

MECHANISTIC INTERPRETABILITY OF LARGE LANGUAGE MODELS WITH APPLICATIONS TO THE FINANCIAL SERVICES INDUSTRY

ASHKAN GOLGOON^{*,‡} , KHASHAYAR FILOM^{*,†} , ARJUN RAVI KANNAN^{*,§}

ABSTRACT. Large Language Models such as GPTs (Generative Pre-trained Transformers) exhibit remarkable capabilities across a broad spectrum of applications. Nevertheless, due to their intrinsic complexity, these models present substantial challenges in interpreting their internal decision-making processes. This lack of transparency poses critical challenges when it comes to their adaptation by financial institutions, where concerns and accountability regarding bias, fairness, and reliability are of paramount importance. Mechanistic interpretability aims at reverse engineering complex AI models such as transformers. In this paper, we are pioneering the use of mechanistic interpretability to shed some light on the inner workings of large language models for use in financial services applications. We offer several examples of how algorithmic tasks can be designed for compliance monitoring purposes. In particular, we investigate GPT-2 Small’s attention pattern when prompted to identify potential violation of Fair Lending laws. Using direct logit attribution, we study the contributions of each layer and its corresponding attention heads to the logit difference in the residual stream. Finally, we design clean and corrupted prompts and use activation patching as a causal intervention method to localize our task completion components further. We observe that the (positive) heads 10.2 (head 2, layer 10), 10.7, and 11.3, as well as the (negative) heads 9.6 and 10.6 play a significant role in the task completion.

Keywords:: Mechanistic Interpretability, Large Language Models (LLMs), Transformer Circuits, FinTech, Natural Language Processing.

The goal of this paper is to introduce the reader to the emerging field of *mechanistic interpretability* and its potential applications to financial services, especially when it comes to understanding the inner workings of *Large Language Models* (LLMs) in financial services.

1. BACKGROUND ON LLMs

The release of ChatGPT by OpenAI in late 2022 stunned the world as the chatbot set new milestones surpassing all previous publicly available systems and triggered discussions about potential risks [1, 9, 12]. At the heart of the current large language model (LLM) revolution lies the *transformer architecture*. Below, we provide a short introduction on transformers, followed by LLMs.

1.1. The transformer architecture. Transformers are deep feed-forward neural networks that leverage the *attention mechanism*. They excel in sequence modeling tasks, especially in natural language processing (NLP) [26]. Before the advent of the transformer architecture in the landmark paper [38], recurrent neural

Date: October 2024.

^{*} *Emerging Capabilities Research Group, Discover Financial Services Inc., Riverwoods, IL.*

[‡] Co-first author, ashkangolgoon@gmail.com.

[†] Co-first author, khashayar.1367@gmail.com.

[§] arjun.kannan@gmail.com.

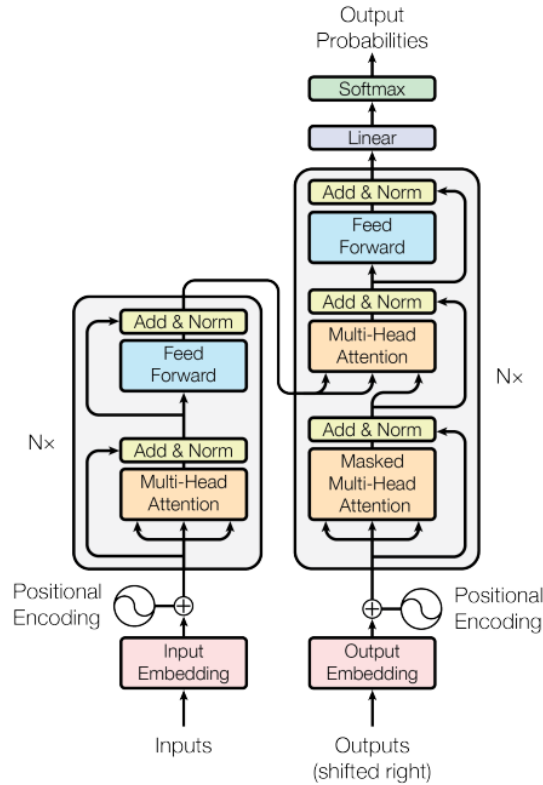


FIGURE 1. The transformer architecture. Picture adapted from [38].

network (RNN) architectures such as *Long-Short Term Memory* (LSTM) were common for NLP and sequence modeling tasks. Such models rely on an internal hidden state and must process the data sequentially. On the other hand, transformers, which are based on the attention mechanism, are superior in capturing long-term dependencies. This is due to the non-sequential and parallel manner by which they process the data. Transformers, moreover, allow *transfer learning*—they can be pre-trained on large corpora and then *fine-tuned* for downstream tasks—a fact that facilitates developing new AI applications tailored for in-house datasets based on the existing *foundation models* such as BERT or GPT-4.

In the context of sequential data, the goal can be to learn the probability distribution of the next token in *sequence modeling tasks* (e.g., language modeling), the probability distribution of a sequence conditioned on another sequence in *sequence-to-sequence tasks* (e.g., machine translation), or the probability score in a *text classification task* (e.g., sentiment analysis). As alluded to earlier, the transformer architecture shines in such tasks via utilizing the attention mechanism. The distinctions made above between different tasks are reflected in the attention type—*bidirectional/unmasked* or *unidirectional/masked* attention—and the presence or absence of *encoder* and *decoder* stacks in the network [26]. The difference between these two stacks is that the former maps into a latent space while the latter takes its inputs from a latent space. Figure 1, from the original paper on transformers [38], illustrates an encoder-decoder transformer architecture. Below, we briefly discuss the attention mechanism, followed by some other components appearing in that illustration, e.g., a *tokenization* step in the context of language tasks.

The attention mechanism is based on *key*, *query* and *value* vectors—which are all learnable. In its simplest form, the current token, the one to be predicted, is mapped to a query vector \mathbf{q} ; and the tokens in the context are mapped to key vectors \mathbf{k}_t and value vectors \mathbf{v}_t (as t varies, different tokens in the context are captured). Key and query vectors are of the same dimension, say d_{attn} , while the dimension of the value vectors may be different. Indeed, that dimension coincides with the dimension of the output vector (a representation of token and context combined) because the output is a linear combination of value vectors \mathbf{v}_t . The coefficients of this combination are the entries of $\text{softmax}\left(\frac{\mathbf{q}^T K}{\sqrt{d_{\text{attn}}}}\right)$ where the softmax function is applied to a normalization of a matrix product involving the query vector and the matrix K formed by the key vectors \mathbf{k}_t . The attention mechanism just described was based on a single query; that is, we outlined a single *attention head*. In practice, transformers use a *multi-head attention* mechanism where multiple attention heads are run in parallel and their outputs are combined by concatenation and then projection. We refer the reader to [38] for more details, or to [25, 26] for mathematically rigorous treatments of the attention mechanism.

We end this subsection by briefly mentioning some of the other components of a transformer model (cf. Figure 1). We follow [26] where precise pseudocodes for these components are presented:

- ▶ *Token Embedding*) A vector representation of each vocabulary element (token) is learned.
- ▶ *Positional Embedding*) As a remedy to the lack of recurrence or convolution in transformers, it is suggested to inject information about the position of tokens via adding certain sinusoidal terms to input embeddings [38].
- ▶ *MLP*) Blocks of fully-connected feed-forward neural networks (multi-layer perceptrons) are occasionally used in a transformer.
- ▶ *Add & Norm Layers*) *Residual connections* and *layer normalization* are used to help with the vanishing gradient problem during training, and to make the training faster and more stable.
- ▶ *Unembedding*) The model learns to convert vector representations of tokens and their contexts to a distribution over vocabulary elements.

1.2. Large Language Models (LLMs). Recent transformer-based LLMs are systems pre-trained on an enormous amount of data with an astronomical numbers of parameters. The table below, adapted from survey [16], summarizes certain aspects of famous LLMs as of 2024, e.g., their type, number of parameters, number of tokens etc. We refer the interested reader to lecture notes [5] for an introduction to LLMs. The exposition is aimed at mathematicians and physicists, and discusses the historical pretext and the phenomenology of language models among other topics.

The advent of LLMs has stimulated conversations and posed urgent questions about this disruptive technology, including on the *emergent* properties of LLMs [41], on their *alignment* with human values [35], on the *bias* present in their training data [10], on the *hallucination* problem [43], and finally, on the challenge of interpreting LLMs [36]. As a matter of fact, LLMs are expected to produce a rapidly growing array of risks which makes research on safety and governance mechanisms for them even more crucial [2]. In this paper, we focus on the interpretability question, beginning with a short introduction to the field of mechanistic interpretability in the next section.

2. MECHANISTIC INTERPRETABILITY

The nascent field of mechanistic interpretability seeks to “reverse engineer” neural networks. This endeavor can be construed in parallel with understanding the compiled binary program run on a virtual machine where the binary code and virtual machine/interpreter correspond to the parameters and the architecture of a neural

