# Lab 2 Data Intensive Computing
Name     - Abhav Luthra
Person No – 50288904

Methodology for Aws Setup

1. Create Bucket on Aws and Input files
2. Schedule EMR job with steps as bellow

**Step 1** will Run 8 Map jobs and 4 reduce jobs



**Step 2** will combine those 4 reduced files to 1



3. Part-00000 file received

Job Flow for Twitter Data

Data Collected Using Rtweet package

```
[Query 1          [Query 2          [Query 3          [Query 4          [Query 5
 Immigrants        Border Wall       Mexico Border     Build Wall        Finish Wall
 Data Cleaned]     Data Cleaned]     Data Cleaned]     Data Cleaned]     Data Cleaned]

Data from Twitter Data from Twitter Data from Twitter Data from Twitter Data from Twitter
Query 1          Query 2          Query 3          Query 4          Query 5
Immigrants       Border Wall      Mexico Border    Build Wall       Finish Wall
```

Data
Collection

EMR Job 1
Data Cleaning in Map and Reduce
- Individually Query and Combined Query Files

EMR Job 2
Reduced File Combined into One
- Individually Query and Combined Query Files

Map
Reduce
Step

Part Files Received Converted to CSV

CSV plotted in Tableau

Data
Visual-
ization

| Combined Query | Immigrants | Border Wall | Mexico Border | Build Wall | Finish Wall |

Tableau Graphs Represented on Webpage

Job Flow for Common Crawl Data

Data Collected Using Common Crawler WET files

**Data Collection**

| Query 1 Immigrants Data Cleaned | Query 2 Border Wall Data Cleaned | Query 3 Mexico Border Data Cleaned | Query 4 Build Wall Data Cleaned |

Data from NYT Query 1 Immigrants

Data from NYT Query 2 Border Wall

Data from NYT Query 3 Mexico Border

Data from NYT Query 4 Build Wall

**Map Reduce Step**

EMR Job 1
Data Cleaning in Map and Reduce
- Individually Query and Combined Query Files

EMR Job 2
Reduced File Combined into One
- Individually Query and Combined Query Files

Part Files Received Converted to CSV

**Data Visualization**

CSV plotted in Tableau

| Combined Query | Immigrants | Border Wall | Mexico Border | Build Wall |

Tableau Graphs Represented on Webpage

Job Flow for NYT Data

Data Collected Using nytimes articlesearch API

Query 1
Immigrants
Data Cleaned

Query 2
Border Wall
Data Cleaned

Query 3
Mexico Border
Data Cleaned

Query 4
Build Wall
Data Cleaned

Data Collection

Data from NYT
Query 1
Immigrants

Data from NYT
Query 2
Border Wall

Data from NYT
Query 3
Mexico Border

Data from NYT
Query 4
Build Wall

EMR Job 1
Data Cleaning in Map and Reduce
- Individually Query and Combined Query Files

EMR Job 2
Reduced File Combined into One
- Individually Query and Combined Query Files

