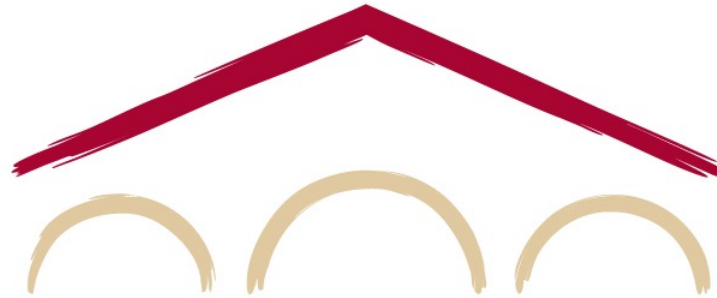


Natural Language Processing with Deep Learning

CS224N/Ling284



Tatsunori Hashimoto
Open problems and discussion

Lecture plan + Logistics

Announcements:

- Final project report + Code: Due March 14th, EOD
 - No late submissions past March 18th 7:00PM
- Poster session: March 18th 7:00 PM
 - Location: AOERC
- Final project poster PDF: March 18th, EOD
 - Will not be accepted late (due to tight grading timelines!)

Lecture plan

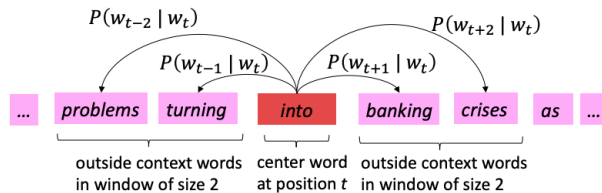
1. Important ideas from 224n for you to remember
2. A brief, broad (and incomplete!) survey of some open problems.
3. Post-lecture QA

Major ideas in 224n

Word vectors

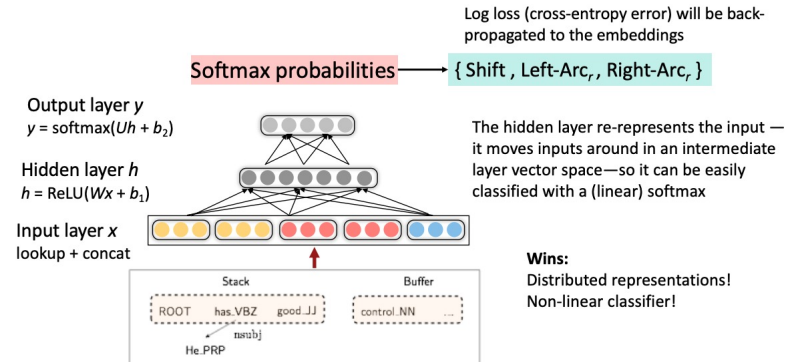
Word2Vec Overview

Example windows and process for computing $P(w_{t+j} | w_t)$



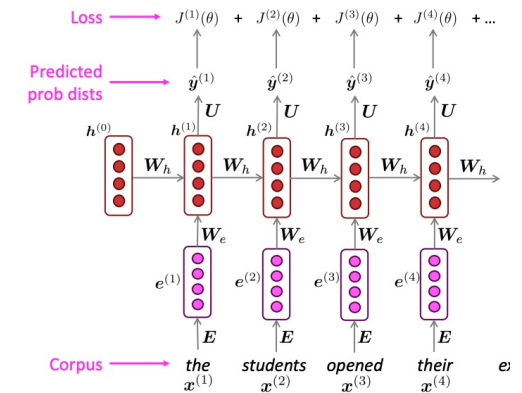
Neural NLP

(A simple feed-forward neural network multi-class classifier)



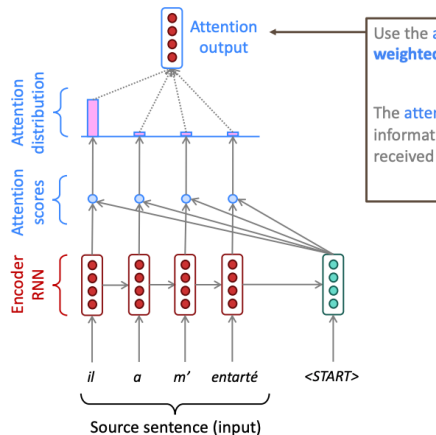
RNNs/LSTMs

Training an RNN Language Model



Transformers, LMs

Sequence-to-sequence with attention



Pretraining, Post-training

What kinds of things does pretraining teach?

There's increasing evidence that pretrained models learn a wide variety of things about the statistical properties of language. Taking our examples from the start of class:

- *Stanford University is located in _____, California.* [Trivia]
- *I put ____ fork down on the table.* [syntax]
- *The woman walked across the street, checking for traffic over ____ shoulder.* [coreference]
- *I went to the ocean to see the fish, turtles, seals, and ____.* [lexical semantics/topic]
- *Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ____.* [sentiment]
- *Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the ____.* [some reasoning – this is harder]
- *I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, ____* [some basic arithmetic; they don't learn the Fibonacci sequence]
- Models also learn – and can exacerbate racism, sexism, all manner of bad biases.

... and more

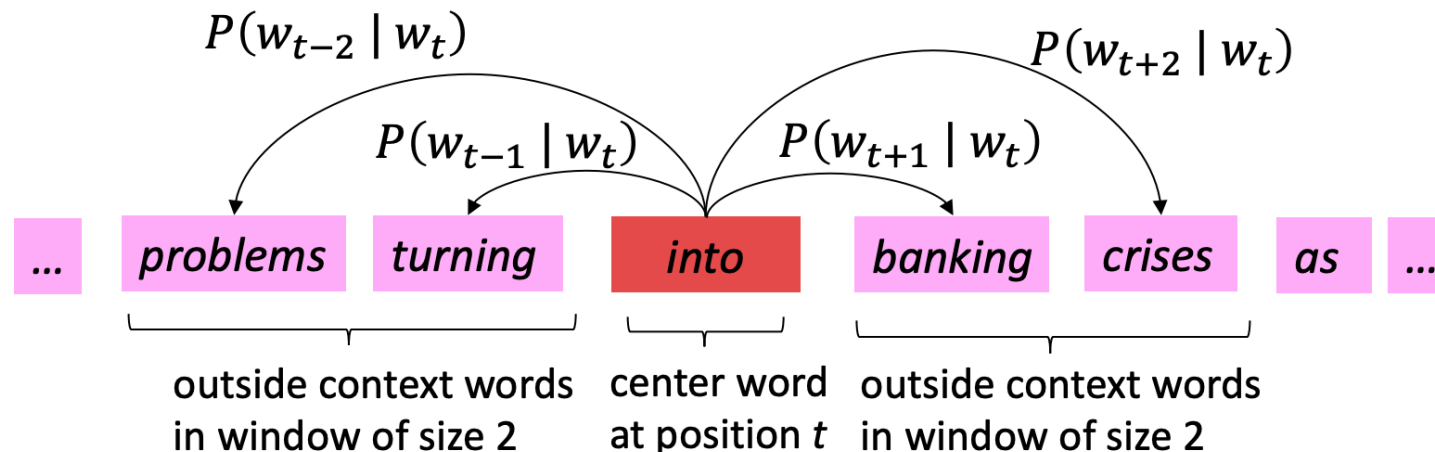
Idea 1: Dense representations and distributional semantics

Representing words by their context



- **Distributional semantics:** A word's meaning is given by the words that frequently appear close-by
 - *"You shall know a word by the company it keeps"* (J. R. Firth 1957: 11)
 - One of the most successful ideas of modern statistical NLP!

Instantiation: predict neighboring words using a word vector representation



Idea 2: Depth and neural networks

The challenges of training large, deep nets

Many tricks like residual connections

But building large neural networks isn't easy or obvious

Greedy Layer-Wise Training of Deep Networks

Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle
Université de Montréal
Montréal, Québec
{bengioy, lamblinp, popovicd, larocheh}@iro.umontreal.ca

much less expressive than deep ones.

However, until recently, it was believed too difficult to train deep multi-layer neural networks. Empirically, deep networks were generally found to be not better, and often worse, than neural networks with one or two hidden layers (Tesauro, 1992). As this is a negative result, it has not been much reported in the machine learning literature. A reasonable explanation is that gradient-based optimization starting from random initialization may get stuck near poor solutions. An approach that has been explored with some success in the past is based on *constructively* adding layers. This was previously done using a supervised criterion at each stage (Eckman & LeBlond, 1990; Lencz & Deneys, 1996). Winter

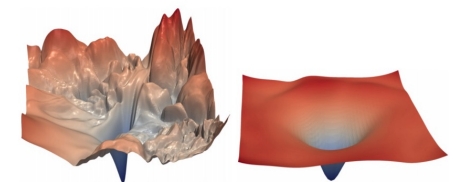
- **Residual connections** are a trick to help models train better.
 - Instead of $X^{(i)} = \text{Layer}(X^{(i-1)})$ (where i represents the layer)

$$X^{(i-1)} \longrightarrow \text{Layer} \longrightarrow X^{(i)}$$

- We let $X^{(i)} = X^{(i-1)} + \text{Layer}(X^{(i-1)})$ (so we only have to learn “the residual” from the previous layer)

$$X^{(i-1)} \longrightarrow \text{Layer} \oplus \longrightarrow X^{(i)}$$

- Gradient is **great** through the residual connection; it's 1!
- Bias towards the identity function!



