# Stream2Graph: Dynamic Knowledge Graph for Online Learning Applied in Large-scale Network
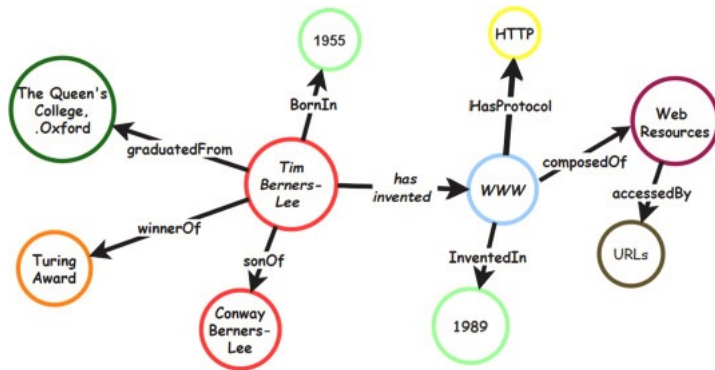
데이터사이언스학과 빅데이터 관리 및 응용 연구실

석사과정 김민선

2023-07-05

# Preliminaries
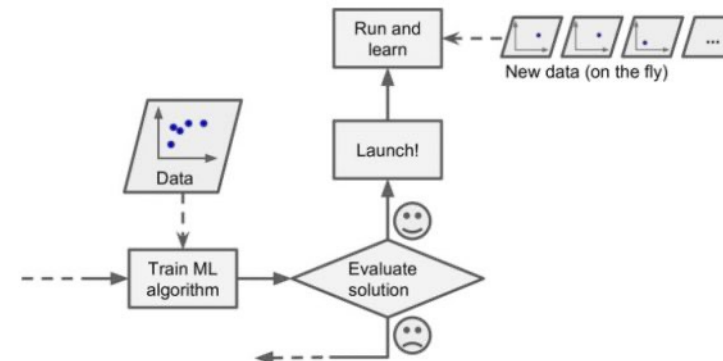
- ## Knowledge Graph



"
*…a graph of data intended to accumulate and convey knowledge of the real world, whose nodes represent entities of interest and whose edges represent potentially different relations between these entities.*
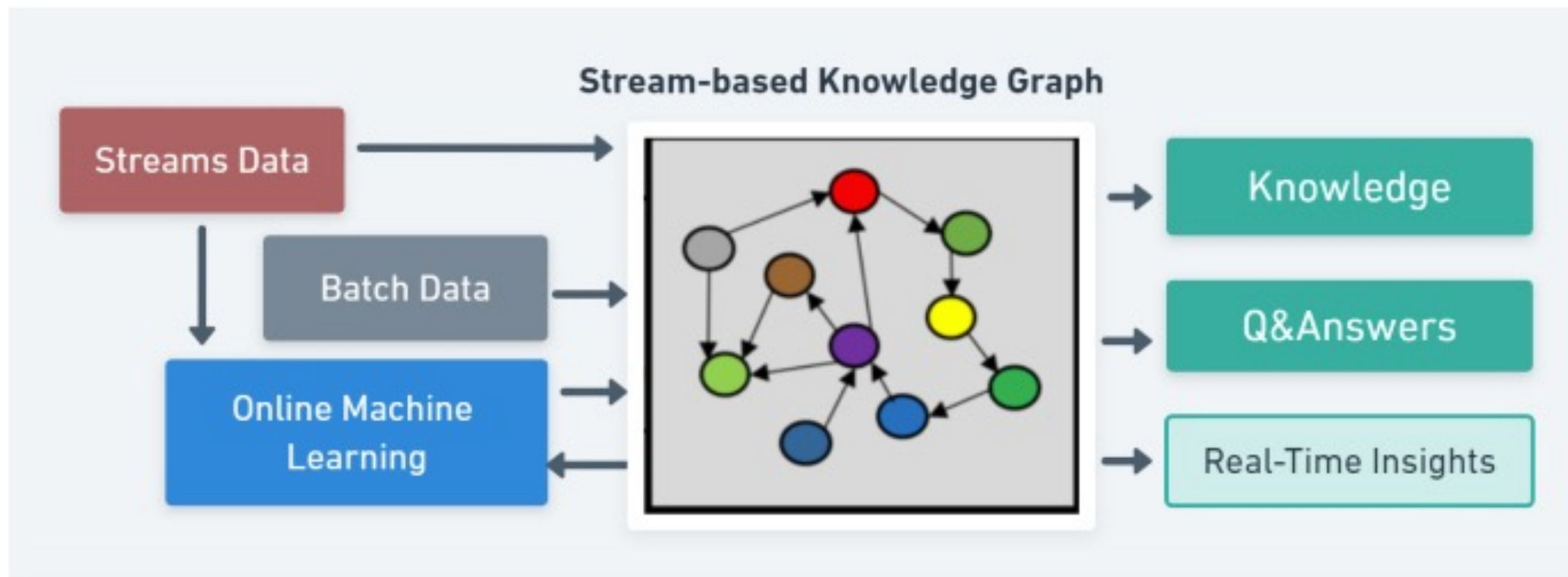
- ## Online Learning



