



MSH

MySmartHome AG

Our Big Data Architecture | Adding value to IoT streaming data

Niels Humbeck (CDE)

20.03.2021 Köln

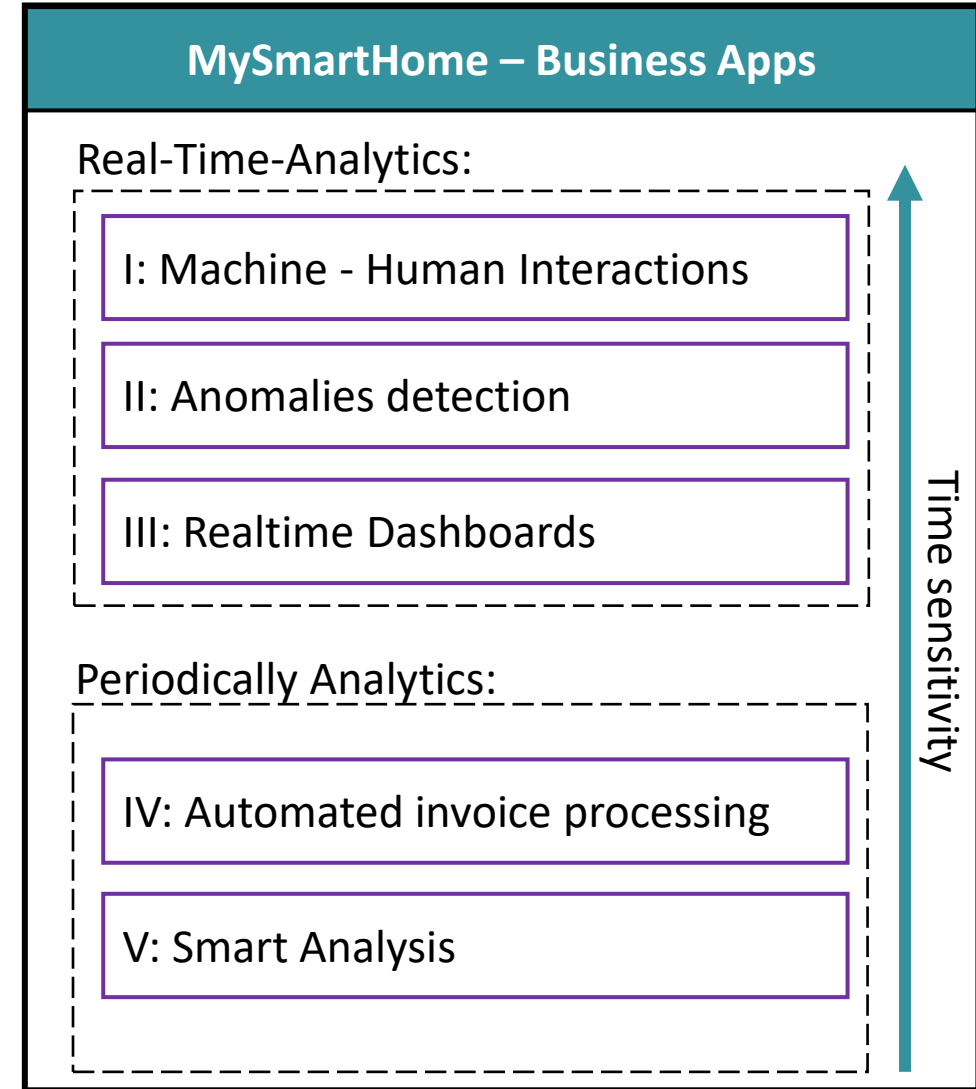
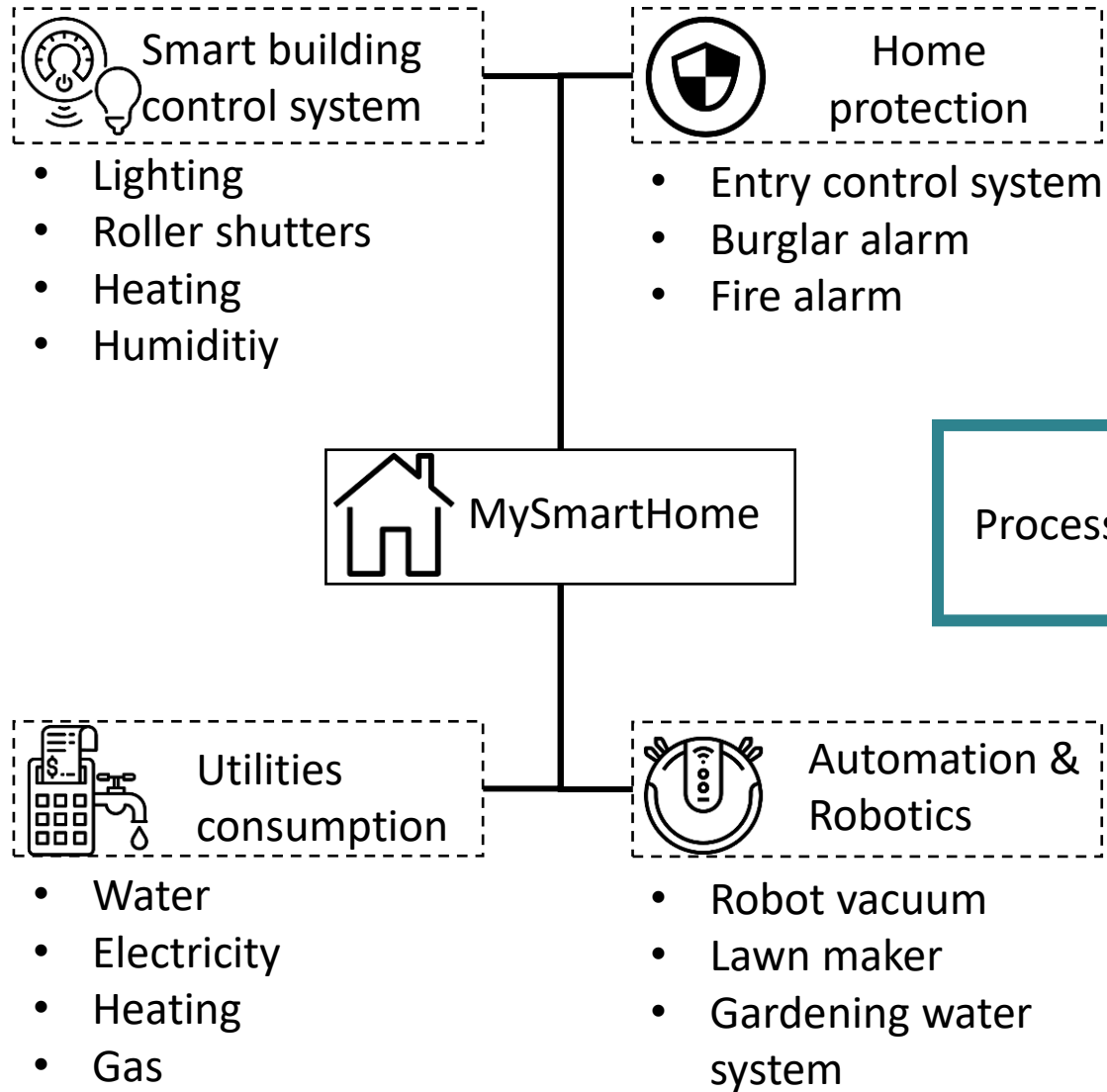
Agenda