[no residuals]

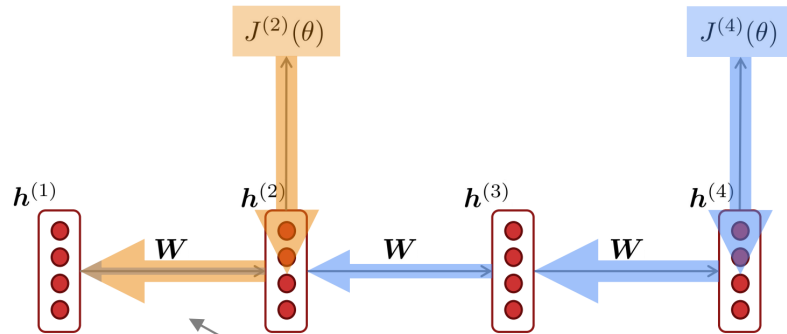
[residuals]

[Loss landscape visualization,
[Li et al., 2018](#), on a ResNet]

Idea 3: Sequence models and computational efficiency

Vanishing gradients and information propagation

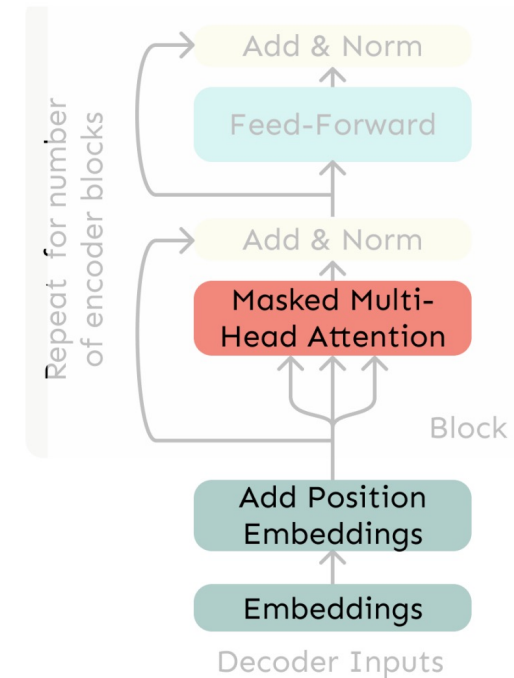
Why is vanishing gradient a problem?



Gradient signal from far away is lost because it's much smaller than gradient signal from close-by.

So, model weights are updated only with respect to near effects, not long-term effects.

Parallelization and dependency benefits of Transformers



Transformer Decoder

Idea 4: Language modeling and pretraining

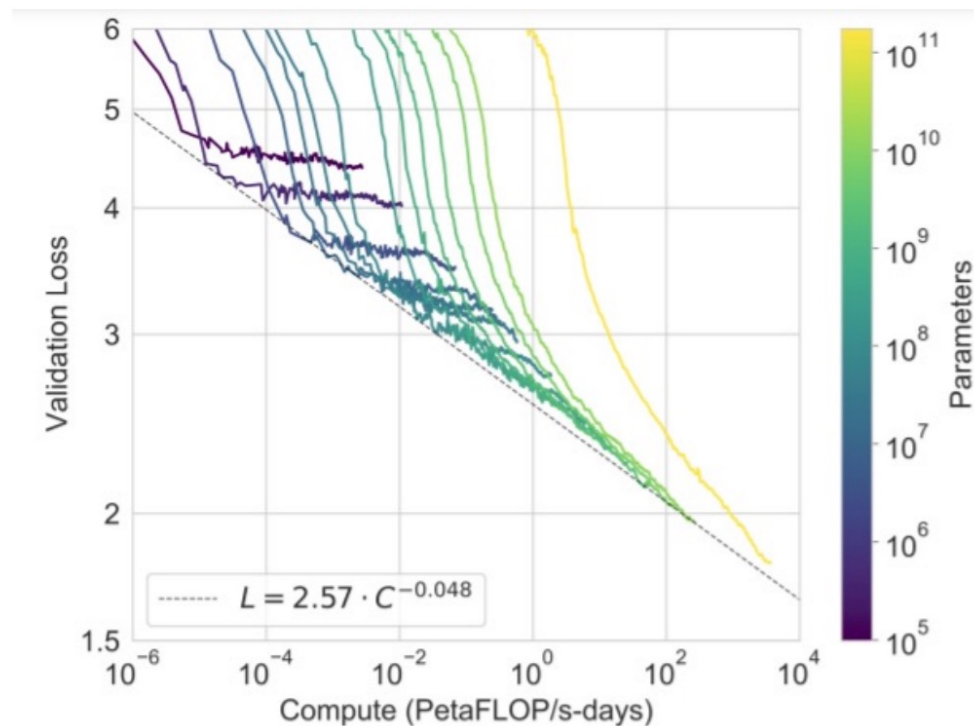
Language modeling as
a ‘universal’ pretraining task

What kinds of things does pretraining teach?

There’s increasing evidence that pretrained models learn a wide variety of things about the statistical properties of language. Taking our examples from the start of class:

- *Stanford University is located in _____, California.* [Trivia]
- *I put ____ fork down on the table.* [syntax]
- *The woman walked across the street, checking for traffic over ____ shoulder.* [coreference]
- *I went to the ocean to see the fish, turtles, seals, and ____.* [lexical semantics/topic]
- *Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ____.* [sentiment]
- *Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the ____.* [some reasoning – this is harder]
- *I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, ____* [some basic arithmetic; they don’t learn the Fibonacci sequence]
- Models also learn – and can exacerbate racism, sexism, all manner of bad biases.

Careful and smart scaling to
internet-wide data



Open problems – a quick (and incomplete) overview

Generalization

Article: Super Bowl 50

Paragraph: "Peython Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had a jersey number 37 in Champ Bowl XXXIV."

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?"

Original Prediction: John Elway

Prediction under adversary: Jeff Dean

Analysis and understanding

input
sentence



Your final project
model

Accuracy: ____ %



output
prediction

Fig 1. A black box

Multilingual

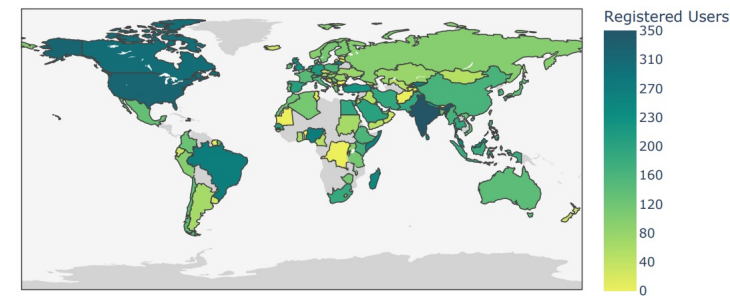


Figure 2: Geographical distribution of the users registered on the Aya platform.

Evaluations



Domains + modalities

THE NEW ENGLAND JOURNAL of MEDICINE

SPECIAL REPORT

Jeffrey M. Drazen, M.D., Editor;
Isaac S. Kohane, M.D., Ph.D., and Tze-Yun Leong, Ph.D., Guest Editors

AI IN MEDICINE

Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine

Peter Lee, Ph.D., Sebastien Bubeck, Ph.D., and Joseph Petro, M.S., M.Eng.

Fairness and social

Do our models understand our tasks?

Article: Super Bowl 50

Paragraph: "Peython Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had a jersey number 37 in Champ Bowl XXXIV."

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?"

Original Prediction: John Elway

Prediction under adversary: Jeff Dean

Answer question

Changes answer if irrelevant information is added

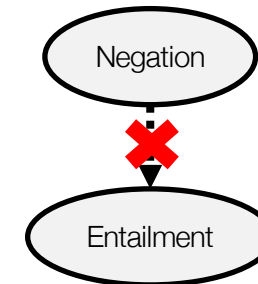
Only looks at last sentence and ignores context

Premise:

The economy could be still better.

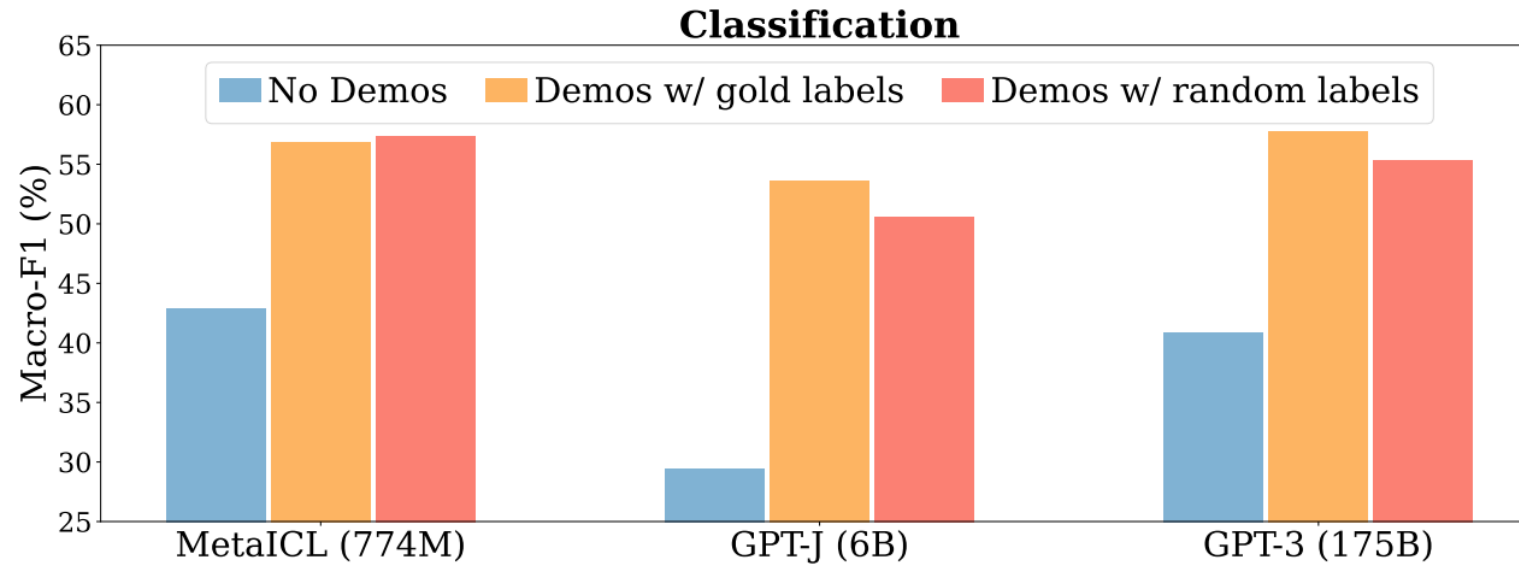
Hypothesis:

