

# ENRON EMAIL CORPUS

A K-MEANS CLUSTER ANALYSIS



# ENRON

THE ORIGINAL EVIL CORP

# BACK STORY

- HOUSTON BASED COMPANY FOUNDED IN 1985
- BEFORE BANKRUPTCY ON DEC 3, 2001
  - HAD 29,000 EMPLOYEES
  - MAJOR ELECTRIC, NATURAL GAS, COMMUNICATIONS & PULP & PAPER CO.
  - REVENUE OF 101 BILLION IN 2000
  - NAMED “AMERICA’S MOST INNOVATIVE Co.” 6 CONSECUTIVE YEARS - FORTUNE.

# EVIL CORP.

- THEY COOKED THEIR BOOKS
- WHICH CAUSED THE ENRON SCANDAL
- SARBANES-OXLEY ACT OF 2002
- KILLED ARTHUR ANDERSON
- INTENTIONALLY ENCOURAGED REMOVAL OF POWER DURING CA ENERGY CRISIS
- BUT MOST IMPORTANTLY...



# THE PEOPLE

- 144 SENIOR MANAGERS
- 29,000 STAFF
- THOUSANDS LOST THEIR 401K'S AND PENSIONS
- YOUR LIFE'S WORK





File Edit View Insert Cell Kernel Widgets Help

Not Trusted

Python 3



In [2]: 1 df = pd.read\_csv('./data/enron\_email\_20\_percent.csv')

In [3]: 1 df.shape

Out[3]: (103480, 52)

## Cleaning Data

In [193]: 1 # show 55 columns and rows to see my column data  
2  
3 pd.set\_option('display.max\_columns', 55)  
4 pd.set\_option('display.max\_rows', 55)

In [5]: 1 df.head(1)

Out[5]:

|   | Unnamed: | Message-ID | Date |
|---|----------|------------|------|
| 0 | 0        |            |      |

THE PROJECT

|   |        |   |            |                                      |   |
|---|--------|---|------------|--------------------------------------|---|
| 0 | 403470 | <117238.1075846810566.JavaMail.evans@thyme> | 2000-12-11 | frozenset({'susan.scott@enron.com'}) | frozenset({'alicia.perkins@enron.com'}) |
|   |        |   | 09:22:00   |                                      |   |

# ENRON EMAIL ANALYSIS

- DATA - ENRON EMAILS
- CULPABILITY VS PROBABLE CAUSE VS SAVE TIME ???
- RESEARCHED & ACQUIRED OUR DATA
- CLEANED THE DATA
- EXPLORED THE DATA
- MODELED DATA THROUGH K-MEANS CLUSTER
- PRESENT DATA



# Enron Email Dataset

DATASET BY BRIAN RAY

Comment

631

Explore this dataset

[Overview](#) [Contributors](#) [Discussion](#) [Activity](#)

## Overview

### DESCRIPTION

Enron Email Dataset converted to tabular format: From, To, Subject, and Content. Some records labeled by CMU students.

### SUMMARY

The Enron email dataset contains approximately 500,000 emails generated by employees of the Enron Corporation. It was obtained by the Federal Energy Regulatory Commission during its investigation of Enron's collapse.

This is the May 7, 2015 Version of dataset, as published at <https://www.cs.cmu.edu/~enron/>

Those 1,700 records that were labeled by students are marked with `labeled` set were labelled by CMU students.

There are up to 12 categories per email:

- Cat [1-12] level\_1 = top-level category
- Cat [1-12] level\_2 = second-level category

Show more ▾

## About this dataset

SHARED WITH Everyone

CREATED 2 years ago by  
 @brianray

SIZE 2.91 GB · Download

TAGS enron, emails,  
unstructured text, nlp

LICENSE Other

DICTIONARY 14 files, 377 columns ·  
View

## THE DATA

### Queries (1)

[enron\\_05\\_17\\_2015\\_with\\_labels\\_v2](#)

### Related projects (5)



# ENRON EMAIL CORPUS

- ORIGINALLY 1.7 MM EMAILS
  - MOST HAVE BEEN REMOVED BECAUSE OF PRIVACY ISSUES
- 500,000 REMAIN PUBLICLY AVAILABLE
- SEVERAL KAGGLE PROJECTS ON JUST CLEANING THE DATA
- I USED THE DATA SET BY BRIAN RAY FOUND AT DATA.WORLD
- REMOVED UNNECESSARY NAN'S, DUPLICATES, ALSO MISC. COLUMNS

File Edit View Insert Cell Kernel Widgets Help

Not Trusted

Python 3



## Countvectorize

In [156]:

```
1 stopwords = ENGLISH_STOP_WORDS
2 my_stopwords = ENGLISH_STOP_WORDS.union(['ect', 'hou', 'com', 'recipient', 'sent', 'enron',
3                                         'forwarded', 'corp', 'said', 'attached', 'ect', \
4                                         'recipient', 'email', 'original', 'doc', 'pm', 'ma
5                                         'enronxgate', 'na', 'year'])
6
7
8 # tf = TfidfVectorizer(analyzer='word', stop_words=my_stopwords, min_df=2, max_df=0.5)
9 cv = CountVectorizer(analyzer='word', stop_words=my_stopwords, max_features=2000)
10
11 # vectorizer = TfidfVectorizer(stop_words='english')
12 X = cv.fit_transform(df['content']) # this is my sparse matrix
13
14
```