1. Our Mission - How to add value to smart home IoT-data?
2. Core requirements of big data processing systems
3. Stream versus Batch processing of big data
4. MSH - Our big data architecture
5. Conclusion



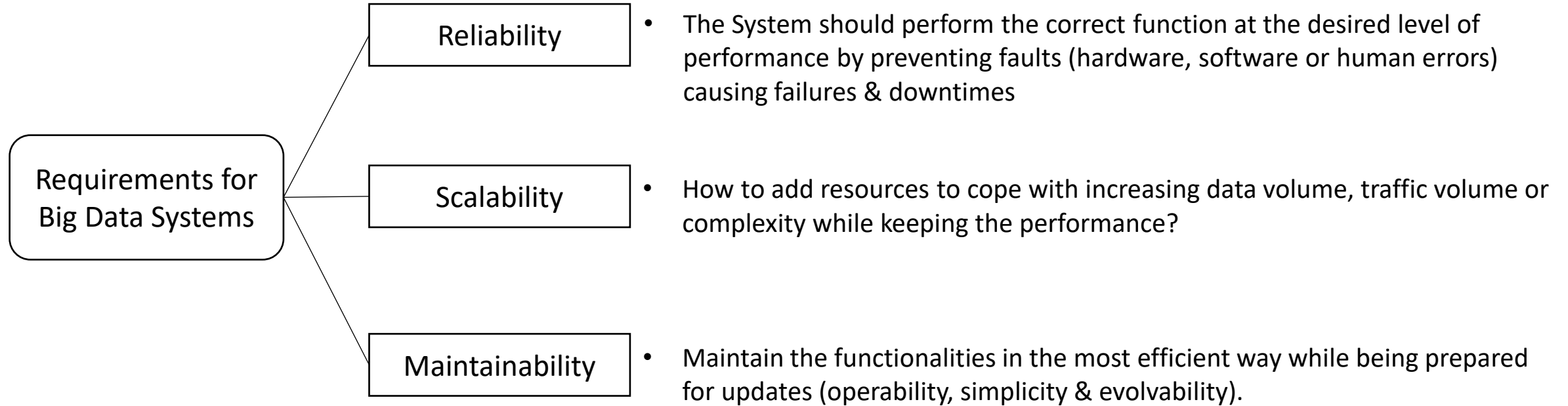
Save peoples time by minimizing their manually daily chores by the means of IoT devices and big data technology.

