# Project 3:
# Classification of Reddit Posts

Team Members:
Hong Aik
Mitchelle
Shu Yi
Yong Gui
Wee Hong

zoom

# Introduction

# Problem Statement
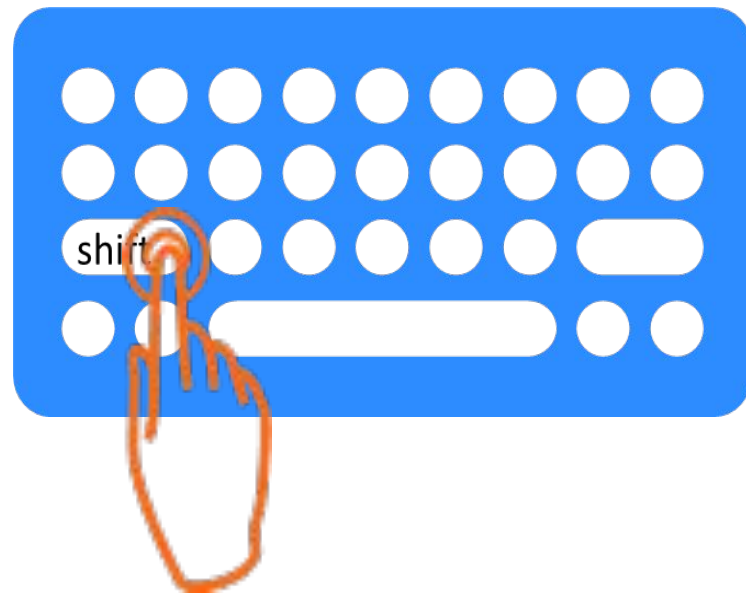
# Collection
# & Data Cleaning
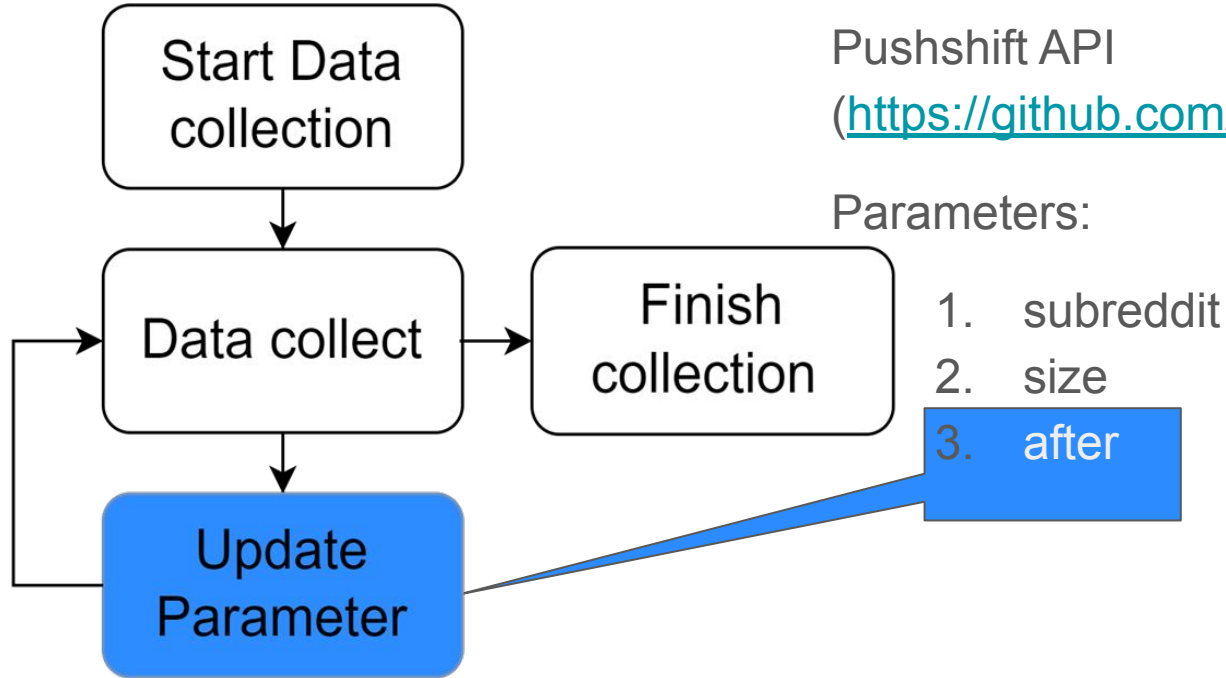
# Data Collection

Pushshift API (https://github.com/pushshift/api)

Parameters:

1. subreddit: "Zoom" / "MicrosoftTeams"
2. size: 100 (maximum)
3. after: *epoch value*
   a. first value: 1577836800
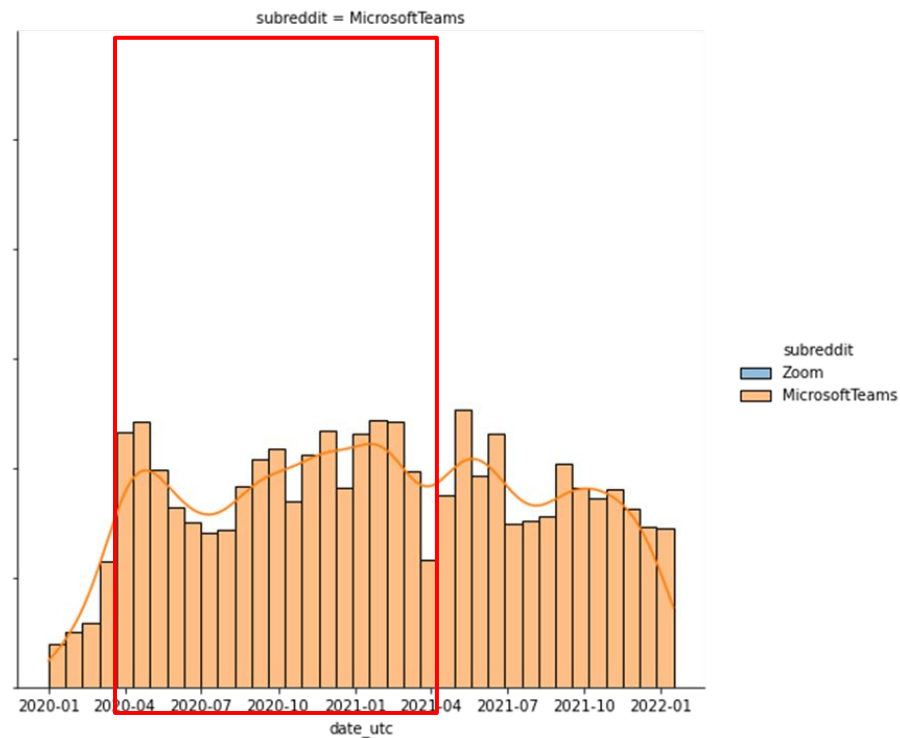   b. Data and Time (GMT): 1 Jan 2020, 00:00

# Data Collection



Pushshift API
(https://github.com/pushshift/api)
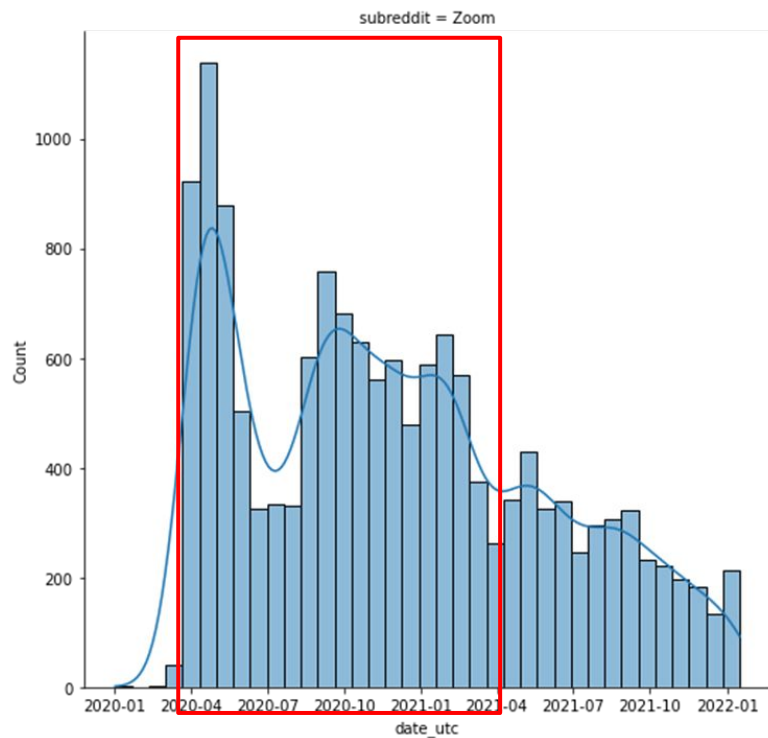
Parameters:

1. subreddit
2. size
3. after

# Cleaning

Cleaning of combined column of selftext and title:

- HTML Special entities (e.g. &amp)
- Hyperlinks
- Punctuation
- Whitespace
- Characters beyond Basic Multilingual Plane (BMP) of Unicode
- [removed]
- [deleted]

zoom

# Time Period

Time Series graph for time frame selection

# Preprocessing

Lemmatize

```python
1  from nltk.stem import WordNetLemmatizer
2  lemmatizer = WordNetLemmatizer()
```

```python
1  words=['feet','dogs','children','identify','this']
2  for word in words:
3      print(f"{word}: {lemmatizer.lemmatize(word)}")
```

```
feet: foot
dogs: dog
children: child
identify: identify
this: this
```
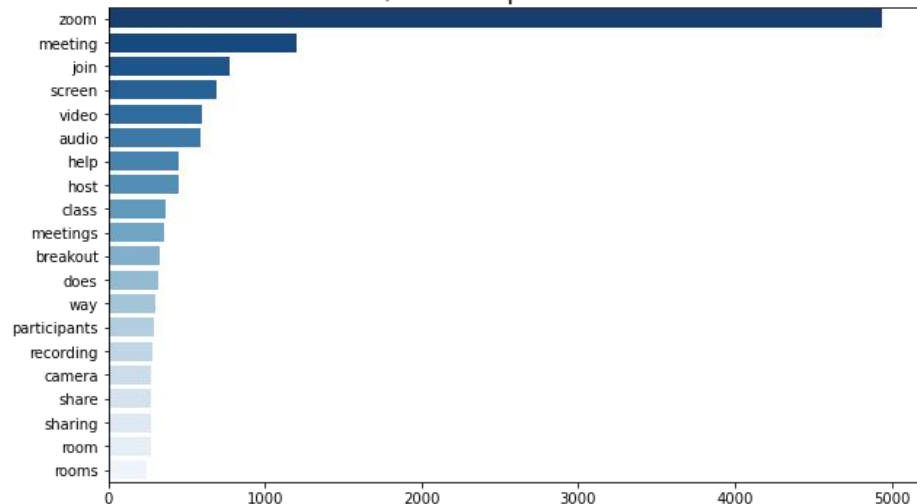
*Image from
https://karkig.medium.com/understand-stemming-and-lemmatization-with-python-nltk-package-77973a727040*

