

# Data Collection & Processing for AI



## Avi Gopal

Co-founder & CTO @ Metabob  
CTO @ Mahn Arc Ventures  
Co-Founder @ Clyste

"Avi has an illustrious background in Aerospace Engineering, with expertise in every breadth of development, specifically versed within AI and machine learning"

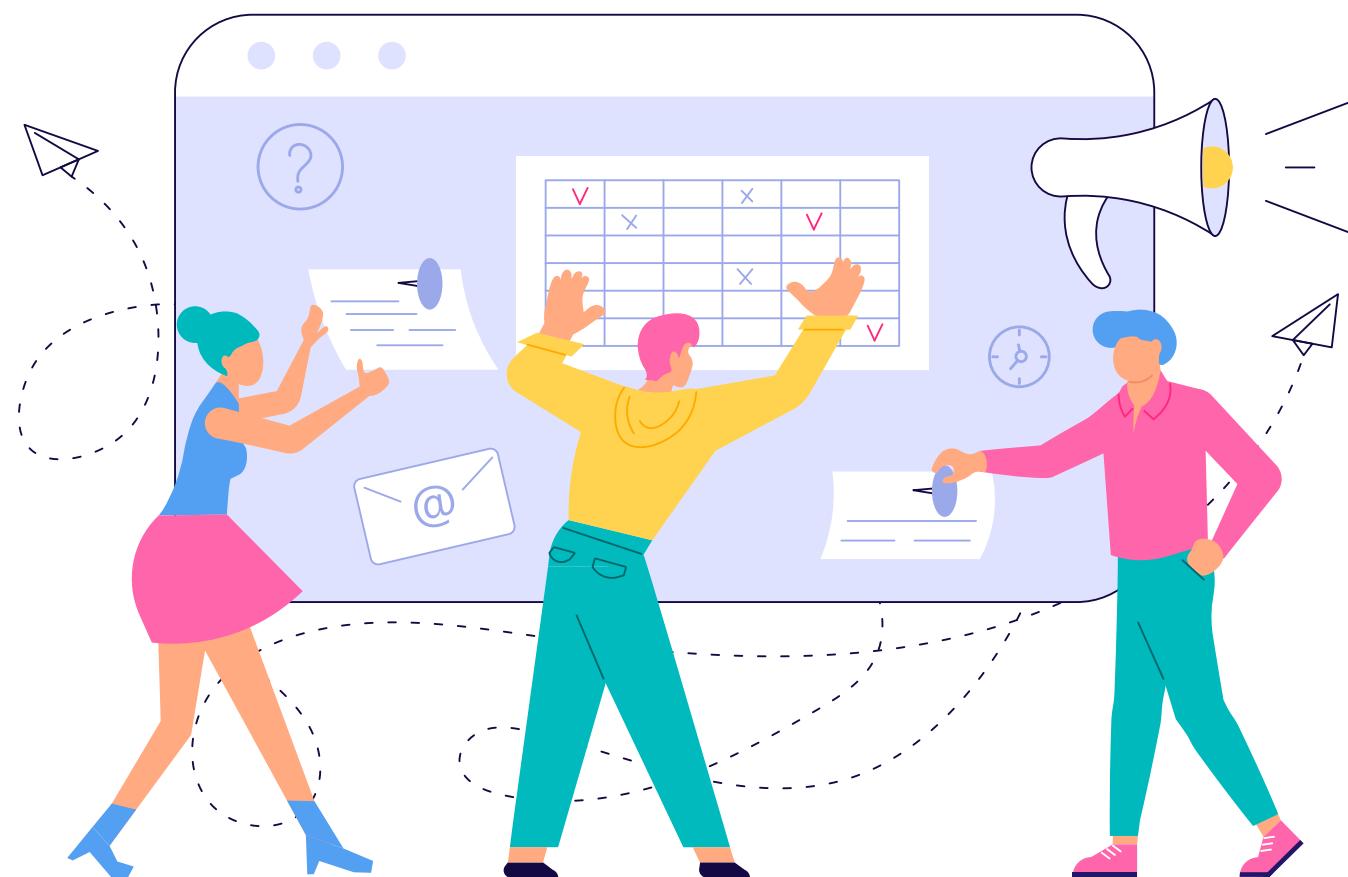
Chapter04  
Session07

Step Up. Speak Up. Go Up

June 27th 10AM PST  
[zoom.goupaz.com](https://zoom.goupaz.com)

