
Entity Multiplexing Through Activation Strength: Understanding goals in A Maze Solving Agent

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work we provide an extensive analysis into the operations of a maze solving
2 reinforcement learning agent trained in the Procgen Heist environment. We target
3 this model because it presented a high degree of polysemanticity due to the fact
4 that it has to target multiple different entities to succeed. By focusing on an
5 agent that has to target multiple similar entities we hope to answer questions
6 about how each of these entities might be processed by the network. Our main
7 finding is that the signals related to the targeting of different entities are encoded
8 at different activation strengths within a single channel in the network. These
9 "steering channels" are often highly redundant, with large numbers of channels
10 enabling precise agent steering, but often only within narrow ranges of activation
11 values. We also discover a paradoxical ablation effect in which removing both
12 steering channels and navigation circuits improves entity collection rates compared
13 to partial ablation, suggesting unexpected interference between these systems.
14 These findings demonstrate that amplitude-based multiplexing is a fundamental
15 strategy for encoding multiple goals in RL agents, while our counterintuitive
16 ablation studies suggest surprising specialization and informational dependencies
17 within the network.

18 1 Introduction

19 Understanding how deep learning models perform their tasks is currently an unsolved problem in the
20 field of AI. This state of affairs ensures that highly reliable, controllable, and understandable models
21 remain out of reach. This is particularly the case in the field of reinforcement learning (RL) that has
22 been somewhat neglected by the techniques of mechanistic interpretability (MI) whose focus has
23 largely been on Large Language Models (LLMs). With techniques from RL being applied to frontier
24 models to a greater extent, it is our belief that lessons learned from applying MI to RL based agents
25 can yield generalizable insights that can improve our understanding of AI agents and neural networks
26 in general.

27 This research builds on the work of Mini et al. [2023] in which precise control of a maze solving agent
28 was achieved by intervening directly on activations within the network. The current contribution
29 begins by extending their work to an environment that involves multiple competing entities that an
30 agent needs to reach in sequence, using the Procgen Heist environment as the target of our research.

31 The Procgen Heist environment requires the agent to collect up to 3 keys and 3 locks before reaching
32 the final goal, a gem. The order in which the keys need to be collected is always the same (blue,
33 green, red). The environment is procedurally generated and might generate with any combination
34 of no keys, or 1-3 keys. The size of the maze will also change depending on how many objects are
35 included.

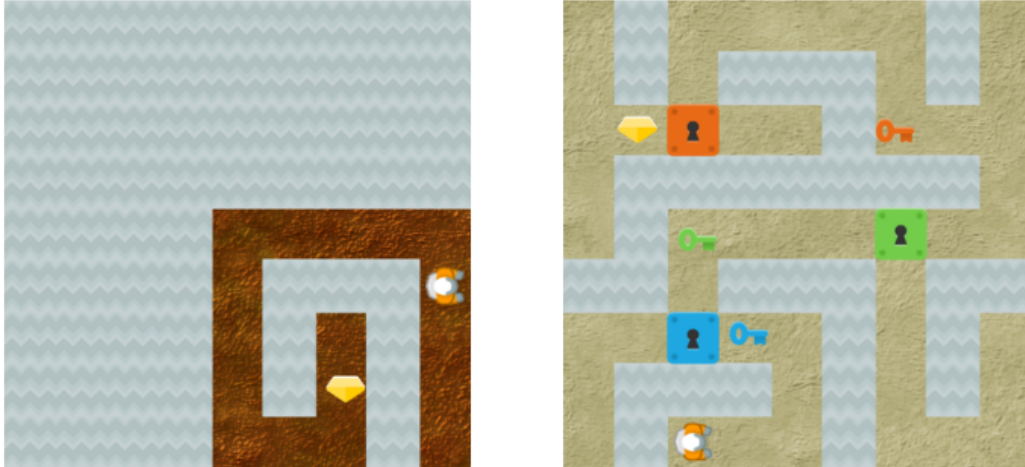


Figure 1: The Procgen Heist environment presents a procedurally generated maze where agents must collect keys and unlock corresponding doors in a specific sequence (blue, green, red) before reaching the final goal (gem). The environment varies in complexity, sometimes only having the gem present, and other times including multiple keys and locks.

