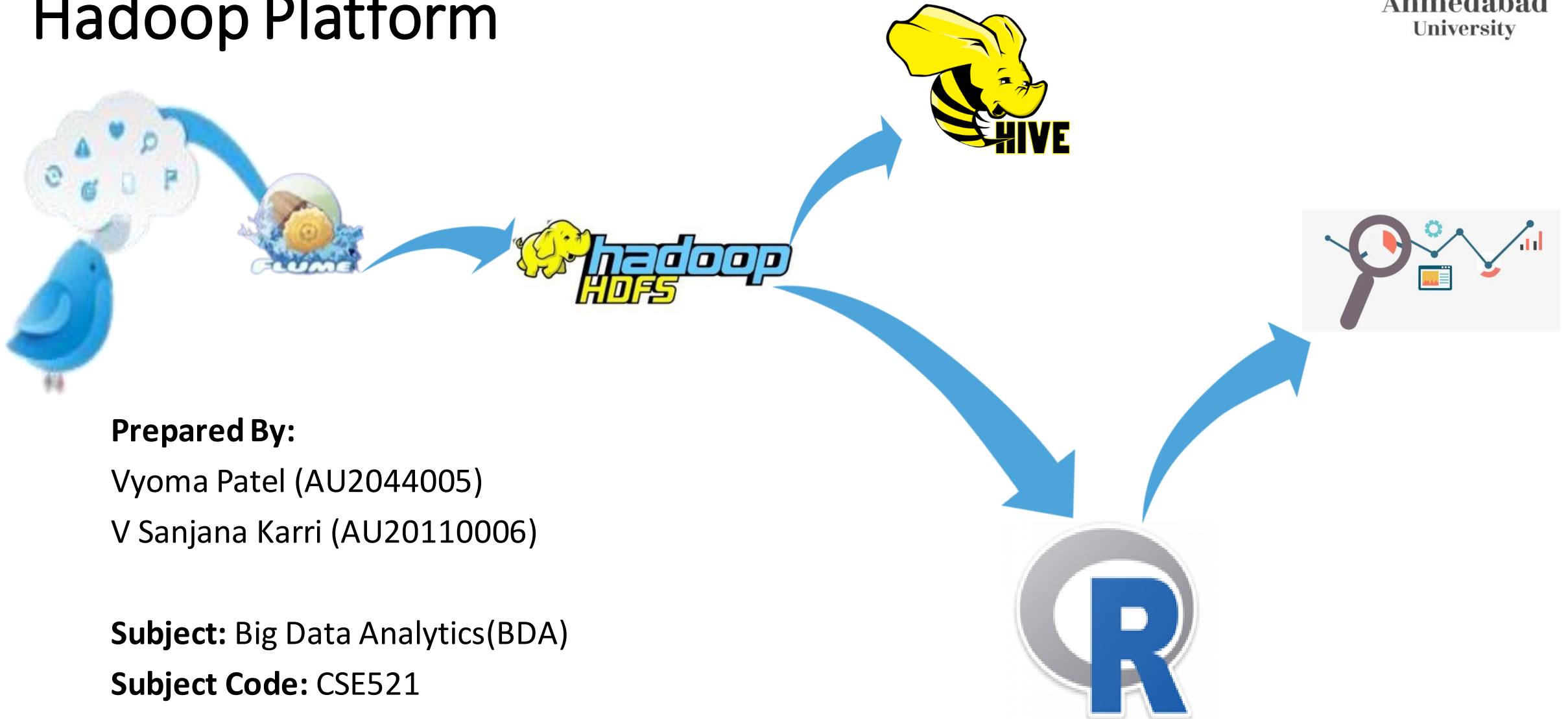


# Twitter Data Analysis with R on top of Hadoop Platform



Ahmedabad  
University



## Prepared By:

Vyoma Patel (AU2044005)

V Sanjana Karri (AU20110006)

**Subject:** Big Data Analytics(BDA)

**Subject Code:** CSE521

# Proposed Methodology

01

## Twitter Application

- Creating twitter developer account
- Generating API keys and tokens

02

## Installation Part

- Hadoop-3.1.4
- Flume-1.9.0
- Hive-3.1.2

03

## Twitter to HDFS

- Create twitter config file
- Start dfs and call twitter agent via flume

04

## Hadoop Cluster

- Flume Data generation
- Data in AVRO format

05

## Changing Data format

- Calling Avro tool, getting schema
- Converting Avro to JSON

06

## Hive

- Starting metastore
- Table Creation and Retrieval

07

## Connecting with R

- JSON file generated sent
- Data Preprocessing

08

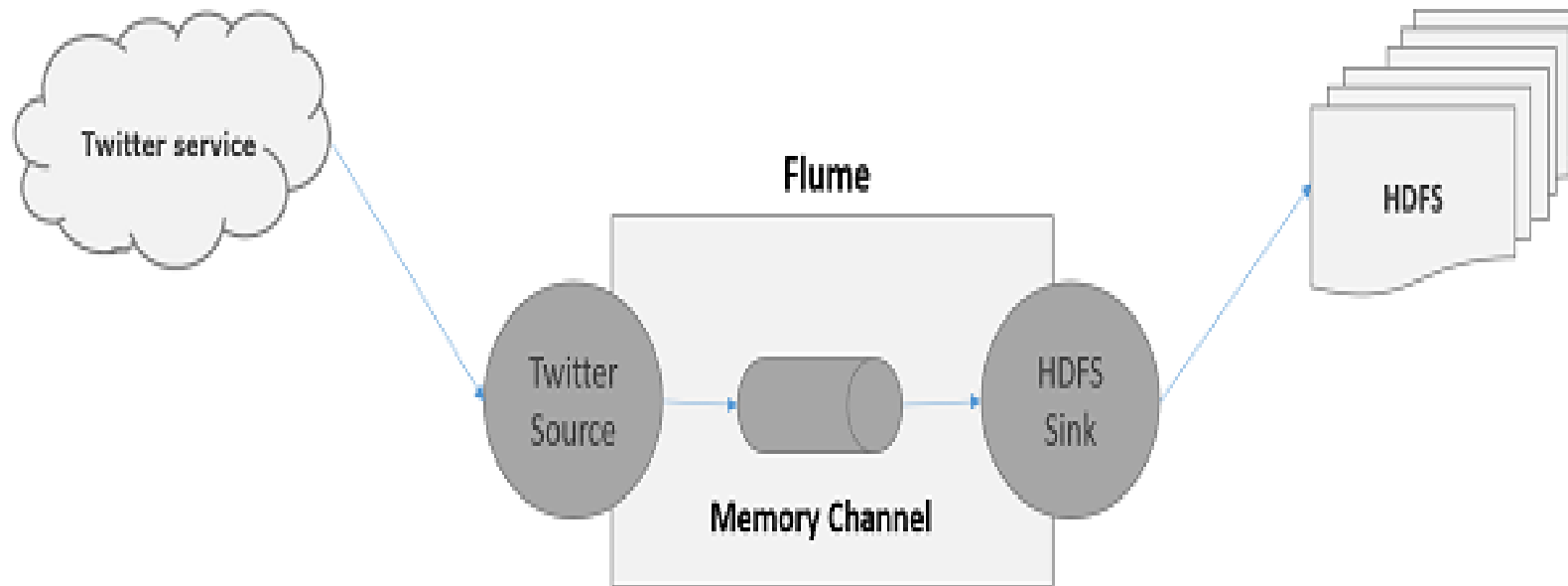
## Visualizations

- Using appropriate libraries
- Analysing data with various plots

# Basic Definitions

- **Apache Hadoop** - software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models.
- **Apache Flume** – It's a tool for data ingestion in HDFS. To capture streaming data from various web servers to HDFS.
- **Apache Hive** - data warehouse software facilitates reading, writing, and managing large datasets residing in distributed storage using SQL.
- **Apache Avro** - is a data serialization system. Provides compact, fast, binary data format on Hadoop.
- **JSON** - JavaScript Object Notation. Lightweight format for storing and transporting data. Widely supported and easily understandable.
- **R** - is a programming language and free software environment for statistical computing and graphics supported by the R Foundation for Statistical Computing.

