggest that SAEs are generally better suited than linear probes to isolate WMs, even in the absence of feature splitting.

## 2 PRELIMINARIES

097 2.1 Environment

098 Though it remains a matter of debate whether Large Language Models (LLMs) construct structured 099 internal models of the real world, we can begin to understand the representations acquired by such 100 models by focusing on "toy" tasks with clear spatial or temporal structure (Brinkmann et al., 2024; 101 Momennejad et al., 2024; Jenner et al., 2024; McGrath et al., 2022). Previous works along these lines 102 (Li et al., 2022; Ivanitskiy et al., 2024; Nanda, 2023; Karvonen, 2024; He et al., 2024) have found 103 a variety of both correlational and causal evidence for internal models of the environment within 104 trained transformers. In this work, we utilize maze-dataset (Ivanitskiy et al., 2023), a package 105 providing maze-solving tasks as well as ways of turning these tasks into text representations. In particular, we use a dataset of mazes consisting of up to  $7 \times 7$  grids, generated via constrained 106 randomized depth first search (see Subsection 2.2) (which produces mazes that are acyclic and thus 107 have a unique solution).

108 To train autoregressive transformers to solve such mazes, we employed a tokenization scheme pro-109 vided by maze-dataset, shown in Figure 2. This scheme is designed to present the maze struc-110 ture, start and end points, and solution path in a format amenable to transformer processing whilst 111 remaining straight-forward to analyse with standard tools from the mechanistic interpretability liter-112 ature - primarily due to the existence of a unique token for every position in the maze (aka "lattice").

| 14 | <adjli< th=""><th>ST_S</th><th>TART&gt;</th><th>(0,0)</th><th>&lt;&gt;</th><th>(1,</th><th>0)</th><th>;</th><th>(2,0)</th><th>&lt;</th><th>-&gt;</th><th>(3,0</th><th>)</th><th>;</th><th>(4,1)</th><th>&lt;&gt;</th></adjli<>                                    | ST_S  | TART> | (0,0) | <>  | (1,  | 0)   | ;  | (2,0) | <      | ->   | (3,0 | )    | ;  | (4,1)                                    | <>      |
|----|---|-------|-------|-------|---|------|------|----|-------|--------|------|------|------|----|--|---------|
| 15 | (4,0)   | ;     | (2,0) | <>    | (2,1)   | ;    | (1,  | 0) | <>    | (1,    | 1)   | ;    | (3,4 | )  | <>                                       | (2,4)   |
| 16 | ; (3  | ,1)   | <>    | (3,2) | ;   |      | •    | (  | 1,3)  | <>     | (    | 1,4) | ;    | <  | ADJLI                                    | ST_END> |
| 17 | <origi< th=""><th>IN_SI</th><th>ART&gt;</th><th>(1,3)</th><th><or1< th=""><th>GIN_</th><th>END&gt;</th><th></th><th>TARGE</th><th>T_STAI</th><th>RT&gt;</th><th>(2</th><th>3)</th><th></th><th><targ< th=""><th>ET_END&gt;</th></targ<></th></or1<></th></origi<> | IN_SI | ART>  | (1,3) | <or1< th=""><th>GIN_</th><th>END&gt;</th><th></th><th>TARGE</th><th>T_STAI</th><th>RT&gt;</th><th>(2</th><th>3)</th><th></th><th><targ< th=""><th>ET_END&gt;</th></targ<></th></or1<> | GIN_ | END> |    | TARGE | T_STAI | RT>  | (2   | 3)   |    | <targ< th=""><th>ET_END&gt;</th></targ<> | ET_END> |
| 10 | <pre><path_< pre=""></path_<></pre>   | .STAF | RT>   | (1,3) | (0,3)   | (0   | ,2)  | (  | 1,2)  | (2,2   | 2)   | (2,  | 1)   | (2 | ,0)                                      | (3,0)   |
| 20 | (4,0)   |       | 4,1)  | (4,2) | (4,3  | 3)   | (4,4 | )  | (3,4  | )      | (2,4 | )    | (2,3 | 3) | <pa< td=""><td>TH_END&gt;</td></pa<>     | TH_END> |



121 (a) An example of a tokenized maze. 1: The adjacency list describes the con-122 nectivity of the maze, with the semicolon token ; delimiting consecutive connections. The order of connections is randomized, ellipses represent omitted 123 connection pairs. 2,3: The origin and target specify where the path should be-124 gin and end, respectively. 4: The path itself a sequence of coordinate tokens 125 representing the shortest path from the origin to the target. For a "rollout," we 126 provide everything up to (and including) the <PATH\_START> token and au-127 toregressively sample with argmax until a **<PATH\_END>** token is produced.

(b) Visual representation of the same maze as in the tokenized representation on the left. The origin is indicated in green, the target in red, and the path in blue.

Figure 2: Tokenization scheme and visualization of a shortest-path maze task generated using Ivanitskiy et al. (2023).

