



# Supply data processing evolution

@ HomeToGo

Darius Kasiulevičius

Vilnius, 2022-03-03



# About me



- **Darius Kasiulevičius**
- **+10 years PHP developer**
- **~21 years since wrote first code line**
- **Team lead & Software architect at HomeToGo**



# Our Challenges

## Supply Data is critical for business

- Missing inventory – missed sales opportunity
  - Including search filters, e.g. destination, dates, pool, wifi, pets, etc.
- Data accuracy (price, availabilities) issues leads to poor conversion rates
- Missing data points leads to poor conversion rates. E.g. cancellation policies, room plan, etc.

# First integration

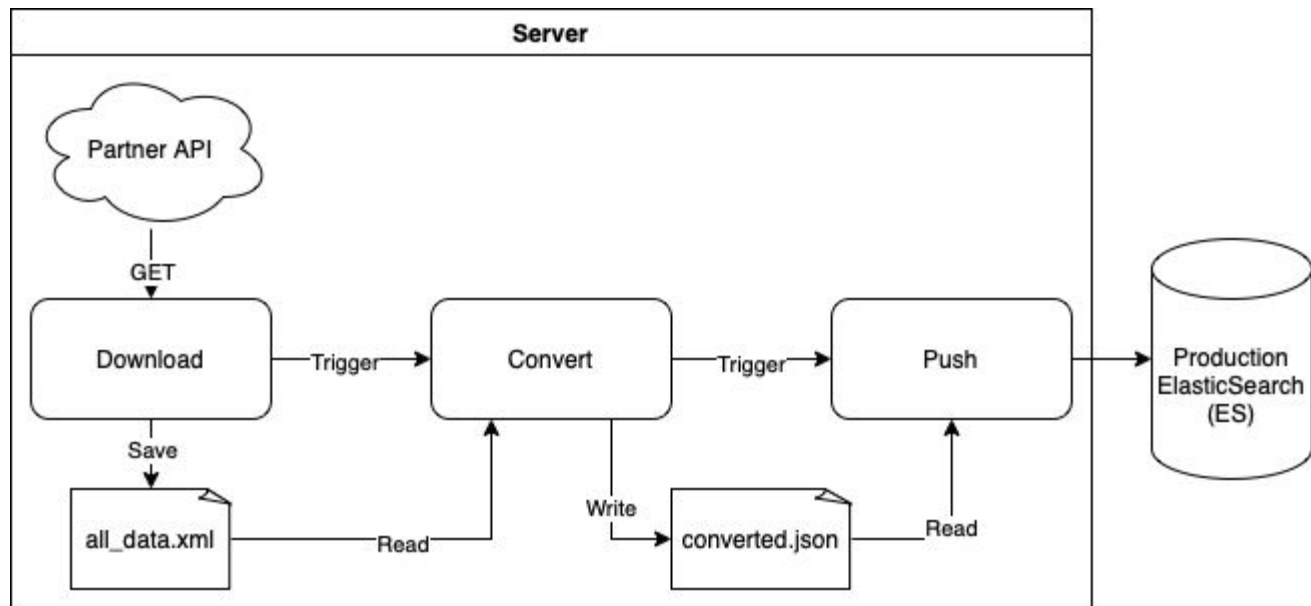
First partner

# First partner

- Small one – less than 1000 offers
- One file with full data
- Data file type XML
- Simple download over HTTP



