

Mockstagram crawler documentation

Table of Contents

Understanding problem statement.....	1
Proposed Solution.....	3
Implementation design (single node).....	4
Choice of platform for implementation design.....	5
Distributed crawler:.....	5
Distributed Timeseries database (TSDB):.....	5
NoSQL.....	7
Infrastructure installation.....	9
Running the code.....	10
Crawling logic.....	11
Updating aggregated metrics.....	12
Updating suspicious status.....	13
Accessing data for API server.....	14
Performance Evaluation.....	15
Summary.....	17

Understanding problem statement

<https://www.notion.so/V2-Affable-Data-Engineering-Task-5d4bec24edbf42ebbe2a285caa699b26>

1) Design crawler for Mockstagram for tracking time series data of 1 million accounts.

Assumption:

- Mockstagram endpoint has no crawling protection and impossible to DDos.

Comment:

- Requires use of efficient, distributed time series database for data with high cardinality.
- Requires distributed crawler design to distribute load across multiple servers
- Strong consistency is not a priority (I assume if followers are in ranges of 100k+, daily changes will be relatively insignificant to 1) affect the performance of machine learning models used, 2) affect the decision making of the data users)

2) Provide some easily accessible aggregated metrics (eg. averageFollowerCount, most recent data) for many users

Assumption:

- The aggregated metrics are globally aggregated (over the entire time series) rather than sliding window (moving averages).

Comment:

- The averageSomeMetrics can be downsampled without too much implication on accuracy while increasing performance speed greatly (eg. Of 1,000 time series data points, sample every 100th interval for calculation)
- Recommended to have database for storing these pre-computed metrics
- Provide a distributed cache with multiple read-only slaves storing these aggregated values.

3) Design feature for updating “suspicious” status, which is computationally heavy

Assumption:

- Server for machine learning service cannot be modified

Comment:

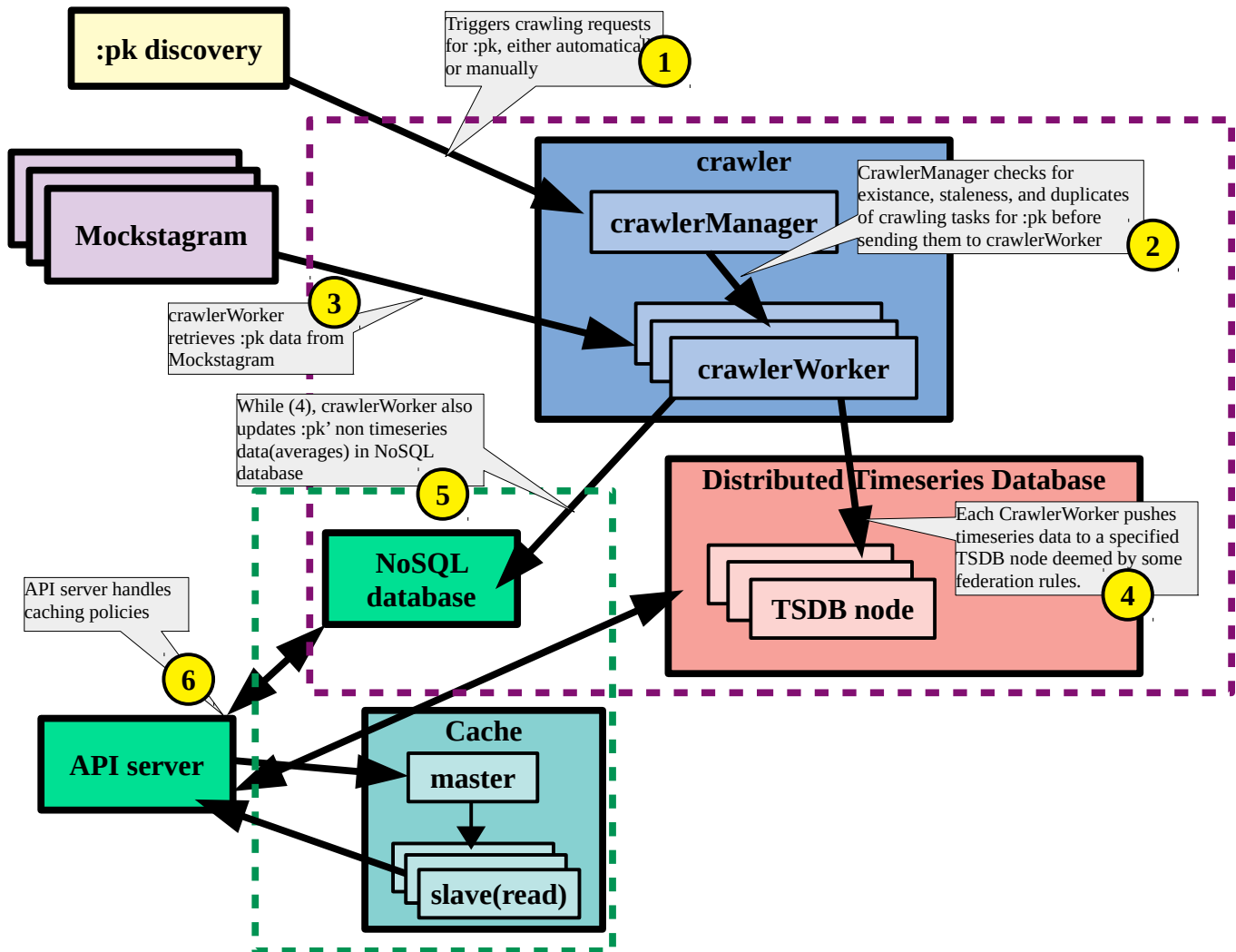
- Updates are not time critical (once daily)
- If results can be cached, we will use a cache store for such information (eg. Same input that has been calculated before)
- If results cannot be cached (eg. Function is always changing), then need a buffer to queue these requests as the consumer (‘/api/v1/influencers/is_suspicious’) is much slower than the producer (it is much easier to send 1000 requests to the consumer than to receive the same amount of replies in a short time)