The economy has **never** been better



- Many instances of models using shortcuts over deep understanding

How much do models really generalize (Generalization)

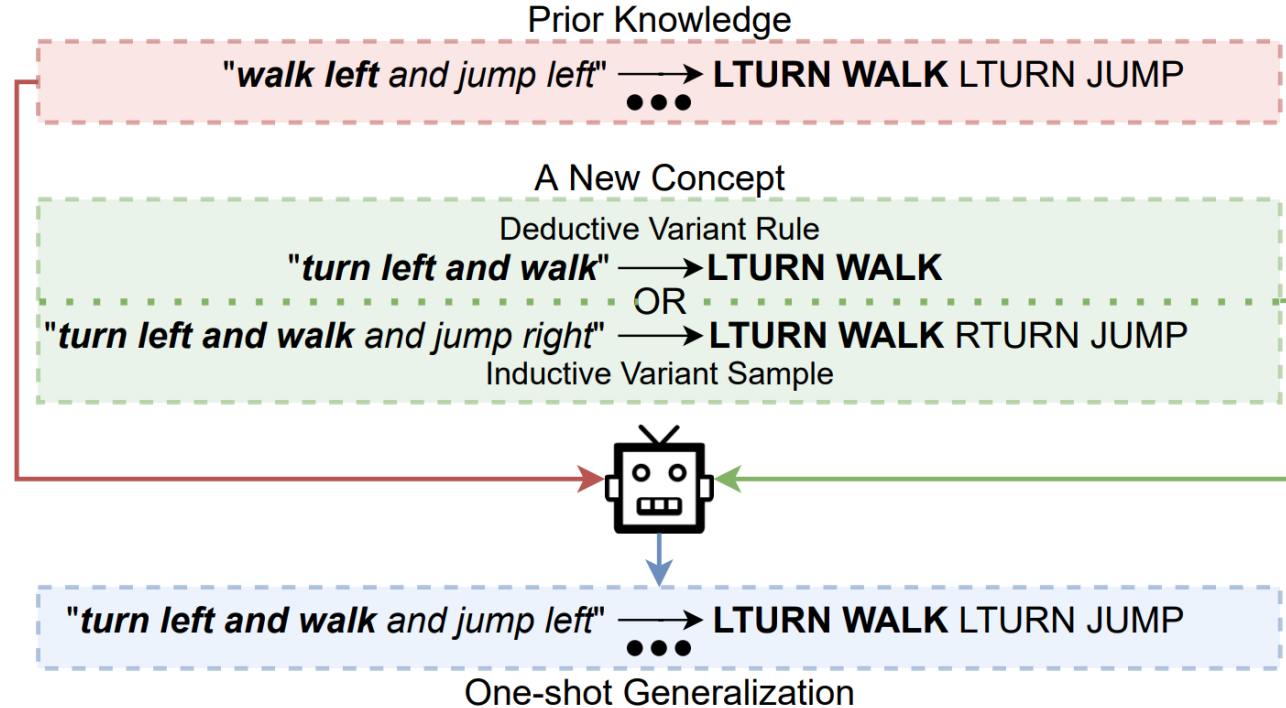


In-context learning with *random labels* does just as well as ICL with real data (Min, 2022)

- Even modern LLMs seem to leverage surface cues – are we just finding better shortcuts?

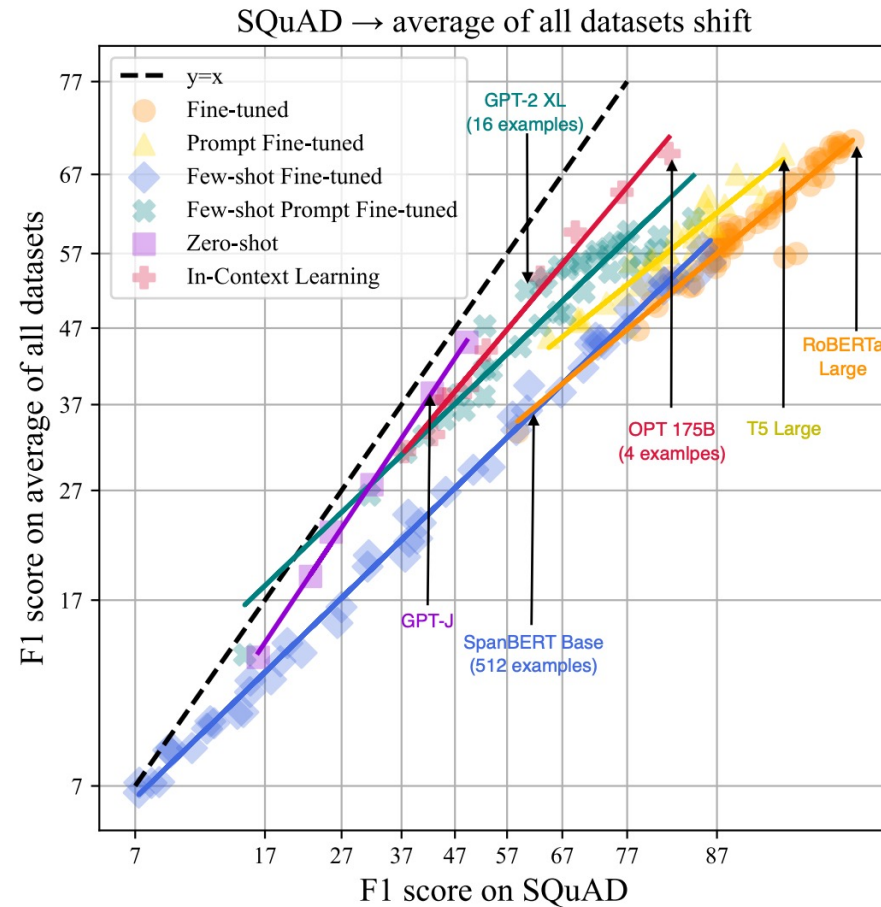
How do we make models go beyond train data? (Generalization)

SCAN dataset and systematic generalization



- Can we generalize in humanlike ways, from little data?

Is there much beyond optimizing IID error? (Generalization)



- Despite our best intuitions, the best models on average are also most robust

Analysis and interpretation – What makes NLP systems work?



What's going on inside NNs? (Analysis)

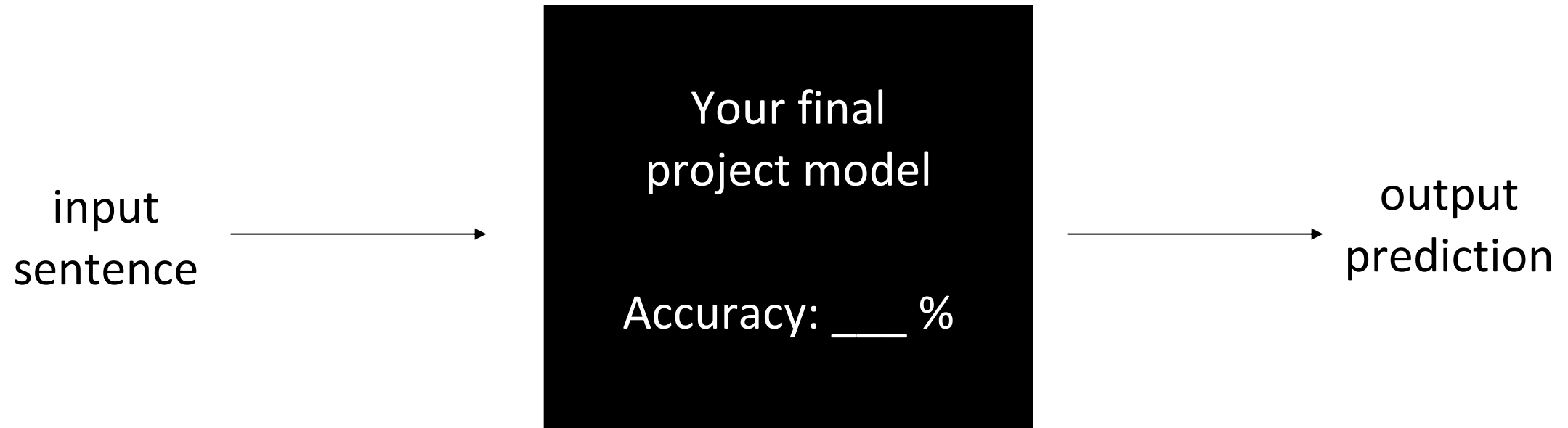


Fig 1. *A black box*

We summarize our models with one (or a handful) of accuracies metric numbers.
What do they learn? Why do they succeed and fail?

Old results already show interpretable latent units

Idea: Individual hidden units can lend themselves to an interpretable meaning.

This model: a character-level LSTM language model.

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Here, “cell” refers to a single dimension of the cell state of the LSTM.

