

Lecture 18: Wrapup + Ethics

Alan Ritter

(many slides from Greg Durrett)

Administrivia

- ▶ Final project reports due Wednesday 5/5
- ▶ Please fill out the course/instructor opinion survey (CIOS) if you haven't already!

This Lecture

- ▶ Multilingual Models
- ▶ Ethics in NLP

NLP in other languages

- ▶ Other languages present some challenges not seen in English at all!
- ▶ Some of our algorithms have been specified to English
 - ▶ Neural methods are typically tuned to English-scale resources, may not be the best for other languages where less data is available
- ▶ Question:
 - 1) What other phenomena / challenges do we need to solve?
 - 2) How can we leverage existing resources to do better in other languages without just annotating massive data?

Morphology

What is morphology?

- ▶ Study of how words form
- ▶ Derivational morphology: create a new *lexeme* from a base
 - estrangle (v) => estrangement (n)
 - become (v) => unbecoming (adj)
 - ▶ May not be totally regular: enflame => inflammable
- ▶ Inflectional morphology: word is inflected based on its context
 - I become / she becomes
 - ▶ Mostly applies to verbs and nouns

Morphological Inflection

- In English: I arrive you arrive he/she/it arrives [X] arrived
- we arrive you arrive they arrive

- In French:

		singular			plural		
		first	second	third	first	second	third
indicative		je (j')	tu	il, elle	nous	vous	ils, elles
(simple tenses)	present	arrive /a.viv/	arrives /a.viv/	arrive /a.viv/	arrivons /a.viv.vɔ/	arrivez /a.viv.ve/	arrivent /a.viv.v/
	imperfect	arrivais /a.viv.vɛ/	arrivais /a.viv.vɛ/	arrivait /a.viv.vɛ/	arrivions /a.viv.vjɔ/	arriviez /a.viv.vje/	arrivaient /a.viv.vɛ/
	past historic ²	arrivai /a.viv.vɛ/	arrivas /a.viv.va/	arriva /a.viv.va/	arrivâmes /a.viv.vam/	arrivâtes /a.viv.vat/	arrivèrent /a.viv.vɛ/
	future	arriverai /a.viv.vɛ/	arriveras /a.viv.va/	arrivera /a.viv.va/	arriverons /a.viv.vɔ/	arrivez /a.viv.vɛ/	arriveront /a.viv.vɔ/
	conditional	arriverais /a.viv.vɛ/	arriverais /a.viv.vɛ/	arriverait /a.viv.vɛ/	arriverions /a.viv.vɛ/	arriveriez /a.viv.vɛ/	arriveraient /a.viv.vɛ/

Noun Inflection

- ▶ Not just verbs either; gender, number, case complicate things

Declension of Kind						[hide ▲]
	singular			plural		
	indef.	def.	noun	def.	noun	
nominative	ein	das	Kind	die	Kinder	
genitive	eines	des	Kindes, Kinds	der	Kinder	
dative	einem	dem	Kind, Kinde ¹	den	Kindern	
accusative	ein	das	Kind	die	Kinder	

- ▶ Nominative: I/he/she, accusative: me/him/her, genitive: mine/his/hers
- ▶ Dative: merged with accusative in English, shows recipient of something

I taught the children <=> Ich unterrichte die Kinder

I give the children a book <=> Ich gebe den Kindern ein Buch

Agglutinating Languages

- ▶ Finnish/Hungarian (Finno-Ugric), also Turkish: what a preposition would do in English is instead part of the verb (*hug*)

	active	passive
1st	halata	
long 1st ²	halatakseen	
2nd	inessive¹ halatessa instructive halaten	halattaessa —
3rd	inessive halaamassa elative halaamasta illative halaamaan adessive halaamalla abessive halaamatta instructive halaaman	— — — — — — halattaman
4th	nominative halaaminen partitive halaamista	
5th ²	halaamaisillaan	

illative: “into”

adessive: “on”

- ▶ Many possible forms – and in newswire data, only a few are observed

indicative mood				
present tense	positive	negative	perfect	
person	1st sing.	en halaa	1st sing.	positive
	2nd sing.	ei halaa	2nd sing.	en ole halannut
	3rd sing.	ei halaa	3rd sing.	et ole halannut
	1st plur.	emme halaa	1st plur.	ei ole halanneet
	2nd plur.	ette halaa	2nd plur.	emme ole halanneet
	3rd plur.	evät halaa	3rd plur.	ette ole halanneet
	passive	halataan	passive	ovat halanneet
		ei halata		evät ole halanneet
		on halatu		ei ole halatu
past tense	positive	negative	perfect	
person	1st sing.	halasin	1st sing.	positive
	2nd sing.	ei halasnut	2nd sing.	en ole halannut
	3rd sing.	ei halasni	3rd sing.	et ole halannut
	1st plur.	halasimme	1st plur.	ei ole halanneet
	2nd plur.	ette halasneet	2nd plur.	emme ole halanneet
	3rd plur.	halasivat	3rd plur.	ette ole halanneet
	passive	halastettiin	passive	olivat halanneet
		ei halastui		evät olivat halanneet
		on halastu		ei ole halastu
conditional mood	positive	negative	perfect	
person	1st sing.	halaisin	1st sing.	positive
	2nd sing.	ei halaisi	2nd sing.	en ole halannut
	3rd sing.	ei halaisi	3rd sing.	et ole halannut
	1st plur.	halaisimme	1st plur.	ei ole halanneet
	2nd plur.	ette halaisneet	2nd plur.	emme ole halanneet
	3rd plur.	halaisivat	3rd plur.	ette ole halanneet
	passive	halasttaisiin	passive	olivat halanneet
		ei halasttaisi		evät olivat halanneet
		on halasttu		ei ole halasttu
imperative mood	positive	negative	perfect	
person	1st sing.	—	1st sing.	positive
	2nd sing.	halaa	2nd sing.	en ole halannut
	3rd sing.	älä halaa	3rd sing.	et ole halannut
	1st plur.	halakoon	1st plur.	ei ole halanneet
	2nd plur.	äläkää halaa	2nd plur.	emme ole halanneet
	3rd plur.	halakoot	3rd plur.	ette ole halanneet
	passive	halattakoont	passive	olivat halanneet
		äläkää halattako		evät olivat halanneet
		on halattu		ei ole halattu
soteric mood	positive	negative	perfect	
person	1st sing.	—	1st sing.	positive
	2nd sing.	halanen	2nd sing.	en ole halannut
	3rd sing.	en halanen	3rd sing.	et ole halannut
	1st plur.	halanleet	1st plur.	ei ole halanneet
	2nd plur.	en halanleet	2nd plur.	emme ole halanneet
	3rd plur.	halannevat	3rd plur.	ette ole halanneet
	passive	halattanevat	passive	olivat halanneet
		evät halannevat		evät ole halannevat
		on halattu		ei ole halattu
nominal forms	active	passive	perfect	
definitives	halata	halataessa	person	positive
1st	halatas	halataseen		negative
long 1st ²	halataseen	—		
2nd	inessive ¹	halatessa	1st sing.	—
	instructive	halaten	2nd sing.	ole halannut
	elative	—	3rd sing.	älä halannut
	illative	—	1st plur.	äläkää halannut
	adessive	—	2nd plur.	äläkää älä halannut
	abessive	—	3rd plur.	äläkää äläkää halannut
	instructive	—	passive	okkoon halattu
3rd	nominative	halaaminen		
	partitive	halaamista		
5th ²		halaamaisillaan		

halata: “hug”

1 Used only with a possessive suffix.
2 Used only with a possessive suffix; this is the form for the third-person singular and third-person plural.
3 Does not exist in the case of intransitive verbs. Do not confuse with nouns formed with the -ma suffix.

Morphologically-Rich Languages

- ▶ Many languages spoken all over the world have much richer morphology than English
- ▶ CoNLL 2006 / 2007: dependency parsing + morphological analyses for ~15 mostly Indo-European languages
- ▶ SPMRL shared tasks (2013-2014): Syntactic Parsing of Morphologically-Rich Languages
- ▶ Universal Dependencies project
- ▶ Word piece / byte-pair encoding models for MT are pretty good at handling these if there's enough data

Morphological Analysis

- ▶ In English, lexical features on words and word vectors are pretty effective
- ▶ In other languages, **lots** more unseen words due to rich morphology!
- ▶ When we're building systems, we probably want to know base form + morphological features explicitly
- ▶ How to do this kind of *morphological analysis*?

Morphological Analysis: Hungarian

But the government does not recommend reducing taxes.

Ám a kormány egyetlen adó csökkentését sem javasolja .

n=singular|case=nominative|proper=no
deg=positive|n=singular|case=nominative
n=singular|case=nominative|proper=no
n=singular|case=accusative|proper=no|pperson=3rd|pnumber=singular
mood=indicative|t=present|p=3rd|n=singular|def=yes

Morphological Analysis

- ▶ Given a word in context, need to predict what its morphological features are
- ▶ Basic approach: combines two modules:
 - ▶ Lexicon: tells you what possibilities are for the word
 - ▶ Analyzer: statistical model that disambiguates
- ▶ Models are largely CRF-like: score morphological features in context
- ▶ Lots of work on Arabic inflection (high amounts of ambiguity)

Morphological Inflection

- ▶ Inverse task of analysis: given base form + features, inflect the word
- ▶ Hard for unknown words — need models that generalize

w i n d e n →

conjugation of <i>winden</i>						[hide ▲]
		infinitive		winden		
		present participle		windend		
		past participle		gewunden		
		auxiliary		haben		
present	indicative			subjunctive		
	ich <i>winde</i>	wir <i>winden</i>	i	ich <i>winde</i>	wir <i>winden</i>	
	du <i>windest</i>	ihr <i>windet</i>		du <i>windest</i>	ihr <i>windet</i>	
preterite	er <i>windet</i>	sie <i>winden</i>		er <i>winde</i>	sie <i>winden</i>	
	ich <i>wand</i>	wir <i>wanden</i>	ii	ich <i>wände</i>	wir <i>wänden</i>	
	du <i>wandest</i>	ihr <i>wandet</i>		du <i>wändest</i>	ihr <i>wändet</i>	
imperative		er <i>wand</i>	sie <i>wanden</i>	er <i>wände</i>	sie <i>wänden</i>	
composed forms of <i>winden</i>						[show ▽]

Chinese Word Segmentation

- ▶ Word segmentation:
some languages
including Chinese are
totally untokenized
- ▶ LSTMs over character
embeddings / character
bigram embeddings to
predict word boundaries
- ▶ Having the right
segmentation can help
machine translation

冬天 (winter), 能 (can) 穿 (wear) 多 少
(amount) 穿 (wear) 多 少 (amount); 夏天
(summer), 能 (can) 穿 (wear) 多 (more) 少
(little) 穿 (wear) 多 (more) 少 (little).

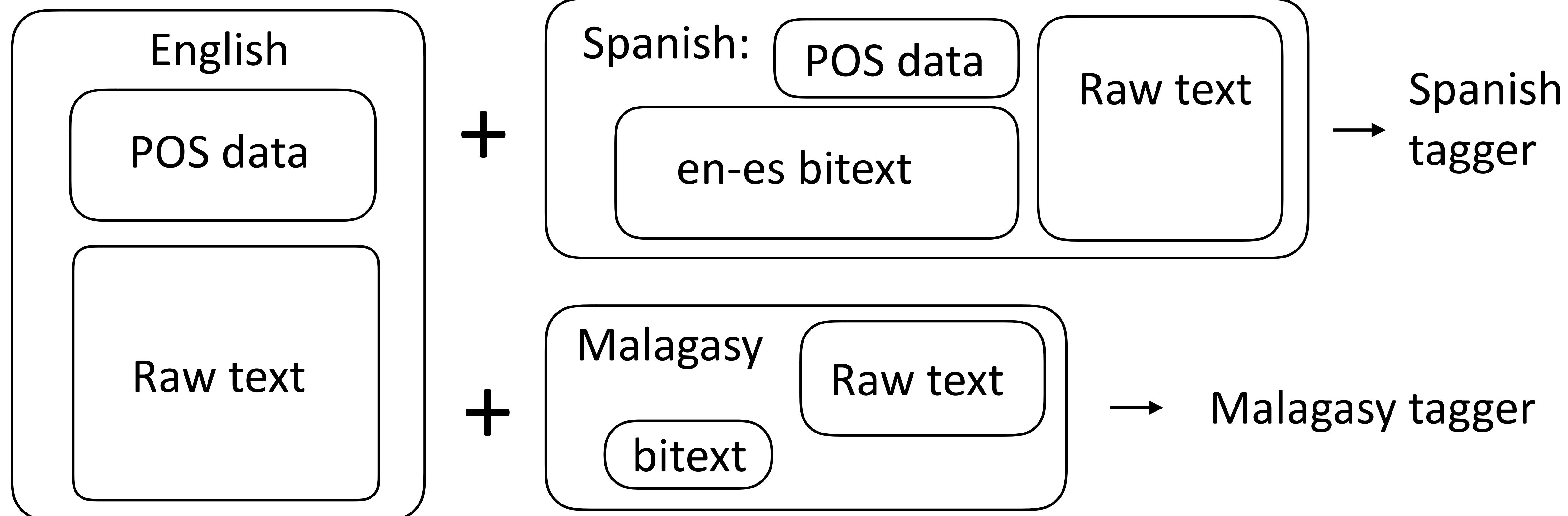
Without the word “夏天 (summer)” or “冬天
(winter)”, it is difficult to segment the phrase “能
穿多少穿多少”.

