





TEAM DELPHI



Ivo Canjels



Dennis Dams



Richard Doornbos



Emile van Gerwen



Christos Kitsanelis



Anca Lichiardopol



Micha Lipplaa



JeanLuc Oudshoorn



Monika Srivastava



Erik Teesink

hilips Image Guided Therapy

Philips, our purpose is to improve people's health and well-being through meaningful novation. We aim to improve 2.5 billion lives per year by 2030, including 400 million in derserved communities.





PHILIPS - MARKET LEADER IN INTERVENTIONAL X-RAY SYSTEMS

- Philips IGT (Image Guided Therapy) is worldwide market leader in high-end interventional X-ray systems
- Used in patient treatment in minimally invasive procedures:
 - Cardiology
 - Neurology
 - Oncology
 - Other clinical domains
- Installed base > 30,000 systems



Azurion System

QUALITY CONTROL UNDER REGULATORY OVERSIGHT IS EXPENSIVE

The current reality in managing customer feedback:

- Installed base generates >200,000 user feedback reports (yearly)
 - Large teams needed for analysis → Millions €
- Analysis must be systematic to ensure correction and prevention
 - Handle information overload: large amounts of unstructured text
 - Manage human interpretation: analysis tends to be subjective

Examples of quality processes that are subject to regulatory oversight:

- Customer feedback and complaints
- Enhancement Requests
- CAPA (Corrective and Preventive Actions)
- PPMRR (Production/Post-Market Risk Reviews)



CAN AI HELP?

Opportunity: Improve quality while reducing cost

- Can save > 1 Million per year
- By AI support during analysis of feedback, complaint and CAPA data

<u>Challenge</u>: How to ensure trustworthiness of Al's answers?

- Correctness?
- Reproducibility?
- Validation?