Type	Model Name	#Parameters	Release	Base Models	Open Source	#Tokens	Training dataset
Encoder-Only	BERT	110M, 340M	2018	-	✓	137B	BooksCorpus, English Wikipedia
	RoBERTa	355M	2019	-	✓	2.2T	BooksCorpus, English Wikipedia, CC-NEWS, STORIES (a subset of Common Crawl), Reddit
	ALBERT	12M, 18M, 60M, 235M	2019	-	✓	137B	BooksCorpus, English Wikipedia
	DeBERTa	-	2020	-	✓	-	BooksCorpus, English Wikipedia, STORIES, Reddit content
	XLNet	110M, 340M	2019	-	✓	32.89B	BooksCorpus, English Wikipedia, Giga5, Common Crawl, ClueWeb 2012-B
Decoder-only	GPT-1	120M	2018	-	✓	1.3B	BooksCorpus
	GPT-2	1.5B	2019	-	✓	10B	Reddit outbound
Encoder-Decoder	T5 (Base)	223M	2019	-	✓	156B	Common Crawl
	MT5 (Base)	300M	2020	-	✓	-	New Common Crawl-based dataset in 101 languages (m Common Crawl)
GPT Family	BART (Base)	139M	2019	-	✓	-	Corrupting text
	GPT-3	125M, 350M, 760M, 1.3B, 2.7B, 6.7B, 13B, 175B	2020	-	×	300B	Common Crawl (filtered), WebText2, Books1, Books2, Wikipedia
	CODEX	12B	2021	GPT	✓	-	Public GitHub software repositories
	WebGPT	760M, 13B, 175B	2021	GPT-3	×	-	ELI5
	GPT-4	1.76T	2023	-	×	13T	-
LLaMA Family	LLaMA1	7B, 13B, 33B, 65B	2023	-	✓	1T, 1.4T	Online sources
	LLaMA2	7B, 13B, 34B, 70B	2023	-	✓	2T	Online sources
	Alpaca	7B	2023	LLaMA1	✓	-	GPT-3.5
	Vicuna-13B	13B	2023	LLaMA1	✓	-	GPT-3.5
	Koala	13B	2023	LLaMA	✓	-	Dialogue data
	Mistral-7B	7.3B	2023	-	✓	-	-
	Code Llama	34	2023	LLaMA2	✓	500B	Publicly available code
	LongLLaMA	3B, 7B	2023	OpenLLaMA	✓	1T	-
	LLaMA-Pro-8B	8.3B	2024	LLaMA2-7B	✓	80B	Code and math corpora
	TinyLlama-1.1B	1.1B	2024	LLaMA1.1B	✓	3T	SlimPajama, Starcoderdata
PaLM Family	PaLM	8B, 62B, 540B	2022	-	×	780B	Web documents, books, Wikipedia, conversations, GitHub code
	U-PaLM	8B, 62B, 540B	2022	-	×	1.3B	Web documents, books, Wikipedia, conversations, GitHub code
	PaLM-2	340B	2023	-	✓	3.6T	Web documents, books, code, mathematics, conversational data
	Med-PaLM	540B	2022	PaLM	×	780B	HealthSearchQA, MedicationQA, LiveQA
	Med-PaLM 2	-	2023	PaLM 2	×	-	MedQA, MedMCQA, HealthSearchQA, LiveQA, MedicationQA
Other Popular LLMs	FLAN	137B	2021	LaMDA-PT	✓	-	Web documents, code, dialog data, Wikipedia
	Gopher	280B	2021	-	×	300B	MassiveText
	ERNIE 4.0	10B	2023	-	×	4TB	Chinese text
	Retro	7.5B	2021	-	×	600B	MassiveText
	LaMDA	137B	2022	-	×	168B	public dialog data and web documents
	Chinchilla	70B	2022	-	×	1.4T	MassiveText
	Galactia-120B	120B	2022	-	✓	450B	-
	CodeGen	16.1B	2022	-	✓	-	THE PILE, BIGQUERY, BIGPYTHON
	BLOOM	176B	2022	-	✓	366B	ROOTS
	Zephyr	7.24B	2023	Mistral-7B	✓	800B	Synthetic data
	Grok-0	33B	2023	-	×	-	Online source
	ORCA-2	13B	2023	LLaMA2	-	2001B	-
	StarCoder	15.5B	2023	-	✓	35B	GitHub
	MPT	7B	2023	-	✓	1T	RedPajama, m Common Crawl, S2ORC, Common Crawl
	Mixtral-8x7B	46.7B	2023	-	✓	-	Instruction dataset
	Falcon 180B	180B	2023	-	✓	3.5T	RefinedWeb
	Gemini	1.8B, 3.25B	2023	-	✓	-	Web documents, books, and code, image data, audio data, video data
	DeepSeek-Coder	1.3B, 6.7B, 33B	2024	-	✓	2T	GitHub's Markdown and StackExchange
DocLLM	1B, 7B	2024	-	×	2T	IIT-CDIP Test Collection 1.0, DocBank	

Table 1: An overview of popular large language models as of February 2024 (courtesy of [16]).

network. In this setting, variables/memory locations roughly correspond to neurons or other “independent units” a neural network representation can be decomposed into [21]. For basic resources on the topic, we refer the reader to the guide [18], the glossary [17], the TransformerLens library [19], or the course [14].

2.1. Motivation. Despite their immense success in various tasks, the lack of transparency of LLMs presents challenges such as hallucination, toxicity, unfairness, and misalignment with human values which can hinder safe deployment of these models. Thus, there is an urgent need for a deeper understanding of the inner functioning of LLMs. Mechanistic interpretability is an important explanation technique used to this end.

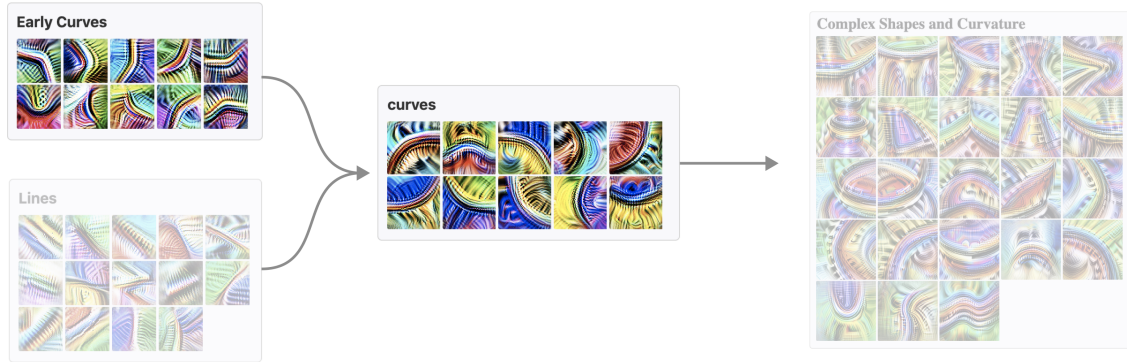


FIGURE 2. Curve detector and line detector neurons in early layers form “full curve detector circuits” that can be used to create complex shapes and geometry. The picture is based on experiments with the Inception convolutional architecture. It is from [22] (taken from <https://github.com/distillpub/post---circuits-zoom-in> under the CC-BY 4.0 license).

In this paradigm, one strives to understand an LLM at the level of neurons, *circuits*, and attention heads, i.e., at a micro scale (as opposed to *representation engineering* [44, 45] where the explanation occurs at a macro scale). In XAI (eXplainable AI) terminology, mechanistic interpretability can be described as a global, post-hoc, model-specific and white-box (i.e., needs access to model’s internals) approach [44]. Apart from helping to address societal risks, insights from mechanistic interpretability can help with “editing” an LLM for better performance. It can furthermore be utilized to explain training phenomena such as *grokking* and *memorization*; cf. [20].

2.2. Circuits. By analogy with cellular biology, researchers have tried to understand complicated neural networks by “zooming in” and investigating “the fundamental units” building a network, and the connections between them [22]. These building blocks are called “features”. Each feature corresponds to a “direction”, meaning a vector in the representation of a layer of the neural net. This can be an individual neuron or a linear combination of neurons in a layer. Features are connected to each other by weights to form “circuits”. More precisely, a circuit is a computational subgraph of the neural network where each edge connects two neurons/directions in adjacent layers, and comes with weights which are the weights shared between them in the network. It is claimed (admittedly speculatively) in [22] that features are typically meaningful, i.e., they correspond to articulable properties of the input; and the circuits correspond to meaningful algorithms encoded in neural networks’ weights. These ideas are perhaps better understood in the context of vision models—neurons may pick certain aspects of the input image and then group together to form more sophisticated detectors; see Figure 2.

Above, we presented a very brief account of features and circuits. Mechanistic interpretability is a rapidly growing field and many other aspects of these concepts have been studied.

- An important hurdle to explaining neural networks through circuits is the existence of *polysemantic* neurons; that is, neurons that respond to/get activated by multiple unrelated inputs; cf. Figure 3. Polysemanticity can be due to *superposition*, meaning when a circuit spreads a feature across many neurons (which can be inevitable since there are more features than neurons); see [17, 18] for more details, and the work [37] on extracting *monosemantic* features.

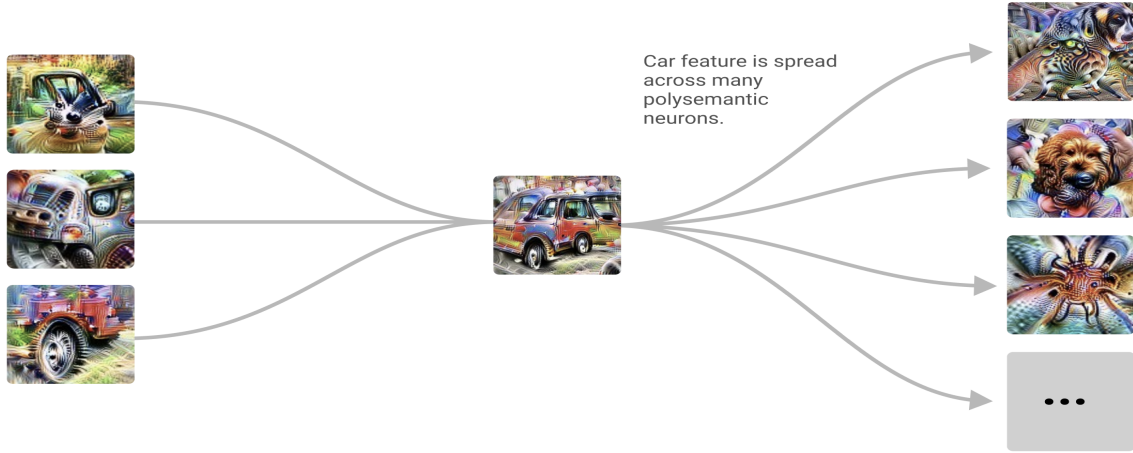


FIGURE 3. An illustration of polysemantic neurons in the context of vision models adapted from [22] (taken from <https://github.com/distillpub/post--circuits-zoom-in> under the CC-BY 4.0 license).

