

# Streaming vs Batching

Conundrum, May 2020

[mark.andreev@gmail.com](mailto:mark.andreev@gmail.com) — Mark Andreev

# Agenda

- About data processing pipeline
- Batching approach
- Streaming approach
- Architecture
- Tools
- Conclusion



<https://clck.ru/VQL9W>



# Data pipeline

- Event := (timestamp, payload)
- Pipeline :=  $F(\text{Event}_1 \dots \text{Event}_n; \text{state})$

## Examples:

- Web Banner CTR
- Forecast Ads Budget Consumption
- Fraud prevention
- Compute analytics over session windows

### Batch

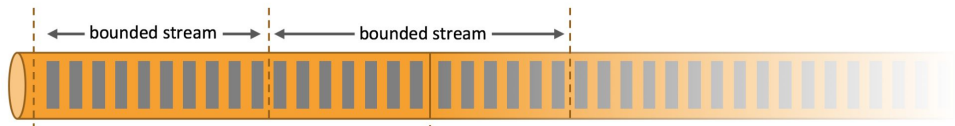
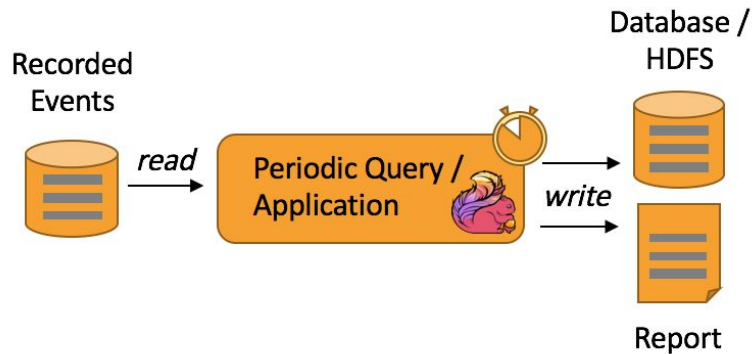
Latency - Hours

### Streaming

Latency - Minutes/Seconds

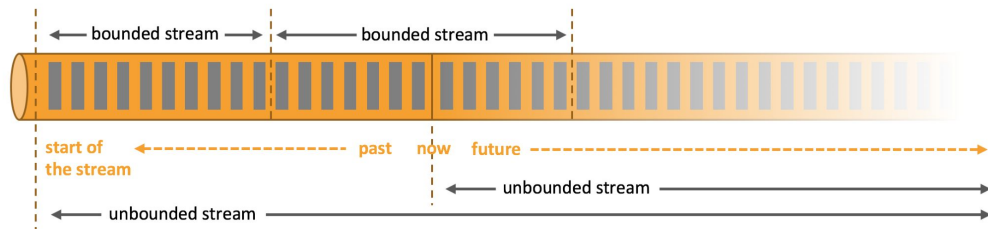
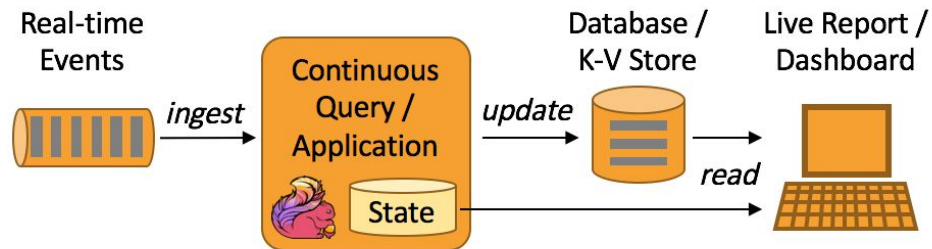
# Batching approach

- Fetch periodically
- Bounded stream approach
- Easy to write (look like SQL select/insert)
- Scheduled by standalone tool (like Airflow)

