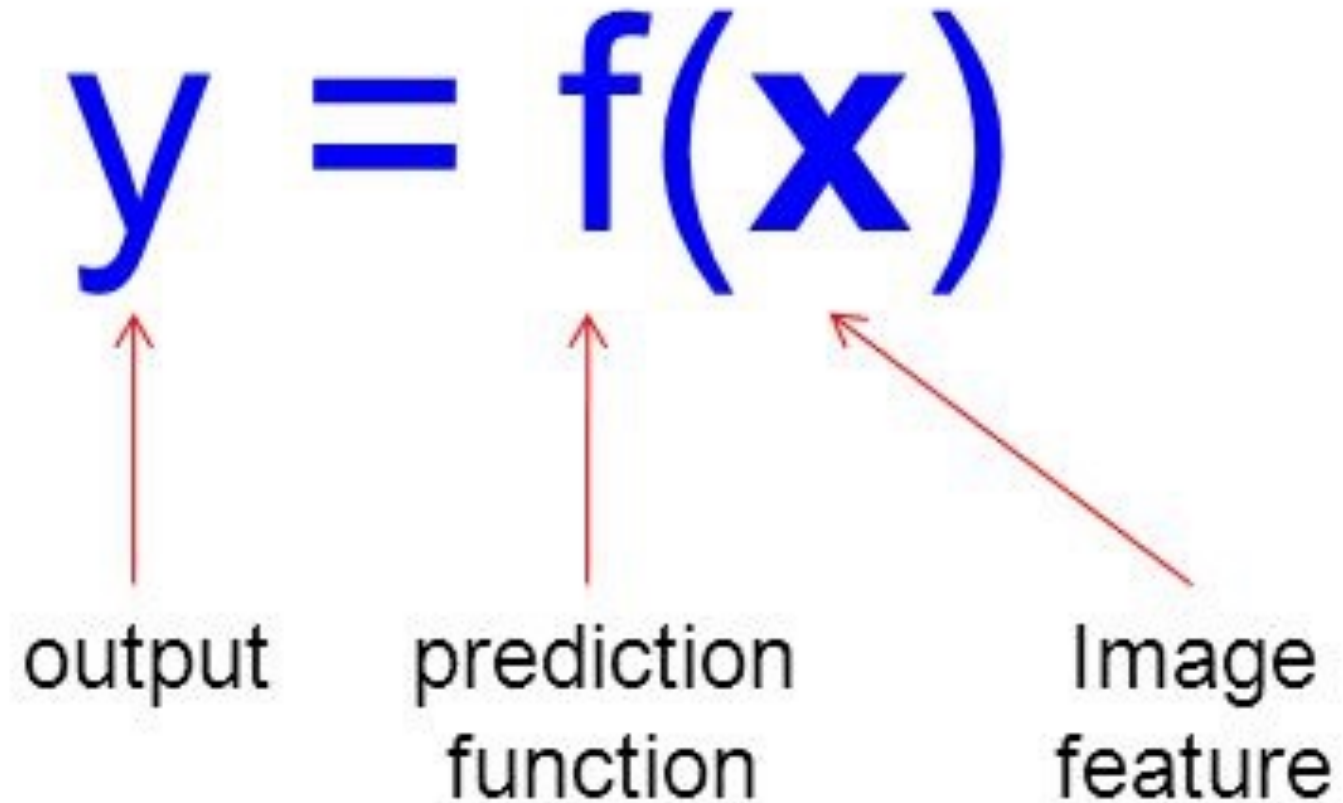# Data Labeling

ICS 491

# For Final Project Checkin #1 on Sep. 26, be prepared to discuss:

- General project idea

- Specific datasets you will use or recruitment strategy

- How your project goes beyond prior work

I will provide feedback in class for other students to learn. (Roughly 2.5-3 minutes for each student's discussion)

# Why do we need data labels?

Machine Learning Setup:

# In machine learning, we learn from data AND labels

- You need to provide a bunch of input/output pairs to train a machine learning model

- The model learns to predict the correct output based on the input

- We talked about acquiring data last class, but for machine learning, we need to be able to make predictions from those data