# First pipeline

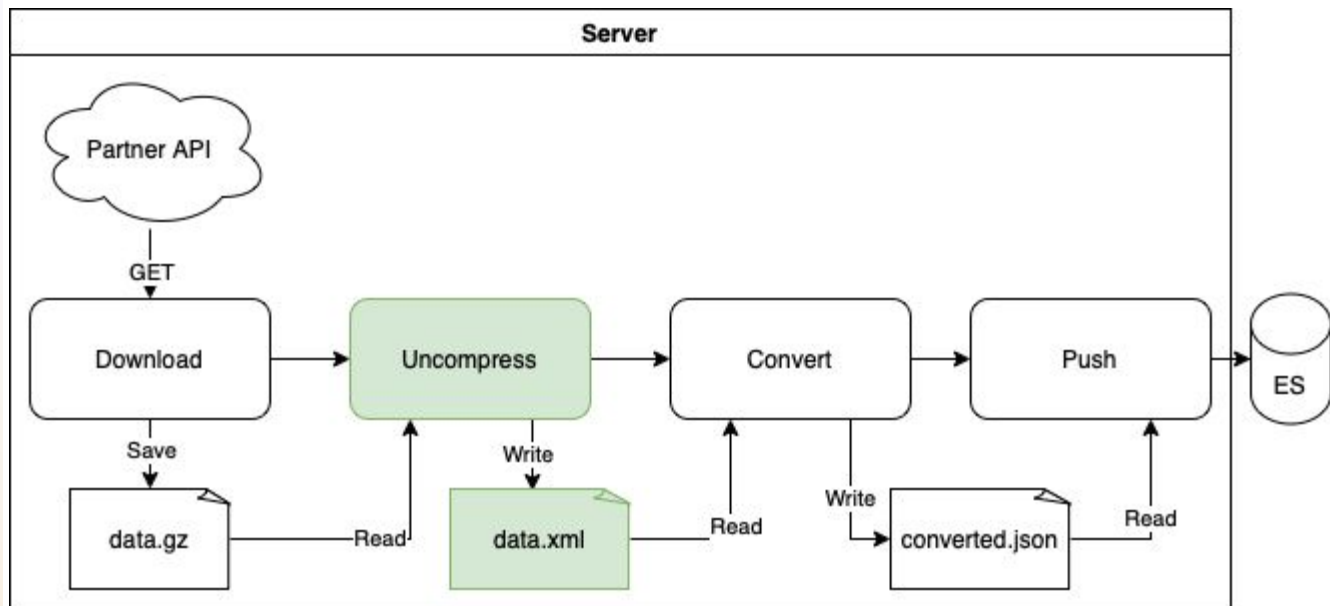


## Next partners

- Compressed file
- Download over FTP, SFTP, SOAP
- File types CSV, JSON
- To get all data need make X GET requests



# Next pipeline



# First improvement

## Custom pipelines per partner

Download more than one file (units.xml.gz, locations.xml, descriptions.xml.gz)

Uncompress

Validation

# Task Manager

- Custom pipeline per partner
- Monitor pipeline status
- Resume pipeline from failed task
- More than one pipeline per partner  
(locations)

# Task Manager configuration

```
12 vrt_providers_adria_gate.task_manager.download_manager.config:
13 |
14 |     task: 'download_manager'
15 |     depends: { task: ['prepare_db'] }
16 |
17 vrt_providers_adria_gate.task_manager.healer.config:
18 |
19 |     task: 'heal'
20 |     params: { source: 'original_source/raw_*.xml' }
21 |     depends: { task: ['download_manager'] }
22 |
23 vrt_providers_adria_gate.task_manager.split_prepare_data.config:
24 |
25 |     task: 'prepare_data'
26 |     params:
27 |         source: 'adria_gate'
28 |         aggregationConfig: '%vrt_providers_adria_gate.aggregation_config%'
29 |     depends: { container: ['vrt_providers_adria_gate.task_manager.healer'] }
30 |
31 vrt_providers_adria_gate.task_manager.after_prepare_data.config:
32 |
33 |     task: 'after_prepare'
34 |     depends: { container: ['vrt_providers_adria_gate.task_manager.split_prepare_data'] }
35 |
36 vrt_providers_adria_gate.task_manager.template.config:
37 |
38 |     template: 'common_from splitted_convert.yml'
39 |     params:
40 |         inheritableTaskId: 'after_prepare'
41 |         masterTable: SessionUnit
42 |         suffix: 'adria_gate'
```

