

Дослідження технологій структуризації неструктурованих текстових даних

Виконав: Домілевський О.С, ІПСА, ДА-62

Керівник: к.т.н. Сергеев-Горчинський О. О.

Об'єкт дослідження

- Неструктуровані дані

Предмет дослідження

- Структурування та аналіз текстових даних
- Платформа Elastic Stack

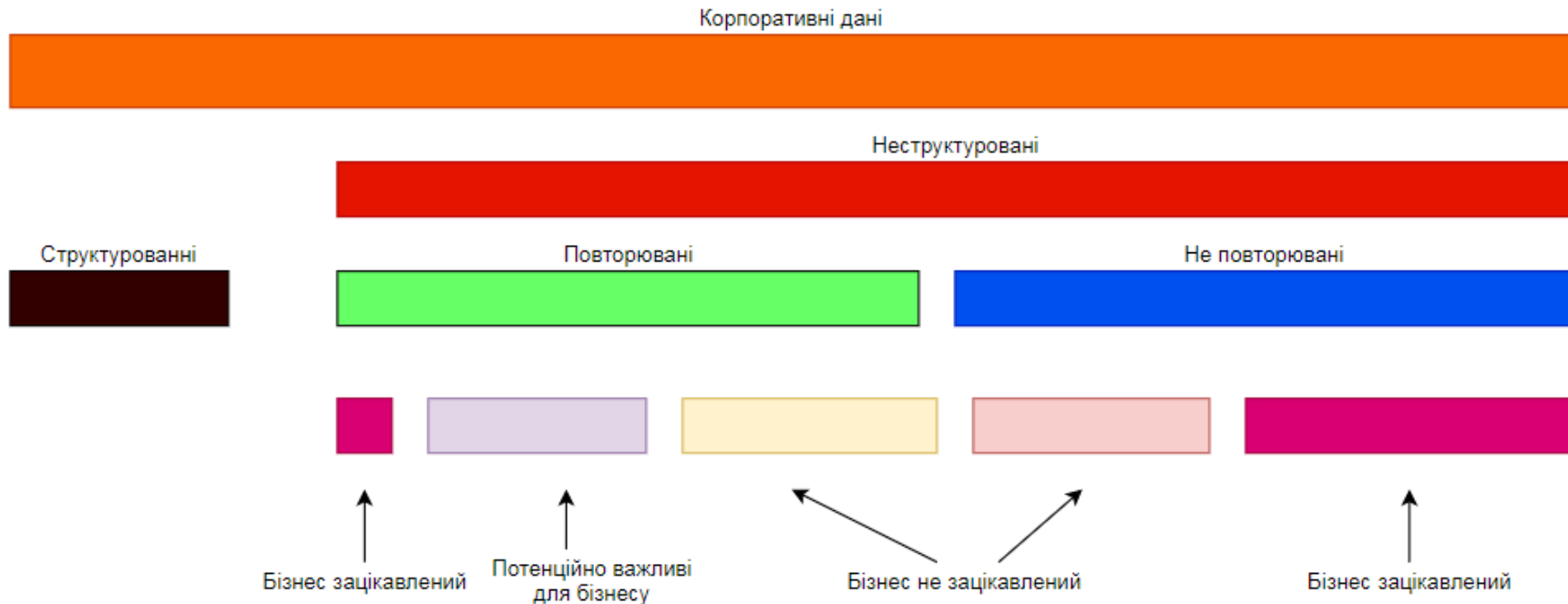
Актуальність

- 80% текстової інформації – неструктурована
- Неможливість вручну обробляти величезні обсяги інформації
- Більшість сучасних ІТ компаній використовують аналіз текстових даних

Мета

- Розглянути методи структуризації та аналізу текстових даних
- Дослідити компоненти системи Elastic Stack
- Провести аналіз сучасних NLP бібліотек
- Створити додаток для аналізу неструктурованих текстових даних та їх анотування

Дослідження даних в корпорації



Методи аналізу тексту

Основні методи

- Частота слова
- Словосполучення
- Відповідність

Передові методи

- **Класифікація**
 - Аналіз почуттів
 - Тематичний аналіз
 - Виявлення намірів
- **Анотування тексту**
 - Виявлення ключових слів
 - Виявлення особи
 - Розбіжність сенсу слова
- **Кластеризація**

