

# BUSINESS DASHBOARDS

using Bonobo, Airflow and Grafana

**Romain Dorgueil**

romain@makersquad.fr

Building software from Zero to Market

 **rdorgueil**



makersquad.fr

# Content

Intro. Product

1. **Plan.**



2. **Implement.**



3. **Visualize.**



4. **Monitor.**

Outro. References, Pointers

# DISCLAIMERS

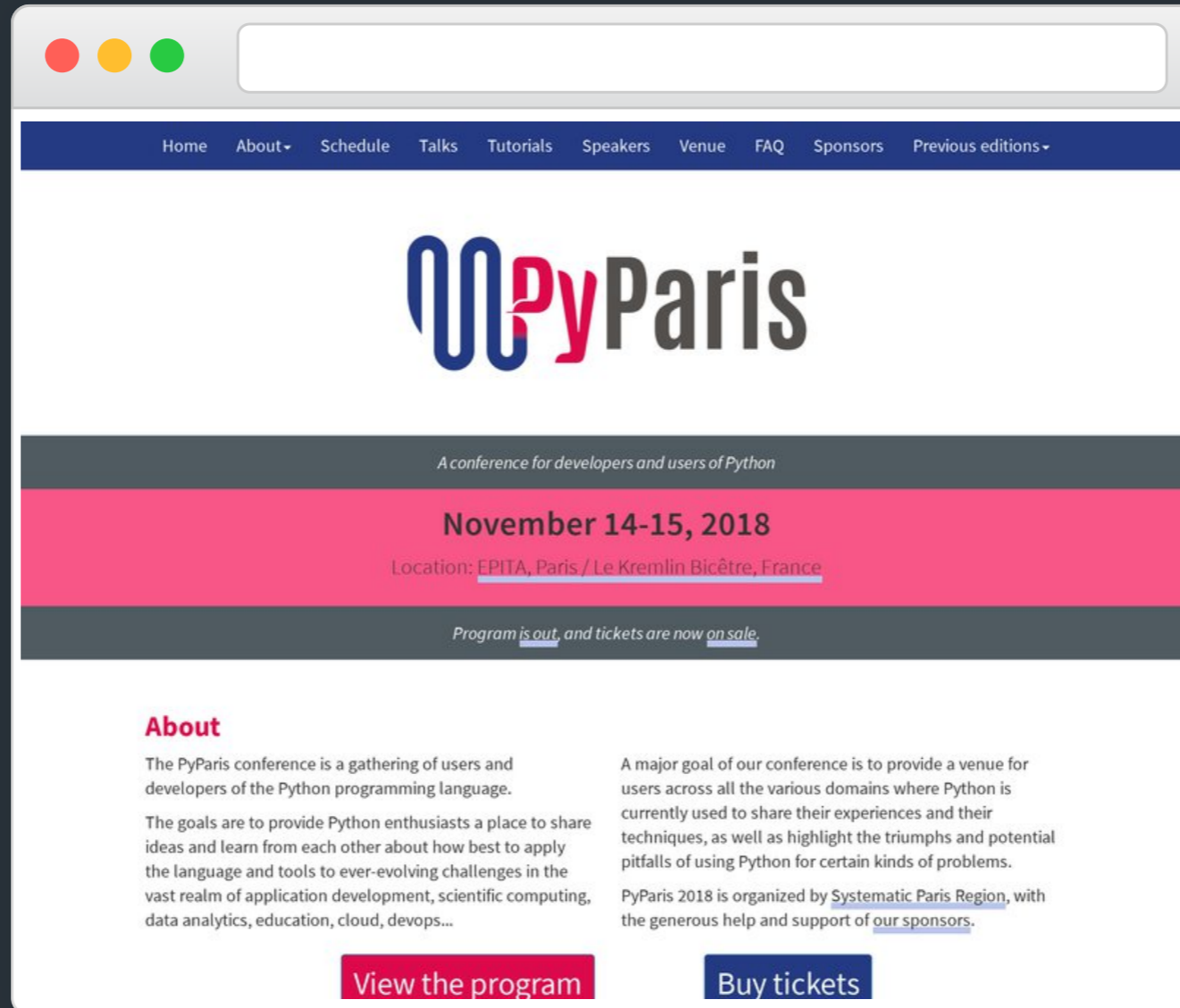
If you build a product, **do your own research.**  
Take time to learn and understand tools you consider.

I don't know nothing, and **I recommend nothing.**  
I assume things and change my mind when I hit a wall.

I'm the creator and main developer of Bonobo ETL.  
I'll try to be objective, but **there is a bias here.**

PRODUCT

GET `https://aprc.it/api/800x600/http://pyparis.org/`



The image shows a screenshot of a web browser displaying the PyParis website. The browser window has a white address bar and three colored window control buttons (red, yellow, green) on the left. The website has a dark blue navigation bar with links: Home, About, Schedule, Talks, Tutorials, Speakers, Venue, FAQ, Sponsors, and Previous editions. The main content area features the PyParis logo, which consists of a stylized 'Py' in blue and red followed by 'Paris' in black. Below the logo, there is a dark grey banner with the text 'A conference for developers and users of Python'. A pink banner below that contains the dates 'November 14-15, 2018' and the location 'Location: EPITA, Paris / Le Kremlin Bicêtre, France'. Another dark grey banner below that says 'Program is out, and tickets are now on sale.' The main content area is divided into two columns. The left column has a red heading 'About' and two paragraphs of text. The right column has two paragraphs of text. At the bottom, there are two buttons: 'View the program' in a pink box and 'Buy tickets' in a blue box.

Home About Schedule Talks Tutorials Speakers Venue FAQ Sponsors Previous editions

# PyParis

A conference for developers and users of Python

**November 14-15, 2018**  
Location: [EPITA, Paris / Le Kremlin Bicêtre, France](#)

Program [is out](#), and tickets are now [on sale](#).

## About

The PyParis conference is a gathering of users and developers of the Python programming language.

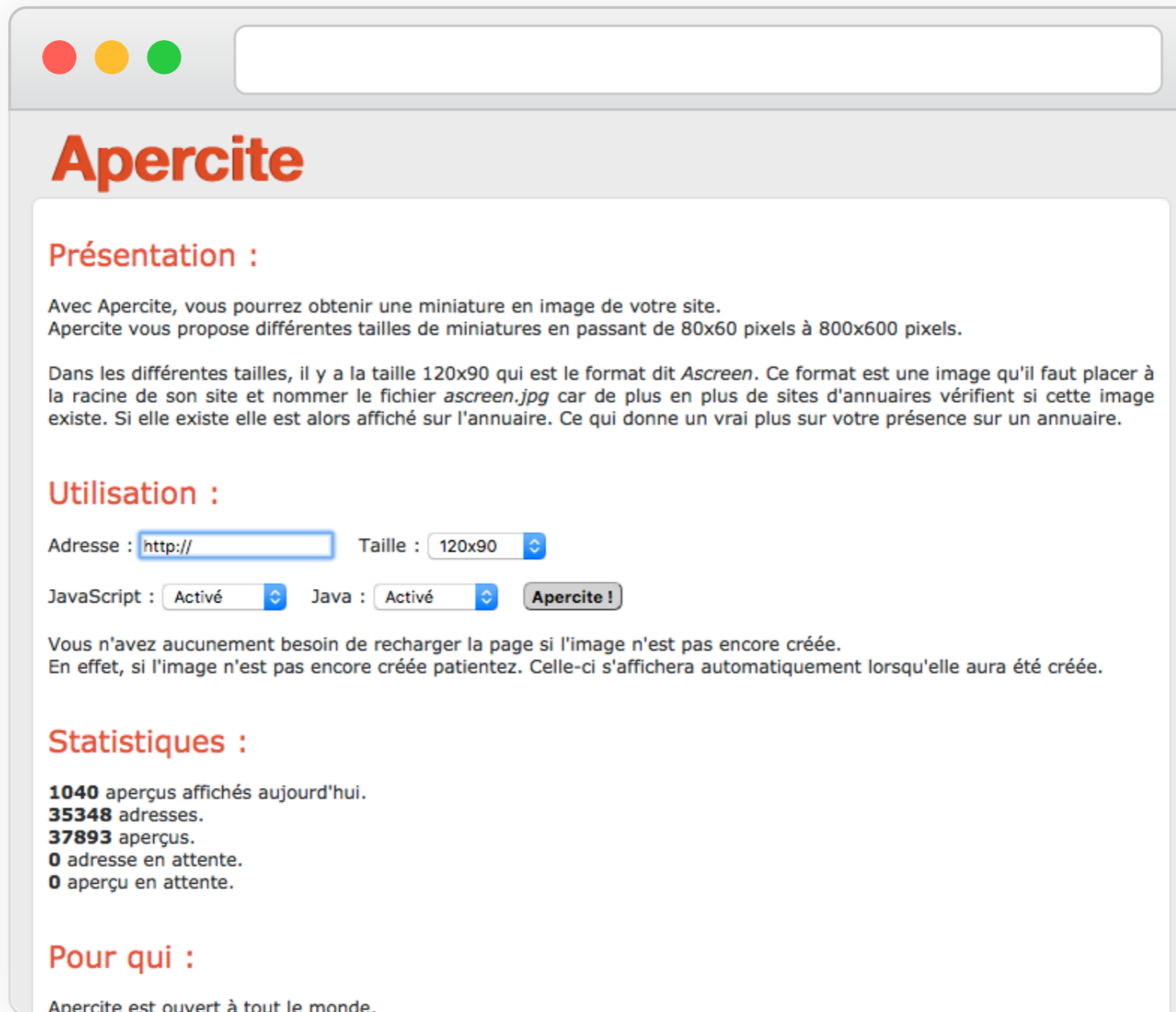
The goals are to provide Python enthusiasts a place to share ideas and learn from each other about how best to apply the language and tools to ever-evolving challenges in the vast realm of application development, scientific computing, data analytics, education, cloud, devops...

A major goal of our conference is to provide a venue for users across all the various domains where Python is currently used to share their experiences and their techniques, as well as highlight the triumphs and potential pitfalls of using Python for certain kinds of problems.

PyParis 2018 is organized by [Systematic Paris Region](#), with the generous help and support of [our sponsors](#).

[View the program](#) [Buy tickets](#)

# January 2009



The screenshot shows a web browser window with a title bar containing three colored buttons (red, yellow, green) and an empty address bar. The page content is as follows:

## Apercite

### Présentation :

Avec Apercite, vous pourrez obtenir une miniature en image de votre site.  
Apercite vous propose différentes tailles de miniatures en passant de 80x60 pixels à 800x600 pixels.

Dans les différentes tailles, il y a la taille 120x90 qui est le format dit *Ascreen*. Ce format est une image qu'il faut placer à la racine de son site et nommer le fichier *ascreen.jpg* car de plus en plus de sites d'annuaires vérifient si cette image existe. Si elle existe elle est alors affichée sur l'annuaire. Ce qui donne un vrai plus sur votre présence sur un annuaire.

### Utilisation :

Adresse :  Taille :

JavaScript :  Java :

Vous n'avez aucunement besoin de recharger la page si l'image n'est pas encore créée.  
En effet, si l'image n'est pas encore créée patientez. Celle-ci s'affichera automatiquement lorsqu'elle aura été créée.

### Statistiques :

**1040** aperçus affichés aujourd'hui.  
**35348** adresses.  
**37893** aperçus.  
**0** adresse en attente.  
**0** aperçu en attente.

### Pour qui :

Apercite est ouvert à tout le monde.

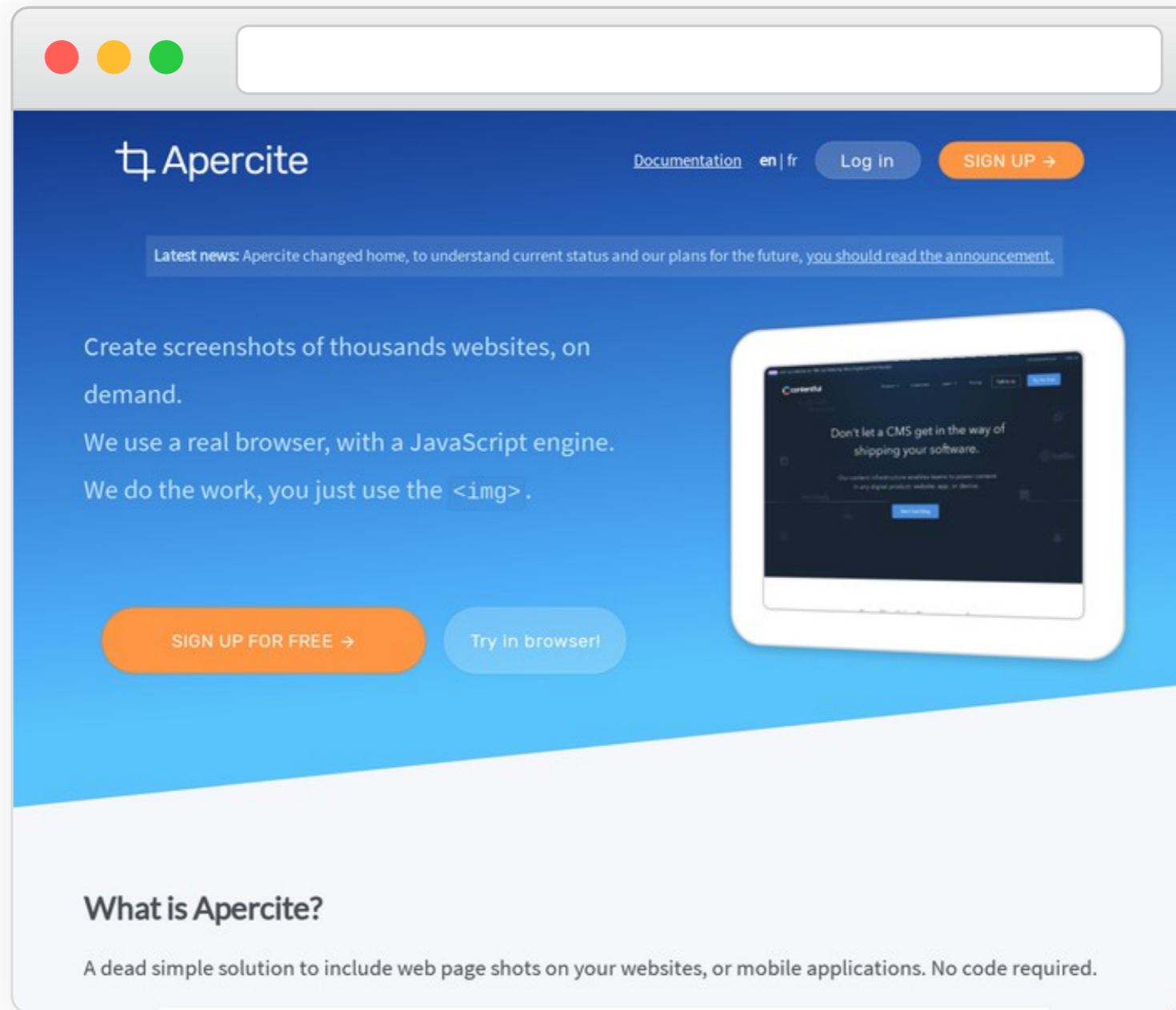




timedelta (years=9)



# April 2018



The image shows a browser window displaying the Apercite website. The page has a blue header with the Apercite logo, navigation links for 'Documentation', 'en | fr', 'Log In', and a 'SIGN UP →' button. A blue banner contains a 'Latest news' announcement. The main content area features a blue background with text describing the service: 'Create screenshots of thousands websites, on demand. We use a real browser, with a JavaScript engine. We do the work, you just use the <img>.' To the right is a tablet displaying a screenshot of a website. Below the text are two buttons: 'SIGN UP FOR FREE →' and 'Try in browser!'. The bottom section has a white background with the heading 'What is Apercite?' and the text 'A dead simple solution to include web page shots on your websites, or mobile applications. No code required.'

