



Anne Castille Buisson



Hayden Ratliff

WIZER PFIZER



Connecting People To Knowledge Faster

○ Faculty Advisor: Yu Ma

○ Company Advisors: Kim Adler, Abby Freeman

Project Overview

Problem Statement

At Pfizer, it is always a priority to **reduce the time required to get drugs to market.**

A key opportunity to improve time to market is "Knowledge Transfer," where relevant R&D and acquired documents are transitioned to the manufacturing department.



9 months per molecule to classify documents



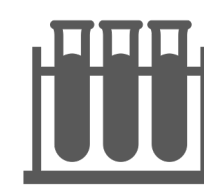
Not possible for Q&A with document contents

Data



33.6k

Documents



7

Molecules



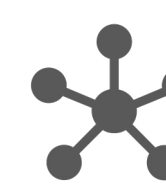
10

File Formats



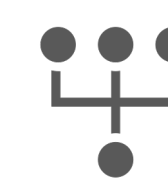
20

Q&A Pairs



633

Nodes



1.6k

Relationships

Product 1: Identifying Relevant Documents

Methods

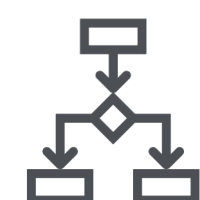
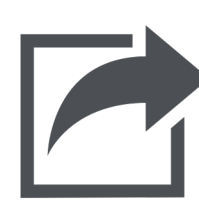
Documents

Text Extraction

Embedding

Classification

Relevance



- Confidential data
- Unbalanced labels
- Clinical, batches, investigations

- Computer Vision w/ Optical Character Recognition
- Limit first 2 pages

- Open-source LLM: BioMed RoBERTa
- Transfer Learning
- Semantic meaning

- HistXGBoost
- Weighted loss
- Custom metric: recall & AUC

- Classification probability
- Threshold = 0.5

Results

Performance

83.5%

Recall

0.845

AUC

71%

Filtering

Exceeds 80% target

Processing Time per Molecule (>1k Docs)

9 Months

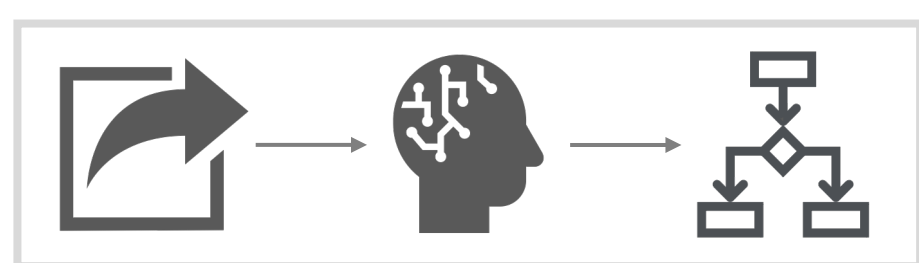
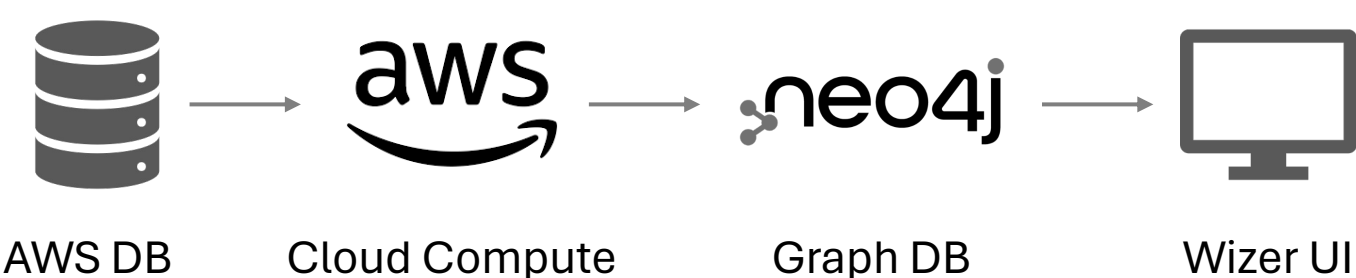
Manual Approach