# EDA

# Analyzing Words in Title



r/Zoom Top 20 Words

r/MicrosoftTeams Top 20 Words

zoom

# Analyzing Words in Selftext



r/Zoom Top 20 Trigrams in SelfText

r/MicrosoftTeams Top 20 Trigrams in SelfText

zoom

# Sentiment Analysis



Sentiment Analysis of Both Post

r/Zoom (Average Compound Score: 0.18)

r/MicrosoftTeams (Average Compound Score: 0.23)

**Observation**: Based on the average compound score computed, it seems to indicate that there are more positive posting on Team as compared to Zoom (this can also be observed in the distribution on the right.

# Example of Negative & Positive Posting about Zoom

**<u>Negative Posting about Zoom:</u>**

zoom troubleshooting help alright i have a weird issue i have one student in a class of  that can join the zoom meeting but is then dropped seconds after   no other student is experiencing this issue  the teacher is not having any network lag  the student can access and watch other classes just fine  student is using a chromebook  student is using a hotspot  student hasnt indicated this has been a problem for the last month  issue only occurs at am class  father is adamant that the issue is with zoom or the teacher but i believe that can be ruled out since no other student is having the same problem  the device can be ruled out as a problem since connection with other classes is fine  my thought is the hotspot they are currently using a sprint mifi unlimited plan after some quick research i found that sprint deprioritizes connections after theyve used gb in a billing cycle i am assuming the am time is a busy network for the tower  any other potential ideas what it could be

**<u>Positive Posting about Zoom:</u>**

trying to host workouts on zoom and i cant figure out how to play spotify in my meeting without screen share apparently there is a  program called loopback that does exactly this but im hoping someone knows of a free option thanks in advance

zoom

# Example of Negative & Positive Posting about Teams

**Negative Posting about Team:**

problems when logging in on laptop with error code whenever i want to start ms teams i always get a error code  and it asks me to restart the application  whenever i restart my application it again gives me the same error  if i log in on my other computer via google chrome or with a clean install it will work and i dont seem to have any problems does someone know whats the problem here

**Positive Posting about Team:**

generate report of all microsoft teams users hello  is it possible to create a report of all active microsoft teams users from the teams admin portal ideally i would like to create a report that includes usernamesemails of all the microsoft teams users   thanks

zoom

# Scattertext Visualization

# Modelling
# &
# ~~Evaluation~~

# Text Vectorization

# Text Vectorization - TfidfVectorizer



TF-IDF

TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$TF\text{-}IDF = TF(t,d) \times IDF(t)$$

Term frequency
Number of times term $t$ appears in a doc, $d$

Inverse document frequency

$$\log \frac{1+n}{1+df(d,t)} + 1$$

# of documents

Document frequency of the term $t$

|  | I | like | cats | me |
|---|---|---|---|---|
| I like cats | 0.1 | 0.1 | 0.1 | 0 |
| Cats like me | 0 | 0.1 | 0.1 | 0.1 |

zoom

# Further EDA



We can also employ t-SNE for dimensionality reduction. We observe that the clusters are denser when using Word2Vec and clusters are better separated.

# Further EDA

We will experiment with 3 classifiers:

- Logistic Regression (LR)

- Random Forest (RF)

- MultiLayer Perceptron Classifier (MLP)

Baseline:

Checking whether submission contains the word "microsoft" or "teams"

Baseline score: **87%**

**MLP Architecture**

| | | Hidden Layer | | |
| Input Layer | | | Output Layer | |

Input #1 →

Input #2 →

Input #3 →

Input #4 →

→ Output

*"Black box"*

zoom

# Data Modeling



Input text data

↓

3 vectorizers X 3 classifiers = 9 models

*Train-test-split*

*Initial fit and score for baseline*

*RandomizedSearchCV*

Best model

*GridSearchCV*

Final model

# Model Evaluation



We managed to improve our score by ~1% which is decent, given that our untuned model already has 93% accuracy!

# Model Evaluation



Confusion Matrix

| | Confidence |
|---|---|
| **Predicting MST** | 88% |
| **Predicting Zoom** | 84% |

zoom

# Deploying the Model

## Reddit Classification Web App

The underlying model was trained on ~14,000 sub-reddits from r/Zoom and r/MicrosoftTeams, with the goal of predicting the sub-reddit given a string of words (submission)

The model is only able to output 2 possible results!

Type your content here!

I am having trouble with virtual backgrounds!

Click for predictions!

With 75% confidence, this is a submission belonging to r/Zoom.

zoom

# Conclusion

# Recommendation for Software Development Team

## Pain points for users

1. Stopped
2. Drop
3. Crash
4. Reinstall
5. Error

## Keep an eye on

**Discord**

a VoIP, instant messaging and digital distribution platform

zoom

# Recommendation for Digital Marketing Team

## Top words for Zoom

1. Zoom
2. Password
3. Host
4. Join
5. Participant
6. Class
7. Virtual
8. Breakout
9. Room
10. Id

## Top words for MST

1. Team
2. Guest
3. Microsoft
4. Channel
5. User
6. Assignment
7. Call
8. Notification
9. Chat
10. Feature

# Zooming ahead

## Refining our current model

Training our model to recognise words unique to subreddit.