Apercite

Documentation en | fr Log In SIGN UP →

Latest news: Apercite changed home, to understand current status and our plans for the future, [you should read the announcement.](#)

Create screenshots of thousands websites, on demand.

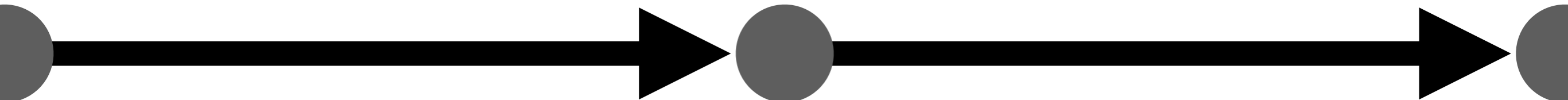
We use a real browser, with a JavaScript engine.

We do the work, you just use the `<img>`.

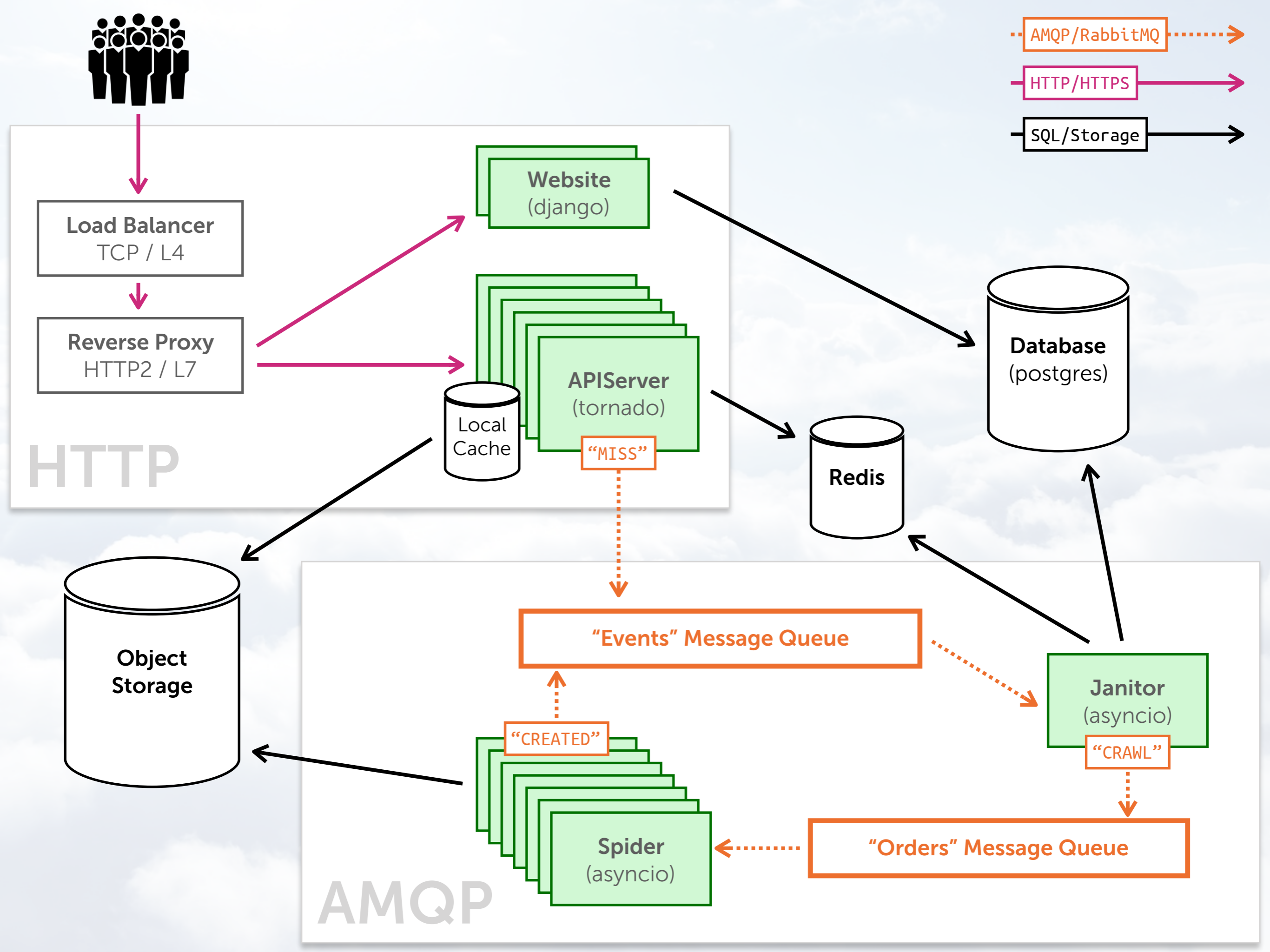
SIGN UP FOR FREE → Try in browser!

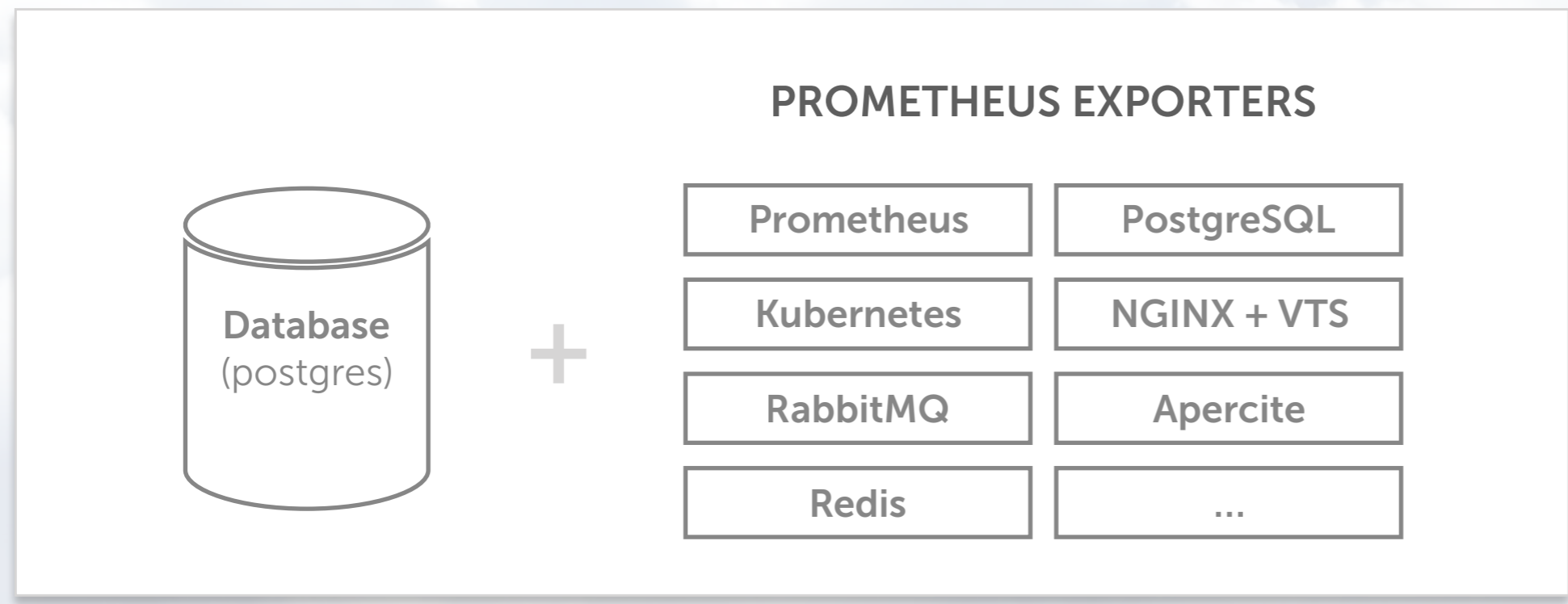
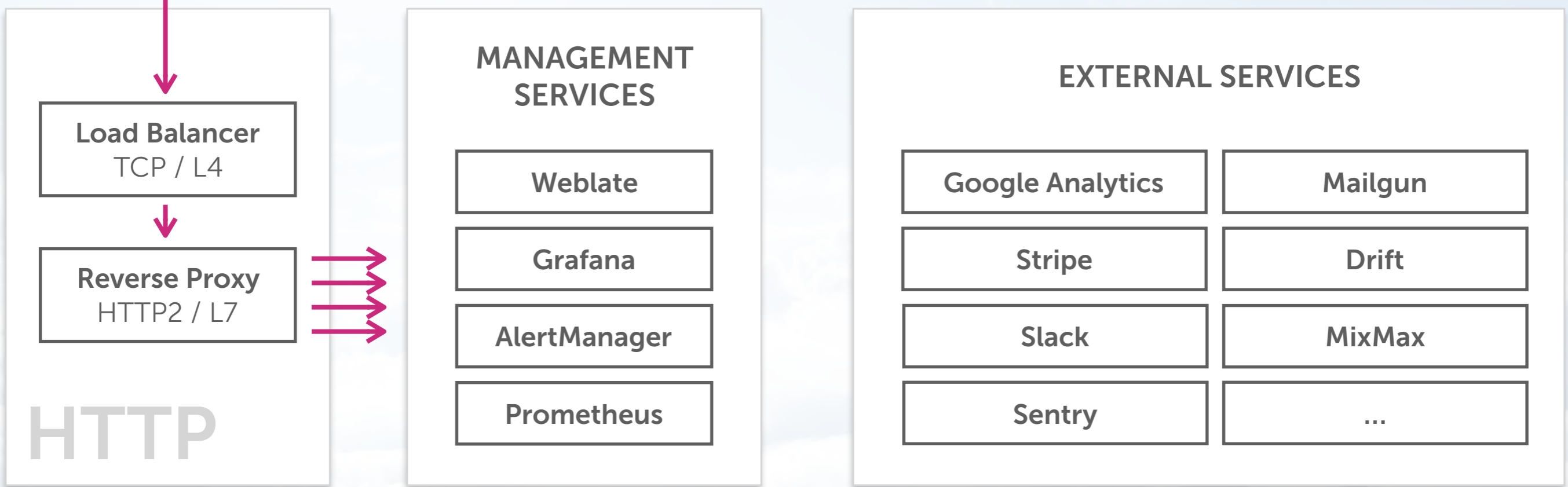
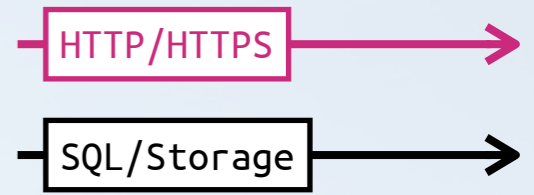
### What is Apercite?