# Download Manager configuration

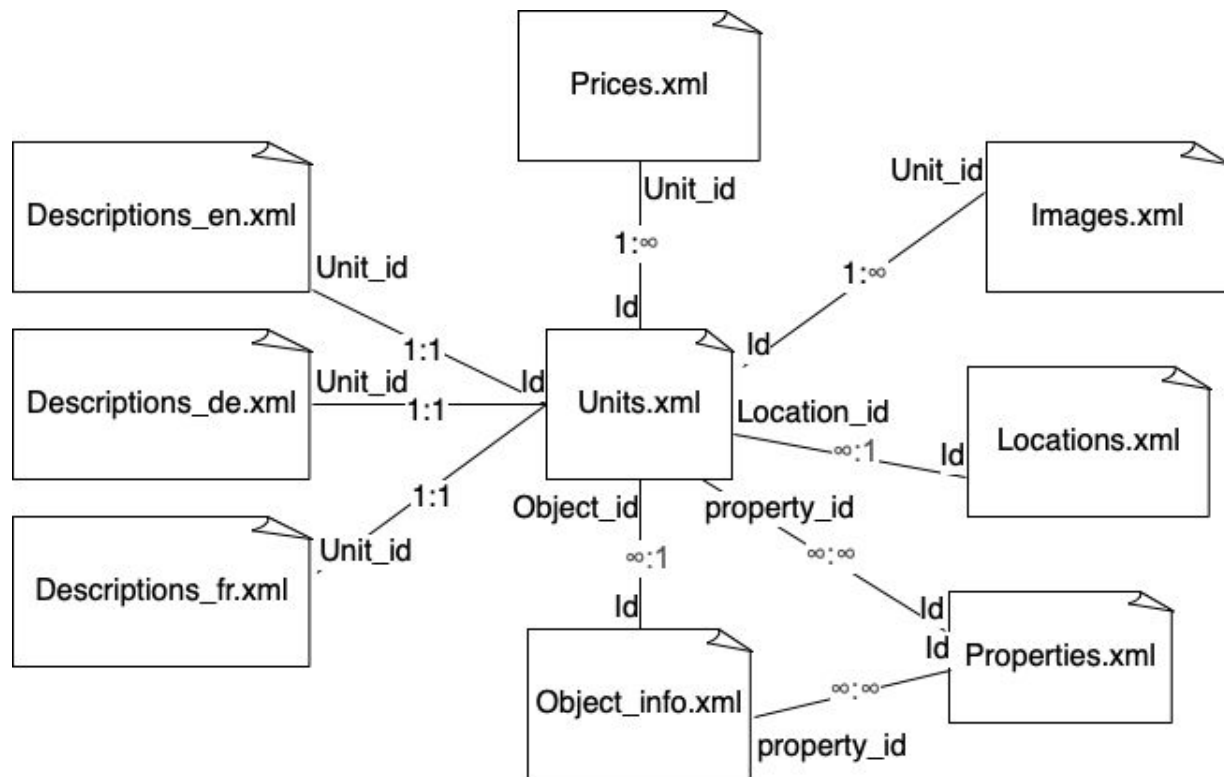
```
82 vrt_providers_adria_gate.download_manager_config:
83   _params:
84     concurrency: 20
85     curl:
86       CURLOPT_HTTPAUTH: 1 # CURLAUTH_BASIC
87       CURLOPT_USERPWD: XXXXXXXXXXXXX
88   areas:
89     source: 'https://ws.adriagate.com/agws/v03/xml/GetAreas?language=en'
90     target: 'areas_en.xml'
91   index:
92     source: 'https://ws.adriagate.com/agws/v03/xml/GetListingIndex'
93     target: 'index.xml'
94   objects:
95     parent: 'index'
96     source: '{$url}'
97     target: ''
98     generator:
99       -
100         class: 'Htg\Backend\DownloadManagerBundle\Request\Generators\Ids'
101         idPattern: '(?<=<url>)(.+?)(?=</url>)'
102         pattern: '{$url}'
103         ignoreEmpty: true
104       -
105         class: 'Vrt\Providers\AdriaGateBundle\DownloadManager\DecodeSourceUrl'
106       -
107         class: 'Vrt\Providers\AdriaGateBundle\DownloadManager\ModifyTarget'
108         search: '/(.+id=)(\d+-\d+)(.*)/'
109         replace: 'object_$.xml'
110   availabilities:
111     parent: 'objects'
112     source: 'https://ws.adriagate.com/agws/v03/AdriagateService.svc/xml/GetUnitAvailability?i'
113     target: 'availabilities_{$id}.xml'
```

# Second improvement

More than one file to get full data

Images.xml  
Descriptions\_de.xml  
Descriptions\_en.xml  
Descriptions\_fr.xml  
Prices.xml  
Units.xml

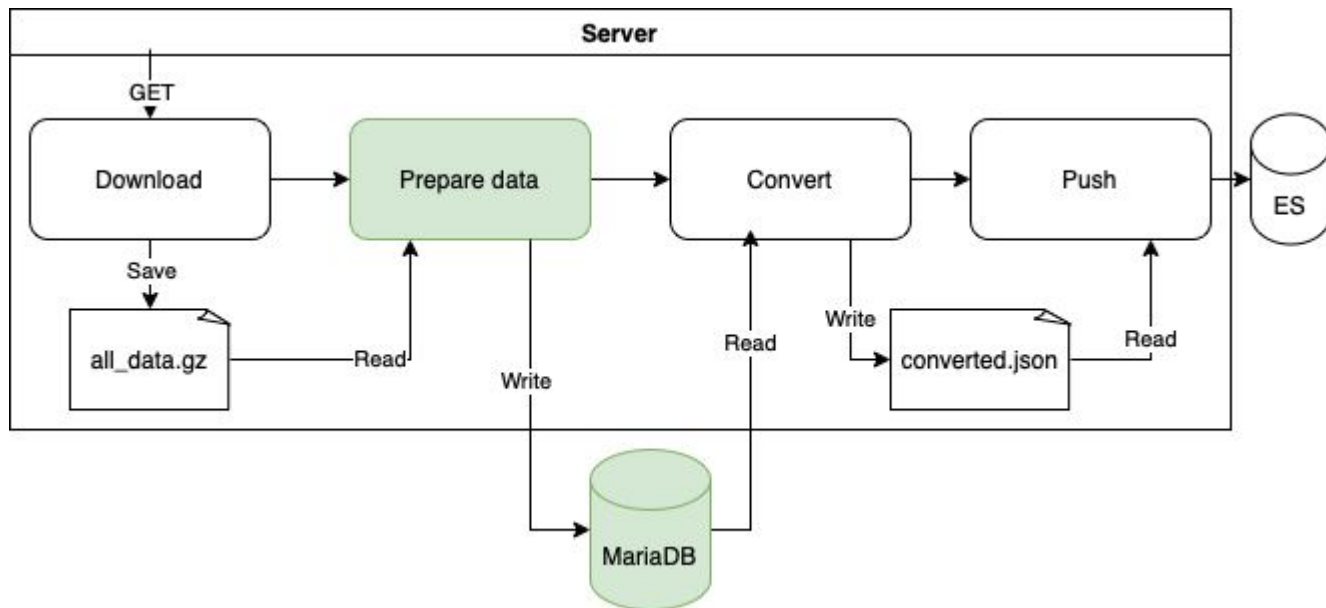
# Data relationship





# Permanent DB

- Separate DB per partner
- Separate table per file type (prices, images, description, etc.)
- Uniq table structure per partner



