# Impact of AI Decal:
## *The Data Revolution*

Gokul Swamy & Brenton Chu

Quiz: https://tinyurl.com/impactsp19q3

*Legolas and Gimli advanced on the orcs, raising their weapons with a harrowing war cry.*

The orcs' response was a deafening onslaught of claws, claws, and claws; even Elrond was forced to retreat. "You are in good hands, dwarf," said Gimli, who had been among the first to charge at the orcs; it took only two words before their opponents were reduced to a blood–soaked quagmire, and the dwarf took his first kill of the night. The battle lasted for hours until two of the largest Orcs attempted to overwhelm Aragorn. When they finally stopped, they lay defeated and lifeless for miles and miles.

"I take nothing," said Aragorn. "But I give my word, at my peril and mine, that I will never forget this day of horror. None of us will forget. Ever!"

"I'll never forget it!" cried Gimli, who had been in the thick of the battle but hadn't taken part in it. One of the wounded orcs he had carried off, he was the only one of the survivors who remained uninjured. "We'll keep the memory of that day of evil, and the war with it, alive as long as we live, my friends!"

"Then we'll keep it alive as long as we live," added Legolas. "And we won't forget the first great battle of the night, even if we may have forgotten the final defeat."

"I agree," Gandalf said, "but we will all remember it as the last battle in Middle–earth, and the first great battle of the new day."

Aragorn drew his sword, and the Battle of Fangorn was won. As they marched out through the thicket the morning mist cleared, and the day turned to dusk.

The Two Rings were returned to Rivendell. Frodo and Sam woke up alone in their room, and Frodo found a note on his pillow. He opened it and read:

# Current AI Systems Can...

Write a sensible Harry Potter chapter

# FALSE

# Remember this?

# OpenAI's Realistic Text-Generating AI Triggers Ethics Concerns

**This AI is so good at writing that its creators won't let you use it**

Too scary? Elon Musk's OpenAI company won't release tech that can generate fake news
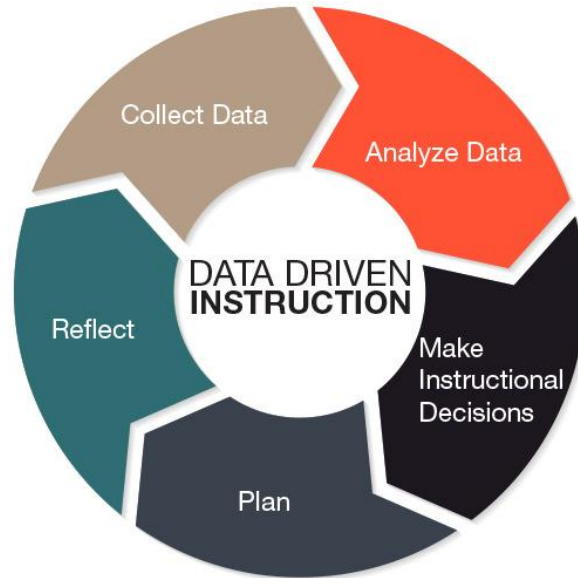
# What's the secret?

➔ Lots of computation power and data!
➔ Used 40GB of text scraped from the internet
  ◆ Equivalent to nearly 100,000 full length books
➔ Previous iteration used only ~5GB
➔ Ethical concerns about this AI (specifically in relation to generating fake news) will be discussed in a later lecture

# Context

# Why is data important?

➔ A: The dominant paradigm in modern ML is supervised learning
➔ B: Supervision is a data-driven approach
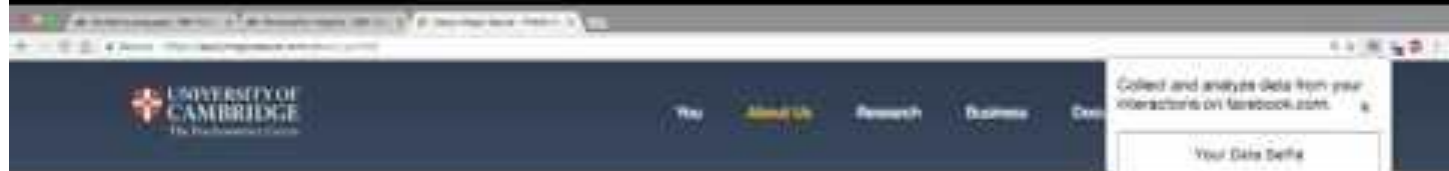➔ Thus, we need data to drive the engines of modern ML
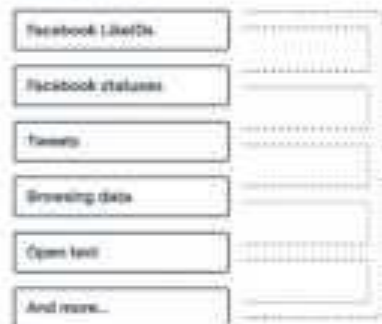
# Q:

*How is data collected?*

# How is data collected?

➔ Data can be collected explicitly from users by from users filling out surveys
  ◆ This requires some sort of incentive for the user to complete so is harder to implement in practice

➔ Most data is collected without notifying the user
  ◆ Facebook can see what kinds of posts you like
  ◆ Amazon can see what kinds of products you view or purchase
  ◆ Google can see what links you visit

# Attention Merchants

➔ Term used to describe corporations who make much of their profits based on serving the attention of users to advertisers.
➔ When new features are being developed at Facebook, metrics related to amount of time spent on the website are considered.

# How is data used?

➔ Data allows the collectors to better understand the consumers.

➔ The collectors can use this insight for whatever they want

➔ Prioritizing results
   ◆ Google PageRank used to generate results

➔ Recommender systems:
   ◆ Spotify playlist generations (built using [this](#))

➔ Psychographics:
   ◆ Using likes, can predict black/white with 90% accuracy and straight/queer with 88% accuracy
   ◆ Used by Cambridge Analytica

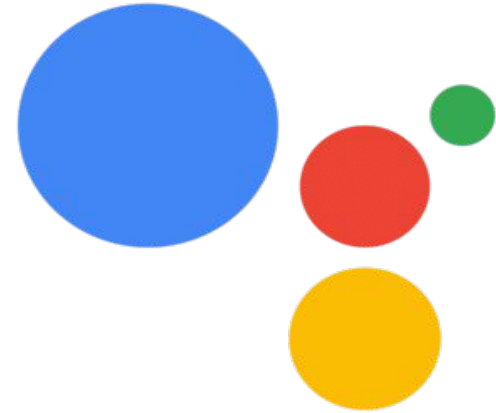# Facebook Like Psychographics (1):

# Facebook Like Psychographics (2):

# Data Use Case Study: Google vs. Apple

**Siri**
Introduced by Apple for iOS in 2011

**Google Assistant**
Introduced by Google for Android in 2016

GOOGLE ASSISTANT **VS** SIRI

# Responses: Differential Privacy

➔ Strength in numbers
➔ Security through randomization and then aggregation of data
   ◆ Uses statistics to still arrive at correct answers without compromising on one's privacy
➔ Currently used heavily by Apple/Google (in areas like typing suggestions)
➔ Only works when there is a large amount of data and learning from aggregates is a feasible strategy

# Differential Privacy: Coin Flipping

➜ In your head, think of your answer to the question "do you like pineapple on pizza?"

➜ Google "coin flip"
  ◆ If heads, say yes
  ◆ If tails, vote your true feeling

➜ In the limit, this should match the true value.
  ◆ This is called RAPPOR

# Responses: Blockers

➜ **Adblock**
  ◆ Block general ads from being show as you browse
➜ **Privacy Badger**
  ◆ Blocks ads that track you between sites and gather info without you knowing
  ◆ Something similar implemented by default in Safari

# Responses: GDPR

➔ Law passed in the EU that defines rules for protecting data and penalties for noncompliance.
➔ Applies to everything from name to sexual orientation to political opinions
➔ Fail-safe default: use highest privacy option by default
➔ Design security in from the start, don't just try and obscure data
➔ Affirmative consent about what data will be used for, how long it will be used, and who it will be shared with
   ◆ Right to be forgotten in limited cases

# Impact of AI Decal:
*Activity*

# What data should be given?

➔ You will be given a series of scenarios about collecting data for certain purposes.

➔ Separate into two sides of the room based on whether you agree with the collection of data in the scenario

➔ A short discussion will follow after each prompt

Taking Reddit posts and comments to train a text based ML model

Using DNA gathered by sequencing companies to identify genes and genetic diseases

Tracking browser history to serve you more relevant products

Looking at photos on your phone to automatically create albums for you

Collecting location data to provide live traffic information

# Q:

*How do we balance privacy with the need for data in ML systems?*

# Writing Assignment One

➜ Due 11:59pm next Friday

➜ Two pages, double-spaced

➜ Helpful to do your own research to assist you

◆ Please cite your sources!

➜ Address one of these prompts:

◆ Given past developments in AI, what are some things we may expect AI to be capable of doing in the next few years? Please discuss at least two of the subfields talked about in week 2.

◆ What are some ways that we can protect ourselves from the negative effects of data collection (compromised privacy, malicious use of data, etc)? This can be policy proposal, awareness campaign, suggestions for researchers, or any other form of solution.

◆ Take any topic/issue that has been discussed in class up to this point and elaborate on it. What is a potential solution, and why would it work?

# Impact of AI Decal:
*Next: A World of Pure Automation?*