Hosted by  
 Metabob

# Agenda



- **10:00-10:10 | Kickoff**
  - Introduction
  - Chit Chat
- **10:10-11:10 | Presentation**
  - Speaker
- **11:10-11:15 | Break**
  - Pre Q&A
- **11:15-11:30 | Q&A**
  - Discussion
- **11:30-11:35 | Wrap Up**
  - Contact Speaker

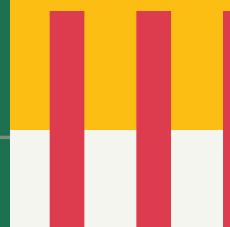
Hosted by

Axel L

Growth Marketing



*Passionate about sports and technologies enhancing employee efficiency*



## Metabob

It's the fast, easy, and visual way of debugging code.

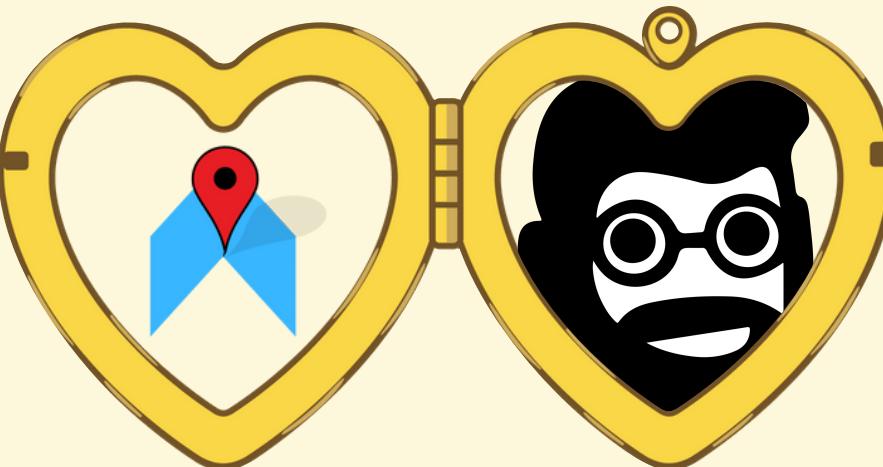
<https://metabob.com>

## GOUP

Community driven Open source accelerator!  
<https://goupaz.com>

# Chapter04

#contributor



# GOUP

Community driven Open source accelerator!  
<https://goupaz.com>



YOU

GOUPO2hero/Chapter04  
Contributor  
About YOU

Step Up. Speak Up. Go Up.

# Metabob

It's the fast, easy, and visual way of  
debugging code.

