

Building a **modern** data pipeline: Lessons learned



Saulius Valatka



Adform

Full stack online advertising platform

Serving banners, real-time auctions, data management, etc.



Data driven products

- Analytics
- Algorithm optimization
- Fraud detection
- Forecasting

>1,000,000 QPS

bids, impressions, clicks, requests, ...

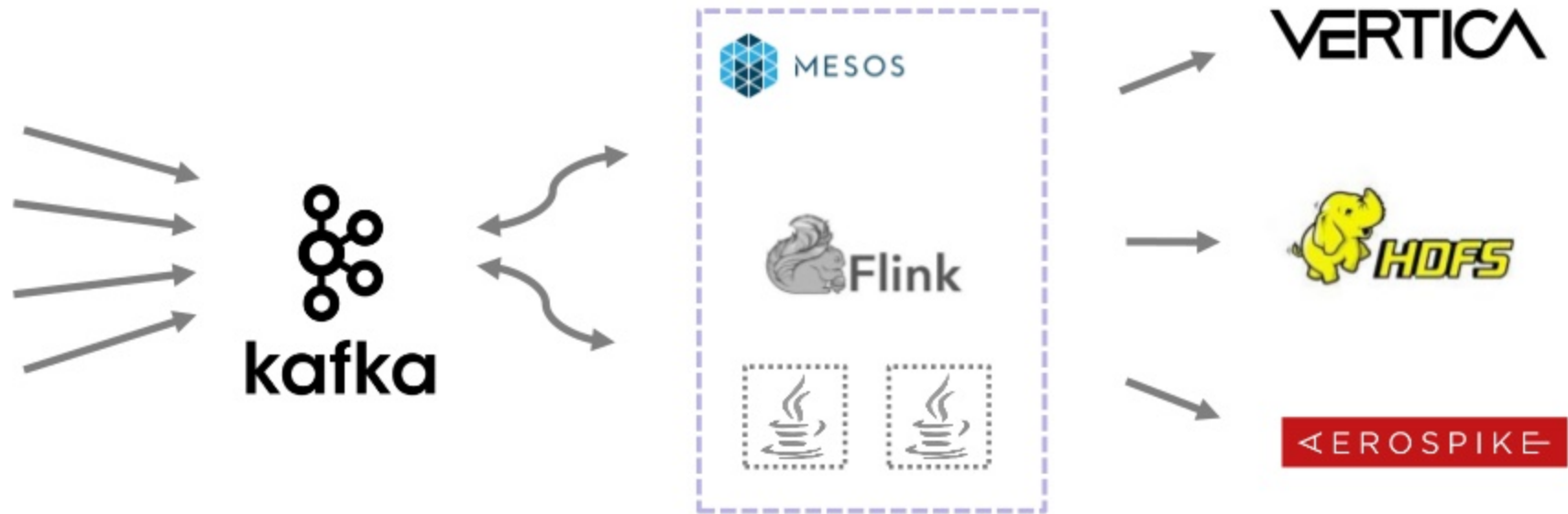
>500 MB/s

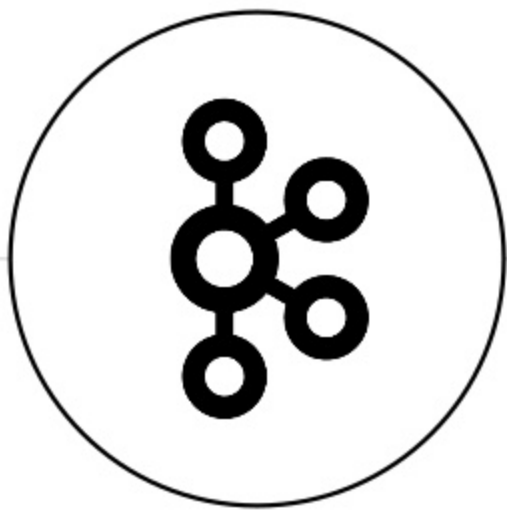
data incoming

>500 TB

data stored







Kafka

the backbone of any modern data pipeline

Topic

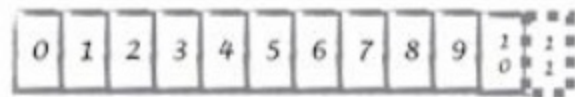
Partition 1



Partition 2



Partition 3



Writes

Old



New



Kafka Delivery Guarantees

Producing

Transactional “exactly-once”