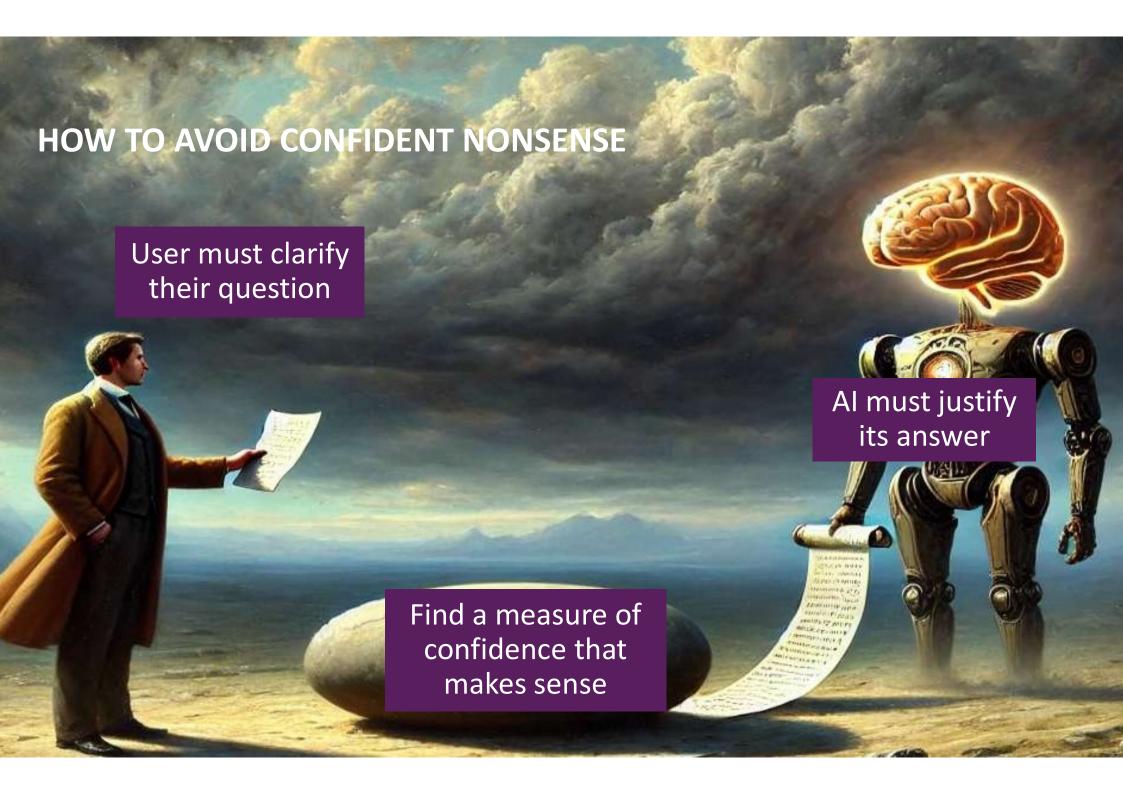




Project DELPHI – Joint initiative Philips IGT & TNO-ESI

- Experiment e.g. which aspects of the quality process can be supported by AI?
- Measure e.g. what is the ground truth to compare to?
- Validate e.g. how do we convince ourselves and the authorities?





LET'S LOOK AT A RECENT USE CASE

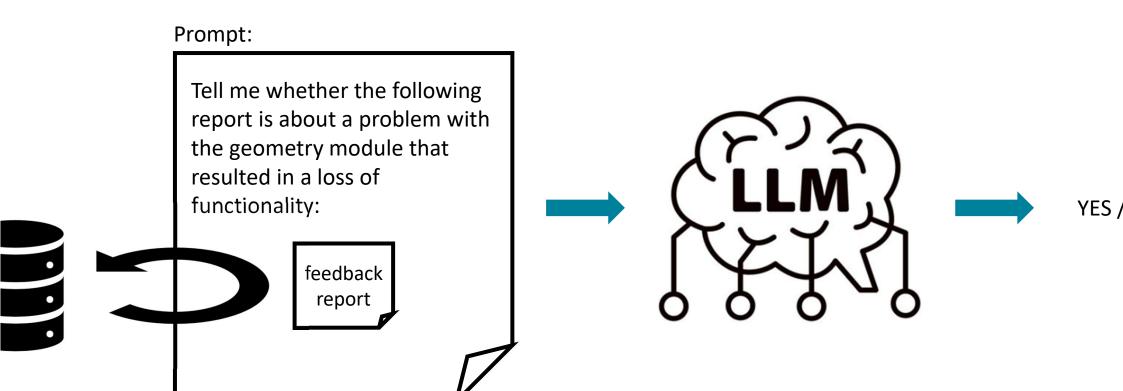


Geometry Module

"How many feedback reportance about a problem with geometry module that resin a loss of functionality?"



WE LET AN LLM LOOK AT EACH REPORT



"Complex classification"

THE LLM NEEDS TO ACCOMPANY ITS ANSWER BY A RATIONALE

Prompt:

Tell me whether the following report is about a problem with the geometry module that resulted in a loss of functionality:

feedback report 4879421



YES, because:

Malfunctioning of shutter button

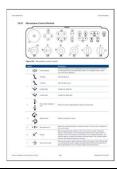
YOU NEED TO TEACH THE LLM SOME DOMAIN-SPECIFIC KNOWLEDGE

Prompt:

Tell me whether the following report is about a problem with the geometry module that resulted in a loss of functionality:

feedback report 4879421

Documentation of the module:





YOU NEED TO EXPLAIN TO THE LLM HOW TO READ THE REPORT

Work Order #: WO-11739790

12-Nov-2024

The system meets the specification for the performed service and is returned to use

Resolution: Patient entry option not available.

Resolution: Reload DFI and Reset SCPU. Reconnect subsystem.

