# Разработка интеллктуальной системы анализа патентов химической отрасли для представления данных в структурированном виде

Студент 2-го курса: Кайда Анатолий Сергеевич
Научный руководитель: Глинский Андрей Владимирович

МФТИ

# Формулировка проблемы

▶ Количество ежегодно публикуемых статей и патентов в химии растет экспоненциально

▶ Процесс работы с в патентной и литературной информацией по-прежнему остается в значительной степени ручным

▶ Сложность навигации в большом объеме литературы приводит к тому, что важные научные открытия остаются незамеченными в течение длительного времени
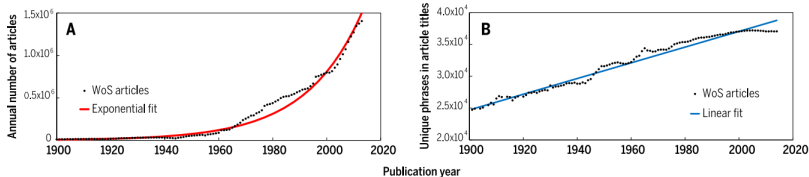


Рис. 1: (A) Годовой выпуск научных статей, индексированных в базе данных WoS. (B) Рост идей, охватываемых статьями, индексированными в WoS. Это было определено путем подсчета уникальных заглавных фраз (концепций) в фиксированном количестве статей.

# Формулировка проблемы

1. Ananikov V. Top 20 Influential AI-Based Technologies in Chemistry. Chemistry, 2024.

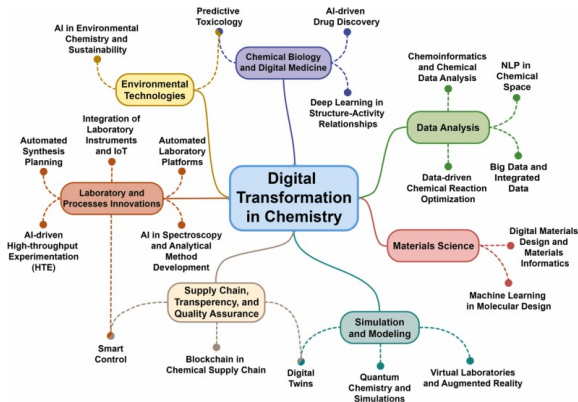В последние годы набирает обороты «цифровизация» химии.



Рис. 2: Связь технологий на основе ИИ с более широкими темами в зависимости от их применения.

# Формулировка проблемы



Рис. 3: Примеры входных данных

# Цели и задачи

Цель работы: Создать цифрового «ассистента» на основе большой языковой модели для извлечения структурированных данных из патентной документации

Задачи:

- Выбрать домен для проведения исследования и провести релевантный патентный поиск и создать БД документов для дальнейшего извлечения информации
- Провести обзор современных фреймворков и технологий для работы с БЯМ
- Создать агентов на основе БЯМ способных решать следующие задачи:
    - Сегментация и фильтрация текста
    - Классификация текста и выделение информации о синтетических процедурах
    - Запрос к БЯМ
    - Формирование датасета

∧МФТИ

# Обзор литературы

2. Ramos M.C., Collison C.J., White A.D. A Review of Large Language Models and Autonomous Agents in Chemistry: arXiv:2407.01603. arXiv, 2024.
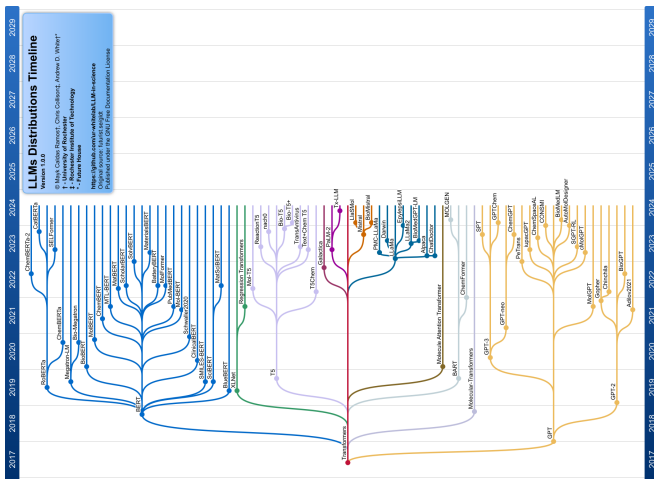


Рис. 4: Иллюстрация хронологической эволюции больших языковых моделей.

