

Report 1.2

Lynx: System for Financial News Analysis Using NLP
and Graph Databases

Andrei Zhdanov Mikhail Sofin

October 7, 2025

A. Team

Team Akira:

Andrei Zhdanov (an.zhdanov@innopolis.university)

Mikhail Sofin (m.sofin@innopolis.university)

B. Project topic

Lynx: System for Financial News Analysis Using NLP and Graph Databases.

C. Links

Link to Github repository with all source code and notebooks: [Lynx](#)

Link to Miro: [Miro](#)

D. What has been done so far

Over the past three weeks, our work has been focused primarily on two main areas: developing part of the backend infrastructure and conducting experiments in the machine learning component of the project. In the architecture diagram below, the parts we concentrated on are highlighted in red to clearly indicate our scope of work during this period.

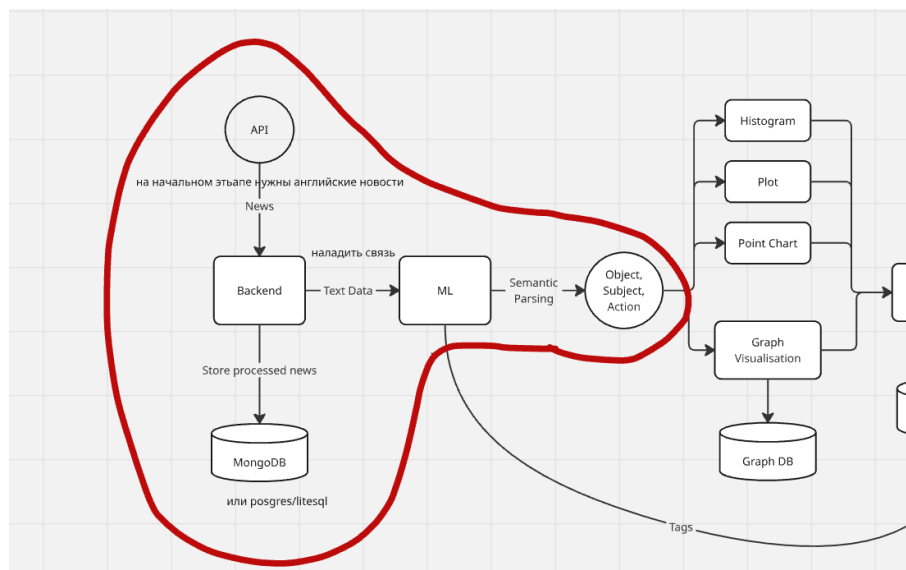


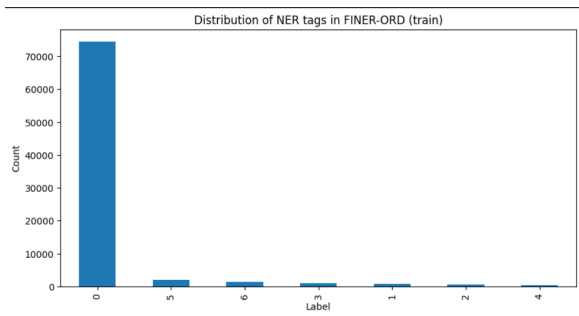
Fig. 1. Architecture

- **Backend Development:** We initiated the development of the backend system. This included the initialization of a MongoDB database for storing raw news articles and a Qdrant vector database for storing news text embeddings to enable semantic search capabilities.

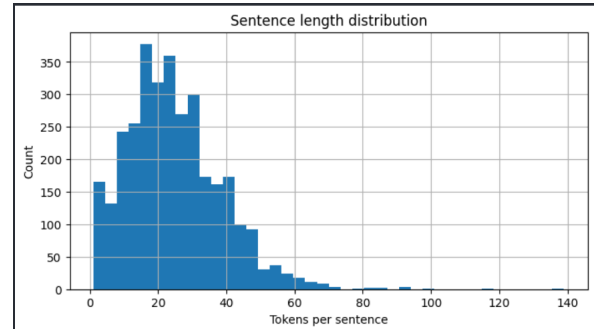


Fig. 2. backend

- **Exploratory Data Analysis (EDA):** We conducted thorough exploratory data analysis on two datasets:
 - **FINER-ORD** for Named Entity Recognition (NER) tasks.