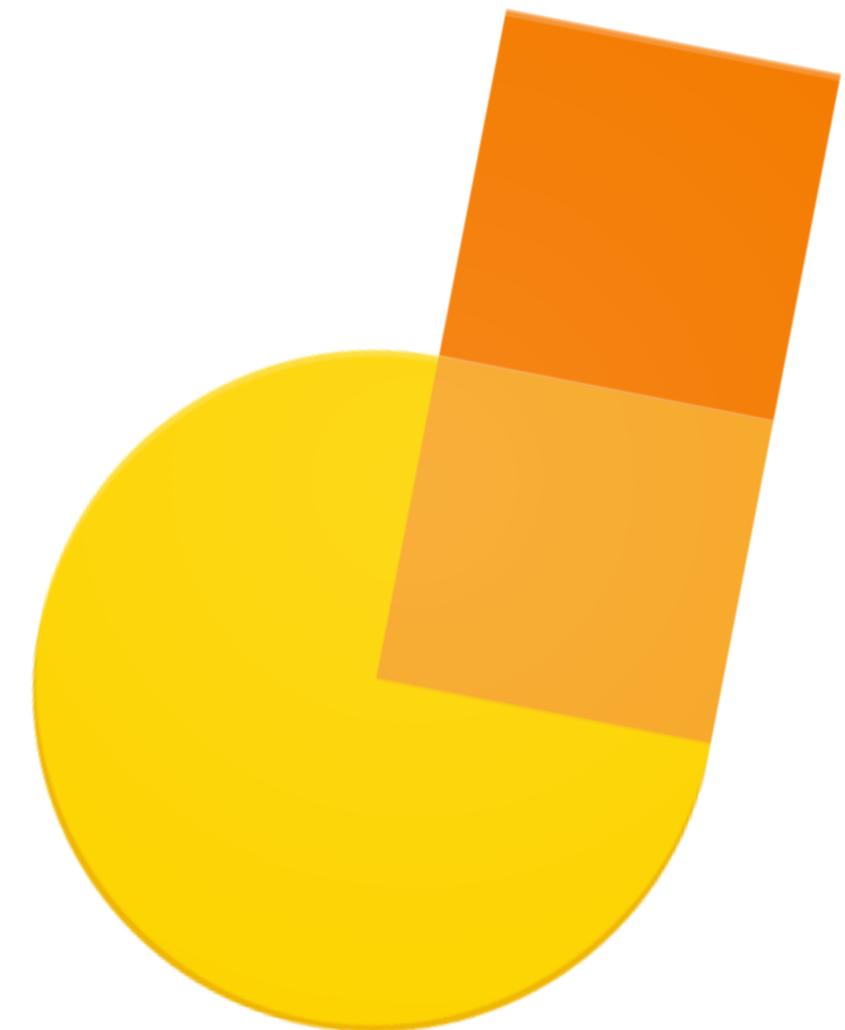
<https://metabob.com>

Are you ready?

Let's Begin!

CHITCHAT

# JAMBOARD



# What do you wish to learn from this webinar?

I want learn more about the importance of collecting data for AI

**As A Beginner How Can I Learn AI - Akhil**

I want to learn how are data and AI (which I believe are mostly algorithms) are related.



Fiza: I want to learn how can we collect effective data for AI.