10 Minutes

Automatic Approach

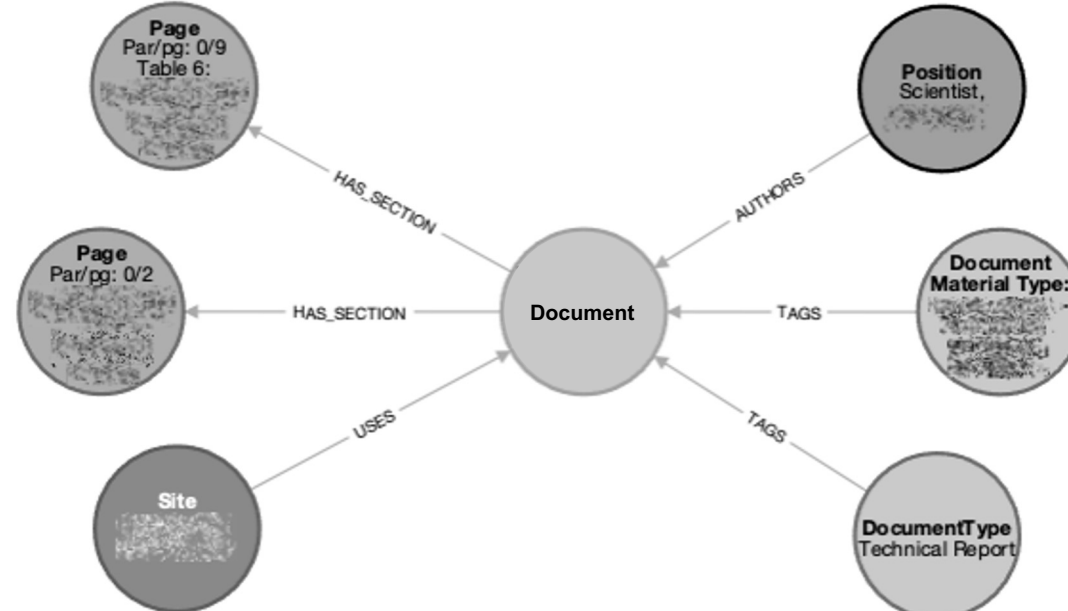
Implementation

Document Classification Workflow



Document Classification Pipeline

Knowledge Graph Database: Neo4j



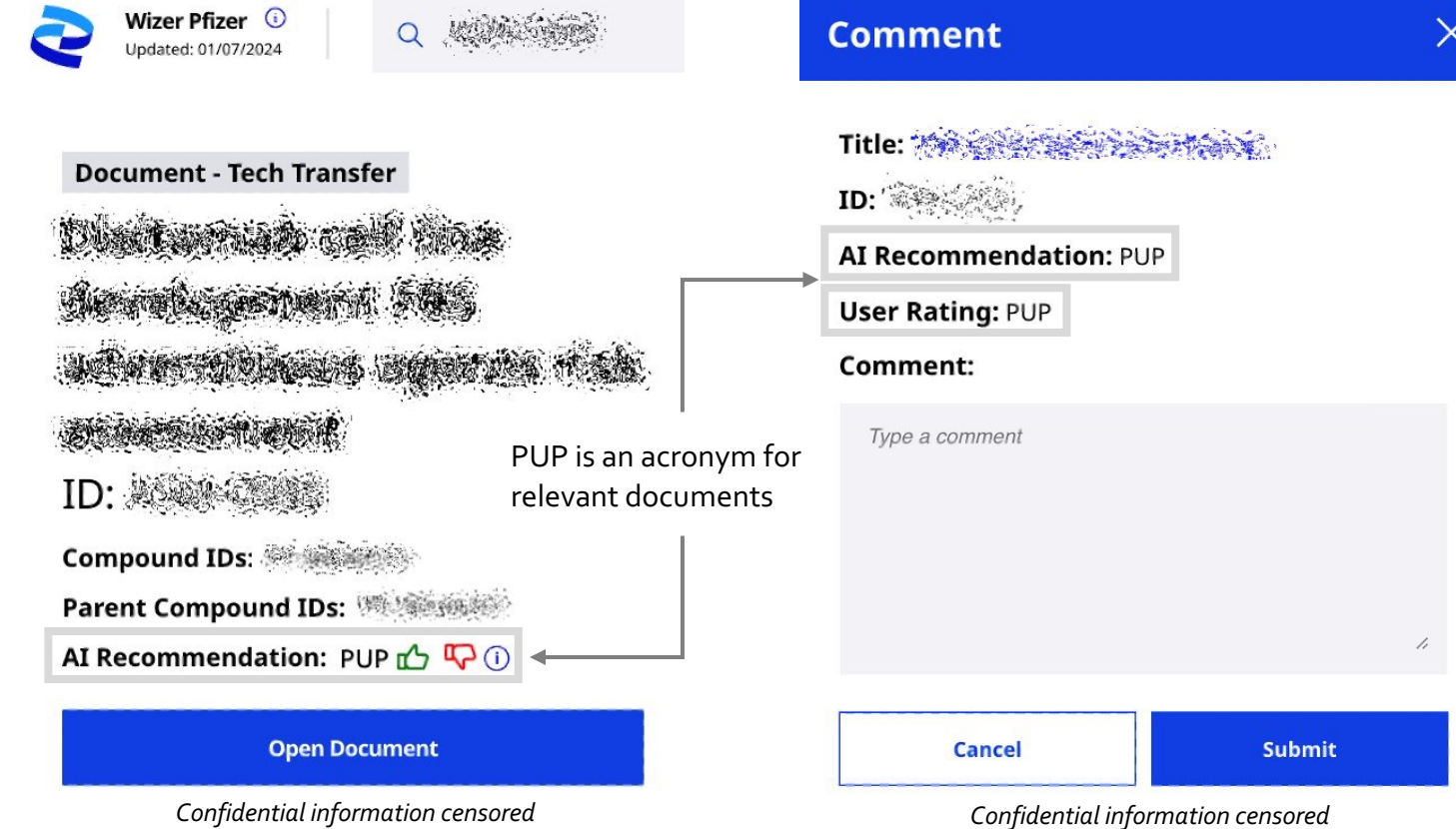
Node properties

author	[REDACTED]
compound_ids	[REDACTED]
document_id	[REDACTED]
document_title	[REDACTED]
pup_probability	0.640241417228176
pup_probability_updated_on	"2024-06-24T20:52:25.285000000Z"

Confidential information censored

Wizer Pfizer: User Interface

Human-In-The-Loop



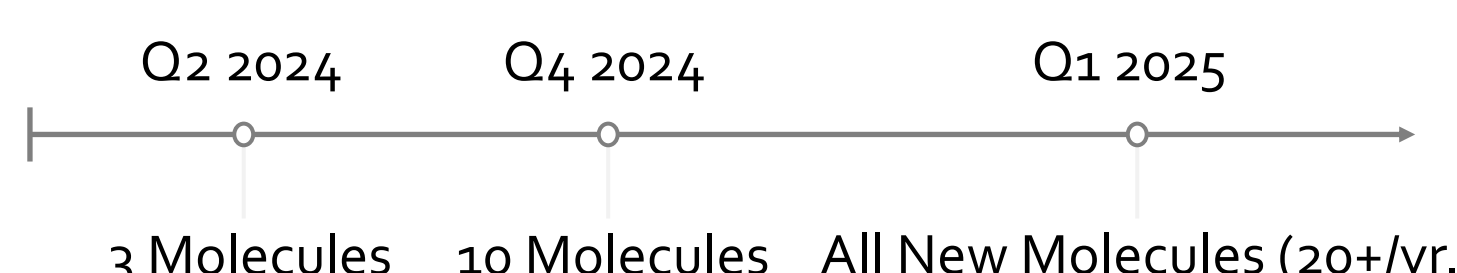
Confidential information censored

Confidential information censored

Performance and Specifications

- Processes each document in <1 second
- Workflow runs in parallel with 8 AWS EC2 nodes
- Results written to Knowledge Graph with ETL script

Product Adoption



Product 2: Answering User Questions

Methods

Question

Query Generation

Context Retrieval

