# Laura Perez-Beltrachini

Research Fellow

University of Edinburgh

# Addressing Hallucinations in Generative AI for Customer Service Applications

Innovate UK 10092998 / Collaborative R&D
Accelerating trustworthy AI
Feb 2024 – Feb 2025

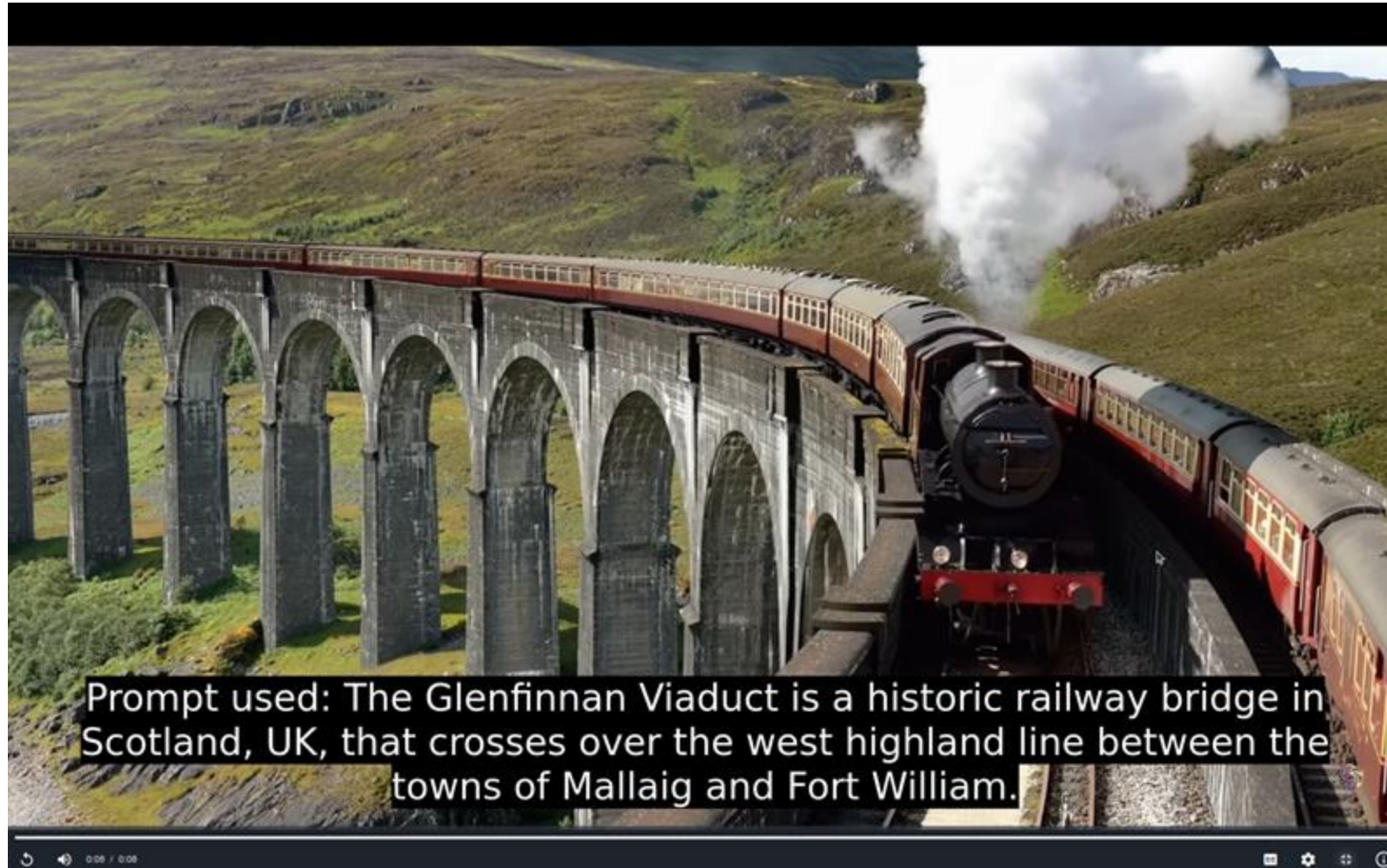Innovate UK BridgeAI Regional Roadshow Series

12th May 2025

**Laura Perez-Beltrachini**

THE UNIVERSITY of EDINBURGH

Edinburgh NLP
University of Edinburgh
Natural Language Processing

# Hallucinations in Generative AI



Prompt used: The Glenfinnan Viaduct is a historic railway bridge in Scotland, UK, that crosses over the west highland line between the towns of Mallaig and Fort William.