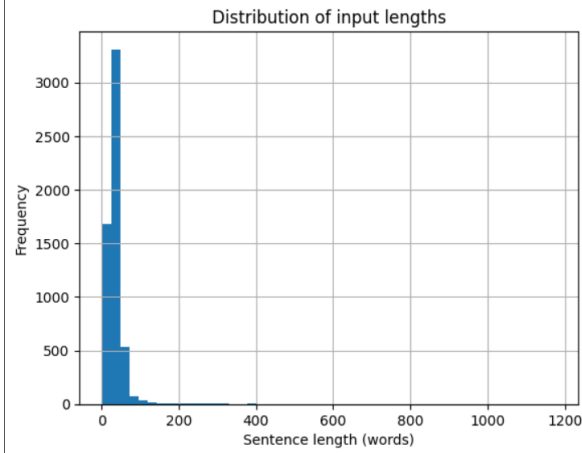
(a) Distribution of NER tags



(b) Sentence length statistics

Fig. 3. Statistical analysis of the FINER-ORD dataset

- **FINRED** for Relation Extraction (RE) tasks.



(a) Sentence length statistics

```
[('product_or_material_produced', 1373),
 ('industry', 1167),
 ('owned_by', 673),
 ('headquarters_location', 602),
 ('parent_organization', 524),
 ('employer', 482),
 ('owner_of', 450),
 ('subsidiary', 351),
 ('developer', 267),
 ('location_of_formation', 238),
 ('manufacturer', 198),
 ('position_held', 186),
 ('founded_by', 174),
 ('chief_executive_officer', 160),
 ('stock_exchange', 159),
 ('operator', 157),
 ('chairperson', 154),
 ('legal_form', 97),
 ('brand', 91),
 ('creator', 84)]
```

(b) Classes distribution

Fig. 4. Statistical analysis of the FINRED dataset

The EDA included statistical analysis of data distribution, class imbalance checking, and inspection of entity and relation labels to ensure data quality and consistency.

- **NER Model Fine-tuning:** For the Named Entity Recognition (NER) task, we selected the pre-trained model `dslim/bert-base-NER` as the foundation for fine-tuning. This model is based on the BERT architecture and has been pre-trained specifically for NER tasks using the CoNLL-2003 dataset, which includes common entity types such as organizations, locations, and persons. By fine-tuning `dslim/bert-base-NER` on the FINER-ORD dataset, we were able to effectively adapt the model to the financial domain, achieving an accuracy of 90% on the validation set, demonstrating strong performance in identifying named entities in text.

Начинаем обучение...

[612/612 01:28, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Precision	Recall	F1
1	0.047500	0.044900	0.883607	0.885057	0.884331
2	0.022600	0.041599	0.912252	0.904762	0.908491
3	0.017600	0.042218	0.907743	0.904762	0.906250

обучение завершено.

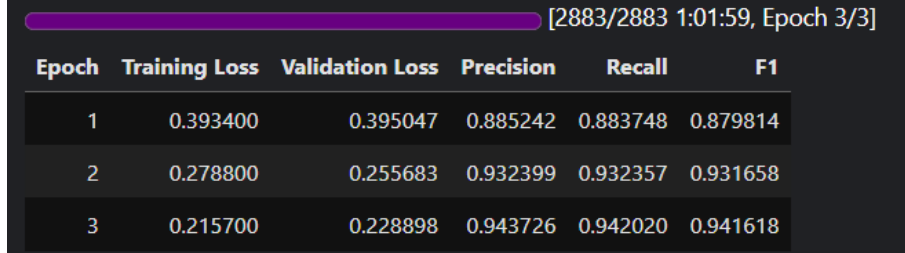
Fig. 5. Training NER model

```
Entity: Donald Trump, Label: PER, Score: 0.9904
Entity: Russia, Label: LOC, Score: 0.9963
```

Fig. 6. Example of NER

- **RE Model Fine-tuning:** For the Relation Extraction (RE) task, we fine-tuned the `microsoft/deberta-v3-small` model. This model belongs to the DeBERTa (Decoding-enhanced BERT with Disentangled Attention) family, which improves upon BERT and RoBERTa by separating content and positional embeddings and using an enhanced mask decoder mechanism. We selected this model because it

achieves strong performance across various NLP benchmarks while maintaining relatively low computational cost due to its smaller size compared to larger DeBERTa variants. Its strong contextual understanding makes it particularly effective for identifying semantic relationships between entities in complex financial texts. After fine-tuning on the FINRED dataset, the model achieved an accuracy of 96%, demonstrating its ability to accurately capture relational structures in financial news.



