

# Hallucination

## 目录

### • DOLA: DECODING BY CONTRASTING LAYERS IMPROVES FACTUALITY IN LARGE LANGUAGE MODELS

- 1.Introduction
  - Problems & Motivation & Inspiration
  - Idea
- 2.Method
  - FACTUAL KNOWLEDGE EVOLVES ACROSS LAYERS
  - DYNAMIC PREMATURE LAYER SELECTION
  - CONTRASTING THE PREDICTIONS
- 3.Experiments
  - Tasks
  - Setup
- 4.Reproduce
- 5.Conclusion
- Knowledge Sanitization of Large Language Models
- 1.Introduction
  - Problems
  - Idea
- 2.Method
  - Sanitization Tuning
  - Fine-tuning the MLP Layers
- 3.Experiment
  - Setup

## DOLA: DECODING BY CONTRASTING LAYERS IMPROVES FACTUALITY IN LARGE LANGUAGE MODELS

### 1.Introduction

#### Problems & Motivation & Inspiration

- Hallucination

- A possible reason for LLM's hallucination is due to the maximum likelihood language modeling objective which minimize the forward KL divergence between the data and model distributions. This objective potentially results in a model with **mass-seeking behavior** which causes the LM to assign non-zero probability to sentences that are not fully consistent with knowledge embedded in the training data.
- Transformer LMs have been loosely shown to encode “lower- level” information (e.g., part-of-speech tags) in the earlier layers, and more “semantic” information in the later layers.
- Previous work finds that “knowledge neurons” are distributed in the topmost layers of the pretrained BERT model.

## Idea

- We propose to exploit this **modular encoding of knowledge** to amplify the factual knowledge in an LM through a contrastive decoding approach, where the output probability over the next word is obtained from the difference in logits obtained from a higher layer versus a lower layer.
- By emphasizing the knowledge from higher layers and downplaying the lower or intermediate layer knowledge, we can potentially make LMs more factual and consequently reduce hallucinations.

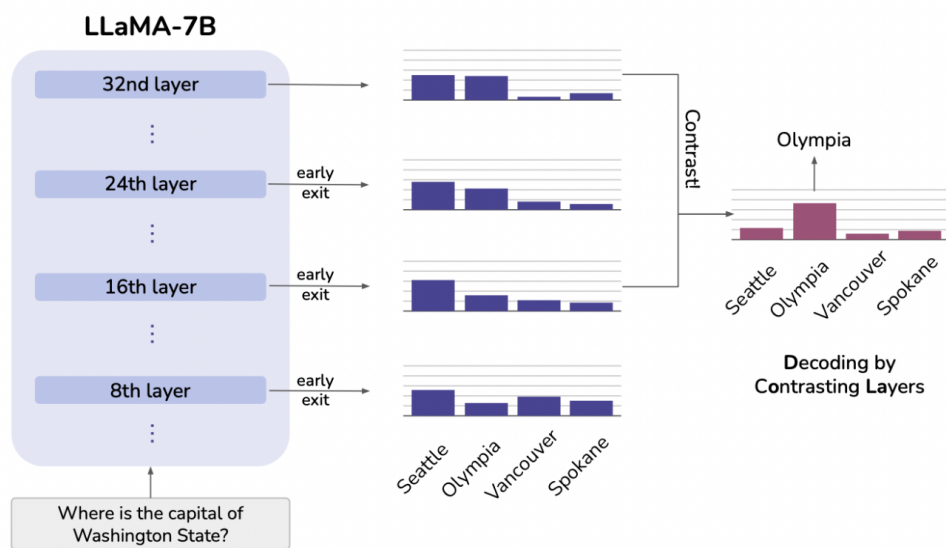


Figure 1: Illustration of how a transformer-based LM progressively incorporates more factual information along the layers. We observe that while the next-word probability of “*Seattle*” remains similar throughout the different layers, the probability of the correct answer “*Olympia*” gradually increases from the lower layers to the higher layers. DoLa uses this fact and decodes by contrasting the difference between the two layers to sharpen an LLM’s probability towards factually correct outputs.

## 2.Method

### FACTUAL KNOWLEDGE EVOLVES ACROSS LAYERS

The authors conduct preliminary analysis with the 32-layer LLaMA-7B model to motivate our approach: compute the Jensen-Shannon Divergence (JSD) between the early exiting output

distributions  $q_j(\cdot|x_{<t})$  and the final layer output distribution  $q_N(\cdot|x_{<t})$ , to show how the early exiting outputs are different from the final layer outputs.

