**Events**

| UserID. 1009 | UserID: 901 | UserID: 1408 |
|---|---|---|
| Login | Update picture | Add skill set |
| IP, Browser, etc | IP, Browser, etc | IP, Browser, etc |

**Server 1**
- Accept
- Process

**Server 2**
- Accept
- Process

. . . .

**Events Topic**

| UserID 1-1000 | UserID 1001-2000 | . . . . . . . . | UserID n-n+log | Alerts Topic |
|---|---|---|---|---|

**Kafka**

→ The figure above depicts the use of kafka to implement LinkedIn user profile experience. The normal operations possible here are login, update profile information, add skill set, add new job, etc.

→ We create 2 topics here for

example:

### 1) Events topic

This will service all the requests for profile updates or operations. As shown in figure there are separate partitions for users, which * sends re different requests from same user to the same partition.

### 2) Alerts topic

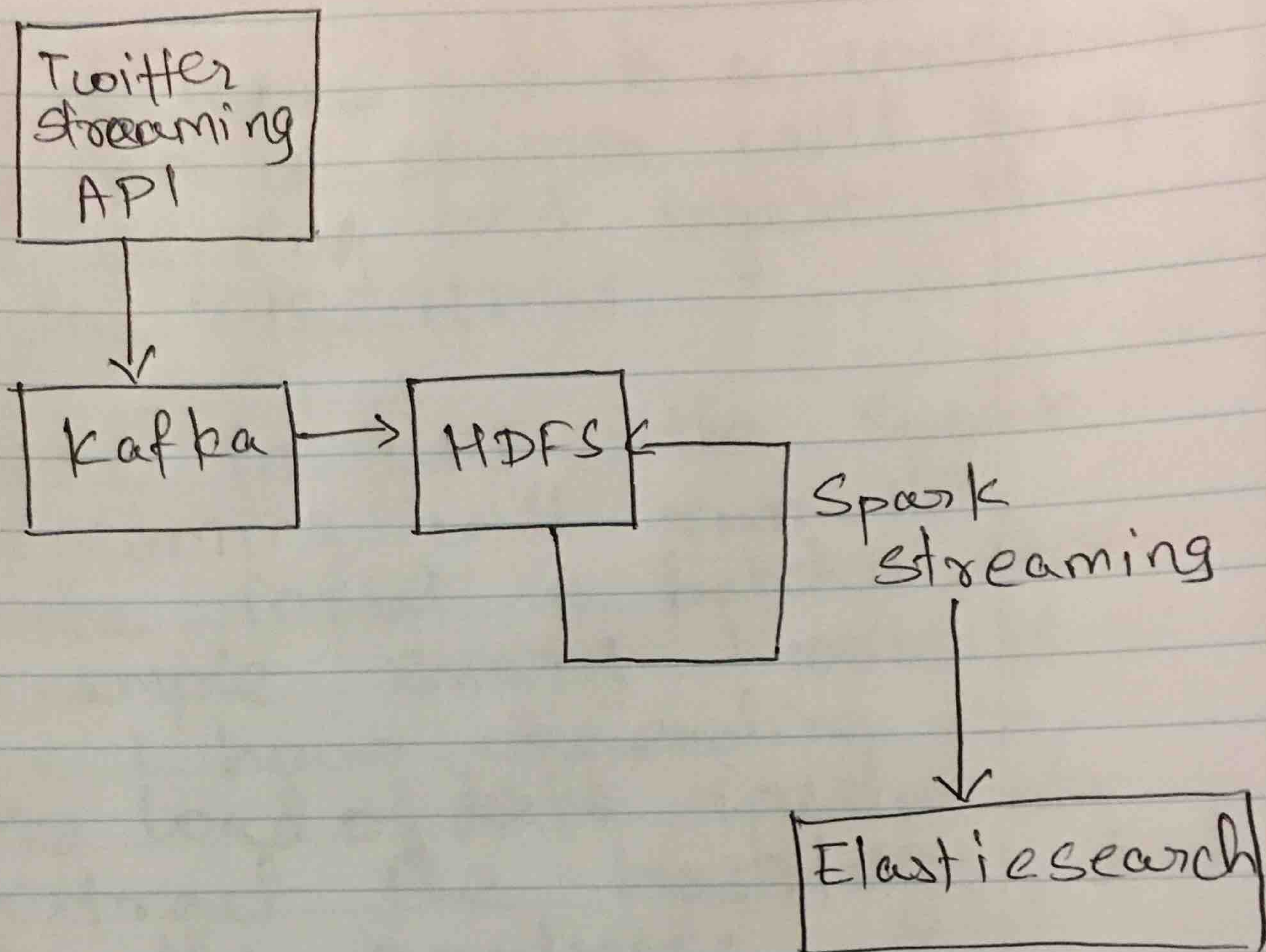This topic is basically for event logging and ingesting data into Hadoop for analytics.

→ As shown in the figure as soon as user sends a request by performing an action, it goes to the nearest server first which then creates an event and adds into Kafka.

→ The events gets processed in

events topic. In each server there is a publisher which publishes events & there is consumer which takes responses. If All the log data, these consumer sends to Alerts topic where there can be another consumer to take those log and add into HDFS for analytics.

⇒ Since kafka is scalable, this design is highly scalable

```
┌──────────────┐
│  Twitter     │
│  streaming   │
│  API         │
└──────────────┘
        │
        ▼
┌──────────┐      ┌──────────┐
│  kafka   │ ───→ │  HDFS ◄──┼──────┐
└──────────┘      └──────┬───┘      │   Spark
                         │          │   streaming
                      ┌──┴───────┐  │
                      │          │  │
                      └──────────┘  │
                                    ▼
                        ┌──────────────────┐
                        │  Elastic search  │
                        └──────────────────┘
```

→ As shown in the figure the twitter streaming API will be used to fetch the real time data. Suppose we use flask web app to download tweets using twitter streaming API.

→ This Flask App will keep adding tweets into kafka as events. Kafka will then publish

to a topic which a consumer
no on a server will keep
consuming and ingest the
data into HDFS.

→ Here in HDFS the spark
streaming will run. i.e.
take input in batches, for
example every 1 minute
or 1 hour depending on
the load of data incoming
extract the hashtags, and
do the analytic

→ Finally the tweets will
be stored in elastic
search.

→ links = # RDD of (url, neighbors)
This will contain all the
urls and its neighbors
i.e. A, C = A → C

→ ranks = # RDD of (url, rank)
This will contain the rank
of each url.
i.e. A, 1

→ [url, [links, rank]] = pair
return [(dest, rank/len(links)) for
dest in links]

On each iteration, have page P
send a contribution of
rank(P)/num_of_neighbors(P)
to its neighbors.

↑
→ Contribs = links.join(ranks).flatMap
(comput_contrib

→ ranks = contribs. reduceByKey
  (lambda x, y: x+y) \
  .mapValues (lambda x: 0.15+0.85 **#
                                   **#)

ranks. SaveAsTextFile (...)
↑
→ Set each page's rank to
  0.15 + 0.85 × Contributions received

→ Then save the result in
  text file.

→ The RDD generated

  ranks = (url, rank).