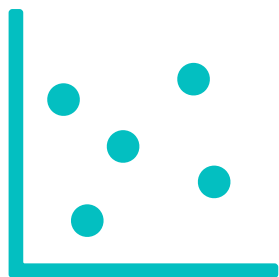
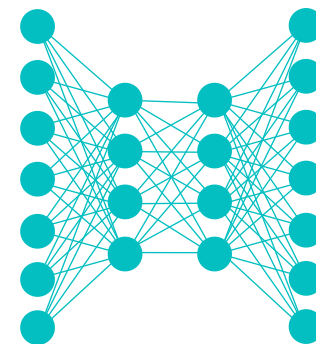


Lecture Notes for **Neural Networks and Machine Learning**



A Practical Example of
Ethically “Aware” NLP Practices

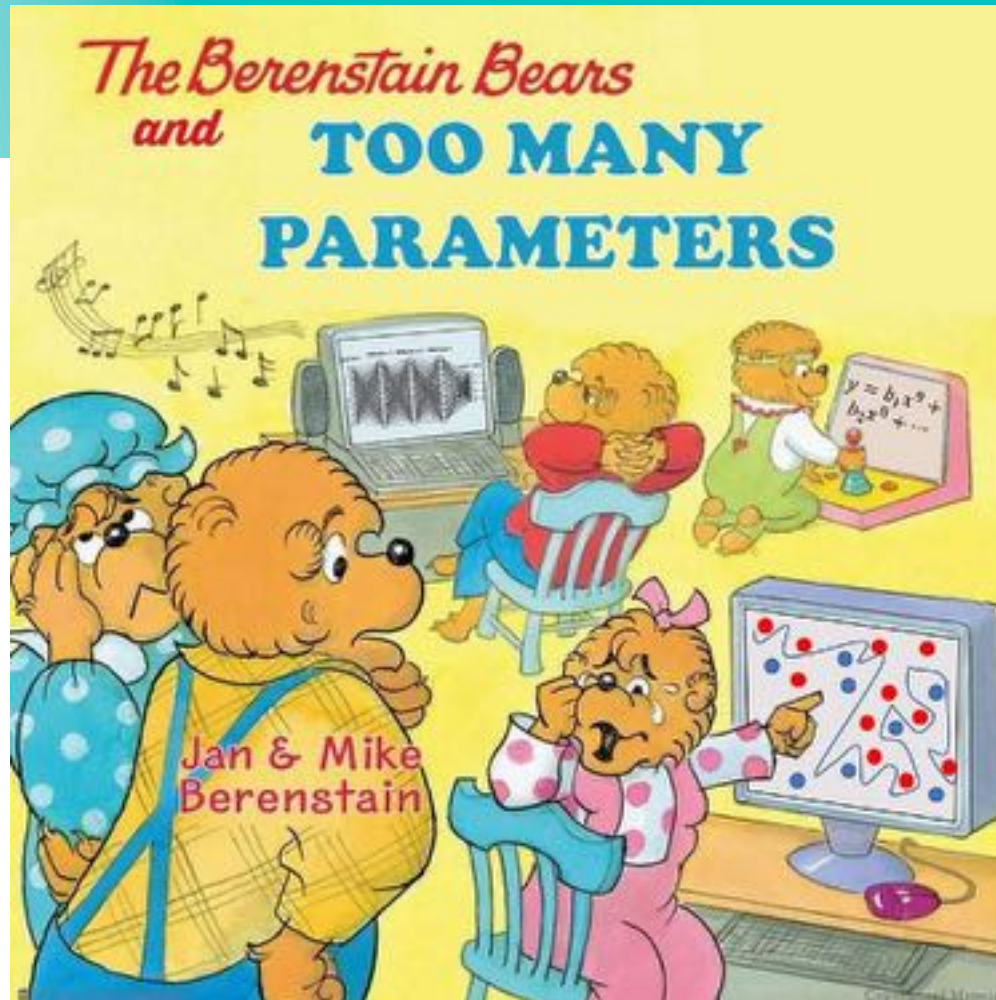


Logistics and Agenda

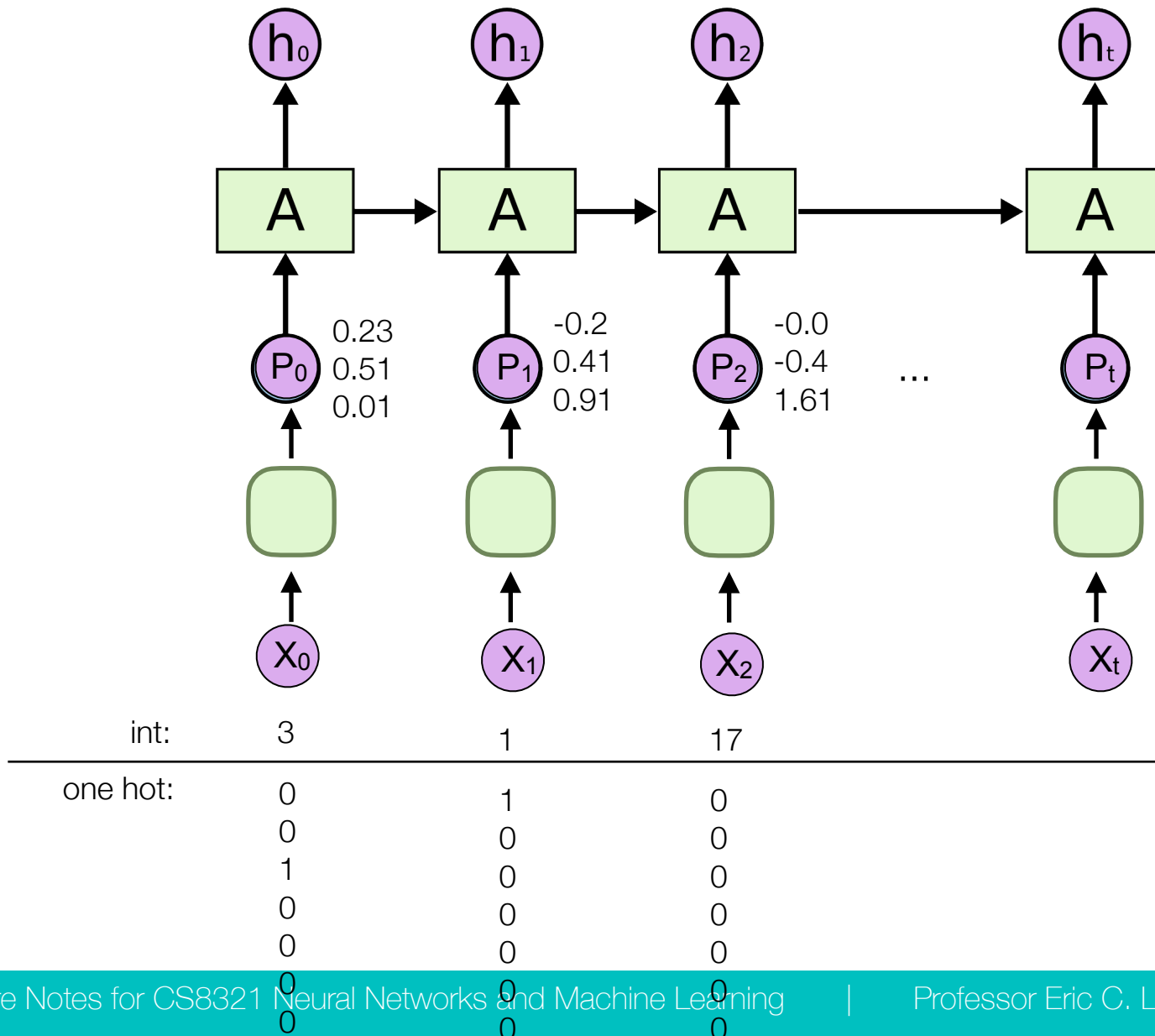
- Logistics
 - Lecture discussion assignments
 - Office hours
- Last Time:
 - Ethical Guidelines
 - Case Studies
- Agenda
 - Word Embedding Review
 - Implicit additive de-biasing
 - Subtractive de-biasing



NLP Embeddings Review

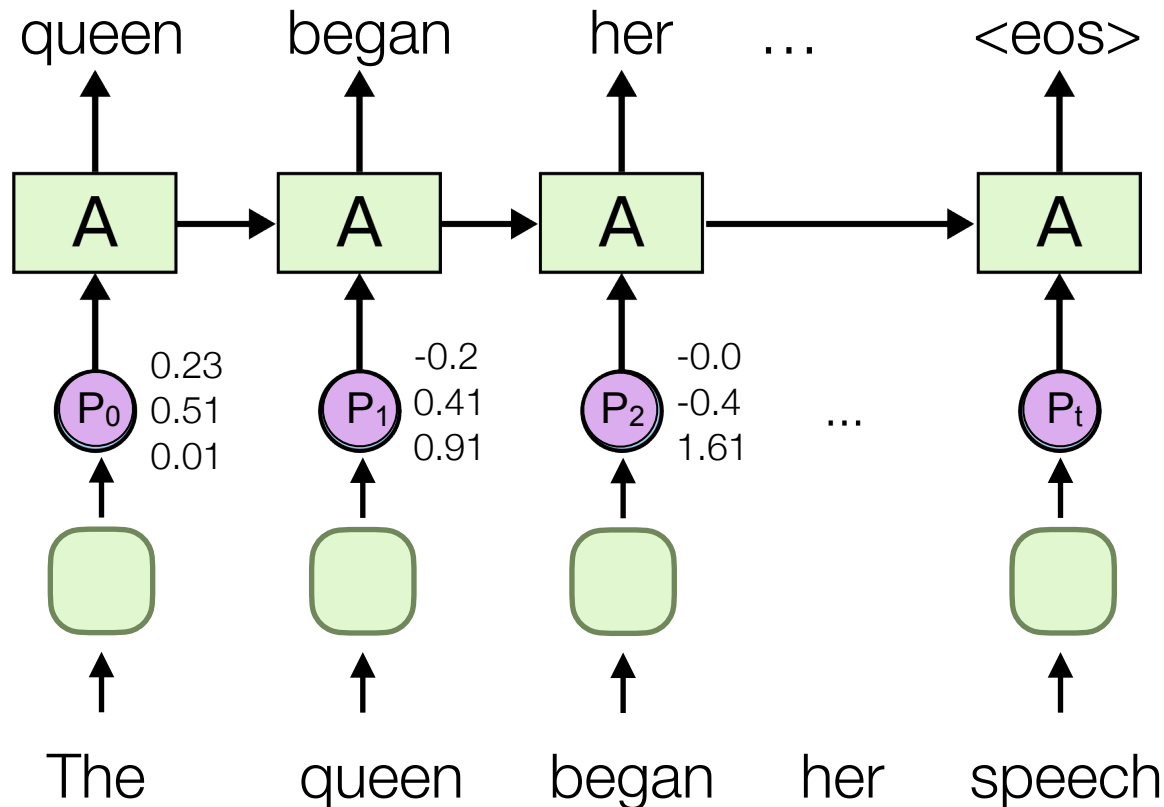


Word Embeddings Review



Word Embeddings: Training Review

- many training options exist
 - a popular option, next word prediction



GloVe Review

GloVe

Global Vectors for Word Representation

Highlights

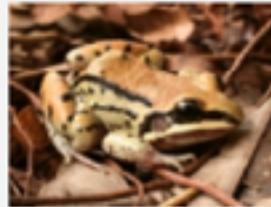
1. Nearest neighbors

The Euclidean distance (or cosine similarity) between two word vectors provides an effective method for measuring the linguistic or semantic similarity of the corresponding words. Sometimes, the nearest neighbors according to this metric reveal rare but relevant words that lie outside an average human's vocabulary. For example, here are the closest words to the target word *frog*:

0. *frog*
1. frogs
2. toad
3. litoria
4. leptodactylidae
5. rana
6. lizard
7. eleutherodactylus



3. litoria



4. leptodactylidae

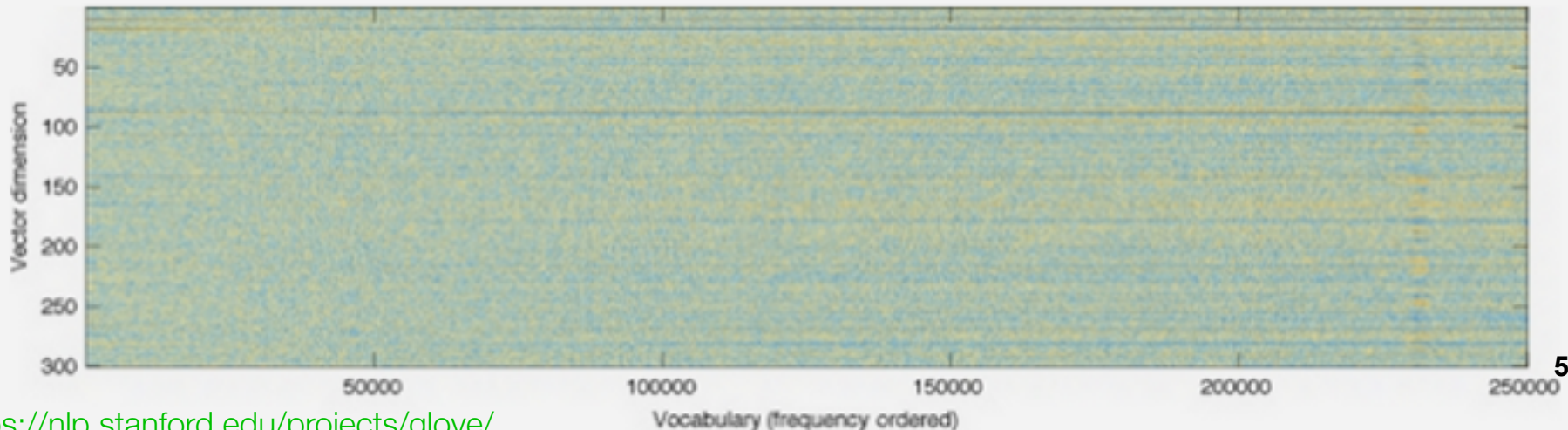


5. rana



7. eleutherodactylus

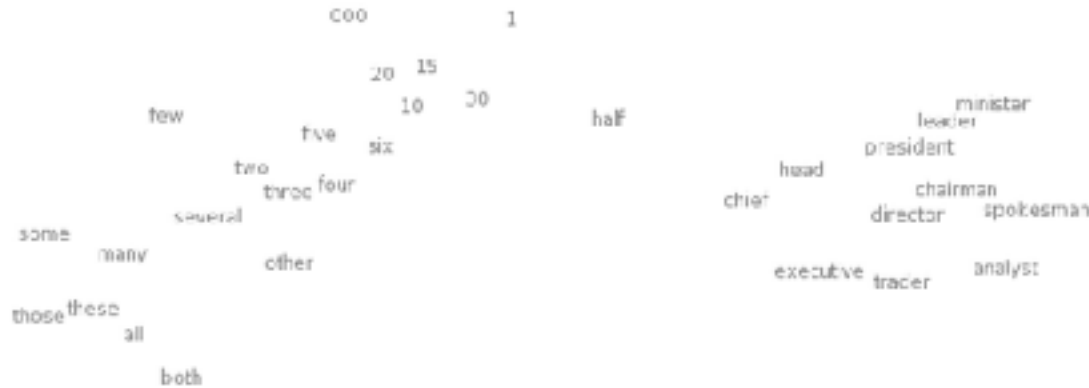
GloVe produces word vectors with a marked banded structure that is evident upon visualization:



Word Embeddings: proximity

GloVe Review

Global Vectors for Word Representation



t-SNE visualizations of word embeddings. Left: Number Region; Right: Jobs Region. From Turian *et al.* (2010), see complete image.

FRANCE	JESUS	XBOX	REDDISH	SCRATCHED	MEGABITS
AUSTRIA	GOD	AMIGA	GREENISH	NAILED	OCTETS
BELGIUM	SATI	PLAYSTATION	BLuish	SMASHED	MB/S
GERMANY	CHRIST	MSX	PINKISH	PUNCHED	BIT/S
ITALY	SATAN	IPOD	PURPLISH	POPPED	DAUD
GREECE	KALI	SEGA	BROWNISH	CRIMPED	CARATE
SWEDEN	INDRA	PSNUMBER	GREYISH	SCRAPED	KBIT/S
NORWAY	VISHNU	HD	GRAYISH	SCREWED	MEGAHERTZ
EUROPE	ANANDA	DREAMCAST	WHITISH	SECTIONED	MEGAPIXELS
HUNGARY	PARVATI	GEFORCE	SILVERY	SLASHED	GBIT/S
SWITZERLAND	GRACE	CAPCOM	YELLOWISH	RIPPED	AMPERES

What words have embeddings closest to a given word? From Collobert *et al.* (2011)

<http://colah.github.io/posts/2014-07-NLP-RNNs-Representations/>

The **chairman** called the **meeting** to order.