Performed calibration, tested and handed over machine to user in working condition

Problem Description by Engineer: Problem description by engineer:

system is not booting

Tools & Equipment: MULTIMETER,2025-02-19 00:00:00

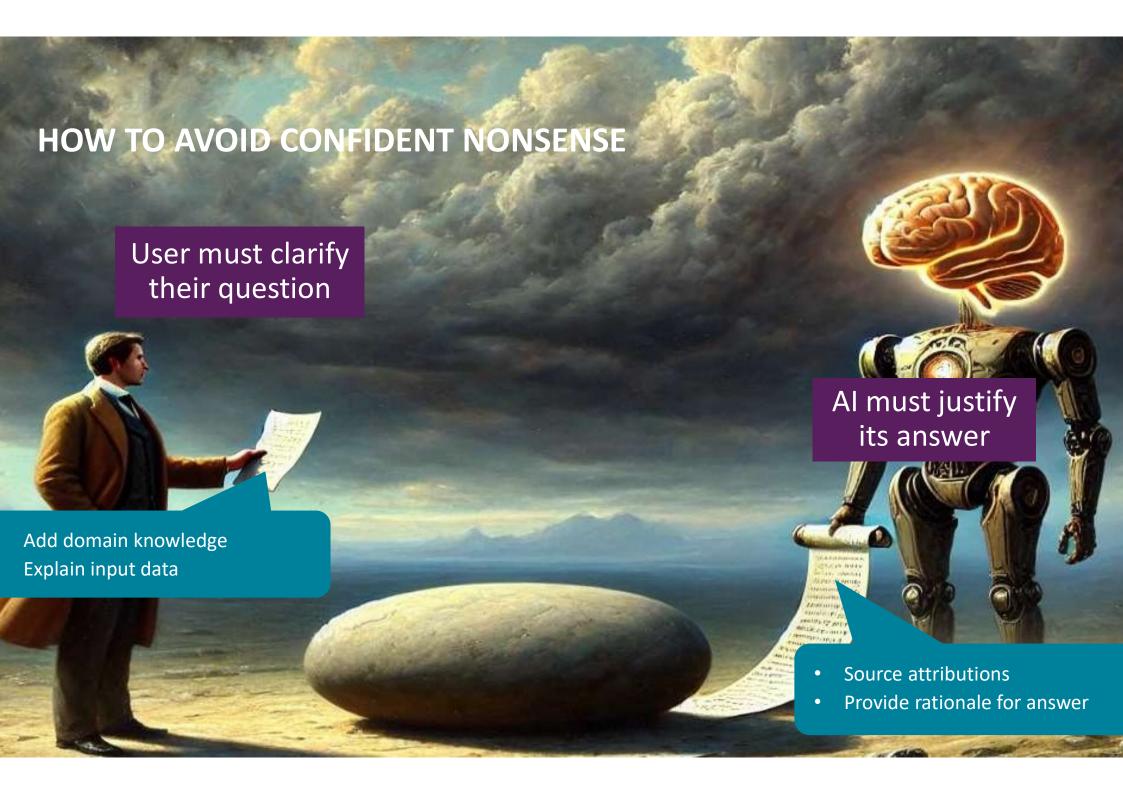
Work Order #: WO-11739412

The system does not meet the specification for the performed service and is not returned to use

Problem Description by Engineer: Problem description by engineer:

system is not booting

Diagnostic performed by Engineer: Checked and found DFI is not booting.





LOOKING FOR A MEASURE OF CONFIDENCE THAT MAKES SENSE

• Idea 1: Let the LLM to answer "DONTKNOW" in case it's not really sure

• Idea 2: Ask the LLM how sure it is about its answer (low/medium/high confidence)

THE LLM WAS TAUGHT TO GUESS RATHER THAN TO ADMIT UNCERTAINTY

• Idea 1: Let the LLM to answer "UNKNOWN" in case it's not really sure

Idea 2: Ask the LLM how sure it is about its answer (low/medium/high confidence)

Why Language Models Hallucinate

Adam Tauman Kalai* Ofir Nachum Santosh S. Vempala[†] Edwin Zhang OpenAI OpenAI Georgia Tech OpenAI