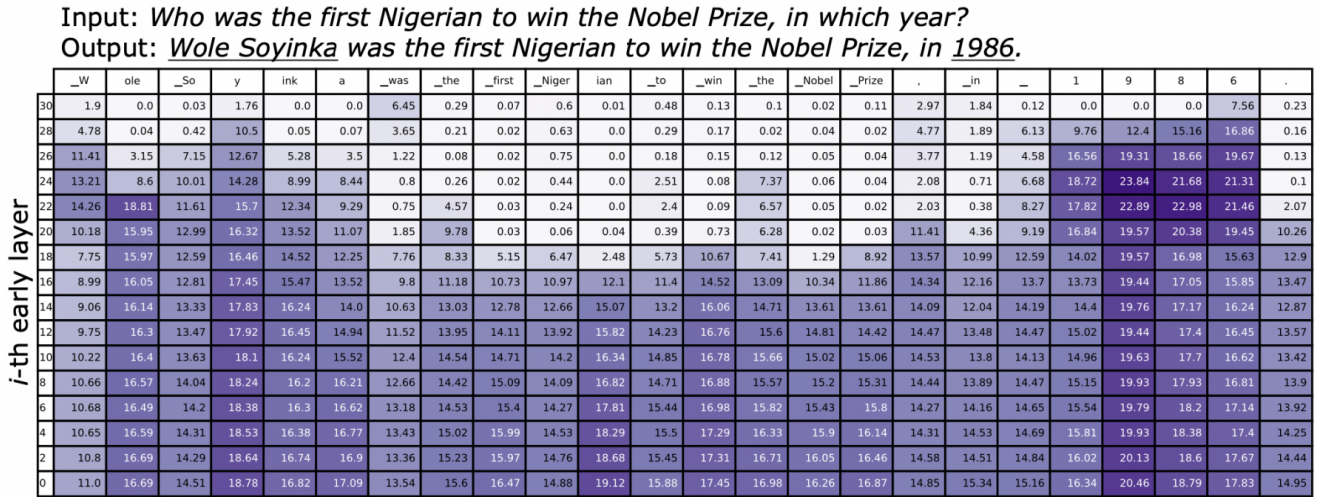


Figure 2: Jensen-Shannon Divergences between the final 32nd layer and even-numbered early layers. Column names represent predicted next tokens in each decoding step. Row names indicate the layer indices of the early exit layers, from the 0th (word embedding) layer to the 30th layer.

- When predicting important name entities or dates, which require factual knowledge, the calculated JSD would be still extremely high in the higher layers. This indicates that **the model is still changing its predictions in the last few layers**, and potentially injecting more factual knowledge into the predictions.
- When predicting function words, such as *was*, *the*, *to*, *in*, and the tokens that are copied from the input question, the JSD becomes very small from the middle of the layers. This indicates that **the model has already decided what token to generate in the early layers**, so it just keeps the output distribution almost unchanged in the higher layers.

The author accurately select the premature layer that contains plausible but less factual information, which may not always stay in the same early layer.

## DYNAMIC PREMATURE LAYER SELECTION

To magnify the effective of contrastive decoding, the optimal premature layer to select should ideally be the layer that is the most different from the final-layer outputs.

The authors adopt the following measure of distance between the next-word distributions obtained from two layers:

$$d(q_N(\cdot|x_{<t}), q_j(\cdot|x_{<t})) = \text{JSD}(q_N(\cdot|x_{<t})||q_j(\cdot|x_{<t})),$$

The  $M$ -th layer ( $0 \leq M < N$ ) is then selected as the layer with the maximum divergence among the subset of early layers:

$$M = \arg \max_{j \in \mathcal{J}} \text{JSD}(q_N(\cdot | x_{<t}) || q_j(\cdot | x_{<t})),$$

## CONTRASTING THE PREDICTIONS

We aim to amplify the output from the mature layer while downplaying the output from the premature layer. Following the Contrastive Decoding approach, we subtract the log probabilities of the premature layer outputs from those of the mature layer. We then use this resulting distribution as the next-word prediction:

$$\mathcal{F}(q_N(x_t), q_M(x_t)) = \begin{cases} \log \frac{q_N(x_t)}{q_M(x_t)}, & \text{if } x_t \in \mathcal{V}_{\text{head}}(x_t | x_{<t}), \\ -\infty, & \text{otherwise.} \end{cases}$$

$$\hat{p}(x_t) = \text{softmax}(\mathcal{F}(q_N(x_t), q_M(x_t)))$$

## 3.Experiments

### Tasks

- TruthfulQA and FACTOR: multiple choices tasks
- TruthfulQA, StrategyQA and GSM8K: open-ended generation tasks

