# The speech we miss:
# How keyword-based data collection obscures youth participation in online political discourse

**Sarah Shugars**
they/them/theirs
Assistant Professor, Rutgers University

RUTGERS

# Detecting and correcting bias in linked data sources

**The speech we miss:** How Keyword-Based Data Collection Obscures Youth Participation in Online Political Discourse

Adina Gitomer, Sarah Shugars, Ryan J. Gallagher, Stefan McCabe, Brooke Foucault Welles

*Computational Communication Research*, 5(1). 2023.
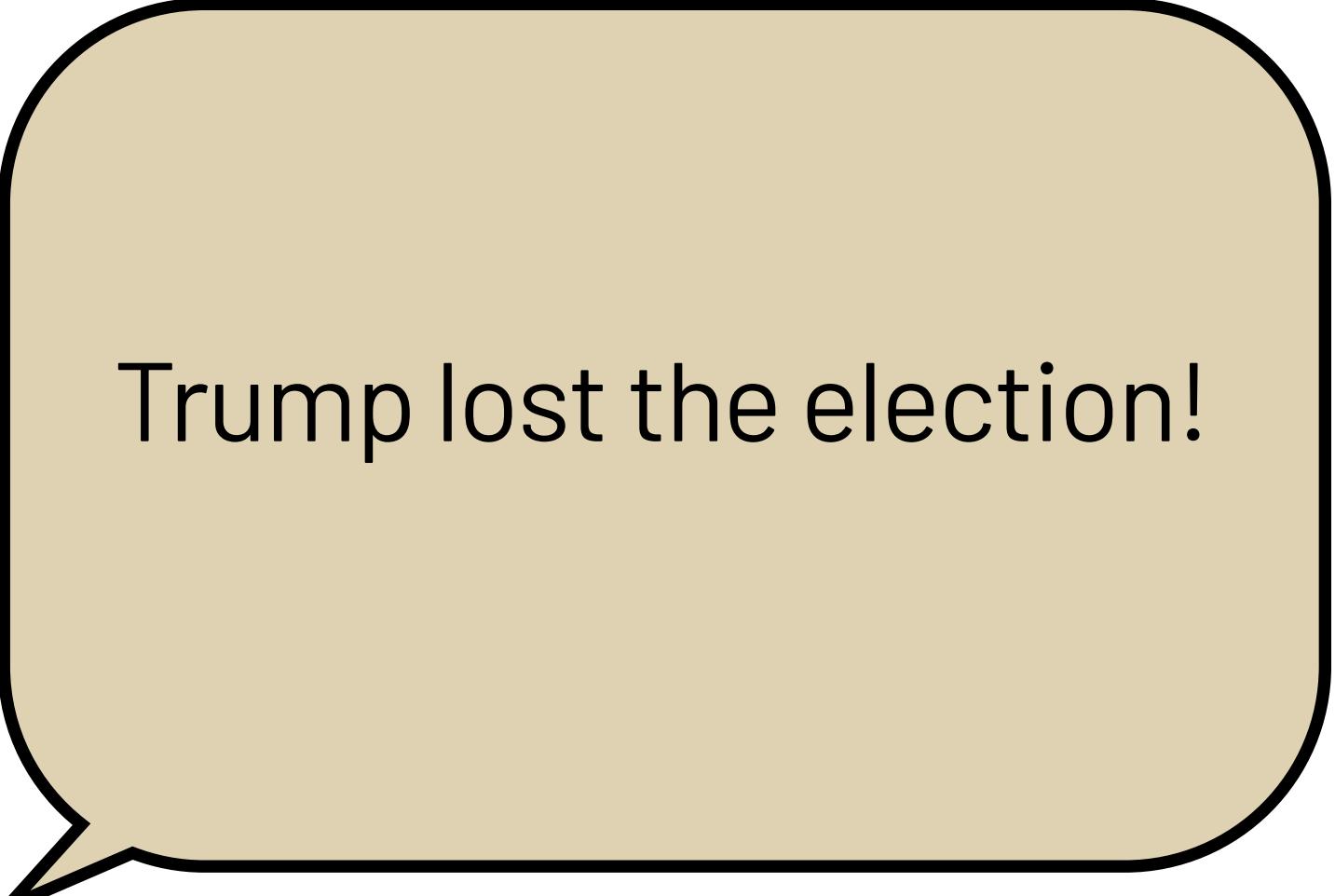
# RIP Twitter API 🪦

```
url = 'https://api.twitter.com

response = requests.request(
    "GET", url,
    params= {
    query = 'Trump OR election OR fascism'
    })
```

List of "political"
keywords

# Capturing "Political" Speech

Keyword search implicitly assumes
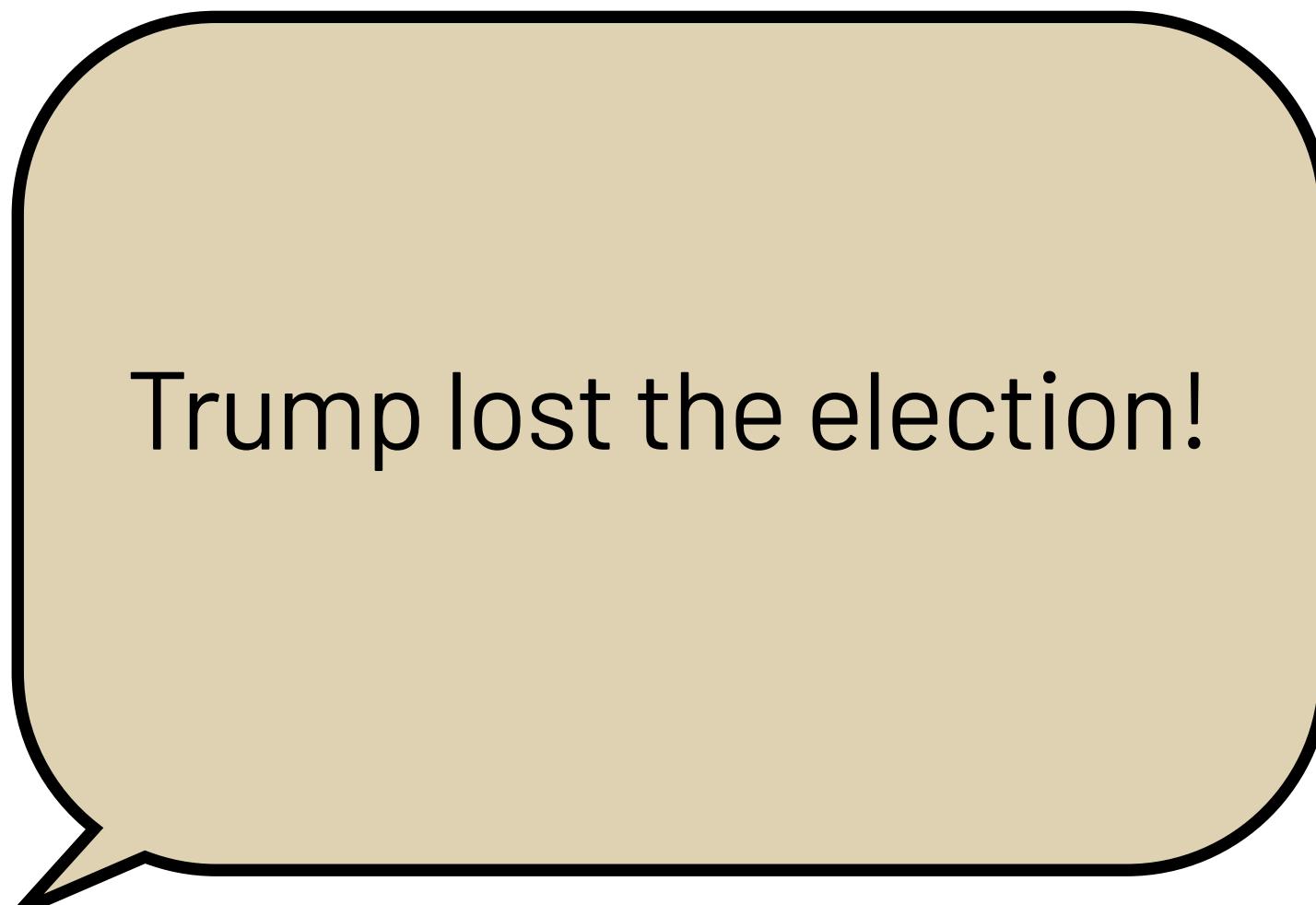political content looks like:

Trump lost the election!

# Capturing "Political" Speech

Keyword search implicitly assumes
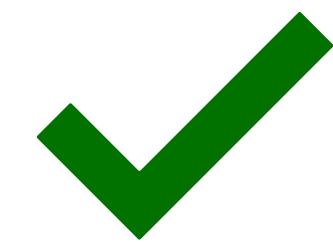political content looks like:

What political speech
actually looks like



Trump lost the election!



orange was ejected

# Capturing "Political" Speech

# The Speech We Miss



Keyword classifier

**Trump** lost the **election**!

✅ Yes, election-related

❌ Not election-related

**RQ1:**
How big of a problem is this?

# The Speech We Miss

Keyword classifier

**Trump** lost the **election**!

✔️ Yes, election-related

❌ Not election-related

orange was ejected

**RQ1:**
How big of a problem is this?

**RQ2:**
Is there variation by **age**?

# Data Matching

Hughes et al. (2020)
Shugars et al. (2021)

# Data Matching

✦ Panel of 1.6 million Twitter users matched to US
voting records

Hughes et al. (2020)
Shugars et al. (2021)

# Data Matching

✦ Panel of 1.6 million Twitter users matched to US voting records

✦ Matching procedure:

Hughes et al. (2020)
Shugars et al. (2021)

# Data Matching

✦ Panel of 1.6 million Twitter users matched to US voting records

✦ Matching procedure:

➡ Collect Twitter profiles from 10% sample (2014-2017)

Hughes et al. (2020)
Shugars  et al. (2021)

# Data Matching

✦ Panel of 1.6 million Twitter users matched to US voting records

✦ Matching procedure:

➡ Collect Twitter profiles from 10% sample (2014-2017)

➡ Match users to voting records:

Hughes et al. (2020)
Shugars et al. (2021)

# Data Matching

✦ Panel of 1.6 million Twitter users matched to US voting records

✦ Matching procedure:

➡ Collect Twitter profiles from 10% sample (2014-2017)

➡ Match users to voting records:

  ✦ Name + (city, state) must be unique match

Hughes et al. (2020)
Shugars  et al. (2021)

# Data Matching

- ✦ Panel of 1.6 million Twitter users matched to US voting records

- ✦ Matching procedure:
  - ➡ Collect Twitter profiles from 10% sample (2014-2017)
  - ➡ Match users to voting records:
    - ✦ Name + (city, state) must be unique match



**Sarah Shugars**
@Shugars

Assistant Professor @RutgersCommInfo. CSS & PolCom. Previously: @NYUDataScience, @NUnetsi. They/them. #FirstGen

Bluesky: shugars.bsky.social

New Brunswick, NJ    sarahshugars.com    Joined February 2009

1,140 Following    3,720 Followers

Match if I'm the **only** "Sarah Shugars" registered to vote in "New Brunswick, NJ"

Hughes et al. (2020)
Shugars  et al. (2021)

# Data Matching

- ✦ Panel of 1.6 million Twitter users matched to US voting records

- ✦ Matching procedure:
  - ➡ Collect Twitter profiles from 10% sample (2014-2017)
  - ➡ Match users to voting records:
    - ✦ Name + (city, state) must be unique match

- ✦ Demographically representative of Twitter users overall

**Sarah Shugars**
@Shugars

Assistant Professor @RutgersCommInfo. CSS & PolCom. Previously: @NYUDataScience, @NUnetsi. They/them. #FirstGen

Bluesky: shugars.bsky.social

📍 New Brunswick, NJ    🔗 sarahshugars.com    🗓 Joined February 2009

**1,140** Following    **3,720** Followers

Match if I'm the **only** "Sarah Shugars" registered to vote in "New Brunswick, NJ"

Hughes et al. (2020)
Shugars et al. (2021)

