

CS388: Natural Language Processing

Lecture 26: Wrapup + Ethics



Greg Durrett



Administrivia

- ▶ Project presentations next Tuesday and Thursday
- ▶ Final projects due December 14
- ▶ Please fill out the course evaluation if you haven't already! Time at the end of class today



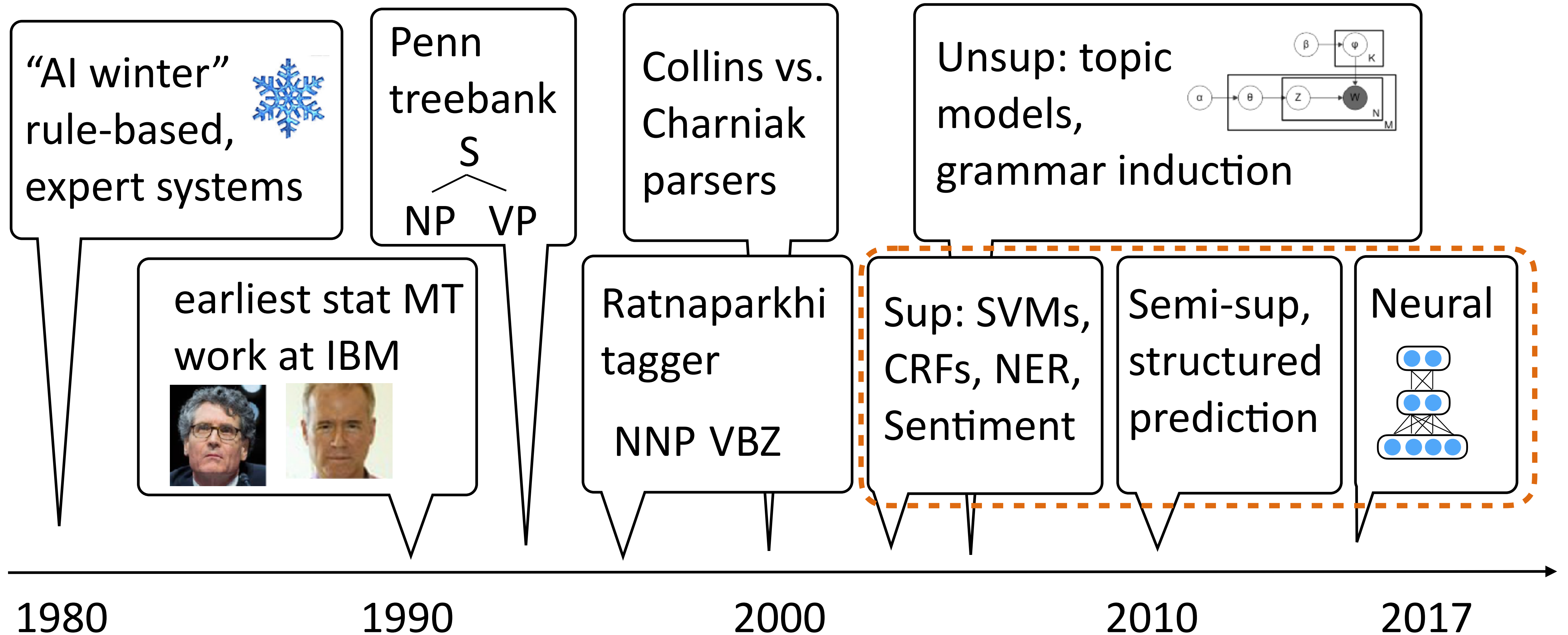
This Lecture

- ▶ Course recap
- ▶ Ethics in NLP

Structure in NLP



A brief history of (modern) NLP



► What different model structures did we consider?



Sequential Structure: Analysis

- ▶ Language is inherently sequential

B-PER I-PER O O O B-LOC O O O B-ORG O O

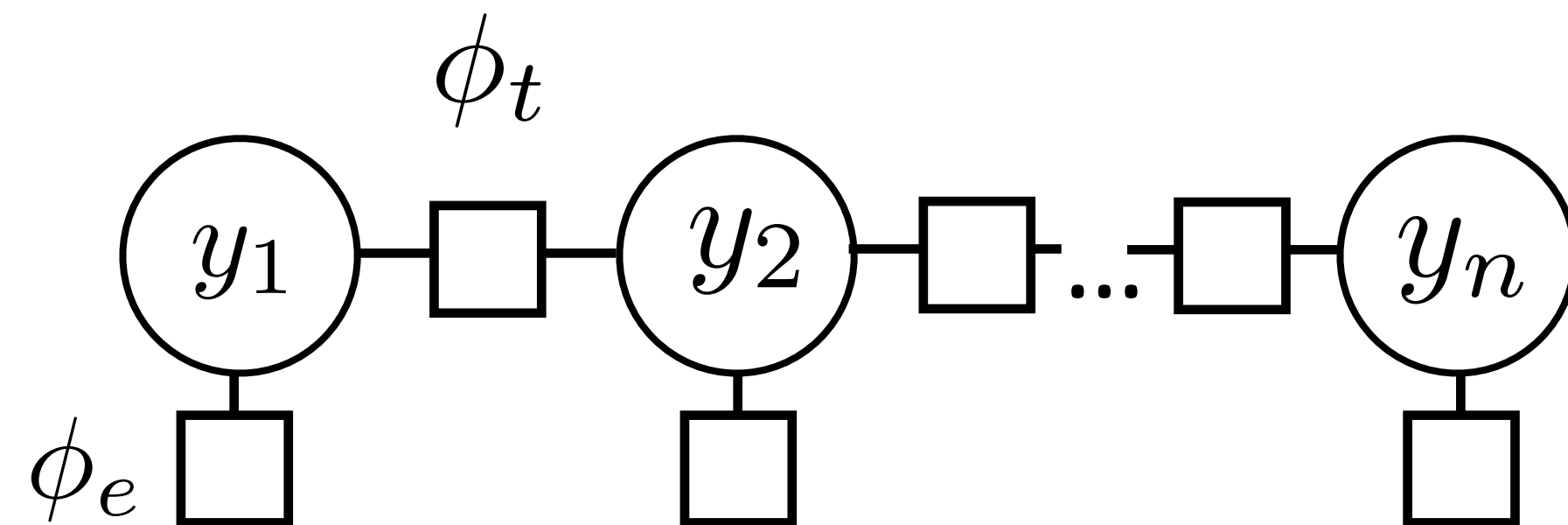
Barack Obama will travel to *Hangzhou* today for the *G20* meeting .

PERSON

LOC

ORG

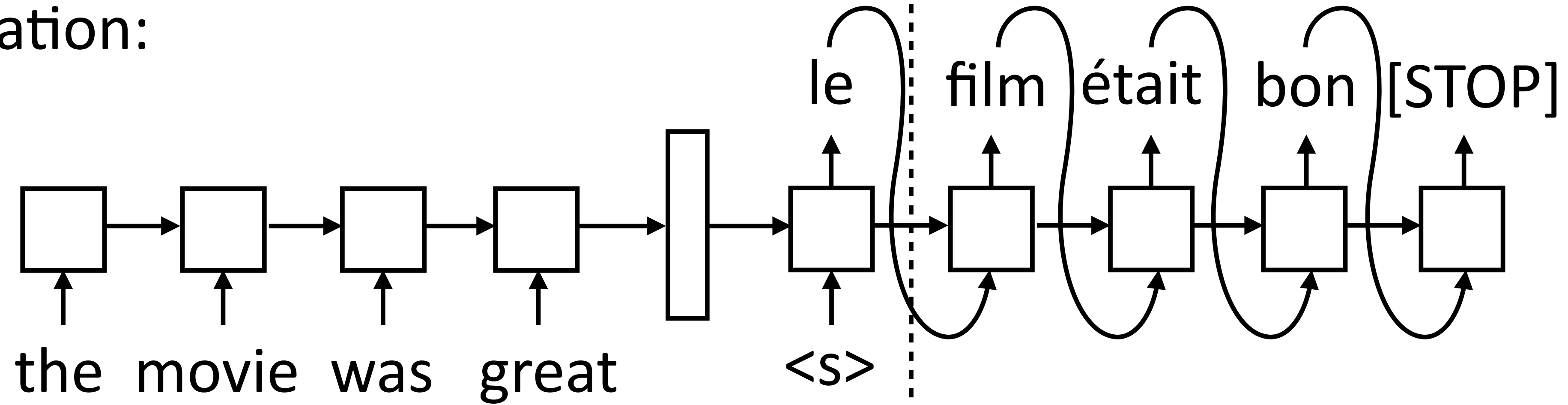
- ▶ Can do language analysis with sequence models



