



The Future of Data Pipelines

John Hammink

Developer Advocate, Evangelist @ Aiven

john@aiven.io

All Things Open 2019, Raleigh NC

The Future of Data Pipelines

- Let's start with where the **future looked** from **about a year ago**
- Let's assume that **pub-sub** and particularly technologies like **Kafka** have created a **pivot point**
- Let's follow up with a **discussion of trends** we're all seeing

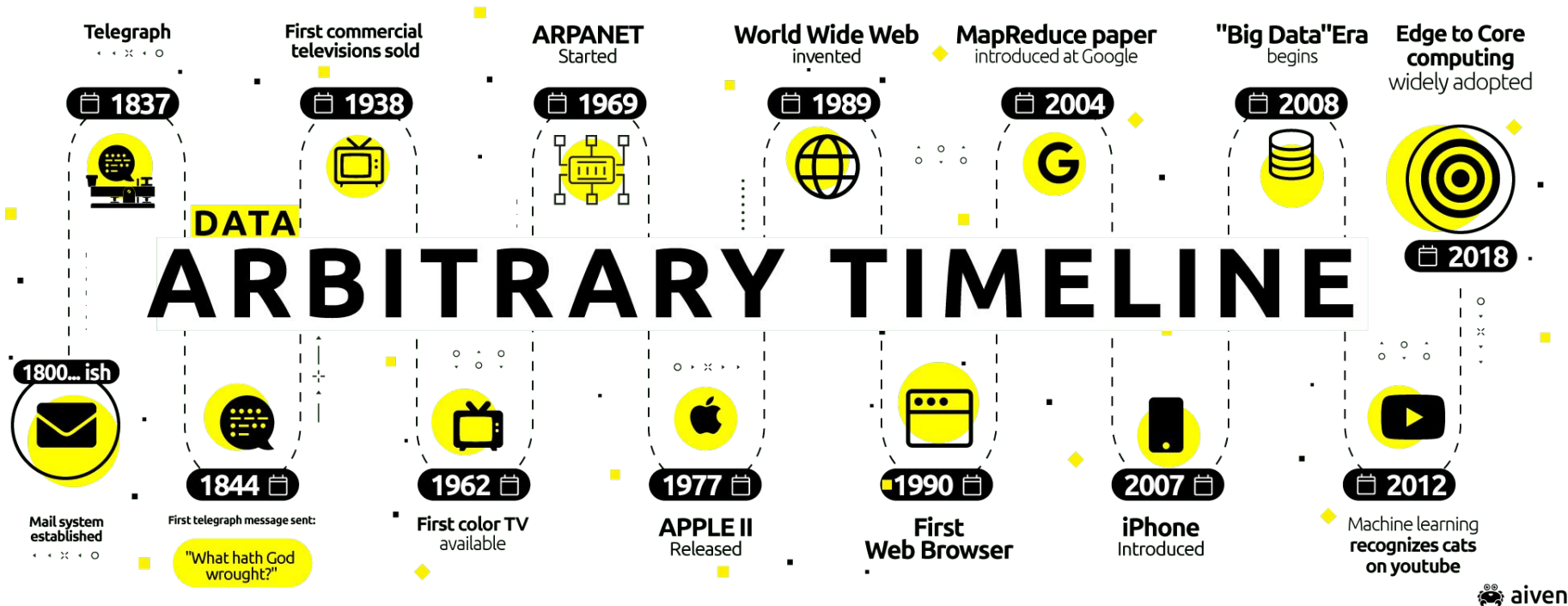


The Future of Data Pipelines

- Background
- Learnings
- Some numbers on data
- Functionality
- Design
- Compliance
- Reliability/Performance/Usability/Other
- Conclusions
- Q & A