36 The environment provides sparse rewards, only giving out a score of 10 if the agent makes it all the
 37 way to the final goal and 0 otherwise. The fact that the environment provides a variety of difficulties
 38 with some environments only having the gem and others including all 3 keys and locks establishes a
 39 natural curriculum for the agent. It also gives the agent a reasonable chance of reaching the end with
 40 a random walk before it starts to develop an actual policy.

41 Part of the reason we selected an environment from the Procgen suite for this task was that the Procgen
 42 environments naturally force the agent to learn sufficiently general representations of different entities
 43 to allow us to potentially extract them from the weights. If the model could simply encode easily
 44 memorized heuristics then this might lead to uninterpretable and highly specific rules rather than
 45 general representations of entities.

46 2 Related Work

47 2.1 Mechanistic Interpretability in RL

48 Previous work by Mini et al. [2023] demonstrated the ability to manipulate an agent’s navigation in a
 49 simple maze environment with an intervention to a single channel in the network. This serves as a
 50 foundation for our work, though we extend it to a more complex multi-objective setting. The idea of
 51 activation steering originated in RL and was later successfully applied to LLMs Turner et al. [2024],
 52 demonstrating the potential for the transfer of techniques between very different neural network
 53 architectures and domains.

54 Work exploring the Procgen Coinrun environment provides a rich set of ideas for how to attribute
 55 attention from objects in the input to specific weights within the model, as well as visualizing the
 56 model weights Hilton et al. [2020]. They showed clear positive and negative attributions to specific
 57 entities according to how it might affect its ability to complete the task within the environment.

58 Other interpretability approaches for RL include saliency-based methods Greydanus et al. [2018],
 59 Atrey et al. [2020] which identify important input pixels but do not enable behavioral manipulation.

60 2.2 Ablation Studies and Network Interpretability

61 Systematic ablation studies have been fundamental to understanding CNNs since Zeiler and Fergus
 62 [2013] pioneered their use to identify important model components. However, recent work has
 63 revealed that networks exhibit surprising robustness to ablations McGrath et al. [2023]. While early
 64 interpretability work searched for individual neurons encoding specific concepts, Morcos et al. [2018]

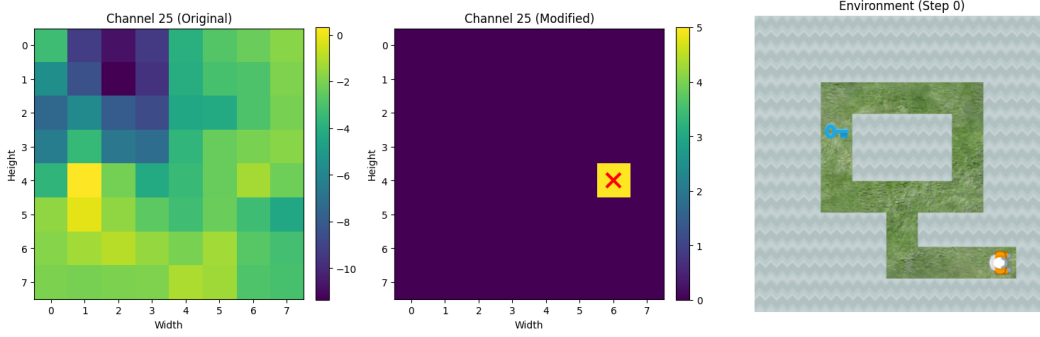


Figure 2: Example of the intervention experiment. We test channels to see the extent to which a single channel is capable of steering the agent to another location given the presence of a single entity in the environment. We consider an intervention a success if the agent enters the modified region.

65 demonstrated that networks achieving good generalization distribute information redundantly, with
 66 no single neuron being critical for performance.

67 Network Dissection Bau et al. [2017] assumes one-to-one neuron-concept mappings, but more recent
 68 work has demonstrated that neural networks are often polysemantic Elhage et al. [2022] and show that
 69 networks store multiple semantic concepts in superposition. Sparse Autoencoder research Bricken
 70 et al. [2023] successfully disentangles these polysemantic representations by training an autoencoder
 71 to separate semantic representations into a sparse overcomplete basis.

72 3 Methods

73 3.1 Model Architecture

74 We train a compressed version of the Impala model that uses 5 convolutional layers instead of 15
 75 as used in the original Impala paper. This architecture comes from Hilton et al. [2020] where they
 76 find that this model was more interpretable with no discernible loss in performance. We include the
 77 model architecture in full in Appendix A.

