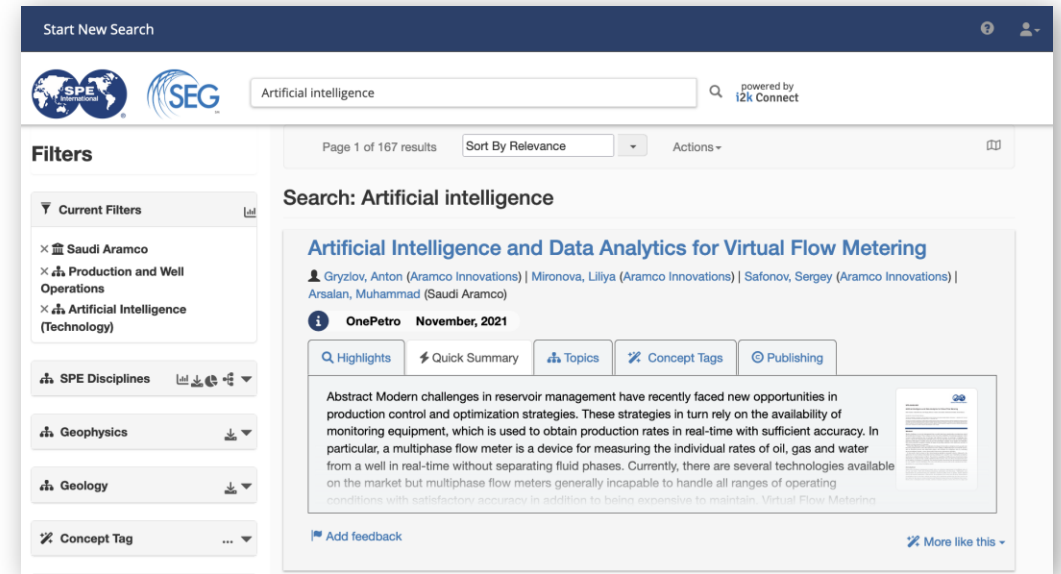
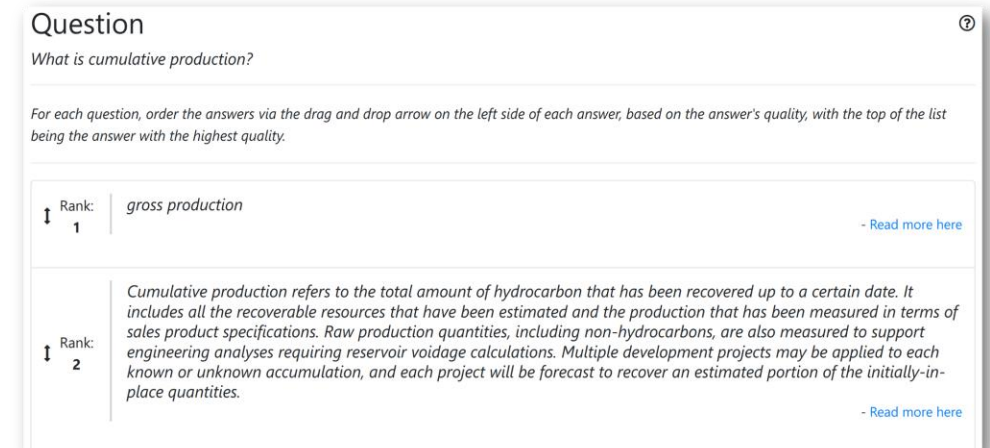


# i2k Connect

- Government experience
  - Cage Code: 70GV3 and Small Business
  - DARPA HR00112190013: **Fine-Grained Knowledge Delivery in Communities of Practice** (Seedling)
  - DARPA BAA HR001121S0034: **Knowledge Management at Scale and Speed** (SRI Prime Contractor)
  - IARPA BAA W911NF-23-S-0007 **REASON** Proposer
- Relevant capabilities
  - Experience applying SOTA ML techniques to information management and discovery
  - Extensive generative AI R&D
  - Experience deploying LLMs in security-sensitive contexts
  - Advanced document processing
- Contact information
  - Becky Thomas, [bthomas@i2kconnect.com](mailto:bthomas@i2kconnect.com)



