

Дослідження методів обробки потоків даних у Big Data

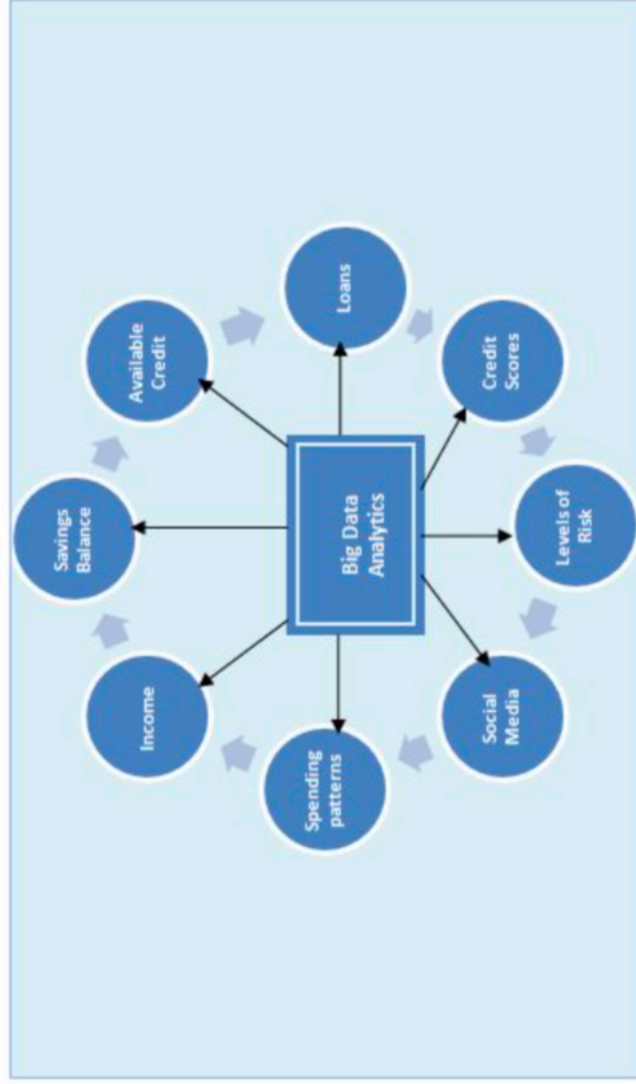
виконав
ст. гр. ІПЗм-17-1 Рукавиця А.С.
керівник
проф. каф. ПІ Руткас А.Г.