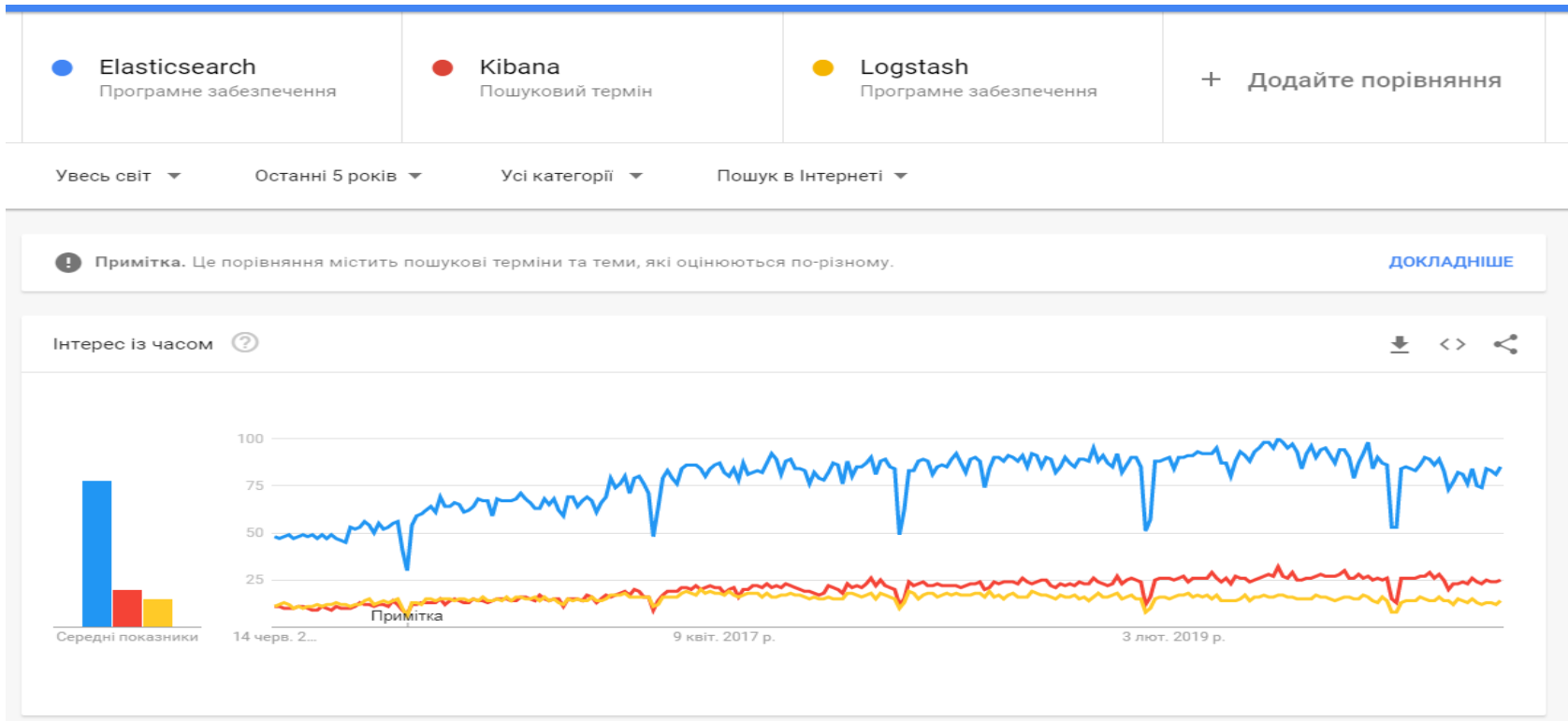
Обробка текстових даних

Слово	Potrer	Lancaster	Snowball	Лематизація	Без стоп слів
Big	big	big	big	Big	Big
Data	data	dat	data	Data	Data
analytics	analyt	analys	analyt	analytic	analytics
can	can	can	can	can	
help	help	help	help	help	help
organizations	organ	org	organ	organization	organizations
better	better	bet	better	better	better
understand	understand	understand	understand	understand	understand
the	the	the	the	the	
information	inform	inform	inform	information	information
contained	contain	contain	contain	contain	contained
within	within	within	within	within	
the	the	the	the	the	
data	data	dat	data	datum	Data