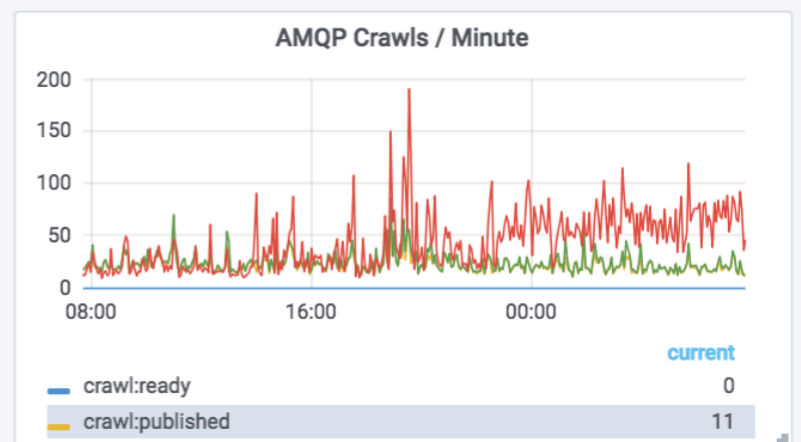
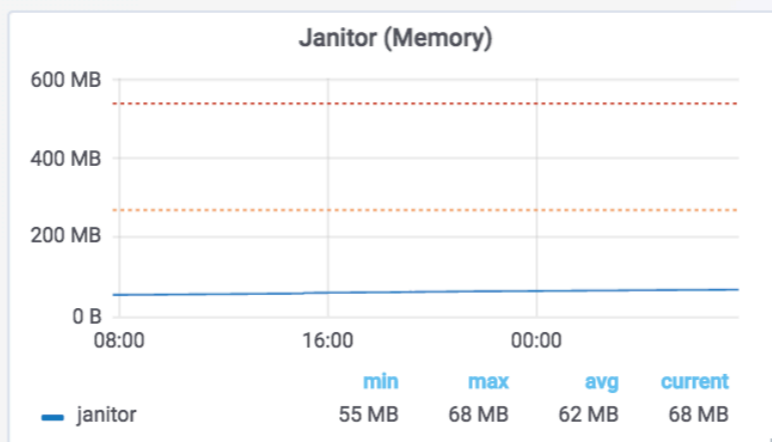
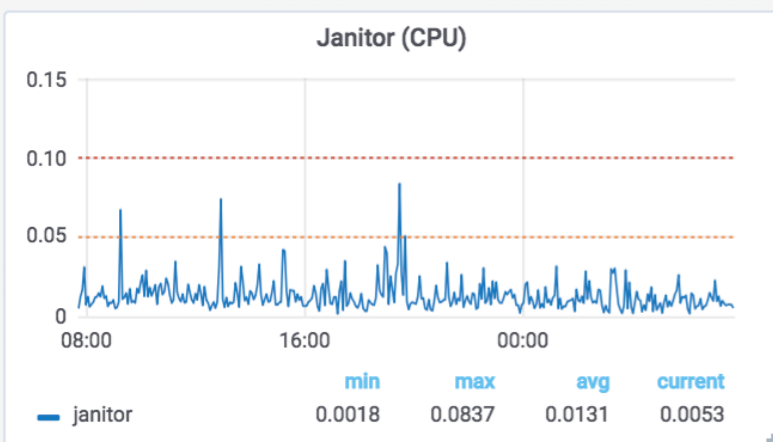
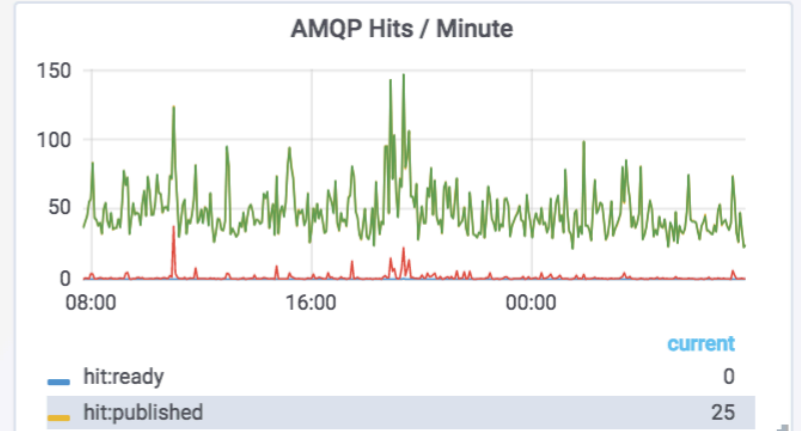
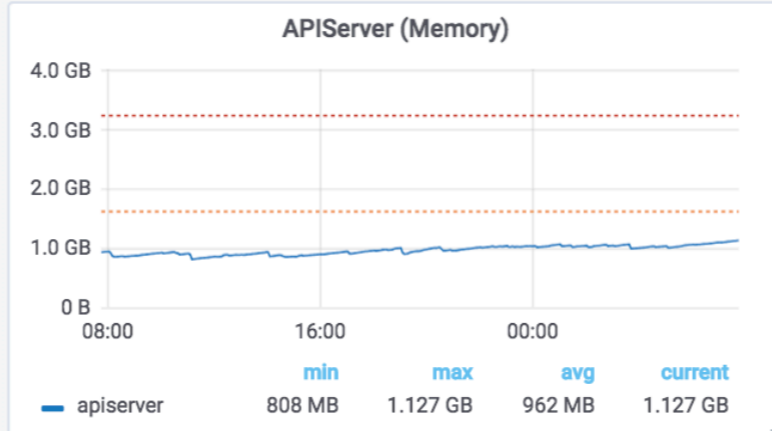
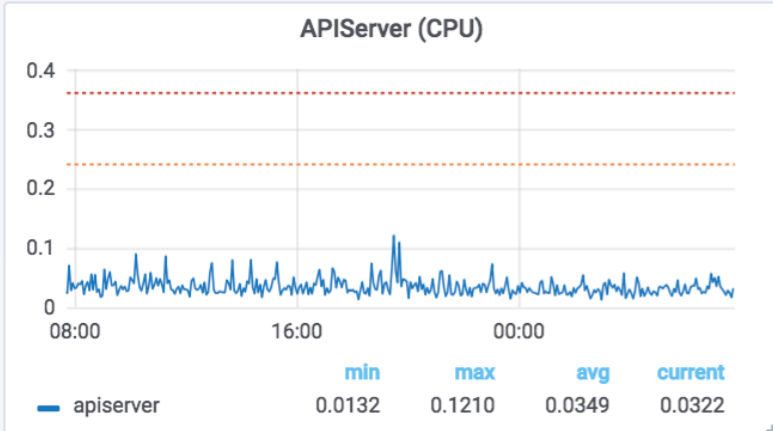
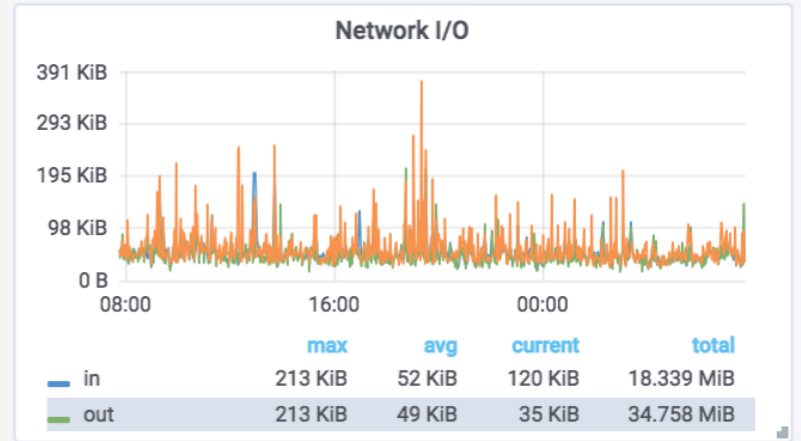
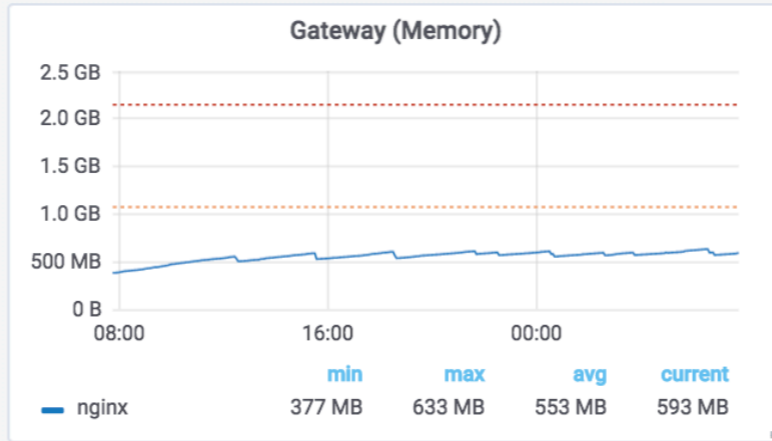
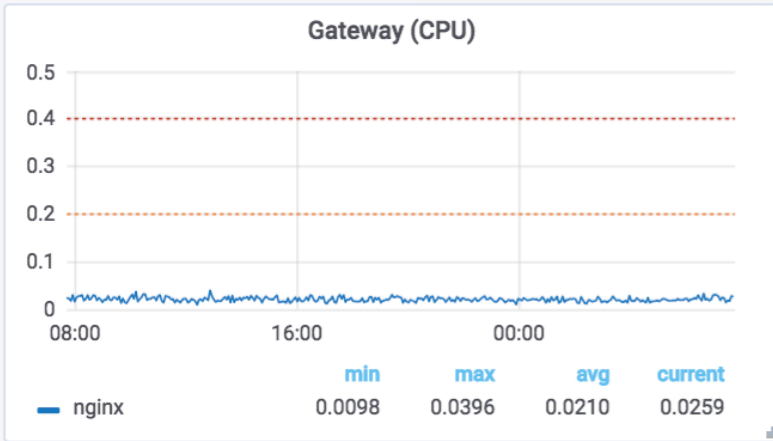
A dead simple solution to include web page shots on your websites, or mobile applications. No code required.

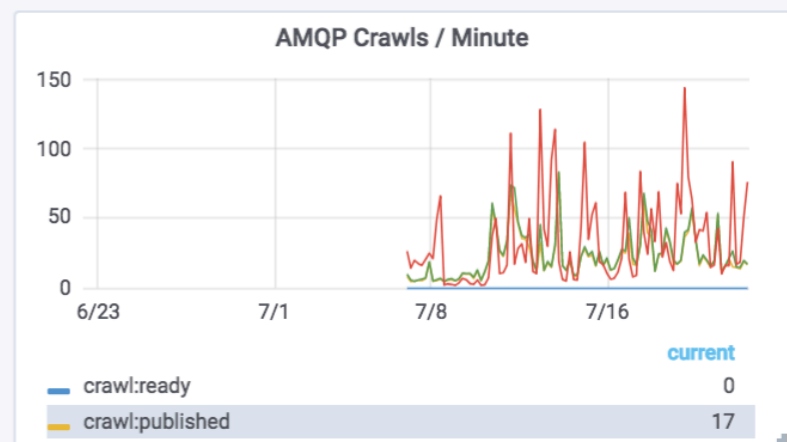
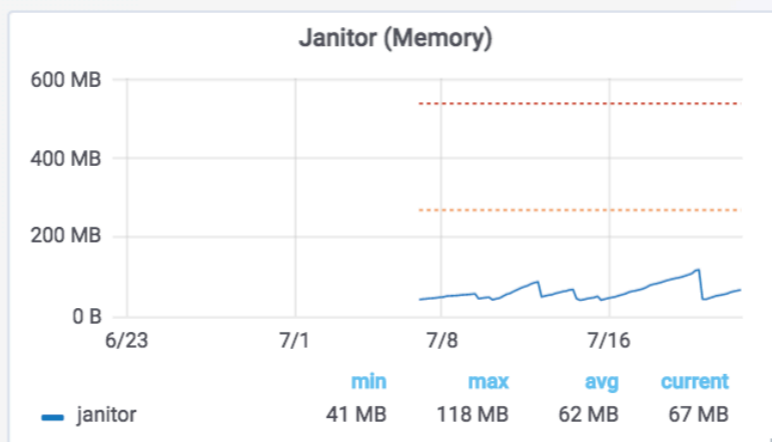
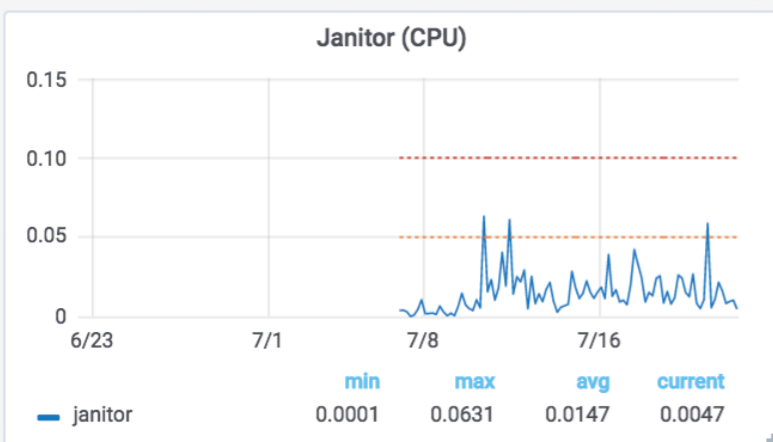
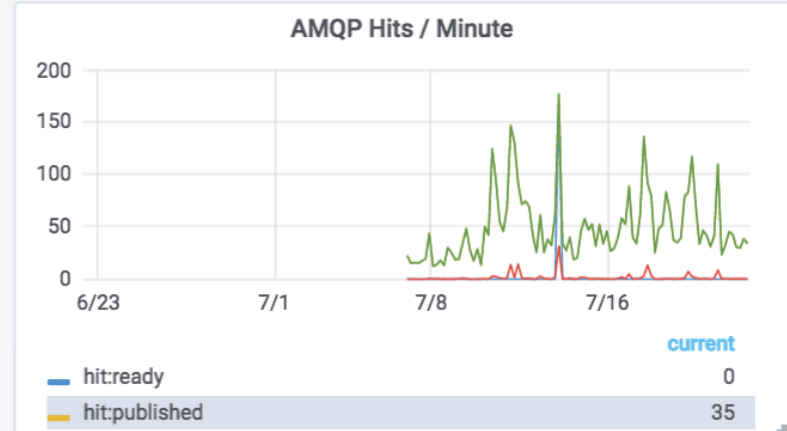
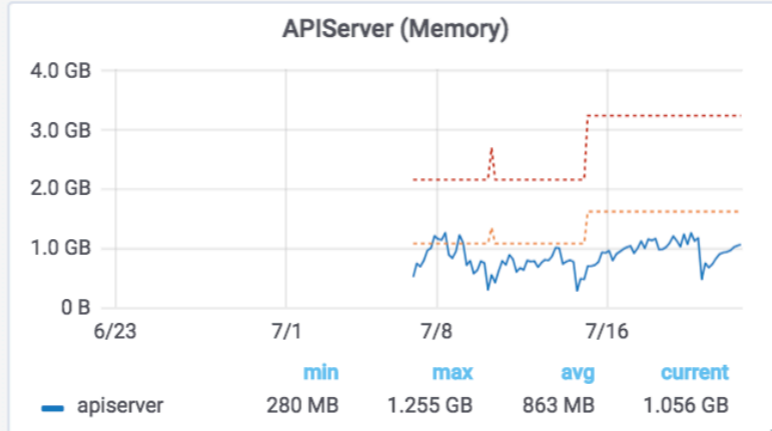
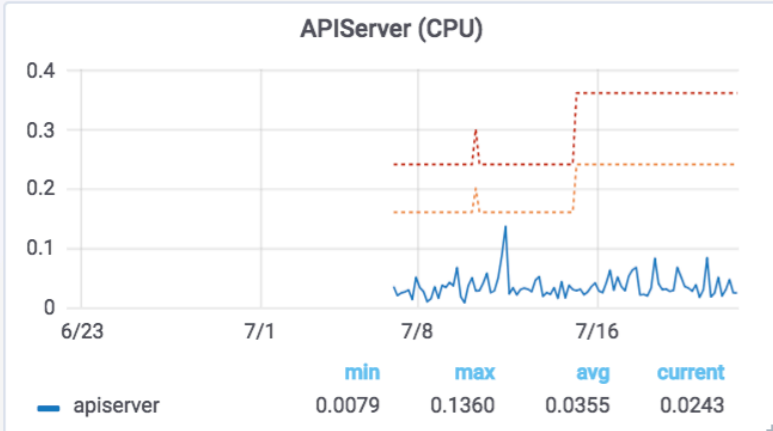
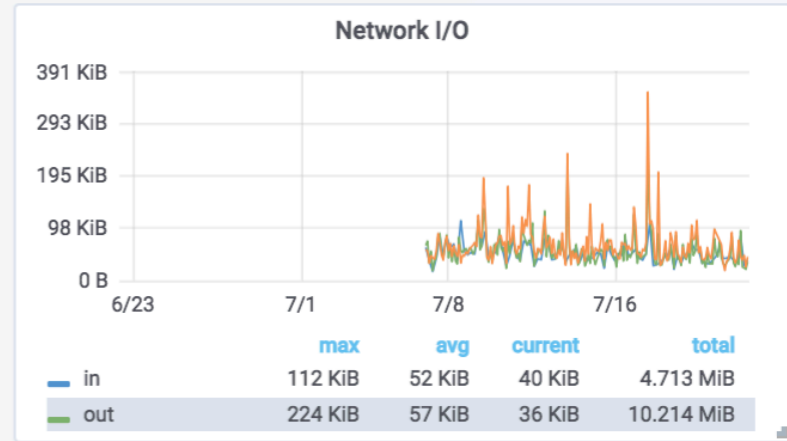
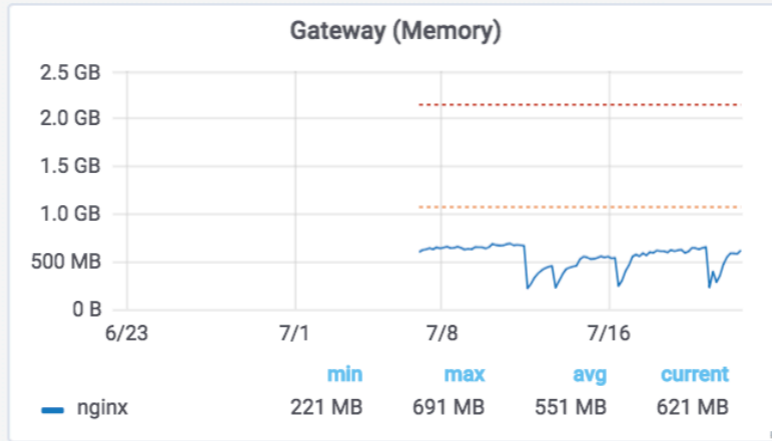
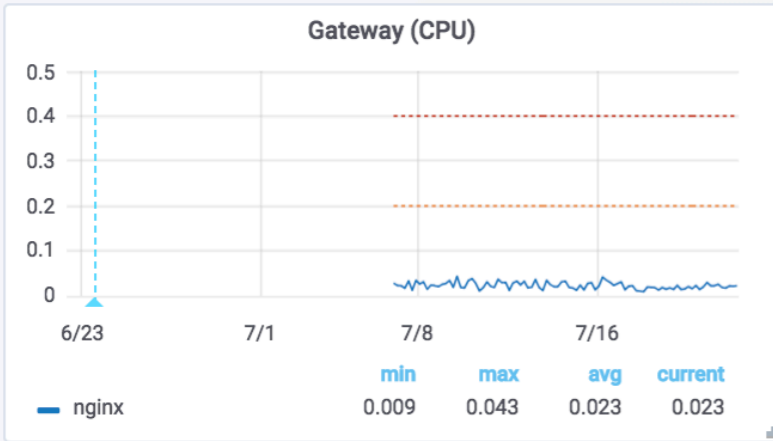


