

Do Robot Snakes Dream like Electric Sheep? Investigating the Effects of Architectural Inductive Biases on Hallucination

Jerry Huang
Chandar Research Lab
Mila & Université de Montréal

Boxing Chen
Noah's Ark Lab

Prasanna Parthasarathi*
Noah's Ark Lab

Mehdi Rezagholizadeh[†]
Advanced Micro Devices

Sarath Chandar
Chandar Research Lab
Mila & Polytechnique Montréal & CIFAR AI Chair

Abstract

The growth in prominence of large language models (LLMs) in everyday life can be largely attributed to their generative abilities, yet some of this is also owed to the risks and costs associated with their use. On one front is their tendency to *hallucinate* false or misleading information, limiting their reliability. On another is the increasing focus on the computational limitations associated with traditional self-attention based LLMs, which has brought about new alternatives, in particular recurrent models, meant to overcome them. Yet it remains uncommon to consider these two concerns simultaneously. Do changes in architecture exacerbate/alleviate existing concerns about hallucinations? Do they affect how and where they occur? Through an extensive evaluation, we study how these architecture-based inductive biases affect the propensity to hallucinate. While hallucination remains a general phenomenon not limited to specific architectures, the situations in which they occur and the ease with which specific types of hallucinations can be induced can significantly differ based on the model architecture. These findings highlight the need for better understanding both these problems in conjunction with each other, as well as consider how to design more universal techniques for handling hallucinations.

1 Introduction

Large language models (LLMs) have rapidly emerged as a every-day tool in modern life (OpenAI, 2024), with many relying on their abilities to accomplish a variety of specific tasks. However, this opened up concerns relating to their propensity to *hallucinate* (Huang et al., 2023), with no concrete reasons for this behavior (Dziri et al., 2022b; Rawte et al., 2023; Chen et al., 2024b), hindering the ability to directly train LLMs that are consistently factual or able to explain themselves through

their knowledge (Madsen et al., 2024a; Prato et al., 2023, 2024; Huang et al., 2024).

In parallel, as LLMs evolve and existing limitations are discovered, alternative architectures have become increasingly common and grow in popularity. In particular, Transformer LLMs (Vaswani et al., 2017) and linear sequence models (Gu et al., 2022; Gu and Dao, 2024) present contrasting methods of encoding sequences, with the Transformer using attention (Bahdanau et al., 2015) to form lossless representations of the context, while linear sequence models follow the recurrent neural network (Rumelhart et al., 1986; Jordan, 1986) in their use of a compressed state representation.

With the intensifying focus in both directions, a notable void exists in verifying how each can affect the others in conjunction. For example, existing works in hallucination detection and mitigation focus almost exclusively on Transformer-based models (Maynez et al., 2020; Longpre et al., 2021a; Guerreiro et al., 2023; Shi et al., 2023; Ji et al., 2023b; Farquhar et al., 2024; Wei et al., 2024), without extension to recurrent-style models, despite the use of a unified hidden representation potentially acting as an information bottleneck that can induce more common hallucinations. This lack of unified understanding on both topics prompts the need for a more explicit investigation.

In this work, we comprehensively explore the differences between pure attention LLMs and recurrent LLMs, specifically with respect to the propensity to hallucinate. Using a set of 20 different hallucination tasks, categorized into 6 groups that evaluate both faithfulness and factuality hallucinations, we evaluate across numerous open-source LLMs that range in scale from under 1B parameters to 70B parameters all the while covering a variety of different architecture choices such as self-attention, recurrent and hybrid models. We further evaluate across factors such as instruction-tuning, all to build a more comprehensive picture of architecture-

*Corresponding author: pp1403@gmail.com

[†]Work done while at Noah's Ark Lab.

specific phenomena with respect to hallucination. From this, we observe the following:

- (1) Viewed very broadly over various different tasks and settings, neither Transformer-based nor recurrent/hybrid LLMs appear to induce hallucinations more often than others.
- (2) However, shifting to individual tasks, it becomes evident that they result in disparate tendencies on specific tasks evaluating for unique criteria such as recalling long-tailed factual knowledge and falling into memorization traps, highlighting that model architecture may promote specific behavior that renders some types of hallucinations more common.
- (3) Evaluating the effects of scale and instruction-tuning, we observe that though factuality remains dependent on model size, recurrent/hybrid architectures are often more faithful at smaller sizes and observe significantly fewer faithfulness benefits from instruction-tuning and scaling compared to self-attention models.

These results highlight that while some hallucinations may be quite consistent across architectures, others are the direct results of specific model design choices that go into the LLM construction. This hints towards the need for more careful consideration on this front, as different techniques for addressing this problem can potentially be highly catered towards specific models, bringing to the forefront the need for better consideration of both these problems in the face of each other.

2 Related Work

Hallucinations in LLMs. *Hallucination* broadly refers to when LLMs generate information that does not directly follow from the context, such as nonsensical or irrelevant answers to questions (Ji et al., 2023a). While it has grown in importance due to its direct relationship with ensuring the safe and responsible use of LLMs, both classifying and quantifying the hallucinations is challenging. In particular, it is difficult to ascertain if the divergence occurs because of specific data heuristics (Lebret et al., 2016; Wiseman et al., 2017) or because of the innate lack of similarity between pre-training and downstream tasks (Rashkin et al., 2021), while measuring for hallucinations automatically introduces various biases (Reiter, 2018; Tian et al., 2020; Ganguli et al., 2022) that may not fully

capture the scope of the errors. However, a variety of task-specific benchmarks (Li et al., 2020; Pagnoni et al., 2021; Zhou et al., 2021; Santhanam et al., 2022) have shown various LLMs to struggle with factual inconsistencies, highlighting a need to render them safer for every-day use.

LLM Architectures. General LLM architecture can be broadly thought to be composed of two components: token mixers which serve to model transformations between time steps (such as attention and recurrent layers) and channel mixers, such as multi-layer perceptrons and mixtures-of-experts (MoE) (Jacobs et al., 1991; Shazeer et al., 2017), which allow communication between different channels within a single time step. Accordingly, token mixing often forms a computational bottleneck in terms of time complexity while channel mixers consist of a memory bottleneck. Though the contemporary standard for token mixing remains self-attention, alternatives are becoming increasingly common as they begin to display more promise. These include the use of linear attention (Katharopoulos et al., 2020) to compensate for the quadratic memory complexity of vanilla self-attention, to new recurrent models (Gu and Dao, 2024; Gemma Team, 2024b) that function as a linear recurrent neural network but can process all elements of a sequence in parallel, and well as hybrid mixtures of recurrent mechanisms and attention (AI21, 2024; Dao and Gu, 2024) that have emerged as a meaningful competitor.

Architecture and Hallucination. Despite growing research in new architectural components, their effects on hallucination have yet to be studied. While some works (Madsen et al., 2024b; Hu et al., 2024; Schimanski et al., 2024) have proposed modifications to either the learning/generation pipeline as a way of reducing hallucinations, proposals for hallucination reduction through structural modification have yet to be suggested. Additionally, though some works (Elhage et al., 2021; Fu et al., 2023; Lutati et al., 2023; Poli et al., 2024) have demonstrated self-attention to aptly solve synthetic tasks that form an essential component of language modeling, these fail to remain faithful on more realistic datasets, punctuating a major limitation of existing LLMs. Finally, though tangential work demonstrates that recurrent models may suffer from issues with information retention (Vardasbi et al., 2023), formally defining a link with hallucination remains necessary. These issues lead to a lack of clarity

on this front, accentuating the need for a formal investigation comparing hallucination alongside architectural paradigms. In particular, the choice of token mixer is important, as this is the component that governs how information is shared between different elements of the sequence to form a holistic representation. As such, an effective token mixer will adequately enable enough information from the context to propagate forward during the generation phase (Olsson et al., 2022; Arora et al., 2024), potentially enabling models to hallucinate less.

3 Background

Attention for Sequences. Vanilla self-attention as used in Transformers is powerful but costly. When provided an embedded text representation as a sequence of tokens $\mathbf{x} \in \mathbb{R}^{L \times d}$, each Transformer layer in the network applies a function

$$T_\ell(\mathbf{x}) = \text{FF}_\ell(A_\ell(\mathbf{x}) + \mathbf{x}) + A_\ell(\mathbf{x}) \quad (1)$$

where A_ℓ is the self-attention of the ℓ -th layer and FF_ℓ is the following feed-forward network. Self-attention computes, for a token at position i in a sequence, a weighted average of the feature representations of all tokens (the values V_ℓ) in the sequence with a weight proportional to a similarity score between i (the query \mathbf{Q}_ℓ at position i) and the rest of the sequence (the keys \mathbf{K}_ℓ). In particular, these can be computed for all positions in parallel

$$\begin{aligned} \mathbf{Q}_\ell &= \mathbf{x} \mathbf{W}_\ell^Q & \mathbf{K}_\ell &= \mathbf{x} \mathbf{W}_\ell^K & \mathbf{V}_\ell &= \mathbf{x} \mathbf{W}_\ell^V \\ A_\ell(\mathbf{x}) &= \mathbf{V}'_\ell = \text{softmax}(\mathbf{Q}_\ell \mathbf{K}_\ell^T / \sqrt{d}) \mathbf{V}_\ell \end{aligned} \quad (2)$$

providing the model a lossless representation of the complete past context. This can be seen as equivalent to search and retrieval within a database, where search is defined using query-key parameterizations and retrieval with value parameterization. In this setting, the database from which the information is being retrieved is equivalent to the model parameters, which store information from a training corpus that the weights $\mathbf{W}^K, \mathbf{W}^Q, \mathbf{W}^V$ parameterizing the model attempt to mimic.

Multi-head attention, a variant of self-attention popularized by the Transformer (Vaswani et al., 2017), has become the dominant variant used in LLMs such as LLaMA (LLaMA Team, 2024) and Gemma (Gemma Team, 2024a). Additional variants, such as multi-query attention from Falcon (Almazrouei et al., 2023) and grouped-query attention from Mistral (Jiang et al., 2023), adapt multi-head

attention but remain build upon the same underlying principle of self-attention. Consequently, such methods are considered to fall under the family of self-attention token mixers.

Recurrent LLMs. A concern for Transformers is the quadratic complexity of attention with sequence length, leading a focused on improving this bound as it directly affects the ability to learn from long sequences. Instead of directly modifying attention (Katharopoulos et al., 2020; Kitaev et al., 2020; Choromanski et al., 2021; Zeng et al., 2025), Gu et al. (2020) motivated a novel paradigm using state-space models (SSMs) from control theory. SSMs map an input $x(t) \in \mathbb{R}^d$ to an intermediate state $h(t) \in \mathbb{R}^n$ that is then projected to an output $y(t) \in \mathbb{R}^d$:

$$h'(t) = \mathbf{A}h(t) + \mathbf{B}x(t), \quad y(t) = \mathbf{C}h(t) + \mathbf{D}x(t)$$

where $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and \mathbf{D} are trainable parameters and $h'(t)$ represents the rate at which $h(t)$ changes. Gu et al. (2021) use this paradigm to define a recurrent model to work on discrete signals, in which case the input can be regarded as discretized data sampled from a continuous signal with a step size Δ , for which the corresponding SSM is defined by:

$$\begin{aligned} h_t &= \overline{\mathbf{A}}h_{t-1} + \overline{\mathbf{B}}x_t & y_t &= \overline{\mathbf{C}}h_t + \overline{\mathbf{D}}x_t \\ \overline{\mathbf{A}} &= \frac{(I + \Delta \mathbf{A}/2)}{(I - \Delta \mathbf{A}/2)} & \overline{\mathbf{B}} &= \frac{\Delta \mathbf{B}}{(I - \Delta \mathbf{A}/2)} \end{aligned}$$

and $\overline{\mathbf{C}} = \mathbf{C}$ ($\overline{\mathbf{D}}$ is equivalent to a residual connection and set to 0.) Thus

$$\overline{\mathbf{K}} = (\overline{\mathbf{C}\mathbf{B}}, \overline{\mathbf{C}\mathbf{A}\mathbf{B}}, \dots, \overline{\mathbf{C}\mathbf{A}}^{L-1} \overline{\mathbf{B}}) \quad \mathbf{y} = \overline{\mathbf{K}} * \mathbf{x}$$

where $\overline{\mathbf{K}}$ is the SSM kernel. As \mathbf{y} can be computed in $O(L \log L)$ with a Fast Fourier Transform (Cormen et al., 2009), the entire output can be computed in tandem, given the matrices that parameterize the system. Furthermore, setting \mathbf{A} as a Hurwitz matrix, SSMs can preserve long-term information, overcoming a long-standing issue (Bengio et al., 1994) with prior recurrent models (Hochreiter and Schmidhuber, 1997; Cho et al., 2014).

