



COMMUNITY DAY

— **TEL AVIV** —

Serverless data pipeline and ETL

Michael Haberman | 2018





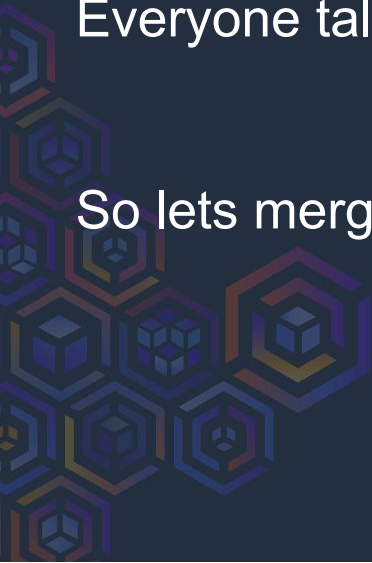
COMMUNITY DAY

Feeling the buzz?

Everyone talks about serverless

Everyone talks about data

So lets merge them together!





COMMUNITY DAY

Our Goal

Create a clickstream **data pipeline** with **ETL** process which is **serverless!**





COMMUNITY DAY



Save the events

Basic transformation

Some Storage

Data is ready / after some time

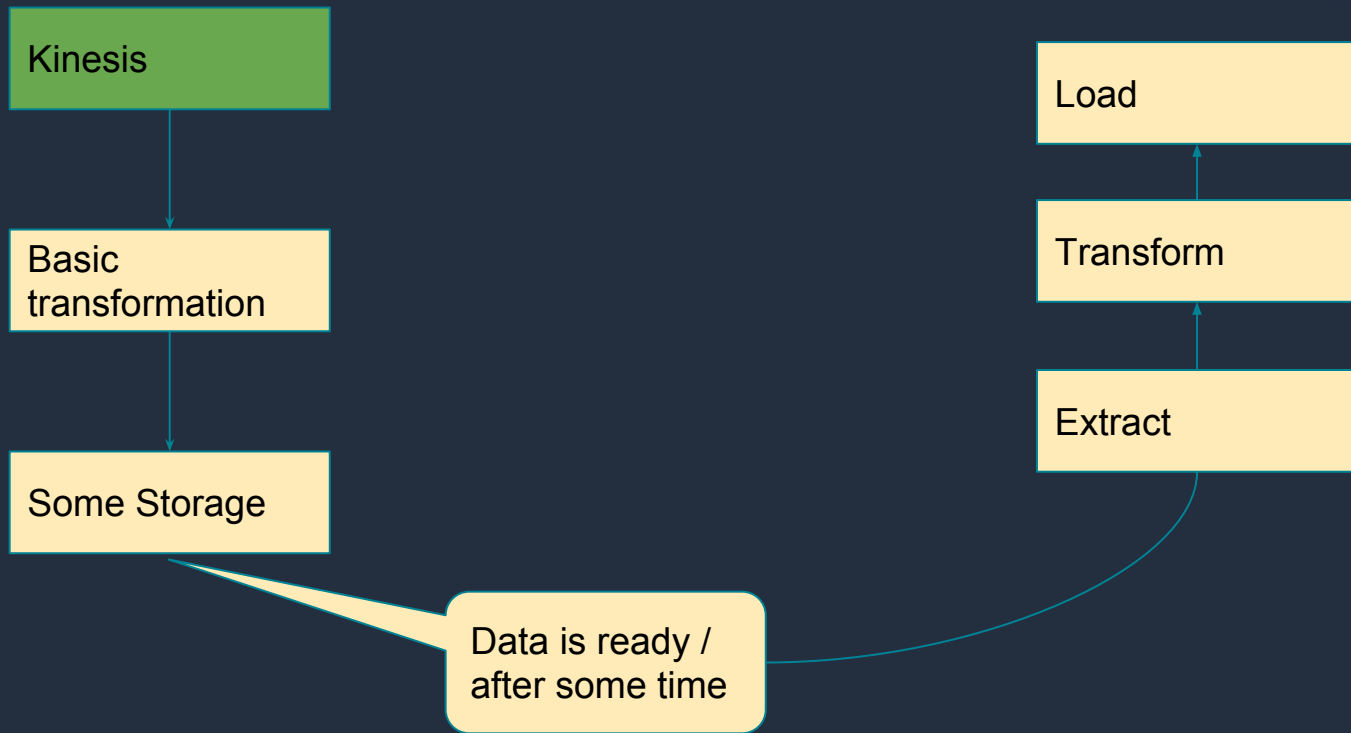
Load

Transform

Extract



COMMUNITY DAY





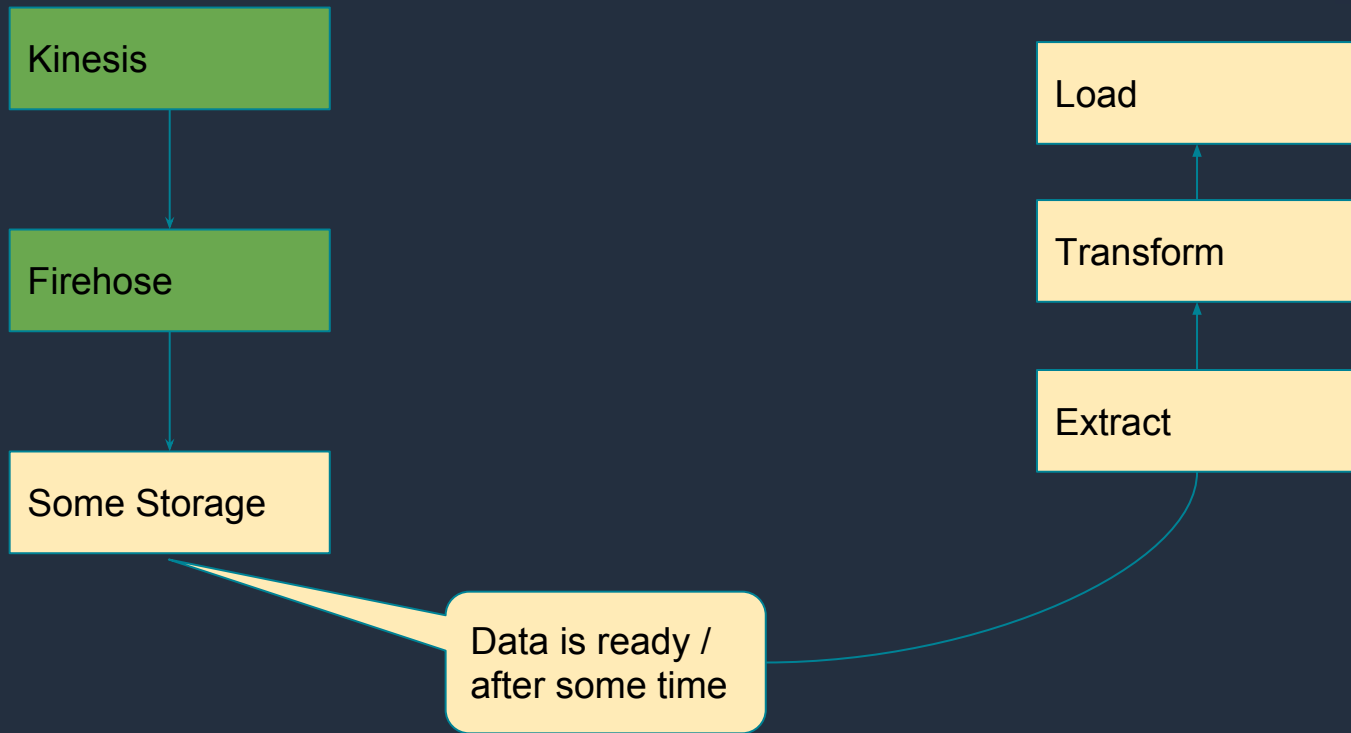
COMMUNITY DAY

Kinesis

- Ingest big data - sharded
- Multiple producers and consumers
- Data retention (7 days max)
- Sub products
 - Data analytics - SQL queries over windows
 - Firehose - transformation (Lambda, json to parquet), really fast
- Scale defined by numbers of shard