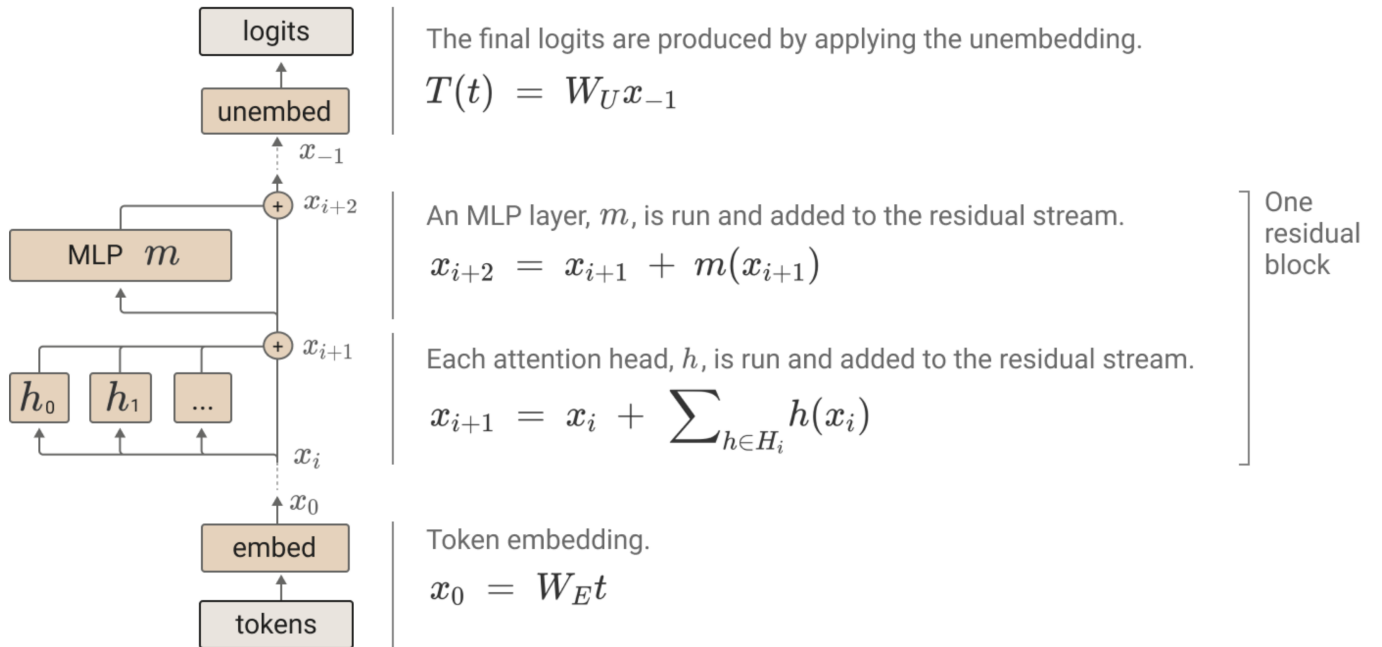


FIGURE 4. A depiction of the residual stream in a transformer (courtesy of [6]).

- *Universality* (or *convergent learning*) is the claim that analogous features and circuits form across models and tasks. This hypothesis, inspired by cell theory, clearly shapes the mechanistic interpretability research by suggesting focusing on “universal” neurons, features or circuits. This line of research has been investigated in [3, 7].

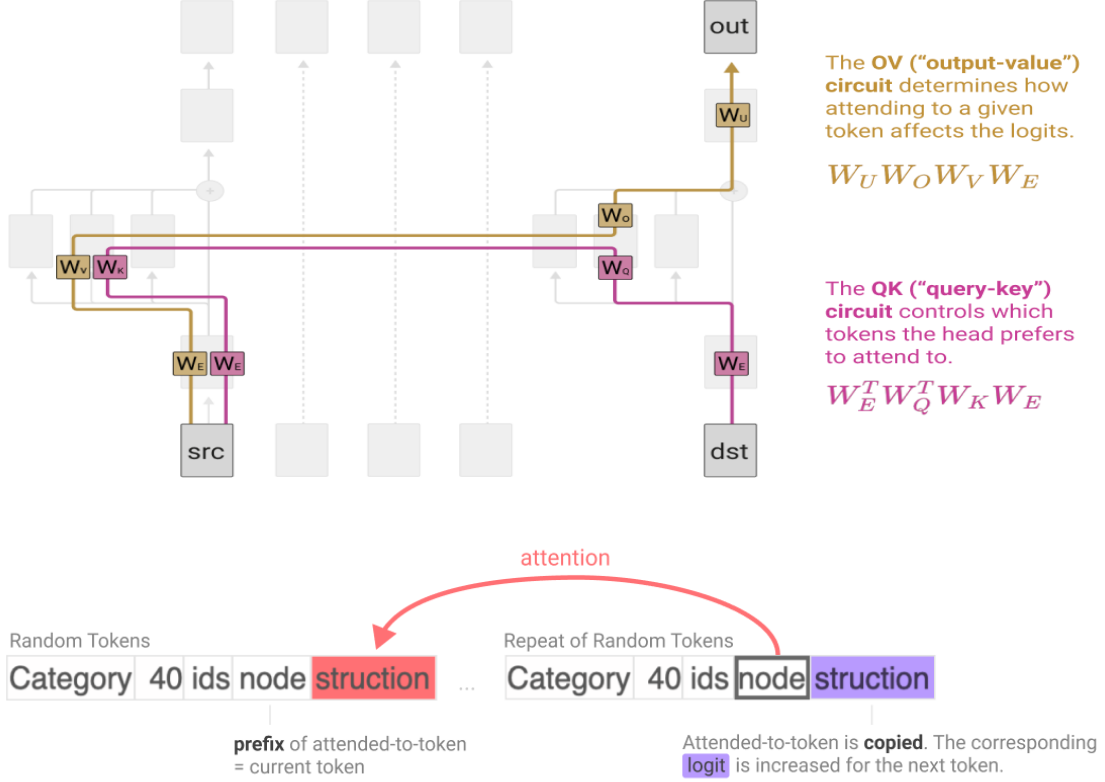


FIGURE 5. Top, an illustration of QK (query-key) and OV (output-value) circuits (courtesy of [6]). Bottom, an illustration of induction heads' behavior (courtesy of [23]).

We now delve into a more precise treatment of circuits following [6, 14]. In transformers, a layer's output is added to the *residual stream*, i.e., the sum of outputs of the previous layers and the original embedding. The residual stream can be considered as a “communication channel” between different transformer components such as MLP layers and attention heads; see Figure 4. The function of attention heads can then be understood as moving the information: reading from the residual stream of one token, and writing to the residual stream of another token. Following [14, chap. 1.1 §2, chap. 1.2 §4], and [6], we argue that each attention head consists of two circuits (Figure 5 (top)):

- A QK (query-key) circuit “determines where to move information to and from” [14, chap. 1.1 §2]. For attention head h , this can be captured by the matrix product $W_E^T (W_Q^h)^T W_K^h W_E$.
- An OV (output-value) circuit “determines what information to move” [14, chap. 1.1 §2]. For attention head h , this can be captured by the matrix product $W_U W_O^h W_V^h W_E$.

In what appeared above, matrices W_E and W_U capture embedding and unembedding of tokens; W_K^h , W_Q^h and W_V^h denote the key, query and value matrices of the attention head h ; and W_O^h captures the attention head's output weights. A thorough study of these circuits for transformers with no more than two layers

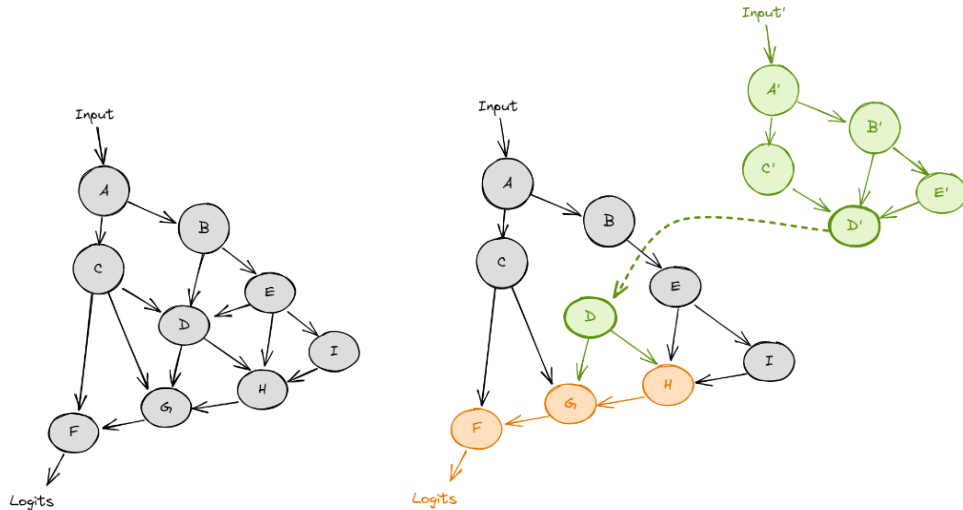


FIGURE 6. In activation patching, a clean run (e.g., "When Mary and John went to the store, John gave a drink to ...") in the IOI task is considered along with a corrupted run for which the answer is flipped (e.g., "When Mary and John went to the store, Mary gave a drink to ..."). The clean run is captured on the left. On the right, an activation from the corrupted run is patched in the corresponding one from the clean run. (Pictures courtesy of [14].)

and no MLP layer can be found in [6]. Great insights can be gained from investigating these toy models: when there is one transformer layer, or none, there are concrete descriptions in terms of n -gram models, whereas two-layer transformers are much more capable since now attention heads can compose (in three different ways corresponding to keys, queries, and values), and this results in creation of *induction heads* (Figure 5 (bottom)). Going beyond toy examples, paper [23] studies induction heads for larger models: they are implemented by circuits consisting of attention heads in different layers that work together to copy or complete patterns, and thus, contributing to *in-context learning*.

We focus on *circuit discovery* in the next subsection, especially the idea of *algorithmic tasks* which is essential to the financial services applications proposed in the last section.

2.3. Circuit discovery. There is a growing body of literature on circuit discovery for transformers; cf. [4]. As a concrete example, the phenomenon of grokking, i.e., sudden transition from overfitting to generalization after numerous training steps, can be explained through circuit formation. This is the content of paper [20] where grokking is reverse engineered for transformer models trained for simple arithmetic tasks. Also, see [8] where grokking and "emergent" abilities in LLMs are studied through the lens of circuits.