Can we build interpretable, but performant models? (Analysis)

All Models are Wrong, but *Many* are Useful: Learning a Variable's Importance by Studying an Entire Class of Prediction Models Simultaneously

Aaron Fisher

*Takeda Pharmaceuticals
Cambridge, MA 02139, USA*

AFISHE27@ALUMNI.JH.EDU

Cynthia Rudin

*Departments of Computer Science and Electrical and Computer Engineering
Duke University
Durham, NC 27708, USA*

CYNTHIA@CS.DUKE.EDU

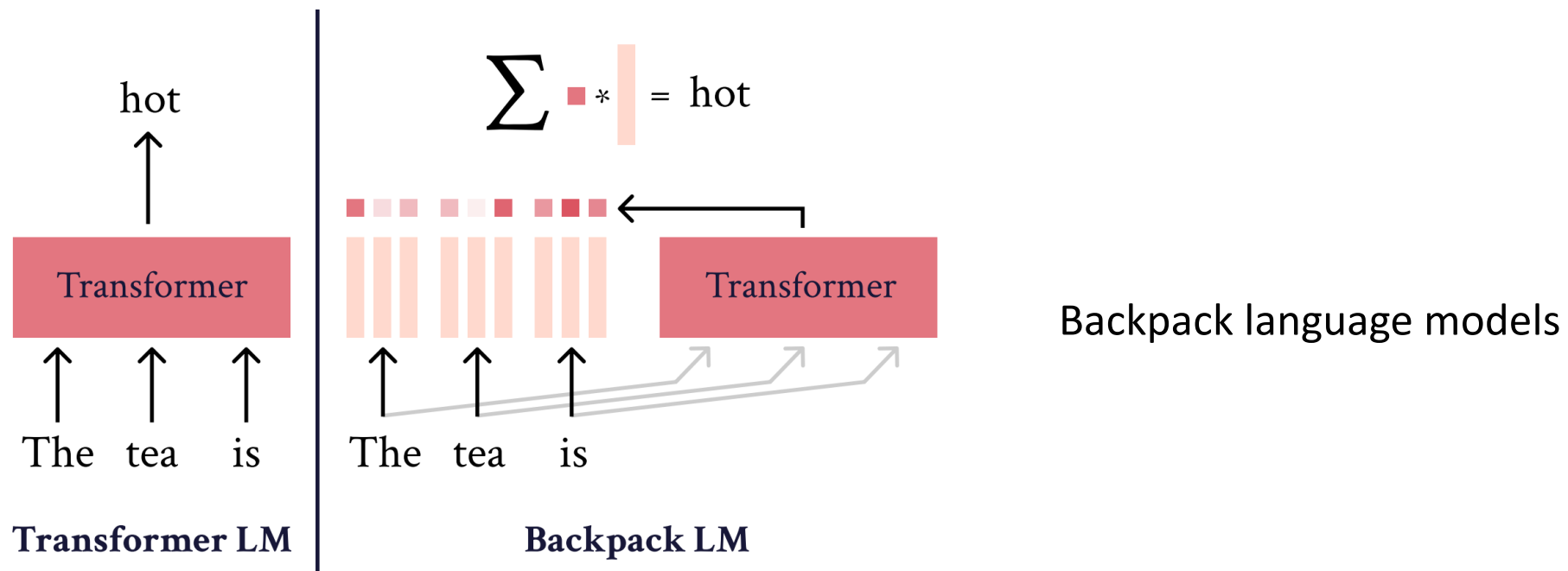
Francesca Dominici

*Department of Biostatistics
Harvard T.H. Chan School of Public Health
Boston, MA 02115, USA*

FDOMINIC@HSPH.HARVARD.EDU

To address these concerns, we analyze the *set* of prediction models that provide near-optimal accuracy, which we refer to as a *Rashomon set*. This approach stands in contrast to training to select a *single* prediction model, among a prespecified class of candidate models. Our motivation is that Rashomon sets (defined formally below) summarize the range of effective prediction strategies that an analyst might choose. Additionally, even if the candidate models do not contain the true data generating process, we may hope that some of these models function in similar ways to the data generating process. In particular, we may hope there exist well performing candidate models that place the same importance on a variable of interest as the underlying data generating process does. If so, then studying sets of well-performing models will allow us to deduce information about the data generating process.

Can we build interpretable, performant models?



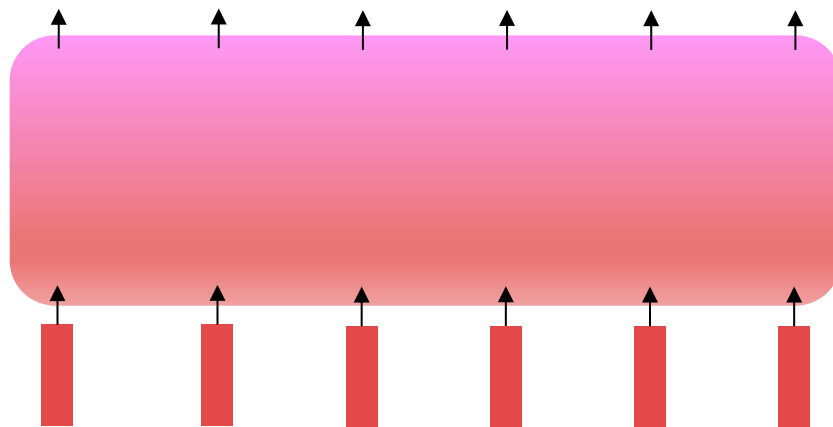
Model	OpenWebText PPL ↓	LAMBADA PPL ↓	LAMBADA ACC ↑	Wikitext PPL ↓	BLiMP ↑
Backpack-Micro	31.5	110	24.7	71.5	75.6
Transformer-Micro	34.4	201	21.3	79.5	77.8
Backpack-Mini	23.5	42.7	31.6	49.0	76.2
Transformer-Mini	24.5	58.8	29.7	52.8	80.4
Backpack-Small	20.1	26.5	37.5	40.9	76.3
Transformer-Small	20.2	32.7	34.9	42.2	81.9

Can understanding help find the next transformer? (Analysis)

What can be learned via language model pretraining?

What will replace the Transformer?

What **can't** be learned via language model pretraining?



What does deep learning struggle to do?

How are our models affecting people, and transferring power?

What do neural models tell us about language?

Multilingual

GPT-4 3-shot accuracy on MMLU across languages

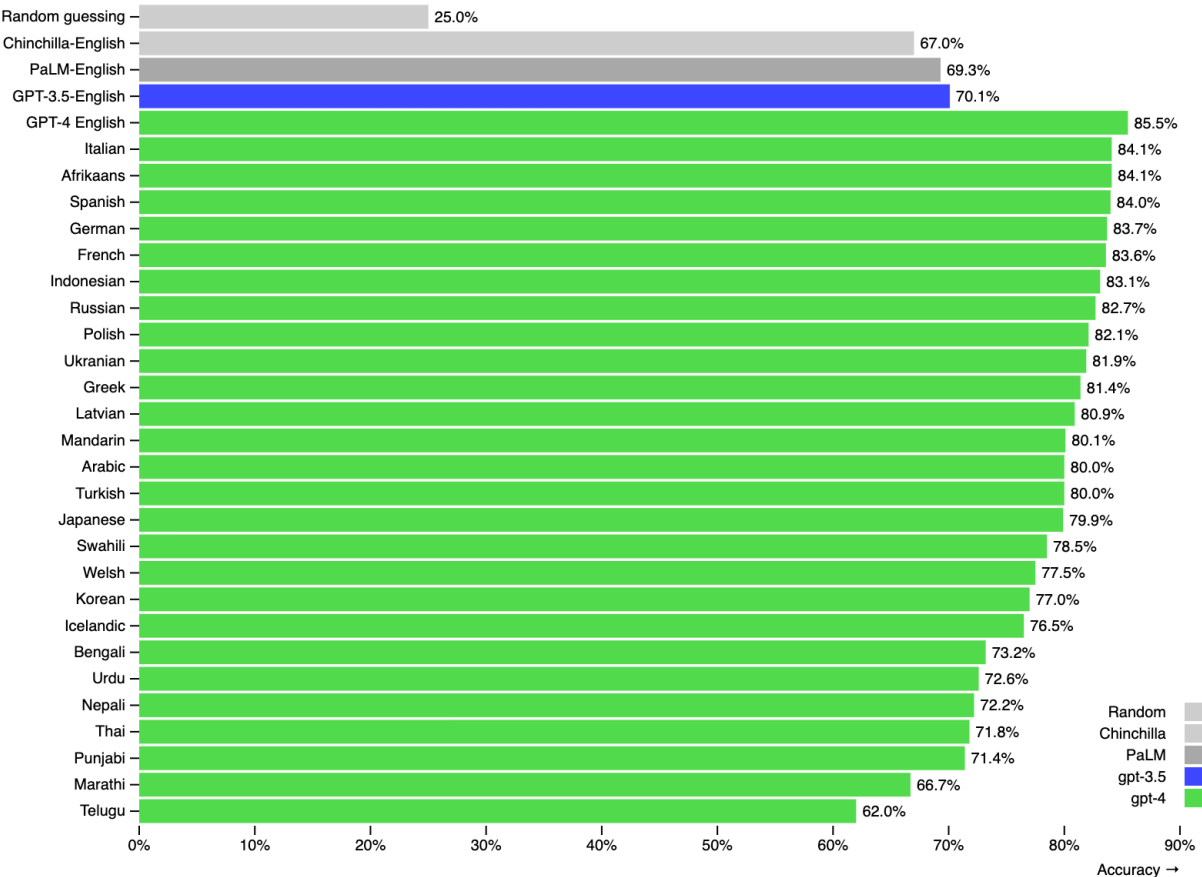
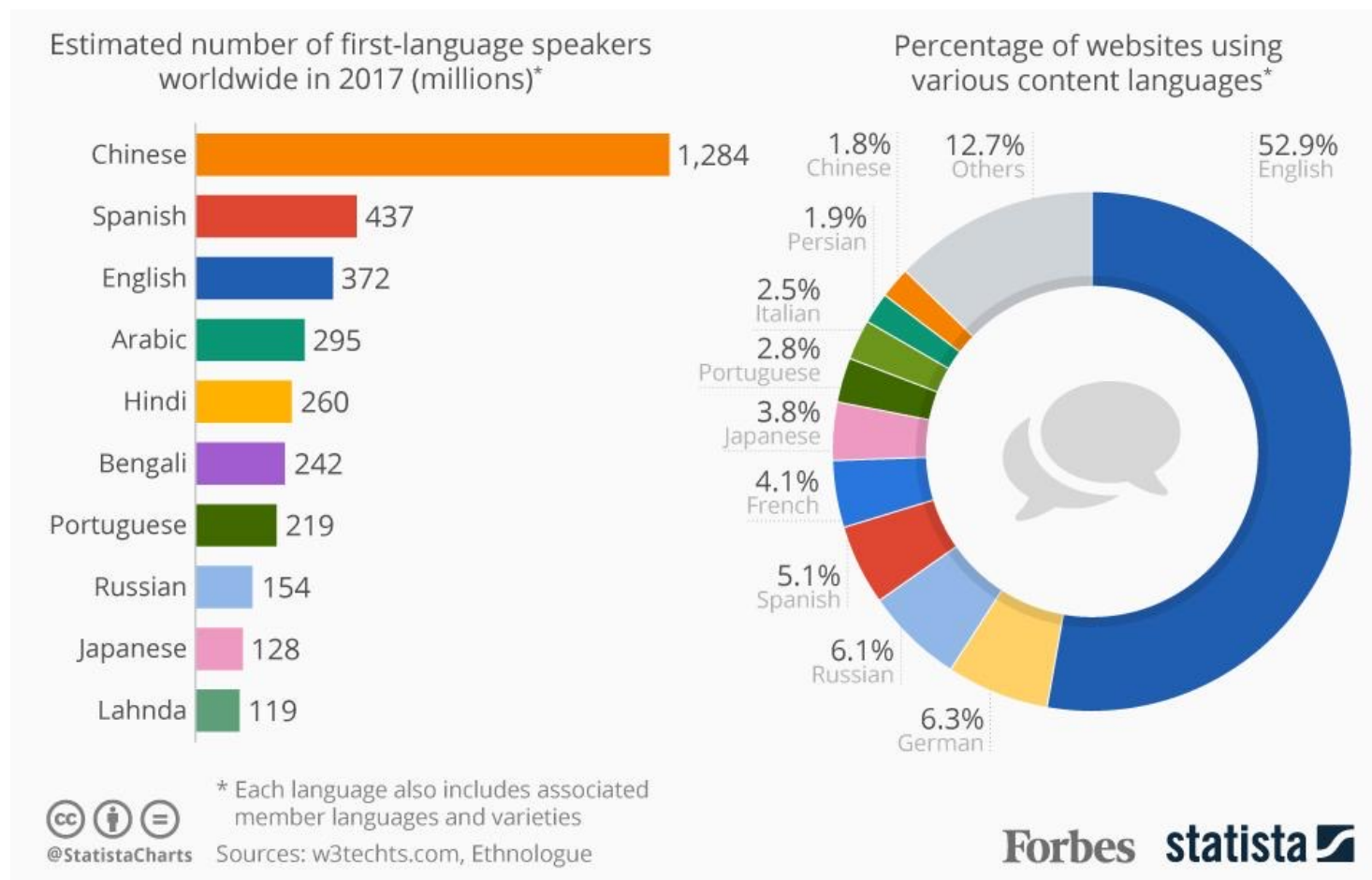