# Data Matching

- Panel of 1.6 million Twitter users matched to US voting records

- Matching procedure:

  ➡ Collect Twitter profiles from 10% sample (2014–2017)

  ➡ Match users to voting records:

    - Name + (city, state) must be unique match

- Demographically representative of Twitter users overall

  - Youngest users were 17 in 2017

**Sarah Shugars**
@Shugars

Assistant Professor @RutgersCommInfo. CSS & PolCom. Previously: @NYUDataScience, @NUnetsi. They/them. #FirstGen
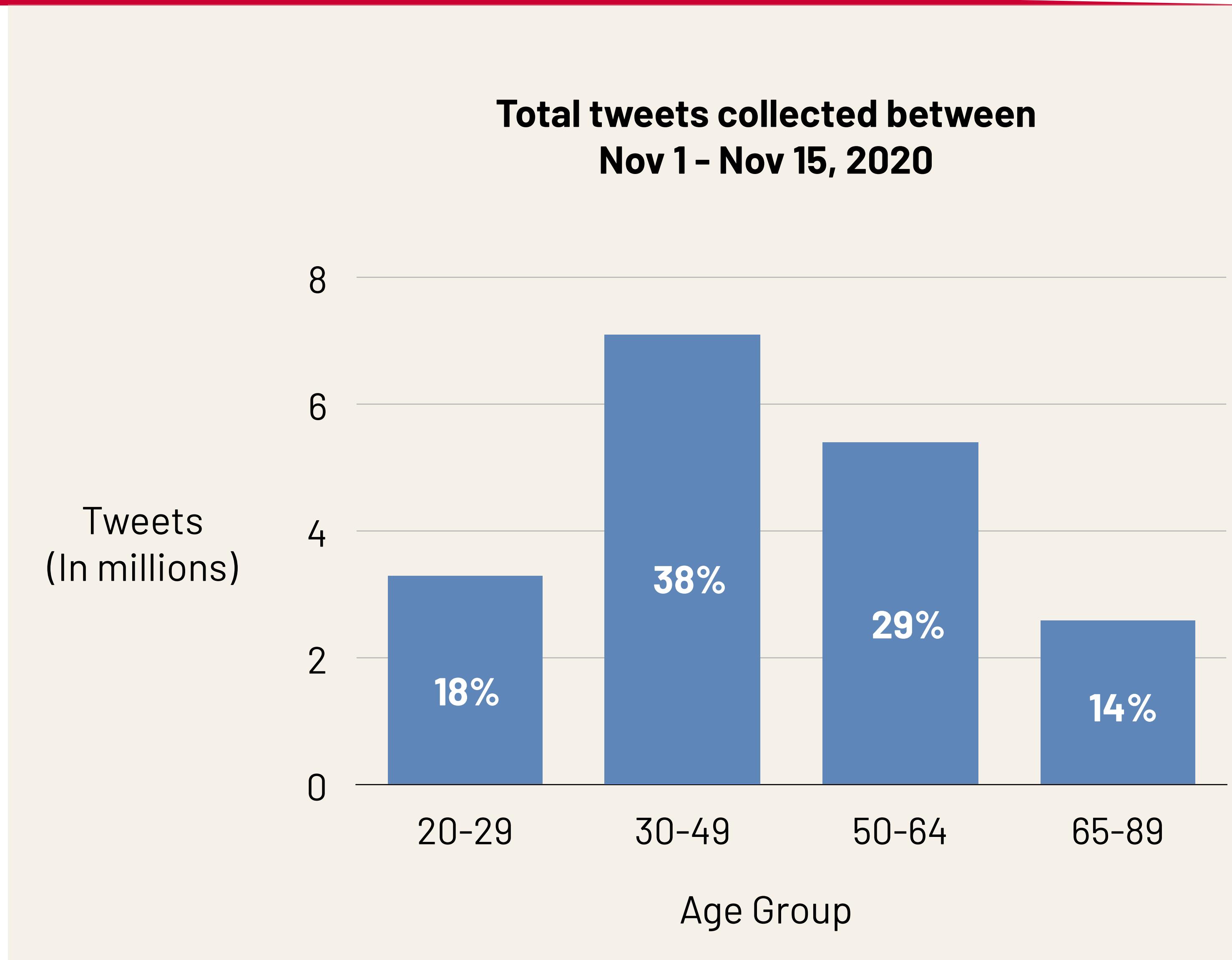
Bluesky: shugars.bsky.social

New Brunswick, NJ   sarahshugars.com   Joined February 2009

1,140 Following   3,720 Followers

Edit profile

Match if I'm the **only** "Sarah Shugars" registered to vote in "New Brunswick, NJ"
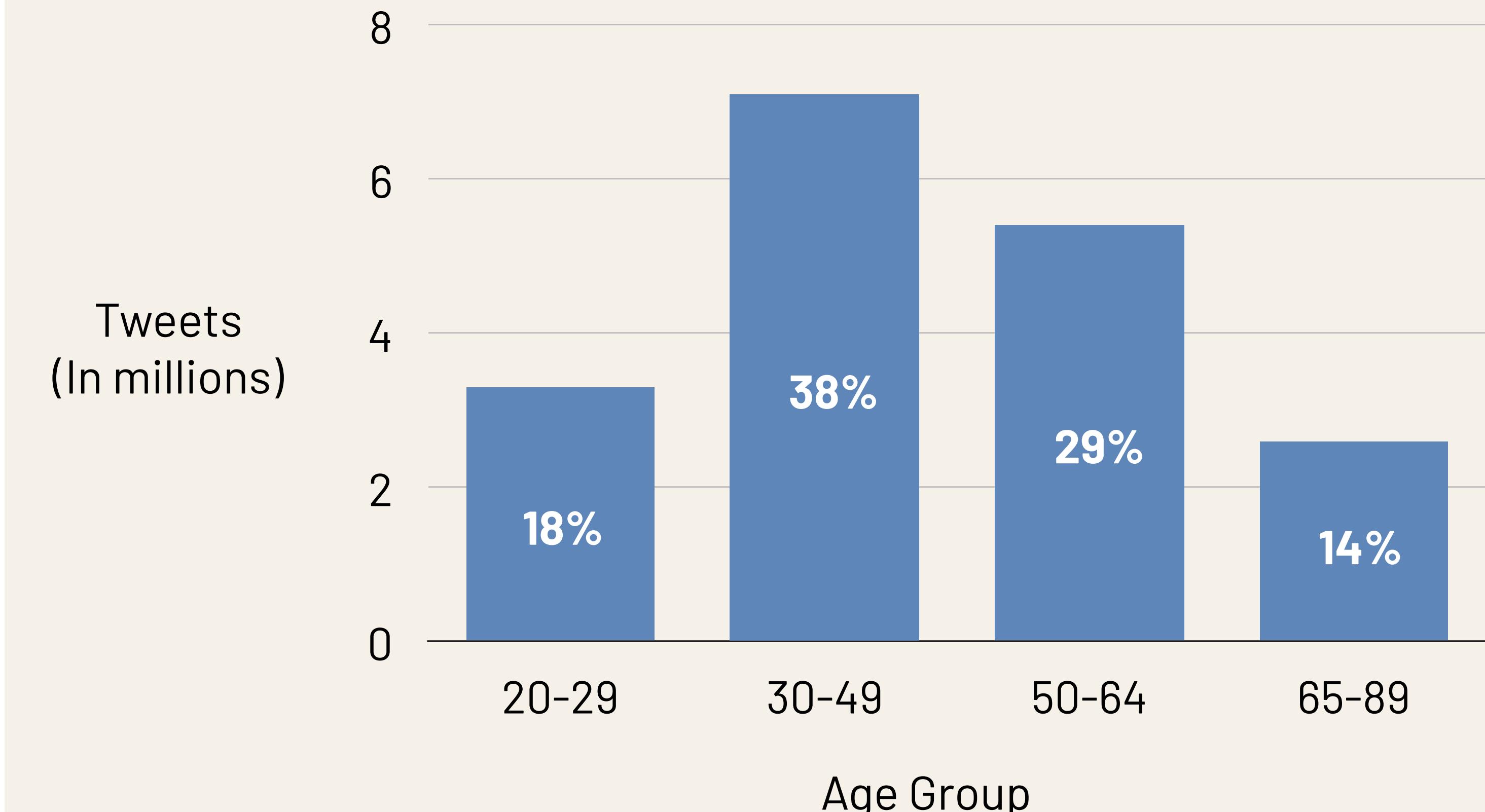
Hughes et al. (2020)
Shugars  et al. (2021)

# Data



Total tweets collected between
Nov 1 – Nov 15, 2020

Tweets
(In millions)

18%

38%

29%

14%

20-29    30-49    50-64    65-89

Age Group

# Data

We collect <u>all</u> posts made by panelists between
Nov 1 - Nov 15 2020

**Total tweets collected between
Nov 1 - Nov 15, 2020**



Tweets
(In millions)

8

6

4

2

0

18%

38%

29%

14%

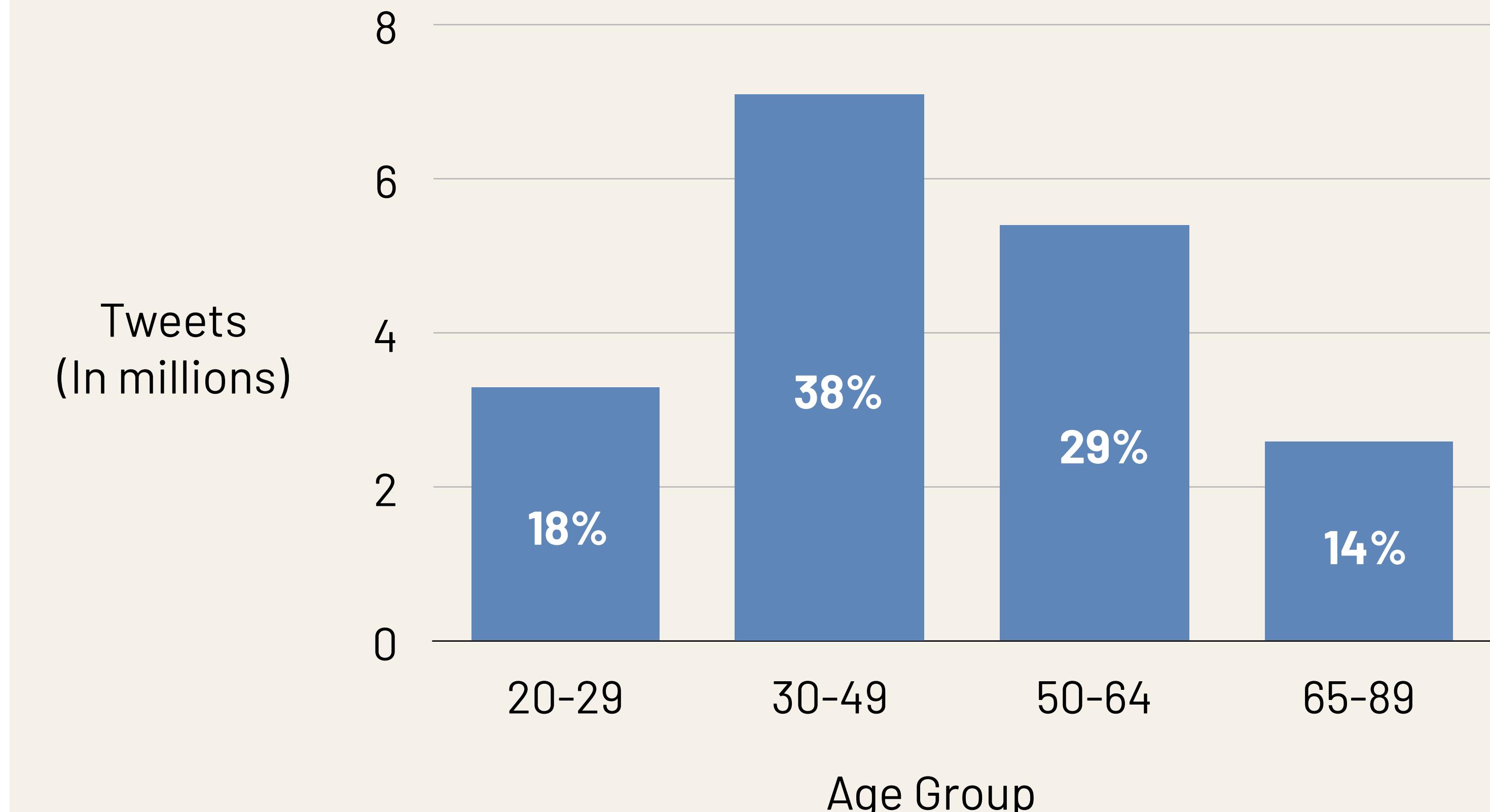20-29    30-49    50-64    65-89

Age Group

# Data

We collect <u>all</u> posts made by panelists between
Nov 1 - Nov 15 2020

➡ Which tweets would be retrieved using keyword-based search?