**under the hood ...**









PLAN

« What gets measured  
gets improved. »

— *Peter Drucker*



# Planning

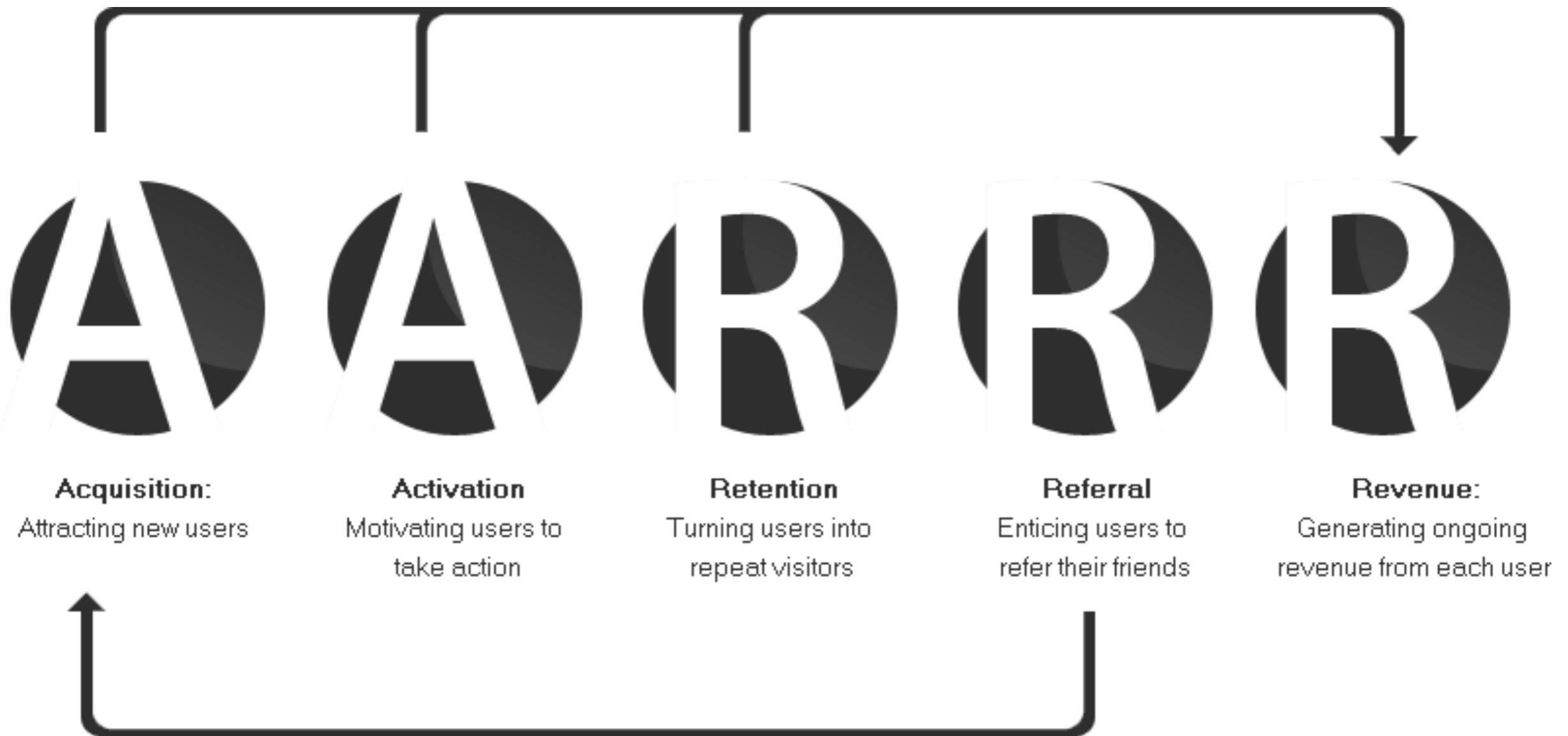
- Take your time to **choose metrics** wisely.
- **Cause vs Effect.**
- Less is More.
- **One at a time.** Or one by team.
- There is **not one answer** to this question.

Vanity metrics will  
**waste your time**

# Planning

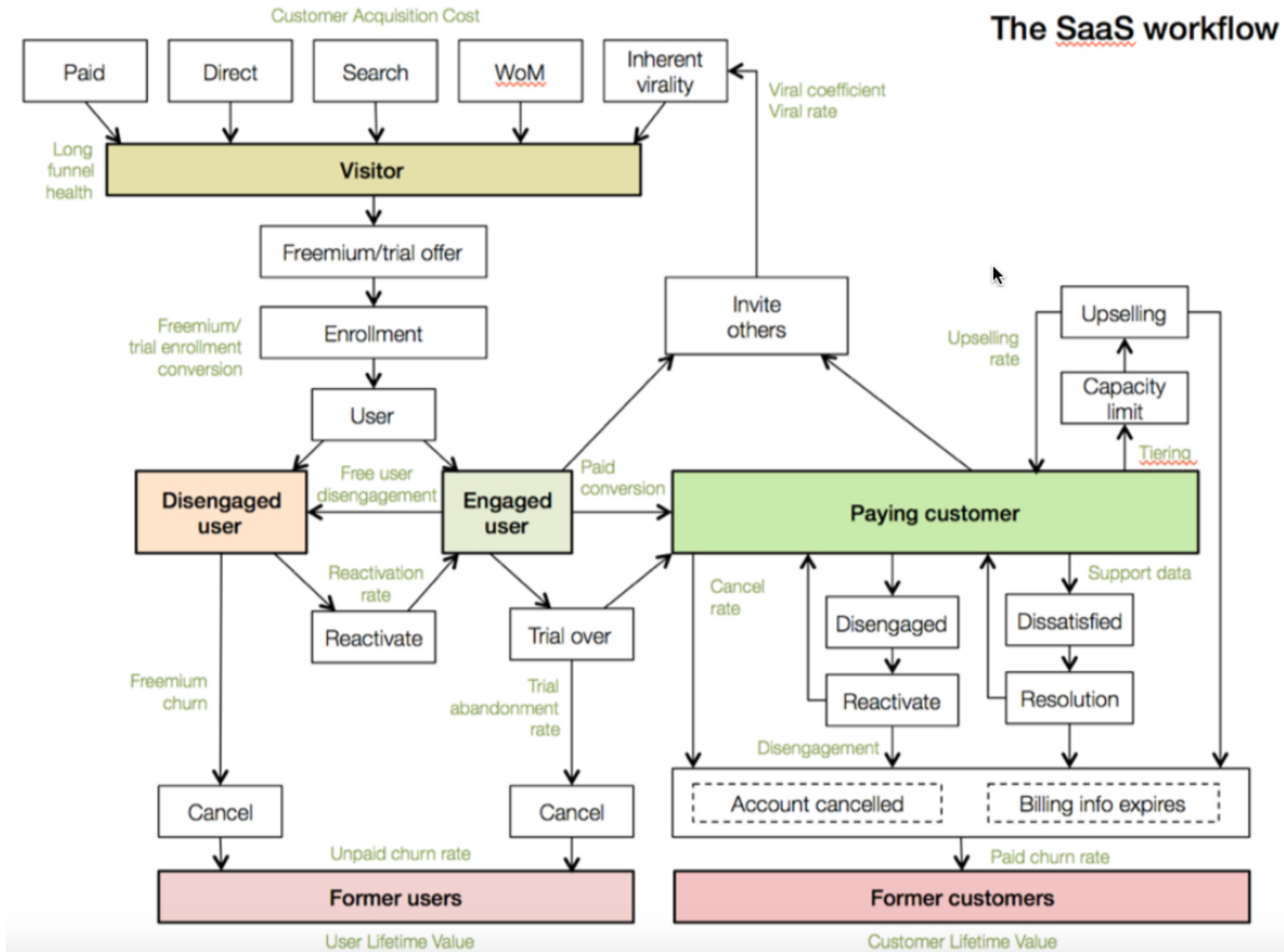
- Start with a **framework**.
- You may build your own, later.

# Pirate Metrics

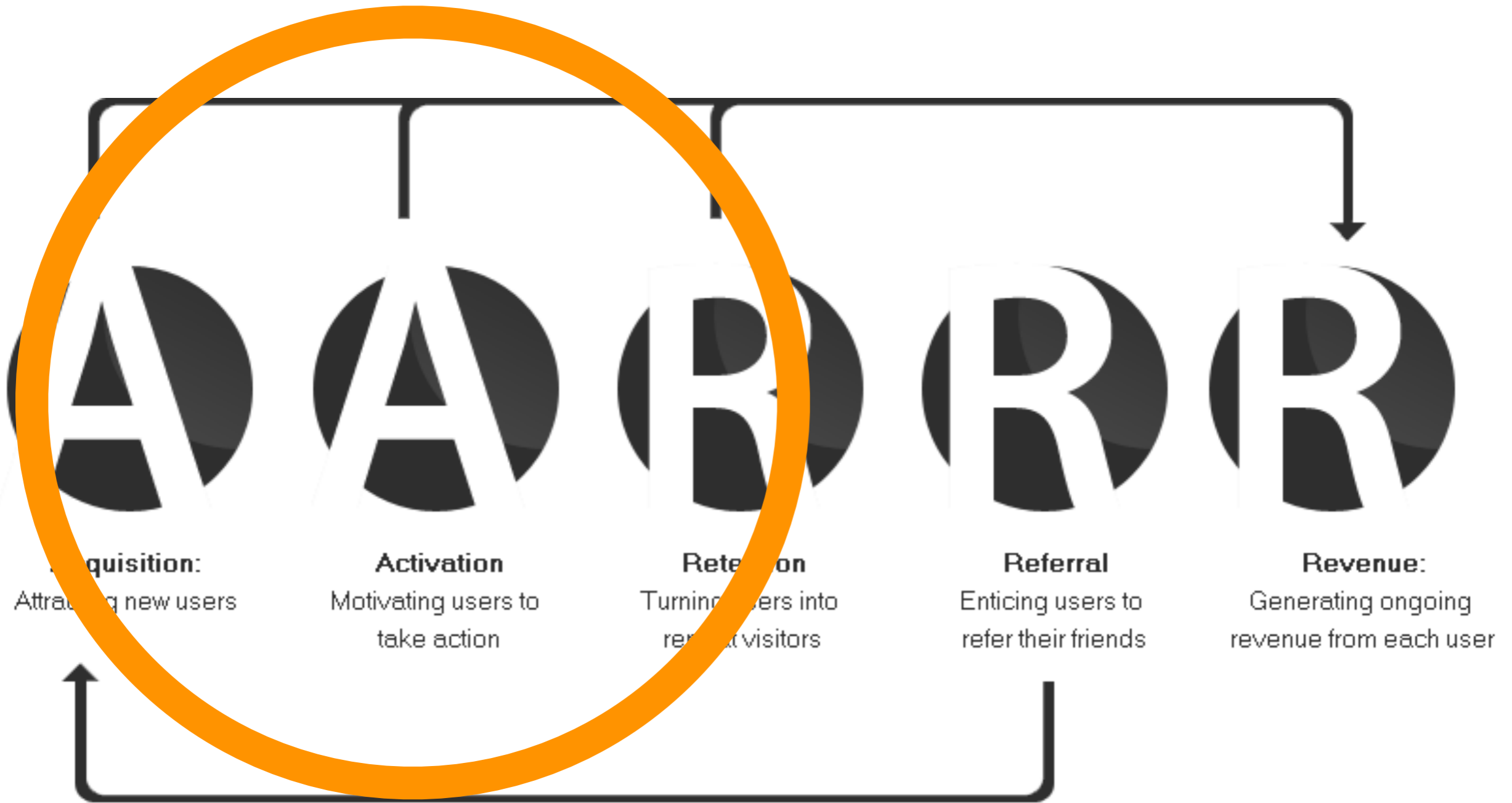


# Lean Analytics

(book by Alistair Croll & Benjamin Yoskovitz)



# Plan A



# Plan A

- What **business?** **Software as a Service**
- What **stage?** **Empathy / Stickyness**
- What **metric** matters?
  - Rate from acquisition to activation.
  - QOS (both for display and measure improvements).

IMPLEMENT

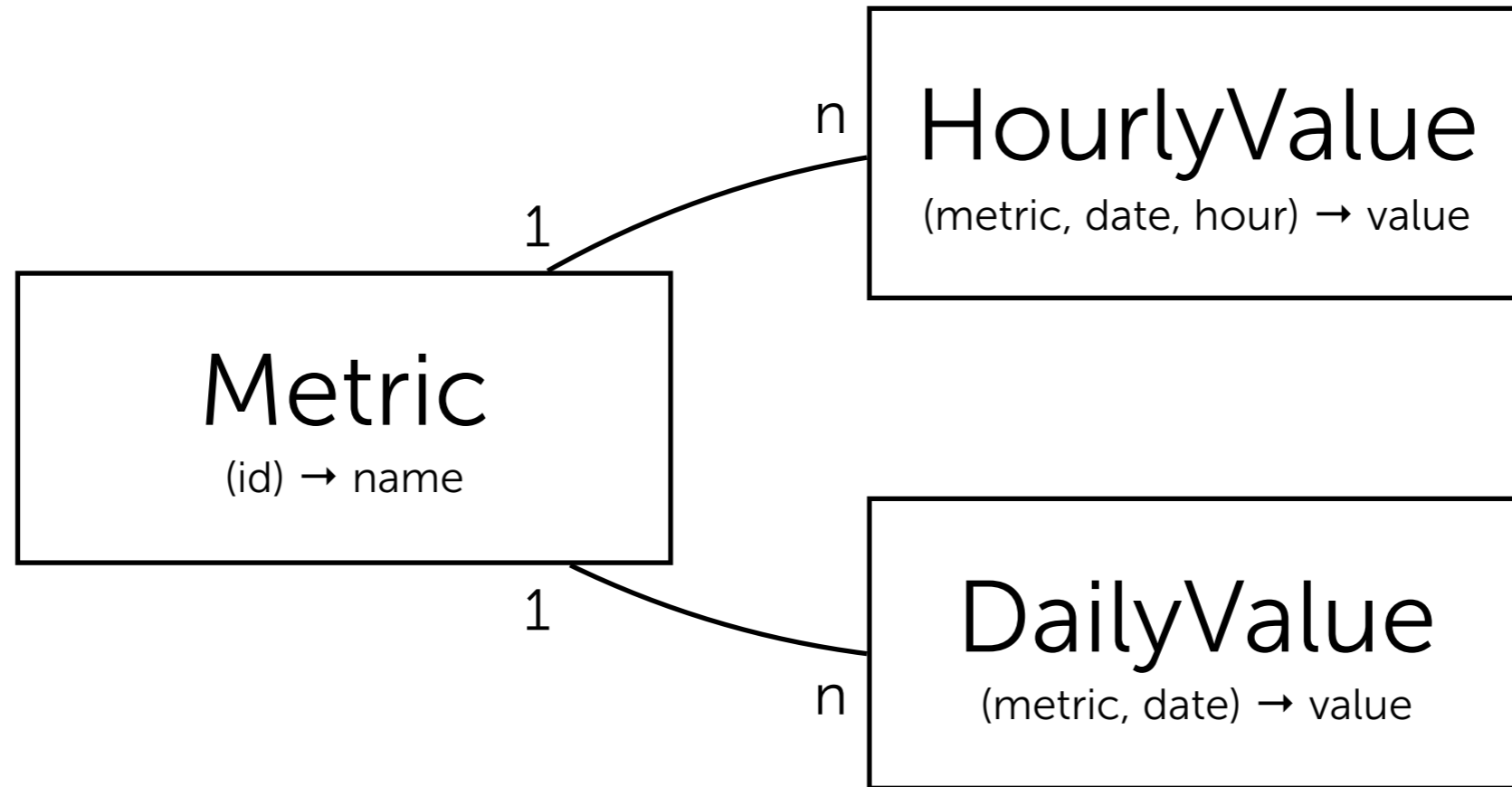




# Idea



# Model



Quick to write.