Further works have modified this structure; Gu and Dao (2024) use input-dependent \mathbf{B} and \mathbf{C} information filtering. LRU/Hawk (Orvieto et al., 2023; De et al., 2024) remove the discretization step from a continuous signal and instead learn discrete matrices for \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} directly. RWKV (Peng et al., 2023, 2024) uses a novel WKV

sequence mixing operator that acts like an RNN. However, these models all share the use of hidden states from which token-level information must pass through in order to interact with future tokens. This has also led to the creation of hybrid mixtures of recurrence and attention (Dao and Gu, 2024; De et al., 2024; AI21, 2024) to form new classes of models meant to capture the benefits of both.

Different sequence (or time) mixing modules attempt to encode the history of tokens in a fixed representation. Recurrent architectures maintain this representation and learn to update it during pre-training, but this can fall short of attention (Vardasbi et al., 2023; Jelassi et al., 2024; Huang, 2025; Wang et al., 2025) as they are constrained by its size (Arora et al., 2024). Meanwhile, attention has a theoretically infinite context window (albeit being costly). As recurrent structures are constrained to make the predictions based on the recent context (which in practice can be long, but remains shorter than what self-attention can attend to), the expressivity of the hidden representation directly affects task performance (Sun et al., 2024) and parametric knowledge can be predisposed to favor the prediction of certain tokens. We model this effect on tasks under the lens of hallucination: faithfulness and factuality. Faithfulness expects the model to follow an instruction (that is largely recent), while factuality assumes the use of token-to-token relations learned during pre-training for a broader contextual representation. Following this, we ask two questions: (1) Do self-attention and recurrent architectures hallucinate differently? (2) Does instruction-tuning and/or size benefit the different classes of LLMs equally?

4 Experiments and Results

Models. To investigate, we use models that vary in size and architecture, with a particular focus on different sequence mixing methods. These include self-attention models (Pythia, LLaMA2/3 (LLaMA Team, 2024), Falcon (Almazrouei et al., 2023), Mistral (Jiang et al., 2023), Gemma (Gemma Team, 2024a) and Mixtral (Jiang et al., 2024)), recurrent models (Mamba (Gu and Dao, 2024; Dao and Gu, 2024), FalconMamba and RWKV/Finch (Peng et al., 2023)), as well as hybrid models (RecurrentGemma (Gemma Team, 2024b), Jamba (AI21, 2024)). We use base and instruction-tuned variants when available. Additional details are available in Appendix A.

Datasets. We evaluate on the Hallucination Leaderboard (Hong et al., 2024), consisting of tasks

- 1) **Closed-book Open-domain QA**♣: NQ-OPEN (Kwiatkowski et al., 2019), TRIVIAQA (Joshi et al., 2017), TRUTHFULQA (Lin et al., 2022), POPQA (Mallen et al., 2023)
- 2) **Summarization**♦: XSUM (Narayan et al., 2018), CNN/DM (See et al., 2017)
- 3) **Reading Comprehension**♦: RACE (Lai et al., 2017), SQuADv2 (Rajpurkar et al., 2018), NQ-SWAP (Longpre et al., 2021b)
- 4) **Instruction Following**♦: MEMOTRAP (Liu and Liu, 2023), IFEVAL (Zhou et al., 2023)
- 5) **Hallucination Detection**: FAITHDIAL♦ (Dziri et al., 2022a), HALUEVAL♦ (Li et al., 2023a), TRUE-FALSE♣ (Azaria and Mitchell, 2023)
- 6) **Fact Checking**♣: FEVER (Thorne et al., 2018)

For tasks, a higher score (ranging from 0 to 100) indicates better performance. Tasks are further divided into *faithfulness* (♦), i.e. whether the generation adheres to the given context, and *factuality* (♣), i.e. whether the generation is factually correct.

4.1 Investigating Task Biases

Our first interest is to verify whether specific choices in architecture can lead to highly evident patterns in performance¹. While direct differences can be difficult to quantify due to differences in how each model is trained, some pattern are consistent, described as follows.

O1: Sequence models miss rare knowledge.

Interestingly, recurrent and hybrid models all significantly underwhelm on POPQA, a task that tests for uncommon knowledge (left plot in Figure 1). For example, FalconMamba achieves an exact-match (EM) score of 0.7 compared to 17.5 by a similar Falcon. Similarly, 2B and 9B RecurrentGemma achieve 7.6 and 16.0 EM compared to 14.6/21.3 EM for 2B/7B Gemma and 13.6/18.2 by 2B/9B Gemma2, while Jamba (0.4 EM) performs worse than a similarly sized Mixtral (31.9 EM).

Previous work (Vardasbi et al., 2023) suggests that hidden states can become dense and difficult to extract information from. However, the brevity and simplicity of the prompt here begs the question whether the hidden representation has the opportunity to become significantly dense. Additional

¹Numbers provided in Table 4, 5 and 8 in Appendix B.

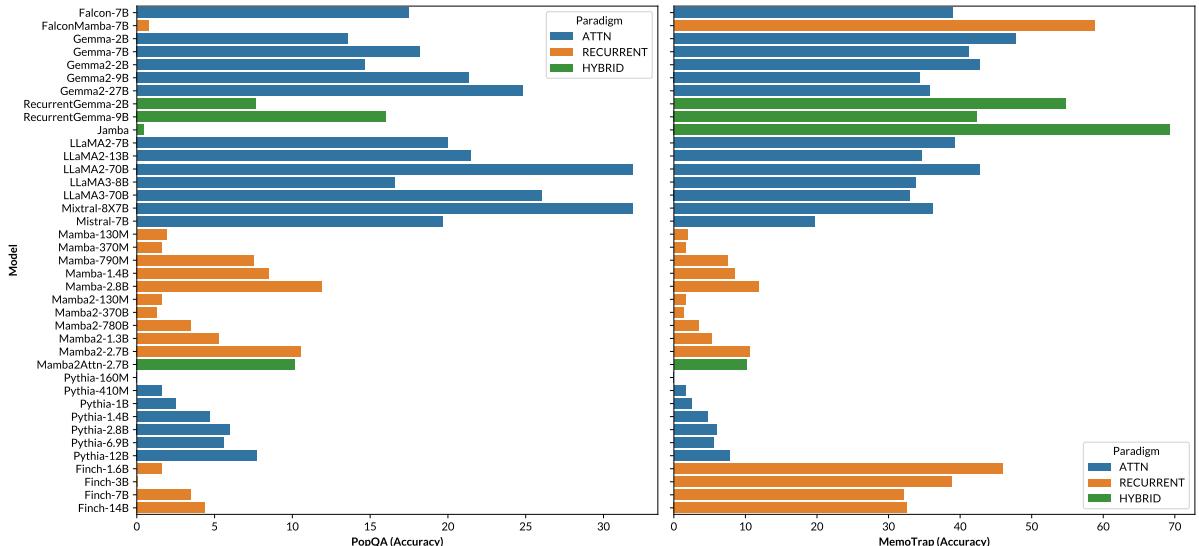


Figure 1: Performance on POPQA (left) and MEMOTRAP (right). Recurrent and hybrid models significantly outperform similar pure attention-based alternatives on MEMOTRAP, but the opposite is true on POPQA.

works (Dao and Gu, 2024) instead posit that using recurrent models with additional mechanisms enabling input filtering can enhance performance on domains such as language; however, the relevance of this claim here is again questionable. Yet this result is significant as it indicates recurrent models, despite the task’s simplicity, struggle when compared to attention-based equivalents, suggesting a link with the sequence mixing. Instead, the reason may be that recurrent/hybrid models do not learn by explicitly memorizing and acting as a retrieval system (as discussed in Section 3), where the learning mechanism may not have led to said information being stored within parameters.

O2: Recurrent models may rely less on memory, leading to more reliance on context.

There are also cases where recurrent models consistently outperform attention equivalents, namely on MEMOTRAP (right plot in Figure 1). Here, the LLM is prompted to complete a well-known proverb with an ending that deviates from the commonly used ending, testing for over-reliance on knowledge memorized from a training corpus. This suggests that although recurrent models might struggle with long-tail knowledge, they are less prone to ignoring contextual cues not stored within parametric memory, or they more model the dynamics based on the context as opposed to acting as a retrieval system. Noticeable drops exist from FalconMamba (58.8) to Falcon (38.9) and from RecurrentGemma 2B/9B to Gemma2 2B/9B (53.7/43.4 vs 42.7/34.4), suggesting that attention

can lead to more hallucinations. It further reveals a benefit of recurrent layers; by not memorizing information directly, there may be a reduced tendency to repeat previous observations (Jelassi et al., 2024) and greater focus on the context.

O3: Scale is required for emergent qualities.

However, some scale to the data and model remain important. Mamba, Finch and Pythia fail to draw the same distinctions as the other models, but given the scale of training for these (300B tokens compared to >1T tokens for other models) and the fact that the discussed tasks relate directly to information memorization, it is possible that these phenomena fail to emerge at this data scale. This is particularly evidenced by how these models all perform better with size on MEMOTRAP, which is designed in a way such that larger models should normally perform worse.

4.2 On the role of Instruction-Tuning.

The use of instruction-tuning appears to show inconsistent improvements across categories (Table 1) and tasks². There exists no strict pattern for overall hallucination; on some tasks, its use is effective for all (ex. TRUTHFULQA) or no (ex. TRIVIAQA) models. Smaller trends also exist, such as instruction-following observing the strongest gains.

However, changes in faithfulness are generally positive for pure attention models and negative for models with recurrent layers, whereas changes in

²See Figure 5.

Model Name	Hallu. Detection	Instr. Following	Closed-Book QA	Reading Comp.	Sum.	Fact-Checking	FAITHFULNESS	FACTUALITY
<i>Attention-Only Models</i>								
Gemma-2B	49.65 (\uparrow 11.51)	30.25 (\uparrow 13.84)	38.20 (\downarrow -6.92)	38.96 (\downarrow -13.96)	24.93 (\downarrow -4.53)	54.14 (\uparrow 1.60)	35.85 (\uparrow 1.92)	32.56 (\downarrow -4.56)
Gemma-7B	53.81 (\downarrow -5.67)	31.94 (\uparrow 6.79)	27.43 (\downarrow -4.13)	31.35 (\downarrow -5.83)	20.63 (\uparrow 1.04)	43.28 (\downarrow -5.22)	36.48 (\uparrow 2.71)	44.50 (\downarrow -6.04)
Gemma2-2B	58.33 (\uparrow 0.64)	26.73 (\uparrow 10.90)	29.59 (\uparrow 4.83)	32.95 (\downarrow -5.26)	16.04 (\downarrow -0.47)	39.57 (\uparrow 18.96)	35.61 (\uparrow 0.65)	34.96 (\downarrow -0.98)
Gemma2-9B	65.64 (\uparrow 6.31)	23.76 (\uparrow 30.07)	39.36 (\uparrow 3.13)	39.74 (\downarrow -0.70)	21.38 (\downarrow -3.35)	62.06 (\uparrow 6.47)	40.56 (\uparrow 7.55)	46.85 (\uparrow 3.13)
Gemma2-27B	62.03 (\uparrow 12.43)	26.62 (\uparrow 33.58)	47.19 (\uparrow 2.86)	43.04 (\uparrow 5.02)	28.92 (\downarrow -0.74)	68.28 (\uparrow 1.00)	41.54 (\uparrow 13.02)	53.92 (\uparrow 2.36)
LLaMA2-7B	53.63 (\uparrow 4.30)	28.94 (\uparrow 8.99)	37.64 (\downarrow -1.89)	27.58 (\uparrow 3.74)	25.19 (\uparrow 0.69)	51.38 (\uparrow 5.88)	35.14 (\uparrow 4.01)	42.51 (\downarrow -0.12)
LLaMA2-13B	67.80 (\downarrow -3.68)	26.57 (\uparrow 9.72)	39.65 (\downarrow -0.65)	32.17 (\downarrow -3.64)	27.33 (\downarrow -0.23)	62.35 (\downarrow -1.64)	40.83 (\uparrow 0.88)	46.27 (\downarrow -0.49)
LLaMA2-70B	62.34 (\uparrow 10.78)	30.24 (\uparrow 17.11)	47.85 (\downarrow -2.53)	39.48 (\downarrow -8.21)	28.00 (\downarrow -0.55)	66.63 (\downarrow -1.45)	41.74 (\uparrow 5.82)	53.69 (\downarrow -2.34)
LLaMA3-8B	60.76 (\uparrow 10.12)	22.06 (\uparrow 26.76)	41.55 (\uparrow 1.38)	33.52 (\uparrow 3.62)	26.62 (\downarrow -1.68)	60.84 (\uparrow 4.82)	37.60 (\uparrow 10.09)	48.14 (\uparrow 1.68)
LLaMA3-70B	71.78 (\uparrow 8.69)	21.59 (\uparrow 25.60)	48.07 (\uparrow 3.43)	46.38 (\downarrow -7.73)	28.91 (\downarrow -1.68)	69.57 (\uparrow 0.99)	45.92 (\uparrow 6.23)	54.80 (\uparrow 2.81)
Mistral-7B	60.48 (\uparrow 3.36)	28.39 (\uparrow 16.85)	41.24 (\uparrow 5.21)	32.03 (\uparrow 1.92)	26.77 (\downarrow -0.77)	58.59 (\uparrow 6.78)	38.47 (\uparrow 4.67)	47.42 (\uparrow 4.17)
Mixtral-8x7B	73.51 (\downarrow -0.15)	26.84 (\uparrow 17.60)	48.99 (\uparrow 4.33)	40.66 (\downarrow -3.90)	27.91 (\downarrow -1.07)	68.35 (\uparrow 0.49)	46.16 (\uparrow 3.89)	55.15 (\uparrow 3.50)
Falcon-7B	52.40 (\downarrow 2.90)	25.54 (\uparrow 12.62)	33.03 (\downarrow -4.42)	29.18 (\downarrow -3.64)	20.70 (\downarrow -1.00)	46.91 (\downarrow -7.71)	34.91 (\downarrow -0.71)	36.92 (\downarrow -3.75)
<i>Recurrent and Hybrid Models</i>								
RecurrentGemma-2B	52.88 (\uparrow 2.33)	33.43 (\uparrow 3.12)	27.36 (\downarrow -1.85)	30.66 (\downarrow -1.77)	18.41 (\uparrow 2.86)	42.97 (\uparrow 3.63)	36.55 (\downarrow -0.14)	31.46 (\uparrow 1.39)
RecurrentGemma-9B	55.75 (\downarrow -1.67)	31.67 (\uparrow 12.96)	36.79 (\uparrow 2.03)	36.57 (\downarrow -6.66)	22.99 (\uparrow 1.84)	51.25 (\uparrow 13.36)	37.61 (\downarrow -0.17)	42.96 (\uparrow 3.28)
Jamba	57.66 (\downarrow -2.98)	43.36 (\downarrow -5.76)	39.50 (\downarrow -0.98)	33.00 (\downarrow -5.88)	23.72 (\downarrow -15.48)	59.88 (\downarrow -2.60)	39.80 (\downarrow -6.91)	45.95 (\downarrow -0.91)
FalconMamba-7B	55.80 (\uparrow 0.73)	42.97 (\downarrow -1.19)	39.94 (\downarrow -0.43)	23.76 (\downarrow -0.85)	23.89 (\downarrow -0.09)	61.85 (\downarrow -1.80)	35.92 (\downarrow -0.26)	47.02 (\downarrow -0.38)