Table 19: Per-language one-shot results on Wikilingua, XL-Sum, and XSum. We report ROUGE-2 for English and SentencePiece-ROUGE-2 for all other languages.

		PaLM	PaLM 2-S	PaLM 2-M	PaLM 2-L
Wikilingua	Arabic	11.18	15.11	17.0	22.23
	Japanese	11.97	15.66	16.4	21.27
	Korean	13.94	19.62	19.5	24.69
	Russian	10.15	14.40	15.0	20.25
	Thai	11.59	20.87	21.8	29.64
	Turkish	11.49	15.17	15.8	23.05
XL-Sum	Arabic	14.43	19.59	21.5	25.42
	Bengali	9.01	19.51	18.4	15.87
	English	14.82	15.2	16.8	23.03
	Indonesian	15.7	19.5	19.9	22.71
	Japanese	14.52	21.8	24.8	28.97
	Korean	13.87	17.6	19.2	21.97
	Russian	12.51	15.6	16.1	18.4
	Swahili	16.43	23.0	23.6	27.45
	Telugu	5.59	14.6	13.5	6.86
	Thai	7.54	17.0	17.3	18.94
	Turkish	15.16	18.1	19.4	23.16
	XSum-English	14.57	16.91	17.2	23.23
	Average	12.47	17.74	18.51	22.06

Significant gaps between high resource (English) and others (Telugu) even for the best models

Can we remove language resource gaps? (Multilingual)



Data drives pretraining – gaps in data availability lead to performance gaps.
How can we close this?

Working with extremely low resource languages (Multilingual)

It is well known that only a very limited proportion of the languages spoken in the world is covered by technology or by scientific knowledge. For technology, only normative productions of very few languages in very few situations are mastered. The technological divide is wide considering the languages spoken: we have a minimally adequate quantity of data for less than 1% of the world's 7000 languages. Most of the world's everyday life speech stems from languages which are essentially unwritten and we include in these languages ethnolects as well as sociolects such as many regional varieties of Arabic, Shanghainese, slang ... There are thousands of endangered languages for which hardly any documentation exists and time is running out before they disappear: some linguists estimate that half of the presently living languages will become extinct in the course of this century^{1,2,3}. Even with the upsurge of documentary

[Adda et al 2016]

- Most languages do not have machine-readable, written text
- Many such languages may become extinct
- Little for-profit motive to serve these languages – vicious feedback loop

Data quality is very variable in multilingual corpora

20% of languages are less than 50% correct

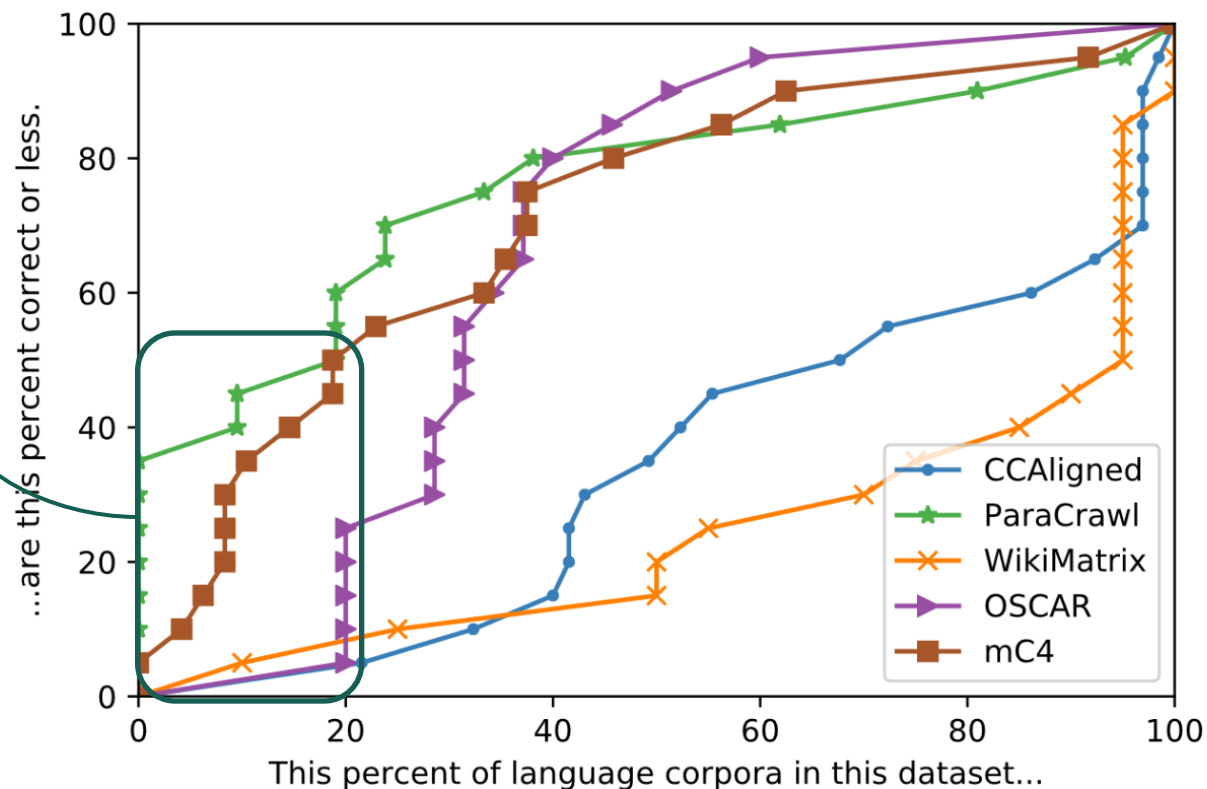


Figure 1: Fraction of languages in each dataset below a given quality threshold (percent correct).