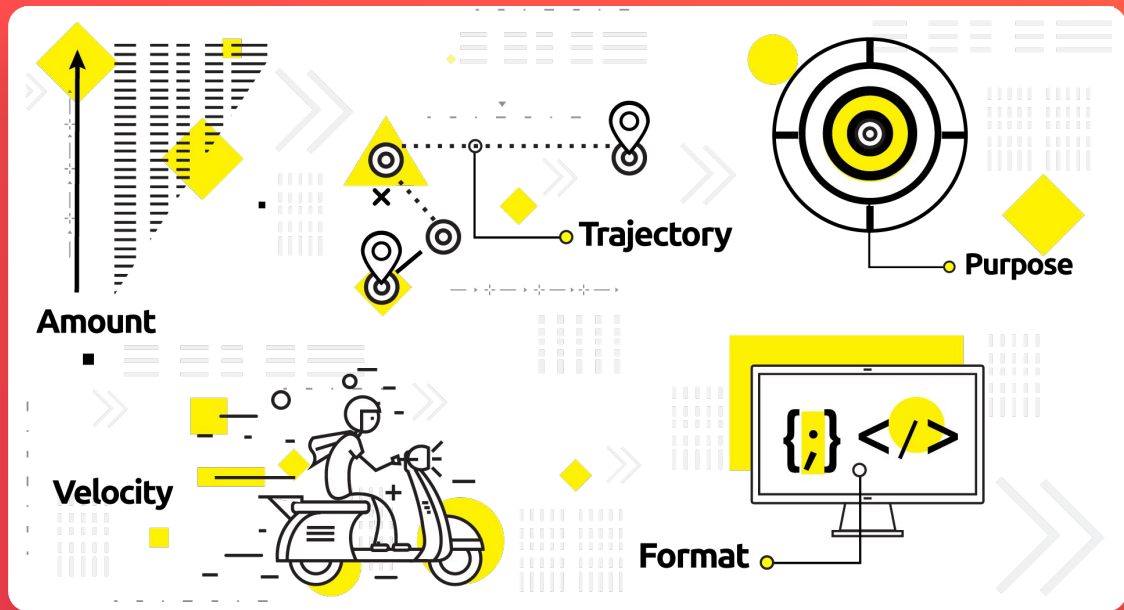


An arbitrary timeline of data



What have we learned?

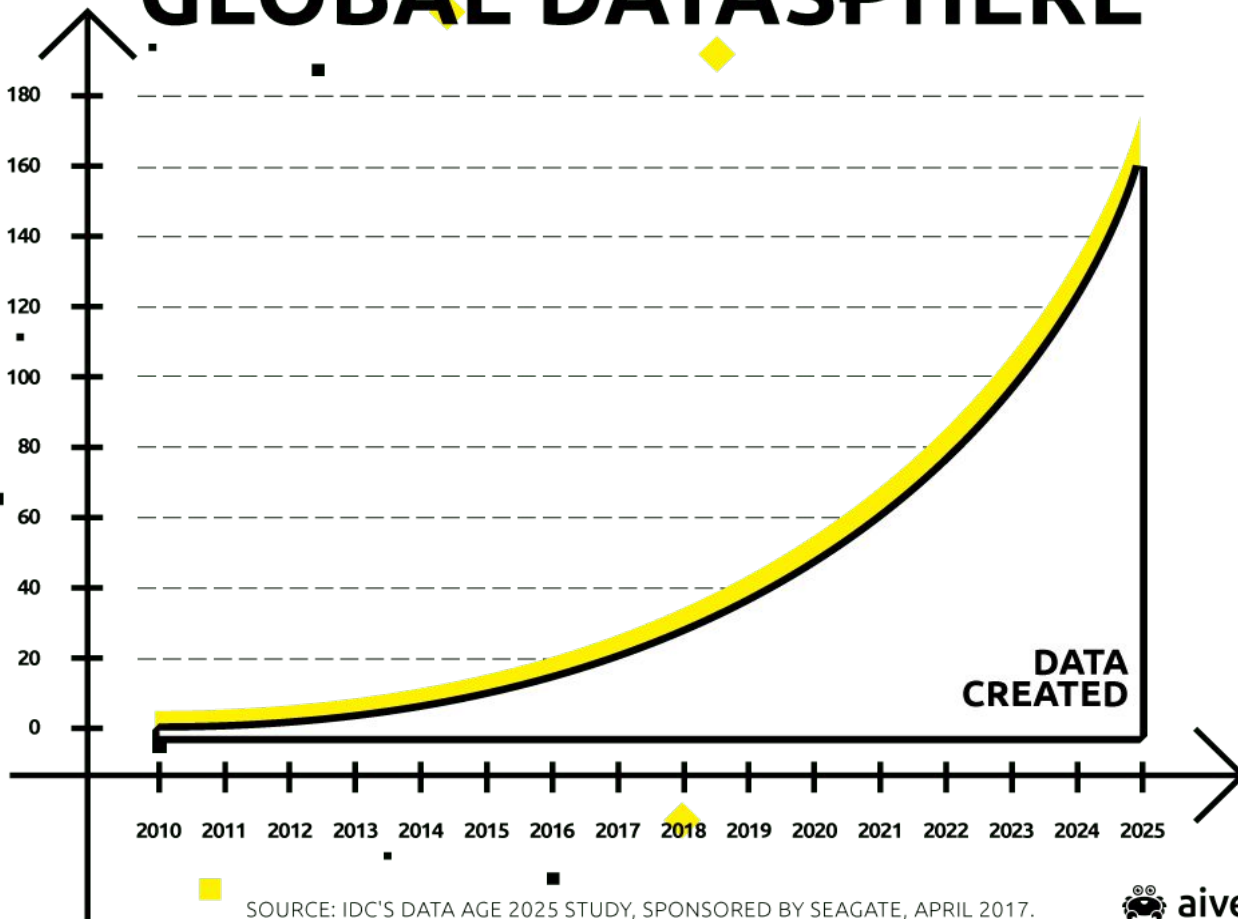
Pipeline data has
grown/transformed
in terms of:



ANNUAL SIZE OF THE

GLOBAL DATASPHERE

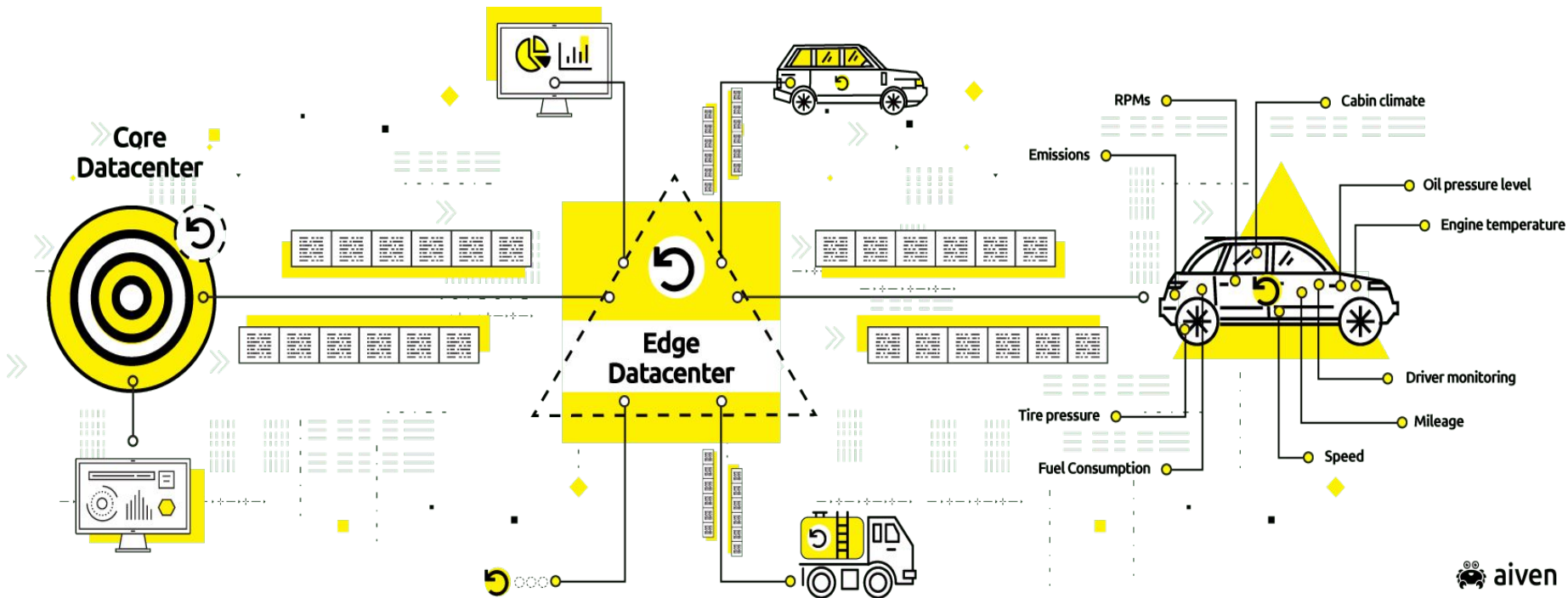
ZETTABYTES



SOURCE: IDC'S DATA AGE 2025 STUDY, SPONSORED BY SEAGATE, APRIL 2017.