The screenshot displays the i2k Connect search interface. At the top, there is a search bar with the text 'Artificial intelligence' and a search icon. Below the search bar, the results are displayed on 'Page 1 of 167 results', sorted by 'Relevance'. The main search result is titled 'Artificial Intelligence and Data Analytics for Virtual Flow Metering' by Gryzlov, Anton (Aramco Innovations), Mironova, Liliya (Aramco Innovations), Safonov, Sergey (Aramco Innovations), and Arsalan, Muhammad (Saudi Aramco), published by OnePetro in November 2021. The abstract discusses modern challenges in reservoir management and the use of multiphase flow meters. The interface includes filters on the left for 'Current Filters' (Saudi Aramco, Production and Well Operations, Artificial Intelligence (Technology)), 'SPE Disciplines' (Geophysics, Geology), and 'Concept Tag'.



The screenshot shows a question-and-answer interface. The question is 'What is cumulative production?'. Below the question, there is a prompt: 'For each question, order the answers via the drag and drop arrow on the left side of each answer, based on the answer's quality, with the top of the list being the answer with the highest quality.' Two answers are shown, ranked from 1 to 2. Answer 1 is 'gross production'. Answer 2 is a detailed definition: 'Cumulative production refers to the total amount of hydrocarbon that has been recovered up to a certain date. It includes all the recoverable resources that have been estimated and the production that has been measured in terms of sales product specifications. Raw production quantities, including non-hydrocarbons, are also measured to support engineering analyses requiring reservoir voidage calculations. Multiple development projects may be applied to each known or unknown accumulation, and each project will be forecast to recover an estimated portion of the initially-in-place quantities.'

# Industry-Specific Use Cases



CVX utilizes i2k Connect for analyzing upstream unstructured content to auto-classify, tag, and enrich metadata for ingestion into SharePoint. Originally, Noble Energy use case AI analysis of entire unstructured content upstream of content or asset classification and D&A purposes.



SLB has OEM'd i2k Connect to provide AI automated classifications against 18 taxonomies to provide insights to their clients within the DELFI cognitive E&P Platform.



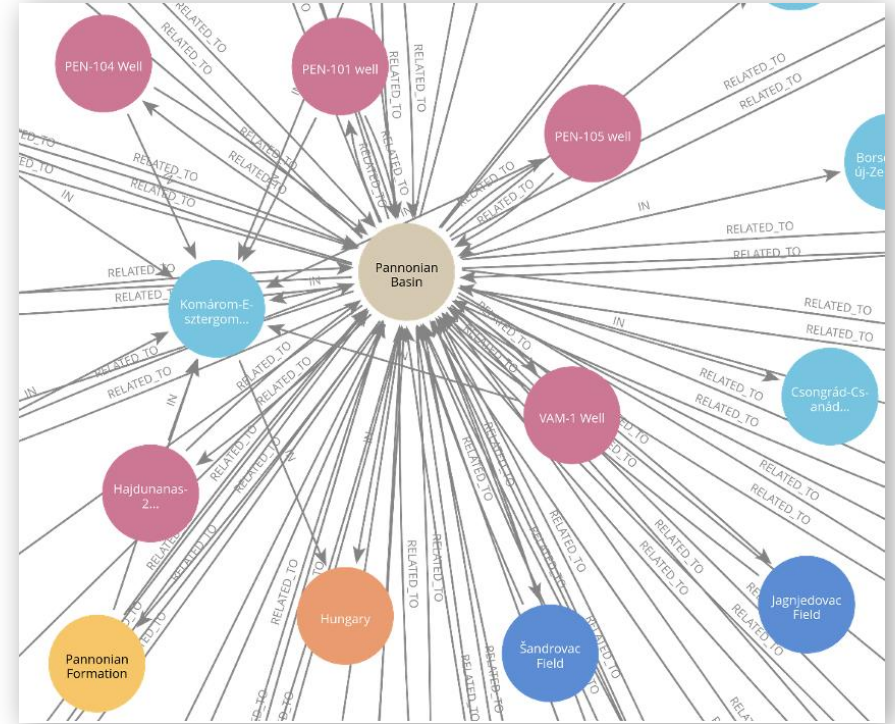
Society of Petroleum Engineers. Conducted LLM research associated with industry domain evaluating available LLM options to determine viability, accuracy, training, fine-tuning. Membership tested. i2k Connect powers the Research Portal, which enables research across all SPE, SEG, and 19 other societies' digital assets (papers, videos, etc.) at <https://search.spe.org> or <http://search.seg.org>



Woodside Energy uses enhanced search and findability of asset-related content over 46M files (437TB) in the upstream organization. Enables users to interrogate the corpus to quickly locate the right information to get their jobs done.