September 4, 2025

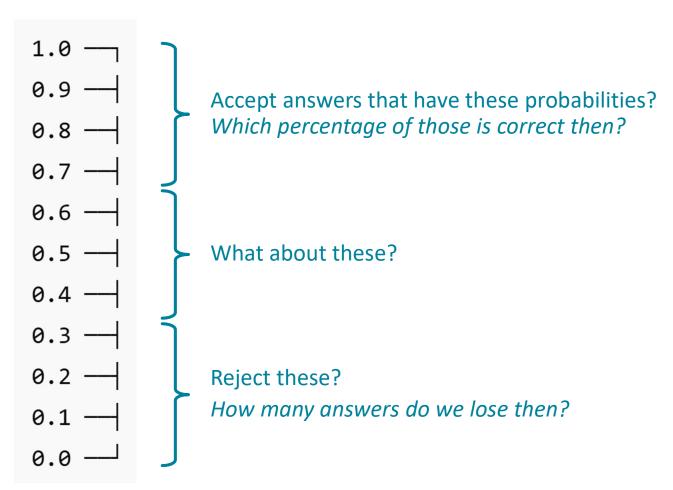
Abstract

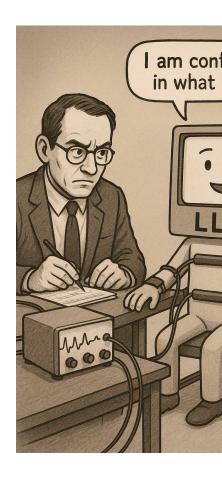
Like students facing hard exam questions, large language models sometimes guess when uncertain, producing plausible yet incorrect statements instead of admitting uncertainty. Such "hallucinations" persist even in state-of-the-art systems and undermine trust. We argue that language models hallucinate because the training and evaluation procedures reward guessing over acknowledging uncertainty, and we analyze the statistical causes of hallucinations in the modern training pipeline. Hallucinations need not be mysterious—they originate simply as errors in binary classification. If incorrect statements cannot be distinguished from facts, then hallucinations in pretrained language models will arise through natural statistical pressures. We then argue that hallucinations persist due to the way most evaluations are graded—language models are optimized to be good test-takers, and guessing when uncertain improves test performance. This "epidemic" of penalizing uncertain responses can only be addressed through a socio-technical mitigation: modifying the scoring of existing benchmarks that are misaligned but dominate leaderboards, rather than introducing additional hallucination evaluations. This change may steer the field toward more trustworthy AI systems.

Introduction

SO WE HOOKED UP THE LLM TO A POLYGRAPH

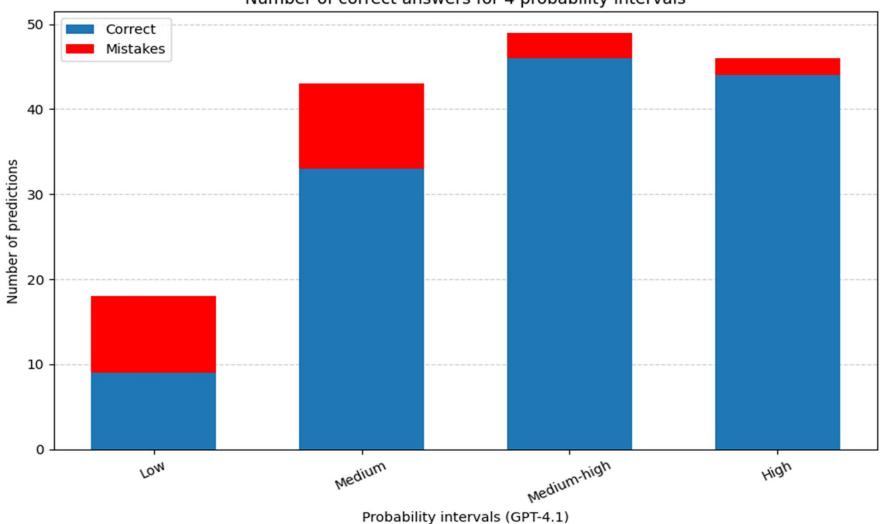
Idea 3: Measure the <u>token-level probability</u> of each answer





