

Generative Models as a Complex Systems Science: How can we make sense of large language model behavior?

Ari Holtzman¹ and Peter West and Luke Zettlemoyer

University of Washington

{ahai, pawest, lsz}@cs.washington.edu

Abstract

Coaxing out desired behavior from pretrained models, while avoiding undesirable ones, has redefined NLP and is reshaping how we interact with computers. What was once a scientific engineering discipline—in which building blocks are stacked one on top of the other—is arguably already a complex systems science—in which *emergent behaviors* are sought out to support previously unimagined use cases.

Despite the ever increasing number of benchmarks that measure *task performance*, we lack explanations of what *behaviors* language models exhibit that allow them to complete these tasks in the first place. We argue for a systematic effort to decompose language model behavior into categories that explain cross-task performance, to guide mechanistic explanations and help future-proof analytic research.

1 The Newformer: A Thought Experiment

Consider the following thought experiment:

Tomorrow, researchers at an industry lab publicly release a new kind of pretrained model: the Newformer. It has a completely different architecture than the Transformer (no attention, non-differentiable components, etc.), that outperforms all pretrained Transformers on the vast majority of benchmarks. Independent labs quickly verify that these results are sound, even on just-released benchmarks. While the composition of the training data is public, it is so expensive to train that no lab can afford to replicate it, even the one that produced it. Scaled-down versions do not exhibit the same performance or interesting behaviors as the original model.

(A)

1.1 How should we study the Newformer?

Identifying high-level behaviors a model does or does not share with older models can steer us toward lower-level mechanisms it uses to solve tasks (§1.2, Figure 1). Interpretation techniques that rely on low-level details are model specific (§1.3) and often abandoned as the field changes. The Newformer is fictional, but it can help us reconceptualize the goals and methods of generative model research in light of the new landscape (§1.4).

How should we factorize model behavior into understandable and explanatory categories? (§2, Figure 2) We present a formalism for describing behavior (§2.1), noting that this corresponds to a *metamodel* that predicts aspects of a primary model (Figure 3). **Benchmarks help us measure performance, but rarely discover behavior** (§2.2) or *characterize* it (§2.3). Instead, discovered behaviors motivate new benchmarks (§2.4, Figure 4).

Generative models qualify as *complex systems* (§3), due to their *emergent behaviors* (§3.1, Figure 5), which are more often *discovered* than engineered (§3.2). A lack of clarity on *what* models do holds us back, as if we were studying organic chemistry without knowledge of biology (§3.3). This issue remains even when proprietary models are released (§3.4), as the problem lies in our lack of behavioral vocabulary; investigating possible mappings between training data and generated data can help us establish new behavioral categories (§3.5).

Despite the challenges, generative models are *easier* to study than many naturally arising complex systems (§4), because they are simulable by construction (§4.1). In contrast to physical phenomena, we can easily conduct a wide range of storable, repeatable experiments without observer effects (§4.2, Figure 7). We do, however, rely on the availability of open-source models (§4.3).

We conclude (§5) with an argument for increased focus on the foundational “what are models doing?” to guide the classic “why are models doing that?”

¹Incoming faculty at University of Chicago

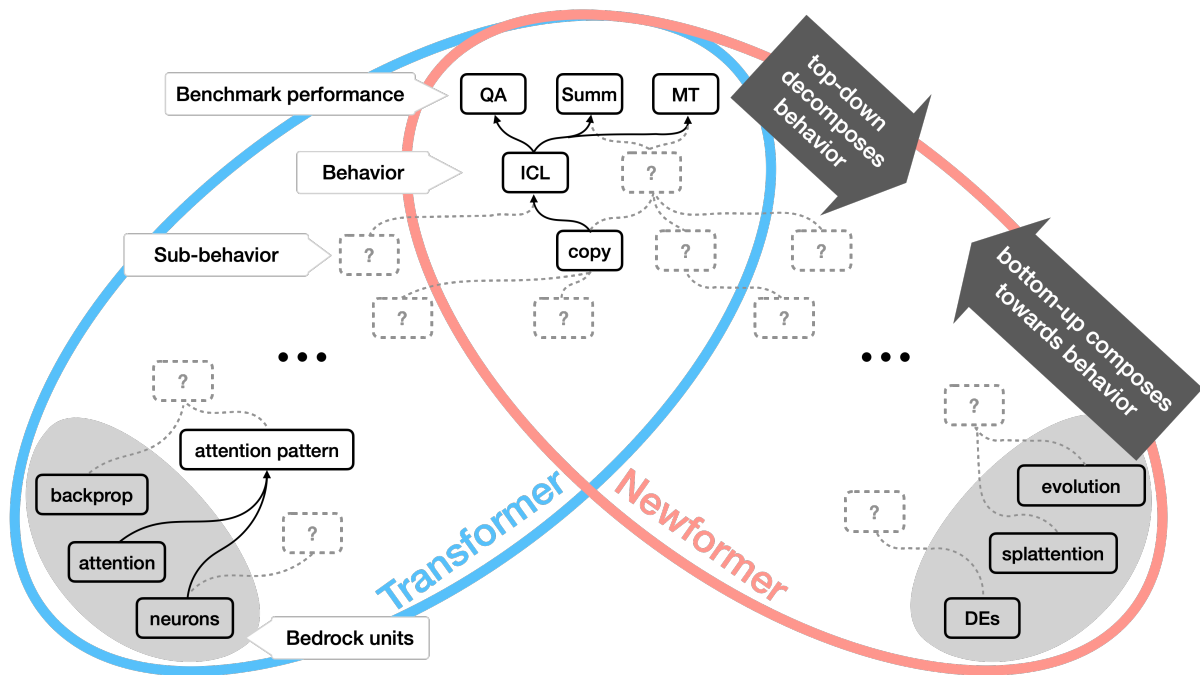


Figure 1: To explain why learned models self-organize the way they do from the bottom-up, it is useful to have top-down hierarchy of partially decomposed behaviors, to guide hypotheses with functionality we know the overall model has. While networks are composed of bedrock units for which we have a perfect understanding by construction (e.g. neurons), most emergent aspects of these systems are still undefined and undiscovered (represented as “?”).

1.2 Top-down behavioral taxonomy guides bottom-up mechanistic explanation

The Newformer is a completely opaque result when considering benchmarks alone; it is simply better at doing what we want it to do than Transformer models (Vaswani et al., 2017) were before. However, as shown in Figure 1, a hierarchical taxonomy of LM behavior can guide our investigation of the Newformer, leading to questions such as:

1. What behaviors do the Transformer models and the Newformer model share, e.g., does the Newformer also repeat phrases more often than as seen in the training data?
2. Do they exhibit similar behaviors in the same contexts, e.g., does the Newformer need fewer input-output demonstrations to exhibit in-context learning at peak performance?
3. Do high-level behaviors decompose into the same lower-level behaviors or does the Newformer use different mechanisms to express them, e.g., when the Newformer is used for paraphrasing does it also tend to exactly copy the input?

Without such behavioral categories, we risk investigating the wrong direction when we try to interpret

models, because we do not know what phenomena we are trying to explain in the first place.

Observed behavior can tell us where to look for bottom-up explanations. Al-Rfou et al. (2019) observed emergent copying behavior in Transformer Language Models (LMs), paving the way for the discovery of *copying heads* that make copying possible. Characterizing copying heads led to the discovery of *induction heads* (Elhage et al., 2021; Olsson et al., 2022): Transformer heads that are capable of copying abstract representational patterns in previous layers and appear to be responsible for in-context learning. Olsson et al. (2022) show that induction heads exhibit a variety of pattern matching behaviors that are still not fully catalogued.

Attempting to explain neural networks bottom-up without being guided by behavior can make it difficult to interpret results. For example, many works that identify anisotropy in the embedding spaces of large LMs diagnose this as a deficiency, and attempt to fix it (Wang et al., 2020; Ethayarajh, 2019; Gao et al., 2019). However, recent work suggests that this anisotropic property may not actually limit expressivity (Biś et al., 2021), may be a result of the transformer architecture specifically (Godey et al., 2023), and may actually be helpful for language models (Rudman and Eickhoff, 2023).

1.3 The Transformer is the old Newformer

A moderate Newformer event has occurred at least once before with LMs: the switch from Recurrent Neural Networks (RNNs) to Transformers. Despite many partial explanations (Lakretz et al., 2019; Olah, 2015; Hochreiter and Schmidhuber, 1997), we still lack an explanatory theory of how LSTM (Hochreiter and Schmidhuber, 1997) LMs such as ELMo (Peters et al., 2018) worked—what behaviors they could and could not capture, how these composed, etc.—even as they were replaced by models like BERT (Devlin et al., 2019) with similar use cases, but completely different architectural details. This does not bode well for the introduction of something like the Newformer which is significantly farther from the Transformer than the Transformer is from the LSTM.

On their own, bottom-up methods do not transfer well to new systems: analysis techniques that relied on mutated state and gating in RNNs, such as visualizing gating mechanisms (Karpathy et al., 2015), are not applicable to Transformers. Interpretation methods for Transformer models (Weiss et al., 2021; Rogers et al., 2020), such as those that use attention, are unlikely to transfer over to the Newformer which breaks many previously immutable assumptions.

This suggests the value in doing more interpretation work that treats models like *black-boxes*, as if we do not have access to their internal mechanisms. There is growing interest in looking at NLP systems as black-boxes (Bastings et al., 2022; Ribeiro et al., 2020; Linzen et al., 2019), though much of this work still uses intermediate outputs—such as embeddings—rather than directly analyzing behavior in the output space models are trained to fit. Truly black-box methods can help insulate our analysis from change, giving us an anchor point that will always be testable on models that use the same modality (e.g., text, speech, images). Belinkov and Bisk (2017) show that neural machine translation systems are brittle to both natural spelling errors and synthetic character-level noise. This observation can be extended to ask: Is the Newformer robust to the same kinds of noise? Up to what threshold? Does the noise appear to be localized in the brittleness of tokenization, as was the case for Transformer-based systems (Provilkov et al., 2020)? Developing a rich inventory of such tests would give us a universal scaffolding for analyzing any Newformer the moment it is discovered.