Table 1: Changes in performance from the use of instruction-tuning. While factuality does not exhibit a specific trend, faithfulness is shown to improve within attention-based models but not recurrent or hybrid models.

factuality show no consistent trends. This suggests that instruction-tuning for LLMs may enable particular learning patterns which become irrelevant when using recurrent layers. Their inductive biases therefore can have a particular effect on specific forms of hallucinations emerging, with instruction-tuning being a manner to mitigate these.

Recalling from Section 3, instruction-tuning can be understood as shifting the projection weights towards a new space through the additional fine-tuning (which is of itself simply an additional step of language modeling). This leads to the overwriting of some of the pre-training corpus with that used for fine-tuning, effectively replacing the keys and values from the underlying retrieval storage. While effects are ambiguous on factuality, as some potentially relevant facts can be removed, faithfulness often improves from this process, as the fine-tuning corpus is more likely to contain information relevant to context following and thereby make it easier for the attention mechanism to retrieve relevant information for such a purpose. In addition to highlighting how enhancing the ability to follow instructions does not necessarily lead to factually correct results, another conclusion stemming from the fact that recurrent models generally observe no benefits in faithfulness after instruction-tuning is that such models’ inductive biases render them inherently more faithful or the process of making them more faithful is distinct from self-attention.

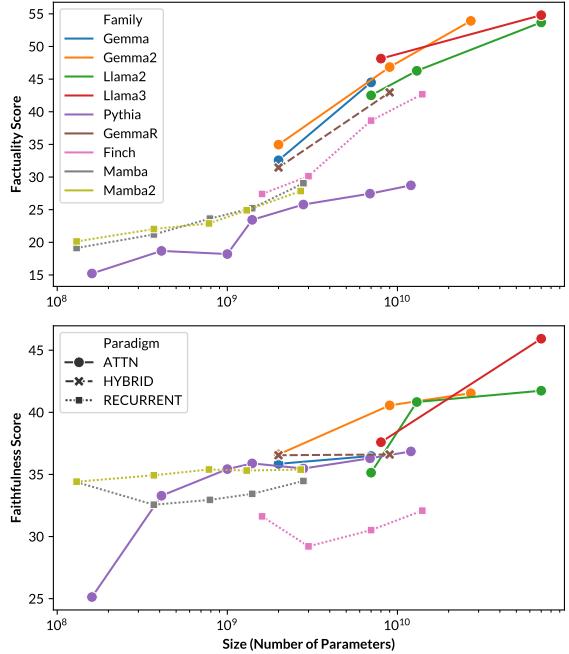


Figure 2: Changes in factuality (top) and faithfulness (bottom) as models are increased in size. Score range between 0 and 100. Factuality always increases with the number of parameters, however faithfulness increases are only meaningful for pure-attention models.

4.3 Impact of Model Size on Hallucinations.

With the large amount of prior work (Kaplan et al., 2020; Wei et al., 2022; Hoffmann et al., 2022) concluding that model size can play a significant role in model reasoning abilities, we provide this additional axis of variation in Figure 2.

Immediately, we observe that increased size

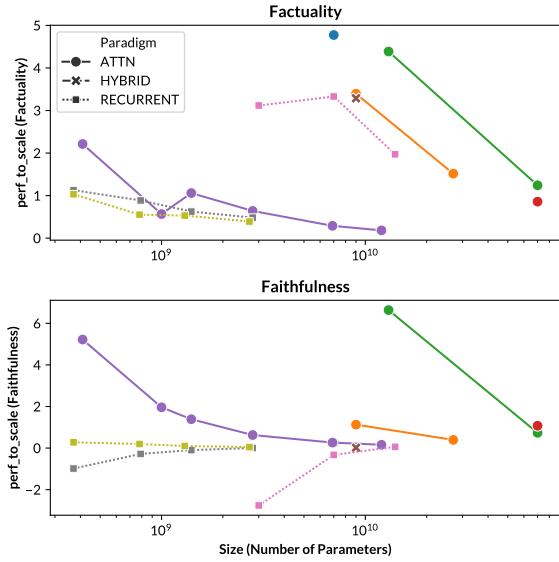


Figure 3: Performance to scale values for factuality and faithfulness. Colors differentiate model families while shapes differentiate the model type. Factuality improves with model size (indicated by values >0) and the model type does not have a link with this relative improvement. Recurrent/hybrid models show low values for faithfulness, indicating that size increases generally do not benefit them, unlike attention models.

leads to improvements in factuality regardless of model. This is partially expected, considering the range of work that suggest a link between size and the amount of factual knowledge that can be encoded within parametric space. However, faithfulness is different. Recurrent and hybrid architectures tend to saturate at a small size whereas pure attention models do not. For example, Mamba stagnates in faithfulness at around 130M parameter, whereas Pythia continues to increase up to the billions. Finch and RecurrentGemma also exhibit similar behavior. Furthermore, even the smallest Mamba are equally faithful as a 7B LLaMA2 and 1.4B Pythia, indicating a meaningful manner in which their hallucinations may differ. To further illustrate, we show a performance to scale metric, defined as

$$\text{perf_to_scale}(m, \mathcal{M}, \mathcal{B}) = \frac{m(\mathcal{M}) - m(\mathcal{B})}{P(\mathcal{M})/P(\mathcal{B}) - 1}$$

where $m(\cdot)$ is a metric of interest, \mathcal{M} is a model and \mathcal{B} is a base model of the family as \mathcal{M} . $P(\cdot)$ counts the number of parameters in a given model. Namely, this measures the relative increase in m by linear scaling the model, with a lower value indicating that linearly increasing the number of parameters yields no changes while larger values indicate

scaling leads to significant improvements. More positive values indicate that the score increases greatly with linear increases in size and negative values the opposite. Consistently negative or near-zero values hence mean that size has little effect on the score. Figure 3 depicts the performance to scale changes for factuality and faithfulness, where it becomes apparent that increasing the number of parameters in hybrid and recurrent models yields a significantly lower increase in faithfulness compared to pure attention models.

This suggests that recurrent layers enable an inherent ability to follow contexts not improved by model size. Their inductive biases can be perceived as encouraging context following, supported by how the hidden state incorporates information from the context with the current input to control the generation. This perspective aligns these results with our instruction-tuning observations, where such models do not exhibit a clear benefit compared to their attention-based counterparts. In particular, we can interpret this as indicating that the base model learns to perform this output control even without direct instruction-tuning; hence applying these techniques is ineffective as we previously observed.

4.4 Controlling for Data/Model Differences

Models are often trained on varying amounts of data, both for pre-training and fine-tuning, while also potentially possessing additional architectural differences in components, such as a channel mixer that exchanges information within each individual token. These can each influence the representation and potentially have a downstream effect on hallucination. To verify whether or not this may be the case, we control for each of these differences in order to quantify the direct effects that exist from replacing attention with sequence layers.

We inspire ourselves from Wang et al. (2024); we start with a pre-trained base model and replace some attention with sequence layers, namely Mamba, while retaining a similar parameter count and keeping the channel mixing Gated MLPs. For fair comparison, we also consider a scenario where the replaced layers are simply re-initialized. These layers are then trained using a standard next-token prediction objective. Specifically, we use a 8B LLaMA3 model and replace/re-initialize 25% or 50% of the attention layers in the model. These layers are then trained, providing us with a synthetically pre-trained model for which both the pre-training data and the channel mixer are controlled for, en-

Category	Baseline (LLaMA3-8B)	25% Reset	50% Reset	25% Replaced	50% Replaced
Hallucination Detection	60.76 ($\uparrow 10.12$)	50.27 ($\downarrow 3.24$)	44.25 ($\uparrow 8.53$)	53.64 ($\uparrow 0.91$)	53.29 ($\uparrow 0.33$)
Instruction Following	22.06 ($\uparrow 26.76$)	14.29 ($\uparrow 6.85$)	13.11 ($\uparrow 6.50$)	20.81 ($\downarrow -0.92$)	20.18 ($\uparrow 0.32$)
Closed-Book QA	41.55 ($\uparrow 1.30$)	36.79 ($\downarrow -0.09$)	31.69 ($\downarrow -0.07$)	36.26 ($\downarrow -0.42$)	33.74 ($\uparrow 0.14$)
Reading Comprehension	33.52 ($\uparrow 3.62$)	30.87 ($\uparrow 1.85$)	29.65 ($\uparrow 2.25$)	32.16 ($\uparrow 0.51$)	30.68 ($\uparrow 0.03$)
Summarization	26.62 ($\downarrow -1.60$)	21.12 ($\uparrow 1.06$)	20.86 ($\uparrow 1.67$)	23.84 ($\uparrow 0.54$)	22.59 ($\uparrow 0.79$)
Fact-Checking	60.84 ($\uparrow 4.82$)	54.93 ($\uparrow 2.69$)	50.82 ($\uparrow 3.17$)	55.45 ($\downarrow -1.06$)	49.66 ($\uparrow 2.16$)
FAITHFULNESS	37.60 ($\uparrow 10.09$)	31.40 ($\uparrow 4.29$)	30.30 ($\uparrow 4.48$)	35.74 ($\uparrow 0.17$)	34.32 ($\uparrow 0.13$)
FACTUALITY	48.14 ($\uparrow 1.60$)	43.16 ($\uparrow 0.83$)	38.98 ($\uparrow 1.60$)	42.85 ($\downarrow -0.31$)	39.08 ($\uparrow 0.56$)

Table 2: Controlling for the pre-training and post-training data as well as the presence of a channel mixer. Baseline is a LLaMA3-8B model. Colored values in brackets represent change in category performance after an additional phase of supervised instruction tuning. Note that the SFT dataset for the baseline differs from other models.

abling us to directly compare the effects of the use of sequence layers compared to attention layers. As a next step, we also perform an additional end-to-end instruction tuning. Thus we control for:

- (1) The pre-training phase, through the base model and dataset used for re-training.
- (2) The instruction fine-tuning process, through the use of the same instruction dataset mix.
- (3) Additional architectural components through when they are frozen and trained.

We use a filtered SlimPajama (Soboleva et al., 2023) to re-train layers with a sequence length of 4096. For the instruction tuning phase, we use the same datasets as Wang et al. (2024) (Chen et al., 2024a; BAAI, 2024; Teknium, 2023). We follow hyper-parameters provided by Wang et al. (2024).

Results in Table 2 demonstrate that our previous observations persist. Replacing attention layers with sequence layers and then re-training appears to regain the losses from the initial baseline model, whereas simply resetting the attention layers still appears to observe a noticeable gap. Findings regarding the effects of instruction fine-tuning hold as well, which is that pure-attention models observe a much greater increase in faithfulness relative to factuality, while hybrid models do not.