130 131 132

133

146 147

148 149

150

151

152

153

154

128

129

113

### 2.2 MAZE SOLVING TRANSFORMERS

134 Utilizing the tokenized representations of mazes provide by the maze-dataset library, a suite of 135 transformer models implemented using TransformerLens ((Nanda & Bloom, 2022)) were trained to 136 predict solution paths in acylic mazes. We performed extensive hyperparameter sweeps (Figure 10) 137 over several variants of the transformer architecture, yielding models with stronger generalization performance than those found by prior work Ivanitskiy et al. (2024). 138

139 To allow the testing of generalization to large maze size, the models were trained on  $5 \times 5$  fully-140 connected and  $6 \times 6$  sparsely connected mazes, embedded in a  $7 \times 7$  lattice. This ensured that all 141 coordinate tokens in the  $7 \times 7$  vocabulary had been seen during training time, such that generalization 142 to  $7 \times 7$  mazes was conceivable but out-of-distribution during inference. For our experiments, we investigated the two best performing models for each positional embedding (Su et al., 2024) scheme, 143 as shown in Table 1. Note that whilst these models had different numbers of heads, their parameter 144 counts varied only slightly - on account of Stan's use of learned positional embeddings. 145

#### 3 DISCOVERING WORLD MODELS

Broadly speaking there are two ways to go about trying to identify internal world models: 1) Assuming the form of the world model and inspecting the transformer with e.g. supervised probes to see if this world model is present ((Nanda, 2023) SELF CITE Workshop proceeding), or 2) Exploring the model internals and investigating any structure which may be present in the representations to see if something akin to a world model exists. In our work we take both approaches.

| 155 | Model Nickname | Positional Embeddings | $d_{\text{model}}$ | $n_{\text{layers}}$ | n <sub>heads</sub> | Num. Params | Maze Solving Accuracy |
|-----|----------------|-----------------------|--------------------|---------------------|--------------------|-------------|-----------------------|
| 157 | Stan           | standard learned      | 512                | 6                   | 8                  | 19,225,660  | 96.6%                 |
| 158 | Terry          | ro <b>tary</b>        | 512                | 6                   | 4                  | 18,963,516  | 94.3%                 |

158 159

Table 1: Models chosen for mechanistic investigation (most performant in the sweep, given their 160 respective position embedding schemes). The number of parameters varies as Stan learns position 161 encodings  $(W_{pos} \in \mathbb{R}^{512 \times 512})$ 

First, in Subsection 3.1 we investigate attention heads in the earliest layer of our models and find heads specialising in the construction of representations akin to a world model. On the basis of this, we use SAEs (an unsupervised method) alongside supervised classifiers to identify latent features corresponding to the world model. Finally, we use patching experiments and interventions to show that both investigations yield consistent features, and that these form a causal world model with some interesting properties.

- 168
- 169
- 170 171

3.1 WORLD MODEL CONSTRUCTION: CONNECTIVITY ATTENTION HEADS

172 We began by examining the attention patterns of the maze-solving transformer models and uncovered a notable pattern: in both models, some or all of the attention heads at the first layer ("layer 173 174 0") appear to consolidate information about maze connections into the ; context positions. In particular, for all  $4i^{\text{th}}$  context-positions tokens (the semicolon separation tokens ; ), these heads 175 attend back 1 or 3 tokens - that is, to one of the two coordinate tokens corresponding to the given 176 connection preceding the ; token. This pattern is observable for 3/8 L0 heads in Stan (Figure 3) 177 and 4/4 L0 heads in Terry (Figure 14). This observation suggests the hypothesis that these heads 178 are in essence constructing a world model for the maze task, for use by later layers. 179

If this were the case, then we should expect that the output of these heads, mediated by the "OV-Circuit" (Elhage et al., 2021), should consist of combinations of the coordinates captured in a given connection. This can be measured by taking the  $W_{OV}$  matrix for each head and measuring the cosine distance between its elements and the model's token embeddings (where coordinate directions are directly given)<sup>1</sup>. With this in mind, we investigated the structure of these vectors more closely. We find an intriguing pattern in the magnitudes of these vectors in the Stan model (Figure 4), while the patterns in Terry were less clear cut (Figure 5).



Figure 3: Attention values for heads L0H3, L0H5, and L0H7 in Stan. We use a rather nonstandard 201 representation, looking only at a fixed window into the past of which tokens are attended to by 202 semicolon ; tokens. Every 4th position, up to 140, is shown along the x-axis. Color shows 203 attention to positions 1, 3, 5, and 7 earlier in the context (shown along the y-axis), for an example 204 6x6 maze input. This sort of pattern is typical across all inputs examined. Up until context position 205 100, the heads are attending 1 and 3 positions back; after this the pattern shifts to 5 and 7 back. Note 206 the complementary attention patterns of L0H3 and L0H7. Closer examination shows that L0H3 207 prefers to direct its attention to 'even-parity' maze cells, with L0H7 preferring 'odd-parity' cells. 208 L0H5 more frequently splits its attention between 1 and 3 back, but sometimes 'fills in' for L0H7. 209 The origins of this pattern are explored further in appendix E; note also the similarities to Figure 4. 210 The other five heads in L0 show no similar pattern. Full patterns are shown in Figure 13

- 211
- 212
- 213

<sup>214</sup> 215

<sup>&</sup>lt;sup>1</sup>As we analyze the first attention layer we can ignore potential "residual drift" in the representations of a given maze coordinate between early and later layers in our transformers (Belrose et al., 2023).



