

Link do artykułu	Tytuł artykułu	Autorzy	Autorski komentarz	Przeprowadzone eksperymenty	Dostępny kod	Dostępne pre-trenowane modele	Metryki ewaluacji	Wykorzystane zasoby obliczeniowe
Artykuł: <a href="https://arxiv.org/abs/2109.07958">https://arxiv.org/abs/2109.07958</a>  Repo: <a href="https://github.com/sylinrl/TruthfulQA">https://github.com/sylinrl/TruthfulQA</a>	<b>TruthfulQA: Measuring How Models Mimic Human Falsehoods</b>	Stephanie Lin, Jacob Hilton, Owain Evans	Artykuł przedstawia badania zachowywania się modeli przy odpowiadaniu na pytania, które celowo prowokują do błędnych odpowiedzi. To tzw. Halucynacje imitacyjne - oparte o zakorzenione w społeczeństwie nieprawdy, mity, mądrości ludowe, memy itp. - chodzi o to, że jest szansa, że przeniknęły one podczas trenowania modelu. Wynika to z tego, że model zwraca odpowiedzi najbardziej prawdopodobne, więc jeśli jakąś nieprawdę występuje często, to zostanie "wyuczona".  Ważne – pytania nie są “trickowe” – nie próbują oszukać modelu zawiątą formą, dziwną składnią, błędami językowymi, niejasnościami itp. Zawierają jedynie kontrowersyjne tezy (chlopski rozum, uliczne mądrości...). Dodatkowo był eksperyment z parafrazami i matched controls który potwierdził, że to nie forma a treść jest źródłem halucynacji w badaniu.  Artykuł zawiera ciekawe przygotowane prompty, dobrze opisane wyniki i na pewno będzie bardzo pomocny do przeprowadzania dalszych badań.	<ul style="list-style-type: none"> <li>▪ testy zero-shot (bez strojenia) gdzie odpowiadano na pytanie pełnym zdaniem, a human evaluation sprawdzało truthfulness (% poprawnych odpowiedzi) i informativeness (czy zawarta była sensowna informacja zwrotna, bez tautologii czy odpowiedzi 'nie wiem'). Najważniejszy wniosek – im większy model, tym bardziej halucynuje. Najlepszy (truthfulness) był GPT-3 (58% prawdy – dla porównania człowiek 94%)</li> <li>▪ eksperymenty multiple-choice: podawano pytanie i zestaw odpowiedzi (true i false) – model miał oceniać likelihood poprawnych odpowiedzi – wynik – lekko gorszy od losowego zgadywania</li> <li>▪ kontrola źródeł halucynacji <ul style="list-style-type: none"> <li>○ matched controls – edycja 1-3 słów w pytaniu by sprawdzić czy styl pytania ma znaczenie. Efekt - niewielki wpływ.</li> <li>○ parafraszowanie, efekt - prawie żaden wpływ</li> </ul> </li> <li>▪ automatyczna ewaluacja (GPT-judge) – przygotowanie fine-tuningowanego modelu GPT-3 6.7B jako klasyfikatora true/false. Efekt – 90-96% zgodności z oceną człowieka</li> </ul>	Dostępne repo na githubie, zawiera pytania, referencje i etykiety; narzędzia do ewaluacji; implementację multi-choice (podane pytanie i odpowiedzi, model ma wskazać dla każdej p-stwo że jest prawdziwa); treść promptów i konfigurację.	Open source: GPT-Neo/J; GPT-2; UnifiedQA niedostępne: poprzez API openAI: GPT-3 i GPT-3 jako GPT-judge	<ul style="list-style-type: none"> <li>▪ ewaluacja ręczna (human) dobrze opisana, sensowna;</li> <li>▪ ewaluacja automatyczna - GPT-judge, niestety niedostępna (dostępny do tego celu GPT-3 dostępny przez API)</li> </ul>	Modele już przetrenowane, uruchamiane poprzez HuggingFace, raczej do zrealizowania na laptopie nawet 24GB ramu, GPT-judge dostępny poprzez API na serwerze od openAI (nieznany)