# Prepare data configuration

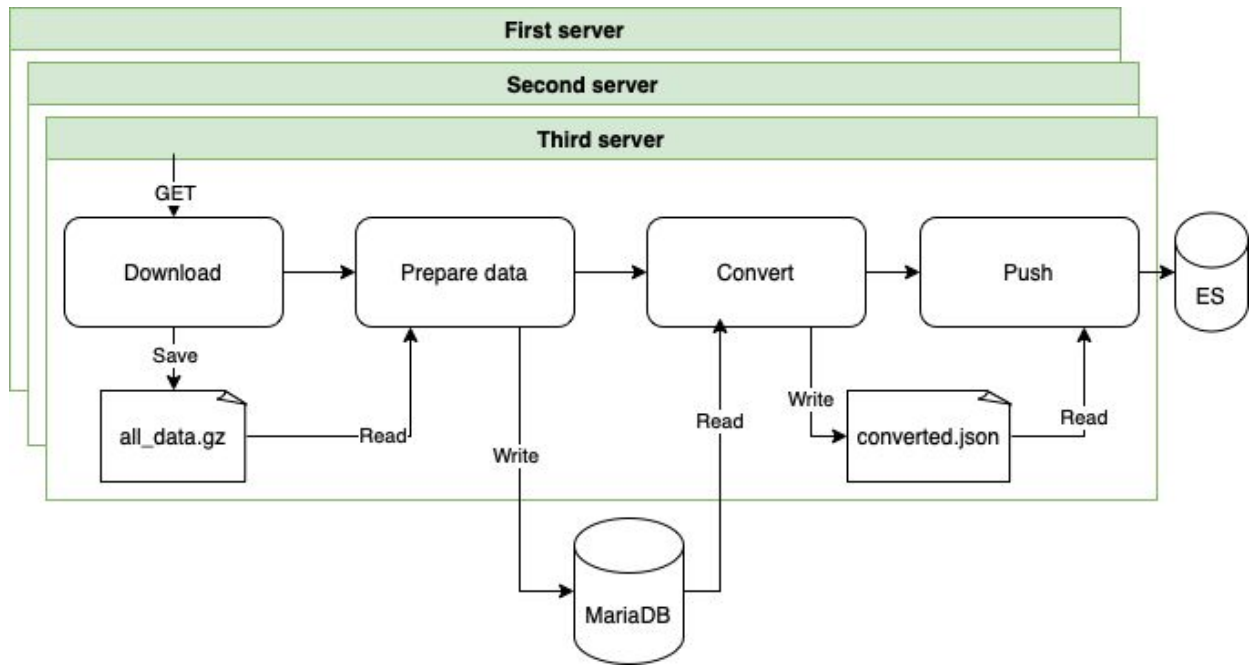
```
5  vrt_providers_adria_island.aggregation_config:
6
7      objects:
8          tag: Property
9          xsl: unit.xsl
10         entity: Units
11         master: true
12         types: [default]
13         file: 'inventory_*.xml'
14         type: csv2
15         params: <1 key>
16         dir: original_source
17         source: XmlParserFileSource
18     prices:
19         tag: Prices
20         xsl: price.xsl
21         entity: Prices
22         types: [default]
23         file: 'prices_*.xml'
24         type: csv2
25         params:
26             duplicate: REPLACE
27             file_name_parts: 'prices_(\d+)_.*.xml'
28             file_name_parts_mapping: 'year'
29         dir: original_source
30         source: XmlParserFileSource
31     vacancies:
32         tag: Availability
33         xsl: availability.xsl
34         entity: Availabilities
35         types: [default]
36         file: 'vacancies_*.xml'
37         type: csv2
38         params:
```