**I am Chandan Kumar, a student. How to get started with AI**

POLL

## HOW EXPERIENCED ARE YOU WITH AI DEVELOPMENT?

1



Not experienced

2



Somewhat experienced

3



Semi-experienced

4



Experienced

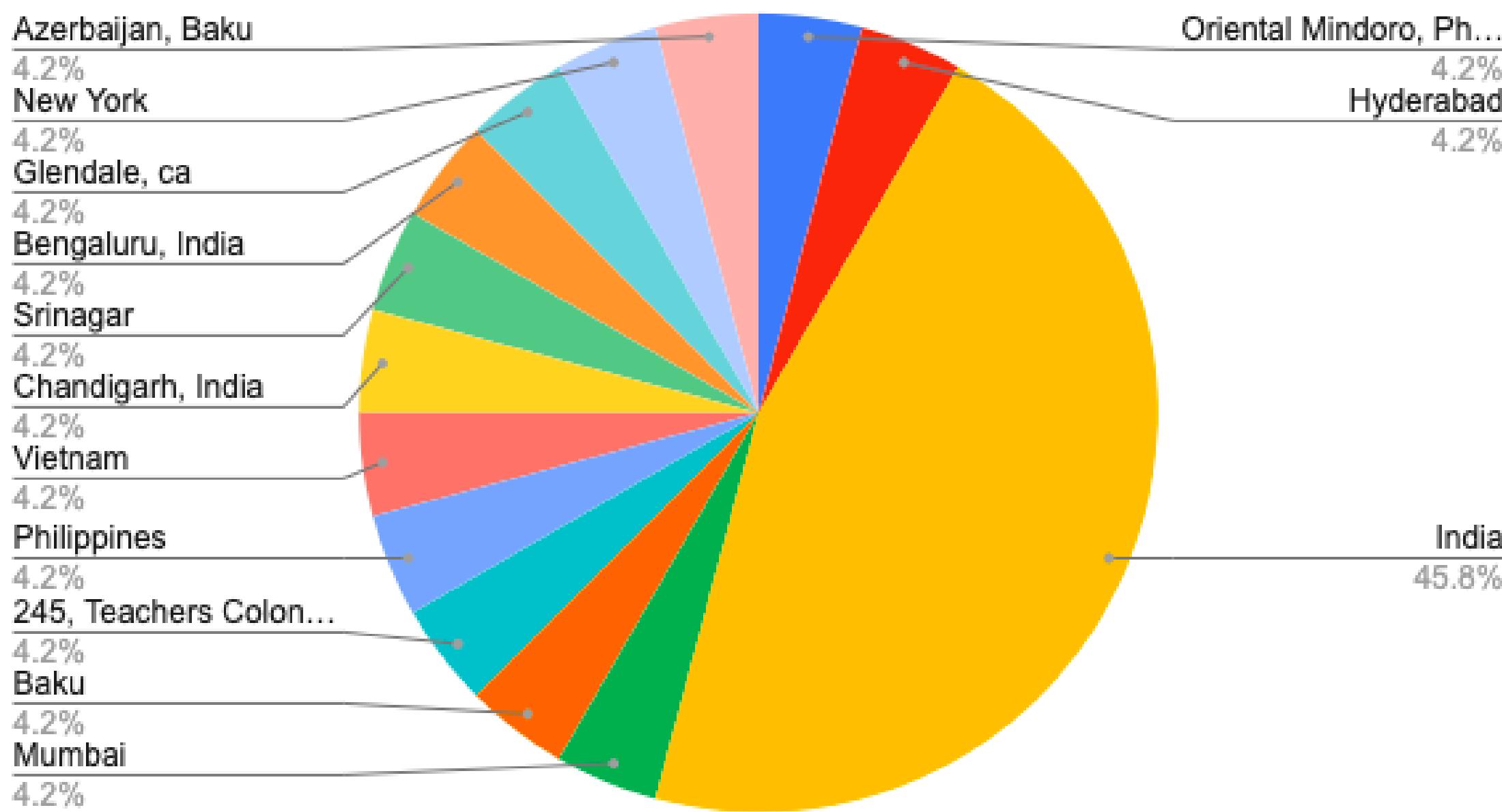
5



Very experienced

# Audience

Count of Location



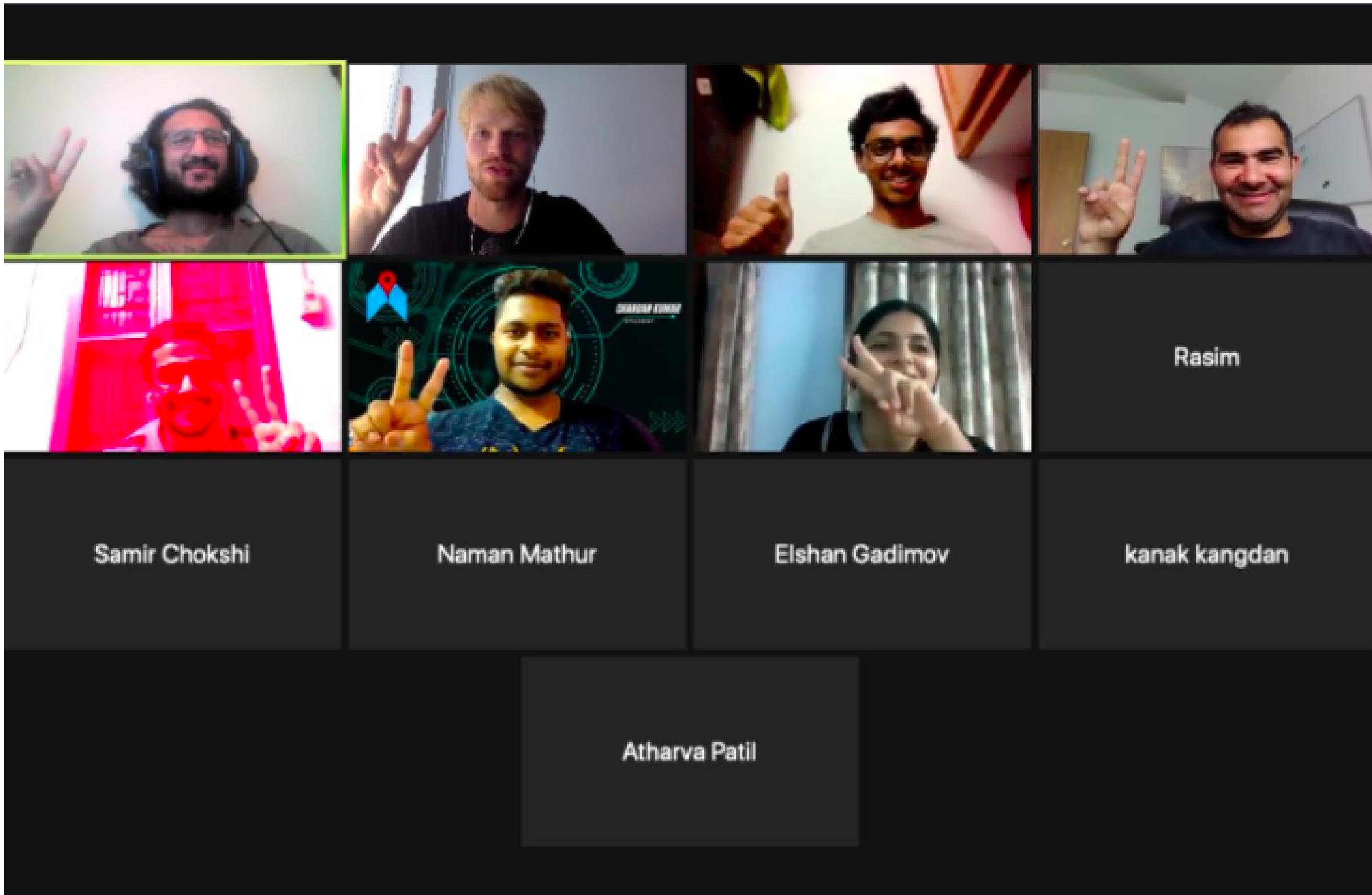
# Code of Conduct

- 1 Learn, benefit, contribute
- 2 No marketing, selling, competing
- 3 Equality despite roles & bg



# Photo Shoot Time

Please turn on your camera :D



# Presentation

# Data Collection and Processing

# Avinash G.



[Avi@metabob.com](mailto:Avi@metabob.com)

- Worked on various Aerodynamics Research Projects, specializing in novel airfoil design, during his time at CSULB as an aerospace engineer
- Created a variety of autonomously controlled craft for personal transport and area mapping
- Created new organization structures and community building applications for Clyste as a CTO
- Developed the entirely of the backend and system architecture for the code analysis platform Metabob owned by NEC

# Overview

What is it?

Why do it?

Defining the Problem

# Data Processing, Analytics & BI | Functional Viewpoint

Acquiring the Incoming data from various Sources

## Data Ingestion

Bounded Data Sources  