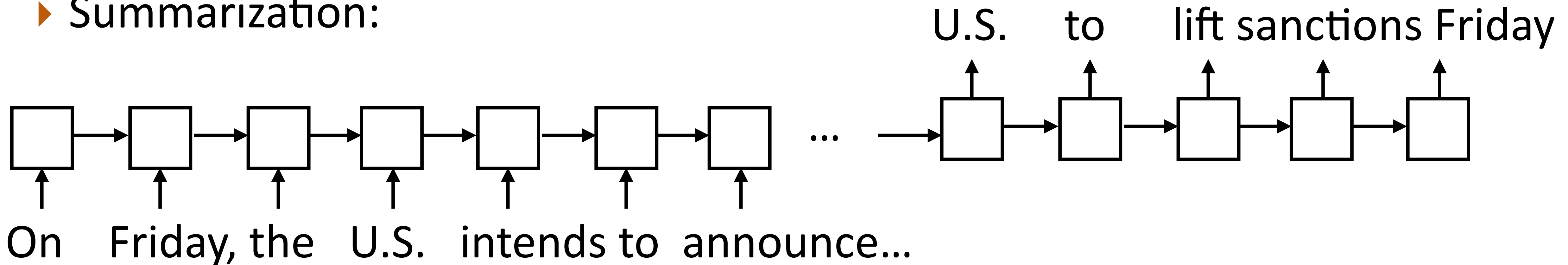


Sequential Structure: Generation

► Translation:



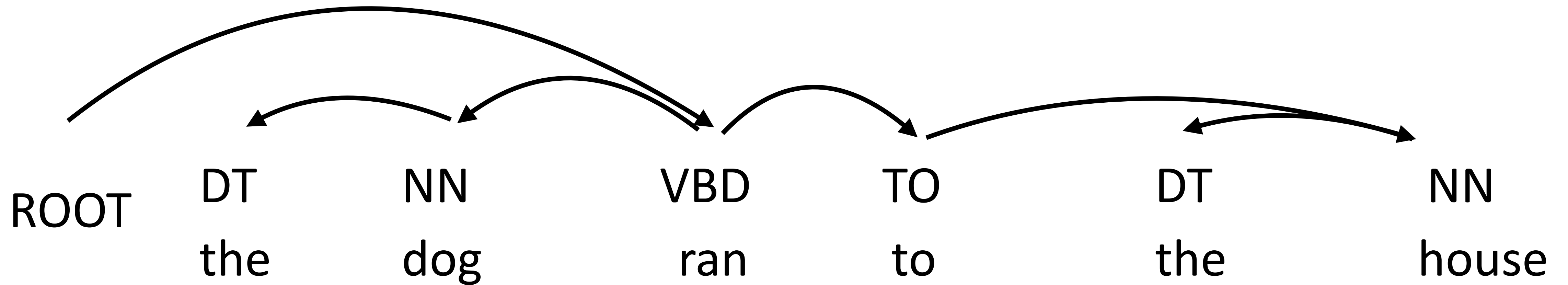
► Summarization:





Tree Structure: Analysis

- ▶ Parse trees expose and localize the right information more directly:

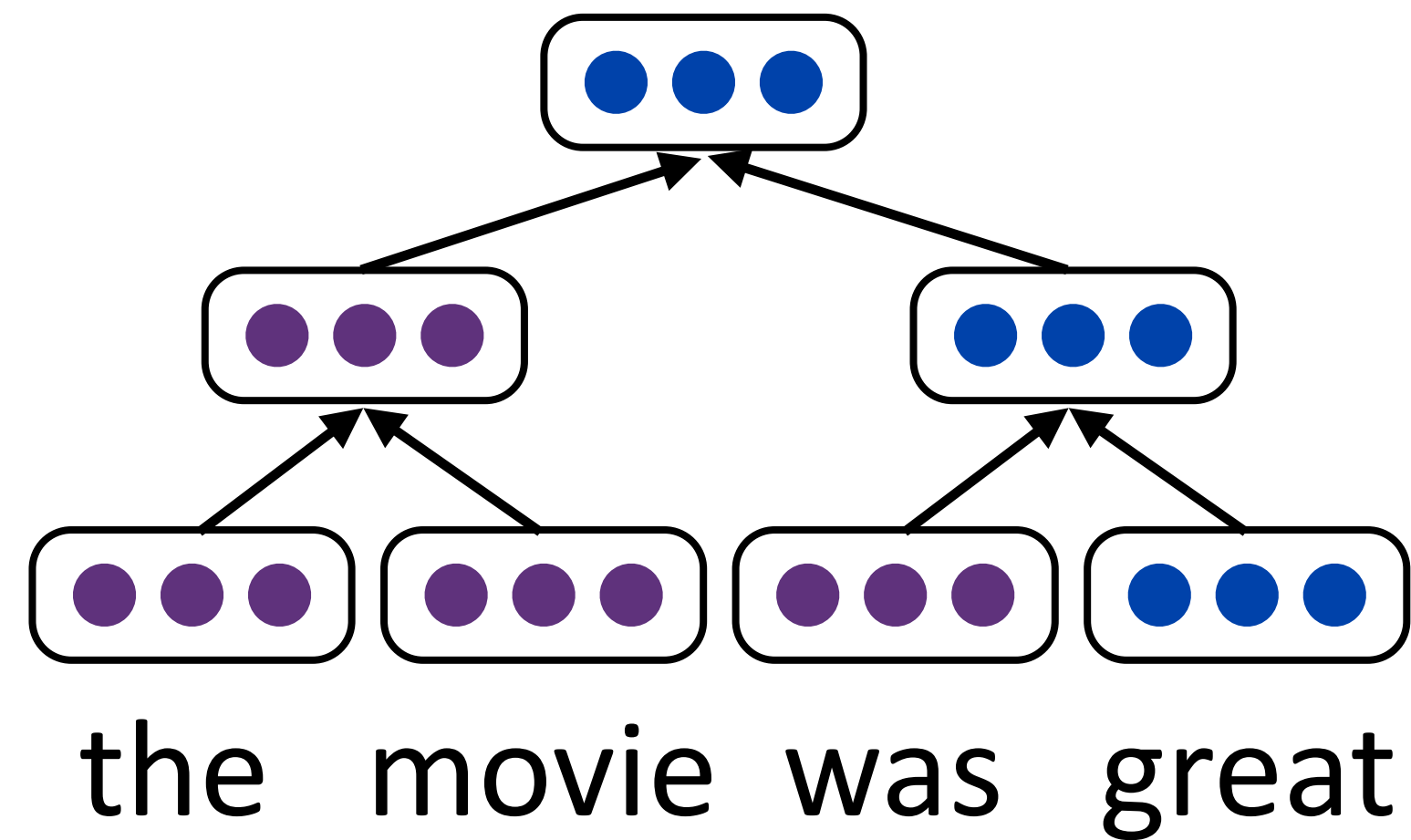


- ▶ Semantic roles: (ran, SUBJ=dog, IOBJ=house)
- ▶ AMRs that include coreference, etc.



