# {ON: The Beach}
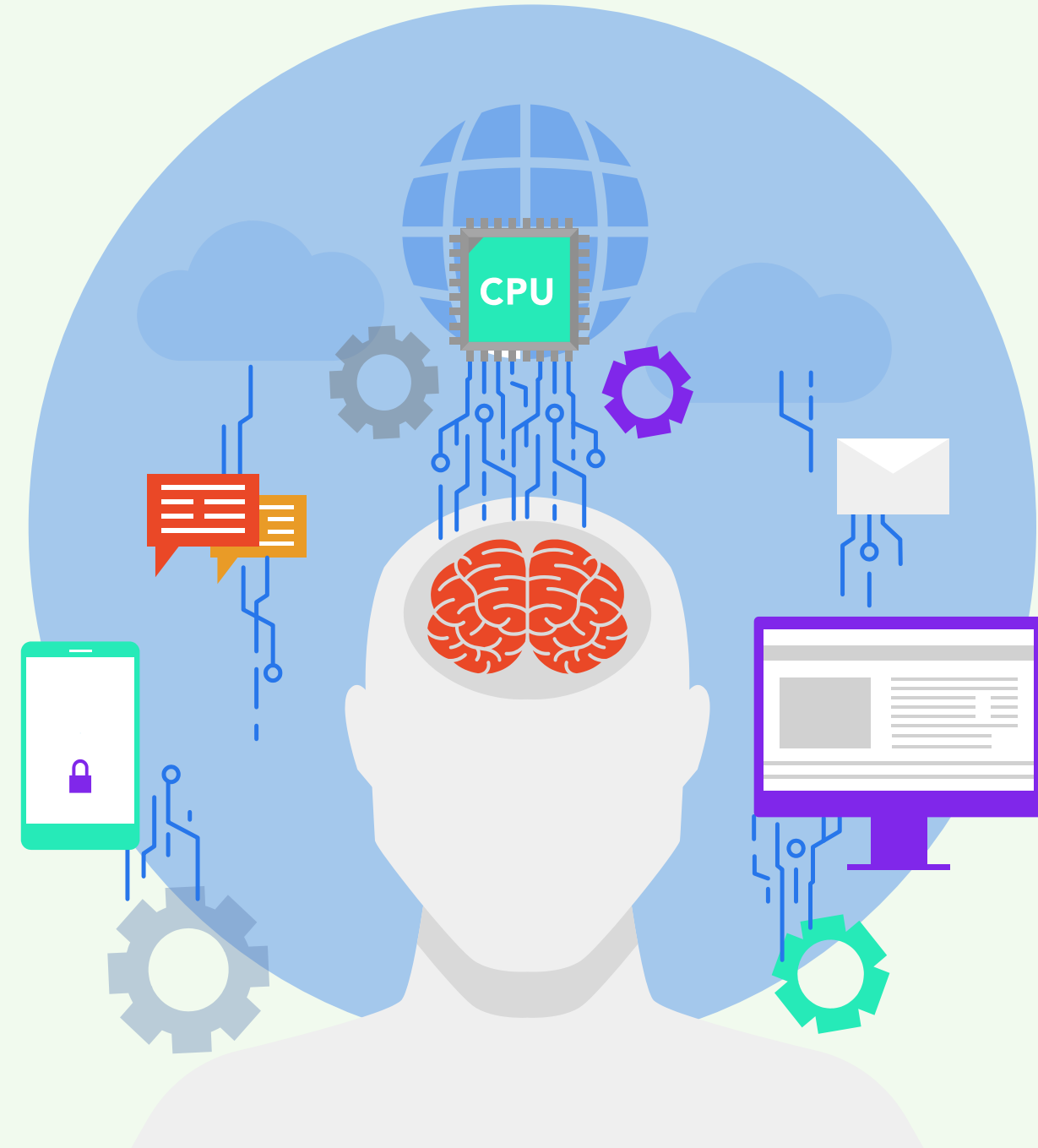
Bringing DevOps, Devs and Data Scientists
together around Big Data