COMMUNITY DAY





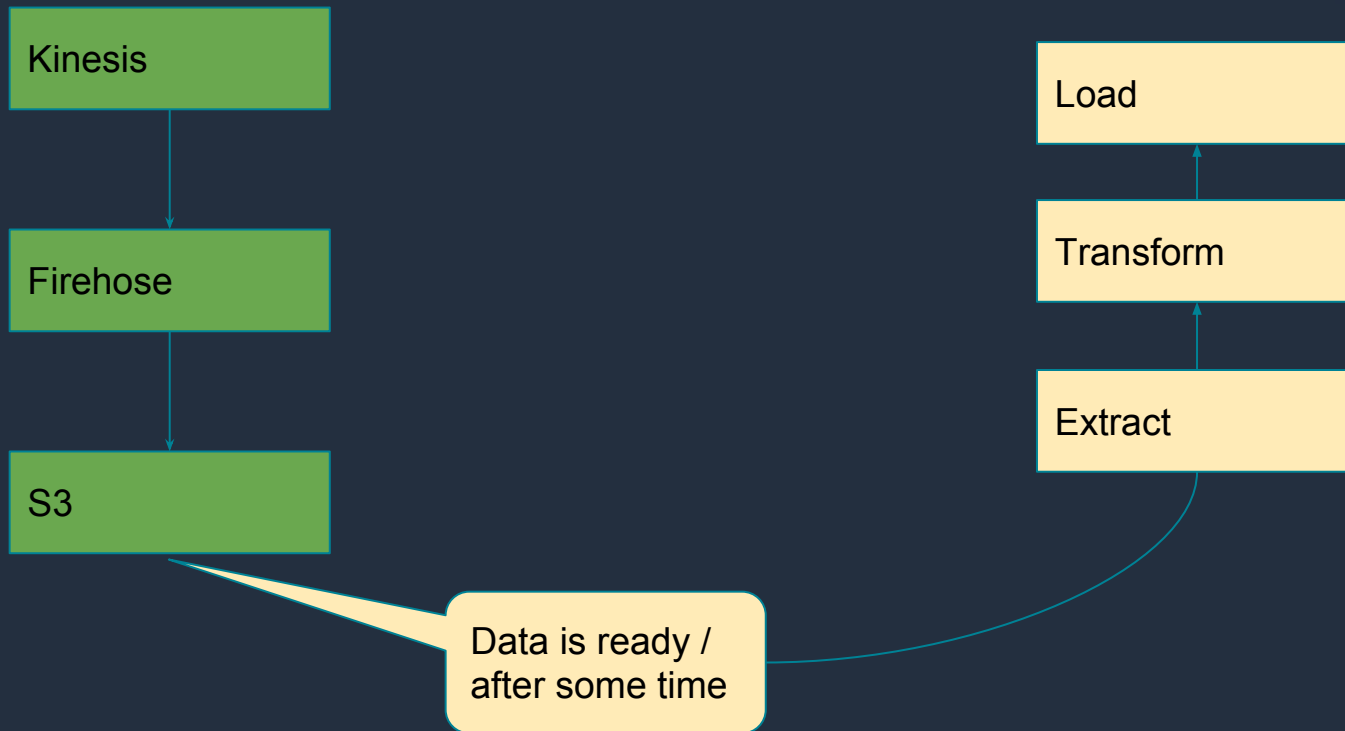
COMMUNITY DAY

Firehose

- Input
 - Kinesis
 - Direct
- Transformation
 - Using Lambda
 - Changes schema and type using AWS Glue
- Output
 - S3
 - Redshift
 - Elastic search
- No scaling definition
- Paying per event + Lambda execution