Додаток А
Слайди презентації

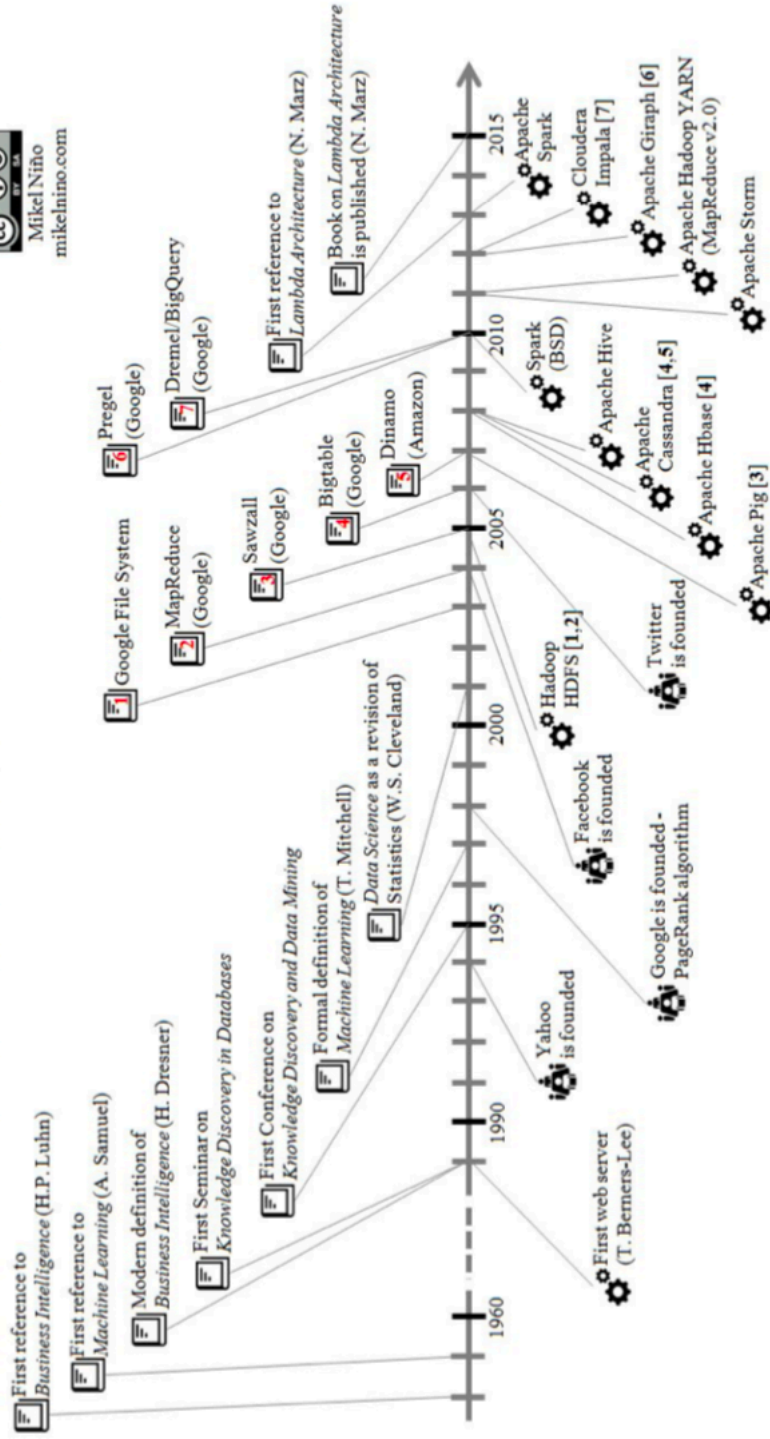
Мета роботи. Об'єкт і предмет дослідження

- Дослідження алгоритмів обробки потоків даних у режимі близькому до реального часу в хмарному середовищі із застосуванням паралельних та розподілених обчислень
- Дослідження способів забезпечення відмовостійкості при обробці потоків даних
- Дослідження способів обробки даних що приходять із запізненням
- Дослідження Big Data екосистем Spark, Flink та Kafka