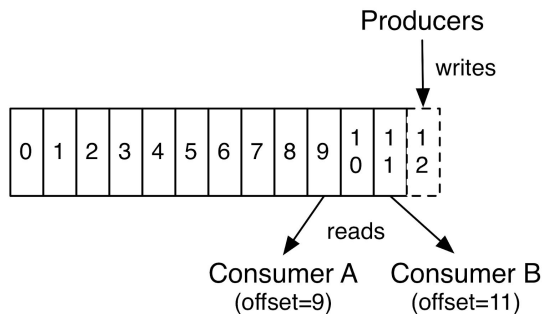


Core-to-edge computing

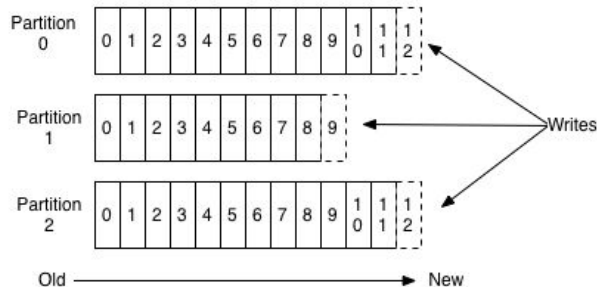


Functionality

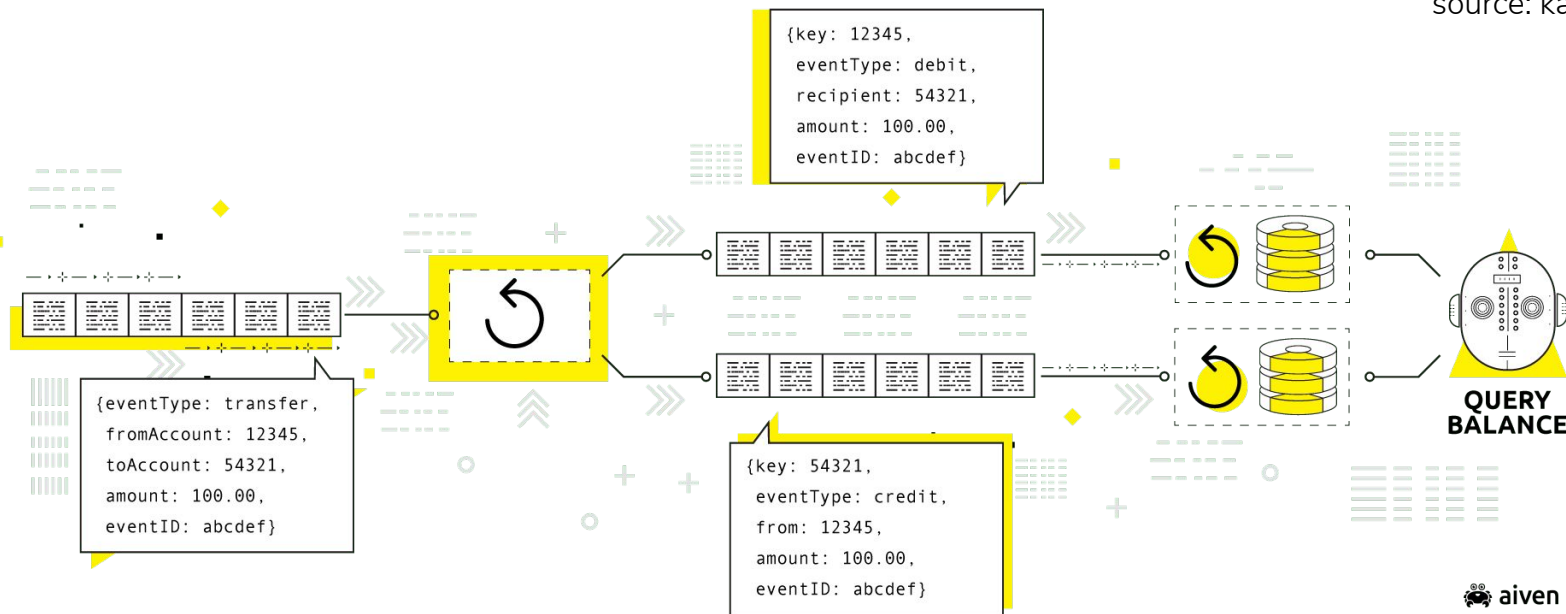
<ul style="list-style-type: none">● can autoscale, shard, tolerate partitions	<ul style="list-style-type: none">● can capture, fix and requeue errored events
<ul style="list-style-type: none">● can be troubleshot and configured on fly	<ul style="list-style-type: none">● data can be used for AI/ML
<ul style="list-style-type: none">● agnostic: accomodates all possible formats	<ul style="list-style-type: none">● can self-perpetuate and automate continuous improvements



Anatomy of a Topic



source: kafka.apache.org



Other design considerations



Kill switch - *when things go wrong; a way to stem the data flow.*



Query on the streams - like KSQL, Apache Flink, SQL on Amazon Kinesis, etc. Ultimately: one query interface for all data.