In [157]:

```
1 # Additional EDA on number only data that I don't believe is mission critical
2 pd.set_option('display.max_columns', 500)
```

In [158]:

```
1 # New DataFrame created to remove vectorized numbers
2 X_df = pd.DataFrame(X.toarray(), columns=cv.get_feature_names()) # this is my dense matrix
```

In [159]:

```
1 # Looking at shape of new dataframe
```

# DATA MODELING

# ENRON EMAIL ANALYSIS - MODELING

- ❶ TERM FREQUENCY-INVERSE DOCUMENT FREQUENCY WAS INITIALLY USED
- ❷ I ENDED UP USING COUNTVECTORIZOR
  - ❸ STOPWORDS  
MYSTOPWORDS
  - ❹ COUNTVECTORIZER  
MAX\_FEATURES = 2000

# ENRON EMAIL ANALYSIS - MODELING

- K-MEANS CLUSTER ALGORITHM
- CLUSTERED GROUPS OF WORDS MORE CLOSELY RELATED
  - $K = 5$  CLUSTERS

# ENRON EMAIL ANALYSIS - INITIAL RESULTS

## WHOLE CORPUS FREQUENTLY OCCURRING WORDS / NUMBER OF OCCURRENCES

|          |       |            |       |         |       |          |       |            |       |
|----------|-------|------------|-------|---------|-------|----------|-------|------------|-------|
| power    | 61521 | ees        | 35715 | just    | 25315 | make     | 20947 | services   | 18356 |
| energy   | 59937 | california | 31829 | houston | 24148 | question | 20779 | think      | 18186 |
| new      | 55393 | need       | 31753 | john    | 24042 | use      | 20611 | work       | 18158 |
| time     | 43476 | state      | 30168 | mark    | 23802 | today    | 20337 | million    | 17956 |
| compan   | 42883 | business   | 29407 | agreem  | 23129 | jeff     | 20170 | report     | 17948 |
| gas      | 42459 | like       | 28488 | meeting | 22331 | td       | 19901 | monday     | 17805 |
| thanks   | 40003 | day        | 28418 | image   | 22274 | contact  | 19495 | service    | 17636 |
| know     | 39739 | let        | 27730 | deal    | 21771 | don      | 19424 | credit     | 17589 |
| market   | 37091 | week       | 27554 | group   | 21493 | date     | 19286 | electricit | 17371 |
| informat | 36272 | price      | 26530 | trading | 21250 | font     | 19209 | friday     | 17235 |

# ENRON EMAIL ANALYSIS - MODELING

- K-MEANS CLUSTER ALGORITHM
- CLUSTERED GROUPS OF WORDS MORE CLOSELY RELATED
  - $K = 5$  CLUSTERS

# ENRON EMAIL ANALYSIS - CLUSTERS

## CLUSTER 1

Businesses  
option  
Defense  
manager  
linda  
evening  
known  
assignment  
attend  
analysis  
season  
crude  
hot  
sellers  
business  
efforts  
debbie  
students  
charge  
dale

## CLUSTER 2

sources  
engine  
stories  
role  
applicationes  
30  
body  
trades  
quickly  
75  
doug  
south  
follow  
site  
annual  
computer  
person  
light  
break  
alternative