Since Kafka 0.11



Kafka Delivery Guarantees

Producing

Transactional “exactly-once”
Since Kafka 0.11

Consuming

Consumer must implement
“atomic” offset committing

1

Long term storage loading

HDFS, Vertica and the like



Loading **projects** topics

_topic	_partition	_offset	cookie_id	url
<i>impressions</i>	1	401	123412341	http://www.reddit.com/
<i>impressions</i>	1	402	23412354	http://lwn.net/Kernel
<i>impressions</i>	2	390	234123566	http://join.adform.com/
<i>impressions</i>	2	391	64342453	http://www.ufc.com/rankings
<i>impressions</i>	1	403	341235346	http://nasza-klasa.pl/
<i>impressions</i>	3	943	465344354	http://www.facebook.com/



VERTICA



AEROSPIKE



VERTICA



AEROSPIKE

Multiple data stores

All the **benefits** of CQRS

But consistency is still **hard**



Kafka is awesome

We can achieve **real-time** data flow with **exactly-once** guarantees and have **statically-typed** data at every step of the way



Kafka is awesome

We can achieve **real-time** data flow with **exactly-once** guarantees and have **statically-typed** data at every step of the way

Kafa is our **source of truth**

2

Stream Processing

Flink, Spark, Kafka-Streams, oh my!

The world beyond batch: Streaming 101

A high level tour of modern data-processing concepts

By Tyler Akidau. August 5, 2015

The world beyond batch: Streaming 102

The what, where, when, and how of unbounded data processing.

By Tyler Akidau. January 20, 2016



Stream processing

Modern frameworks like spark (structured streaming),
flink and kafka streams allow for **exactly-once** **stateful**
data processing



Stream processing

Modern frameworks like spark (structured streaming), flink and kafka streams allow for **exactly-once** **stateful** data processing

Streaming should be **preferred** over batching

3

Batch Processing

... or what we don't want to do with streams



What do we do in batch ?

- ⦿ Algorithm training
- ⦿ Attributions with wide windows
- ⦿ Data cubing



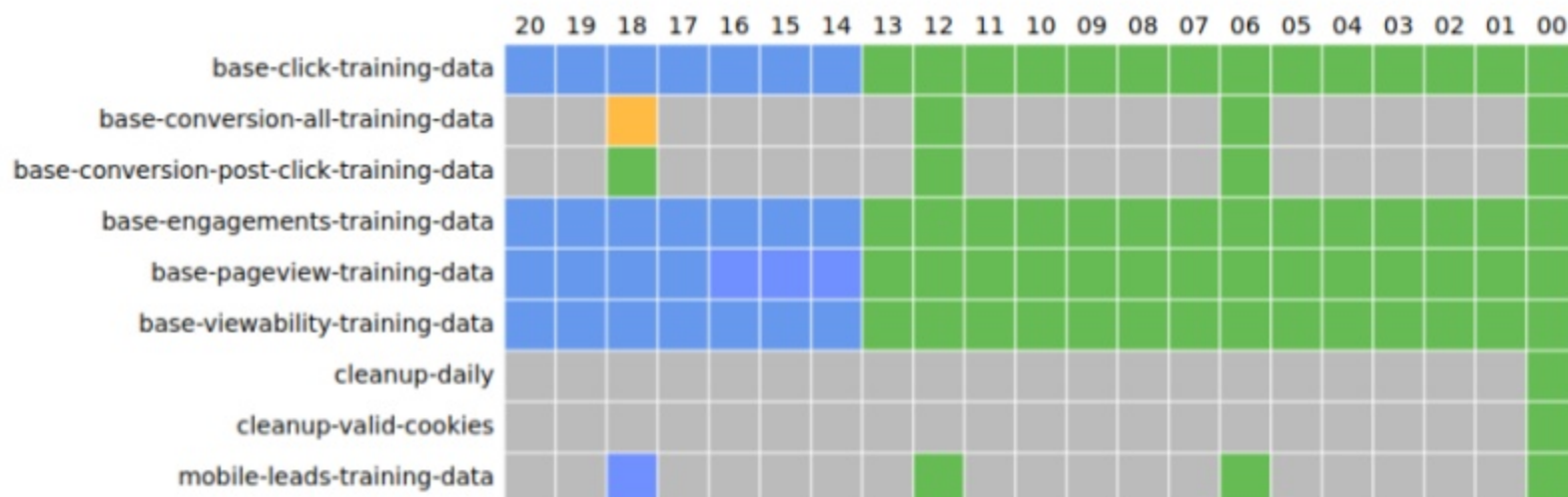
How do you even batch ?

azkaban, oozie, luigi, airflow ...

why do people keep **reinventing** this ?