Our Mission - How to add value to smart home IoT-data?



Core requirements of big data processing systems

Reliability, scalability and maintainability

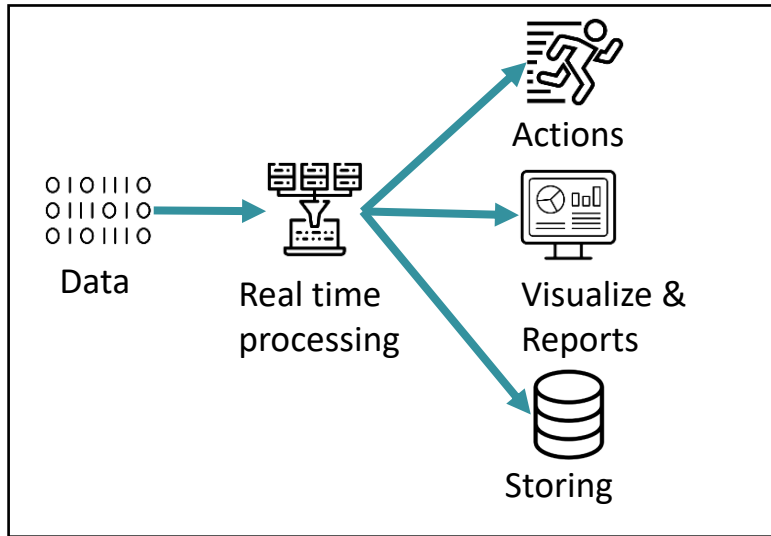


Focus on:

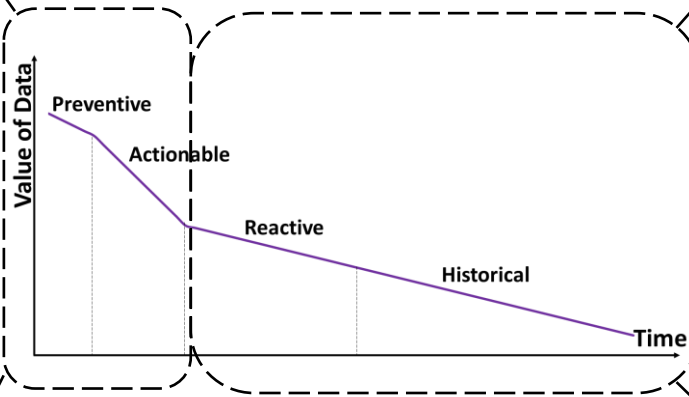
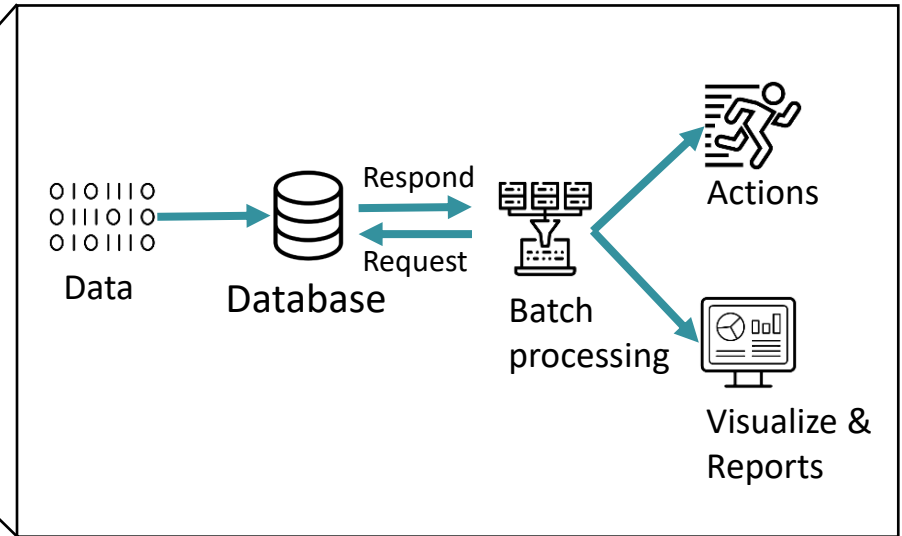
Optimized for availability, high throughput and low latency

