

StreamFlow: A System for Summarizing and Learning Over Industrial Big Data Streams

Mariam BARRY, Saad El JAOUHARI, Albert BIFET,
Jacob MONTIEL, Eric GUERIZEC, Raja CHIKY

Contributions

StreamFlow - proposed data pipeline

Use sliding windows to summarize high-velocity data.

StreamFlow result - feature vector for machine learning tasks.

System deployed in a banking system.

Speedups in training time with good predictive performance.

Data

- network traffic data streams from BNP Paribas
- telecommunication logs

Related work

	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]	StreamFlow
Streams Aggregates					✓	✓					✓
Multiple Streams Fusion	✓	✓	✓			✓	✓		✓		✓
Categories Time-alignment					✓						
Real-world Industrial Data		✓	✓	✓	✓	✓		✓	✓	✓	✓
Industrial Big Data summurization					✓		✓	✓	✓	✓	✓
Online Machine Learning Applications		✓						✓	✓		✓

Challenges



Big Data Processing

Collect asynchronous events from multiple data sources to get a structured information

Collecting High-velocity streams & Fine-tuning Logs pipelines, windowing strategy

Dealing with resources allocation across distributed platform computing partial tasks



Big Data Summarizing

Optimizing big data clusters & partitioning resources AI/ops & IT process in Real-time

Industrialization of a online learning model which process events summaries incrementally

Updating a model with asynchronous batch & streams events using big data summaries

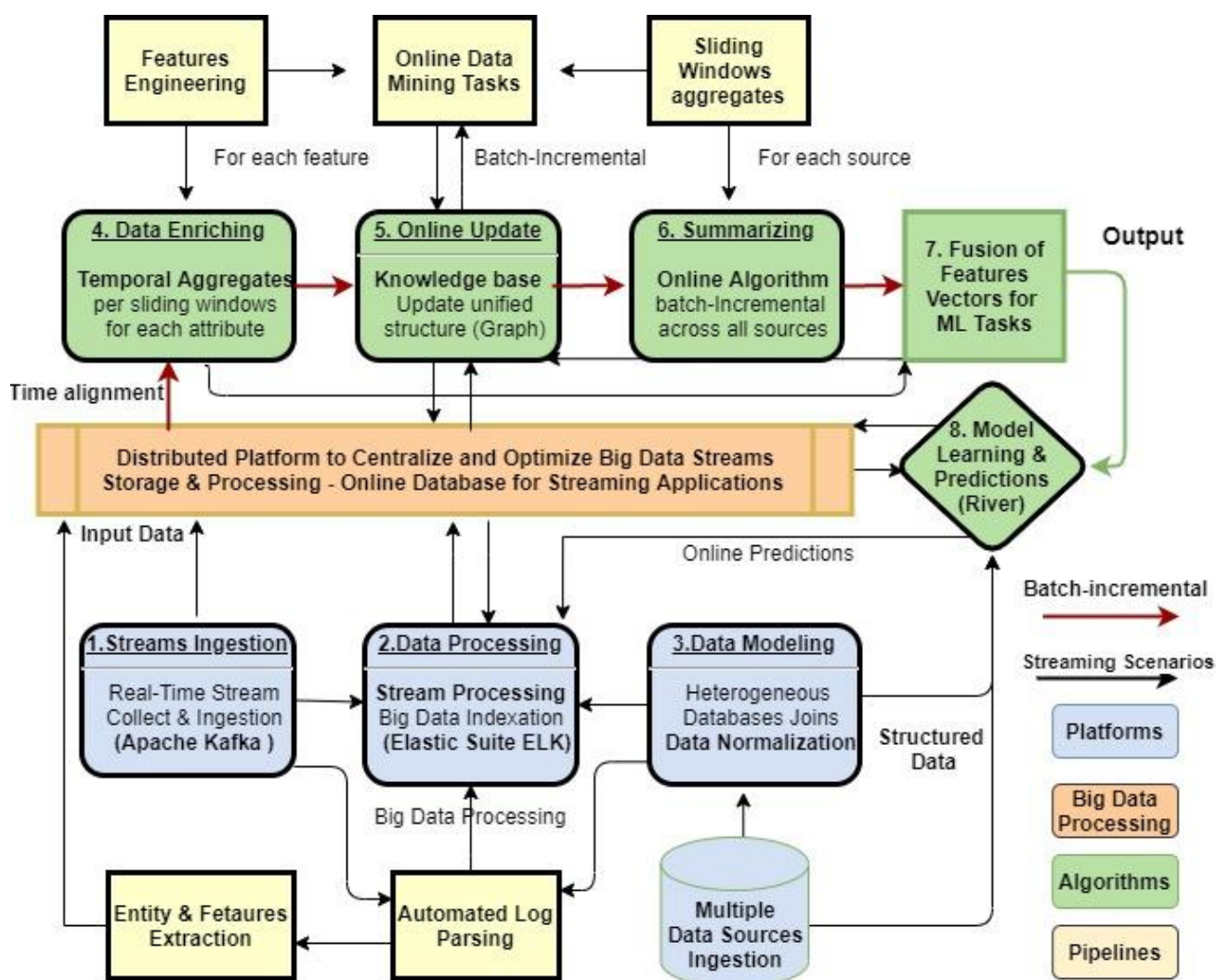


Big Data Interoperability

Building & scaling automated Data Pipelines connecting Big Data tools (ELK, Kafka, Flink..)