Вплив Big Data



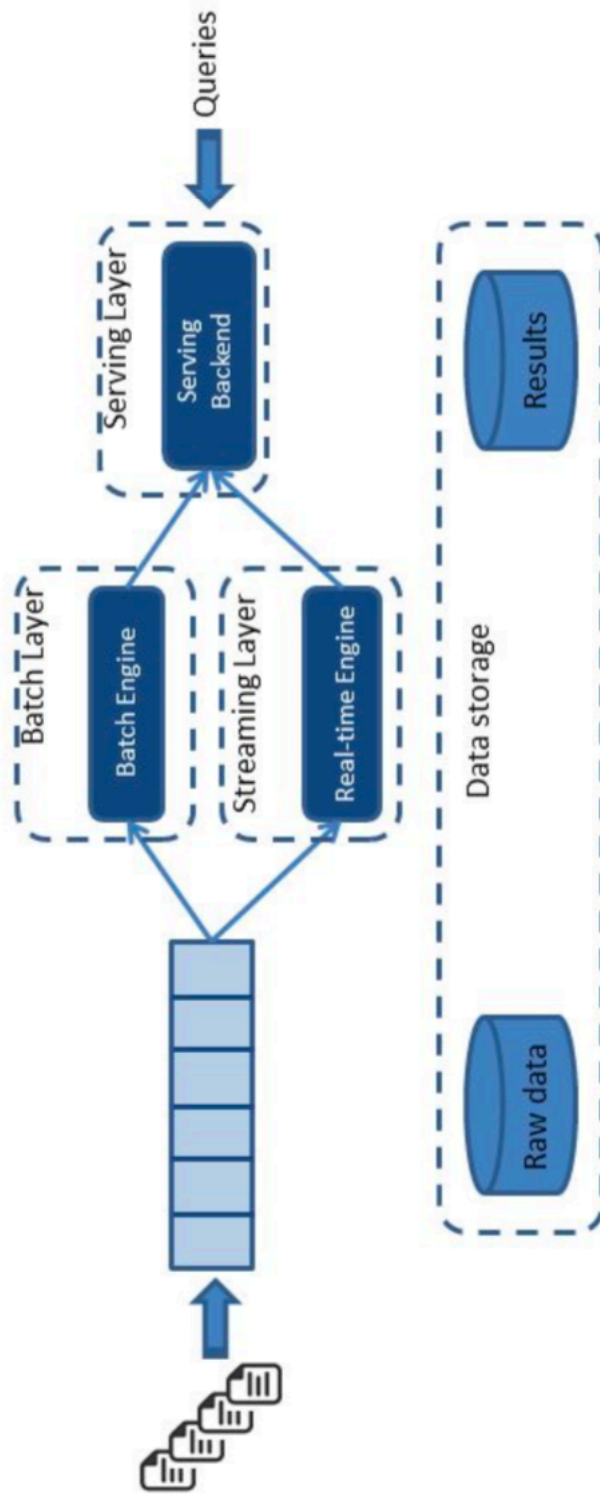
CHRONOLOGY OF ANTECEDENTS, ORIGIN AND DEVELOPMENT OF BIG DATA



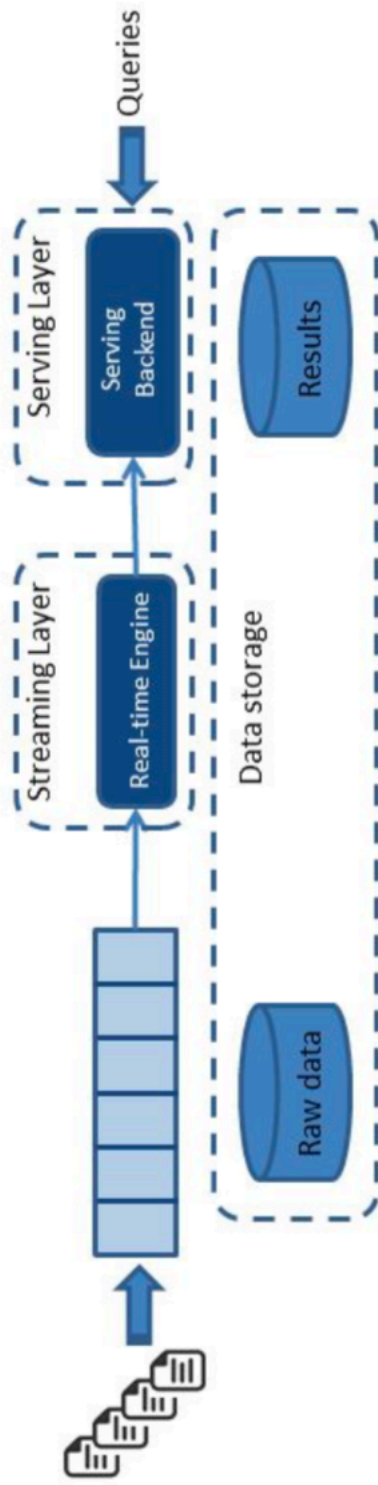
Dissemination
 Company
 Tool
 [n] Based on dissemination

v1.0 (150924)