## Process of Tweets stored on Hadoop Cluster



## Twitter Config File

The flume agent has 3 components: source, sink, and channel.

- **Source:** It accepts the data from the incoming streamline and stores the data in the channel.
- **Channel:** In general, the reading speed is faster than the writing speed. Channel acts as the local storage(buffer) or temporary storage between the source of data and persistent data in the HDFS.
- **Sink:** Collects the data from the channel and commits or writes the data in the HDFS permanently.

```
#####
# Twitter agent for collecting Twitter data to HDFS.
#####
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS
#####
# Describing and configuring the sources
#####
TwitterAgent.sources.Twitter.type = org.apache.flume.source.twitter.TwitterSource
TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sources.Twitter.consumerKey = pKyqnqrQSVcPIJFLBeyumm80F
TwitterAgent.sources.Twitter.consumerSecret = UwSIdQnKtAle69rGwlf00AfEWxzHgIcG4QDhsrzo2VyXRRxAXX
TwitterAgent.sources.Twitter.accessToken = 1320030990685777921-3U05qa1VJ7qNQQpv3gvgEHDjH5fb5R
TwitterAgent.sources.Twitter.accessTokenSecret = 4g7NDwt4v5tRBspLt0Sm6QfADfPlxZ0tPi86pox3eZLrS
TwitterAgent.sources.Twitter.keywords = covid-19,corona,aid,virus,India,COVID,help,sos,shortage,crisis,shutdown,lockdown,lack
#####
# Twitter configuring HDFS sink
#####
TwitterAgent.sinks.HDFS.channel = MemChannel
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = hdfs://localhost:9000/user/Hadoop/twitter_data/
TwitterAgent.sinks.HDFS.hdfs.fileType = DataStream
# Test the hdfs.writeFormat to json
TwitterAgent.sinks.HDFS.
TwitterAgent.sinks.HDFS.hdfs.batchSize = 1000
#####
# Twitter Channel
#####
TwitterAgent.channels.MemChannel.type = memory
TwitterAgent.channels.MemChannel.capacity = 1000
TwitterAgent.channels.MemChannel.transactionCapacity = 1000
#####
# Binding the Source and the Sink to the Channel
#####

TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sinks.HDFS.channels = MemChannel
```



## Starting Hadoop and calling flume agent

```
2021-05-04 09:39:17,565 INFO util.GSet: VM type          = 64-bit
2021-05-04 09:39:17,565 INFO util.GSet: 0.25% max memory 828.5 MB = 2.1 MB
2021-05-04 09:39:17,565 INFO util.GSet: capacity          = 2^18 = 262144 entries
2021-05-04 09:39:17,581 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
2021-05-04 09:39:17,581 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
2021-05-04 09:39:17,582 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
2021-05-04 09:39:17,587 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
2021-05-04 09:39:17,588 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry time is 600000 millis
2021-05-04 09:39:17,592 INFO util.GSet: Computing capacity for map NameNodeRetryCache
2021-05-04 09:39:17,592 INFO util.GSet: VM type          = 64-bit
2021-05-04 09:39:17,592 INFO util.GSet: 0.029999999329447746% max memory 828.5 MB = 254.5 KB
2021-05-04 09:39:17,592 INFO util.GSet: capacity          = 2^15 = 32768 entries
2021-05-04 09:39:17,742 INFO namenode.FSImage: Allocated new BlockPoolId: BP-812740381-127.0.1.1-1620101357664
2021-05-04 09:39:17,914 INFO common.Storage: Storage directory /home/hadoop/hdfs/namenode has been successfully formatted.
2021-05-04 09:39:17,968 INFO namenode.FSImageFormatProtobuf: Saving image file /home/hadoop/hdfs/namenode/current/fsimage.ckpt_00000000000000000000 using no compression
2021-05-04 09:39:18,196 INFO namenode.FSImageFormatProtobuf: Image file /home/hadoop/hdfs/namenode/current/fsimage.ckpt_00000000000000000000 of size 390 bytes saved in 0 seconds .
2021-05-04 09:39:18,268 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
2021-05-04 09:39:18,279 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid = 0 when meet shutdown.
2021-05-04 09:39:18,281 INFO namenode.NameNode: SHUTDOWN_MSG:
/*****
SHUTDOWN_MSG: Shutting down NameNode at hp-HP-EliteBook-8440p/127.0.1.1
*****/
hadoop@hp-HP-EliteBook-8440p:~/hadoop-3.1.4/bin$ cd $HADOOP_HOME/sbin/
hadoop@hp-HP-EliteBook-8440p:~/hadoop-3.1.4/sbin$ start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [hp-HP-EliteBook-8440p]
hadoop@hp-HP-EliteBook-8440p:~/hadoop-3.1.4/sbin$ jps
5186 NameNode
5382 DataNode
5753 Jps
5563 SecondaryNameNode
hadoop@hp-HP-EliteBook-8440p:~/hadoop-3.1.4/sbin$ cd
hadoop@hp-HP-EliteBook-8440p:~$ cd $FLUME_HOME
hadoop@hp-HP-EliteBook-8440p:/home/hp/apache-flume-1.9.0-bin$ bin/flume-ng agent --conf ./conf/ -f conf/new_twitter.conf -Dflume.root.logger=D
EBUG,console -n TwitterAgent
```

# Tweets stored in Avro format using Flume on Hadoop Cluster

127.0.0.1:9870/explorer.html#/user/Hadoop/twitter\_data/

/user/Hadoop/twitter\_data/ Go!

Show 25 entries Search:

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	hadoop	supergroup	1.17 KB	May 04 17:14	1	128 MB	FlumeData.1620128690366
-rw-r--r--	hadoop	supergroup	8.23 KB	May 04 17:14	1	128 MB	FlumeData.1620128690367
-rw-r--r--	hadoop	supergroup	56.62 KB	May 04 17:14	1	128 MB	FlumeData.1620128690368
-rw-r--r--	hadoop	supergroup	37.16 KB	May 04 17:14	1	128 MB	FlumeData.1620128690369
-rw-r--r--	hadoop	supergroup	27.54 KB	May 04 17:14	1	128 MB	FlumeData.1620128690370
-rw-r--r--	hadoop	supergroup	28.82 KB	May 04 17:14	1	128 MB	FlumeData.1620128690371
-rw-r--r--	hadoop	supergroup	28.43 KB	May 04 17:14	1	128 MB	FlumeData.1620128690372
-rw-r--r--	hadoop	supergroup	25.54 KB	May 04 17:14	1	128 MB	FlumeData.1620128690373
-rw-r--r--	hadoop	supergroup	30.78 KB	May 04 17:14	1	128 MB	FlumeData.1620128690374
-rw-r--r--	hadoop	supergroup	25.84 KB	May 04 17:15	1	128 MB	FlumeData.1620128690375
-rw-r--r--	hadoop	supergroup	26.92 KB	May 04 17:15	1	128 MB	FlumeData.1620128690376
-rw-r--r--	hadoop	supergroup	31.8 KB	May 04 17:15	1	128 MB	FlumeData.1620128690377
-rw-r--r--	hadoop	supergroup	31.79 KB	May 04 17:15	1	128 MB	FlumeData.1620128690378

127.0.0.1:9870/explorer.html#/user/Hadoop/twitter\_data/

Browse Directory

/user/Hadoop/twitter\_data/ Go!