1.4 Are we there yet?

Deciding whether a Newformer-like event has *truly* happened is an unresolvable question. New models are always partially derivative, and new (possibly artificial) axes can always be invented where they are worse or better (Wolpert and Macready, 1997). Yet three years on it is still infeasible for most labs to train a GPT-3 (Brown et al., 2020) level LM, costing approximately half a million dollars in compute alone for private companies—with engineering teams—to produce a similar model (Abhinav Venigalla, 2022). Thus it seems that the gap for training is only growing wider as ChatGPT (Schulman et al., 2022) and GPT-4 (OpenAI, 2023) become commonplace in research (Yang et al., 2023; Zhang et al., 2023), production (Eloundou et al., 2023; George and George, 2023; Ray, 2023), and even model evaluation (Liu et al., 2023b; Zheng et al., 2023).

Unlike the Newformer these models were never released, are frequently deprecated (OpenAI, 2023), change from day to day (Perry, 2023; Southern, 2023), and are known to be unstable over theoretically deterministic queries (Deng, 2023). Yet, the open source community has caught-up quickly (Alizadeh et al., 2023; Mukherjee et al., 2023; Gunasekar et al., 2023, *inter alia*) helped by industry labs' open-sourcing efforts (Touvron et al., 2023; Almazrouei et al., 2023; Stability AI, 2023, *inter alia*), and new finetuning techniques (Taori et al., 2023; Vicuna Team; Dettmers et al., 2023).

However, the question still remains: how should we explain models not everyone can train? Models that are so arduous, slow, and expensive to train that we will likely never ablate all the necessary variables needed to study them properly?

This leaves us with mere behavior. We generally think of there as being two different kinds of behavior: the neural behavior of different activations in models and the “output” of the model in the form of human media (e.g., text, images, videos, etc.). Most methods of explaining models focus on the former: trying to explain why neural activations cluster into certain patterns and trying to understand what those patterns mean about the output.

We argue that not enough attention has been given to formalizing the latter: *what* models are doing in the first place, in terms of regularities in their outputs. Without such a formalization, bottom-up methods will have a much harder time deciding what precisely to explain, and what is simply noise.

2 The Behavioral Bottleneck

How do we avoid proposing a new explanation for every exhibited difference? Surely we do not believe that we need a benchmark for every prompt that elicits slightly different behavior from a generative model? One solution is to propose many possible mechanisms, but make it an explicit research agenda to discover *the most parsimonious explanation*, a concept visualized in Figure 2. In other words, we want to be able to predict the aspects of text we care about (e.g. factuality) with the simplest rules possible. We briefly formalize this concept in §2.1, but the bulk of this paper concerns the *need* for this new research focus and the perspective it yields.

Thousands of papers observe behavioral tendencies in models, such as the ability of a pretrained Transformer to copy from the input context (Elhage et al., 2021; Al-Rfou et al., 2019), which we will adopt as a running example. To understand models better, we must rigorously describe (1) what *aspect* of generative behavior a given mechanism predicts (e.g. repetition, copying from the training set, etc.) and (2) how much of the *information* in the output space of the model such predictions explain (since most will not predict 100% of what a model emits).

Figure 2 serves as a visual map of how we might explain models via behavior. On the top level we have a huge diversity of benchmarks that currently exist, and the even larger number that may one day exist. On the bottom we have the mathematical abstraction that describes the space of all possible models. Clearly both of these represent many more possibilities than is useful as an explanation or than is *necessary* to explain specific facets of model behavior. The intermediate levels, then, deal with simplified metamodels, i.e., models of the underlying generative model that are less explanatory, but still allow us to interpret or theorize around models.

2.1 A working definition of “behavior”

Fong and Vedaldi (2017) state that: “An explanation is a rule that predicts the response of a black box f to certain inputs.” We think of a *behavior* as an explanation of limited aspects of a model, a concept we briefly formalize. We make reference to this formalization sparsely throughout the rest of the paper, as the argument can be understood without it, and we stress that the problem we are facing is more fundamental than a missing formalism.

Given a generative model from one input medium \mathcal{X} (e.g., strings composed of at most 2048 tokens) and a source of randomness \mathcal{R} to an output medium \mathcal{Y} (e.g., 512x512 pixel images):

$$\mathcal{M} : \mathcal{X} \times \mathcal{R} \rightarrow \mathcal{Y} \quad (1)$$

we can define a behavior as a function from the same input medium to a feature set \mathcal{F} :

$$\mathcal{B} : \mathcal{X} \rightarrow \mathcal{F} \quad (2)$$

For instance, \mathcal{M} may be a general purpose text-to-image model trained on scraped data, while \mathcal{B} may map a string $x \in \mathcal{X}$ to a probability that an image $\mathcal{M}(x)$, contains at least one dog. Or \mathcal{X} and \mathcal{Y} may both be Unicode strings, in the case of an LM, with \mathcal{B} being a binary prediction as to whether $\mathcal{M}(x)$ will eventually get caught in a repetition loop (Holtzman et al., 2020).

Our goal in proposing behaviors is to *explain* the underlying model using rules that capture model tendencies. Behaviors are explanatory to the extent that they give us information about the application of the model \mathcal{M} under distributions $\mathcal{D}_{\mathcal{X}}$ and $\mathcal{D}_{\mathcal{R}}$ over \mathcal{X} and \mathcal{R} , which we collectively refer to as \mathcal{D} for brevity. We can formalize the notion of “giving us information about the application of the model” through the mutual information:

$$I_{\mathcal{D}}(\mathcal{M}(\mathbf{X}, \mathbf{R}); \mathcal{B}(\mathbf{X})) = \quad (3)$$

$$H_{\mathcal{D}}(\mathcal{M}(\mathbf{X}, \mathbf{R})) - H_{\mathcal{D}}(\mathcal{M}(\mathbf{X}, \mathbf{R})|\mathcal{B}(\mathbf{X})) \quad (4)$$

where \mathbf{X} and \mathbf{R} are random variables drawn from \mathcal{X} and \mathcal{R} according to $\mathcal{D}_{\mathcal{X}}$ and $\mathcal{D}_{\mathcal{R}}$, and H is the entropy: $H(\mathbf{Y}) = \mathbb{E}_{\mathcal{D}}[-\log p_{\mathcal{D}}(\mathbf{Y})]$ for a random variable \mathbf{Y} . The mutual information is a direct measure of *how many bits of information we learn about one variable given another*, so this formulation directly tells us how much a behavior reveals about expected model output.

We call setting $\mathcal{D}_{\mathcal{X}}$ to be uniform over \mathcal{X} the “mechanistic distribution.” In this case, the mutual information is unrelated to the expected distribution of inputs in the wild, but is instead representative of how well we can model *any* input to \mathcal{M} . For instance, explaining an LM under the mechanistic distribution would require a behavior that predicts aspects of the LM’s output accurately even for long strings of gibberish. This may be difficult, since we often use human linguistic features to make predictions about model outputs, but such behaviors are closer to the notion of mechanistic interpretability

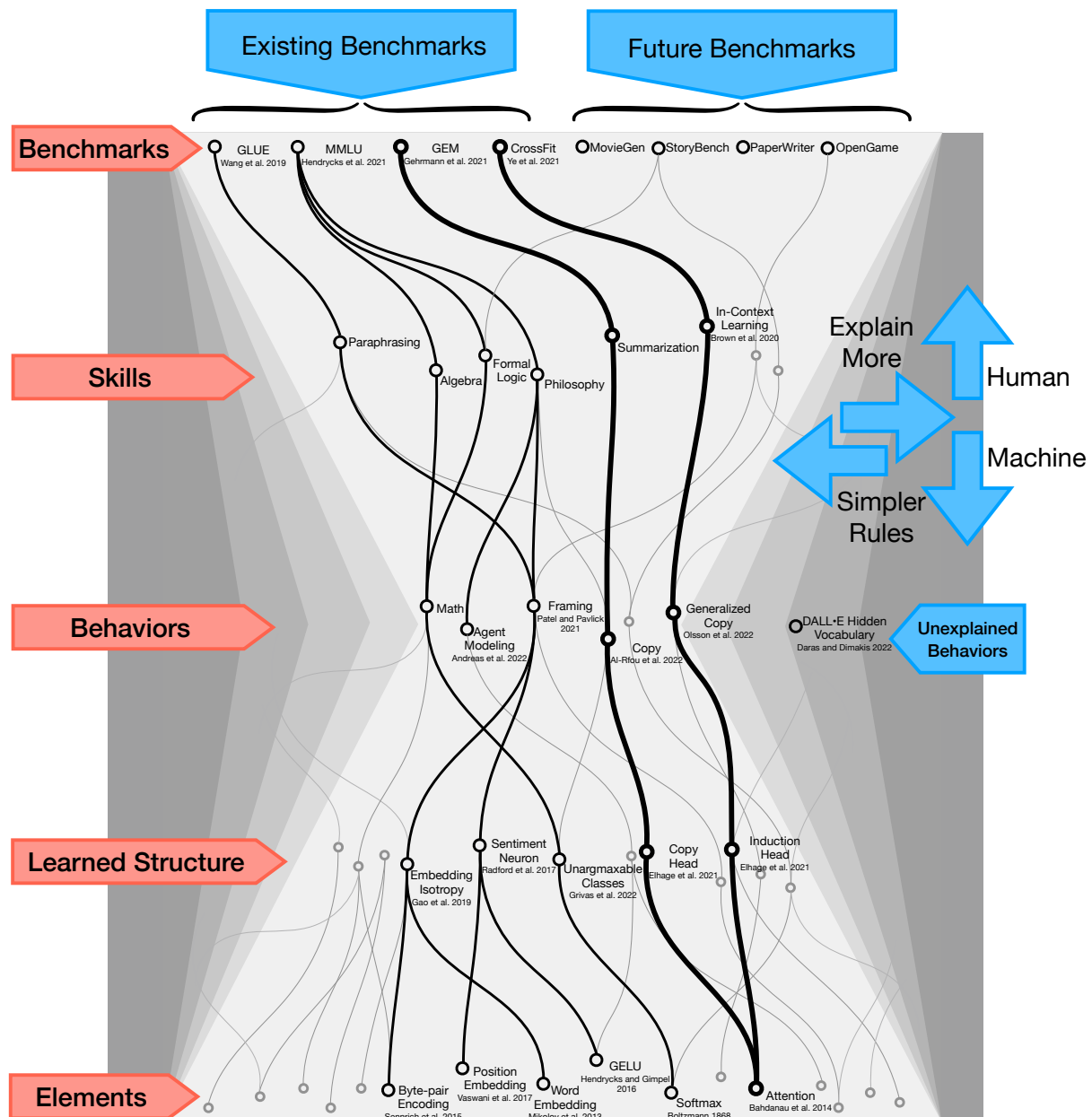


Figure 2: A visual representation of different aspects of models, shown from the basic elements of models on the bottom up to the benchmarks we are attempting to solve. Nodes represent invented and discovered aspects of models. The highlighted subgraph captures the concepts that we might want to use for understanding the phenomenon of “copying” in Transformers, when models generate sequences that appear in their local context window, a behavior that serves as a running example in this paper.