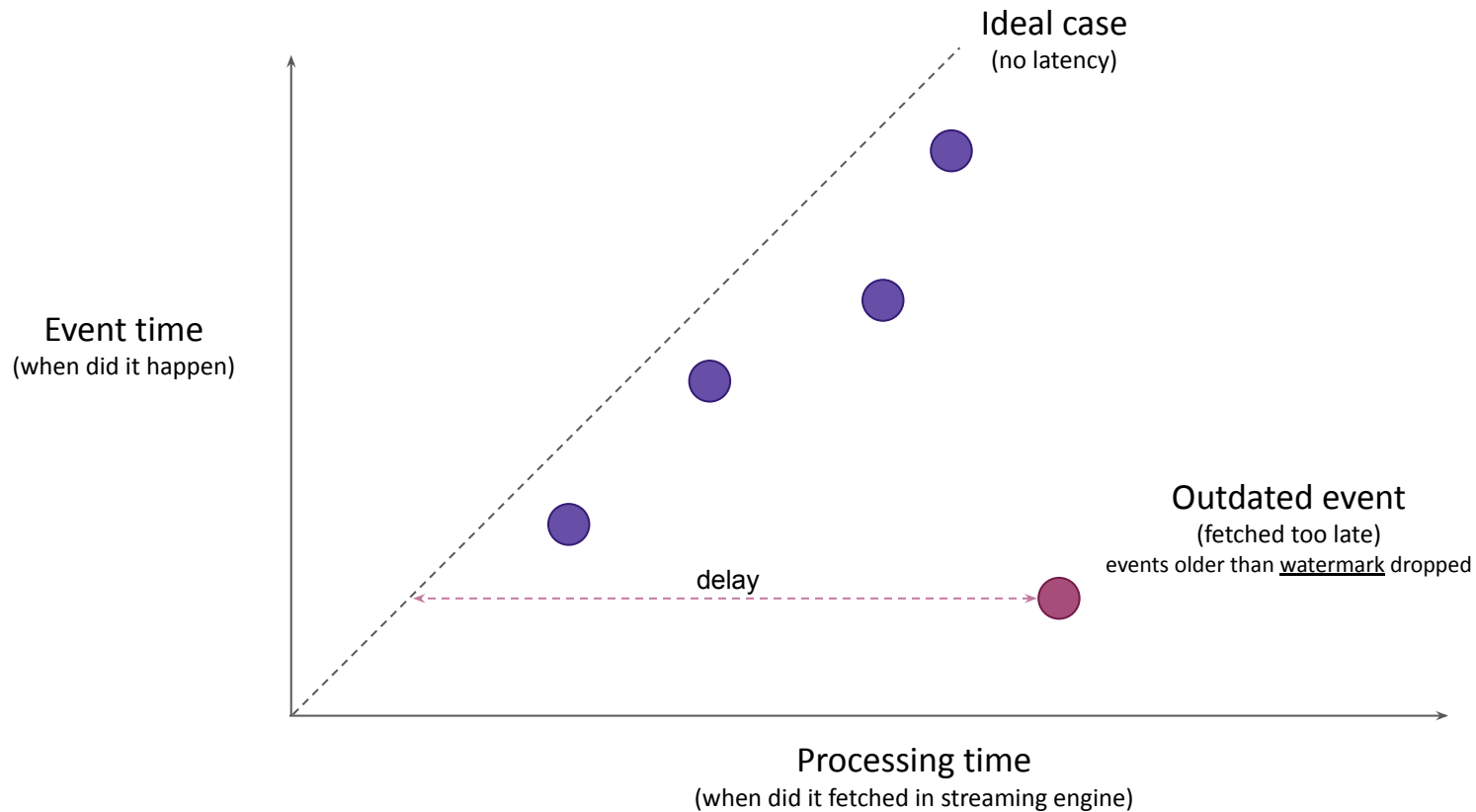


# Streaming approach

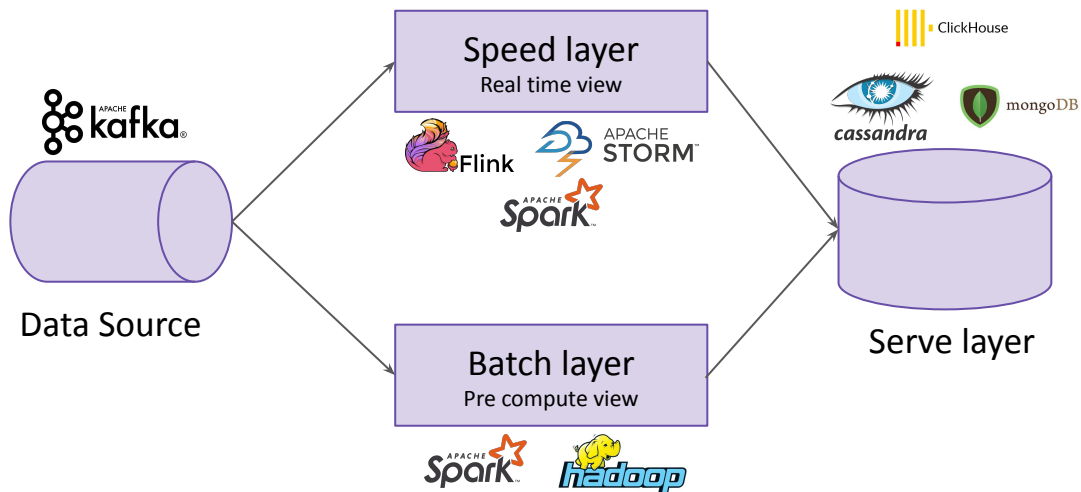
- Fetch as soon as possible