COMMUNITY DAY





COMMUNITY DAY

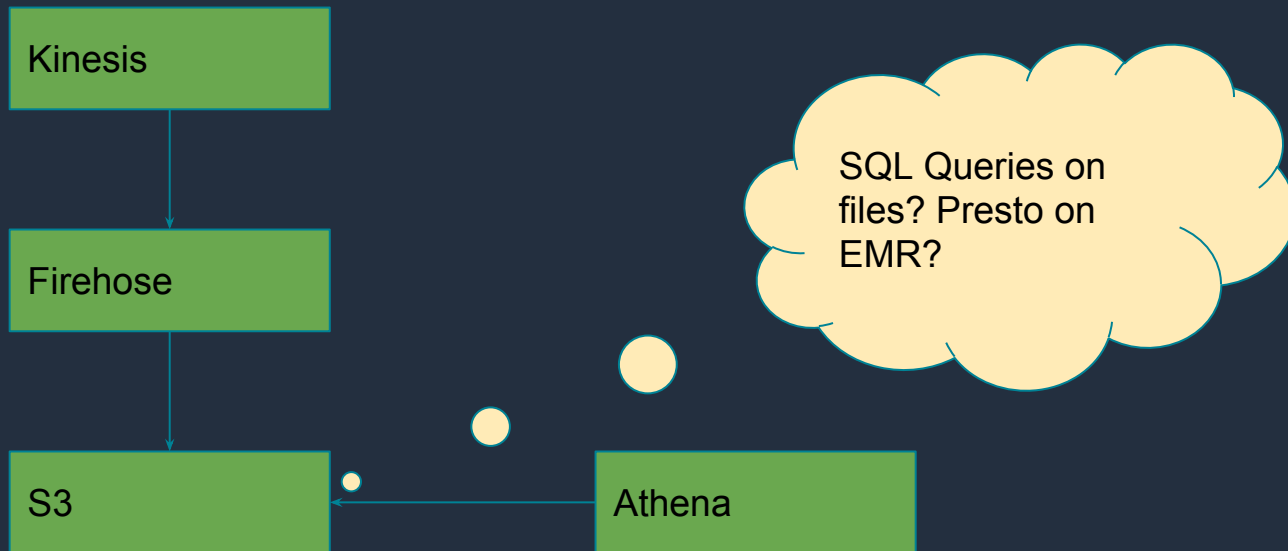
S3

- I guess you know S3...
- What if I want to query S3?