Serving (MLOps) online models to continuously learn from incremental vectors

Set up real-time data pipelines from streams collection, storage, modeling to predictions



Elasticsearch overview



Last 15 minutes

Show dates



[Overview](#) [Nodes](#) [Indices](#) [Machine learning jobs](#) [CCR](#)

Status

● Green

Nodes

21

Indices

235

JVM Heap

46.4 GB / 84.0 GB

Total shards

928

Unassigned shards

0

Documents

50,129,460,773

Data

21.5 TB

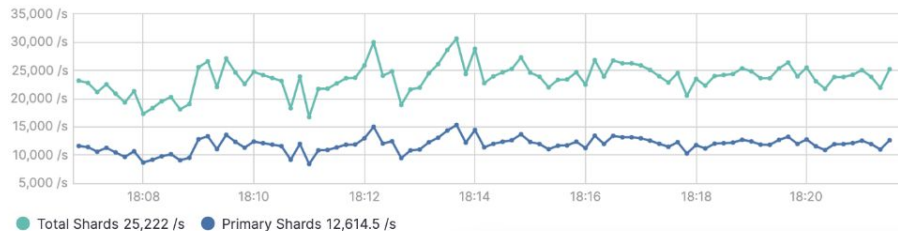
Search Rate (/s) ⓘ



Search Latency (ms) ⓘ

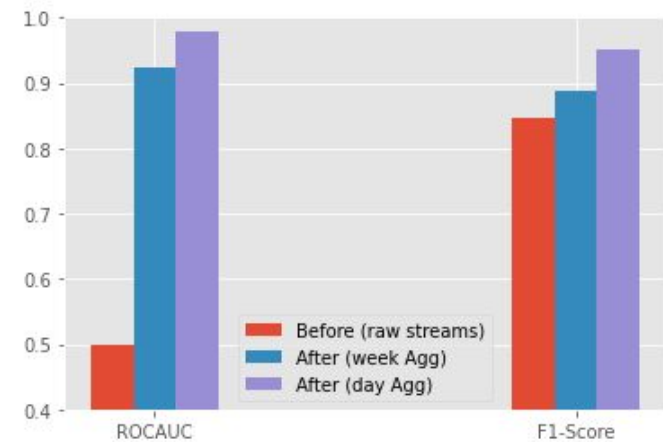


Indexing Rate (/s) ⓘ

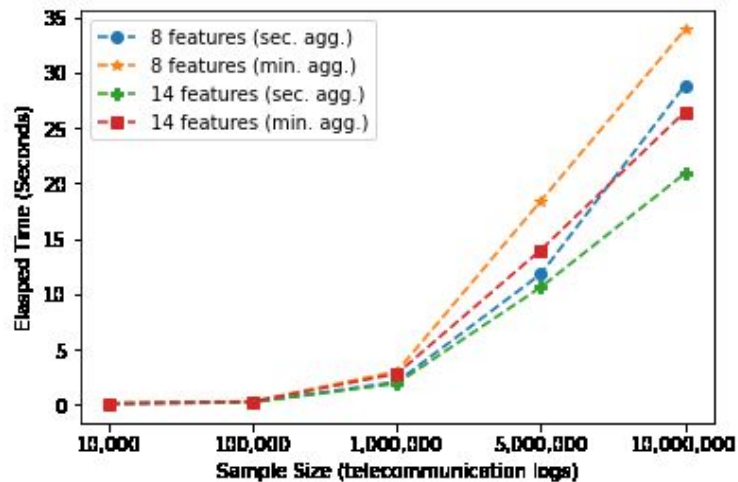


Indexing Latency (ms) ⓘ

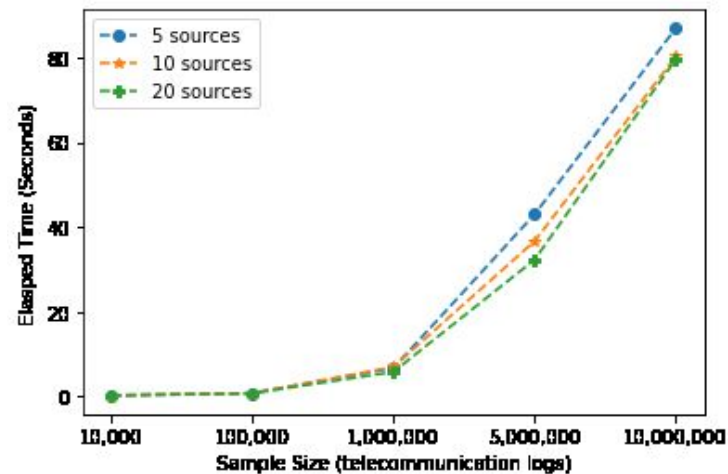




Predictive performance improvement of linear regression model.



Impact of number of features on time.



Impact of number of sources on time.

	Linear Regression			CART / Decision Trees			Random Forest			ARF	HT
	ROCAUC	F1	Time (s)	ROCAUC	F1	Time (s)	ROCAUC	F1	Time (s)	F1	F1
BEFORE - Raw big data[~1.5M]	0.500	84.53	1.70	0.987	99.23	4.22	0.987	99.23	245.61	93.84	96.33
AFTER - Minute (Week) [~90K]	0.922	88.86	1.18	0.981	96.59	0.30	0.983	97.31	6.73	98.23	96.17
AFTER - Second (Day) [~20K]	0.978	95.18	0.54	0.992	97.99	0.06	0.995	98.58	1.57	94.20	97.42