

Course project - Big Data Processing

Team

- [Dmytro Lutchyn](#)
- [Yarema Mishchenko](#)

System design

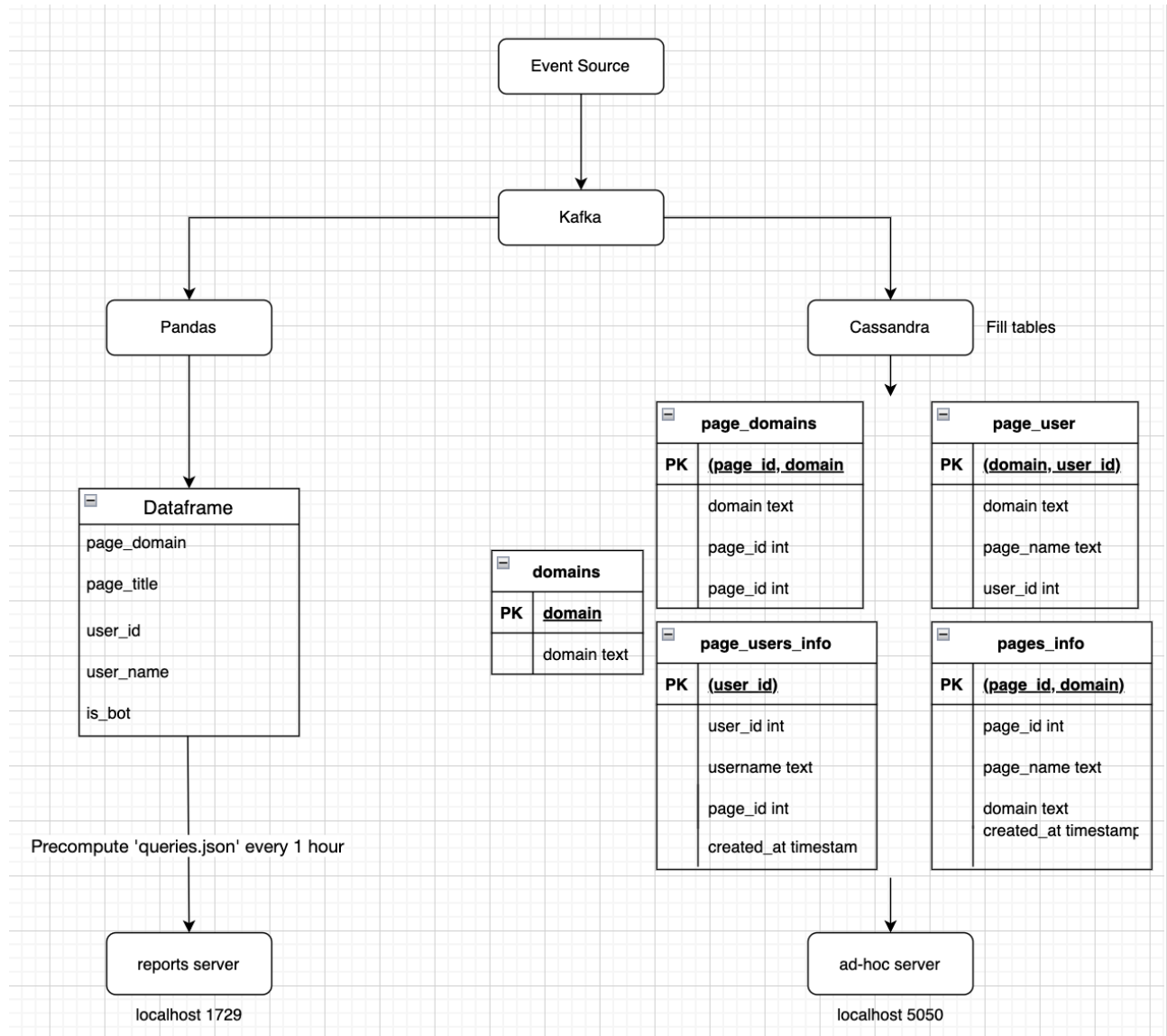
Our system consists of 4 independent services, as well as Kafka, Zookeeper, and Cassandra instances.