**Total tweets collected between
Nov 1 - Nov 15, 2020**



Tweets
(In millions)

Age Group

# Data

We collect <u>all</u> posts made by panelists between
Nov 1 - Nov 15 2020

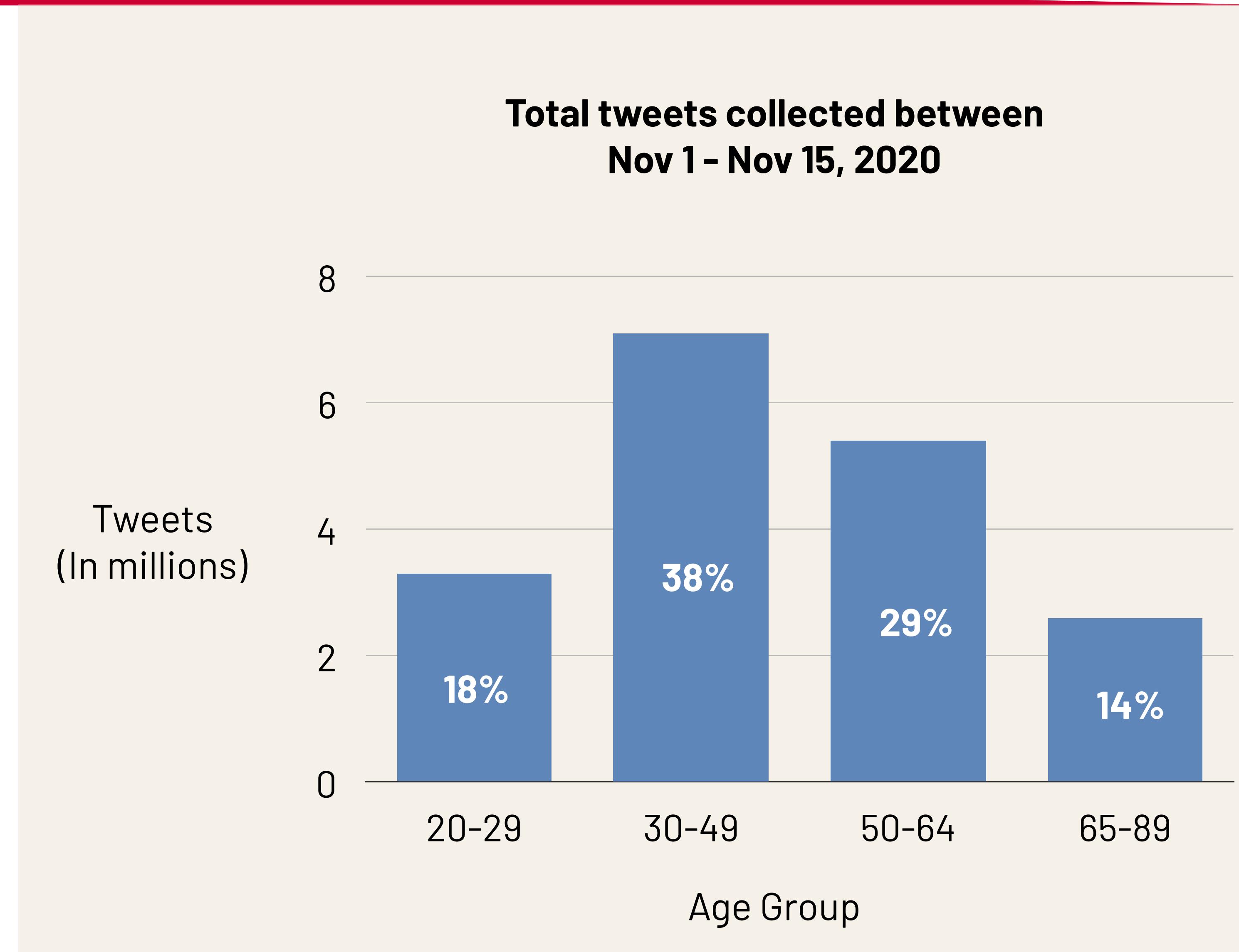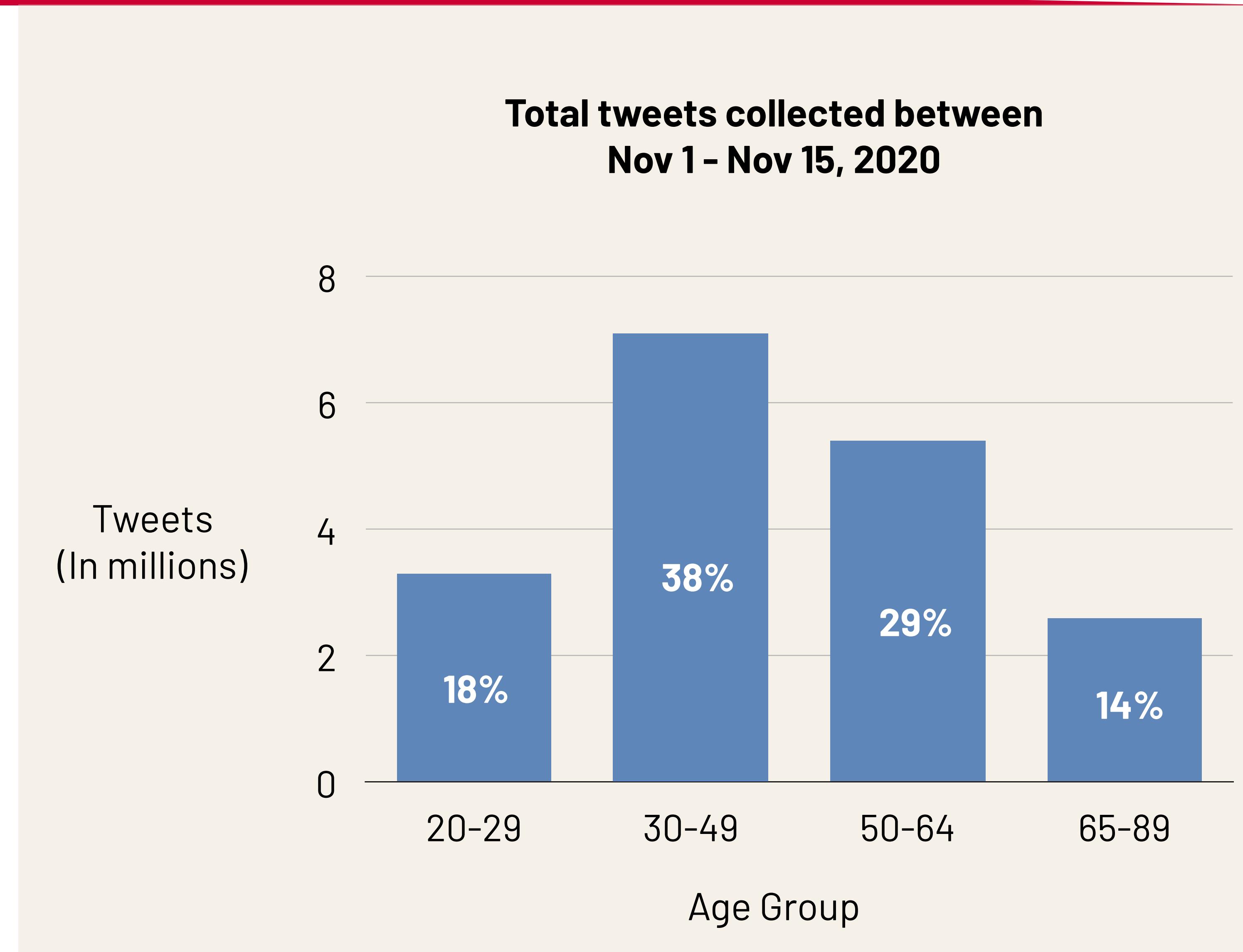➡ Which tweets would be retrieved using keyword-based search?

➡ Which tweets are actually election-related? (Via handcoding)

**Total tweets collected between
Nov 1 - Nov 15, 2020**



Tweets
(In millions)

| Age Group | Percentage |
| --- | --- |
| 20-29 | 18% |
| 30-49 | 38% |
| 50-64 | 29% |
| 65-89 | 14% |

Age Group

# Data

We collect <u>all</u> posts made by panelists between
Nov 1 - Nov 15 2020

➡ Which tweets would be retrieved using keyword-based search?

➡ Which tweets are actually election-related? (Via handcoding)

➡ Stratify sample by age

**Total tweets collected between
Nov 1 - Nov 15, 2020**

Tweets
(In millions)



Age Group

# Method

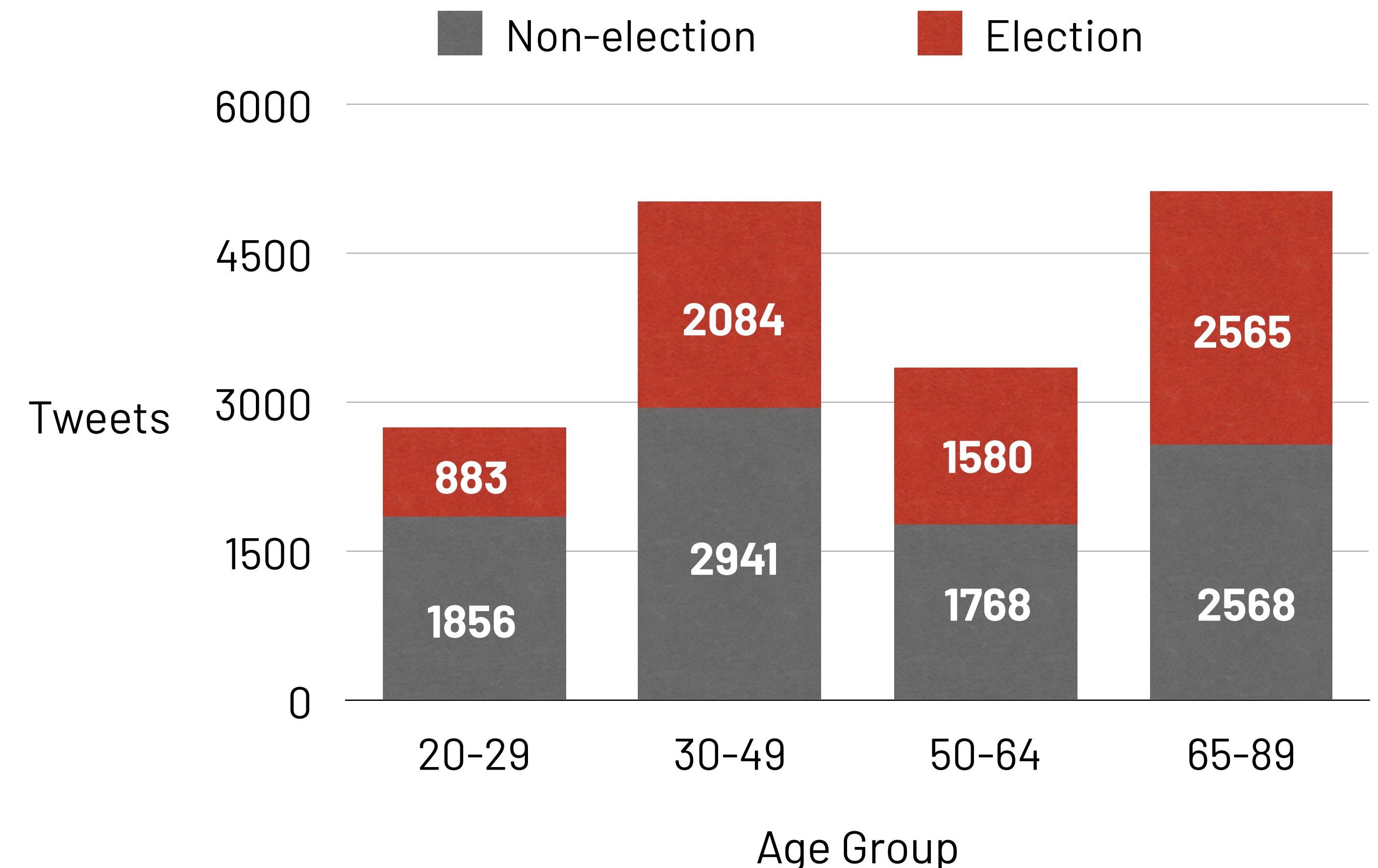**Step 1:** Use an initial keyword classifier to take an informed sample for hand coding

**Step 2:** Hand code 16,245 tweets for election relevance

# Method

**Step 1:** Use an initial keyword classifier to take an informed sample for hand coding

**Step 2:** Hand code 16,245 tweets for election relevance

**Tweets sampled for handcoding and keyword-based classification**

# Keyword Classifier Accuracy

How frequently did the keyword classifier
disagree with our handcoding?

