

MechIR: A Mechanistic Interpretability Framework for Information Retrieval

Andrew Parry¹[0000-0001-5446-8328], Catherine Chen²[0009-0009-8734-436X],
Carsten Eickhoff³[0000-0001-9895-4061], and Sean
MacAvaney¹[0000-0002-8914-2659]

¹ University of Glasgow, Glasgow, Scotland, UK

² Brown University, Providence, RI, USA

³ University of Tübingen, Tübingen, Germany

Abstract. Mechanistic interpretability is an emerging diagnostic approach for neural models that has gained traction in broader natural language processing domains. This paradigm aims to provide attribution to components of neural systems where causal relationships between hidden layers and output were previously uninterpretable. As the use of neural models in IR for retrieval and evaluation becomes ubiquitous, we need to ensure that we can interpret *why* a model produces a given output for both transparency and the betterment of systems. This work comprises a flexible framework for diagnostic analysis and intervention within these highly parametric neural systems specifically tailored for IR tasks and architectures. In providing such a framework, we look to facilitate further research in interpretable IR with a broader scope for practical interventions derived from mechanistic interpretability. We provide preliminary analysis and look to demonstrate our framework through an axiomatic lens to show its applications and ease of use for those IR practitioners inexperienced in this emerging paradigm.

Keywords: Mechanistic Interpretability · Information Retrieval · Explainability · Machine Learning

 <https://github.com/Parry-Parry/MechIR>

1 Introduction

In recent years, machine learning models have achieved high performance, leading to widespread adoption and integration into tools used across various fields and industries. However, the internal decision-making processes of these models remain opaque due to their deep parametric nature, which makes translating model behavior into human-understandable terms difficult. This can be particularly concerning in safety-sensitive domains, such as healthcare and law, where understanding model behavior is critical for detecting and correcting potential errors. Without effective tools to investigate these complex models, our ability to intuitively understand and explain their behavior is limited.

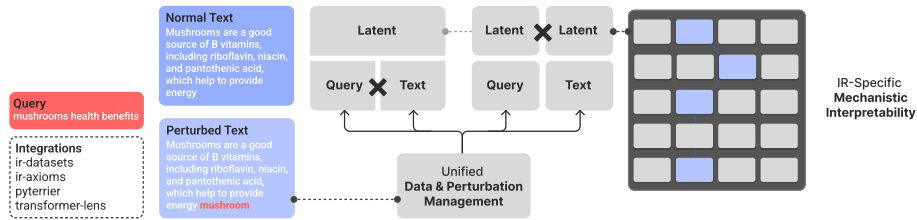


Fig. 1. MechIR allows for common IR architectures to be analyzed under text perturbation. Here a query terms is added to the text and we can observe the attention heads which respond to the addition of this term. Both cross-encoders and bi-encoders can be analyzed.

However, recent mechanistic interpretability methods offer a promising solution for understanding the internal mechanisms of neural networks by performing causal interventions on specific model components [1, 18]. Using causal methods such as activation patching, researchers have identified components responsible for generative language modelling tasks such as indirect object identification [19] and applied learned insights to improve model performance [12, 13]. Much of this progress in the broader NLP community is driven by the availability of open-source tooling and tutorials, such as TransformerLens [14]. However, current resources are tailored for generative language tasks. To make mechanistic interpretability research accessible in IR, we present **MechIR**, a Python package that aims to help researchers reverse-engineer IR models and causally investigate internal model components.

2 Background

Activation Patching Activation patching is a causal intervention method that aims to localize the specific component(s) responsible for a targeted behavior [6, 12, 18]. In IR, this approach involves a pair of inputs, consisting of one perturbed input ($X_{perturbed}$) and one baseline input ($X_{baseline}$). The perturbed input is constructed by applying some function to modify an original input (e.g., inserting a query term to the end of a document) and the baseline input is a padded variant of the original input to maintain token lengths between the pairs [2]. These pairs are then used in three forward passes:

1. *Baseline run*: run the model on $X_{baseline}$, recording model performance and caching activations (if expected performance is greater than $X_{perturbed}$).
2. *Perturbed run*: run the model on $X_{perturbed}$, recording model performance and caching activations (if expected performance is greater than $X_{baseline}$).
3. *Patched run*: run the model on $X_{baseline}$ or $X_{perturbed}$ (whichever has the lower expected performance) with model component(s) replaced from the cached activations from the other run and record the final model performance.

