

Machine-learned NLP Systems	Broad Areas	
<ul> <li>Aggregate textual information to make predictions</li> <li>Hard to know why some predictions are made</li> <li>More and more widely use in various applications/sectors</li> </ul>	<ul> <li>Bias amplification: systems exacerbate real-world bias rather than correct for it</li> </ul>	
<ul> <li>What are the risks here?</li> <li>of certain applications?</li> </ul>	Exclusion: underprivileged users are left behind by systems	
<ul> <li>IE / QA / summarization?</li> <li>MT?</li> </ul>	<ul> <li>Dangers of automatic systems: automating things in ways we don't understand is dangerous</li> </ul>	
<ul><li>Dialog?</li><li>of machine-learned systems?</li></ul>	Unethical use: powerful systems can be used for bad ends	
of deep learning specifically?		

# **Bias Amplification**

 Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias

.

- Can we constrain models to avoid this while achieving the same predictive accuracy?
- Place constraints on proportion of predictions that are men vs. women?



Zhao et al. (2017)









۲	Dangers of Automatic Systems	
1HE VERO	📕 TECH - SCIENCE - CULTURE - CARS - REVIEWS - LONGFORM VIDEO MORE - 🛛 🗲 🛩 🔊 ᆂ Q	
Faceb Palest	ook apologizes after wrong translation sees	
Facebook tro	anslated his post as 'attack them' and 'hurt them' ayOng   Oct 24, 2017, 10:43am EDT	

# Dangers of Automatic Systems

Translations of gay	
adjective	
homosexual	homosexual, gay, camp
alegre	cheerful, glad, joyful, happy, merry, gay
brillante	bright, brilliant, shiny, shining, glowing, glistening
vivo	live, alive, living, vivid, bright, lively
vistoso	colorful, ornate, flamboyant, colourful, gorgeous
jovial	jovial, cheerful, cheery, gay, friendly
gayo	merry, gay, showy
noun	
el homosexual	homosexual, gay, poof, queen, faggot, fagot > Offensive terms
<ul> <li>el jovial</li> </ul>	gay
	Slide credit: <u>allout.org</u>

Slide credit: The Verge

# 

# Dangers of Automatic Systems

"Instead of relying on algorithms, which we can be accused of manipulating for our benefit, we have turned to machine learning, an ingenious way of disclaiming responsibility for anything. Machine learning is like money laundering for bias. It's a clean, mathematical apparatus that gives the status quo the aura of logical inevitability. The numbers don't lie."

## - <u>Maciej Cegłowski</u>

### Slide credit: Sam Bowman

# Dangers of Automatic Systems

- "Amazon scraps secret AI recruiting tool that showed bias against women"
  - "Women's X" organization was a negative-weight feature in resumes
  - Women's colleges too

Was this a bad model? May have actually modeled downstream outcomes correctly...but this can mean learning humans' biases

> Slide credit: https://www.reuters.com/article/us-amazon-com jobs-automation-insight/amazon-scraps-secret-ai-recruitingtool-that-showed-bias-against-women-idUSKCN1MK08G

Unethical Use			al Use	Unethical Use	
Surveillance a     Generating co     Comment lo:     1060075606075     Dear Commissioners:     Hi, I'd like to comment on     net neutrality regulations.     i want to     implore     the government to     repeal     Barack Obama's     decision to     regula     interret access.     individuals,	pplications? prvincing fak FCC comment ID: 10000135300724 Dear Chairman Pai, I'm a voter worried about Internet freedom. I'd like to ask Ajit Pai to repeal President Obama's order to regulate broadband. people like me,	E Decys / fak FCC Comment ID: 1060773200312 Infoemation Infoemation ISTROIGNA ISTR	e comments? • What if these were undetectable?	<ul> <li>Sophia: "chatbot" that the creators make incredible claims about</li> <li>Creators are actively misleading people into thinking this robot has sentience</li> <li>Most longer statements are scripted by humans</li> <li>"If I show them a beautiful smiling robot face, then they get the feeling that 'AGI' (artificial general intelligence) may indeed be nearby and viable None of this is what I would call AGI, but nor is it</li> </ul>	
rather than	rather than	rather than		simple to get working" 2018/10/12/sophia-modern-marvel-or- mindless-marketing/	



# Final Thoughts You will face choices: what you choose to work on, what company you choose to work for, etc. Tech does not exist in a vacuum: you can work on problems that will fundamentally make the world a better place or a worse place (not always easy to tell) As Al becomes more powerful, think about what we *should* be doing with it to improve society, not just what we *can* do with it