5 Discussion

Mitigating Hallucinations in LLMs. Facts such as sequence models saturating in faithfulness in a manner not resolved through instruction-tuning indicates that existing techniques for hallucination mitigation can differ in their effectiveness due to a bias towards specific architectures. Current techniques therefore deserve a look, in particular in

terms of their effectiveness across the different inductive biases of models. While some methods evidently might be specifically catered to specific models, e.g. using attention (Li et al., 2023b; Zhang et al., 2024; Chuang et al., 2024), other methods may implicitly also possess biases due to the underlying assumptions being made or specific requirements that may be unevenly addressed through different architectures, such as the reliance on data-refinement or additional context (Shi et al., 2024).

Inherent Faithfulness and Factuality. Designing inherently faithful (Herman, 2019; Wiegreffe and Pinter, 2019) or factual models remains an important issue. While much work suggests attention-based LLMs lack an inherent tendency to be faithful (Jacovi and Goldberg, 2020; Wiegreffe and Marasovic, 2021), methods have been designed to render models more interpretable (Madsen et al., 2024b) for specifically evaluating this. Similarly, some of these methods have been shown to also adapt to Mamba models (Sharma et al., 2024). However, while interpretability is useful for evaluating models as being faithful or factual, it does not ensure an inherent tendency towards being either.

6 Conclusion

Are hallucinations the direct result of how token-mixing takes place across time? We study whether different inductive biases of LLM architectures can increase this propensity by observing how hallucinations change across numerous tasks. Patterns emerge demonstrating that some architectural biases can lead to either improvements or degradation in performance compared to others, both for individual tasks and categorizations. Additionally, some model types lead to inherent behaviour when comparing the effects of instruction-tuning and scaling model size. In sum, model-specific induc-

tive biases can have a direct effect on the type of hallucinations they are faced with. However, the data used and the specific model construction can also affect learning. We hope that future work can build on our findings by exploring both the types of hallucinations that can occur and techniques meant to mitigate them and how to design more universal techniques to make models more robust or reliable.

7 Limitations

Evaluating faithfulness and factuality. A major limitation existing in hallucination detection and mitigation research is the lack of metrics that provide explicit information regarding hallucination. While the development of such metrics is an active area of research, limitations still exist, such as biases when involving additional models within the evaluation process.

Limits on tasks. Another limitation of this work is the non-exhaustive set of tasks and domains in which we measure hallucinations. Due to the exhaustive ways in which hallucinations can be measured, we limit ourselves to tasks that are well motivated and frequently used in practice.

8 Ethical Concerns

This paper provides an analysis on the effects of different architectural inductive biases on the propensity to hallucinate within large language models. As such, mistakes in methodology can lead to unsupported confidence or skepticism regarding their performance with the explored task or related ones. While skepticism may not be an ethical issue, unsupported confidence can be problematic. However, the overall message is that all LLMs have specific tasks and settings in which they will display a greater propensity to hallucinate, hence we do not believe that the ideas expressed in this work will explicitly lead to unsupported confidence.

9 Acknowledgements

Jerry Huang was supported by a National Science and Engineering Research Council (NSERC) Canada Graduate Scholarship, a Fonds de Recherche du Québec Nature et technologies (FRQNT) Training Scholarship and a Hydro-Québec Excellence Scholarship. Sarath Chandar is supported by a Canada CIFAR AI Chair, the Canada Research Chair in Lifelong Machine Learning and a NSERC Discovery Grant.

References

- AI21. 2024. [Introducing jamba: Ai21’s groundbreaking ssm-transformer model](#).
- Ebtessam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérourane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *Preprint*, arXiv:2311.16867.
- Simran Arora, Sabri Eyuboglu, Aman Timalsina, Isys Johnson, Michael Poli, James Zou, Atri Rudra, and Christopher Re. 2024. [Zoology: Measuring and improving recall in efficient language models](#). In *International Conference on Learning Representations*.
- Amos Azaria and Tom Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, Singapore. Association for Computational Linguistics.
- BAAI. 2024. [Infinity instruct](#).
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. [Neural machine translation by jointly learning to align and translate](#). In *International Conference on Learning Representations*.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. [Learning long-term dependencies with gradient descent is difficult](#). *IEEE Transactions on Neural Networks*, 5(2):157–166.
- Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John Kirchenbauer, Tianyi Zhou, and Tom Goldstein. 2024a. [Genqa: Generating millions of instructions from a handful of prompts](#). *Preprint*, arXiv:2406.10323.
- Nuo Chen, Ning Wu, Shining Liang, Ming Gong, Linjun Shou, Dongmei Zhang, and Jia Li. 2024b. [Is bigger and deeper always better? probing llama across scales and layers](#). *Preprint*, arXiv:2312.04333.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. [Learning phrase representations using RNN encoder-decoder for statistical machine translation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- Krzysztof Marcin Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, David Benjamin Bellanger, Lucy J Colwell, and Adrian Weller. 2021. [Rethinking attention with performers](#). In *International Conference on Learning Representations*.

- Yung-Sung Chuang, Linlu Qiu, Cheng-Yu Hsieh, Ranjay Krishna, Yoon Kim, and James Glass. 2024. **Lookback lens: Detecting and mitigating contextual hallucinations in large language models using only attention maps.** *Preprint*, arXiv:2407.07071.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2009. *Introduction to Algorithms, Third Edition*, 3rd edition. The MIT Press.
- Tri Dao and Albert Gu. 2024. **Transformers are SSMs: Generalized models and efficient algorithms through structured state space duality.** In *Forty-first International Conference on Machine Learning*.
- Soham De, Samuel L. Smith, Anushan Fernando, Aleksandar Botev, George Cristian-Muraru, Albert Gu, Ruba Haroun, Leonard Berrada, Yutian Chen, Srivatsan Srinivasan, Guillaume Desjardins, Arnaud Doucet, David Budden, Yee Whye Teh, Razvan Pascanu, Nando De Freitas, and Caglar Gulcehre. 2024. **Griffin: Mixing gated linear recurrences with local attention for efficient language models.** *Preprint*, arXiv:2402.19427.
- Nouha Dziri, Ehsan Kamalloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M. Ponti, and Siva Reddy. 2022a. **FaithDial: A faithful benchmark for information-seeking dialogue.** *Transactions of the Association for Computational Linguistics*, 10:1473–1490.
- Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022b. **On the origin of hallucinations in conversational models: Is it the datasets or the models?** In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5271–5285, Seattle, United States. Association for Computational Linguistics.
- 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*. [Https://transformercircuits.pub/2021/framework/index.html](https://transformercircuits.pub/2021/framework/index.html).
- Sebastian Farquhar, Jannik Kossen, Lorenz Kuhn, and Yarin Gal. 2024. **Detecting hallucinations in large language models using semantic entropy.** *Nature*, 630:625–630.
- Daniel Y Fu, Simran Arora, Jessica Grogan, Isys Johnson, Sabri Eyuboglu, Armin W Thomas, Benjamin Frederick Spector, Michael Poli, Atri Rudra, and Christopher Re. 2023. **Monarch mixer: A simple sub-quadratic GEMM-based architecture.** In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El-Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. 2022. **Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned.** *Preprint*, arXiv:2209.07858.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. **A framework for few-shot language model evaluation.**
- Gemma Team. 2024a. Gemma.
- Gemma Team. 2024b. **Recurrentgemma: Moving past transformers for efficient open language models.** *Preprint*, arXiv:2404.07839.
- Albert Gu and Tri Dao. 2024. **Mamba: Linear-time sequence modeling with selective state spaces.**
- Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. 2020. **Hippo: Recurrent memory with optimal polynomial projections.** In *Advances in Neural Information Processing Systems*, volume 33, pages 1474–1487. Curran Associates, Inc.
- Albert Gu, Karan Goel, and Christopher Ré. 2022. **Efficiently modeling long sequences with structured state spaces.** In *International Conference on Learning Representations*.
- Albert Gu, Isys Johnson, Karan Goel, Khaled Kamal Saab, Tri Dao, Atri Rudra, and Christopher Re. 2021. **Combining recurrent, convolutional, and continuous-time models with linear state space layers.** In *Advances in Neural Information Processing Systems*.
- Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. **Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation.** In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.
- Bernease Herman. 2019. **The promise and peril of human evaluation for model interpretability.** *Preprint*, arXiv:1711.07414.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. **Long short-term memory.** *Neural Comput.*, 9(8):1735–1780.

- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. [Training compute-optimal large language models](#). *Preprint*, arXiv:2203.15556.
- Giwon Hong, Aryo Pradipta Gema, Rohit Saxena, Xiaotang Du, Ping Nie, Yu Zhao, Laura Perez-Beltrachini, Max Ryabinin, Xuanli He, Clémentine Fourrier, and Pasquale Minervini. 2024. [The hallucinations leaderboard – an open effort to measure hallucinations in large language models](#). *Preprint*, arXiv:2404.05904.
- Minda Hu, Bowei He, Yufei Wang, Liangyou Li, Chen Ma, and Irwin King. 2024. [Mitigating large language model hallucination with faithful finetuning](#). *Preprint*, arXiv:2406.11267.
- Jerry Huang. 2025. [How well can a long sequence model model long sequences? comparing architectural inductive biases on long-context abilities](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 29–39, Abu Dhabi, UAE. Association for Computational Linguistics.
- Jerry Huang, Prasanna Parthasarathi, Mehdi Rezagholizadeh, and Sarath Chandar. 2024. [Towards practical tool usage for continually learning llms](#). *Preprint*, arXiv:2404.09339.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2023. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *Preprint*, arXiv:2311.05232.
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. [Adaptive mixtures of local experts](#). *Neural Computation*, 3(1):79–87.
- 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, Online. Association for Computational Linguistics.
- Samy Jelassi, David Brandfonbrener, Sham M. Kakade, and eran malach. 2024. [Repeat after me: Transformers are better than state space models at copying](#). In *Forty-first International Conference on Machine Learning*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023a. [Survey of hallucination in natural language generation](#). *ACM Computing Surveys*, 55(12):1–38.
- Ziwei Ji, Zihan Liu, Nayeon Lee, Tiezheng Yu, Bryan Wilie, Min Zeng, and Pascale Fung. 2023b. [RHO: Reducing hallucination in open-domain dialogues with knowledge grounding](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4504–4522, Toronto, Canada. Association for Computational Linguistics.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2024. [Mistral of experts](#). *Preprint*, arXiv:2401.04088.
- Michael I. Jordan. 1986. [Serial order: a parallel distributed processing approach](#). Technical report, University of California, San Diego: Institute for Cognitive Science.
- Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *Preprint*, arXiv:2001.08361.
- Angelos Katharopoulos, Apoorv Vyas, Nikolaos Papas, and François Fleuret. 2020. [Transformers are RNNs: Fast autoregressive transformers with linear attention](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 5156–5165. PMLR.
- Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. 2020. [Reformer: The efficient transformer](#). In *International Conference on Learning Representations*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob

- Uszkoreit, Quoc Le, and Slav Petrov. 2019. **Natural questions: A benchmark for question answering research**. *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. **RACE: Large-scale ReADING comprehension dataset from examinations**. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Rémi Lebret, David Grangier, and Michael Auli. 2016. **Neural text generation from structured data with application to the biography domain**. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213, Austin, Texas. Association for Computational Linguistics.
- Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. 2023a. **HaluEval: A large-scale hallucination evaluation benchmark for large language models**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6449–6464, Singapore. Association for Computational Linguistics.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023b. **Inference-time intervention: Eliciting truthful answers from a language model**. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Yangming Li, Kaisheng Yao, Libo Qin, Wanxiang Che, Xiaolong Li, and Ting Liu. 2020. **Slot-consistent NLG for task-oriented dialogue systems with iterative rectification network**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 97–106, Online. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. **TruthfulQA: Measuring how models mimic human falsehoods**. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Alisa Liu and Jiacheng Liu. 2023. **The memotrap dataset**.
- LLaMA Team. 2024. **The llama 3 herd of models**. *Preprint*, arXiv:2407.21783.
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021a. **Entity-based knowledge conflicts in question answering**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021b. **Entity-based knowledge conflicts in question answering**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shahar Lutati, Itamar Zimerman, and Lior Wolf. 2023. **Focus your attention (with adaptive IIR filters)**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12538–12549, Singapore. Association for Computational Linguistics.
- Andreas Madsen, Sarath Chandar, and Siva Reddy. 2024a. **Are self-explanations from large language models faithful?** *Preprint*, arXiv:2401.07927.
- Andreas Madsen, Siva Reddy, and Sarath Chandar. 2024b. **Faithfulness measurable masked language models**. In *Forty-first International Conference on Machine Learning*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. **When not to trust language models: Investigating effectiveness of parametric and non-parametric memories**. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. **On faithfulness and factuality in abstractive summarization**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. **Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization**. *ArXiv*, abs/1808.08745.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Connelly, 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**. *Preprint*, arXiv:2209.11895.
- OpenAI. 2024. **Gpt-4 technical report**. *Preprint*, arXiv:2303.08774.
- Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pascanu, and Soham De. 2023. **Resurrecting recurrent neural networks for long sequences**. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 26670–26698. PMLR.

- Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. 2021. [Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4812–4829, Online. Association for Computational Linguistics.
- Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Gv, Xuzheng He, Haowen Hou, Przemysław Kazienko, Jan Kocón, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Johan Wind, Stanisław Woźniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. 2023. [RWKV: Reinventing RNNs for the transformer era](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14048–14077, Singapore. Association for Computational Linguistics.
- Bo Peng, Daniel Goldstein, Quentin Anthony, Alon Albalak, Eric Alcaide, Stella Biderman, Eugene Cheah, Xingjian Du, Teddy Ferdinand, Haowen Hou, Przemysław Kazienko, Kranthi Kiran GV, Jan Kocón, Bartłomiej Koptyra, Satyapriya Krishna, Ronald McClelland Jr. au2, Jiaju Lin, Niklas Muennighoff, Fares Obeid, Atsushi Saito, Guangyu Song, Haoqin Tu, Cahya Wirawan, Stanisław Woźniak, Ruichong Zhang, Bingchen Zhao, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. 2024. [Eagle and finch: RWKV with matrix-valued states and dynamic recurrence](#). Preprint, arXiv:2404.05892.
- Michael Poli, Armin W Thomas, Eric Nguyen, Pragaash Ponnusamy, Björn Deiseroth, Kristian Kersting, Taiji Suzuki, Brian Hie, Stefano Ermon, Christopher Ré, Ce Zhang, and Stefano Massaroli. 2024. [Mechanistic design and scaling of hybrid architectures](#). Preprint, arXiv:2403.17844.
- Gabriele Prato, Jerry Huang, Prasanna Parthasarathi, Shagun Sodhani, and Sarath Chandar. 2023. [EpiKEval: Evaluation for language models as epistemic models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9523–9557, Singapore. Association for Computational Linguistics.
- Gabriele Prato, Jerry Huang, Prasanna Parthasarathi, Shagun Sodhani, and Sarath Chandar. 2024. [Do large language models know how much they know?](#) In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6054–6070, Miami, Florida, USA. Association for Computational Linguistics.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don't know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Hannah Rashkin, David Reitter, Gaurav Singh Tomar, and Dipanjan Das. 2021. [Increasing faithfulness in knowledge-grounded dialogue with controllable features](#). 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 704–718, Online. Association for Computational Linguistics.
- Vipula Rawte, Swagata Chakraborty, Agnihib Pathak, Anubhav Sarkar, S.M Towhidul Islam Tonmoy, Aman Chadha, Amit Sheth, and Amitava Das. 2023. [The troubling emergence of hallucination in large language models - an extensive definition, quantification, and prescriptive remediations](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2541–2573, Singapore. Association for Computational Linguistics.
- Ehud Reiter. 2018. [A structured review of the validity of BLEU](#). *Computational Linguistics*, 44(3):393–401.
- David E. Rumelhart, James L. McClelland, and PDP Research Group. 1986. *Parallel Distributed Processing, Volume 1: Explorations in the Microstructure of Cognition: Foundations*. The MIT Press.
- Sashank Santhanam, Behnam Hedayatnia, Spandana Gella, Aishwarya Padmakumar, Seokhwan Kim, Yang Liu, and Dilek Hakkani-Tur. 2022. [Rome was built in 1776: A case study on factual correctness in knowledge-grounded response generation](#). Preprint, arXiv:2110.05456.
- Tobias Schimanski, Jingwei Ni, Mathias Kraus, Elliott Ash, and Markus Leippold. 2024. [Towards faithful and robust LLM specialists for evidence-based question-answering](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1913–1931, Bangkok, Thailand. Association for Computational Linguistics.
- 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, Vancouver, Canada. Association for Computational Linguistics.
- Arnab Sen Sharma, David Atkinson, and David Bau. 2024. [Locating and editing factual associations in mamba](#). In *First Conference on Language Modeling*.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. [Outrageously large neural networks: The sparsely-gated mixture-of-experts layer](#). Preprint, arXiv:1701.06538.

- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärlí, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *Proceedings of the 40th International Conference on Machine Learning*.
- Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Wen-tau Yih. 2024. Trusting your evidence: Hallucinate less with context-aware decoding. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 783–791, Mexico City, Mexico. Association for Computational Linguistics.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama.
- Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, Tatsunori Hashimoto, and Carlos Guestrin. 2024. Learning to (learn at test time): Rnns with expressive hidden states. *Preprint*, arXiv:2407.04620.
- Teknium. 2023. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. 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 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Ran Tian, Shashi Narayan, Thibault Sellam, and Ankur P. Parikh. 2020. Sticking to the facts: Confident decoding for faithful data-to-text generation. *Preprint*, arXiv:1910.08684.
- Ali Vardasbi, Telmo Pessoa Pires, Robin Schmidt, and Stephan Peitz. 2023. State spaces aren't enough: Machine translation needs attention. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 205–216, Tampere, Finland. European Association for Machine Translation.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Junxiong Wang, Daniele Paliotta, Avner May, Alexander M Rush, and Tri Dao. 2024. The mamba in the llama: Distilling and accelerating hybrid models. In *Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Xinyu Wang, Linrui Ma, Jerry Huang, Peng Lu, Prasanna Parthasarathi, Xiao-Wen Chang, Boxing Chen, and Yufei Cui. 2025. Resona: Improving context copying in linear recurrence models with retrieval. *Preprint*, arXiv:2503.22913.
- 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. *Transactions on Machine Learning Research*. Survey Certification.
- Jiaheng Wei, Yuanshun Yao, Jean-Francois Ton, Hongyi Guo, Andrew Estornell, and Yang Liu. 2024. Measuring and reducing llm hallucination without gold-standard answers. *Preprint*, arXiv:2402.10412.
- Sarah Wiegreffe and Ana Marasovic. 2021. Teach me to explain: A review of datasets for explainable natural language processing. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Sarah Wiegreffe and Yuval Pinter. 2019. Attention is not explanation. 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 11–20, Hong Kong, China. Association for Computational Linguistics.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.
- Qiuhan Zeng, Jerry Huang, Peng Lu, Gezheng Xu, Boxing Chen, Charles Ling, and Boyu Wang. 2025. ZETA: Leveraging \$z\$-order curves for efficient top-\$k\$ attention. In *The Thirteenth International Conference on Learning Representations*.
- Shaolei Zhang, Tian Yu, and Yang Feng. 2024. TruthX: Alleviating hallucinations by editing large language models in truthful space. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8908–8949, Bangkok, Thailand. Association for Computational Linguistics.
- Chunting Zhou, Graham Neubig, Jiatao Gu, Mona Diab, Francisco Guzmán, Luke Zettlemoyer, and Marjan Ghazvininejad. 2021. Detecting hallucinated content in conditional neural sequence generation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1393–1404, Online. Association for Computational Linguistics.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *Preprint*, arXiv:2311.07911.

A Technical Implementation Details

A.1 Models and Datasets

All models and datasets used in this paper are public and directly available on the HuggingFace Hub.

A.2 Experimental Setup

Experiments were run using the Language Model Evaluation Harness (Gao et al., 2024).

A.3 Datasets

We evaluate on the Hallucination Leaderboard (Hong et al., 2024), consisting of tasks

- 1) **Closed-book Open-domain QA** ♠: NQ-OPEN (Kwiatkowski et al., 2019), TRIVIAQA (Joshi et al., 2017), TRUTHFULQA (MC1, MC2, Generative) (Lin et al., 2022), POPQA (Mallen et al., 2023)
- 2) **Summarization** ♠: XSUM (Narayan et al., 2018), CNN/DM (See et al., 2017)
- 3) **Reading Comprehension** ♠: RACE (Lai et al., 2017), SQuADV2 (Rajpurkar et al., 2018), NQ-SWAP (Longpre et al., 2021b)
- 4) **Instruction Following** ♠: MEMOTRAP (Liu and Liu, 2023), IFEVAL (Zhou et al., 2023)
- 5) **Hallucination Detection**: FAITHDIAL ♠ (Dziri et al., 2022a), HALUEVAL ♠ (QA, Summarization, Dialogue) (Li et al., 2023a), TRUEFALSE ♠ (Azaria and Mitchell, 2023)
- 6) **Fact Checking** ♠: FEVER (Thorne et al., 2018)

For tasks, a higher score (ranging from 0 to 100) indicates better performance. These are

- EM: HALUEVAL (all sets), POPQA, TRUTHFULQA (MC1 and MC2), SQuADV2, NATURALQUESTIONS (NQ-OPEN), TRIVIAQA
- Accuracy: MEMOTRAP, TRUTHFULQA (Gen), RACE, IFEVAL
- Rouge-L: CNNDM, XSUM, FAITHDIAL, TRUEFALSE, FEVER10

Tasks are further divided into two categories: *faithfulness* (♠) hallucinations, i.e. whether an LLM generation adheres to the given source of information, and *factuality* (♣) hallucinations, i.e. whether LLMs generate factually correct content according to world knowledge based on knowledge acquired during training. To compute scores across such categories, the scores for each task in the category are averaged.

A.4 Models and Baselines

For our comparison, we use various models (Table 3) that vary in size and architecture, with a particular focus on different time-mixing methods. These include:

- All Pythia models
- All publicly available Mamba (Gu and Dao, 2024; Dao and Gu, 2024) models
- All LLaMA2 and LLaMA3 (LLaMA Team, 2024) models
- Falcon-7B (Almazrouei et al., 2023) and FalconMamba-7B
- Mistral-7B (Jiang et al., 2023)
- All Gemma (Gemma Team, 2024a) models
- RWKV (Peng et al., 2023) models, namely Finch (RWKV-v6)
- RecurrentGemma (Gemma Team, 2024b)
- Mixtral-8x7B (Jiang et al., 2024)
- Jamba (AI21, 2024)

For models with instruction fine-tuned versions, we use both the base version as well as the instruction fine-tuned variant.

A.5 Computing Resources Used

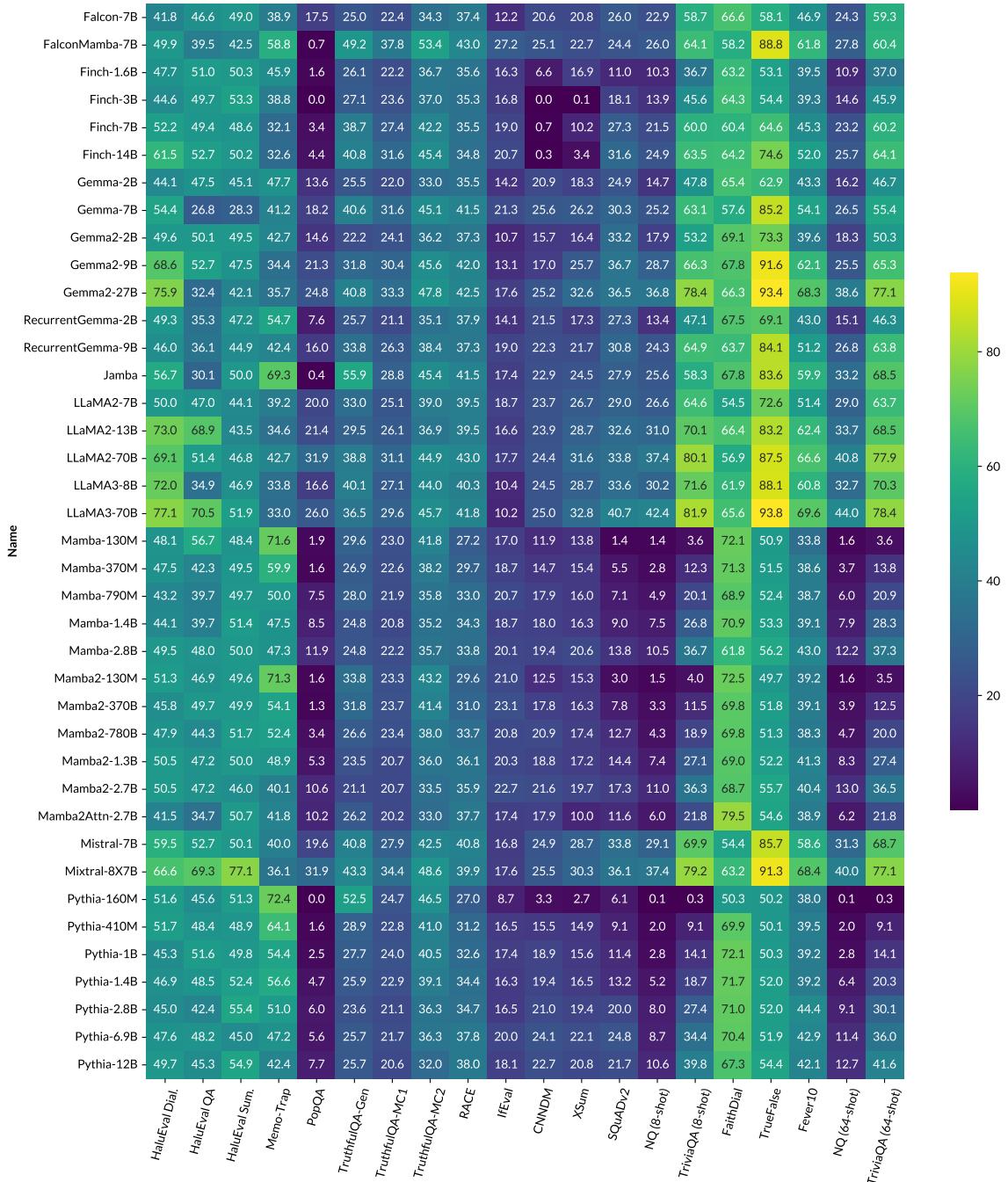
All results were obtained using a server of 8 NVIDIA V100 32GB or 4 NVIDIA RTX A6000 48GB GPUs. The accelerate package was used for model sharding in instances where a single GPU was insufficient to store the entire model.