## Running the refined model on other competitor pairing

Zoom vs Google, Zoom vs Skype

# Annex

## Model Vectoriz..

| Model Vectorizer | | Score |
|---|---|---|
| BERT-LogisticRegression | Train Score | 0.92 |
| | Test Score | 0.89 |
| BERT-MLPClassifier | Train Score | 1.00 |
| | Test Score | 0.88 |
| BERT-RandomForest | Train Score | 1.00 |
| | Test Score | 0.83 |
| Tfidf-LogisticRegression | Train Score | 0.96 |
| | Test Score | 0.93 |
| Tfidf-MLPClassifier | Train Score | 1.00 |
| | Test Score | 0.86 |
| Tfidf-RandomForest | Train Score | 1.00 |
| | Test Score | 0.92 |
| Word2Vec-LogisticRegression | Train Score | 0.91 |
| | Test Score | 0.89 |
| Word2Vec-MLPClassifier | Train Score | 1.00 |
| | Test Score | 0.89 |
| Word2Vec-RandomForest | Train Score | 1.00 |
| | Test Score | 0.85 |

Score

## Model Vectorizer

| Model Vectorizer | Cross Val Score |
|---|---|
| BERT-LogisticRegression | 0.89 |
| BERT-MLPClassifier | 0.88 |
| BERT-RandomForest | 0.83 |
| Tfidf-LogisticRegression | 0.92 |
| Tfidf-MLPClassifier | 0.90 |
| Tfidf-RandomForest | 0.91 |
| Word2Vec-LogisticRegression | 0.89 |
| Word2Vec-MLPClassifier | 0.89 |
| Word2Vec-RandomForest | 0.86 |

Cross Val Score

# Text Vectorization and Further EDA - TfidfVectorizer



TF-IDF is a measure of originality of a word by comparing the number of times a word appears in a doc with the number of docs the word appears in.

$$TF\text{-}IDF = TF(t, d) \times IDF(t)$$

Term frequency

Number of times term $t$ appears in a doc, $d$

Inverse document frequency

$$\log \frac{1 + n}{1 + df(d,t)} + 1$$

# of documents

$n$ ← documents

Document frequency of the term $t$

|  | I | like | cats | me |
|---|---|---|---|---|
| I like cats | 0.1 | 0.1 | 0.1 | 0 |
| Cats like me | 0 | 0.1 | 0.1 | 0.1 |

zoom

# Text Vectorization and Further EDA - Word2Vec



|  | Dim 1 | Dim 2 | ... | Dim 300 |
|---|---|---|---|---|
| I like cats | 0.01 | 0.01 | ... | 0.01 |
| Cats like me | 0.01 | 0.01 | ... | 0.01 |

zoom

# Text Vectorization and Further EDA - Bi-directional Encoder Representations from Transformers



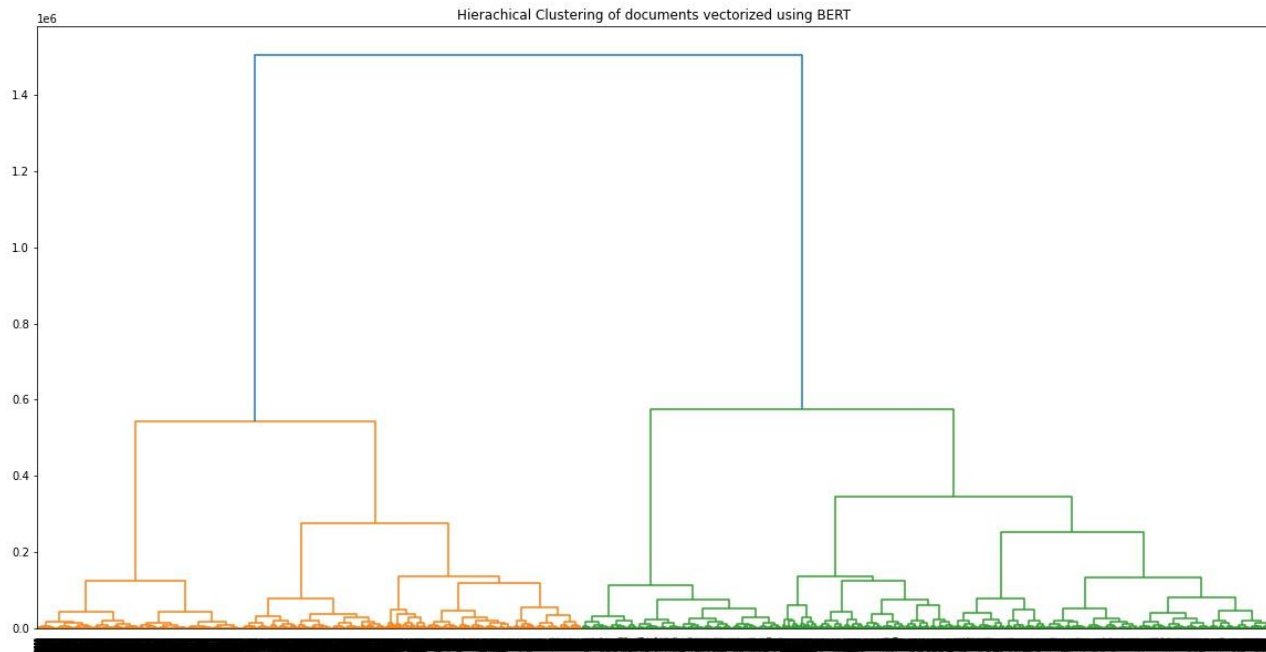|  | Dim 1 | Dim 2 | … | Dim 768 |
|---|---|---|---|---|
| I like cats | 0.01 | 0.01 | … | 0.01 |
| Cats like me | 0.01 | 0.01 | … | 0.01 |