# Third improvement

Update data for all partners at same time

Partner count increased

# Increasing server count to 3

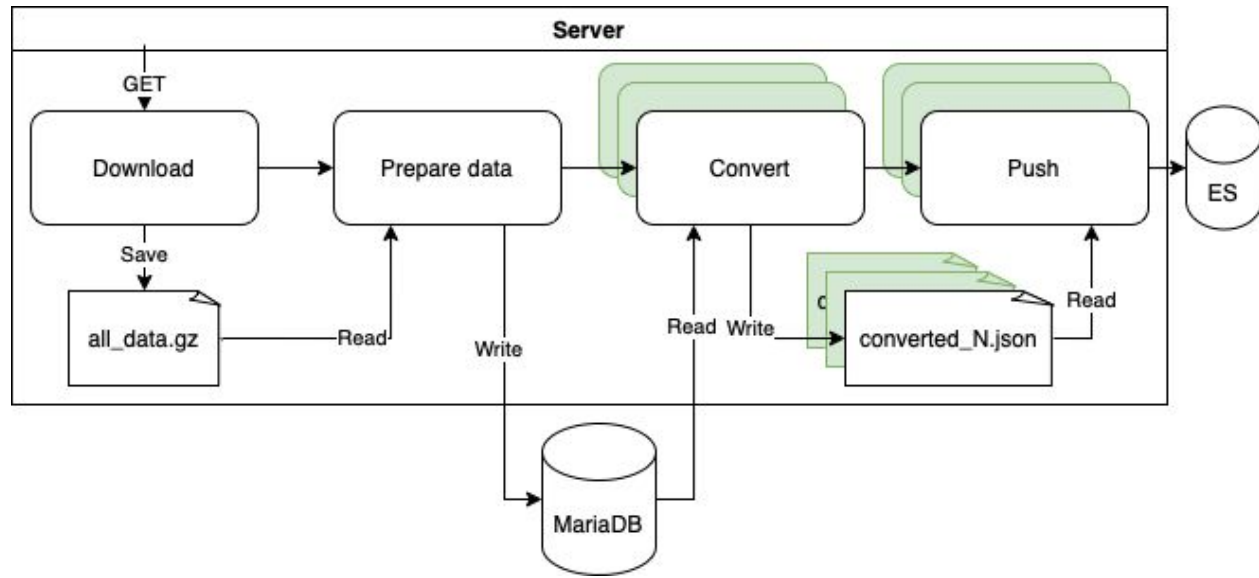


# Fourth improvement

## Process run time optimization

Partner with 1 million offers

# Split convert



`SELECT select_list FROM table_name LIMIT offset, row_count;`

First convert run time << second convert run time <<< third convert run time

`SELECT select_list FROM table_name WHERE autoincrement_id BETWEEN 5000 AND 10000;`