■ Classified as
election-related

■ Classified as not
election-related

20-29          30-49          50-64          65-89

# Keyword Classifier Accuracy

How frequently did the keyword classifier
disagree with our handcoding?

■ Classified as
election-related

2084

1580

2565

883

■ Classified as not
election-related

20-29  30-49  50-64  65-89

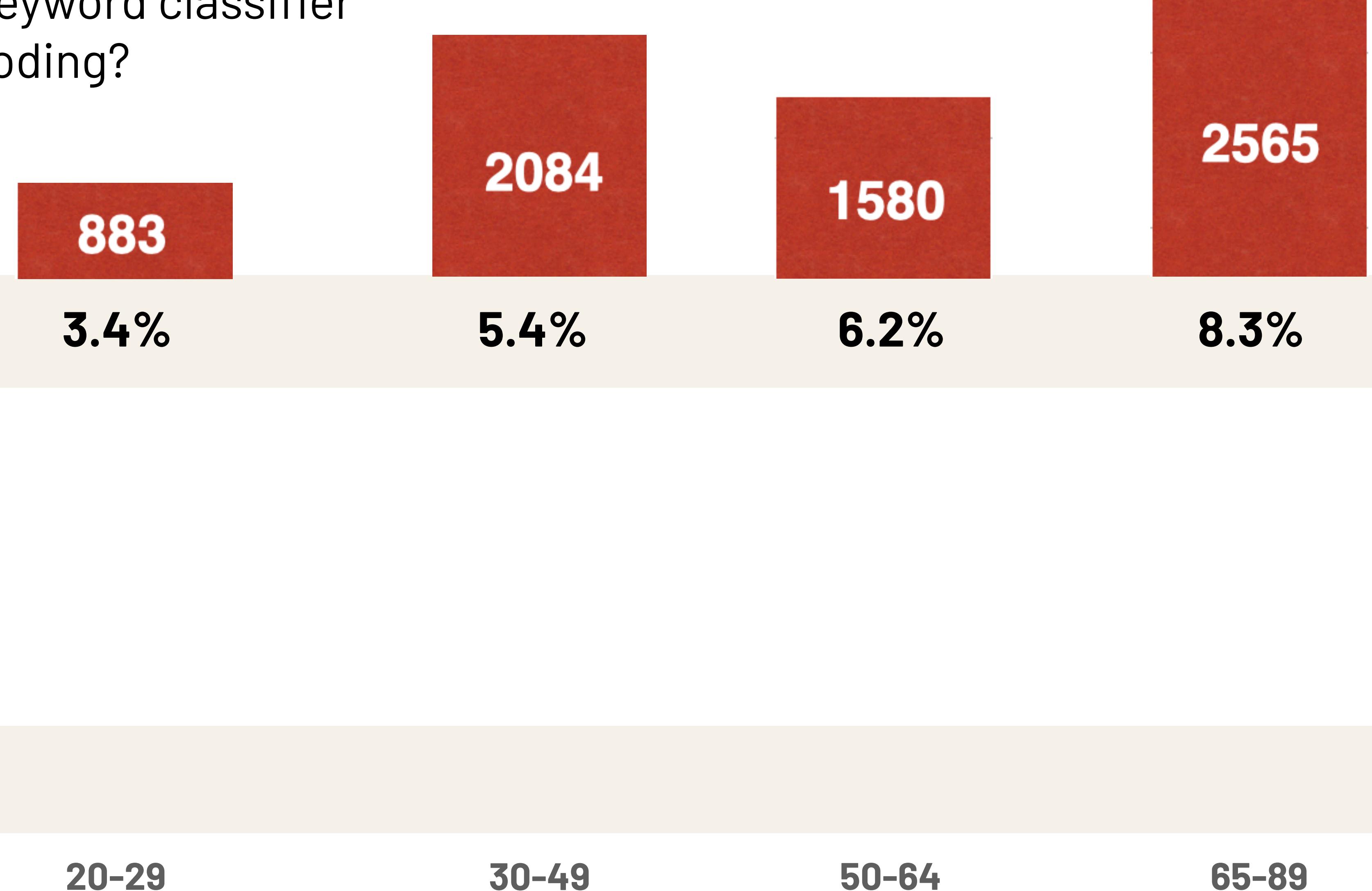# Keyword Classifier Accuracy

How frequently did the keyword classifier
disagree with our handcoding?

■ Classified as
  election-related

**2565**

**2084**

**1580**

**883**

**False Positive
Rate**

**3.4%**          **5.4%**          **6.2%**          **8.3%**

■ Classified as not
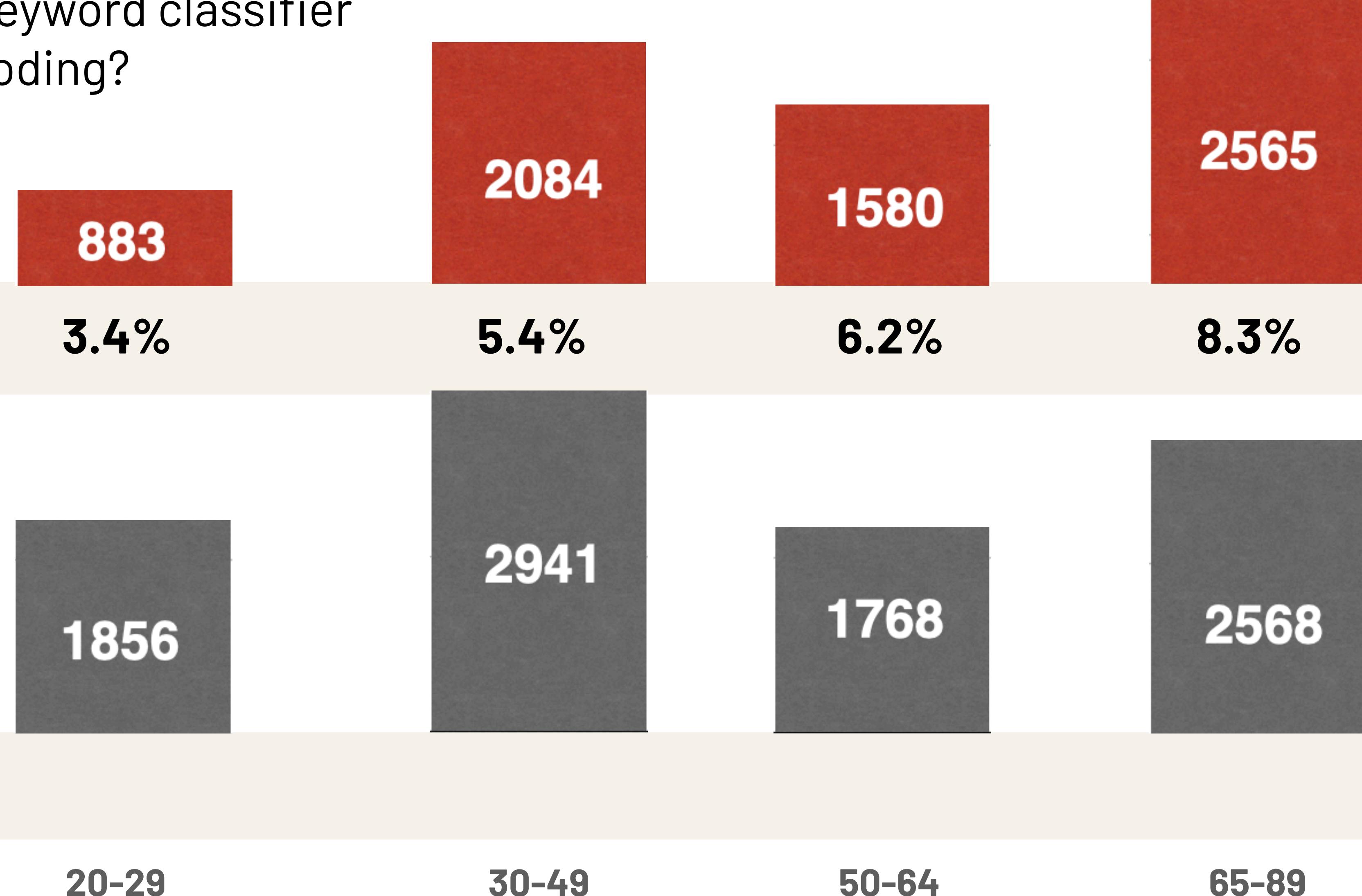  election-related

20-29          30-49          50-64          65-89

# Keyword Classifier Accuracy

How frequently did the keyword classifier disagree with our handcoding?

Classified as election-related

**883**

**2084**

**1580**

**2565**

**False Positive Rate**

**3.4%**

**5.4%**

**6.2%**

**8.3%**

Classified as not election-related

**1856**

**2941**

**1768**

**2568**

20-29

30-49

50-64

65-89

# Keyword Classifier Accuracy

How frequently did the keyword classifier disagree with our handcoding?