Proposed Solution

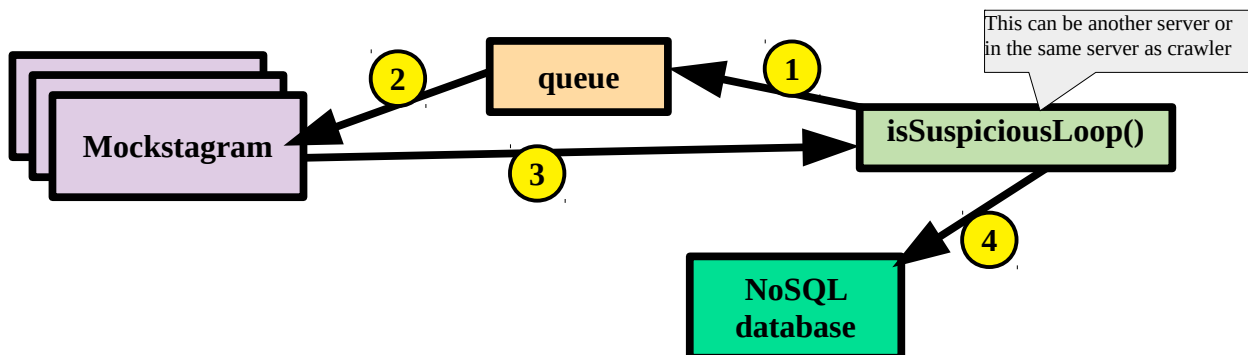
(Color of dotted boxes matches the problem it is solving)

Purple box: Solution for (1), Design crawler for Mockstagram with 1 million accounts to track.

Green box: Solution for (2), Provide some easily accessible metrics (eg. averageFollowerCount, most recent data) for many users

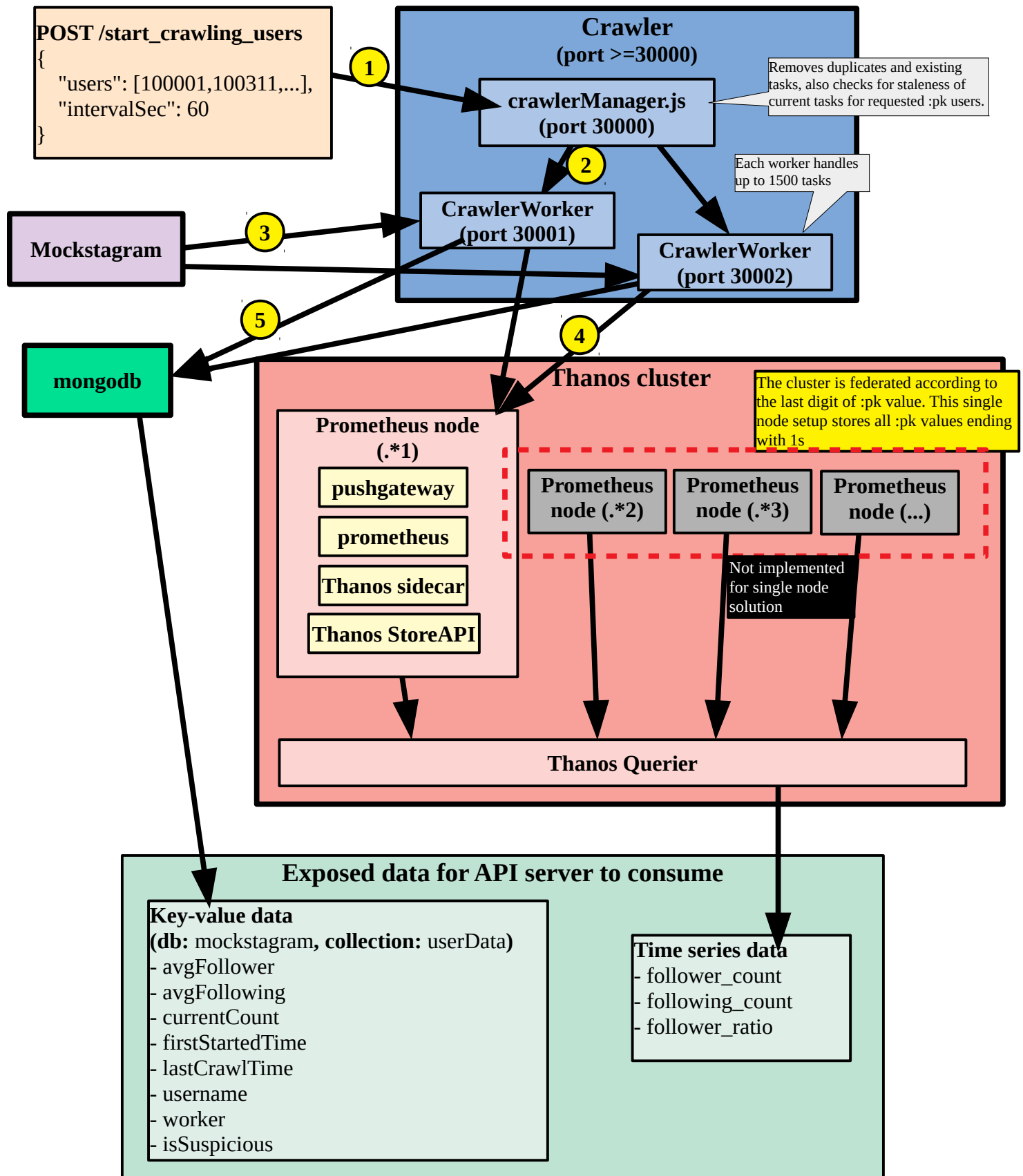


For (3), Design feature for updating “suspicious” status, which is computationally heavy



Implementation design (single node)

Due to time constraints, and a desirably simplified solution for a single node setup, the implementation design is a watered down version of the proposed solution.