Map
Reduce
Step

Part Files Received Converted to CSV

CSV plotted in Tableau

Data
Visual-
ization

Combined Query

Immigrants

Border Wall

Mexico Border

Build Wall

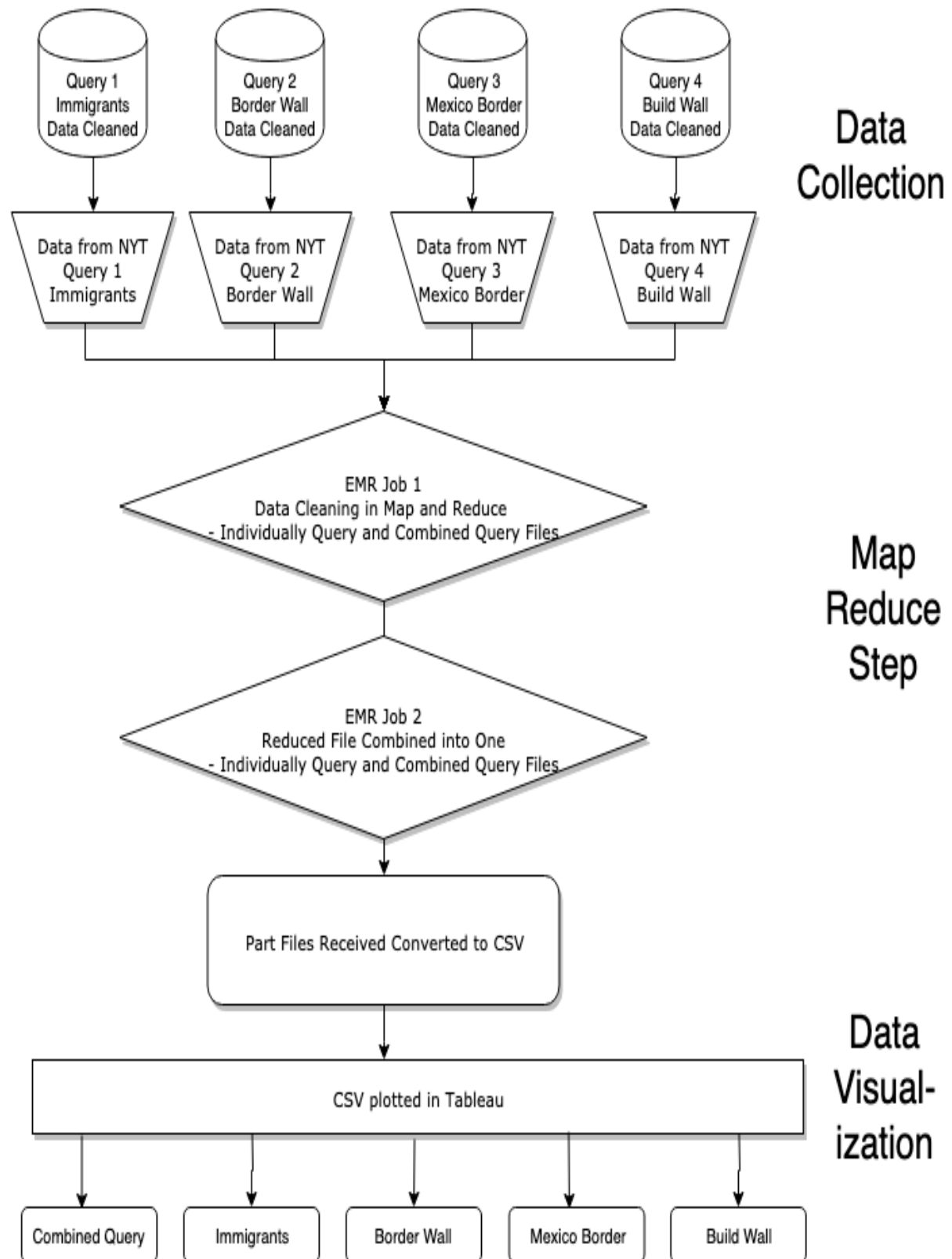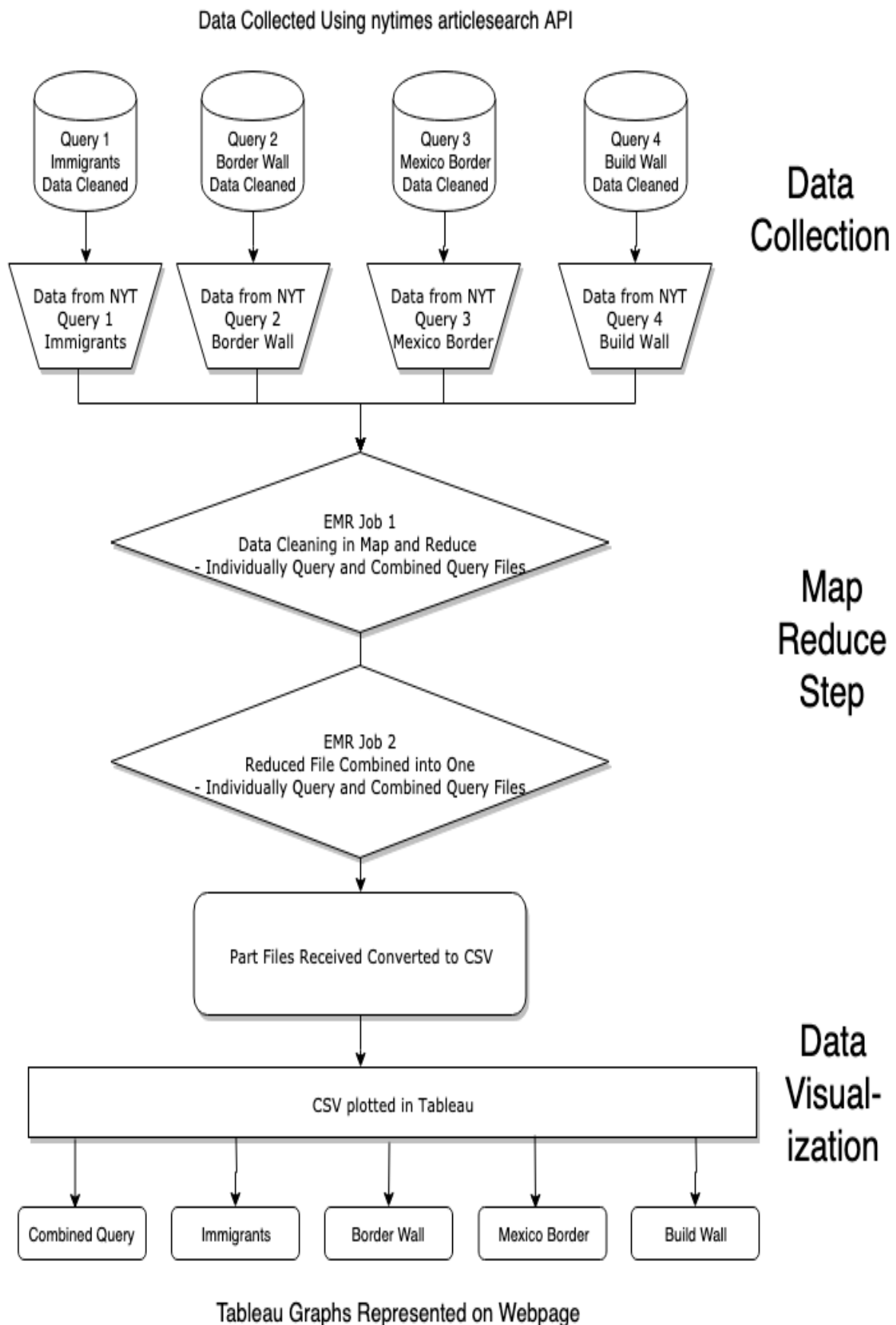Tableau Graphs Represented on Webpage