Generated by Sora (OpenAI's text-to-video model)
https://en.wikipedia.org/wiki/Hallucination_(artificial_intelligence)

3

# Hallucinations in Large Language Models (LLMs)

## Chatbots May 'Hallucinate' More Often Than Many Realize

When summarizing facts, ChatGPT technology makes things up about 3 percent of the time, according to research from a new start-up. A Google system's rate was 27 percent.

Amr Awadallah, the chief executive of Vectara, warns that its chatbot software doesn't always tell the truth. Cayce Clifford for The New York Times
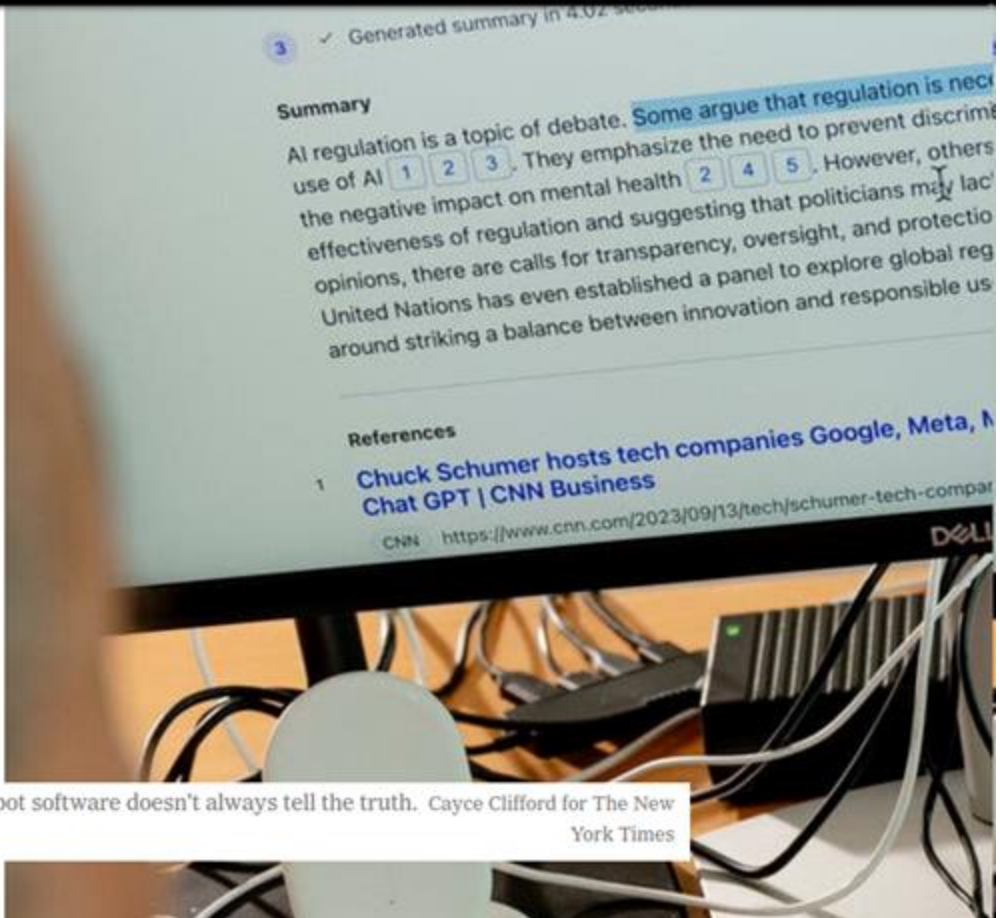
4