# Popular Data Labeling Platforms

- Amazon Mechanical Turk

- Snorkel

- Labelbox

- Hive

- …

# Paid Data Labeling

# Amazon Mechanical Turk

## Crowdsourcing platform

# Amazon Mechanical Turk: Example Task

# Amazon Mechanical Turk: Example Task

# Amazon Mechanical Turk: Example Task



(a)                (b)                (c)

# Amazon Mechanical Turk

All HITs    Your HITs Queue

## HIT Groups (1-20 of 640)

Show Details    Hide Details    | Items Per Page: 20

| Requester | Title | HITs | Reward | Created | Actions | |
|-----------|-------|------|--------|---------|---------|---|
| ScoutIt | Classify Receipt | 151 | $0.03 | 14s ago | Preview | 🔒 Qualify |
| Crowdsurf Support | Full Text Review - Earn up to $... | 53 | $0.17 | 3m ago | Preview | 🔒 Qualify |
| Laura A. King | Personality, Information Proce... | 1 | $0.15 | 4m ago | Preview | Accept & Work |
| Crowdsurf Support | Review, edit, and score the tra... | 1,091 | $0.02 | 5m ago | Preview | 🔒 Qualify |
| Erica Fissel | Quick Demographic Survey!(~... | 1 | $0.01 | 6m ago | Preview | Accept & Work |
| ScoutIt | Extract summary information fr... | 1 | $0.05 | 9m ago | Preview | Accept & Work |
| Crowdsurf Support | Transcribe up to 35 Seconds o... | 1,042 | $0.05 | 10m ago | Preview | 🔒 Qualify |
| Ben Stevens | Help Pick a Book Cover! | 1 | $0.10 | 12m ago | Preview | Accept & Work |
| ScoutIt | Extract summary information fr... | 1 | $0.05 | 12m ago | Preview | Accept & Work |
| Michael Busseri, PhD | Answer survey (10 minutes) a... | 1 | $1.00 | 12m ago | Preview | 🔒 Qualify |
| Amy Minnikin | Feedback Seeking Motives Pr... | 111 | $0.25 | 15m ago | Preview | 🔒 Qualify |
| SEO BrainTrust | Summarize and write three ke... | 23 | $0.35 | 25m ago | Preview | 🔒 Qualify |

# Amazon Mechanical Turk: Worker Rights

Platforms like Turkopticon allow for rating Requesters rather than Workers

# Amazon Mechanical Turk: Worker Rights

**TIME**

SPOTLIGHT MAHSA AMINI'S DEATH STILL HAUNTS THE IRANIAN REGIME

TECH • ARTIFICIAL INTELLIGENCE

## Gig Workers Behind AI Face 'Unfair Working Conditions,' Oxford Report Finds

ARTIFICIAL INTELLIGENCE

## Amazon Mechanical Turk Pays Less Than 40% of US Minimum Wage, Research Suggests

Updated on December 9, 2022
By **Martin Anderson**

## Crowdsourcing Platform Market Will Show the Highest Growth Rates & Incredible Demand By 2029 | 99designs, Amazon Mechanical Turk, Crowdcube, CrowdSource, Crowdspring

**WIRED** BACKCHANNEL BUSINESS CULTURE GEAR IDEAS SCIENCE SECURITY MERCH

MORGAN MEAKER BUSINESS SEP 11, 2023 6:00 AM

## These Prisoners Are Training AI

In high-wage Finland, where clickworkers are rare, one company has discovered a novel labor force—prisoners.

## Clickwork in Brazil: A mom balancing endless gigs with childcare

Through platforms like UHRS, Amazon Mechanical Turk and Appen, a clickworker takes on a job consisting of hundreds of quick, repetitive tasks.

THE RISE OF AI

## Clickworkers are turning against each other

Brazil's online gig workers, who power features like social media moderation and AI training, are increasingly wary of newcomers.