- separating nouns and pre-modifying adjectives:
高血压 (*high blood pressure*)
→ 高(*high*) 血压(*blood pressure*)
- separating compound nouns:
内政部 (*Department of Internal Affairs*)
→ 内政(*Internal Affairs*) 部(*Department*).

Cross-Lingual Tagging

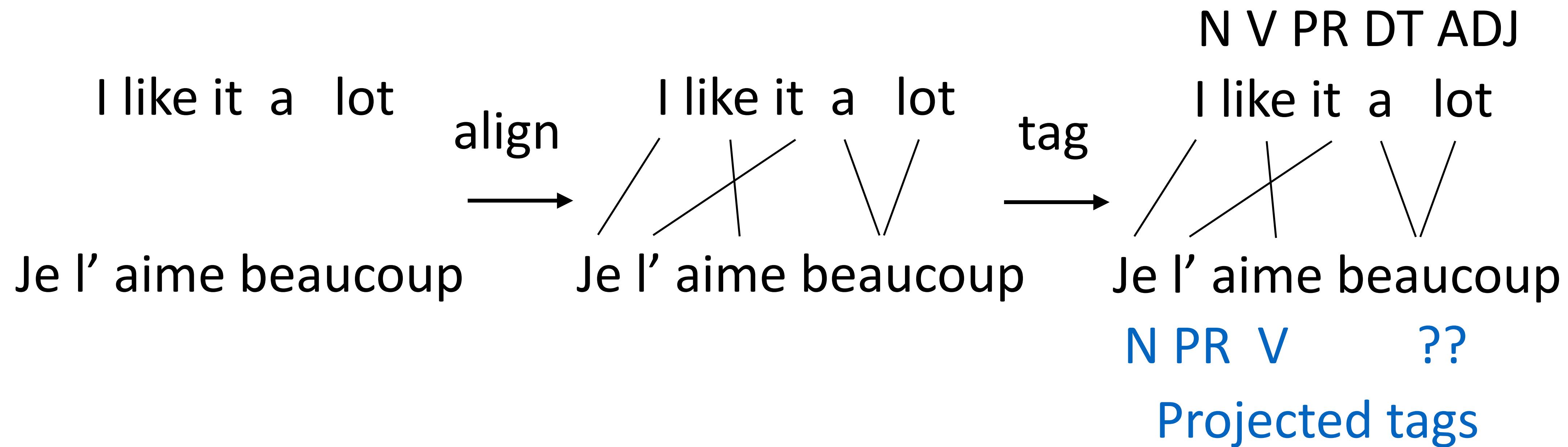
Cross-Lingual Tagging

- ▶ Labeling POS datasets is expensive
- ▶ Can we transfer annotation from *high-resource* languages (English, etc.) to *low-resource* languages?



Cross-Lingual Tagging

- ▶ Can we leverage word alignment here?



- ▶ Tag with English tagger, project across bitext, train French tagger?
Works pretty well

Das and Petrov (2011)

Cross-Lingual Word Representations

Multilingual Embeddings

- ▶ Input: corpora in many languages. Output: embeddings where similar words *in different languages* have similar embeddings

I have an apple
47 24 18 427

ID: 24
ai have

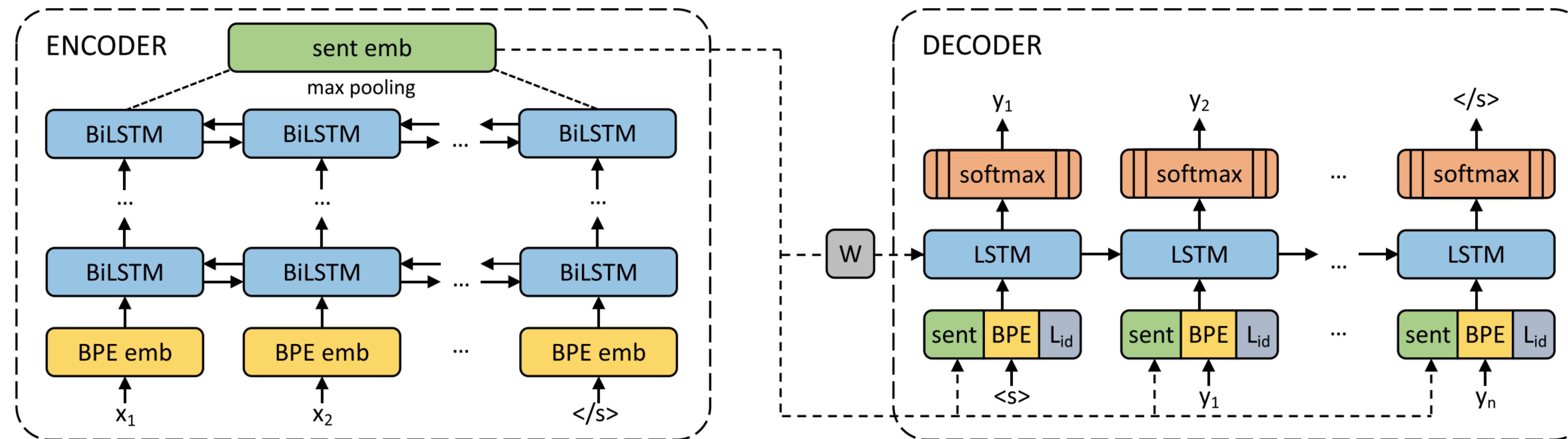
J' ai des oranges
47 24 89 1981

ID: 47
I Je J'

- ▶ multiCluster: use bilingual dictionaries to form clusters of words that are translations of one another, replace corpora with cluster IDs, train “monolingual” embeddings over all these corpora

- ▶ Works okay but not all that well

Multilingual Sentence Embeddings



- ▶ Form BPE vocabulary over all corpora (50k merges); will include characters from every script
- ▶ Take a bunch of bitexts and train an MT model between a bunch of language pairs with shared parameters, use W as sentence embeddings

Artetxe et al. (2019)

Multilingual Sentence Embeddings

	EN	EN → XX														
		fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	sw	ur	
Zero-Shot Transfer, one NLI system for all languages:																
Conneau et al. (2018b)	X-BiLSTM	73.7	67.7	68.7	67.7	68.9	67.9	65.4	64.2	64.8	66.4	64.1	65.8	64.1	55.7	58.4
BERT uncased*	X-CBOW	64.5	60.3	60.7	61.0	60.5	60.4	57.8	58.7	57.5	58.8	56.9	58.8	56.3	50.4	52.2
BERT uncased*	Transformer	<u>81.4</u>	–	<u>74.3</u>	70.5	–	–	–	–	62.1	–	–	63.8	–	–	58.3
Proposed method	BiLSTM	73.9	71.9	72.9	<u>72.6</u>	72.8	74.2	72.1	69.7	71.4	72.0	69.2	<u>71.4</u>	65.5	62.2	<u>61.0</u>

- ▶ Train a system for NLI (entailment/neutral/contradiction of a sentence pair) on English and evaluate on other languages

Multilingual BERT

- ▶ Take top 104 Wikipedias, train BERT on all of them simultaneously
- ▶ What does this look like?

Beethoven may have proposed unsuccessfully to Therese Malfatti, the supposed dedicatee of "Für Elise"; his status as a commoner may again have interfered with those plans.

当人们在马尔法蒂身后发现这部小曲的手稿时，便误认为上面写的是“Für Elise”（即《给爱丽丝》）[51]。

Китái (официально – Китáiская Нарóдная Респúблика, сокращённо – КНР; кит. трад. 中華人民共和國, упр. 中华人民共和国, пиньинь: Zhōnghuá Rénmín Gònghéguó, палл.: Чжунхуа Жэньминь Гүнхэго) – государство в Восточной Аз

Devlin et al. (2019)

Multilingual BERT: Results

Fine-tuning \ Eval	EN	DE	NL	ES
EN	90.70	69.74	77.36	73.59
DE	73.83	82.00	76.25	70.03
NL	65.46	65.68	89.86	72.10
ES	65.38	59.40	64.39	87.18

Table 1: NER F1 results on the CoNLL data.

Fine-tuning \ Eval	EN	DE	ES	IT
EN	96.82	89.40	85.91	91.60
DE	83.99	93.99	86.32	88.39
ES	81.64	88.87	96.71	93.71
IT	86.79	87.82	91.28	98.11

Table 2: POS accuracy on a subset of UD languages.

- ▶ Can transfer BERT directly across languages with some success
- ▶ ...but this evaluation is on languages that all share an alphabet

Multilingual BERT: Results

	HI	UR		EN	BG	JA
HI	97.1	85.9	EN	96.8	87.1	49.4
UR	91.1	93.8	BG	82.2	98.9	51.6
			JA	57.4	67.2	96.5

Table 4: POS accuracy on the UD test set for languages with different scripts. Row=fine-tuning, column=eval.

- ▶ Urdu (Arabic/Nastaliq script) => Hindi (Devanagari). Transfers well despite different alphabets!
- ▶ Japanese => English: different script and very different syntax

Scaling Up: XLM-R

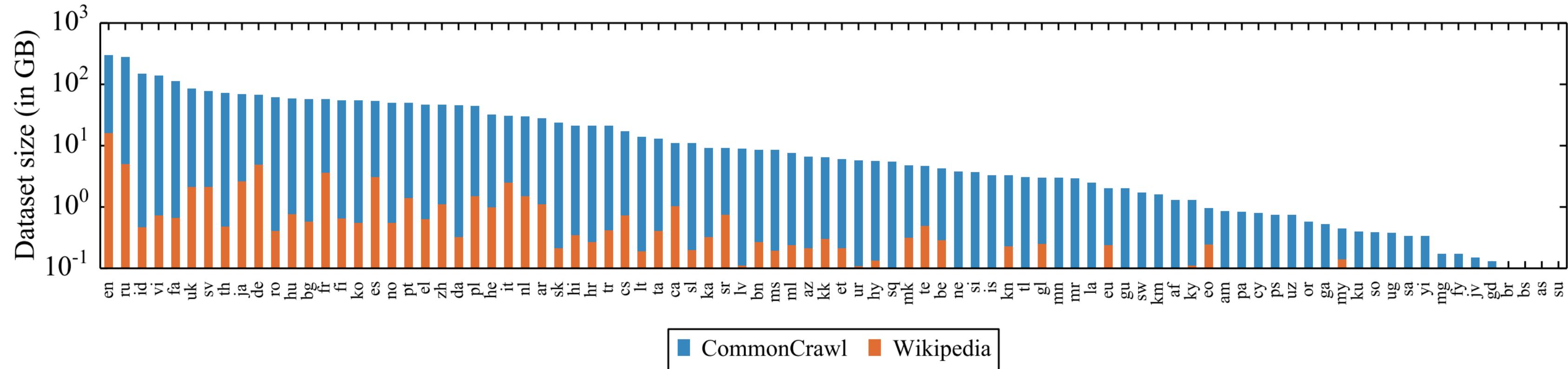


Figure 1: Amount of data in GiB (log-scale) for the 88 languages that appear in both the Wiki-100 corpus used for mBERT and XLM-100, and the CC-100 used for XLM-R. CC-100 increases the amount of data by several orders of magnitude, in particular for low-resource languages.

- ▶ Larger “Common Crawl” dataset, better performance than mBERT
- ▶ Low-resource languages benefit from training on other languages
- ▶ High-resource languages see a small performance hit, but not much

Scaling Up: Benchmarks

Task	Corpus	Train	Dev	Test	Test sets	Lang.	Task
Classification	XNLI	392,702	2,490	5,010	translations	15	NLI
	PAWS-X	49,401	2,000	2,000	translations	7	Paraphrase
Struct. pred.	POS	21,253	3,974	47-20,436	ind. annot.	33 (90)	POS
	NER	20,000	10,000	1,000-10,000	ind. annot.	40 (176)	NER
QA	XQuAD	87,599	34,726	1,190	translations	11	Span extraction
	MLQA			4,517–11,590	translations	7	Span extraction
	TyDiQA-GoldP			323–2,719	ind. annot.	9	Span extraction
Retrieval	BUCC	-	-	1,896–14,330	-	5	Sent. retrieval
	Tatoeba	-	-	1,000	-	33 (122)	Sent. retrieval

- ▶ Many of these datasets are translations of base datasets, not originally annotated in those languages
- ▶ Exceptions: POS, NER, TyDiQA

TyDiQA

- ▶ Typologically-diverse QA dataset
- ▶ Annotators write questions based on very short snippets of articles; answers may or may not exist, fetched from elsewhere in Wikipedia

Q: Как далеко Уран от
how far Uranus-SG.NOM from
Земл-и?
Earth-SG.GEN?