Stream versus Batch processing of big data

Stream processing:



Batch processing:



- Low latency
- Up to date data
- Expensive
- Higher complexity
- Normally less complex analysis

Use Cases:

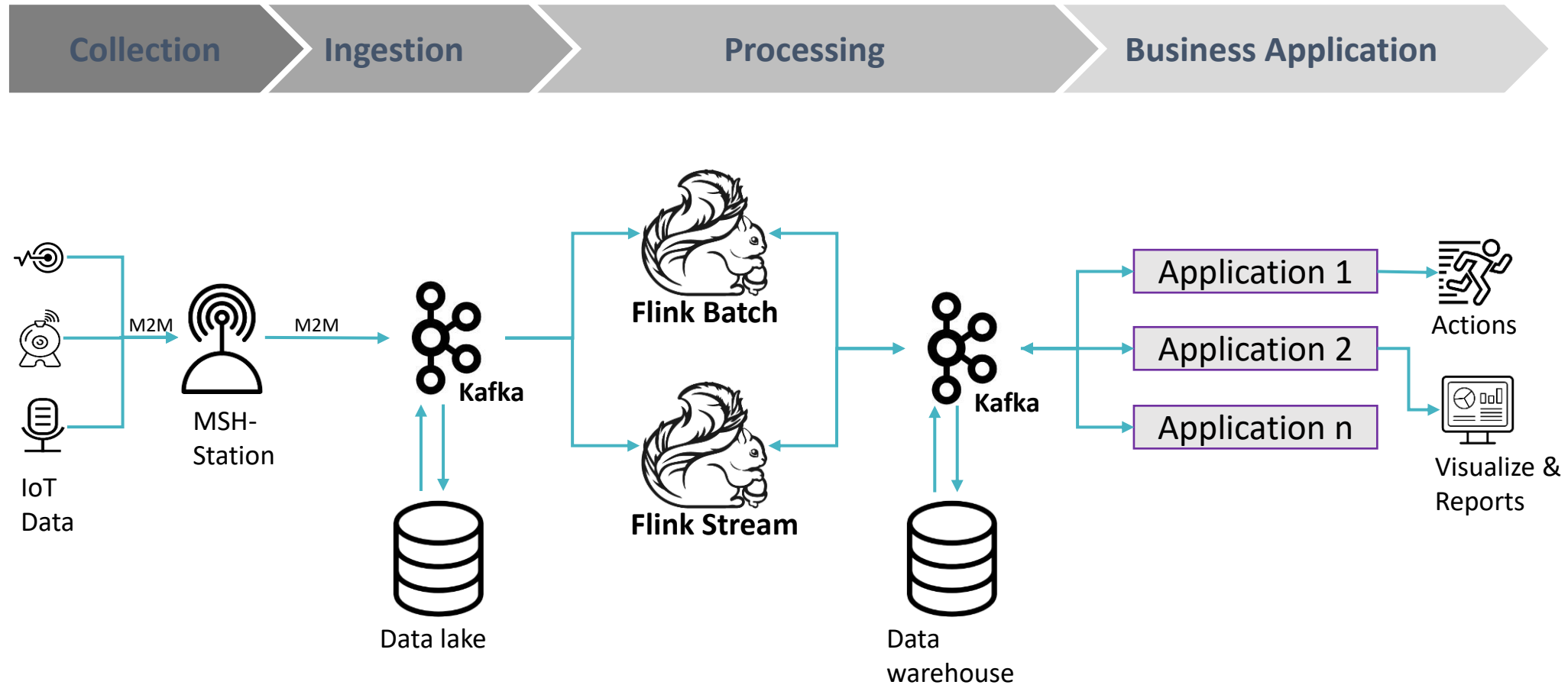
- Machine - Human Interactions
- Anomalies detection
- Realtime Dashboards

- Large batches of data
- Complex analytics & independency
- Efficient, cost effective
- High latency

Use Cases:

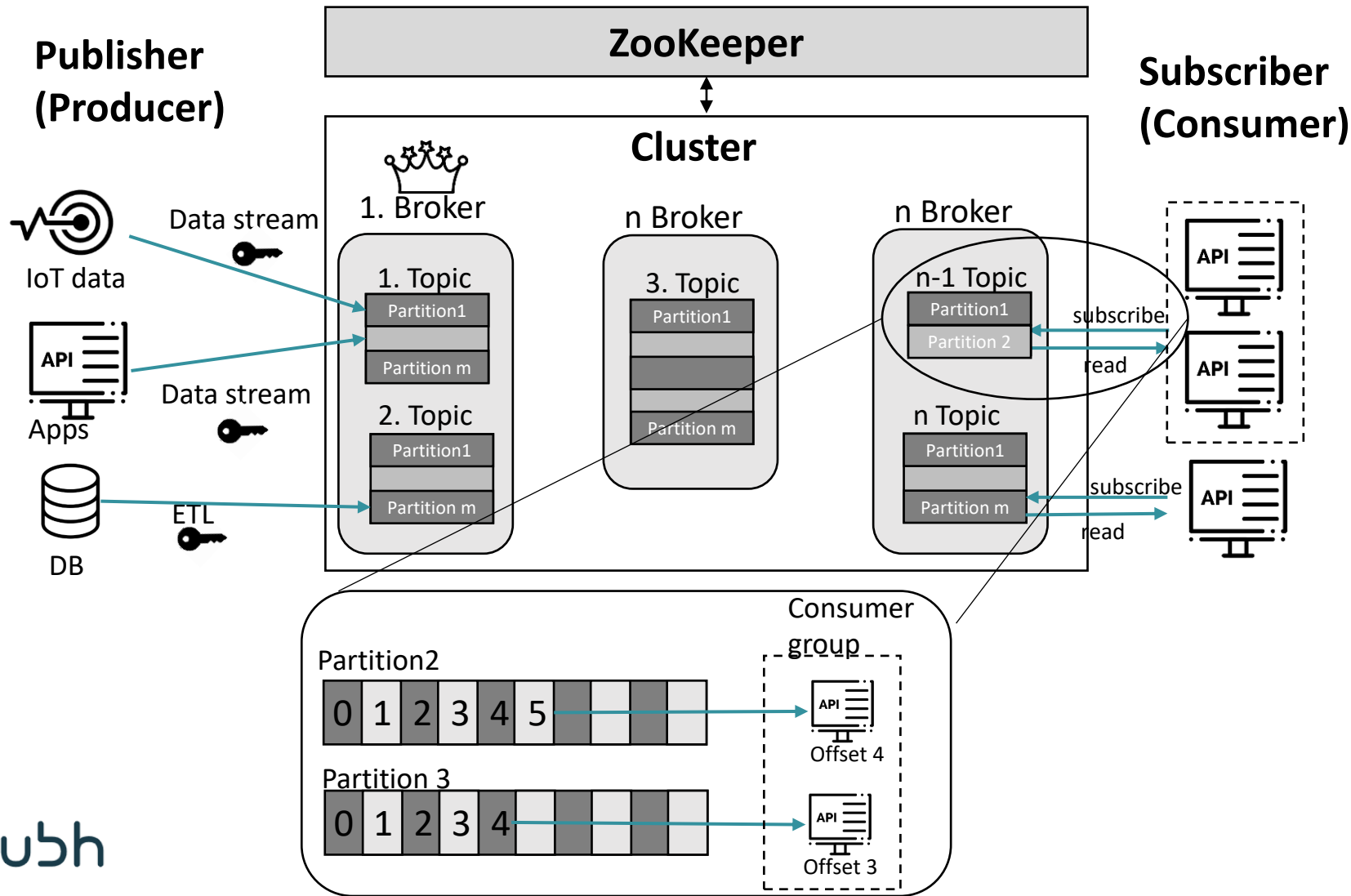
- Automated invoice processing
- Smart Analysis

MSH - Our big data architecture



MSH - Our big data architecture

Kafka | Pub/Sub messaging system with distributed immutable commit log



Reliability guaranties:

- Guaranties the order in a partition
- At-least once messages guarantee

Scalability:

- Highly scalable (horizontal)