"
*… to maximize the accuracy of the sequence of predictions made By the online learner, given the knowledge of correct answers to Previous prediction tasks and possibly additional information*

# Stream2Graph

A domain-agnostic system to easily build and operationalize stream-based knowledge graphs and combine them with online learning applications



Stream-based Knowledge Graph

Streams Data

Batch Data

Online Machine Learning

Knowledge

Q&Answers

Real-Time Insights

# Desiderata

1. Heterogeneous data
2. Training and deployment of predictive models on evolving data
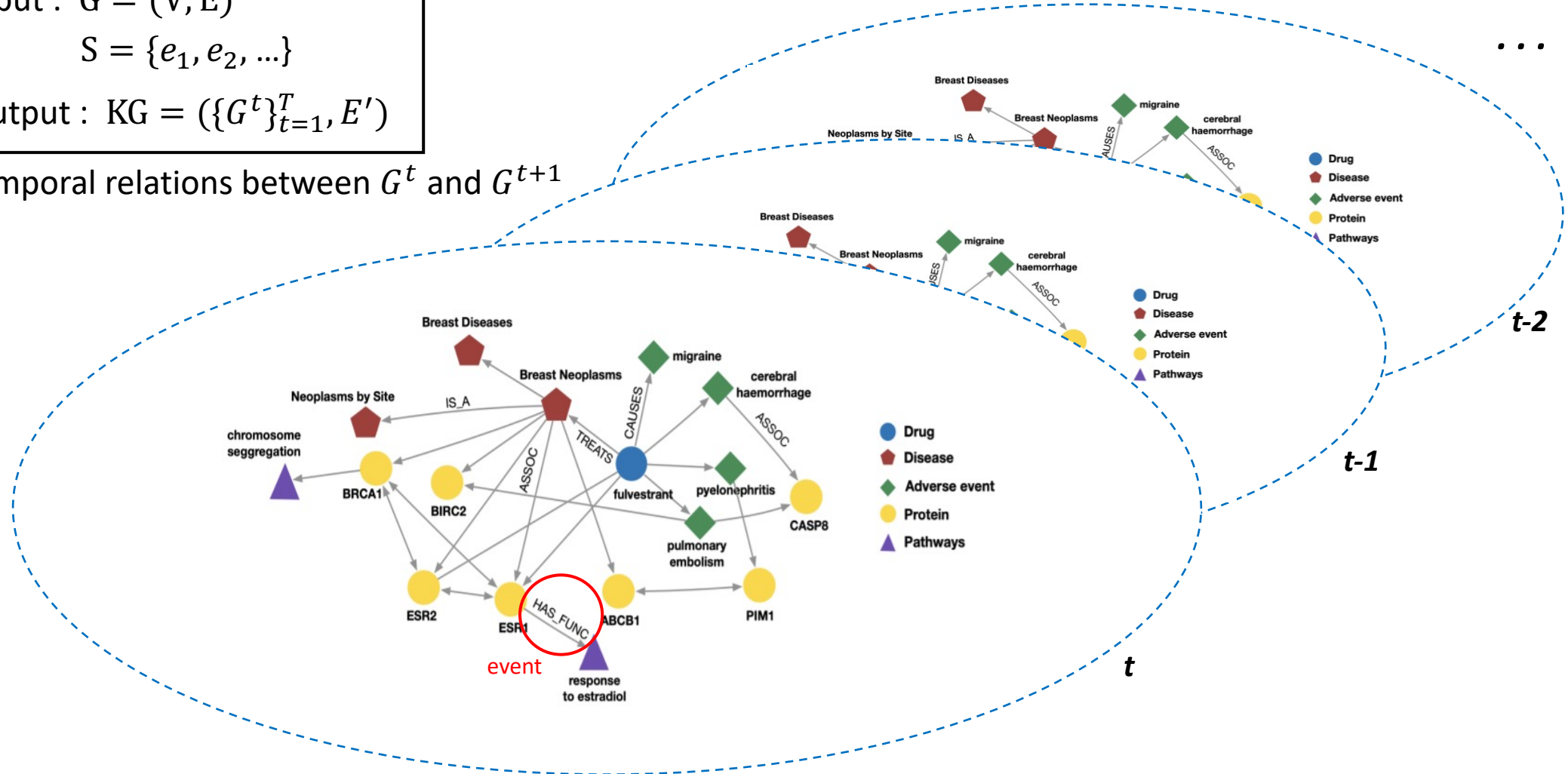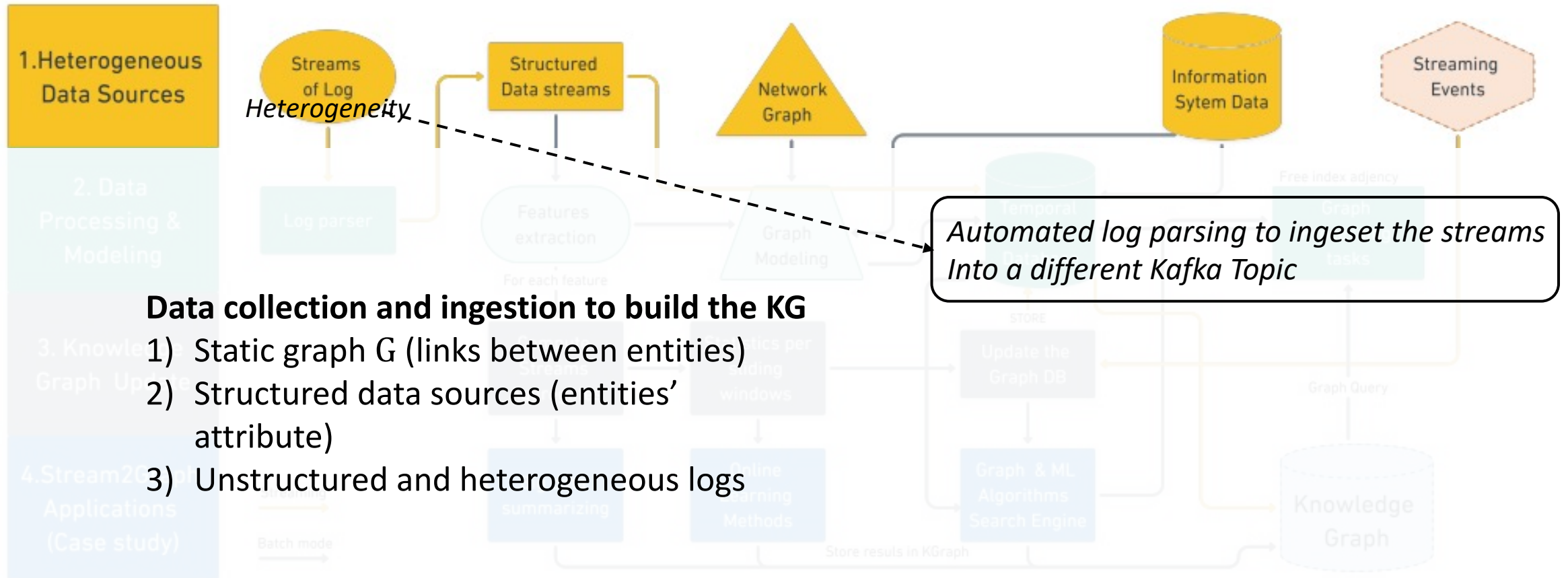3. Data pipelines for updating and maintaing the KG