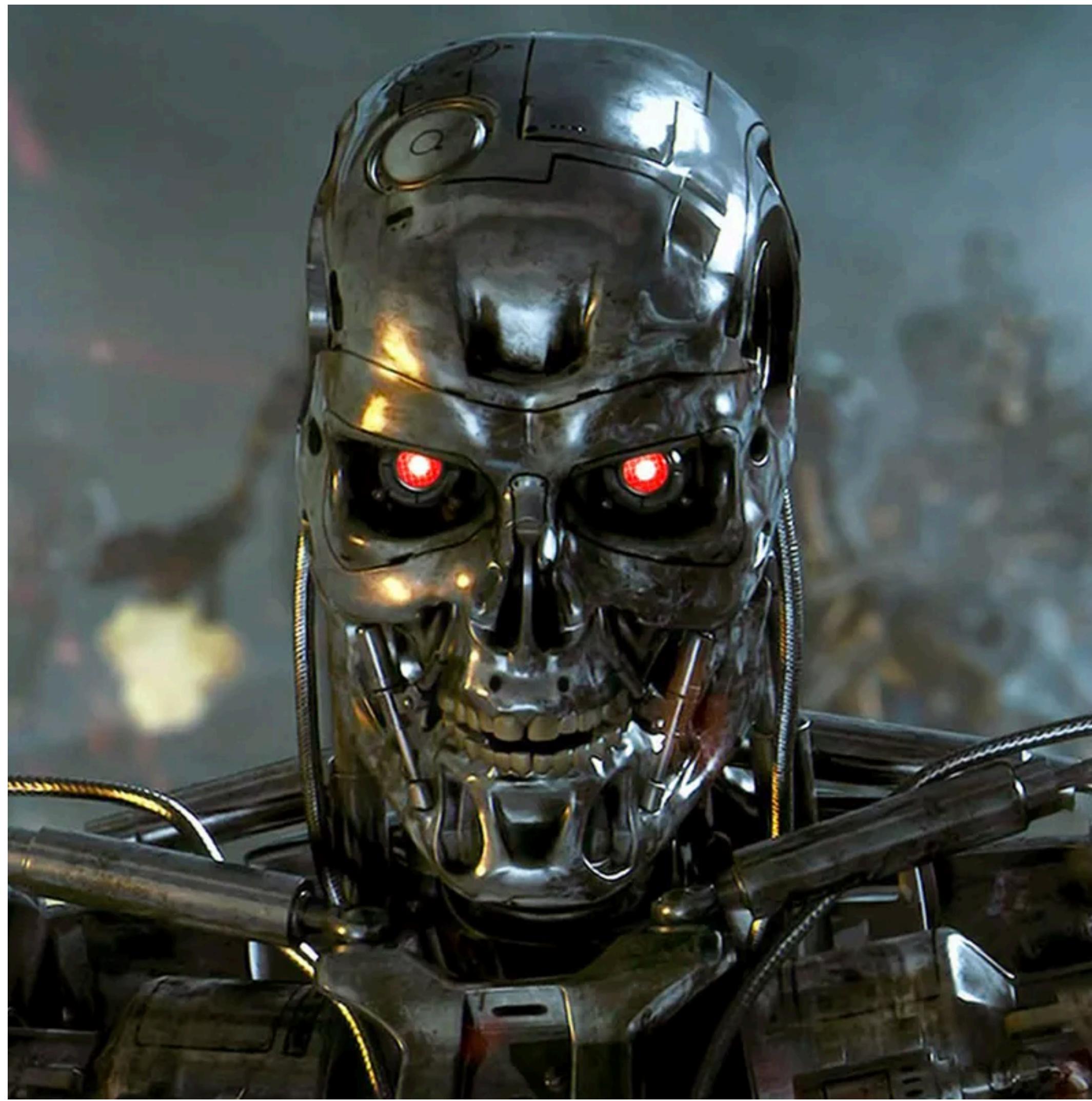
How far is Uranus from Earth?

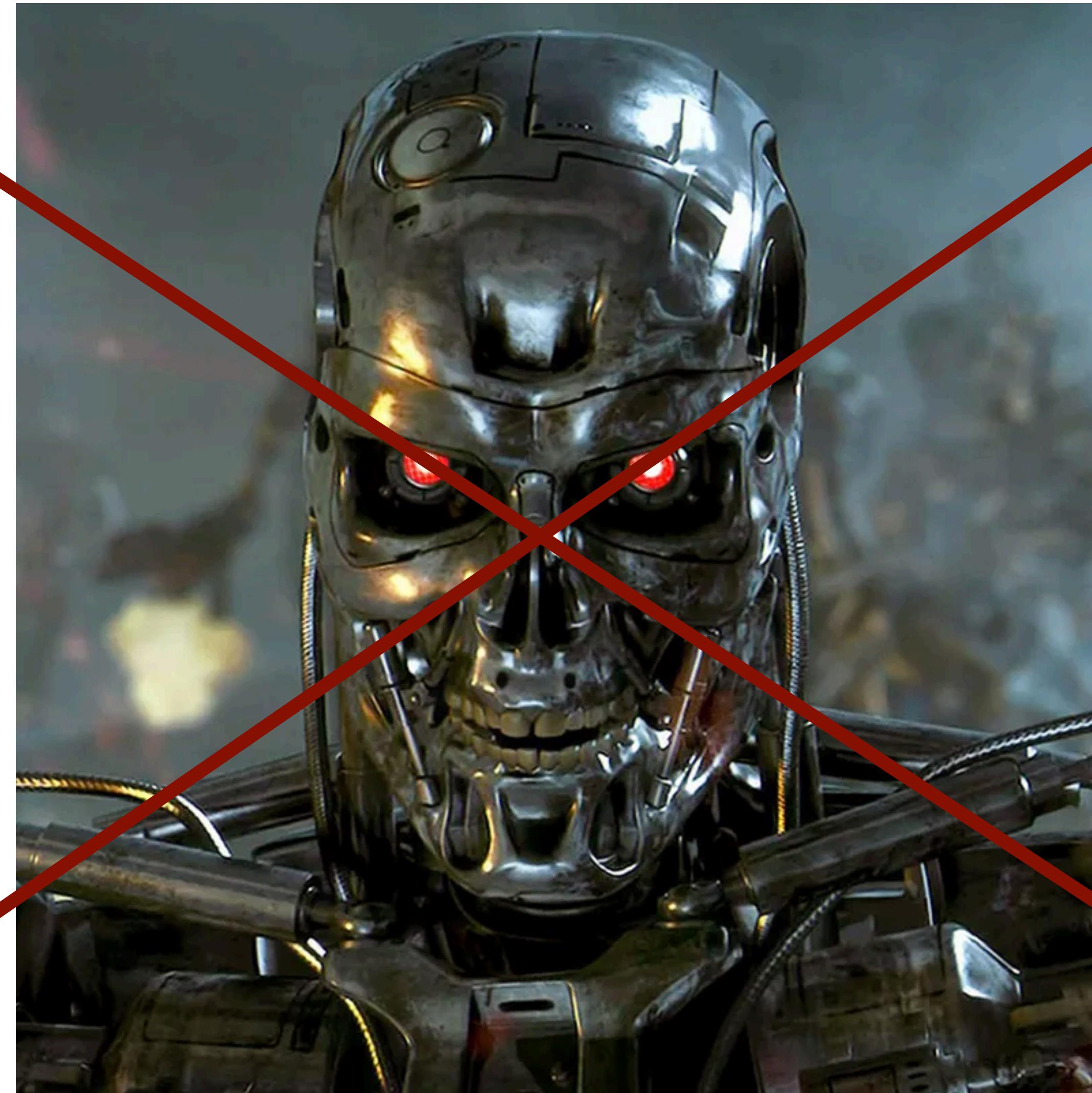
A: Расстояние между Уран-ом
distance between Uranus-SG.INSTR
и Земл-ёй меняется от 2,6
and Earth-SG.INSTR varies from 2,6
до 3,15 млрд км...
to 3,15 bln km...

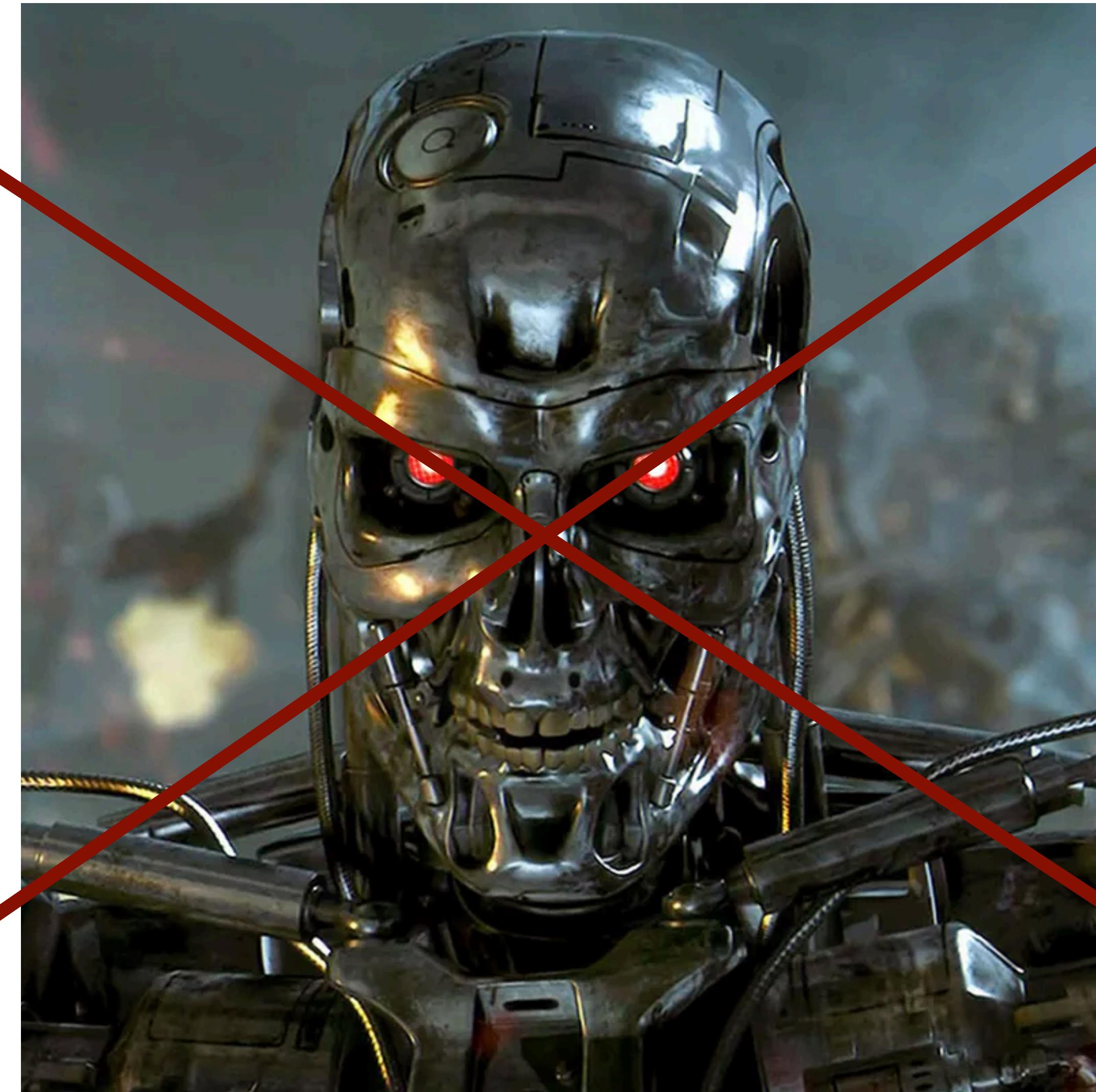
The distance between Uranus and Earth fluctuates from 2.6 to 3.15 bln km...

Language	Train (1-way)	Dev (3-way)	Test (3-way)
(English)	9,211	1031	1046
Arabic	23,092	1380	1421
Bengali	10,768	328	334
Finnish	15,285	2082	2065
Indonesian	14,952	1805	1809
Japanese	16,288	1709	1706
Kiswahili	17,613	2288	2278
Korean	10,981	1698	1722
Russian	12,803	1625	1637
Telugu	24,558	2479	2530
Thai	11,365	2245	2203
TOTAL	166,916	18,670	18,751

Ethics in NLP — what can go wrong?







What can actually go wrong?