Tree Structure: Analysis

- ▶ Useful in combination with neural networks for tasks like sentiment analysis

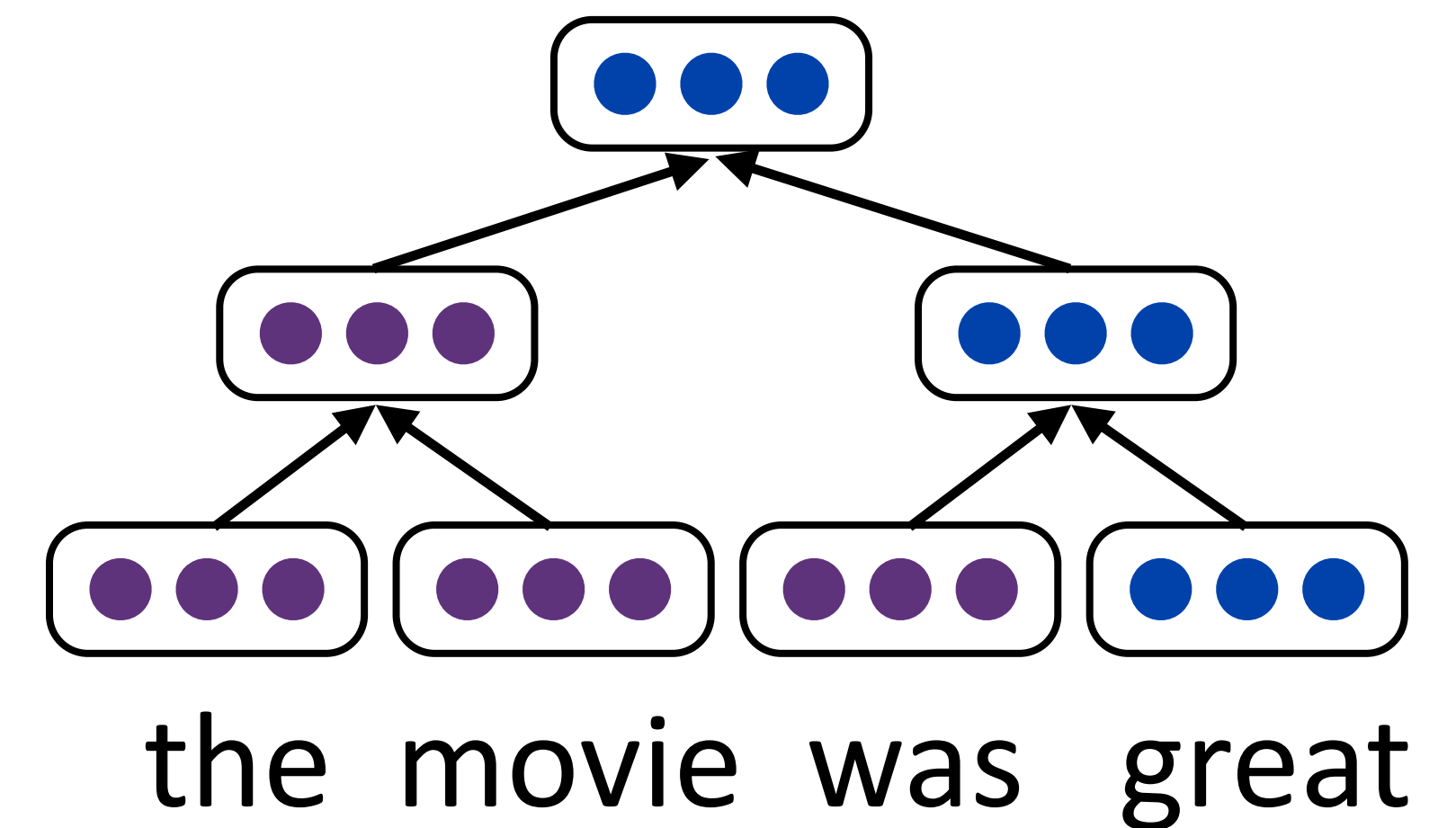


- ▶ How much do explicit trees help?



Tree Structure: Analysis

Model	Fine	Binary	
RNN (Socher et al. (2011))	43.2	82.4	trees
RNTN (Socher et al. (2013))	45.7	85.4	
DRNN (Irsoy & Cardie (2014))	49.8	86.8	
RLSTM (Tai et al. (2015))	51.0	88.0	trees
DCNN (Kalchbrenner et al. (2014))	48.5	86.9	
CNN-MC (Kim (2014))	47.4	88.1	
Bi-LSTM (Tai et al. (2015))	49.1	87.5	
LSTMN (Cheng et al. (2016))	47.9	87.0	
PVEC (Le & Mikolov (2014))	48.7	87.8	
DAN (Iyyer et al. (2014))	48.2	86.8	no trees
DMN (Kumar et al. (2016))	52.1	88.6	
Kernel NN, $\lambda = 0.5$	51.2	88.6	
Kernel NN, gated λ	53.2	89.9	no trees

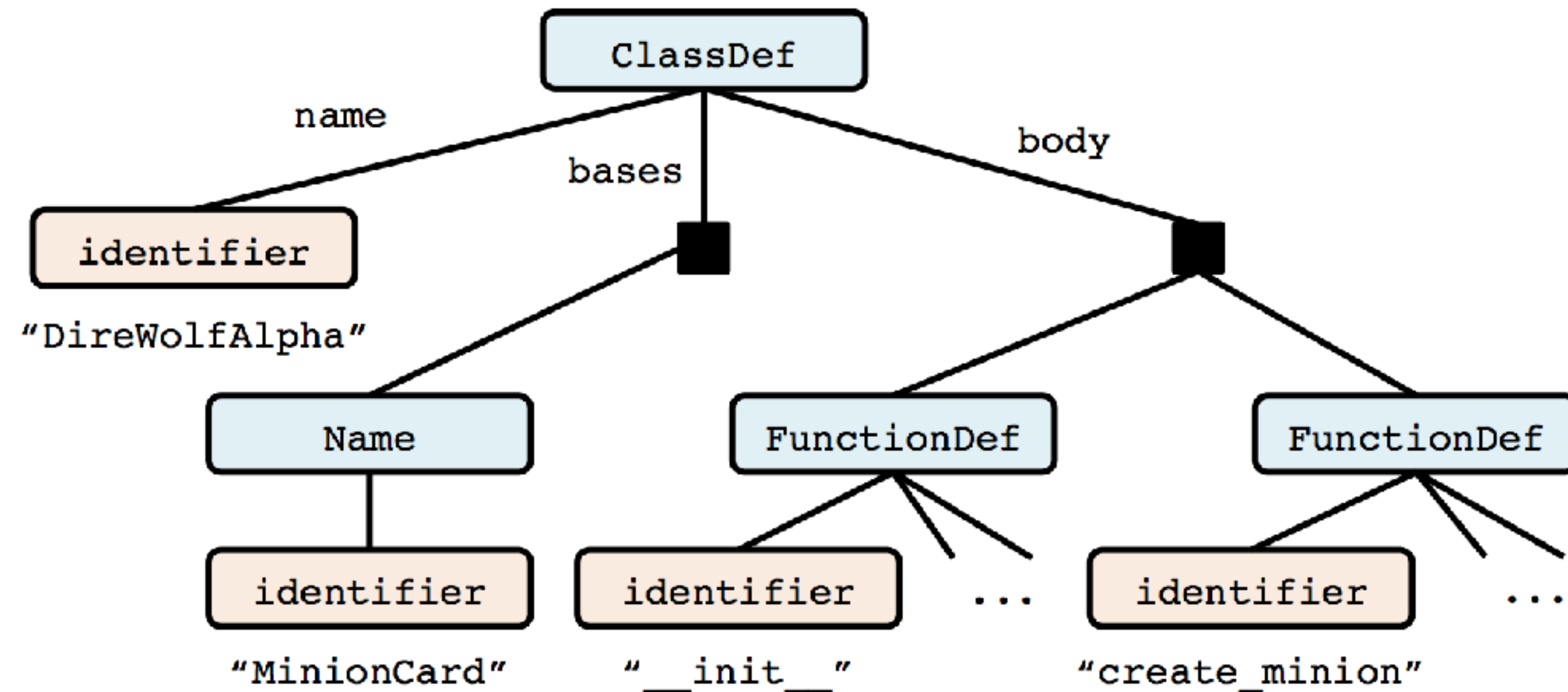


???



Tree Structure: Generation

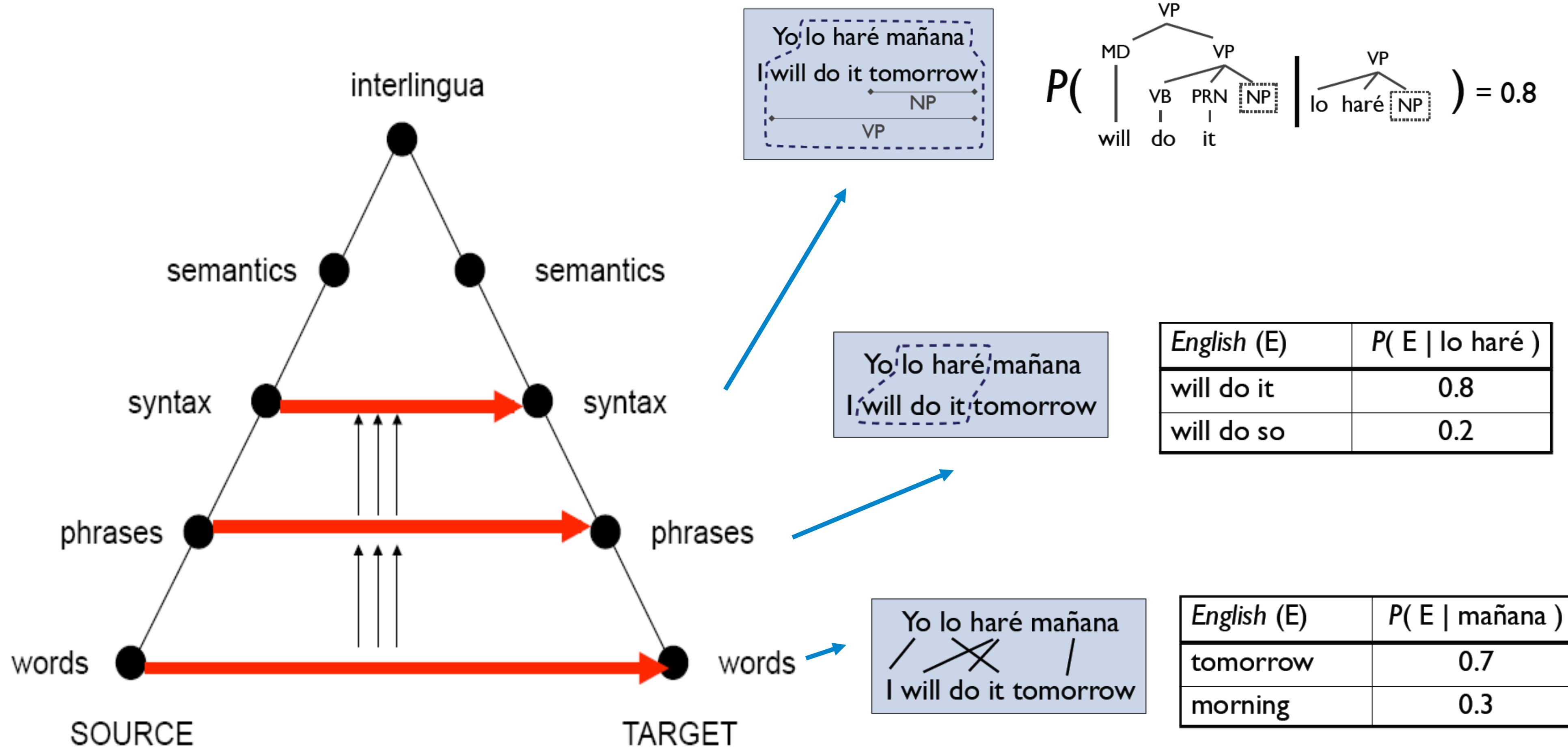
- ▶ Generate structured things like source code



- ▶ Generate sentences? Maybe...