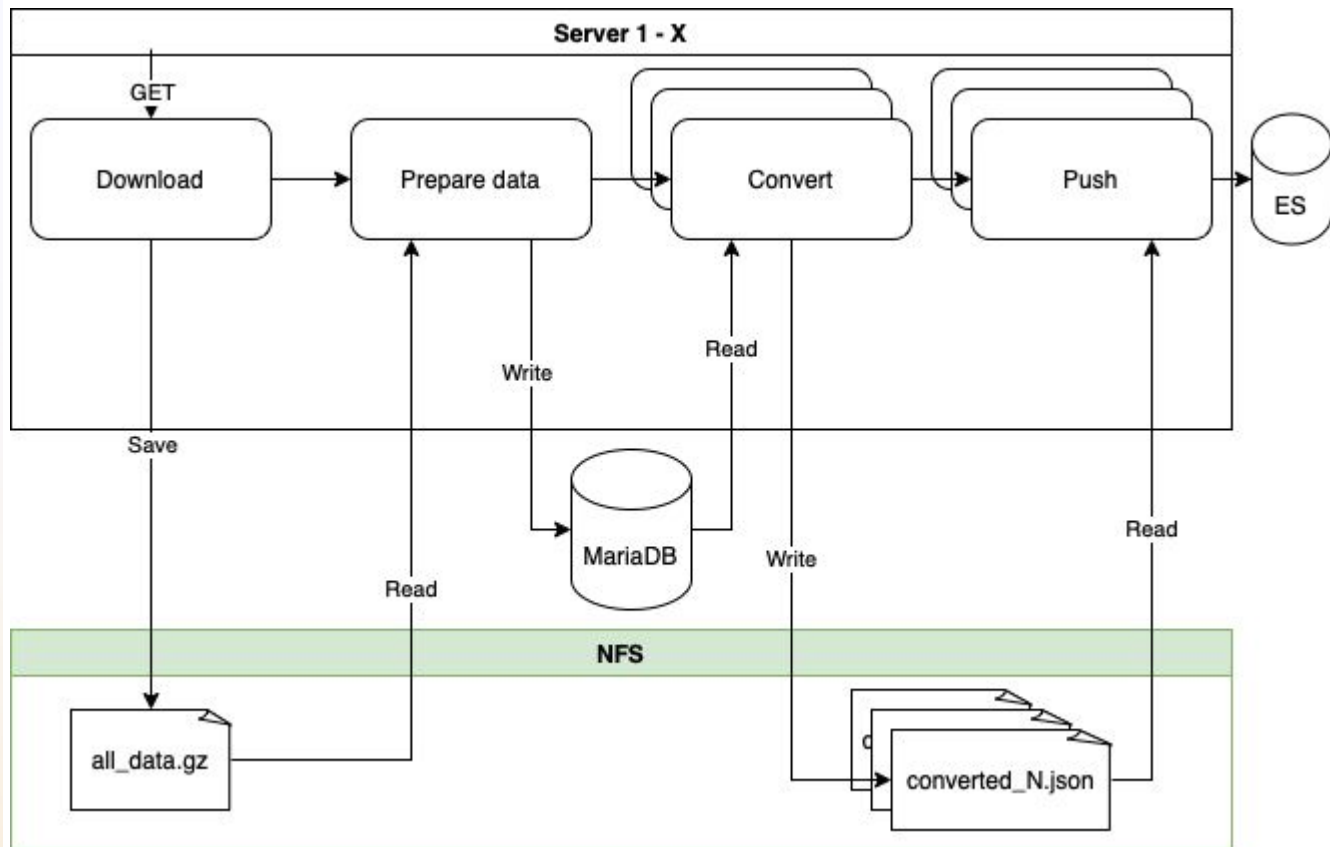
# Fifth improvement

## Equal use of server resources

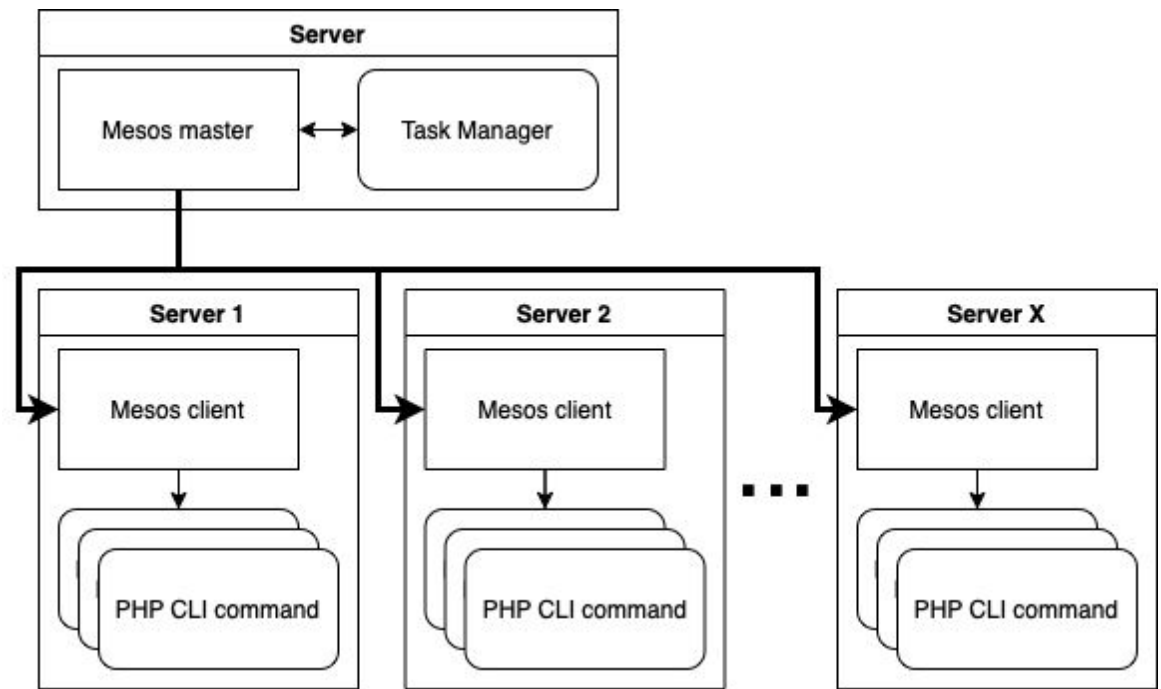
Big partners competes with the small ones