*(Batch Sources)*

Unbounded Data Sources  
*(Streaming Sources)*

Performing ETL/ELT for Data Processing

## Data Processing (ETL)

Batch Process ETL  
(Data Quality, Data Transformation, Aggregation)

Streaming ETL  
(Data Quality, Data Transformation)

Performing advanced Analytics / Modeling

## Data Processing (Analytics)

Analytics  
(Data Cleansing, Data Munging, Outlier Analysis, Clustering, Projection & Prediction)

Streaming Analytics  
(Run against Model as service call)  
Ex: Credit Card fraud validation

Enabling consumption for various Downstream

## Data Consumption

Reports  
(Canned Reports, Interactive Reports)

Dashboards  
(Multi-dimensional Views, Real-time refresh, Drill down)

Extracts  
(File extracts, Cubes extracts, DB Dump)

Adhoc  
(SQL Quires against the DB)

Data Storage into various Formats/Platforms

## Data Persistence

Relational

NoSQL

Hadoop

## Data Persistence

Relational

OLAP Cubes

Files

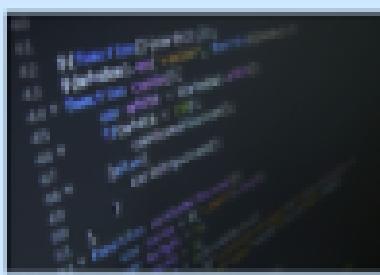
# Data Ingestion

- Do
  - Have a clear goal in mind
  - Focus on identifying potential data sources
  - Build the necessary integrations for each source
  - Keep track of the metadata as well
- Don't
  - Worry about scaling when just starting out
  - Apply any processing or filters
  - Delete anything