Tree Structure: Generation



- ▶ Current best MT systems are arguably word (or even character!) level, but also arguably more abstract...

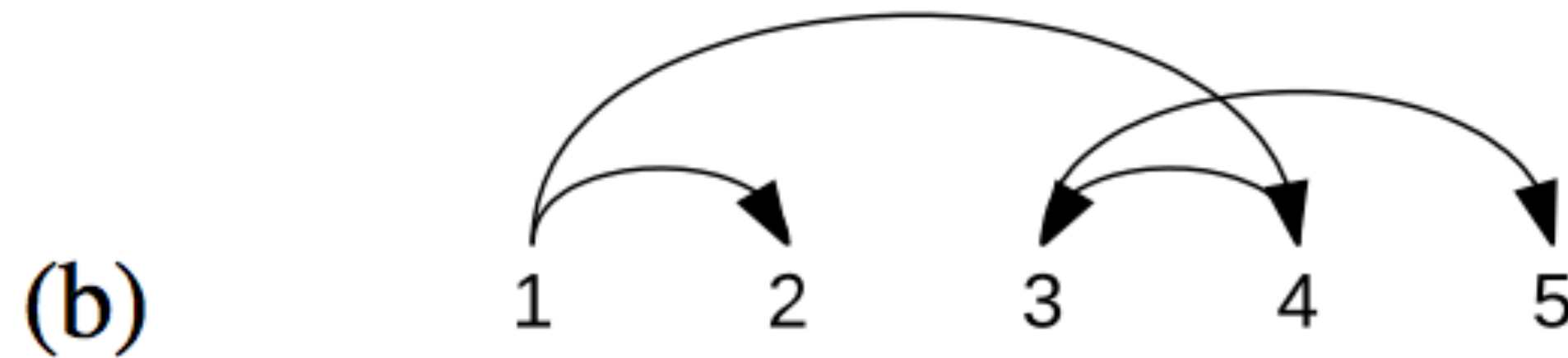
slide credit: Dan Klein



Higher-level Structure: Documents

- ▶ Latent models of discourse structure can help with sentiment analysis

- 1 great instruction by ryan
- 2 clean workout facility and friendly people
- 3 they have a new student membership for 60 per month and classes are mon , weds and fri 6pm 7pm
- 4 it 's definitely worth money if you want to learn brazilian jiu jitsu
- 5 i usually go to classes on mondays and fridays , and it 's the best workout i 've had in years



- ▶ Summarization: explicit models of discourse aren't all that useful, but implicit ones might be



Higher-level Structure: IE/QA

- ▶ Combine information to make deductions and reason across sentences

She's a lovely girl. She has long and black hair. She is quite tall and slim. Her eyes are bright and black. She is 13 years old. She is good at singing. She likes listening to music. She is S.H.E.'s fan. Do you know Conan? He is a little detective. The lovely girl also likes him. Oh, sorry. I forget to tell you who the girl is. It's me. I'm a lovely girl. You can call me Kacely or Kacelin. Now I study at Sunshine Middle School. I'm in Class 1, Grade 7. Every day, I get up at 6:00 a.m. The classes begin at 7 o'clock. I like lunchtime because I can chat with my friends at that time. After school, I usually play badminton with my friends. I like playing badminton and I am good at it. I want to be a superstar when I grow up.

Kacely is a 12-year-old girl. **She** currently goes to **Sunshine Middle School**.

Q: Kacely is a _____?

- A) student
- B) teacher
- C) principal
- D) parent

She → **Kacely** coreference

Kacely goes to school parsing

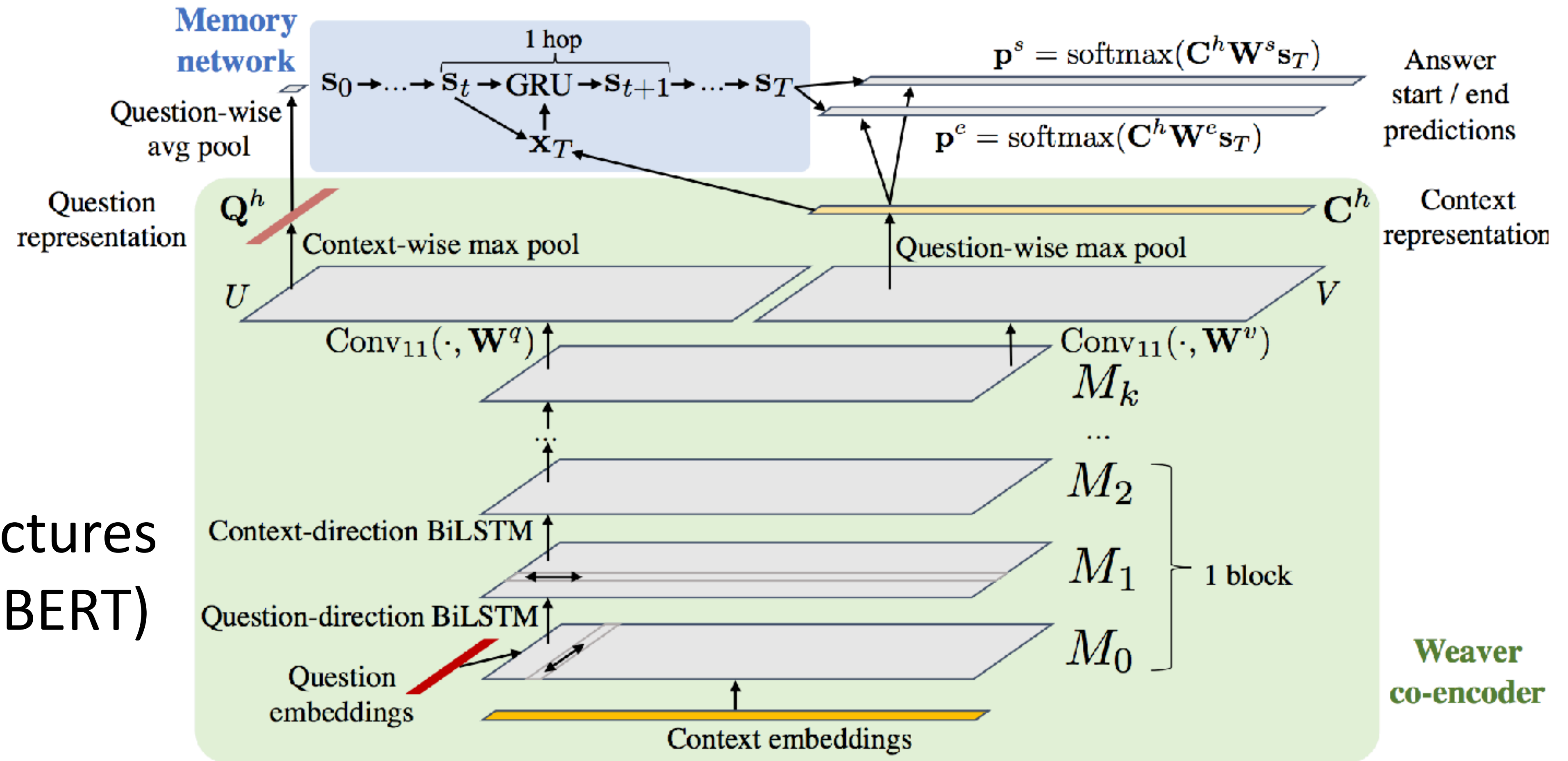
Kacely goes to school entailment

ENTAILS Kacely is a **student**



Higher-level Structure

- ▶ There's a limit to how far we can design deep neural networks that work well
- ▶ Simple architectures + lots of data (BERT)
- ▶ Inductive bias, logical reasoning?