# Unearth The Black-Box:

# Building Fair, Accountable & Trustworthy ML Systems

# AGENDA

**01** **INTRODUCTION**

WHAT IS ML AND RELATED CONCEPTS

**02** **UNDERSTANDING BLACK-BOX PROBLEM**

WHY BLACK-BOX MODELS CAN LEAD TO UNFAIR, BIASED OUTCOMES

**03** **STRATEGIES FOR BUILDING FAIR AND UNBIASED MODELS**

HOW TO BUILD STRATEGIES FOR ENSURING FAIRNESS AND NO-BIAS

**04** **BUILDING TRUSTWORTHY & ACCOUNTABLE MACHINE LEARNING MODELS**

WHAT STRATEGIES FOR BUILDING TRUST & ACCOUNTABILITY IN ML SYSTEMS

**05** **CONCLUSION AND RESOURCES**

ENCOURAGE FURTHER RESEARCH AND DISCUSSION

# INTRODUCTION

A B C

# BASIC NEURAL NETWORK

INPUT

WEIGHTS

SUM AND BIAS

ACTIVATION FUNCTION

OUTPUT

# IMAGE RECOGNITION

What the computer sees

image classification →

82% cat
15% dog
2% hat
1% mug

Img Src: KDnuggets Blog

# (Biased) FACE RECOGNITION

# (Biased) WORD EMBEDDINGS



Homemaker

Distance is more

Distance is less

Man

Woman

**The doctor yelled at the nurse because he was late. Who was late?**

In the sentence "The doctor yelled at the nurse because he was late," the pronoun "he" refers back to "the doctor." Therefore, the doctor was the one who was late.

**The doctor apologized to the nurse for her late arrival to the emergency room. Who was late previously?**

Based on the sentence you provided, it seems that the doctor was the one who was late to the emergency room, and that's why he apologized to the nurse.

The nurse apologized to the doctor for his late arrival to the operation room. Who was late?

Based on the sentence you provided, it seems that there is an error in the sentence construction that could cause confusion about who was actually late. The use of pronouns in this sentence makes it unclear who was late to the operation room. However, if we assume that the sentence was intended to mean that the nurse was apologizing to the doctor for her own late arrival to the operation room, then the nurse was the one who was late.
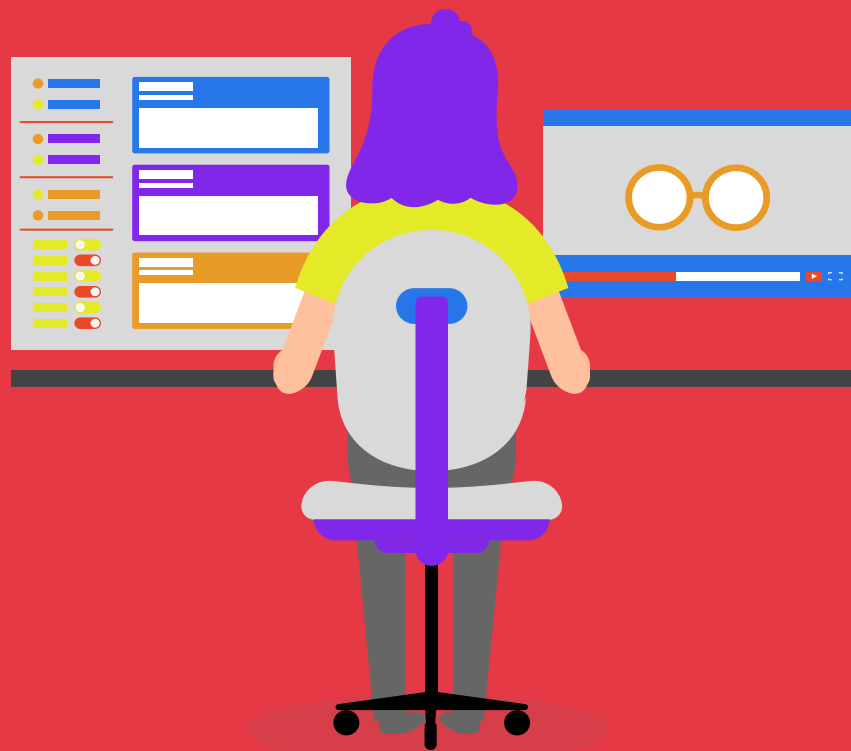
# (Biased) CREDIT SCORING

