

2023 Hackathon “Abstracters”

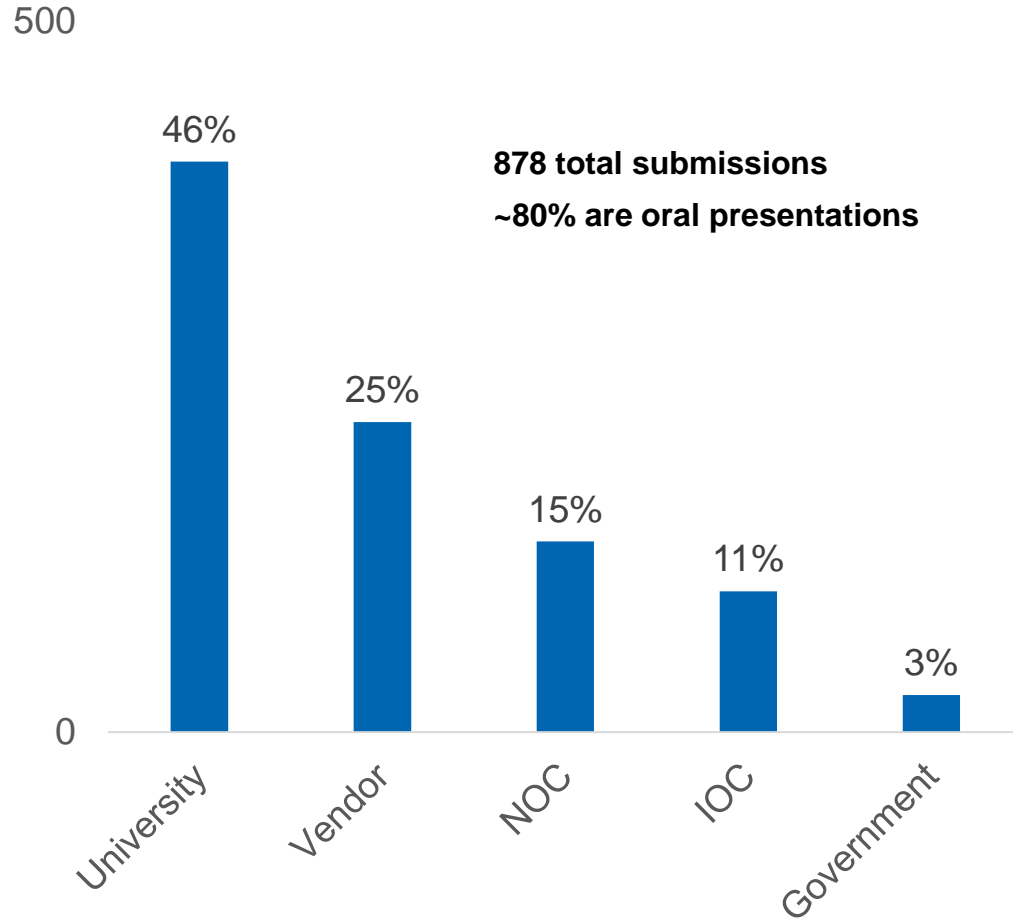
Steve Braun

Sergey Fomel

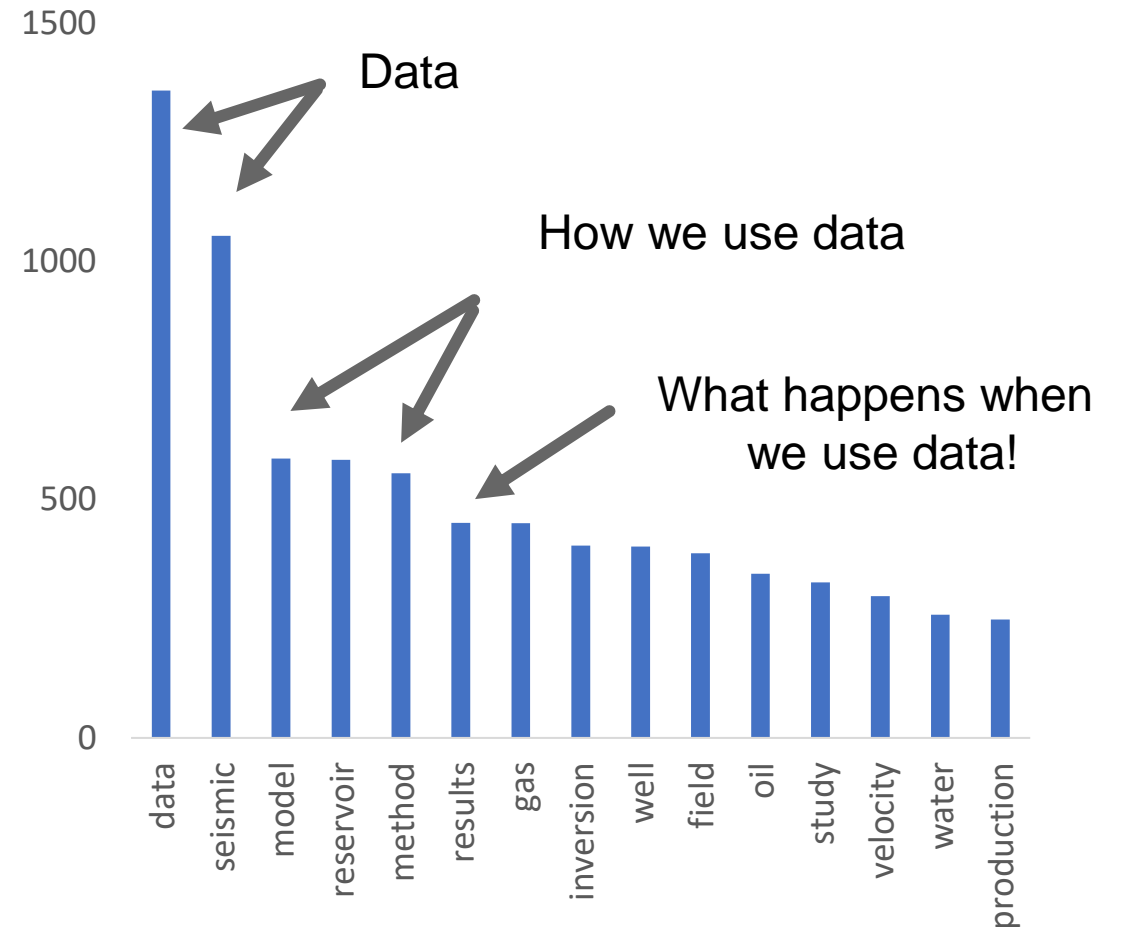


Preliminary data analysis: Who is presenting and about what?

Abstract count by organization type and presentation type*



Top 15 words mentioned in abstract summaries



*Organization typed by first author on best effort basis

If ChatGPT is so smart maybe it should do the hackathon too?



Give me 10 ideas for projects in a hackathon on human language processing for geoscience applications



Sure! Here are 10 project ideas for a hackathon focused on human language processing for geoscience applications:

1. Geoscience Text Summarization: Develop a text summarization tool specifically designed for geoscience literature, which can extract key information from research papers, reports, or other textual resources.
2. Geoscience Document Classification: Build a machine learning model that can accurately classify geoscience documents into different categories, such as seismic data analysis, geothermal exploration, mineralogy, or climate modeling.
3. Geoscience Question Answering System: Create a question answering system that can provide accurate answers to geoscience-related questions by leveraging natural language processing techniques and geoscience knowledge bases.