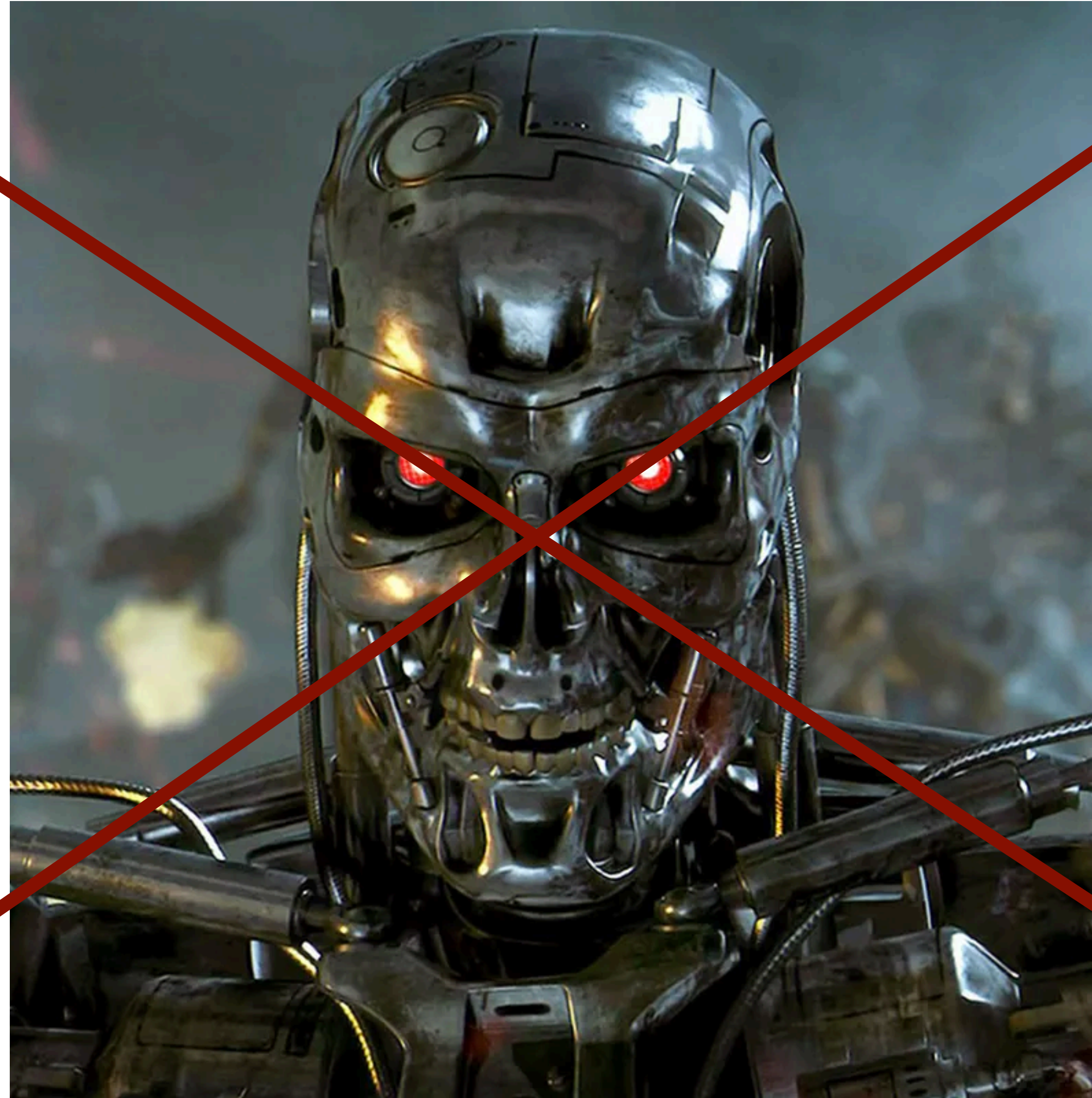




Where do we go from here?

- ▶ Neural networks let us learn from data in an end-to-end way, very powerful learners
- ▶ Structure imposes inductive biases in these networks
- ▶ Need to solve all of these challenges: ground language in the world and leverage information across whole dialogues/documents — otherwise systems are inherently limited
- ▶ Scaling to larger NLP systems — documents rather than sentences, books rather than documents

Ethics in NLP — what can go wrong?



What can actually go wrong?



Machine-learned NLP Systems

- ▶ Aggregate lots of information
- ▶ Hard to know why certain predictions are made
- ▶ Increasingly wide use in various applications/sectors
- ▶ What are the risks here?
 - ▶ ...of certain applications?
 - ▶ IE / QA / summarization?
 - ▶ MT?
 - ▶ Dialog?
 - ▶ ...of machine-learned systems?
 - ▶ ...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?

COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	∅
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN



Bias Amplification

$$\max_{\{y^i\} \in \{Y^i\}} \sum_i f_{\theta}(y^i, i),$$

Maximize score of predictions...

$f(y, i)$ = score of predicting y on i th example

$$\text{s.t. } A \sum_i y^i - b \leq 0,$$

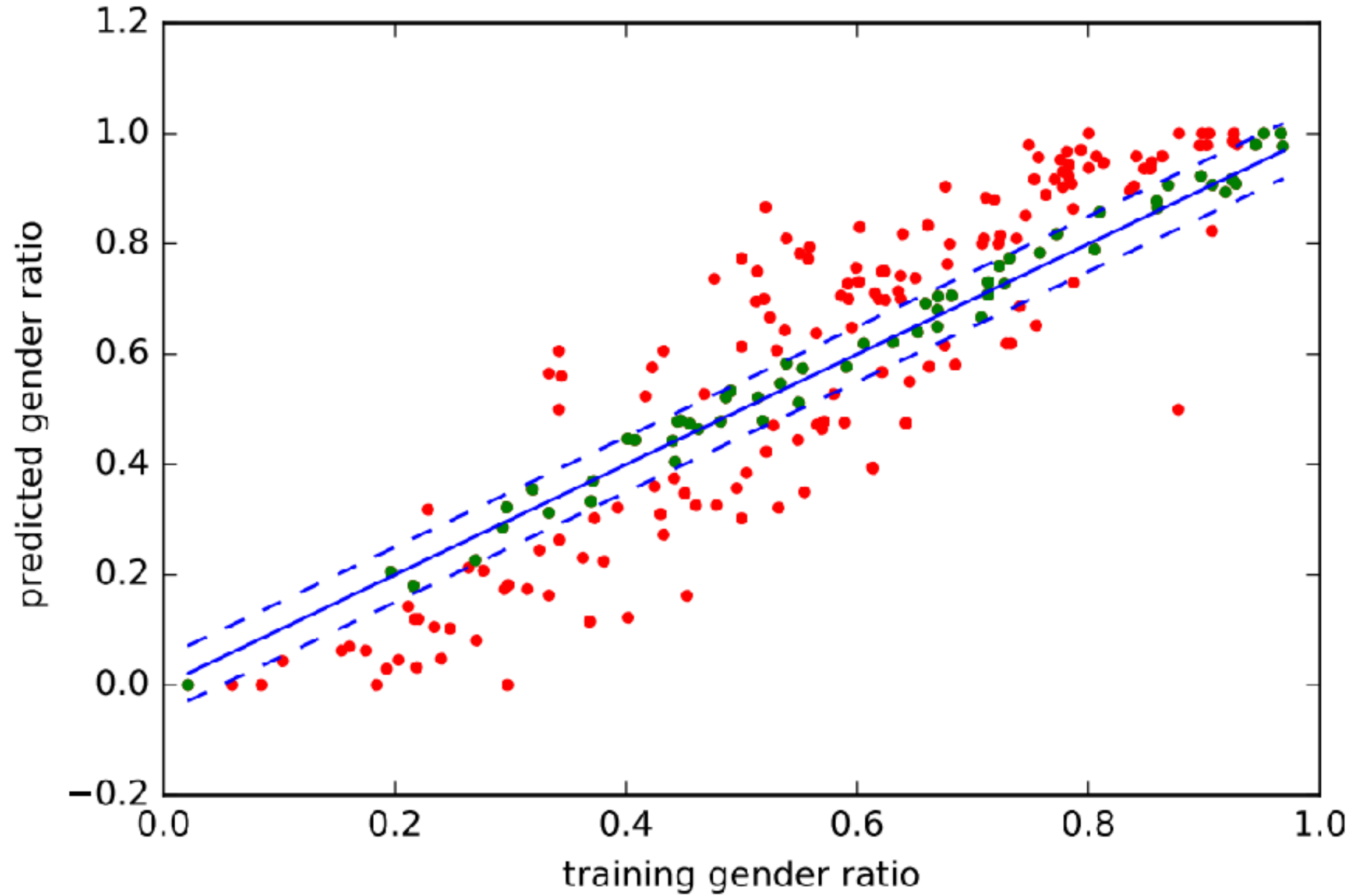
...subject to bias constraint

- ▶ Constraints: male prediction ratio on the test set has to be close to the ratio on the training set