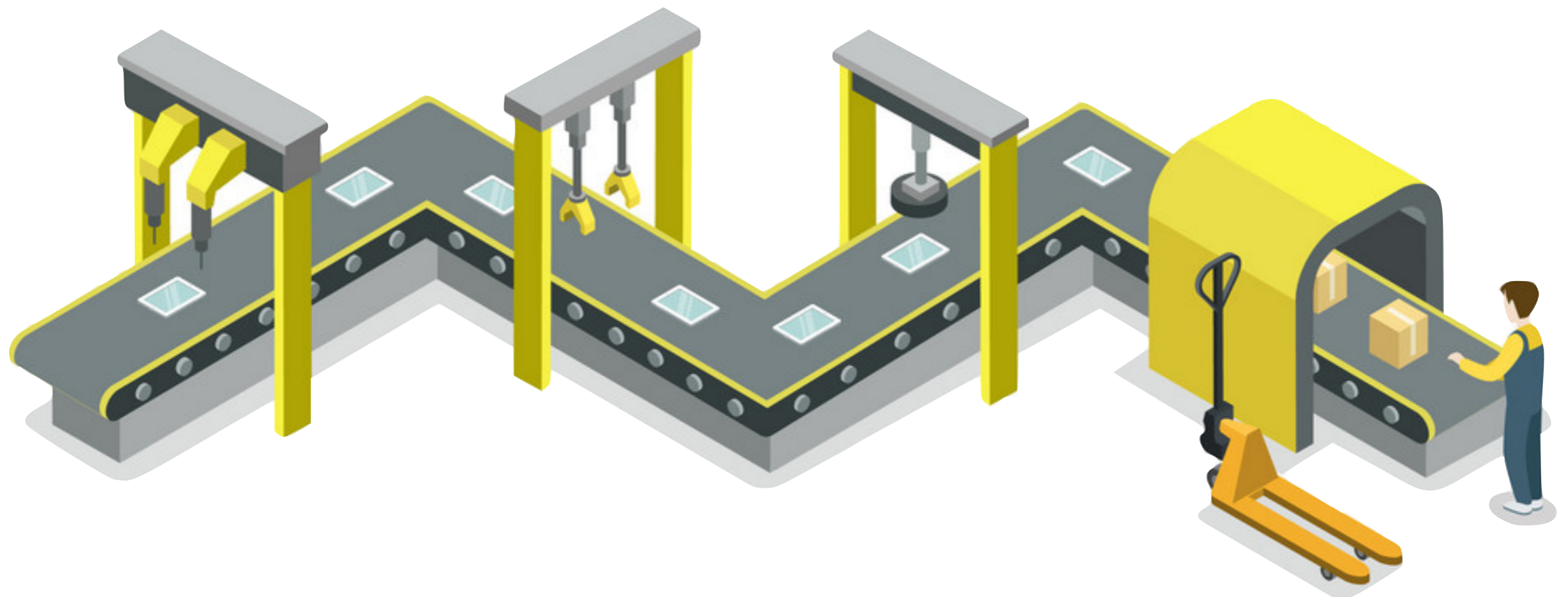
Not the best.

*Keywords to read more: Star and Snowflake Schemas*



# Bonobo

## Extract Transform Load





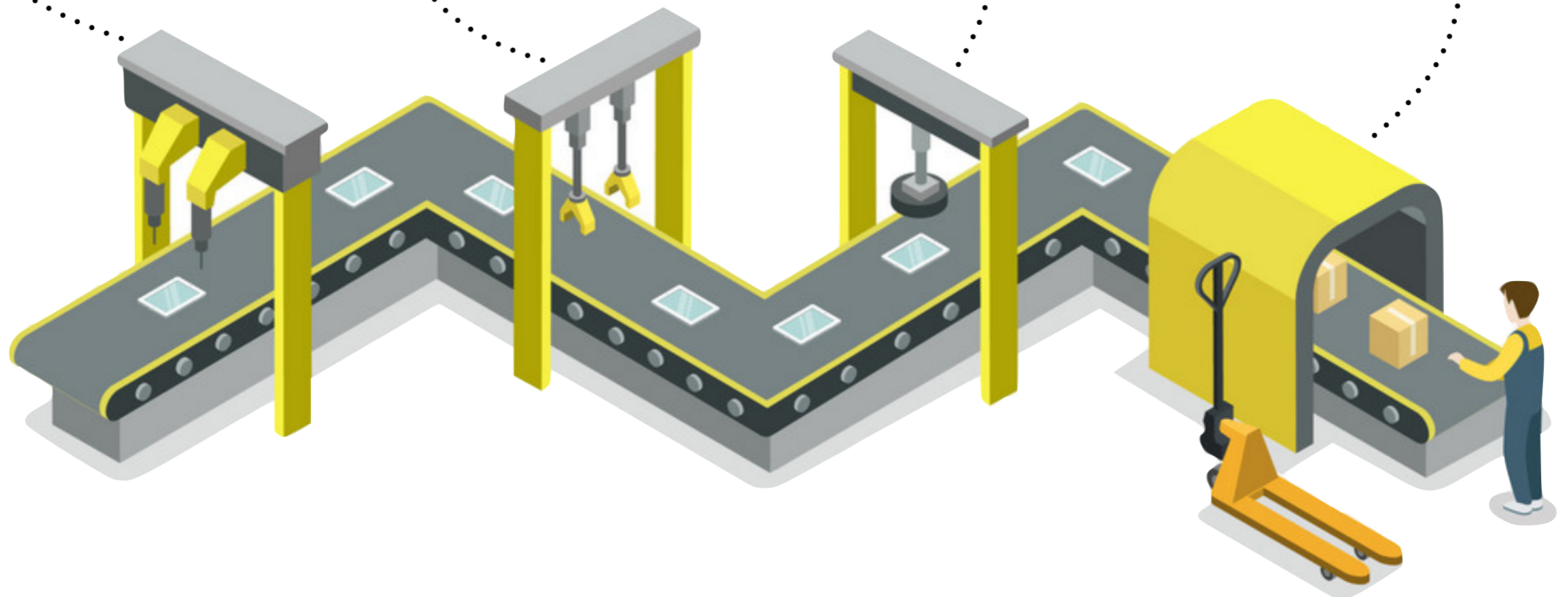
# Bonobo

```
Select('''  
SELECT *  
FROM ...  
WHERE ...  
...''')
```

```
def qualify(row):  
    yield (  
        row,  
        'active' if ...  
        else 'inactive'  
    )
```

```
Join('''  
SELECT count(*)  
FROM ...  
WHERE uid = %(id)s  
...''')
```

```
def report(row):  
    send_email(  
        render(  
            'email.html', row  
        )  
    )
```





# Bonobo

- Independent **threads**.
- Data is passed **first in, first out**.
- Supports any kind of **directed acyclic graphs**.
- Standard Python **callable and iterators**.
- Getting started still fits in here (ok, barely)

```
$ pip install bonobo  
$ bonobo init somejob.py  
$ python somejob.py
```

Let's write our jobs.

# Extract

... counts from website's database

```
from bonobo.config import use_context, Service
from bonobo_sqlalchemy.readers import Select

@use_context
class ObjectCountsReader(Select):
    engine = Service('website.engine')
    query = '''
        SELECT count(%(0)s.id) AS cnt
        FROM %(0)s
    ...
    output_fields = ['dims', 'metrics']

    def formatter(self, input_row, row):
        now = datetime.datetime.now()
        return ({
            'date': now.date(),
            'hour': now.hour,
        }, {
            'objects.{}.count'.format(input_row[1]): row['cnt']
        })
```



# Extract

... counts from website's database

```
TABLES_METRICS = {
    AsIs('apercite_account_user'): 'users',
    AsIs('apercite_account_userprofile'): 'profiles',
    AsIs('apercite_account_apikey'): 'apikeyes',
}

def get_readers():
    return [
        TABLES_METRICS.items(),
        ObjectCountsReader(),
    ]
```

# Normalize

All data should look the same

```
bonobo.SetFields(['dims', 'metrics'])
```

# Load

```
class AnalyticsWriter(InsertOrUpdate):
    dims = Option(required=True)
    filter = Option(required=True)

    @property
    def discriminant(self):
        return ('metric_id', *self.dims)

    def get_or_create_metrics(self, context, connection, metrics):
        ...

    def __call__(self, connection, table, buffer, context, row, engine):
        dims, metrics = row

        if not self.filter(dims, metrics): return

        # Get database rows for metric objects.
        db_metrics_ids = self.get_or_create_metrics(context, connection, metrics)

        # Insert or update values.
        for metric, value in metrics.items():
            yield from self._put(table, connection, buffer, {
                'metric_id': db_metrics_ids[metric],
                **{dim: dims[dim] for dim in self.dims},
                'value': value,
            })
```

# Compose

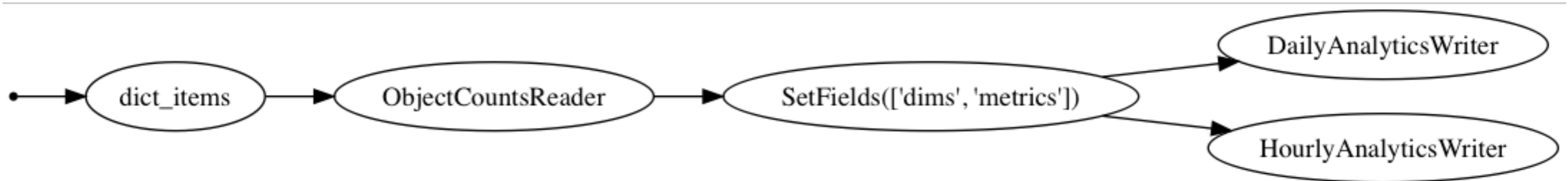
```
def get_graph():
    normalize = bonobo.SetFields(['dims', 'metrics'])
    graph = bonobo.Graph(*get_readers(), normalize)

    graph.add_chain(
        AnalyticsWriter(
            table_name=HourlyValue.__tablename__,
            dims=('date', 'hour'),
            filter=lambda dims, metrics: 'hour' in dims,
            name='Hourly',
        ), _input=normalize
    )
    graph.add_chain(
        AnalyticsWriter(
            table_name=DailyValue.__tablename__,
            dims=('date',),
            filter=lambda dims, metrics: 'hour' not in dims,
            name='Daily',
        ),
        _input=normalize
    )

    return graph
```

# Inspect

```
bonobo inspect --graph job.py | dot -o graph.png -T png
```



# Configure

```
def get_services():  
    return {  
        'sqlalchemy.engine': EventsDatabase().create_engine(),  
        'website.engine': WebsiteDatabase().create_engine(),  
    }
```

# Run

```
$ python -m apercite.analytics read objects --write  
- dict_items in=1 out=3 [done]  
- ObjectCountsReader in=3 out=3 [done]  
- SetFields(['dims', 'metrics']) in=3 out=3 [done]  
- HourlyAnalyticsWriter in=3 out=3 [done]  
- DailyAnalyticsWriter in=3 [done]
```

Got it.

Let's add readers.

*We'll run through, you'll have the code.*



# Google Analytics

```
@use('google_analytics')
def read_analytics(google_analytics):
    reports = google_analytics.reports().batchGet(
        body={...}
    ).execute().get('reports', [])

    for report in reports:
        dimensions = report['columnHeader']['dimensions']
        metrics = report['columnHeader']['metricHeader']['metricHeaderEntries']
        rows = report['data']['rows']

        for row in rows:
            dim_values = zip(dimensions, row['dimensions'])
            yield (
                {
                    GOOGLE_ANALYTICS_DIMENSIONS.get(dim, [dim])[0]:
                    GOOGLE_ANALYTICS_DIMENSIONS.get(dim, [None, IDENTITY])[1](val)
                    for dim, val in dim_values
                },
                {
                    GOOGLE_ANALYTICS_METRICS.get(metric['name'], metric['name']):
                    GOOGLE_ANALYTICS_TYPES[metric['type']](value)
                    for metric, value in zip(metrics, row['metrics'][0]['values'])
                }
            )
    )
```

# Prometheus

```
class PrometheusReader(Configurable):
    http = Service('http')
    endpoint = 'http://{host}:{port}/api/v1'.format(PROMETHEUS_HOST, PROMETHEUS_PORT)
    queries = [...]

def __call__(self, *, http):
    start_at, end_at = self.get_timerange()

    for query in self.queries:
        for result in http.get(...).json().get('data', {}).get('result', []):
            metric = result.get('metric', {})
            for ts, val in result.get('values', []):
                name = query.target.format(**metric)
                _date, _hour = ...
                yield {
                    'date': _date,
                    'hour': _hour,
                }, {
                    name: float(val)
                }
```

# Spider counts

```
class SpidersReader(Select):
    kwargs = Option()
    output_fields = ['row']

    @property
    def query(self):
        return '''
            SELECT spider.value AS name,
                   spider.created_at AS created_at,
                   spider_status.attributes AS attributes,
                   spider_status.created_at AS updated_at
            FROM
                spider
            JOIN ...
            WHERE spider_status.created_at > %(now)s
            ORDER BY spider_status.created_at DESC
            ...

    def formatter(self, input_row, row):
        return (row, )
```

# Spider counts

```
def spider_reducer(self, left, right):  
    result = dict(left)  
  
    result['spider.total'] += len(right.attributes)  
  
    for worker in right.attributes:  
        if 'stage' in worker:  
            result['spider.active'] += 1  
        else:  
            result['spider.idle'] += 1  
  
    return result
```

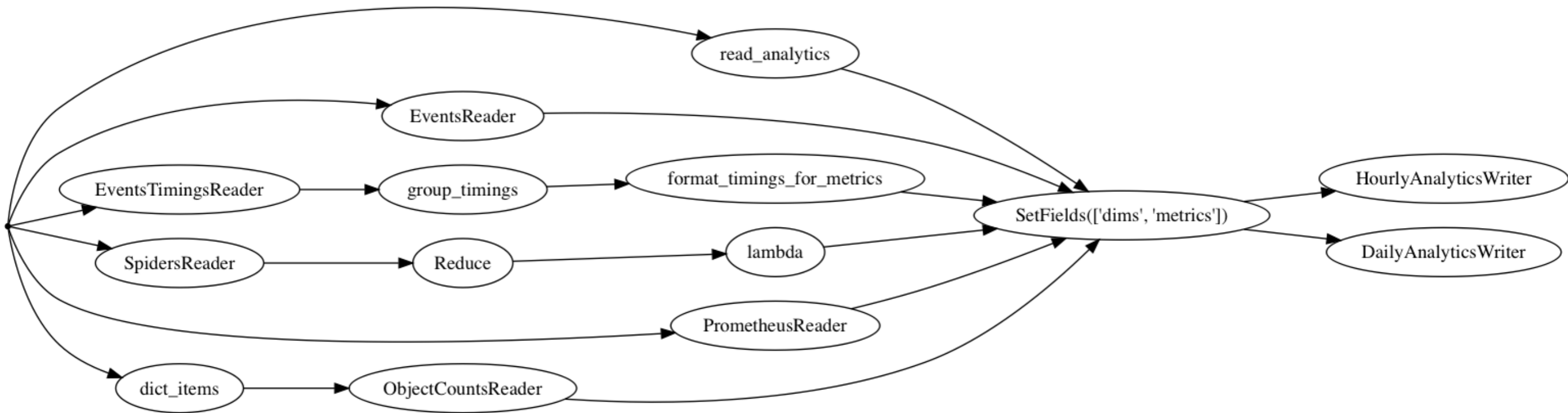
# Spider counts

```
now = datetime.datetime.utcnow() - datetime.timedelta(minutes=30)
```

```
def get_readers():  
    return (  
        SpidersReader(kwargs={'now': now}),  
        Reduce(spider_reducer, initializer={  
            'spider.idle': 0,  
            'spider.active': 0,  
            'spider.total': 0,  
        }),  
        (lambda x: ({'date': now.date(), 'hour': now.hour}, x))  
    )
```

etc.