## CLUSTER 3

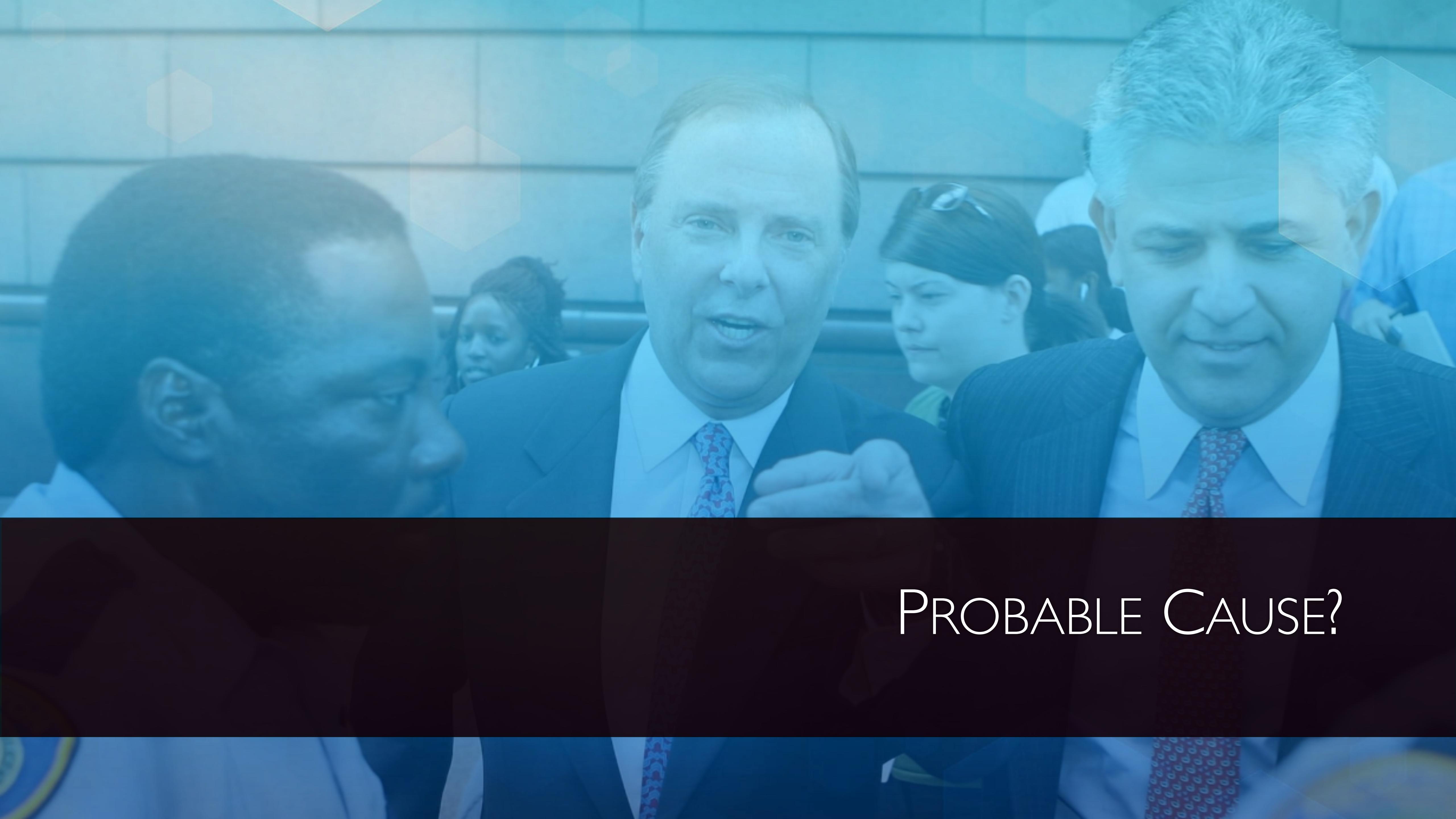
manager  
stake  
station  
including  
defense  
option  
day  
gov  
evening  
main  
known  
investigation  
involved  
commodity  
businesses  
tools  
holding  
assignment  
hub  
knowledge

## CLUSTER 4

meter  
bpa  
impacted  
quantity  
heat  
software  
assistance  
mkt\_type  
ancillary  
appear  
option  
members  
shapiro  
kenneth  
industry  
sense  
holding  
read  
judy  
dear

## CLUSTER 5

option  
season  
defense  
attend  
debbie  
commodities  
evening  
team  
pager  
page  
notice  
analysis  
dave  
red  
manager  
known  
clean  
teams  
brown  
office

A photograph showing several men in business attire, including suits and ties, standing in what appears to be a hallway or lobby. One man in the center foreground is gesturing with his right hand while speaking. The background is slightly blurred, showing other people and architectural details.

PROBABLE CAUSE?

# ENRON EMAIL ANALYSIS - CLUSTER 3

|           |               |
|-----------|---------------|
| manager   | known         |
| stake     | investigation |
| station   | involved      |
| including | commodity     |
| defense   | businesses    |
| option    | tools         |
| day       | holding       |
| gov       | assignment    |
| evening   | hub           |
| main      | knowledge     |

A group of diverse business professionals in a modern office setting, smiling and engaged in conversation.

# PRACTICAL APPLICATIONS

# PRACTICAL APPLICATIONS

- 1.7 MILLION EMAILS
- AVG EMAIL READ RATE - 20 SECONDS
- 34,000,000 SECONDS (THAT'S 34 MILLION SECONDS)
- 393 - 24 HOUR DAYS TO READ THROUGH 1.7 MM EMAILS
- 1,180 - 8 HOUR WORK DAYS (3.23 YEARS)

# PRACTICAL APPLICATIONS

- SORT THROUGH A CLUSTER OF EMAILS
- LESS TIME CONSUMING
- POTENTIALLY MORE PRODUCTIVE
- FOR EXAMPLE...

# CONSULTATION GIG

- ◆ PULL OUT MY TRUSTY TITANIUM PowerBook G4
  - ◆ CIRCA 2001
  - ◆ 128 MB RAM
  - ◆ 128 MB GRAPHICS CARD
  - ◆ 10 GIGS OF DISK SPACE



A group of diverse business professionals in a modern office setting, smiling and engaged in conversation.