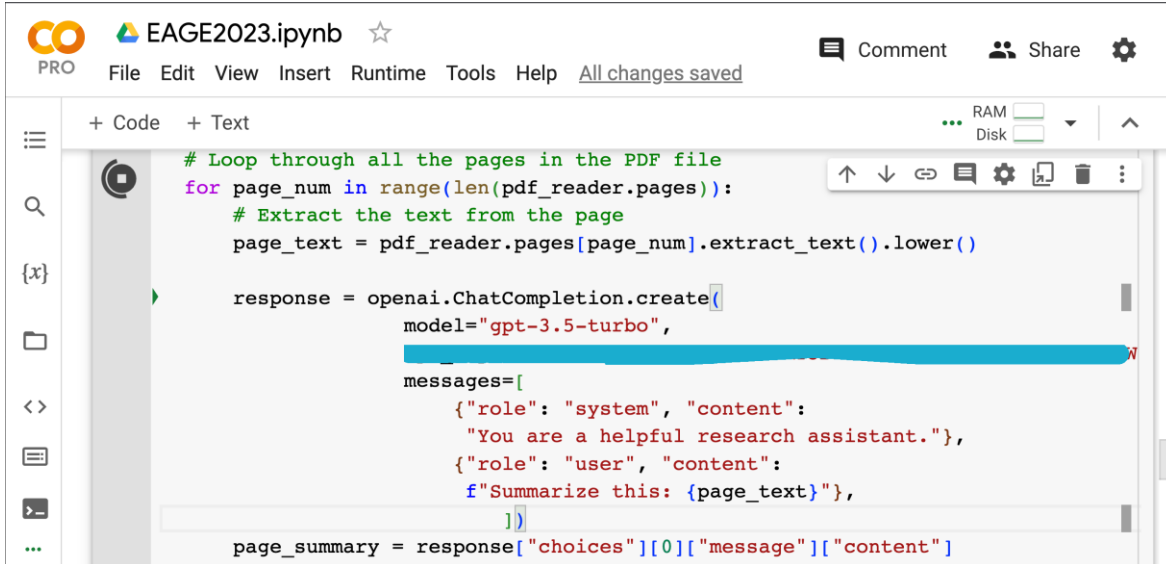
NRCAN, NPD have both released tools doing this in the last year



We decided to try these, with differing levels of success

Followed templates from openAI for working with ChatGPT

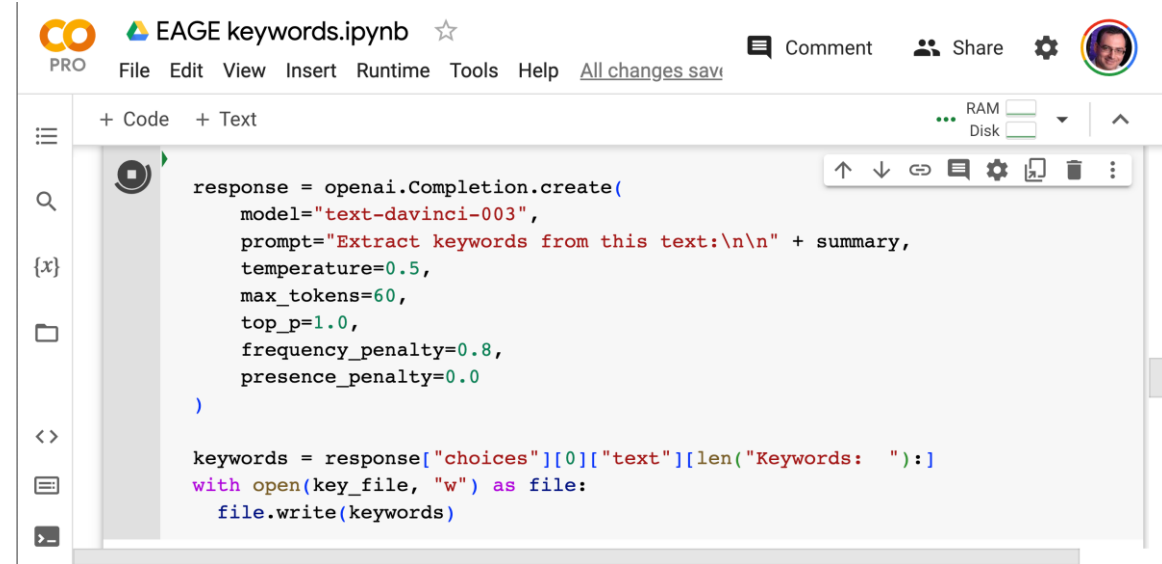
Automating GPT3.5 summaries



```
# Loop through all the pages in the PDF file
for page_num in range(len(pdf_reader.pages)):
    # Extract the text from the page
    page_text = pdf_reader.pages[page_num].extract_text().lower()

    response = openai.ChatCompletion.create(
        model="gpt-3.5-turbo",
        messages=[
            {"role": "system", "content":
              "You are a helpful research assistant."},
            {"role": "user", "content":
              f"Summarize this: {page_text}"},
        ]
    )
    page_summary = response["choices"][0]["message"]["content"]
```