Choice of platform for implementation design

Distributed crawler:

NodeJS is used for the backend here because of convenience and the ease of handling asynchronous events.

Distributed Timeseries database (TSDB):

This is a difficult choice because of a few reasons.

1) Key-value store vs TSDB:

Key-value NoSQL databases are widely used and hence there is a greater variety of choices and community support for production use. However, they are not optimized for time series data hence may not perform better than TSDB for the problem statement.

TSDB, on the other hand, are only handful in terms of variety (Opentsdb, InfluxDB, TimescaledDB, Prometheus/Thanos). But some of them have been used for production with good community support hence these choices are definitely preferred over a key-value database.

User reviews:

“Cassandra/MongoDB (NoSQL) has been disqualified, since it is [much slower than Timescale](#)(TSDB)” [5]

2) Choice of TSDB

This is an even more difficult choice to make. Between the various TSDBs, each of them have their pros and cons in a way that doesn't make any of them stood out as an obvious choice (although there are easy rejects).

The pros/cons are summarized briefly below:

Opentsdb:

Pros:

- Based on Hadoop and Hbase

Cons: ?

InfluxDB:

Pros:

- Developed from PostgresDB

Cons:

- Expensive distributed solution
- Performance issues

- “Influxdb suck for collecting large point in time metrics. Also the memory it uses is huge.”[1]
- “If you want clustering for HA or for horizontal scaling, you need the enterprise version of InfluxDB. “[1]
- “InfluxDB performance dropped ... This is significant performance loss comparing to other competitors(VictoriaMetrics, TimescaleDB)” [6]
- “InfluxDB: didn’t finish because it required more than 60GB of RAM.”[6]

TimescaledDB:

Pros:

- Developed from PostgresDB

Cons:

- Inefficient storage, “Timescale data occupies **whopping 29GB** on HDD. That’s 50x more than InfluxDB and 75x more than VictoriaMetrics”[5]

Prometheus/Thanos:

Thanos as a platform to manage Prometheus clusters.

Pros:

- More scalable than some other platforms,
 - “Yea, the Prometheus + Thanos combo for larger deployments is crazy awesome. It allows for nearly infinite storage, while keeping the deployment simple and robust against failure.”[1]
 - “Prometheus open source is more scalable than influxdb”[1]

Cons:

- Not a conventional choice for TSDB, mainly for monitoring
- No strong consistency, data stored may not be accurate.
- Higher number of moving parts (more complex to setup)
- By nature, ingests data by ‘pulling’. But a *pushgateway* can be used to change the ingestion to a ‘push’ one.
- Backfilling is a pain to implement, I did not consider this feature before, but now it is too late to change.

Prometheus/VictoriaMetrics

VictoriaMetrics as a platform to manage Prometheus clusters.

Pros:

- Better query speed, “VictoriaMetrics wins InfluxDB and Timescale in all the queries by a margin of up to 20x. It especially excels at heavy queries, which scan many millions of datapoints across thousands of distinct timeseries.”[5]

- Better insert speed and performance on low cardinality, “VictoriaMetrics wins in insert performance and in compression ratio(with respect to InfluxDB, TimescaleDB)” [6]

Cons:

- Same as Prometheus/Thanos

Conclusion:

The obvious choice is to use a TSDB, but which one? Unfortunately, time does not permit for setting up a test environment to benchmark these various TSDBs. So judgement can only be based on the respective platform features and user reviews (which may be biased). The reasons for deciding on the choice of TSDB are as follows:

- InfluxDB is a strong reject based on multiple sources on performance issues
- TimescaleDB may have some performance issues
- Not enough user review on Opentsdb
- Prometheus/Thanos and Prometheus/VictoriaMetrics seems comparable
- There are seemingly more users/community support for Prometheus/Thanos
- I am already familiar with Prometheus

Hence, Prometheus/Thanos became the final choice for distributed TSDB.

References

- [1] https://www.reddit.com/r/devops/comments/8qvpz7/prometheus_or_influxdb_tick/
- [2] <https://medium.com/faun/comparing-thanos-to-victoriametrics-cluster-b193bea1683?>
- [3] https://www.reddit.com/r/devops/comments/941n2k/tsdbs_at_scale_part_one/
- [4] <https://blog.timescale.com/blog/timescaledb-vs-influxdb-for-time-series-data-timescale-influx-sql-nosql-36489299877/>
- [5] <https://medium.com/@valyala/when-size-matters-benchmarking-victoriametrics-vs-timescale-and-influxdb-6035811952d4>
- [6] <https://medium.com/@valyala/high-cardinality-tsdb-benchmarks-victoriametrics-vs-timescaledb-vs-influxdb-13e6ee64dd6b>

NoSQL

This is used for storing meta-data for :pk users. No nested structure or other complications involved so a simple key-value NoSQL database is enough. MongoDB is an obvious choice.

Distributed design may not be needed as

1) data is presumed to grow very slowly, or since :pk users is capped at 1 million, it may not grow at all. Only involves updating the existing 1million data points.