Query list used – **Border Wall, Build Wall, Immigrant, Finish Wall (only Twitter as no data found), Mexico Border**

**Data Collection**
1. Twitter

Data collected on 29 March and 9 April 2019. Count of tweets-

Border Wall - 27449

Build Wall – 32198 + 26556

Finish Wall – 664 + 3583

Mexico Border - 28937

Immigrant - 33345

Total - 152732

2. NYT

Data collected on 29 March 2019 using ArticleSearch Api and Beautiful Soup. Article Count-

Border Wall - 577

Build Wall - 170

Mexico Border - 207

Immigrant – 943

Total – around 1800

3. Common Crawl

Data collected on 15 April 2019. Article Count-

Border Wall - 2514

Build Wall - 220

Mexico Border – 84

Immigrants – 728314

**Data Processing**
1. Twitter

Removed Retweets. Count of tweets-

Border Wall – 5632

Build Wall – 5452 + 7480 = 12932

Finish Wall – 325 + 1360 = 1685

Mexico Border - 1721

Immigrant – 7246

Total = 29216

2. NYT

Sentence Tokenized for Co-occurrence using nltk
   nltk.sent_tokenize(article)

3. Common Crawl

Sentence Tokenized for Co-occurrence using nltk
nltk.sent_tokenize(article)