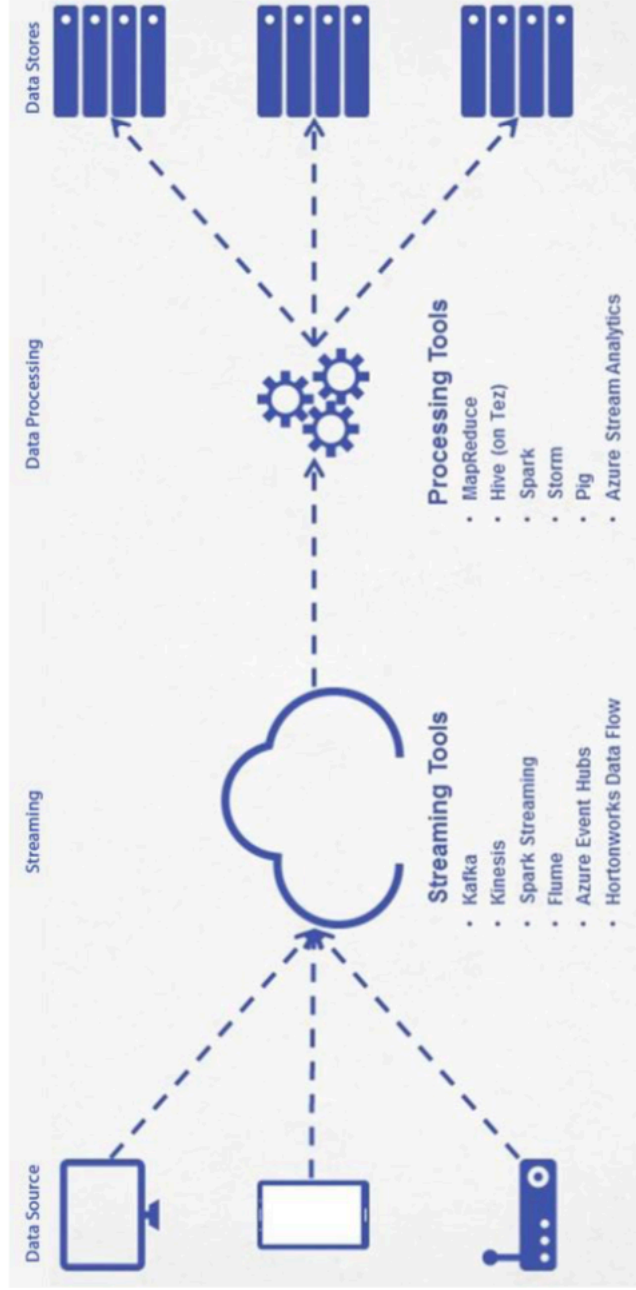
Lambda architecture



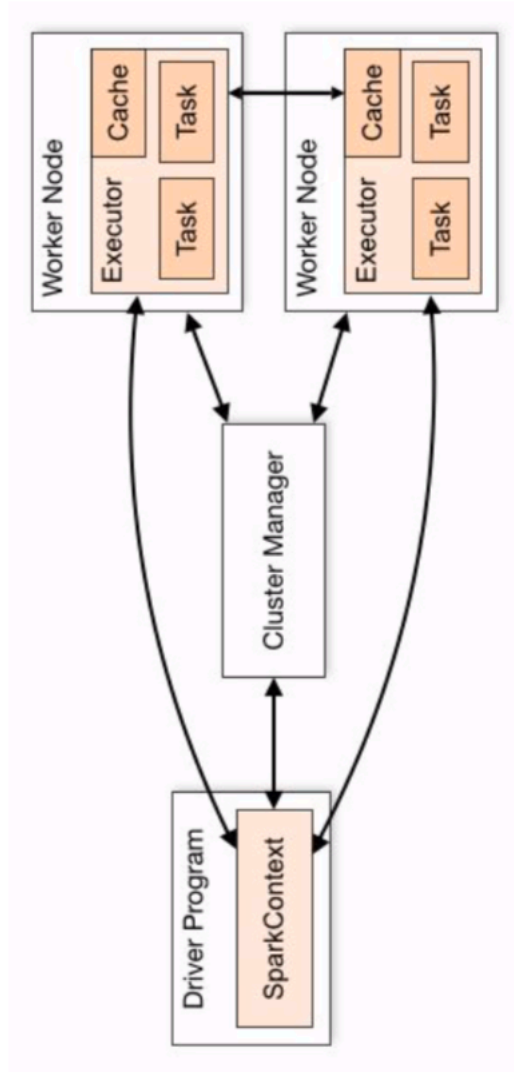
Карра архітектура