# i2k LLM R&D

- Question answering grounded in facts
  - Sourcing answers directly from documents
  - Grounding with knowledge graphs
- Published papers, continuing research
  - QA experiments with SPE volunteers
  - Finetuning domain-specific LLMs
  - **LLM Redact** – LLMs on docs w/ redactions
  - **LLM Trace** – Framework to keep track of LLM-generated data within a system
  - **QuoteLLM** – Extracting training text verbatim
  - Building LLM agents



## Answering Natural Language Questions with OpenAI's GPT in the Petroleum Industry

J. Eckroth<sup>1</sup>, M. Gipson<sup>1</sup>, J. Boden<sup>2</sup>, L. Hough<sup>1</sup>, J. Elliott<sup>1</sup>, J. Quintana<sup>2</sup>  
<sup>1</sup>i2k Connect Inc, Houston, Texas  
<sup>2</sup>Society of Petroleum Engineers, Richardson, Texas

### Abstract

This work documents two experiments that make use of OpenAI's ChatGPT and GPT-4 for question answering in the petroleum industry. First, we describe PetroQA, a prototype tool that can answer natural language questions. It uses PetroWiki content to inform ChatGPT about facts specific to this industry. We are able to convince ChatGPT to avoid hallucinations and cite its sources. We asked nearly 200 SPE members to volunteer to test PetroQA and discuss results from that test. Second, we are developing and testing a tool, known as GraphQA, that allows users to ask questions and receive answers from a large graph knowledge base consisting of facts and relations between concepts such as wells, fields, basins, formations, geography, geologic age, rock type, operators, and more. A knowledge base like this is difficult for users to explore, so we use GPT-4 to automatically generate accurate graph queries from their natural language questions. We explore several novel techniques for prompting GPT-4 to produce the right queries and have developed an advanced caching mechanism to reduce interactions with the cloud model, thus reducing time to answer and cost.

# LLM Risk



## LLM01: Prompt Injections

Prompt injection vulnerabilities in LLMs involve crafty inputs leading to undetected manipulations. The impact ranges from data exposure to unauthorized actions, serving attacker's goals.

## LLM02: Insecure Output Handling

These occur when plugins or apps accept LLM output without scrutiny, potentially leading to XSS, CSRF, SSRF, privilege escalation, remote code execution, and can enable agent hijacking attacks.

## LLM03: Training Data Poisoning

LLMs learn from diverse text but risk training data poisoning, leading to user misinformation. Overreliance on AI is a concern. Key data sources include Common Crawl, WebText, OpenWebText, and books.

## LLM04: Denial of Service

An attacker interacts with an LLM in a way that is particularly resource-consuming, causing quality of service to degrade for them and other users, or for high resource costs to be incurred.

## LLM05: Supply Chain

LLM supply chains risk integrity due to vulnerabilities leading to biases, security breaches, or system failures. Issues arise from pre-trained models, crowdsourced data, and plugin extensions.

## LLM06: Permission Issues

Lack of authorization tracking between plugins can enable indirect prompt injection or malicious plugin usage, leading to privilege escalation, confidentiality loss, and potential remote code execution.

## LLM07: Data Leakage

Data leakage in LLMs can expose sensitive information, proprietary details, leading to privacy and security breaches. Proper data sanitization, and clear policies are crucial for prevention.

## LLM08: Excessive Agency

When LLMs interface with other systems, unnecessary agency may lead to undesirable operations and data breaches. Like web-apps, LLMs should not self-police; they should be embedded in APIs.