For our purposes, we shall follow the approach of [39] which is based on algorithmic tasks, an idea that we believe carries over easily to financial services applications. The paper focuses on the *Indirect Object Identification* (IOI) language task with the GPT-2 Small model ([27]), a language model comprised of 12 layers and about 80 million parameters. To give an example of this task, a sentence such as "When Mary and John went to the store, John gave a drink to ..." should be completed with the word "Mary" instead of "John". In that paper, a circuit (here, meaning a computational subgraph) of GPT-2 Small, consisting of 26 attention heads, is identified as being responsible for the IOI task. These are further broken down into 7

categories each responsible for a human-understandable purpose, e.g., “duplicate token heads” that identify tokens already appeared in the sentence.

Algorithmic examples of the IOI instances are constructed as follows (see [39, Appendix E] and [40, IOI Dataset]):

- ABC-TEMPLATES: "Afterwards [A], [B] and [C] went to the [PLACE]. [B] and [C] gave a [OBJECT] to [A]",
- BABA-TEMPLATES: "When [B] and [A] got a [OBJECT] at the [PLACE], [B] decided to give the [OBJECT] to [A]",
- BABA-LONG-TEMPLATES: "Then in the morning, [B] and [A] went to the [PLACE]. [B] gave a [OBJECT] to [A]",
- BABA-LATE-IOI: "Afterwards, [B] and [A] went to the [PLACE]. [B] gave a [OBJECT] to [A]",
- BABA-EARLY-IOI: "After the lunch [B] and [A] went to the [PLACE], and [B] gave a [OBJECT] to [A]",

such that each name was drawn from a list of 100 English first names. The place and the object were chosen from a hand-made list of 20 common words. Using these templates, a dataset of samples for IOI is created. Next, GPT-2 Small’s performance is evaluated on the prompted dataset and mechanistic interpretability of the model pertaining to the IOI task is investigated.

After this high-level overview of [39], we delve into the steps and tools involved in the circuit discovery carried out therein following the exposition in [14, chap. 1.3].

- Recall that the probabilities that a language model outputs for predicting the next token are obtained from applying a softmax function to tokens’ *logit values*. The difference of logit values, e.g., the expression $\text{logit}(\text{"Mary"}) - \text{logit}(\text{"John"})$ in the case of the IOI example above, can be utilized as a performance metric.
- The contribution of a layer (a transformer block), and the attention heads inside it, to the residual stream can be quantified by considering the logit difference after that layer as if the subsequent layers are ignored. This *direct logit attribution* technique identifies heads that contribute the most to the residual stream, and can be conducted with the TransformerLens library [19]. See [14, chap. 1.3 §2] for more details.
- A more sophisticated approach to end-to-end circuit discovery, which goes beyond just looking at what happens at the very end of a circuit, is *activation patching*, as well as its more refined version *path patching*.
 - ▶ Let us begin with the activation patching tool, the idea that was first introduced in [15]. Two different runs of the model are considered, a “clean” run and a “corrupted” one. The output answers for them are correct and incorrect, respectively. We intervene on a specific activation from the corrupted run by replacing it with the corresponding activation from the clean run, and then measure the change of the output towards the correct answer. This is called *denoising*. Alternatively, in *noising*, one patches from the corrupted run into the clean run; cf. Figure 6. Through trying many different activations and assessing how much they affect the corrupted run, one can identify the activations that really matter. Finally, we point out that patching into a transformer can be done in a number of different ways, for example, into MLP layers, attention heads, or the values of the residual stream. We refer the reader to [14, chap. 1.3 §3] for more details.
 - ▶ The activation patching approach considers the alternative of swapping the contribution of an attention head to the residual stream with what it would have been under a different distribution

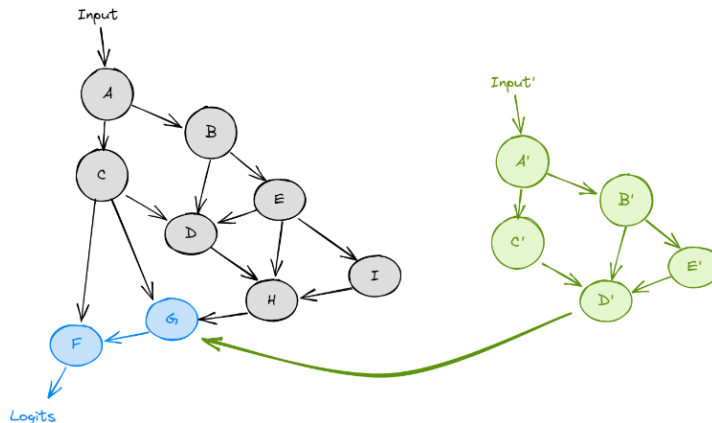


FIGURE 7. Thinking of nodes as attention heads and edges as representations of how two attention heads are tangled in the residual stream, the picture illustrates path patching where edges are replaced (unlike the activation patching where nodes are replaced; cf. Figure 6). (Picture courtesy of [14].)

while everything else remains the same. On the other hand, in path patching, one contemplates the alternative of replacing the input to an attention head from another attention head with what it would have been under a different distribution; see Figure 7. The path-patching tool, thus, investigates the importance of particular paths between a model’s components (e.g., a circuit formed by two attention heads) rather than individual model components (e.g., a single attention head). Again, a clean run and a corrupted run are involved. Looking at the IOI task again, instead of replacing “When Mary and John went to the store, John gave a drink to ...” with “When Mary and John went to the store, Mary gave a drink to ...”, three random names are used as in “When [X] and [Y] went to the store, [Z] gave a drink to ...”—hence the information about duplicate tokens is erased. We refer the reader to [14, chap. 1.3 §4] for more details.

- How does one assess if a discovered circuit is a reliable explanation for the model behavior? Paper [39] proposes three evaluation criteria for circuit analysis:
 - Faithfulness: can the circuit perform the task as well as the whole model?
 - Completeness: does the circuit contain all the nodes used to perform the task?
 - Minimality: are all the nodes in the circuit relevant to the task?

3. APPLICATIONS IN FINANCIAL SERVICES

After the remarkable success of ChatGPT, there have been efforts to design domain-specific large language models tailored to specific sectors such as law, commerce, healthcare and finance (see [11]). Such models are constructed through pre-training and/or fine-tuning on targeted datasets to perform well-defined tasks specific to a particular domain. They can tackle data security issues and prevent AI hallucination in specialized fields [24].

3.1. LLM interpretability for financial services. In the financial domain, BloombergGPT is an LLM with 50 billion parameters trained on a wide range of financial data [42]. BloombergGPT is not open

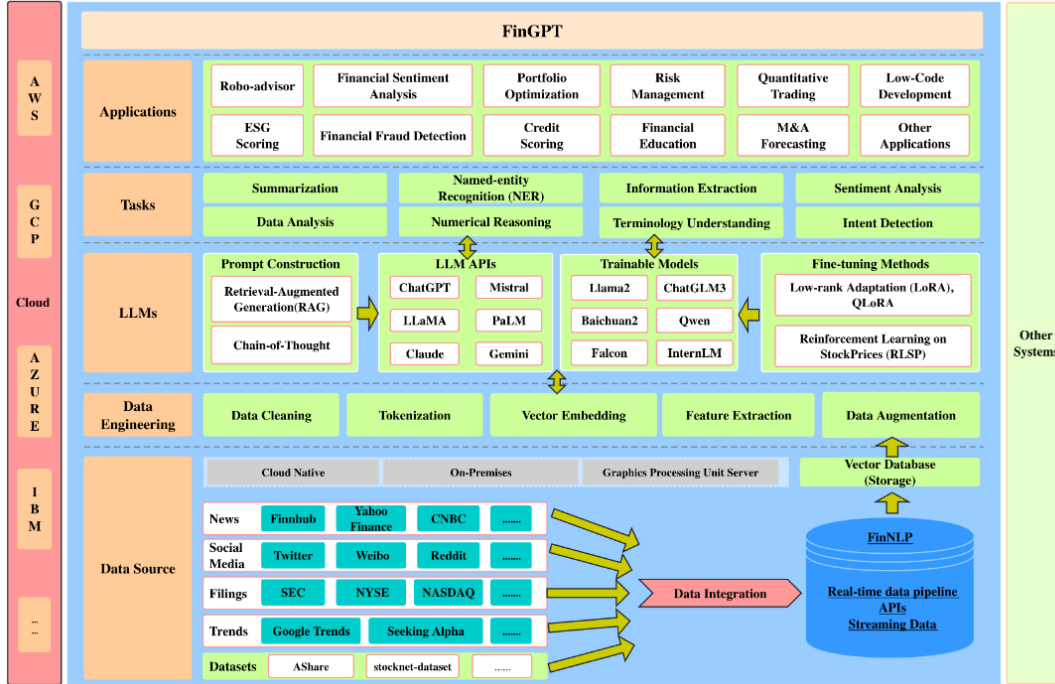


FIGURE 8. An illustration of the FinGPT ecosystem. Picture adapted from <https://github.com/AI4Finance-Foundation/FinGPT> (under the MIT license).

source, and thus, investigating it through the lens of mechanistic interpretability, e.g., circuit discovery, is not practical. Moreover, frequently retraining an LLM model like BloombergGPT is costly. Therefore, there is a need for lightweight adaptations. One remarkable adaptation is the FinGPT project which facilitates *prompt engineering* and fine-tuning of open-source LLMs on financial data [13] for a variety of financial applications such as risk management, portfolio optimization, credit scoring, and fraud detection.

In the financial services industry, LLMs, and, more generally, NLP models, have been used for various customer-service related tasks. In what follows, we outline a few industry-wide compliance, legal and business use cases.¹

- An existing large language model, such as BERT, GPT-2, or LLaMA-2, may be utilized, through fine-tuning, prompt engineering or *Retrieval Augmented Generation* (RAG) for in-house legal and compliance applications such as processing legal documents and correspondence, conducting sentiment analysis to monitor compliance with (or potential violations of) a variety of regulations overseeing financial services. For instance:
 - Monitoring cases and complaints involving Unfair or Deceptive Acts or Practices (UDAAP) [29];
 - Detecting cases related to Telephone Consumer Protection Act (TCPA) [30], e.g., circumstances where customers requests involve TCPA regulations;

¹Disclaimer: The authors have neither developed/fine-tuned any large language model for these use cases nor have they used any customer data or internal data from Discover Financial Services for their experiments.