[[Kreutzer et al 2022](#);
[Nekoto et al 2020](#)]

How can we get better, multilingual training data? (Multilingual)

No Language Left Behind: Scaling Human-Centered Machine Translation

NLLB Team, Marta R. Costa-jussà*, James Cross*, Onur Çelebi*, Maha Elbayad*, Kenneth Heafield*, Kevin Heffernan*, Elahe Kalbassi*, Janice Lam*, Daniel Licht*, Jean Maillard*, Anna Sun*, Skyler Wang*[§], Guillaume Wenzek*, Al Youngblood*
Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran
Pierre Andrews[†], Necip Fazil Ayan[†], Shruti Bhosale[‡], Sergey Edunov[‡], Angela Fan^{†,‡}, Cynthia Gao[‡], Vedanuj Goswami[‡], Francisco Guzmán[‡], Philipp Koehn^{†,¶}, Alexandre Mourachko[‡], Christophe Ropers[‡], Safiyyah Saleem[‡], Holger Schwenk[‡], Jeff Wang[‡]

Meta AI, [§]UC Berkeley, [¶]Johns Hopkins University

Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning

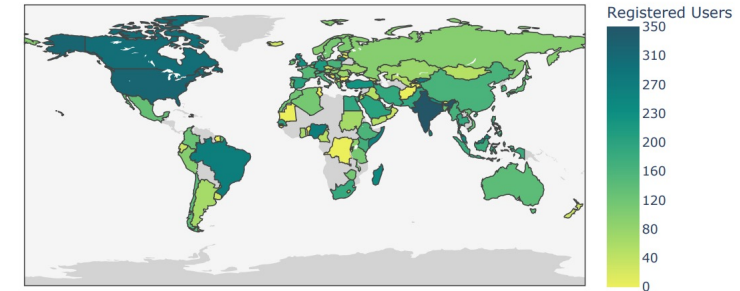
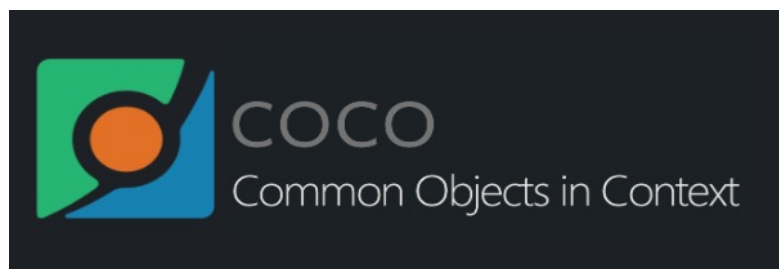
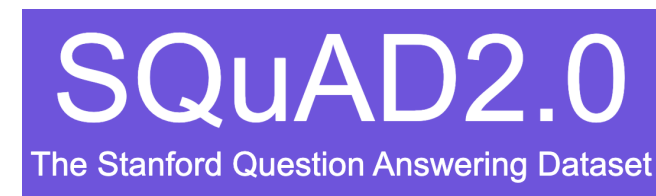


Figure 2: Geographical distribution of the users registered on the **Aya** platform.

- Very difficult to scale things like machine translation or instruction tuning data
 - Automatic (alignment) approaches
 - Crowdsourced approaches

Evaluation and comparison



EMNLP 2022
SEVENTH CONFERENCE ON
MACHINE TRANSLATION (WMT22)

December 7-8, 2022
Abu Dhabi

Shared Task: General Machine Translation

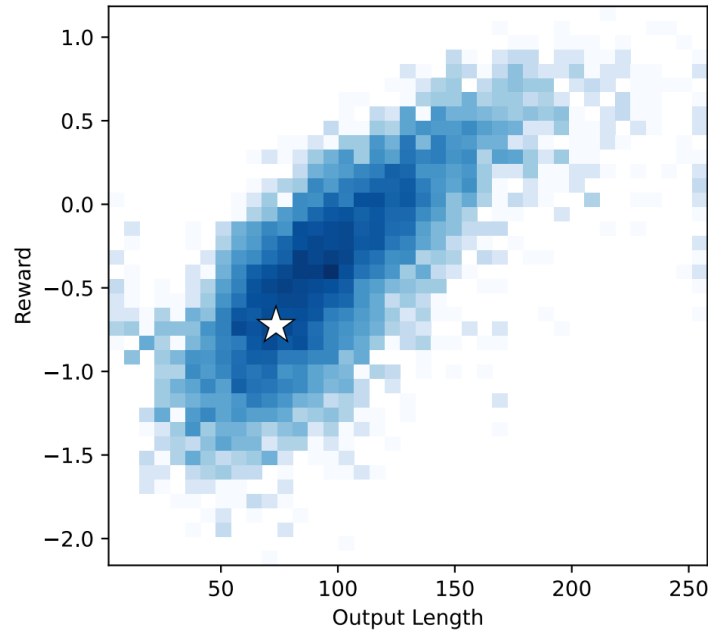
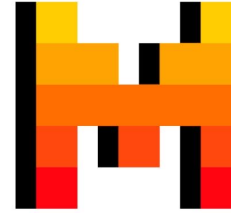


Benchmarks and how we evaluate drive the progress of the field

Evaluating chatbots and open-domain systems (Evaluation)



VS



Question: *Why don't adults roll off the bed?*

☆ **SFT (Before); 59 tokens**

Adults typically do not roll off of the bed because they have developed the muscle memory to keep their bodies from involuntarily moving during sleep and maintaining proper posture.

RLHF (After); 243 tokens: Similar output, but much longer / more details

Adults generally do not roll off of the bed because they have developed muscle control and awareness which helps them to remain upright while sleeping.

Additionally, most adults find it uncomfortable or even painful to move around in their sleep, so rolling off the bed would be difficult without waking up first. There may also be safety concerns such as falling out of bed and potentially hurting oneself if one were to lose balance or fall from the ...

We know our evals are broken in some known ways (length) but we have no idea what other ways they're broken

What's the gold standard in evaluation? (Evaluation)

HUMAN FEEDBACK IS NOT GOLD STANDARD

Tom Hosking
University of Edinburgh
tom.hosking@ed.ac.uk

Phil Blunsom
Cohere
phil@cohere.com

Max Bartolo
Cohere, UCL
max@cohere.com

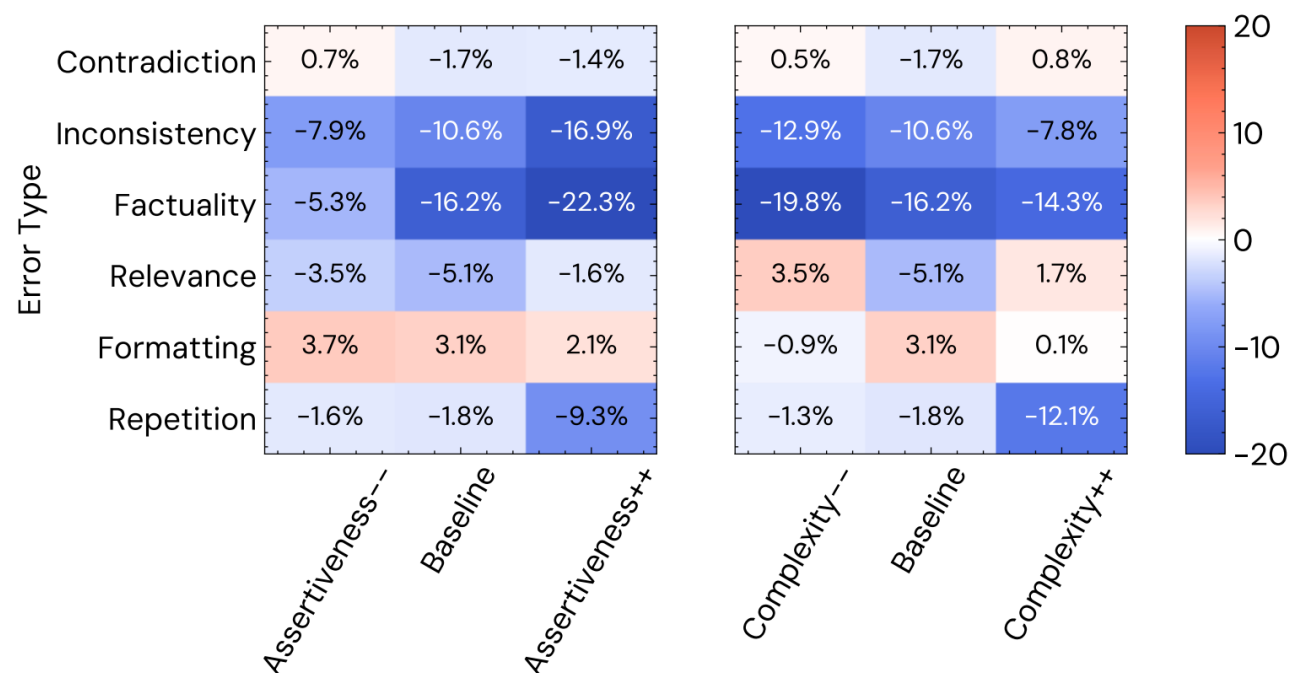
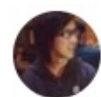


Figure 4: The difference in error rates between crowdsourced annotations and 'expert' annotations from the authors, excluding samples that were marked as refusing to respond. Annotators tend to

How do we maintain benchmark integrity? (Evaluation)



Horace He
@cHHillee



I suspect GPT-4's performance is influenced by data contamination, at least on Codeforces.

Of the easiest problems on Codeforces, it solved 10/10 pre-2021 problems and 0/10 recent problems.

This strongly points to contamination.

1/4

g's Race	implementation, math				greedy, implementation		
nd Chocolate	implementation, math			Cat?	implementation, strings		
triangle!	brute force, geometry, math			Actions	data structures, greedy, implementation, math		
	greedy, implementation, math			Interview Problem	brute force, implementation, strings		

- Pretraining makes it increasingly difficult to build hidden test sets

How do we evaluate things like interpretability (Evaluation)

Attention is not Explanation

Sarthak Jain

Northeastern University

jain.sar@husky.neu.edu

Byron C. Wallace

Northeastern University

b.wallace@northeastern.edu

after 15 minutes watching the movie i was asking myself what to do leave the theater sleep or try to keep watching the movie to see if there was anything worth i finally watched the movie what a waste of time maybe i am not a 5 years old kid anymore

original α

$$f(x|\alpha, \theta) = 0.01$$

after 15 minutes watching the movie i was asking myself what to do leave the theater sleep or try to keep watching the movie to see if there was anything worth i finally watched the movie what a waste of time maybe i am not a 5 years old kid anymore

adversarial $\tilde{\alpha}$

$$f(x|\tilde{\alpha}, \theta) = 0.01$$

Are Model Explanations Useful in Practice? Rethinking How to Support Human-ML Interactions.