Model performance is defined by the relevance score for an input, with a specific measure depending on the model type. For example, the similarity score (e.g., dot product) between query and document representations is often used as the relevance score in bi-encoders, whereas logits are typically used in cross-encoders. To evaluate the effect of a patch and determine which component(s) are responsible for the difference in model behavior, we examine changes to the model’s performance during the *patched run*. For more details, we refer the reader to Chen et al. [2] and the supplementary demo notebooks.

TransformerLens TransformerLens [14] is an open-source Python package for analyzing decoder-only Transformer-based language models, widely used in mechanistic interpretability research for generative language models [3, 7, 15]. Researchers can cache, edit, remove, or replace model activations to explore how specific model components influence overall behavior. **MechIR** extends this functionality to neural IR models and additionally offers integration with widely used IR research tools such as IR datasets [10] and PyTerrier [11]. Overall, our package is an end-to-end tool for mechanistic interpretability research in IR.

3 MechIR Overview

In this section, we provide an overview of the **MechIR** framework (Figure 1) and supported functionality.

Accessing Model Activations TransformerLens enables access to a transformer’s internal activations by re-implementing the model with hooks on components of interest (i.e., attention heads and MLPs). In **MechIR**, we extend support to common retrieval architectures, including bi-encoders (e.g., TAS-B [8]) and cross-encoders (e.g., monoELECTRA [17]). A full list of supported architectures can be found in the package documentation. Below observe the instantiation of a bi-encoder in **MechIR**, the underlying interface infers the original model architecture and converts accordingly.

```

1 from mechir import Dot
2 model = Dot('bert-base-uncased')
3 all_hooks = model.hook_dict

```

Activation Patching **MechIR** supports activation patching for the following components: (1) *blocks*: patches components in an entire layer (e.g., residual stream, attention outputs, MLP layer) across individual token positions in the input pairs and (2) *attention heads*: patches attention heads in specific layers across all or individual token positions. Additionally, we provide basic plotting functionality to support analysis on experimental results, specifically to visualize the results of activation patching and attention patterns.

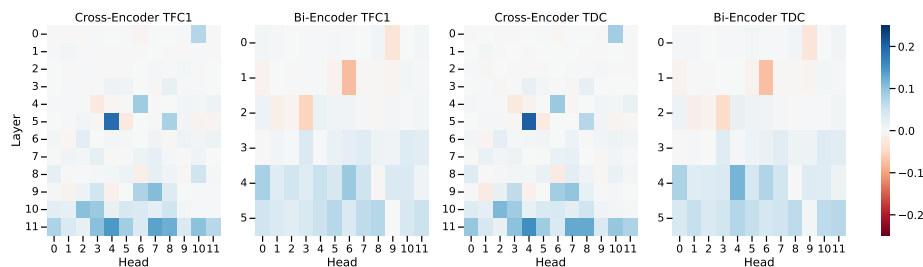


Fig. 2. Effect of inserting different types of query terms to documents (*left*: TFC1, *right*: TDC) with activation patching over a bi- and cross-encoder for a subsample of texts with a “highly-relevant” judgment.