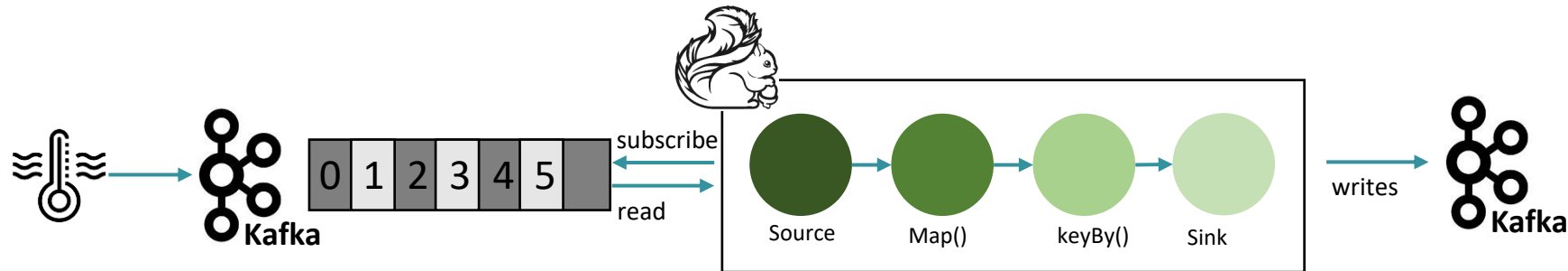
Maintainability:

- Highly configurable + retention
- Reduces integration complexity
- Fault tolerant
- Multiple producers/ consumers

- + ETL & Big Data ingestion
- High Throughput
- Fairly low latency
- Training
- Not for real low latency

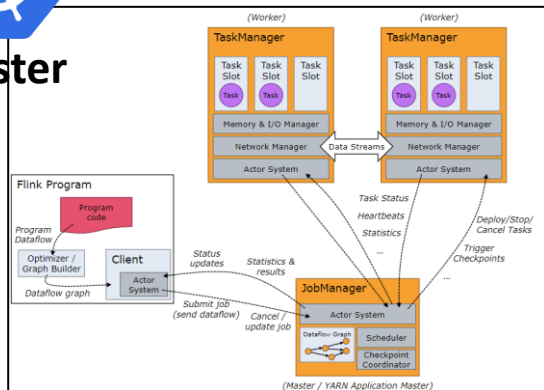
MSH - Our big data architecture

Flink | Distributed processing engine for stateful computations over (un-) & bounded data streams



Cluster

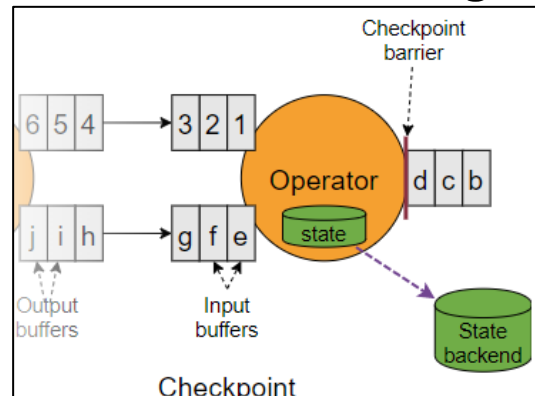
Distributed



Scalability:

- Horizontal scaling
- High throughput

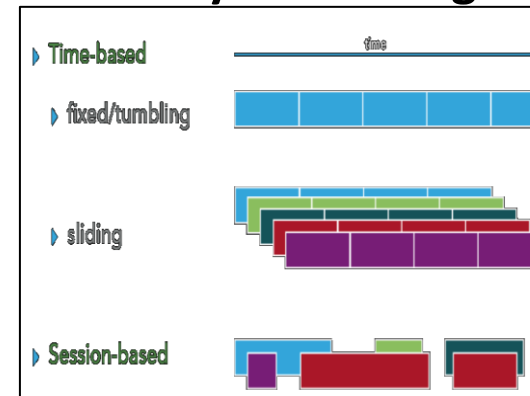
Stateful Processing



Reliability:

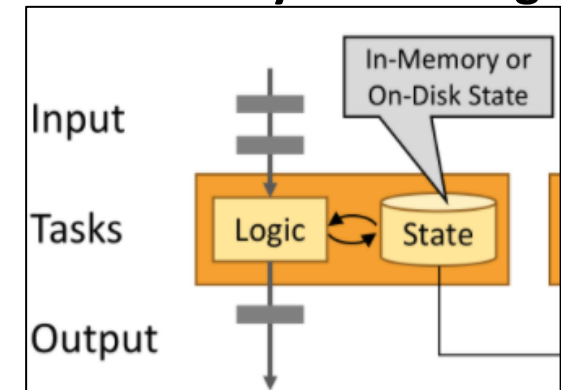
- Exactly once guarantee
- Fault tolerant