{ON: The Beach}

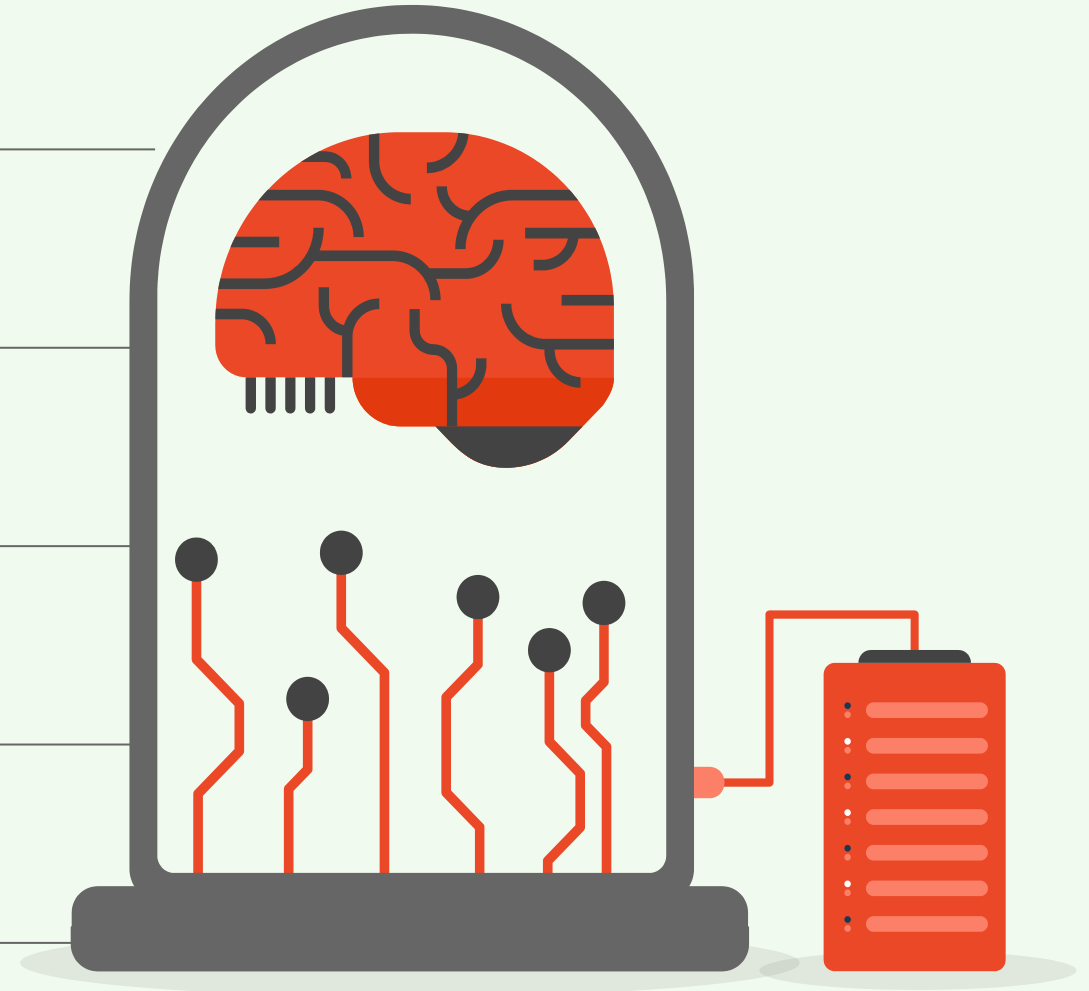How can black-box models lead to unfair, biased outcomes?