[2883/2883 1:01:59, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Precision	Recall	F1
1	0.393400	0.395047	0.885242	0.883748	0.879814
2	0.278800	0.255683	0.932399	0.932357	0.931658
3	0.215700	0.228898	0.943726	0.942020	0.941618

Fig. 7. Training

- **Dataset Construction for Russian News:** We have begun collecting news articles in Russian and preparing a dataset with annotations using large language models (LLMs). This dataset will be used to train and evaluate models for NER and RE in Russian-language news, enabling multilingual capabilities for our system.

```
our_feeds = {
    'Kommersant Economy': 'https://www.kommersant.ru/RSS/economy.xml',
    'Vedomosti Politics': 'https://www.vedomosti.ru/rss/politics',
    'Vedomosti Finance': 'https://www.vedomosti.ru/rss/finance',
    'RBC News': 'https://rssexport.rbc.ru/rbcnews/news/30/full.rss',
    'Interfax': 'https://www.interfax.ru/rss.asp',
    'RIA Novosti': 'https://ria.ru/export/rss2/index.xml',
    'Lenta Politics': 'https://lenta.ru/rss/news/politics/',
    'TASS': 'https://tass.ru/rss/v2.xml',
    'Rossiyskaya Gazeta': 'https://rg.ru/xml/index.xml',
    'Gazeta Ru': 'https://www.gazeta.ru/export/rss/first.xml'
}
```

(a) News RSS collection pipeline

Title,Description,Link,Publication Date
 Орбан заявил о нежелании вводить евро в Венгрии из-за «развала» Евросоюза, "Использование евро
 сентябрь в Москве второй раз с начала века оказался солнечнее августа,"Сентябрь в Москве оказа
 Tesla подали в суд из-за аварии с Cybertruck студентами, "В США родители двух
 Семин назвали неожиданным уход гендиректора «Спартак», "Бывший тренер сборной России похвалит
 Уличные музыканты «парализовали» работу музыкальной школы в Петербурге, "Администрация Централь
 Минкультуры предложило критерии восприятия традиционных ценностей в кино, "Критерии применят дл
 В Ростове заочно приговорили к пожизненному организатора покушения из СВУ, "Новый окружной воен
 "Сильно узнала, что Сырский расформировал бывшую группировку «Хортица», "Группировка «Днепр» (быв
 В Москве наградили лауреатов НК-премии «за кадры», "В Москве подвели итоги премии правительства
 Оценивая гражданство Касаткина досрочно завершил сезон, "Теннисистка в этом сезоне впервые
 Умерла автор цикла книг «Хроники Ратмира» Дарья Кулер, "Британская писательница Дилли Кулер ук
 Гладков сообщил о погибшем в результате ракетного удара по Белгороду, "Один человек погиб и еще
 В Волгограде эвакуировали ТЦ из-за задымления, "Посетителей торгового центра «Город«
 ВКС Украины впервые выступили в роли условного противника НАТО на учениях, "Военно-морские силы
 Путин назначил двух заместителей директора ФСВН, "Заместителем Аркадия Гостева назначен генера
 «Радиостанция Судного дня» повторила сообщение от февраля 2022 года, "Радиостанция УВБ-76, более
 "Залуцкий заявил, что Украине пора посмотреть в космос", "Украине необходимо достичь технологич
 Сборная Узбекистана поедет на ЧМ-2026 в Фабью Каннаваро, "Итальянец будет готовить сборную узбе
 В Подольске задержали избившего собаку инвалида мужчину, "Полиция задержала мужчину, жестоко из
 Евразийский суд ЕС в плане Еврокомиссии найти деньги для производства дронов, "Сейчас ЕС ищет спосо
 В Урале опровергли фейк о запрете выезда из д/р из-за холеры, "Информация о распространении хо
 Красноварские власти назвали временные ограничения на отделеях ЕСВ, "Ситуация с поставками то
 болельщик за «Спартак» ректор «Бауманин» извинился за пост про ЦСКА, "Спартак» 5 октября проигр
 Hindustan Times узнала о планах Индии закупить еще пять российских Б-400, "Индия планирует заку
 СВН обвинила Лондон в планах провокации ЕС «злостными агентами Кремля», "Власти Британии расстр
 Мужчина в ножах ранил двух человек в стенах Цоя на Арбате, "Мужчина в ножах ранил двух человек в

(b) Example of collected news item

Fig. 8. Dataset creation process: collecting and preparing Russian financial news

E. Work Distribution

- Andrei Zhdanov: Data EDA, fine-tune models, ml architecture
- Mikhail Sofin: backend architecture, data collection

F. Plan for the Next Weeks

In the next stages of the project, we plan to focus on completing the integration of the backend and machine learning pipelines into a unified system. This includes connecting the data ingestion, processing, and model inference components.

Additionally, we aim to start developing the frontend part of the application and to experiment with graph-based visualizations of entities and relations using Neo4j.

Another important direction will be the continuation of dataset development — expanding the corpus with new annotated examples and improving data quality. We also intend to conduct further experiments with machine learning models, including the possibility of designing and training our own custom model for financial text understanding.