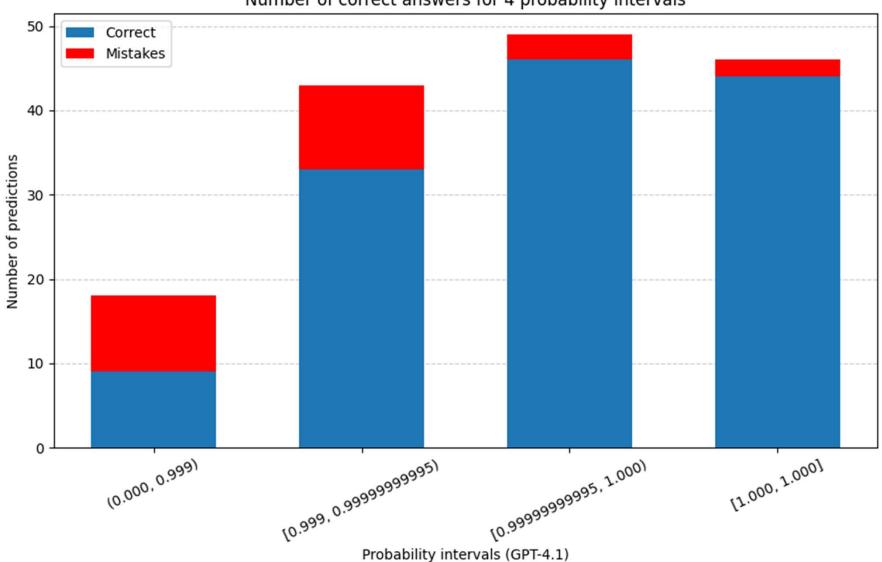
AT WHICH PROBABILITY VALUES CAN WE TRUST THE ANSWER?



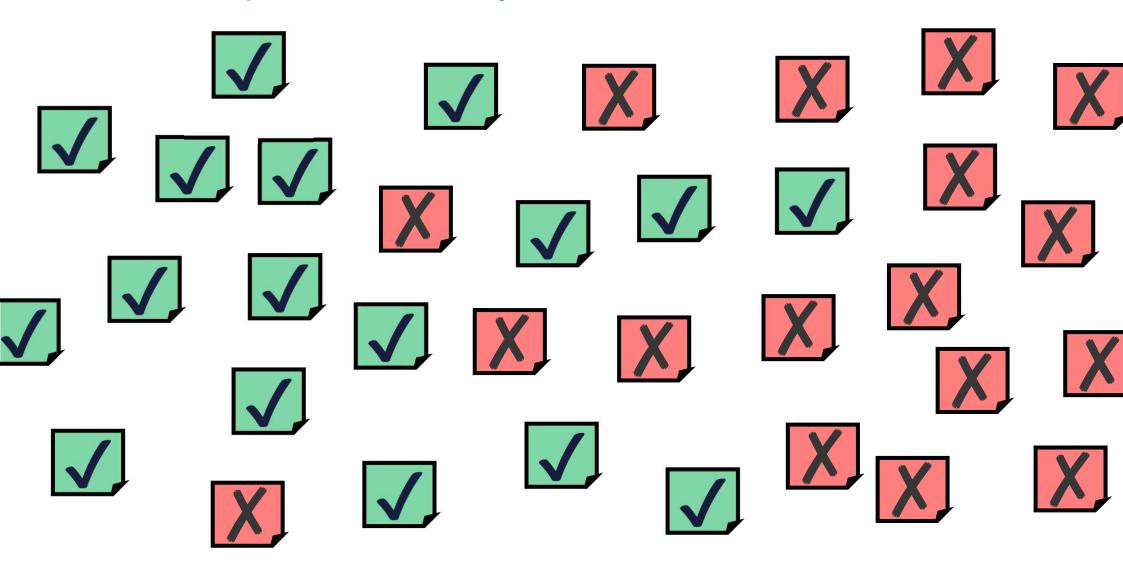


AT WHICH PROBABILITY VALUES CAN WE TRUST THE ANSWER?