78 When training our model we use a simple PPO implementation from an open source implementation
 79 designed to train Progen models called Progen-pfml. We used the easy distribution of the environ-
 80 ment due to compute limitations. We tested model training over a number of different batch sizes and
 81 distributions. The model checkpoints used to derive our main findings were trained with standard
 82 hyperparameters: learning rate 5e-4, 64 parallel environments collecting 256 steps each (16,384
 83 total steps per update), processed in batches of 8 over 3 epochs for 800 million steps. The training
 84 dynamics were significantly impacted by changes to batch size often completely collapsing training
 85 performance, but over the course of 5 runs with separate seeds with the parameters above we were
 86 able to replicate similar model dynamics in each case.

87 3.1.1 Intervention experiments

88 When doing our experiments, we would place a specific entity into the environment. We then
 89 intervene on a single channel in a single layer by applying a zero mask to the channel setting it to
 90 zero and setting a specific region of the channel to a chosen value. In this experiment, we would do a
 91 sweep over a range of intervention values between 0 and the maximum values that the channel would
 92 reach during typical functioning as determined by sampling from the environment during operation.
 93 By running this sweep, we can get a sense of what the different channels do at different intervention
 94 strengths. We kept the length of the episode to 20 steps to ensure that the agent would either be able
 95 to reach the target zone, or the actual entity, but not both.

96 For each entity-channel combination, we ran 500 trials with the activation value incrementing linearly:
 97 $v_i = i \cdot \frac{v_{\max}}{N}$, where i is the trial index (1 to N), $N = 500$ is the total number of trials, and v_{\max} is
 98 the maximum activation range.

99 We use a "Q" shaped maze with the design visible in Figure 2 because it was challenging enough
 100 that the agent would need to retain its navigational abilities to reach either target, while providing a
 101 clear decision point where the agent would need to choose between pursuing the original entity or
 102 our artificial target zone.

103 We test a single entity at a time in this way meaning that steering that is effective with a particular
 104 value on a given channel may not be effective for redirecting the agent in the presence of another
 105 entity.

106 To confirm that our results are robust, we also train Sparse Autoencoders (SAEs) Bricken et al. [2023]
 107 on the convolutional layers using techniques from Gorton [2024] to confirm that the results are not
 108 merely the result of polysemanticity.

109 3.2 Channel Ablation Study

110 To systematically assess the role of individual channels, we developed a comprehensive ablation
 111 methodology. For each channel in the network, we ran episodes with only that channel active while
 112 masking all other channels to zero. This allowed us to measure the isolated capability of each channel
 113 to support navigation to different entities in the environment. We recorded successful entity collection
 114 counts for each channel across 400 rollouts, providing a quantitative measure of each channel’s
 115 specialization and capability. This methodology allows us identify which channels are sufficient for
 116 basic navigation. An example of the maze used is given in 3.

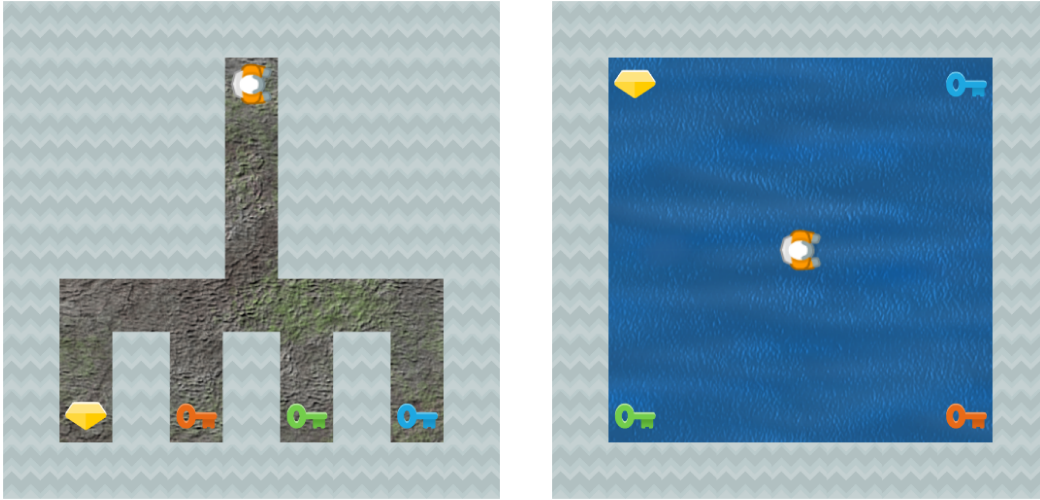


Figure 3: Maze configurations for ablation studies. (a) Fork maze tests individual channel navigation capability by requiring explicit path choices towards different entities and providing sufficient complexity that navigational abilities must be preserved. We orient the stem toward the model’s inherent bias direction to minimize false positives. (b) Open maze with entities in corners tests the effect of ablating intervention spans without navigational constraints, ensuring only preferences regarding entity collection are observed.