# Add NFS



# Start using



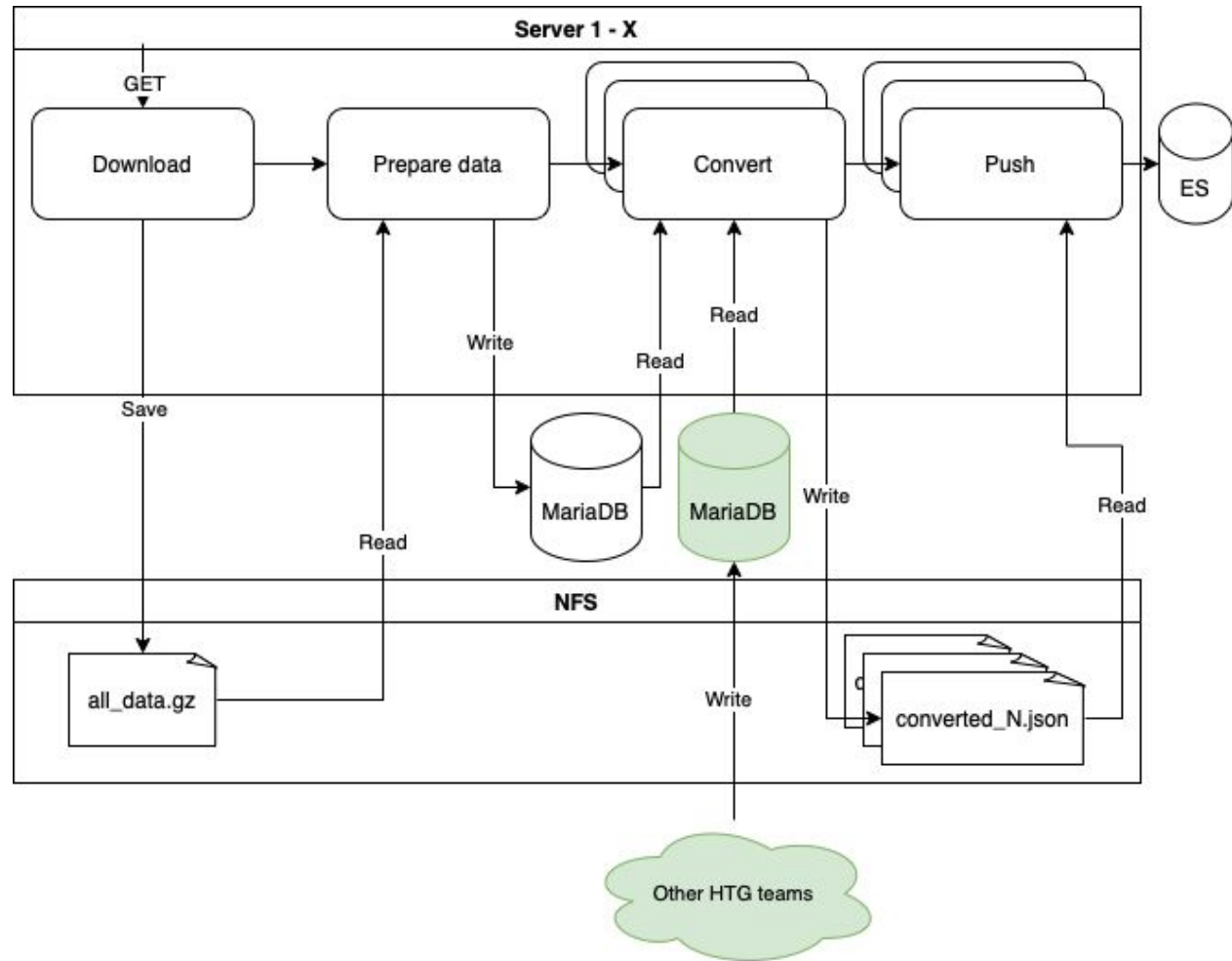
- Task Manager changes
  - Communicate with Mesos API
  - Support X servers
  - Partner commands run in any server
  - Support auto scaling