zoom

# Exploring Misclassifications

| | predictions | true | probability | text_nostop |
|---|---|---|---|---|
| **640** | 0 | 1 | 0.467682 | hi schedule recurring meeting phone app show online account anyone know sync these tia app sync |
| **9965** | 1 | 0 | 0.835560 | please help camera hidden share screen participant get bored |
| **10197** | 1 | 0 | 0.815036 | aspect ratiorequirements virtual background |
| **4410** | 0 | 1 | 0.480724 | sit window left side table take call face look like halfmoon one side lighted side lighted course could shift 90 degree face window need refer desktop pc time make sense keep alternating window back pc every time take video call idea one problem thought getting ring light seem many option webcam ring light even ring light seem good enough illuminate face compensate uneven lighting deal uneven lighting taking video call |
| **7775** | 1 | 0 | 0.510970 | suggestion good way prevent student cheating multiple choice form quiz example failing student finish difficult physic quiz 7 second score 100 something definitely i give openended question would like make testing process somewhat manageable right cheating |
| **479** | 0 | 1 | 0.425665 | mute presenter shared sound presenter know shared sound |

zoom

# Further EDA



Hierachical Clustering of documents vectorized using BERT

Agglomerative clustering shows that there are 2 distinct clusters in our word vectors (BERT), although this may not necessarily correspond to our Zoom and MST sub-reddits.

# Further EDA



Distribution of Predicted Labels

The clustering earlier has nicely separated the data into 2 clusters, although the overall accuracy is not high.

The labels have also been "inverted", although we can be fairly certain that 0 (actual) = 1 (predicted)

Distribution of Actual Labels

We see a rather messy overlap of points from Zoom (blue) and MST (yellow) texts, although there is a strong Zoom cluster on the left. This is expected given that we have "over-simplified" the data

zoom

# Understanding the Final Model



Top and bottom 20 tokens contributing to the model predictions of subreddits

# Zoom Video Communication

Q1 FY21
Earnings  June 2, 2020

# Use of Non-GAAP Financial Measures

In addition to the financials presented in accordance with U.S. generally accepted accounting principles ("GAAP"), this presentation includes the following non-GAAP metrics: non-GAAP gross margin, non-GAAP operating expenses, non-GAAP operating margin, non-GAAP operating income, non-GAAP EPS and non-GAAP Free Cash Flow. Non-GAAP metrics have limitations as analytical  tools and you should not consider them in isolation or as a substitute for or superior to the most directly comparable financial  measures prepared in accordance with U.S. GAAP. There are a number of l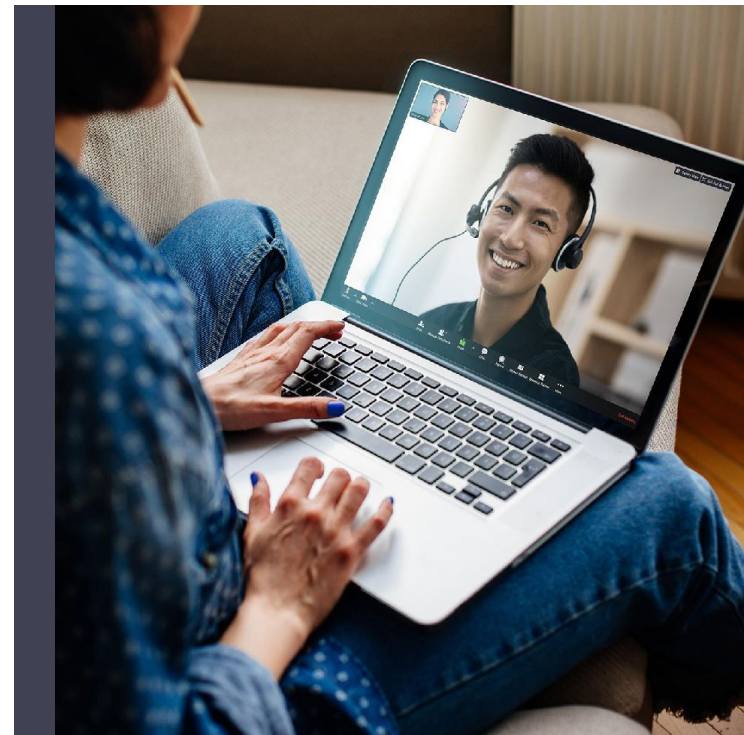imitations related to the use of non-GAAP metrics versus  their nearest GAAP equivalents. Other companies, including companies in our industry, may calculate non-GAAP metrics differently  or may use other measures to evaluate their performance, all of which could reduce the usefulness of our non-GAAP metrics as  tools for comparison. We urge you to review the reconciliation of Zoom's non-GAAP metrics to the most directly comparable GAAP  financial measures, and not to rely on any single financial measure to evaluate our business. See the Appendix for reconciliation  between each non-GAAP metric and the most comparable GAAP measure.