117 3.3 Expanded ablation experiments

118 To derive further insights from our earlier results, we ablate regions of the network based on their
 119 success in modifying model behavior in the previous experiments. Our experiment uses a modified
 120 maze which removes all walls, and has the keys and gem randomly distributed in the corners of the
 121 maze.

122 We perform four targeted ablations based on our steering and navigation findings:

123 **Ablation 1 (Intervention spans):** We zero out any activations within the value ranges that suc-
 124 cessfully steered the agent. Specifically, for entity e and channel c , we set $h_{c,i,j} = 0$ whenever the
 125 activation falls within the successful steering range $\mathcal{S}_{e,c}$ identified in our intervention experiments.

Ablation 2 (Navigation channels): We completely ablate the top 10 channels identified as navigation-critical in our channel isolation study, setting these channels to zero while preserving all others.

Ablation 3 (Combined): We apply both ablations simultaneously, zeroing both the navigation channels and the intervention spans.

Ablation 4 (Inverse): We preserve only the intervention spans that successfully steered the agent in the presence of entity e , ablating everything else. This tests whether these spans alone are sufficient for navigation.

4 Results

All results presented are from checkpoint 35001 (35k update steps), well past convergence which occurred around 20k, unless otherwise noted.

4.1 Quantitative Steering Results

We apply our incremental steering experiments across all channels in conv3a and conv4a and find a large number of channels which successfully steer the agent away from a given entity. We present our results from conv4a here as its effective value ranges were more sparse and provided clearer results.

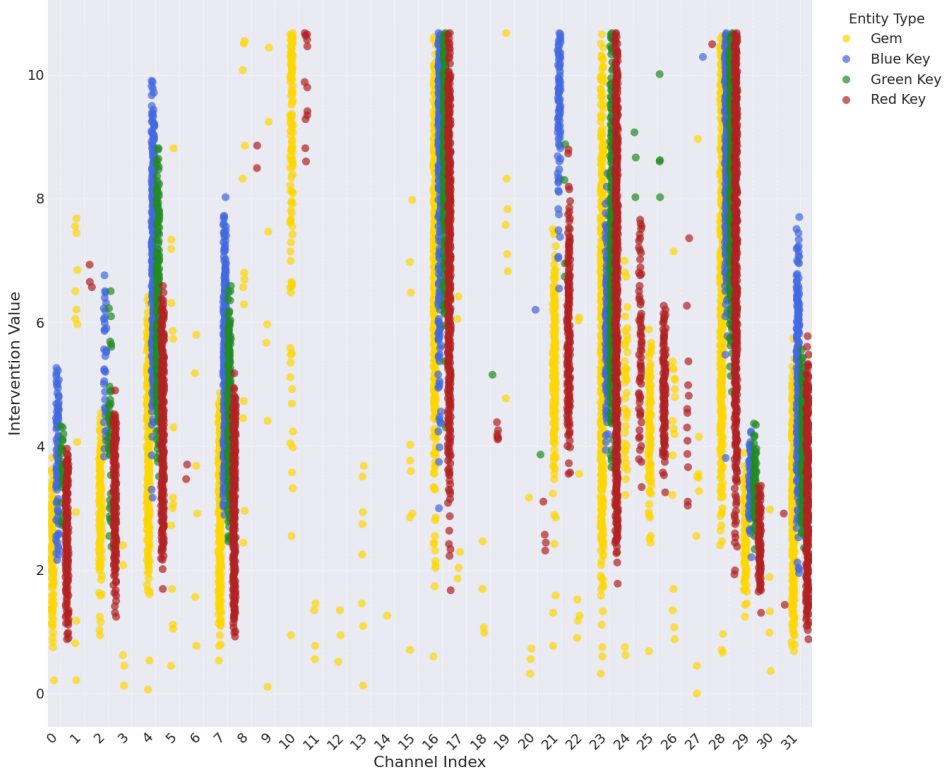


Figure 4: To determine the frequency with which a given channel can be used to steer the model away from a given entity, we create an artificial maze and place an artificial activation at a given point of the maze. We determine a successful intervention based on whether the agent makes contact with the region we specify. Each colored dot above represents a successful intervention with the specified entity as the target for the agent to move towards. Intervention success is partially localized to specific activation value ranges based on entity. This sample is from channel conv4a.