■ Classified as election-related

**883** | **2084** | **1580** | **2565**

**False Positive Rate** | 3.4% | 5.4% | 6.2% | 8.3%

■ Classified as not election-related

**1856** | **2941** | **1768** | **2568**

**False Negative Rate** | 36.6% | 39.3% | 25.1% | 16.1%

20-29 | 30-49 | 50-64 | 65-89

# Keyword Classifier Accuracy
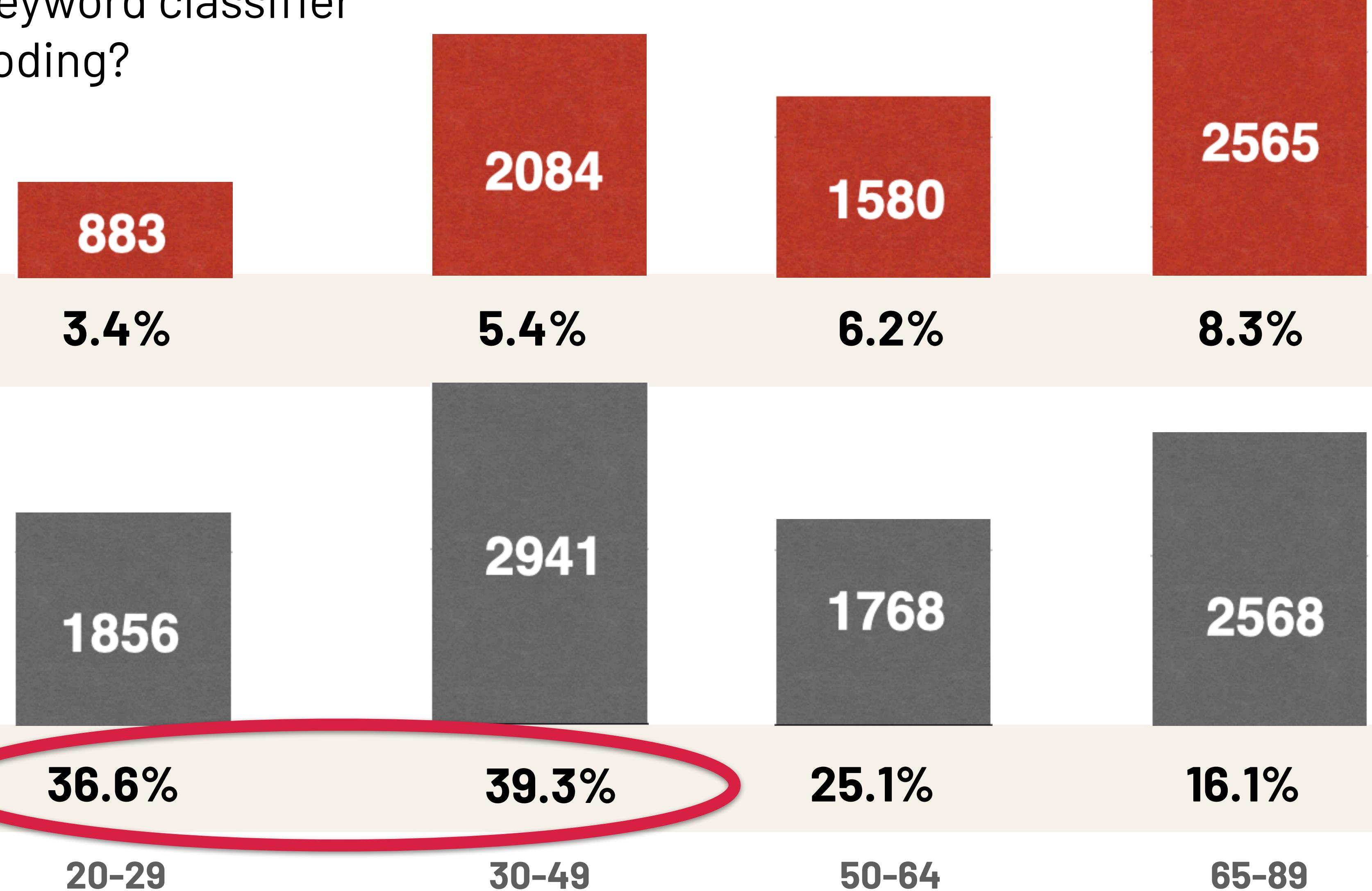
How frequently did the keyword classifier
disagree with our handcoding?

■ Classified as
   election-related

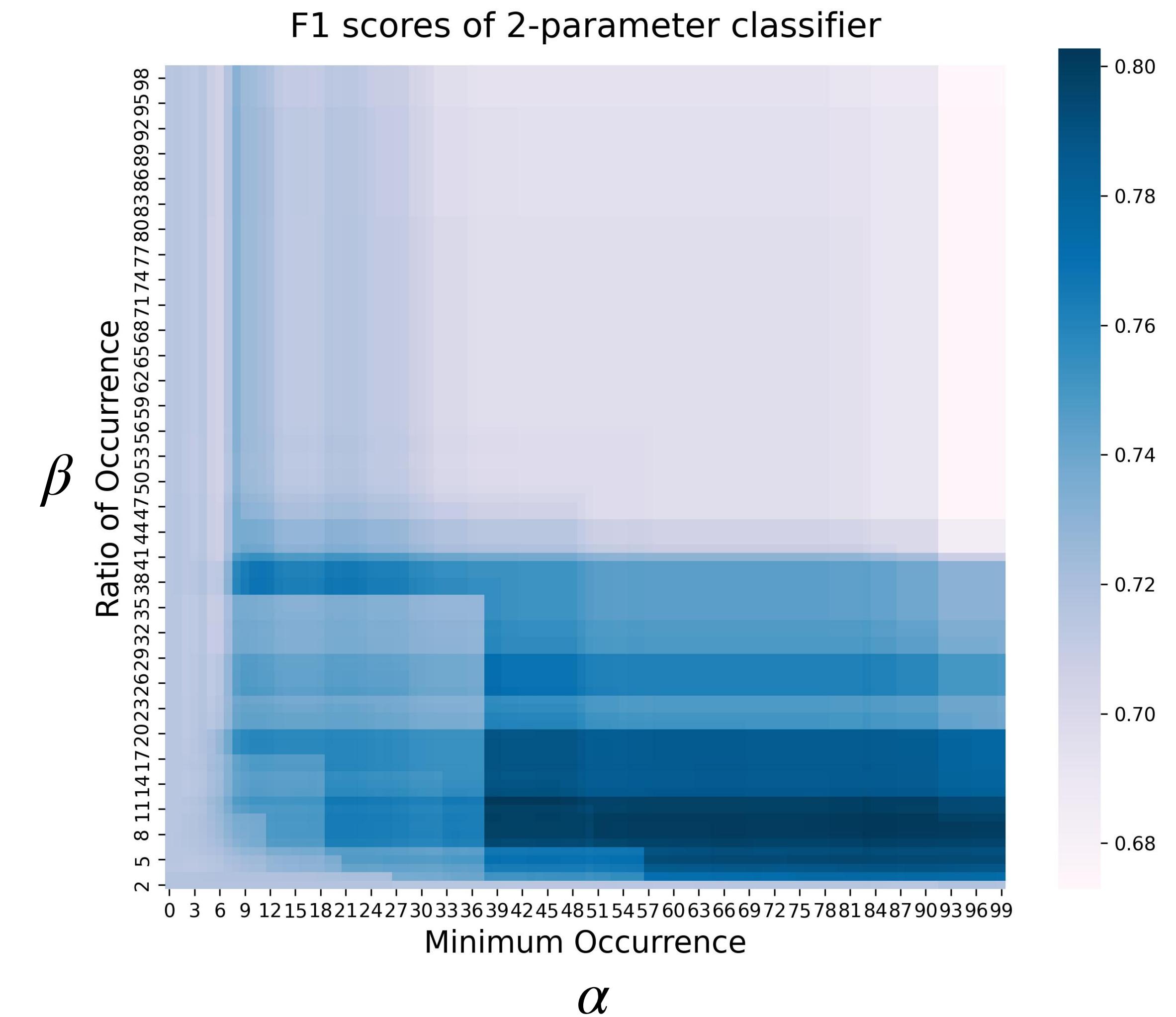**False Positive Rate**

| | 20–29 | 30–49 | 50–64 | 65–89 |
|---|---|---|---|---|
| Count | 883 | 2084 | 1580 | 2565 |
| Rate | 3.4% | 5.4% | 6.2% | 8.3% |

■ Classified as not
   election-related

**False Negative Rate**

| | 20–29 | 30–49 | 50–64 | 65–89 |
|---|---|---|---|---|
| Count | 1856 | 2941 | 1768 | 2568 |
| Rate | 36.6% | 39.3% | 25.1% | 16.1% |