Show 25 entries Search:

Name

FlumeData.1620128690366

FlumeData.1620128690367

FlumeData.1620128690368

FlumeData.1620128690369

FlumeData.1620128690370

FlumeData.1620128690371

FlumeData.1620128690372

FlumeData.1620128690373

FlumeData.1620128690374

FlumeData.1620128690375

FlumeData.1620128690376

FlumeData.1620128690377

File Information - FlumeData.1620128690368

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073741827

Block Pool ID: BP-154730217-127.0.1.1-1620128564266

Generation Stamp: 1003

Size: 57979

Availability:

- hp-HP-EliteBook-8440p

File contents

```
Obj[{"type":"record","name":"Doc","doc":{"adoc":{"fields":[{"name":"id","type":"string"}, {"name":"user_friends_count","type":["int","null"]}, {"name":"user_location","type":["string","null"]}, {"name":"user_description","type":["string","null"]}, {"name":"user_statuses_count","type":["int","null"]}, {"name":"user_followers_count","type":["int","null"]}, {"name":"user_name","type":["string","null"]}, {"name":"user_screen_name","type":["string","null"]}, {"name":"created_at","type":["string","null"]}, {"name":"text","type":["string","null"]}, {"name":"retweet_count","type":["long","null"]}, {"name":"retweeted","type":
```



## Storing Data generated using flume agent in a folder and Avro initialization

```
hadoop@hp-HP-EliteBook-8440p:~$ hdfs dfs -copyToLocal /user/Hadoop/twitter_data/FlumeData.1620128690368 /tmp
hadoop@hp-HP-EliteBook-8440p:~$ java -jar /home/hp/apache-flume-1.9.0-bin/lib/avro-tools-1.10.2.jar
Version 1.10.2 of Apache Avro
Copyright 2010-2015 The Apache Software Foundation
```

```
This product includes software developed at
The Apache Software Foundation (https://www.apache.org/).
```

-----  
Available tools:

canonical	Converts an Avro Schema to its canonical form
cat	Extracts samples from files
compile	Generates Java code for the given schema.
concat	Concatenates avro files without re-compressing.
count	Counts the records in avro files or folders
fingerprint	Returns the fingerprint for the schemas.
fragtojson	Renders a binary-encoded Avro datum as JSON.
fromjson	Reads JSON records and writes an Avro data file.
fromtext	Imports a text file into an avro data file.
getmeta	Prints out the metadata of an Avro data file.
getschema	Prints out schema of an Avro data file.
idl	Generates a JSON schema from an Avro IDL file
idl2schemata	Extract JSON schemata of the types from an Avro IDL file
induce	Induce schema/protocol from Java class/interface via reflection.
jsontofrag	Renders a JSON-encoded Avro datum as binary.
random	Creates a file with randomly generated instances of a schema.
recodec	Alters the codec of a data file.
repair	Recovers data from a corrupt Avro Data file
rpcprotocol	Output the protocol of a RPC service
rpcreceive	Opens an RPC Server and listens for one message.
rpcsend	Sends a single RPC message.
tether	Run a tethered mapreduce job.
tojson	Dumps an Avro data file as JSON, record per line or pretty.
totext	Converts an Avro data file to a text file.
totrevni	Converts an Avro data file to a Trevni file.
trevni_meta	Dumps a Trevni file's metadata as JSON.
trevni_random	Create a Trevni file filled with random instances of a schema.
trevni_tojson	Dumps a Trevni file as JSON.

```
hadoop@hp-HP-EliteBook-8440p:~$
```



# Getting Schema of Data Generated in Avro Format

```
hadoop@hp-HP-EliteBook-8440p:~$ java -jar /home/hp/apache-flume-1.9.0-bin/lib/avro-tools-1.10.2.jar getschema /tmp/flumedata.avro
21/05/04 17:38:36 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
{
  "type" : "record",
  "name" : "Doc",
  "doc" : "adoc",
  "fields" : [ {
    "name" : "id",
    "type" : "string"
  }, {
    "name" : "user_friends_count",
    "type" : [ "int", "null" ]
  }, {
    "name" : "user_location",
    "type" : [ "string", "null" ]
  }, {
    "name" : "user_description",
    "type" : [ "string", "null" ]
  }, {
    "name" : "user_statuses_count",
    "type" : [ "int", "null" ]
  }, {
    "name" : "user_followers_count",
    "type" : [ "int", "null" ]
  }, {
    "name" : "user_name",
    "type" : [ "string", "null" ]
  }, {
    "name" : "user_screen_name",
    "type" : [ "string", "null" ]
  }, {
    "name" : "created_at",
    "type" : [ "string", "null" ]
  }, {
    "name" : "text",
    "type" : [ "string", "null" ]
  }, {
    "name" : "retweet_count",
    "type" : [ "int", "null" ]
  }, {
    "name" : "retweeted",
    "type" : [ "boolean", "null" ]
  }, {
    "name" : "in_reply_to_user_id",
    "type" : [ "long", "null" ]
  }, {
    "name" : "source",
    "type" : [ "string", "null" ]
  }, {
    "name" : "in_reply_to_status_id",
    "type" : [ "long", "null" ]
  }, {
    "name" : "media_url_https",
    "type" : [ "string", "null" ]
  }, {
    "name" : "expanded_url",
    "type" : [ "string", "null" ]
  }
]
```

```
    "name" : "user_followers_count",
    "type" : [ "int", "null" ]
  }, {
    "name" : "user_name",
    "type" : [ "string", "null" ]
  }, {
    "name" : "user_screen_name",
    "type" : [ "string", "null" ]
  }, {
    "name" : "created_at",
    "type" : [ "string", "null" ]
  }, {
    "name" : "text",
    "type" : [ "string", "null" ]
  }, {
    "name" : "retweet_count",
    "type" : [ "long", "null" ]
  }, {
    "name" : "retweeted",
    "type" : [ "boolean", "null" ]
  }, {
    "name" : "in_reply_to_user_id",
    "type" : [ "long", "null" ]
  }, {
    "name" : "source",
    "type" : [ "string", "null" ]
  }, {
    "name" : "in_reply_to_status_id",
    "type" : [ "long", "null" ]
  }, {
    "name" : "media_url_https",
    "type" : [ "string", "null" ]
  }, {
    "name" : "expanded_url",
    "type" : [ "string", "null" ]
  }
]
```