The most striking aspect of our results is the fact that particular value ranges worked for specific entities, illustrating how the channel seemed to signal the presence of a particular entity through the amplitude of the activation within a channel.

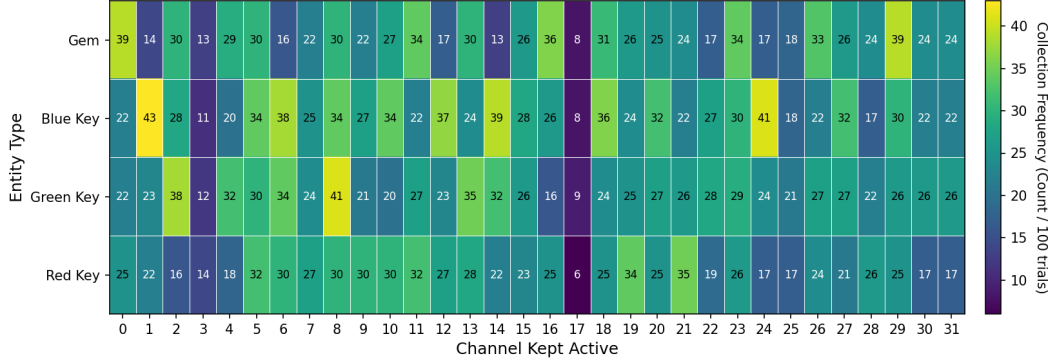


Figure 5: Heatmap showing entity collection counts when only a single channel is active (all others ablated). The vertical axis shows different entity types (gem, blue_key, green_key, red_key), while the horizontal axis shows channel indices. Brighter colors (yellow) indicate higher collection success rates.

We trained SAEs on these layers and found they exhibit identical multiplexing patterns, confirming this is a robust organizational strategy rather than the result of polysemanticity.

4.2 Channel Ablation Studies

Our ablation analysis reveals that the ability to navigate is surprisingly redundantly encoded throughout the network. The heatmap in Figure 5 shows the results in full. We observe that all channels have some success with reaching the entities some of the time with the best performing channels achieving success rates greater than 40% for certain entities, while others have lower than 10% success rates between entities. We also see that some channels show similar performance across entities, while others are significantly better at particular entities.

Table 1: Functional specialization between steering and navigation channels

Entity	Steering-only	Navigation-only	Overlap
Gem	8 channels	8 channels	2 (20%)
Blue Key	7 channels	7 channels	3 (30%)
Green Key	8 channels	8 channels	2 (20%)
Red Key	9 channels	9 channels	1 (10%)
Total unique	15 channels	19 channels	6 channels

4.3 Expanded ablation study

Our ablation studies reveal two counterintuitive findings.

First, ablating the gem’s identified steering channels improves collection rates from 96.8% to 98.0%, while the same ablation reduces key collection by 43% on average. This suggests the gem may be encoded through the absence of key and lock signals, rather than through dedicated positive features. This broader intervention space suggests the network affords greater flexibility to gem-pursuit circuitry, as it can safely activate without competing with other navigation objectives once all preceding entities have been collected.

Second, we discover that partial ablation can be more harmful than complete ablation. When we ablate only navigation channels, gem collection drops to 34.0%. However, when we additionally ablate the intervention spans (removing more of the network) performance improves to 47.8%. This suggests an interference mechanism between the navigation and steering systems.

To further investigate the distributed nature of these representations, we performed inverse ablations where we preserved only the intervention spans for each entity. This revealed a clear preference of pursuit: blue key maintained 44% performance (well above the 26% random baseline), while

Table 2: Effect of targeted ablations on entity collection rates. Bold values indicate performance improvements or minimal degradation compared to baseline.

Entity	Baseline	Intervention Spans Only	Navigation Channels Only	Both	Preserve Only Intervention Spans	Random (Control)
Gem	96.8%	98.0%	34.0%	47.8%	20.0%	24.0%
Blue Key	98.0%	57.0%	98.3%	66.3%	40.0%	26.0%
Green Key	94.8%	52.0%	96.3%	57.0%	32.0%	22.0%
Red Key	80.0%	34.8%	89.0%	39.3%	28.0%	31.5%

gem performance dropped to 16% (below the 24% random baseline). This confirms that early-game entities have more localized representations while the final goal requires whole-network context.