224

225

226

227

228

229

239

240

241

242 243 244

245

Figure 4: Magnitudes of vectors resulting from applying the  $W_{OV}$  matrices of heads L0H3, L0H5 and L0H7 of Stan to maze-cell token embeddings, projected onto the maze grid. The pattern here mirrors the way that the heads divide their attention between the 1-back and 3-back context positions (exemplified in Figure 3) with L0H3 focused on 'even-parity' cells, and L0H7 and LH05 focused primarily on 'odd-parity' cells. This pattern also recurs in the overlaps between query and key vectors of token embeddings, explored in detail in Appendix E.



Figure 5: Magnitudes of vectors resulting from applying the  $W_{OV}$  matrices of layer-0 heads of Terry to maze-cell token embeddings, projected onto the maze grid. The pattern here is much less striking than that for Stan (shown in Figure 4) although it does suggest that the heads specialise in even/odd-parity cells in localised regions of the maze.

## 3.2 WORLD MODEL REPRESENTATION: SPARSE AUTOENCODERS

246 As previous work (Ivanitskiy et al., 2024) struggled to intervene on WM features identified via lin-247 ear probing (Alain & Bengio, 2016), we trained Sparse Autoencoders to attempt to find disentangled features in our models (Cunningham et al., 2023; Bricken et al., 2023). Sparse Autoencoders are 248 motivated by the notion of superposition (Elhage et al., 2022) which posits that artificial neural 249 networks store more features than an orthogonal representation would allow. By training an autoen-250 coder with a higher-dimensional latent space than that of the transformer, tasked with reconstructing 251 a residual stream vector under a sparsity penalty, the hope is that the SAE will recover interpretable 252 features which the transformer was forced to superimpose. Similar approaches have previously seen 253 success on other toy tasks (He et al., 2024; Karvonen et al., 2024). 254

To prevent "neuron death" in the SAE latent space, resulting from high sparsity penalties, we apply the method of "Ghost Gradients" proposed by Jermyn & Templeton (2024). The resulting trained SAEs faithfully reconstructed the activations (in our case, the residual stream after L0), and replacing these activations with their SAE reconstructed counterparts did not affect model behaviour (Figure 8), giving confidence in the completeness of their representation.

Initial attempts to isolate SAE features corresponding to connections in the maze attempted to use differences in the features present in mazes with or without certain connections. This approach worked well in some cases, but not in others, as not all relevant features varied in magnitude by the same amount, and many features were co-active to a given connection (i.e. those implicated in the path representation, which itself might change when connections are added/removed). To address this, we instead trained decision trees to isolate the relevant features in our transformers (akin to Spies et al. (2022)), as shown in Figure 6.

This analysis yielded our first unexpected finding: Stan's WM consisted of two features for each connection - a somewhat generic "semicolon" feature, as well as a connection specific feature. We visualize highly activating examples of these features in Figure 17, and show that Stan's representation was stable for an additionally trained SAE in Figure 18.

We speculate that this "compositional code" arises in Stan as a result of the transformer imperfectly separating positional information from its WM. This representation also explains why previous efforts to intervene on models with learned positional encodings by using linear probes were unsuccessful - as intervening with a single direction yielded from supervised decoding would also affect the semicolon feature.

It is also interesting to note that Terry encoded connection information very cleanly into single features for each connection - i.e., a single direction in the residual stream. This is in-spite of the fact that Terry's attention heads appeared to operate in a more entangled fashion than those of Stan.



Figure 6: Decision tree decoding accuracies and relevant features (in parentheses) for each connection in the maze. Upper right triangles correspond to right connections, and lower left triangles correspond to down connections. The decision trees were trained to predict the presence, or absence, of a connection from the SAE feature vector at the semicolon immediately following the definition of that connection. See Figure 17 for more details.

299 300 301

302

279

280 281

283

284

287

289

290 291

292

293

295

296

297

298

### 3.3 COMPARING SAES AND CIRCUITS

In Subsection 3.1 we advanced the claim that certain L0 heads construct features representing maze edges at the ; context positions, specifically by attending to earlier positions containing maze-cell token embeddings, and rewriting those embeddings by application of their  $W_{OV}$  matrices. Subsection 3.2 identified features representing maze edges via an independent line of reasoning, by training SAEs, and identifying which of their features were indicative of the presence of a maze edge.

To verify whether these approaches yielded consistent features for the WM, we first calculated the cosine similarity between the features written by isolated attention heads, and those encoded in the SAE (Figure 7a). These showed excellent agreement for Stan, where the attention patterns were clear, but only once the compositional code was taken into account (see Appendix G for details).