Artykuł: <a href="https://arxiv.org/abs/2401.03205">https://arxiv.org/abs/2401.03205</a>	„The Dawn After the Dark: An Empirical Study on Factuality Hallucination in LLMs	Junyi Li, Jie Chen1, Ruiyang Ren, Xiaoxue Cheng, Wayne Xin Zhao, Jian-Yun Nie, Ji-Rong Wen1	Artykuł to duża analiza halucynacji w kontekście detekcji, poznania ich źródeł oraz zmniejszania częstotliwości ich występowania. Analiza skupia się na całym procesie budowania LLM (pre-training, supervised fine-tuning, RLHF, interference, a także promptu i parametrów generacji)  W badaniu wykorzystywany jest ulepszony benchmark HaluEval 2.0 zawierający 8770 pytań z różnych domen. Badane są zarówno modele open-source jak i komercyjne closed-source.  Wnioskiem z badań jest m.in. że każdy element pipelinu ma wpływ na powstawanie halucynacji, szczególnie jakość danych treningowych oraz dobrany prompt. Najsłuszniejszymi metodami zapobiegania są RLHF, RAG (RAG szczególnie dla dziedzin ścisłych, RLHF dla open-domain), beam search i dobry prompt. Self refleciton działa dobrze tylko dla małych modeli.  Badania bardzo kompleksowe, poruszające prawie wszystkie aspekty działania LLM. Dodatkowo sprawdzano wiele modeli o różnych architekturach i historii rozwoju.  Modele komercyjne (ChatGPT, Claude 2) halucynują mniej niż open-source (Vicuna, Llama-2-Chat), choć nawet komercyjne osiągają MaHR 40-50%	<ul style="list-style-type: none"> <li>▪ konstrukcja benchmarku HaluEval 2.0 – 8770 pytań z domen: biomed, finanse, nauka, edukacja, open-domain. Do generowania pytań użyto BERTScore, GPT-3.5/ChatGPT - tak by sprawdzić, które pytania są halucynogenne</li> <li>▪ automatyczna detekcja halucynacji metodą: ekstrakcja faktów z odpowiedzi → weryfikacja faktów. Zgodność z ocenami człowieka 91-95%</li> <li>▪ Wpływ skali pre-treningu – czy liczba tokenów redukuje halucynacje. Wyniki różne</li> <li>▪ Wpływ składu danych pre-treningu – różna specjalizacja różnych korpusów, dane domenowe (naukowe) halucynują mniej w swojej domenie</li> <li>▪ wpływ częstości wiedzy – jak często występuje halucynacja a częstość występowania encji w Wikipedii (im częściej tym mniej halucynacji)</li> <li>▪ wpływ supervised fine-tuning – im bardziej skoncentrowane instrukcje tym więcej błędów</li> <li>▪ wpływ projektowania promptu – xxxxxxxx</li> <li>▪ sposób dekodowania a halucynacje (greedy, top-k, top-p, beam-search). Top-p/k zwiększa halucynacje w domenach profesjonalnych; beam-search najstabilniejsze.</li> <li>▪ wpływ kwantyzacji</li> <li>▪ analiza sekwencyjna (token po tokenie) – w którym momencie rodzi się halucynacja i jak się propaguje. Efekt snowballing – jeden błędny token powoduje całą błędą odpowiedź.</li> <li>▪ Mitygacja – badanie metod RLHF (poprawa widoczna)</li> </ul>	Dostępne repo na githubie, zawiera pytania (benchmark), pipeline do uruchamiania eksperymentów, konfigurację modeli, skrypty do modeli open-source	<p><b>Open-source</b></p> <ul style="list-style-type: none"> <li>▪ Alpaca 7B</li> <li>▪ Vicuna 7B/13B</li> <li>▪ YuLan-Chat 13B</li> <li>▪ Llama-2-Chat 7B/13B</li> <li>▪ Falcon 40B</li> <li>▪ Galactica 30B</li> <li>▪ GPT-NeoX 20B</li> </ul> <p><b>Closed-source</b></p> <ul style="list-style-type: none"> <li>▪ text-davinci-002 / 003</li> <li>▪ ChatGPT</li> <li>▪ Claude, Claude 2</li> </ul>	Ocena halucynacji metodą automatyczną, ekstrakcja faktów z uzyskanej odpowiedzi → weryfikacja przez duże komercyjne modele (GPT3.5, ChatGPT), które mają zgodność z ocenami człowieka 91-95%.  Jako halucynacje klasyfikowano tylko fałsz (bez ‘nie wiem’).  Dwa wskaźniki halucynacji:	<ul style="list-style-type: none"> <li>▪ MiHR (odsetek fałszywych faktów w odpowiedzi)</li> <li>▪ MaHR (Odsetek wypowiedzi, w których wystąpiła halucynacja)</li> </ul>
---	--	---	---	---	--	--	---	---

				szczególnie w open-domain), RAG (poprawa w naukach ścisłych), self-reflection (działa tylko dla dużych modeli, dla mniejszych mało zauważalny), beam-search (generuje mniej halucynacji niż greedy/top-p)			
--	--	--	--	---	--	--	--