# PRACTICAL APPLICATIONS

# ENRON EMAIL ANALYSIS - CLUSTER 3

|           |               |
|-----------|---------------|
| manager   | known         |
| stake     | investigation |
| station   | involved      |
| including | commodity     |
| defense   | businesses    |
| option    | tools         |
| day       | holding       |
| gov       | assignment    |
| evening   | hub           |
| main      | knowledge     |

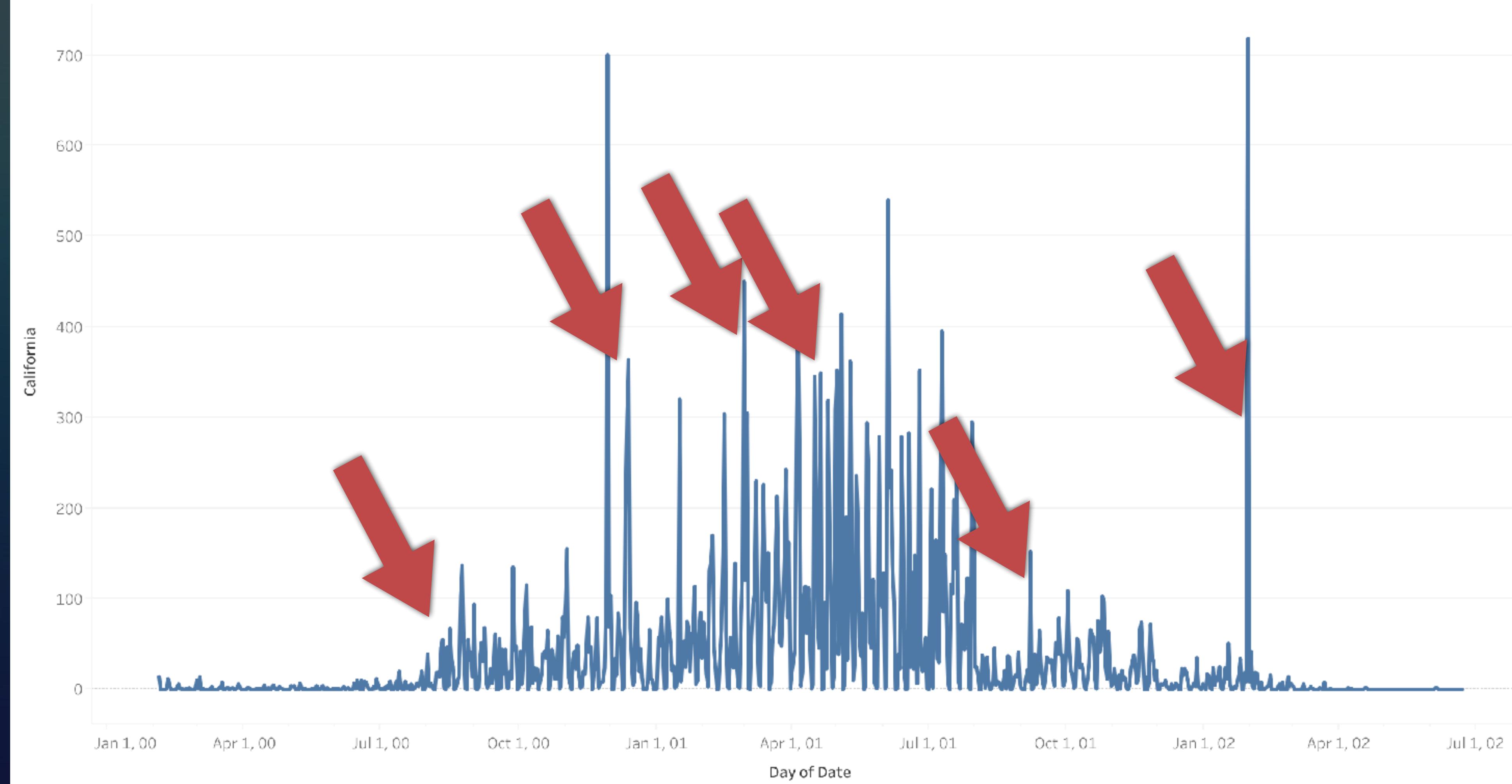


Q & A

A photograph of a group of people in a professional environment. In the foreground, a man with light-colored hair and a blue suit is looking directly at the camera with a slight smile. Behind him, another man with grey hair and a dark suit is also smiling. To the left, a person's head is partially visible, and in the background, there are other individuals, some wearing glasses. The scene is overlaid with several large, semi-transparent white geometric shapes, including triangles and hexagons.

