

# Consistency text similarity on the example of the task of recognizing hallucinations of language models

Kseniia Petrushina

Moscow Institute of Physics and Technology,  
Skolkovo Institute of Science and Technology

*Scientific supervisor:* Alexander Panchenko

2023

# Introduction

## Explanation

**Hallucination** of the language model is a grammatically correctly generated response, which, however, contains incorrect information.

## Examples

Paraphrase generation & Machine translation – different meaning.  
Definition modeling – deviation from the database.

## Problem statement

*Language model* is a function

$$\mathbf{f} : \mathcal{P}(\mathbf{T}_s^{L_s}) \rightarrow \mathcal{P}(\mathbf{T}_h^{L_h}),$$

$\mathbf{s}_i$  is called *source sentence* and  $\mathbf{h}_i$  is called *model hypothesis*.

We can define function

$$\mathbf{f}^{-1} : \mathcal{P}(\mathbf{T}_h^{L_h}) \rightarrow \mathcal{P}(\mathbf{T}_s^{L_s})$$

Then it is said that  $\mathbf{h} = \mathbf{f}(\mathbf{s})$  is a *hallucination* of the language model  $\mathbf{f}$  with the input  $\mathbf{s}$  if

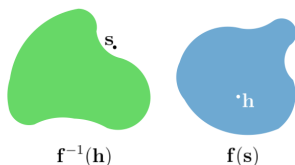
$$p(\mathbf{f}^{-1}(\mathbf{f}(\mathbf{s})) = \mathbf{s}) = 0.$$

# Problem statement

The task of recognizing hallucinations is to find a function  $\text{sim} : \mathbf{T}_s^{L_s} \times \mathbf{T}_h^{L_h} \rightarrow [0, 1]$ , such that

$$\mathbb{E}_{\mathbf{s}_i \sim \mathbf{T}_s^{L_s}, \mathbf{h}_i \sim f(\mathbf{s}_i)} \{ \mathbb{I}[\text{sim}(\mathbf{s}_i, \mathbf{h}_i) \geq \text{thr}] = y_i \} \rightarrow \max_{\text{sim}, \text{thr}},$$

where  $y_i$  denotes the presence of a hallucination.



**Figure:** An illustration of a model's hallucination.  $\mathbf{s}$  does not belong to the set of possible outputs of  $\mathbf{f}^{-1}(\mathbf{h})$

# Existing solutions

## 1. Words or characters n-grams

$$\text{sim}_{\text{BLEU}}(\mathbf{s}, \mathbf{h}) = \frac{|N_s \cap N_h|}{|N_h|}$$

## 2. Similarity between static embeddings

$$\text{sim}_{\text{cos}}(\mathbf{s}, \mathbf{h}) = \cos(\mathbf{v}_s, \mathbf{v}_h)$$

## 3. Similarity between contextualized embeddings

$$R = \frac{1}{L_s} \sum_{v_i \in \mathbf{v}_s} \max_{\hat{v}_j \in \mathbf{v}_h} v_i^T \hat{v}_j \quad P = \frac{1}{L_h} \sum_{\hat{v}_j \in \mathbf{v}_h} \max_{v_i \in \mathbf{v}_s} v_i^T \hat{v}_j$$

$$\text{BERTScore} = 2 \frac{PR}{P + R}$$

## 4. Similarity between embeddings from bi-encoders

$$sim_{bi-enc}(\mathbf{s}, \mathbf{h}) = \cos(\mathbf{enc}_s(\mathbf{s}), \mathbf{enc}_h(\mathbf{h}))$$

## 5. Symmetric and asymmetric cross-encoders

$$sim_{cross-enc}(\mathbf{s}, \mathbf{h}) = \mathbf{clf}(\mathbf{enc}(\mathbf{s}, \mathbf{h}))$$

In the general case, the similarity function should be defined for objects from different spaces  $\mathbf{T}_s^{L_s}$  and  $\mathbf{T}_h^{L_h}$ .

The existing methods do not investigate whether there is enough information in  $\mathbf{h}$  to restore  $\mathbf{s}$ .

# Computational experiment

We are given the dataset

$$\mathcal{D} = \{(\mathbf{s}_i, \mathbf{h}_i, y_i)\}_{i=1}^N, \quad \mathbf{h}_i \in \mathbf{f}(\mathbf{s}_i), \quad y_i \in \{0, 1\}$$

The target variable  $y_i$  indicates the occurrence of a hallucination in the  $\mathbf{f}$  model at the input of  $\mathbf{s}_i$  and the output of  $\mathbf{h}_i$ .

1. The proportion of correct predictions:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[\hat{y}_i \geq \text{thr}] = y_i$$

2. Spearman's rank correlation coefficient:

$$r_s = \rho_{R(Y), R(\hat{Y})} = \frac{\text{cov}(R(Y), R(\hat{Y}))}{\sigma_{R(Y)} \sigma_{R(\hat{Y})}}$$

# Results

## ► Paraphrase generation task

| Method         | Accuracy $\uparrow$ | $r_s$ $\uparrow$ |
|----------------|---------------------|------------------|
| $sim_{bi-enc}$ | 0.808               | 0.153            |
| $sim_C$        | <b>0.824</b>        | <b>0.186</b>     |

Table: Hallucination recognition results in the PG task

## ► Machine translation task

| Method            | Accuracy $\uparrow$ | $r_s$ $\uparrow$ |
|-------------------|---------------------|------------------|
| $sim_{LaBSE}$     | 0.786               | 0.592            |
| $sim_{BLASER-QE}$ | <b>0.802</b>        | <b>0.605</b>     |

Table: Hallucination recognition results in the MT task



# Conclusion

- ▶ Analysis of the existing measures of textual similarity.
- ▶ New method that corrects the disadvantages of the previous ones.
- ▶ Hypotheses about properties of the *consistency similarity measure*.
- ▶ Computational experiments with different measures.

## **Future work:**

- ▶ Theoretical justification of hypotheses.
- ▶ Extend the method to work with an external database
- ▶ Conduct comprehensive ablation study.

# Literature

1. David Dale et al. 2023. [HalOmi: A Manually Annotated Benchmark for Multilingual Hallucination and Omission Detection in Machine Translation](https://arxiv.org/pdf/2305.11746.pdf).  
URL: <https://arxiv.org/pdf/2305.11746.pdf>
2. Nikolay Babakov et al. 2022. [A large-scale computational study of content preservation measures for text style transfer and paraphrase generation](https://aclanthology.org/2022.acl-srw.23.pdf).  
URL: <https://aclanthology.org/2022.acl-srw.23.pdf>
3. Ashish Vaswani et al. 2017. [Attention is All you Need](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)  
URL: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)