* Heterogeneous graph

# Problem definition: Graph

- Input : $G = (V, E)$
  $S = \{e_1, e_2, ...\}$

- Output : $KG = (\{G^t\}_{t=1}^{T}, E')$

$E'$ : temporal relations between $G^t$ and $G^{t+1}$

$e_i \in E$

# 1. Knowledge Collcetion from multiple sources

**1.Heterogeneous Data Sources**

Streams of Log

*Heterogeneity*

Structured Data streams

Network Graph

Information Sytem Data

Streaming Events

*Automated log parsing to ingeset the streams Into a different Kafka Topic*

**Data collection and ingestion to build the KG**
1) Static graph G (links between entities)
2) Structured data sources (entities' attribute)
3) Unstructured and heterogeneous logs

2. Data Processing & Modeling

3. Knowledge Graph Update

4.Stream2Graph Applications (Case study)

Log parser

Features extraction

For each feature

Graph Modeling

Streams

Sliding windows

streaming

summarizing

Batch mode

Online Learning Methods

Temporal

STORE

Update the Graph DB

Graph & ML Algorithms Search Engine

Store resuls in KGraph

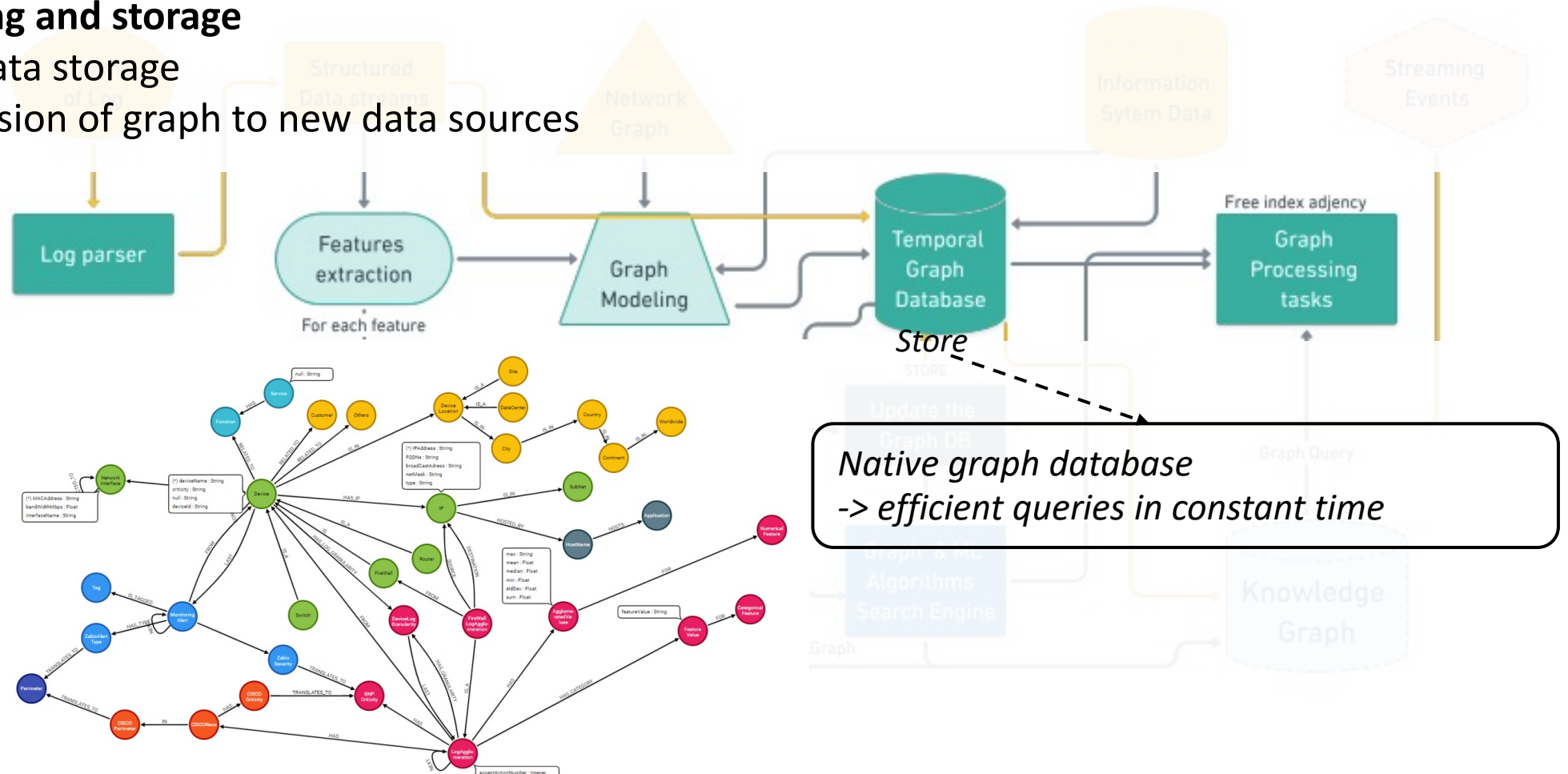Free index adjancy

Graph

tasks

Graph Query

Knowledge Graph

# 2. Knowledge modelling, enrichment, and mapping

**Data modelling and storage**

1) Scalable data storage
2) Easy extension of graph to new data sources



*Native graph database*
*-> efficient queries in constant time*

# 3. Knowledge Graph incremental update



STORE

| Compute Streams Aggregates | → | Statistics per sliding windows | → | Update the Graph DB |

Graph Query

1) Extract features from event stream
2) Compute statistics and trends for each feature type
3) Update current graph with new nodes, edges and recent statistics