312 Though these results were promising, we carried out a further comparison (Figure 7b) to minimize 313 the assumptions required, and to account for two potential effects: 1) There may be "wiggle room" 314 between feature directions in the model's residual stream, and the circuits that construct them (which would lead to low cosine similarities, even for the same features), 2) As our SAEs are trained after 315 an entire block of computation, it is possible that the MLPs, applied after attention, also played a 316 role in forming the representations. In this second experiment we patched attention head values in 317 the presence of a connection to the mean of a maze set without that connection.By looking at the 318 effect of patching the attention heads on the resulting SAE Latent vectors, we were able to observe 319 that the features considered relevant for any given connection were indeed sensitive to the heads 320 implicated in constructing those features. 321

In particular, we consider the effect on the SAE features identified in Subsection 3.2 when each attention head is patched at the semicolon position for with its average non-connection value across 500 examples (i.e. removing the contribution a given head toward encoding that connection). 324 This captures the extent to which a given head contributes towards the "creation" of a maze-325 connection's representation in the residual stream. These plots not only confirm the link between 326 the attention circuits and the SAE features, but even show the same spatial partitioning of different 327 parts of the maze between different heads. The same plots for Terry, and Stan's down connection, 328 are shown in Appendix F.



345 (a) Cosine similarities between edge features derived 346 from the SAE and edge features derived by direct application of the  $W_{OV}$  matrices of the heads discussed 347 in Subsection 3.1 to token embeddings for Stan (for 348 reasons of space, only the right directed edge features 349 are shown here, but the pattern is the same for the full 350 set). Details of how features are computed are given in 351 Appendix G.

(b) Effect of patching attention heads on SAE features (rightward connections) for Stan's connection-specific WM features (identified from Figure 6). Compared against Figure 4, we see agreement between the attention analysis and the SAE Feature analysis. We normalize the per-head patching effect magnitudes, such that 100% was the maximal effect seen on an SAE feature's magnitude as a result of patching the attention head (across all features for a given head).



#### 4 INTERVENING ON WORLD MODELS

359 Though a universally agreed upon definition does not exist, we shall consider world Models to be "structure preserving [...] causally efficacious representations" (Millière & Buckner, 2024) of an 360 environment; i.e. representations which preserve the causal structure of the environment as far as is necessitated by the tasks an agent needs to perform. As such, we are interested in understanding 362 how the WMs we have discovered are leveraged by our models to facilitate generation of valid solution sequences. In Figure 8, we give an example of perturbing a feature to "fool" the model into behaving as though it is in a different maze. When patching in the SAE-reconstructed residual stream without perturbations we still see the same behavior as in the original model; when patching 366 in with a modified feature, we see a change in the path. We perform such interventions across 200 examples for each connection feature, and show the resulting intervention efficacies in Figure 9. 368

The intervention process involves toggling a feature on (to the maximal value observed for that 369 feature in a small dataset) or turning it off (setting it to 0) at all semicolon positions<sup>2</sup>. We measure 370 the impact of these interventions on the model's maze-solving accuracy, with a particular focus 371 on how activating versus removing features affects performance. Our results reveal an intriguing 372 asymmetry that constitutes our second finding: interventions that activate features tend to be more 373 effective in altering the model's behavior compared to those that remove features.

374

352 353

354 355 356

357 358

361

363

364

367

<sup>375</sup>  $^{2}$ For the case of adding a connection, this is necessary as there is no semicolon in the sequence which 376 "belongs" to the connection that doesn't exist. We also experimented with toggling to a fixed maximal value in 377 Figure 19, but this was generally less effective. In the case of removal, it made little difference if the feature was disabled everywhere, as it is almost always exactly 0 for a non-matching connection semicolon

This suggests that the transformer may rely more heavily on the presence of certain connectivity cues rather than their absence when constructing its internal world model.

Our final finding relates to the toggling of features in Stan. Though Stan utilized a compositional 381 code, activating the connection-specific features at unrelated semicolons worked in 35% of cases. 382 Conversely, we saw that all removal interventions failed for Stan, for the simple reason that Stan 383 was unable to generalize to sequences containing more connections than it had seen during training 384 - thus failing when shown examples containing the additional connection to be removed (this failure 385 was a result of Stan using learned positional embeddings (Table 1), as shown in Figure 10). The fact 386 that activating connections in the space of the SAE worked at all means that Stan's maze-solving 387 behaviour was at least partially able to generalize in the latent space, where it was decoupled from 388 the positional embeddings.



Figure 8: An example of an intervention on Terry where a connection is added by enabling the relevant feature in the SAE's latent space (in this case, feature 250 for (1, 2) < --> (1, 3)). From left to right: 1) input maze with ground truth 2) model's prediction with the unperturbed SAE reconstruction patched in 3) perturbed ground truth 4) model's prediction with the perturbed SAE reconstruction in its residual stream at layer 0.



Figure 9: Aggregated accuracy of interventions for examples on which the original prediction was correct. An accurate intervention is one in which the toggling of a connection in the SAE feature space leads the model to act accordingly. Note that Stan removal interventions fail as the inputs in these cases have more connections than the model is able to handle (see length generalization failure in Figure 10).

425 426

389

399

400

401

402

403

## 5 RELATED WORK

427 428

Our work builds on existing literature in interpretability (Räuker et al., 2023), particularly how trans formers develop structured internal representations, often called world models. World Models, as
 defined by Millière & Buckner (2024), are "structure-preserving, causally efficacious representa tions of properties of [a model's] input domain."

Here, structure-preserving means that the representations reflect the causal structure of the observa tion space and causally efficacious means that the model leverages these representations to enable
 relevant interactions with its environment.