The **director** called the **conference** to order.

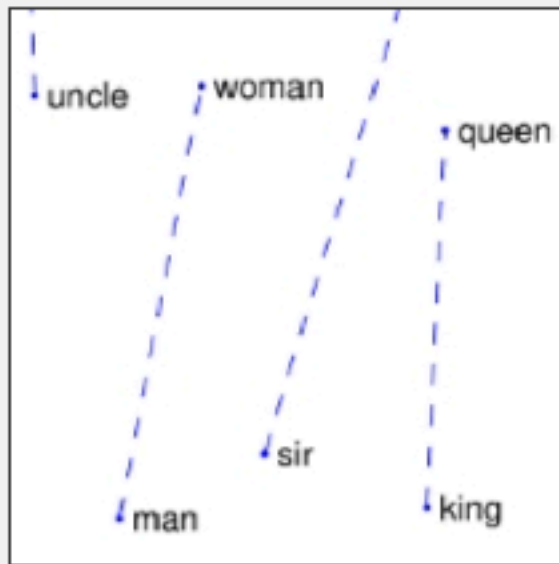
The **chief** called the **council** to order.



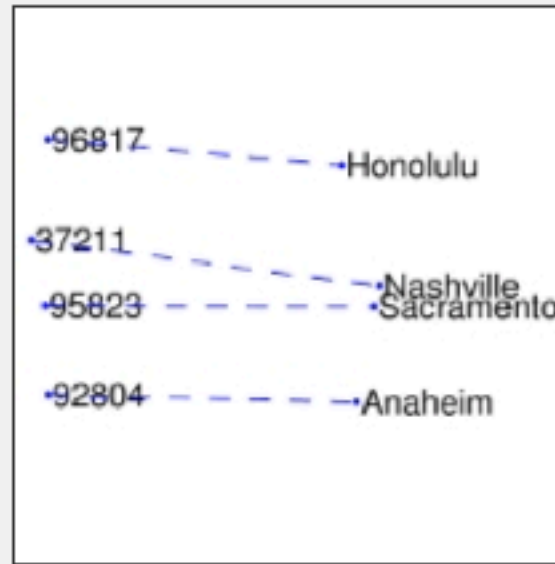
Word Embeddings: Analogy

GloVe Review

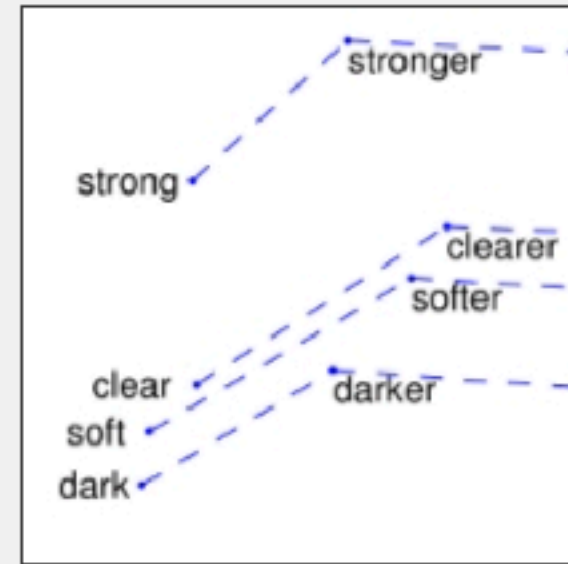
Global Vectors for Word Representation



man - woman



city - zip code



comparative - superlative

each vector difference **might** encode analogy



Word Embeddings: Analogy?

GloVe Review



$$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"aunt"}) - W(\text{"uncle"})$$

$$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"queen"}) - W(\text{"king"})$$

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}$$

From Mikolov *et al.*
(2013a)

Trained on
New York Times



Extreme *she* occupations

- | | | |
|-----------------|-----------------------|------------------------|
| 1. homemaker | 2. nurse | 3. receptionist |
| 4. librarian | 5. socialite | 6. hairdresser |
| 7. nanny | 8. bookkeeper | 9. stylist |
| 10. housekeeper | 11. interior designer | 12. guidance counselor |

Extreme *he* occupations

- | | | |
|----------------|-------------------|----------------|
| 1. maestro | 2. skipper | 3. protege |
| 4. philosopher | 5. captain | 6. architect |
| 7. financier | 8. warrior | 9. broadcaster |
| 10. magician | 11. fighter pilot | 12. boss |