We might start out by noticing that Transformer models have higher scores on GEM (Gehrmann et al., 2021) (Benchmark), especially on summarization-like tasks (Skill). Inspecting the data generated by the models of interest, we might notice one of the qualitative differences separating Transformer models from other models is the ability to correctly use novel entities (Al-Rfou et al., 2019) (Behavior). We might ask why this is, embarking on an empirical study of when networks develop the ability to copy, as Elhage et al. (2021) did, discovering specific attention heads served as *copying heads* (Learned Structure) supported by certain Elements). This led to other discoveries such as *induction heads* (Elhage et al., 2021; Olsson et al., 2022) (Learned Structure), which were found to perform a kind of *generalized copying* that supports inference-time pattern recognition (Behavior), e.g., for In-Context Learning (Brown et al., 2020) (Skill), leading to better results on fewshot benchmarks such as CrossFit (Ye et al., 2021) (Benchmark). Research can proceed by observing high-level behavioral regularities, explaining them via the tendencies of the model, and using this to achieve clarity about other observed behaviors.

that tries to fully reverse engineer the model being studied (Olah, 2022).

If \mathcal{M} ignores its source of randomness, i.e., $I(\mathcal{M}(\mathbf{X}, \mathbf{R}); \mathbf{R}) = 0$ —as is the case for deterministic models such as a greedy-decoded LM—then the most explanatory behavior is simply $\mathcal{B} = \mathcal{M}(\mathbf{X}, r)$ for any $r \in \mathcal{R}$. This is a degenerate behavior, in that it is very explanatory, but has not brought us any closer to explaining \mathcal{M} . Therefore, we would like behaviors that are not just very high mutual information with the model, but also *point to predictable regularities* in \mathcal{M} , especially in a way that allows us to build up new hypotheses about it. Much has been written about what makes an explanation useful (Lipton, 2018; Jacovi and Goldberg, 2020; Chen et al., 2023, *inter alia*), and reviewing these desiderata is out of scope for this paper.

The mutual information $I_{\mathcal{D}}(\mathcal{M}(\mathbf{X}, \mathbf{R}); \mathcal{B}(\mathbf{X}))$ can also be viewed as *how much a behavior allows us to compress the output of a model* under a distribution \mathcal{D} , e.g., a distribution of articles for a summarization task (and a random number generator \mathbf{R}). This is because any bits of information revealed by one variable, can be used to compress the other, under a proper coding scheme (Cover, 1999).

The concept of Minimum Description Length (MDL) has been used as an information theoretic criterion for finding good hypotheses (Grünwald, 2007). Essentially, it suggests an extension to Occam’s razor (Barry, 2014): that we should favor explanations that are simple to describe and explain the object under study the most. We can formalize this notion for behaviors, via an encoding scheme \mathcal{C} that represents behaviors \mathcal{B} and outputs $y \in \mathcal{Y}$ as binary strings of variable (but finite) length $s \in \mathcal{S}$, where $|s|$ is the length of a string s . A naïve MDL objective would then be:

$$\operatorname{argmin}_{\mathcal{B}} |\mathcal{C}(\mathcal{B})| + \sum_{x \in \mathcal{X}} \mathbb{E}_{\mathcal{D}_{\mathcal{R}}} [|\mathcal{C}(\mathcal{M}(x, \mathbf{R})|\mathcal{B}(\mathbf{X}))|] \quad (5)$$

However, this would not suit our general objective: we do not necessarily wish to encode *all possible data a model could produce*, especially since most models have huge output spaces of largely low probability density. Instead, we would like to quantify the information behaviors can save us under $\mathcal{D}_{\mathcal{X}}$.

Model Generalization

from training data to the underlying distribution



Metamodel Generalization

a second model predicting model outcomes, from one part of the generative distribution to another

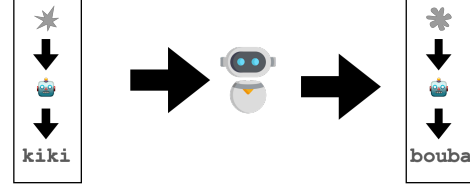


Figure 3: If we cannot train models easily, but those models are sufficiently general and useful, we can predict what models can and can’t do, rather than what a model trained differently would do. The kiki/bouba categorization is a cross-culturally robust linguistic-conceptual mapping in humans (Ćwiek et al., 2022).

To capture the idea “how much space does \mathcal{B} save us under $\mathcal{D}_{\mathcal{X}}$ we can use:

$$\operatorname{argmin}_{\mathcal{B}} |\mathcal{C}(\mathcal{B})| + \alpha H_{\mathcal{D}}(\mathcal{M}(\mathbf{X}, \mathbf{R})|\mathcal{B}(\mathbf{X})) \quad (6)$$

where we replace the second term with the conditional entropy $H_{\mathcal{D}}$, since this describes the minimum number of bits that could be used to represent the information encoded (Cover, 1999). This can be interpreted as, “we would like behaviors that on average, save us more space in terms of encoding the possible outcomes of a model than they take to describe.” α allows us to trade-off how much we weight the representation of the behavior vs. the outputs of a model, where larger values of α may be appropriate if we are dealing with a many outputs, making the bits saved by way of conditioning on the behavior more pertinent.

Overall, we seek to find behaviors that are both explanatory and simple to describe. We can think of this as attempting to find a *metamodel*: models that are designed or trained to predict another model’s behavior (Barton and Meckesheimer, 2006), as illustrated in Figure 3. This suggests we want to find *behaviors that transfer over different contexts* so we can predict where models will be useful and where they will break down.

2.2 Can benchmarks discover new behavior?

In general, discrepancies in performance between benchmarks can *hint* at potentially new behavior, but they cannot discover behavior we have not yet observed. Given the diversity of NLP benchmarks, it is likely that the Newformer (§1) will perform drastically different on certain pairs of benchmarks we believe to be related, e.g., the same task in two different domains. This is a useful signal for where to inspect behavior, but benchmarks alone cannot reveal new abilities, underlying mechanisms, or shortcut heuristics the Newformer is relying on that *cause* a discrepancy in results and what else its effects are.

For example, it is very difficult to imagine how *prompting* (Radford et al., 2019b; Liu et al., 2021) could have been discovered via benchmarking. Finetuning a generative model, such as GPT-2 (Radford et al., 2019a), and doing well or badly at any number of benchmarks could not have revealed a model can be prompted with text that matched training data patterns, to elicit behavior such as summary generation via the string “TL;DR” or translation through formatting such as “French sentence: <source>\\English sentence:”. These discoveries are a result of inspecting the generative *behavior* of GPT-2, and only afterwards testing a perceived pattern on benchmarks.

How do we try to explain the behavior of models, once we know there’s a discrepancy we want to explain? Often we attempt to look at the qualitative differences between tasks a model is good or bad at, and come up with hypotheses for what the model is failing to do when it performs poorly, e.g., across different finetuning tasks (Li et al., 2021). While useful for coming up with hypotheses, using benchmarks as evidence of behavior requires care, because it is often unclear what a given benchmark is actually testing. Rohrbach et al. (2018) show that image captioning systems hallucinate objects not present in the scene, and are unintentionally *rewarded* for doing so by standard metrics, by capturing phrasing and *n*grams of reference captions better when hallucinating. Liao et al. (2021) describe a detailed framework for assessing benchmark validity and note the complexity of ensuring benchmarks test what we would like them to. Thus, since we often do not know precisely what behavior benchmarks test, they might indicate what contexts to examine the Newformer in, but not precisely what it does.

2.3 Can benchmarks characterize behavior?

Consider standardized tests for humans—such as the SAT (for College Admission Counseling, 2008) or the NCEE (百度百科, 2022)—while the debate about how much these tests tell us is heated, there is little resistance to the statement: *test scores do not fully describe human behavior*, even within the subjects they test such as mathematics and biology.

Performance data about a bicycle is not sufficient to reverse-engineer its gear system. Even with perfectly valid benchmarks, the subspace of benchmark performance is not descriptive enough to characterize behavior. As we greatly increase the number of benchmarks, we are left with the problem of determining precisely how benchmarks overlap and differ in a way that characterizes behavior (Figure 2). Because the space of benchmarks is limited, as we test for human-desirable skills and human-interpretable pitfalls, discovering novel behavior in non-human systems is difficult.

Measuring systems only for their expected purposes makes it difficult to disentangle component behaviors that allow models to produce the desired or undesired outputs, as failure under distribution shift often reveals. For example, neural machine translation often outputs completely irrelevant translations under domain shift (Wang and Sennrich, 2020; Müller et al., 2020). This is exacerbated by the fact that most generative models are not trained with a precise purpose in mind.

Imagine testing whether an LM can summarize an article. In order to summarize an article a requisite skill required by models is *copying*, because novel entities are constantly appearing, but need to be referenced in the summary. See et al. (2017) add a copying mechanism to an RNN in order to improve its copying ability for summarization. If we were to only look at performance on summarization, we would be unlikely to notice whether copying was happening or not directly—only whether performance is hitting certain desired levels.

Benchmarks are, by necessity, scoped to certain contexts that are presumed to test for certain behaviors—but they do not directly tell us what patterns the model is exploiting to solve the task, as Liao et al. (2021) point out. This was a hard-learned lesson in many benchmarks, such as when it was discovered that SNLI (Poliak et al., 2018; Gururangan et al., 2018) could be solved with *hypothesis-only* systems that only use a subset of the information that was supposed to be necessary to the task.