# Inspect



We can generate ETL graphs with all readers or only a few.

# Run

```
$ python -m apercite.analytics read all --write
```

```
- read_analytics in=1 out=91 [done]  
- EventsReader in=1 out=27 [done]  
- EventsTimingsReader in=1 out=2039 [done]  
- group_timings in=2039 out=24 [done]  
- format_timings_for_metrics in=24 out=24 [done]  
- SpidersReader in=1 out=1 [done]  
- Reduce in=1 out=1 [done]  
- <lambda> in=1 out=1 [done]  
- PrometheusReader in=1 out=3274 [done]  
- dict_items in=1 out=3 [done]  
- ObjectCountsReader in=3 out=3 [done]  
- SetFields(['dims', 'metrics']) in=3420 out=3420 [done]  
- HourlyAnalyticsWriter in=3420 out=3562 [done]  
- DailyAnalyticsWriter in=3420 out=182 [done]
```



Easy to **build**.

Easy to **add** or **replace** parts.

Easy to **run**.

*Told ya, slight bias.*

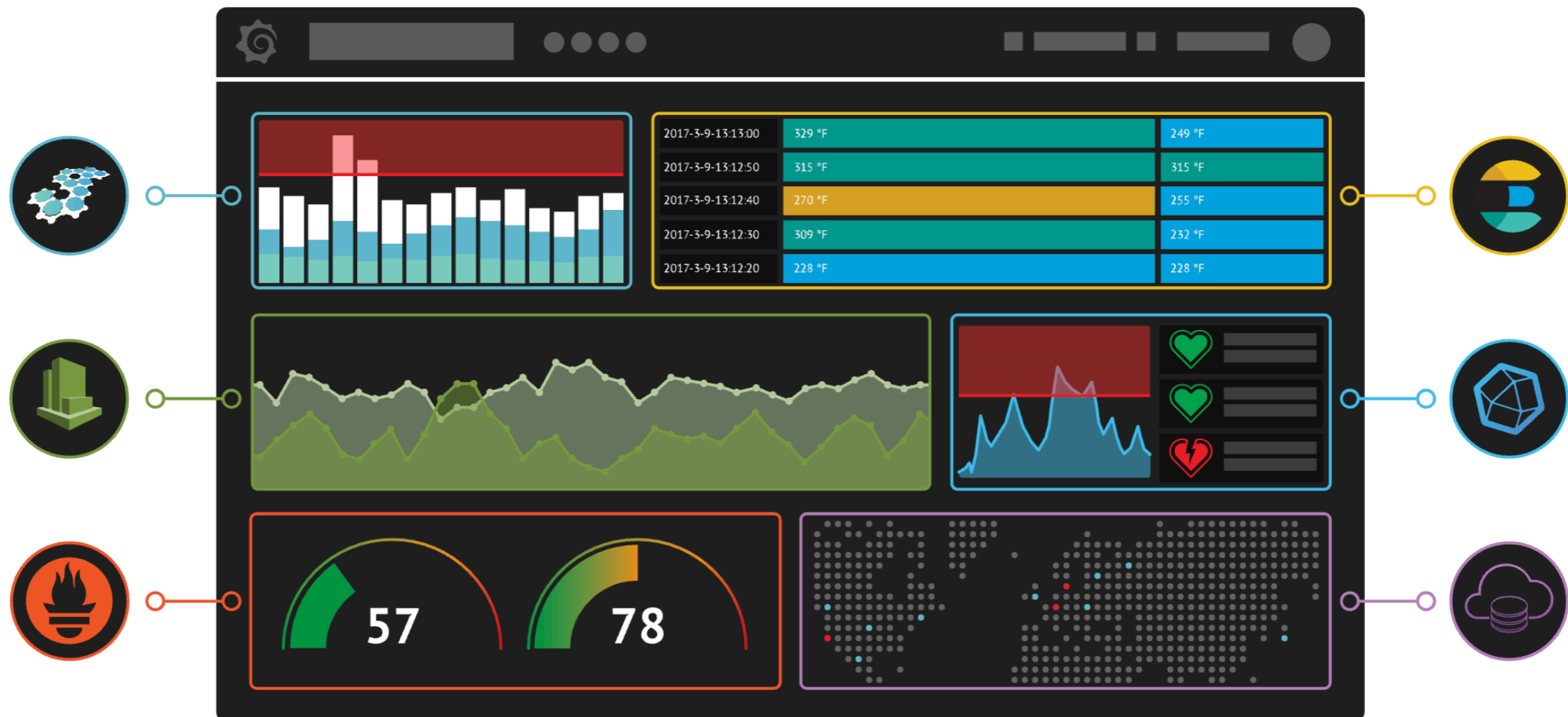
VISUALIZE



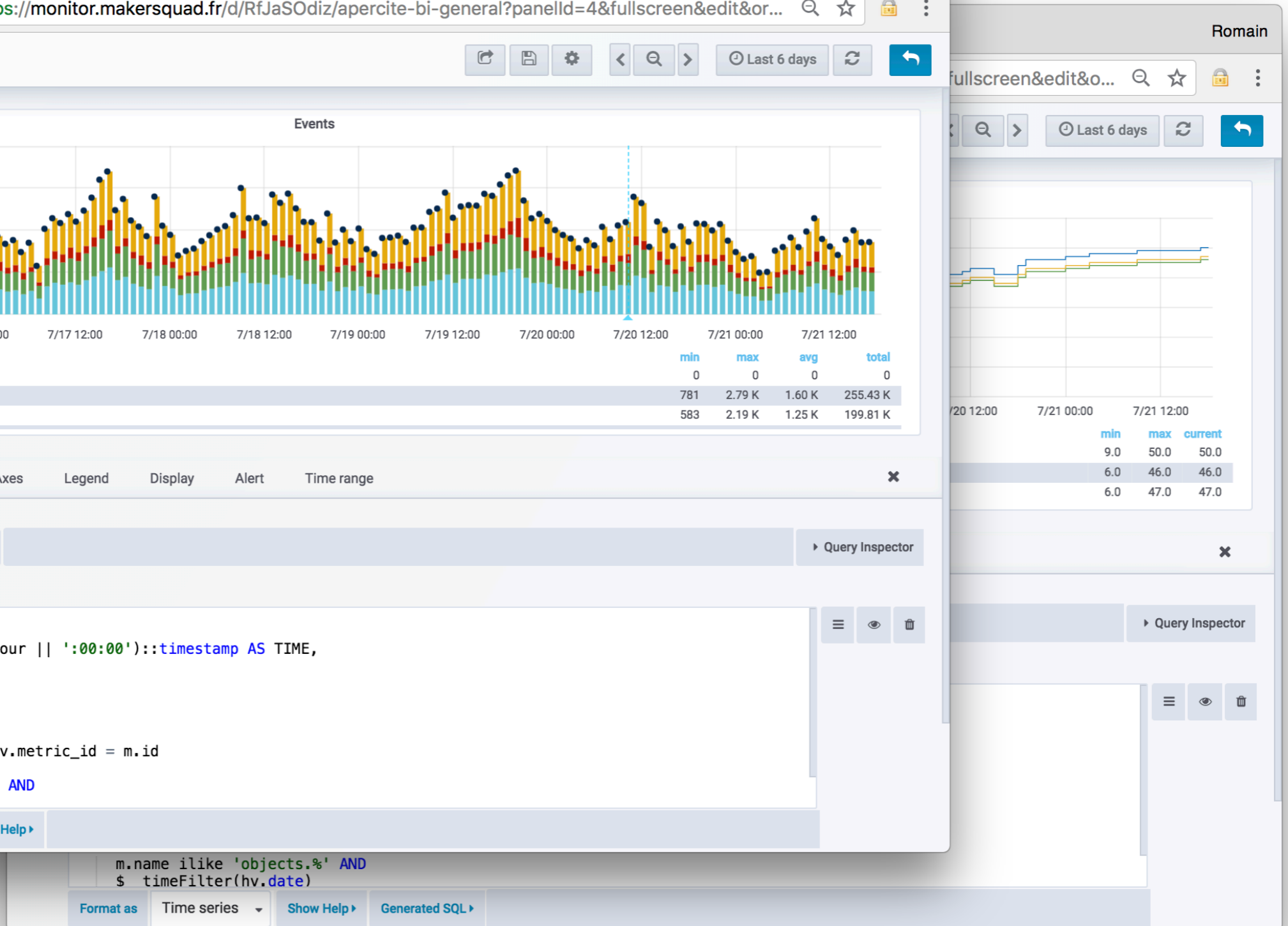
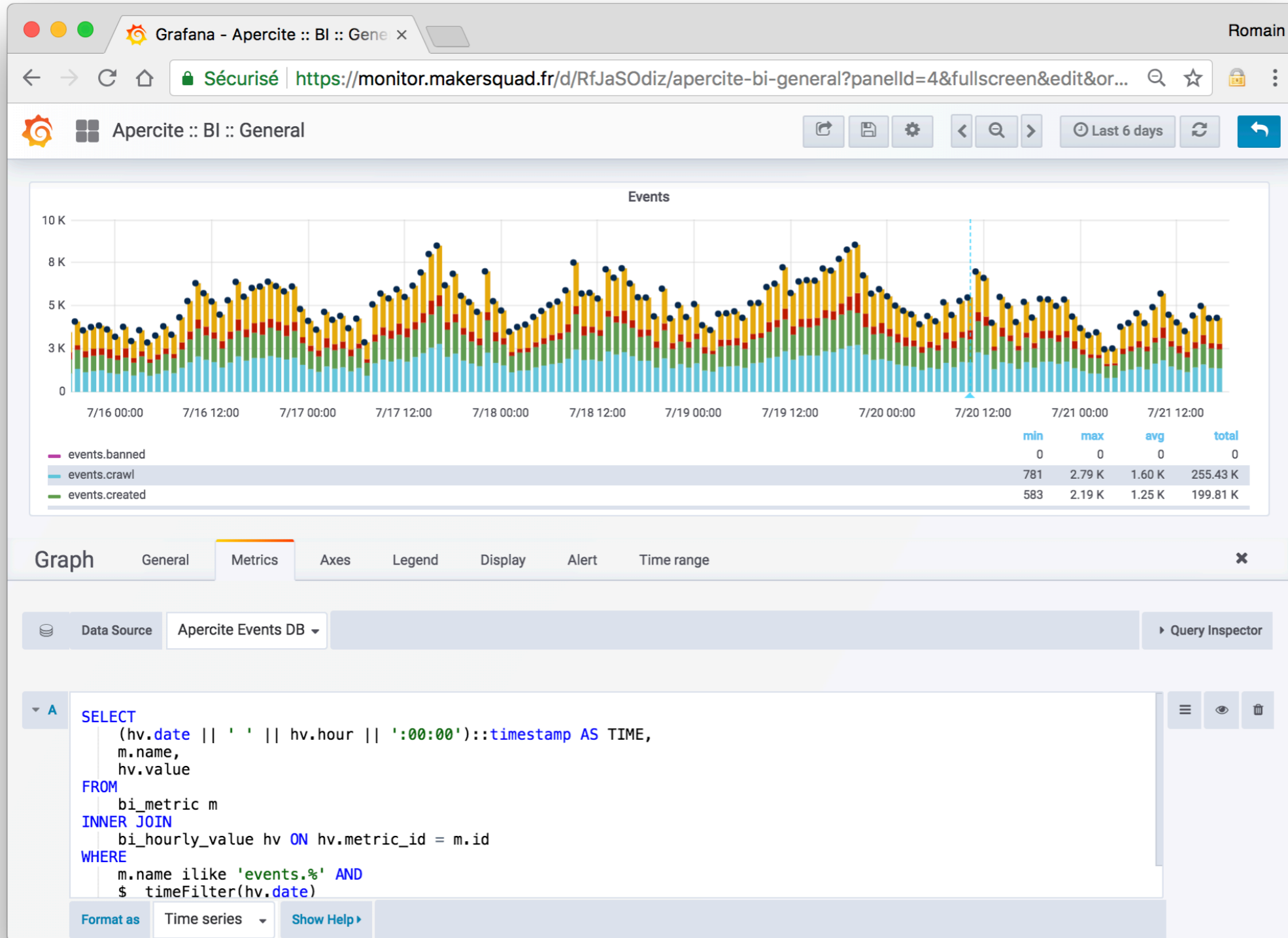