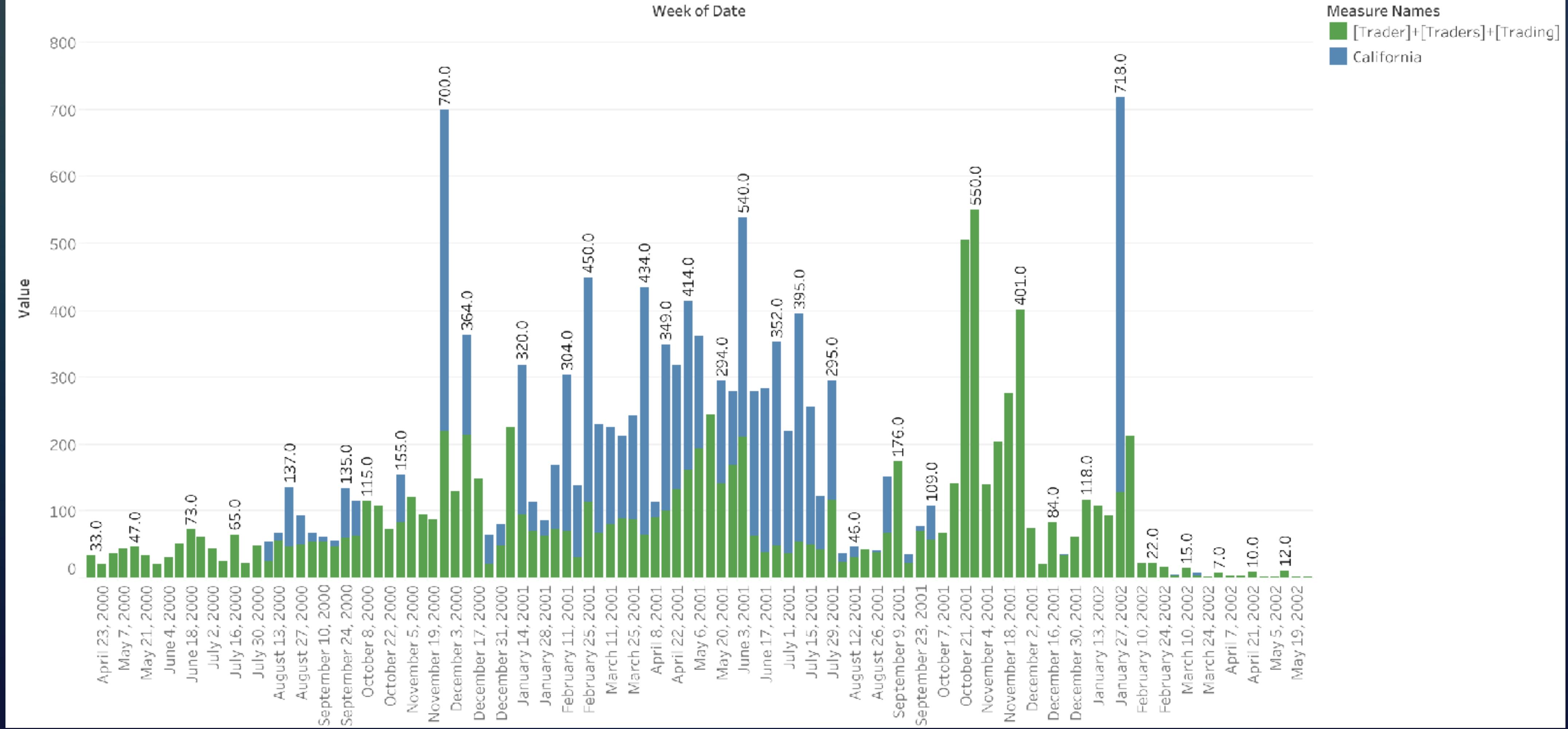
# INTERESTING FINDS

California Bars



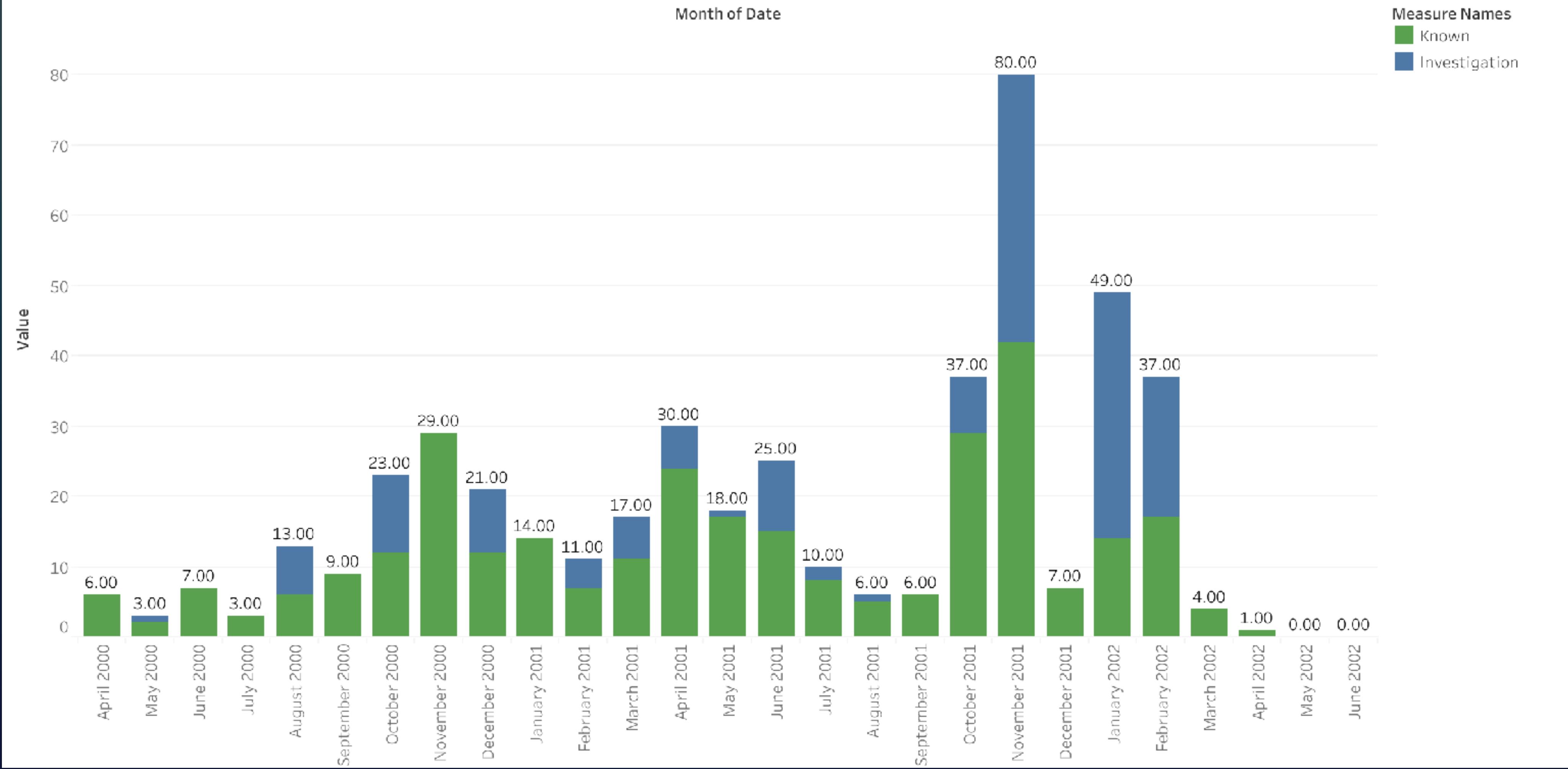
CALIFORNIA WORD COUNT

## California vs Trading



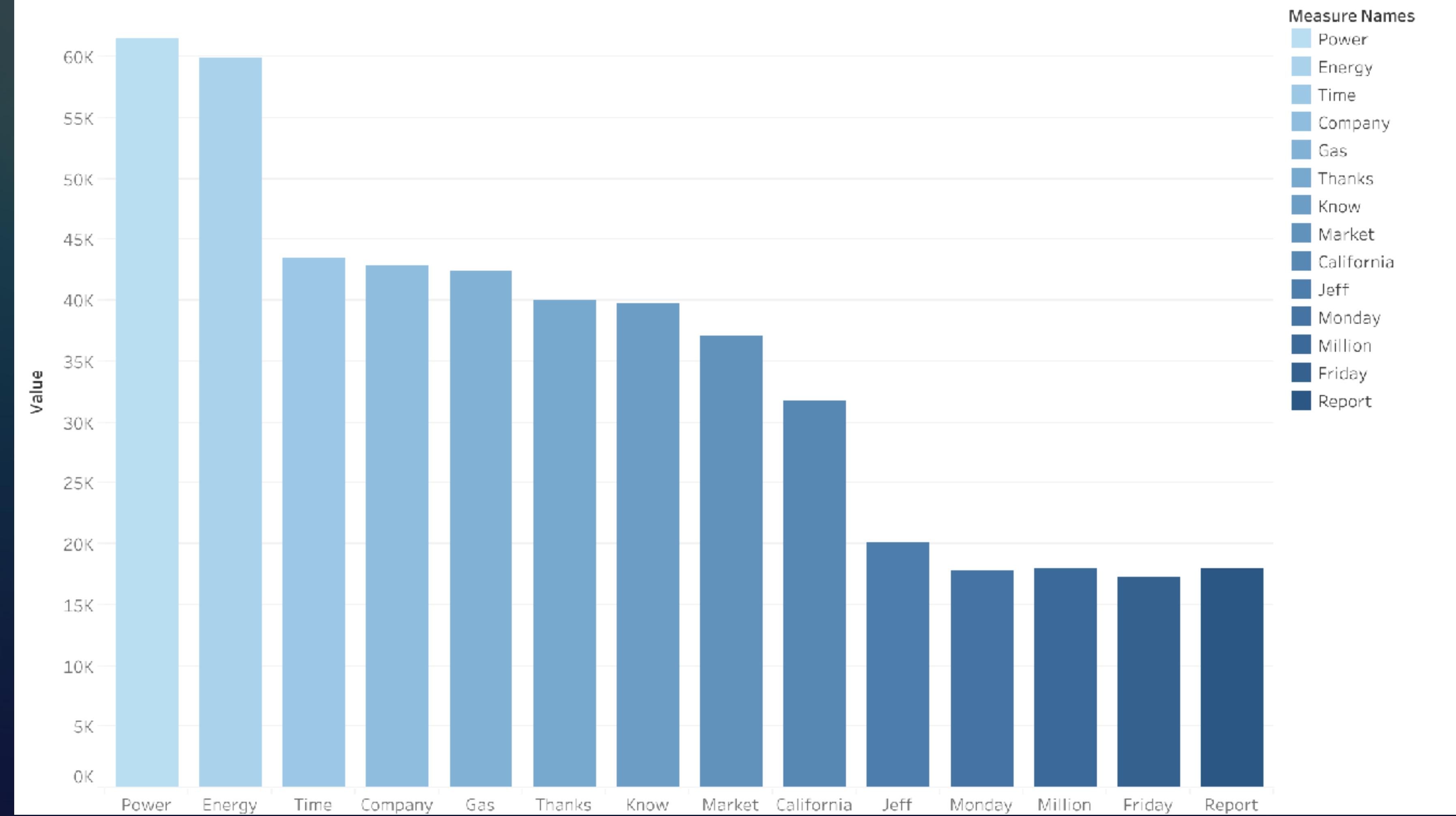
# CALIFORNIA VS TRADING

## Known vs Investigation



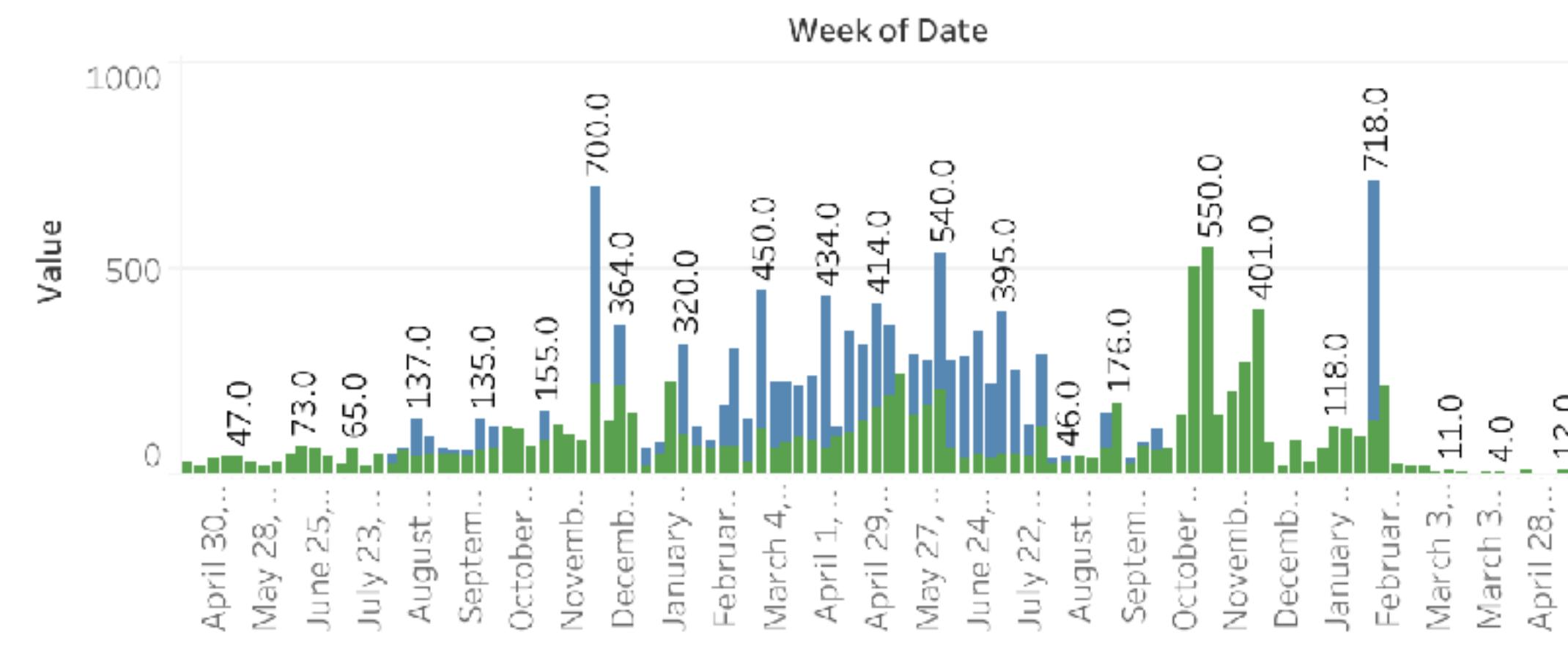
# KNOWN VS INVESTIGATION

## Top Corpus Word Counts

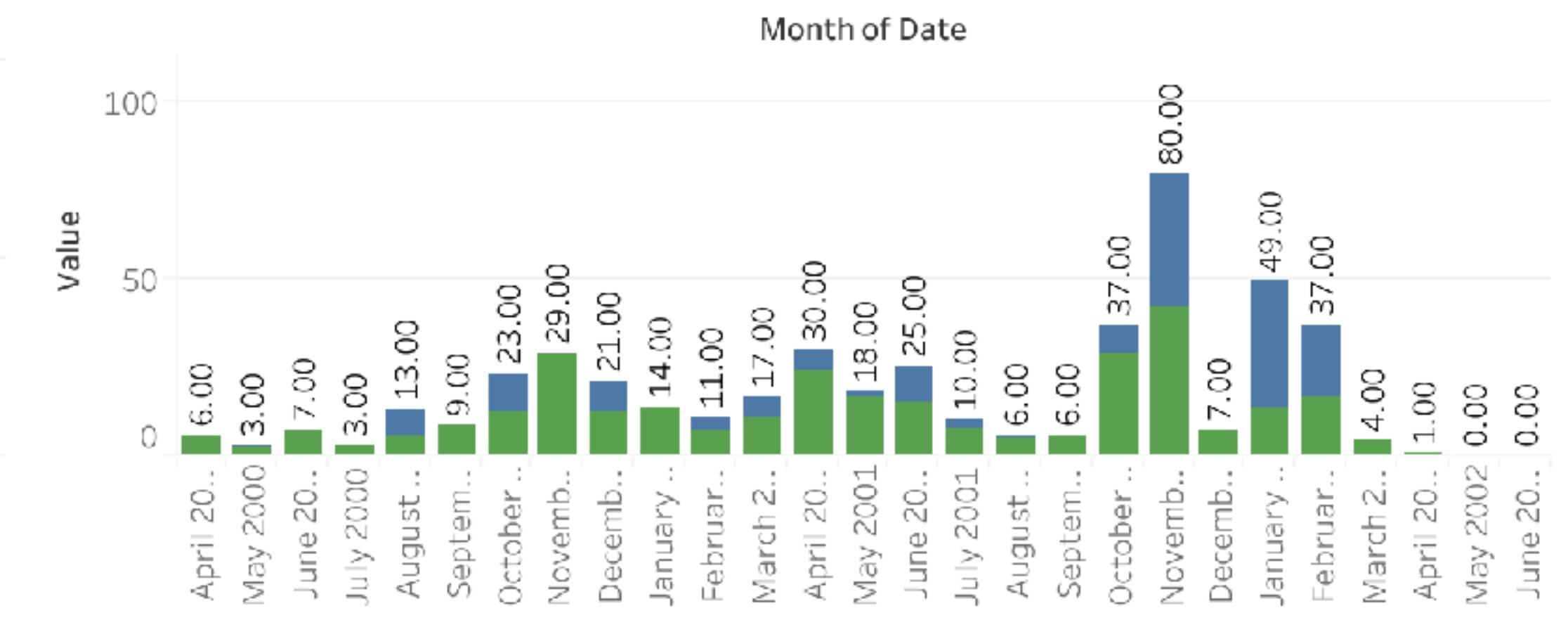