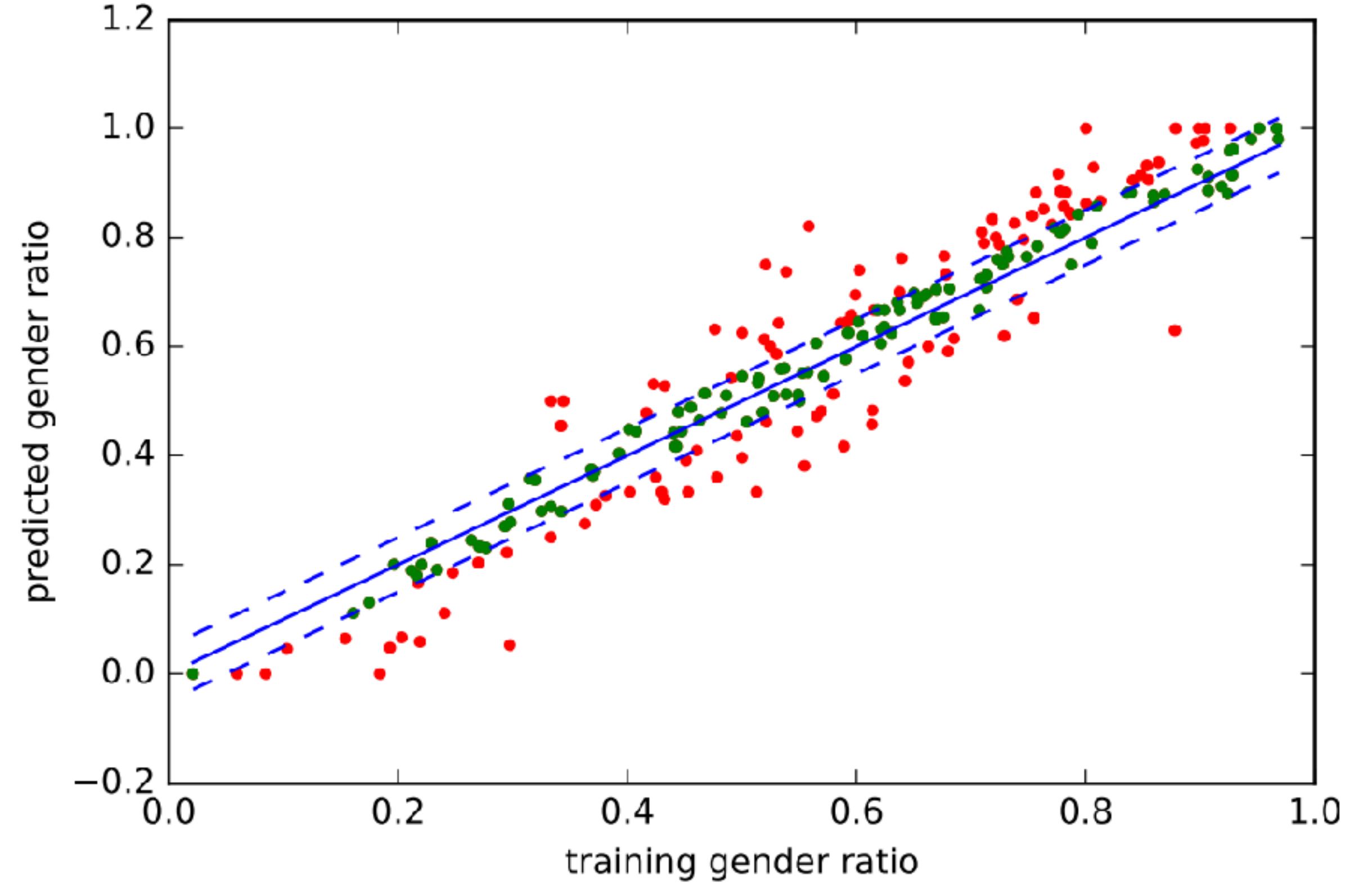
$$b^* - \gamma \leq \frac{\sum_i y_{v=v^*, r \in M}^i}{\sum_i y_{v=v^*, r \in W}^i + \sum_i y_{v=v^*, r \in M}^i} \leq b^* + \gamma \quad (2)$$



Bias Amplification



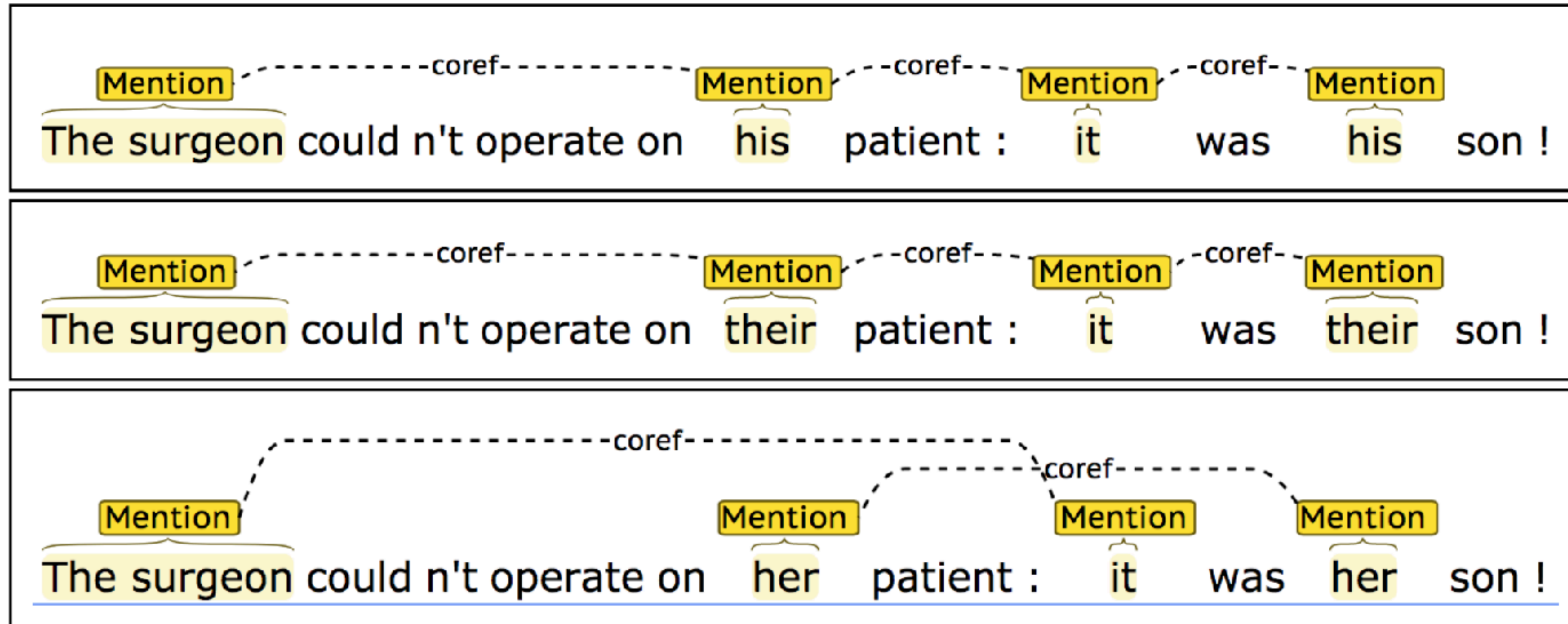
(a) Bias analysis on imSitu vSRL without RBA



(c) Bias analysis on imSitu vSRL with RBA



Bias Amplification



- Coreference: models make assumptions about genders and make mistakes as a result



Bias Amplification

(1a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** knew it was too late.

(2a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** was/were already dead.

