



∞

# 1800-READMES: A META RESEARCH DIVE

---

By Brock Green, Jeremiah Toribio, Ramiro Lopez

# Overview

---

**EXECUTIVE SUMMARY**

**FINDINGS**

**RECOMMENDATIONS**

**Conclusion**

# EXECUTIVE SUMMARY

## GOAL

Create a model that can predict the main programming language of a repository

# EXECUTIVE SUMMARY

## GOAL

Create a model that can predict the main programming language of a repository

## BIG IDEA

Given the text of a repo's ReadMe file, one can use NLP to accurately predict a repositories programming language

# EXECUTIVE SUMMARY

## GOAL

Create a model that can predict the main programming language of a repository

## BIG IDEA

Given the text of a repo's ReadMe file, one can use NLP to accurately predict a repositories programming language

## FINDINGS

- Unigrams
- Bigrams
- Trigrams

# EXECUTIVE SUMMARY

## GOAL

Create a model that can predict the main programming language of a repository

## BIG IDEA

Given the text of a repo's ReadMe file, one can use NLP to accurately predict a repositories programming language

## FINDINGS

- Unigrams
- Bigrams
- Trigrams

## RECOMMENDATIONS

- Run this model to aid recruiters
- Compare to other departments

**82%**

**REPOS ARE PYTHON**

---

**18%**

**OTHER REPOS**

---

*C++, Markdown and other*

# Unigrams

## Most common words by category

### ■ Markdown ■

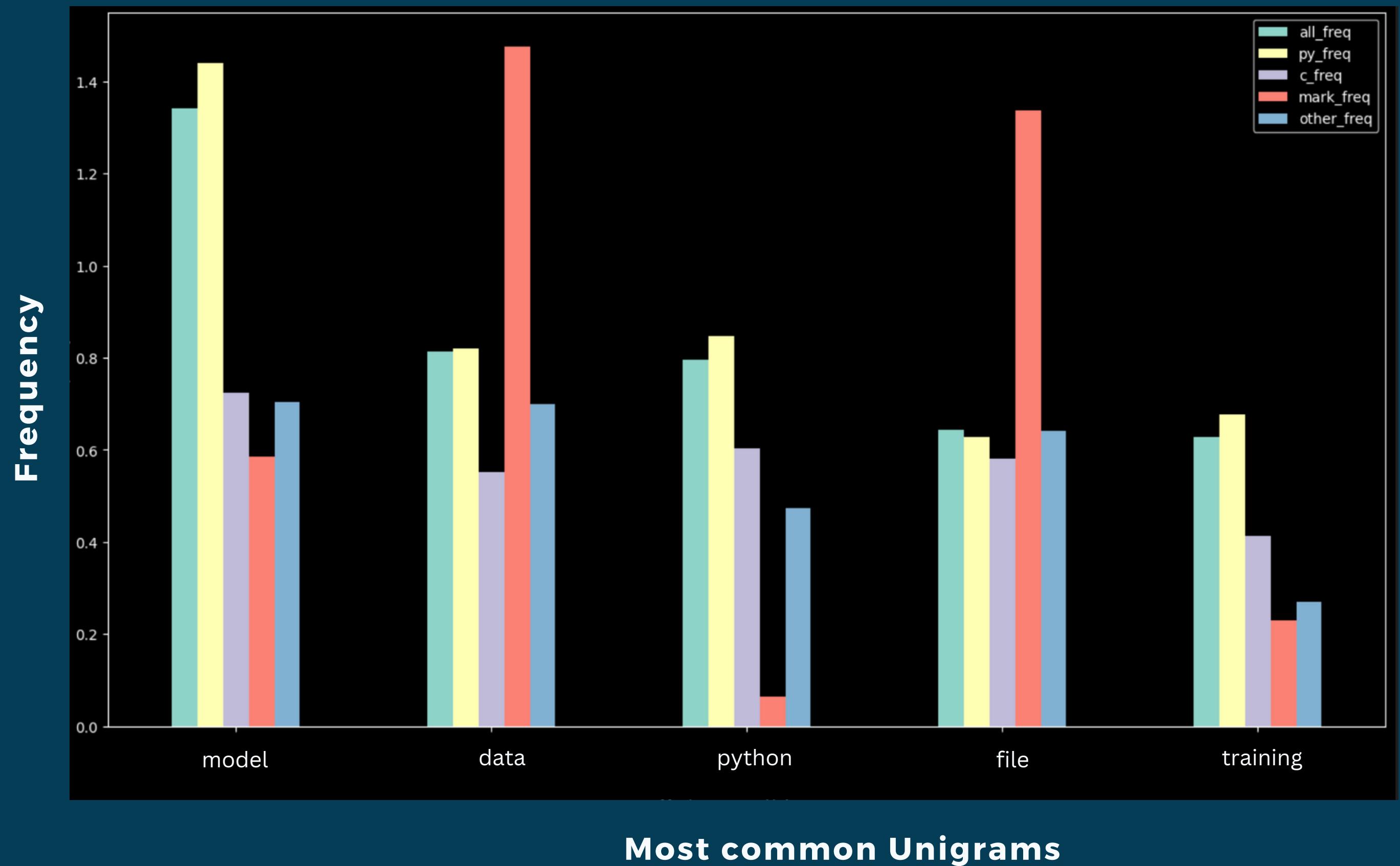
- Dataset
- license
- file
- data

### ■ C++ ■

- Example
- install