These 4 services are, as named in the project:

- `producer/producer.py` - responsible for fetching the Wikipedia data from the stream and passing it forward
- `consumer/consumer.py` - responsible for processing the data, filling Cassandra tables and maintaining a Pandas dataframe
- `precomputed-reports/reports_server.py` - responsible for maintaining a JSON file with precomputed outputs and serving it upon request
- `ad-hoc-queries/ad_hoc_server.py` - responsible for maintaining a Cassandra client connection and serving query results upon request.

In the diagram below, you can see how all those components create a coherent system.

Diagram



Kafka - producer & consumer

As events come into the system, they pass through a Kafka MQ. This is done for scalability as well as for separating services with different purposes. Preprocessing events in the producer allows the consumer to focus on the important logic, which is writing records to a Cassandra DB and precomputing reports.

For reports from Category A, we are using Pandas instead of Apache Spark. The reason is twofold:

- The amount of data is small enough to be suitable for a non-distributed system. We are keeping track of all records from the last 8 hours, which (at the rate of ~2/sec) there are about 57k. Given that one entry usually weighs no more than 500 bytes (pessimistic estimate, for long titles), we can easily fit our records in less than 30 MB.
- We are more comfortable with Pandas :D While often overlooked, the amount of developer's experience in a framework is often extremely important for a streamlined development process.

Cassandra - DB for storage

Our Cassandra storage has 5 different tables for each ad-hoc queries.

- domains -> (domain text), PRIMARY KEY (domain)
- page_user -> (domain text, page_name text, user_id int, page_id int,) PRIMARY KEY (user_id, domain, page_id)
- page_domains -> (domain text, page_id int) PRIMARY KEY (domain, page_id)
- pages_info -> (page_id int, page_name text, domain text, created_at timestamp) PRIMARY KEY (page_id, domain)
- page_users_info -> (user_id int, username text, page_id int, created_at timestamp) PRIMARY KEY (user_id, created_at, page_id)

This way the use of tables becomes easier and more efficient when doing SELECT queries.

You can check out `cassandra-client/create-tables.cql` for keyspace initialization.

Precomputed reports server - Category A

Every hour, the consumer recomputes records for the last 6 full hours, excluding the last hour. These requests, in a JSON form, are then sent via a POST request to our precomputed server. Here, the server stores `queries.json` to give out on a GET request, which you can do as a client via port forwarding.

Ad hoc server - Category B

Consumer creates `CassandraClient` that updates all existing tables as soon as the new message comes to Kafka. `CassandraClient` has separate methods for inserting data to different tables (`insert_to_page_user` , etc.), and methods for getting answers to all queries (`select_from_pages_info` , etc.)

To get SELECT results we have separate script that gets requests from the server and executes corresponding SELECT queries.

`client_demo.py` sends these requests via POST request to the server.

Results

All of our results are stored in the `project-results` directory. These results are generated using the `client_demo.py` script. As per the duration requirement, the system was running for ~8h 27m before making the requests.

Category A

For category A, the JSON file with query results is saved to `project-results/queries.json` .

The file is extremely large because of the 3rd query: the top author created over 20000 pages.

First query:

```
▼ "first_query": [  
  ▼ {  
    "time_start": "01:00",  
    "time_end": "02:00",  
    ▼ "statistics": {  
      "ar.wikipedia.org": 197,  
      "arz.wikipedia.org": 99,  
      "ban.wikipedia.org": 1,  
      "bg.wikipedia.org": 2,  
      "bjn.wikipedia.org": 1,  
      "ca.wikipedia.org": 1,  
      "ceb.wikipedia.org": 1,  
      "commons.wikimedia.org": 1010,  
      "cy.wikisource.org": 1,  
      "da.wikipedia.org": 14,  
      "de.wikipedia.org": 4,  
      "en.wikibooks.org": 1,  
      "en.wikinews.org": 4,  
      "en.wikipedia.org": 179,  
      "en.wikiquote.org": 1,  
      "en.wikisource.org": 121,  
      "en.wiktionary.org": 34,  
      "es.wikipedia.org": 34,  
      "es.wikiquote.org": 1,  
      "es.wikisource.org": 4,  
      "es.wiktionary.org": 1,  
      "eu.wikipedia.org": 1,  
      "fa.wikipedia.org": 34,  
      "fi.wikipedia.org": 2,  
      "fr.wikipedia.org": 14,  
      "fr.wikisource.org": 18,  
      "gu.wikipedia.org": 1,  
      "hi.wikipedia.org": 1,  
      "hr.wikipedia.org": 1,  
      "ht.wikipedia.org": 2,  
      "hu.wikipedia.org": 2,  
      "id.wikipedia.org": 4,  
      "incubator.wikimedia.org": 7,  
      "is.wikipedia.org": 2.
```