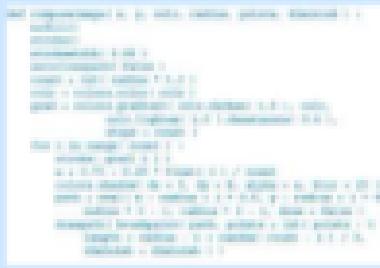
# What we did at Metabob



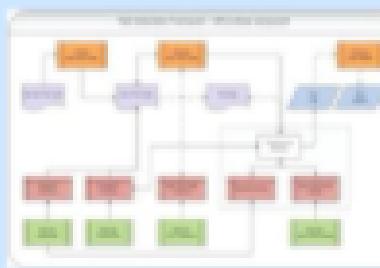
high-quality  
open source  
repos



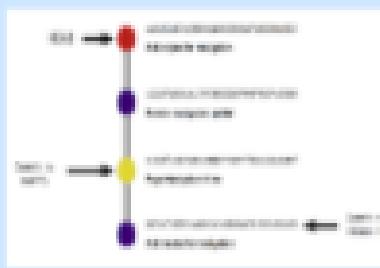
company  
SOPs



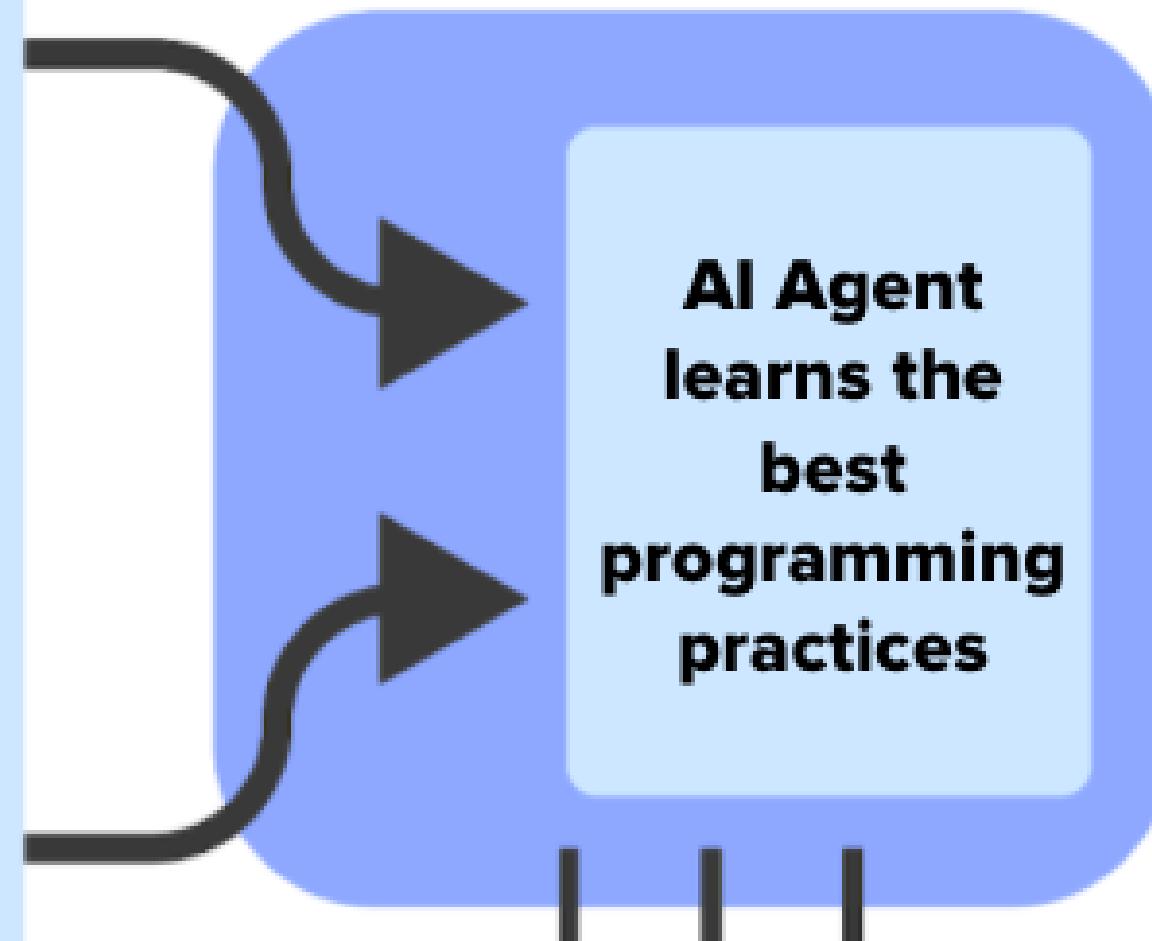
company  
design  
standard



program  
analysis



manual  
review  
historical  
results



A set of  
recommended  
code changes

AI agent continues  
updating and learning  
through interaction



# Data Processing

Only build data processing as needed

Focus initially on regularizing input and/or aggregating items from different sources

Store intermediary steps to save on processing time down the road