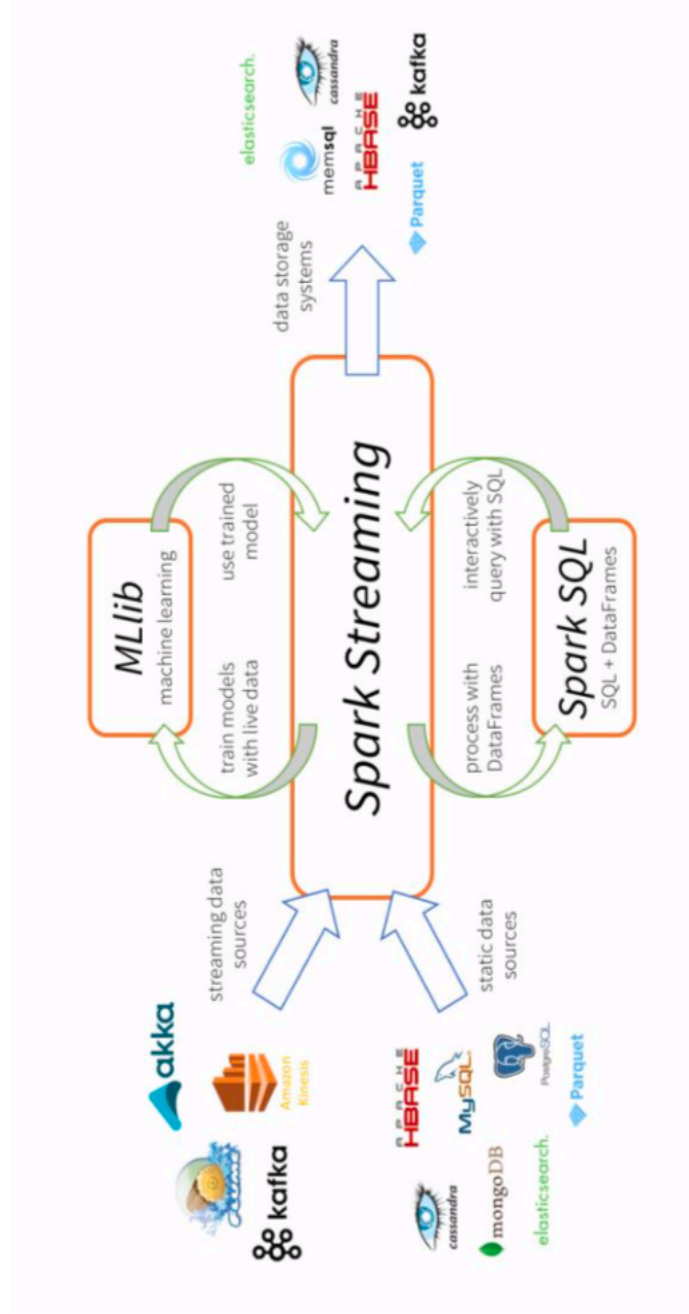
Потоки даних



Spark Cluster



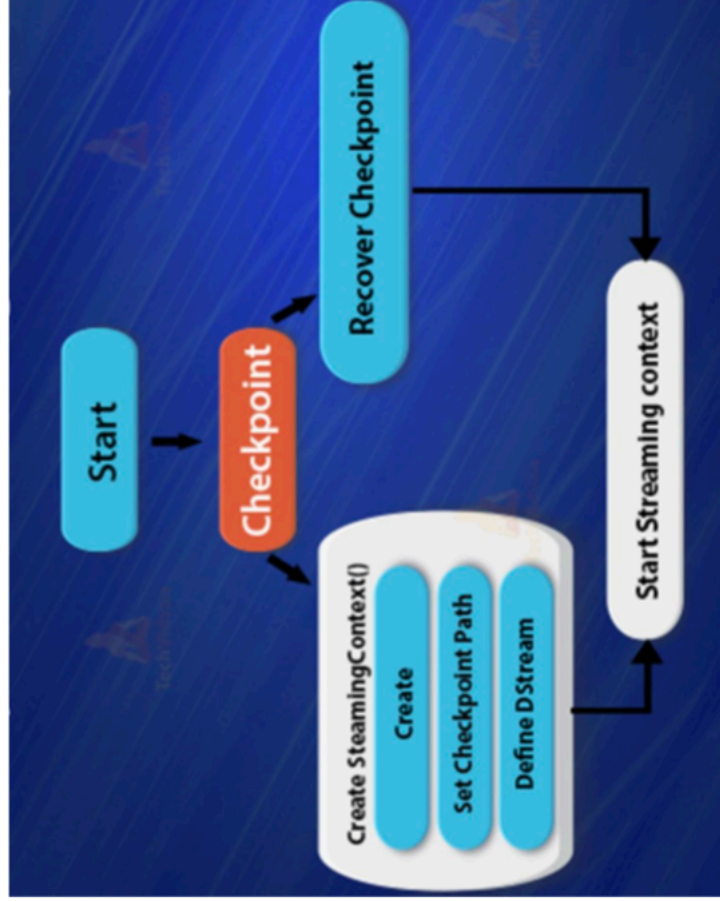