Answer Extraction

Answer



- About data in Knowledge Graph
- Documents, sites, people, teams

- Translation LLM: TinyLLaMA
- Question to Cypher
- Few-shot learning w/ rules, schema

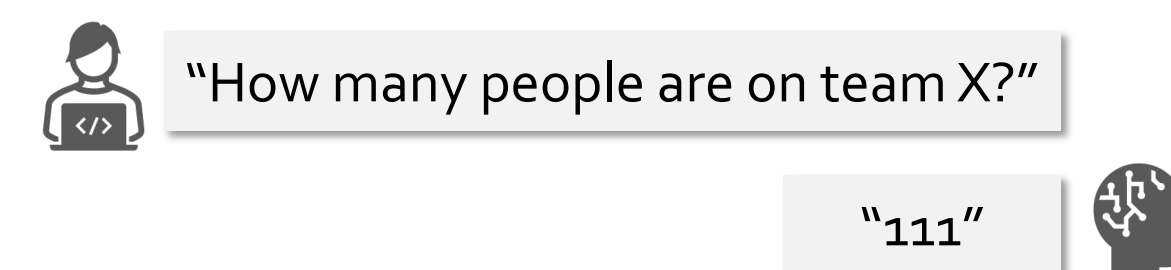
- Run query on Knowledge Graph
- Retrieve data required for answer

- Answer LLM: RoBERTa-Squadz
- Retrieval Augmented Generation (RAG)

- Informed answer if context found
- "I don't know" if no context found

Results

Sample Conversation



Performance

- Infers answer from database
- No hallucination during testing
- 15-20 seconds to answer on laptop
- Concept proven

Business Impact

Identifying Documents



1,200

Hours Saved



\$5.3 M

Savings Per Year



↓90%

Manual Work



1,000

Users

Current: Q2 2024

Projected: Q1 2025

Retrieving Information



14

Mins Saved / Search



150

Users



420+

Hours Saved / Year



300

Users

Current: Q2 2024

Projected: Q1 2025