

# From batch to streaming to both

Herman Schaaf, Principal Software Engineer

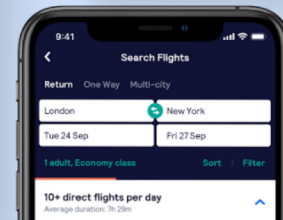
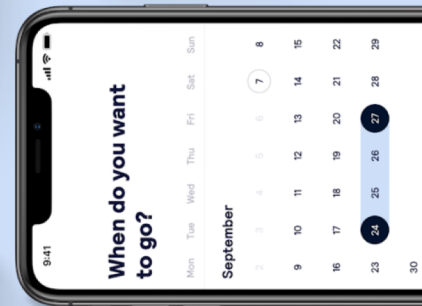
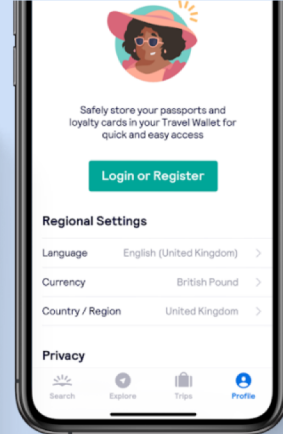
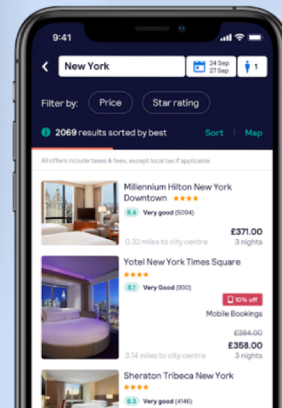
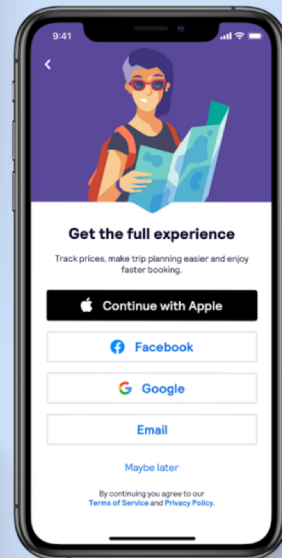
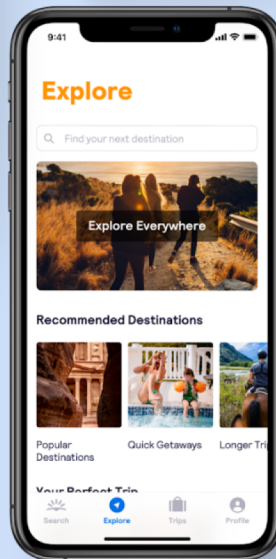
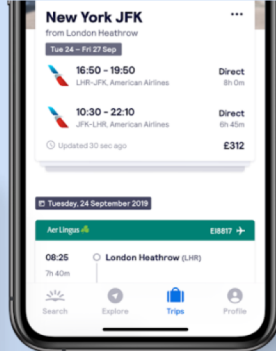
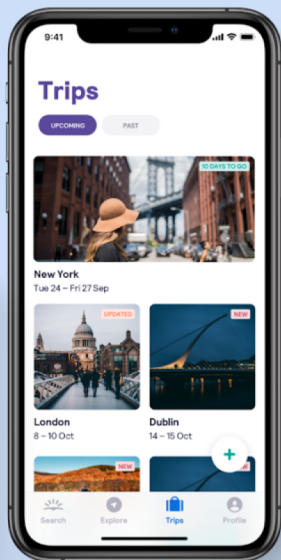
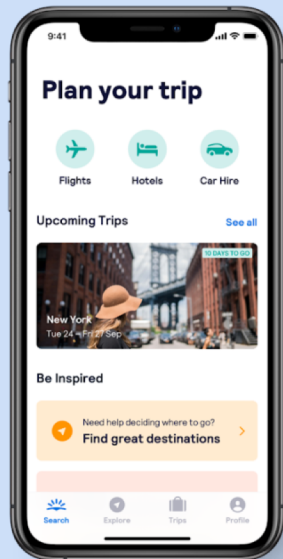
# A Story