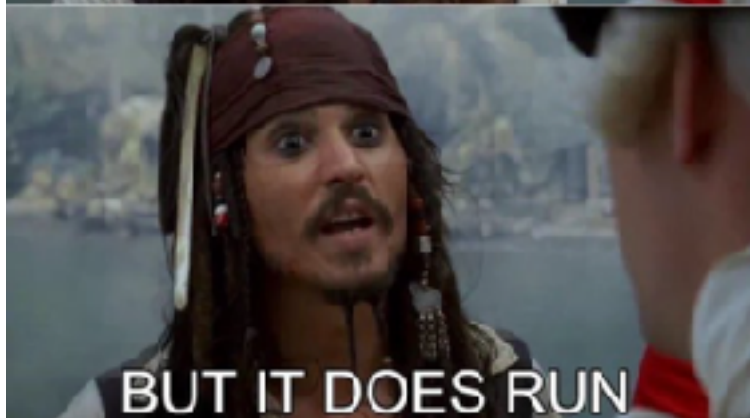
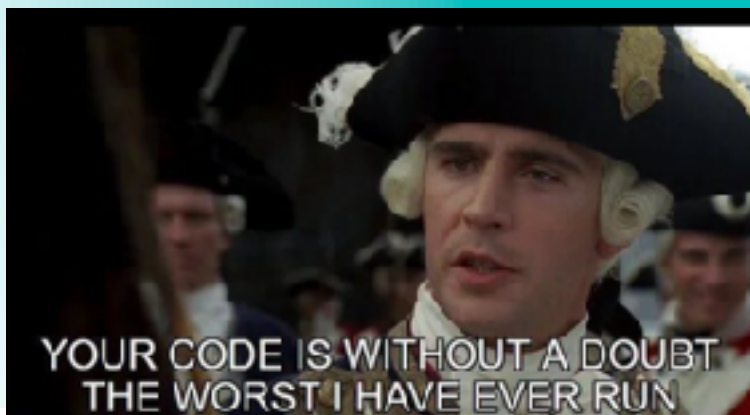
Bolukbasi et al., NeurIPS 2016

<https://arxiv.org/pdf/1607.06520.pdf>

<https://nlp.stanford.edu/projects/glove/>



Practical Example in NLP



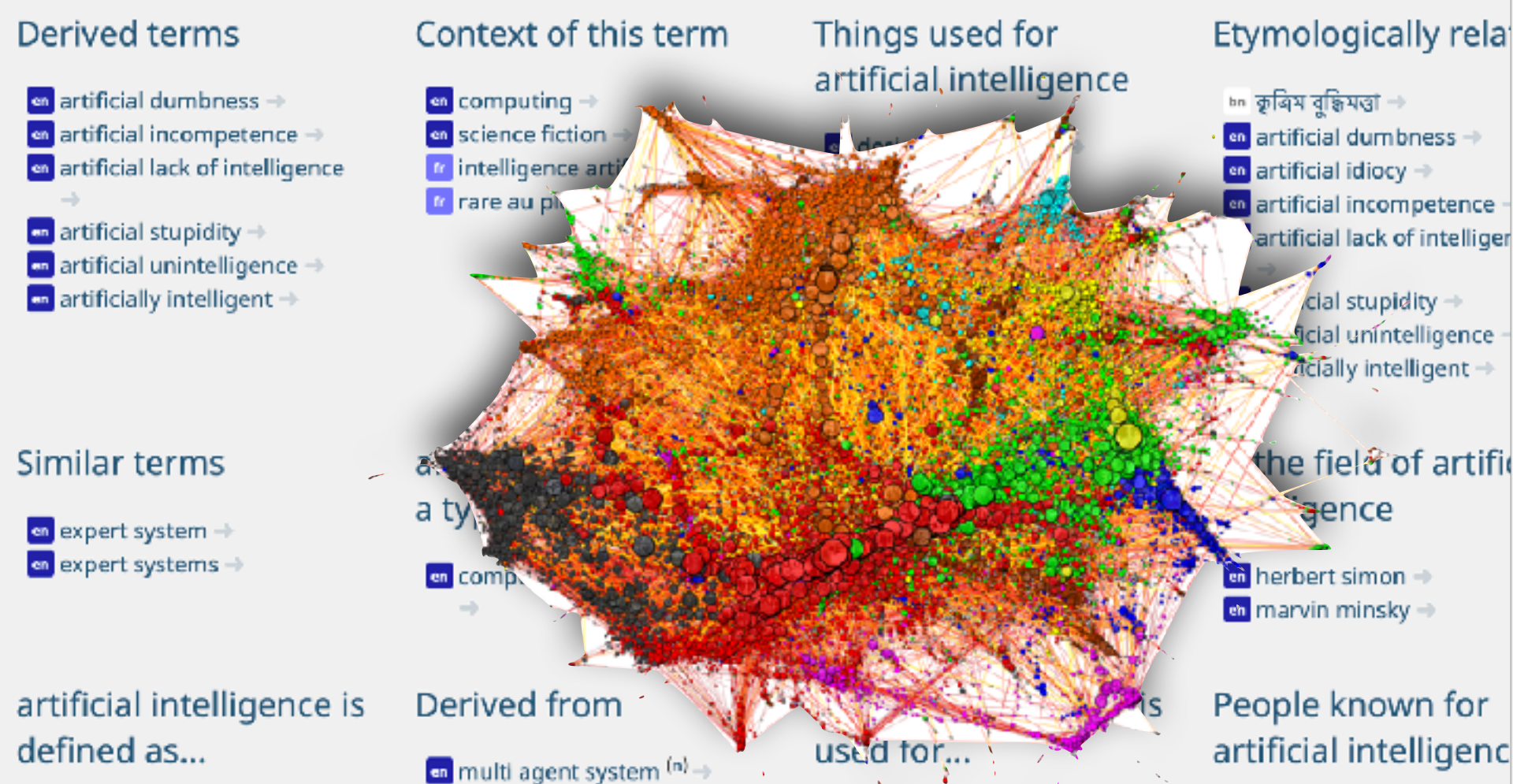
De-biasing Strategies

- Explicit:
 - Incorporate fairness in loss function
 - Incorporate other constraints during training
- Implicit:
 - Additive: Incorporate additional knowledge sources that could mitigate bias
 - ◆ More structured relationships should help to lessen biased data correlations
 - Subtractive: Identify features within model that influence bias and shift/eliminate
 - ◆ Prune out offending weights, while tracking performance trade off (if any)