Scheduling reinvented `_(\ツ)_/'





Scheduling reinvented `~_(\ツ)_/~`

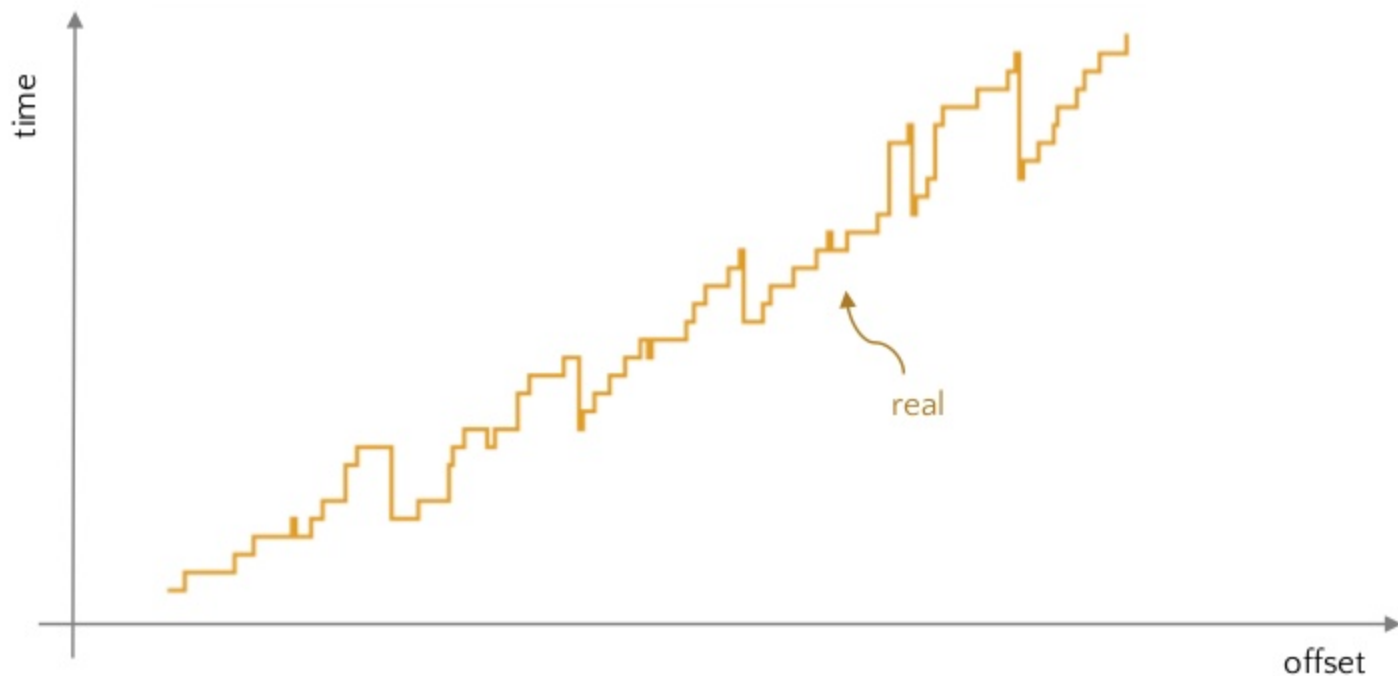
- Declarative configuration
- Runs containers (<https://github.com/adform/sprint>)
- Not (yet) open-sourced

Time based scheduling cannot **in principle** guarantee data consistency because of late data