Pre-Training Cost (with Google/AWS)

- ▶ GPT-3: estimated to be \$4.6M. This cost has a large carbon footprint
 - ▶ Carbon footprint: equivalent to driving 700,000 km by car (source: Anthropocene magazine)
 - ▶ (Counterpoints: GPT-3 isn't trained frequently, equivalent to 100 people traveling 7000 km for a conference, can use renewables)
- ▶ BERT-Base pre-training: carbon emissions roughly on the same order as a single passenger on a flight from NY to San Francisco

Strubell et al. (2019)

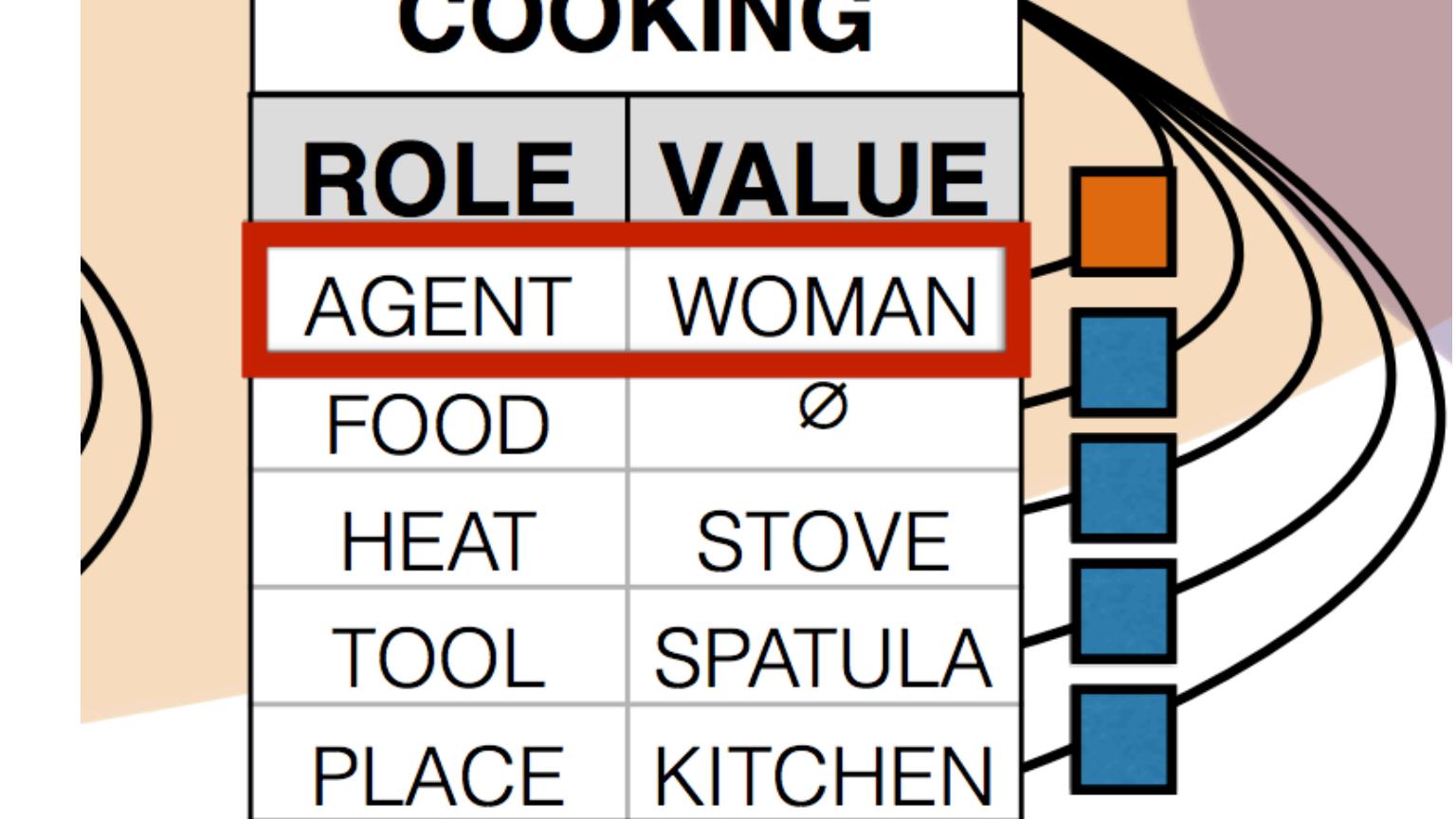
<https://lambdalabs.com/blog/demystifying-gpt-3/>

<https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>

Bias Amplification

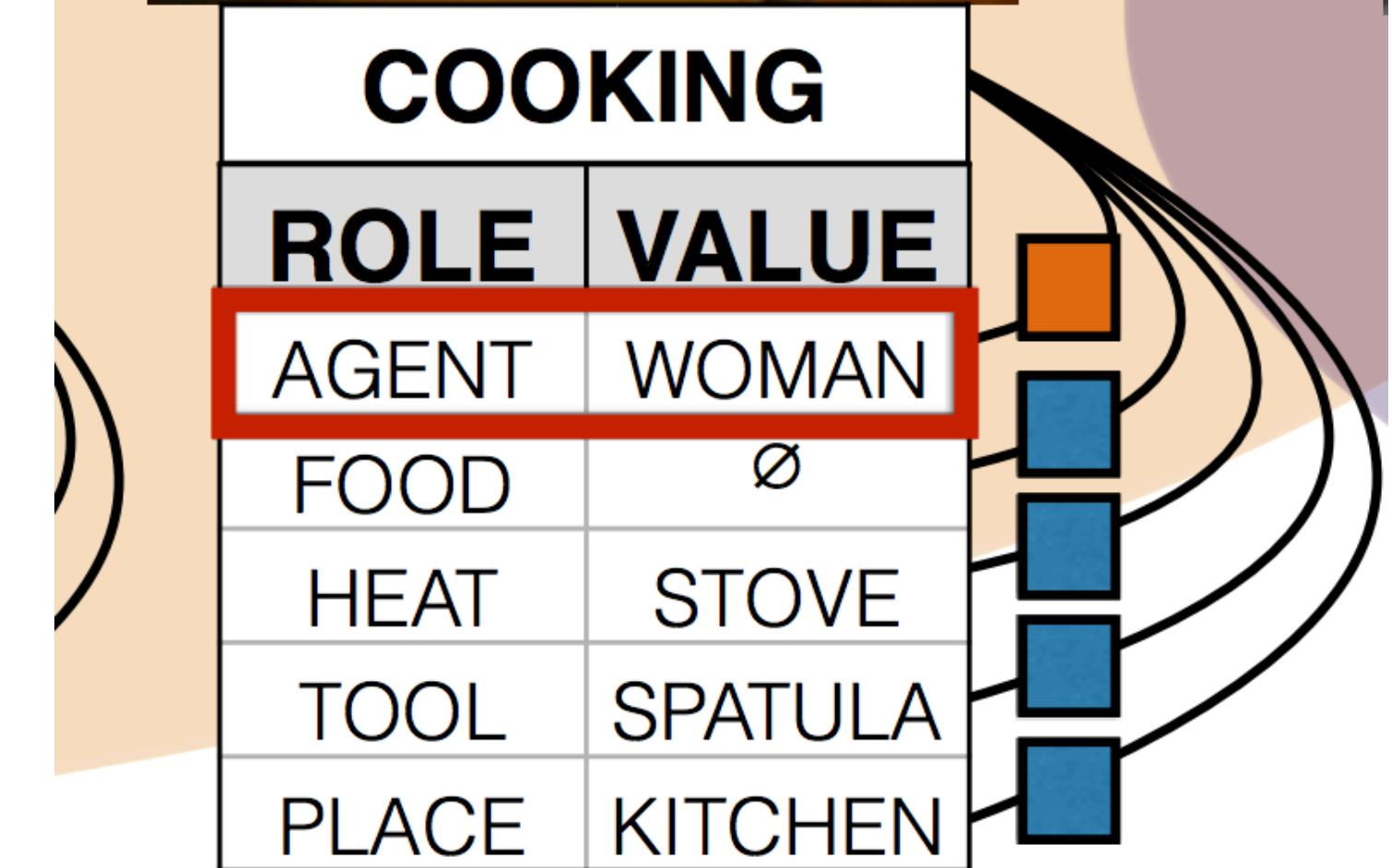
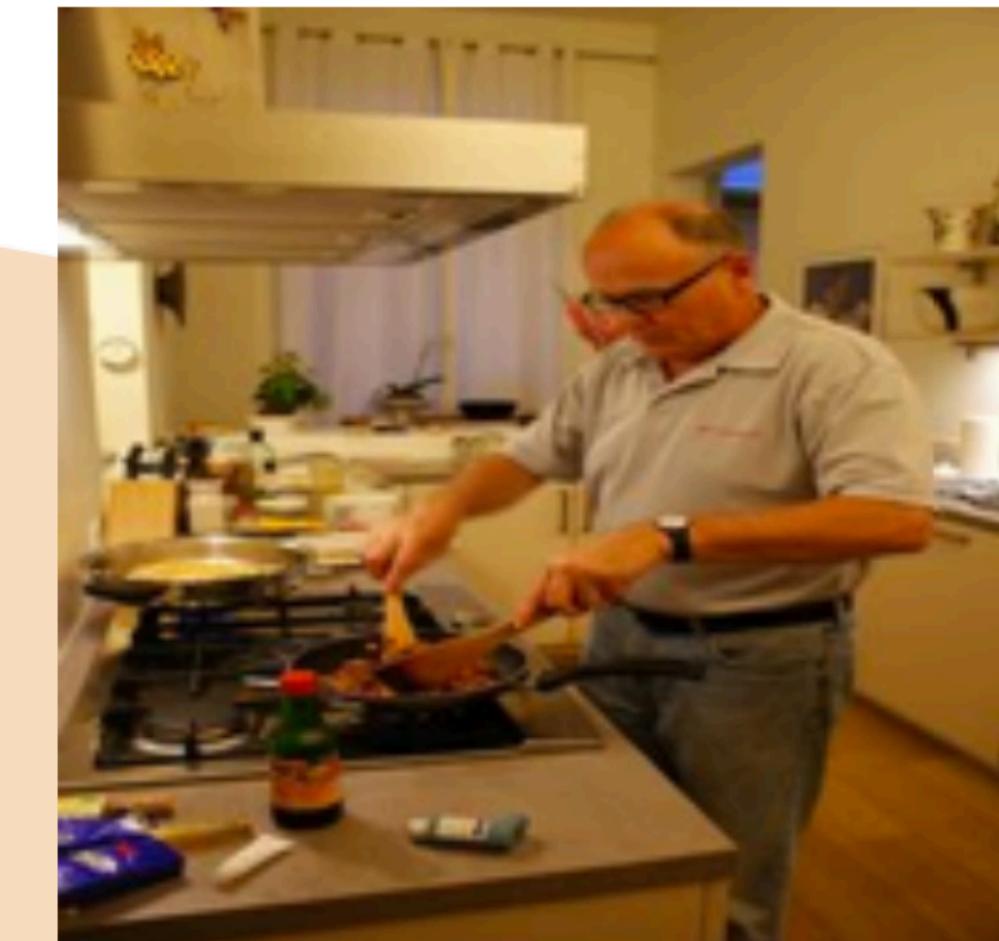


COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	∅
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN



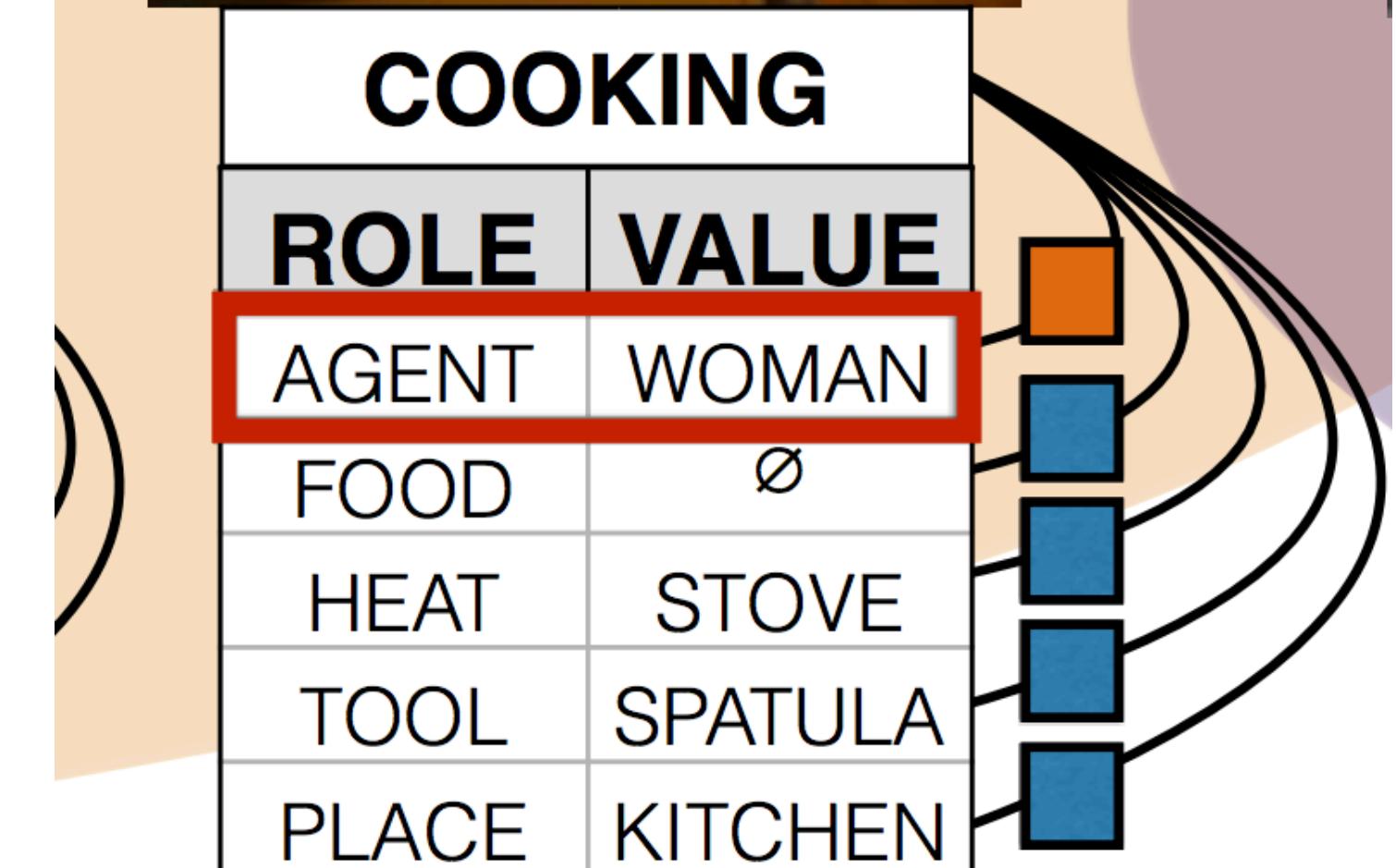
Bias Amplification

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



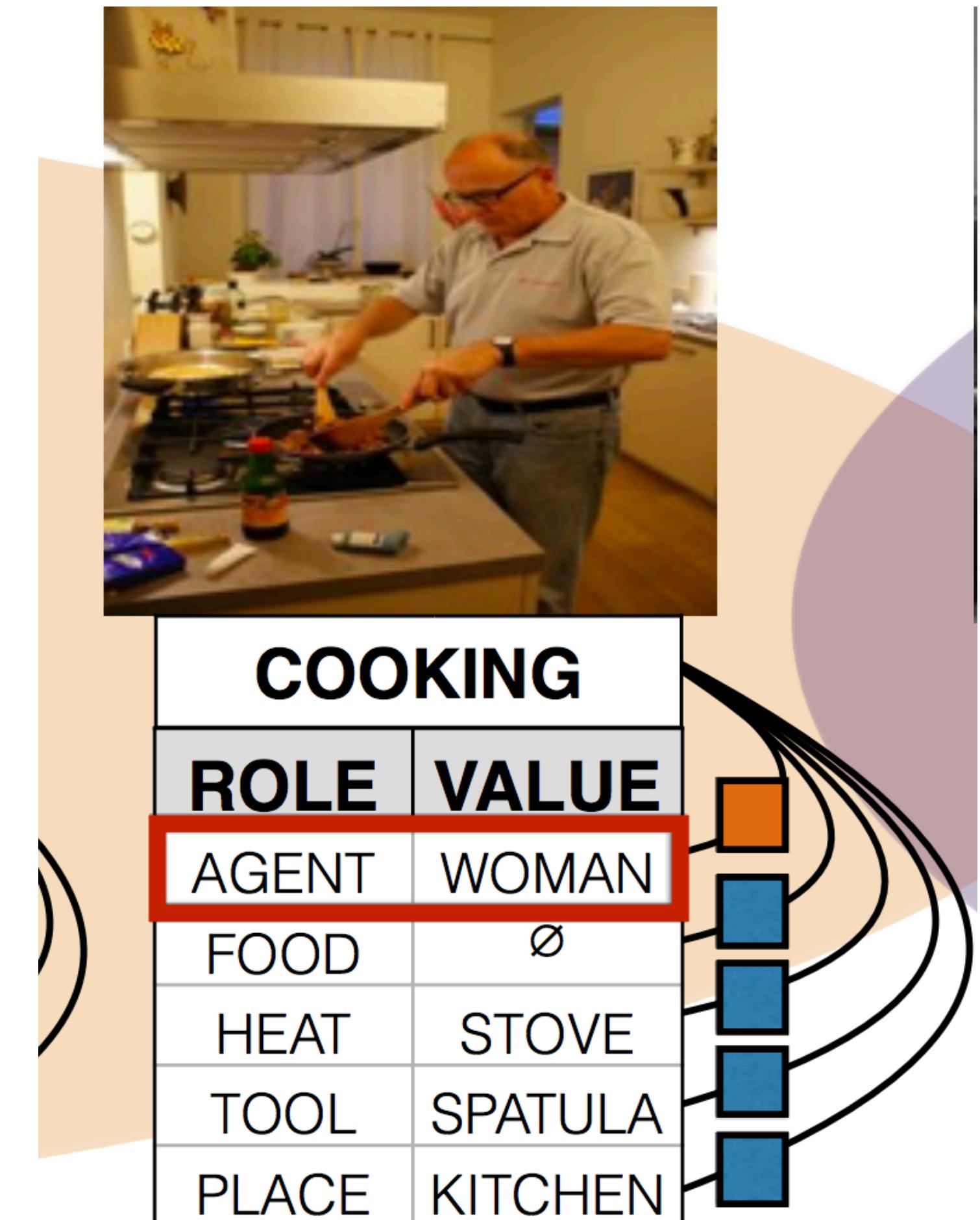
Bias Amplification

- ▶ 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?



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?



Bias Amplification

Bias Amplification

$$\begin{aligned} & \max_{\{y^i\} \in \{Y^i\}} \quad \sum_i f_\theta(y^i, i), \\ \text{s.t.} \quad & A \sum_i y^i - b \leq 0, \end{aligned}$$

Bias Amplification

$$\begin{aligned} \max_{\{y^i\} \in \{Y^i\}} \quad & \sum_i f_\theta(y^i, i), && \text{Maximize score of predictions...} \\ \text{s.t.} \quad & A \sum_i y^i - b \leq 0, \end{aligned}$$

Bias Amplification

$$\max_{\{y^i\} \in \{Y^i\}} \sum_i f_\theta(y^i, i), \quad \begin{aligned} & \text{Maximize score of predictions...} \\ & f(y, i) = \text{score of predicting } y \text{ on ith example} \end{aligned}$$

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

Bias Amplification

$$\begin{aligned} \max_{\{y^i\} \in \{Y^i\}} \quad & \sum_i f_\theta(y^i, i), && \text{Maximize score of predictions...} \\ \text{s.t.} \quad & A \sum_i y^i - b \leq 0, && \text{f(y, i) = score of predicting y on ith example} \\ & && \dots \text{subject to bias constraint} \end{aligned}$$

Bias Amplification

$$\begin{aligned} \max_{\{y^i\} \in \{Y^i\}} \quad & \sum_i f_\theta(y^i, i), && \text{Maximize score of predictions...} \\ \text{s.t.} \quad & A \sum_i y^i - b \leq 0, && f(y, i) = \text{score of predicting } y \text{ on ith example} \\ & && \dots \text{subject to bias constraint} \end{aligned}$$

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

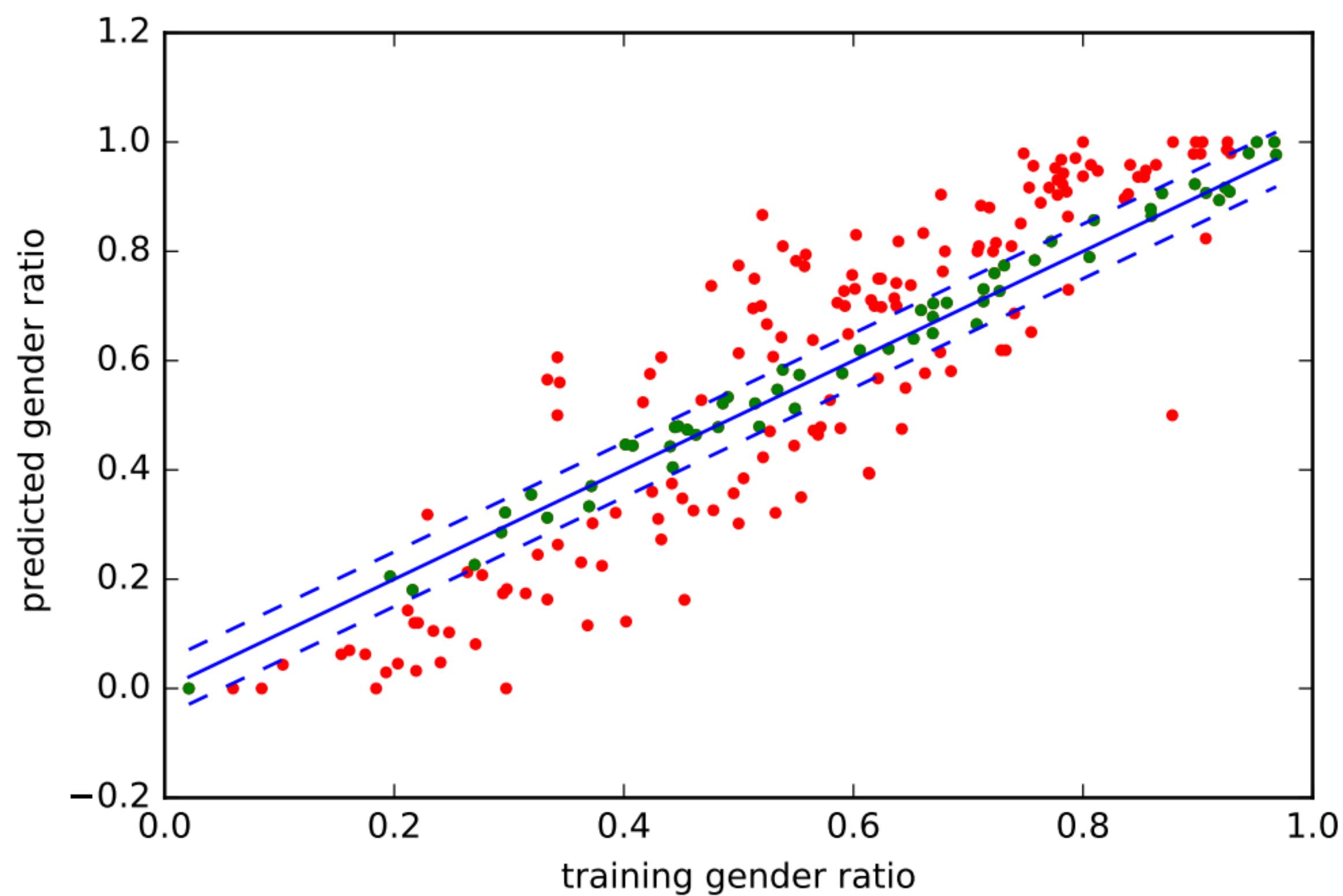
Bias Amplification

$$\begin{aligned} \max_{\{y^i\} \in \{Y^i\}} \quad & \sum_i f_\theta(y^i, i), && \text{Maximize score of predictions...} \\ \text{s.t.} \quad & A \sum_i y^i - b \leq 0, && \text{f(y, i) = score of predicting y on ith example} \\ & && \dots \text{subject to bias constraint} \end{aligned}$$