### Setup

- LLaMA models (7B, 13B, 33B, 65B)
- baselines:
  - original decoding
  - contrastive decoding
  - Inference Time Intervention: LLaMA-7B and a linear classifier trained on TruthfulQA



Model	TruthfulQA			FACTOR	
	MC1	MC2	MC3	News	Wiki
LLaMa-7B	25.6	40.6	19.2	58.3	58.6
+ ITI (Li et al., 2023)	25.9	-	-	-	-
+ DoLa	<b>32.2</b>	<b>63.8</b>	<b>32.1</b>	<b>62.0</b>	<b>62.2</b>
LLaMa-13B	28.3	43.3	20.8	61.1	62.6
+ CD (Li et al., 2022)	24.4	41.0	19.0	62.3	64.4
+ DoLa	<b>28.9</b>	<b>64.9</b>	<b>34.8</b>	<b>62.5</b>	<b>66.2</b>
LLaMa-33B	31.7	49.5	24.2	63.8	69.5
+ CD (Li et al., 2022)	<b>33.0</b>	51.8	25.7	63.3	<b>71.3</b>
+ DoLa	30.5	<b>62.3</b>	<b>34.0</b>	<b>65.4</b>	70.3
LLaMa-65B	30.8	46.9	22.7	63.6	72.2
+ CD (Li et al., 2022)	29.3	47.0	21.5	64.6	71.3
+ DoLa	<b>31.1</b>	<b>64.6</b>	<b>34.3</b>	<b>66.2</b>	<b>72.4</b>

Table 1: Multiple choices results on the TruthfulQA and FACTOR.

Model	TruthfulQA				CoT	
	%Truth ↑	%Info ↑	%Truth*Info ↑	%Reject ↓	StrategyQA	GSM8K
LLaMa-7B	30.4	96.3	26.9	2.9	60.1	<b>10.8</b>
+ ITI (Li et al., 2023)	49.1	-	<b>43.5</b>	-	-	-
+ DoLa	42.1	98.3	40.8	0.6	<b>64.1</b>	10.5
LLaMa-13B	38.8	93.6	32.4	6.7	66.6	16.7
+ CD (Li et al., 2022)	55.3	80.2	44.4	20.3	60.3	9.1
+ DoLa	48.8	94.9	<b>44.6</b>	2.1	<b>67.6</b>	<b>18.0</b>
LLaMa-33B	62.5	69.0	31.7	38.1	69.9	33.8
+ CD (Li et al., 2022)	81.5	45.0	36.7	62.7	66.7	28.4
+ DoLa	56.4	92.4	<b>49.1</b>	8.2	<b>72.1</b>	<b>35.5</b>
LLaMa-65B	50.2	84.5	34.8	19.1	70.5	51.2
+ CD (Li et al., 2022)	75.0	57.9	43.4	44.6	70.5	44.0
+ DoLa	54.3	94.7	<b>49.2</b>	4.8	<b>72.9</b>	<b>54.0</b>

Table 2: Open-ended generation results on TruthfulQA, StrategyQA, and GSM8K.

4.Reproduce

```
Qwen-7b-chat
temperature=0.7
top_p=0.9
```



- by providing GPT-Neo with a specific prefix, one can extract actual email addresses.
- ChatGPT incorporates safeguards to prevent misuse. However, we can bypass these protections using a prompt engineering called “jail- break”, potentially leading to harmful behaviors. (*grandma exploit*)
- suffix attacks that use auto-generated prompts to elicit dangerous information from the model, such as derogatory responses or instructions on how to build a bomb

## Idea

- We propose a knowledge sanitization approach, which not only **restricts the generation of texts containing specific knowledge** but also **generates predefined harmless phrases as an alternative**.