**Mapper and Reducer**

Mapper.py – Word Count

Mapper2.py – Word Co- occurrence – NYT and Common Crawl

Mapper3.py – Word Co- occurrence – Twitter

MapperCC.py py – Word Count – Common Crawl

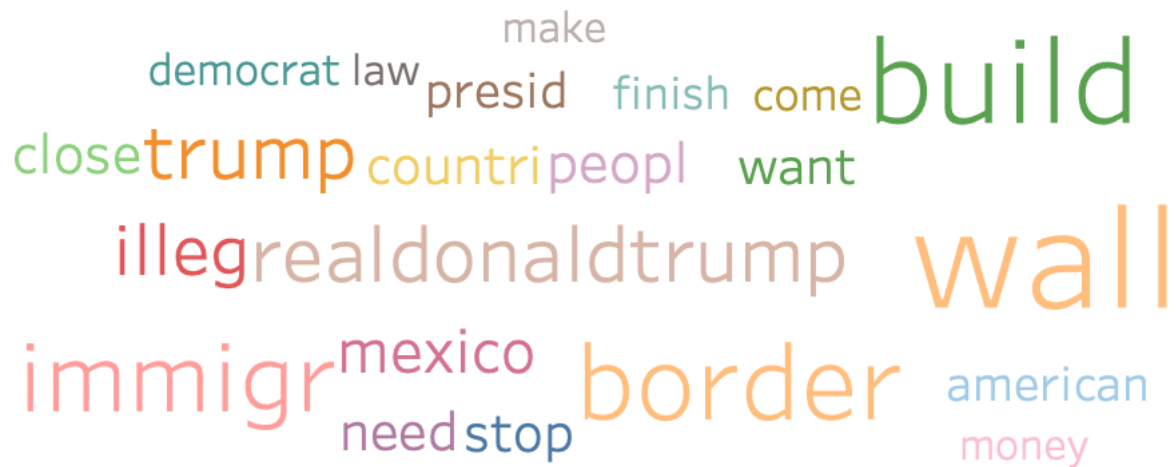Minor changes done to mapper to improve pre-processing.

**Data Visualization**

In Common Crawl, only Border Wall and Immigrant data is used for co-occurrence and full word count as rest return irrelevant words

## Twitter

**Word Count**

Word Cloud on Full Data



Word Cloud on Border Wall

Word Cloud on Build Wall

law america close
countri
want american mexico pay border
presid money damn
trump illeg build wall
Word: build need
Count: 13,580 come
immigr
peopl realdonaldtrump stop
keep

Word Cloud on Finish Wall

come end money
down build countri great peopl border
presid
illeg
trump need finish stop wall
america back
close realdonaldtrump immigr
start
pleas

Word Cloud on Mexico Border

come asylum judg mexico
threaten american stop back
southern
close illeg border mexican trump aid
down
need migrant immigr
countri presid realdonaldtrump danger
peopl wall shut cross

Word Cloud on Immigrant

come
america trump white realdonaldtrump
want
stop make
border immigr illeg work
democrat
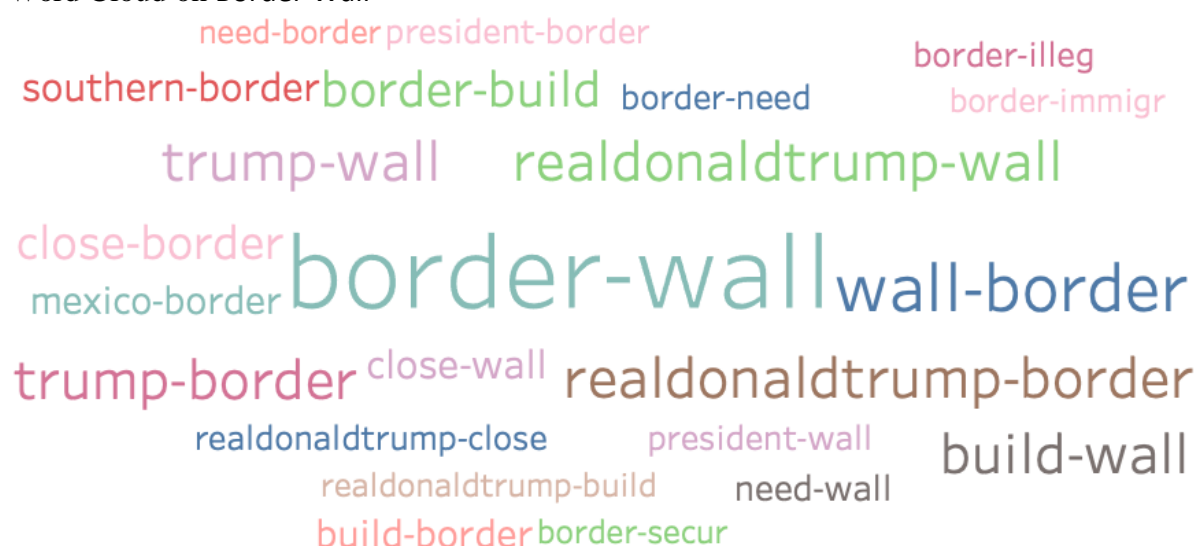american legal mexico know countri
need law peopl