Second query:

```
1,
▼ "second_query": {
  "time_start": "01:00",
  "time_end": "07:00",
  ▼ "statistics": [
    ▼ {
      "domain": "ar.wikipedia.org",
      "created_by_bots": 392
    },
    ▼ {
      "domain": "as.wikipedia.org",
      "created_by_bots": 17
    },
    ▼ {
      "domain": "bn.wikipedia.org",
      "created_by_bots": 2
    },
    ▼ {
      "domain": "ca.wikipedia.org",
      "created_by_bots": 1
    },
    ▼ {
      "domain": "commons.wikimedia.org",
      "created_by_bots": 3407
    },
    ▼ {
      "domain": "de.wikipedia.org",
      "created_by_bots": 10
    },
    ▼ {
      "domain": "en.wikipedia.org",
      "created_by_bots": 65
    },
    ▼ {
      "domain": "es.wikipedia.org",
      "created_by_bots": 2
    },
    ▼ {
      "domain": "eu.wikipedia.org",
      "created_by_bots": 7
    },
    ,
  ]
}
```

Third query:

```
    },
    "third_query": {
      "time_start": "01:00",
      "time_end": "07:00",
      "top_users": [
        {
          "uid": 963971,
          "uname": "Vojtěch Dostál",
          "num_created": 21460,
          "titles": [
            "Q112453917",
            "Q112453918",
            "Q112453919",
            "Q112453920",
            "Q112453921",
            "Q112453922",
            "Q112453923",
            "Q112453924",
            "Q112453925",
            "Q112453926",
            "Q112453927",
            "Q112453928",
            "Q112453929",
            "Q112453930",
            "Q112453931",
            "Q112453932",
            "Q112453933",
            "Q112453934",
            "Q112453935",
            "Q112453936",
            "Q112453937",
            "Q112453938",
            "Q112453939",
            "Q112453940",
            "Q112453941",
            "Q112453942",
            "Q112453943",
            "Q112453944"
```


Category B

As for category B, the output is displayed directly when running `client_demo.py`.