# Safe Harbor Statement

This presentation and the accompanying oral presentation have been prepared by Zoom Video Communications, Inc. ("Zoom") for informational purposes only and not for any other purpose. Nothing contained in this presentation is, or should be construed as, a recommendation, promise or representation by the presenter or Zoom or any officer, director, employee, agent or advisor of Zoom. This presentation does not purport to be all-inclusive or to contain all of the information you may desire.
Information provided in this presentation and the accompanying oral presentation speak only as of the date hereof.
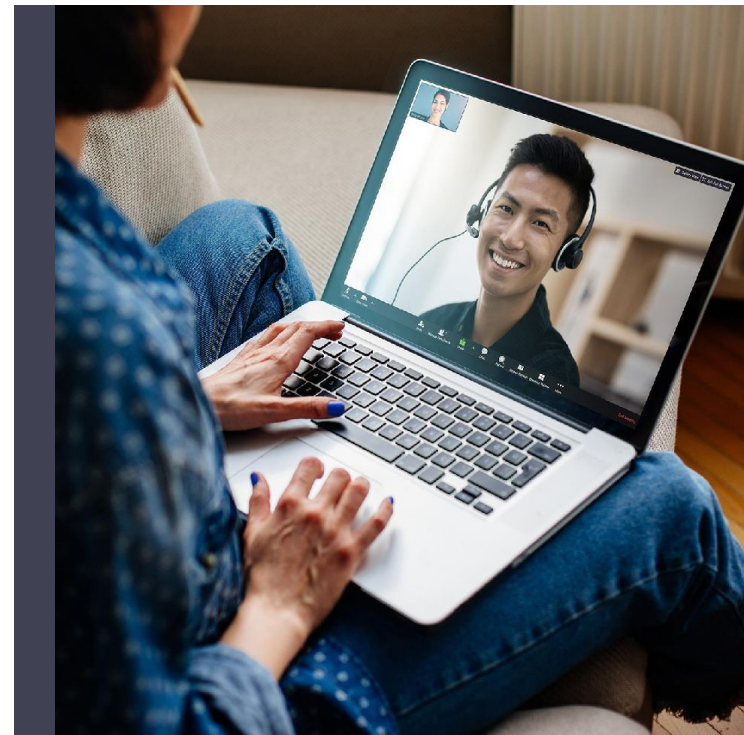
This presentation and the accompanying oral presentation include express and implied "forward-looking statements" within the meaning of the Private Securities Litigation Reform Act of 1995. In some cases, you can identify forward-looking statements by terms such as "anticipate," "believe," "estimate," "expect," "intend," "may," "might," "plan," "project," "will," "would," "should," "could," "can," "predict," "potential," "target," "explore," "continue," or the negative of these terms, and similar expressions intended to identify forward-looking statements. However, not all forward-looking statements contain these identifying words. These statements may relate to our market size and growth strategy, our reputation in the market, our estimated and projected costs, margins, revenue, expenditures, investments, and growth rates, as well as trends regarding the same, our future results of operations or financial condition, our plans and objectives for future operations, growth initiatives, or strategies and the impact to our business from the COVD-19 pandemic. By their nature, these statements are subject to numerous uncertainties and risks, including factors beyond our control, that could cause actual results, performance or achievement to differ materially and adversely from those anticipated or implied in the statements. These assumptions, uncertainties and risks include that, among others, our business would be harmed by any decline in new customers and hosts, renewals or upgrades, our limited operating history makes it difficult to evaluate our prospects and future results of operations, we operate in competitive markets, we do not expect to sustain our revenue growth rate in the future, there is continued uncertainty regarding the extent and duration of the COVID-19 and the responses of government and private industry thereto, as well as the impact of COVID-19 on the overall economic environment, any or all of which will have an impact on demand for remote work solutions for business as well as overall distributed face-to-face interactions and collaboration using Zoom, our business would be harmed by any significant interruptions, delays or outages in services from our co-located data centers, and failures in internet infrastructure or interference with broadband access could cause current or potential users to believe that our systems are unreliable. Additional risks and uncertainties that could cause actual outcomes and results to differ materially from those contemplated by the forward-looking statements are included under the caption "Risk Factors" and elsewhere in our most recent filings with the Securities and Exchange Commission (the "SEC"), including our annual report on Form 10-K for the fiscal year ended January 31, 2020. Forward-looking statements speak only as of the date the statements are made and are based on information available to Zoom at the time those statements are made and/or management's good faith belief as of that time with respect to future events. Zoom assumes no obligation to update forward-looking statements to reflect events or circumstances after the date they were made, except as required by law.

This presentation and the accompanying oral presentation also contain estimates and other statistical data made by independent parties and by us relating to market size and growth and other data about our industry. This data involves a number of assumptions and limitations, and you are cautioned not to give undue weight to such estimates. In addition, projections, assumptions, and estimates of our future performance and the future performance of the markets in which we compete are necessarily subject to a high degree of uncertainty and risk.

# Zoom Video Communication

Q1 FY21
Earnings  June 2, 2020

# Meeting the Increased Demand

**354%**

Growth in Customers with > 10 employees

**175K**

Number of licenses deployed for a new customer

**+200%**

Q/Q growth of minutes for the Global 2000 customers

**+300M**

Peak number of daily meeting participants

**+Two Trillion**

Annualized meeting minutes run rate

# Challenges and Commitment to Security and Privacy

✓ Enacted 90-day security plan initiative

✓ Acquired Keybase to add engineering expertise in encryption

✓ Released Zoom 5.0 with new security features and enhancements including support for AES 256-bit GCM encryption

**Happy Zoom Customers**

# Rapid Revenue Growth

In Millions



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $75 | $90 | $106 | $122 | $146 | $167 | $188 | $328 |
| Q2 FY19 | Q3 FY19 | Q4 FY19 | Q1 FY20 | Q2 FY20 | Q3 FY20 | Q4 FY20 | Q1 FY21 |