COMMUNITY DAY

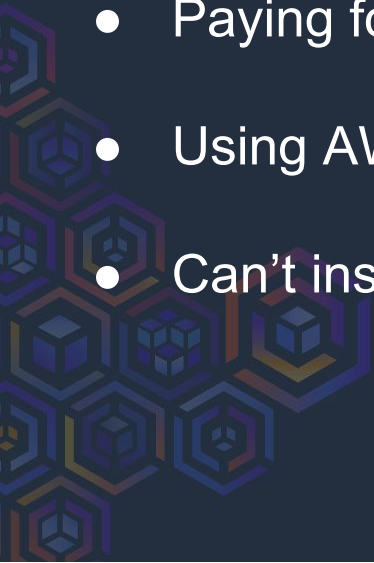




COMMUNITY DAY

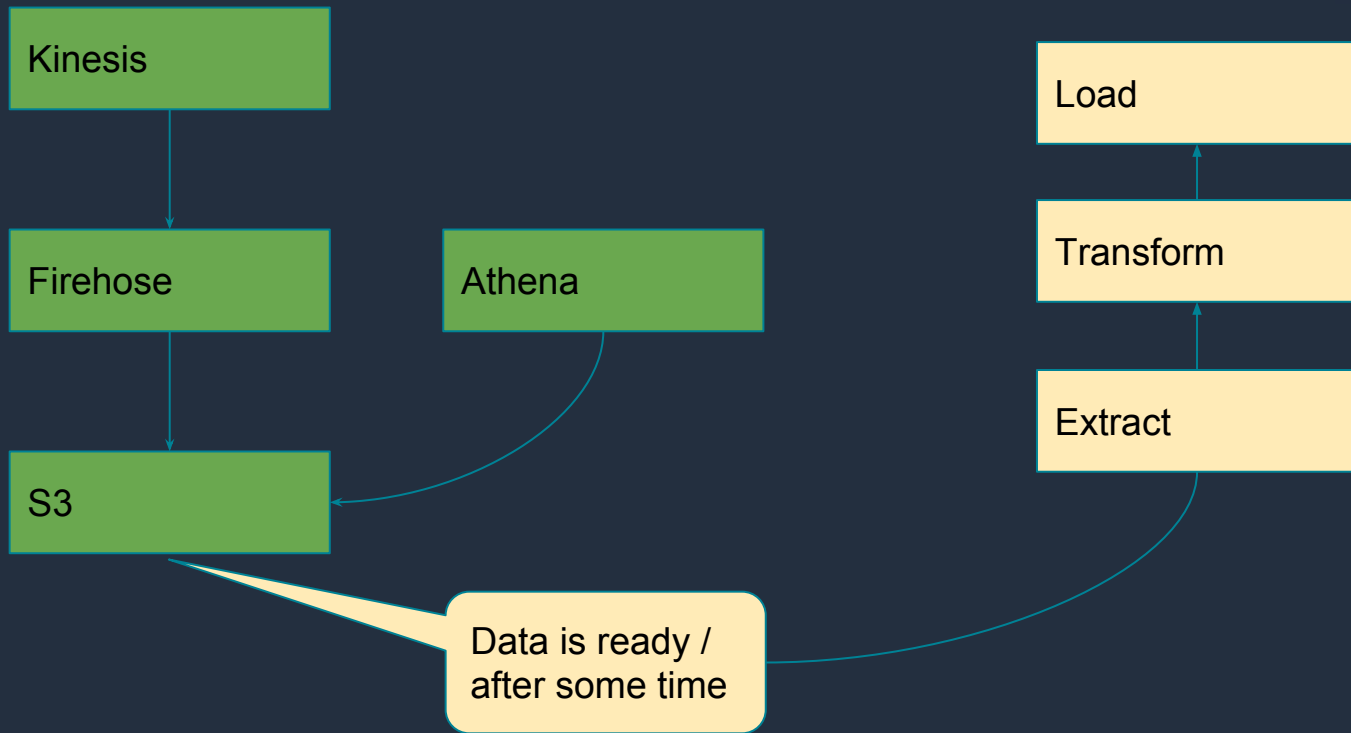
Athena

- Basically Presto
- Paying for scanned data (1TB = 5\$)
- Using AWS Glue for metastore
- Can't insert data



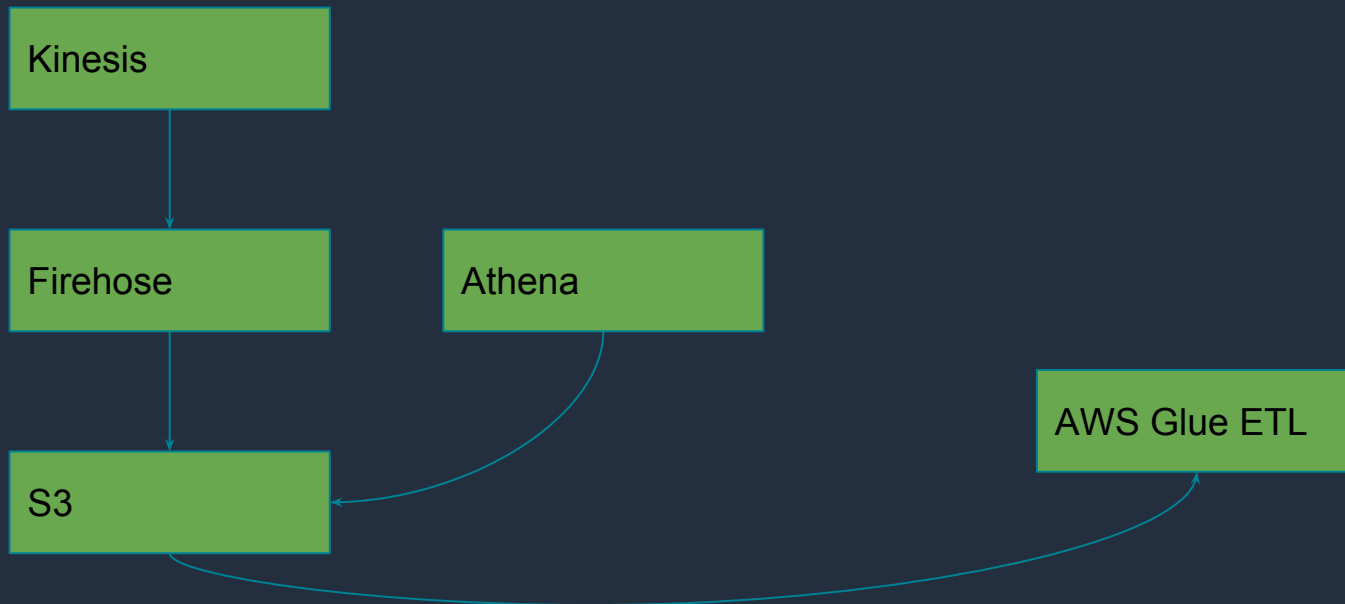


COMMUNITY DAY





COMMUNITY DAY





COMMUNITY DAY

Glue

- Extract - using metastore
- Transform - python code
- Load - S3, Redshift





COMMUNITY DAY

Quicksight

- Analytical dashboards





COMMUNITY DAY

Thanks!

michael@topsight.io