Elastic Stack



Популярність Elasticsearch, Kibana, Logstash



Elastic Search

Аналіз

Масштабованість

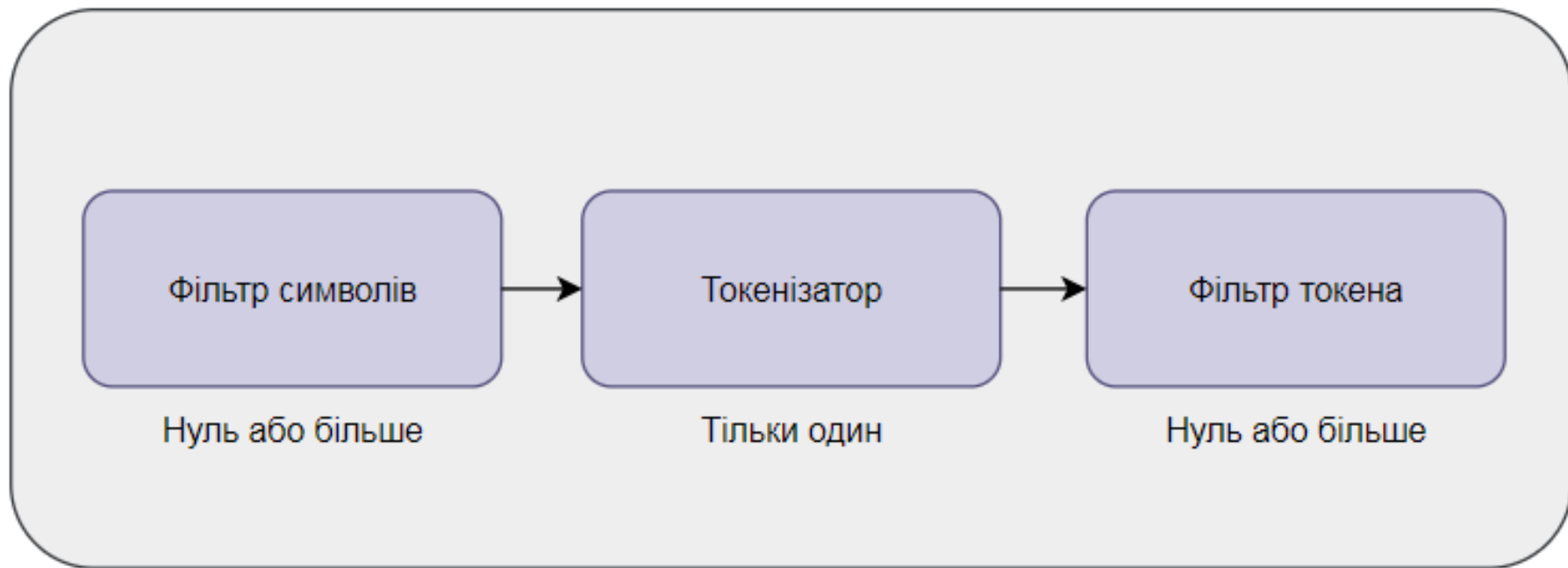
REST API

Відмовостійкість

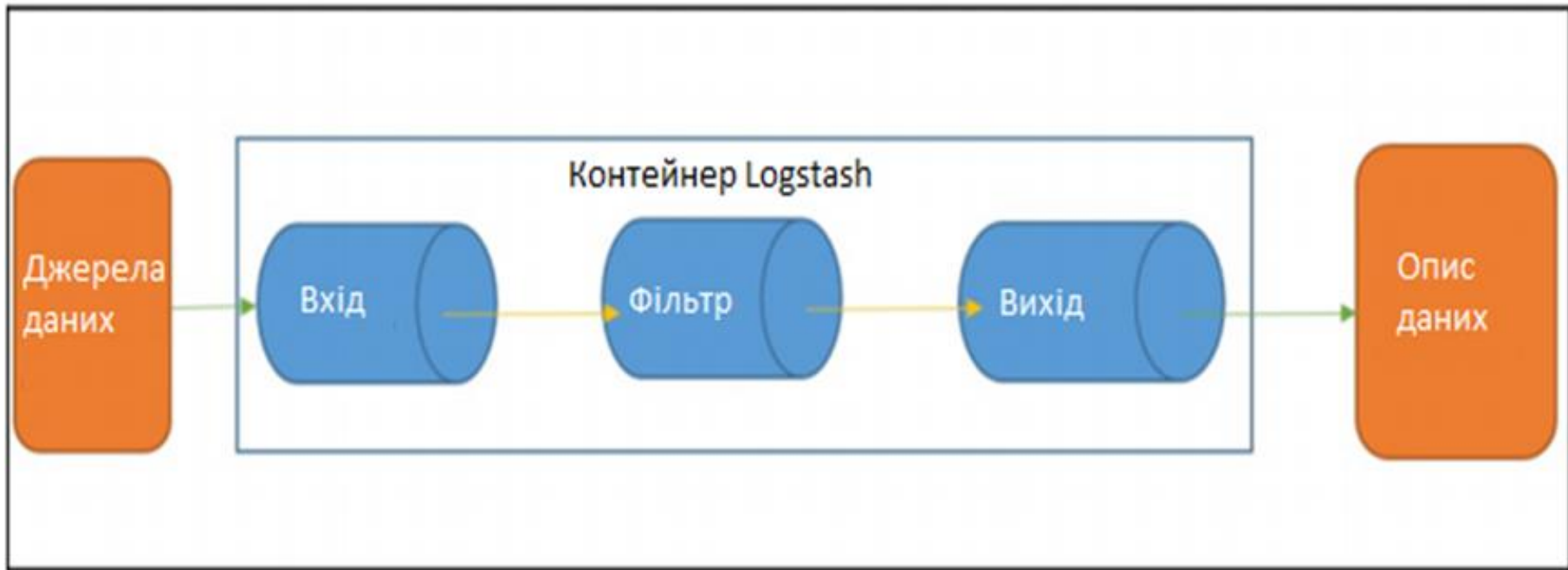
Неструктурованість



Загальний вигляд аналізатора



Logstash



Порівняння можливостей основних бібліотек для NLP

Ім'я	Spark NLP	spaCy	NLTK	CoreNLP
Виявлення речень	+	+	+	+
Токенізація	+	+	+	+
Стематизація	+	+	+	+
Лематизація	+	+	+	+
Позначення часткової мови	+	+	+	+
Іменовані сутності	+	+	+	+
Виявлення залежностей	+	+	+	+
Збіг тексту	+	+	-	+
Збіг дати	+	-	-	+
Поділ на фрагменти	+	+	+	+
Перевірка орфографії	+	-	-	-
Аналіз настроїв	+	-	-	+
Можливість перетренувати моделі	+	+	+	+
Тренувальні моделі	+	+	+	+

Работа програми

Выберите файл: Файл не выбран

Имя файла 1.txt

Тип файла text/plain

Размер в байтах 622

```
{'text': 'In modern conditions a comprehensive reconstruction of all spheres of public life arises fundamentally new requirements for training future information security professionals in higher education institutions, which differ not only in a high level of professional skills, but also in a harmoniously developed interests, in ability to continuously improve the educational level and react to changes in life. Therefore, a special place in the training of future professionals takes a generation of professionally significant qualities that is the basis for the formation of a coherent worldview of existing reality in youth.\n', 'chain': {'life': {'arise': {'requirement': [1, 13]}}, 'requirement': {'train': {'information': [1, 18]}}, 'institution': {'differ': {'level': [1, 28]}}, 'skill': {'develop': {'interest': [1, 43]}}, 'ability': {'improve': {'level': [1, 49]}}, 'level': {'react': {'change': [1, 54]}}, 'professional': {'take': {'generation': [2, 11]}}, 'worldview': {'exist': {'reality': [2, 30]}}}}
```