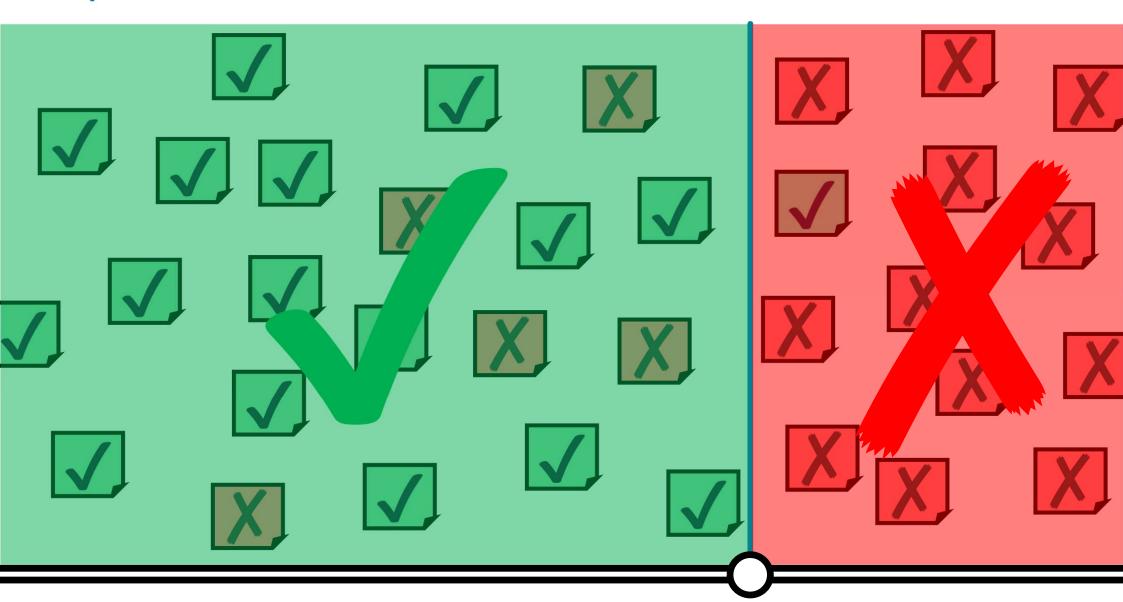




THE TEST SET (\checkmark = YES, X = NO)



YES/NO ASSESSMENT BY A SINGLE PROMPT



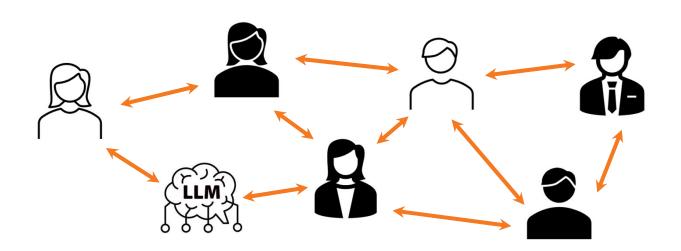
IDEA 4: CONFIDENCE LEVELS BY USING 2 CONSERVATIVE PROMPTS





SO WHAT DID WE LEARN?

- LLM-based analysis of customer feedback reports
 - reduces cost because fewer human specialists need to be called in
 - increases consistency because an LLM can tirelessly repeat the same analysis over 1000s of reports
- Trustworthiness of LLM output can be brought to an acceptable level
 - by applying the suggestions from this presentation...
 - ... and comparing the LLM to what human experts can realistically achieve



AND WHAT'S NEXT FOR US?

- Support trust building via conversational Prompt Refinement Interface (demo in the innovation market!)
- Bring the current prototype to a larger user group for further feedback
- Turn our experimental platform to the level of a validated tool
- And add device logs into the mix!