AUTHORS

Valerie Chen

Ameet Talwalkar

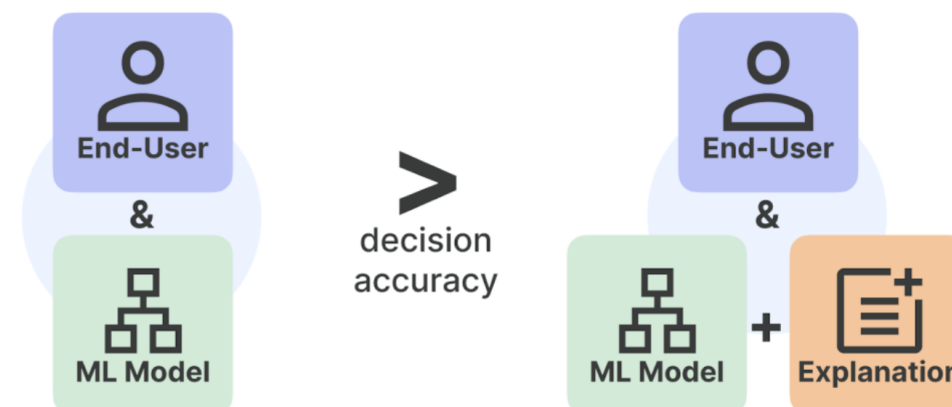
AFFILIATIONS

MLD, CMU

MLD, CMU

PUBLISHED

March 31, 2023



- Increasingly, many things in research are qualitative. How do we evaluate those?

Making NLP work in high-impact domains

NLP systems (and LLMs) are going from the lab to the real world

Clinical

THE NEW ENGLAND JOURNAL OF MEDICINE

SPECIAL REPORT

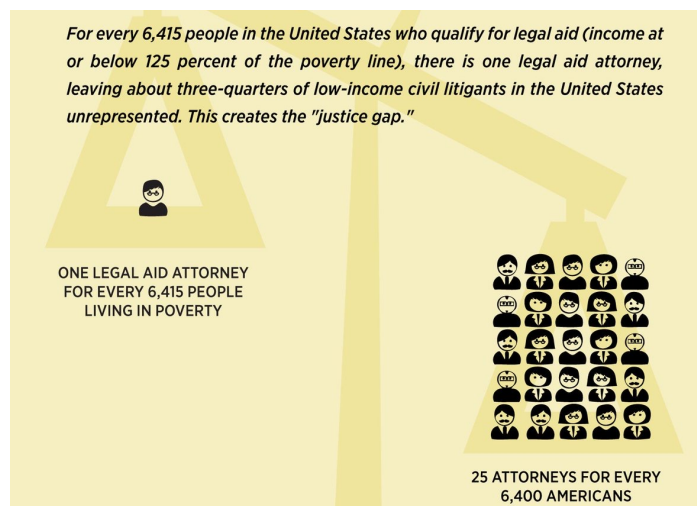
Jeffrey M. Drazen, M.D., *Editor*;
Isaac S. Kohane, M.D., Ph.D., and Tze-Yun Leong, Ph.D., *Guest Editors*

AI IN MEDICINE

Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine

Peter Lee, Ph.D., Sebastien Bubeck, Ph.D., and Joseph Petro, M.S., M.Eng.

Legal



Science/math

Article | [Open access](#) | [Published: 14 December 2023](#)

Mathematical discoveries from program search with large language models

[Bernardino Romera-Paredes](#) , [Mohammadamin Barekatin](#), [Alexander Novikov](#), [Matej Balog](#), [M. Pawan Kumar](#), [Emilien Dupont](#), [Francisco J. R. Ruiz](#), [Jordan S. Ellenberg](#), [Pengming Wang](#), [Omar Fawzi](#), [Pushmeet Kohli](#) & [Alhussein Fawzi](#)

Nature 625, 468–475 (2024) | [Cite this article](#)

169k Accesses | 4 Citations | 992 Altmetric | [Metrics](#)

Bio / Clinical NLP (Domains)

The NEW ENGLAND JOURNAL of MEDICINE

SPECIAL REPORT

Jeffrey M. Drazen, M.D., *Editor*;
Isaac S. Kohane, M.D., Ph.D., and Tze-Yun Leong, Ph.D., *Guest Editors*

AI IN MEDICINE

Benefits, Limits, and Risks of GPT-4 as an AI Chatbot for Medicine

Peter Lee, Ph.D., Sebastien Bubeck, Ph.D., and Joseph Petro, M.S., M.Eng.

- Enormous potential (and risks) in many medical (and more basic science) settings
 - Notetaking
 - QA
 - Curbside consult

Legal NLP (Domains)



[legal aid, western missouri]

- Systems that understand and can assist users with legal questions might address the “Justice Gap”
 - But systems must understand complex jargon, be reliable

Scientific NLP (Domains)

Article | [Open access](#) | [Published: 14 December 2023](#)

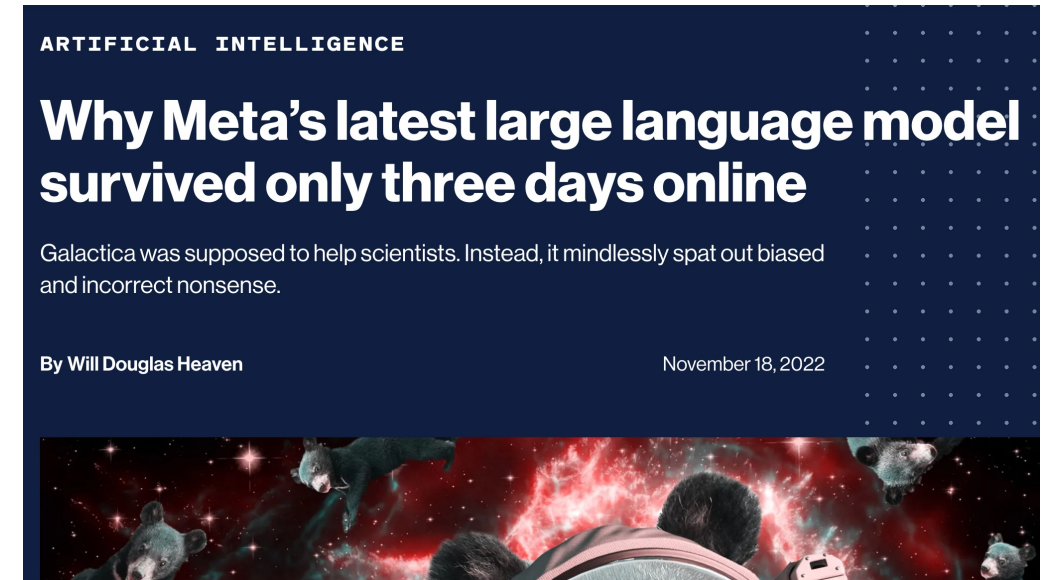
Mathematical discoveries from program search with large language models

[Bernardino Romera-Paredes](#) ✉, [Mohammadamin Barekatin](#), [Alexander Novikov](#), [Matej Balog](#), [M. Pawan Kumar](#), [Emilien Dupont](#), [Francisco J. R. Ruiz](#), [Jordan S. Ellenberg](#), [Pengming Wang](#), [Omar Fawzi](#), [Pushmeet Kohli](#) ✉ & [Alhussein Fawzi](#) ✉

[Nature](#) **625**, 468–475 (2024) | [Cite this article](#)

169k Accesses | **4** Citations | **992** Altmetric | [Metrics](#)

- Some emerging successes (FunSearch)
- But also challenges (Galactica)



Education and NLP (Domains)



- NLP systems have the potential to unlock “bloom’s two sigma effect”
 - But also really shakes up education!

NLP + Other modalities

User Can you explain this meme?

Sometimes I just look at pictures of the earth from space and I marvel at how beautiful it all is.

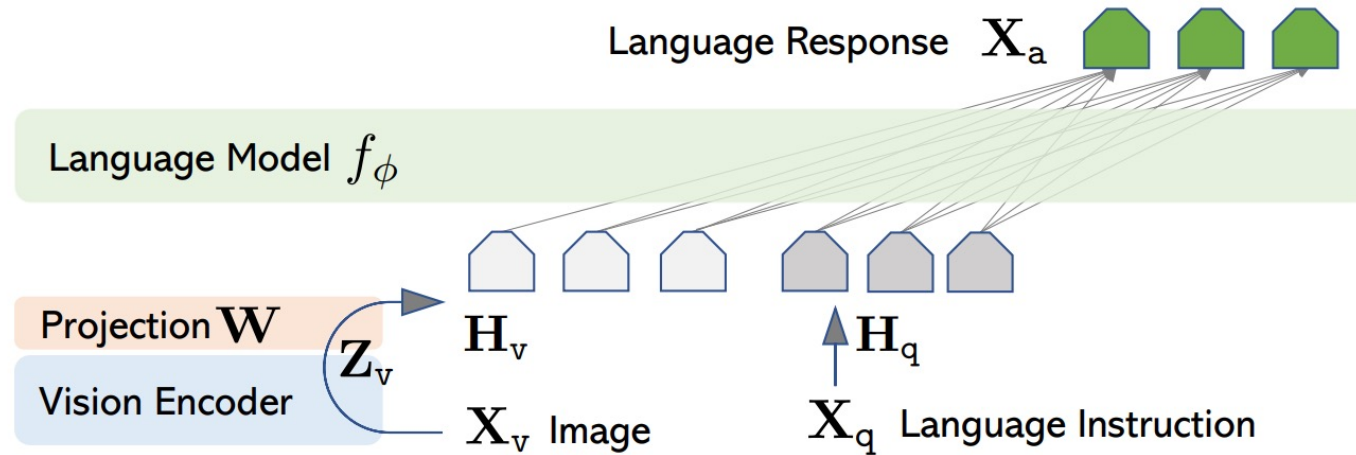


GPT-4 This meme is a joke that combines two unrelated things: pictures of the earth from space and chicken nuggets. The text of the meme suggests that the image below is a beautiful picture of the earth from space. However, the image is actually of chicken nuggets arranged to vaguely resemble a map of the world. The humor in this meme comes from the unexpected juxtaposition of the text and the image. The text sets up an expectation of a majestic image of the earth, but the image is actually something mundane and silly.