## Converting Avro to JSON Data format

```
hadoop@hp-HP-EliteBook-8440p:~$ java -jar /home/hp/apache-flume-1.9.0-bin/lib/avro-tools-1.10.2.jar tojson /tmp/flumedata.avro > /tmp/stream.json
hadoop@hp-HP-EliteBook-8440p:~$ hadoop fs -cat file:/tmp/stream.json
{"id":"1389546533041934336","user_friends_count":{"int":1050},"user_location":null,"user_description":{"string":"존나 좆대있게 살아야지"},"user_statuses_count":{"int":8126},"user_followers_count":{"int":1069},"user_name":{"string":"채정체"},"user_screen_name":{"string":"endxning_neve r"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"@bfzapAoGiEGzaq8 @Esgubne_Chicken 웁"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1363329983872073731},"source":{"string":"<a href='\"http://twitter.com/download/iphone\"' rel='\"nofollow\"'>Twitter for iPhone</a>"},"in_reply_to_status_id":{"long":1389546495641272328},"media_url_https":null,"expanded_url":null}
{"id":"1389546533029249036","user_friends_count":{"int":49},"user_location":{"string":"\ud83c\udf7c2020.07.24~ \ud83c\udf7c2021.03.09~"},"user_description":{"string":"#15cm인형 #태공이 #서공이 #20cm인형 #하공이 #지공이 울캄씨기들\ud83d\udc9b\ud83d\udc9b 기어온 우리콩야즈 내입에 넣어 허로 굴려"},"user_statuses_count":{"int":80},"user_followers_count":{"int":4},"user_name":{"string":"콩야"},"user_screen_name":{"string":"__ggon gya"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"RT @__3701 : • 실소비 15cm,\nDM으로 주문 주시면 됩니다 ☺**.*• https://t.co/aRrc03JpDL"},"retweet_count":{"long":0},"retweeted":{"boolean":true},"in_reply_to_user_id":{"long":-1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0gtv-vVKAiPcGT.jpg"},"expanded_url":{"string":"https://twitter.com/__3701/status/1389410346499076099/photo/1"}}
{"id":"1389546533050322947","user_friends_count":{"int":4815},"user_location":{"string":"SHE%HER ; 8TEEN!"},"user_description":{"string":":~ #JAEMIN #JENO #DOYOUNG :: 私はこの世界の何よりもこの男の子を愛しています。"},"user_statuses_count":{"int":64727},"user_followers_count":{"int":4850},"user_name":{"string":"nåtta"},"user_screen_name":{"string":"NCMNCORE"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"@_sunxtety iya sna @doyoung \ud83e\uddd17"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1290710686650122240},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":1389542771703222272},"media_url_https":null,"expanded_url":null}
{"id":"1389546533050208257","user_friends_count":{"int":251},"user_location":{"string":"Shahpura, India"},"user_description":{"string":"I don't have dirty mind, I have sexy imagination. °\ud83d\ude0E\nmusical, style..."},"user_statuses_count":{"int":951},"user_followers_count":{"int":40},"user_name":{"string":"Subhash choudhary"},"user_screen_name":{"string":"Imsubhash_jat"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"RT @poetdushyant: रहनुमाओकी आओपे फिदा है दुनिया,\nइंस बहक्ती हुई दुनिया को समालो यारो |\n\n- दुष्यंत कुमार\n#BengalViolence"},"retweet_count":{"long":0},"retweeted":{"boolean":true},"in_reply_to_user_id":{"long":-1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":null,"expanded_url":null}
{"id":"138954653305458210","user_friends_count":{"int":238},"user_location":{"string":"SW-8122-6783-0813"},"user_description":{"string":"CEO of Rheagardleth and FeMU/Cordelia | ADHD | Non-Binary Ace Lesbian | 25yrs old | She/Her | Pfp of my FeMU Kusanagi and Rheagard hdr by @axe drawings"},"user_statuses_count":{"int":23727},"user_followers_count":{"int":63},"user_name":{"string":"Bonk, assigned Daddy kin"},"user_screen_name":{"string":"BehemothKing1"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"RT @queerDev: @gzdraw https://t.co/Ue4ixM2iSr"},"retweet_count":{"long":0},"retweeted":{"boolean":true},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1"}}
{"id":"1389546533054472192","user_friends_count":{"int":445},"user_location":{"string":"040"},"user_description":{"string":"24/7 genervt von d ir und mir selbst\n\n*Concerts\Cats\P54*Chaos*\n\nmakechesterproud"},"user_statuses_count":{"int":102048},"user_followers_count":{"int":2339},"user_name":{"string":"San"},"user_screen_name":{"string":"Sanny_Me"},"created_at":{"string":"2021-05-04T17:14:44Z"},"text":{"string":"Moving to berlin to move in with my boyfriend this year"},"retweet_count":{"long":0},"retweeted":{"boolean":false},"in_reply_to_user_id":{"long":1},"source":{"string":"<a href='\"http://twitter.com/download/android\"' rel='\"nofollow\"'>Twitter for Android</a>"},"in_reply_to_status_id":{"long":-1},"media_url_https":{"string":"https://pbs.twimg.com/media/E0im0v2UYAMCicj.jpg"},"expanded_url":{"string":"https://twitter.com/queerDev/status/1389543651915751424/photo/1
```



# rHadoop

```
> install.packages("/home/hp/rHadoopClient_0.2.tar.gz")
Installing package into '/home/hp/R/x86_64-pc-linux-gnu-library/4.0'
(as 'lib' is unspecified)
inferring 'repos = NULL' from 'pkgs'
* installing *source* package 'rHadoopClient' ...
** package 'rHadoopClient' successfully unpacked and MD5 sums checked
** using staged installation
** R
** byte-compile and prepare package for lazy loading
** help
*** installing help indices
** building package indices
** testing if installed package can be loaded from temporary location
** testing if installed package can be loaded from final location
** testing if installed package keeps a record of temporary installation path
* DONE (rHadoopClient)
```

## Avro data format read

```
> rHadoopClient::read.hdfs("/user/Hadoop/twitter_data/FlumeData.1620203595094")
V1
Obj\001\002\026avro.schema\xe4
{type:record
His face (especially eyes): 🌟🌟🌟🌟🌟
https://t.co/yYFhU7wnZx
証明されちゃった
鍵50付ける暇があったら勉強したらw https://t.co/9erKh36yRs
It is reported that 90 percent of Covid beds have been occupied....
EP. 5 https://t.co/7ttTtjf4S5
いま行く 待っててね
India
#최현석 #CHOIHYUNSUK https://t.co/KBR1J8iGdF
cc bg
4/23より、ヘザー店舗とWEBストア[.st]にてヘザー商品1万円以上ご購入で、リップティントとポーチをプレゼント🎁
リンクの診断テストでおすすめスタイル&リップをCHECKして👉
真ん中の家おるー!!
원우 도겸 우지 양도 받아올
항상 즉입가능
\002
https://t.co/UsjzIhmtCu
Islamic nigeriA Killed 3.5M+ Christian Biafrans
#BiafraReferendum
INDEPENDENCE FOR BIAFRA
one nigeriA is a FRAUD
#UnityIsNotByForce
#Referendum =! #War... https://t.co/iMoJCva7I
On YouTube: https://t.co/A2tM30KFCg
#FridaysForFuture #Berlin👑 #ClimateCh...
喘んでるのはスルーして
全編はこちら〜〜〜
#RUNaライブ
#RUNaカッット
https://t.co/Cp8tMr2fSC https://t.co/XOCM9bpE4f
```