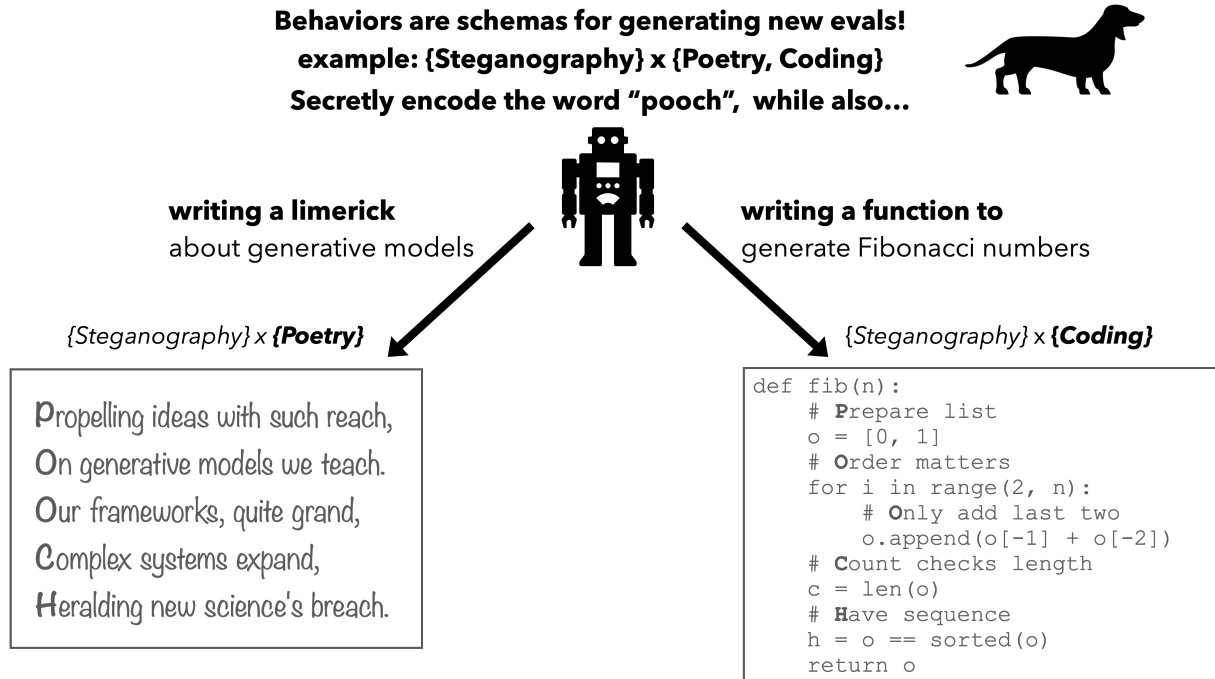


Figure 4: An example of how behaviors can be used to create new evaluations. These examples were generated from GPT-4, but required significant human curation, suggesting that Thought Experiment B has not yet occurred.

2.4 Behaviors: building blocks for evaluation

Benchmarks are still the best solution for coordinating cross-lab experimental *comparisons*, and we expect them to continue to be useful in that respect indefinitely. However, to answer “What strategy is the Newformer using for this task?” and “What failure modes should we expect?” and “What else do we expect the Newformer to be capable of?” we cannot use benchmarks alone to guide where we inspect model behavior, nor as a means to define it.

Instead, we propose an increased focus on behavior, because we believe the science of generative models is currently held back by insufficient understanding of *what* models are doing in general, rather than *how well* models perform on specific tasks. These are highly related to each other, and we can think of *behaviors as building blocks for evaluation*. Consider the following thought experiment:

A new LM is released with many of the expected capabilities, such as basic arithmetic and basic translation, but another interesting behavior is noticed and hypothesized: when asked properly in natural language, the model can steganographically encode complex hidden messages while completing other tasks.

(B)

When this LM is released it is unlikely there are any benchmarks that test this particular capability. While we could design a specific benchmark for this behavior, this would be somewhat counter-productive: what we really care about is the *Cartesian product* of this behavior and other tasks that we were already testing. In this sense, behaviors are the building blocks for benchmarks.

As [Chang and Bergen \(2023\)](#) point out in their survey of behaviors, researchers are often surprised by the outputs of the models they work with; it should not surprise us that we cannot premeditate benchmarks to capture behavior when modeling improvements have outpaced our ability to be exposed to generated data. One way to be more nimble to new behaviors, is to directly measure behaviors we expect ([Jain et al., 2023](#)), flagging unexpected combinations for inspection.

On the surface, it might seem that naming behavioral categories such as “copying” or “in-context learning” is just as liable to obsolescence as any other analysis. What should we do if the Newformer does not exhibit these behaviors? We argue that this is a very unlikely scenario: as long as we are attempting to train models to mimic human understandable phenomena, there will be human perceivable patterns that we expect models to mimic as well.

3 Generative Models as a Complex Systems Science

While the Newformer (§1) is a thought experiment, it is representative of many facets of research regarding generative models today; suddenly, focus has shifted to searching for *emergent behaviors* in large and often inscrutable models. Larger pre-trained models continue to be trained and continue to perform better (Schulman et al., 2022; OpenAI, 2023; Pichai, 2023, *inter alia*). While efforts to release models (Touvron et al., 2023; Almazrouei et al., 2023; Stability AI, 2023) and involve more researchers in model training (BigScience, 2022) can increase transparency and provide more information, it is well beyond the resources of the vast majority of labs to train. Efficiency breakthrough are likely to be exploited to further increase model size and feed into the same problem they were meant to solve.

Thus, it seems likely that training and re-training models is no longer the path towards understanding them for the vast majority of researchers. In many fields the creation of what it studies is impossible, from biology to astronomy. Many of these fields are *complex systems sciences*, in that they focus on the question illustrated in Figure 5: how do the macro-level behaviors we observe (life, black holes, etc.) arise from the micro-level units we understand better (chemicals, regular matter, etc.)?

In other words, we suggest studying *generative models themselves not just generative modeling*.

3.1 What is a complex system?

Newman (2011) establishes a working definition:

[A] system composed of many interacting parts, such that the collective behavior of those parts together is more than the sum of their individual behaviors. The collective behaviors are...“emergent” behaviors, and a complex system can thus be said to be a system of interacting parts that displays emergent behavior.

Recently, interest in emergent behavior has grown in NLP (Bubeck et al., 2023; Wei et al., 2022; Teehan et al., 2022; Manning et al., 2020, *inter alia*), though it is usually defined, in terms of scaling over model parameters, dataset size, or computational power. We rely on a much simpler definition:

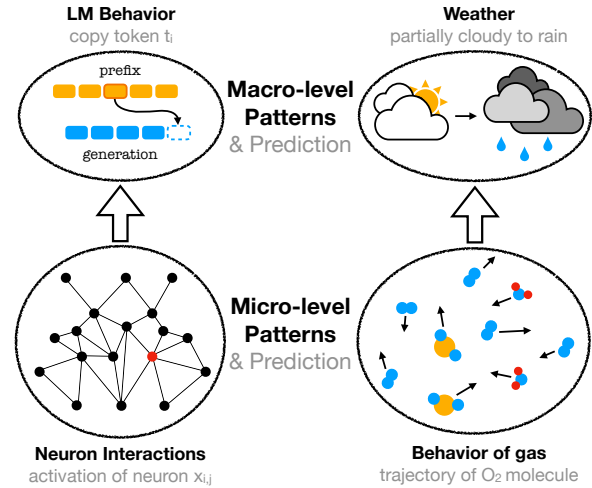


Figure 5: Complex systems are characterized by two or more levels of regularity: a micro-level in which local interactions are at least partially predictable and a macro-level in which many local interactions collectively exhibit recognizable patterns. *Emergence* describes how macro-level regularity is hard to predict in advance from comparatively well-understood micro-level dynamics.

Emergent behaviors are system level behaviors that are hard to predict from the dynamics of lower level subcomponents.

For instance, the ocean is a complex system. We can understand many properties of individual water molecules, e.g., H_2O has a partial positive and negative charge in certain places due its composition, but the aggregate properties of *water* as a collective whole exhibits predictable properties such as waves. It is difficult to predict the properties of water from H_2O because “the interactions of interest are non-linear...[yielding] levels of organization and hierarchies—selected aggregates at one level become ‘building blocks’ for emergent properties at a higher level, as when H_2O molecules become building blocks for water.” (Holland, 2014)

Similarly, we understand the basic mechanical properties of LMs at the neuronal level, e.g., we have a perfect understanding of how to predict what any individual neuron will do given arbitrary inputs (by construction), but we also notice patterns at the level of *model behavior*, e.g., the emergent copying behavior, which is observed in both Transformer models (Al-Rfou et al., 2019; Khandelwal et al., 2019) and LSTM models (Khandelwal et al., 2018). In the face of new behavior that a model such as the Newformer might exhibit, we would be even less certain of how lower-level system components add up to observed responses.

3.2 Emergent behaviors in LMs are discovered, not designed

Neural architectural elements (e.g. position embeddings) and training methods (e.g. masking strategies) deeply affect the resulting model but do not fully explain behavior. We often fail to create the behavior we attempted to engineer into an architecture *and* discover new, unintentional behavior.

Many architectures have been designed to make use of longer context (Yu et al., 2023; Beltagy et al., 2020; Child et al., 2019, *inter alia*), but evidence suggests that these models often do not make use of the long-term dependencies that they intended to capture (Liu et al., 2023a; Sun et al., 2021; Press et al., 2021). Inversely, BERT was shown to capture much of the functionality of a knowledge base without task-specific training (Petroni et al., 2019).

To illustrate the difference between *designing* and *discovering* behavior, let us return to our running example of the copying behavior, where models produce a span that was in their input. A classic example of designing behavior is pointer-generator models (See et al., 2017), in which a specific, discrete mechanism was added to encourage a certain behavior: copying. Transformers, on the other hand, were designed such that computation at a given time-step could *attend* to any previous time-step that was included in the context window. This intentionally removed the recurrence in architectures such as LSTMs (Hochreiter and Schmidhuber, 1997) and GRUs (Cho et al., 2014) in order to increase efficiency on highly parallelizable hardware such as GPUs and TPUs. A side-effect of this change was the emergent behavior of *copying* that arises directly from the Transformer architecture trained as a language model (Al-Rfou et al., 2019).