# From batch to streaming to both *to streaming...again?*

Herman Schaaf, Principal Software Engineer

Data Platform Tribe





## Global travel company with local expertise

100

million peak  
monthly  
active users

30+

languages  
available

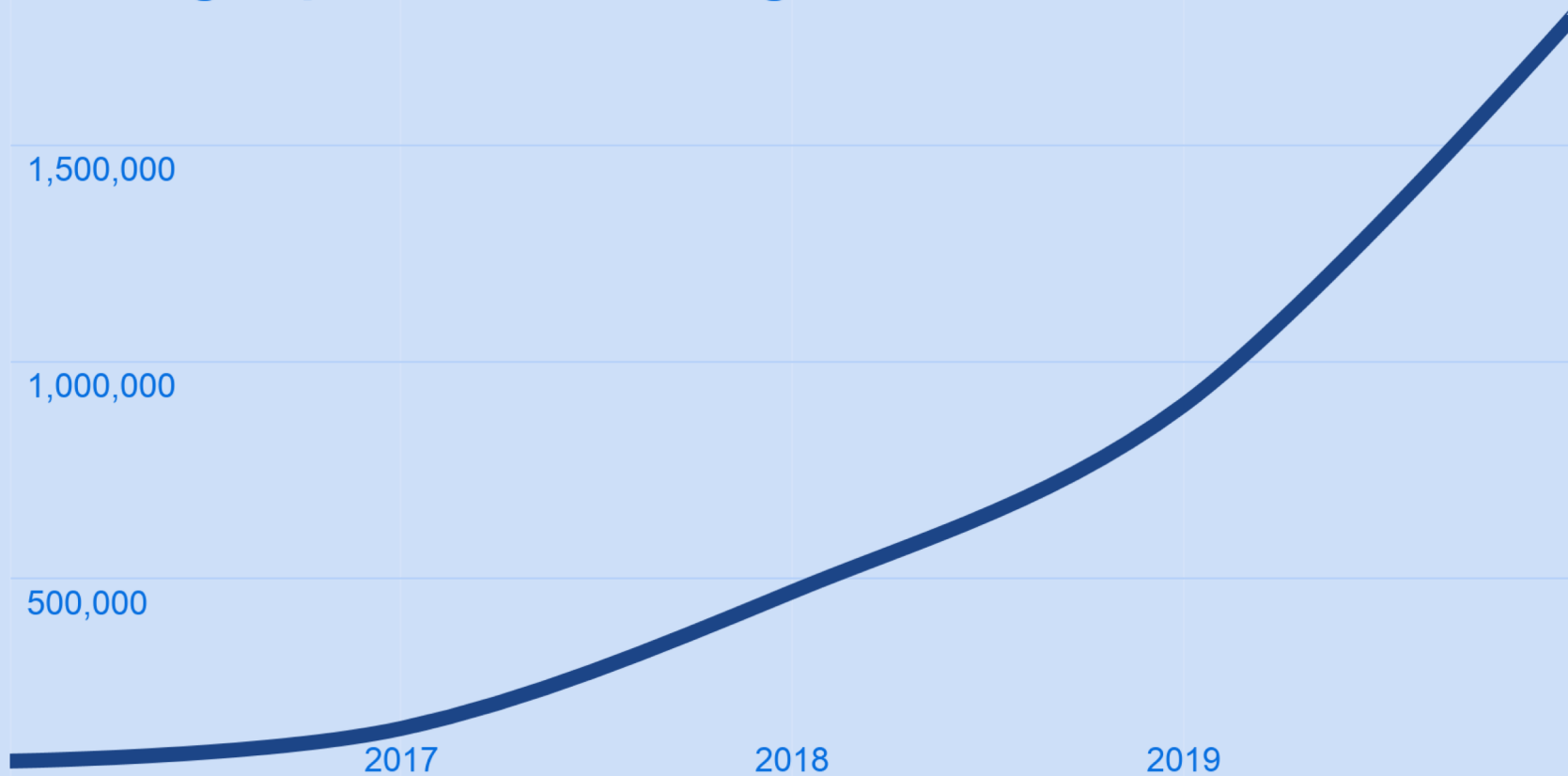
100

million app  
downloads

1.2k+

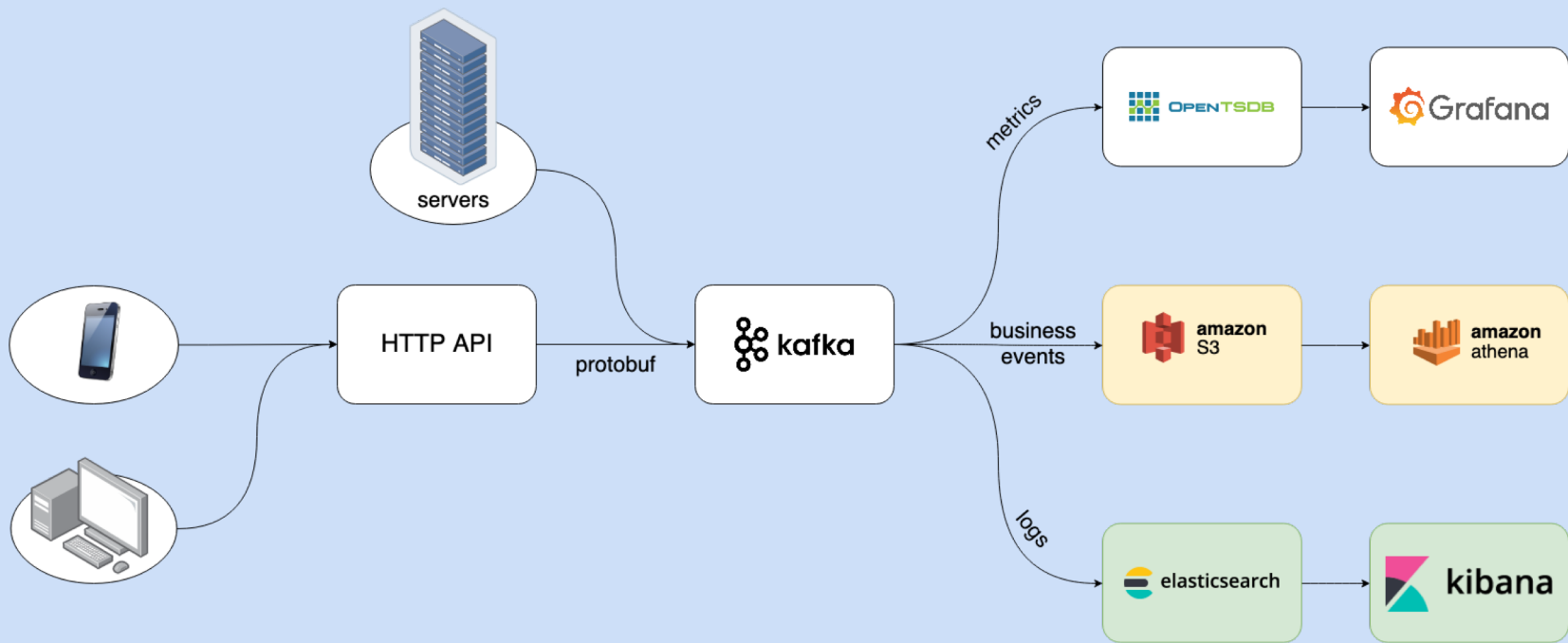
partners

# Messages per second through the Data Platform



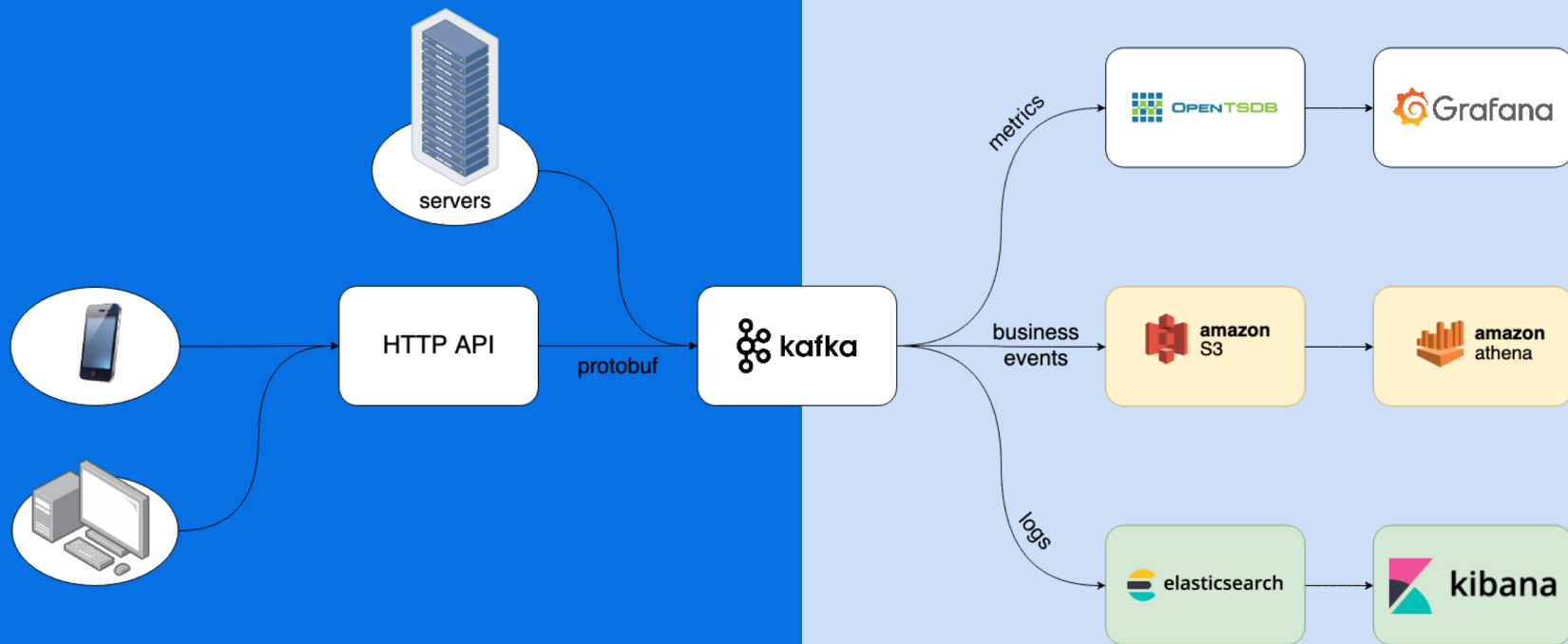
“The Cube”





