# Hybrid Summarization of Chemistry Patents

Darby Brown
W266: Natural Language Processing
Dec 2023
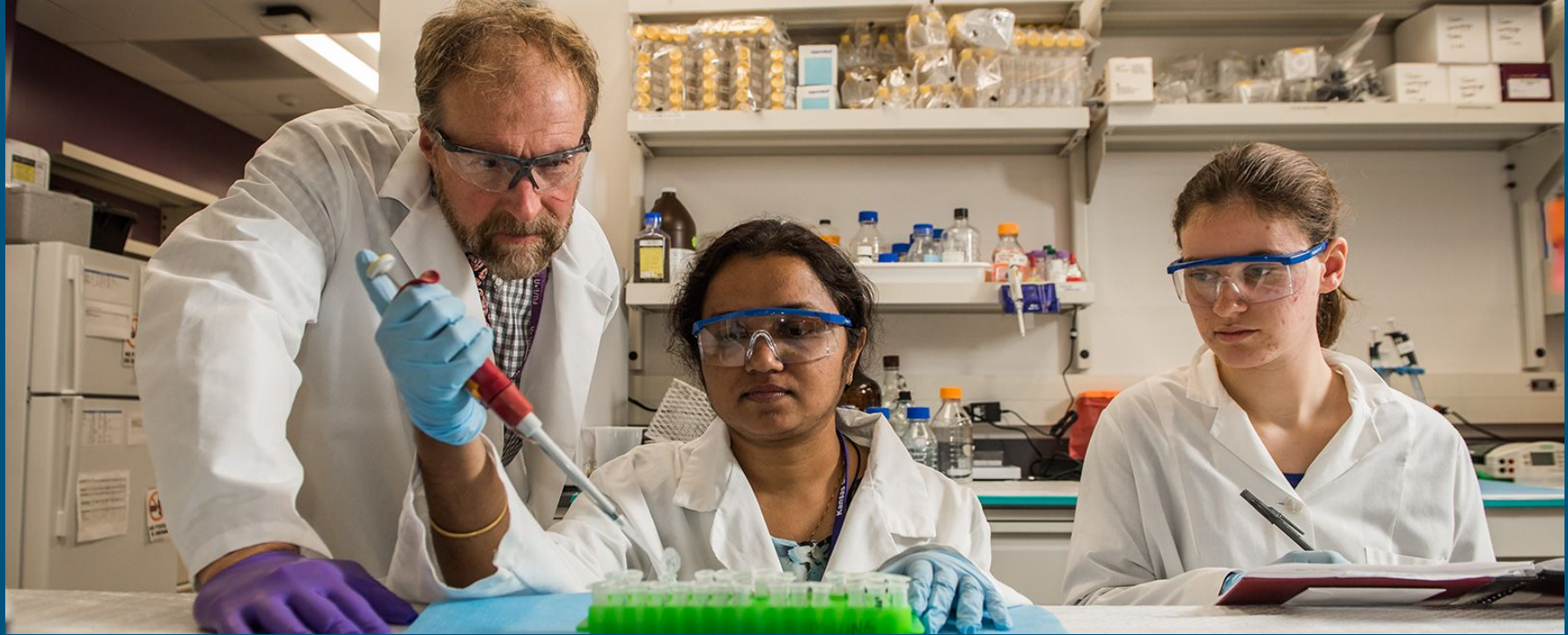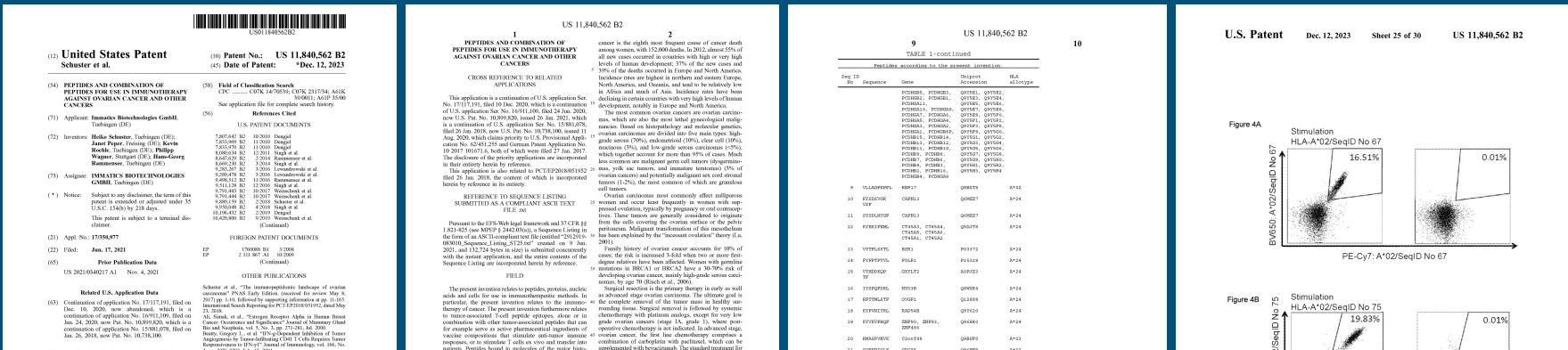
# Why patent summarization?

## 71%

of new technology published in patent literature
is not published elsewhere.