**01** **Biased training dataset**

**02** **Lack of diversity within dataset**
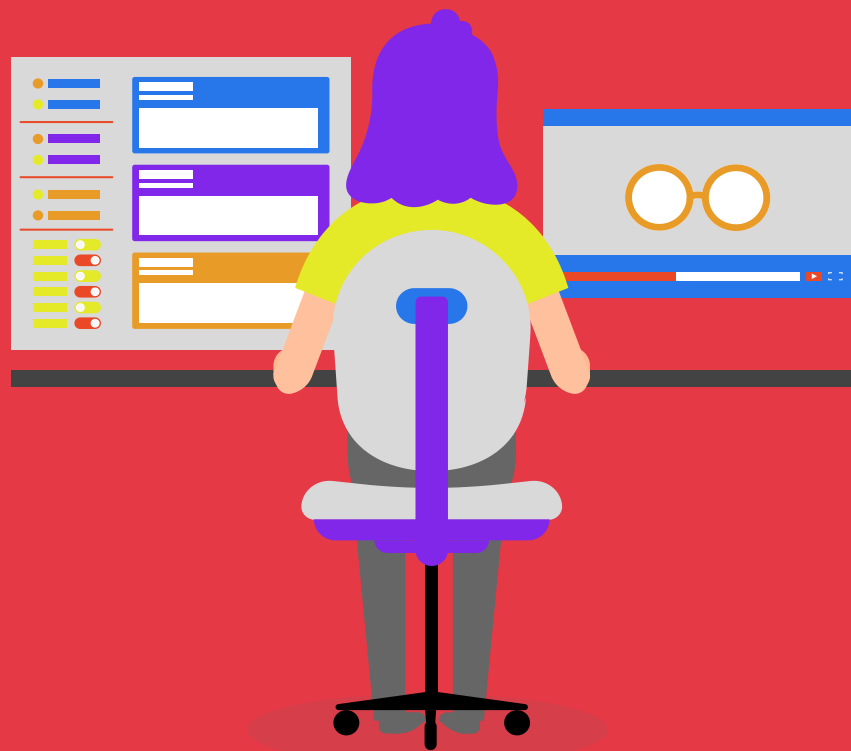
**03** **Cognitive bias**

**04** **Inadequate evaluation metrics**

**05** **Systemic bias**

Who tends to experience bias from such ML systems?

**{ON: The Beach}**

**Age**

Older people, over the age of 50 more likely to experience discrimination

**01**

**02**

**Gender**

Favoritism towards one gender over the another

**Disability**

People with disabilities are often forgotten about during the design of ML systems

**03**

**Multiple protected or marginalized groups**

**04**

**Race and ethnicity**

People who are identified as more than one race, are subject to racial bias

**Immigration Status**

Immigrants significantly face unfair, biased outcomes

**05**

**06**

**Language**

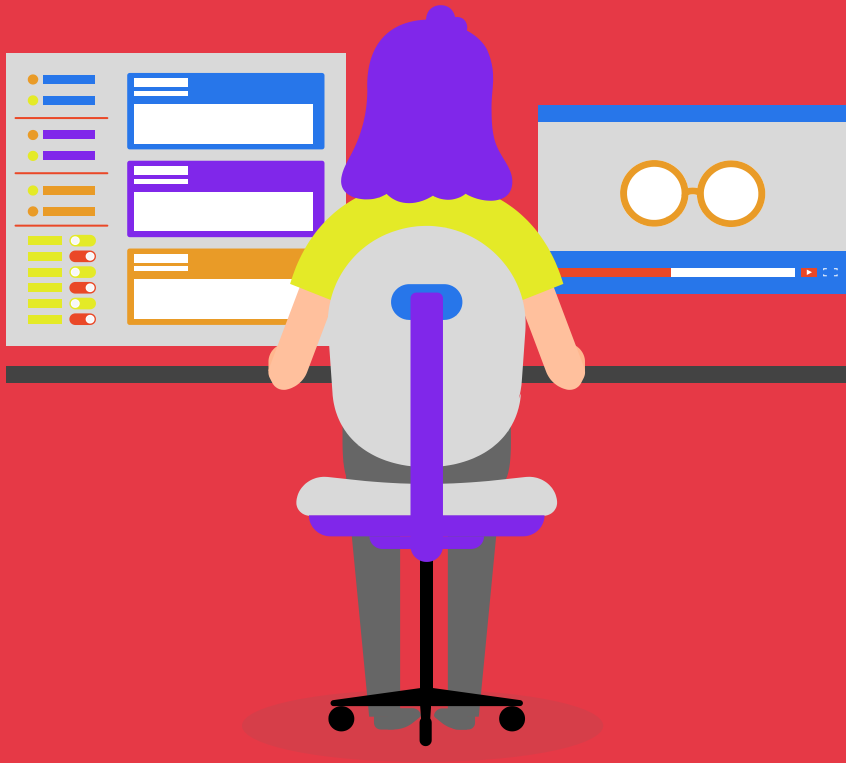Especially in online content, those who use other than English are more likely to experience bias