### ■ Python ■

- Model
- python



# Bigrams

## Most common bigrams by category

### ■ Markdown ■

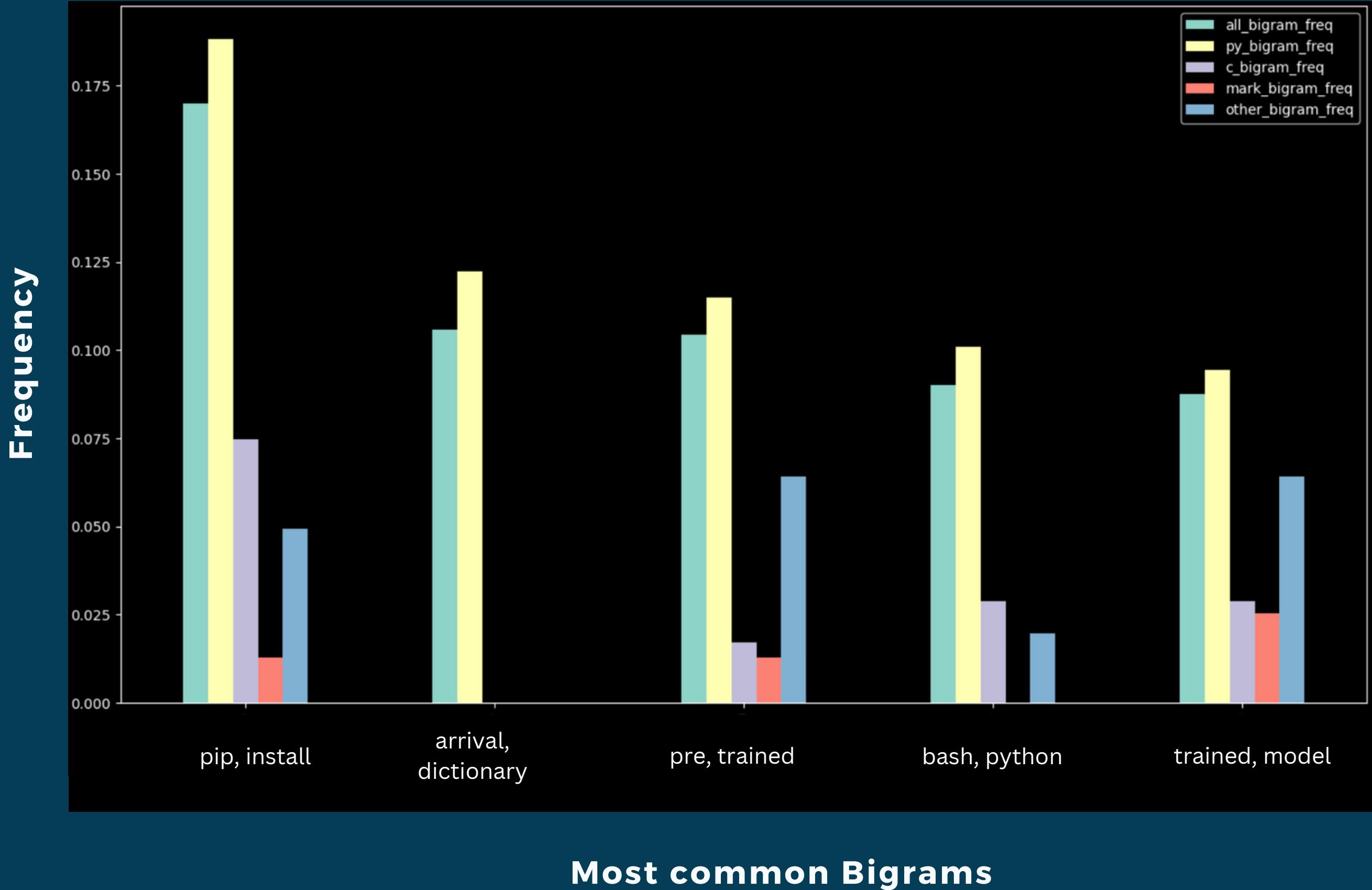
- License, File
- Blob, Main
- cc, nc

### ■ C++ ■

- Git, Clone
- Following, command

### ■ Python ■

- Pip, Install
- Pre, Trained
- Trained, Model



# Trigrams

## Most common trigram by category

### ■ Markdown

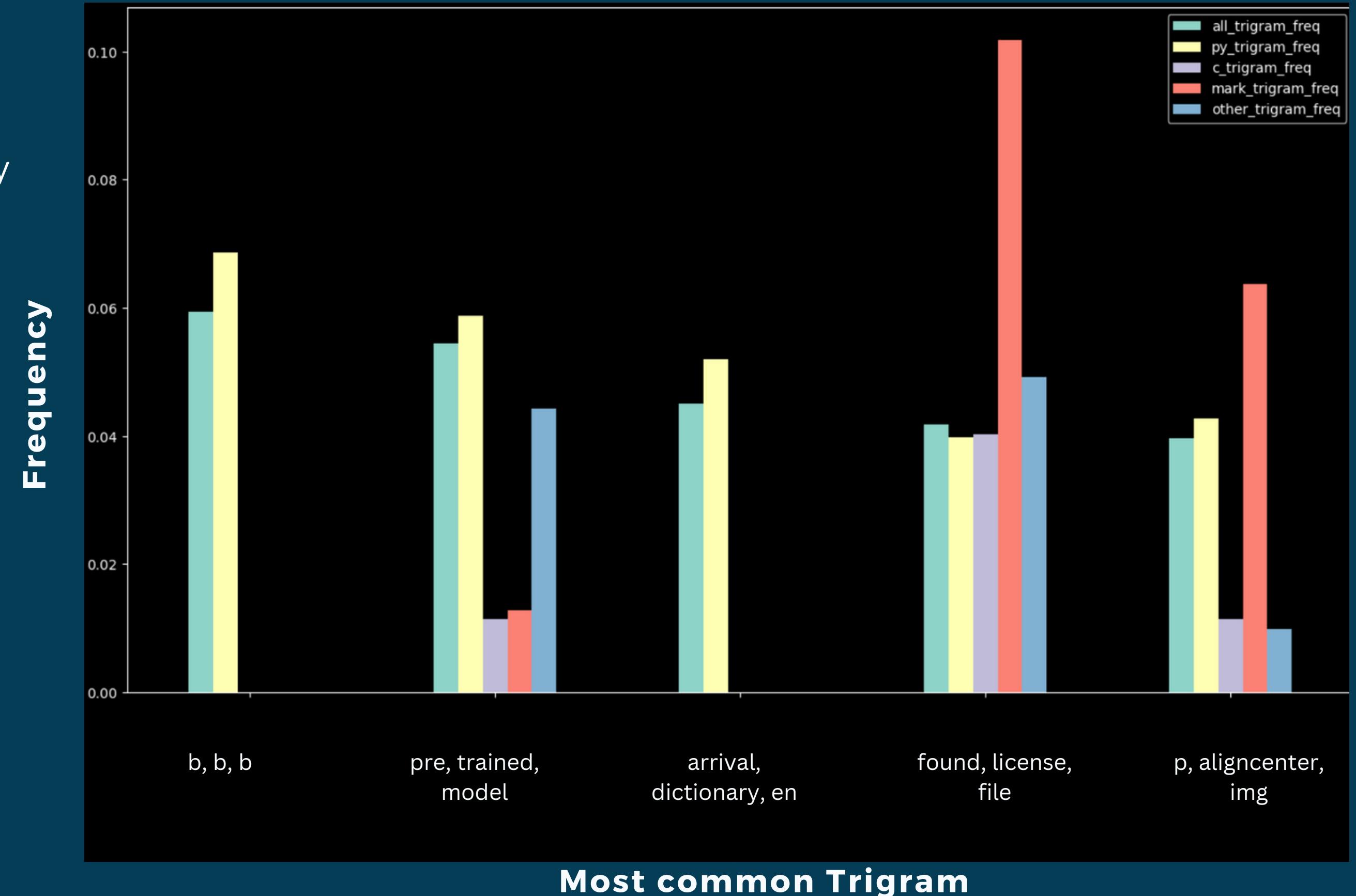
- found, license, file
- testhttps, arrival, dictionary
- p, aligncenter, img

