An aerial photograph of an LNG terminal. A large ship with three prominent orange spherical storage tanks is docked at a pier. Other smaller ships and industrial structures are visible in the background.

# Python, SQLAlchemy and Scrappy at Kpler

Jean Maynier  
CTO

# Business : commodities markets

- **Domain:** trading / maritime transport
- **Customers:** physical traders / analysts
- **Needs:** follow cargoes flows and prices impact

The first “**Cargo-Tracking**” system,  
bridging the gap between LNG **Shipping**  
and **Trading**.



Methodology : **4** data layers

“**Meta**” AIS (10 networks)

**Ports** Authorities

Gas **Inventories**

**Commercial** Database



## Why python ?

- Simple, fast to prototype
- Data libraries (sqlAlchemy, Pandas, Scikit-learn)
- Data oriented community
- Scrapy: best scraping framework

## Scrapy

- Simplify crawling
- Reusable extractors (xpath, css selectors)
- PaaS: Scrapinghub
- Blacklist: Tor, crawlera for IP rotation

## Process raw data

- Clean and normalize data
  - Positions
  - Port calls
  - Contracts
  - Prices
- Python PubSub system
- Persist data : sqlalchemy & Postgresql

## Enrich data

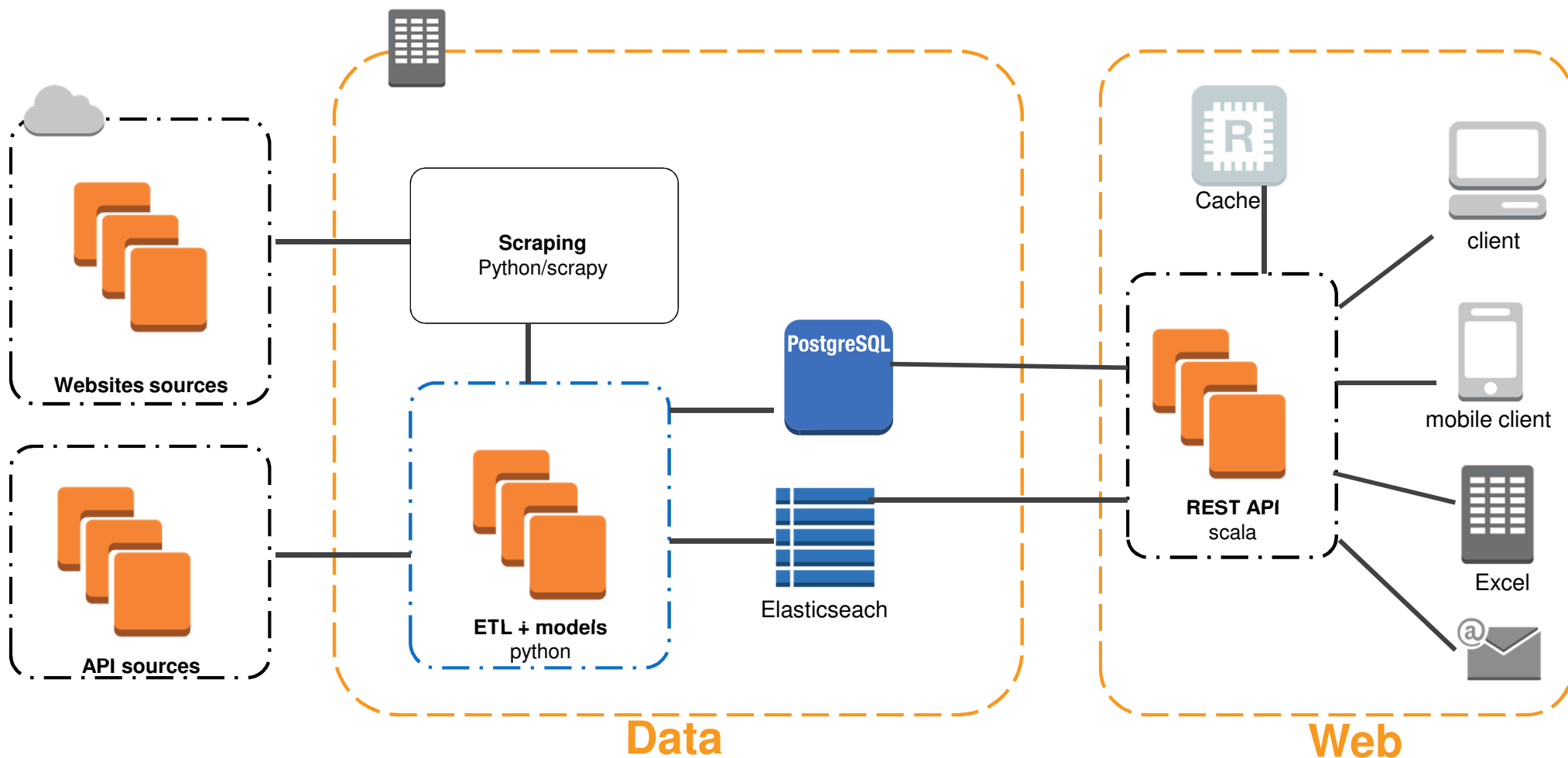
### Done

- Vessel dock/undock events
- Predict next destination
- Find buyer/seller
- Basic routing