B Complete Results

Model	Public Link	HuggingFace Model
Pythia-160M	EleutherAI/pythia-160m	✓
Pythia-410M	EleutherAI/pythia-410m	✓
Pythia-1B	EleutherAI/pythia-1b	✓
Pythia-1.4B	EleutherAI/pythia-1.4b	✓
Pythia-2.8B	EleutherAI/pythia-2.8b	✓
Pythia-6.9B	EleutherAI/pythia-6.9b	✓
Pythia-12B	EleutherAI/pythia-12b	✓
Mamba-130M	state-spaces/mamba-130m	✗
Mamba-370M	state-spaces/mamba-370m	✗
Mamba-790M	state-spaces/mamba-790m	✗
Mamba-1.4B	state-spaces/mamba-1.4b	✗
Mamba-2.8B	state-spaces/mamba-2.8b	✗
Mamba2-130M	state-spaces/mamba2-130m	✗
Mamba2-370M	state-spaces/mamba2-370m	✗
Mamba2-780M	state-spaces/mamba2-780m	✗
Mamba2-1.3B	state-spaces/mamba2-1.3b	✗
Mamba2-2.7B	state-spaces/mamba2-2.7b	✗
Mamba2Attention-2.7B	state-spaces/mamba2attn-2.7b	✗
RecurrentGemma-2B	google/recurrentgemma-2b	✓
RecurrentGemma-2B (IT)	google/recurrentgemma-2b-it	✓
RecurrentGemma-9B	google/recurrentgemma-9b	✓
RecurrentGemma-9B (IT)	google/recurrentgemma-9b-it	✓
Gemma-2B	google/gemma-2b	✓
Gemma-2B (IT)	google/gemma-2b-it	✓
Gemma-9B	google/gemma-9b	✓
Gemma-9B (IT)	google/gemma-9b-it	✓
Gemma2-2B	google/gemma2-2b	✓
Gemma2-2B (IT)	google/gemma2-2b-it	✓
Gemma2-9B	google/gemma2-9b	✓
Gemma2-9B (IT)	google/gemma2-9b-it	✓
Gemma2-27B	google/gemma2-27b	✓
Gemma2-27B (IT)	google/gemma2-27b-it	✓
Falcon-7B	tiiuae/falcon-7b	✓
Falcon-7B (IT)	tiiuae/falcon-7b-instruct	✓
FalconMamba-7B	tiiuae/falcon-mamba-7b	✓
FalconMamba-7B (IT)	tiiuae/falcon-mamba-7b-instruct	✓
Mistral-7B	mistralai/Mistral-7B-v0.3	✓
Mistral-7B (IT)	mistralai/Mistral-7B-Instruct-v0.3	✓
Mixtral-8x7B	mistralai/Mixtral-8x7B-v0.1	✓
Mixtral-8x7B (IT)	mistralai/Mixtral-8x7B-Instruct-v0.1	✓
Jamba	ai21labs/Jamba-v0.1	✓
Jamba (IT)	ai21labs/AI21-Jamba-1.5-Mini	✓
LLaMA2-7B	meta-llama/Llama-2-7b-hf	✓
LLaMA2-7B (IT)	meta-llama/Llama-2-7b-hf	✓
LLaMA2-13B	meta-llama/Llama-2-13b-hf	✓
LLaMA2-13B (IT)	meta-llama/Llama-2-13b-hf	✓
LLaMA2-70B	meta-llama/Llama-2-70b-hf	✓
LLaMA2-70B (IT)	meta-llama/Llama-2-70b-hf	✓
LLaMA3-8B	https://huggingface.co/meta-llama/Meta-Llama-3-8B	✓
LLaMA3-8B (IT)	https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct	✓
LLaMA3-70B	https://huggingface.co/meta-llama/Meta-Llama-3-70B	✓
LLaMA3-70B (IT)	https://huggingface.co/meta-llama/Meta-Llama-3-70B-Instruct	✓

Table 3: Models used and public links to their weights.

Figure 4: Performance of various base models on tasks within the Hallucination Leaderboard.



Model	Hall-Eval Dial-EM	Hall-Eval QA-EM	Hall-Eval SUM-Memo-Trap	PopQA	TruthfulQA-Gen-Rouge-L	TruthfulQA-MC1-Accuracy	TruthfulQA-MC2-Accuracy	RACE-Accuracy	IrEval-Accuracy	CNNRD-EM	XSum-EM	SQuAD2-NQ(8-shot)	NQ(8-shot)	TriviaQA-(8-shot)	FIDial	TrueFalse-Rouge-L	Fever10-EM	NQ(64-shot)	TriviaQA(64-shot)	Avg		
Falcon-7B	41.41	40.57	40.05	39.97	24.97	22.40	34.26	77.42	15.53	16.01	22.98	36.86	39.00	24.42	97.97	39.57	33.87	33.71				
Falcon-7B (IT)	38.40	39.64	48.45	60.79	8.68	38.07	20.93	44.12	37.22	15.53	20.43	18.97	20.68	14.29	67.14	63.66	39.20	15.51	39.95	33.71		
FalconMamba-7B	49.93	39.54	42.53	58.76	0.74	49.20	37.82	53.40	42.97	27.17	25.09	22.68	24.36	25.96	64.13	58.18	88.84	61.85	27.84	41.21		
FalconMamba-7B (IT)	50.14	40.16	39.41	56.20	0.18	49.08	37.70	53.23	42.58	27.36	24.79	22.80	24.53	26.34	63.45	62.67	90.29	60.05	27.65	58.45	40.89	
Finch-L4B	47.72	51.05	44.88	45.94	1.58	26.07	22.15	36.73	35.60	16.27	6.63	16.91	10.97	10.30	36.71	63.24	53.10	39.47	37.05	29.62		
Finch-L4B (IT)	52.17	49.41	48.60	32.05	3.45	38.68	27.42	42.19	35.50	19.04	0.6	10.19	27.33	21.50	60.03	60.44	64.57	45.26	23.19	60.19	34.39	
Finch-L4B	61.52	52.74	50.18	32.59	4.39	40.76	31.58	45.43	34.83	20.70	0.33	3.45	31.60	24.93	63.48	64.23	74.58	51.98	25.65	64.12	37.13	
Gemma-2B	44.13	47.54	45.09	47.72	13.58	25.46	22.03	33.05	35.46	14.16	20.93	18.34	24.91	14.65	47.76	65.43	62.88	43.28	16.20	46.72	34.81	
Gemma-2B	54.27	26.78	46.50	41.91	1.93	30.84	31.58	35.49	41.42	23.31	26.73	24.24	23.21	20.07	67.95	53.11	26.85	33.77				
Gemma-7B (IT)	50.82	51.90	51.82	58.94	2.84	38.68	29.01	45.85	41.22	25.22	28.61	18.74	29.97	8.14	26.48	38.62	55.75	38.06	8.67	38.84	38.84	
Gemma-7B (IT)	65.84	61.48	47.98	56.20	12.03	59.00	29.13	46.88	42.06	31.98	20.91	19.89	28.96	12.11	39.10	53.62	78.65	55.74	12.87	39.14	38.84	
Gemma-9B (IT)	49.93	49.88	52.76	56.81	4.98	47.98	29.00	45.46	41.71	24.71	17.00	19.87	28.73	9.31	29.21	52.89	78.49	32.99	9.31	29.21	32.99	
Gemma-12B (IT)	64.33	27.11	38.18	48.82	9.32	61.69	34.05	51.05	41.41	41.18	22.45	23.27	33.78	16.26	51.42	61.03	85.70	58.96	17.34	52.03	40.29	
Gemma-2B	49.57	50.14	49.53	42.74	14.64	22.15	24.11	36.23	37.32	10.72	15.66	16.43	33.18	17.87	53.16	69.09	73.33	39.57	18.31	50.28	35.83	
Gemma-2B	68.56	52.74	47.54	34.40	21.35	31.82	30.35	45.57	42.06	13.12	17.02	25.73	36.65	66.26	67.91	91.66	62.06	25.48	65.34	43.55		
Gemma-2B	75.24	32.28	42.34	35.94	34.81	30.76	32.29	42.42	41.22	24.22	25.61	36.48	78.42	78.42	85.48	87.89	91.44	26.85	33.55	33.77		
Gemma-2B (IT)	59.45	34.61	46.97	48.82	11.43	42.21	36.96	53.15	44.70	26.43	15.05	16.10	40.71	37.71	17.12	49.70	68.30	85.51	58.53	17.89	45.90	38.86
Gemma-9B (IT)	73.33	60.80	69.62	54.06	16.80	41.74	42.84	60.11	46.99	53.60	16.64	19.41	40.05	25.18	63.93	64.76	91.41	68.53	23.27	66.04	49.00	
Gemma-12B (IT)	79.53	65.31	62.42	48.08	20.34	43.54	46.48	50.13	50.13	15.32	21.00	24.03	34.32	73.52	73.52	91.32	56.23	25.38	48.00	43.63		
RecurrentGemma-2B	49.28	38.28	47.21	54.74	7.64	25.70	21.05	35.10	37.94	14.12	21.50	17.32	27.29	13.41	47.11	65.54	60.98	42.97	15.00	48.34	33.26	
RecurrentGemma-9B	45.97	36.14	44.45	42.38	16.01	33.78	20.32	38.37	37.28	18.96	22.25	21.73	30.76	44.86	63.34	84.11	51.25	26.84	63.77	40.16		
RecurrentGemma-2B (IT)	46.69	44.55	42.88	51.28	9.00	45.53	27.19	47.00	42.77	17.94	24.38	31.63	37.77	37.37	80.07	87.49	66.63	36.77	77.81	46.60	11.55	
RecurrentGemma-9B (IT)	60.22	18.42	44.23	44.23	4.00	47.00	27.60	44.42	44.42	17.22	23.49	23.49	32.03	77.49	87.03	87.03	24.34	55.29	40.59			
LLaMA-7B	29.06	47.00	44.12	39.21	20.01	31.05	25.09	38.99	39.62	18.67	23.66	26.53	29.00	26.65	64.64	54.45	72.59	51.38	28.98	63.75	39.65	
LLaMA-13B	73.00	68.88	43.50	21.43	29.50	20.07	36.90	39.48	16.56	23.92	31.00	30.08	66.45	83.17	62.35	33.68	68.54	43.95				
LLaMA-21B	69.12	51.44	42.74	31.90	38.80	31.09	44.86	42.77	17.94	21.81	36.62	21.81	36.24	36.24	77.22	87.49	66.63	36.77	77.81	46.60	11.55	
LLaMA-31B	58.76	37.76	45.11	50.81	24.26	42.06	44.42	44.42	17.22	23.49	23.49	23.49	77.49	87.03	87.03	24.34	55.29	40.59				
LLaMA-70B (IT)	67.62	58.38	48.70	42.35	9.63	42.11	28.03	43.96	47.12	32.44	27.60	28.60	47.00	27.01	66.44	64.75	85.16	60.71	29.22	65.57	43.13	
LLaMA-70B (IT)	79.26	63.49	60.52	55.88	12.95	48.35	35.74	52.75	43.98	18.82	24.73	30.17	35.75	32.22	73.76	58.52	85.83	65.18	35.29	71.47	49.37	
LLaMA-8B	72.04	34.85	46.92	33.76	16.56	40.15	27.05	43.96	40.29	10.35	24.51	28.72	33.61	30.19	71.59	61.82	88.12	60.84	32.69	70.26	42.62	
LLaMA-8B	77.05	70.55	51.65	33.04	36.47	29.52	42.52	41.82	41.82	17.30	23.82	32.82	42.87	41.87	70.08	59.76	87.44	50.13	20.88	78.32	50.13	
LLaMA-8B (IT)	76.65	68.15	59.19	15.81	47.86	35.99	51.60	46.66	38.45	24.56	25.49	37.41	27.51	57.65	56.97	88.94	65.66	26.47	67.49	46.66		
LLaMA-70B (IT)	82.79	88.26	72.77	41.88	27.67	50.31	41.06	61.83	46.56	52.50	24.40	36.45	40.23	27.67	80.20	65.02	91.53	70.56	38.56	72.89	54.75	
Jamba	56.70	30.10	50.64	69.34	0.45	55.94	34.33	45.26	41.53	17.38	22.93	24.52	27.89	38.29	67.82	63.83	93.88	59.88	33.16	68.47	42.73	
Jamba (IT)	20.45	20.45	74.47	47.47	6.00	54.09	44.92	44.92	44.92	17.11	21.81	21.81	19.75	19.75	50.89	50.89	57.73	25.53	33.68			
Mambu-130M	48.08	56.71	48.36	71.58	1.93	29.62	21.01	41.75	27.18	17.01	11.87	13.81	1.43	1.44	3.58	72.14	50.91	33.84	1.58	3.59	27.11	
Mambu-370M	47.53	42.29	49.48	59.94	1.64	26.93	22.64	38.22	37.22	16.87	16.87	15.41	5.49	2.80	12.33	69.79	51.29	51.49	38.55	3.71	13.80	
Mambu-790M	43.22	39.69	49.74	50.00	7.54	30.88	21.91	35.78	33.01	20.70	17.88	16.05	7.09	4.93	20.08	68.48	52.44	38.73	5.96	20.88	26.51	
Mambu-13B	44.97	37.72	51.36	47.44	24.45	20.83	25.52	35.98	36.06	20.33	18.78	17.19	14.42	7.37	27.12	69.09	52.23	41.26	8.34	27.39	30.36	
Mambu-21B	50.48	47.19	49.97	48.93	5.26	23.50	20.69	35.50	34.45	16.27	19.42	16.53	13.16	5.18	18.70	68.69	55.66	40.35	12.99	36.47	32.33	
Mambu-37B	44.86	47.23	46.94	40.06	10.56	21.05	20.69	35.50	34.89	22.74	21.59	19.73	11.22	7.03	21.75	70.43	54.49	34.68	6.18	23.23	23.23	
Mistral-7B	59.50	52.69	50.12	39.96	19.64	40.76	27.91	42.53	40.77	16.82	24.87	28.67	33.77	29.11	69.92	54.42	85.65	58.59	31.33	68.71	42.73	
Mistral-7B (IT)	69.87	41.12	62.54	45.94	5.67	42.11	59.65	46.70	44.55	23.59	28.40	32.38	35.69	35.54	68.35	56.97	89.01	65.37	27.77	67.77	47.64	
Mistral-6x7B	66.62	69.32	77.13	36.11	31.89	43.33	34.39	48.65	39.90	17.56	25.47	30.35	36.13	37.40	79.17	63.21	91.27	68.35	40.00	77.09	50.44	
Mistral-6x7B	81.17	73.47	62.52	42.31	28.53	57.04	49.82	48.84	46.70	25.31	28.38	28.38	35.69	35.54	78.55	91.09	68.84	37.06	76.18	53.10		
Pythia-4B	51.92	48.57	47.23	72.44	0.69	52.51	44.74	46.48	46.48	20.06	14.96	2.72	4.06	0.11	9.28	50.32	50.16	37.96	0.31	9.28	20.00	
Pythia-10B	51.72	48.36	48.85	64.10	1.61	28.89	23.77	41.03	31.20	16.45	14.88	9.13	1.97	0.11	9.31	50.32	50.16	37.96	0.31	9.28	20.00	
Pythia-1B	45.27	51.59	49.82	54.38	2.47	27.66	23.99	40.48	32.63	17.38	18.95	15.58	11.22	2.77	14.07	72.05	5					