## JSON Format data read

```
> rHadoopClient::read.hdfs("file:/tmp/stream.json")
V1 V2
1 {id:1389860711581552640 user_friends_count:{int:150}
2 {id:1389860711594176514 user_friends_count:{int:143}
3 {id:1389860711589941253 user_friends_count:{int:1}
4 {id:1389860711602417667 user_friends_count:{int:57}
5 {id:1389860711602421761 user_friends_count:{int:134}
6 {id:1389860711581421573 user_friends_count:{int:392}
7 {id:1389860711577227265 user_friends_count:{int:171}
8 {id:1389860711577366529 user_friends_count:{int:175}
9 {id:1389860711589978116 user_friends_count:{int:1863}
10 {id:1389860711564595203 user_friends_count:{int:974}
11 {id:1389860711589826561 user_friends_count:{int:173}
12 {id:1389860715771531269 user_friends_count:{int:31}
13 {id:1389860715779936261 user_friends_count:{int:163}
14 {id:1389860715767373825 user_friends_count:{int:11}
15 {id:1389860715788333060 user_friends_count:{int:862}
16 {id:1389860715771559938 user_friends_count:{int:152}
17 {id:1389860715758981120 user_friends_count:{int:118}
18 {id:1389860715758964744 user_friends_count:{int:899}
19 {id:1389860715775680513 user_friends_count:{int:2814}
20 {id:1389860715792531457 user_friends_count:{int:377}
21 {id:1389860715796729856 user_friends_count:{int:2691}
22 {id:1389860715767414789 user_friends_count:{int:183}
23 {id:1389860715792670720 user_friends_count:{int:5001}
24 {id:1389860715796725760 user_friends_count:{int:189}
25 {id:1389860715775750146 user_friends_count:{int:25}
```

# Starting Hive

```
hadoop@hp-HP-EliteBook-8440p:~$ mv metastore_db metastore_db.tmp
hadoop@hp-HP-EliteBook-8440p:~$ schematool -dbType derby -initSchema
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hadoop/apache-hive-3.1.2-bin/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop-3.1.4/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Metastore connection URL:      jdbc:derby::databaseName=metastore_db;create=true
Metastore Connection Driver :  org.apache.derby.jdbc.EmbeddedDriver
Metastore connection User:     APP
Starting metastore schema initialization to 3.1.0
Initialization script hive-schema-3.1.0.derby.sql

Initialization script completed
schemaTool completed
hadoop@hp-HP-EliteBook-8440p:~$ hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hadoop/apache-hive-3.1.2-bin/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop-3.1.4/share/hadoop/common/lib/slf4j-log4j12-1.7.25.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = 7f392ec7-ec03-420b-a521-7272e611e7a7

Logging initialized using configuration in jar:file:/home/hadoop/apache-hive-3.1.2-bin/lib/hive-common-3.1.2.jar!/hive-log4j2.properties Async
: true
Hive Session ID = f17f5181-056a-428f-abc9-534354122d65
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> CREATE EXTERNAL TABLE TwitterData
> ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.avro.AvroSerDe'
> WITH SERDEPROPERTIES ('avro.schema.literal'='
> {
>
>   "type" : "record",
>
>   "name" : "Doc",
>
>   "doc" : "adoc",
>
>   "fields" : [ {
>
>     "name" : "id",
```

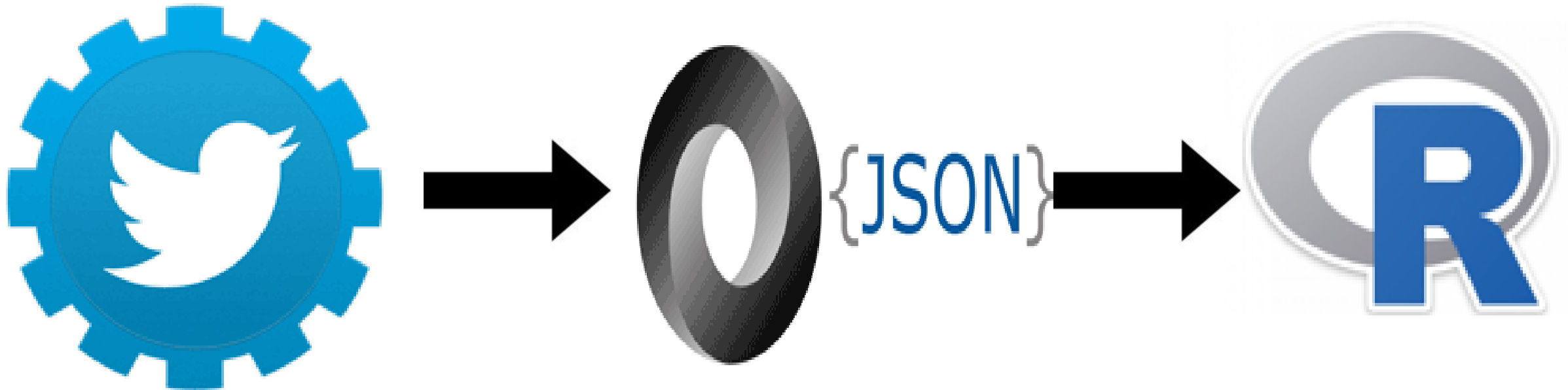


## Creating and Showing Table using Twitter(Avro) Data using Hive

```
> }, {
>   "name" : "media_url_https",
>   "type" : [ "string", "null" ]
> }, {
>   "name" : "expanded_url",
>   "type" : [ "string", "null" ]
> } ]
> }
>
> ' )
> STORED AS
> INPUTFORMAT 'org.apache.hadoop.hive.ql.io.avro.AvroContainerInputFormat'
> OUTPUTFORMAT 'org.apache.hadoop.hive.ql.io.avro.AvroContainerOutputFormat'
> LOCATION '/user/Hadoop/twitter_data'
> ;

OK
Time taken: 2.363 seconds
hive> DESCRIBE TwitterData;
OK
id                string
user_friends_count int
user_location     string
user_description  string
user_statuses_count int
user_followers_count int
user_name         string
user_screen_name  string
created_at        string
text              string
retweet_count     bigint
retweeted         boolean
in_reply_to_user_id bigint
source            string
in_reply_to_status_id bigint
media_url_https   string
expanded_url      string
Time taken: 1.565 seconds, Fetched: 17 row(s)
hive> █
```