- Unbounded stream approach
- Require KV storage for state



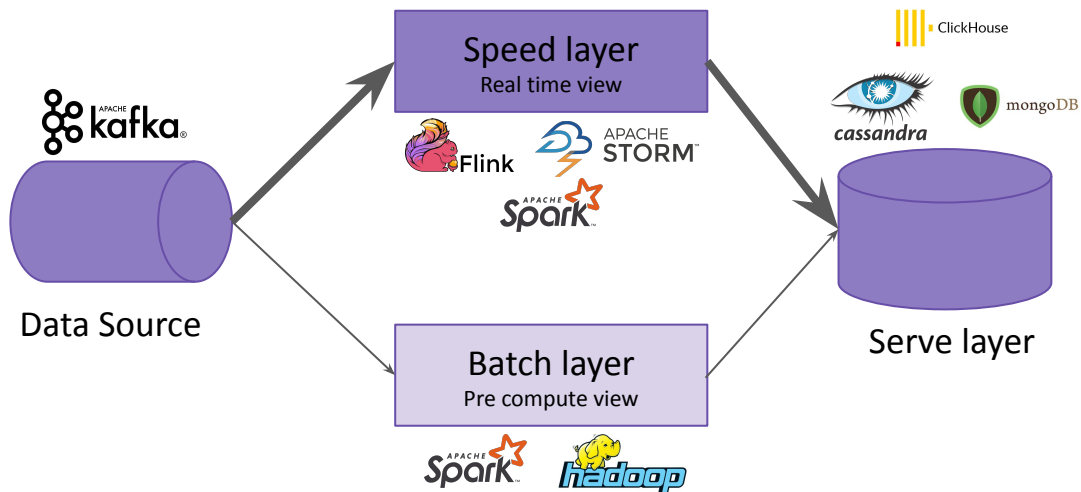
# Streaming approach [Watermarks]