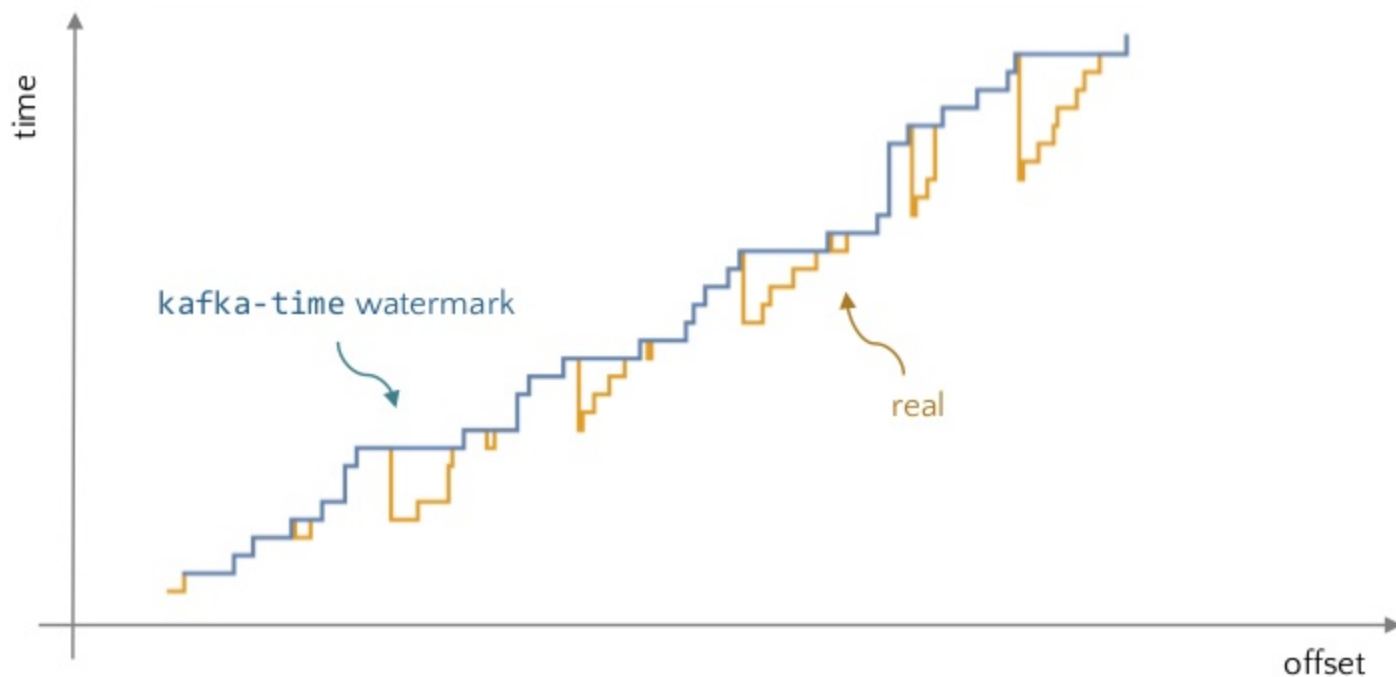


When is it safe to run ?





kafka-time to the rescue!



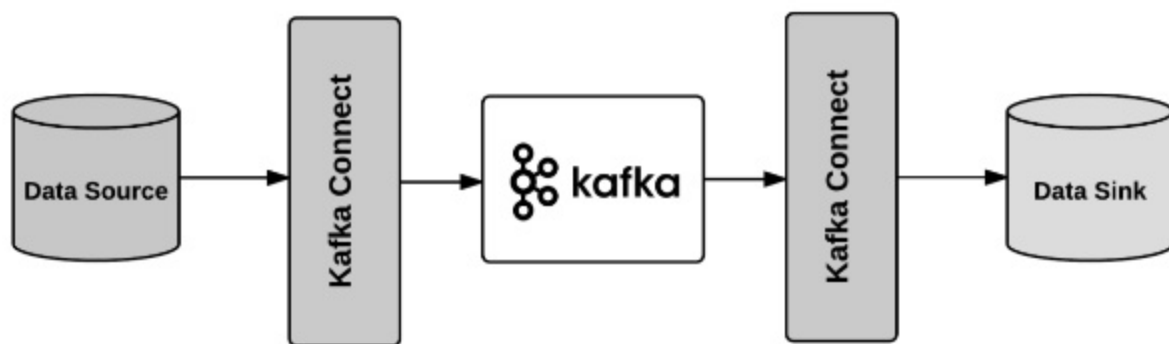
4

What about “dimensions” ?

a.k.a. metadata: clients, campaigns, etc.



Data Distribution Platform





A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?



A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?

Immutable

Reproducible, testable,
auditable, other good words



A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?

Real Time

Ain't nobody got time to wait for batches ...

Immutable

Reproducible, testable, auditable, other good words



A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?

Real Time

Ain't nobody got time to wait for batches ...

Immutable

Reproducible, testable, auditable, other good words

Performant

Seriously, ain't nobody got time to wait for queries ...



A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?

Real Time

Ain't nobody got time to wait for batches ...

Statically Typed

Yes, I said it, so sue me ...

Immutable

Reproducible, testable, auditable, other good words

Performant

Seriously, ain't nobody got time to wait for queries ...



A modern data pipeline is ...

Consistent

Otherwise, what are you even doing?

Real Time

Ain't nobody got time to wait for batches ...

Statically Typed

Yes, I said it, so sue me ...

Immutable

Reproducible, testable, auditable, other good words

Performant

Seriously, ain't nobody got time to wait for queries ...

Accessible

Clearly understood, easily used and extended



Thanks!

Any *questions* ?

You can find me at

- @saulius_vl
- github.com/sauliusvl