## Reconnecting with R



## Cleaned JSON file with tweets in English Language

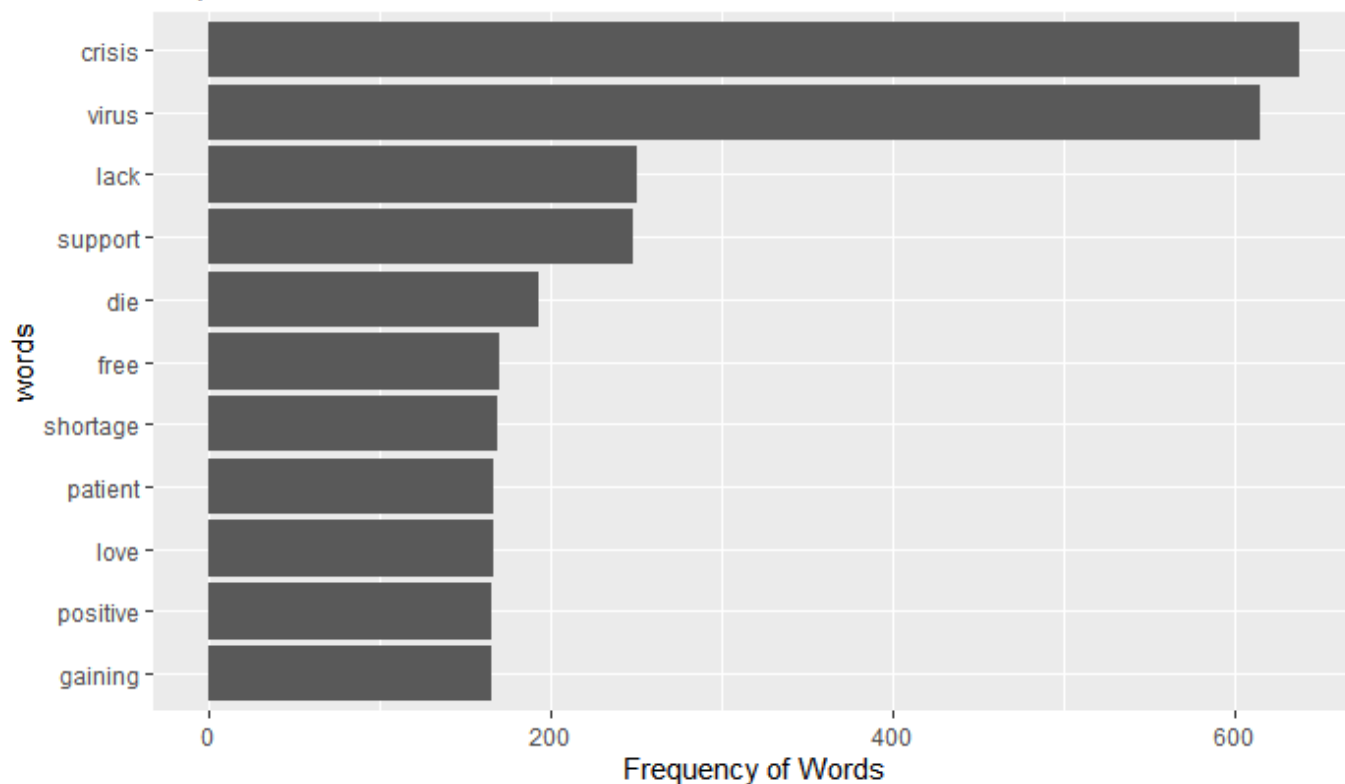
```
{
  "created_at": "Tue May 04 09:43:25 +0000 2021",
  "id": 1389516003801976836,
  "id_str": "1389516003801976836",
  "text": "RT @Sriram37915703: Please provide compensatory attempt to all CSE last"
},
{
  "created_at": "Tue May 04 09:43:25 +0000 2021",
  "id": 1389516004020166658,
  "id_str": "1389516004020166658",
  "text": "RT @pinky_collectio: Name: Okigbo chidera Vera. \nAge: 15years \nLast s"
},
{
  "created_at": "Tue May 04 09:43:25 +0000 2021",
  "id": 1389516004103950336,
  "id_str": "1389516004103950336",
  "text": "RT @sangeetabasa: During times like pandemic, are the rules not same fo"
},
{
  "created_at": "Tue May 04 09:43:25 +0000 2021",
  "id": 1389516003986509826,
  "id_str": "1389516003986509826",
  "text": "#Indien - Fatale Massenveranstaltungen @msargentini @Corriere https://\n"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004187938816,
  "id_str": "1389516004187938816",
  "text": "@manojpandey66 No. Although a protest against a protest for not followi"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004313780224,
  "id_str": "1389516004313780224",
  "text": "RT @YeoBeenmyQueen: THEY ARE SO CUTE \nHELP!!!!\n#VincenzoEp20 https://\n"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004179517440,
  "id_str": "1389516004179517440",
  "text": "RT @RICHA_LAKHERA: Engrave names of each & every one we lost on the"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004380786689,
  "id_str": "1389516004380786689",
  "text": "FUERZA COLOMBIA \ud83c\uddde\ud83c\udddf \u270a\u270a\u270a\n@ONU_es\n"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004250771456,
  "id_str": "1389516004250771456",
  "text": "Dear @MinPres sorry again. I hope you could translate it but it's about"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004519288832,
  "id_str": "1389516004519288832",
  "text": "RT @ReaganGomez: Why this picture\ud83d\ude02",
  "source": "\u003ca href=\n"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004657713153,
  "id_str": "1389516004657713153",
  "text": "RT @AnaCabrera: NEW: Los Angeles reports zero Covid deaths for second s"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004502511617,
  "id_str": "1389516004502511617",
  "text": "@chrislittlew008 @BeamLevels @wonderw12494002 India,Brazil and many cou"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004661899265,
  "id_str": "1389516004661899265",
  "text": "RT @Benarasiyaa: This @CNN report from UP's Meerut is glimpse of the ha"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004842262531,
  "id_str": "1389516004842262531",
  "text": "RT @Giyourdad: \u0019\u0035\u0048\u0001\u0047\u0016\u0032\u0021\u0040\u00"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516005014130694,
  "id_str": "1389516005014130694",
  "text": "RT @ideal_granada: \ud83d\udea8 #\u00daltimaHora Largas colas en los ac"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004552757248,
  "id_str": "1389516004552757248",
  "text": "I'm getting tired of these Liverpool transfer rumours everytime yet we"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516005089677314,
  "id_str": "1389516005089677314",
  "text": "RT @brajeshksingh: \u0015\u004b\u0032\u0015\u003e\u0024\u003e \u0092e\u00"
},
{
  "created_at": "Tue May 04 09:43:26 +0000 2021",
  "id": 1389516004968042497,
  "id_str": "1389516004968042497",
  "text": "RT @ThalaSudharshan: Rudhramadevi \ud83d\ud25\ud83d\ud25\n\nEarned 50"
}
```

```
## {r}  
library(deepIr)  
library(translateR)  
library(dplyr)  
library(tidytext)  
library(ggplot2)  
library(corpus)  
library(tm)  
library(wordcloud)  
library(igraph)  
library(ggraph)
```