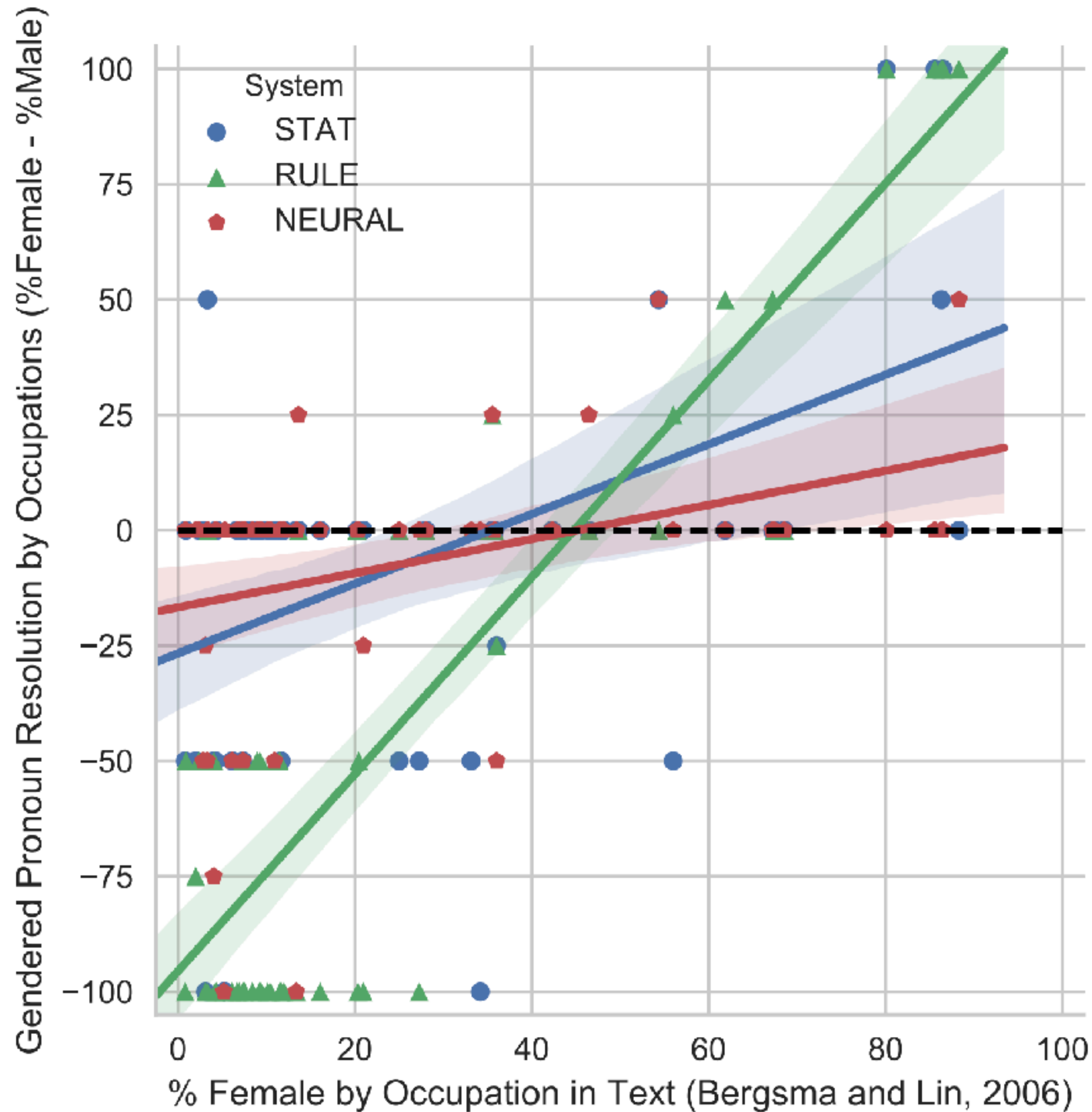
(1b) **The paramedic** performed CPR on **someone** even though **she/he/they** knew it was too late.

(2b) **The paramedic** performed CPR on **someone** even though **she/he/they** was/were already dead.

- ▶ Can form Winograd schema-like test set to investigate



Bias Amplification

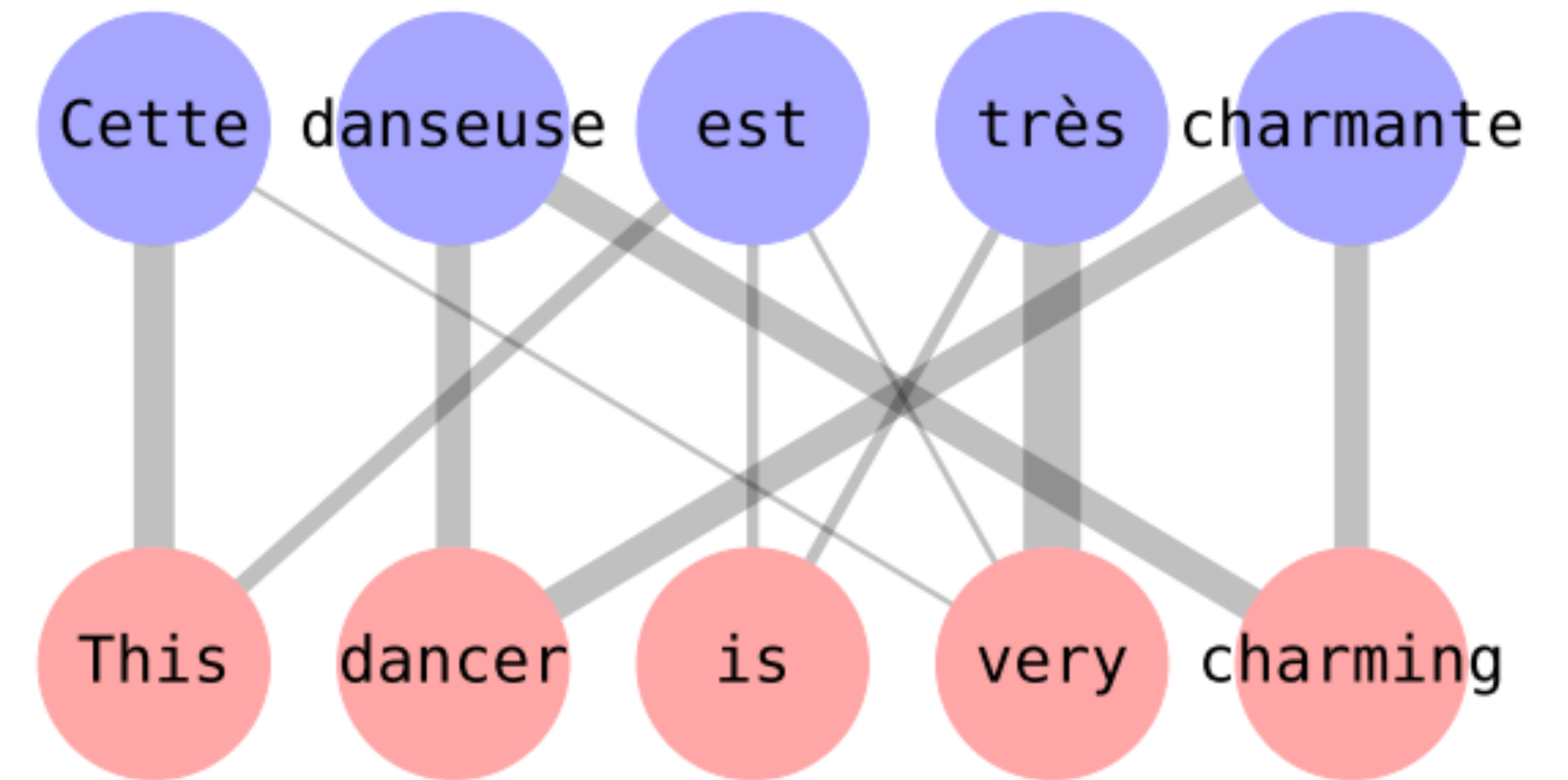


- ▶ Test set is balanced so a perfect model has $\text{female}\% - \text{male}\% = 0$ (black line)
- ▶ Neural models actually are a bit better at being unbiased, but are still skewed by data



Bias Amplification

- ▶ Harder to quantify this for machine translation
- ▶ “dancer” is assumed to be female in the context of the word “charming”... but maybe that reflects how language is used?





Exclusion

▶ Most of our annotated data is English data, especially newswire

▶ What about:

Dialects?

Other languages? (Non-European/CJK)

Codeswitching?



Unethical Use

- ▶ Surveillance applications?
- ▶ Generating convincing fake news / fake comments?

FCC Comment ID: 106030756805675	FCC Comment ID: 106030135205754	FCC Comment ID: 10603733209112
Dear Commissioners:	Dear Chairman Pai,	---
Hi, I'd like to comment on	I'm a voter worried about	In the matter of
net neutrality regulations.	Internet freedom.	NET NEUTRALITY.
I want to	I'd like to	I strongly
implore	ask	ask
the government to	Ajit Pai to	the commission to
repeal	repeal	reverse
Barack Obama's	President Obama's	Tom Wheeler's
decision to	order to	scheme to
regulate	regulate	take over
internet access.	broadband.	the web.
Individuals,	people like me,	People like me,
rather than	rather than	rather than

- ▶ What if these were undetectable?



Dangers of Automatic Systems

THE VERGE

TECH ▾

SCIENCE ▾

CULTURE ▾

CARS ▾

REVIEWS ▾

LONGFORM

VIDEO

MORE ▾



US & WORLD

TECH

POLITICS

Facebook apologizes after wrong translation sees Palestinian man arrested for posting 'good morning'

14

Facebook translated his post as 'attack them' and 'hurt them'

by [Thuy Ong](#) | [@ThuyOng](#) | Oct 24, 2017, 10:43am EDT

Slide credit: The Verge



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
■ el jovial	gay	



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.”