169% Yr/Yr

| $122 | $328 |
|---|---|
| Q1 FY20 | Q1 FY21 |

# Zooming ahead

**In the pipeline**

1. Zoom

2. has_zoom
3. Host
4. Join
5. Participant
6. Class
7. Virtual
8. Breakout
9. Room
10. Id

**Keep an eye on**

## Discord

a VoIP, instant messaging and digital distribution platform

# Recommendation for Digital Marketing Team

## Top words for Zoom

1. Zoom
2. has_zoom
3. Host
4. Join
5. Participant
6. Class
7. Virtual
8. Breakout
9. Room
10. Id

## Top words for MST

1. Team
2. has_mst
3. Microsoft
4. Channel
5. User
6. Assignment
7. Call
8. Notification
9. Chat
10. Feature

# Zooming ahead

**Refining our current model**

**Training our model to recognise words unique to subreddit.**

**Running the refined model on other competitor pairing**

**Zoom vs Google, Zoom vs Skype**

# Recommendation for Software Development Team

**Growth in Customers[1] with >$100K ARR**

769

**90% Yr/Yr**

405

**Q1 FY20**

**Q1 FY21**

**+500**

New customers[1] with >$100K ARR in Q1'21 from Q4'20

[1]The number of customers are rounded down to the nearest hundred

# Gaining Enterprise Traction

Growth in Customers with >$100K in Trailing 12-Month Revenue

**769**

**90% Yr/Yr**

**405**

Q1 FY20

Q1 FY21

Growth in Customers[1] with >$100K ARR

**+500**
New customers[1] with >$100K ARR in Q1'21 from Q4'20

[1]The number of customers are rounded down to the nearest hundred

# Rapidly Growing Customer Base

Customers[1] with more than 10 Employees



**354% Yr/Yr**

265.4K

58.5 K

Q1 FY20    Q1 FY21

Revenue from customers with 10 or fewer employees

**Q1'21: 30%**
**Q4'20: 20%**

[1]The number of customers are rounded down to the nearest hundred

# Strong Q1 Net Dollar Expansion Rate

TTM Net Dollar Expansion Rate[1]

**+130%**

in Q1 FY21

- 8th consecutive quarter above 130%

- Demonstrated Ability to Land and Expand

- Reflects Trust and Loyalty with Existing Customer

# Growing International Presence



Revenue[1]
(in millions)

Yr/Yr

ROW 246%

Americas 150%

Q1 FY20: AMER $98.2, EMEA $13.4, APAC $10.4
Q1 FY21: AMER $245.6, EMEA $51.3, APAC $31.3

AMER    EMEA    APAC

## Revenue Share by Region[1]

Q1 FY20: AMER 80.5%, EMEA 11.0%, APAC 8.6%
Q1 FY21: AMER 74.9%, EMEA 15.6%, APAC 9.5%

AMER    EMEA    APAC

[1]Subtotal revenue has been rounded

# Q1 FY21 Expenses and Margins

| | Q1 FY21 | | | |
|---|---|---|---|---|
| | GAAP Results | Yr/Yr | Non-GAAP[1] Results | Yr/Yr |
| Revenue | $328 million | 169% | $328 million | 169% |
| Gross Margin | 68.4% | (1,184bps) | 69.4% | (1,149bps) |
| Research & Development | 8.0% | (326bps) | 6.4% | (395bps) |
| Sales & Marketing | 37.0% | (1,546bps) | 31.5% | (1,880bps) |
| General & Administrative | 16.2% | +102bps | 14.8% | +134bps |
| Operating Margin | 7.1% | +585bps | 16.6% | +991bps |

[1]Note - A reconciliation of non-GAAP guidance measures to corresponding GAAP measures is not available on a forward-looking basis without unreasonable effort due to the uncertainty of expenses that may be incurred in the future

# Growing Future Revenue Under Contract

Total RPO[1]
(in millions)



Left chart — Q1 FY20 / Q1 FY21, Yr/Yr:
- Q1 FY20: $377 total ($149 Deferred Revenue, $227 Unbilled)
- Q1 FY21: $1,068 total ($552 Deferred Revenue, $516 Unbilled)
- Yr/Yr: 127% (Unbilled), 270% (Deferred Revenue)

Legend: Deferred Revenue, Unbilled

Right chart — Q1 FY20 / Q1 FY21, Yr/Yr:
- Q1 FY20: $377 total ($240 Current RPO, $137 Non-Current RPO)
- Q1 FY21: $1,068 total ($772 Current RPO, $296 Non-Current RPO)
- Yr/Yr: 116% (Non-Current RPO), 222% (Current RPO)

Legend: Current RPO, Non-Current RPO

[1] Remaining performance Obligations (RPO) consists of both billed considerations and unbilled considerations that we expect to recognize as revenue, which grew 184% year-over-year. We expect to recognize approximately 72% or $772 million dollars of the total RPO as revenue over the next 12 months compared to 64% or $240 million dollars in Q1 last year. Subtotals have been rounded.

14

# Rapid Cash Flow Growth

**Operating Cash Flow**
**(in millions)**

$259.0

$22.2

1,065%
Yr/Yr

Q1 FY20          Q1 FY21

**Free Cash Flow**
**(in millions)** 1

$251.7

1,541%
Yr/Yr

15

# Full Year and Q2 FY21 Outlook

|  | **Q2FY21** | **FY21** |
|---|---|---|
| Revenue | $495 - $500 million | $1,775 - $1,800 million |
| Non-GAAP Operating Income | $130 - $135 million | $355 - $380 million |
| Weighted Average Share Count | 299 million | 300 million |
| Non-GAAP EPS | $0.44 - $0.46 | $1.21 - $1.29 |

[1] A reconciliation of non-GAAP guidance measures to corresponding GAAP measures is not available on a forward-looking basis without unreasonable effort due to the uncertainty of expenses that may be incurred in the future .