# Grafana

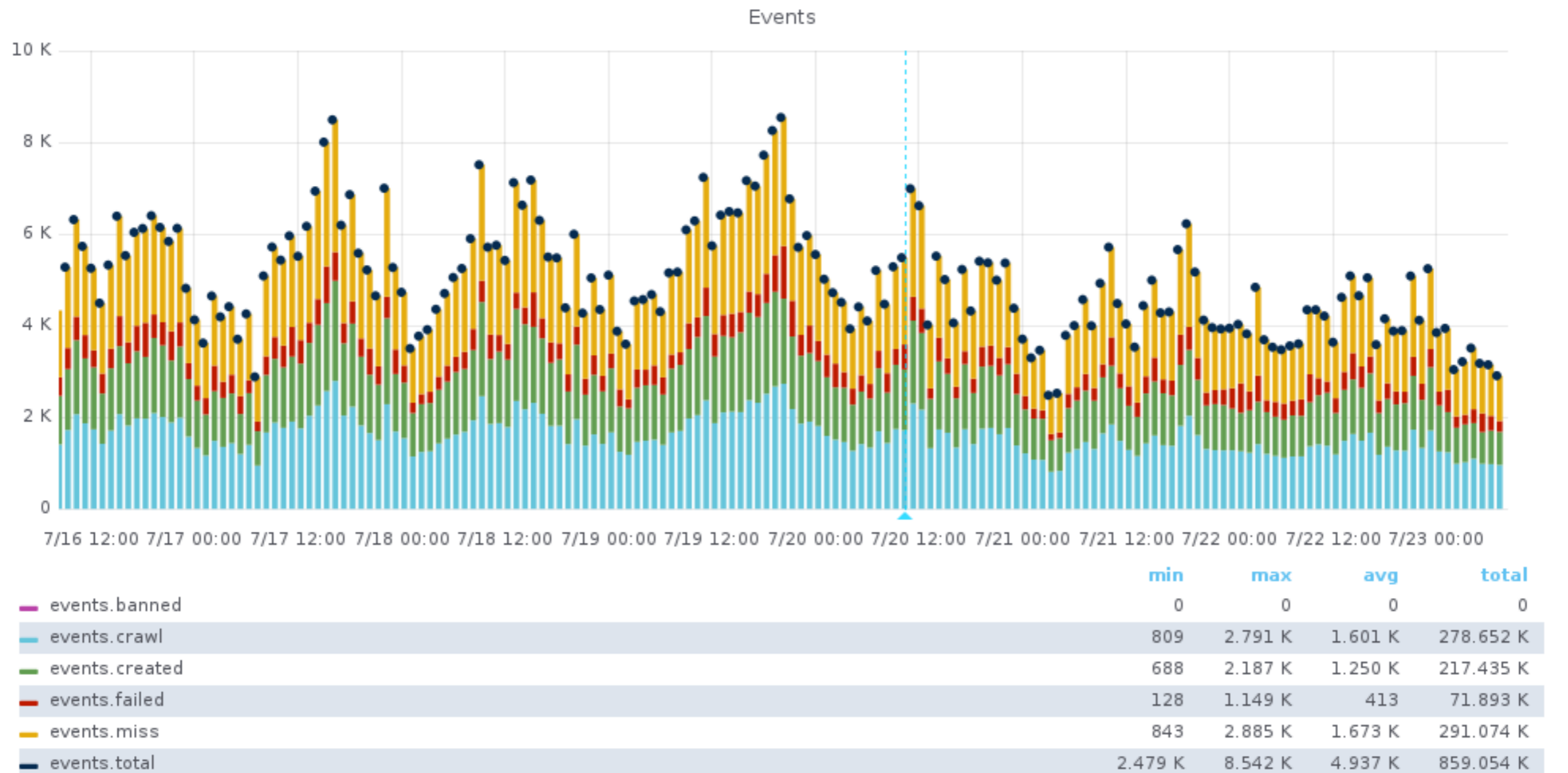
## Analytics & Monitoring



# Dashboards

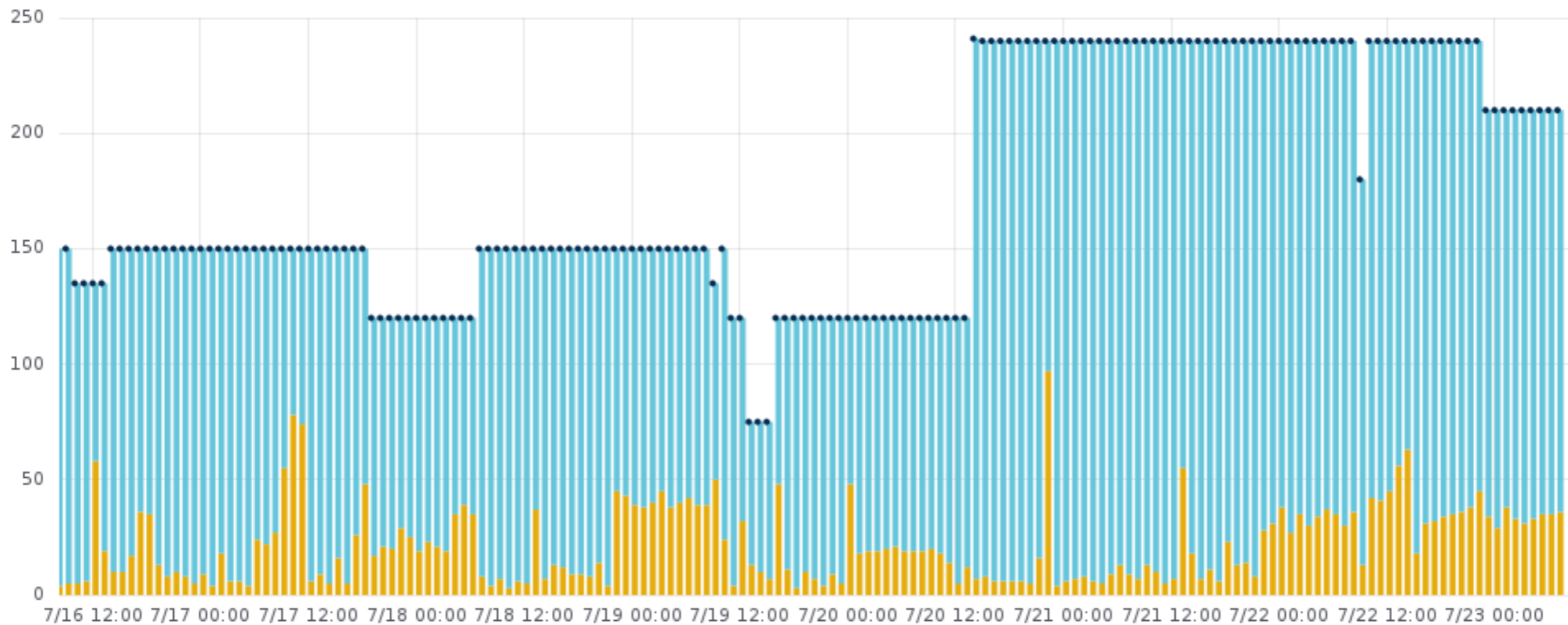


# Quality of Service



# Quality of Service

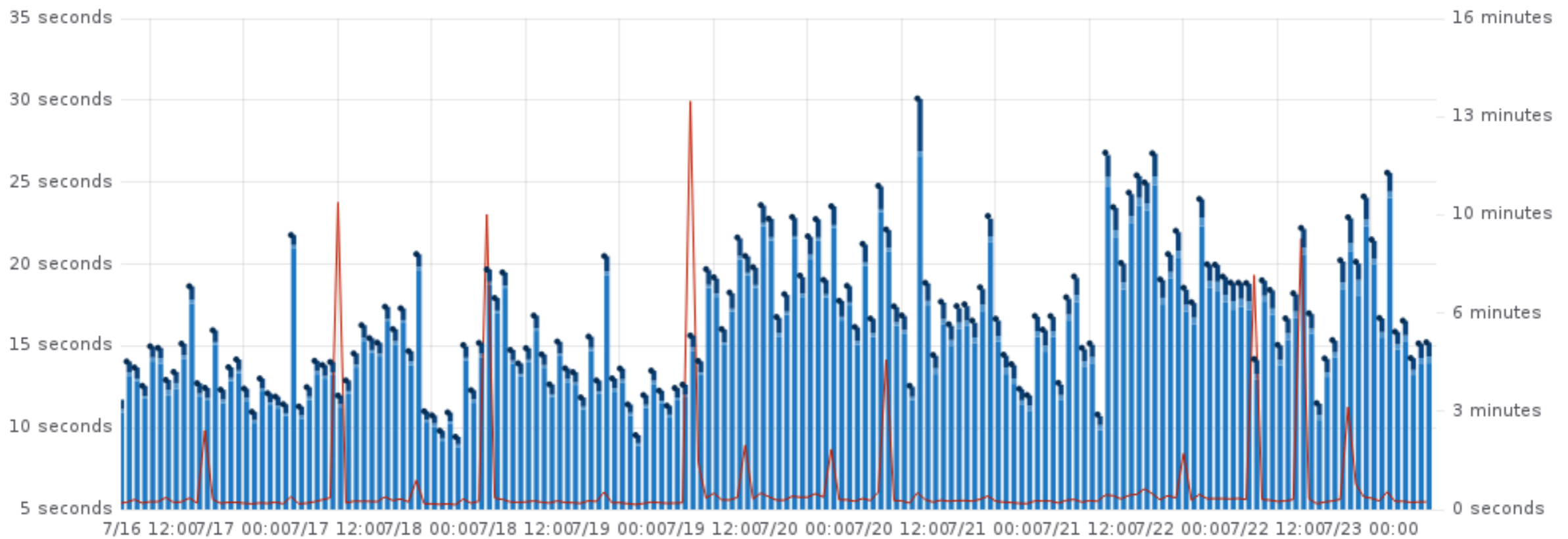
Spider Statuses



	min	max	avg	current	total
spider.active	3	97	22	36	3.7910 K
spider.idle	62	236	151	174	26.3000 K
spider.total	75	241	173	210	30.0910 K

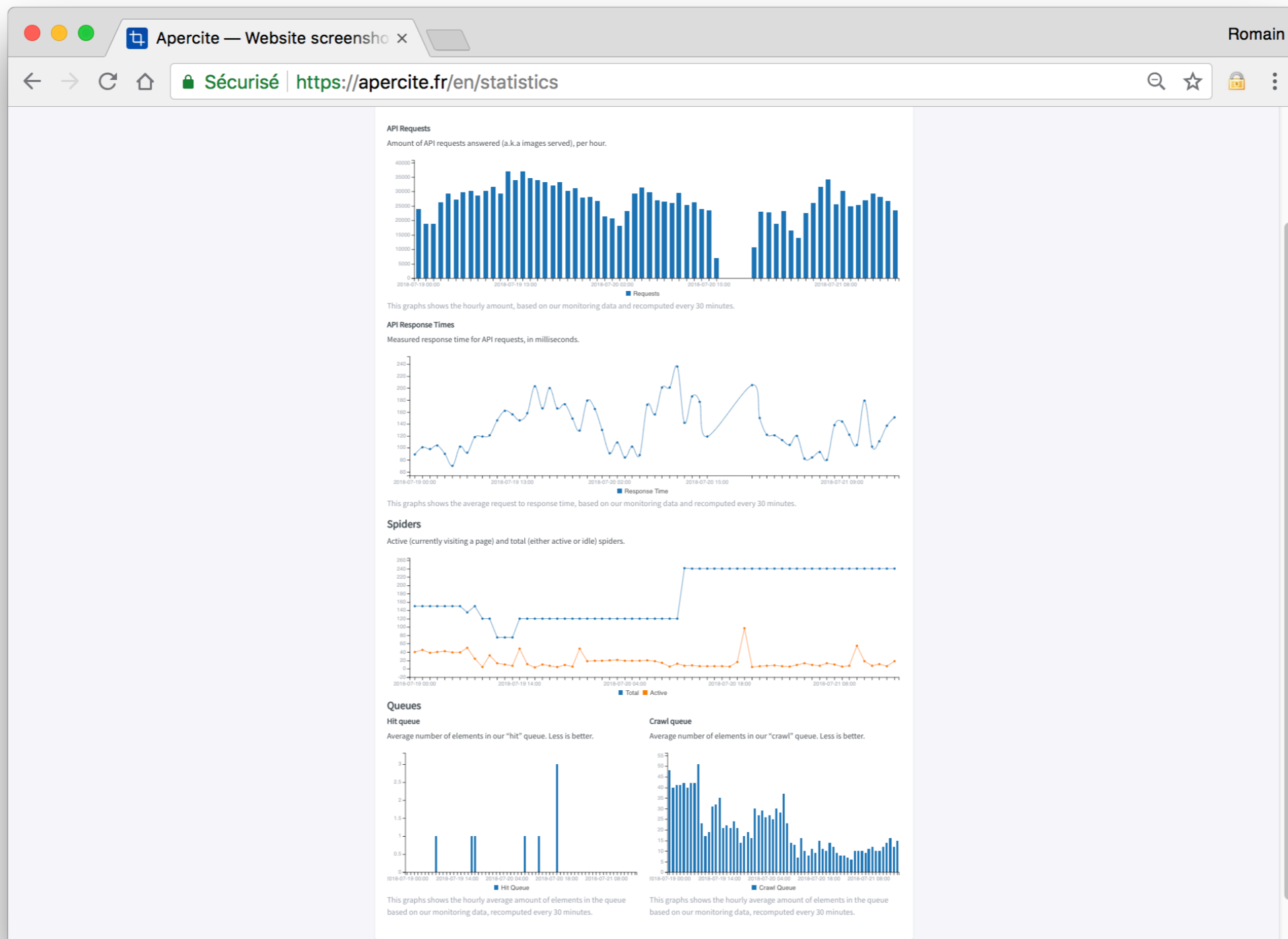
# Quality of Service

Spider Timings



	min	max	avg	current
spider.timings.browser	8 seconds	26 seconds	15 seconds	13 seconds
spider.timings.checks	140 milliseconds	933 milliseconds	265 milliseconds	381 milliseconds
spider.timings.delta (right-y)	9 seconds	13 minutes	42 seconds	15 seconds
spider.timings.thumbs	451 milliseconds	3 seconds	829 milliseconds	855 milliseconds
spider.timings.total	9 seconds	30 seconds	16 seconds	15 seconds