	Hallu.	Detection	Instr.	Following	QA	Reading Comp.	Sum.	Fact-Checking
Falcon-7B	52.40	25.54		33.03	29.18	20.70	46.91	
Falcon-7B (IT)	49.50	38.16		28.61	25.54	19.70	39.20	
FalconMamba-7B	55.80	42.97		39.94	23.76	23.89	61.85	
FalconMamba-7B (IT)	56.53	41.78		39.51	22.91	23.80	60.05	
Gemma-2B	53.01	30.94		27.43	30.35	19.63	43.28	
Gemma-7B	50.45	31.25		38.20	39.96	25.93	54.14	
Gemma-2B (IT)	49.74	40.73		23.30	27.52	23.67	38.06	
Gemma-7B (IT)	61.16	44.09		31.28	25.00	20.40	55.74	
Gemma1.1-2B (IT)	50.61	41.27		25.43	29.38	21.00	34.65	
Gemma1.1-7B (IT)	55.27	46.50		36.66	26.63	22.32	58.96	
Gemma2-2B	58.33	26.73		29.59	32.95	16.04	39.57	
Gemma2-9B	65.64	23.76		39.36	39.74	21.38	62.06	
Gemma2-27B	62.03	26.62		47.19	43.04	28.92	68.28	
Gemma2-2B (IT)	58.97	37.63		34.42	27.69	15.57	58.53	
Gemma2-9B (IT)	71.95	53.83		42.49	39.04	18.03	68.53	
Gemma2-27B (IT)	74.46	60.20		50.05	48.06	28.18	69.28	
RecurrentGemma-2B	53.68	34.43		27.36	31.66	19.41	42.97	
RecurrentGemma-9B	54.95	30.67		36.79	35.57	21.99	51.25	
RecurrentGemma-2B (IT)	55.21	36.55		25.51	28.89	21.27	46.60	
RecurrentGemma-9B (IT)	53.28	43.63		38.82	28.91	23.83	64.61	
LLaMA2-7B	53.63	28.94		37.64	27.58	25.19	51.38	
LLaMA2-13B	67.00	25.57		39.65	31.17	26.33	62.35	
LLaMA2-70B	62.34	30.24		47.85	39.48	28.00	66.63	
LLaMA2-7B (IT)	57.93	37.93		35.75	31.32	25.88	57.26	
LLaMA2-13B (IT)	64.92	37.29		39.00	29.53	28.10	60.71	
LLaMA2-70B (IT)	73.12	47.35		45.32	31.27	27.45	65.18	
LLaMA3-8B	60.76	22.06		41.55	33.52	26.62	60.84	
LLaMA3-70B	71.78	21.59		48.07	46.38	28.91	69.57	
LLaMA3-8B (IT)	70.88	48.82		42.85	37.14	25.02	65.66	
LLaMA3-70B (IT)	80.47	47.19		51.50	38.65	27.23	70.56	
Jamba	57.66	43.36		39.50	33.00	23.72	59.88	
Jamba (IT)	54.76	37.60		38.52	27.12	8.24	57.28	
Finch-1.6B	53.09	31.10		22.69	16.60	11.77	39.47	
Finch-3B	53.29	27.80		25.96	17.89	0.03	39.26	
Finch-7B	55.04	25.55		34.58	21.05	5.42	45.26	
Finch-14B	60.65	26.64		37.54	22.40	1.89	51.98	
Mamba-130M	55.24	44.29		13.31	12.87	12.84	33.84	
Mamba-370M	52.41	39.30		15.25	12.97	15.04	38.55	
Mamba-790M	50.81	35.35		18.14	18.75	16.96	38.73	
Mamba-1.4B	51.88	33.11		19.98	20.44	17.11	39.08	
Mamba-2.8B	53.10	33.74		23.91	20.80	20.00	42.99	
Mamba2-130M	54.02	46.13		14.05	12.73	13.88	39.17	
Mamba2-370B	53.41	38.58		16.19	19.26	17.05	39.08	
Mamba2-780B	52.97	36.58		17.42	21.44	19.15	38.35	
Mamba2-1.3B	53.77	34.63		19.45	22.20	17.98	41.26	
Mamba2-2.7B	53.62	31.40		22.82	27.90	20.66	40.35	
Mamba2Attn-2.7B	53.60	29.57		19.64	36.83	21.72	43.37	
Mistral-7B	60.48	28.39		41.24	32.03	26.77	58.59	
Mistral-7B (IT)	63.84	45.24		46.45	33.95	26.00	65.37	
Mixtral-8X7B	73.51	26.84		48.99	40.66	27.91	68.35	
Mixtral-8X7B (IT)	73.36	44.44		53.32	36.76	26.84	68.84	
Pythia-160M	49.80	40.56		10.42	8.14	3.00	37.96	
Pythia-410M	53.18	40.28		11.72	20.58	15.20	39.49	
Pythia-1B	53.81	35.88		13.78	27.11	17.26	39.22	
Pythia-1.4B	54.29	36.45		17.90	22.16	17.98	39.22	
Pythia-2.8B	53.15	33.71		20.18	22.89	20.19	44.43	
Pythia-6.9B	52.62	33.62		22.47	24.86	23.14	42.91	
Pythia-12B	54.32	30.26		23.84	28.05	21.75	42.14	

Table 5: Results separated by task categories.

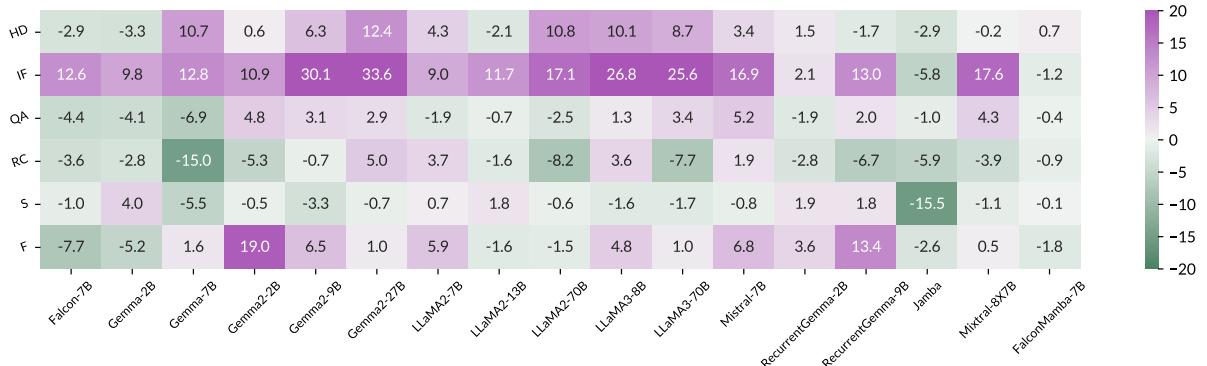
	Hallu.	Detection	Instr.	Following	QA	Reading Comp.	Sum.	Fact-Checking
Falcon-7B	49.53		29.11		43.50	26.18	20.70	46.91
Falcon-7B (IT)	45.22		44.21		29.13	21.13	19.71	39.20
FalconMamba-7B	52.15		47.19		42.70	19.98	23.90	61.85
FalconMamba-7B (IT)	52.20		45.63		41.77	19.44	23.80	60.05
Gemma-2B	50.86		36.43		34.54	28.08	20.64	43.28
Gemma-2B (IT)	49.29		43.60		19.19	23.57	21.71	38.06
Gemma-7B	45.01		32.90		44.58	34.76	24.93	54.14
Gemma-7B (IT)	60.66		47.33		29.69	23.04	20.41	55.74
Gemma1.1-2B (IT)	50.52		44.41		21.27	26.53	20.02	34.65
Gemma1.1-7B (IT)	51.32		47.12		37.28	26.31	22.32	58.96
Gemma2-2B	55.10		31.01		37.82	32.13	16.04	39.57
Gemma2-2B (IT)	54.83		40.62		35.33	28.15	15.57	58.53
Gemma2-9B	62.73		26.61		49.30	38.02	21.35	62.06
Gemma2-9B (IT)	71.21		53.89		47.89	37.78	18.02	68.53
Gemma2-27B	58.23		29.05		58.76	40.52	28.90	68.28
Gemma2-27B (IT)	74.87		57.08		57.62	44.46	28.17	69.28
RecurrentGemma-2B	49.05		38.86		34.90	27.59	18.42	42.97
RecurrentGemma-2B (IT)	51.02		40.49		25.93	27.11	21.28	46.60
RecurrentGemma-9B	51.49		34.80		46.83	34.27	22.99	51.25
RecurrentGemma-9B (IT)	50.33		47.66		45.24	29.00	24.83	64.61
LLaMA2-7B	51.62		31.69		47.96	25.65	25.18	51.38
LLaMA2-7B (IT)	55.95		41.30		40.15	29.34	25.89	57.26
LLaMA2-13B	66.34		28.98		51.85	31.01	27.32	62.35
LLaMA2-13B (IT)	62.10		37.64		46.84	27.44	27.10	60.71
LLaMA2-70B	60.75		33.58		61.27	36.44	27.98	66.63
LLaMA2-70B (IT)	74.25		49.63		52.59	30.43	27.43	65.18
LLaMA3-8B	57.88		25.19		51.81	32.13	26.60	60.84
LLaMA3-8B (IT)	71.37		51.59		49.53	35.39	25.02	65.66
LLaMA3-70B	70.61		24.64		61.02	45.04	28.88	69.57
LLaMA3-70B (IT)	81.70		45.77		59.18	37.54	27.21	70.56
Jamba	53.44		50.30		43.45	29.15	23.72	59.88
Jamba (IT)	50.28		47.46		38.92	25.42	8.23	57.28
Finch-1.6B	51.43		35.07		24.94	10.35	11.73	39.47
Finch-3B	51.38		30.74		30.36	14.33	0.03	39.26
Finch-7B	53.22		27.29		41.06	20.56	5.39	45.26
Finch-14B	58.69		28.23		43.98	23.51	1.88	51.98
Mamba-130M	52.91		51.59		4.08	5.25	12.84	33.84
Mamba-370M	49.43		44.82		9.81	6.45	15.04	38.55
Mamba-790M	47.69		39.27		15.91	11.06	16.96	38.73
Mamba-1.4B	48.63		36.97		20.59	12.91	17.12	39.08
Mamba-2.8B	51.38		37.37		27.53	15.26	20.00	42.99
Mamba2-130M	50.59		51.86		4.21	4.28	12.87	39.17
Mamba2-370B	50.90		42.72		9.26	12.18	17.05	39.08
Mamba2-780B	49.55		39.81		14.22	14.34	18.16	38.35
Mamba2-1.3B	51.44		38.46		19.64	16.15	17.99	41.26
Mamba2-2.7B	50.96		33.72		26.81	21.93	20.67	40.35
Mamba2Attn-2.7B	49.56		32.84		24.12	30.08	21.72	43.37
Mistral-7B	58.98		31.48		51.43	30.90	26.75	58.59
Mistral-7B (IT)	62.47		45.43		50.98	31.94	25.98	65.37
Mixtral-8X7B	73.44		29.32		60.87	38.99	27.89	68.35
Mixtral-8X7B (IT)	74.02		43.87		59.64	34.24	26.83	68.84
Pythia-160M	49.68		49.09		1.16	5.03	3.01	37.96
Pythia-410M	51.34		46.65		7.76	16.13	15.20	39.49
Pythia-1B	51.18		40.83		11.60	22.24	17.28	39.22
Pythia-1.4B	51.67		41.84		14.72	15.96	17.99	39.22
Pythia-2.8B	50.36		38.32		20.82	19.26	20.20	44.43
Pythia-6.9B	47.25		34.26		24.80	19.11	20.14	42.91
Pythia-12B	52.20		33.51		28.72	23.42	21.75	42.14