Spark Streaming



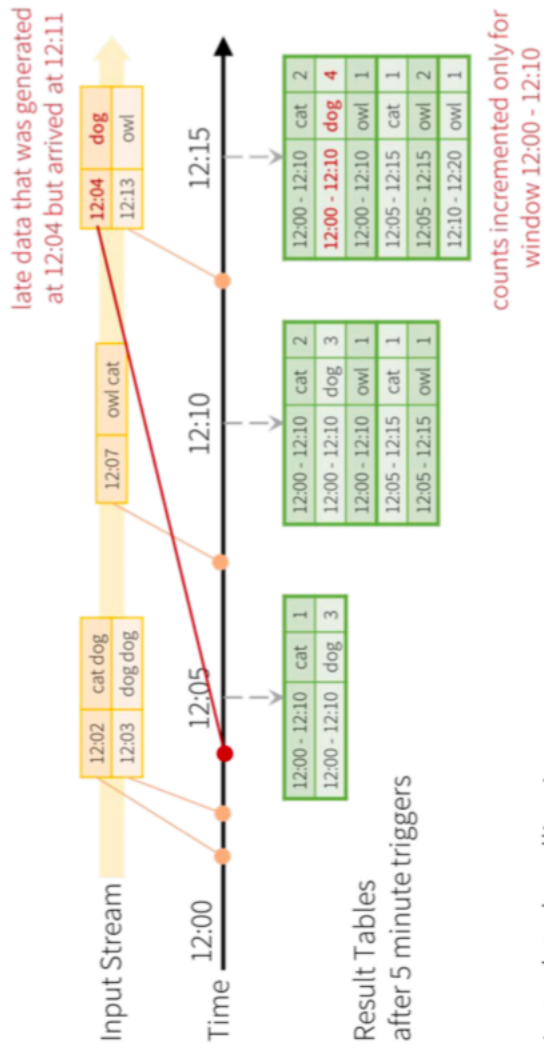
Spark Streaming використовує батч інтервали для обробки подій даних