Загрузка ещё одного файла

Выберите файл: Файл не выбран

Приклад зберігання в базі Elastic Search

```
Result Source
{
  "_index": "diplom",
  "_type": "text_to_check",
  "_id": "a0CTmXIBDmFsjpQr1It9",
  "_version": 1,
  "_score": 1,
  "_source": {
    "text": "Text analysis is the automated process of understanding and sorting unstructured text, making it easier to manage. Word cloud tools, for example, are used to perform very basic text analysis techniques, like detecting keywords and phrases that appear most often in your your data. However, to sort your data into custom categories, you'll need to use more advanced text analysis tools, which you can try out right away with MonkeyLearn. Text analysis tools, like t s free online sentiment analyzer, are often used to unearth valuable insights in social media conversations, survey responses, online reviews, and more. Maybe you're new to artificial intelligence and work in customer support, sales or product. You might even be a data-savvy analyst or software developer. Either way, this guide offers a comprehensive introduction to text analysis with machine learning ",
    "chain": {
      "understanding": { "sort": { "text": [ 1,10] } },
      "example": { "use": { "text": [ 2,7], "perform": { "text": [ 2,9] } },
      "technique": { "detect": { "keyword": [ 2,16] } },
      "phrase": { "appear": { "datum": [ 2,21] } },
      "category": { "will": { "text": [ 3,10], "need": { "text": [ 3,11], "use": { "text": [ 3,13] } },
      "analyzer": { "use": { "insight": [ 4,14] } },
      "guide": { "offer": { "introduction": [ 7,5] } }
    }
  }
}
```

Підсумки виконаної роботи

1. Проаналізовано методи аналізу тексту
2. Проаналізовано систему Elastic Stack
3. Порівняно бібліотеки NLP
4. На основі бібліотеки spacy створено веб-додаток, який аналізує текстову інформацію та зберігає в Elastic Search

Майбутні напрямки роботи та досліджень

- Проаналізувати варіанти оптимізації ресурсів, необхідних для роботи стеку
- Реалізувати Веб-сервіс у вигляді плагіну
- Забезпечити функцію створення шаблонів для аннотування текстових даних
- Забезпечити функцію автоматичного завантаження тексту з cloud-платформ

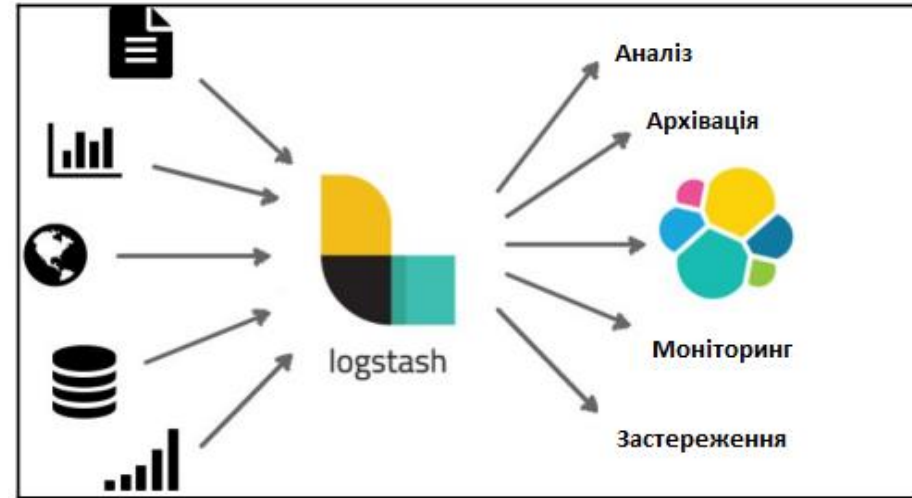
Дякую за увагу!