# Architecture [Lambda Architecture]

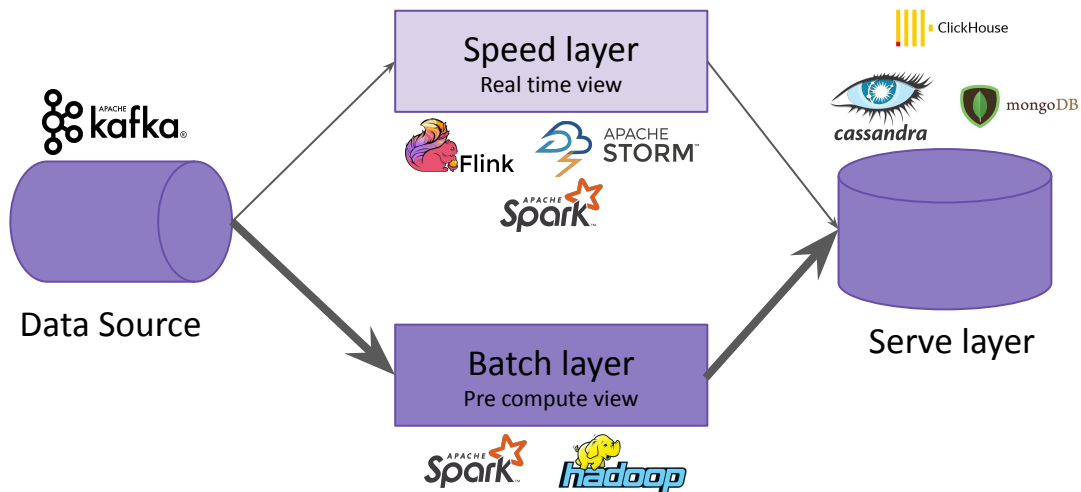


# Architecture [Lambda Architecture]

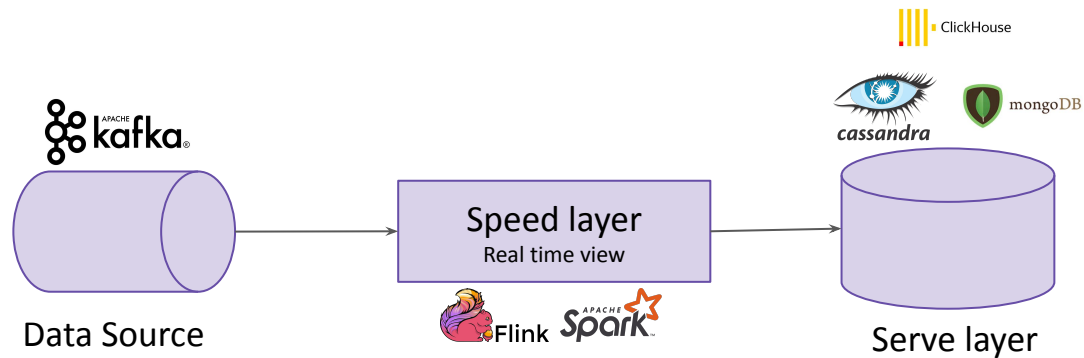




# Architecture [Lambda Architecture]



# Architecture [Kappa Architecture]



# Architecture [Lambda vs Kappa]

Lambda	Kappa
<b>Batch + Streaming</b>	<b>Streaming</b>
<b>Two scripts for both approach</b>	<b>Single script</b>
<b>Query all data</b>	<b>Incremental algorithms on deltas</b>
<b>Batch is reliable Streaming is approximate</b>	<b>Streaming with consistency</b>

# Architecture [Lambda vs Kappa]

Lambda	Kappa
Batch + Streaming	Streaming
Two scripts for both approach	Single script
Query all data	Incremental algorithms on deltas
Batch is reliable Streaming is approximate	Streaming with consistency

# Architecture [Lambda vs Kappa]

Lambda	Kappa
Batch + Streaming	Streaming
Two scripts for both approach	Single script
Query all data	Incremental algorithms on deltas
Batch is reliable Streaming is approximate	Streaming with consistency

# Architecture [Lambda vs Kappa]

Lambda	Kappa
Batch + Streaming	Streaming
Two scripts for both approach	Single script
Query all data	Incremental algorithms on deltas
Batch is reliable Streaming is approximate	Streaming with consistency

## References

- Flink Concepts [<https://clck.ru/VQLMU>]
- Spark Streaming [<https://clck.ru/VQLPA>]
- Streaming Systems: The What, Where, When, and How of Large-Scale Data Processing
- Stream Processing with Apache Flink: Fundamentals, Implementation, and Operation of Streaming Applications