Timely Processing



- Enabling stream processing

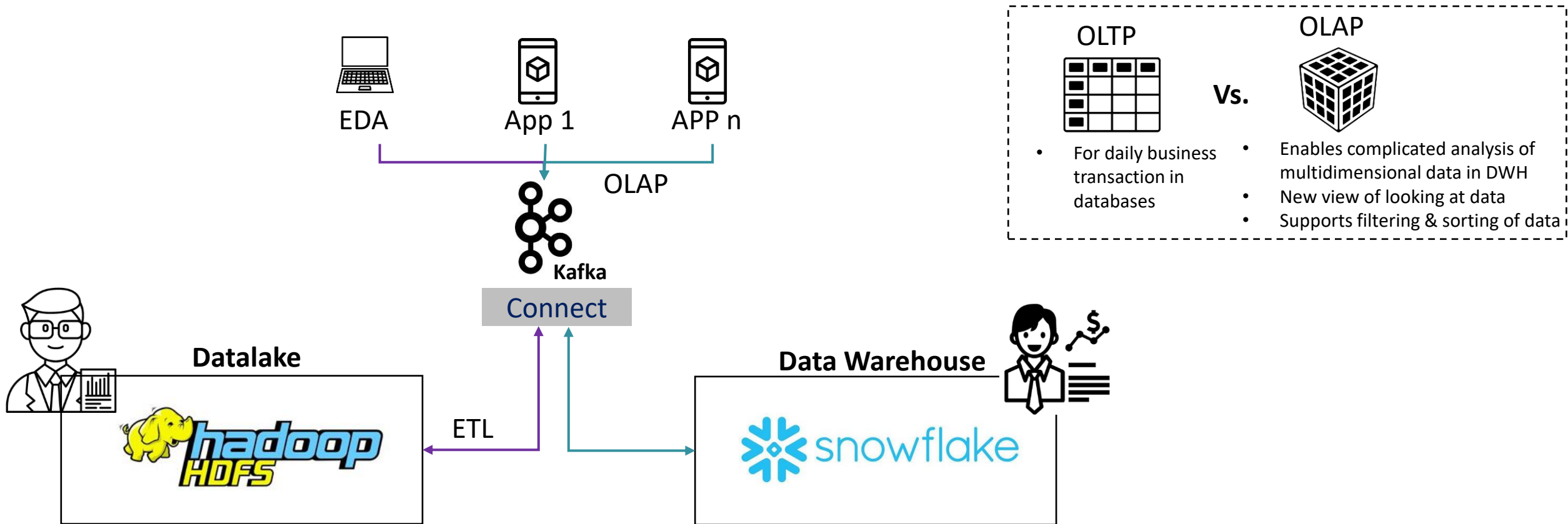
In-Memory Processing



- Low latency

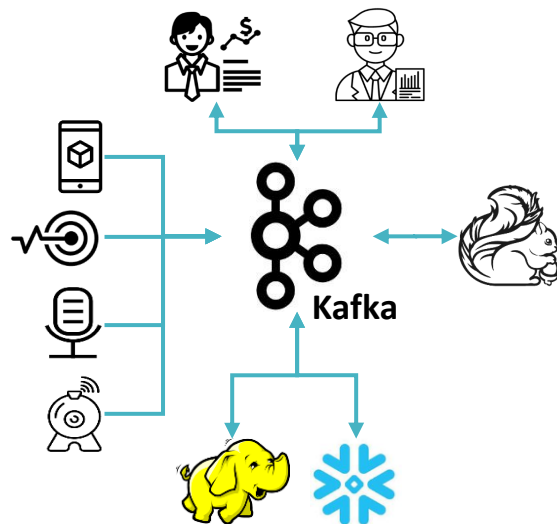
MSH - Our big data architecture

Data storage & business application layer



Conclusion

- Reliable, maintainable and scalable real time streaming pipeline
 - Kafka as our data backbone enabling multi data ingestion & low latency connectivity
 - Real time analysis
- Next generation cloud data warehouse



MySmartHome – Business Apps

Real-Time Analytics:

I: Machine - Human Interactions

II: Anomalies detection

III: Realtime Dashboards

Periodically Analytics:

IV: Automated invoice processing

V: Smart Analysis

Time sensitivity



Flink Stream



Flink Batch

Library (1/2)