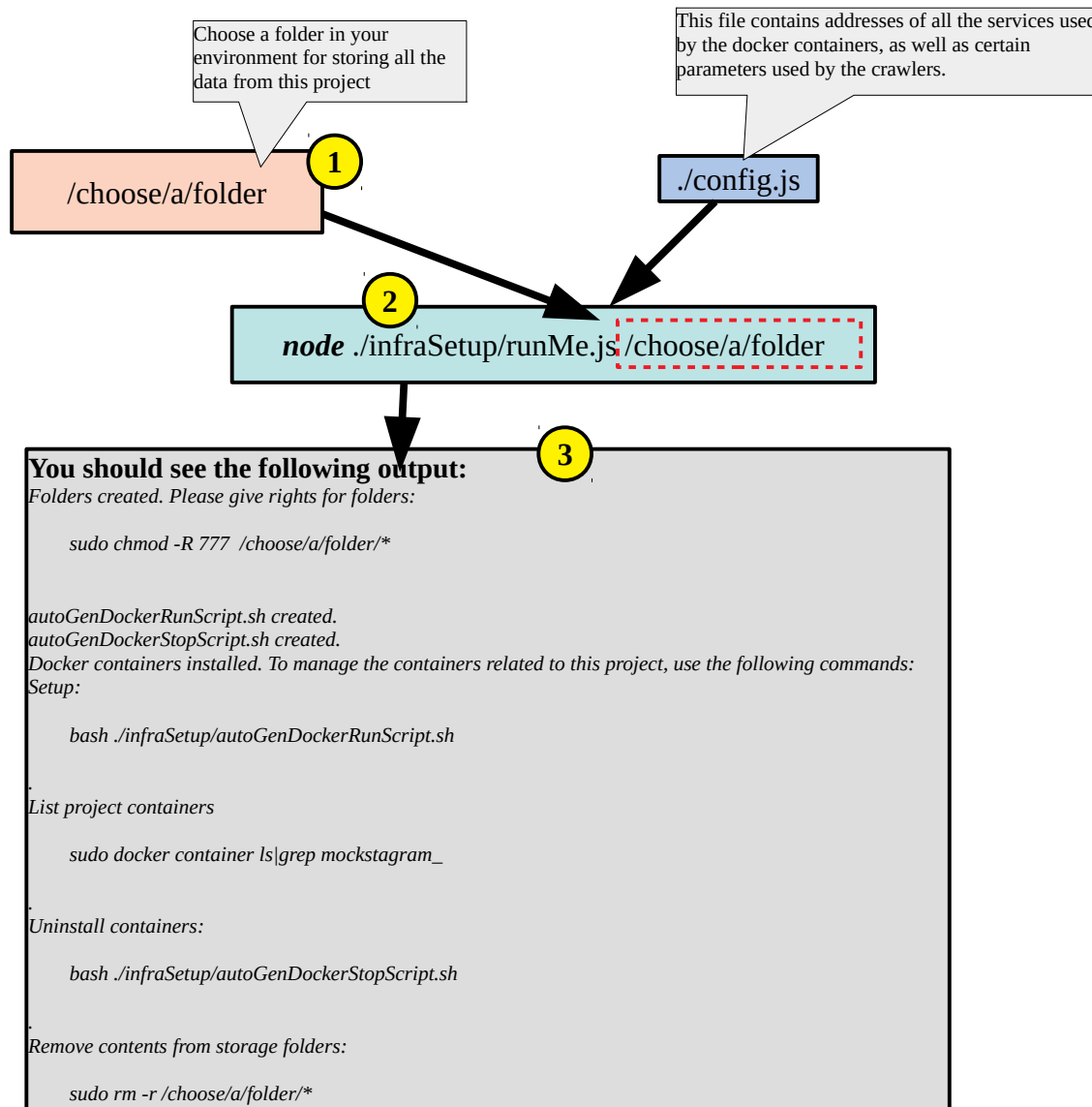
2) Taking fully saturated storage into account, each field storing 64 bytes (worst case) , total storage requirements : $1,000,000 \text{ users} \times 7 \text{ fields(see implementation design)} \times 64 \text{ bytes} < \mathbf{0.5Gb}$ which is quite manageable.

Infrastructure installation

All infrastructure platforms are ran in docker containers to avoid polluting and potentially conflicting the test environment.

The setup steps are as below:

- cd into repository
- Run “node ./infraSetup/runMe.js /choose/a/folder”



There are 7 containers used for this project. They are

- mockstagram_pushgateway
- mockstagram_prometheus
- mockstagram_thanos_sidecar
- mockstagram_thanos_store

- mockstagram_thanos_querier
- mockstagram_mongodb
- mockstagram_pushwiper

You should `sudo docker container ls|grep mockstagram_` make sure they are running successfully.

Optional

You may find some Docker/Kubernetes files in `/mockstagram_scalable` folder for increasing Mockstagram instances so as to handle crawling load.

Installing dependencies

- run `'npm i --save'`

Running the code

- cd into the repo

1) Starting crawler manager

- run `'node ./crawler/crawlerManager.js'`

Output:

```
crawlerManager listening on port 30000
```

To start crawling requests for :pk values, send HTTP request with the following format

```
POST localhost:30000/start_crawling_users
{
  'users': [1996161,1787381],
  'intervalSec': 5
}
```

2) Sending some requests to crawler with auto-generated :pk values

This is to help with auto-generating large number of random :pk values to help with load testing.