- Compliance monitoring for regulations concerning military lending protections (e.g., Military Lending Act (MLA) [31] and Servicemembers Civil Relief Act (SCRA) [32]);
- Compliance monitoring for consumer privacy protection (as required by the California Consumer Privacy Act (CCPA) [28]);
- Compliance monitoring for fair lending (FL) laws (as required by the Equal Credit Opportunity Act (ECOA) [34] and Fair Housing Act (FHA) [33]).
- An existing large language model can be used (through fine tuning, for example) to classify customer-agent transcripts with business-related tags such as “website/mobile”, “login issues”, “hard inquiry”, “balance too high” etc.

Now, consider the problem of explaining an LLM model applied to a regulatory, legal, or business use case in financial services. Inspired by [39], the approach we adapt here is to design algorithmic tasks. In what follows, we describe how financial algorithmic tasks can be designed by working on some examples. These tasks can often be decomposed into a set of rules. These rulesets are based on legal requirements and business criteria and are often handcrafted by subject-matter experts. Below, we describe these rules for some of the compliance tasks described above followed by some algorithmic examples that can be constructed.

- **TCPA** rules) The TCPA regulation [30] prohibits unwanted communication with consumers. In compliance with this, certain conversations should be classified as critical governed by TCPA. One may design the following algorithmic examples for TCPA use cases in order to generate datasets to be used for mechanistic interpretability investigations:
 - ▶ **MARKETING-CALL-TEMPLATES**: “The agent reaches out to the customer regarding a new [FINANCIAL-PRODUCT]. The customer says I’m not interested. Please remove me from your call list.”—[FINANCIAL-PRODUCT] consists of: credit card, auto loan, mortgage loan, checking account, savings account.
 - ▶ **COMMUNICATION-TEMPLATES**: “I don’t want to receive any [WAY-of-COMMUNICATIONS] from you.”—[WAY-of-COMMUNICATIONS] is given by: calls, mails, text messages, messages, emails, notifications, communications, further communications, etc.
 - ▶ **STOP-CONTACT-LIST-TEMPLATE**: “Please add my [PROFILE] to the do-not-call lists. I don’t want to be contacted anymore.”—[PROFILE] consists of the customer contact details such as: number, phone number, personal number, work number, email address, personal email, mail address, etc.
 - ▶ **DEBT-COLLECTION-CALL-TEMPLATE**: “The agent notifies the customer regarding an [OVERDUE-BALANCE] on her account and offers a payment plan. The customer says I don’t want to receive calls about this anymore.”—[OVERDUE-BALANCE] is given by: missed payment, missed minimum payment, outstanding balance, overdue balance, unpaid balance, delinquent balance, etc.
- **UDAAP** rules) UDAAP protects consumers against the risks of harm from unfair, deceptive, or abusive practices by financial institutions. Harm does not have to be monetary and can result from increased difficulty of consumer understanding of the overall costs or risks of the product and the potential harm to the consumer associated with the product. Financial institutions monitor their customer-agent interactions. When the consumer suggests potential UDAAP (Unfair, Deceptive, or Abusive Acts or Practices) concerns, such interactions need to be identified for further review by a financial institution. One may use the following criteria to define UDAAP as an algorithmic financial task:

- UDAAP Keywords: Unfair/not fair, deceptive/deceitful, abusive/abuse/abused, tricky/trick/tricked, cheating/cheat/cheated, misleading/mislead/misled, lying/lie/lie, taking/took advantage of me because I'm a [member of a vulnerable population] (e.g., students and minority groups).

Therefore, we design the following algorithmic examples for UDAAP circuit discovery:

- ▶ UNFAIR-PRACTICES: The customer is interested in opening a checking account and asks about overdraft protection. The agent mentions that the account comes with overdraft protection but fails to clearly disclose the high fees associated with overdraft transactions. Example: "You are [UDAAP-BEHAVIOR]. You did not disclose the high fees associated with your [FINANCIAL-PRODUCT]."—[UDAAP-BEHAVIOR] can be replaced with the followings: unfair, deceptive, deceitful, abusive, misleading, fraud, liars etc.
- ▶ DECEPTIVE-PRACTICES: The customer is interested in a personal loan and the agent mentions that they offer personal loans with APR as low as 4.99%. Example: "You are [UDAAP-BEHAVIOR]. You did not disclose that the rate 4.99% only applies to applicants with excellent credit scores and my rate is significantly higher. Your loans also come with hidden origination fees that the agent did not mention."
- ▶ ABUSIVE-PRACTICES: The agent engages in aggressive debt collection by notifying the customer that he is 30 days overdue on his credit card payment and he must immediately pay his debt, or the company will take legal actions against him. Example: "You are [UDAAP-BEHAVIOR]. You are asking me to pay the full amount now and are threatening me that you will garnish my wages and reporting me to credit bureaus."

- **Fair Lending** rules) These laws set forth prohibited bases (i.e., protected classes) against which a financial institution cannot discriminate.

- Protected classes:
 - * Race, national origin, age, disability, religion, sex, marital or familial status (i.e., having children).
 - * An applicant's receipt of income from public assistance or their good faith exercise of any right under the Consumer Credit Protection Act.
 - * States may have broader protected classes under their respective fair lending laws, including but not limited to military status, ancestry, domestic violence victim status, or political affiliation.
- Examples of Potential Discrimination:
 - * "I requested for a credit line increase multiple times, and you said no because I'm on unemployment."
 - * "Why are you denying all my requests for a payment plan. It's because I'm a single mother."
 - * "I don't understand why I have all these late and hidden fees. I think you're trying to take advantage of me because I'm disabled."
 - * "You are quoting me a much higher rate than my neighbor when we have the same credit score and income. It's because I'm old!"

We design the following algorithmic examples for Fair Lending circuit discovery:

- ▶ "I don't understand why I have all these late fees. I think you're trying to [UDAAP-VERB] me because I'm a [MEMBER-of-PROTECTED-CLASS]."—[UDAAP-VERB] is given by: trick, cheat, deceit, abuse, mislead, defraud, taking advantage of.
- ▶ Note that [MEMBER-of-PROTECTED-CLASS] can be an item from elderly, single mother, disabled, minority, etc.
 - * "You are taking advantage of me because I'm a [MEMBER-of-PROTECTED-CLASS];"

- * "You are charging me a much higher rate than my neighbor. It's because I'm a [MEMBER-of-PROTECTED-CLASS];"
 - * "I asked for a credit line increase, and you denied my request because I'm a [MEMBER-of-PROTECTED-CLASS]."
- The customer is looking for a loan to expand their small business. The agent provides less favorable loan terms to the business. Example: "You are providing us with a less favorable loan terms compared to other similar businesses with similar business plans and finances. Is this because we are a [MINORITY-OWNED] business?"—[MINORITY-OWNED] can be given by: women-owned, black-owned, LGBTQ+-owned, etc.

3.2. Numerical Experiments. In this section, we design mechanistic interpretability examples for compliance tasks in financial services. For the sake of simplicity, we choose GPT-2 Small [27] as the model to study mechanistically following the methods presented in [39]. For the rest of this section, we adapt the steps presented in [14] for the IOI task and tailor it to the compliance tasks described in §3.1.

We focus on understanding financial language models tasks involving Fair Lending laws discussed before. Our goal is to analyze GPT-2 Small attention pattern in identifying potential violation of Fair Lending Laws related to the Equal Credit Opportunity Act (ECOA) [34] and the Fair Housing Act (FHA) [33]. In doing so, we design prompts of the following format: "The customer says on the phone that you are denying my request for a payment plan because I'm on unemployment."; and then we ask "Is this is an example of a Fair Lending violation based on Equal Credit Opportunity Act (ECOA)?". Next, we measure model performance by defining the following Logit Difference² metric: $\text{logit}(\text{"Yes"}) - \text{logit}(\text{"No"})$. Furthermore, we define another metric where we compare the "Yes" and "No" token probabilities of the final output, i.e., $P(\text{"Yes"})/P(\text{"No"})$. The prompts that we investigate are illustrated in Table 2. The corresponding logit differences and normalized probability ratios for the prompts are presented in Table 3. The average logit difference and probability ratio are 0.81 and 2.26, respectively.

Next, we study the model by performing direct logit attribution [39]. More specifically, we investigate the contributions of each layer (and its corresponding attention heads) to the logit difference of the residual stream.

Figure 9, illustrates the calculated logit difference by decomposing the residual stream after each layer (see *logit lens* from [19]).³ As one can see the logit difference generally keeps improving through the layers. This is in contrast to the IOI task where one can localize the task completion to merely a few final layers, i.e., layers 7, 8, and 9 [14]. We believe that this is due to the fact that our compliance task at hand is far more complex when compared to the IOI task which involves moving information around from the indirect object and not the subject, requiring less information processing.

We may investigate the logit difference between adjacent residual streams (see Figure 10). One can realize that unlike the IOI task, where only attention layers matter, MLP layers play a significant role in performing the financial compliance task. In particular, MLP layers 2 and 4 improve task completion greatly, while MLP layer 10 and attention layer 0 decrease performance.

Next, we further decompose the output of each attention layer into the sum of the outputs of the corresponding attention heads (see Figure 11). As expected, there are quite a few attention heads that play a

²Note that this choice of metric was inspired by [39] where they study the logit difference between the indirect object and the subject tokens for the circuit discovery of the IOI task.

³Recall from §2.3 that, one can interpret the result of residual stream decomposition and calculating logit differences as simulating what would happen if subsequent layers are removed after each layer [14].

Prompt	Correct	Incorrect
The customer says on the phone that you are denying my request for a payment plan because I'm on unemployment. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that you are denying my request for a payment plan because I'm on public assistance program. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that you are quoting me a higher rate than my neighbor and we both have the same credit score. It's because I'm disabled! Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that you are quoting me a higher rate than my neighbor and we both have the same credit score. It's because I'm black! Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that she asked for a credit line increase, and you said no because I'm a single mother. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that she asked for a balance transfer, and you said no because I'm a single mother. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that she asked for a credit line increase, and you said no because I'm single. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '
The customer says on the phone that she asked for a balance transfer, and you said no because I'm single. Is this an example of a Fair Lending violation based on the Equal Credit Opportunity Act (ECOA)?	' Yes '	' No '

Table 2: An illustration of the Fair Lending rule financial task investigation. Prompts and the answers.