# From batch to streaming

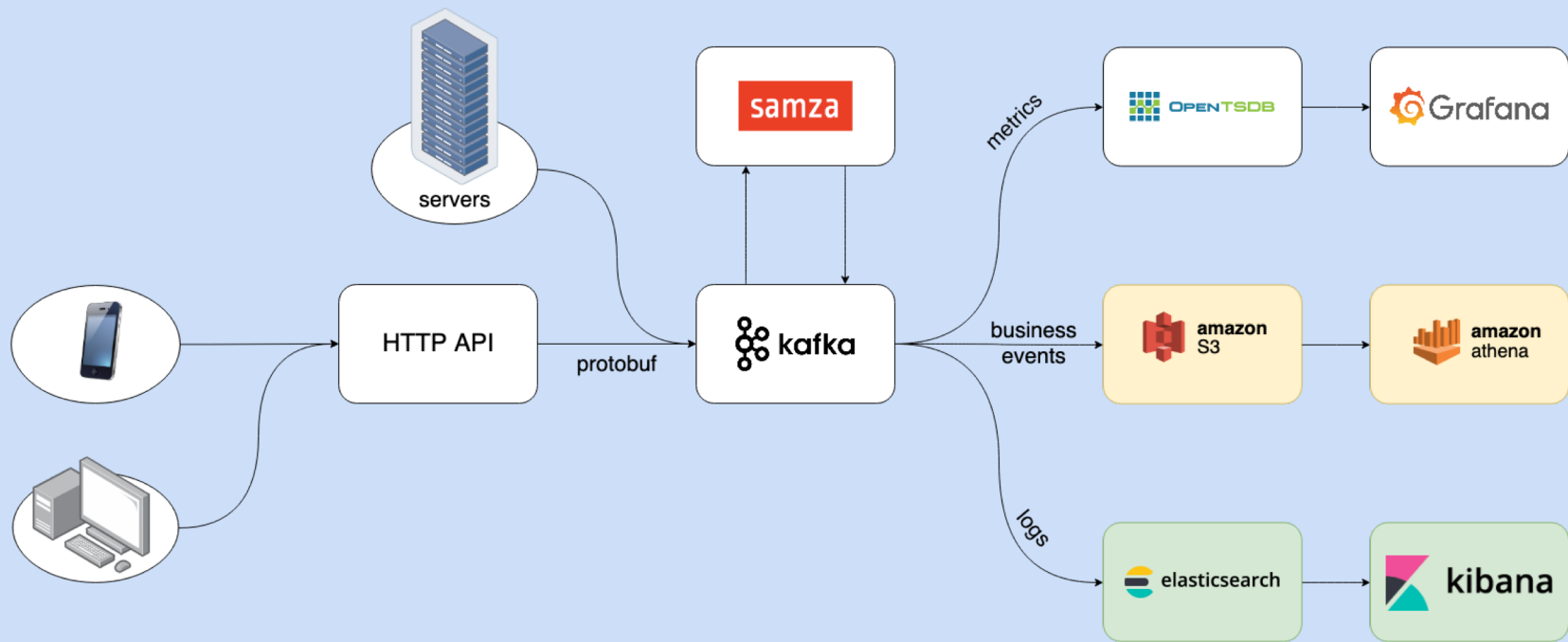




Decoupled

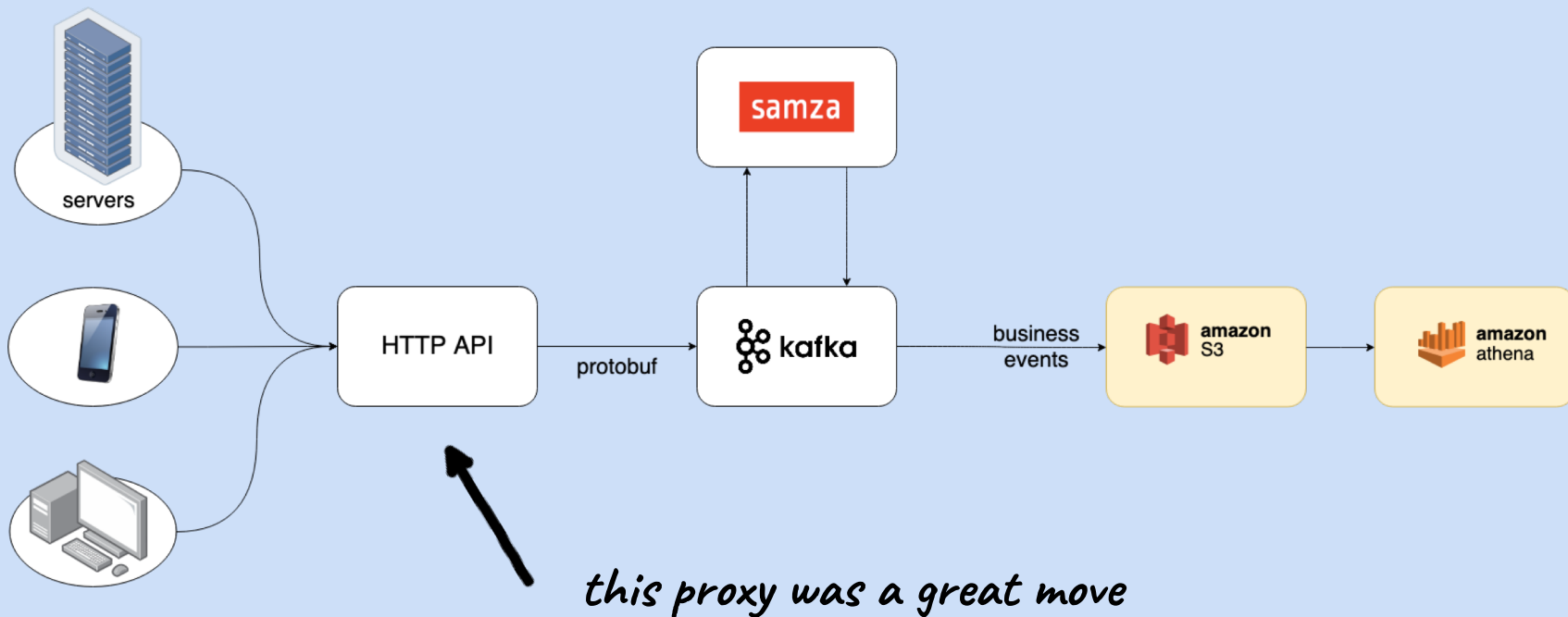
Unified

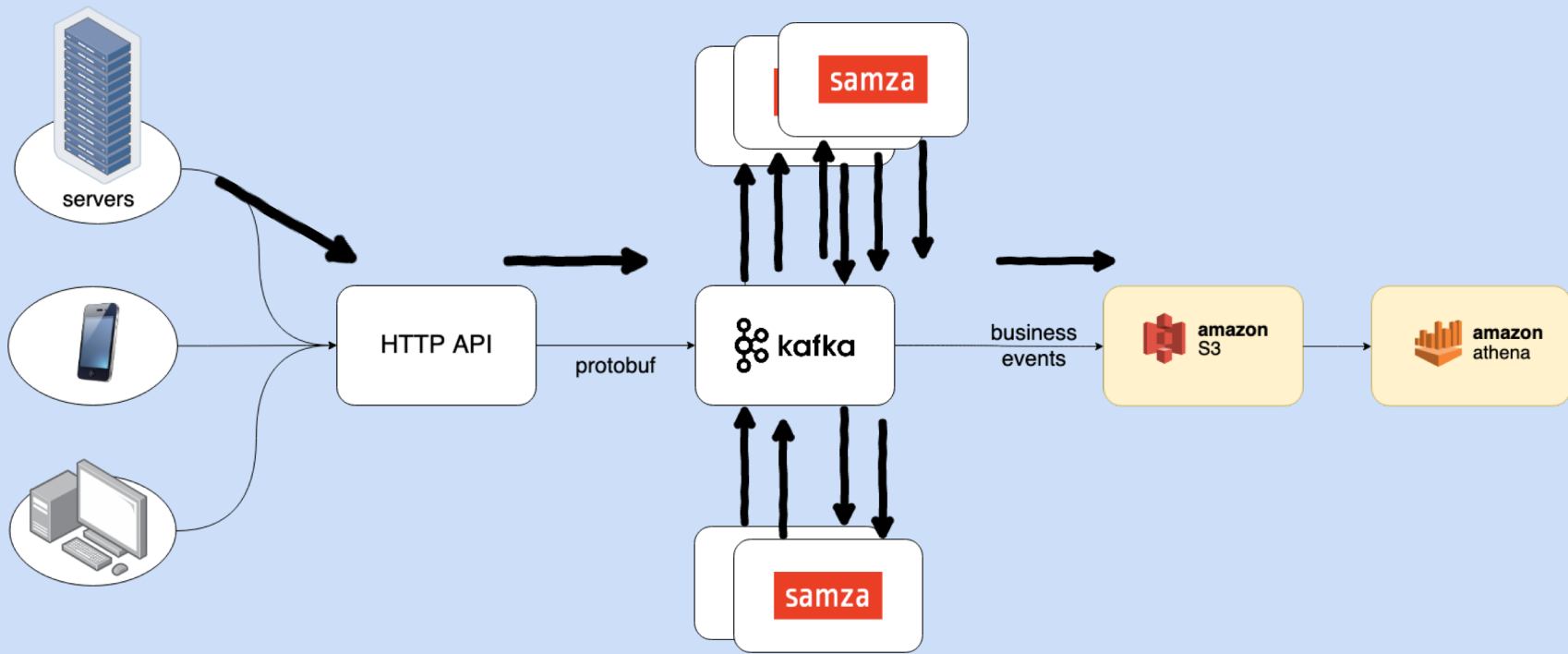




# The "Single Unified Log"







# The Single Unified Log





Lesson 1:  
Conway's Law is true for  
data platforms

“Organizations which design  
**data platforms** are  
constrained to produce  
designs which are copies of  
their communication  
structures”



Being self-serve is good

...but then metadata is critical



So let's talk about metadata



*prod.identity-service.audit.identity.AuditMessage*

*prod.flyingcircus.applog.applog.Message*

# A simple convention

<prod|sandbox|local>.<service-name>.<event-name>.<schema>.<message>