# Обзор литературы

3. Lála J. et al. PaperQA: Retrieval-Augmented Generative Agent for Scientific Research: arXiv:2312.07559. arXiv, 2023.



Рис. 5: PaperQA — это агент, который преобразует вопрос в ответ с указанием источников. Агент использует три инструмента: поиск, сбор данных и ответ на вопрос. Инструменты позволяют ему находить и анализировать соответствующие полнотекстовые исследовательские работы, определять конкретные разделы в работе, которые помогают ответить на вопрос, суммировать эти разделы с контекстом вопроса (называемые доказательствами), а затем генерировать ответ на основе доказательств.

# Обзор литературы

4. Zheng Z. et al. ChatGPT Chemistry Assistant for Text Mining and Prediction of MOF Synthesis.
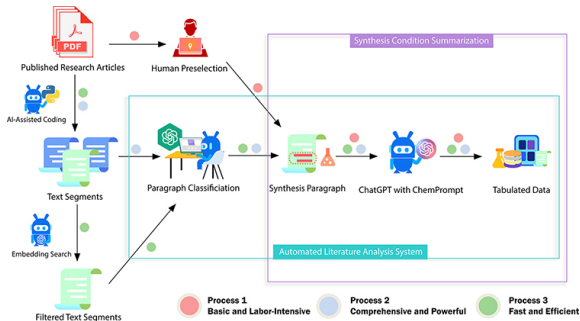


Рис. 6: ChatGPT Chemistry Assistant ChatGPT и ChemPrompt для эффективного анализа текста и обобщения условий синтеза MOF из разнообразного набора опубликованных исследовательских статей.

# Текущие результаты и планы

| Catalyst type: | PE: ZN on SiO2 |
|---|---|
| KEY WORDS: | T/A/C/D: (((((Ziegler-Natta catalyst+) or (silica?supported Ziegler-Natta catalyst+)) and (titanium +chloride)) and gas-phase) and (PE or polyethylene)) |
| RESTRICTION: | ALL |
| number of documents before relevance is determined | There were about 2491 documents (ORBIT) |
| DATE | Priority date from 01/01/2002 |
| Date of research | 13.03.2022 |
| 2491 | patented inventions |
| 0,5 | owned by top 10 players |