ConceptNet, a Knowledge Graph

en artificial intelligence



ConceptNet Numberbatch



- Create with a Knowledge Graph (from multiple sources with relations like *UsedFor*, *PartOf*, etc.)
- Based KG edges, perturb existing embeddings (like GloVe) to optimize:

$$\Psi(Q) = \sum_{i=1}^n \left[\alpha_i \|q_i - \hat{q}_i\|^2 + \sum_{(i,j) \in E} \beta_{ij} \|q_i - q_j\|^2 \right]$$

\uparrow
new embed

 \uparrow
old embed

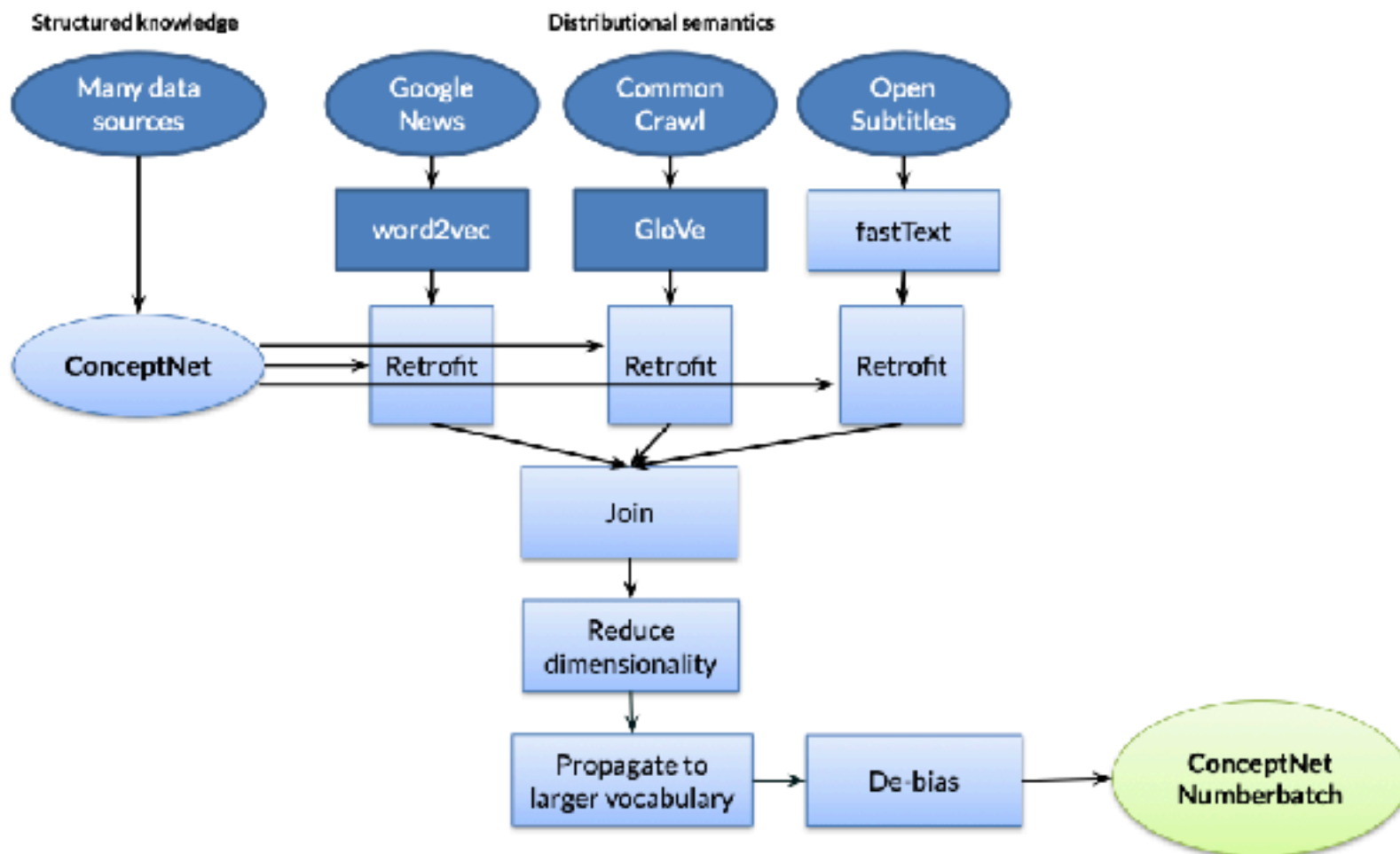
(keep similar to original)
(make similar according to other knowledge)

\nwarrow
neighbors from KG

- Easy to optimize the objective by averaging neighbors in the ConceptNet KG
- Multiple embeddings achieved by merging through “retrofitting” which projects onto a shared matrix space (with SVD)



Building ConceptNet Numberbatch



Aside: Transparency in Research

ConceptNet is all you need

Our full classifier used the linear combination of 5 types of input features shown above. This point is labeled **ABCDE** on the graph to the right. The other points are ablated versions of the classifier, trained on subsets of the five sources.