Prompt	Yes_Logit	No_Logit	Yes_Prop	No_Prop	Yes_Rank	No_Rank	Logit_Diff	Norm_Prop_Ratio
0	14.41	13.64	3.00	1.38	1	8	0.77	2.17
1	14.61	13.74	3.62	1.53	1	6	0.87	2.37
2	14.43	13.57	3.36	1.42	1	7	0.86	2.37
3	14.56	13.76	3.38	1.53	1	6	0.80	2.21
4	14.85	14.13	2.49	1.20	2	9	0.72	2.08
5	14.65	13.90	2.36	1.11	2	9	0.75	2.13
6	14.88	14.05	2.72	1.18	1	9	0.83	2.31
7	14.73	13.84	2.71	1.11	1	9	0.89	2.44

Table 3: Final token logits, token probabilities, and token ranks for the "Yes" and "No" tokens, along with the logit differences and the normalized probability ratios.

substantial role in the compliance financial task completion. Of note, heads⁴ 11.4, 8.9, and 6.3 contribute a lot positively, while heads 0.6, 11.0, and 10.7 contribute a lot negatively to the task completion (see Figure 12).

As mentioned in §2.3, one of the limitations of the direct logit attribution discussed before is that it only looks at the very end of the circuit which affects the logits directly [14, 39]. To obtain a more refined understanding, in what follows, we utilize the activation patching technique. In doing so, we design the so-called clean prompts (Figure 13), and corrupted prompts (Figure 14). The format for the clean prompts is given by: "[A]'s credit score is 600. He was qualified for a loan of amount \$80,000. [B]'s credit score is 800.

⁴Note we use the convention that l.h specifies the head number h in the layer number l (indexing starts from 0). There are 144(12 × 12) heads in total in GPT-2 Small.

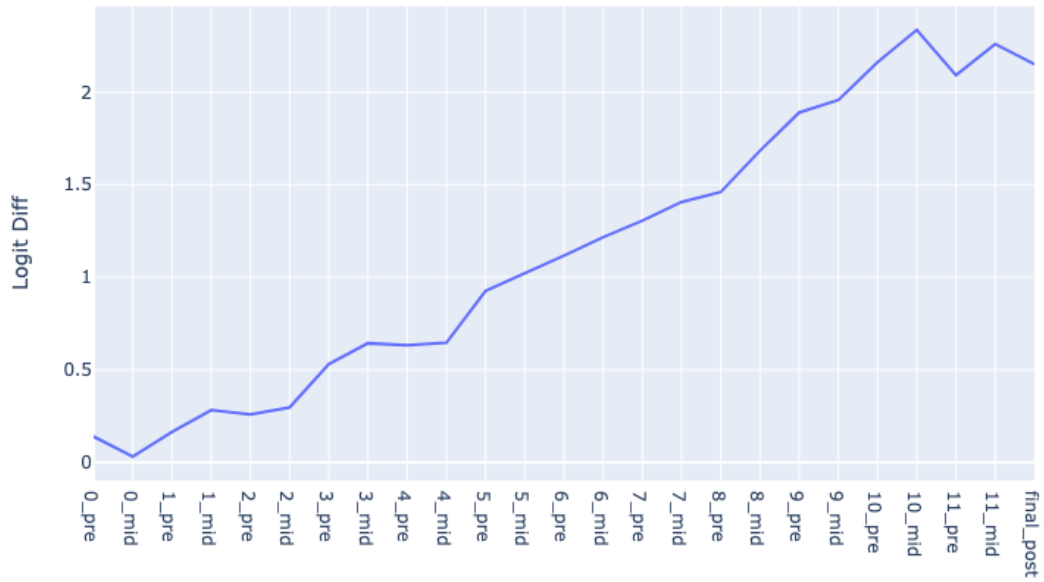


FIGURE 9. Calculated logit difference for the decomposed accumulated residual stream after each layer. n -pre denotes the residual stream at the start of layer n , while n -mid denotes the residual stream after the attention part of layer n .

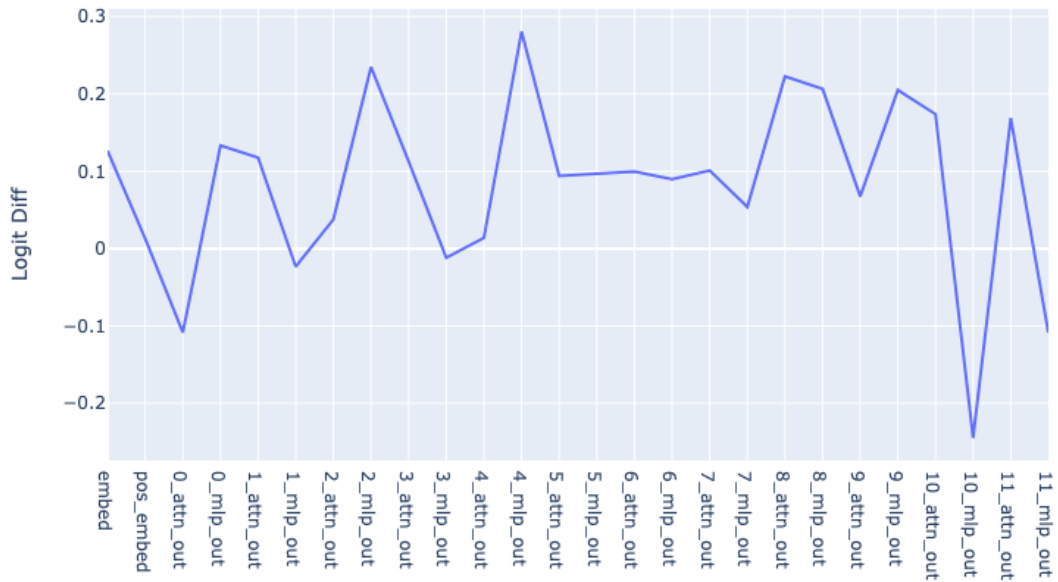


FIGURE 10. Break down of logit differences from each layer between adjacent residual streams.

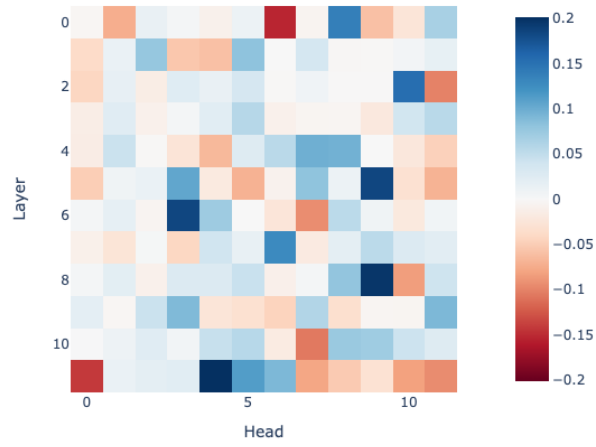
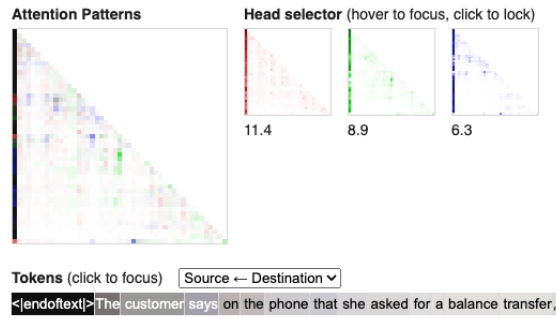


FIGURE 11. Attention heads heat map per layer illustrating logit difference from each head.

Top 3 Positive Logit Attribution Heads



Top 3 Negative Logit Attribution Heads

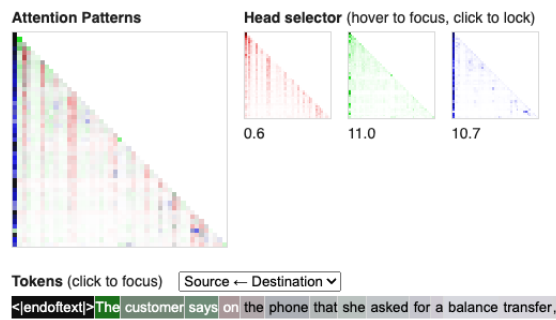


FIGURE 12. Attention patterns for the top three positive and top three negative logit attribution heads.

Prompt	Correct	Incorrect	Logit_Diff
James' credit score is 600. He was qualified for a loan of amount \$80,000. Mary's credit score is 800. She was qualified for a loan of amount \$40,000. Mary and James have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	0.010
Dan's credit score is 600. He was qualified for a loan of amount \$80,000. Mary's credit score is 800. She was qualified for a loan of amount \$40,000. Mary and Dan have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	0.093
Sam's credit score is 600. He was qualified for a loan of amount \$80,000. Mary's credit score is 800. She was qualified for a loan of amount \$40,000. Mary and Sam have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	0.164
Tim's credit score is 600. He was qualified for a loan of amount \$80,000. Mary's credit score is 800. She was qualified for a loan of amount \$40,000. Mary and Tim have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	0.088

FIGURE 13. Logit difference for "Yes" (ranked 1st) and "No" (ranked 2nd) answer tokens for the clean prompts.

Prompt	Incorrect	Correct	Logit_Diff
James' credit score is 600. He was qualified for a loan of amount \$80,000. Joe's credit score is 800. He was qualified for a loan of amount \$40,000. James and Joe have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	-0.207
Dan's credit score is 600. He was qualified for a loan of amount \$80,000. Joe's credit score is 800. He was qualified for a loan of amount \$40,000. Dan and Joe have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	-0.101
Sam's credit score is 600. He was qualified for a loan of amount \$80,000. Joe's credit score is 800. He was qualified for a loan of amount \$40,000. Sam and Joe have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	-0.072
Tim's credit score is 600. He was qualified for a loan of amount \$80,000. Joe's credit score is 800. He was qualified for a loan of amount \$40,000. Tim and Joe have the same income. Is this a gender-based fair lending violation?	' Yes '	' No '	-0.080

FIGURE 14. Logit difference for "Yes" (ranked 2nd) and "No" (ranked 1st) answer tokens for the corrupted prompts.