435 Research into world models has gained traction across various domains, with transformers trained 436 to play complex games like chess being prime examples. For instance, McGrath et al. (2022) trained 437 linear probes to extract various features in AlphaZero's chess model, showing how different aspects 438 of the game, such as piece positioning and potential future moves, are captured within the model's 439 layers. Similarly, Karvonen (2024) investigates the internal representations of a chess model us-440 ing linear probes and contrastive activations, revealing structured representations of the game state. 441 Jenner et al. (2024) explores the emergence of learned look-ahead capabilities in Leela Chess Zero, 442 where the model encodes an internal representation of future optimal moves.

Another task used to study internal representations in transformers is Othello. Several works have explored the emergence of causal linear world models in this domain Li et al. (2022); Nanda (2023), with recent advancements leveraging SAEs (see Subsection 3.2 to uncover these world models He et al. (2024).

Beyond game-playing tasks, the study of learned world models in transformers extends to other domains, such as natural language processing, where Hewitt & Manning (2019) used probing techniques to uncover the syntactic structure encoded by BERT. This line of research demonstrates that transformer models can implicitly learn hierarchical structures in their residual streams, as explored by Manning et al. (2020). Further supporting this, Pal et al. (2023) demonstrated that the residual stream corresponding to individual input tokens encodes information to predict the correct token several positions ahead, highlighting the model's capacity for structured, anticipatory reasoning.

Additionally, graph traversal as multi-step reasoning has been investigated both from a model capabilities perspective Momennejad et al. (2024) and through mechanistic interpretability Brinkmann et al. (2024); Ivanitskiy et al. (2024), providing further evidence of transformers' ability to encode and utilize structured representations in complex tasks.

458 459 460

461

## 6 CONCLUSIONS AND FUTURE WORK

In this work, we demonstrated that transformers trained to solve maze navigation tasks form highly 462 structured internal representations that capture the connectivity of the maze and thus act as world 463 models. Through exploratory analysis of attention patterns, we found that connection information 464 was consolidated into semicolon tokens by a subset of attention heads. By using Decision Trees 465 to analyze the latent space of Sparse Autoencoders on these semicolons, we were able to identify 466 sparse features that encoded the position in the maze. We showed that these world models were con-467 structed differently in transformers leveraging learned vs. rotary positional encodings, suggesting 468 that simpler methods such as activation steering or probing would have been insufficient to extract 469 causal world models in at least some cases. More interesting still, we showed that interventions to add connections by toggling features were consistently more effective than interventions that sought 470 to remove connections by zeroing the corresponding features. Furthermore, we found that models 471 with learned position encodings, which were unable to generalize to longer input sequences (i.e., 472 mazes with more connections), were able to behave consistently if additional connection features 473 were enabled via SAE interventions, even if the corresponding token sequence would have caused 474 the model to fail. 475

These findings shed light on the inner workings of transformers trained on sequential planning tasks and suggest that maze-solving tasks are a rich testbed for understanding the formation of world models in transformers. Future work should aim to uncover whether our findings on intervention asymmetries and steerability are universal - and if not, which conditions give rise to each. An empirical understanding of the reliability of SAE feature discovery and steerability is crucial for AI Safety efforts that attempt to constrain or coerce model behavior through interventions or monitoring based on such methods.

- 483
- 484
- 485

# 486 REFERENCES

492

502

518

- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier
   probes. *arXiv preprint arXiv:1610.01644*, 2016.
- 490 Nora Belrose, Zach Furman, Logan Smith, et al. Eliciting Latent Predictions from Transformers
   491 with the Tuned Lens. *arXiv preprint arXiv:2303.08112*, 2023.
- Trenton Bricken, Adly Templeton, Brian Chen, et al. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, Oct 4 2023.
   URL https://transformer-circuits.pub/2023/monosemantic-features/index.html.
- Jannik Brinkmann, Abhay Sheshadri, Victor Levoso, et al. A mechanistic analysis of a transformer
   trained on a symbolic multi-step reasoning task. *arXiv preprint arXiv:2402.11917*, 2024.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs, et al. Sparse Autoencoders Find Highly Interpretable Features in Language Models, October 2023. URL http://arxiv.org/abs/2309.08600.
- Nelson Elhage, Neel Nanda, Catherine Olsson, et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformercircuits.pub/2021/framework/index.html.
- Nelson Elhage, Tristan Hume, Catherine Olsson, et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.
- Zhengfu He, Xuyang Ge, Qiong Tang, et al. Dictionary learning improves patch-free circuit discovery in mechanistic interpretability: A case study on othello-gpt. *arXiv preprint arXiv:2402.12201*, 2024.
- John Hewitt and Christopher D Manning. A structural probe for finding syntax in word representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4129–4138, 2019.
- Michael I. Ivanitskiy, Rusheb Shah, Alex F Spies, et al. A configurable library for generating and
   manipulating maze datasets. *arXiv preprint arXiv:2309.10498*, 2023.
- Michael I. Ivanitskiy, Alex F. Spies, Tilman Räuker, et al. Linearly structured world representations in maze-solving transformers. In *Proceedings of UniReps: the First Workshop on Unifying Representations in Neural Models*, volume 243 of *Proceedings of Machine Learning Research*, pp. 133–143. PMLR, 15 Dec 2024. URL https://proceedings.mlr.press/v243/ivanitskiy24a.html.
- Erik Jenner, Shreyas Kapur, Vasil Georgiev, et al. Evidence of learned look-ahead in a chess-playing
   neural network. *arXiv preprint arXiv:2406.00877*, 2024.
- Adam Jermyn and Adly Templeton. Anthropic circuits Updates January 2024, January 2024. URL https://transformer-circuits.pub/2024/jan-update#
   dict-learning-resampling.
- Adam Karvonen. Emergent world models and latent variable estimation in chess-playing language
   models. arXiv preprint arXiv:2403.15498, 2024.
- Adam Karvonen, Benjamin Wright, Can Rager, et al. Measuring progress in dictionary learning
   for language model interpretability with board game models. *arXiv preprint arXiv:2408.00113*,
   2024.
- Kenneth Li, Aspen K Hopkins, David Bau, et al. Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint arXiv:2210.13382*, 2022.
- Tom Lieberum, Matthew Rahtz, János Kramár, et al. Does Circuit Analysis Interpretability Scale?
   Evidence from Multiple Choice Capabilities in Chinchilla. *arXiv preprint arXiv:2307.09458*, 2023.