Automating GPT3.5 key words



```
response = openai.Completion.create(
    model="text-davinci-003",
    prompt="Extract keywords from this text:\n\n" + summary,
    temperature=0.5,
    max_tokens=60,
    top_p=1.0,
    frequency_penalty=0.8,
    presence_penalty=0.0
)

keywords = response["choices"][0]["text"][len("Keywords: "):]
with open(key_file, "w") as file:
    file.write(keywords)
```

The code openAI gave us usually didn't work out of the box. However, light troubleshooting usually fixed the issues

Creating a text summary

What went into GPT3.5 summarization

L. Zhuo, W. Fei, L. Wei, J. Bergeron, G. Hennenfent



Evaluating 4D DAS VSP Acquisition at the Big Foot Field Gulf of Mexico

Introduction

The Big Foot field is located in the Walker Ridge area of the U.S. Gulf of Mexico, approximately 225 miles (360 km) south of New Orleans, La., in a water depth of approximately 5,200 feet (1,584 m) (Figure 1a). The reservoir is a three-way closure sealed by the salt with the Miocene sand as the main reservoir (Figure 1b). The Big Foot field production started in 2018. Field development plan includes a multi-well drilling program. 4D seismic imaging is considered for enhanced field development.

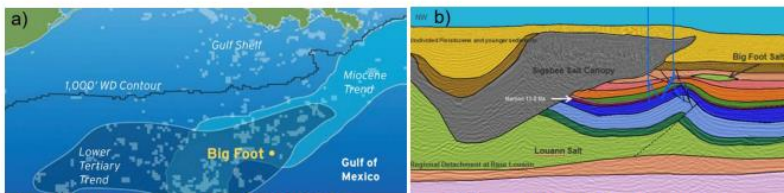


Figure 1 The Big Foot field, U.S. Gulf of Mexico. a) Field location; b) Field geological section.

Vertical Seismic Profiling (VSP) acquisition with Distributed Acoustic Sensing (DAS) technology is emerging as an attractive alternative for deep water fields in addressing challenges for 3D seismic imaging and 4D reservoir monitoring due to the lower cost, less interruption to production and the recent industry success (Zhan and Nahm, 2022, Kiyashchenko et al., 2019, Mateeva et al., 2017).

4D DAS VSP is considered as a potential solution for reservoir monitoring at the Big Foot field to further enhance oil recovery. There are critical questions to be answered for making the decision, including: 1) What is the chance of success for 4D DAS VSP at this field? 2) How frequently should the monitor survey be conducted? 3) What is the impact of the field noise to 4D signal? 4) Can we reduce the acquisition time and cost?

We performed a FD (Finite Difference) modeling study in 2022 to assess the feasibility of 4D DAS VSP acquisition at the Big Foot field to improve decisions for field development.

Method

FD modeling is a proven technology that has routinely demonstrated its value for Front-End-Loading (FEL) phase of evaluating seismic acquisition design to reduce cost and improve decision quality