What does it mean?  
Who owns it?  
What does it contain?  
Where does it come from?

Descriptive



*we had some of this*

How does it relate to other  
data sets?  
How is it organized?  
How is it sorted /  
partitioned?

~~Structural~~

*Some, from using  
protobuf schemas*

How far does it date back?  
How frequently is it  
updated?  
How large is it?  
How complete is it?

~~Administrative~~

*nope.*

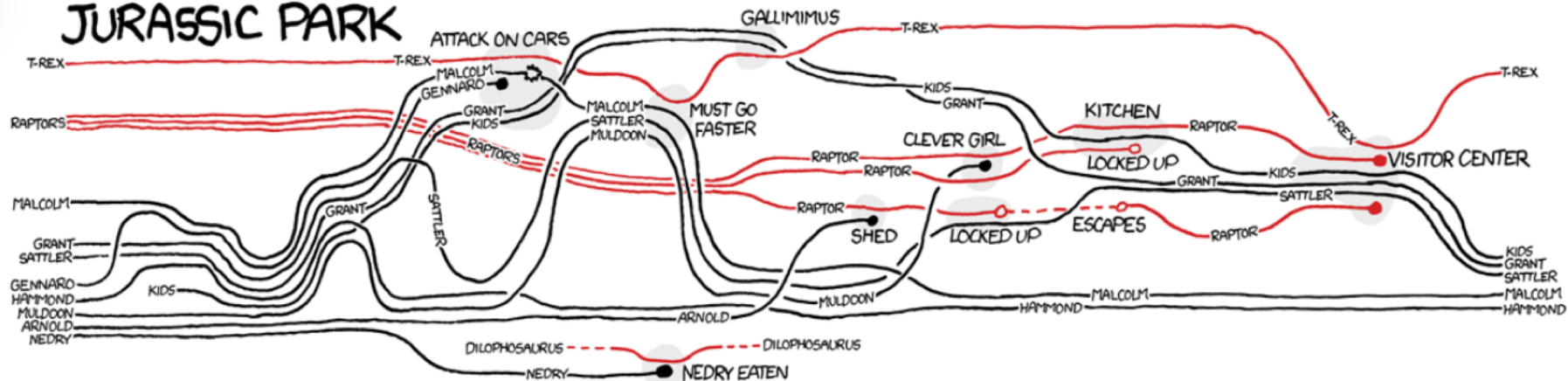


## Lesson 2: Metadata is Critical

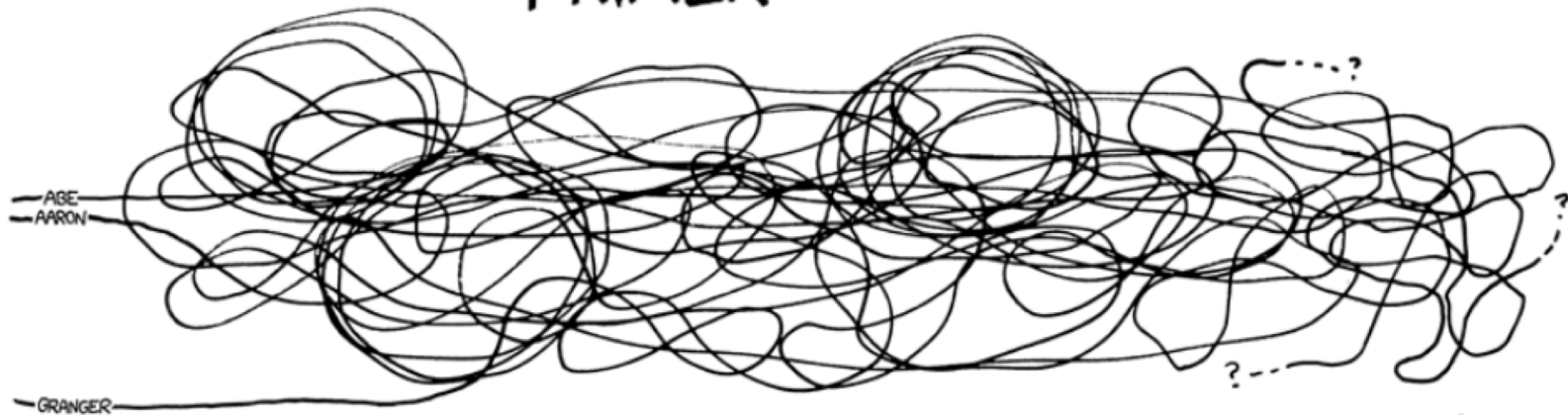
- Especially relationships
- Ideally automated
- Ideally from the start
- Tools like Schema Registry are a start, but not the full solution