# Could we build a better keyword classifier?



F1 scores of 2-parameter classifier

# Could we build a better keyword classifier?

✦ Test 10,000 potential classifiers through a two-parameter model



F1 scores of 2-parameter classifier
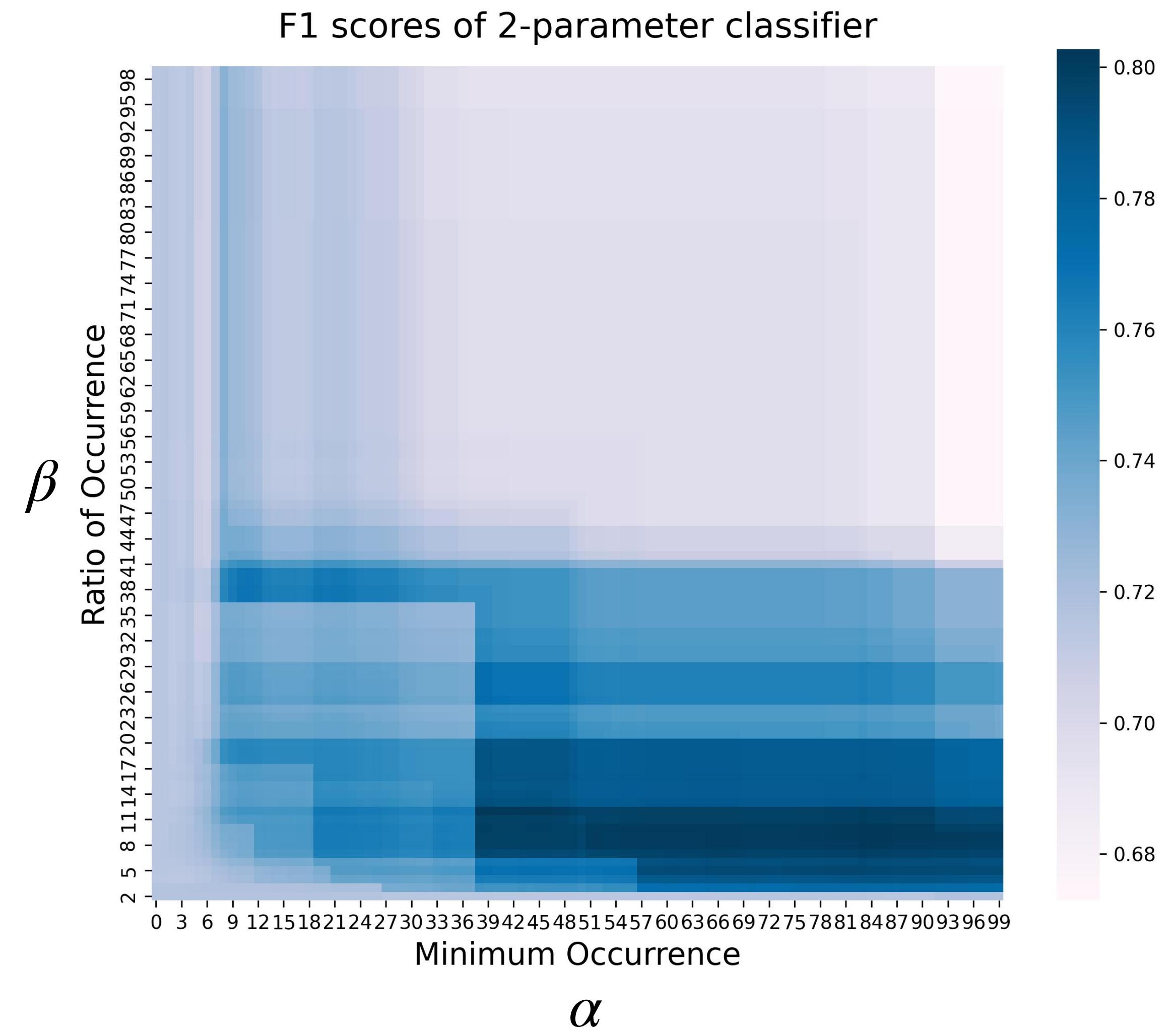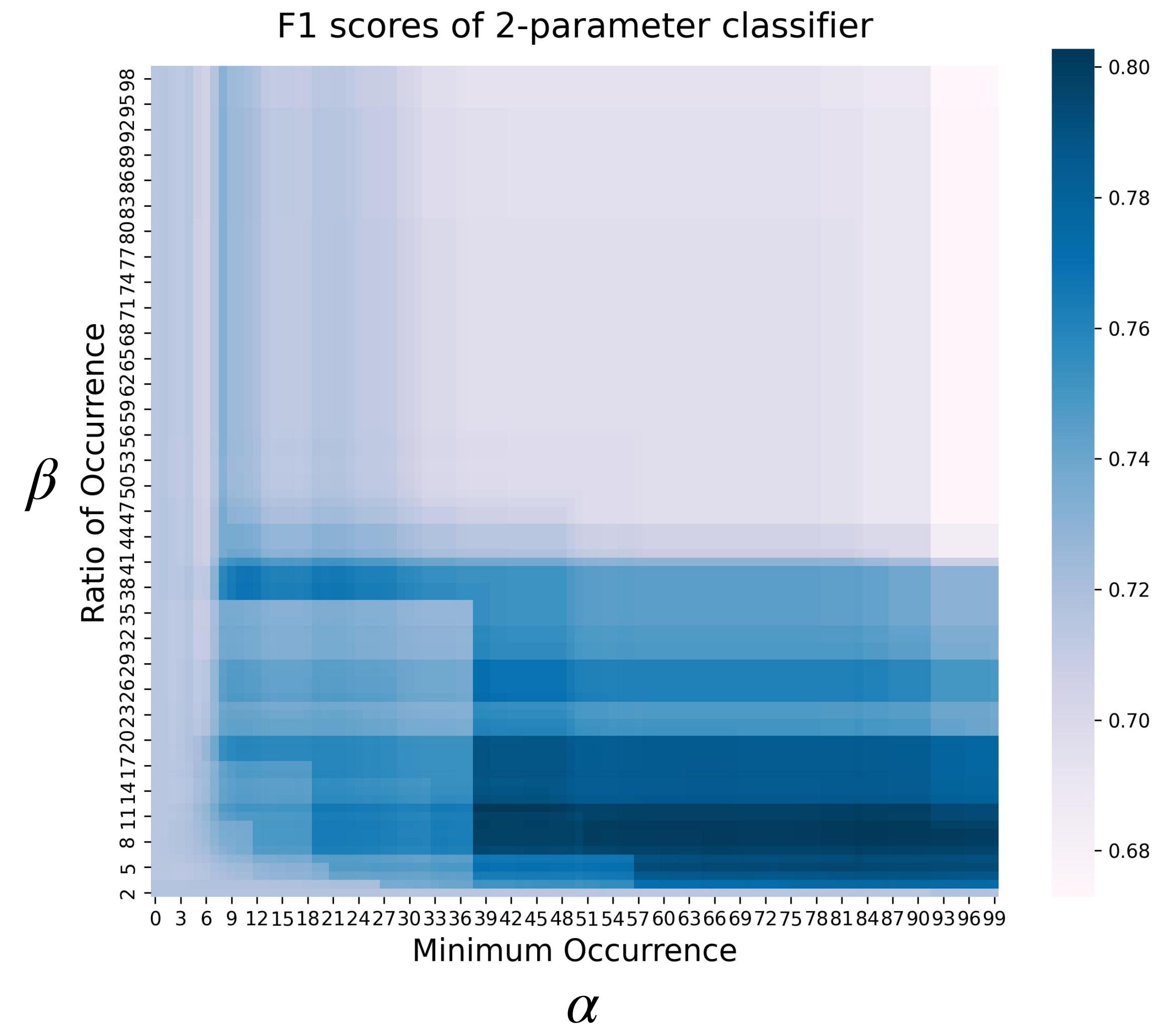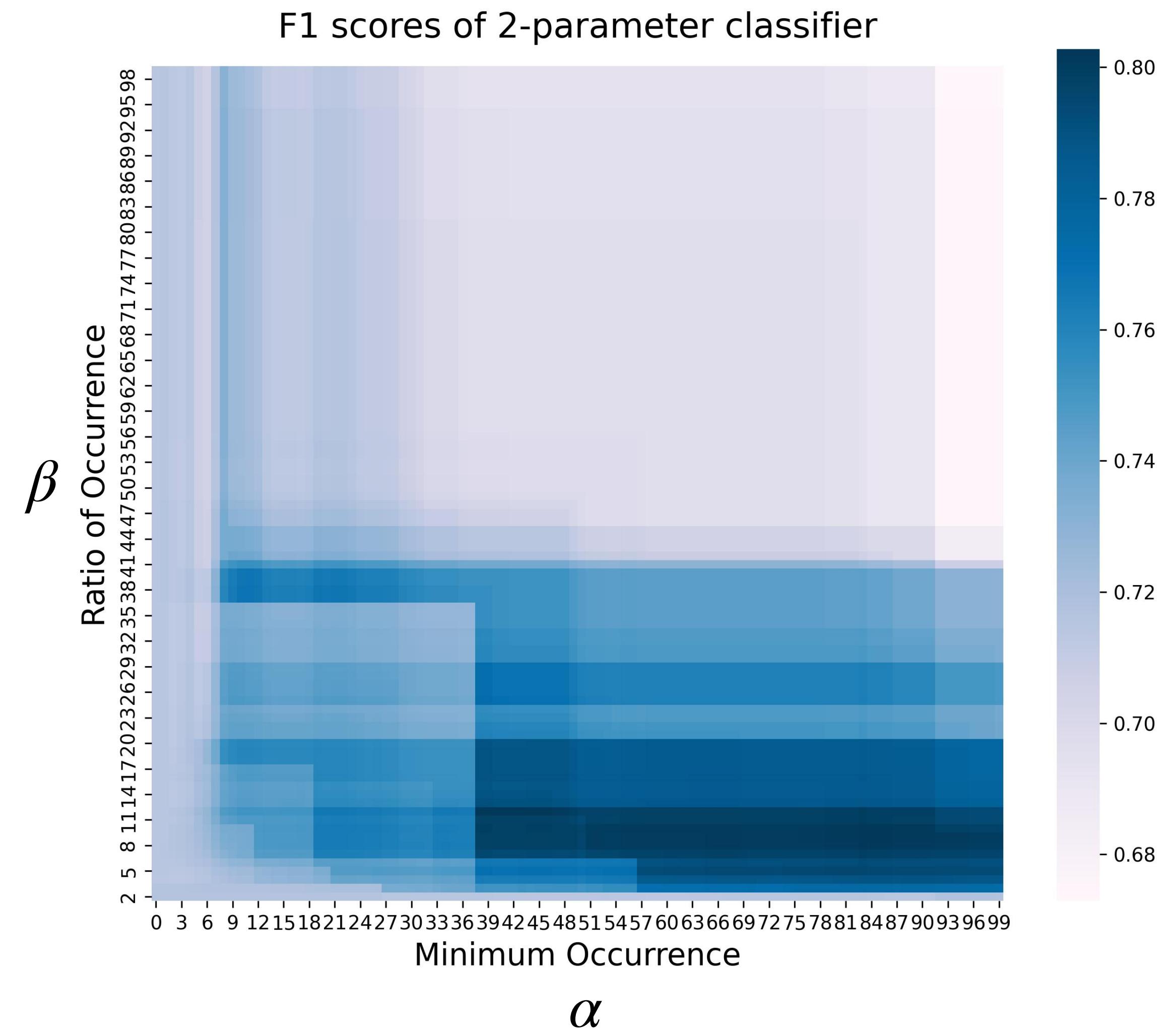
# Could we build a better keyword classifier?

✦ Test 10,000 potential classifiers through a two-parameter model

✦ Terms included in keyword list if:



F1 scores of 2-parameter classifier

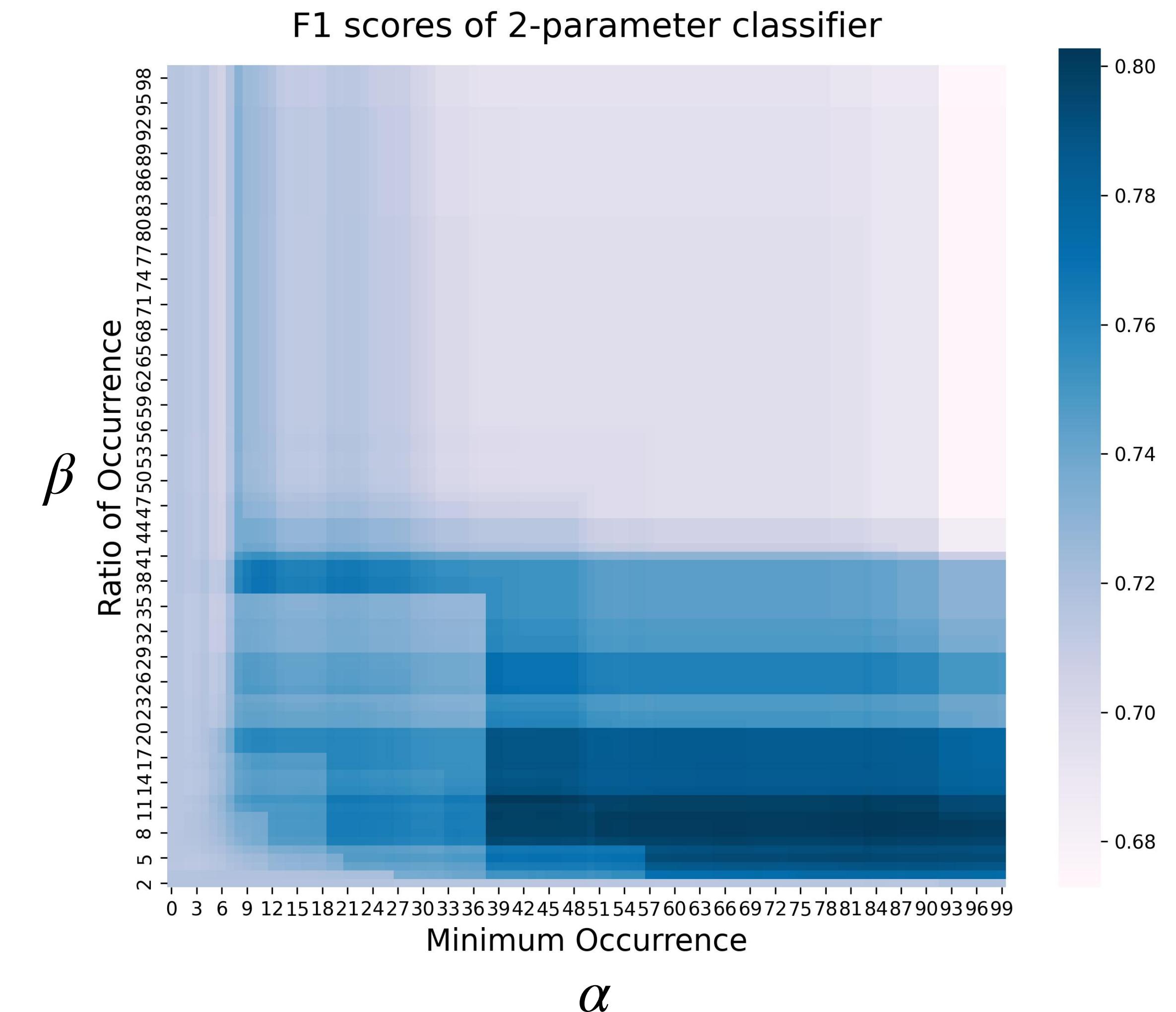# Could we build a better keyword classifier?