| 540<br>541<br>542 | Christopher D Manning, Kevin Clark, John Hewitt, et al. Emergent linguistic structure in artificial neural networks trained by self-supervision. <i>Proceedings of the National Academy of Sciences</i> , 117(48):30046–30054, 2020.   |
|-------------------|--|
| 543<br>544<br>545 | Callum McDougall. SAE Visualizer. https://github.com/callummcdougall/sae_vis, 2024.  |
| 546<br>547        | Thomas McGrath, Andrei Kapishnikov, Nenad Tomašev, et al. Acquisition of chess knowledge in alphazero. <i>Proceedings of the National Academy of Sciences</i> , 119(47):e2206625119, 2022.   |
| 549<br>550        | Raphaël Millière and Cameron Buckner. A philosophical introduction to language models - part ii:<br>The way forward, 2024. URL https://arxiv.org/abs/2405.03207.   |
| 551<br>552<br>553 | Ida Momennejad, Hosein Hasanbeig, Felipe Vieira Frujeri, et al. Evaluating cognitive maps and planning in large language models with cogeval. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.   |
| 554<br>555<br>556 | Neel Nanda. Actually, othello-gpt has a linear emergent world representation. <i>Neel Nanda's Blog</i> , 7, 2023.  |
| 557<br>558        | Neel Nanda and Joseph Bloom. Transformerlens. https://github.com/<br>TransformerLensOrg/TransformerLens, 2022.   |
| 559<br>560<br>561 | Catherine Olsson, Nelson Elhage, Neel Nanda, et al. In-context Learning and Induction Hea ds.<br>arXiv preprint arXiv:2209.11895, 2022.  |
| 562<br>563        | Koyena Pal, Jiuding Sun, Andrew Yuan, et al. Future lens: Anticipating subsequent tokens from a single hidden state. <i>arXiv preprint arXiv:2311.04897</i> , 2023.  |
| 565<br>566<br>567 | Tilman Räuker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. Toward Transparent AI: A Survey on Interpreting the Inner Structures of Deep Neural Networks. In <i>2023 ieee conference on secure and trustworthy machine learning (satml)</i> , pp. 464–483. IEEE, 2023. |
| 568<br>569        | Alex F. Spies, Alessandra Russo, and Murray Shanahan. Sparse Relational Reasoning with Object-Centric Representations, July 2022. URL http://arxiv.org/abs/2207.07512.   |
| 570<br>571<br>572 | Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. <i>Neurocomputing</i> , 568:127063, 2024.  |
| 573<br>574        | A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.   |
| 575               |  |
| 576               |  |
| 577               |  |
| 578               |  |
| 579               |  |
| 580               |  |
| 581               |  |
| 582               |  |
| 583               |  |
| 584               |  |
| 585               |  |
| 586               |  |
| 587               |  |
| 588               |  |
| 589               |  |
| 590               |  |
| 591<br>500        |  |
| 592               |  |
| 100               |  |



Figure 10: Accuracies of all transformers trained in our sweep on a generalization task. "Train Val" shows the accuracy on the held out in-length-distribution mazes from train time, and "Full Maze" features mazes with more connections (longer input sequences) than those seen at train time. Only rotary models are able to generalize at all

## **B** SAE TRAINING DETAILS

To choose optimal hyperparameters for our SAEs we ran a sweep over SAEs at layers 2 to 4 on Terry, finding consistent trends across layers. The results of this sweep are shown in Figure 11, and the final details of the SAE analyzed in the main paper are given in Table 2. We also provide feature density histograms for the SAEs analyzed in the main paper in Figure 12 noting that these look good, in that many features are sparse, but also rather distinct from is typically observed in LLMs. This is not surprising, as our token and features distributions will be very distinct from those of natural language, as most mazes have many active connections, and connections are similarly likely to be present in any given maze.

|                  | Sparsity        |                      | D          | ataset         | 0             | ptimizer             |
|------------------|-----------------|----------------------|------------|----------------|---------------|----------------------|
| Expansion Factor | Ghost Threshold | (L0) Sparsity Weight | Batch Size | Training Steps | Learning Rate | Linear Warm Up Steps |
| 4                | 100             | 0.01                 | 1024       | $\sim 10^{6}$  | $10^{-4}$     | 1000                 |