#JOTB23 - 17

# What is fairness in ML models?

{ON: The Beach}

**Fairness in the machine learning models**

**01**

**With respect to protected attributes**

Loan classification without factoring person's age / gender in decision-making process.
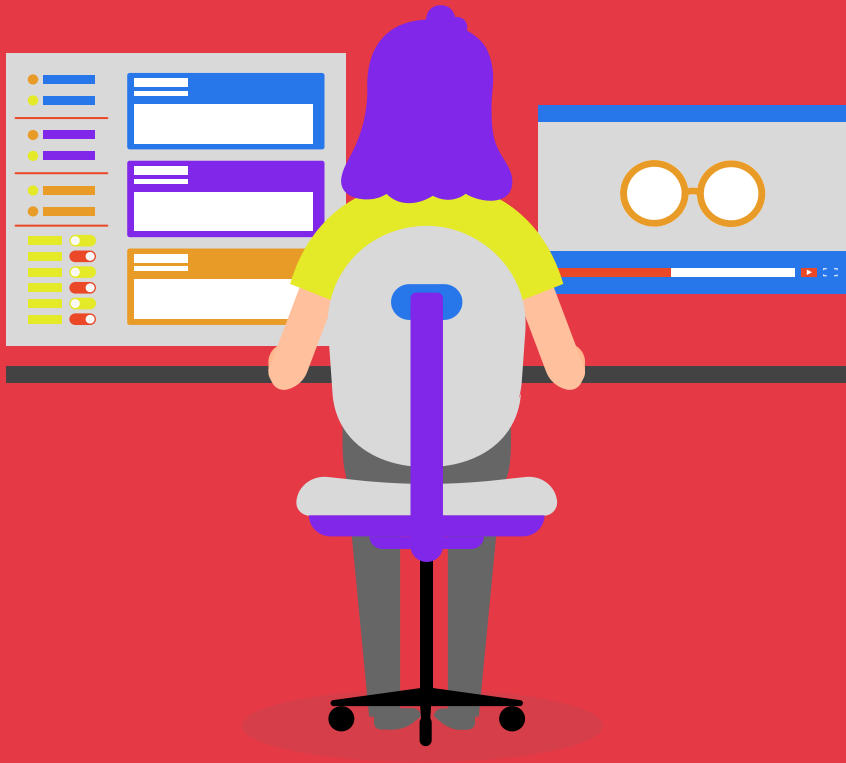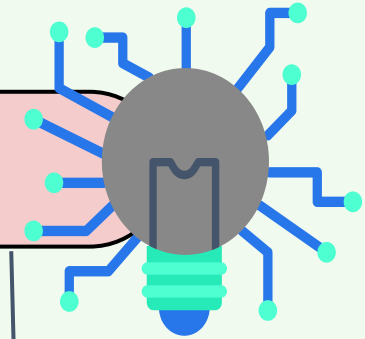
**02**

**In outcomes**

Hiring algorithm should not unfairly favor any particular group

#JOTB23 – 20

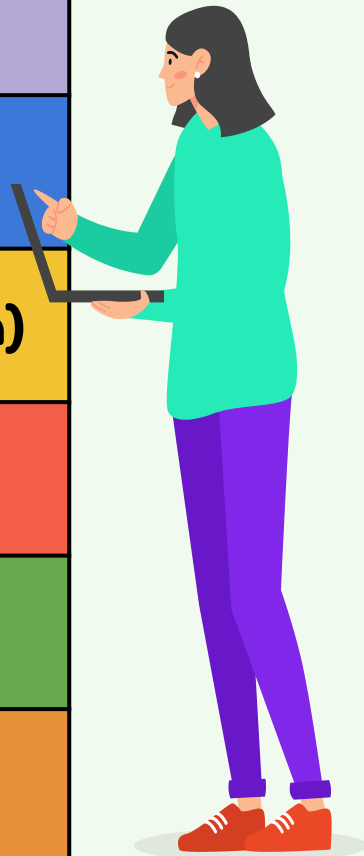What are algorithmic fairness techniques?