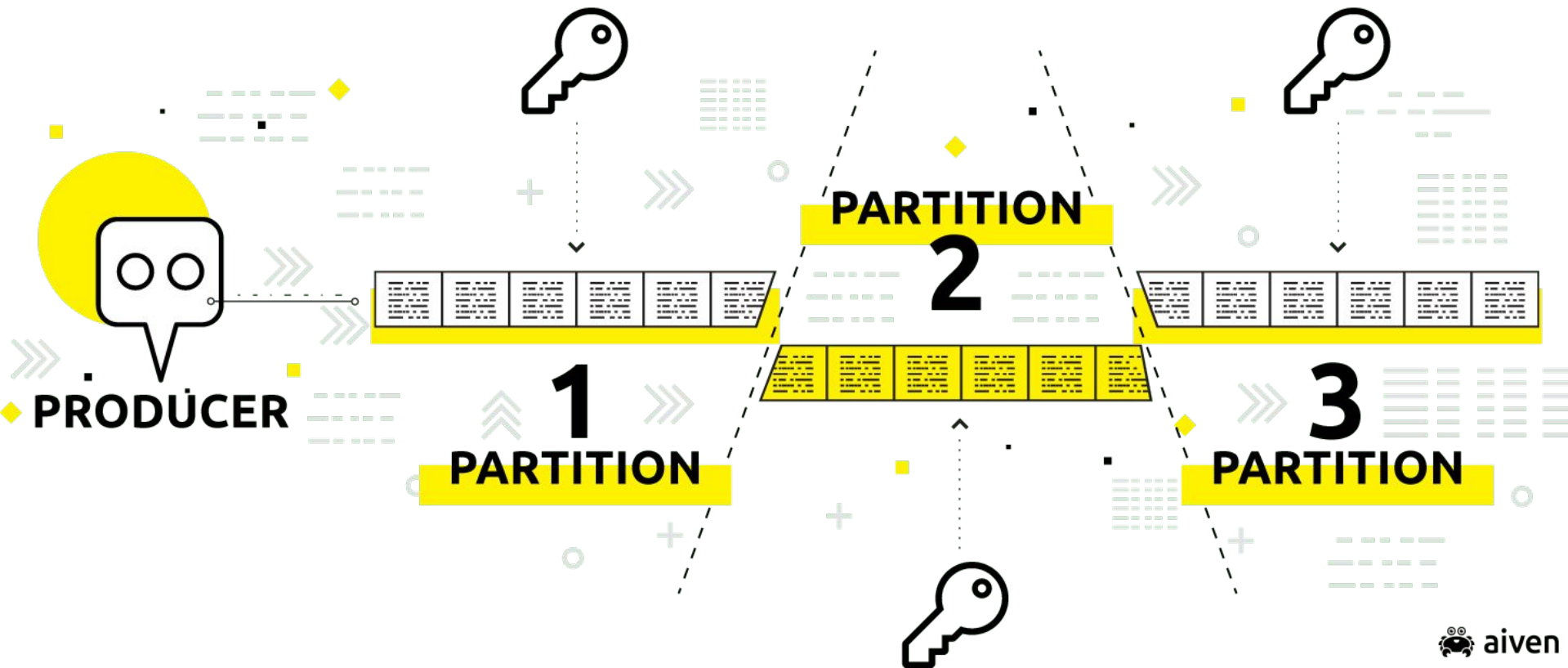


Need to accomodate - core-to-edge, growing data volume, near real-time velocity, bi-directionality



Distributed ledgers - pipelines can support these as tunably consistent distributed datastores

Compliance



Reliability/Performance/Usability/Other

- $MTTF \rightarrow \infty$



Components: from *interpreted* to *native*?



Metadata: will continue to transform, to offer even more space savings.



Pipelines: from *hard-coded* to *drag and drop* designs?



Which **data** is left *volatile* and which is *stored*?

Wrapping up

Check out **The Future of Data Pipelines** at
<https://aiven.io/blog/the-future-of-data-pipelines/>

Discussion

