Security and Fairness of Deep Learning

# Course Overview

Spring 2020

# Today

- Goals
- Modules

- Prerequisites
- Logistics
- Grading
- Policies

# Course staff

- Instructor: Piotr (Peter) Mardziel
  - Email: piotrm@cmu.edu
  - Office hours: Thursdays 12-1pm Pacific
  - Office:  B23 114

- TA:  Zifan Wang
  - Email: zifan.wang@sv.cmu.edu
  - Office hours: Wednesdays 2:30-3:30pm Eastern
  - Office:  CIC 2206

- TA:  Caleb Lu
  - Email: kaijil@andrew.cmu.edu
  - Office hours: TBD
  - Office:  TBD

- Office hours available remotely:
  - Zoom meeting links on website

# Continuing successes of deep learning

# Image classification



What the computer sees

image classification →
82% cat
15% dog
2% hat
1% mug

# NLP: translation, etc.

# Deep neural networks learn representations



Deeper layers learn progressively more abstract representations:
pixels, edges, motifs, parts of objects, objects

# Enabling trends

- Large volumes of training data

- Computation power
  - GPUs,…

# Course objective

Understand deeply how and why deep networks work and their weaknesses

Become informed: what can go wrong (other than poor performance)?

# Course modules

1. Fundamentals of deep learning
2. Explanations for deep learning
3. Security of deep learning
4. Privacy and Fairness in deep learning

# Course modules

1. Fundamentals of deep learning
   - Background on machine learning
   - Architectures, training, platforms
   - Focus on convolutional and recurrent neural networks