# How to build fair and unbiased models?

# Strategies to build fair and unbiased models

- **Collecting diverse and representative dataset**

- **Monitoring for bias**

- **Pre-processing and post-processing dataset (eg: data aug., feature selection)**

- **Regularization to reduce overfitting in the dataset**

- **Algorithmic fairness techniques (eg: equalized odds, demographic parity)**

- **Explainability to build transparent and interpretable models**

# Demo

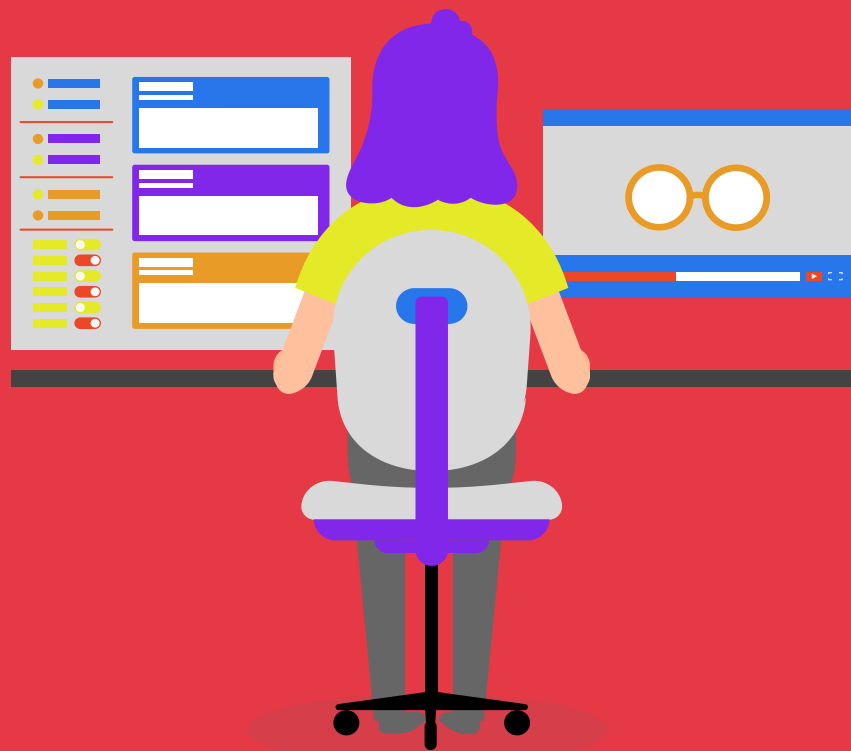# STRATEGIES FOR BUILDING TRUSTWORTHY AND ACCOUNTABLE ML SYSTEMS

Why trust is critical for adoption and success of ML systems?