Таблица 1: результаты поиска по ключевым словам

∧МФТИ

| | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CHINA PETROLEUM & CHEMICAL | 3 | 6 | 2 | 2 | 4 | 2 | | 27 | 27 | 30 | 18 | 30 | 21 | 11 | 56 | 27 | 22 | 10 | 8 | | |
| JAPAN POLYPROPYLENE | 4 | 6 | 6 | 4 | 9 | 18 | 15 | 19 | 15 | 9 | 14 | 23 | 14 | 13 | 6 | 8 | 9 | 7 | 5 | 3 | |
| CHINA SINOPEC BEIJING RESEARCH INSTITUTE OF CHEMIC... | | 2 | 1 | 1 | 1 | 1 | | 5 | 14 | 14 | 1 | 26 | 12 | 10 | 51 | 21 | 16 | 2 | 7 | | |
| BOREALIS | | | | | | | 6 | 16 | 7 | 10 | 10 | 14 | 14 | 9 | 12 | 19 | 12 | 11 | 19 | 3 | |
| BASELL POLYOLEFINE | 8 | 16 | 9 | 6 | 6 | 3 | 6 | 8 | 2 | 10 | 10 | 9 | 1 | 8 | 8 | 11 | 9 | 7 | 7 | 5 | |
| SUMITOMO CHEMICAL | 26 | 30 | 12 | 17 | 6 | 10 | 13 | 14 | 4 | 7 | 3 | | 3 | | | 1 | 1 | | 1 | | |
| DOW GLOBAL TECHNOLOGIES | 10 | 3 | 4 | 10 | 9 | 7 | 10 | 3 | 9 | 3 | 2 | 1 | 2 | 1 | 1 | 4 | 1 | 7 | 6 | 1 | |
| EXXONMOBIL CHEMICAL PATENTS | 2 | 3 | 2 | 23 | 3 | | 1 | | | | | 3 | 4 | 1 | 10 | 7 | 16 | 2 | 6 | 5 | |
| SINOPEC | | 2 | 1 | | 1 | 2 | | 16 | 8 | 10 | 16 | 2 | | | 4 | 6 | 1 | 6 | | | |
| MITSUI CHEMICALS | 5 | 4 | 6 | 5 | 4 | 6 | 2 | 6 | 1 | | | 2 | 1 | 12 | 1 | 1 | 2 | 3 | 2 | | |
| SABIC GLOBAL TECHNOLOGIES | | | | | | | | | | | | 1 | 11 | 11 | 17 | 14 | 2 | 6 | 2 | | |
| ASAHI KASEI | | | | | | | | | | | 2 | 6 | 1 | 8 | 5 | 14 | 8 | 5 | 1 | 4 | |
| LOTTE CHEMICAL | 1 | 1 | | | | 1 | | 2 | 3 | | 2 | 1 | 5 | 4 | 6 | 7 | 8 | 9 | | | |
| LG CHEM | | | 2 | 3 | 5 | 3 | | 7 | | 2 | 1 | 1 | 7 | 5 | 8 | 1 | | | 1 | | |
| NOVA CHEMICALS | | 2 | 1 | 1 | 1 | 1 | | 4 | 1 | | | 4 | 3 | 2 | 13 | 3 | 2 | 2 | 6 | | |
| PRIME POLYMER | 2 | 1 | 1 | 2 | 6 | 3 | 5 | 2 | 1 | | 2 | 2 | 4 | | 1 | 3 | 2 | 5 | 2 | | |
| PETROCHINA | | | 1 | 1 | | 3 | 2 | 4 | 1 | 2 | 7 | 2 | 6 | 3 | 2 | 4 | 5 | 1 | | | |
| UNIVATION TECHNOLOGIES | 1 | 11 | 5 | 4 | 1 | 2 | 1 | | 1 | | 4 | 1 | 1 | | 3 | 4 | | | | | |
| WR GRACE | 7 | 1 | 1 | 1 | | 3 | | 3 | 1 | 4 | 2 | 1 | 4 | | | 3 | 4 | 4 | 1 | | |
| BOREALIS TECHNOLOGY | 4 | 3 | 6 | 2 | 5 | 11 | 2 | | | | | | | 1 | 1 | | | | | | |

Рис. 7: Тепловая карта распределения патентов

# Текущие результаты и планы

План работ на третий семестр

- ▶ Создать набор агентов и инструментов для решения задачи извлечения информации на базе одной БЯМ (ChatGPT-3.5Turbo)
- ▶ Собрать небольшой датасет (около 50-100 наблюдений) для оценки эффективности извлечения данных
- ▶ Оценить эффективность обработки данных с помощью метрики F1-score

МФТИ

# Список литературы

1. Ananikov V. Top 20 Influential AI-Based Technologies in Chemistry. Chemistry, 2024.

2. Ramos M.C., Collison C.J., White A.D. A Review of Large Language Models and Autonomous Agents in Chemistry: arXiv:2407.01603. arXiv, 2024.

3. Lála J. et al. PaperQA: Retrieval-Augmented Generative Agent for Scientific Research: arXiv:2312.07559. arXiv, 2023.

4. Zheng Z. et al. ChatGPT Chemistry Assistant for Text Mining and Prediction of MOF Synthesis.