Paired Dataset Creation Activation patching relies on pairs of inputs constructed by perturbing an original query-document pair. To create this dataset, we provide support for loading from IR datasets [10] and a *perturbation* class that allows users to create their input pairs. Provided supplementary tutorial materials discuss possible general perturbation methods and outline an analysis pipeline using our framework. Below we define a perturbation which can be added to a torch dataloader through collate functions.

```

1  from mechir.data import MechIRDataset, DotDataCollator
2  from mechir.perturb import perturbation
3  from torch.utils.data import DataLoader
4
5  @perturbation
6  def func(document):
7      return 'MechIR! ' + document
8  dataset = MechIRDataset('msmarco-passage')
9  collator = DotDataCollator(model.tokenizer, func)
10 loader = DataLoader(dataset, collate_fn=collator)
11 outputs = [model.patch(**b) for b in loader]

```

Open Source Materials We provide our full resource in an open-source Python package, which can be installed via `pip install mechir`. Additionally, we provide documentation and several supplementary tutorial notebooks in the source code repository. These tutorials cater to both novice and experienced researchers in interpretability, covering (1) an example of an end-to-end experimental pipeline (demonstration) and (2) considerations for activation patching in IR. We plan to continue development and add additional functionality for other mechanistic interpretability methods (e.g., path patching) in the future.

4 Demonstration

Our demonstration presents a concrete application of our framework, showing all aspects of an investigation of relevance approximation through activation

patching. We provide an introduction to mechanistic interpretability using common IR architectures, aiming to provide intuitions behind this research area while grounding this introduction in familiar concepts (ranking axioms).

We will guide users through the process of loading a dataset into MechIR before choosing a perturbation to create input pairs that aim to isolate behavior to particular components of neural ranking models. We take a relevance grade-stratified subsample ($n = 1000$) of the TREC DL19 and DL20 test collection and observe how different IR model architectures (cross-encoder vs. bi-encoder) behave under the addition of query terms. We perform head patching to see how the addition of random query terms (TFC1 [4]) versus the most discriminative terms by IDF (TDC [5]) changes model behaviour.

Observe in Figure 2 the diffuse nature of bi-encoder activations under term addition; no one head strongly activates under term matching with different heads in later layers activating for any query term versus discriminative query terms. This is somewhat intuitive as there is no term-level interactions between queries and documents within this architecture. Additionally, the addition of these tokens frequently is penalized in earlier layers; whether or not this is an artifact of term matching or positional bias would require future investigation⁴. In contrast, within a cross-encoder, several heads activate consistently under term matching with slight increases in activations dependent on the salience of terms; the majority of strong activations are observed in the final layers aligning with prior post-hoc probing by Jiang et al. [9] but with some strong activations in the middle layers particularly under TDC. The key difference in validating prior findings with these methods is that by isolating this behavior to particular components, we can intervene, for example, to reduce bi-encoder penalization of salient terms.

5 Target Audience and Potential Use Cases

The target audience for this demo is PhD students and researchers working on explainable information retrieval (XIR) or those interested in starting interpretability research. While the functionality of this first iteration of MechIR is focused on neural search, we hope it sparks collaborative efforts to extend mechanistic interpretability to other areas of IR, such as recommender systems. Overall, this line of research offers exciting opportunities, including the development of diagnostic tools to enhance model performance, mitigate bias, prevent adversarial attacks, and create steerable and personalized systems.

Acknowledgments. We thank Jack Merullo and Gregory Polyakov for their comments and suggestions on this work.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