Pass a number argument, N, to the script to send crawling requests for N randomly generated :pk values to the crawlerManager. Eg: the following example sends crawling requests for 20 random :pk values

- run `'node ./crawler/crawlerStressTest.js 20'`

3) Tests

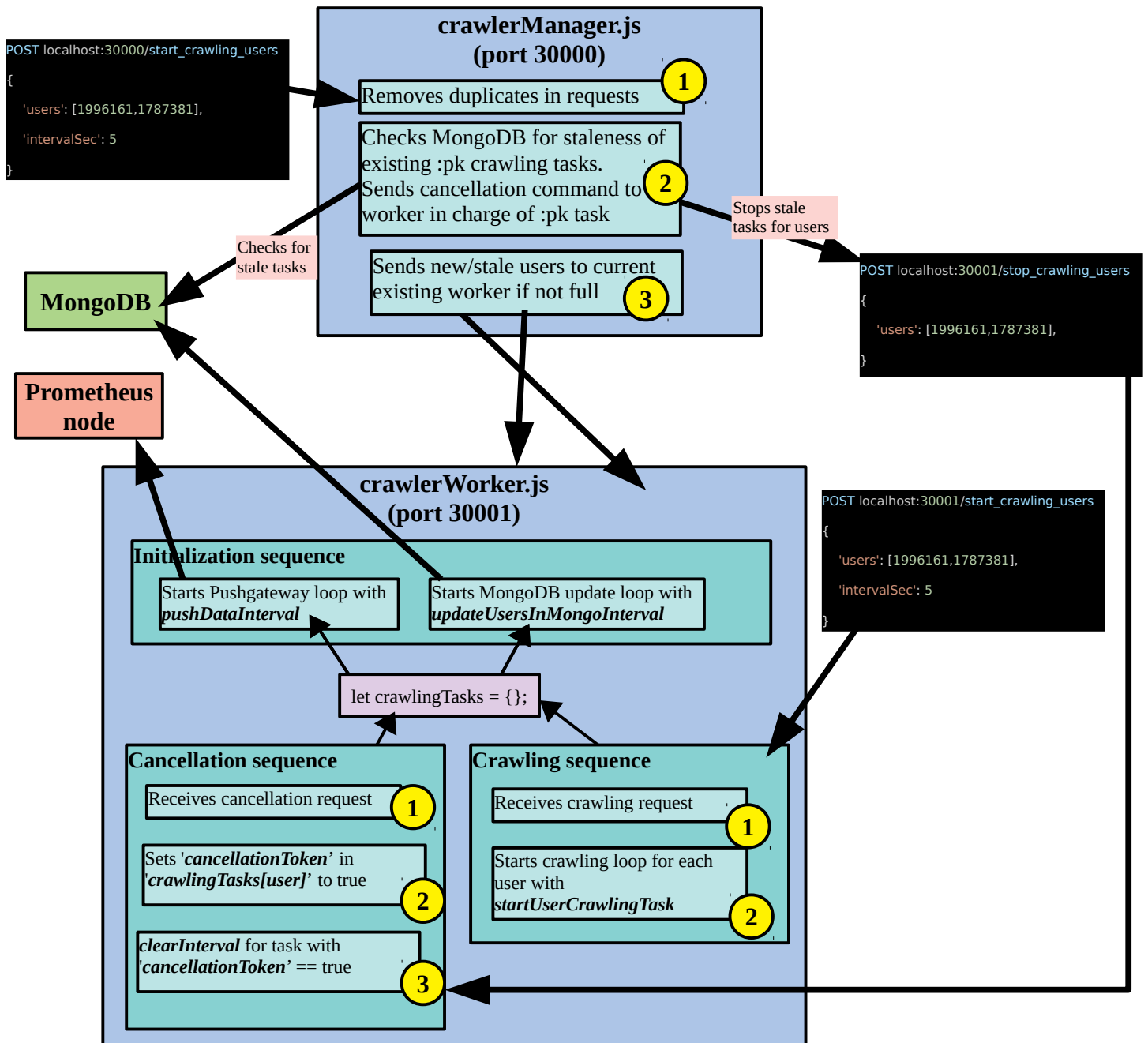
Some basic unit tests are setup in `./crawler/crawlerSpec.js`. To use them,

- run `'./node_modules/jasmine/bin/jasmine.js ./crawler/crawlerSpec.js'`

Crawling logic

The crawlerManager does a series of checks before sending these :pk values to workers.

If there are existing workers with available slots, manager will assigns these :pk tasks to them. If not, the manager creates new workers with ports >30000 and sends these tasks to them.