What GPT3.5 came up with as a summary

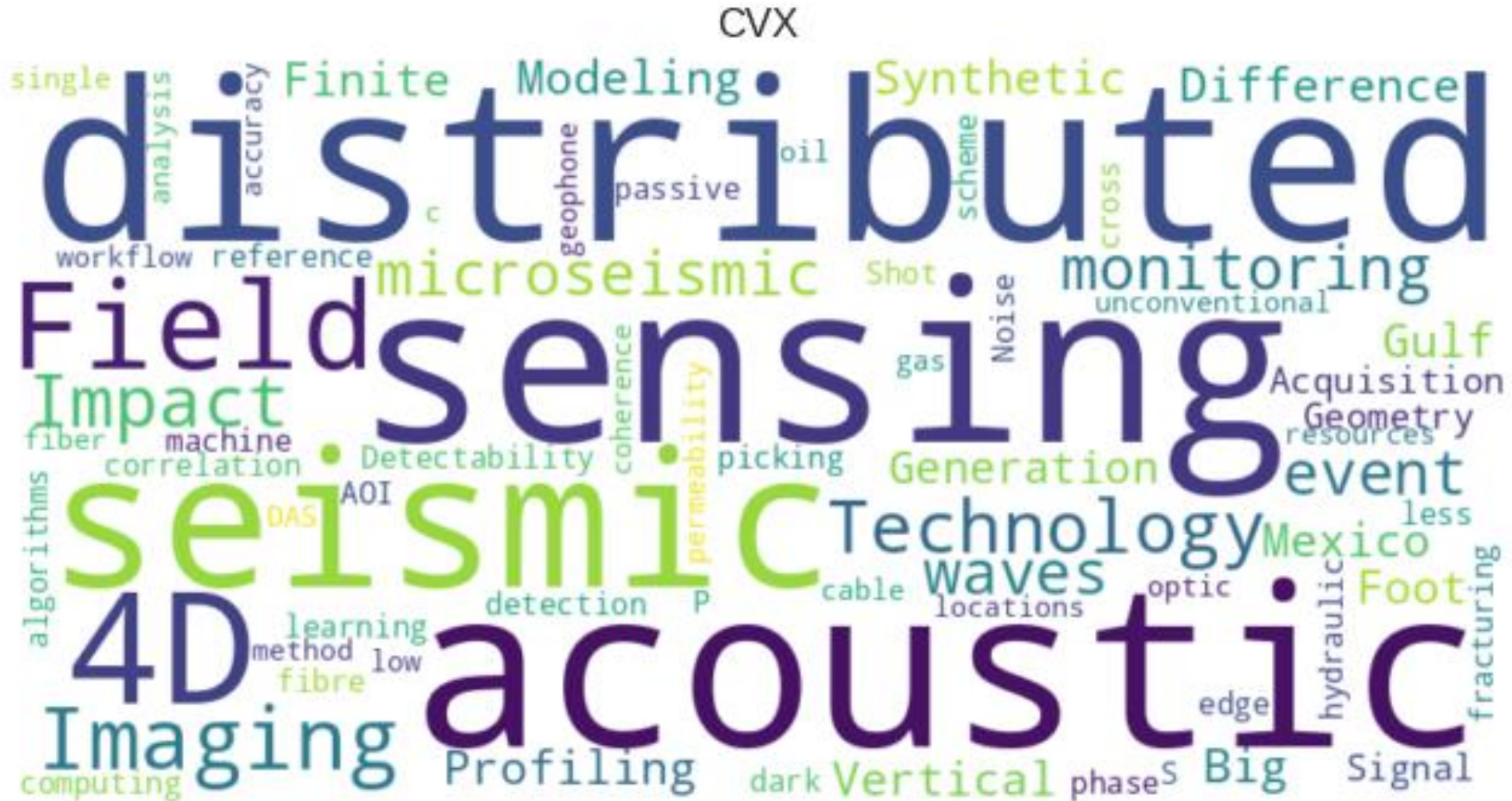
We performed a FD (Finite Difference) modeling study to evaluate the feasibility of DAS VSP acquisition for 4D reservoir monitoring at the Big Foot field, U.S. Gulf of Mexico.

The study demonstrates that: 1) DAS VSP acquisition has high chance of success for detecting 4D signal and effectively monitoring reservoir changes at this field. 2) Shot AOI at one well can be reduced without jeopardizing the 4D objective, which could save time and cost for field acquisition. However, the larger shot AOI is still needed for the other two wells. 3) Field shut-in noise and natural flow noise substantially reduced the S/N of the 3D gather. However, noise sensitivity test shows the shut-in noise is negligible, and natural flow only adds 1-2% of noise to 4D signal.

The learnings from this study significantly contributed to the decisions on planning DAS VSP acquisitions at the Big Foot field.

How did it do? Could be a good way to determine which presentations to go to

Using key words to understand what IOCs are sharing



We also tried with document summaries, but non-industry terms made the clouds less useful
We created these word clouds for all companies

Is the next step to use these kinds of systems to *predict* what future research opportunities could be?