*Incremental update*
*-> online learning settings*

# Online Knowledge Graph Update

- Compute feature vector for each node based on <u>categorical and numerical features</u> from raw log data
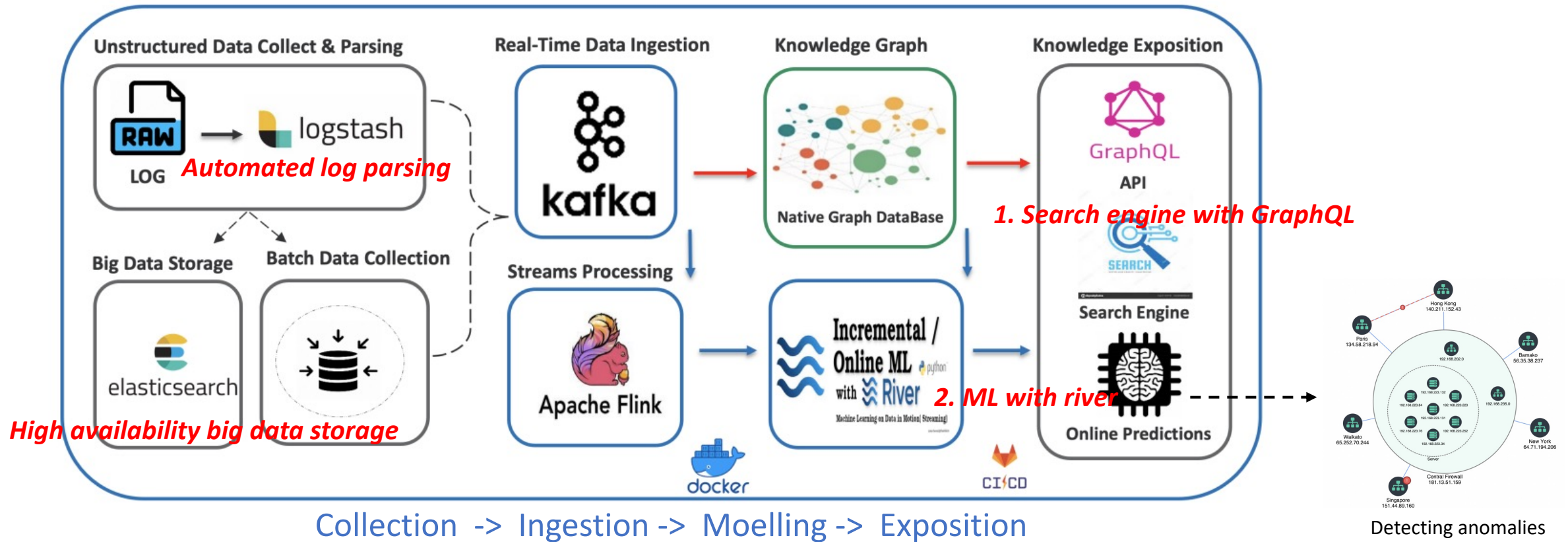


*Log patterns related to network traffic on devices*

**Algorithm 1:** Online Knowledge Graph Update

**Input:** Event stream $S = \{e_1, ..., e_n ...\}$ over time $t$
Knowledge Graph $G^t = (V^t, E^t)$, multivariate vector $r_i$
    $u_i, v_i$ source and destination node of event $e_i$
    $\eta_i$ set of numerical attributes of $e_i \in G^t$,
    $\tilde{\eta}_{ij}$ set of statistics of attribute $\eta_{ij}$,
    $\varsigma_i$ set of categorical attributes of $e_i$,
    $\tilde{\varsigma}_{ij}$ count of attribute $\varsigma_{ij}$, $r_{ij}$ value of feature $j$ in $r_i$

**Output:** Updated Knowledge Graph $\{G^t\}_{t=1}^T = (\mathcal{V}, \mathcal{E})$

1  **while** *new record* $e_i = (u_i, v_i, r_i, t_i)$ *arrives* **do**
2     $V^{t_i} \leftarrow V^t \bigcup u_i \bigcup v_i$
3     **for** $w_{ij}$ *in* $e_i$ **do**
4         **if** ISNUMERICALATTRIBUTE$(r_{ij})$
5             $\eta_i \leftarrow \eta_i \bigcup r_{ij}$
6         **else if** ISCATEGORICALATTRIBUTE$(w_{ij})$
7             $\varsigma_i \leftarrow \varsigma_i \bigcup r_{ij}$
8     ▷ Update $V$ with numerical statistics
9     **for** $\eta_{ij}$ *in* $\eta_i$ **do**
10        $\tilde{\eta}_{ij} \leftarrow$ COMPUTESTATS$(\eta_{ij})$
11        $V^{t_i} \leftarrow V^{t_i} \bigcup \eta_{ij}$
12     ▷ Update $V$ with categorical counts
13     **for** $\varsigma_{ij}$ *in* $\varsigma_i$ **do**
14        $\tilde{\varsigma}_{ij} \leftarrow$ UPDATECOUNT$(\varsigma_{ij})$
15        $V^{t_i} \leftarrow V^{t_i} \bigcup \varsigma_{ij}$
16     $E^{t_i} \leftarrow E^t \bigcup e_i$
17     $G^{t_i} \leftarrow$ GRAPH$(V^{t_i}, E^{t_i})$
18     $G^T \leftarrow G^T \bigcup G^{t_i}$
19     ▷ **Output updated graph** $G^T$