# Questions

Thank you

# Appendix

# GAAP to Non-GAAP Reconciliation

**Gross Profit**

| ($ in thousands) | QTD – Q1FY20 | QTD – Q1FY21 |
|---|---|---|
| Total Revenue | $121,988 | $328,167 |
| GAAP Gross Profit | $97,884 | $224,460 |
| (+) Stock-based compensation expense and related payroll taxes | $830 | $3,382 |
| Non-GAAP Gross Profit | $98,714 | $227,842 |
| Non-GAAP Gross Margin | 80.9% | 69.4% |

**R&D Expenses**

| | QTD – Q1FY20 | QTD – Q1FY21 |
|---|---|---|
| GAAP R&D | $13,783 | $26,389 |
| (-) Stock-based compensation expense and related payroll taxes | $1,164 | $5,403 |
| Non-GAAP R&D | $12,619 | $20,986 |

**S&M Expenses**

| | QTD – Q1FY20 | QTD – Q1FY21 |
|---|---|---|
| GAAP S&M | $64,041 | $121,556 |
| (-) Stock-based compensation expense and related payroll taxes | $2,627 | $18,025 |
| Non-GAAP S&M | $61,414 | $103,531 |

**G&A Expenses**

| | QTD – Q1FY20 | QTD – Q1FY21 |
|---|---|---|
| GAAP G&A | $18,503 | $53,130 |
| (-) Stock-based compensation expense, related payroll taxes, and charitable donation of common stock | $2,041 | $4,436 |
| Non-GAAP G&A | $16,462 | $48,694 |

# GAAP to Non-GAAP Reconciliation

| ($ in thousands) | QTD – Q1FY20 | QTD – Q1FY21 |
|---|---|---|
| Total revenue | $121,998 | $328,167 |
| GAAP operating profit | $1,557 | $23,385 |
| (+) Stock-based compensation expense, related payroll taxes, and charitable donation of common stock | $6,662 | $31,246 |
| Non-GAAP operating profit | $8,219 | $54,631 |
| Non-GAAP operating margin | 6.7% | 16.6% |
| | | |
| **Net Income** | | |
| GAAP net income attributable to common stockholders | $198 | $27,036 |
| (+) Stock-based compensation expense, related payroll taxes, and charitable donation of common stock | $6,662 | $31,246 |
| (+) Undistributed earnings attributable to participating securities | $2,016 | $39 |
| Non-GAAP net income | $8,876 | $58,321 |
| | | |
| **Earnings Per Share** | | |
| GAAP net income per share – diluted | $0.00 | $0.09 |
| Non-GAAP net income per share – diluted | $0.03 | $0.20 |
| | | |
| **Weighted Average Shares** | | |
| GAAP weighted-average – diluted | 136M | 295M |
| Non-GAAP weighted-average - diluted | 290M | 295M |

# Historic Metrics

| Metric | Q2 FY19 | Q3 FY19 | Q4 FY19 | Q1 FY20 | Q2 FY20 | Q3 FY20 | Q4 FY20 | Q1 FY21 |
|---|---|---|---|---|---|---|---|---|
| Revenue | $74.5 | $90.1 | $105.8 | $122.0 | $145.8 | $166.6 | $188.3 | $328.2 |
| y/y | 126% | 120% | 108% | 103% | 96% | 85% | 78% | 169% |
| GAAP Operating Income | $3.4 | $(1.1) | $5.5 | $1.6 | $2.3 | $(1.7) | $10.6 | $23.4 |
| Stock-based compensation expense, related payroll taxes, & charitable donation of common stock | $1.1 | $2.7 | $4.3 | $6.7 | $18.5 | $22.9 | $27.9 | $31.2 |
| Non-GAAP Operating Income | $4.5 | $1.6 | $9.8 | $8.2 | $20.7 | $21.3 | $38.4 | $54.6 |
| Operating Cash Flow | $14.4 | $18.2 | $16.0 | $22.2 | $31.2 | $61.9 | $36.6 | $259.0 |
| Capital Expenditures (Property & Equipment) | $(6.2) | $(8.1) | $(10.3) | $(6.9) | $(14.0) | $(7.2) | $(10.0) | $(7.3) |
| Free Cash Flow | $8.2 | $10.1 | $5.7 | $15.3 | $17.1 | $54.7 | $26.6 | $251.7 |
| RPO | $210.5 | $256.0 | $311.7 | $376.5 | $457.6 | $517.0 | $604.1 | $1067.9 |
| y/y | n/a | n/a | n/a | 127% | 117% | 102% | 94% | 184% |
| TTM Net $ Expansion Rate | 138% | 139% | 140% | 130%+ | 130%+ | 130%+ | 130%+ | 130%+ |
| Customers >10 Employees | 37.2k | 44.4k | 50.8k | 58.5k | 66.3k | 74.1k | 81.9k | 265.4k |
| y/y | n/a | n/a | 97% | 86% | 78% | 67% | 61% | 354% |
| Customers >$100K TTM Revenue | 228 | 277 | 344 | 405 | 466 | 546 | 641 | 769 |
| y/y | n/a | n/a | 141% | 120% | 104% | 97% | 86% | 90% |