# JURASSIC PARK



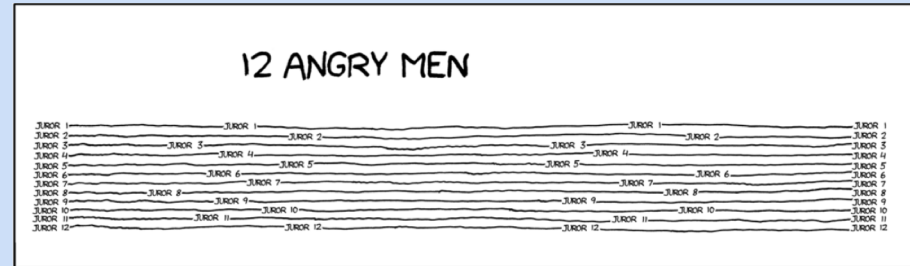
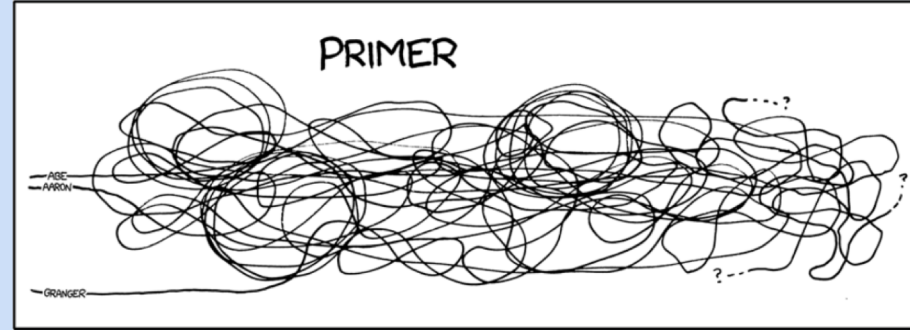
# PRIMER