✦ Test 10,000 potential classifiers through a two-parameter model

✦ Terms included in keyword list if:

➡ Occur at least $\alpha$ times



F1 scores of 2-parameter classifier

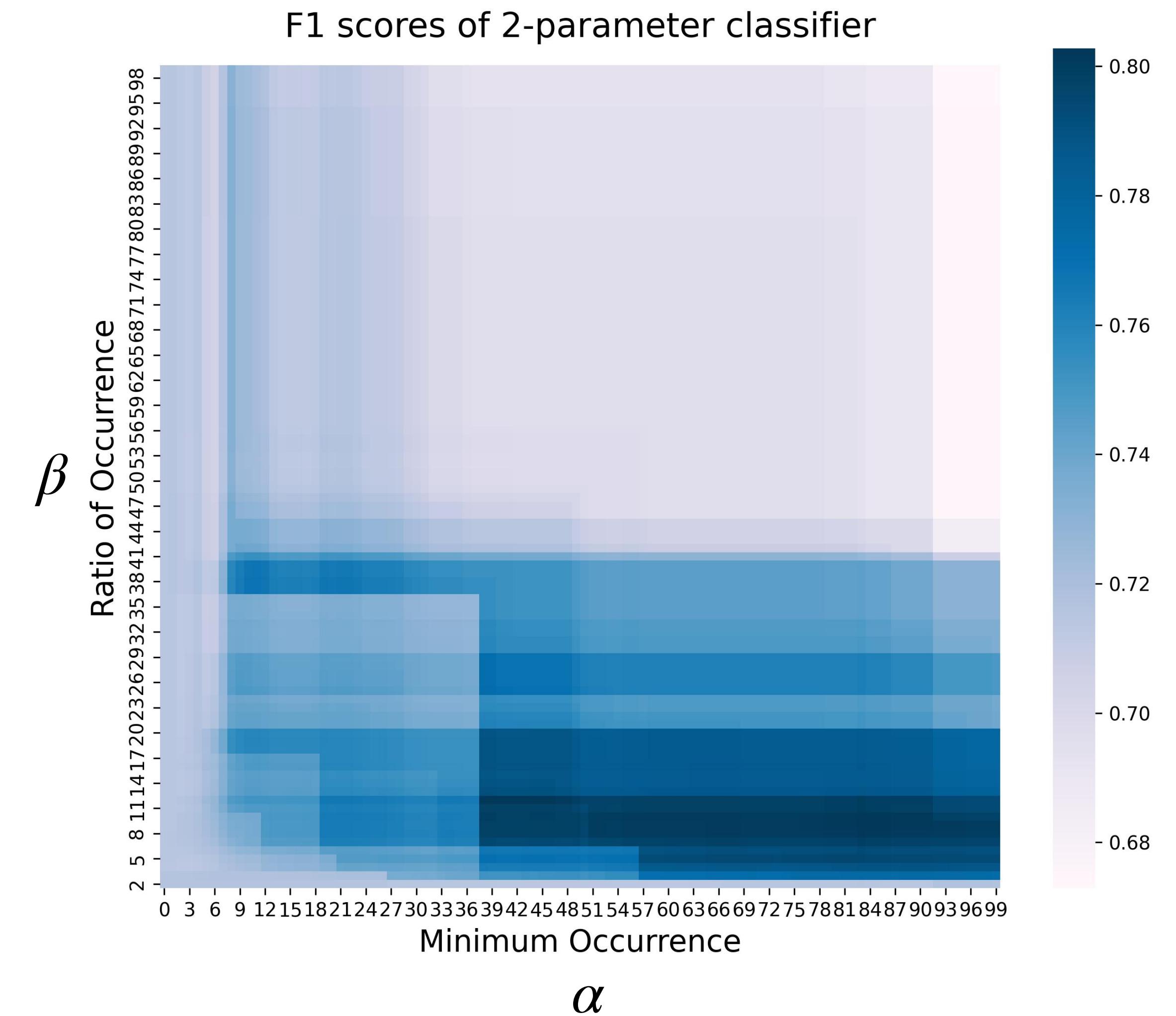$\beta$ — Ratio of Occurrence

$\alpha$ — Minimum Occurrence

# Could we build a better keyword classifier?

- ✦ Test 10,000 potential classifiers through a two-parameter model

- ✦ Terms included in keyword list if:

  - ➡ Occur at least $\alpha$ times

  - ➡ Occur $\beta$ times more in election tweets (defined from handcoding)
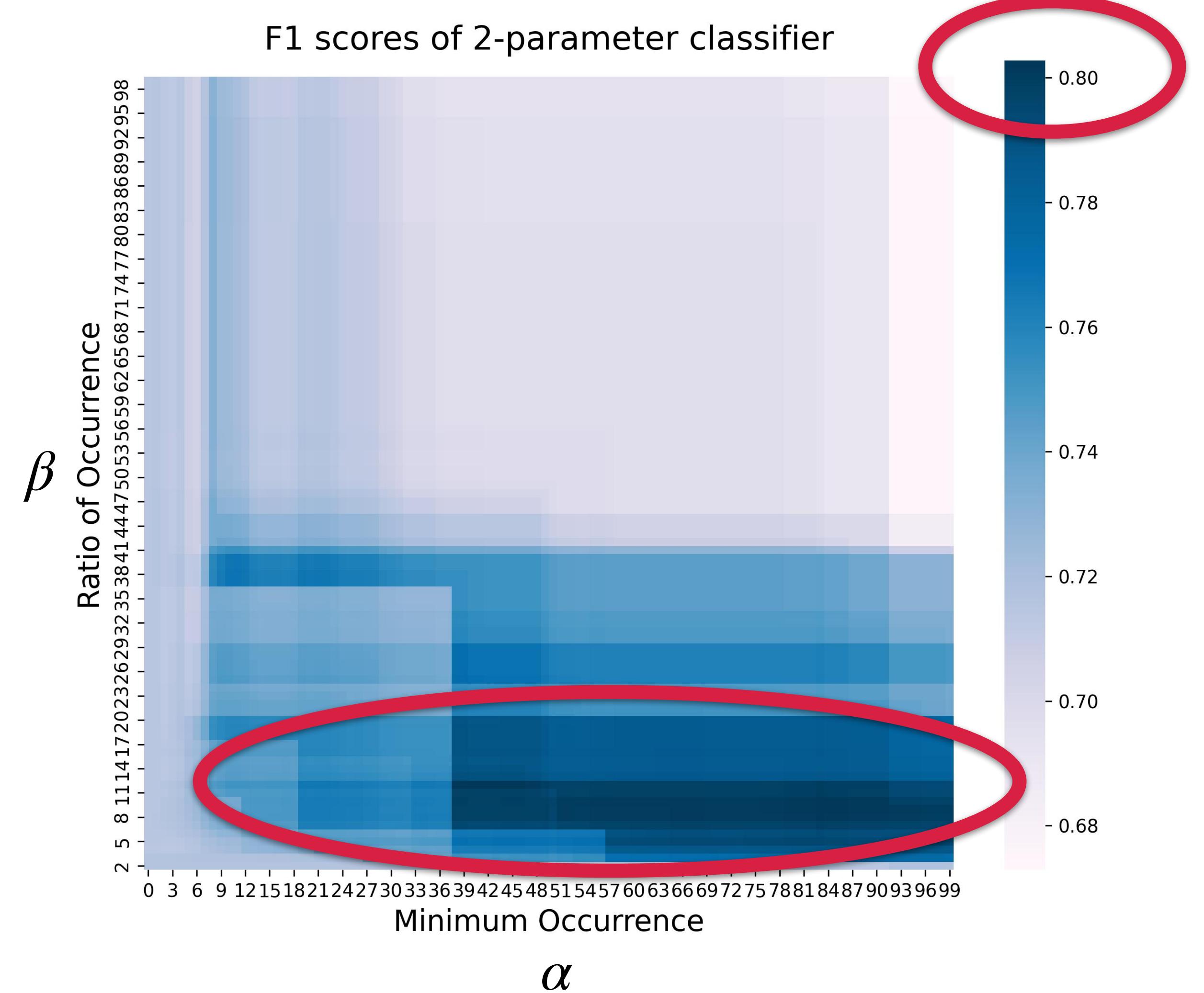


F1 scores of 2-parameter classifier

# Could we build a better keyword classifier?

- ✦ Test 10,000 potential classifiers through a two-parameter model

- ✦ Terms included in keyword list if:
  - ➡ Occur at least $\alpha$ times
  - ➡ Occur $\beta$ times more in election tweets (defined from handcoding)

- ✦ Calculate accuracy for each resulting keyword list
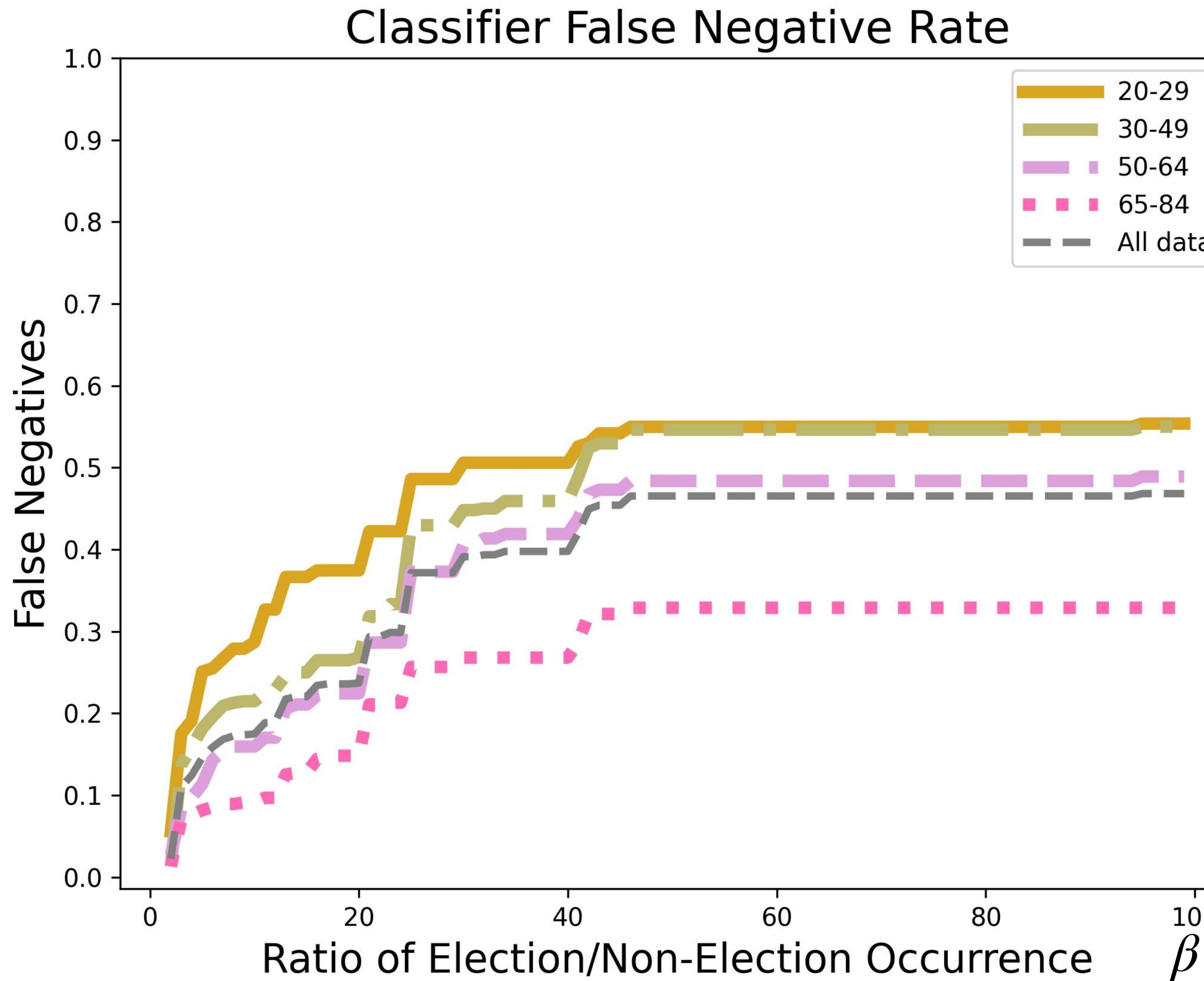