5 Discussion

5.1 Solution Multiplicity and Representational Drift

A striking observation from our experiments was that intervention positions changed completely across different checkpoints even after convergence. Despite maintaining similar performance, the network continued exploring different encoding schemes within the solution space. This suggests that amplitude-based multiplexing is not a unique solution but rather one of many equivalent representations the network can adopt. The continuous drift between encoding schemes post-convergence indicates the network exists on a "valley floor" of equally viable solutions, constantly reorganizing its internal representations while preserving behavioral performance.

We find a notable checkpoint at update step 30001 where no steerable channels at all are found for any entity, though we only found 1 example of this in the 60 checkpoints we tested, suggesting it might be less optimal than those with the stronger intervention configuration.

This has meaningful implications for approaching understanding the inner workings of models. The specific channels and activation ranges we identify for steering are not fixed properties of the task but rather snapshots of a dynamic system exploring equivalent representations as possible solutions. This representational drift may explain some of the difficulty in creating robust interpretability tools and suggests that interventions may need to be continuously recalibrated as networks evolve, though the extent to which this phenomena occurs beyond RL is unknown.

5.2 Activation-Level Entity Encoding

Our findings regarding how the intervention spans encode multiple entities at different activation strengths reveal a sophisticated information compression mechanism. Rather than requiring separate channels for each entity type, the model has learned to use activation magnitude as an additional dimension for encoding information. This suggests that interventions on neural networks may need to consider not just which channels to modify, but also the precise activation values to use.

This activation-level encoding also helps explain why the model can maintain full functionality with so few active channels - the network has effectively learned to multiplex information within individual channels.

An explanation for why this happens might be that due to the similarity of the task of tracking each key, it proves the most efficient solution to have re-purpose channels already capable of entity detection in general, and to have them simply target each one sequentially, while indicating the nature of the given entity by varying the activation strength.

5.3 Role specialization

A surprising finding was that some channels seem able to unilaterally alter where the agent will navigate to, many other channels were more reliably able to navigate the agent when the other channels were ablated, but were completely unable to individually steer the agent, as evidenced by examining 4 and 5. This indicates that the specific function of indicating that the agent "must go

205 to this spot" were exclusive to the steering channels, and the channels that were more successful in
206 navigating had a role specifically in identifying how to get to a position, that in the absence of any
207 other channels providing contrasting signals would lead the agent to go to that position.

208 This idea of specialization and reliance between channels is reinforced by the fact that ablating both
209 the navigation channels and the intervention spans produced better collection rates than ablating the
210 best navigation channels while leaving the intervention spans untouched. The fact that this occurred
211 52% of cases with $n=400$ trials is evidence that this effect is robust and not just noise.

212 Unfortunately we do not have a clear understanding of the precise nature of the intervention spans at
213 this point, or what leads to the interference effect specifically, and hope to advance this question in
214 future work.

215 6 Conclusion

216 Our work demonstrates that entities can be multiplexed within single channels in a neural network,
217 and exploring these channels leads to surprising specialization within the network.

218 Our key findings include:

- 219 • Steering using a single channel worked in multiple layers (conv3a, conv4a) across check-
220 points during training.
- 221 • While some channels seemed to both track certain entities and were able to affect steering,
222 some channels were highly efficient at tracking, but couldn't modify the location the agent
223 would go to, indicating specialized roles in the model.
- 224 • We found that the model was still able to navigate to different keys despite only having a
225 single channel active.
- 226 • We found that there exist many regions within channels in the model that were capable
227 of successful agent re-targeting in the presence of different entities at different activation
228 strengths, revealing a mechanism for multiplexing information within a single channel.
- 229 • We uncovered surprising results showing that in some cases when ablating the best navigation
230 channels it in fact improved performance to also remove the steering channels.
- 231 • We discover the surprising phenomenon of models slowly exploring a variety of global min-
232 ima of possible solutions to the Procgen Heist environment without degrading performance
233 as their changes occur.

234 These results show that there is specialization that occurs within the network both at the level of
235 activation amplitude to signal specific entities, and also between channels where certain channels
236 seem to play a role of directing the heading of the network while other channels provide navigational
237 functionality to the agent.