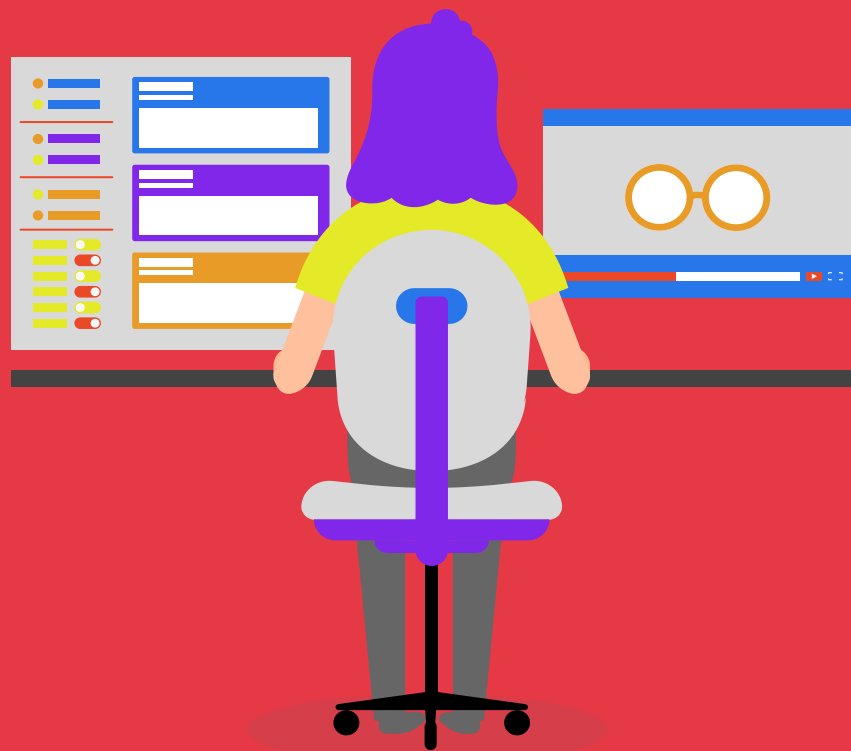Autonomous Driving
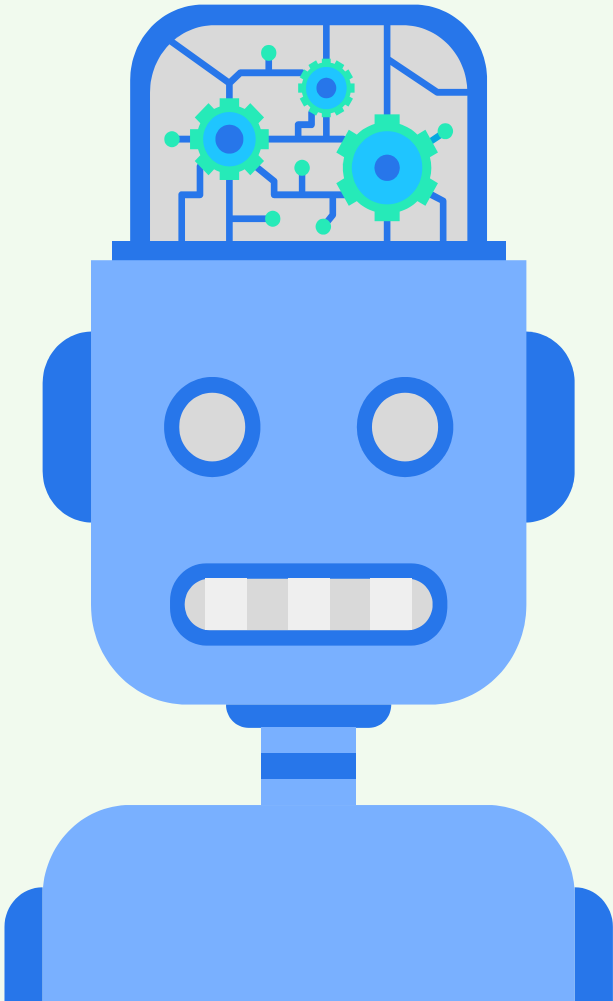


Healthcare

# How to build trustworthy models?

# What is explainability of models?

# Explainable Data

**01**

What data was used to train the model?

# Explainable Predictions

**02**

What features and weights were used for this particular task / prediction?

# Explainable Algorithms

**03**

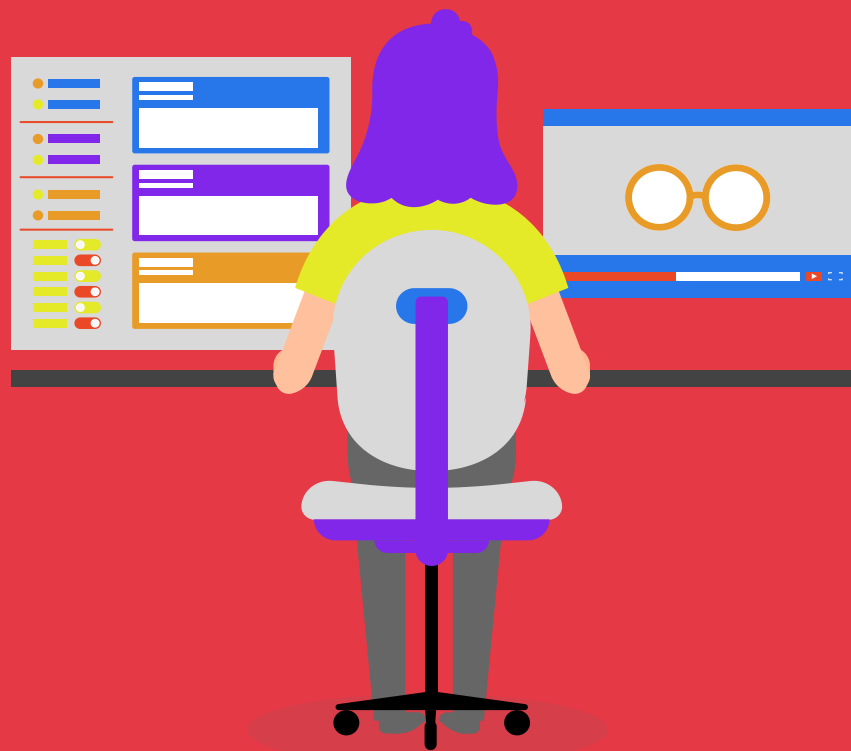What are the individual layers and the thresholds for predictions?