Updating aggregated metrics

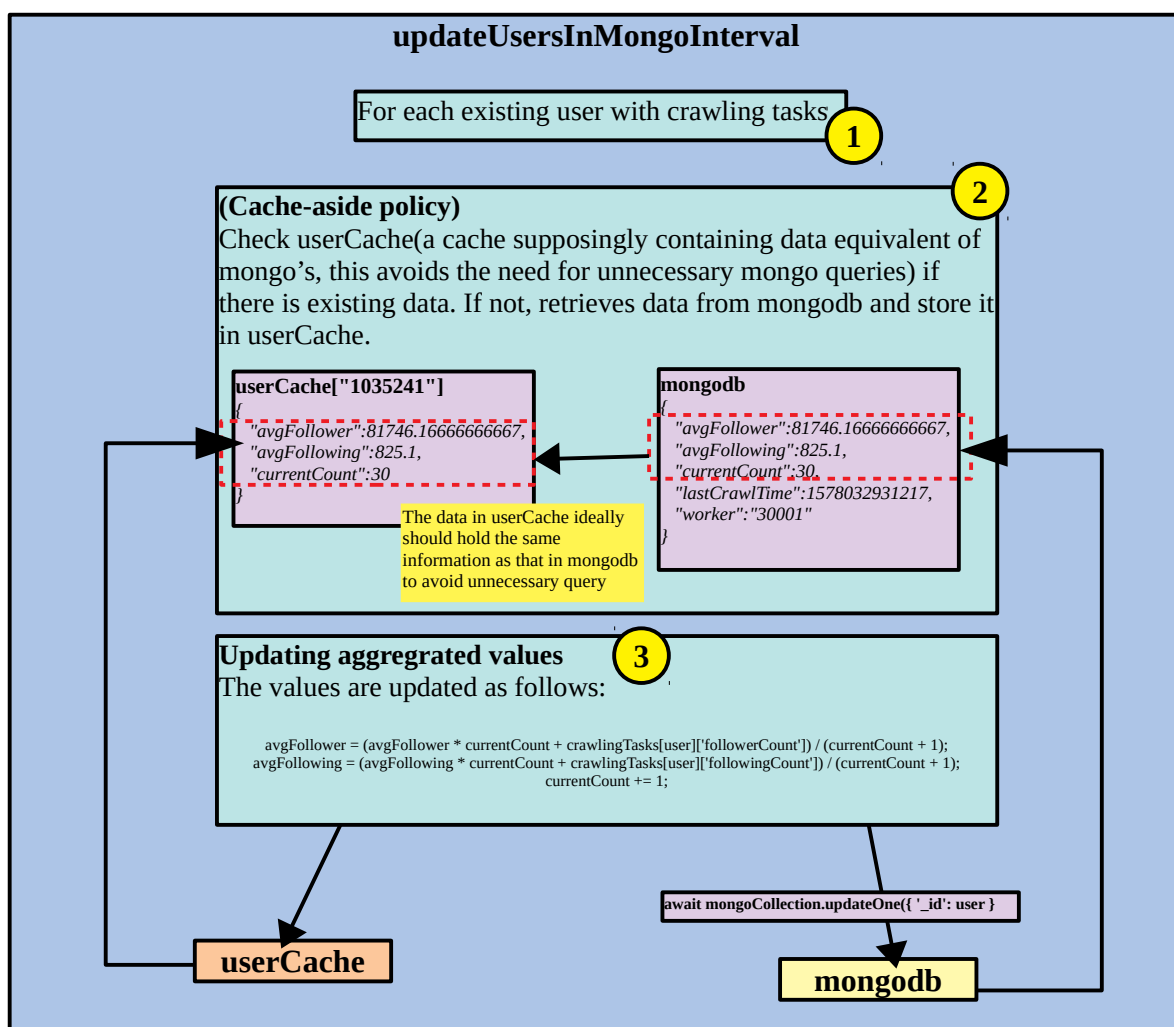
Aggregated metrics refer to metrics like global average follower/following counts.

As seen in the crawling logic above, aggregated metrics is handled by ***updateUsersInMongoInterval*** in ***crawlerWorker***.

Updating of aggregated metrics is separated from the crawling tasks because the decoupling means that we can control the frequency of these 2 different tasks separately.

This means that we can downsample aggregated metrics easily (crawling may be at 1 minute interval, but updating aggregated metrics can be done at hourly or daily interval instead). This helps with improving the performance of overall system.

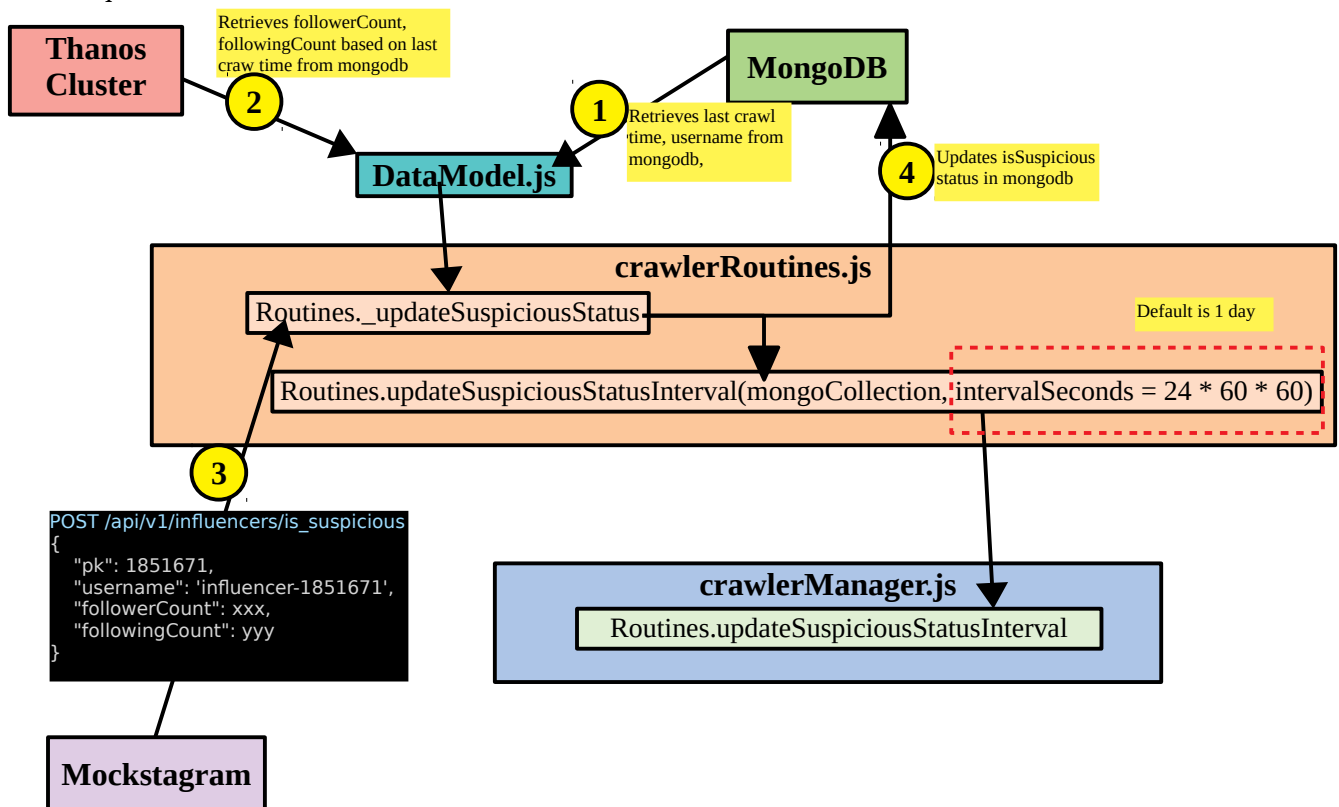
The sequence are as follows:



Updating suspicious status

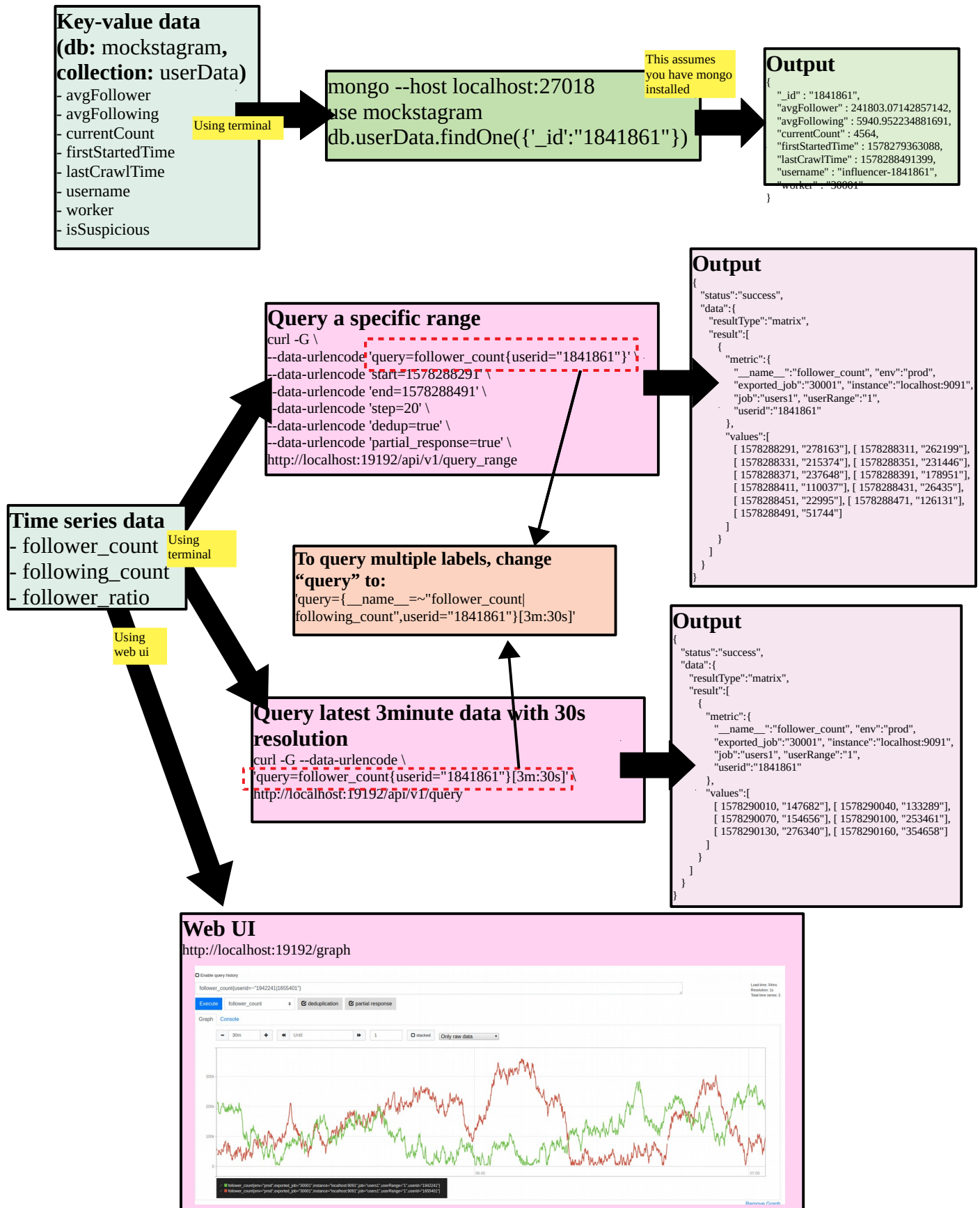
This is implemented minimally due to time constraints. There is no buffer used, so the function will send requests for /is_suspicious for all users concurrently(albeit in a synchronized manner so there is some “queuing” to avoid DDOSing the server) once the interval is triggered.

The sequence are as follows:



Accessing data for API server

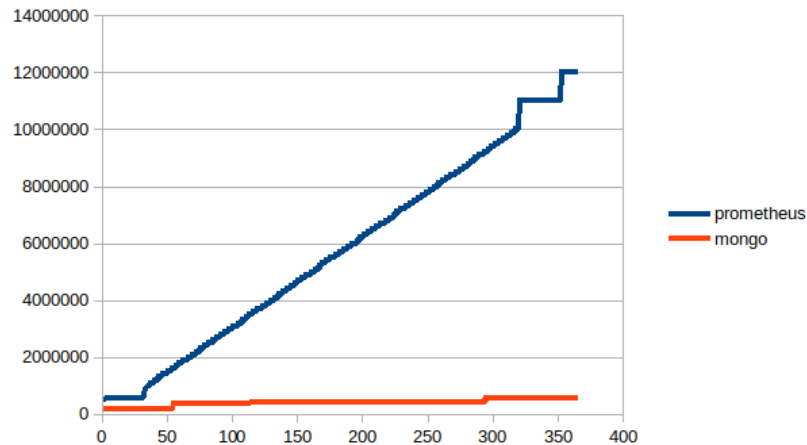
The data from each database can be access as follows:



Performance Evaluation

Storage performance

Testing was done with 1000 writes per second. Storage was profiled using `sudo du -cha --max-depth=1` on user-defined storage folder (as determined when setting up infrastructure) at 1 second interval for a period of ~ 5 minutes.