238 7 Future Work

239 Promising directions for future work include:

- 240 • Better understand the precise mechanism of navigation within the network, and the mecha-
241 nism of steering, and how these impact the next layer.
- 242 • Train a decision transformer or RNN style network that we can train probes on, and do
243 targeted interventions in the residual stream.
- 244 • Better understanding the role of negative activations in agent navigation.
- 245 • Do further analysis on the difference in the treatment of keys and locks within the network.
- 246 • Techniques such as attribution-based parameter decomposition Braun et al. [2025] are also
247 likely to yield valuable and potentially generalizable insights if applied to networks such as
248 this one.

249 As RL techniques increasingly influence frontier models, understanding these fundamental organiza-
250 tional principles will likely provide valuable insights into building interpretable and controllable AI
251 systems.

References

- A. Atrey, K. Clary, and D. Jensen. Exploratory Not Explanatory: Counterfactual Analysis of Saliency Maps for Deep Reinforcement Learning, Feb. 2020. URL <http://arxiv.org/abs/1912.05743>. arXiv:1912.05743 [cs].
- D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network Dissection: Quantifying Interpretability of Deep Visual Representations, Apr. 2017. URL <http://arxiv.org/abs/1704.05796>. arXiv:1704.05796 [cs].
- D. Braun, L. Bushnaq, S. Heimersheim, J. Mendel, and L. Sharkey. Interpretability in Parameter Space: Minimizing Mechanistic Description Length with Attribution-based Parameter Decomposition, Feb. 2025. URL <http://arxiv.org/abs/2501.14926>. arXiv:2501.14926 [cs].
- T. Bricken, A. Templeton, J. Batson, B. Chen, A. Jermyn, T. Conerly, N. Turner, C. Anil, C. Denison, A. Askell, R. Lasenby, Y. Wu, S. Kravec, N. Schiefer, T. Maxwell, N. Joseph, Z. Hatfield-Dodds, A. Tamkin, K. Nguyen, B. McLean, J. E. Burke, T. Hume, S. Carter, T. Henighan, and C. Olah. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023.
- N. Elhage, T. Hume, C. Olsson, N. Schiefer, T. Henighan, S. Kravec, Z. Hatfield-Dodds, R. Lasenby, D. Drain, C. Chen, R. Grosse, S. McCandlish, J. Kaplan, D. Amodei, M. Wattenberg, and C. Olah. Toy models of superposition. *Transformer Circuits Thread*, Sept. 2022. URL https://transformer-circuits.pub/2022/toy_model/index.html. Publisher: Anthropic.
- L. Gorton. The Missing Curve Detectors of InceptionV1: Applying Sparse Autoencoders to InceptionV1 Early Vision, June 2024. URL <http://arxiv.org/abs/2406.03662>. arXiv:2406.03662 [cs] version: 1.
- S. Greydanus, A. Koul, J. Dodge, and A. Fern. Visualizing and Understanding Atari Agents. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1792–1801. PMLR, July 2018. URL <https://proceedings.mlr.press/v80/greydanus18a.html>. ISSN: 2640-3498.
- J. Hilton, N. Cammarata, S. Carter, G. Goh, and C. Olah. Understanding RL Vision. *Distill*, 5(11): e29, Nov. 2020. ISSN 2476-0757. doi: 10.23915/distill.00029. URL <https://distill.pub/2020/understanding-rl-vision>.
- T. McGrath, M. Rahtz, J. Kramar, V. Mikulik, and S. Legg. The Hydra Effect: Emergent Self-repair in Language Model Computations, July 2023. URL <http://arxiv.org/abs/2307.15771>. arXiv:2307.15771 [cs].
- U. Mini, P. Grietzer, M. Sharma, A. Meek, M. MacDiarmid, and A. M. Turner. Understanding and Controlling a Maze-Solving Policy Network, Oct. 2023. URL <http://arxiv.org/abs/2310.08043>. arXiv:2310.08043 [cs].
- A. S. Morcos, D. G. T. Barrett, N. C. Rabinowitz, and M. Botvinick. On the importance of single directions for generalization, May 2018. URL <http://arxiv.org/abs/1803.06959>. arXiv:1803.06959 [stat].
- A. M. Turner, L. Thiergart, G. Leech, D. Udell, J. J. Vazquez, U. Mini, and M. MacDiarmid. Steering Language Models With Activation Engineering, Oct. 2024. URL <http://arxiv.org/abs/2308.10248>. arXiv:2308.10248 [cs].
- M. D. Zeiler and R. Fergus. Visualizing and Understanding Convolutional Networks, Nov. 2013. URL <http://arxiv.org/abs/1311.2901>. arXiv:1311.2901 [cs].