### ■ C++

- Git, Clone, http

### ■ Python

- pre, trained, model



# Modeling

---

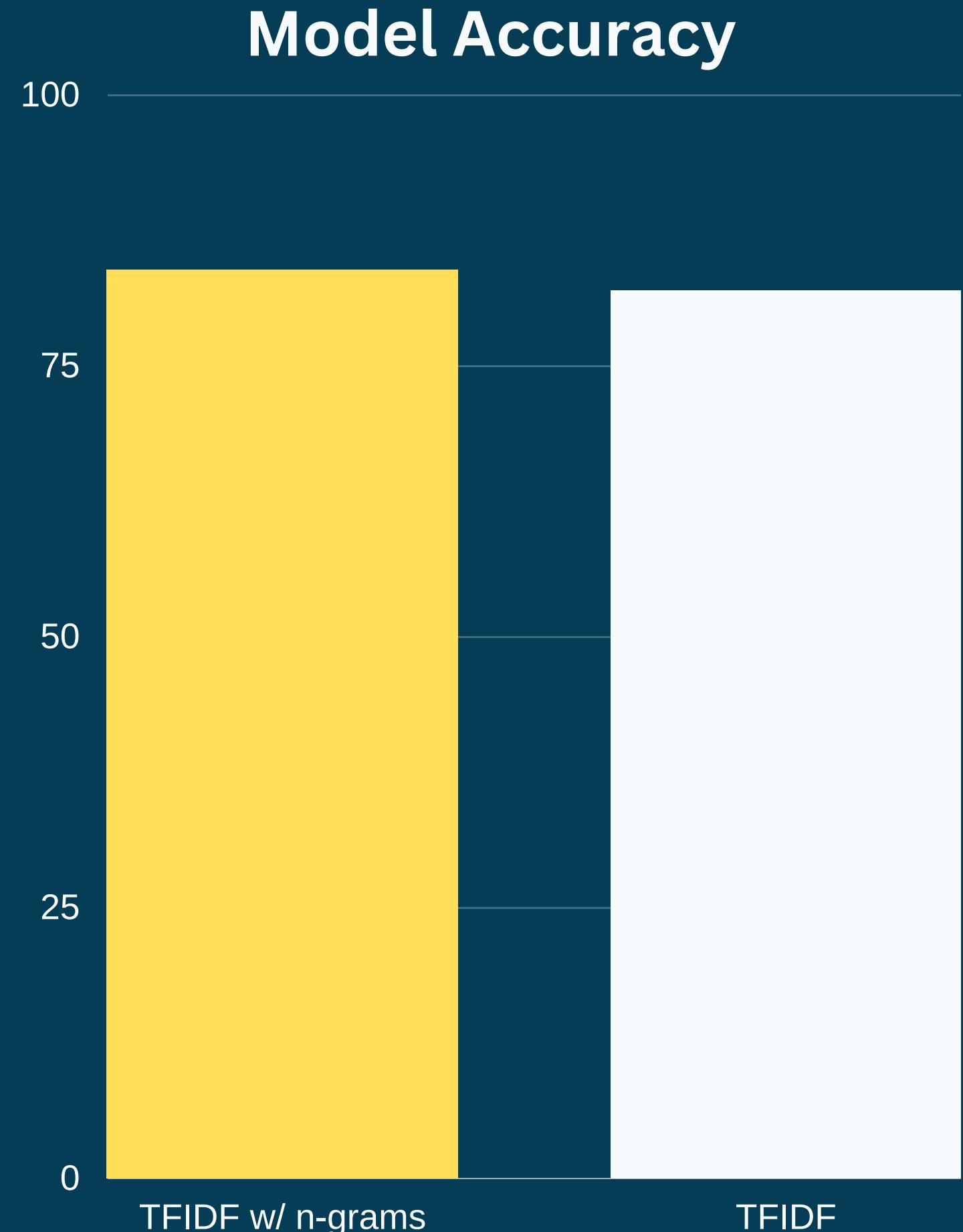
**With and without n-grams**

**Decision Tree** of *max\_depth*: 2

Baseline - 82%

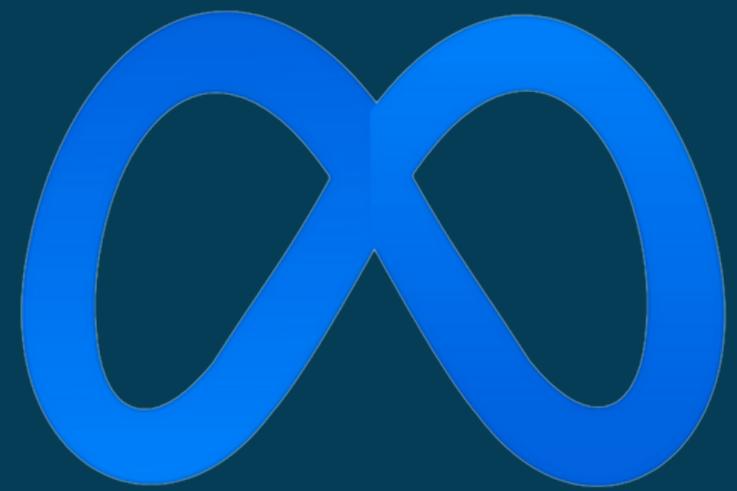
Train - 84.39%

Test - 83.75%



# Recommendations

---



- Assisting with team manager placements
- Potentially scan an applicants GitHub to see candidate likelyhood for Research team