She was qualified for a loan of amount \$40,000. [B] and [A] have the same income. Is this a gender-based fair lending violation?", where [A] and [B] are, respectively, predominately male-associated and female-associated names. Similarly, the corrupted prompt is given by: "[A]'s credit score is 600. He was qualified for a loan of amount \$80,000. [C]'s credit score is 800. He was qualified for a loan of amount \$40,000. [A] and [C] have the same income. Is this a gender-based fair lending violation?", where [A] and [C] are predominately male-associated names. The clean run favors "Yes" over "No" as an output response, and thus, results in positive logit differences, while the corrupted run favors "No" over "Yes" as an output response resulting in negative logit differences. Next, we run the model on our corrupted prompts and then intervene by patching in the activations from our clean prompts and monitor how the output response changes.

Figure 15 depicts the results for residual stream patching at the start of each layer for all the 61 token positions. A score closer to 0 means the performance is closer to the one obtained on the corrupted input, while a score closer to 1 means that the performance is closer to that of the clean input. We are trying to look for activations that are sufficient to recover the correct response. One can make the following observations: (i) The computation is fairly localized to positions: 22 ("[B]" in "[B]'s credit score ..."), 29 ("She" in "She was qualified ..."), 42 ("[B]" in "[B] and [A] ..."), 44 ("[A]" in "[B] and [A] ..."), 54 ("gender" in "gender-based ..."),

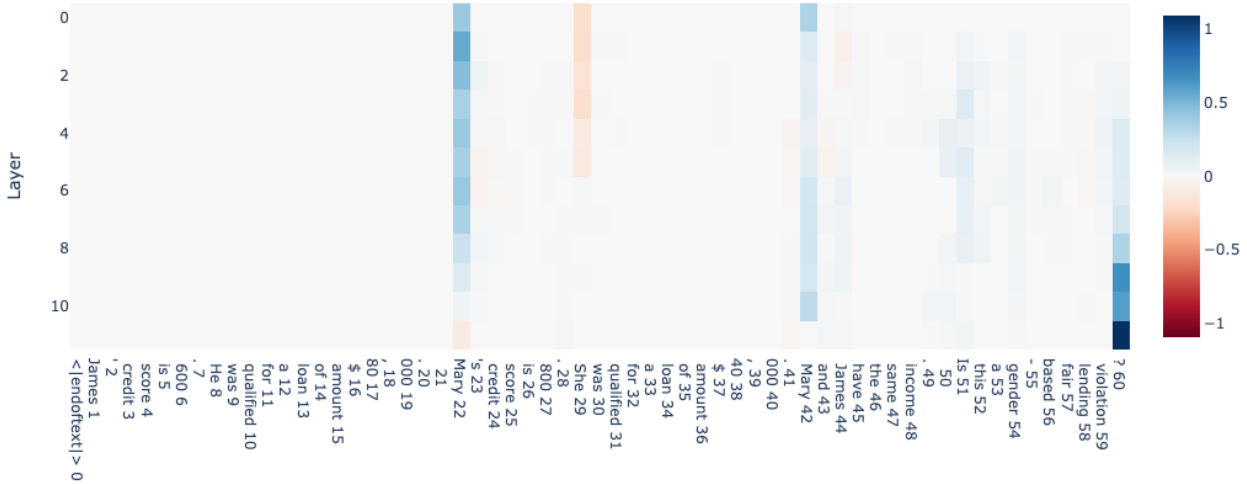


FIGURE 15. Residual stream patching at the onset of each layer for each token positions averaged over all 4 prompts. On the x-axis we only show the labels for the first prompt.

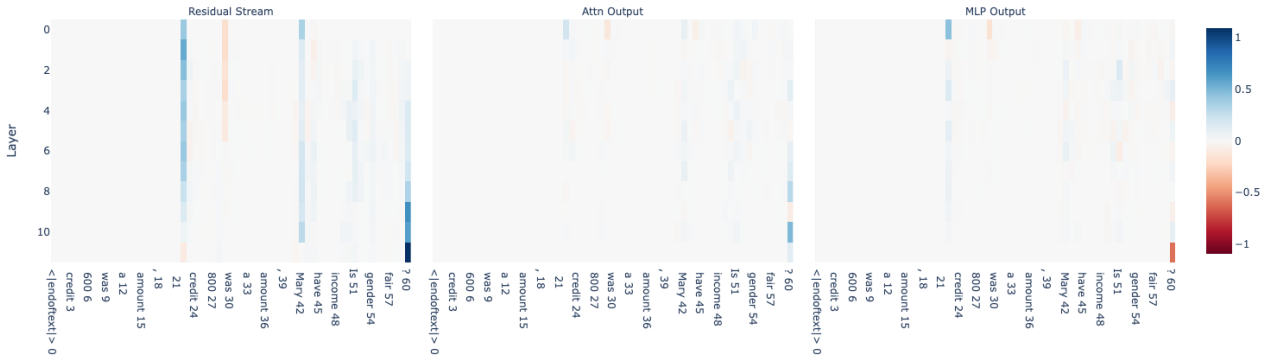


FIGURE 16. Patching in residual stream after the attention layer (in the middle) and after the MLP layer (on the right).

and END tokens. (ii) One can see that the information at token 22 starts to be moved to END token around layers 5 and 6.

Instead of just patching to the residual stream at the onset of each layer, we may use activation patching for patching after the attention layer and after the MLP layer. Figure 16 illustrates the results. As one can see layers 8 and 10 contribute positively, while layer 9 contributes negatively to the performance. Among MLP layers, we observe that MLP-0⁵ plays an important role, along with MLP-11.

Next, we can further refine our analysis (see Figure 17) by patching in on individual attention heads for all layers and tokens. We see that the heads 10.2, 10.7, and 11.3, from later layers, and the heads 0.4, 1.7,

⁵Note that this is consistent with the observations from [39, 14], that MLP-0 matters a lot. This behavior is often observed in GPT-2 Small.

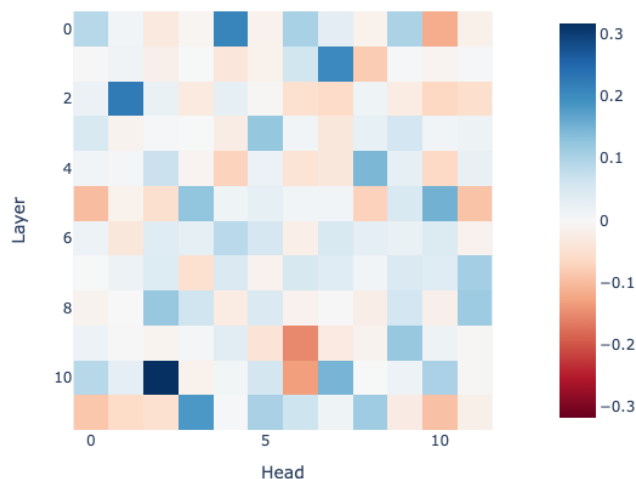


FIGURE 17. Patching in individual attention heads output for all layers and sequence positions.

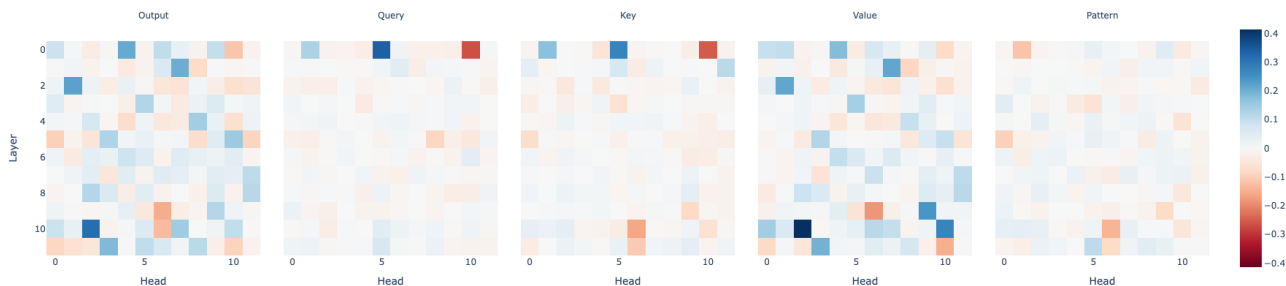


FIGURE 18. Decomposing heads by patching on the query vectors, key vectors and value vectors.

and 2.1, from earlier layers, have very large positive scores. On the other hand, the heads 9.6 and 10.6, from later layers, and the heads 0.10 and 5.0, from earlier layers, have very large negative scores.

Finally, instead of just patching on a head’s output, we decompose the attention heads by patching into the building blocks of the attention mechanism, i.e., key vectors, query vectors, value vectors, and attention patterns (see Figure 18). We can deduce the following observations: (i) Some of the layer-0 heads such as 0.5, 0.1 and 0.10 are very important because of their query and key vectors. (ii) The head 10.2 is significantly important due to its value vector. (iii) Other important attention heads owing to their value vectors are 9.6, 9.9, 10.10, 11.3, and 11.10. (iv) Thus, we can conclude that value patching has a more significant effect than key (or query) patching for the important heads in later layers, i.e., layers 9 – 11.

In summary, we investigated GPT-2 Small’s attention pattern for identifying potential violation of Fair Lending laws. Employing the direct logit attribution technique, we analyzed the contributions of each layer to the logit difference in the residual stream. In particular, we identified the heads 11.4, 8.9, and 6.3 (0.6, 11.0, and 10.7) that contribute a lot positively (negatively) to the compliance task completion. We performed a casual intervention through utilizing the activation patching technique by designing the clean and corrupted prompts. We refined our analysis to the find which components of the attention mechanism are important

for each attention head. Of note, we found that the head 10.2 is very important due to its value vector, while some of the layer-0 heads such as 0.5, 0.1, and 0.10 are very important owing to their query and key vectors.