The increment in storage space of Prometheus worked out to be 25kb per second.

To scale this up to 1 million users:

25Kb/1k Ops/second * 16.7k Ops/second

~ 0.4175 Mb/second

~ 36Gb/day

Running performance

The number of crawler tasks is increased gradually in steps of 500 at 60 seconds interval using the command below:

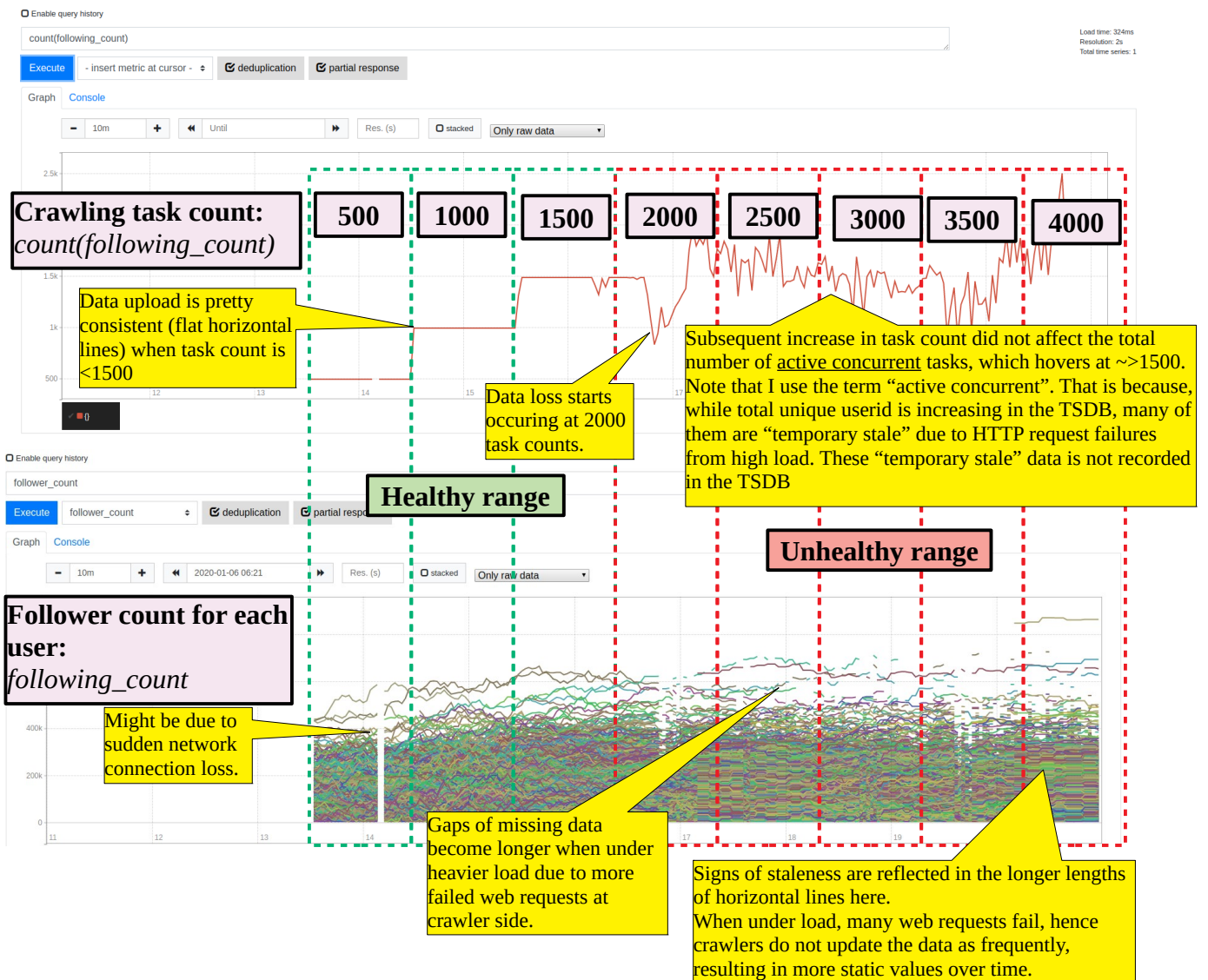
```
while true; do node ./crawler/crawlerStressTest.js 500; sleep 60; done
```

The specs of the hardware is as follows:

CPU: Intel® Core™ i7-6900K CPU @ 3.20GHz × 16

RAM: 64Gb

However, there are other unrelated processes running in the same environment so testing may not be accurate.



One major bottleneck is the number of concurrent HTTP requests that can handle by a single node. This limitation of HTTP requests is due to 2 major factors

- 1) Fully utilized CPU time (Crawler side)
- 2) Suspected fully utilized sockets in either NodeJS or machine itself since the network traffic is not saturated at all. Need more time with this.

Hence, it is highly desirable to operate the crawlers within healthy limits, in this case <1500 per node to avoid data loss.

Summary

This assignment proposes a possible pipeline for crawling Mockstagram. The choice of TSDB(Prometheus) and handling of computational values(pre-computation, downsampling) means that the solution favours higher availability (AP) over stronger consistency (weak CP).