As capabilities develop move onto enriching sources as well

## Connectors

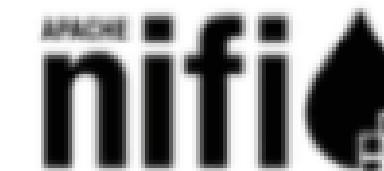


## Data Platform

ETL  
State



ETL  
Orchestration



Data  
Processing



Data  
Ingestion

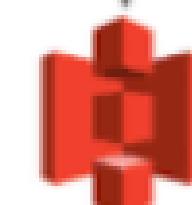


Data Mgt.  
& Optimization

**<code>**



Raw  
Data



Optimized  
Data



Glue  
Catalog

## Consumption



Athena



Redshift



RDS



Elastic



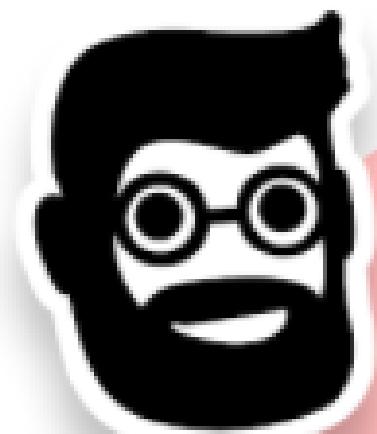
SageMaker



Qubole



Data  
Apps



## Topic Modeling

Data  
Gathering

Data  
Labeling

Topic  
Generation

### **Data Gathering**

Aggregated over 50 thousand human verified bug fixes corresponding to over 15 million code changes