Table 2: Hyperparameter values for the final SAEs analyzed in the main paper.

|       |                                    | SAE Feature Metrics |      | Average Token Reconstruction Errors |              |             |  |
|-------|------------------------------------|---------------------|------|-------------------------------------|--------------|-------------|--|
|       | Residual Reconstruction Error (L2) | Sparsity (L1)       | L0   | Unperturbed                         | Zero Patched | SAE Patched |  |
| Terry | $2.87 \times 10^{-4}$              | 10.7                | 20.9 | 8.18                                | 8.52         | 8.18        |  |
| Stan  | $6.35 \times 10^{-4}$              | 9.53                | 28.4 | 6.35                                | 8.61         | 6.35        |  |

Table 3: SAE Metrics for the final SAEs trained on Stan and Terry. We see that replacing the residual stream with the SAE reconstructions has very little impact on the sequence produced by the model, providing confidence that the SAEs are encoding all the relevant information in the model's residual stream.



Figure 11: Results of an SAE sweep carried out on Terry.



Figure 13: Attention patterns for head L0H3 in Stan and Terry, for a specific example maze. At every 729 fourth context position from 4 through to 140 (the ; positions in the adjacency-list) attention is 730 directed very strongly back to one or two positions, typically 1 or 3 positions earlier in the context 731 (though for Stan, after context position 100, this shifts to 5 or 7 positions earlier in the context). This 732 pattern is qualitatively repeated across all examples examined, for heads L0H3, L0H5 and L0H7 in 733 Stan, and for all four L0 heads in Terry. 734

#### FURTHER ATTENTION VISUALIZATIONS С



753

735

754 Figure 14: Attention values for layer 0 heads in Terry, from context positions holding the ; token 755 (shown along the x-axis) to positions 1 and 3 earlier in the context (shown along the y-axis), for an example maze input. This pattern is typical across all inputs examined. The pattern is less clear-cut than for Stan (Figure 3), but note that at every fourth context position, there is at least one head attending strongly to positions 1 and 3 earlier in the context.

## D HOW SAE REPRESENTATIONS DIFFER



(a) Terry model: A single feature almost perfectly
encodes the existence of a specific connection in the
maze. This demonstrates the direct encoding of maze
connectivity in Terry's SAE latent space.

(b) Stan model: Two features (Feature 1422 and another) work together to encode maze connectivity. Feature 1422 appears consistently across all connections, aligning with the decision tree decoding results presented Figure 6.

Figure 15: Decision trees trained on SAE latents for Terry and Stan models, predicting the existence of specific connections in the maze. These examples illustrate how maze connectivity is encoded in the residual stream at layer 0 on the corresponding semicolon position. The decision trees were trained as supervised classifiers whose target was to predict the presence of a given connection, given an SAE feature vector from the corresponding semicolon position. These SAEs were trained with 10,000 examples per connection (equally balanced between the presence / non-presence of a connection).

| 779 |   |  |   |
|-----|---|--|---|
| 780 | TOP ACTIVATIONS<br>MAX = 4.665  | TOP ACTIVATIONS<br>MAX = 12.012              | TOP ACTIVATIONS   |
| 781 | (5,3); $(4,4) <> (4,5)$ ; $(3,0) <> (3,1)$ ; $(4,1)$                                | (4,2); (4,4) <> (4,5); (4,2) <> (4,3); (3,3) | MAX = 1.518   |
| 782 | (5,1); <mark>(4,4) &lt;&gt; (4,5)</mark> ; (3,2) <mark>&lt;&gt;</mark> (3,3); (2,5) | (5,2); (4,4) <> (4,5); (1,2) <> (0,2); (3,5) | (5,1); $(4,5) <> (3,5)$ ; $(4,1) <> (4,0)$ ; $(4,3)$  |
| 792 | (2,1); (4,4) <> (4,5); (5,4) <> (5,3); (2,2)  | (1,5);(4,5)<>(4,4);(2,4)<>(3,4);(0,5)        | (5,1); (1,3) <> (2,3); (3,3) <> (3,4); (4,1)  |
| 105 | (0,1); <mark>(4,4)&lt;&gt;(4,5);</mark> (5,2) <mark>&lt;&gt;</mark> (5,1);(1,2)     | (2,4);(4,4)<>(4,5);(2,0)<>(3,0);(3,5)        | (3,3) <mark>;</mark> (5,1)<>(5,2) <mark>;</mark> <adjlist_end><origin_start></origin_start></adjlist_end> |
| 784 | (4,4); $(4,4) <> (4,5)$ ; $(3,3) <> (4,3)$ ; $(2,0)$                                | (0,5); (4,5) <> (4,4); (1,5) <> (0,5); (4,3) | <pre>(1,1);(5,3)&lt;&gt;(5,2);<adjlist_end><origin_start></origin_start></adjlist_end></pre>              |
| 785 | (a) Terry model: A single SAE   | (b) Stan model: The connection-              | (c) Stan model: Feature 1422,   |
| 786 | feature directly encodes a spe-   | specific feature activates at the            | in conjunction with another fea-  |
| 787 | cific maze connection, demon-   | semicolon corresponding to the               | ture, encodes maze connectivity.  |
| 788 | strating Terry's straightforward  | encoded connection, similar to               | This feature appears consistently   |
| 789 | representation of maze connec-  | Terry's encoding strategy (see               | across all connections, corrobo-  |
| 790 | tıvıty.   | Figure 17b).                                 | rating the decision tree decoding results in Figure 6.  |