# Other Platforms: Hive

# Other Platforms: Label Studio

# Other Platforms: Label Box

# Other Platforms: CVAT

# Other Platforms: Tag Tot

# Another Solution: Hire Your Own Trained Labelers

**BUSINESS • TECHNOLOGY**

## Exclusive: OpenAI Used Kenyan Workers on Less Than $2 Per Hour to Make ChatGPT Less Toxic

# Another Solution: Hire Your Own Trained Labelers



Data Labeler
Tesla
Palo Alto, CA
via Salary.com

Full-time · Health insurance · Dental insurance
Paid time off

Data Labelling Analyst (Remote)
Staffingly
Anywhere
via ZipRecruiter

Work from home · Full-time

Data Labeler
Acceler8 Talent
via LinkedIn

21 days ago · $ 20–50 an hour · Contractor

Labeling Data Analyst Assistant - Mid Level (Remote)
millenniumsoft
Sparks, MD
via ZipRecruiter

Contractor and Temp work

Data Labeling Specialist, German - REMOTE
Indico Data
Boston, MA
via Lever

Part-time

New job alerts

## Data Labeler

SAVE

Tesla
Palo Alto, CA

**Apply on Salary.com**   **Apply on Tarta.ai**   **Apply on Camino Bluff**

Full-time · Health insurance · Dental insurance · Paid time off

### Job highlights
Identified by Google from the original job post

**Qualifications**

- No previous experience in AI or data labeling required

- The perfect candidate for this role is someone who is adaptable, can apply logic to multiple different scenarios, is attentive to detail, has experience with computers and other software, and enjoys a fast…

- High School diploma or evidence of exceptional ability

- Must have a valid driver's license & knowledge of the roads

- Passionate and curious about technology

- Available to work overtime as needed

7 more items

**Responsibilities**

- You will use in-house tools to label images for the Autopilot team

- Interact with team members to help us improve on the design of an efficient labeling interface

2 more items

**Benefits**

- Pay: $22.00 - $24.00 per hour

- 401(k) matching

6 more items

More job highlights

### Job description

We are looking for a driven team member to contribute to the development of our Full Self Driving software at Tesla. This person labels images & objects that contribute to our deep learning neural network. In this role you will work with in-house tools to label images coming directly from the Tesla fleet.

# Another Solution: Hire Your Own Trained Labelers

# Quality Control

# Quality Control

# Creative Data Labeling Strategies

# Captchas are for crowdsourced data labeling

They know the right label of some but not all figures. Your test is a subset of the figures and the rest is training data.

# Example from my research



Parent Guesses → Child Acts → Data Logged

COMPUTER VISION LIBRARIES

GAME PLAY DATA

Kalantarian, **Washington** et al. *IEEE Conference on Healthcare Informatics.* 2018.

Kalantarian, **Washington** et al. *IEEE Transactions on Games.* 2018.

Kalantarian, ..., **Washington** et al. *Journal of Healthcare Informatics Research.* 2019.

Kalantarian, ..., **Washington** et al. *Artificial Intelligence in Medicine.* 2019.

Penev, ..., **Washington** et al. *Applied Clinical Informatics.* 2021.

