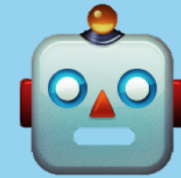# How to Feed Your Robot

## Building and Maintaining Open Machine Learning Datasets 🤖

**Evan Tachovsky**, Rockefeller Foundation

@evantachovsky

2019-05-08

# Supervised ML, grossly simplified

**1. Train** an algorithm with labeled data

| Picture | Show **a label** |
|---|---|
| 🐶 | "This is a **dog**" |
| 🐱 | "This is a **cat**" |
| 🐶 | "This is a **dog**" |
| 🐱 | "This is a **cat**" |
| | *…lots of other examples…* |
| 🐶 | "This is a **dog**" |

**2.** Use the algorithm to **classify** new pictures of dogs and cats

| New **picture** | Algorithm **classifies** |
|---|---|
| 😹 | "I'm 0.79 sure that is a **cat**." 🤖 |
| 🐶 | "I'm 0.87 sure that is a **dog**." 🤖 |

# Massive Bias Problems

**60%**
images in
ImageNet and
Open Images
comes from
**six countries**

1 ▬▬▬▬▬ 642,997

*Source:* No Classification without Representation: Assessing
Geodiversity Issues in Open Data Sets for the Developing World
by Shreya Shankar et al.

# Current Project

Interview ~**20 data set builders, maintainers, and funders** to understand how we can build just and effective datasets for development and humanitarian use

- Semi-structured conversations
- Spread across industry, academic, research institutions
- No attribution unless the info was published
- Most interviewees were working on text or image datasets

| | | | | |
|---|---|---|---|---|
| Alex Ratner | Alistair Johnson | Andrew McCallum | Brandeis Marshall | Celina Lee |
| Courosh Mehanian | Desmond Patton | Donghui Li | George Azzari | Jayme Garcia Arnal Barbedo |
| Leo Anthony Celi | Mutale Nkonde | Mutembesa Daniel | Nick Adams | Rashida Richardson |
| Rediet Abebe | Sebastian Ruder | Sue Márquez | Tariq Khokhar | Tony Hey |

**Thank you** to everyone who took the time to talk! Any errors are on me.

# Five lessons

1. Motivations shape dataset
2. Transactional labels are always worse than you think
3. Find your place on the labeling spectrum
4. Don't ignore shelf life
5. Think infrastructure, not research

# 1. Motivations shape datasets

## Commercial

- I want to build a new product
- This dataset will give us a moat
- This dataset will make our algorithm more robust

## Methodological

- We need a benchmark for the field
- If we have this dataset, we can figure out this method
- I'm curious about _____.

## Applied

- How do I understand these data?
- How to I automate this annoying task?

# 2. Transactional labels are worse than you think

# 3. Find your place on the labeling spectrum

The messy middle

Anyone can label

Experts only

# 3. Find your place on the labeling spectrum

Strategies for the **messy middle**

- Incentives, not just financial
- New tools to make tasks easier
- Change the classification problem

# 4. Don't ignore shelf life

Use Buy:
2020-02-28

```
{
  "users": [
    {
      "id": 0,
      "name": "Jane Smith",
      "work": "ABC",
      "dob": "1978",
      "address": "83 Warner Street",
      "city": "Boston",
      "optedin": true
    },
    {
      "id": 1,
      "name": "Frank Jonses,
      "work": "XYZ",
      "dob": "13/05/1987",
      "address": "9 Coleman Avenue",
      "city": "Toronto",
      "optedin": false
    }
  ],
  "images": [
    "img0.png",
    "img1.png",
    "img2.png"
  ],
  "coordinates": {
    "x": 35.12,
    "y": -21.49
  },
  "price": "$100,000"
}
```

# 5. Think infrastructure, not research

- Ubiquitous, not bespoke
- For a community, not a team
- A separate budget, not a budget line
- Measure and reward contributions

Background Reading:
https://www.are.na/evan-tachovsky/in-training

# Thank you! Get in touch

Evan Tachovsky,
Rockefeller Foundation
@evantachovsky
2019-05-08