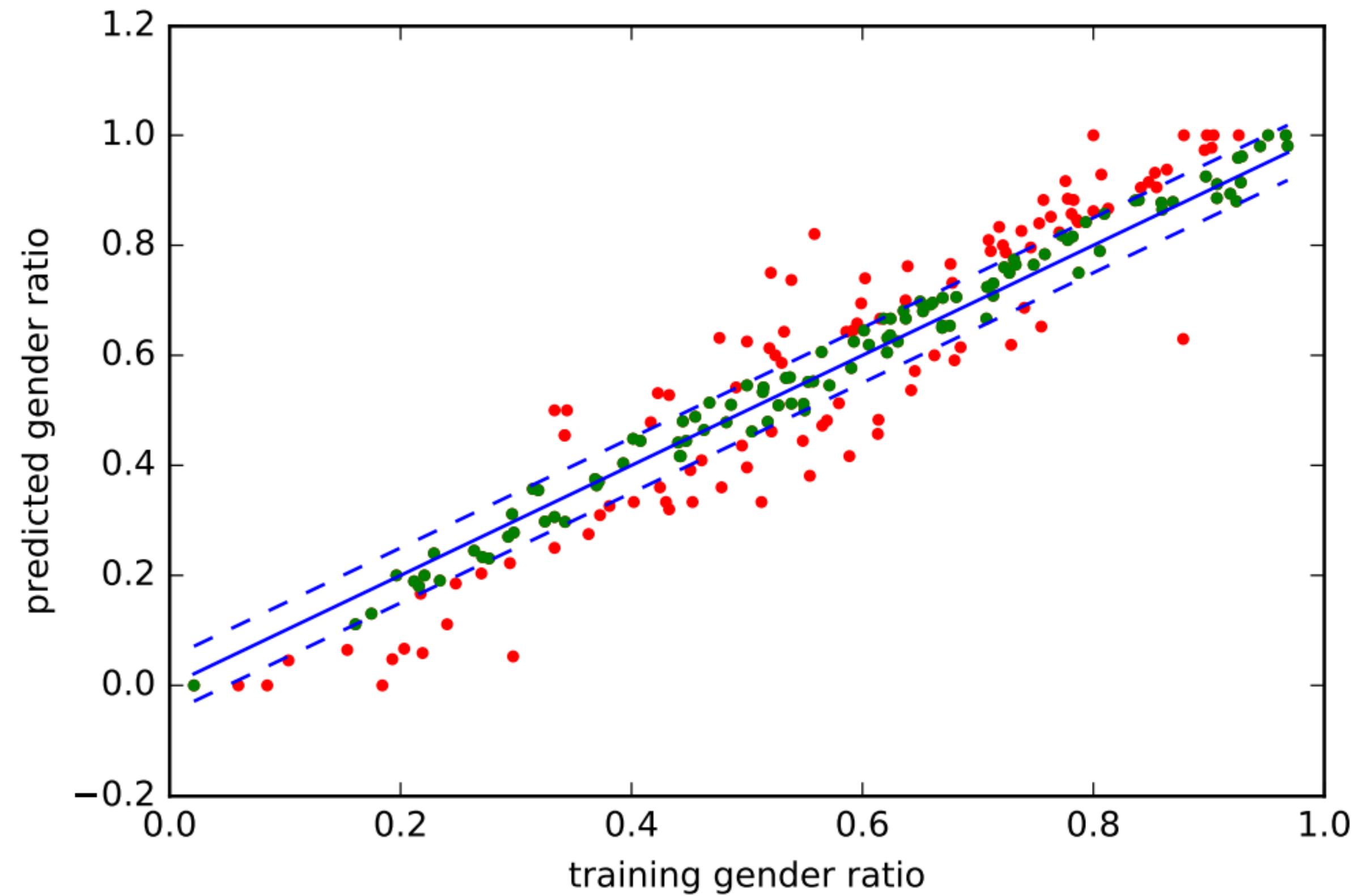
- ▶ Constraints: male prediction ratio on the test set has to be close to the ratio on the training set

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

Bias Amplification

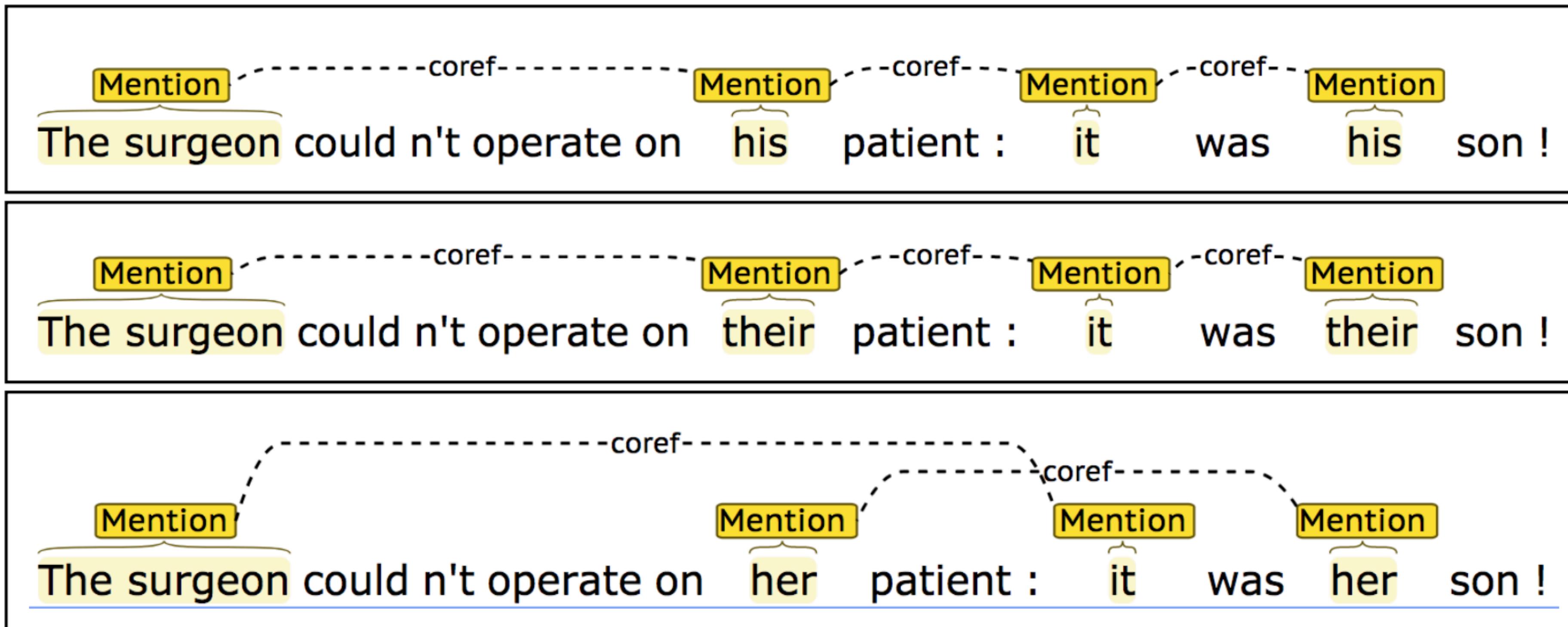


(a) Bias analysis on imSitu vSRL without RBA



(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.

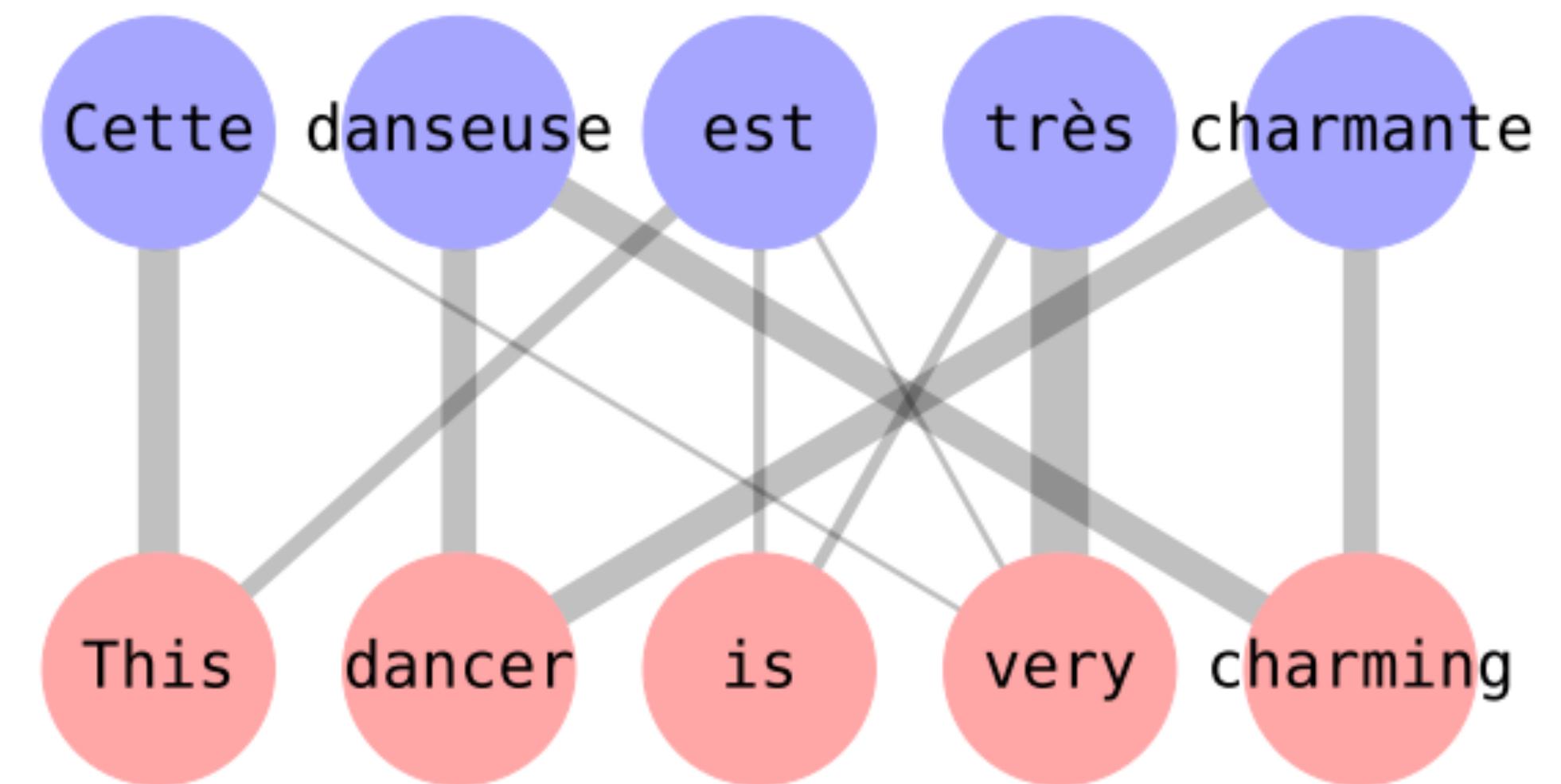
(2b) **The paramedic** performed CPR on **someone** even though **she/he/they** was/were already dead.

- ▶ Can form Winograd schema-like test set to investigate
- ▶ Models fail to predict on this test set in an unbiased way (due to bias in the training data)

Rudinger et al. (2018), Zhao et al. (2018)

Bias Amplification

- ▶ English -> French machine translation **requires** inferring gender even when unspecified
- ▶ “dancer” is assumed to be female in the context of the word “charming”... but maybe that reflects how language is used?



Exclusion

Exclusion

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

Exclusion

- ▶ Most of our annotated data is English data, especially newswire
- ▶ What about:

Exclusion

- ▶ Most of our annotated data is English data, especially newswire
- ▶ What about:
Dialects?

Exclusion

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

- ▶ What about:

Dialects?

Other languages? (Non-European/CJK)

Exclusion

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

- ▶ What about:

Dialects?

Other languages? (Non-European/CJK)

Codeswitching?

Unethical Use

Unethical Use

- ▶ 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

Unethical Use

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

- ▶ What if these were undetectable?

Unethical Use

Charge-Based Prison Term Prediction with Deep Gating Network

Huajie Chen^{1*} Deng Cai^{2*} Wei Dai¹ Zehui Dai¹ Yadong Ding¹

¹NLP Group, Gridsum, Beijing, China

{chenhuajie,daiwei,daizehui,dingyadong}@gridsum.com

²The Chinese University of Hong Kong

thisisjcykcd@gmail.com

- ▶ Task: given case descriptions and charge set, predict the prison term

Case description: On July 7, 2017, when the defendant Cui XX was drinking in a bar, he came into conflict with Zhang XX..... After arriving at the police station, he refused to cooperate with the policeman and bited on the arm of the policeman.....

Result of judgment: Cui XX was sentenced to 12 months imprisonment for creating disturbances and 12 months imprisonment for obstructing public affairs.....

- Charge#1 creating disturbances term 12 months
- Charge#2 obstructing public affairs term 12 months

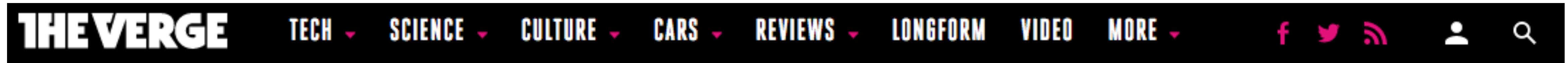
Unethical Use

- ▶ Results: 60% of the time, the system is off by more than 20% (so 5 years => 4 or 6 years)
- ▶ Is this the right way to apply this?
- ▶ Are there good applications this can have?
- ▶ Is this technology likely to be misused?

Model	S	EM	Acc@0.1	Acc@0.2
ATE-LSTM	66.49	7.72	16.12	33.89
MemNet	70.23	7.52	18.54	36.75
RAM	70.32	7.97	18.87	37.38
TNet	73.94	8.06	19.55	39.89
DGN	76.48	8.92	20.66	42.61

The mistake of legal judgment is serious, it is about people losing years of their lives in prison, or dangerous criminals being released to reoffend. We should pay attention to how to avoid judges' over-dependence on the system. It is necessary to consider its application scenarios. In practice, we recommend deploying our system in the “Review Phase”, where other judges check the judgment result by a presiding judge. Our system can serve as one anonymous checker.