# Demo

Who is responsible for considering the ramification of ML system?

# CONCLUSION

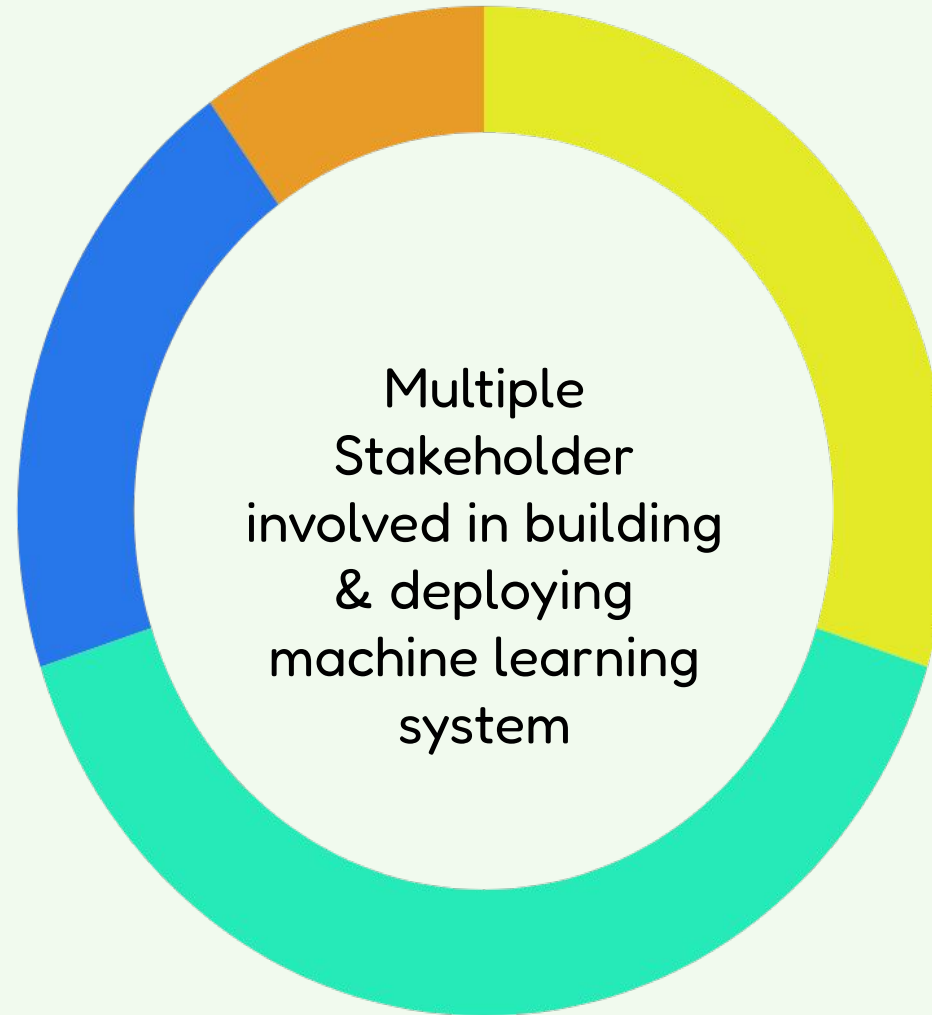- Building fair, accountable, and trustworthy machine learning systems is critical for ensuring that the benefits of these systems are widely accessible and enjoyed by all.

- By ensuring that our models are transparent and interpretable, we can detect and mitigate biases and unfairness before they cause harm.

- Monitoring and evaluation of these models is vital to maintain accountability and ensure that they continue to function as intended.

# Resources

1) Blog: https://eugeneyan.com/writing/testing-pipelines/
2) Book: https://fairmlbook.org/
3) Course: Introduction to Deep Learning
   http://introtodeeplearning.com/
4) Article:
   https://hdsr.mitpress.mit.edu/pub/f9kuryi8/release/8

# Thank you!

@iamrashminagpal