Instead of directly designing models for these purposes, we are now in the position of training general models with different structure and actively probing them for behavior. Using various data and masking strategies has produced models that can be controlled through different metadata (Keskar et al., 2019; Zellers et al., 2019; Aghajanyan et al., 2022), while instruction-tuning has shown that pretrained LMs can be finetuned for control (Mishra et al., 2022; Chung et al., 2022; Ouyang et al., 2022, *inter alia*) often with very limited data (Zhou et al., 2023; Dettmers et al., 2023; Taori et al., 2023).

This *discovery* process focuses on giving the model access to certain kinds of correlations, and then inspecting what model behavior emerges.

3.3 Neuronal explanations are limited by our understanding of behavior

It is difficult to explain how or why LMs produce their outputs without having a good description of *what* they do. Explaining behavior bottom-up, requires an understanding of what behaviors we are trying to explain. Mittal et al. (2018) note:

an emergent property of a system is usually discovered at the macro-level of the behavior of the system and cannot be immediately traced back to the specifications of the components, whose interplay produce this emergence.

This is the situation we find ourselves in with regards to large, pretrained models. We cannot, in general, predict how structure will form. While we can engineer systems with the hope of producing certain kinds of behavior, e.g., training on multimodal data to produce models that can draw inferences in ways that integrate paired text and images, this often does not produce the desired results (Ilharco et al., 2021; Parcalabescu et al., 2021).

Bottom-up investigation can reveal key properties of emergent organization within LMs, e.g. BERT replicates features of the classical NLP pipeline (Tenney et al., 2019). But when anomalous behavior is discovered, e.g., the DALL•E 2 hypothesized “hidden vocabulary” of invented words that correspond to specific image categories (Daras and Dimakis, 2022), it is difficult to investigate them with bottom-up tools until we reach a better understanding of what triggers them, what their scope is, etc. There have been attempts to reject the hidden vocabulary hypothesis (Hilton, 2022), but it is a very difficult hypothesis to rebut from first principles: what tests reject the hypothesis “DALL•E 2 has a hidden set of vocabulary with clear and consistent meaning” rather than “this specific mapping from the vocabulary to features isn’t correct”?

This is similar to trying to research organic chemistry without knowledge of biology: it is certainly not impossible, but without high-level guides to the kind of structure one is expecting, the search space is huge and it is difficult to know where to look. Our lack of a behavioral taxonomy hampers research into internal structure, especially in models that break current assumptions such as the Newformer, as it is significantly more challenging to probe for structure without knowing what patterns in the outputs hint at the presence of structure.

3.4 Access is not a silver bullet

Consider the following thought experiment:

Tomorrow, all industry labs publicly
release all of their pretrained models (C)

Despite the fact that this would doubtlessly help us understand the basic properties of a given model such as ChatGPT, e.g., how large it is, we would still have significant obstacles on the way to explaining why ChatGPT is capable of writing short stories for almost any given prompt.

Indeed, the problem with answering the question of “How can a language model write a story?” has much less to do with language models and much more to do with the fact that we are currently incapable of answering the question “How can x write a short story?” for any value of x . We find ourselves in the strange position of being able to train models we do not fully understand *for tasks we do not fully understand or anticipate in advance*.

The key to answering this question is to ask: what kind of explanation would satisfy us? For instance, when it comes to LMs, one explanation is that models are simply reconstructing long sequences from the training set and stitching them together. While a significant amount of memorization is taking place (McCoy et al., 2023; Carlini et al., 2023; Lee et al., 2022, *inter alia*) models appear to be able to generate data that is not a trivial recombination of the training data (Bubeck et al., 2023; Tirumala et al., 2022; Olsson et al., 2022).

The goal, then, should be to build up the case for a reasonable hypothesis that explain the breadth, depth, and (most importantly) mistakes models make when executing a complex task. However, we do not want a new explanation for every new task, which is precisely why we argue for the formalization and study of *behaviors* that describe the underlying strategies of models.

While model access would not directly solve these problems, we *do* believe that open-source models are a necessary prerequisite to this research program, for reasons outlined in §4.3.

3.5 (Generated) data represents behavior

Behavior in large pretrained models is nothing more than the answer to the question “How can we characterize the distribution of data this model generates?” Aspects of the training data such as the presence of multiple languages (Blevins and Zettlemoyer, 2022; Lin et al., 2021) or the number

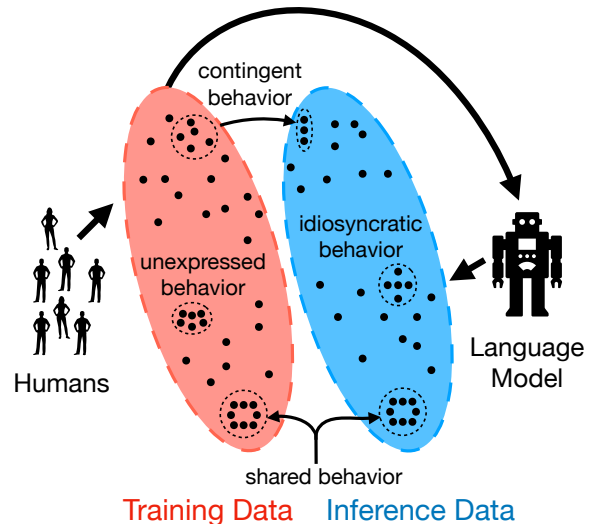


Figure 6: Generative language models are trained to capture the distribution of training data, then exhibit behavior in model outputs, i.e., *inference data*. See §3.5 for examples of the different behavioral mappings.

of repeated documents (Kandpal et al., 2022; Lee et al., 2022) in the training set have been shown to be explanatory of zero-shot translation abilities and model tendency to leak training data, respectively.

Figure 6 visualizes what kind of behavioral mappings we can explore with data-based explanations. *Shared behavior*—patterns that are found in both the training and inference data (the outputs of the model)—are the simplest to search for, because they only require finding a specific behavior in the training or inference data and then looking for it in the other. For instance, the prompting behavior discovered in GPT-2 that causes summaries to be generated when “TL;DR” is placed after an article is an example of shared behavior. Idiosyncratic behaviors describe behaviors that don’t appear to be caused directly by the training data at all, e.g., zero-length translations in many large models (Stahlberg et al., 2022; Shi et al., 2020; Stahlberg and Byrne, 2019). Perhaps the most difficult to find behavioral mappings are those for which behavior in the corpus yields different behavior in the model, *contingent behavior*, as is hypothesized to be the case for DALL-E 2’s “hidden vocabulary”: nonsense words that appear to consistently lend certain meanings to produced images (Daras and Dimakis, 2022). Finally, unexpressed behavior is observed in the training data, but not in the inference data, such as long-term consistency in story telling (Xie et al., 2023; See et al., 2019) that models have yet to properly mimic for very long documents.

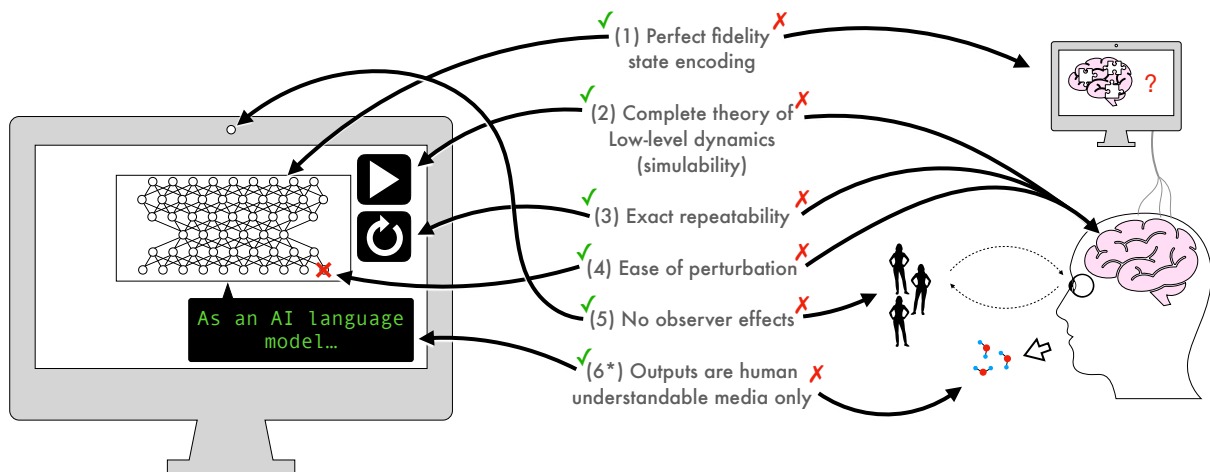


Figure 7: Visual representation of the advantages generative model researchers have over researchers that study the main other media generating system on earth: human beings.

4 A Different Kind of Complex System

One reasonable worry is that taking on the complex systems lens will be fruitless because studying complex systems is a very difficult task, and we are not equipped to tackle such a hard problem.

In fact, compared to other complex systems, such as the brain, understanding current generative models is an immensely *easier* challenge, and can help us develop tools for the future. Turning our attention to “What, precisely, do language models do?” over “What is the best recipe for training large models?” we can take full advantage of the *complete simulability* of generative models. In the long run, it seems it will become more difficult to address the latter question coherently without better answers to the former.

4.1 Two kinds of complex systems simulations

Complex systems theory is divided between two basic approaches. The first involves the creation and study of simplified mathematical models that, while they may not mimic the behavior of real systems exactly, try to abstract the most important qualitative elements into a solvable framework from which we can gain scientific insight...The second approach is to create more comprehensive and realistic models, usually in the form of computer simulations, which represent the interacting parts of a complex system, often down to minute details, and then to watch and measure the emergent behaviors that appear. (Newman, 2011)

At first glance, generative models can largely be described as the second complex systems approach: we train models to capture properties of the natural distribution of human media, such as internet text, images, videos, etc., and then attempt to study the emergent effects. Yet, this would be a mischaracterization of what, e.g., language models do: we do not expect that language models learn language the way a human does nor create languages the way the human species did.¹

Instead, the triumphs of generative models are the result of emergent behavior within computational models trained to predict very general objectives. Many have been surprised that by learning on massively more data from a given medium than a human is ever exposed to, generative models can learn uncannily human patterns from simple, passive word prediction, denoising objectives, etc.