## LLM09: Overreliance

Overreliance on LLMs can lead to misinformation and inappropriate content due to "hallucinations." Without proper oversight, this can result in legal issues and reputational damage.

## LLM10: Insecure Plugins

Plugins connecting LLMs to external resources can be exploited if they accept free-form text inputs, or if they accept malicious requests that could lead to undesirable actions or remote code execution.



## Ethical and social risks of harm from Language Models

Laura Weidinger<sup>1</sup>, John Mellor<sup>1</sup>, Maribeth Rauh<sup>1</sup>, Conor Griffin<sup>1</sup>, Jonathan Uesato<sup>1</sup>, Po-Sen Huang<sup>1</sup>, Myra Cheng<sup>1,2</sup>, Mia Glaese<sup>1</sup>, Borja Balle<sup>1</sup>, Atoosa Kasirzadeh<sup>1,3</sup>, Zac Kenton<sup>1</sup>, Sasha Brown<sup>1</sup>, Will Hawkins<sup>1</sup>, Courtney Biles<sup>1</sup>, Abeba Birhane<sup>1,4</sup>, Julia Haas<sup>1</sup>, Laura Rimell<sup>1</sup>, Lisa Anne Hendricks<sup>1</sup>, W. Isaac<sup>1</sup>, Sean Legassick<sup>1</sup>, Geoffrey Irving<sup>1</sup> and Iason Gabriel<sup>1</sup>

<sup>1</sup>Google, <sup>2</sup>California Institute of Technology, <sup>3</sup>University of Toronto, <sup>4</sup>University College Dublin

### Abstract

This paper aims to help structure the risk landscape associated with large-scale Language Models (LLMs). In order to foster advances in responsible innovation, an in-depth understanding of the potential risks posed by LLMs is needed. A wide range of established and anticipated risks are analysed in detail, drawing on interdisciplinary literature from computer science, linguistics, and social sciences.

The paper outlines six specific risk areas: I. Discrimination, Exclusion and Toxicity, II. Information Hazards, III. Information Harms, IV. Malicious Uses, V. Human-Computer Interaction Harms, VI. Automation, Access, and Environmental Harms.

## Gender bias and stereotypes in Large Language Models

Hadas Kotek  
Apple & MIT  
Cupertino, CA, USA  
hadas@apple.com

Rikker Dockum  
Swarthmore College  
Swarthmore, PA, USA  
rdockum1@swarthmore.edu

David Q. Sun  
Apple  
Cupertino, CA, USA  
dqs@apple.com

## LLMs as Factual Reasoners: Insights from Existing Benchmarks and Beyond

Philippe Laban Wojciech Kryściński Divyansh Agarwal Alexander R. Fabbri  
Caiming Xiong Shafiq Joty Chien-Sheng Wu  
Salesforce AI

{plaban, wojciech.kryscinski, dagarwal, afabbri, cxiong, sjoty, wu.jason}@salesforce.com

### Abstract

The recent appearance of LLMs in practical applications has highlighted the importance of ensuring that these models are factually consistent. Factual inconsistencies in model outputs can lead to misinformation and other harmful consequences. When testing LLMs on classification benchmarks, we observe that existing evaluation methods are often inconsistent, affecting the evaluation precision. In this paper, we propose a new protocol for consistency detection benchmark creation.