Технології, що використовуються

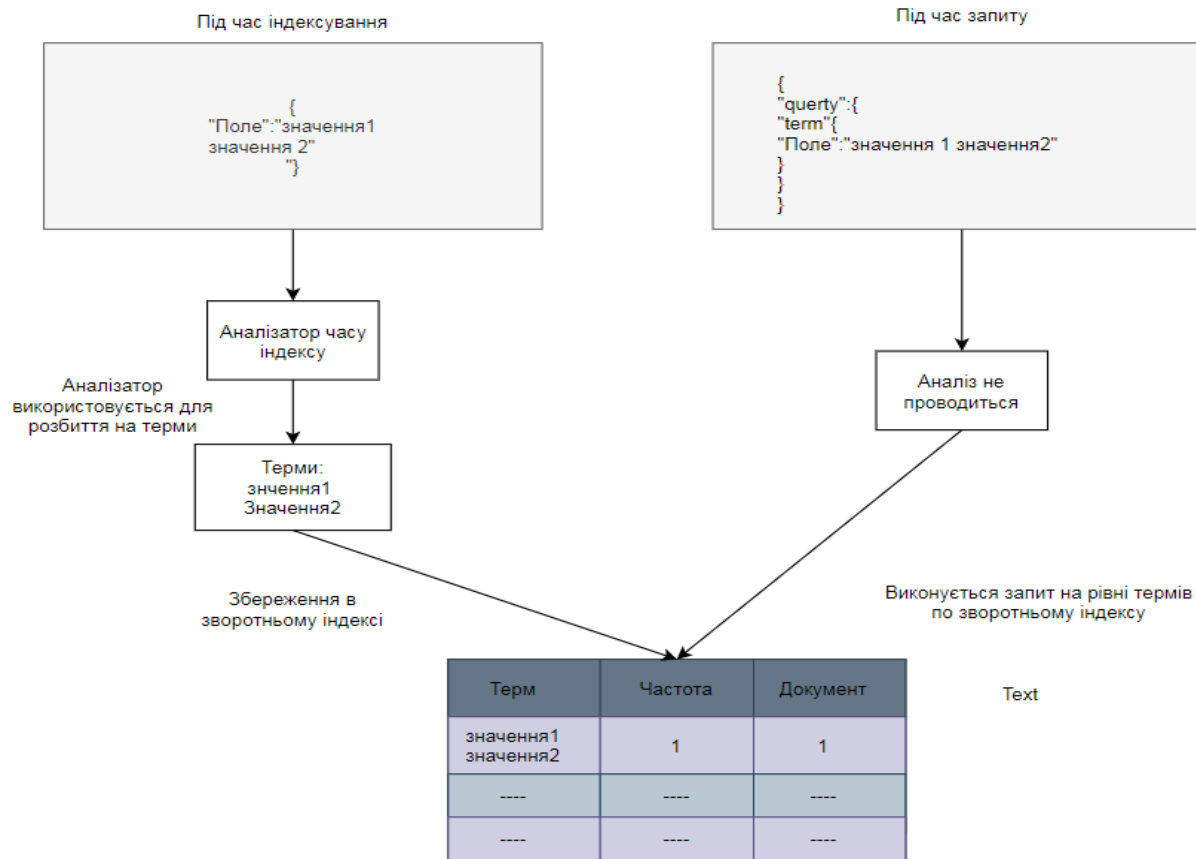


Apache

Типова взаємодія користувача з Kibana та Logstash



Виконання Term запиту



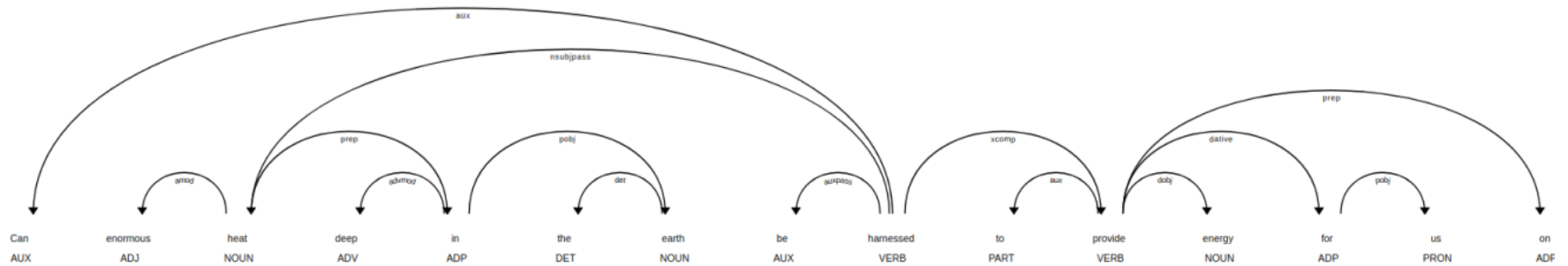
Типи агрегацій, що підтримує ElasticSearch

- Сегментні
- Метричні
- Матричні
- Агрегації контейнерів

Загальна структура індексу



Приклад розбиття на частини мови



Перевірка швидкості бібліотек spaCy і NLTK