# TOP WORDS FOR WHOLE CORPUS

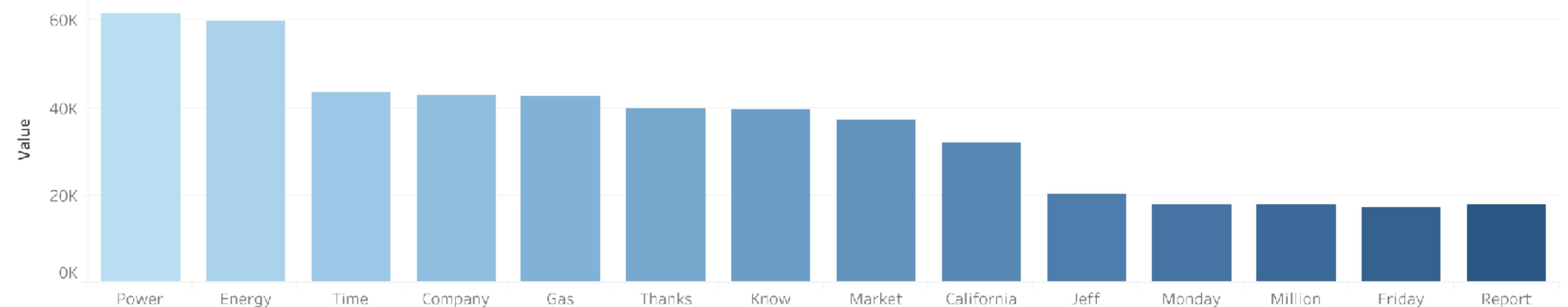
## California vs Trading



## Known vs Investigation



## Top Corpus Word Counts





Q & A

# ORIGINALLY TERMS & CONDITIONS

- HARD TO FIND
- HARD TO SCRAPE
- NOT EASILY ACCESSIBLE
- BEGIN SIGNING UP FOR EVERYTHING
- No, NOT REALLY