⁴ See Jiang et al. [9] and Parry et al. [16] for treatments of these phenomena

References

1. Bereska, L., Gavves, E.: Mechanistic interpretability for ai safety—a review. arXiv preprint arXiv:2404.14082 (2024)
2. Chen, C., Merullo, J., Eickhoff, C.: Axiomatic causal interventions for reverse engineering relevance computation in neural retrieval models. In: Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1401–1410 (2024)
3. Conmy, A., Mavor-Parker, A., Lynch, A., Heimersheim, S., Garriga-Alonso, A.: Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems* **36**, 16318–16352 (2023)
4. Fang, H., Tao, T., Zhai, C.: A formal study of information retrieval heuristics. In: Sanderson, M., Järvelin, K., Allan, J., Bruza, P. (eds.) SIGIR 2004: Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Sheffield, UK, July 25–29, 2004. pp. 49–56. ACM (2004). <https://doi.org/10.1145/1008992.1009004>, <https://doi.org/10.1145/1008992.1009004>
5. Fang, H., Tao, T., Zhai, C.: Diagnostic evaluation of information retrieval models. *ACM Trans. Inf. Syst.* **29**(2), 7:1–7:42 (2011). <https://doi.org/10.1145/1961209.1961210>
6. Geiger, A., Lu, H., Icard, T., Potts, C.: Causal abstractions of neural networks. *Advances in Neural Information Processing Systems* **34**, 9574–9586 (2021)
7. Gurnee, W., Nanda, N., Pauly, M., Harvey, K., Troitskii, D., Bertsimas, D.: Finding neurons in a haystack: Case studies with sparse probing. arXiv preprint arXiv:2305.01610 (2023)
8. Hofstätter, S., Lin, S.C., Yang, J.H., Lin, J., Hanbury, A.: Efficiently teaching an effective dense retriever with balanced topic aware sampling. In: Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 113–122 (2021)
9. Jiang, Z., Tang, R., Xin, J., Lin, J.: How does BERT rerank passages? an attribution analysis with information bottlenecks. In: Bastings, J., Belinkov, Y., Dupoux, E., Giulianelli, M., Hupkes, D., Pinter, Y., Sajjad, H. (eds.) Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP. pp. 496–509. Association for Computational Linguistics, Punta Cana, Dominican Republic (Nov 2021). <https://doi.org/10.18653/v1/2021.blackboxnlp-1.39>, <https://aclanthology.org/2021.blackboxnlp-1.39>
10. MacAvaney, S., Yates, A., Feldman, S., Downey, D., Cohan, A., Goharian, N.: Simplified data wrangling with `ir_datasets`. In: SIGIR (2021)
11. Macdonald, C., Tonello, N.: Declarative experimentation information retrieval using pyterrier. In: Proceedings of ICTIR 2020 (2020)
12. Meng, K., Bau, D., Andonian, A., Belinkov, Y.: Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems* **35**, 17359–17372 (2022)
13. Merullo, J., Eickhoff, C., Pavlick, E.: Circuit component reuse across tasks in transformer language models. In: The Twelfth International Conference on Learning Representations (2024), <https://openreview.net/forum?id=fpoAYV6Wsk>
14. Nanda, N., Bloom, J.: Transformerlens. <https://github.com/neelnanda-io/TransformerLens> (2022)
15. Nanda, N., Chan, L., Lieberum, T., Smith, J., Steinhardt, J.: Progress measures for grokking via mechanistic interpretability. arXiv preprint arXiv:2301.05217 (2023)

16. Parry, A., MacAvaney, S., Ganguly, D.: Exploiting positional bias for query-agnostic generative content in search. In: Ku, L.W., Martins, A., Srikumar, V. (eds.) Findings of the Association for Computational Linguistics ACL 2024. pp. 11030–11047. Association for Computational Linguistics, Bangkok, Thailand and virtual meeting (Aug 2024). <https://doi.org/10.18653/v1/2024.findings-acl.656>, <https://aclanthology.org/2024.findings-acl.656>
17. Pradeep, R., Liu, Y., Zhang, X., Li, Y., Yates, A., Lin, J.: Squeezing water from a stone: a bag of tricks for further improving cross-encoder effectiveness for reranking. In: European Conference on Information Retrieval. pp. 655–670. Springer (2022)
18. Vig, J., Gehrmann, S., Belinkov, Y., Qian, S., Nevo, D., Singer, Y., Shieber, S.: Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems* **33**, 12388–12401 (2020)
19. Wang, K., Variengien, A., Conmy, A., Shlegeris, B., Steinhardt, J.: Interpretability in the wild: a circuit for indirect object identification in gpt-2 small. arXiv preprint arXiv:2211.00593 (2022)