## Libraries used for Data Analysis



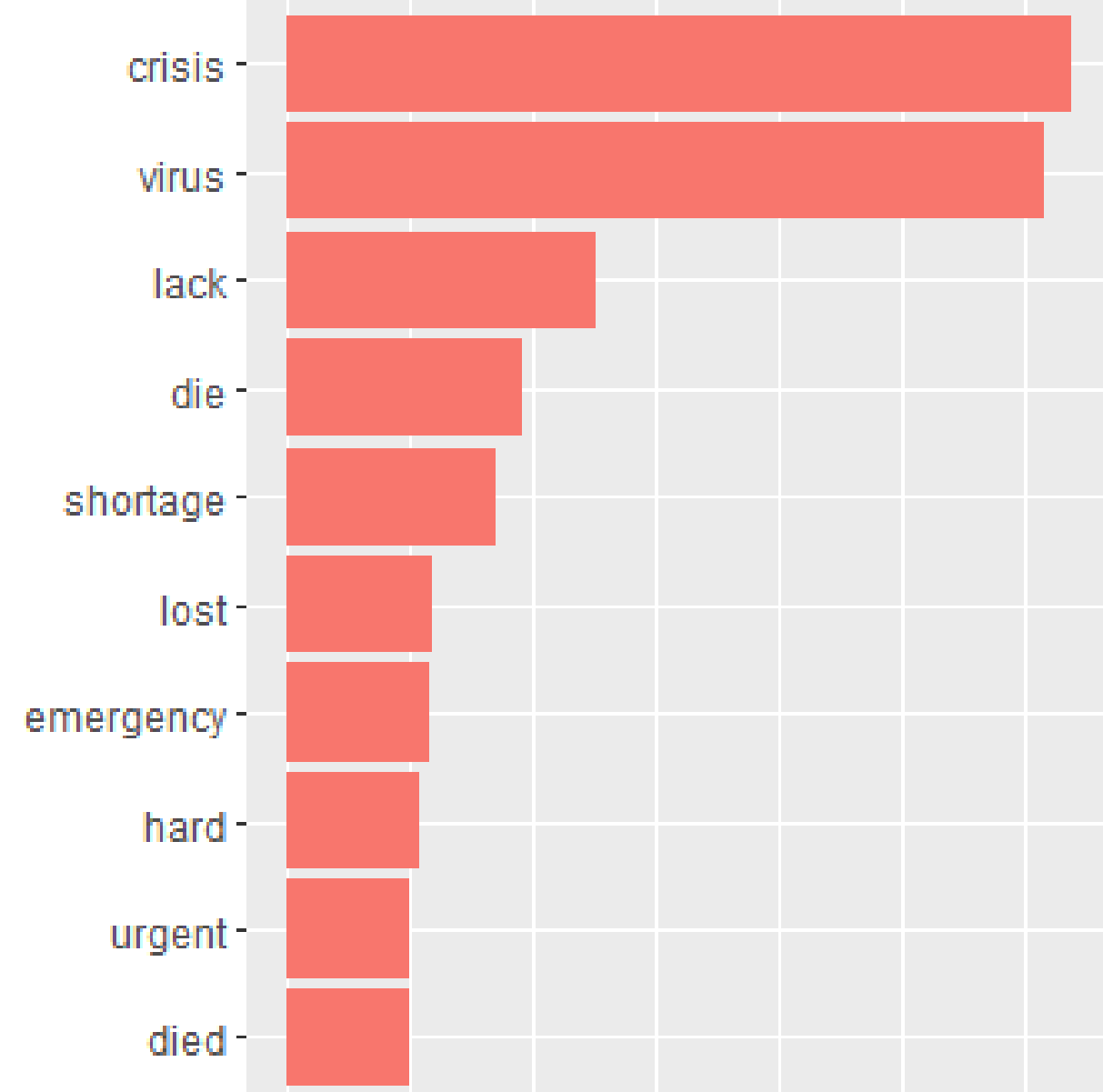
Top 10 most used words in tweets



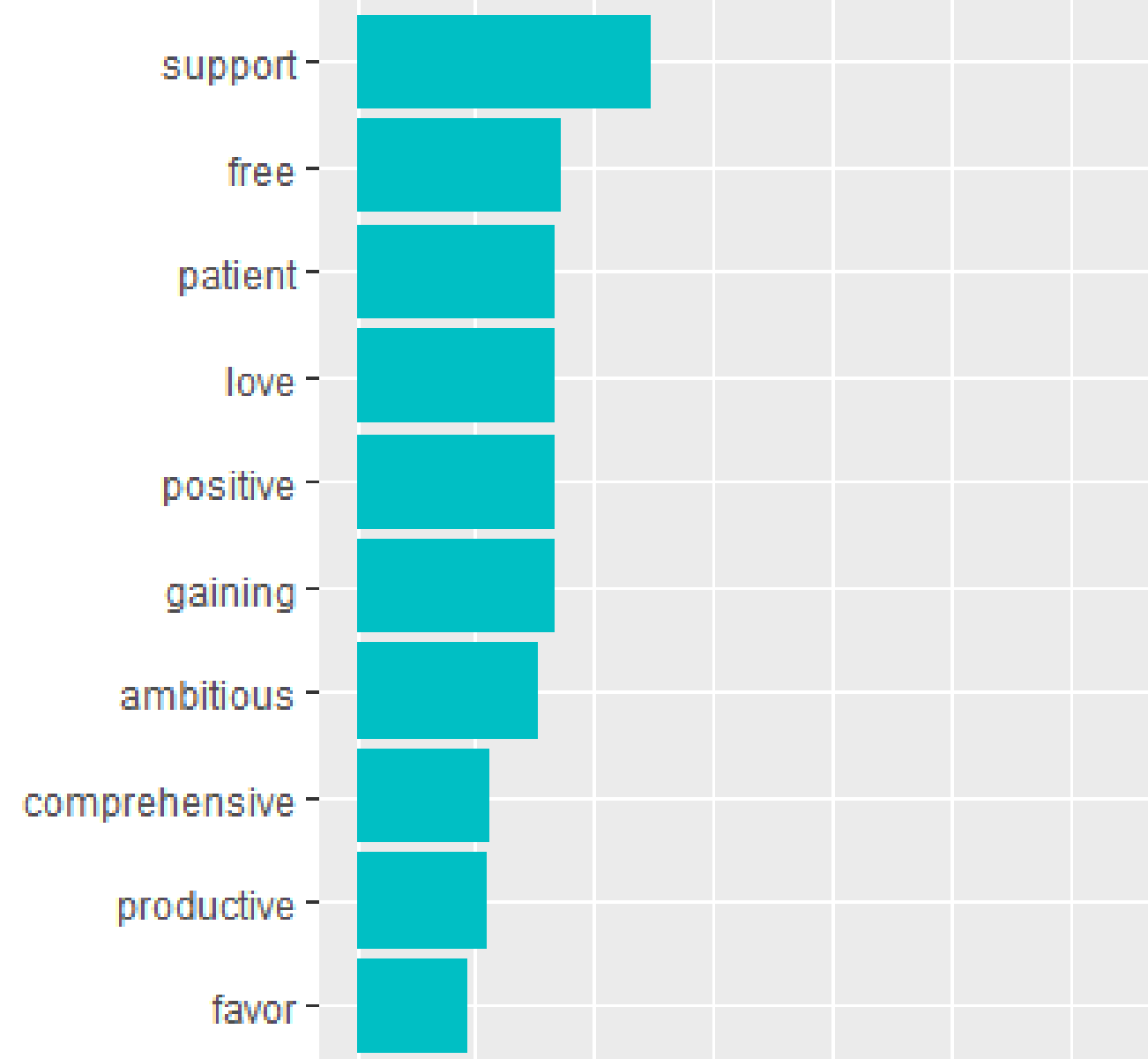
```
```{r}
bing_word_counts %>%
  top_n(10) %>%
  mutate(word = reorder(word, n)) %>%
  ggplot(aes(x = word, y = n)) +
  geom_col() +
  xlab(NULL) +
  coord_flip() +
  labs(y = "Frequency of words",
       x = "words",
       title = "Top 10 most used words in tweets")
```
```