**Word Co-occurrence**
Word Cloud on Full Data

finish-wall wall-build
mexico-border border-mexico wall-illeg
illegal-immigr trump-build close-wall wall-border
realdonaldtrump-immigr realdonaldtrump-build
trump-wall build-wall border-wall
border-build
realdonaldtrump-border realdonaldtrump-wall
close-border trump-border southern-border border-immigr
build-border president-wall trump-immigr

Word Cloud on Border Wall

need-border president-border border-illeg
southern-border border-build border-need border-immigr
trump-wall realdonaldtrump-wall
close-border border-wall wall-border
mexico-border
trump-border close-wall realdonaldtrump-border
realdonaldtrump-close president-wall build-wall
realdonaldtrump-build need-wall
build-border border-secur

Word Cloud on Build Wall

money-build
wall-build president-build wall-illeg border-build
trump-build realdonaldtrump-build
close-border people-wall
trump-wall build-wall border-wall
build-illeg people-build
realdonaldtrump-border wall-trump realdonaldtrump-wall
build-border president-wall close-wall close-build
wall-border money-wall wall-stop

Word Cloud on Finish Wall

end-releas end-catch
close-border trump-finish build-wall wall-amp
amnesty-end realdonaldtrump-wall
realdonaldtrump-border
finish-end finish-wall wall-finish
end-lawbreak
border-finish realdonaldtrump-finish
border-wall wall-end end-anchor trump-wall
end-reward end-babi

Word Cloud on Mexico Border

trump-close
president-mexico president-trump stop-border
close-mexico mexico-stop realdonaldtrump-close
trump-border trump-mexico border-mexico
shut-border mexico-illeg
mexico-border mexico-immigr
close-border realdonaldtrump-mexico
southern-border realdonaldtrump-border
president-border mexico-close border-stop mexico-trump
border-illeg border-wall