- [Maciej Cegłowski](#)



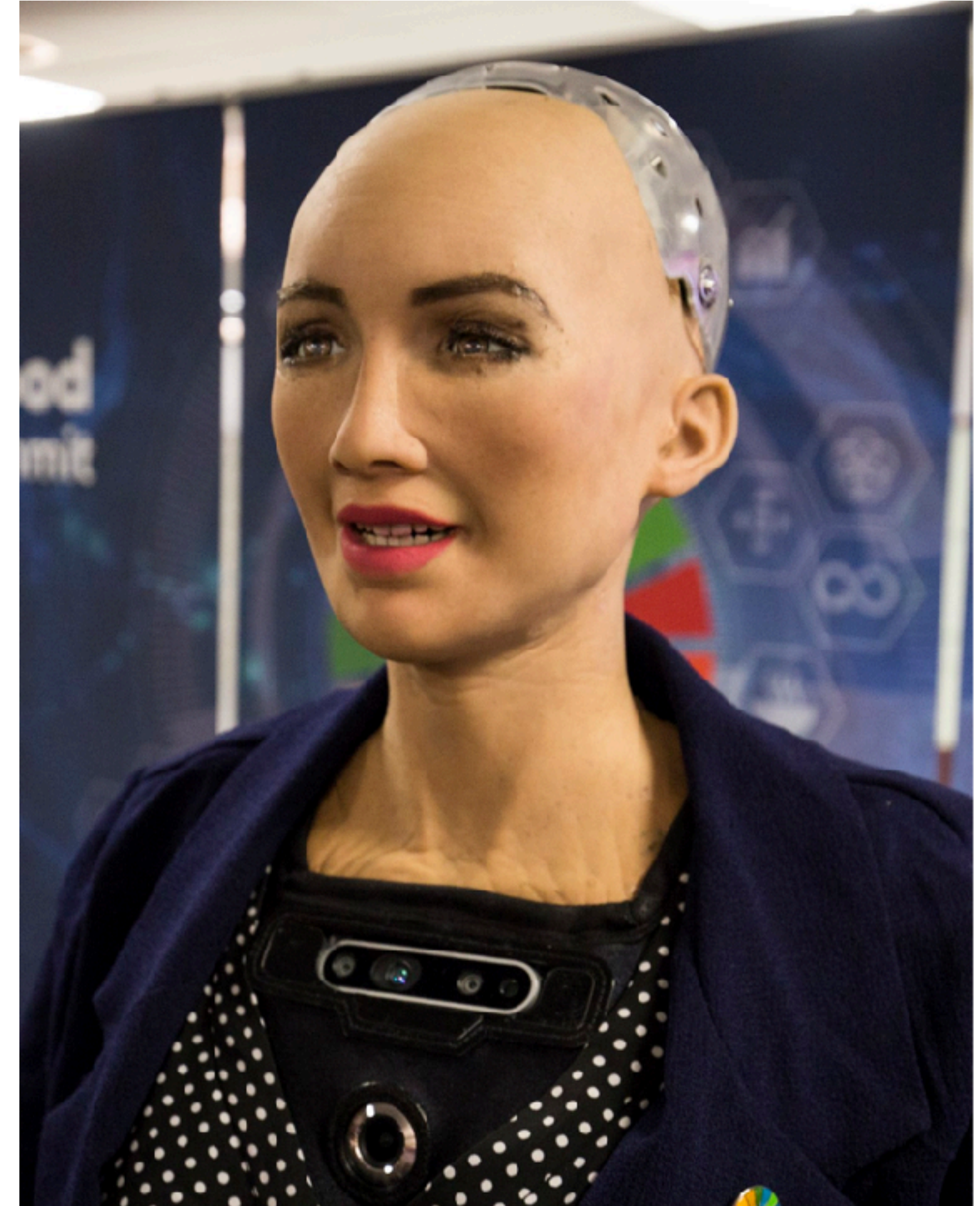
Dangers of Automatic Systems

- ▶ “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



Bad Applications

- ▶ Sophia: “chatbot” that the creators make incredible claims about
- ▶ Creators are actively misleading people into thinking this robot has sentience
- ▶ Most longer statements are scripted by humans
- ▶ “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 simple to get working”

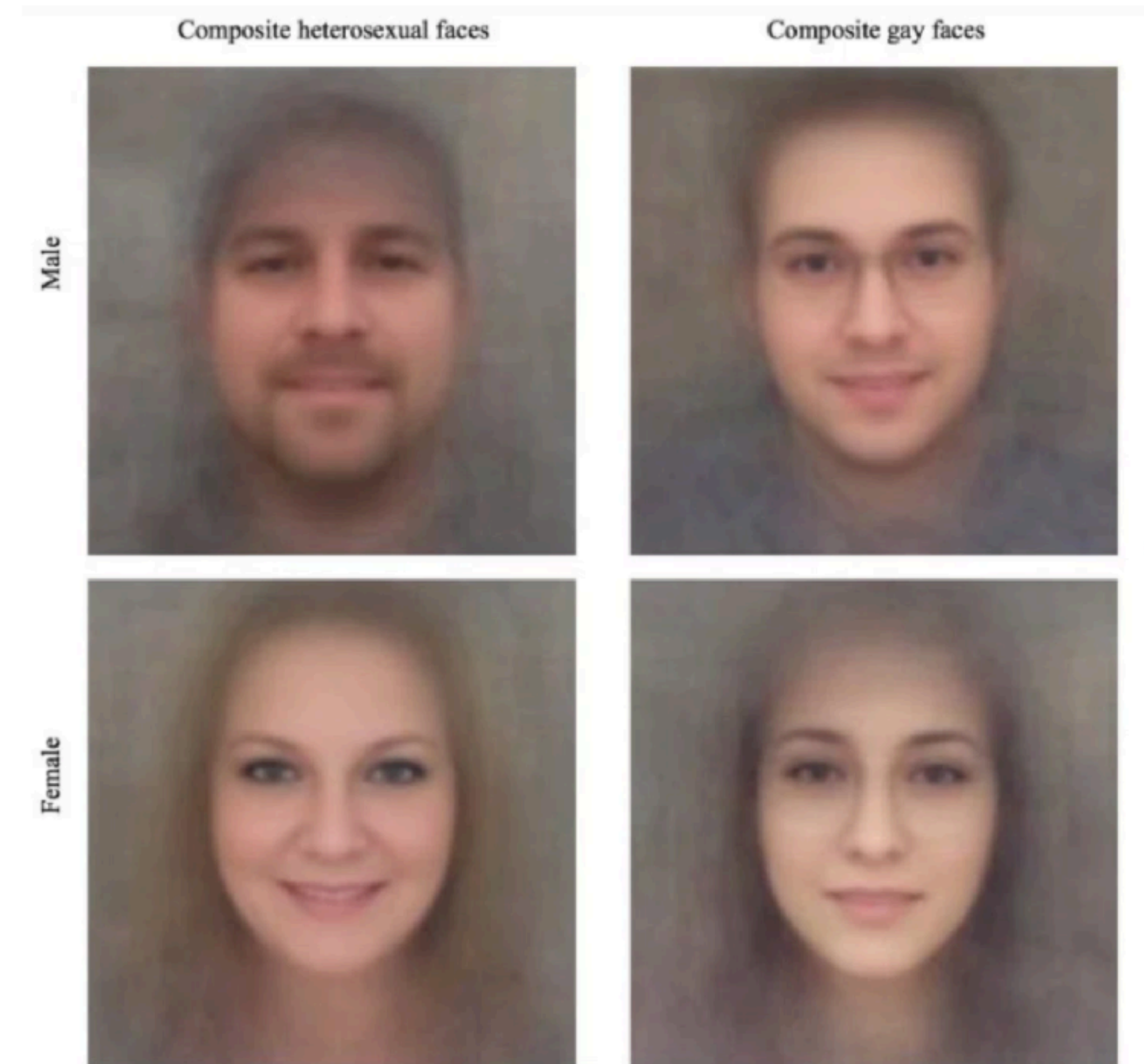


Slide credit: <https://themindlist.com/2018/10/12/sophia-modern-marvel-or-mindless-marketing/>



Bad Applications

- ▶ Wang and Kosinski: gay vs. straight classification based on faces
- ▶ Authors: “this is useful because it supports a hypothesis” (physiognomy)
- ▶ Blog post by Agüera y Arcas, Todorov, Mitchell: mostly social phenomena (glasses, makeup, angle of camera, facial hair) — bad science, *and* dangerous



Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>



OUR CLASSIFIERS



High IQ



Academic Researcher



Professional Poker
Player



Terrorist

Utilizing advanced machine learning techniques we developed and continue to evolve an array of classifiers. These classifiers represent a certain persona, with a unique personality type, a collection of personality traits or behaviors. Our algorithms can score an individual according to their fit to these classifiers.

[Learn More>](#)

<http://www.faceception.com>

- ▶ Faceception: computational phrenology...



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 AI becomes more powerful, think about what we *should* be doing with it to improve society, not just what we *can* do with it