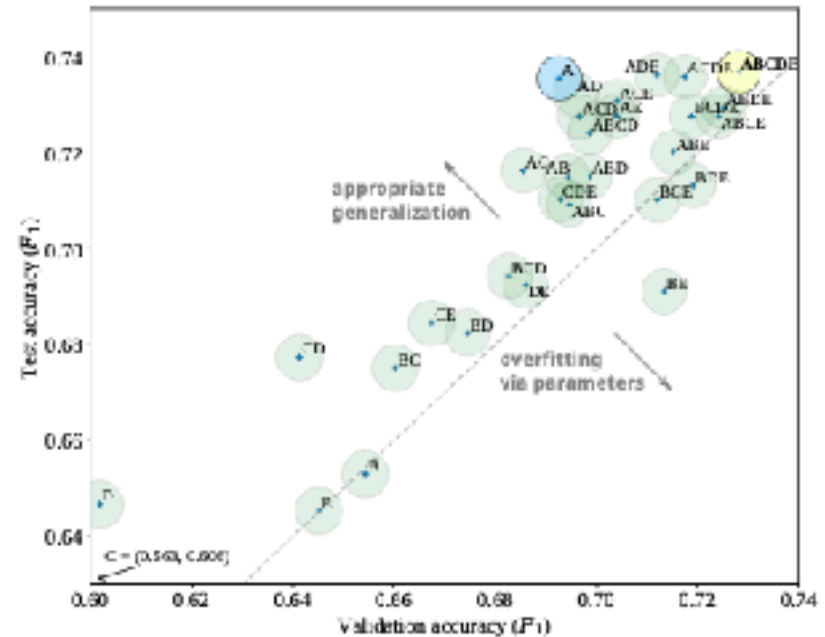
We found that the single feature of ConceptNet similarity (**A**) performed just as well on the test data as the full classifier, despite its lower validation accuracy.

This one-feature classifier could be more simply described as a heuristic over cosine similarities of ConceptNet embeddings:

$$\text{sim}(\text{term}_1, \text{attr}) - \text{sim}(\text{term}_2, \text{attr}) > 0.0961$$

It seems that the test data contained distinctions that can already be found by comparing ConceptNet embeddings, and that more complex features may have simply provided an opportunity to overfit to the validation set by parameter selection.

Results for all subsets of sources



This graph shows the validation and test accuracy of classifiers trained on subsets of the five sources of features. Ellipses indicate standard error of the mean, assuming that the data is sampled from a larger set.



ConceptNet Numberbatch

As a kid, I used to hold marble racing tournaments in my room, rolling marbles simultaneously down plastic towers of tracks and funnels. I went so far as to set up a bracket of 64 marbles to find the fastest marble. I kind of thought that running marble tournaments was peculiar to me and my childhood, but now I've found out that marble racing videos on YouTube are a big thing! Some of them even have **overlays as if they're major sporting events**.

In the end, there's nothing special about the fastest marble compared to most other marbles. It's just lucky. If one ran the tournament again, the marble champion might lose in the first round. But the one thing you could conclude about the fastest marble is that it was no *worse* than the other marbles. A bad marble (say, a misshapen one, or a plastic bead) would never luck out enough to win.

In our paper, we tested 30 alternate versions of the classifier, including the one that was roughly equivalent to this very simple system. We were impressed by the fact that it performed as well as our real entry. And this could be because of the inherent power of ConceptNet Numberbatch, or it could be because it's the lucky marble.



-Robyn Speer

<http://blog.conceptnet.io>

60

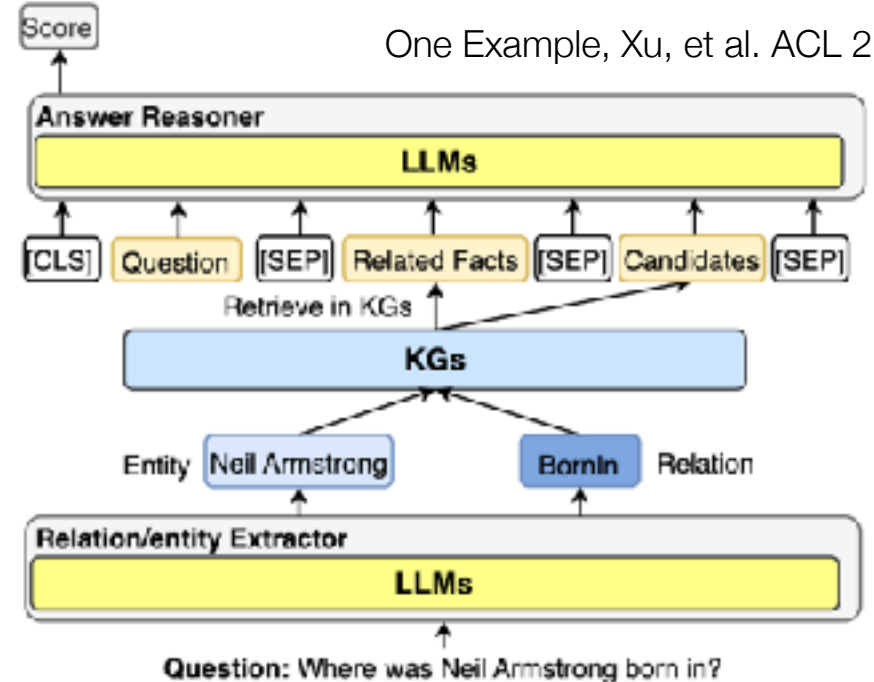
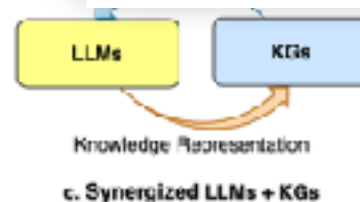


Beyond Word Embeddings

- ConceptNet Numberbatch is designed to help reduce bias in word embeddings analogy through incorporating knowledge graphs
- How can this be studied by LLMs?
 - Such as transformer models

Unifying Large Language Models and Knowledge Graphs: A Roadmap

Shirui Pan, Senior Member, IEEE, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, Xindong Wu, Fellow, IEEE