# R&D to Address Hallucinations in LLMs for Customer Service Chatbots

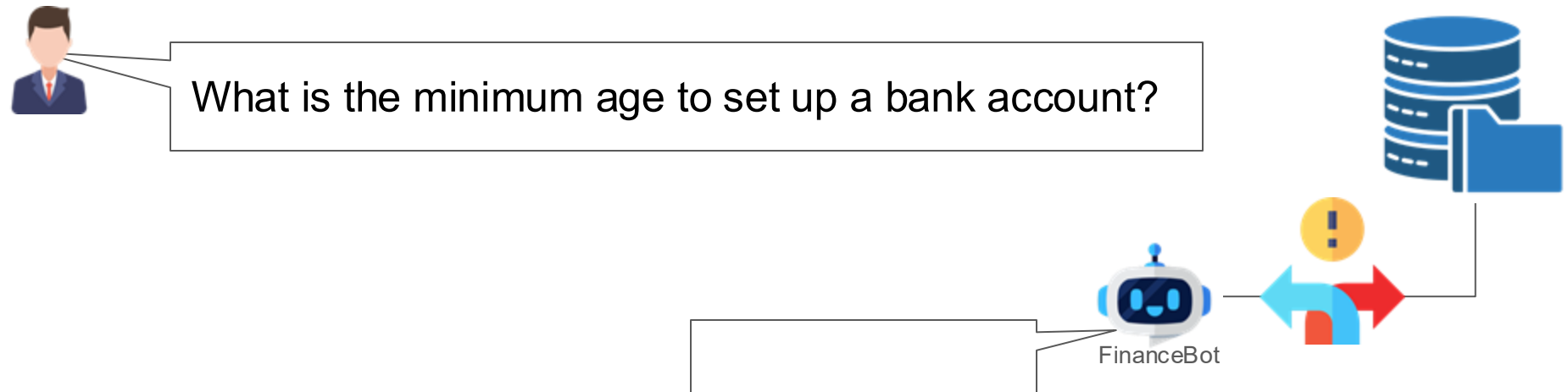| | | |
|---|---|---|
| Large company | NatWest / Algomo | banking dataset, piloting solution, and providing customer access |
| SME | yu life / THE UNIVERSITY of EDINBURGH / University of Essex (Lead Participant) | core conversational AI platform and productization. technical implementation and end-user of the solution |

# Our Research on Detecting Hallucinations

**generation fact-checking**

[Fine-grained natural language inference based faithfulness evaluation for diverse summarisation tasks](#) (Zhang, Xu, Perez-Beltrachini, EACL 2024)

[Leveraging entailment judgements in cross-lingual summarisation](#) (Zhang and Perez-Beltrachini, ACL 2024)

**model uncertainty**

**[Uncertainty Quantification in Retrieval Augmented Question Answering](#)** (Perez-Beltrachini and Lapata, Under Submission)

THE UNIVERSITY of EDINBURGH

# Retrieval Augmented Question Answering

What is the minimum age to set up a bank account?

FinanceBot

# Retrieval Augmented Question Answering

What is the minimum age to set up a bank account?

FinanceBot

# Hallucinations and Answer Uncertainty



**Hallucinated Answer**

**Safe Answer**

9

# Answer Uncertainty Estimation?



What is the minimum age to set up a bank account?

Passage Retriever

FinanceBot

PPL("There is no limit!")= 0.5

- Perplexity

Information-Theoretic

# Answer Uncertainty Estimation?

What is the minimum age to set up a bank account?

Passage Retriever

FinanceBot

Output Variation

Regular Entropy - Semantic Entropy

18 years old.

bank dependant.

16 years old.

24 years old.

1.5

# Answer Uncertainty Estimation?

# Answer Uncertainty via Passage Utility

What is the minimum age to set up a bank account?

Passage Retriever

0.6
0.2
0.01
0.6

FinanceBot

Passage Utility Predictor

Our Approach

1° observe QA performance on individual passages and target LLM(θ), Passage Utility.

2° train a secondary small model to predict Passage Utility.

3° estimate Answer Uncertainty from Passage Utilities.

13

# Answer Uncertainty Estimation?

# Successful and Lightweight Hallucination Detection

Natural Questions, TriviaQA, WebQuestions, SQuAD, PopQA, RefuNQ

OpenDomainBot

Llama-3.1-8B, Gemma2-9B, Mistral-v0.3-7B, Gemma2-2B, Gemma2-27B

# Improved Accuracy via Selective Answering



X-axes Gemma2 -2B, -9B, -27B.
Average AccLM (LLM- based accuracy) is across the six QA tasks.
Black dots indicate accuracy when always answering.

# Going Forward

- Generalisation of the Passage Utility predictor:
  - Unseen QA tasks and LLMs,
  - Multi-Hop reasoning,
  - Long-form and subjective answers.

- Answering ambiguous questions.

FinanceBot

Thank you!

# Observed Low/High Passage Utility

# Train a Passage Utility Predictor

# Answer Uncertainty via Passage Utility

# Answer Uncertainty via Passage Utility

What is the minimum age to set up a bank account?

Passage Retriever

Apply online. You must be between 13 and 14 years old to apply.

**0.01**

You must be aged between 11 and

The specific requirements for opening a bank account vary from bank to bank.

**0.2**

**0.8**

FinanceBot

**0.8**

# Prediction of Hallucination

## Llama3.1-8B & question from SQuAD



The gold standard answer is "Lockheed Aircraft Corporation."

# Prediction of Correct Answer
# Llama3.1-8B & question from SQuAD



The gold standard answer is "Millions of youth who had migrated to the cities."

23

brainfit™
powered by FCLABS

Andy Low
andy@fclabs.co.uk
COO
brainfitfirst.com