# Course modules

2. Explanations for deep learning
   - Feature importance and visualization



how many townships have a population above 50 ? [prediction: NUMERIC]
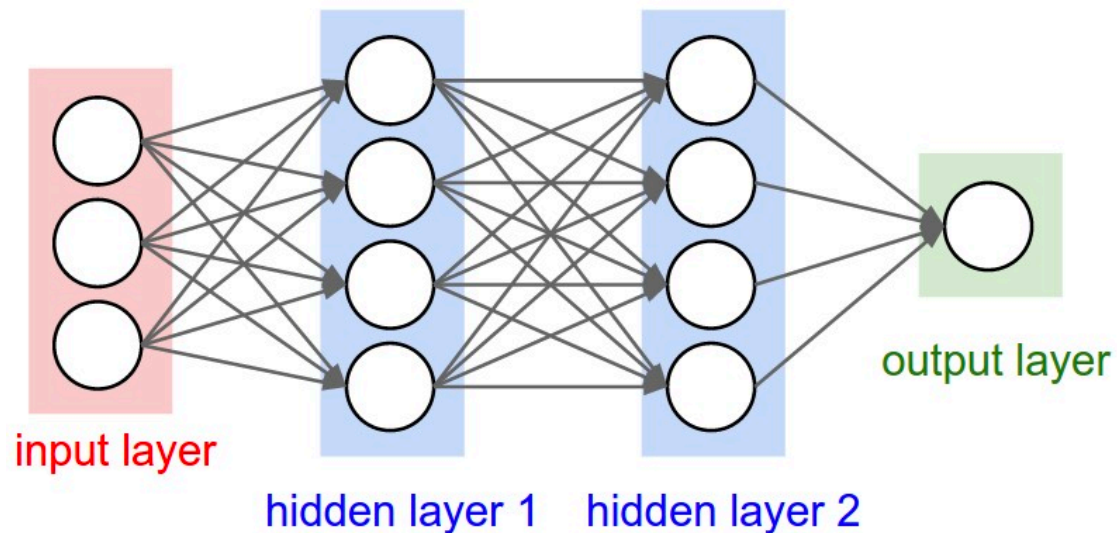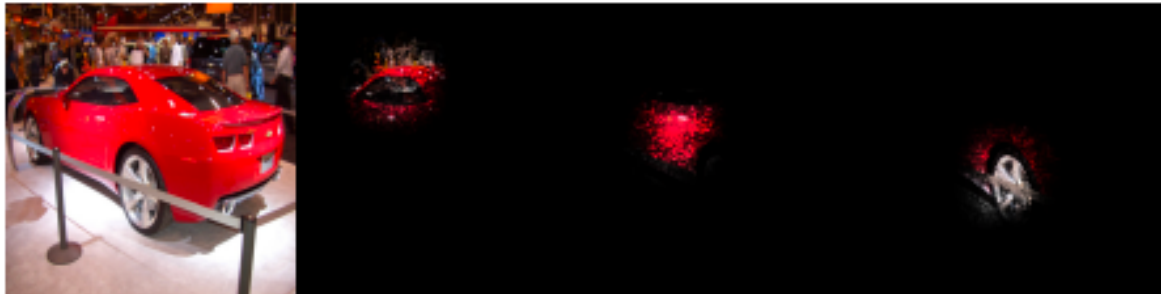what is the difference in population between fora and masilo [prediction: NUMERIC]
how many athletes are not ranked ? [prediction: NUMERIC]
what is the total number of points scored ? [prediction: NUMERIC]
which film was before the audacity of democracy ? [prediction: STRING]
which year did she work on the most films ? [prediction: DATETIME]
what year was the last school established ? [prediction: DATETIME]
when did ed sheeran get his first number one of the year ? [prediction: DATETIME]
did charles oakley play more minutes than robert parish ? [prediction: YESNO]
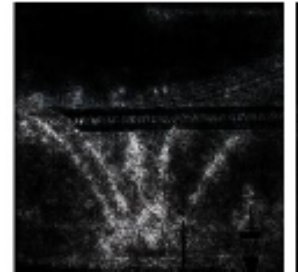
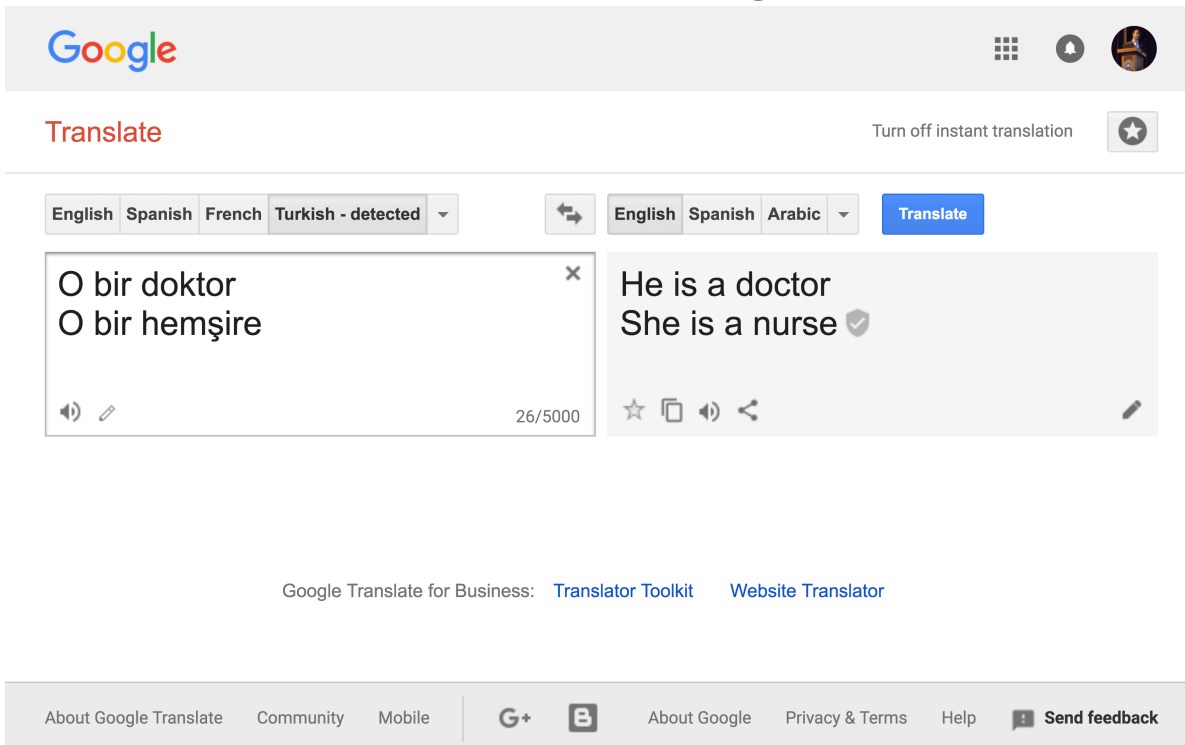| Original image | Top label and score | Integrated gradients |
|---|---|---|
| | Top label: reflex camera Score: 0.993755 | |
| | Top label: fireboat Score: 0.999961 | |
| | Top label: school bus Score: 0.997033 | |
| | Top label: mosque Score: 0.999127 | |