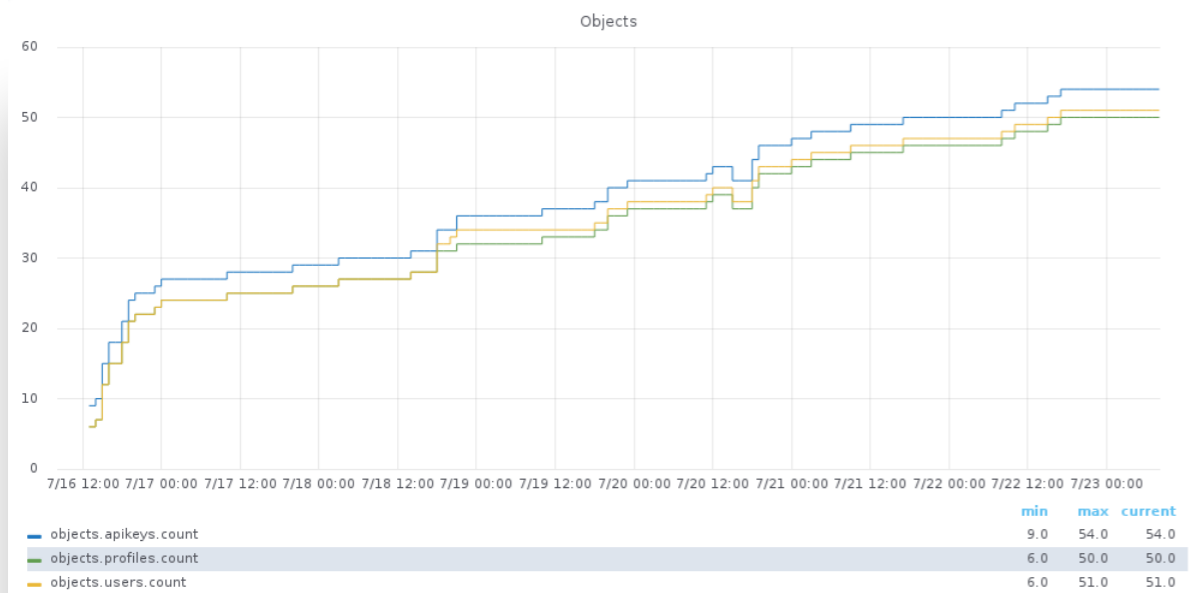
# Public Dashboards





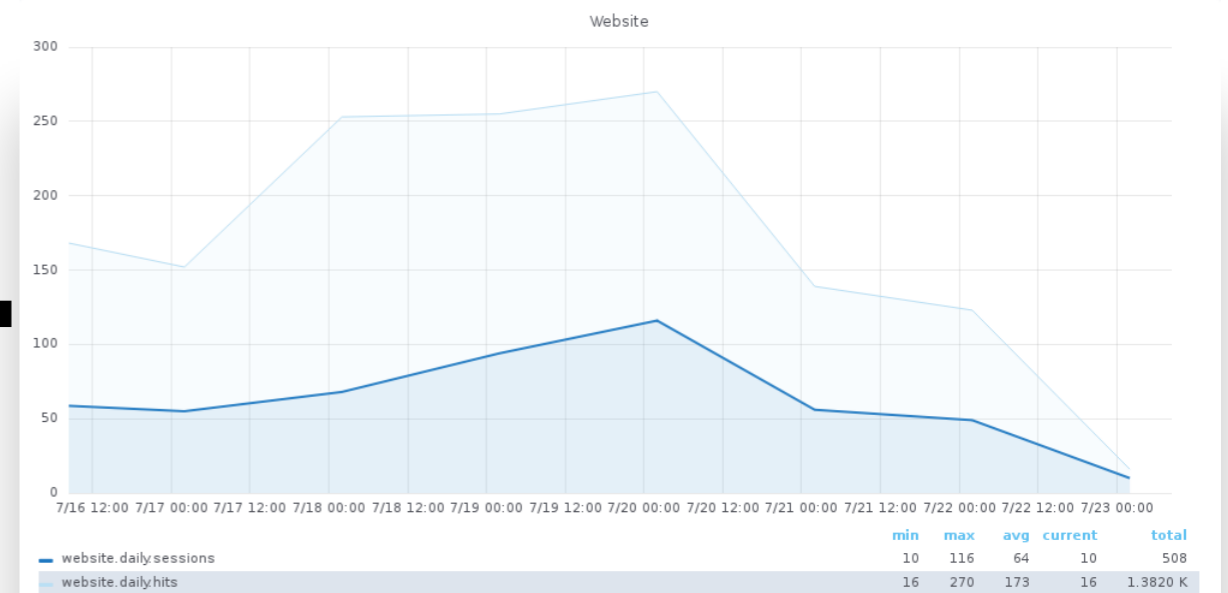
# Acquisition Rate

## User Counts

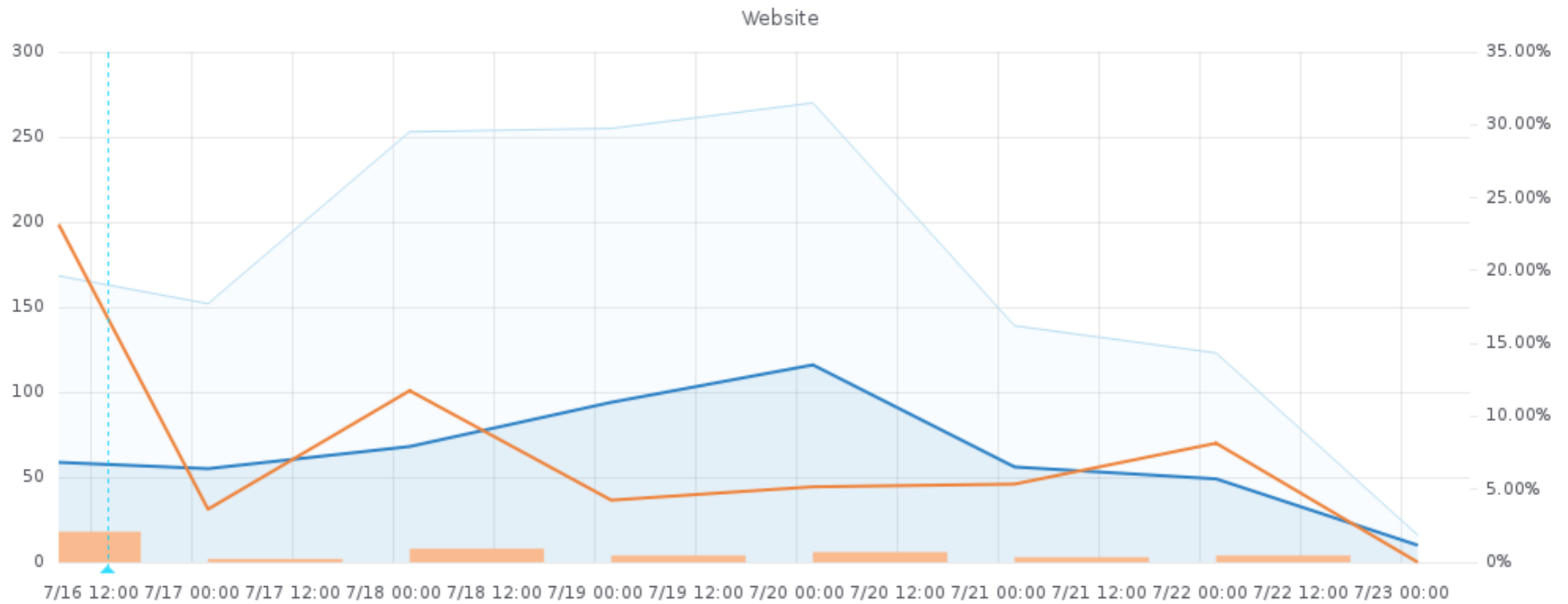


+

## New Sessions



# Acquisition Rate



	min	max	avg	current	total
website.daily.sessions	10	116	64	10	508
website.daily.hits	16	270	173	16	1.3820 K
new users	0	18	6	0	45
activation rate (right-y)	0%	30%	9%	0%	68%

We're just getting started.

MONITOR



# Iteration 0

- Cron job (K8S) runs everything every 30 minutes.
- No way to know if something fails.
- Expensive tasks.
- Hard to run manually.



# Airflow

«**Airflow** is a platform to programmatically **author**, **schedule** and **monitor** workflows.»

- Official docs

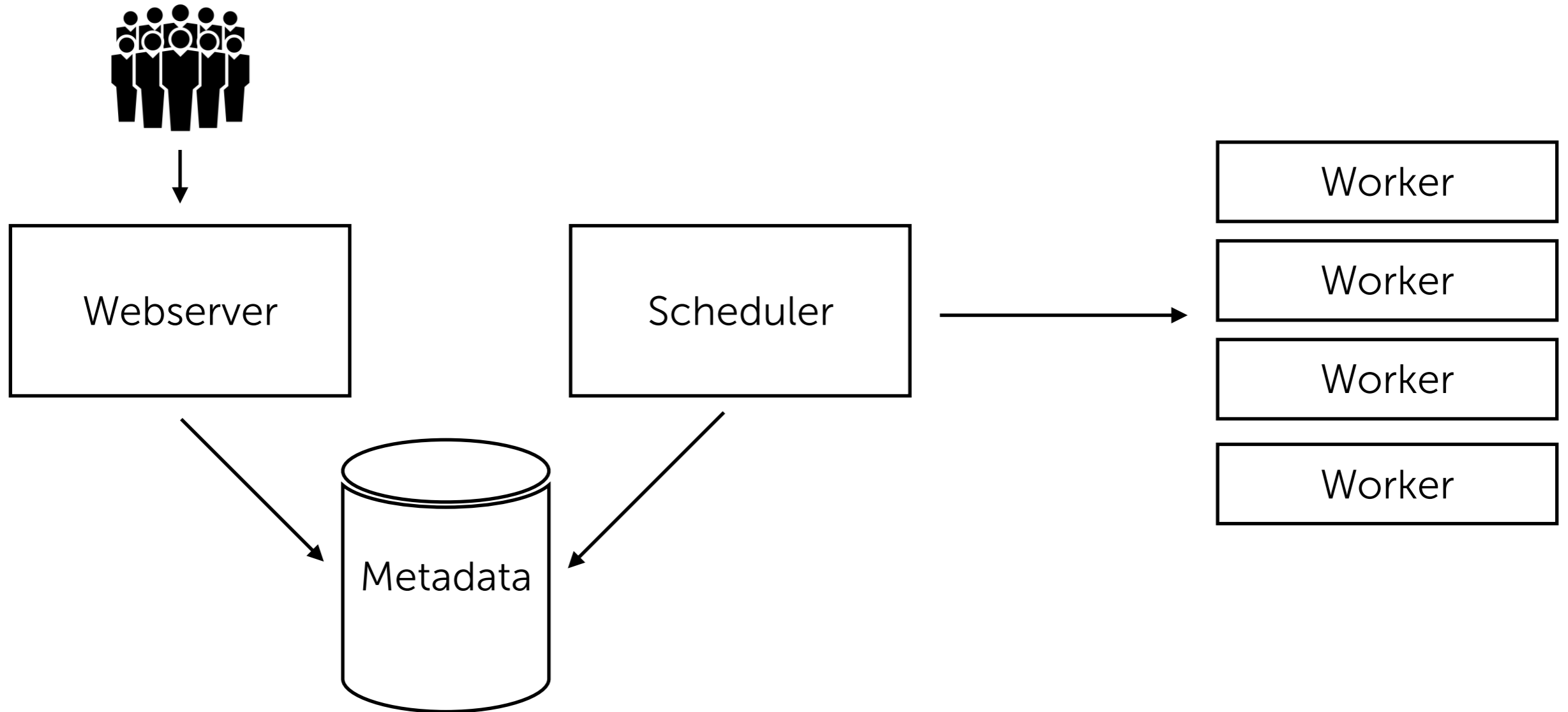


# Airflow

- Created by Airbnb, joined Apache incubation.
- Schedules & monitor jobs.
- Distribute workloads through Celery, Dask, K8S...
- Can run anything, not just Python.



# Airflow



Simplified to show high-level concept.  
Depends on executor (celery, dask, k8s, local, sequential ...)



# DAGs

	<i>i</i>	DAG	Schedule	Owner	Recent Tasks <i>i</i>	Last Run <i>i</i>	DAG Runs <i>i</i>	Links
	<input checked="" type="checkbox"/>	apercite.analytics.cleanup	@daily	apercite	1	2018-07-24 00:00 <i>i</i>	4	
	<input checked="" type="checkbox"/>	apercite.analytics.events	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	71 1	
	<input checked="" type="checkbox"/>	apercite.analytics.events_timings	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	71 1	
	<input checked="" type="checkbox"/>	apercite.analytics.google_analytics	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	71 2	
	<input checked="" type="checkbox"/>	apercite.analytics.objects	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	70 2	
	<input checked="" type="checkbox"/>	apercite.analytics.prometheus	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	71 1	
	<input checked="" type="checkbox"/>	apercite.analytics.spiders	@hourly	apercite	1	2018-07-25 08:00 <i>i</i>	70 2	
	<input type="checkbox"/>	read_object_counts <i>i</i>		apercite				
	<input type="checkbox"/>	sample_dag	@once	airflow	3	2017-09-10 00:00 <i>i</i>	5	

Showing 1 to 9 of 9 entries

Hide Paused DAGs

# DAGs

```
import shlex
from airflow import DAG
from airflow.operators.bash_operator import BashOperator

def _get_bash_command(*args, module='apercite.analytics'):
    return '(cd /usr/local/apercite; /usr/local/env/bin/python -m {} {})' .format(
        module, ' '.join(map(shlex.quote, args)),
    )

def build_dag(name, *args, schedule_interval='@hourly'):
    dag = DAG(
        name,
        schedule_interval=schedule_interval,
        default_args=default_args,
        catchup=False,
    )

    dag >> BashOperator(
        dag=dag,
        task_id=args[0],
        bash_command=_get_bash_command(*args),
        env=env,
    )

    return dag
```

# DAGs

```
# Build datasource-to-metrics-db related dags.
for source in ('google-analytics', 'events', 'events-timings', 'spiders',
'prometheus', 'objects'):
    name = 'apercite.analytics.' + source.replace('-', '_')
    globals()[name] = build_dag(name, 'read', source, '--write')

# Cleanup dag.
name = 'apercite.analytics.cleanup'
globals()[name] = build_dag(name, 'clean', 'all', schedule_interval='@daily')
```

# Connections

List (29) Create With selected ▾

<input type="checkbox"/>		Conn Id	Conn Type	Host	Port	Is Encrypted	Is Extra Encrypted
<input type="checkbox"/>		airflow_ci	mysql	localhost		⊖	⊕
<input type="checkbox"/>		airflow_db	mysql	localhost		⊖	⊖
<input type="checkbox"/>		apercite_events	postgres	cloudsql-proxy.database	5432	⊕	⊖
<input type="checkbox"/>		apercite_website	postgres	cloudsql-proxy.database	5432	⊕	⊖
<input type="checkbox"/>		aws_default	aws			⊖	⊕
<input type="checkbox"/>		beeline_default	beeline	localhost	10000	⊖	⊕
<input type="checkbox"/>		bigquery_default	bigquery			⊖	⊖
<input type="checkbox"/>		databricks_default	databricks	localhost		⊖	⊖
<input type="checkbox"/>		druid_ingest_default	druid	druid-overlord	8081	⊖	⊕
<input type="checkbox"/>		emr_default	emr			⊖	⊕
<input type="checkbox"/>		fs_default	fs			⊖	⊕

# Data Sources

```
from airflow.models import Connection
from airflow.settings import Session
```

```
session = Session()
website = session.query(Connection).filter_by(conn_id='apercite_website').first()
events = session.query(Connection).filter_by(conn_id='apercite_events').first()
session.close()
```

```
env = {}
```

```
if website:
```

```
    env['DATABASE_HOST'] = str(website.host)
    env['DATABASE_PORT'] = str(website.port)
    env['DATABASE_USER'] = str(website.login)
    env['DATABASE_NAME'] = str(website.schema)
    env['DATABASE_PASSWORD'] = str(website.password)
```

```
if events:
```

```
    env['EVENT_DATABASE_USER'] = str(events.login)
    env['EVENT_DATABASE_NAME'] = str(events.schema)
    env['EVENT_DATABASE_PASSWORD'] = str(events.password)
```

Warning: sub-optimal

# Learnings

- Multiple services, not trivial
- Helm charts :-(  
- Astronomer Distro :-)
- Read the Source, Luke

The image shows a hand-drawn diagram on a chalkboard. The word "Outtro" is written in the center in a white, hand-drawn font. The diagram consists of several rectangular boxes arranged in a grid-like structure. There are three arrows pointing to the right, one above the word, one below it, and one below the word. There are also some scribbled lines and a shaded area in the bottom right corner of the diagram.

Outtro

**Bonobo** helps you build assembly lines.  
Does not care about the surrounding factory.

**Airflow** helps you manage the whole factory.  
Does not care about the jobs' content.

**Grafana** let you see the data, with only a few  
queries to write.



romain@makersquad.fr



 rdorgueil