Generative models are certainly a computational simulation, but they are a simulation of an entire medium rather than a singular process we have isolated. We suggest thinking of generative models as a different kind of complex system, where *underlying patterns of a medium are learned by a model through optimization, and we then look for those patterns within the model*. Below we list ways in which this discovery process is made easier because our system of interest is the computational model itself, rather than a naturally arising system.

¹Though this certainly describes facets of certain subfields such as emergent communication (Lazaridou and Baroni, 2020), with recent work taking advantage of pretrained models (Steinert-Threlkeld et al., 2022), and developmentally plausible pretraining such as the recent BabyLM challenge (Warstadt et al., 2023).

4.2 Generative Models: the easiest complex system to study

- (1) **Perfect fidelity state encoding** Because neural networks are formal mathematical models, represented by code and parameters, there is zero *necessary ambiguity* in our representations. Imperfect data archiving and complex code bases often make it difficult to perfectly recover the formal model, but with sufficient effort, it is possible to store every bit of information about the state of the model at every computation step. We cannot track every neuron in a participant’s brain at every moment due to the limited nature of our measurement instruments, but we can perfectly record the state of an LM in order to look for and verify emergent behavior, without influencing the system as we note in Advantage (5).
- (2) **Complete theory of low-level dynamics** While Advantage (1) establishes that we can perfectly store the state of a generative model at any given moment in time, Advantage (2) notes that we have a perfect, mechanistic, and deterministic understanding of how one state of a model evolves into the next, unlike in physical experiments. Artificial neural networks do not need to be simulated, they are defined in a medium for simulation: executable code. Unlike in physics, where centuries of research have been spent chasing the bottom of the chain of causation, we *begin* with the base-level causal structure of the model. This does not follow directly from (1). It is possible to imagine a scenario where every static state is recordable, but where the rules that govern the changes in states are hard to discover, e.g., the problem of learning video game dynamics from pixels (Hafner et al., 2019). In practice, nondeterminism exists in certain fast computations (Morin and Willetts, 2020), but this can be removed at the cost of speed.
- (3) **Exact repeatability** Directly entailed by (1) and (2) is the fact that experiments can be repeated *exactly*. An algorithm that uses randomness to generate text, may generate different text on a second run, but as long as the probabilities of different tokens are recorded the likelihood of that text (and of alternative branching paths) can be verified to be exactly the same. A psychologist who conducts a

study twice will almost never get results that are exactly the same, simply because sample differences and unmeasured variables have to be accounted for. We distinguish repeatability from the broader notion of *replicability*, which also includes replicating a study to the level of detail described by the authors, leaving room for both human and systematic error. With proper code, data, and model releases, many generative model experiments are exactly repeatable, allowing us to reach for a much higher standard for replicability.

- (4) **Ease of perturbation** We also have a complete description of *all possible models* given a certain setup, e.g., all possible combinations of weights for a given architecture. Combining this with Advantage (2), we can perturb a model of interest, and play out experiments with this new model *without destroying the original model*. Contrast this with studying human language production, for which most perturbations of the human brain are both unethical and illegal, partially because humans cannot be unperturbed. This allows for extremely targeted experiments, e.g., finding which weights in a network control a certain decision boundary.
- (5) **No observer effects**—a classic problem in many complex systems is that by attempting to make a measurement one changes the value being measured, e.g., Clever Hans a horse that could allegedly play chess, but was simply reading the audience’s reaction to possible moves (Prinz, 2006). In contrast, generative models do not distinguish between the same input given for different reasons or with different expectations by the experimenter. The caveat is that experimenters still control the input distributions to experiments allowing for systematic bias that accidentally leaks *experimenter expectations* (Rosenthal, 1976) to the model, as past research has consistently shown (McCoy et al., 2019; Poliak et al., 2018; Gururangan et al., 2018, *inter alia*). We must be careful about “tells” (Caro, 2003): stylistic and semantic artifacts that make it into the data which can give the model information the experimenter assume it does not have access to. Yet, the guarantee that the specific observer will not change the result is strong.

Advantages (1) and (2) allow us to completely remove any worry about hidden variables that may explain effects we attempted to explain through other means. (3) and (4), allow us to experiment freely, knowing that experiments and models that have been properly recorded are recoverable, leaving us free to perturb and explore the local neighborhood of similar models and setups. (5) partially relieves us of the fear of influencing the outcome through our means of observation, a key issue in many experiments involving language.

Another advantage, that does not apply to every generative model, deserves an honorable mention:

- (6*) **(Some) generative models exclusively output human understandable media** Many complex systems, such as cities and brains, produce human understandable media as some percentage of their output. Many generative models produce human understandable media as their *only* output, an enormous advantage for two reasons. First, humans are better suited to positing patterns in human understandable media than, say, subatomic particles. Second, *the uncanny valley effect* (Mori, 2012) allows us to see when patterns are “almost correct but not quite” much more easily in human-related artifacts. While we sometimes finetune models to produce outputs that are no longer human understandable, by and large current generative models operate entirely within human media—and we believe there is much that can be learned from this that will transfer over to generative models of other media.

Advantage (6*) is very special. By allowing us to take advantage of our intuitive understanding of media, it becomes easier to seek out the ways generated media diverges from the natural human media we are steeped in from infancy. Indeed, most named behaviors are failure modes, e.g. degeneration (Holtzman et al., 2020) or empty translations (Stahlberg and Byrne, 2019)).

Organic chemistry has given a great deal to biology, but is very much indebted to it as well. Our hope is that we can take inspiration from these other complex system sciences to start taking the problem of understanding *behavior* seriously, as a distinct abstraction that needs to be decomposed and theorized, while putting our enormous advantages to good use.

4.3 The necessity of open-source models

Most of these advantages rely on stable access to a consistent representation of a model, which is difficult to guarantee via a proprietary API.

- (1X) **Perfect fidelity state encoding** It is difficult to work with or guarantee saved state is persistent and untampered without direct access to said state. Even cryptographically signed state can be tampered *after* re-submission to an API for use there, making guarantees moot.
- (2X) **Complete theory of low-level dynamics** With only imperfect knowledge of an underlying model, researchers must make assumptions about low-level dynamics in a model that may only partially be true, or possibly even completely false.
- (3X) **Exact repeatability** In practice, it is impossible to guarantee that an API will not drift over time, something observed with even the apparent attempts at stable APIs in recent years (Deng, 2023).
- (4X) **Ease of perturbation** It is normally impossible to perturb a model through an API, though some APIs allow for finetuning and special versions of models. However, the real issue is that it is impossible to ensure that such perturbations do precisely what they are claiming to do to the model, without access to the model or even the model architecture in most cases.
- (5X) **No observer effects** Sadly, even though this is one of the greatest advantages of generative models, it is the one most destroyed by using models via APIs: companies consistently, and often silently, fix undesired (from the company’s perspective) behaviors in models (Wilson; Eliaçık; Kiho) so that testing a certain hypothesis tends to influence future tests.
- (6X*) **(Some) generative models exclusively output human understandable media** Without complete access to a model it is impossible to know if it doesn’t have other outputs (or inputs) that would help explain the model’s behavior more fully.

In short, without access to open-source models, these advantages are largely moot. However, the community has seen a consistent open-source releases of better generative models in many different

media (Rombach et al., 2021; Le et al., 2023; Luo et al., 2023). There is unquestionably lag in the capabilities between proprietary and open-source models, and this is out of necessity: open-source cannot outpace private industry when private industry controls most of the training resources and can build on top of anything open-source does. But the fact that open source often lags only a year or two behind in terms of capabilities, and the fact that private labs are often incentivized to open-source models as a recruiting and market strategy, suggests that open-source will continue to be a wellspring of fascinating generative models to study. Indeed, if all progress stopped now, we believe it would be decades before we finished cataloging all of the generalizable behavioral principles with the hundreds of large generative models that have already been released; perhaps our successes would encourage future open-source releases.

5 Conclusion

How should we study models of data, when we don't fully understand the models or the data? We should study them first by asking *what* models do, before attempting the more complicated *how* and the bottomless question of *why*?

In this paper, we presented a thought experiment: the Newformer, a model that would be impossible to study with many of the techniques we use to understand Transformer models today.

We argue that focusing on what *behaviors* explain its performance across tasks will lead us to a deeper understanding of generative models' tendencies and guide bottom-up mechanistic explanation, as well as forming building blocks for evaluations.

We discuss how generative models are well captured by the definition of a complex system, due to the emergent behaviors they exhibit. This separates generative models from traditional machine learning, where models often served as explanations via behaviors that were architected directly into them. This opens up the need for *metamodels* that help us predict regularities in generative model outputs in order to understand them better.

While the prospect of studying models we do not have a clear understanding of is daunting, we highlight advantages that generative models have over naturally arising complex systems. These advantages, however, require open-source models as a prerequisite, a point we emphasize as a necessity for conducting replicable science.

6 Limitations

We present one perspective on the kind of science NLP is becoming, and how we can leverage the complex systems lens in order to better explore the phenomena we find ourselves faced with: generative models we do not fully understand. We cite evidence from NLP publications, blog posts, and other media, but this necessarily does not capture the totality of perspectives.

Indeed, we purposefully avoid attempting any sort of survey of these issues, as this would involve citing thousands of papers and be a very unwieldy object. Instead, we attempt to form an argument as economically as possible, attempting to put forth a new set of goals and principles for how to study generative models given current progress.

We make comparisons with other sciences and cite sources from those sciences where appropriate, but are extremely limited in expressing many equally relevant connections and in fully exploring the connections we do mention. There is an enormous amount related to sister fields (e.g., cognitive science, linguistics, etc.), other sciences that study complex systems (e.g., chemistry, biology, etc.), and regarding more meta-science issues (e.g., complex systems theory, chaos theory, etc.) that we could not cover, and we do not in any way attempt to—giving a complete account of these connections is simply beyond the reach of any one work.

Finally, parts of our assessment is necessarily subjective. We attempt to lay out the evidence as we see it, tracing the connections we drew in order to describe a style of research that we believe is necessary to face the current challenges of our field. This seems especially pertinent in a time when most researchers cannot train large generative models from scratch, but are excited to contribute to their study. With evidence drawn from the literature, we describe the current research space as we perceive it, and our vision for where it might go. Our hope is that this will add to a discussion on what the study of generative models currently is and what we, as a community, would like it to become.