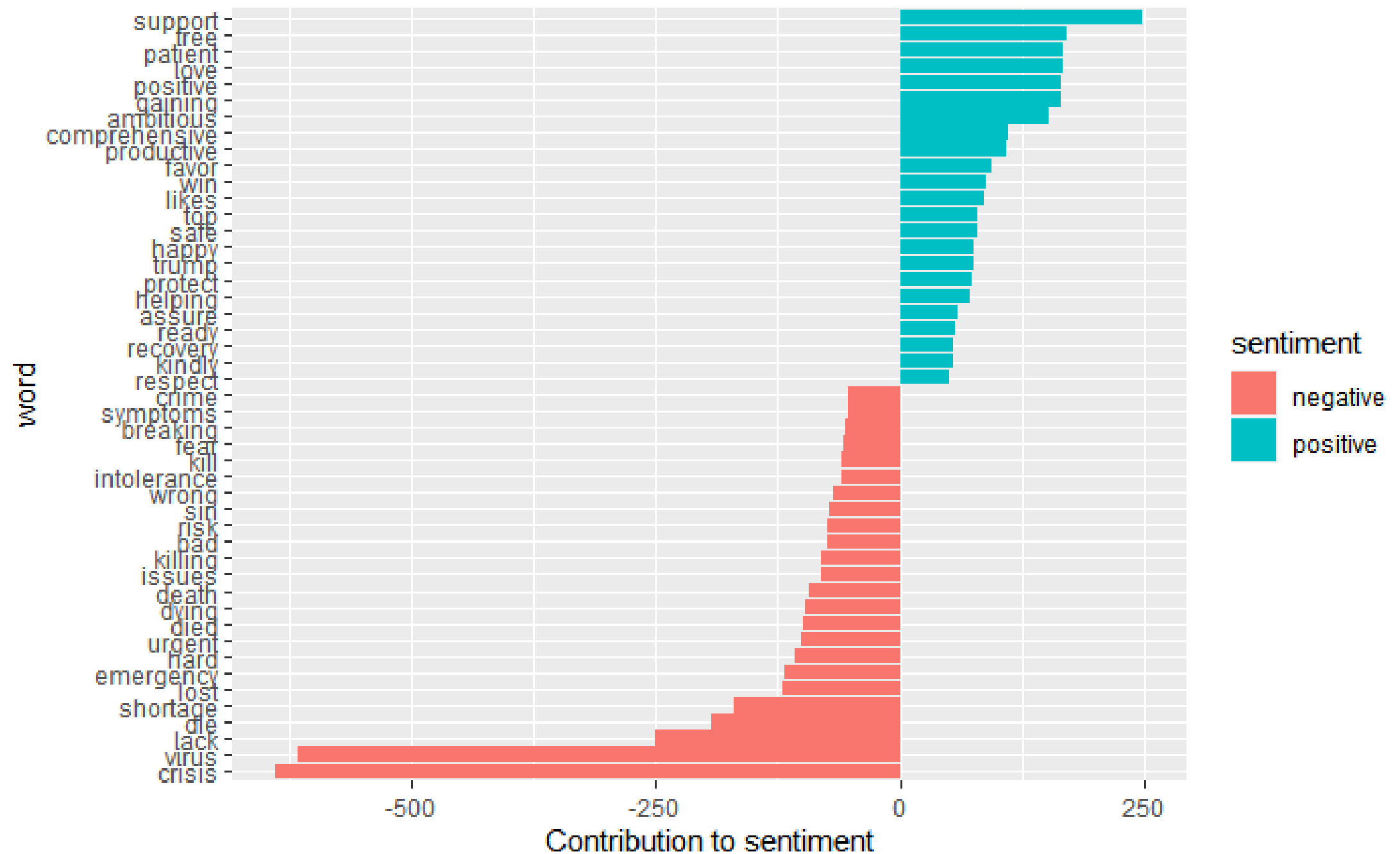
# Sentiment during the COVID.

negative



positive

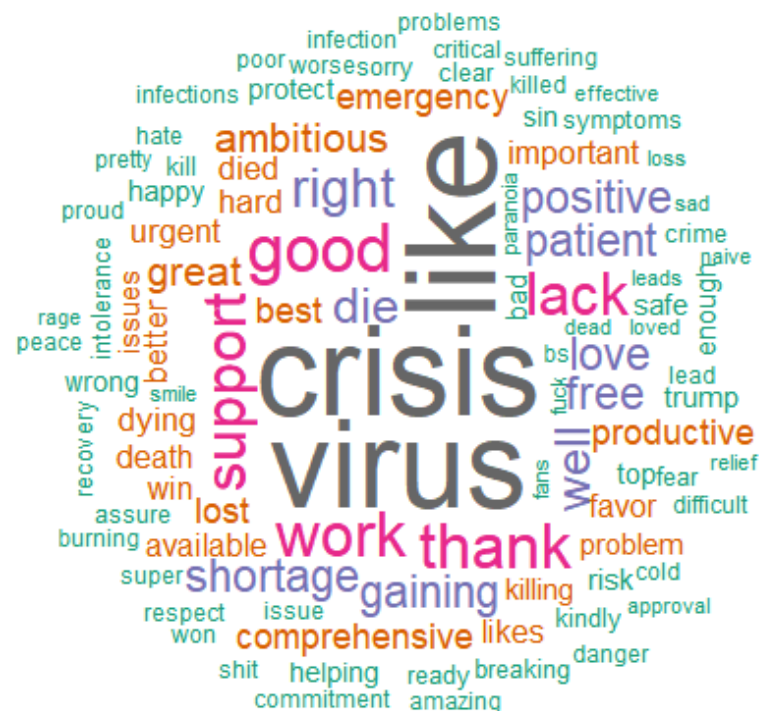




```

```{r}
streamdata_clean %>%
  count(word) %>%
  inner_join(get_sentiments("bing"), by = "word") %>%
  with(wordcloud(word, n, max.words = 100, min.freq = 10,
    random.order = FALSE, colors = brewer.pal(8, "Dark2")))
```

```





## Word Network: Tweets using the hashtag - Covid

Text mining twitter data



crisis virus

strange warning inflated opposition  
funny refuses poor symptoms deadly  
shame lied fake sick vulnerability  
dangerous crime issues issue infected  
worst cold shit hard death fuck refuse  
impose abuse killing suffering problem lost emergency  
shameful sorry risks lack died die  
paranoia false ignore critical sin bs  
threat wrong bs shortage urgent  
severe falling killed problems  
scrap difficult dying worse inaction  
misleading demise cancer intolerance hate stress  
negative pain infections disturbing  
decline worry complaining disaster  
untrue unbelievable scared  
mocked

positive patient free support love gaining ambitious comprehensive

encourage confidence meaningful  
leading effective smile trump safe helping assure ready relief recovery super enjoy worth  
glad trump safe helping assure ready relief recovery super enjoy worth  
commitment appreciated peace loved  
favor won respect leads kindly protect  
win respect leads kindly protect  
productive happy strong helped  
likes fans pretty proud  
trust approval

# References

1. <https://flume.apache.org/FlumeUserGuide.html>
2. <https://www.tutorialspoint.com/apache-flume/fetching-twitter-data.htm>
3. <https://towardsdatascience.com/apache-flume-71ed475eee6d>
4. <https://www.youtube.com/watch?v=PdY31i25SL0>
5. <https://data-flair.training/blogs/apache-hive-installation/>
6. <https://stackoverflow.com/questions/11889261/datanode-process-not-running-in-hadoop>
7. <http://dbmentors.blogspot.com/2017/06/streaming-twitter-data-using-apache.html>
8. <https://medium.com/edureka/apache-flume-tutorial-6f7150210c76>
9. <https://www.h2kinfosys.com/blog/apache-flume-tutorial/>
10. <https://www.confluent.io/blog/avro-kafka-data/#:~:text=Avro%20has%20a%20JSON%20like,in%20a%20compact%20binary%20form.&text=It%20has%20a%20direct%20mapping,inefficient%20for%20high%2Dvolume%20usage>

11. <https://towardsdatascience.com/twitter-sentiment-analysis-and-visualization-using-r-22e1f70f6967>
12. <https://towardsdatascience.com/twitter-data-visualization-fb4f45b63728>
13. <https://codeburst.io/sentiment-analysis-of-twitter-data-359fa9f86bd6>
14. <https://dataaspirant.com/twitter-sentiment-analysis-using-r/>
15. <https://www.r-bloggers.com/2014/04/twitter-sentiment-analysis-with-r/>



# Thank You