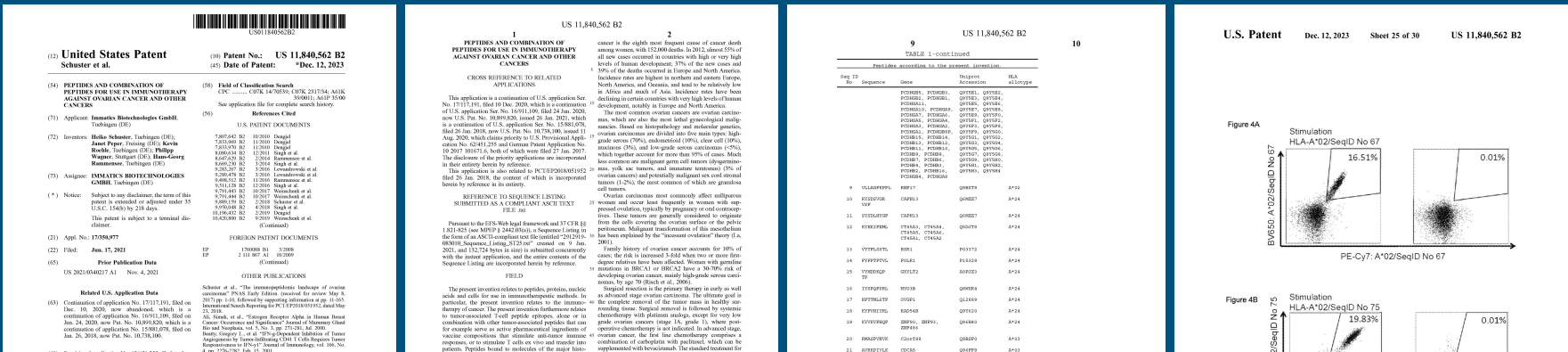**Reduce technology landscaping time to get scientists back into the lab, faster.**

# Working with Patents: Key Challenges



- Average document length: 24000 words
- Technical and legal jargon, often only interpretable by SMEs and attorneys
- Images, tables, and charts are frequently referred to in the text

# Working with Patents: Key Challenges



- **Average document length: 24000 words**
- **Technical and legal jargon,** often only interpretable by SMEs and attorneys
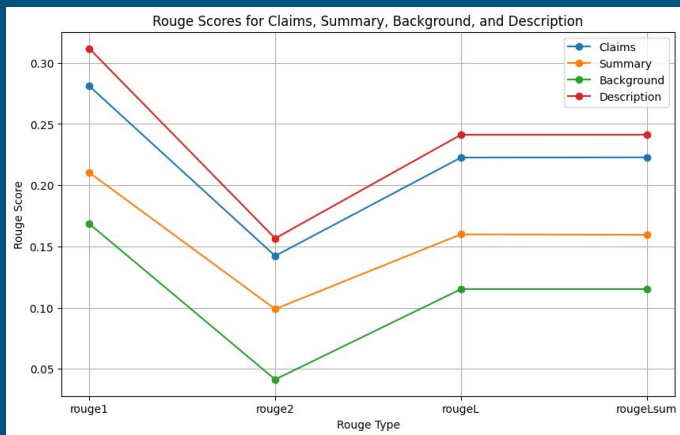- Images, tables, and charts are frequently referred to in the text

# Methods

**Intuition**: **Document length:** Extractive prior to abstractive summarization should mitigate issue.

**Technical jargon:** Domain-specific fine-tuning might improve the model's performance in that space.

| Dataset | <ul><li>Harvard USPTO Patent Dataset</li><li>Organic Chemistry class (2% of overall dataset)</li><li>Inputs: Claims, background, summary, description</li><li>Labels: Abstract</li></ul> |
|---|---|
| Basemodel | <ul><li>Pegasus-large</li></ul> |
| Improvements Tried | <ul><li>Select best sections for input ids</li><li>Hybrid summarization<ul><li>SumBasic → Pegasus-large →</li></ul></li><li>Fine-tuning on domain-specific data</li></ul> |
| Evaluation | <ul><li>ROUGE (esp. ROUGE-L)</li><li>Qualitative review</li></ul> |

# Results

- Best inputs: description and claims
- Hybrid summarization didn't improve quantitative results, but domain-specific fine-tuning did.



| Model | ROUGE | Dataset | |
| | | Description | Extracted |
|---|---|---|---|
| Pegasus -large | R1 | 31.16 | 22.51 |
| | R2 | 15.66 | 8.57 |
| | RL | 24.11 | 15.81 |
| | RLsum | 24.15 | 15.8 |
| Fine tuned Pegasus -large | R1 | **44.68** | 29.66 |
| | R2 | **27.94** | 13.34 |
| | RL | **37.52** | 23.42 |
| | RLsum | **40.42** | 25.09 |

(Baseline highlighted for Description column of Pegasus -large)

# Limitations & Future Research

Limitations:

- Memory constraints limited the dataset to <1000 observations
- Data leakage due to Pegasus' training set
  - Overlap between BigPatent and HUPD

Future work:

- TextRank for extractive summarization
- Leverage newer patents and/or proprietary summaries from industrial stakeholders to reduce data leakage
- Multimodal models to incorporate images

| Data Leakage Example | |
| --- | --- |
| Abstract | The present invention relates to an antibody construct comprising a first human binding domain specific for the extracellular part of the influenza envelope protein M2 (M2e) and a second domain specific for CD3.... |
| Summary of Description | <pad>FIELD OF THE INVENTION The present invention relates to an antibody construct comprising a first human binding domain specific for the extracellular part of the influenza envelope protein M2 and a second domain specific for CD3.... |

# Conclusion

- We established a proof-of-concept for domain-specific fine-tuning of patent summaries, which can serve as a starting point for R&D organizations or patent offices wishing to develop this capability.
- We explored the utility of hybrid summarization applied to patent documents and saw promising qualitative results.
- There is still much opportunity for optimization of patent summaries, both at the single- and multi- document level.