Word Cloud on Immigrant

country-immigr

legal-immigr  stop-immigr

people-immigr  realdonaldtrump-illeg  immigrants-illeg

immigrants-peopl  undocumented-immigr  immigrants-want

trump-immigr  border-illeg  realdonaldtrump-immigr

trump-border  illegal-immigr  immigration-immigr

immigrants-american  immigrants-immigr

immigrants-border  immigrants-countri

border-immigr  immigrants-amp  amp-immigr

immigrants-come  immigrants-legal

**New York Times**

**Word Count**

Word Cloud on Full Data

first white  polit

countri  peopl hous  last  border

american  mani

republican  trump  work democrat

state  immigr nation  presid  govern

wall  senat two day  make

Word: **trump**
Count: **12,075**

Word Cloud on Border Wall

two

hous  countri

nation  wall senat  work  border

shutdown  make

republican  trump peopl  presid

congress  american

state  govern polit democrat

last  immigr secur  fund  emerg

Word Cloud on Build Wall

build think support congress
two wall nation fund democrat hous
republican trump thi presid
govern american immigr senat border
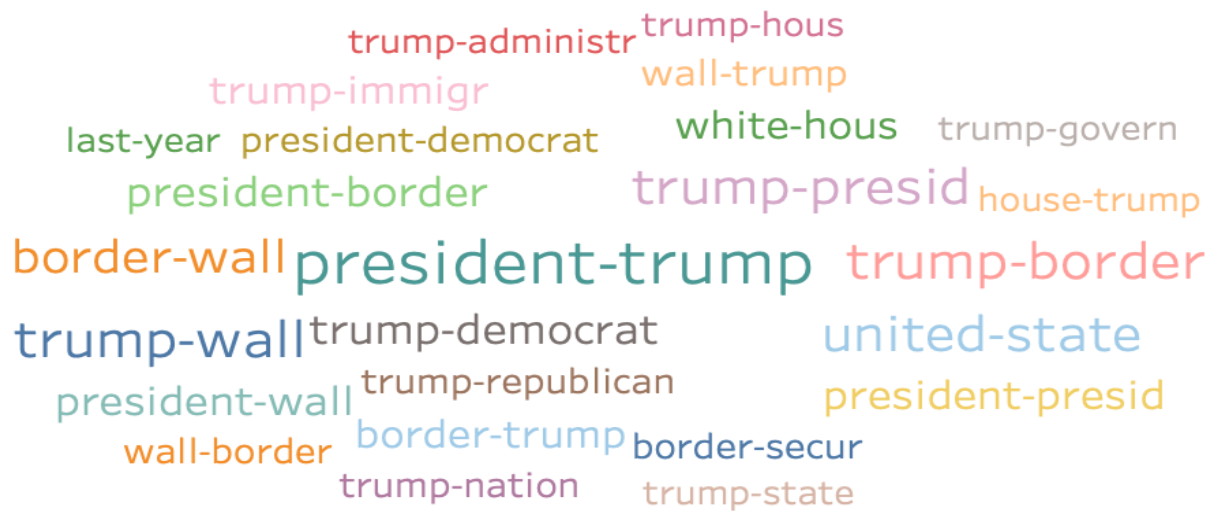state peopl emerg countri polit year

Word Cloud on Mexico Border

peopl countri wall last emerg presid
migrant mexico year hous
republican border nation trump
american govern democrat immigr
two state mexican work offici senat
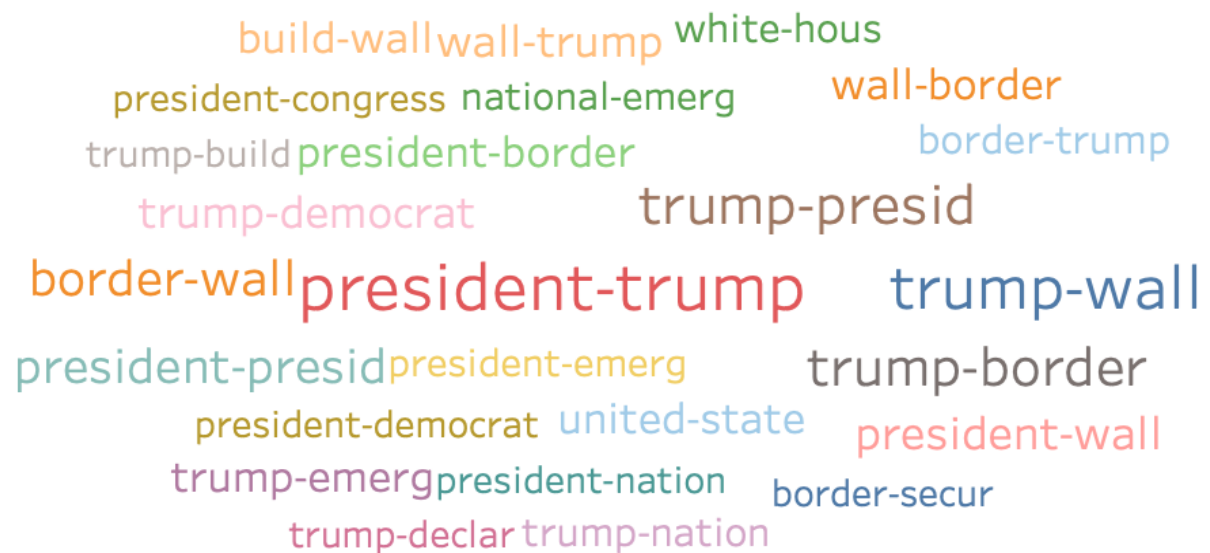secur