Відновлення після падіння

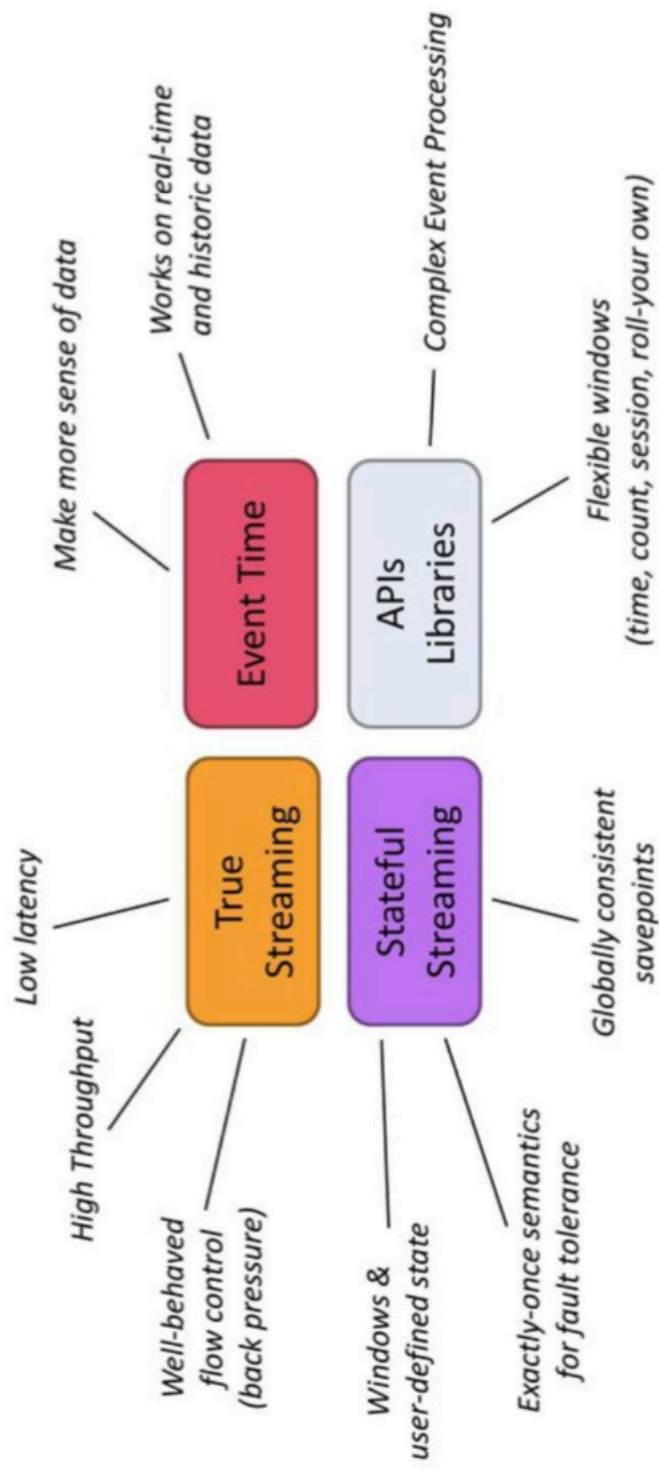


Windowing



Late data handling in Windowed Grouped Aggregation

Apache Flink



Spark vs Flink



Spark is not a true real time processing. It is Near to real time processing framework.

At heart Spark is a Batch processing framework.

Spark's streaming computation model is based on Microbatching.

Spark is implemented in Scala.

Spark does not have an efficient memory manager.

Frequently gets Out of Memory.



Apache Flink

Flink is True processing framework.

At heart Flink is a Stream processing framework.

Flink's Streaming model is based on Windowing and Checkpointing.

Flink is implemented in Java.

Flink has it's own efficient automatic memory manager.

Висновки

- Було досліджені алгоритми обробки потоків даних у режимі близькому до реального часу в хмарному середовищі із застосуванням паралельних та розподілених обчислень
- Були дослідження способи забезпечення відмовостійкості при обробці потоків даних та способи обробки даних що приходять із запізненням з використанням Big Data екосистем Spark, Flink та Kafka