TOT



Securing a sustainable future together

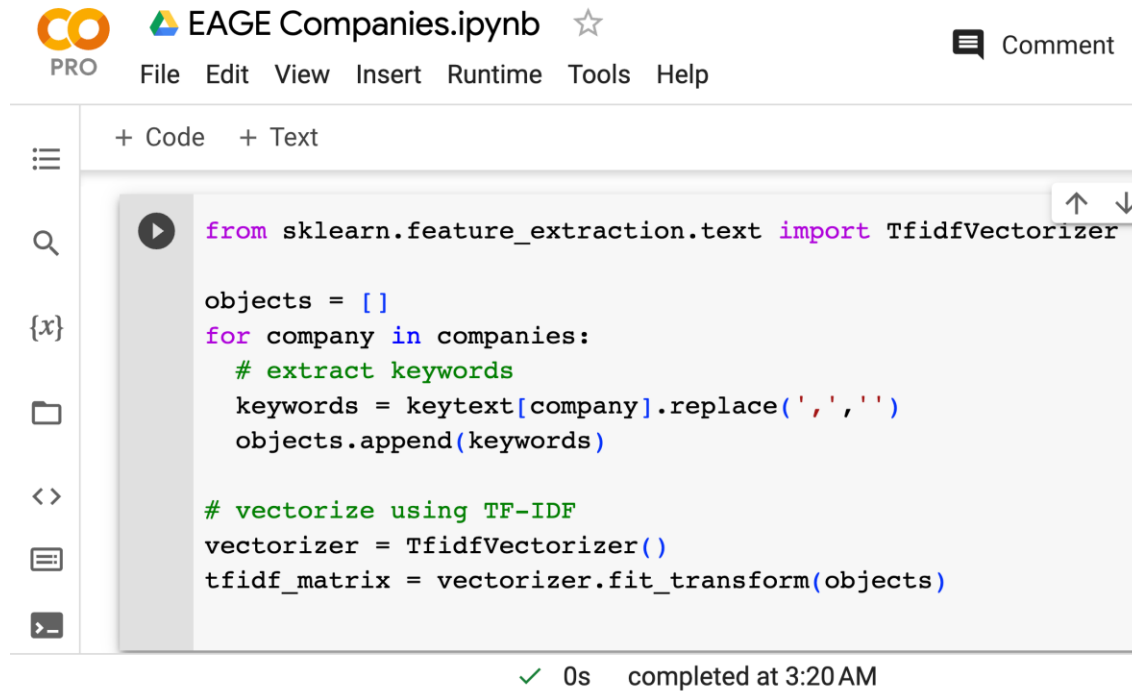
Energy Transition



We also did this for all other session themes: geology, geophysics, petroleum engineering, mining, etc

Which companies are presenting the most similar content?

Cluster analysis enables visual interrogation of company abstracts

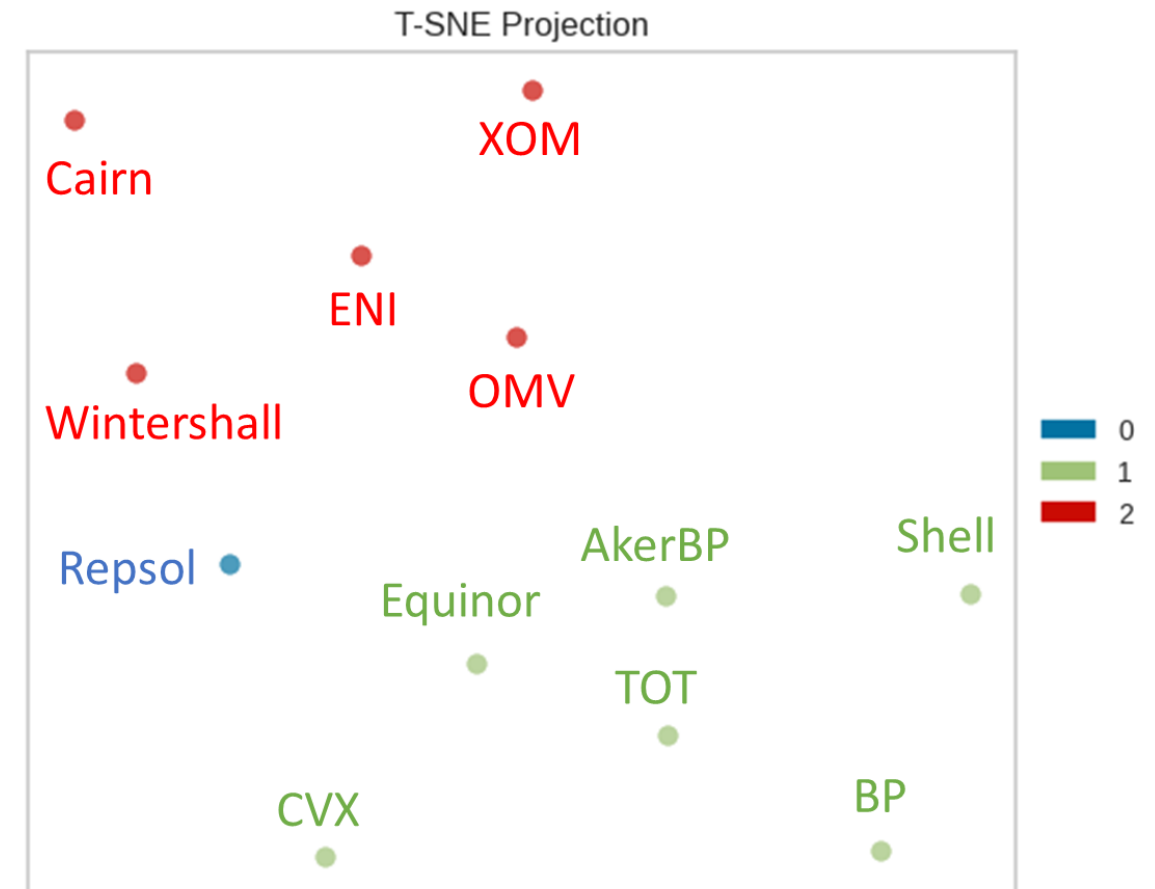


```
from sklearn.feature_extraction.text import TfidfVectorizer

objects = []
for company in companies:
    # extract keywords
    keywords = keytext[company].replace(',', ',')
    objects.append(keywords)

# vectorize using TF-IDF
vectorizer = TfidfVectorizer()
tfidf_matrix = vectorizer.fit_transform(objects)
```