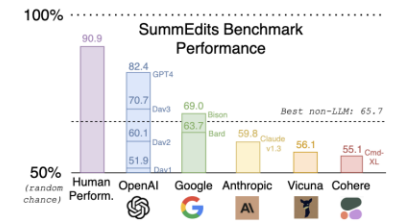
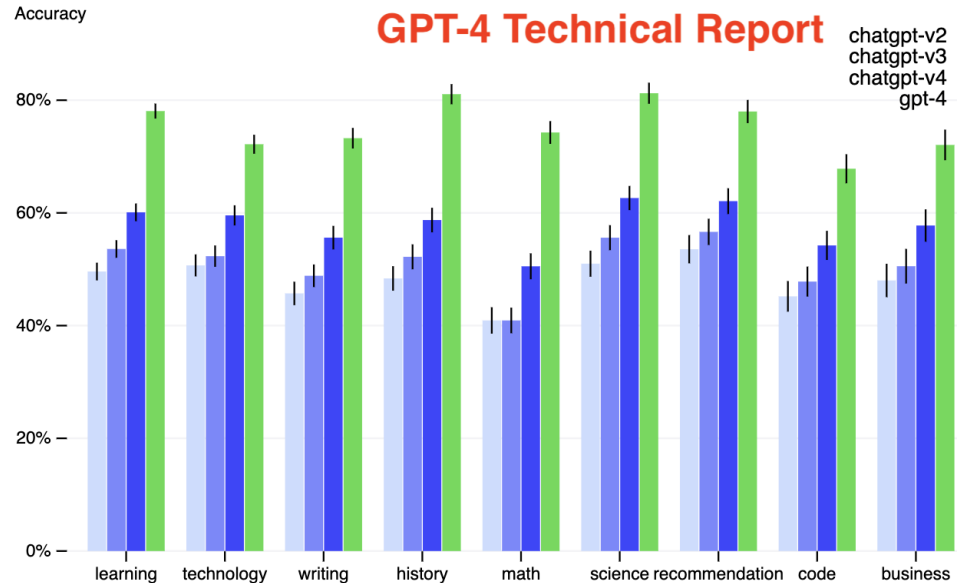


Figure 1: SUMMEDITS is a benchmark to evaluate the factual reasoning abilities of LLMs, measuring if models detect factual inconsistencies when they occur in summaries. Capable detection models can help build more reliable NLG systems.

## Internal factual eval by category





# Mitigating LLM Risk

**Answering Natural Language Questions with OpenAI's GPT in the Petroleum Industry**

J. Eckroth<sup>1</sup>, M. Gipson<sup>1</sup>, J. Boden<sup>2</sup>, L. Hough<sup>1</sup>, J. Elliott<sup>1</sup>, J. Quintana<sup>2</sup>  
<sup>1</sup>i2k Connect Inc, Houston, Texas  
<sup>2</sup>Society of Petroleum Engineers, Richardson, Texas

**LLM01: Prompt Injections**  
Prompt Injection Vulnerabilities in LLMs involve crafty inputs leading to undetected manipulations. The impact ranges from data exposure to unauthorized actions, serving attacker's goals.

**LLM02: Insecure Output Handling**  
These occur when plugins or apps accept LLM output without scrutiny, potentially leading to XSS, CSRF, SSRF, privilege escalation, remote code execution, and can enable agent hijacking attacks.

**LLM03: Training Data Poisoning**  
LLMs learn from diverse text but risk training data poisoning, leading to user misinformation. Overreliance on AI is a concern. Key data sources include Common Crawl, WebText, OpenWebText, and books.

**LLM04: Denial of Service**  
An attacker interacts with an LLM in a way that is particularly resource-consuming, causing quality of service to degrade for them and other users, or for high resource costs to be incurred.

**LLM05: Supply Chain**  
LLM supply chains risk integrity due to vulnerabilities leading to biases, security breaches, or system failures. Issues arise from pre-trained models, crowdsourced data, and plugin extensions.

QuoteLLM Research Project

LLM Redact Research Project

Finetuning Domain-Specific LLMs

LLM Trace Research Project

**LLM06: Permission Issues**  
Lack of authorization tracking between plugins can enable indirect prompt injection or malicious plugin usage, leading to privilege escalation, confidentiality loss, and potential remote code execution.

**LLM07: Data Leakage**  
Data leakage in LLMs can expose sensitive information or proprietary details, leading to privacy and security breaches. Proper data sanitization, and clear terms of use are crucial for prevention.

**LLM08: Excessive Agency**  
When LLMs interface with other systems, unrestricted agency may lead to undesirable operations and actions. Like web-apps, LLMs should not self-police; controls must be embedded in APIs.

**LLM09: Overreliance**  
Overreliance on LLMs can lead to misinformation or inappropriate content due to "hallucinations." Without proper oversight, this can result in legal issues and reputational damage.

**LLM10: Insecure Plugins**  
Plugins connecting LLMs to external resources can be exploited if they accept free-form text inputs, enabling malicious requests that could lead to undesired behaviors or remote code execution.