# Example from my research

Kalantarian, **Washington** et al. *IEEE Conference on Healthcare Informatics.* 2018.
Kalantarian, **Washington** et al. *IEEE Transactions on Games.* 2018.
Kalantarian, …, **Washington** et al. *Journal of Healthcare Informatics Research.* 2019.
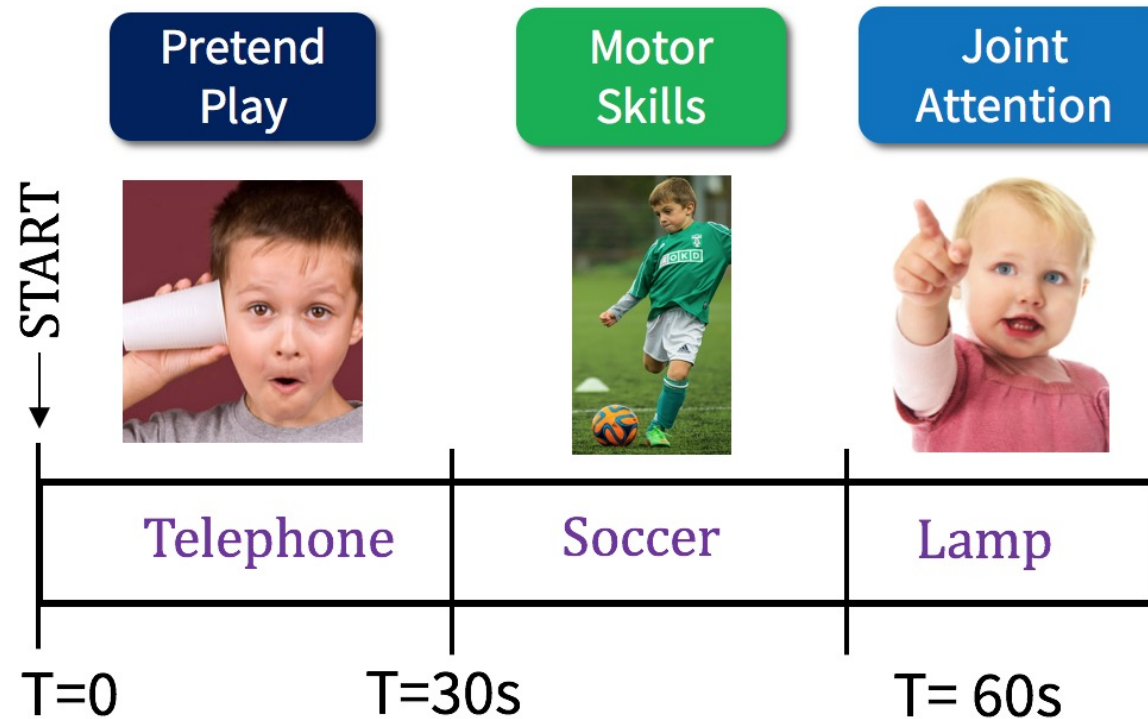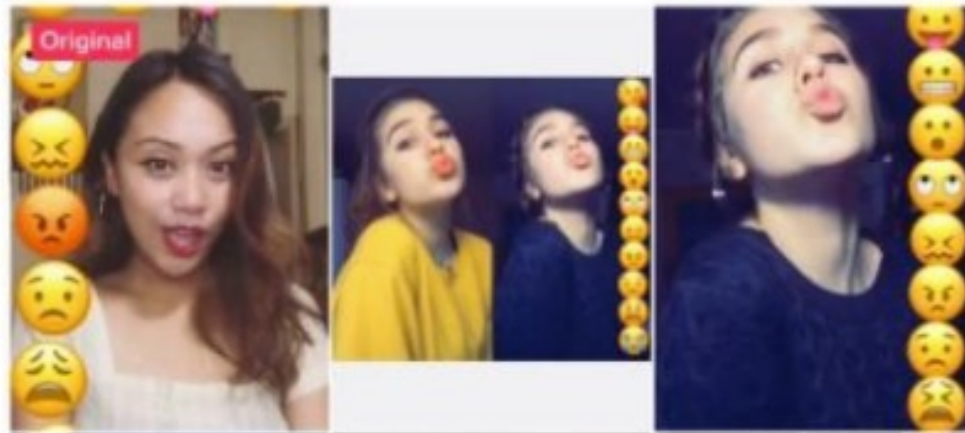Kalantarian, …, **Washington** et al. *Artificial Intelligence in Medicine.* 2019.
Penev, …, **Washington** et al. *Applied Clinical Informatics.* 2021.

# Another example from my research



(a) Best Emoji Face Challenge.

(b) Face Challenge.

Surabhi, …, **Washington** et al. *CVPR Workshop.* 2022.

# Example from my PhD student



| User profile classification | | |
|---|---|---|
| **Word vectorization method** | **Model** | **Metric performance on test set** |
| Keras Embedding | Attention + BiLSTM | Accuracy: 0.87 |
| | | F1 score: 0.805 |
| | | AUC score: 0.78 |

# One more example from my lab's research



Input: Twitter Data

Model: Neural Network

Output: Probability of Cigarette Use

97.23%

Aditi Jaiswal

Dr. Pokhrel

UNIVERSITY OF HAWAI'I CANCER CENTER

Dr. Amin

# Class Exercise

- Spend 5-10 minutes thinking through the (1) recruitment, (2) data collection, (3) data labeling, and (4) data preprocessing strategies that you will need to conduct for your project

- We will discuss people's strategies together as a class