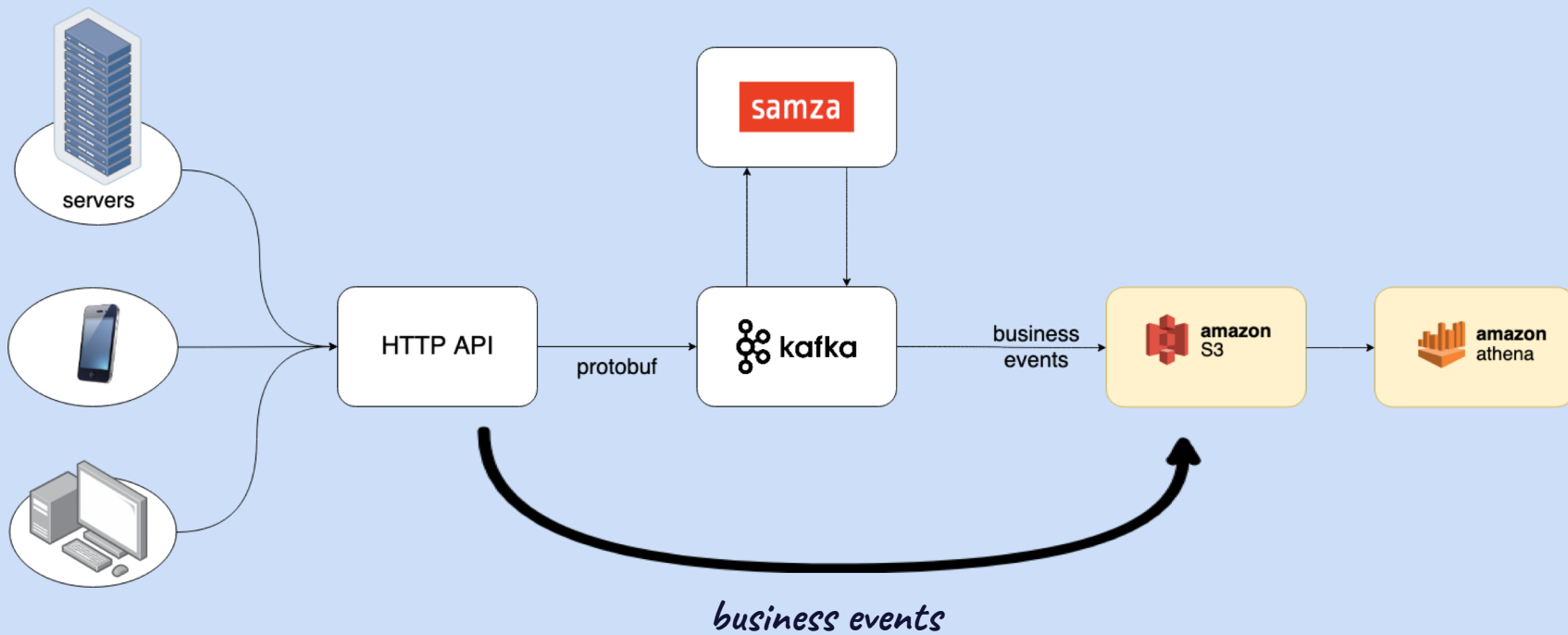


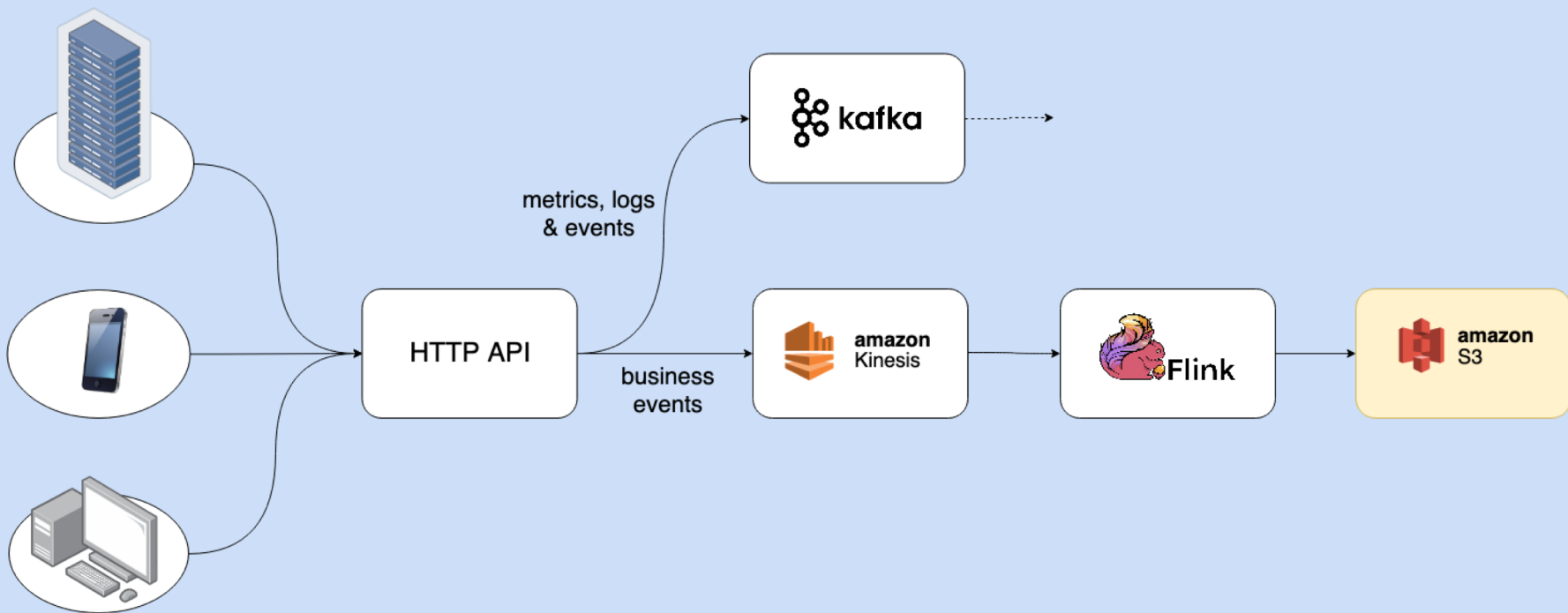


## Lesson 3:

# Data Engineers Control the Plot Line

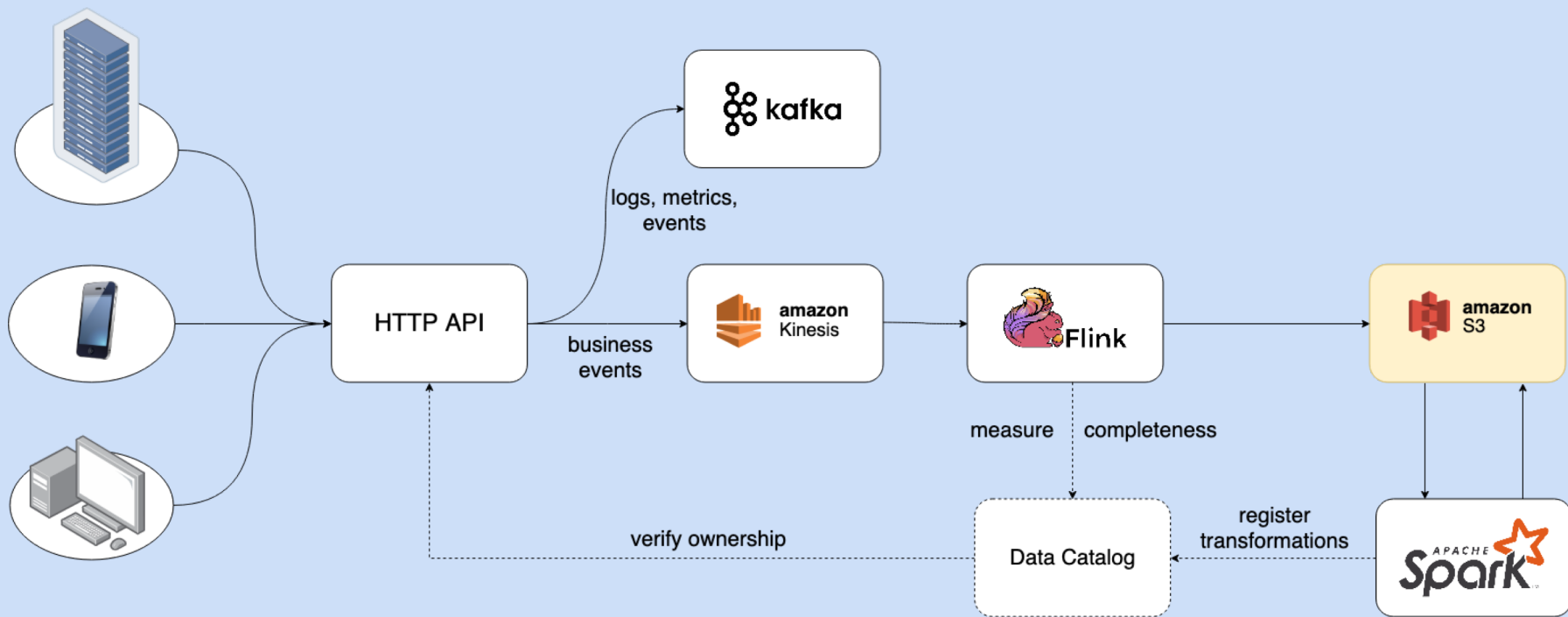






# From streaming to both





## Lesson 4:

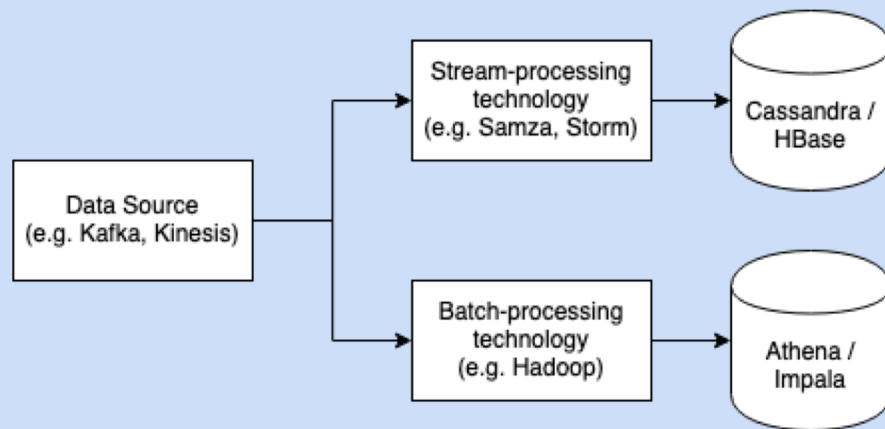
# Repeatability is important

- Streams have to choose between replays and accepting errors as permanent
- Replays can be complicated
- Batch processing can be done again any time
- Going straight to the archive in small batches gets benefits of both