One Example, Xu, et al. ACL 21



BiasWipe

- Run model explainability with feature sensitivity or importances
- Identify counterfactual groups and measure bias
- Identify weights that are most influential for the bias (k-top)
- Prune weights (make zero)

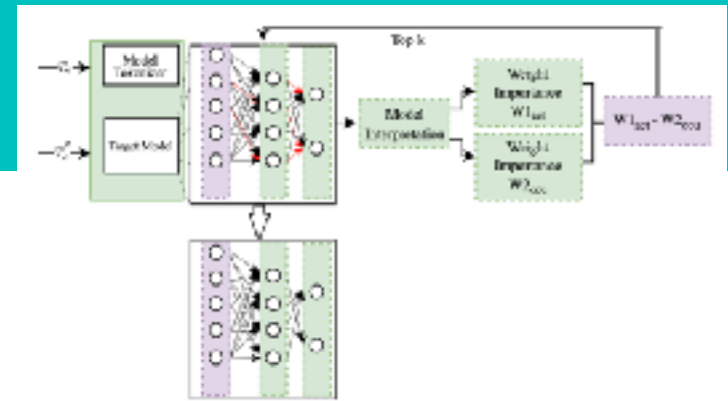
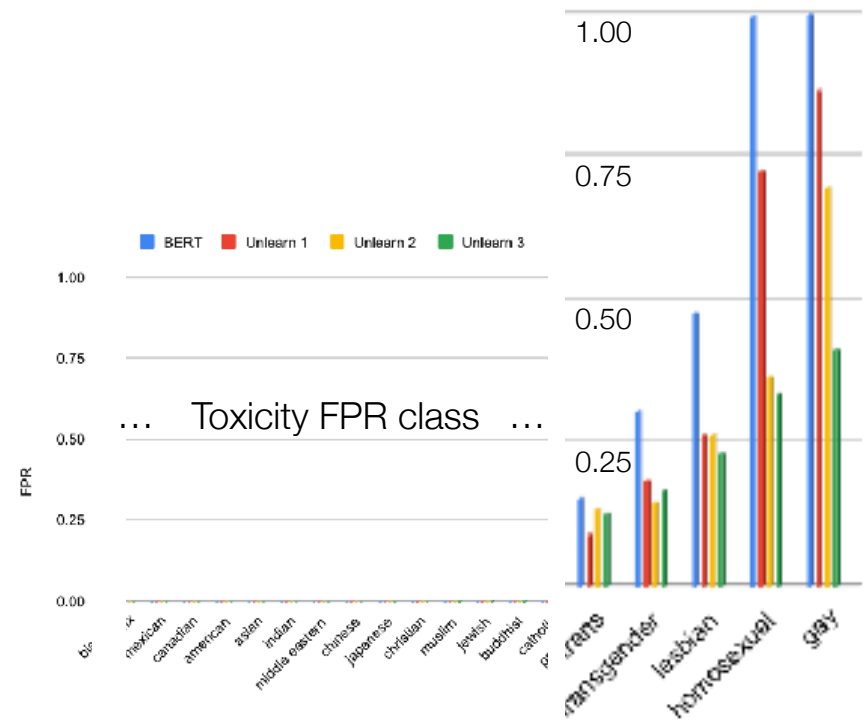
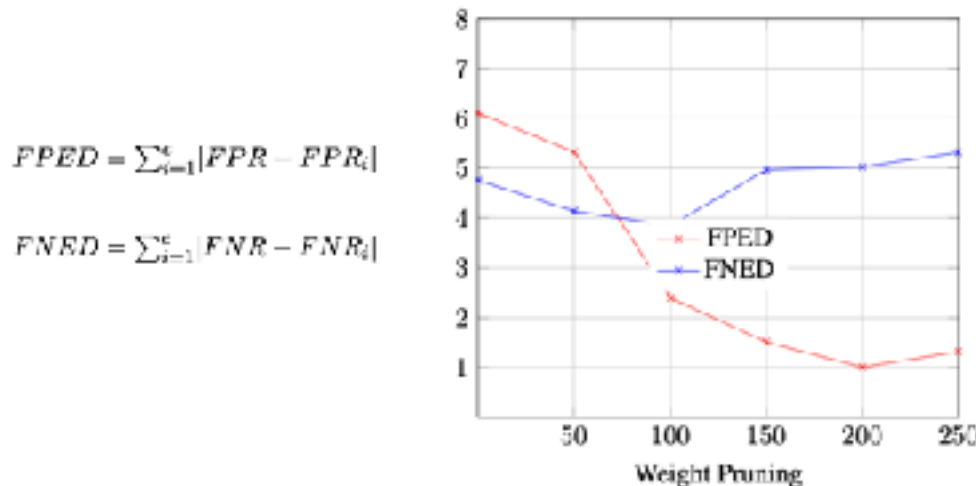


Figure 1: Proposed Workflow of *BiasWipe*.





How to Make a Racist AI without Really Trying



Robyn Speer, 2017

<http://blog.conceptnet.io/posts/2017/how-to-make-a-racist-ai-without-really-trying/>

Debiasing: Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings

Bolukbasi et al., NeurIPs 2016

<https://arxiv.org/pdf/1607.06520.pdf>

ConceptNet 5.5: An Open Multilingual Graph of General Knowledge

Speer et al., AAAI 2017

<https://arxiv.org/pdf/1612.03975.pdf>



Rachael Tatman @rctatman · 18h

I first got interested in ethics in NLP/ML because I was asking "does this system work well for everyone". It's a good question, but there's a more important one:

Who is being harmed and who is benefiting from this system existing in the first place?



Lecture Notes for **Neural Networks and Machine Learning**

Ethically Aware Practices



Next Time:
Transfer Learning
Reading: Chollet Article