### Todo / WIP

- Estimate price
- Routing with weather
- Portfolio optimisation
- Fleet optimisation

# Architecture





# Metrics trends

## Present

- DBs < 100Gb
- 50 sources
- 500 vessels
- Positions every 3min

## Future

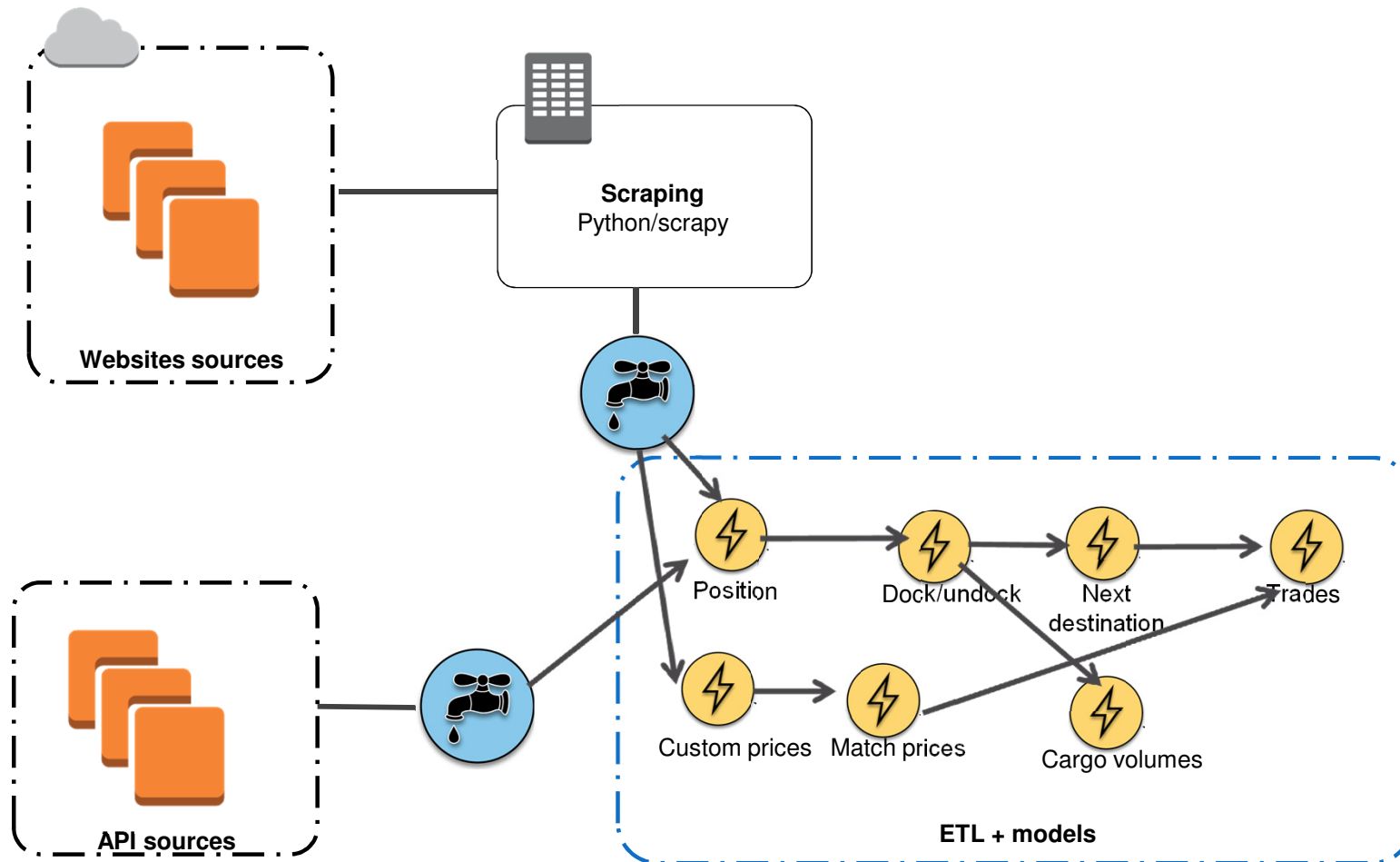
- 1-10 Tb
- 100 sources
- 10k to 100k vessels
- Position every 30s

**Performance problem !**

## Batch to data streaming

- Parallelization
- Granularity : item vs source
- Handle failure, no data loss, monitoring
- Akka streaming, Spark streaming, Celery, **Storm**

# Storm at Kpler: POC



We are recruiting, join us !

[jobs@kpler.com](mailto:jobs@kpler.com)