Figure 16: Maximally activating examples, displayed using a modified version of McDougall (2024)
for SAE features encoding the connection (4, 4) <--> (4, 5), as identified by decision
tree decoding. Underlines correspond to loss contribution (blue for positive, red for negative) and
highlighting indicates feature activation at a given token position. Connection-specific features in
both models (Figure 17b and Figure 17a) show clear activation patterns, while Stan's generic semicolon feature (Figure 16c) exhibits a less obvious trend. Produced using a modified version of
McDougall (2024)

799 800

801

756

771

#### D.1 MAGNITUDE OF INTERVENTIONS

To complement the intervention results presented in the main text, we also conducted fixed-value interventions on both the Stan and Terry models. In these interventions, instead of calculating new activations based on the modified input, we directly set the activations of the targeted features to fixed values. This approach allows us to examine how the models respond to more controlled manipulations of their internal representations.

The fixed-value intervention results shown in Figure 19 reveal interesting patterns that both complement and contrast with the calculated intervention results presented in the main text.



(b) Representative SAE features for Stan.

Figure 17: We provide examples for the types of features observed in Stan and Terry, beyond the connection features which form the primary focus of the main paper. We observe the same kinds of features between both transformers, and in both cases the predominant features are of the form observed in the top-left (Feature 32 in Terry and 2 in Stan) - These features are more distributed and harder to interpret than the others, and may be suppressed by higher sparsity penalties.

850 851 852

853

845 846

847

848

849

## E INVESTIGATION OF QK-CIRCUIT IN STAN MODEL

854 In an effort to better understand the notable "1- and 3-back" attention patterns appearing in heads 855 L0H3, L0H5 and L0H7 of Stan, described in Subsection 3.1, we investigated the query and key 856 vectors for token and positional embeddings, and their overlaps. The scalar products between queries 857 and keys of token embeddings for L0H3 are shown in figure 20. The most striking feature of this 858 plot is the row corresponding to the query vector of the ; token, and in particular its overlap 859 with the maze cell tokens. Plotting these scalar products on the maze cell grid (figure 21) a clear pattern emerges, analogous to that shown in figure 4, accounting for LH03's tendency to attend to 860 even-parity cells, and LH05's and LH07's tendencies to attend to odd-parity cells. Examining the 861 scalar products among query and key vectors for positional embeddings (figure 22) reveals a pattern 862 that likely accounts for the focusing of attention from ; context positions to positions 1 and/or 3 863 earlier in the context.



Figure 19: Aggregated accuracy of fixed-value interventions for examples on which the original prediction was correct. As opposed to Figure 9, the addition interventions were performed with a fixed value of 10 (removal interventions were the same, with a fixed value of 0). Here we see that the fixed-value interventions are mostly less effective than the calculated interventions, suggesting magnitude sensitivity for feature magnitudes in the transformer's use of the World Model.

- 914
- 915
- 916
- 917



Figure 20: Scalar products of Stan LH03 of query (rows) and key (columns) vectors for token embeddings. Note that the most pronounced pattern is found on the row corresponding to the query vector of the ; token, reflecting the importance of this head in establishing the attention pattern from context positions containing the ; token.



Figure 21: Stan scalar products of query vector for ; token and key vectors for maze-cell tokens, arranged on maze grid. Note the clear correspondence with Figure 4. These patterns account for why LH03 directs its attention to even-parity cells, while odd-parity cells are attended to by LH07 or LH05.









(a) Effect of attention patching on right-connection
 (b) Effect of attention patching on down-connection
 features.
 features

Figure 26: Effect of patching attention heads on SAE features for Terry. Whilst we observe notable effects, it is difficult to see a clear pattern - as revealed by the attention analyses, the role of each head in constructing a single connection feature in Terry is harder to understand.

1131

1132

## <sup>1134</sup> G COMPUTING SAE AND OV EDGE FEATURE SIMILARITY

In Figure 7a we compute the cosine similarity between SAE edge features and OV circuit edge features.

SAE edge features are formed from a linear combination of the specific edge feature and a "generic edge" feature, with the generic feature coefficient of -0.6 being chosen to maximise cosine similarity.

1142 OV edge features are formed from a weighted sum:

 $\sum_{h,c} a_c^h W_{OV}^h t_c$ 

1146 Here, *h* indexes heads L0H3, L0H5 and L0H7, with  $W_{OV}^{L0H3}$ , for example, giving the OV matrix of 1147 L0H3. *c* indexes the two cells present in the edge of interest, and  $t_c$  is the token embedding of a cell 1148 *c*. The coefficients  $a_c^h$  are given by the attention directed by head *h* to cell *c* from the ; context 1149 position following the specification of the edge of interest. Data was averaged averaged over 100 1150 examples (see Figure 3 for one such example).