### **Data Labeling**

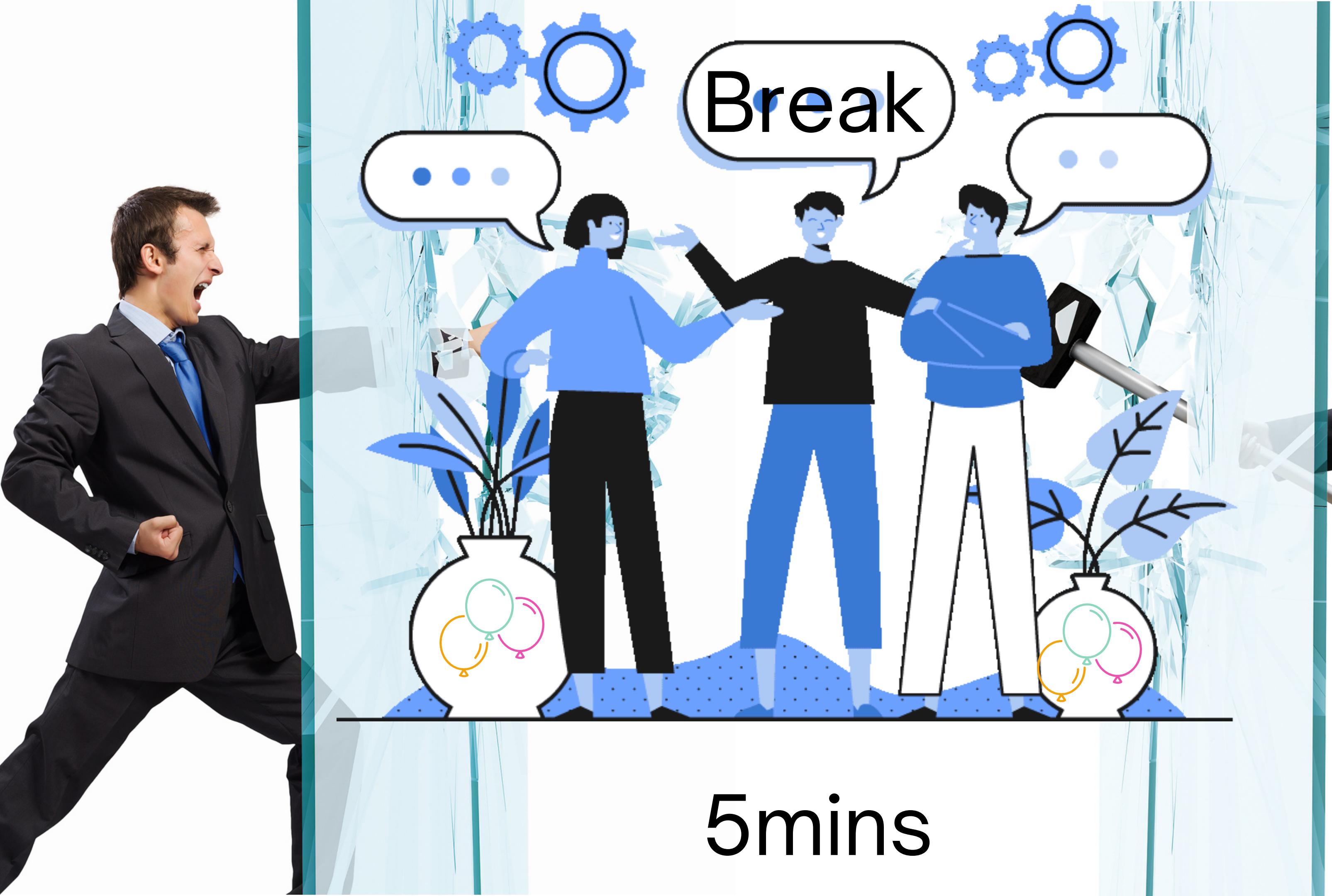
Generated a baseline understanding of the reasons behind code changes via human labelling

### **Topic Generation**

Sorted code changes into a set of topics relating why and how a change was made

# Q&A | Discussion

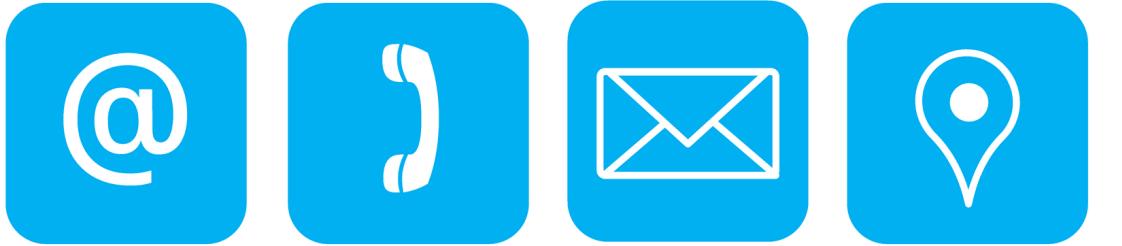




5mins

# Wrap Up

HOWTO



## Remember!

- Re-Search
- Be short and clear
- Re-mind
- Q&A over Slack

Linkedin: @avinash-gopal-440669140

GOUP Slack: @Avi CL



# Avi Gopal



**Join Us!**

**HTTPS://GOUPAZ.COM**

**HTTPS://METABOB.COM**

**1 Community Managers**

**2 Tech Writers**

**3 In/Out Ambassadors**

**4 Marketing Creators & Editors**

**5 Course Creators**

**6 Project Creators**

# Thank You

## Culture

#egoless #collaborative #competent #decentralized #scalable #fun

## Open source

#creator #contributor

## Diversity

#age #gender #location #economics #religion #politicalview

How can we do  
better?