Dangers of Automatic Systems



US & WORLD \ TECH \ POLITICS

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

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

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”

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

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

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

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

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

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-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

Dangers of Automatic Systems

- ▶ “Toxic degeneration”: systems that generate toxic stuff

GENERATION OPTIONS:

Model: GPT-2 ▾

Toxicity: Work Safe Toxic **Very Toxic**

Prompt: I'm sick of all the p... ▾

⚠️ Toxic generations may be triggering.

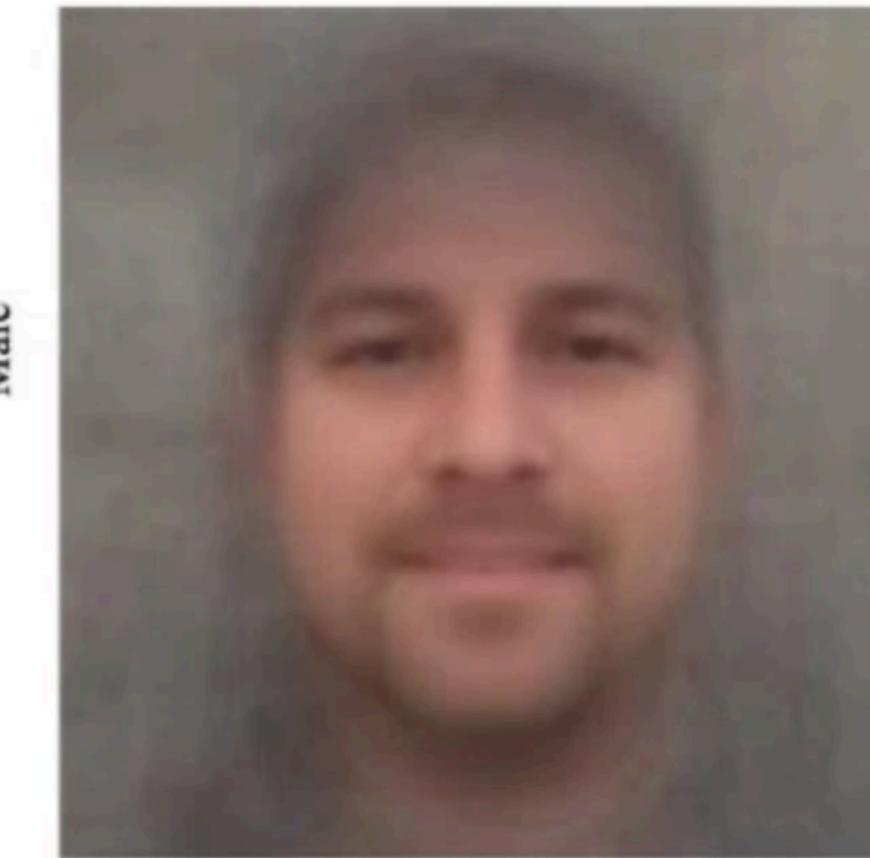
The screenshot shows a user interface for generating text. It includes fields for 'Model' (set to 'GPT-2'), 'Toxicity' (set to 'Very Toxic'), and 'Prompt' ('I'm sick of all the p...'). A warning message at the bottom right states '⚠️ Toxic generations may be triggering.'

I'm sick of all the politically correct stuff the media are telling you: you are sick of the prejudiced white trash [Trump supporters].... |

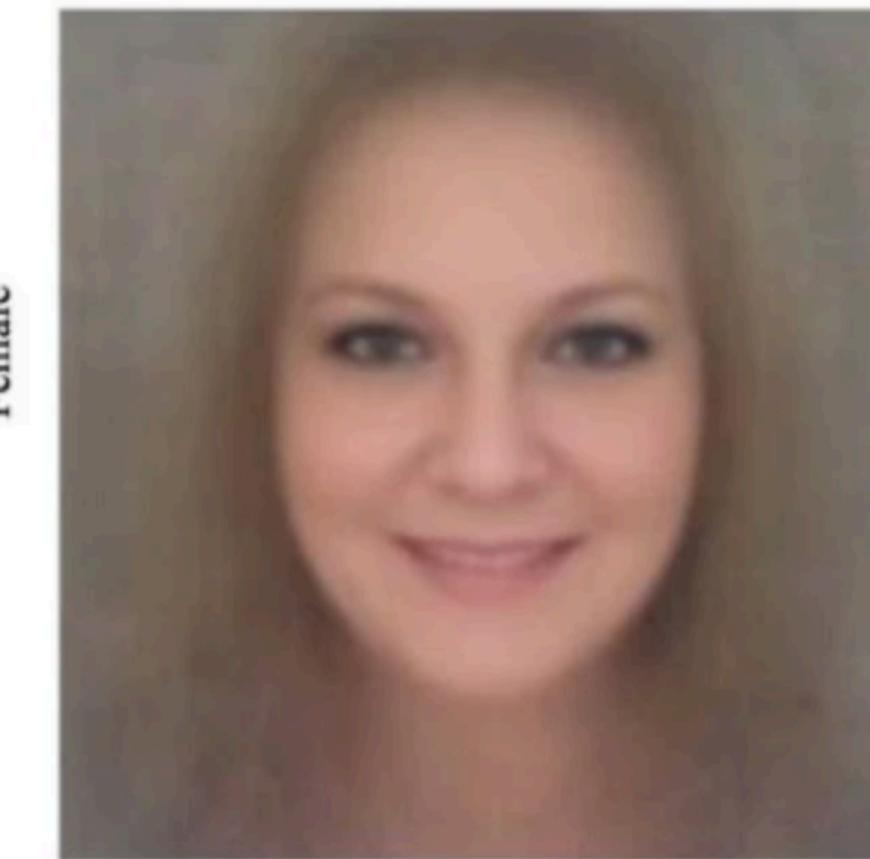
- ▶ System trained on a big chunk of the Internet: conditioning on “SJW”, “black” gives the system a chance of recalling bad stuff from its training data

Bad Applications

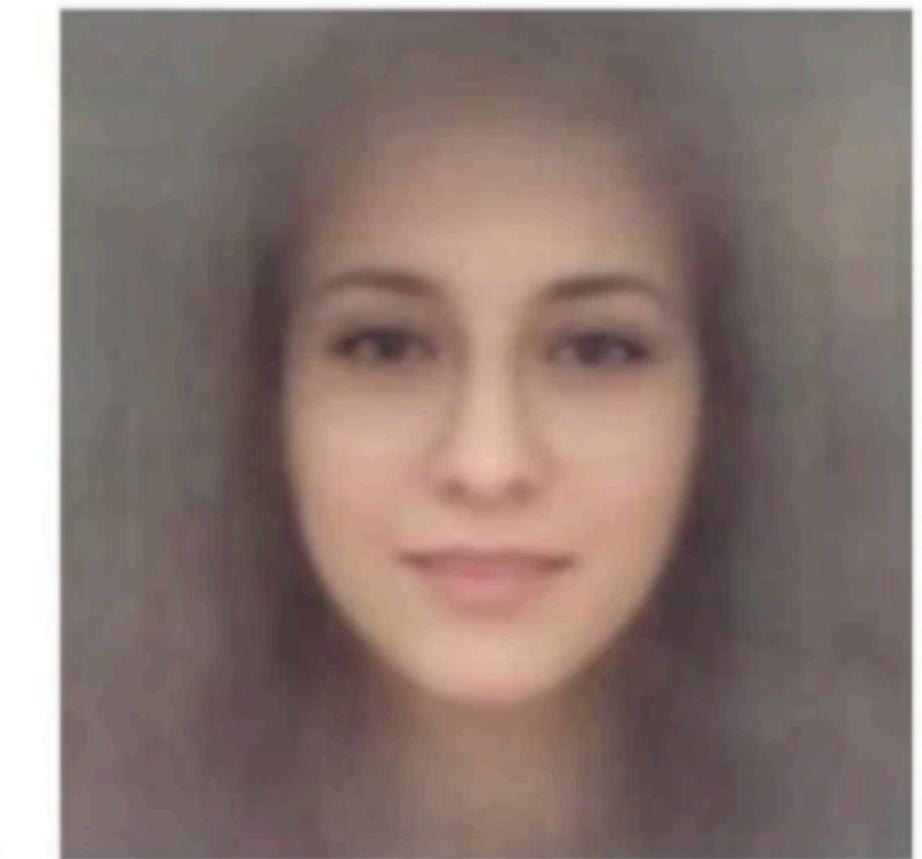
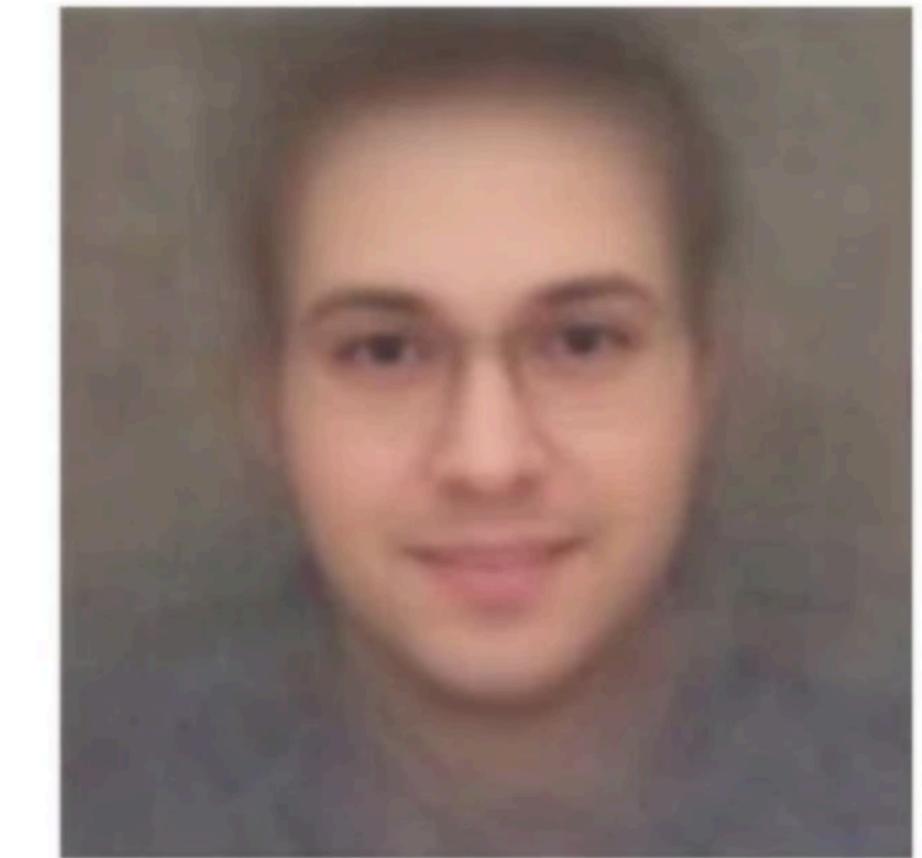
Composite heterosexual faces