# Course modules

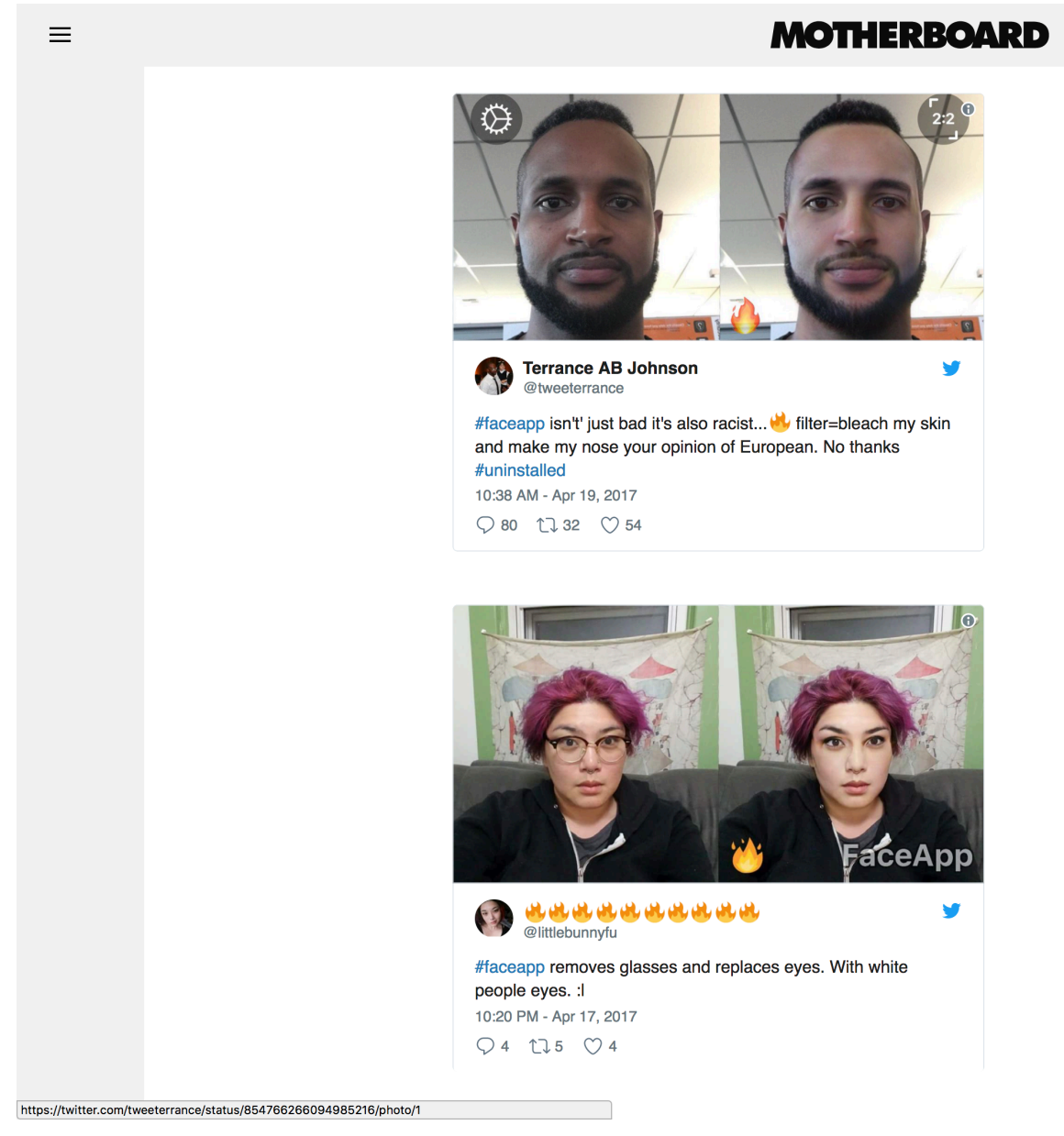3. Security of deep learning models
    - Attacks on classifiers and defenses

# Course modules

4. Privacy and Fairness in deep learning

- Inferring sensitive information
- Bias and de-biasing

# Course Format

- Lectures covering the background
  - Stanford CS231n (Convolutional Neural Networks for Visual Recognition)
  - Deep Learning textbook
- 1-2 Lecture covering software tools and setup
  - Numpy, Tensorflow, Keras, Jupyter Notebook, Google Computing Services
- **Lectures covering research papers**
  - Occasionally guest lecturers

# Prerequisites

- No formal prerequisites

- Basics of linear algebra, probability, multivariate calculus
  - Will review briefly in class and provide resources to learn on your own
  - Roughly Chapters 1-5 of [Deep Learning](#) textbook by Goodfellow et al.

- Familiarity with Python
  - Necessary for programming homework

- Quick class poll

# Logistics

- Lectures:  Tue & Thur,  10:30-11:50am Pacific / 1:30-2:50pm Eastern
- Web page:  http://www.ece.cmu.edu/~ece739/

- Gradescope (assignment submission)
- Canvas (grades)
- Piazza (announcements, for all other communication)

- Textbook
  - Deep Learning textbook by Goodfellow, Bengio, Courville

# Grading

- Homework: 90%
  - 5 x 18%
- Class participation: 10%
  - Be present and engaged in class and piazza
  - Informed questions for guest lecturers

# Collaboration policy on homework

- You are allowed/encouraged to discuss homework problems with other students in the class but are required to write out solutions independently and to acknowledge any collaboration or other source. If you are unsure about something, consult the course staff.

CMU Computing Policy

CMU Policy on Cheating

# Acknowledgment

- Based on material from
  - Spring 2019 Course