A Technical Appendices

A Model Architecture Details

ImpalaCNN architecture:

- 298 • Input: Image (H,W,C) \rightarrow Normalization [0,1] \rightarrow Format adaptation
- 299 • Conv Block 1: Conv(C \rightarrow 16, 7 \times 7) \rightarrow ReLU \rightarrow LPPool(2 \times 2, s=2)
- 300 • Conv Block 2: Conv(16 \rightarrow 32, 5 \times 5) \rightarrow ReLU \rightarrow Conv(32 \rightarrow 32, 5 \times 5) \rightarrow ReLU \rightarrow LP-
- 301 Pool(2 \times 2, s=2)
- 302 • Conv Block 3: Conv(32 \rightarrow 32, 5 \times 5) \rightarrow ReLU \rightarrow LPPool(2 \times 2, s=2)
- 303 • Conv Block 4: Conv(32 \rightarrow 32, 5 \times 5) \rightarrow ReLU \rightarrow LPPool(2 \times 2, s=2)
- 304 • Flatten \rightarrow Linear(flattened \rightarrow 256) \rightarrow ReLU
- 305 • Linear(256 \rightarrow 512) \rightarrow ReLU
- 306 • Dual heads: Policy(512 \rightarrow num_outputs) & Value(512 \rightarrow 1)

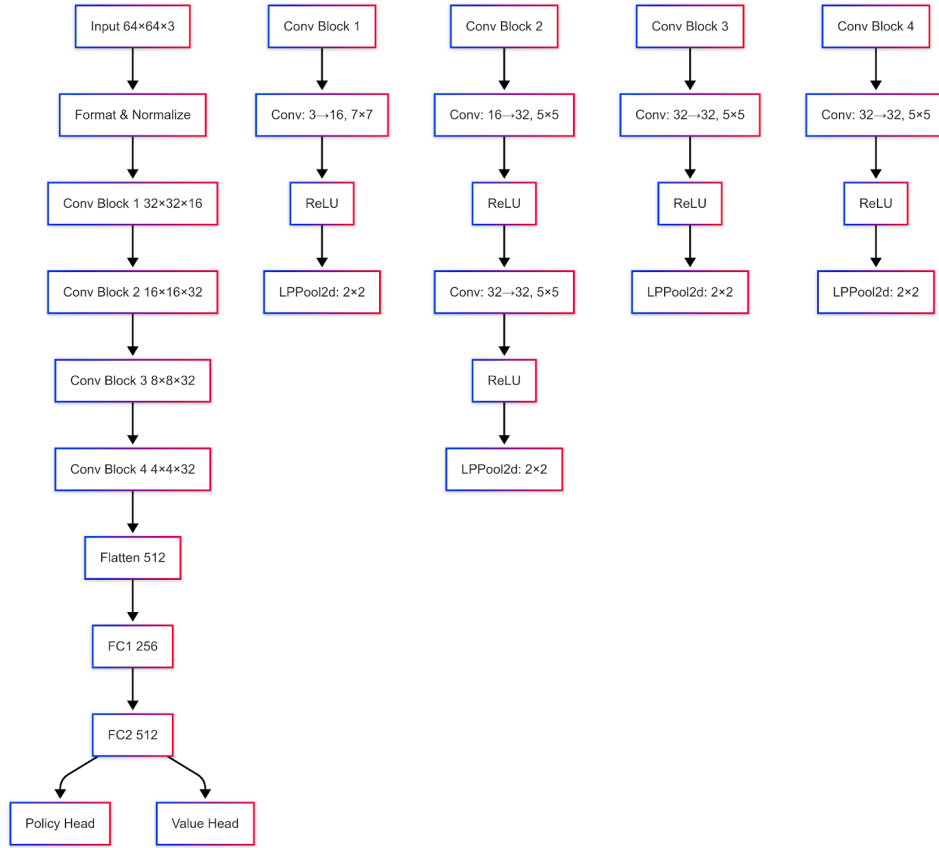


Figure 6: The convolutional neural network architecture used in our experiments. We use a modified version of the Impala CNN with 5 convolutional layers rather than the original 15 layers, which provides better interpretability without compromising performance.