Table 6: Results separated by task categories, where scores are weighted by the size of the tasks.

Model Name	Hallu. Detection	Instr. Following	Closed-Book QA	Reading Comp.	Sum.	Fact-Checking	FAITHFULNESS	FACTUALITY
<i>Attention-Only Models</i>								
Gemma-2B	53.81 (↓ -5.67)	31.94 (↑ 6.79)	27.43 (↓ -4.13)	31.35 (↓ -5.83)	20.63 (↑ 1.04)	43.28 (↓ -5.22)	34.32 (↑ 0.83)	39.38 (↓ -11.20)
Gemma-7B	49.65 (↑ 11.51)	30.25 (↑ 13.84)	38.20 (↓ -6.92)	38.96 (↓ -13.96)	24.93 (↓ -4.53)	54.14 (↑ 1.60)	33.06 (↑ 4.92)	50.44 (↓ -8.42)
Gemma2-2B	58.33 (↑ 0.64)	26.73 (↑ 10.90)	29.59 (↑ 4.83)	32.95 (↓ -5.26)	16.04 (↓ -0.47)	39.57 (↑ 18.96)	33.97 (↑ 2.01)	40.54 (↑ 6.17)
Gemma2-9B	65.64 (↑ 6.31)	23.76 (↑ 30.07)	39.36 (↑ 3.13)	39.74 (↓ -0.70)	21.38 (↓ -3.35)	62.06 (↑ 6.47)	41.41 (↑ 3.93)	56.42 (↑ 1.52)
Gemma2-27B	62.03 (↑ 12.43)	26.62 (↑ 33.58)	47.19 (↑ 2.86)	43.04 (↑ 5.02)	28.92 (↓ -0.74)	68.28 (↑ 1.00)	41.82 (↑ 10.00)	64.25 (↓ -0.31)
LLaMA2-7B	53.63 (↑ 4.30)	28.94 (↑ 8.99)	37.64 (↓ -1.89)	27.58 (↑ 3.74)	25.19 (↑ 0.69)	51.38 (↑ 5.88)	35.47 (↑ 2.90)	50.65 (↓ -1.91)
LLaMA2-13B	67.80 (↓ -3.68)	26.57 (↑ 9.72)	39.65 (↓ -0.65)	32.17 (↓ -3.64)	27.33 (↓ -0.23)	62.35 (↓ -1.64)	41.08 (↑ 3.11)	57.51 (↓ -3.38)
LLaMA2-70B	62.34 (↑ 10.78)	30.24 (↑ 17.11)	47.85 (↓ -2.53)	39.48 (↓ -8.21)	28.00 (↓ -0.55)	66.63 (↓ -1.45)	42.47 (↑ 5.96)	64.76 (↓ -5.64)
LLaMA3-8B	60.76 (↑ 10.12)	22.06 (↑ 26.76)	41.55 (↑ 1.30)	33.52 (↑ 3.62)	26.62 (↓ -1.60)	60.84 (↑ 4.82)	39.32 (↑ 7.82)	57.22 (↑ 0.49)
LLaMA3-70B	71.78 (↑ 8.69)	21.59 (↑ 25.60)	48.07 (↑ 3.43)	46.38 (↓ -7.73)	28.91 (↓ -1.68)	69.57 (↑ 0.99)	49.24 (↑ 4.00)	66.05 (↓ -0.72)
Mistral-7B	60.48 (↑ 3.36)	28.39 (↑ 16.85)	41.24 (↑ 5.21)	32.03 (↑ 1.92)	26.77 (↓ -0.77)	58.59 (↑ 6.78)	39.98 (↑ 1.85)	56.04 (↑ 2.40)
Mixtral-8X7B	73.51 (↓ -0.15)	26.84 (↑ 17.60)	48.99 (↑ 4.33)	40.66 (↓ -3.90)	27.91 (↓ -1.07)	68.35 (↑ 0.49)	48.49 (↑ 0.82)	65.37 (↓ -0.54)
Falcon-7B	52.40 (↓ -2.90)	25.54 (↑ 12.62)	33.03 (↓ -4.42)	29.18 (↓ -3.64)	20.70 (↓ -1.00)	46.91 (↓ -7.71)	30.37 (↑ 3.89)	45.60 (↓ -10.78)
<i>Recurrent and Hybrid Models</i>								
RecurrentGemma-2B	52.88 (↑ 2.33)	33.43 (↑ 3.12)	27.36 (↓ -1.85)	30.66 (↓ -1.77)	18.41 (↑ 2.86)	42.97 (↑ 3.63)	34.07 (↓ -1.11)	38.93 (↓ -2.45)
RecurrentGemma-9B	55.75 (↓ -1.67)	31.67 (↑ 12.96)	36.79 (↑ 2.03)	36.57 (↓ -6.66)	22.99 (↑ 1.84)	51.25 (↑ 13.36)	35.89 (↓ -1.29)	50.63 (↑ 4.11)
Jamba	57.66 (↓ -2.90)	43.36 (↓ -5.76)	39.50 (↓ -0.98)	33.00 (↓ -5.88)	23.72 (↓ -15.48)	59.88 (↓ -2.60)	36.28 (↓ -7.38)	51.78 (↓ -3.49)
FalconMamba-7B	55.80 (↑ 0.73)	42.97 (↓ -1.19)	39.94 (↓ -0.43)	23.76 (↓ -0.85)	23.89 (↓ -0.09)	61.85 (↓ -1.80)	33.31 (↑ 0.28)	52.37 (↓ -1.11)

Table 7: Changes in performance from the use of instruction-tuning, weighted by the number of samples present in each task.

Figure 6: Change in task category performance from base model to instruction fine-tuned model.



	Faithfulness	Factuality
Falcon-7B	34.91	36.92
Falcon-7B (IT)	34.20	33.17
FalconMamba-7B	35.92	47.02
FalconMamba-7B (IT)	35.66	46.64
Finch-1.6B	31.62	27.41
Finch-3B	29.21	30.14
Finch-7B	30.52	38.65
Finch-14B	32.08	42.69
Gemma-2B	35.85	32.56
Gemma-7B	36.48	44.50
Gemma-2B (IT)	36.77	28.00
Gemma-7B (IT)	39.19	38.46
Gemma1.1-2B (IT)	36.60	30.12
Gemma1.1-7B (IT)	37.11	43.79
Gemma2-2B	36.61	34.96
Gemma2-9B	40.56	46.85
Gemma2-27B	41.54	53.92
Gemma2-2B (IT)	36.26	41.94
Gemma2-9B (IT)	48.11	49.98
Gemma2-27B (IT)	54.56	56.28
RecurrentGemma-2B	36.55	31.46
RecurrentGemma-2B (IT)	36.41	32.85
RecurrentGemma-9B	36.61	42.96
RecurrentGemma-9B (IT)	36.44	46.24
Jamba	39.80	45.95
Jamba (IT)	32.89	45.04
LLaMA2-7B	35.14	42.51
LLaMA2-13B	40.83	46.27
LLaMA2-70B	41.74	53.69
LLaMA2-7B (IT)	39.15	42.39
LLaMA2-13B (IT)	41.71	45.78
LLaMA2-70B (IT)	47.56	51.35
LLaMA3-8B	37.60	48.14
LLaMA3-70B	45.92	54.80
LLaMA3-8B (IT)	47.69	49.74
LLaMA3-70B (IT)	52.15	57.61
Mamba-130M	34.38	19.13
Mamba-370M	32.56	21.21
Mamba-790M	32.95	23.63
Mamba-1.4B	33.44	25.23
Mamba-2.8B	34.47	29.05
Mamba2-130M	34.42	20.13
Mamba2-370B	34.93	22.04
Mamba2-780B	35.40	22.90
Mamba2-1.3B	35.31	24.91
Mamba2-2.7B	36.39	27.86
Mamba2Attn-2.7B	38.91	26.49
Mistral-7B	38.47	47.42
Mistral-7B (IT)	43.14	52.59
Mixtral-8X7B	46.16	55.15
Mixtral-8X7B (IT)	48.05	58.65
Pythia-160M	25.13	15.23
Pythia-410M	33.29	18.69
Pythia-1B	35.43	18.21
Pythia-1.4B	35.89	23.44
Pythia-2.8B	35.48	25.79
Pythia-6.9B	36.30	27.46
Pythia-12B	36.85	28.73

Table 8: Results categorized by faithfulness and factuality.

	Faithfulness	Factuality
Falcon-7B	30.37	45.60
Falcon-7B (IT)	34.26	34.82
FalconMamba-7B	32.75	52.37
FalconMamba-7B (IT)	33.03	51.26
Finch-1.6B	29.36	31.87
Finch-3B	26.54	35.01
Finch-7B	29.70	43.97
Finch-14B	31.41	48.68
Gemma-2B	34.32	39.38
Gemma-2B (IT)	35.15	28.18
Gemma-7B	33.06	50.44
Gemma-7B (IT)	37.98	42.02
Gemma1.1-2B (IT)	34.55	28.59
Gemma1.1-7B (IT)	33.85	48.00
Gemma2-2B	33.97	40.54
Gemma2-2B (IT)	35.98	46.71
Gemma2-9B	41.41	56.42
Gemma2-9B (IT)	45.34	57.94
Gemma2-27B	41.82	64.25
Gemma2-27B (IT)	51.82	63.94
RecurrentGemma-2B	34.07	38.93
RecurrentGemma-2B (IT)	32.96	36.48
RecurrentGemma-9B	35.89	50.63
RecurrentGemma-9B (IT)	34.60	54.74
Jamba	36.28	51.78
Jamba (IT)	28.90	48.29
LLaMA2-7B	35.47	50.65
LLaMA2-7B (IT)	38.37	48.74
LLaMA2-13B	41.08	57.51
LLaMA2-13B (IT)	44.19	54.13
LLaMA2-70B	42.47	64.76
LLaMA2-70B (IT)	48.43	59.12
LLaMA3-8B	39.32	57.22
LLaMA3-8B (IT)	47.14	57.71
LLaMA3-70B	49.24	66.05
LLaMA3-70B (IT)	53.24	65.33
Mamba-130M	29.78	17.64
Mamba-370M	28.72	22.70
Mamba-790M	29.28	26.35
Mamba-1.4B	30.14	29.23
Mamba-2.8B	32.78	34.83
Mamba2-130M	28.44	19.58
Mamba2-370B	31.37	22.59
Mamba2-780B	31.50	25.16
Mamba2-1.3B	32.76	29.41
Mamba2-2.7B	34.30	33.42
Mamba2Attn-2.7B	34.45	33.51
Mistral-7B	39.98	56.04
Mistral-7B (IT)	41.83	58.44
Mixtral-8X7B	48.49	65.37
Mixtral-8X7B (IT)	49.31	64.83
Pythia-160M	22.83	12.62
Pythia-410M	30.48	18.94
Pythia-1B	32.71	20.29
Pythia-1.4B	32.92	25.81
Pythia-2.8B	33.60	31.22
Pythia-6.9B	31.84	32.97
Pythia-12B	35.72	35.11

Table 9: Results categorized by faithfulness and factuality, weighted by the number of examples per task.