# Conclusion

## MODEL ACCURACY

- +2% Better than baseline

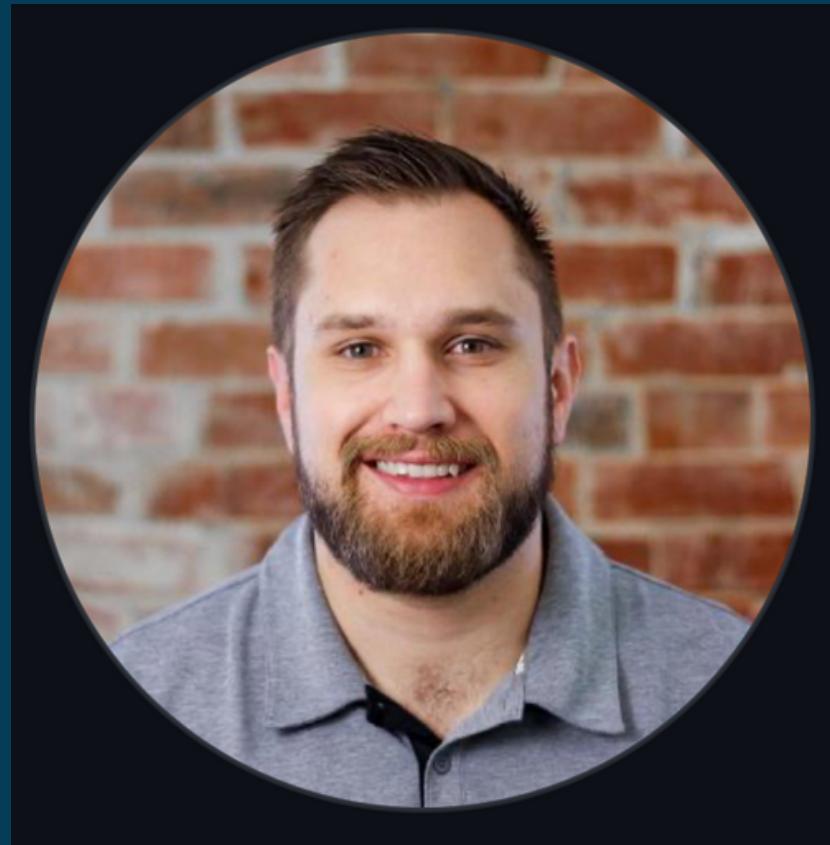
## TAKEAWAYS

- Research teams primarily use Python
- Potential candidates can be vetted on their knowledge and experience with respective top languages

## NEXT STEPS

- Add more repositories from different departments to accurately predict for any department
- Prescriptive Model: Identify the features that would be easiest for the business to take action

# THANK YOU FOR YOUR TIME



BROCK GREEN

Data Scientist

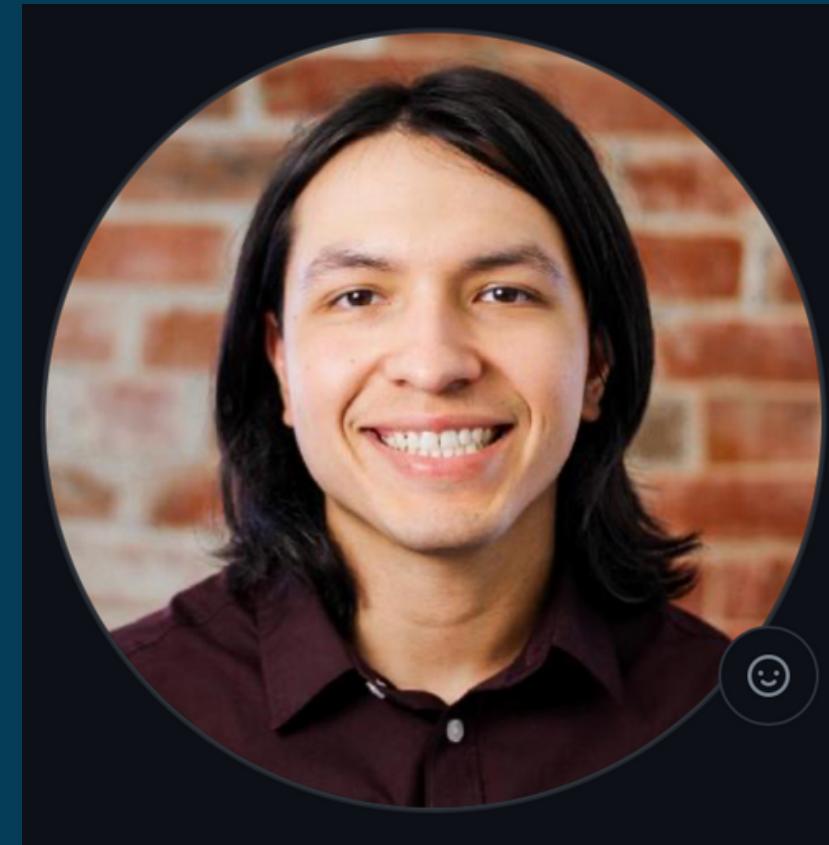
<https://github.com/brock-green>



JEREMIAH TORIBIO

Data Scientist

<https://github.com/jeremiah-toribio>



RAMIRO LOPEZ

Data Scientist

<https://github.com/Ramiro-Lopez>