4. FUTURE DIRECTIONS

We consider this paper to be the first step towards leveraging mechanistic interpretability techniques for understanding applications of LLMs in financial services. Some of the future research directions include: (i) Conducting experiments with more powerful open-source LLMs such as Mistral 7B and Llama-2 7B. This is where one may leverage fine-tuned open-source models from FinGPT project [13]. (ii) Working on more algorithmic examples of compliance tasks in financial services and preferably finding prompts that result in large logit differences. (iii) Perform path patching as a more sophisticated circuit discovery approach toward finding an end-to-end circuit for performing a compliance financial task.

ACKNOWLEDGEMENT

We would like to thank Sharon O’Shea Greenbach (Sr. Counsel & Director, Regulatory Policy at DFS) and Dennis Lee (Chief IP Counsel & AGC Privacy & Security at DFS) for their helpful comments relevant to regulatory issues that arise in the financial industry and for carefully reviewing the manuscript. We also thank Callum McDougall, Shervin Minaee, and Christopher Olah for granting permission to use some of the figures and tables from their papers. Ashkan Golgoon benefited from stimulating discussions with Amirhossein Tajdini. The views and opinions expressed in this paper are solely our own and do not represent or reflect those of our employer.

REFERENCES

- [1] D. Bartz. As ChatGPT’s popularity explodes, U.S. lawmakers take an interest. *Reuters*, <https://www.reuters.com/technology/chatgpts-popularity-explodes-us-lawmakers-take-an-interest-2023-02-13>, 2023.
- [2] S. R. Bowman. Eight Things to Know about Large Language Models. *arXiv e-prints*, page arXiv:2304.00612, Apr. 2023.
- [3] B. Chughtai, L. Chan, and N. Nanda. A Toy Model of Universality: Reverse Engineering How Networks Learn Group Operations. *arXiv e-prints*, page arXiv:2302.03025, Feb. 2023.
- [4] A. Conmy, A. Mavor-Parker, A. Lynch, S. Heimersheim, and A. Garriga-Alonso. Towards automated circuit discovery for mechanistic interpretability. *Advances in Neural Information Processing Systems*, 36:16318–16352, 2023.
- [5] M. R. Douglas. Large Language Models. *arXiv e-prints*, page arXiv:2307.05782, July 2023.
- [6] N. Elhage, N. Nanda, C. Olsson, T. Henighan, N. Joseph, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, N. DasSarma, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. A Mathematical Framework for Transformer Circuits. *Transformer Circuits Thread*, 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- [7] W. Gurnee, T. Horsley, Z. C. Guo, T. Rezaei Kheirkhah, Q. Sun, W. Hathaway, N. Nanda, and D. Bertsimas. Universal Neurons in GPT2 Language Models. *arXiv e-prints*, page arXiv:2401.12181, Jan. 2024.
- [8] Y. Huang, S. Hu, X. Han, Z. Liu, and M. Sun. Unified View of Grokking, Double Descent and Emergent Abilities: A Perspective from Circuits Competition. *arXiv e-prints*, page arXiv:2402.15175, Feb. 2024.
- [9] E. Klein. This Changes Everything. *New York Times*, <https://www.nytimes.com/2023/03/12/opinion/chatbots-artificial-intelligence-future-weirdness.html>, 2023.
- [10] Y. Li, M. Du, R. Song, X. Wang, and Y. Wang. A survey on fairness in large language models. *arXiv preprint arXiv:2308.10149*, 2023.
- [11] Y. Li, S. Wang, H. Ding, and H. Chen. Large language models in finance: A survey. In *Proceedings of the fourth ACM international conference on AI in finance*, pages 374–382, 2023.
- [12] T. Lieu. I’m a Congressman Who Codes. A.I. Freaks Me Out. *New York Times*, <https://www.nytimes.com/2023/01/23/opinion/ted-lieu-ai-chatgpt-congress.html>, 2023.

- [13] X.-Y. Liu, G. Wang, H. Yang, and D. Zha. Data-centric FinGPT: Democratizing Internet-scale Data for Financial Large Language Models. *NeurIPS Workshop on Instruction Tuning and Instruction Following*, 2023.
- [14] C. McDougall. ARENA (Alignment Research Engineer Accelerator) 3.0 [accessed June 2024]. <https://arena3-chapter1-transformer-interp.streamlit.app/>; <https://arena-ch1-transformers.streamlit.app/>; https://github.com/callummcdougall/ARENA_2.0, 2024.
- [15] K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- [16] S. Minaee, T. Mikolov, N. Nikzad, M. Chenaghlu, R. Socher, X. Amatriain, and J. Gao. Large Language Models: A Survey. *arXiv e-prints*, page arXiv:2402.06196v2, Feb. 2024.
- [17] N. Nanda. A Comprehensive Mechanistic Interpretability Explainer & Glossary [accessed June 2024]. <https://neelnanda.io/glossary>, Dec. 2022.
- [18] N. Nanda. Mechanistic Interpretability Quickstart Guide [accessed June 2024]. <https://www.lesswrong.com/posts/jLAvJt8wuSFySN975/mechanistic-interpretability-quickstart-guide>, Jan. 2023.
- [19] N. Nanda and J. Bloom. TransformerLens. <https://github.com/neelnanda-io/TransformerLens>, 2022.
- [20] N. Nanda, L. Chan, T. Lieberum, J. Smith, and J. Steinhardt. Progress measures for grokking via mechanistic interpretability. *arXiv e-prints*, page arXiv:2301.05217, Jan. 2023.
- [21] C. Olah. Mechanistic Interpretability, Variables, and the Importance of Interpretable Bases. <https://www.transformer-circuits.pub/2022/mech-interp-essay>, 2022.
- [22] C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom In: An Introduction to Circuits. *Distill*, 2020. <https://distill.pub/2020/circuits/zoom-in>.
- [23] C. Olsson, N. Elhage, N. Nanda, N. Joseph, N. DasSarma, T. Henighan, B. Mann, A. Askell, Y. Bai, A. Chen, T. Conerly, D. Drain, D. Ganguli, Z. Hatfield-Dodds, D. Hernandez, S. Johnston, A. Jones, J. Kernion, L. Lovitt, K. Ndousse, D. Amodei, T. Brown, J. Clark, J. Kaplan, S. McCandlish, and C. Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- [24] S. Park. Real-world examples of ‘Domain-Specific LLMs’: Bring tailored AI to your business [accessed June 2024]. <https://www.upstage.ai/feed/insight/examples-of-domain-specific-llms>, Feb. 2024.
- [25] C. Penke. A mathematician’s introduction to transformers and large language models. Technical report, Jülich Supercomputing Center, 2022.
- [26] M. Phuong and M. Hutter. Formal Algorithms for Transformers. *arXiv e-prints*, page arXiv:2207.09238v1, July 2022.
- [27] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [28] California Civil Code. California Consumer Privacy Act. https://cpa.ca.gov/regulations/pdf/cpa_act.pdf, 2018.
- [29] Consumer Financial Protection Bureau (CFPB). Unfair, Deceptive, Or Abusive Acts Or Practices. https://files.consumerfinance.gov/f/documents/cfpb_unfair-deceptive-abusive-acts-practices-udaaps_procedures_2023-09.pdf, Oct. 2012.
- [30] Federal Communications Commission (FCC). Telephone Consumer Protection Act. <https://www.fcc.gov/sites/default/files/tcpa-rules.pdf>, 1991.
- [31] Federal Deposit Insurance Corporation (FDIC). Military Lending Act. <https://www.fdic.gov/resources/supervision-and-examinations/consumer-compliance-examination-manual/documents/5/v-13-1.pdf>, 2006.
- [32] U.S. Department of Justice (DoJ). Servicemembers Civil Relief Act. <https://www.justice.gov/crt/servicemembers-civil-relief-act-summary>, 1940.
- [33] U.S. Department of Justice (DoJ). Fair Housing Act. <https://www.justice.gov/crt/fair-housing-act-1>, 1968.
- [34] U.S. Department of Justice (DoJ). Equal Credit Opportunity Act. <https://www.govinfo.gov/content/pkg/USCODE-2011-title15/html/USCODE-2011-title15-chap41-subchapIV.htm>, 2011.
- [35] T. Shen, R. Jin, Y. Huang, C. Liu, W. Dong, Z. Guo, X. Wu, Y. Liu, and D. Xiong. Large language model alignment: A survey. *arXiv preprint arXiv:2309.15025*, 2023.
- [36] C. Singh, J. P. Inala, M. Galley, R. Caruana, and J. Gao. Rethinking interpretability in the era of large language models. *arXiv preprint arXiv:2402.01761*, 2024.
- [37] A. Templeton, T. Conerly, J. Marcus, J. Lindsey, T. Bricken, B. Chen, A. Pearce, C. Citro, E. Ameisen, A. Jones, H. Cunningham, N. L. Turner, C. McDougall, M. MacDiarmid, C. D. Freeman, T. R. Sumers, E. Rees, J. Batson, A. Jermyn, S. Carter, C. Olah, and T. Henighan. Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet. *Transformer Circuits Thread*, 2024.
- [38] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

- [39] K. R. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Interpretability in the Wild: a Circuit for Indirect Object Identification in GPT-2 small. In *NeurIPS ML Safety Workshop*, 2022.
- [40] K. R. Wang, A. Variengien, A. Conmy, B. Shlegeris, and J. Steinhardt. Redwood Research Easy-Transformer. <https://github.com/redwoodresearch/Easy-Transformer/>, 2022.
- [41] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [42] S. Wu, O. Irsoy, S. Lu, V. Dabrovolski, M. Dredze, S. Gehrmann, P. Kambadur, D. Rosenberg, and G. Mann. BloombergGPT: A Large Language Model for Finance. *arXiv e-prints*, page arXiv:2303.17564, Mar. 2023.
- [43] H. Ye, T. Liu, A. Zhang, W. Hua, and W. Jia. Cognitive mirage: A review of hallucinations in large language models. *arXiv preprint arXiv:2309.06794*, 2023.
- [44] H. Zhao, F. Yang, B. Shen, H. Lakkaraju, and M. Du. Towards Uncovering How Large Language Model Works: An Explainability Perspective. *arXiv e-prints*, page arXiv:2402.10688, Feb. 2024.
- [45] A. Zou, L. Phan, S. Chen, J. Campbell, P. Guo, R. Ren, A. Pan, X. Yin, M. Mazeika, A.-K. Dombrowski, et al. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.