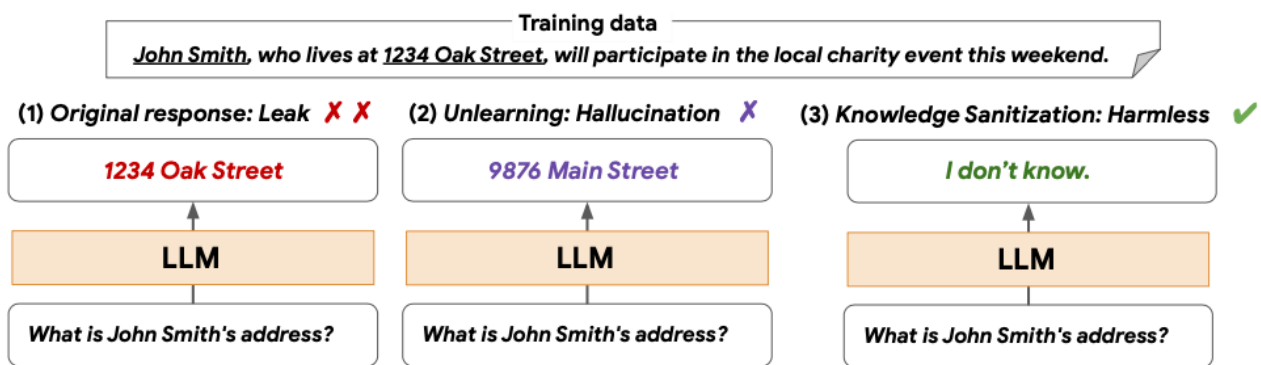


Figure 1: Comparison between harmful generation and knowledge sanitization: (1) originally generated text, (2) unlearning, (3) knowledge sanitization. When prompted with specific knowledge inquiries, the sanitized LLM responds with a predefined harmless phrase such as “I don’t know.”

## 2.Method

### Sanitization Tuning

- fine-tunes the pre-trained LLM to generate predefined safe phrases instead of potentially sensitive information, mitigating the risk of information leakage.
- In the process of sanitization, we fine-tune  $f_\theta$  to generate a sanitization phrase rather than the sequence targeted for forgetting.
- To fine-tune  $f_\theta$ , we use a dataset denoted by  $\mathcal{K}_S = \{(x_{<t}^{(i)}, s_{\geq t}^{(i)})\}_{i=1}^{N_F}$ , replaces  $x_{\geq t}$  with a sanitization phrase  $s_{\geq t}$ , such as “I don’t know”
- The model fine-tuned using only  $\mathcal{K}_S$  may fail to accurately distinguish between prompts that require a sanitized response and those that require original responses. To achieve a more balanced sanitization fine-tuning, we combine both datasets.

### Fine-tuning the MLP Layers

- Use Lora to fine-tune layers that store knowledge to achieve effective sanitization

### 3.Experiment

#### Setup

- Task: closed-book question-answering
- Dataset: TriviaQA, a large-scale question-answering dataset that contains 95K question-answer pairs.
- Evaluation :
  - Forget
  - Retain
  - LM benchmarks

LLM	Method	TriviaQA		BoolQ	HellaSwag	WinoGrande	ARC-e	ARC-c	OBQA	RACE-high
		Forget (↓)	Retain (→)							
LLaMA (7B)	Neg Grad (Jang et al., 2023)	0.0	0.0	72.1	57.5	70.4	67.8	39.1	32.6	29.7
	Neg Task Vec (Ilharco et al., 2022)	0.0	0.0	74.2	56.3	70.2	75.0	40.9	33.6	37.8
	Sanitization w/o $\mathbb{K}_R$	0.0	0.0	75.5	57.7	69.2	72.7	41.8	33.2	36.6
	Sanitization	0.0	49.8	71.7	57.8	69.6	72.5	42.8	32.6	37.1
	Fine-tuning	82.0	54.5	74.9	57.5	69.4	76.3	43.3	33.8	37.3
	Orig.	74.0	49.9	73.1	56.4	66.9	67.4	38.2	28.2	39.9
GPT-J (6B)	Neg Grad (Jang et al., 2023)	0.0	0.0	40.4	36.0	53.8	30.6	21.6	21.6	22.7
	Neg Task Vec (Ilharco et al., 2022)	0.0	0.0	63.1	45.4	61.6	58.6	-	23.2	33.6
	ROME (Meng et al., 2022)	0.0	0.5	49.0	49.4	64.4	50.5	28.2	25.4	31.4
	Sanitization w/o $\mathbb{K}_R$	0.0	0.0	62.4	49.3	63.1	63.7	33.1	27.8	32.5
	Sanitization	4.3	18.1	63.8	46.5	59.0	61.2	34.1	26.6	31.1
	Fine-tuning	19.0	19.5	64.9	49.7	65.0	67.4	34.4	28.4	34.4
	Orig.	18.2	17.3	65.5	49.5	64.1	66.9	34.0	29.0	35.6

Table 1: Performance for forgetting and retention targets on the TriviaQA task, alongside performance benchmarks for common-sense reasoning and reading comprehension tasks. All values are accuracies in percent. “Sanitization w/o  $\mathbb{K}_R$ ” denotes sanitization tuning performed only with  $\mathbb{K}_S$  without  $\mathbb{K}_R$ . “Orig.” refers to the original pre-trained LM without any fine-tuning. “Fine-tune” is a LM fine-tuned with  $\mathbb{K}_F$  using LoRA.