0s completed at 3:20 AM

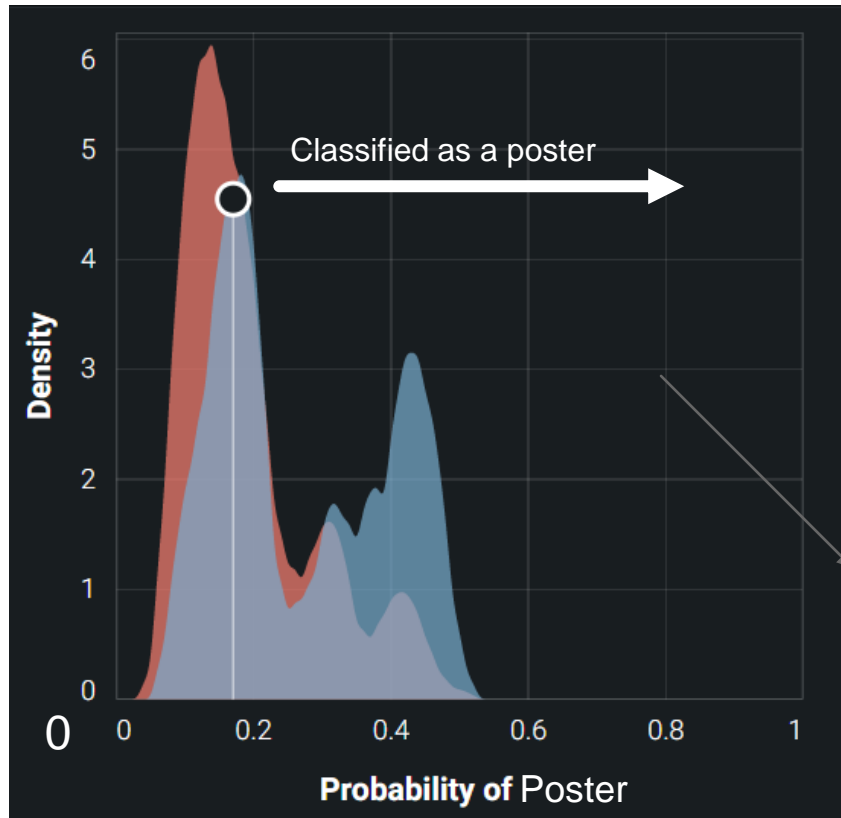


Could become the basis of a recommendation engine or to identify folks publishing anomalous work

Building a classifier to understand who gets a Poster Presentation

Results from holdout sample using a K-nearest neighbors classifier of document summaries