F1 scores of 2-parameter classifier
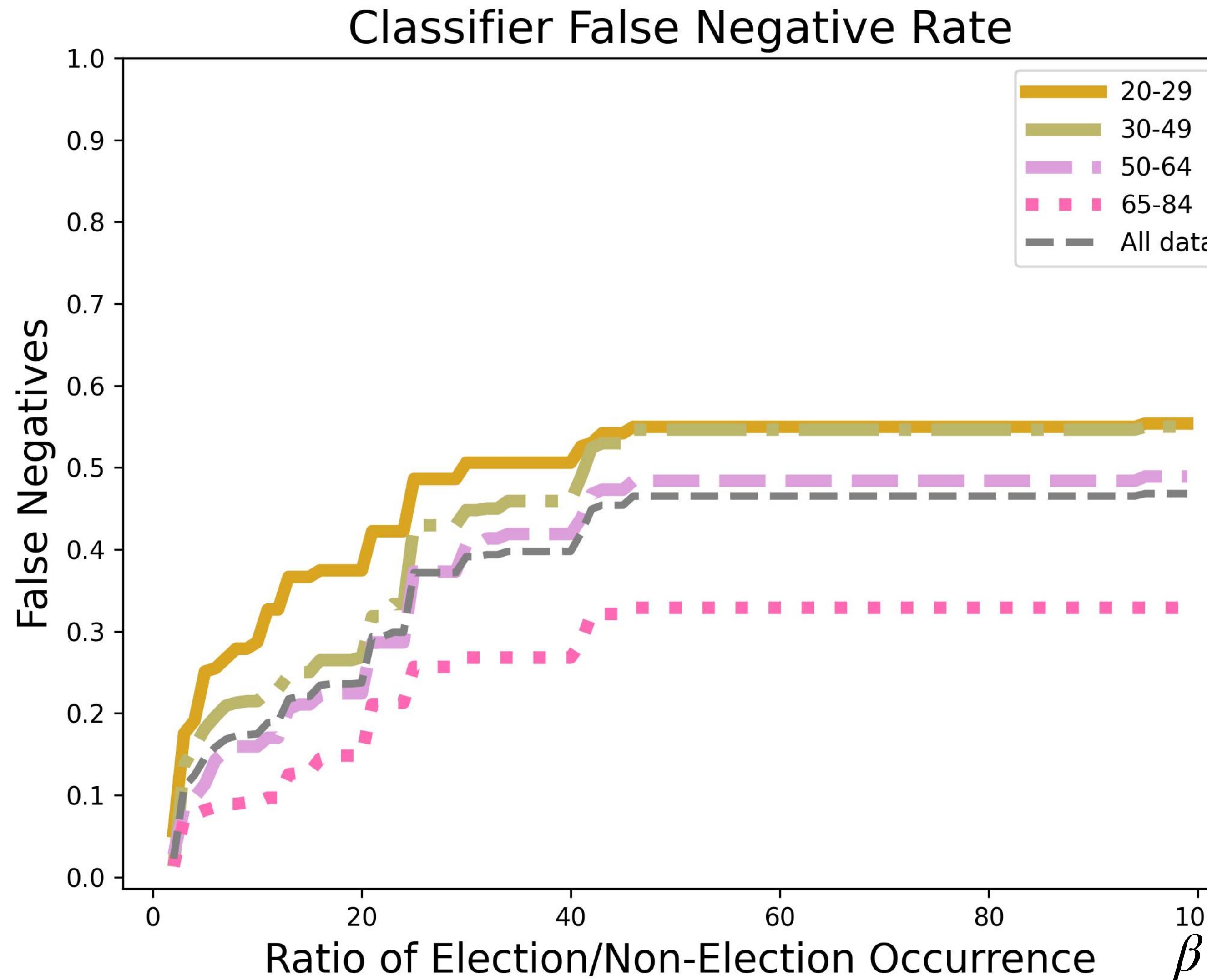
# Could we build a better keyword classifier?

- Test 10,000 potential classifiers through a two-parameter model

- Terms included in keyword list if:

  ➡ Occur at least $\alpha$ times

  ➡ Occur $\beta$ times more in election tweets (defined from handcoding)

- Calculate accuracy for each resulting keyword list

- Best models achieved (only) an F1-score of ~ 0.8



F1 scores of 2-parameter classifier

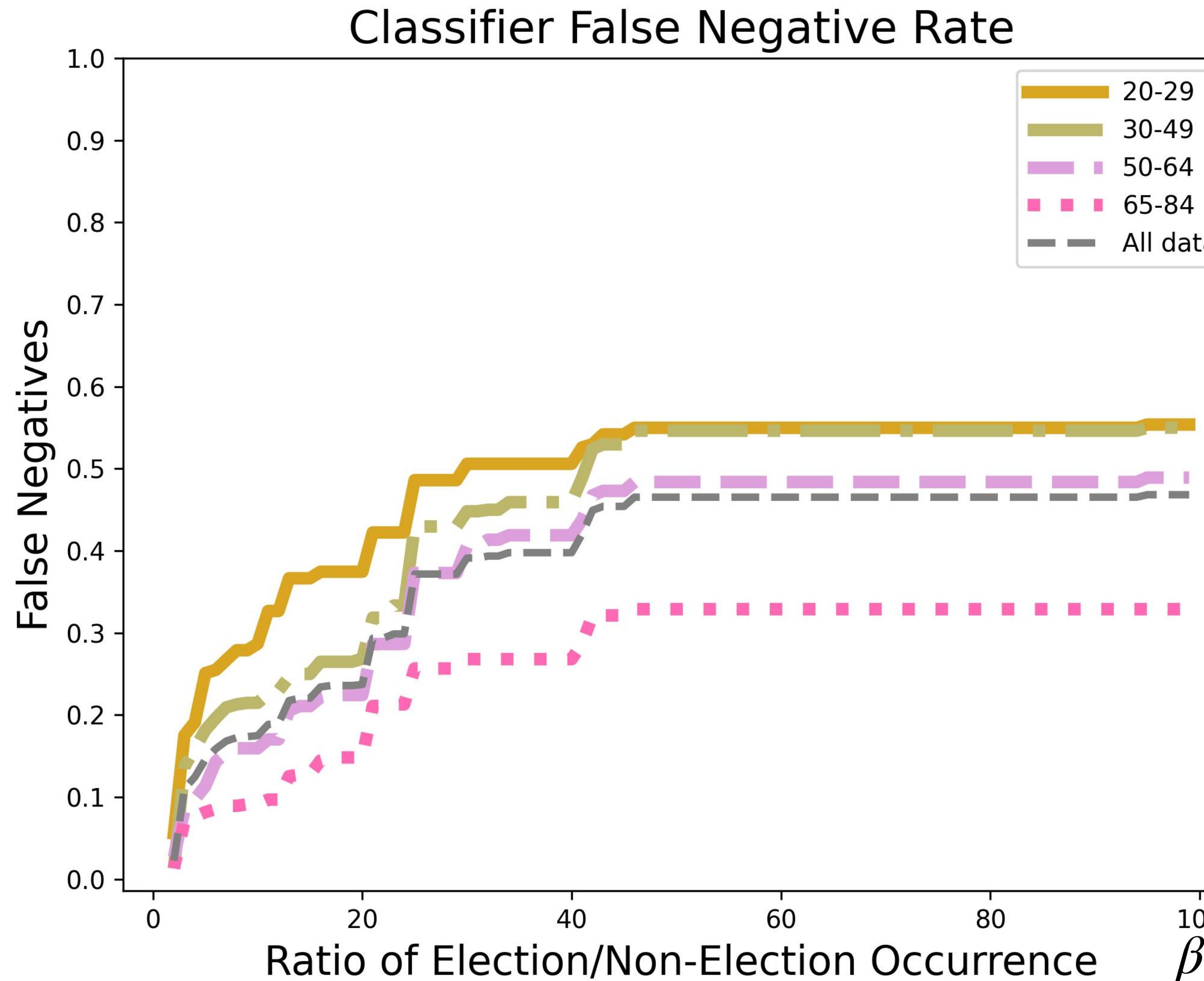# Could we build a better keyword classifier?



Classifier False Negative Rate

# Could we build a better keyword classifier?



Classifier False Negative Rate

Older users consistently have fewer of their election tweets misclassified

# Could we build a better keyword classifier?

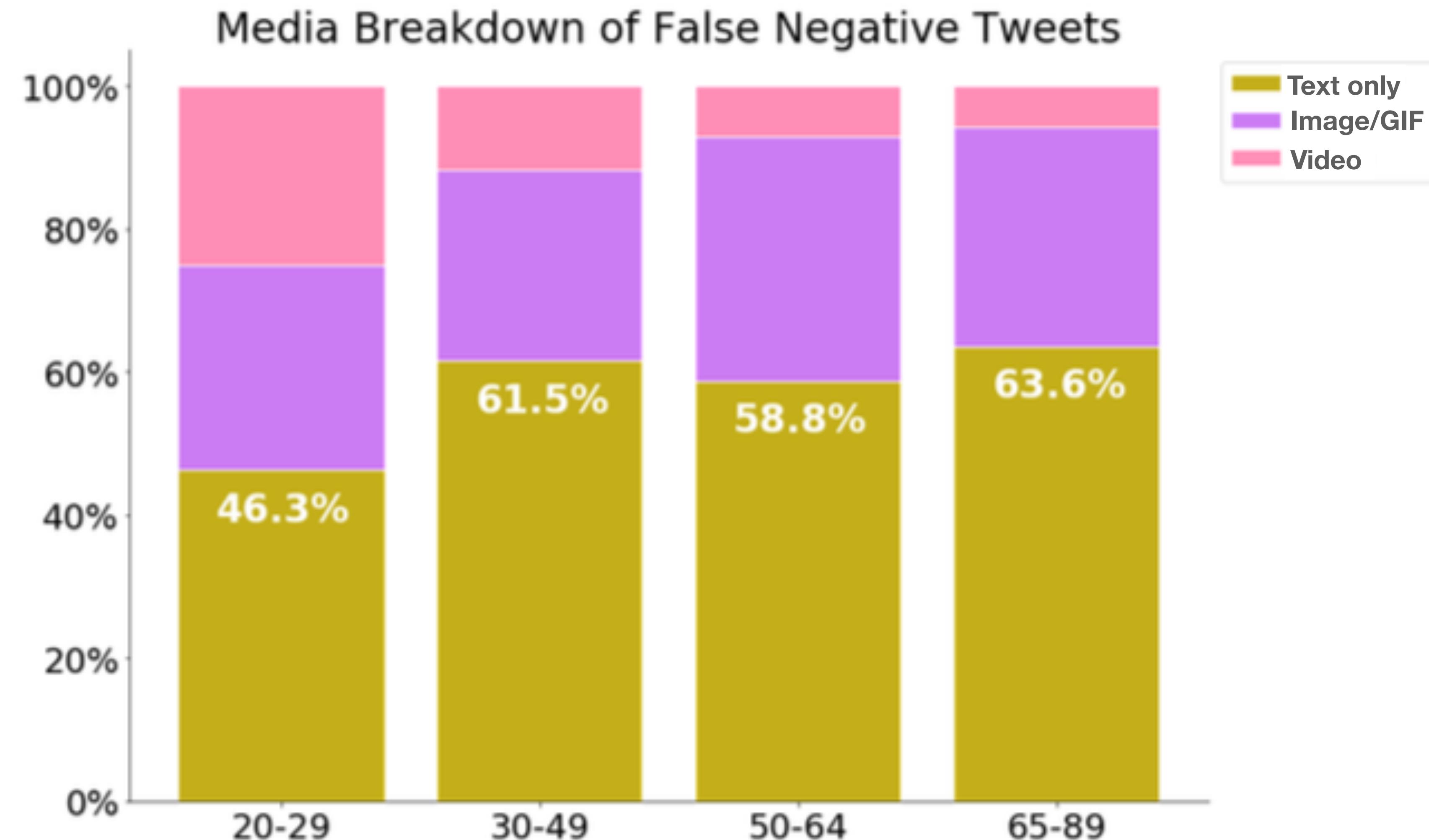## Classifier False Negative Rate



Younger users consistently have more of their election tweets misclassified

Older users consistently have fewer of their election tweets misclassified

# Generational Trends?

Younger users frequently communicate using non-textual media, so keyword classifiers may underrepresent their political speech



Media Breakdown of False Negative Tweets

# Final Thoughts

# Final Thoughts

✦ How we measure "speech" has implications for whose speech we capture

# Final Thoughts

- ✦ How we measure "speech" has implications for whose speech we capture

- ✦ Matched panel data can help us interrogate these measurement biases

# Final Thoughts

✦ How we measure "speech" has implications for whose speech we capture

✦ Matched panel data can help us interrogate these measurement biases

➡ Have meaningful "negative" samples (eg, all users' posts)

# Final Thoughts

- ✦ How we measure "speech" has implications for whose speech we capture

- ✦ Matched panel data can help us interrogate these measurement biases

  - ➡ Have meaningful "negative" samples (eg, all users' posts)

  - ➡ Have age and other demographic info, can compare impact across categories

# Final Thoughts

✦ How we measure "speech" has implications for whose speech we capture

✦ Matched panel data can help us interrogate these measurement biases

➡ Have meaningful "negative" samples (eg, all users' posts)

➡ Have age and other demographic info, can compare impact across categories

## Thank you!

**Sarah Shugars**
*they/them/theirs*

Assistant Professor
Rutgers University

sarah.shugars@rutgers.edu
BlueSky: @Shugars