- BUT TO IDENTIFY WHAT PRIVACY WE MAY OR MAY NOT BE GIVING UP

# ENRON

- ◆ RIGHTS GIVEN UP TERMS AND CONDITIONS LANDED ON THE ENRON DATA SET FROM 2003
- ◆ LEARNED THAT SIRI WAS TRAINED ON IT...
- ◆ YOU'RE PROBABLY THINKING, OH, THAT'S WHY SIRI IS SO BAD
- ◆ ONLY PUBLICLY AVAILABLE CORPUS OF THAT TYPE AVAILABLE TO THE PUBLIC
  - ◆ WHAT DOES THAT MEAN.... HOW MANY OF YOU HAVE GMAIL...
    - IN JULY 2017 GOOGLE PROMISED NOT TO SEARCH YOUR EMAILS - DON'T WORRY THEY WOULD ONLY SEARCH FOR WHAT TO SELL YOU. THEY DON'T DO IT ANY MORE
    - THEY NOW SEARCH YOUR PICTURES THAT YOU UPLOAD TO GOOGLE PHOTOS FOR FREE FOR POTENTIAL MARKETING OPPORTUNITIES. SO AGAIN, YOU'RE GOOD

# ENRON

- IF YOU'RE NOT PAYING FOR THE PRODUCT, YOU ARE THE PRODUCT