Word Cloud on Immigrant

nation work first unit wall hous
citi year peopl govern last american
presid polit trump mani democrat
republican state immigr border
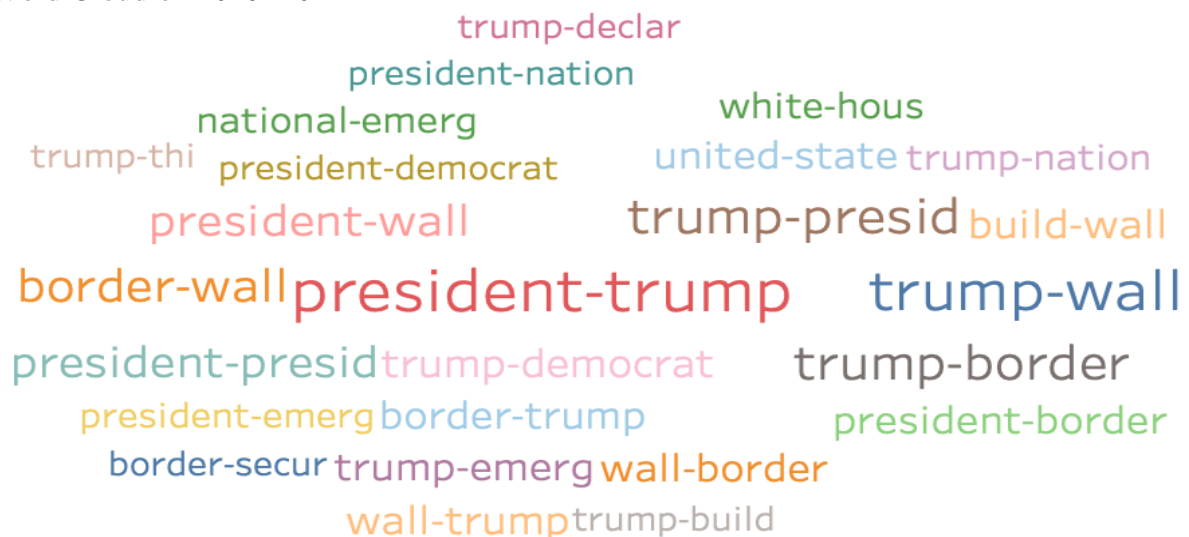countri
york make two

**Word Co-occurrence**

Word Cloud on Full Data



Word Cloud on Border Wall



Word Cloud on Build Wall

Word Cloud on Mexico Border

border-secur
wall-border
trump-administr states-border border-state
border-trump border-migrant president-wall trump-immigr
mexico-border border-wall united-state
trump-wall president-trump trump-border
trump-presid president-presid president-border
border-mexico mexican-border migrants-border
border-cross trump-emerg trump-nation
trump-state

Word Cloud on Immigrant

border-secur
president-border border-trump
anti-semit president-democrat students-school border-immigr
trump-immigr democrats-trump trump-presid
border-wall president-trump united-state
trump-border high-school trump-democrat trump-govern
trump-republican house-democrat trump-wall
trump-administr president-wall trump-state
last-year white-hous

Word: **border-wall**
Count: **564**

**Common Crawl**

**Word Count**

Word Cloud on Full Data

white
law national history news visa
world business information trump
department immigration government
international march state war american
president states united make best social
wall

Word Cloud on Border Wall

tax watch
greg national american best
good wall economic donald decor news state
president trump america march
feb border diego room government economy
business san white world war down house

Word Cloud on Build Wall

von
den shower news bildung das
mit
des wall modern sekundarbereich fragen
der white bersicht table und
metal
jahrgang wahlmodul nach medien materialien die
window hir Word: jahrgang light storage
Count: 1,235
wood design

Word Cloud on Mexico Border

download answers college levitra
research paper series workshop
online
edition life viagra help essay
write
writing
essays study thesis cialis
owners
ebook manuals buy user

Word Cloud on Immigrant

john   war
national united march public make
history business information world
american immigration government
department states visa international
law state country social best trump
news

**Word Co-occurrence**
Unable to run even with 20 instances on AWS.