# 4. Architecture & operationalization of a dynamic knowledge graph



**Unstructured Data Collect & Parsing**
- RAW LOG → logstash
- *Automated log parsing*

**Big Data Storage**
- elasticsearch

**Batch Data Collection**

*High availability big data storage*

**Real-Time Data Ingestion**
- kafka

**Streams Processing**
- Apache Flink

**Knowledge Graph**
- Native Graph DataBase

**Incremental / Online ML** with River
Machine Learning on Data in Motion( Streaming)

*2. ML with river*

**Knowledge Exposition**
- GraphQL API

*1. Search engine with GraphQL*

- SEARCH
- Search Engine
- Online Predictions

docker    CI/CD

Detecting anomalies

Hong Kong 140.211.152.43
Paris 134.58.218.94
Bamako 56.35.38.237
192.168.202.0
192.168.223.132
192.168.235.0
192.168.223.223
Waikato 65.252.70.244
192.168.223.131
192.168.223.76
192.168.223.252
New York 64.71.194.206
192.168.223.34
Server
Central Firewall 181.13.51.159
Singapore 151.44.89.160
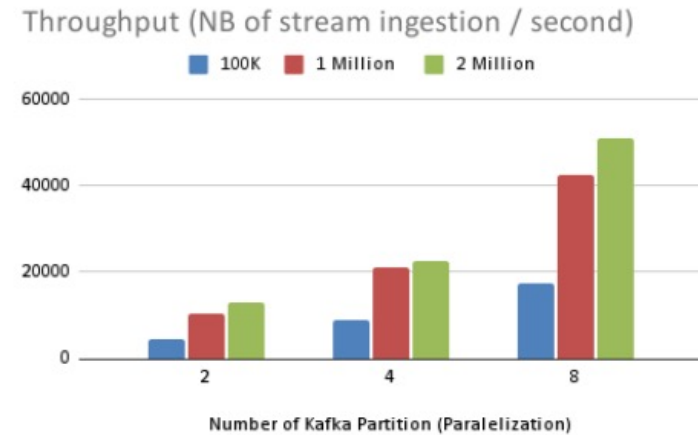
Collection  ->  Ingestion ->  Moelling ->  Exposition

# Experimental evaluation

1. Scalability in knowledge collection

-> Scale the number of instance-varying Kafka partitions from 2 to 8



(a) Knowledge ingestion Time scales linearly with the number of processing instances



(b) Throughput (Number of events/sec) scales linearly with parallelization.

# Experimental evaluation

2. Latency in knowledge ingestion and enrichment to processing billions of events

-> 22.000 events/sec



(c) High-velocity Data (22,000 events/sec) from industrial big data (21 Billions of events emitted)

# Experimental evaluation

3. Improvement of ML performance on high-velocity data streams

1. No redundant information
2. Low-dimensional while keeping relevant information
-> easier to train online ml model with time improvement (while preserving performance)

Table III: time (s) and ROC AUC for events classification on real data. Before (4 million instances) vs After (4,000 records aggregated /seconds) enriched with KG

|  | Raw network data | | Data with KG Features | | |
| --- | --- | --- | --- | --- | --- |
|  | ROC AUC | Time (s) | ROC AUC | Time (s) | Speed up |
| ARF | 0.93 | 10,623 | **0.94** | **5.19** | 2000× |
| AdaBoost | 0.98 | 3,994 | 0.98 | **7.29** | 547× |
| HT | 0.96 | 2,861 | 0.96 | **1.46** | 1900× |