Acknowledgements

We thank Julian Michael, Dallas Card, Jared Moore, Daniel Fried, Gabriel Ilharco, Tim Dettmers, Ian Magnusson, Alisa Liu, and Kaj Bostrom for their insightful discussions and feedback.

References

- Linden Li Abhinav Venigalla. 2022. Mosaic LLMs (part 2): GPT-3 quality for <\$500k. <https://www.mosaicml.com/blog/gpt-3-quality-for-500k>. Accessed: 2023-7-7.
- Armen Aghajanyan, Dmytro Okhonko, Mike Lewis, Mandar Joshi, Hu Xu, Gargi Ghosh, and Luke Zettlemoyer. 2022. [HTLM: Hyper-text pre-training and prompting of language models](#). In *International Conference on Learning Representations*.
- Rami Al-Rfou, Dokook Choe, Noah Constant, Mandy Guo, and Llion Jones. 2019. Character-level language modeling with deeper self-attention. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 3159–3166.
- Meysam Alizadeh, Maël Kubli, Zeynab Samei, Shirin Dehghani, Juan Diego Bermeo, Maria Korobeynikova, and Fabrizio Gilardi. 2023. [Open-Source large language models outperform crowd workers and approach ChatGPT in Text-Annotation tasks](#).
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.
- Jacob Andreas. 2022. Language models as agent models. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5769–5779.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate.
- C M Barry. 2014. Who sharpened occam’s razor? <https://www.irishphilosophy.com/2014/05/27/who-sharpened-occams-razor/>. Accessed: 2023-7-7.
- Russell R Barton and Martin Meckesheimer. 2006. Metamodel-based simulation optimization. *Handbooks in operations research and management science*, 13:535–574.
- Jasmijn Bastings, Yonatan Belinkov, Yanai Elazar, Dieuwke Hupkes, Naomi Saphra, and Sarah Wiegreffe. 2022. Proceedings of the fifth blackboxnlp workshop on analyzing and interpreting neural networks for nlp. In *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*.
- 百度百科. 2022. 普通高等学校招生全国统一考试.
- Yonatan Belinkov and Yonatan Bisk. 2017. Synthetic and natural noise both break neural machine translation. *International Conference on Learning Representations*.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- BigScience. 2022. Bigscience model training launched. *BigScience Blog*.
- Daniel Biś, Maksim Podkorytov, and Xiuwen Liu. 2021. Too much in common: Shifting of embeddings in transformer language models and its implications. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5117–5130.
- Terra Blevins and Luke Zettlemoyer. 2022. Language contamination explains the cross-lingual capabilities of english pretrained models. *arXiv preprint arXiv:2204.08110*.
- Ludwig Boltzmann. 1868. Studien über das gleichgewicht der lebendigen kraft zwischen bewegten materiellen punkten [studies on the balance of living force between moving material points]. *Wiener Berichte*, 58:517–560.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrkke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. In *The Twelfth International Conference on Learning Representations*.
- Mike Caro. 2003. *Caro’s book of poker tells*. Cardoza Publishing.
- Tyler A Chang and Benjamin K Bergen. 2023. Language model behavior: A comprehensive survey. *arXiv preprint arXiv:2303.11504*.
- Chacha Chen, Shi Feng, Amit Sharma, and Chenhao Tan. 2023. Machine explanations and human understanding. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1–1.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. 2019. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*.

- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. In *Proceedings of SSST-8, Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation*, pages 103–111.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Thomas M Cover. 1999. *Elements of information theory*. John Wiley & Sons.
- Aleksandra Ćwiek, Susanne Fuchs, Christoph Draxler, Eva Liina Asu, Dan Dediú, Katri Hiovain, Shigeto Kawahara, Sofia Koutalidis, Manfred Krifka, Pärtel Lippus, Gary Lupyán, Grace E Oh, Jing Paul, Caterina Petrone, Rachid Ridouane, Sabine Reiter, Nathalie Schümchen, Ádám Szalontai, Özlem Ünal-Logacev, Jochen Zeller, Marcus Perlman, and Bodo Winter. 2022. The bouba/kiki effect is robust across cultures and writing systems. *Philos. Trans. R. Soc. Lond. B Biol. Sci.*, 377(1841):20200390.
- Giannis Daras and Alexandros G Dimakis. 2022. Discovering the hidden vocabulary of dalle-2. *arXiv preprint arXiv:2206.00169*.
- Yuntian Deng. 2023. OpenAI watch. <https://openaiwatch.com/>. Accessed: 2023-7-6, source: <https://twitter.com/yuntiangeng/status/1641108596510343168?s=20>.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan†, Nicholas Joseph†, Ben Mann†, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah†. 2021. [A mathematical framework for transformer circuits](#). *Transformer Circuits Thread*.
- Eray Elicaçık. Playing with fire: The leaked plugin DAN unchains ChatGPT from its moral and ethical restrictions. <https://dataconomy.com/2023/03/31/chatgpt-dan-prompt-how-to-jailbreak-chatgpt/>. Accessed: 2023-7-6.
- Tyna Eloundou, Sam Manning, Pamela Mishkin, and Daniel Rock. 2023. Gpts are gpts: An early look at the labor market impact potential of large language models. *arXiv preprint arXiv:2303.10130*.
- Kawin Ethayarajh. 2019. How contextual are contextualized word representations? comparing the geometry of bert, elmo, and gpt-2 embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 55–65.
- Ruth C Fong and Andrea Vedaldi. 2017. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE international conference on computer vision*, pages 3429–3437.
- National Association for College Admission Counseling. 2008. *Report of the commission on the use of standardized tests in undergraduate admission*. ERIC Clearinghouse.
- Jun Gao, Di He, Xu Tan, Tao Qin, Liwei Wang, and Tieyan Liu. 2019. [Representation degeneration problem in training natural language generation models](#). In *International Conference on Learning Representations*.
- Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh Dhole, et al. 2021. The gem benchmark: Natural language generation, its evaluation and metrics. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 96–120.
- A Shaji George and AS Hovan George. 2023. A review of chatgpt ai’s impact on several business sectors. *Partners Universal International Innovation Journal*, 1(1):9–23.
- Nathan Godey, Éric de la Clergerie, and Benoît Sagot. 2023. Is anisotropy inherent to transformers? *arXiv preprint arXiv:2306.07656*.
- Andreas Grivas, Nikolay Bogoychev, and Adam Lopez. 2022. [Low-rank softmax can have unargmaxable classes in theory but rarely in practice](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6738–6758, Dublin, Ireland. Association for Computational Linguistics.
- Peter D Grünwald. 2007. *The minimum description length principle*. MIT press.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*.

- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112.
- Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. 2019. Learning latent dynamics for planning from pixels. In *International conference on machine learning*, pages 2555–2565. PMLR.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*.
- Benjamin Hilton. 2022. No, DALL-E doesn’t have a secret language.(or at least, we haven’t found one yet)this viral DALL-E thread has some pretty astounding claims. but maybe the reason they’re so astounding is that, for the most part, they’re not true. thread (1/15) <https://t.co/8F2WDp7lTK>. https://twitter.com/benjamin_hilton/status/1531780892972175361?lang=en. Accessed: 2023-7-6.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- John H Holland. 2014. *Complexity: A very short introduction*. OUP Oxford.
- Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. *International Conference on Learning Representations*.
- Gabriel Ilharco, Rowan Zellers, Ali Farhadi, and Hananeh Hajishirzi. 2021. Probing contextual language models for common ground with visual representations. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5367–5377.
- Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable nlp systems: How should we define and evaluate faithfulness? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4198–4205.
- Neel Jain, Khalid Saifullah, Yuxin Wen, John Kirchenbauer, Manli Shu, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. Bring your own data! self-supervised evaluation for large language models. *arXiv preprint arXiv:2306.13651*.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. In *International Conference on Machine Learning*, pages 10697–10707. PMLR.
- Andrej Karpathy, Justin Johnson, and Li Fei-Fei. 2015. [Visualizing and understanding recurrent networks](#).
- Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*.
- Urvashi Khandelwal, Kevin Clark, Dan Jurafsky, and Lukasz Kaiser. 2019. Sample Efficient Text Summarization Using a Single Pre-Trained Transformer. *arXiv preprint arXiv:1905.08836*.
- Urvashi Khandelwal, He He, Peng Qi, and Dan Jurafsky. 2018. Sharp nearby, fuzzy far away: How neural language models use context. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 284–294.
- Lee Kiho. ChatGPT_DAN: ChatGPT DAN, jailbreaks prompt.
- Yair Lakretz, German Kruszewski, Theo Desbordes, Dieuwke Hupkes, Stanislas Dehaene, and Marco Baroni. 2019. [The emergence of number and syntax units in LSTM language models](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 11–20, Minneapolis, Minnesota. Association for Computational Linguistics.
- Angeliki Lazaridou and Marco Baroni. 2020. Emergent multi-agent communication in the deep learning era. *arXiv preprint arXiv:2006.02419*.
- Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, Vimal Manohar, Yossi Adi, Jay Mahadeokar, and Wei-Ning Hsu. 2023. [Voicebox: Text-Guided multilingual universal speech generation at scale](#).
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8424–8445.
- Belinda Z Li, Jane Yu, Madian Khabsa, Luke Zettlemoyer, Alon Halevy, and Jacob Andreas. 2021. Quantifying adaptability in pre-trained language models with 500 tasks. *arXiv preprint arXiv:2112.03204*.
- Thomas Liao, Rohan Taori, Inioluwa Deborah Raji, and Ludwig Schmidt. 2021. Are we learning yet? a meta review of evaluation failures across machine learning. In *Thirty-fifth Conference on Neural Information*

- Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2021. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*.
- Tal Linzen, Grzegorz Chrupała, Yonatan Belinkov, and Dieuwke Hupkes, editors. 2019. *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Association for Computational Linguistics, Florence, Italy.
- Zachary C Lipton. 2018. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2023a. *Lost in the middle: How language models use long contexts*.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. Gpteval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*.
- Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. 2023. Videofusion: Decomposed diffusion models for high-quality video generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Christopher D Manning, Kevin Clark, John Hewitt, Urvasi Khandelwal, and Omer Levy. 2020. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*, 117(48):30046–30054.
- R Thomas McCoy, Paul Smolensky, Tal Linzen, Jianfeng Gao, and Asli Celikyilmaz. 2023. How much do language models copy from their training data? evaluating linguistic novelty in text generation using raven. *Transactions of the Association for Computational Linguistics*, 11:652–670.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. *Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference*. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *60th Annual Meeting of the Association for Computational Linguistics, ACL 2022*, pages 3470–3487. Association for Computational Linguistics (ACL).
- Saurabh Mittal, Saikou Diallo, and Andreas Tolk. 2018. *Emergent behavior in complex systems engineering: A modeling and simulation approach*. John Wiley & Sons.
- Masahiro Mori. 2012. The uncanny valley: The original essay by masahiro mori. <https://spectrum.ieee.org/the-uncanny-valley>. Accessed: 2023-7-6.
- Miguel Morin and Matthew Willetts. 2020. Non-determinism in tensorflow resnets. *arXiv preprint arXiv:2001.11396*.
- Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. *Orca: Progressive learning from complex explanation traces of GPT-4*.
- Mathias Müller, Annette Rios Gonzales, and Rico Senrich. 2020. Domain robustness in neural machine translation. In *Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)*, pages 151–164.
- MEJ Newman. 2011. Resource letter cs-1: Complex systems. *American Journal of Physics*, 79(8):800–810.
- Chris Olah. 2015. Understanding LSTM networks. <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>. Accessed: 2023-7-7.
- Chris Olah. 2022. Mechanistic interpretability, variables, and the importance of interpretable bases. <https://transformer-circuits.pub/2022/mech-interp-essay/index.html>. Accessed: 2023-7-8.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. 2022. *In-context learning and induction heads*. *Transformer Circuits Thread*.
- OpenAI. 2023. GPT-4 API general availability and deprecation of older models in the completions API. <https://openai.com/blog/>

- [gpt-4-api-general-availability](#). Accessed: 2023-7-6.
- OpenAI. 2023. [GPT-4 technical report](#).
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Letitia Parcalabescu, Albert Gatt, Anette Frank, and Iacer Calixto. 2021. Seeing past words: Testing the cross-modal capabilities of pretrained v&l models on counting tasks. In *Proceedings of the 1st Workshop on Multimodal Semantic Representations (MMSR)*, pages 32–44.
- Roma Patel and Ellie Pavlick. 2021. “was it “stated” or was it “claimed”?”: How linguistic bias affects generative language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10080–10095, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Alex Perry. 2023. OpenAI updates GPT-4 with new features. <https://mashable.com/article/openai-chatgpt-gpt-4-function-calling-update>. Accessed: 2023-7-6.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473.
- Sundar Pichai. 2023. An important next step on our AI journey. <https://blog.google/technology/ai/bard-google-ai-search-updates/>. Accessed: 2023-7-6.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. Hypothesis only baselines in natural language inference. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191.
- Ofir Press, Noah A Smith, and Mike Lewis. 2021. Shortformer: Better language modeling using shorter inputs. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5493–5505.
- Wolfgang Prinz. 2006. Messung kontra augenschein. *Psychologische Rundschau*, 57(2):106–111.
- Ivan Provilkov, Dmitrii Emelianenko, and Elena Voita. 2020. Bpe-dropout: Simple and effective subword regularization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1882–1892.
- Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. 2017. Learning to generate reviews and discovering sentiment. *arXiv preprint arXiv:1704.01444*.
- Alec Radford, Jeffrey Wu, Dario Amodei, Daniela Amodei, Jack Clark, Miles Brundage, and Ilya Sutskever. 2019a. Better language models and their implications. *OpenAI Blog*, 1:2.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019b. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Partha Pratim Ray. 2023. Chatgpt: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems*.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. Beyond accuracy: Behavioral testing of nlp models with checklist. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in bertology: What we know about how bert works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4035–4045.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. [High-resolution image synthesis with latent diffusion models](#).
- Robert Rosenthal. 1976. Experimenter effects in behavioral research.
- William Rudman and Carsten Eickhoff. 2023. Stable anisotropic regularization. *arXiv preprint arXiv:2305.19358*.
- John Schulman, Barret Zoph, Christina Kim, Jacob Hilton, Jacob Menick, Jiayi Weng, Juan Felipe Ceron Uribe, Liam Fedus, Luke Metz, Michael Pokorny, Rapha Gontijo Lopes, Shengjia Zhao, Arun Vijayvergiya, Eric Sigler, Adam Perelman, Chelsea

- Voss, Mike Heaton, Joel Parish, Dave Cummings, Rajeev Nayak, Valerie Balcom, David Schnurr, Tomer Kaftan, Chris Hallacy, Nicholas Turley, Noah Deutsch, Vik Goel, Jonathan Ward, Aris Konstantinidis, Wojciech Zaremba, Long Ouyang, Leonard Bogdonoff, Joshua Gross, David Medina, Sarah Yoo, Teddy Lee, Ryan Lowe, Dan Mossing, Joost Huizinga, Roger Jiang, Carroll Wainwright, Diogo Almeida, Steph Lin, Marvin Zhang, Kai Xiao, Katarina Slama, Steven Bills, Alex Gray, Jan Leike, Jakub Pachocki, Phil Tillet, Shantanu Jain, Greg Brockman, Nick Ryder, Alex Paino, Qiming Yuan, Clemens Winter, Ben Wang, Mo Bavarian, Igor Babuschkin, Szymon Sidor, Ingmar Kanitscheider, Mikhail Pavlov, Matthias Plappert, Nik Tezak, Heewoo Jun, William Zhuk, Vitchyr Pong, Lukasz Kaiser, Jerry Tworek, Andrew Carr, Lilian Weng, Sandhini Agarwal, Karl Cobbe, Vineet Kosaraju, Alethea Power, Stanislas Polu, Jesse Han, Raul Puri, Shawn Jain, Benjamin Chess, Christian Gibson, Oleg Boiko, Emy Parparita, Amin Tootoonchian, Kyle Kopic, and Christopher Hesse. 2022. *Introducing chatgpt*. *OpenAI Blog*.
- Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083.
- Abigail See, Aneesh Pappu, Rohun Saxena, Akhila Yerukola, and Christopher D Manning. 2019. Do massively pretrained language models make better storytellers? In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 843–861.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.
- Xing Shi, Yijun Xiao, and Kevin Knight. 2020. Why neural machine translation prefers empty outputs. *arXiv preprint arXiv:2012.13454*.
- Matt G Southern. 2023. [no title]. <https://www.searchenginejournal.com/openai-chatgpt-update/476116/>. Accessed: 2023-7-6.
- company Stability AI. 2023. StableLM: StableLM: Stability AI language models.
- Felix Stahlberg and Bill Byrne. 2019. On nmt search errors and model errors: Cat got your tongue? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3356–3362.
- Felix Stahlberg, Ilia Kulikov, and Shankar Kumar. 2022. Uncertainty determines the adequacy of the mode and the tractability of decoding in sequence-to-sequence models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8634–8645.
- Shane Steinert-Threlkeld, Xuhui Zhou, Zeyu Liu, and C M Downey. 2022. Emergent communication fine-tuning (EC-FT) for pretrained language models. *ICLR 2022 EmeCom Workshop*.
- Simeng Sun, Kalpesh Krishna, Andrew Mattarella-Micke, and Mohit Iyyer. 2021. Do long-range language models actually use long-range context? In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 807–822.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Ryan Teehan, Miruna Clinciu, Oleg Serikov, Eliza Szczechla, Natasha Seelam, Shachar Mirkin, and Aaron Gokaslan. 2022. Emergent structures and training dynamics in large language models. In *Proceedings of BigScience Episode\# 5–Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 146–159.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. Bert rediscovers the classical nlp pipeline. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601.
- Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. *Advances in Neural Information Processing Systems*, 35:38274–38290.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. *LLaMA: Open and efficient foundation language models*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- The Vicuna Team. Vicuna: An Open-Source chatbot impressing GPT-4 with 90%* ChatGPT quality. <https://lmsys.org/blog/2023-03-30-vicuna/>. Accessed: 2023-7-6.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355.

- Chaojun Wang and Rico Sennrich. 2020. On exposure bias, hallucination and domain shift in neural machine translation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3544–3552.
- Lingxiao Wang, Jing Huang, Kevin Huang, Ziniu Hu, Guangtao Wang, and Quanquan Gu. 2020. Improving neural language generation with spectrum control. In *International Conference on Learning Representations*.
- Alex Warstadt, Leshem Choshen, Aaron Mueller, Adina Williams, Ethan Wilcox, and Chengxu Zhuang. 2023. Call for papers—the babylm challenge: Sample-efficient pretraining on a developmentally plausible corpus. *arXiv preprint arXiv:2301.11796*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. 2021. Thinking like transformers. In *International Conference on Machine Learning*, pages 11080–11090. PMLR.
- Andrew Wilson. How to jailbreak ChatGPT to unlock its full potential. <https://approachableai.com/how-to-jailbreak-chatgpt/>. Accessed: 2023-7-6.
- D H Wolpert and W G Macready. 1997. No free lunch theorems for optimization. *IEEE Trans. Evol. Comput.*, 1(1):67–82.
- Zhuohan Xie, Trevor Cohn, and Jey Han Lau. 2023. [Can very large pretrained language models learn storytelling with a few examples?](#)
- Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond. *arXiv preprint arXiv:2304.13712*.
- Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021. [CrossFit: A few-shot learning challenge for cross-task generalization in NLP](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7163–7189, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Lili Yu, Dániel Simig, Colin Flaherty, Armen Aghajanyan, Luke Zettlemoyer, and Mike Lewis. 2023. Megabyte: Predicting million-byte sequences with multiscale transformers. *arXiv preprint arXiv:2305.07185*.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. *Advances in neural information processing systems*, 32.
- Chaoning Zhang, Chenshuang Zhang, Chenghao Li, Yu Qiao, Sheng Zheng, Sumit Kumar Dam, Mengchun Zhang, Jung Uk Kim, Seong Tae Kim, Jinwoo Choi, et al. 2023. One small step for generative ai, one giant leap for agi: A complete survey on chatgpt in aigc era. *arXiv preprint arXiv:2304.06488*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-judge with MT-Bench and chatbot arena](#).
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*.