- Akanbi, A., Masinde, M. (2020), A distributed Stream Processing Middleware Framework for Real-Time Analysis of Heterogeneous Data on Big Data Platform: Caes of Environmental Monitoring, Central University of Technology, South Africa
- Chaudhuri, S., Dayal, U. (n.a.), An Overview of Data Warehousing and OLAP Technology
- Cheng, C., Li, S., Ke, H. (2018), Analysis on the Status of Big Data Processing Framework, International Computers, Signals and Systems Conference
- Dendane, Y., Petrillo, F., Mcheick, H., Ben Ali, S. (2019), Quality model for evaluating and choosing a stream processing framework architecture, Universit du Qubec de Chicoutimi
- Finematics (2019), Apache Kafka Explained; <https://finematics.com/apache-kafka-explained/>, last access: 13.03.2021
- Flink 1 (2020), Flink Architecture, <https://ci.apache.org/projects/flink/flink-docs-release-1.12/concepts/flink-architecture.html>, last access 14.03.2021 at 11:31
- Flink 2(2020), Steteful Stream Processing, <https://ci.apache.org/projects/flink/flink-docs-release-1.12/concepts/stateful-stream-processing.html>, last access 14.03.2021 at 11:31
- Flink 3 (2020), What is Apache Flink?-Architecture, <https://flink.apache.org/flink-architecture.html>, last access 14.03.2021 at 11:31
- Foto 1ste Seite: <https://thenounproject.com/photo/pattern-cubes-4dEanb/>
- Gualtiri, M., Curran, R. (2016). Perishable Insights – Stop wasting money on unactionable analytics. Forrester
- Gupta, S. (2020), Architecture for High-Throughput Low-Latency Big Data Pipeline on Cloud, <https://towardsdatascience.com/scalable-efficient-big-data-analytics-machine-learning-pipeline-architecture-on-cloud-4d59efc092b5>, last access 16.03.2021 at 21:14
- Kidd, C. (2020), Data Storage Explained: Data Lake vs Warehouse vs Database, <https://www.bmc.com/blogs/data-lake-vs-data-warehouse-vs-database-whats-the-difference/>, last access 14.03.2021 at 18:24

Library (2/2)

- Kleppmann, M. (2017). Designing data intensive applications: The big ideas behind reliable, scalable, and maintainable systems. Sebastopol, CA: O'Reilly
- Knight, T. (2018), Enabling new retail experiences with Big Data, <https://www.youtube.com/watch?v=-HX-El5uhsQ>, last access 16.03.2021 as 21:01
- Marz, N., Warren, J. (2015), Big Data – Principles and best practises of scalable real time-time data systems, Manning Shelter Island
- Müller-Kett (2020); Course Book: Data Engineer – DLMDSEDE01, IUBH
- Nasiri, H., Nasehi, S., Goudarzi, M. (2019), Evaluation of distributed stream processing frameworks for IoT applications in Smart Cities, Journal of Big Data
- Patil, P. (2018), What is Explorative Data Analysis?, <https://towardsdatascience.com/exploratory-data-analysis-8fc1cb20fd15>, last access 14.03.2021 at 18:34
- Sakr, S. (2020), Big Data 2.0 Processing Systems-A Systems Overview, Springer, Institute of Computer Science, University of Tartu, Estonia
- Waski, A. (2016), Windowing data in Big Data Streams, Spark, Flink Kafka, Akka; <https://softwaremill.com/windowing-in-big-data-streams-spark-flink-kafka-akka/>, last access 14.03.2021 at 11:31

Sources | Icons (1/2)

- <https://thenounproject.com/term/smart-home/1832362/>
- <https://thenounproject.com/search/?q=control+system&i=3340124>
- <https://thenounproject.com/search/?q=Protection&i=396887>
- <https://thenounproject.com/search/?q=utilities&i=2629112>
- <https://thenounproject.com/search/?q=robotic+cleaner&i=3307576>
- <https://thenounproject.com/term/machine-code/1706949/>
- <https://thenounproject.com/term/sensor/93304/>
- <https://thenounproject.com/term/web-camera/3409539/>
- <https://thenounproject.com/term/voice/3747767/>
- <https://thenounproject.com/term/antenna/1905220/>
- <https://thenounproject.com/search/?q=database&i=2321456>
- <https://thenounproject.com/term/action/1396548/>
- <https://thenounproject.com/search/?q=visualization&i=3060300>
- <https://thenounproject.com/search/?q=key&i=2474274>
- <https://thenounproject.com/search/?q=data+scientist&i=3463591>



Sources | Icons (2/2)

- <https://thenounproject.com/search/?q=businessman&i=3346861>
- <https://thenounproject.com/search/?q=app&i=1860579>
- <https://thenounproject.com/search/?q=Laptop&i=3029683>
- <https://thenounproject.com/term/table/250445/>
- <https://thenounproject.com/term/cube/1986590/>
- <https://thenounproject.com/search/?q=chat&i=2644028>
- <https://thenounproject.com/search/?q=excel&i=3267693>
- <https://thenounproject.com/search/?q=vision&i=1852050>
- <https://thenounproject.com/search/?q=sand+clock&i=3741831>
- <https://thenounproject.com/search/?q=check&i=1438093>
- <https://thenounproject.com/search/?q=mission&i=3405804>