The arguments were chosen to fit well in a screen.

```
raven@DESKTOP-78L52C4: ~/bigdata/project
raven@DESKTOP-78L52C4:~/bigdata/project$ python3 client_demo.py
task 1:
['el.wikipedia.org', 'bpy.wikipedia.org', 'an.wikipedia.org', 'fa.wiktionary.org', 'id.wikisource.org', 'ht.wikipedia.org', 'tr.wikisource.org', 'yi.wikipedia.org', 'sw.wikipedia.org', 'skr.wiktionary.org', 'n
ia.wiktionary.org', 'da.wikipedia.org', 'ban.wikipedia.org', 'ms.wiktionary.org', 'te.wikisource.org', 'hr.wikipedia.org', 'br.wikipedia.org', 'ru.wikinews.org', 'es.wikiquote.org', 'nso.wikipedia.org', 'sr.wi
ktionary.org', 'ar.wikipedia.org', 'hi.wikisource.org', 'id.wikipedia.org', 'it.wikisource.org', 'meta.wikimedia.org', 'or.wikipedia.org', 'nl.wikipedia.org', 'ko.wikisource.org', 'zh-min-nan.wiktionary.org',
'uk.wiktionary.org', 'www.mediawiki.org', 'sv.wiktionary.org', 'cs.wiktionary.org', 'lv.wikipedia.org', 'ary.wikipedia.org', 'fa.wikipedia.org', 'mai.wikipedia.org', 'dty.wikipedia.org', 'ps.wikipedia.org', 'm
l.wikipedia.org', 'en.wikisource.org', 'en.wikinews.org', 'gu.wikisource.org', 'tg.wiktionary.org', 'hs.wikipedia.org', 'bug.wikipedia.org', 'commons.wikimedia.org', 'az.wikipedia.org', 'lt.wikipedia.org', 'os
.wikipedia.org', 'cv.wikipedia.org', 'sk.wikipedia.org', 'mr.wikipedia.org', 'ja.wikipedia.org', 'kw.wikipedia.org', 'gu.wikipedia.org', 'he.wiktionary.org', 'su.wikipedia.org', 'pt.wikipedia.org', 'kk.wikiped
ia.org', 'ar.wikisource.org', 'ms.wikipedia.org', 'it.wikiquote.org', 'bs.wikipedia.org', 'te.wikipedia.org', 'pt.wikinews.org', 'fi.wiktionary.org', 'sr.wikipedia.org', 'uz.wikipedia.org', 'nah.wikipedia.org',
'eo.wikipedia.org', 'skr.wikipedia.org', 'de.wikisource.org', 'pl.wiktionary.org', 'oc.wiktionary.org', 'ka.wikipedia.org', 'de.wikivoyage.org', 'sv.wikipedia.org', 'olo.wikipedia.org', 'cs.wikiquote.org', '
it.wikipedia.org', 'sq.wikipedia.org', 'hi.wikipedia.org', 'es.wikipedia.org', 'he.wikisource.org', 'nds.wikipedia.org', 'it.wikivoyage.org', 'fr.wiktionary.org', 'af.wikipedia.org', 'zh-yue.wikipedia.org', 'c
e.wikipedia.org', 'pt.wiktionary.org', 'bcl.wikipedia.org', 'bn.wikipedia.org', 'km.wikipedia.org', 'www.wikidata.org', 'io.wikipedia.org', 'ro.wikipedia.org', 'hif.wikipedia.org', 'de.wiktionary.org', 'la.wik
ipedia.org', 'be-tarask.wikipedia.org', 'ban.wikisource.org', 'udm.wikipedia.org', 'vi.wikipedia.org', 'ta.wikipedia.org', 'lld.wikipedia.org', 'mk.wikipedia.org', 'shn.wikipedia.org', 'de.wikipedia.org', 'ba.
wikipedia.org', 'bcl.wiktionary.org', 'fr.wikipedia.org', 'simple.wikipedia.org', 'et.wikipedia.org', 'sn.wikipedia.org', 'ig.wikipedia.org', 'sv.wikisource.org', 'ceb.wikipedia.org', 'els.wikipedia.org', 'bn.
wiktionary.org', 'my.wikipedia.org', 'arz.wikipedia.org', 'diq.wiktionary.org', 'he.wikipedia.org', 'cy.wikisource.org', 'eu.wikipedia.org', 'ur.wikipedia.org', 'tl.wikipedia.org', 'bd.wikimedia.org', 'ak.wik
ipedia.org', 'pa.wikipedia.org', 'my.wikipedia.org', 'en.wiktionary.org', 'sm.wiktionary.org', 'xmf.wikipedia.org', 'smn.wikipedia.org', 'mn.wikipedia.org', 'pl.wikinews.org', 'ky.wikipedia.org', 'ja.wiktionar