# Sixth improvement

Enrich data from third sources

Distance to water  
Sorting data

# Enrich data

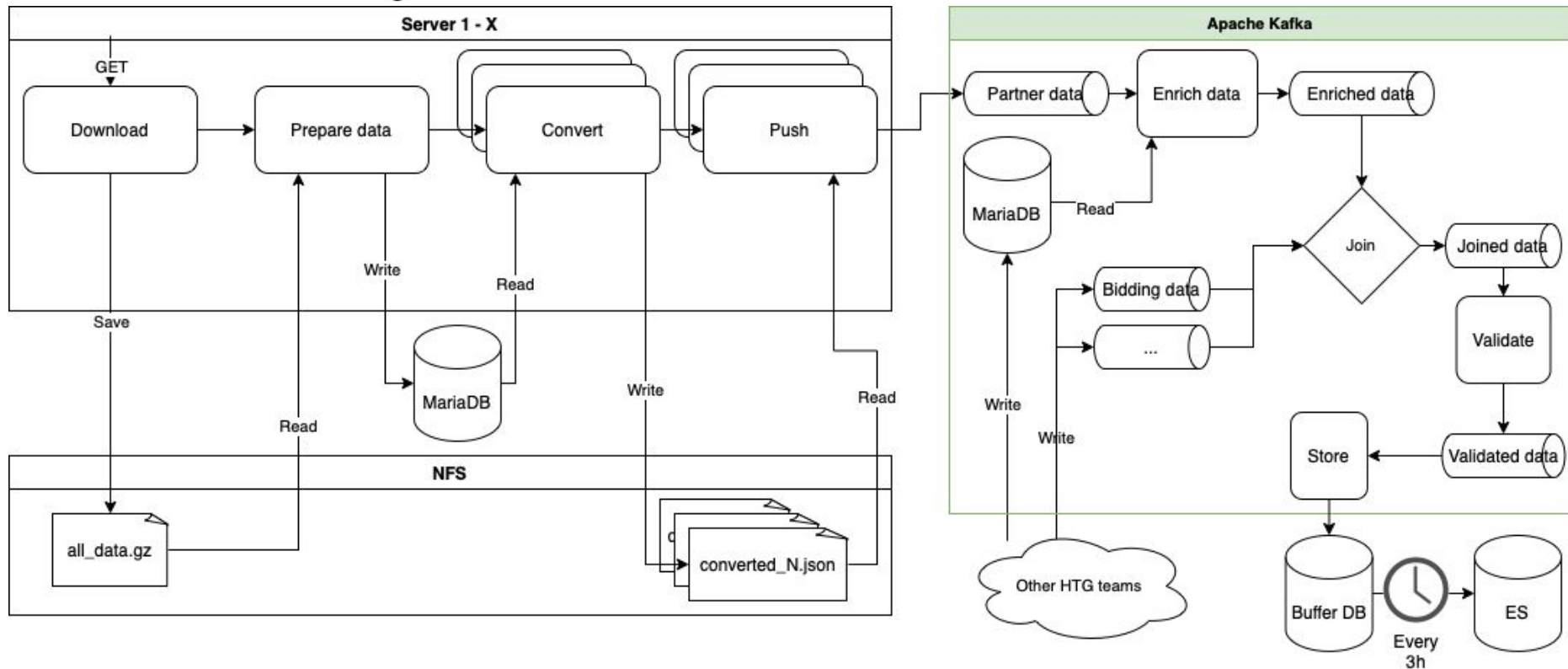


# N improvement

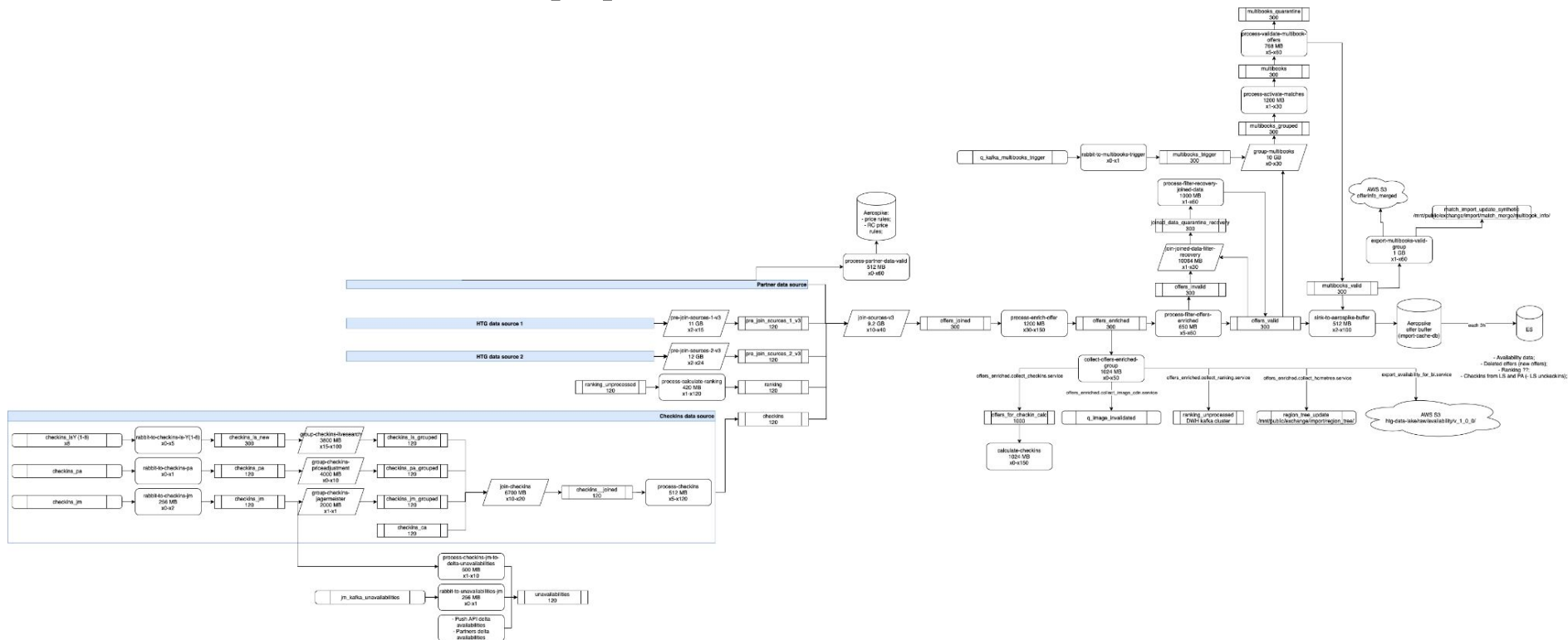
## Parallel data update

Bidding data  
Match and merge offers

# Start using kafka<sup>®</sup>



home  
to go.





# Production Kafka

- Today we have
  - 131 topics
  - 12 brokers
  - 37 consumer groups (scaling 1 – 300 processes)
  - ~6.5 billion messages per day

Rate	Mean	1 min	5 min	15 min
Messages in /sec	46k	32k	34k	40k
Bytes in /sec	482m	297m	344m	416m
Bytes out /sec	3.1b	2.3b	2.5b	3.1b

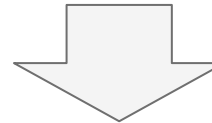
# N+1 improvement

## Dynamic resource pool

Increased user traffic  
Seasons

# Changing

- Mesos
  - Last release 24 Nov 2020 version 1.11.0
  - New OS don't have builded package
- Kubernetes
  - Fast auto scaling
  - Better resource utilization
  - Containerization



**kubernetes**

**You aren't gonna need it**

# YAGNI

- From Wiki
  - is a principle that states a programmer should not add functionality until deemed necessary
- From personal experience
  - Business changes a lot
  - We keep discovering new opportunities
  - Do what is needed only today with the opportunity to expand

# Let's connect

Darius Kasiulevičius

<https://www.linkedin.com/in/darius-kasiulevicius/>