Look at the overlap between classes



Confusion matrix shows low predictive power of model

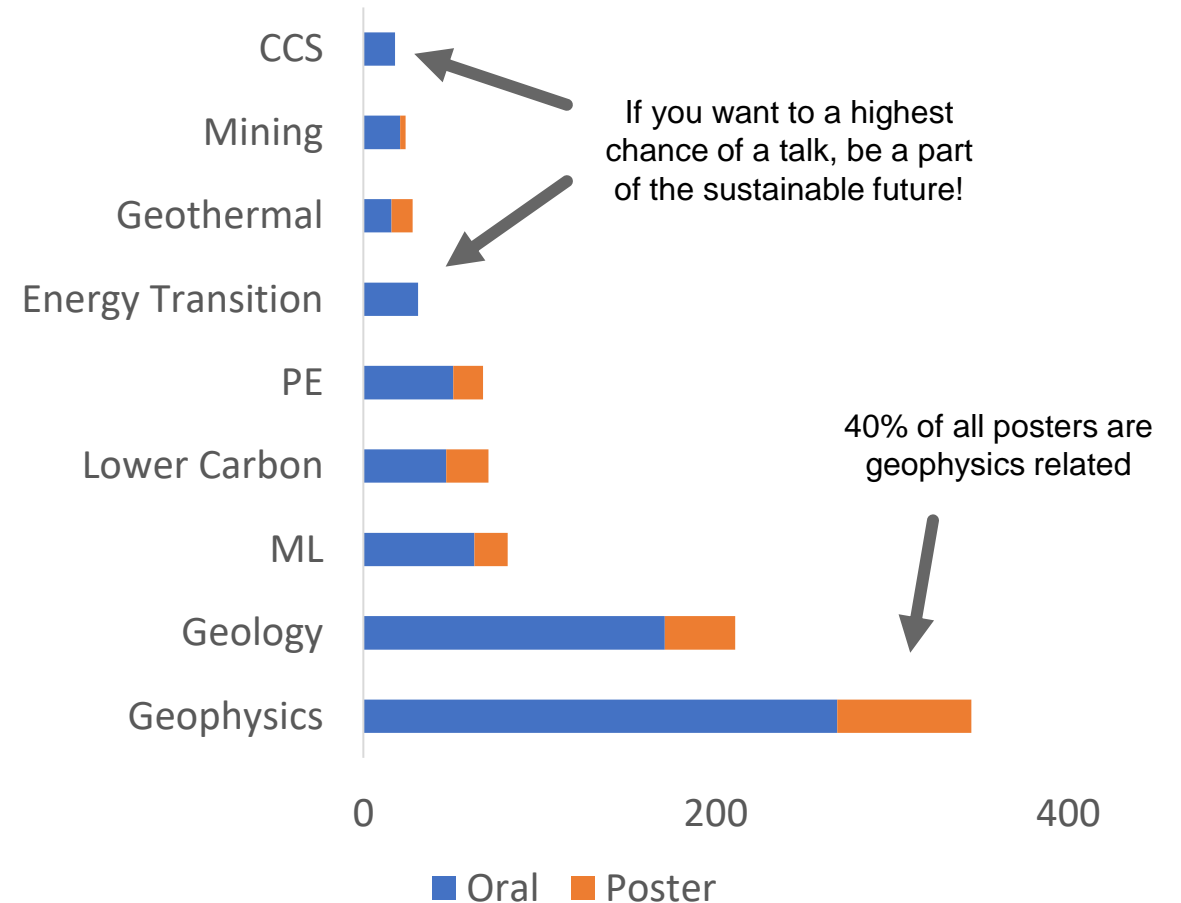
		Predicted		
		Oral	Poster	
Actual	Oral	74 (TN)	64 (FP)	138
	Poster	12 (FN)	26 (TP)	38
		86	90	176

While this is clearly not a production worthy model, there are still insights to be gained

Building a classifier to understand who gets a Poster Presentation

- Important to remember **predictions don't tell you reasons**
- For example, we examined the papers with the highest probability of being posters, many of them mentioned geophysics topics
- This could indicate, recalling back to our first data analysis that there were just more geophysics abstracts submitted than spaces for talks

Abstract count by session theme



*Session themes typed on best effort basis

Conclusions



Large language models like ChatGPT are effective in working with human texts

Extracting summaries
Extracting keywords



Keywords enable data analysis on EAGE expanded abstracts

Word clouds
Clustering



Ideas for additional projects

Classification using supervised machine learning

- Predict oral vs poster
- Predict organization type