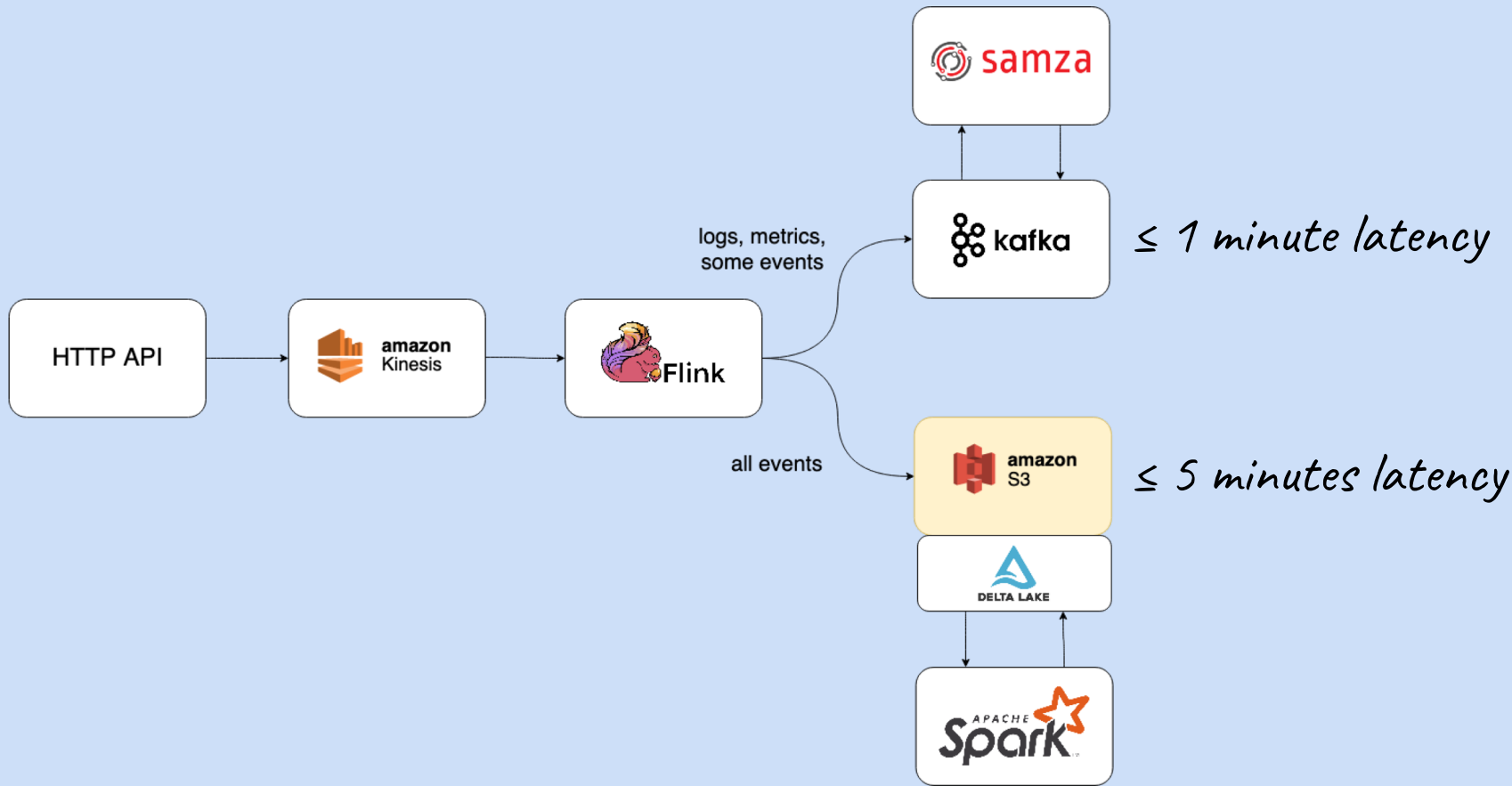


# Lambda architecture?





<https://delta.io/>





## Key Takeaways

- Conway's Law is true for data platforms
- Metadata is Critical
- Data Engineers Control the Plot Line
- Repeatability is important
- Delta Lake + Structured Streaming can bring the benefits of streaming to batch



# Thanks

## Contact

If you have any questions regarding  
Skyscanner please contact:

**Herman Schaaf**

[herman.schaaf@skyscanner.net](mailto:herman.schaaf@skyscanner.net)

**Skyscanner**

**Herman Schaaf**  
**@ironzeb**