y.org', 'bh.wikipedia.org', 'uk.wikipedia.org', 'zh.wikipedia.org', 'gl.wikipedia.org', 'tt.wikipedia.org', 'fi.wikipedia.org', 'ha.wikipedia.org', 'as.wikipedia.org', 'guw.wikipedia.org', 'is.wikipedia.org', '
kn.wikipedia.org', 'bg.wikipedia.org', 'sah.wikiquote.org', 'tr.wikipedia.org', 'pl.wikisource.org', 'wikisource.org', 'zh.wiktionary.org', 'sl.wikipedia.org', 'en.wikibooks.org', 'nn.wikipedia.org', 'mdf.wik
ipedia.org', 'vi.wiktionary.org', 'el.wiktionary.org', 'pam.wikipedia.org', 'zh.wikisource.org', 'as.wikisource.org', 'vep.wikipedia.org', 'sg.wiktionary.org', 'nl.wiktionary.org', 'no.wikipedia.org', 'li.wiki
news.org', 'incubator.wikimedia.org', 'zsh.wikipedia.org', 'en.wikipedia.org', 'en.wikiquote.org', 'bat-smg.wikipedia.org', 'sat.wikipedia.org', 'uk.wikisource.org', 'ar.wiktionary.org', 'ru.wiktionary.org', '
ml.wikipedia.org', 'fa.wikisource.org', 'test.wikipedia.org', 'bg.wiktionary.org', 'bja.wikipedia.org', 'ku.wikipedia.org', 'sr.wikisource.org', 'ca.wikipedia.org', 'de.wikiiversity.org', 'hu.wikipedia.org', 't
h.wikipedia.org', 'pt.wikiiversity.org', 'pnb.wikipedia.org', 'ja.wikisource.org', 'hi.wikiquote.org', 'sa.wikipedia.org', 'en.wikivoyage.org', 'it.wiktionary.org', 'es.wiktionary.org', 'uk.wikinews.org', 'bn.w
ikisource.org', 'ru.wikipedia.org', 'es.wikibooks.org', 'tcy.wikipedia.org', 'be.wikipedia.org', 'species.wikimedia.org', 'ca.wikipedia.org', 'ko.wiktionary.org', 'hy.wikipedia.org', 'be.wikisource.org', 'ku.w
iktionary.org', 'ceb.wikipedia.org', 'si.wikipedia.org', 'avk.wikipedia.org', 'zh-classical.wikipedia.org', 'sk.wiktionary.org', 'es.wikisource.org', 'ast.wikipedia.org', 'fr.wikinews.org', 'ko.wikipedia.org',
'jv.wikipedia.org', 'cy.wikipedia.org', 'pl.wikipedia.org', 'fr.wikisource.org', 'th.wiktionary.org', 'id.wikimedia.org']

-----

task 2:
[['Module:sce-pron', 'en.wiktionary.org'], ['Template:sce-pron', 'en.wiktionary.org']]

-----

task 3:
160

-----

task 4:
['Islamophobia_and_Violence_against_Muslims_in_India', 'en.wikipedia.org', '2022-06-13 00:53:23.660000']

-----

task 5:
[[1121688, 'Denelson83', 105], [3083127, 'DrThneed', 90], [2004310, 'Marcomogollon', 10], [1043913, 'New user message', 1], [193075, 'Doncram', 1], [15078, 'Azmi1995', 2], [3006008, 'Hugo999', 1], [83322, 'Alex
anders', 1], [11260960, 'AskeBot', 20], [604824, 'Alemann', 1], [2572341, 'Benezius', 2], [063971, 'Vojtech Dostál', 119], [30626738, 'Adamtine123', 1], [18185, 'Krutvyuss', 1], [57645, 'Songhongyi', 1], [10
676741, 'Masileann', 1], [11325218, 'Micred', 1], [588882, 'BotMultichill1', 9], [111660, 'New user message', 1], [988148, 'Ixoactus', 1], [14926857, 'Gerald Perer', 1], [4356, 'Stura', 3], [26498, 'B3333 B3
3333 B3333', 1], [644015, 'Quangdat201', 1], [2996010, 'Nkywuong', 2], [41345, 'Atheist Armenian', 1], [85347, 'Apsite', 1], [71355, 'Yanguas', 1], [496876, 'GnuBotmarcoo', 6], [493719, 'whoop whoop pull up'
, 1], [603378, 'Christian Ferrer', 1], [2728610, 'DeltaBot', 2], [57122, 'Dcirovic', 4], [43117470, 'Ishiura', 2], [302461, 'Wikimedia Commons Welcome', 3]]
```