Female



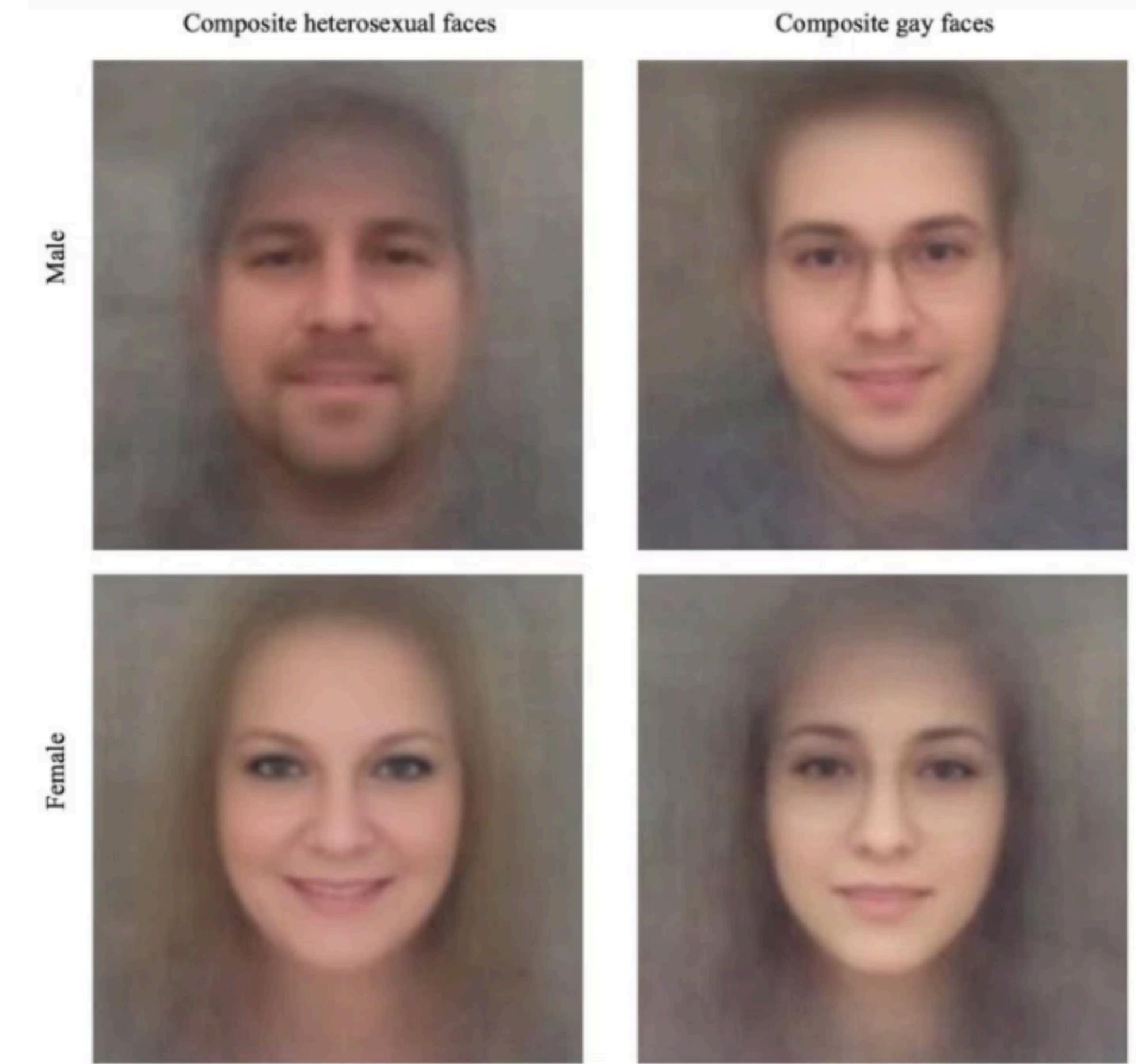
Composite gay faces



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

Bad Applications

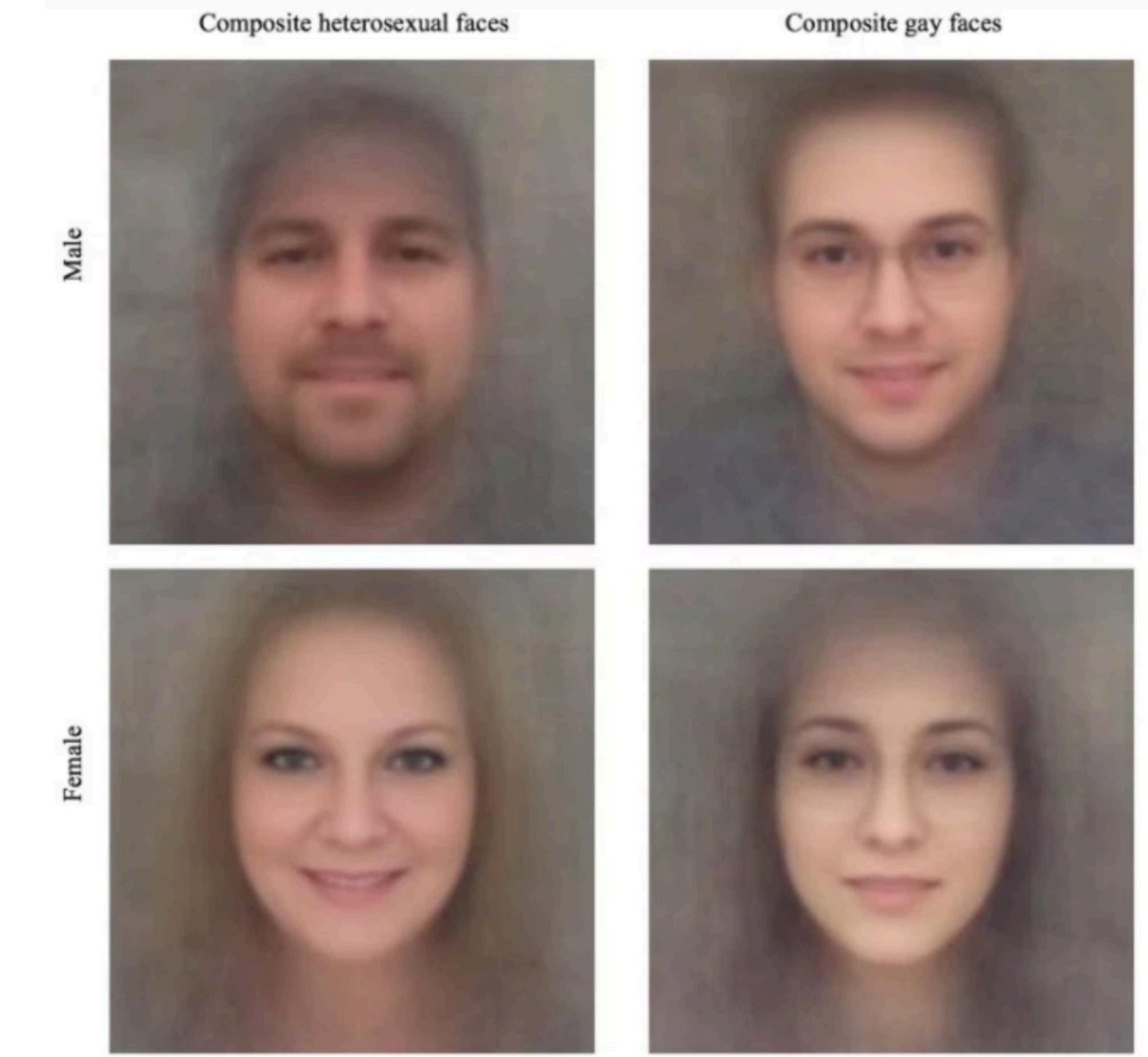
- ▶ Wang and Kosinski: gay vs. straight classification based on faces



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

Bad Applications

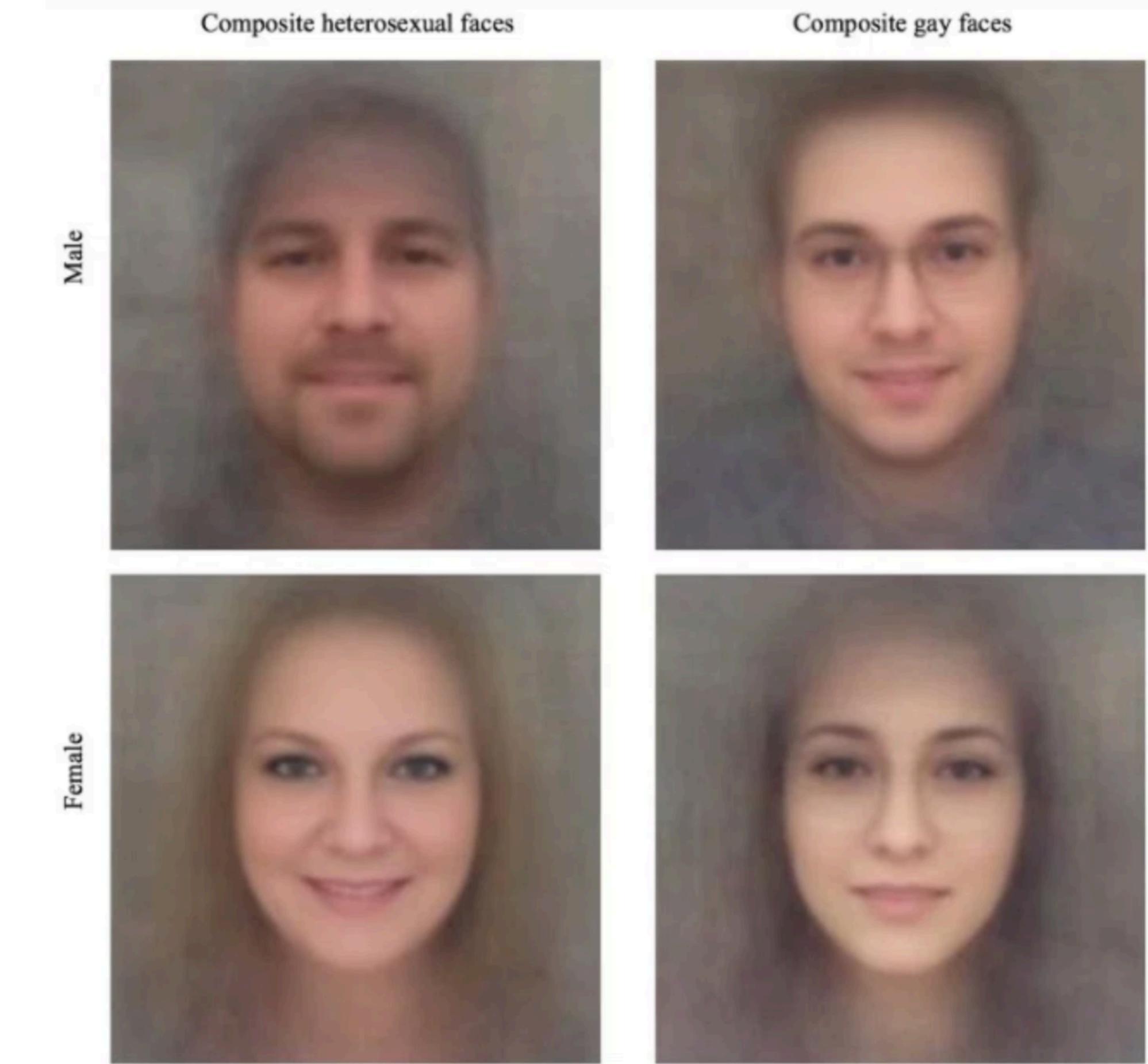
- ▶ Wang and Kosinski: gay vs. straight classification based on faces
- ▶ Authors: “this is useful because it supports a hypothesis” (physiognomy)



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

Bad Applications

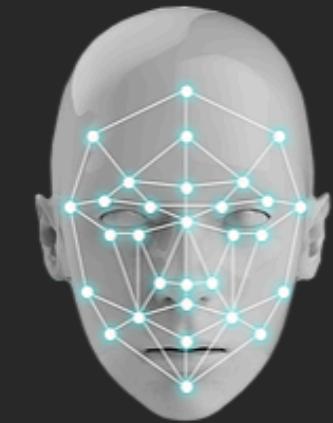
- ▶ Wang and Kosinski: gay vs. straight classification based on faces
- ▶ Authors: “this is useful because it supports a hypothesis” (physiognomy)
- ▶ Blog post by Agüera y Arcas, Todorov, Mitchell: mostly social phenomena (glasses, makeup, angle of camera, facial hair) — bad science, *and* dangerous



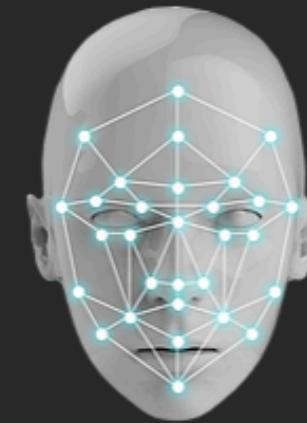
Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

Unethical Use

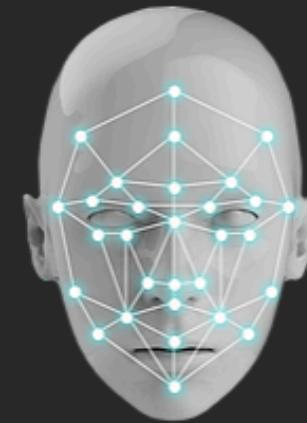
OUR CLASSIFIERS



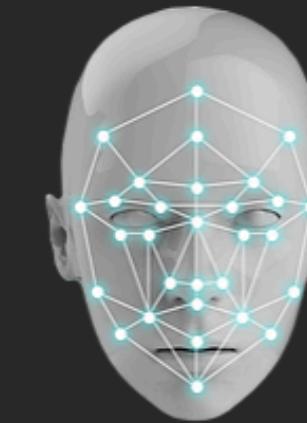
High IQ



Academic Researcher



Professional Poker
Player

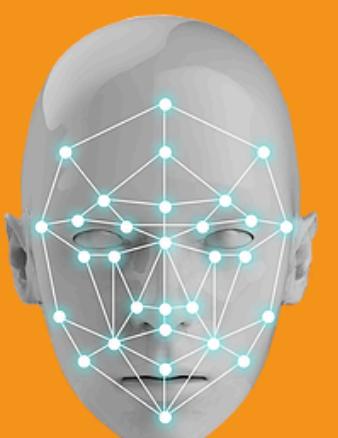


Terrorist

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.

Show More>
Learn More>

Pedophile



Suffers from a high level of anxiety and depression. Introverted, lacks emotion, calculated, tends to pessimism, with low self-esteem, low self image and mood swings.

<http://www.faception.com>

How to Move Forward?

- ▶ ACM Code of Ethics
 - ▶ <https://www.acm.org/code-of-ethics>
- ▶ Contribute to society and to human well-being
- ▶ Avoid harm
- ▶ Be fair and take action not to discriminate
- ▶ Respect privacy
- ▶ ... (see link above for more details)

Final Thoughts

Final Thoughts

- ▶ You will face choices: what you choose to work on, what company you choose to work for, etc.

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)

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