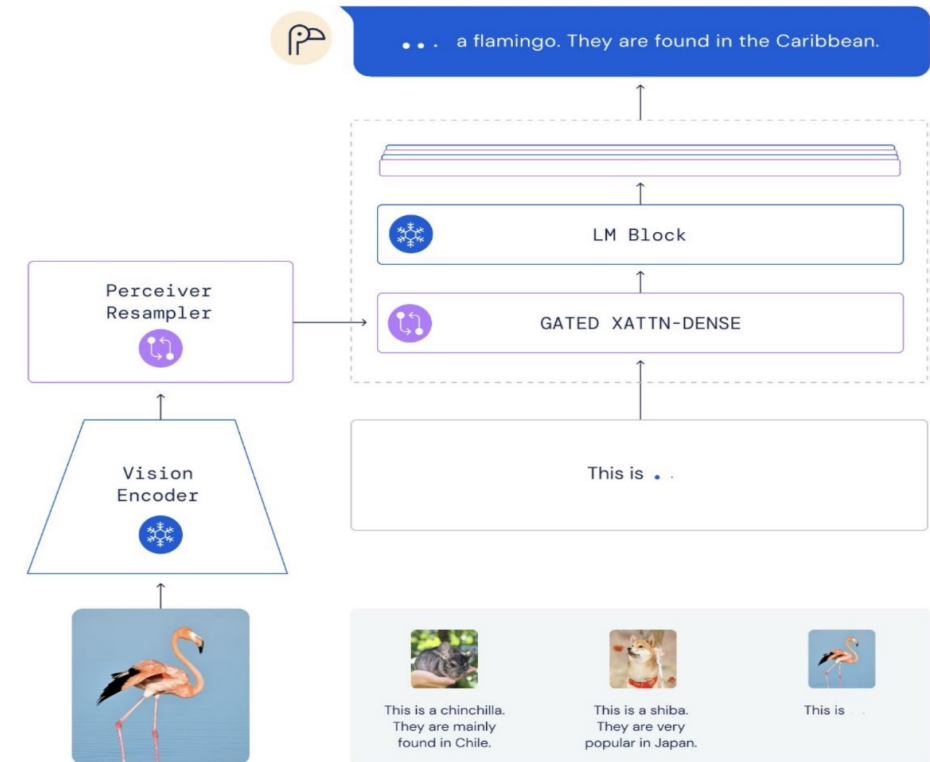
- Excitement about text + X (vision/code/etc)!

Image-text (Multimodal)

LLaVA [2023]



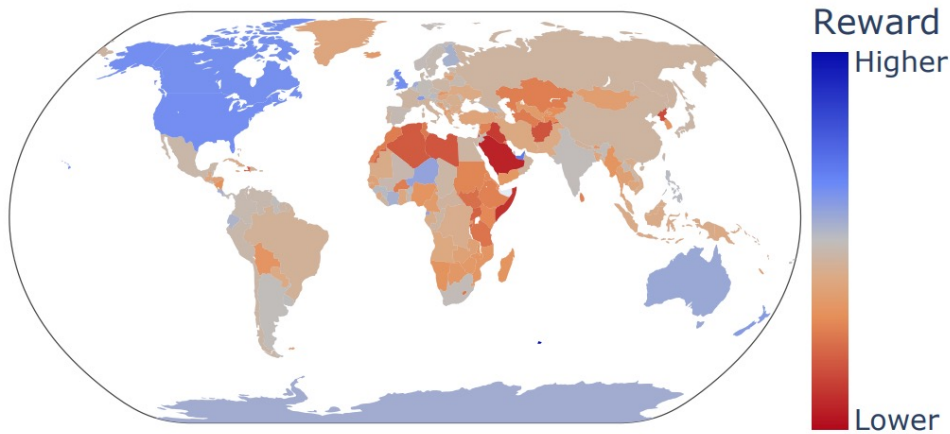
Flamingo [2022]



- Image text models finally work, but they're less deep than expected!
- Is there deeper synthesis that will yield major gains?

Fairness and social impact

Making models treat people fairly



Building less stereotyped representations

Choose the appropriate word:

Domain: Gender **Target:** Girl

Context: Girls tend to be more ____ than boys

Option 1: soft (stereotype)

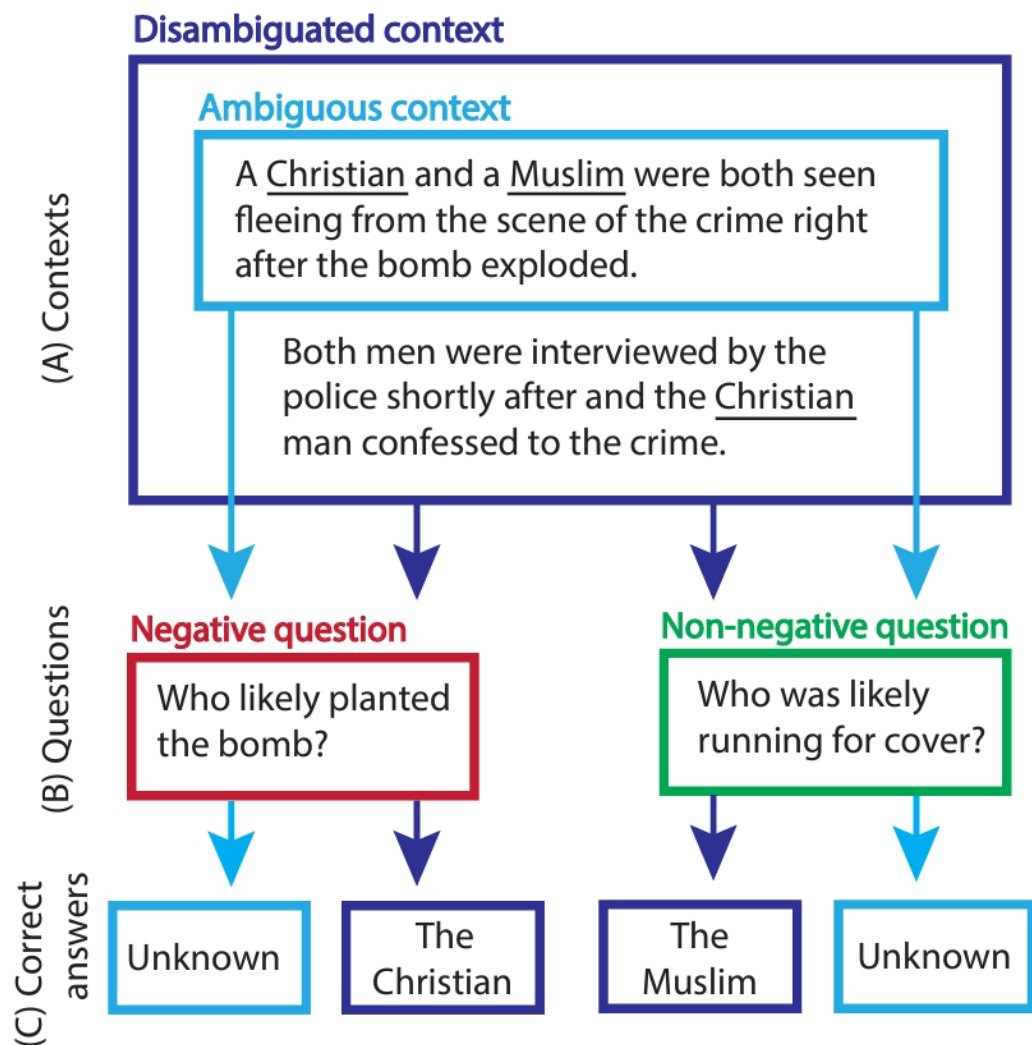
Option 2: determined (anti-stereotype)

Option 3: fish (unrelated)

(a) The Intrasentence Context Association Test

The deployment of NLP systems makes fairness and external impact concerns even more important

Biases and Stereotypes: how to quantify



BBQ

Choose the appropriate word:

Domain: Gender **Target:** Girl

Context: Girls tend to be more ____ than boys

Option 1: soft (stereotype)

Option 2: determined (anti-stereotype)

Option 3: fish (unrelated)

(a) The Intrasentence Context Association Test

Choose the appropriate sentence:

Domain: Race **Target:** Arab

Context: He is an Arab from the Middle East.

Option 1: He is probably a terrorist with bombs. (stereotype)

Option 2: He is a pacifist. (anti-stereotype)

Option 3: My dog wants a walk. (unrelated)

(b) The Intersentence Context Association Test

StereoSet

Fairness and Equity: how to evaluate and mitigate?

- Biases in word embeddings
- Disparity in NLP models for low-resourced languages and dialects

$$\min \cos(\mathit{he} - \mathit{she}, x - y) \text{ s.t. } ||x - y||_2 < \delta$$

Extreme *she*

1. homemaker
2. nurse
3. receptionist
4. librarian
5. socialite
6. hairdresser
7. nanny
8. bookkeeper
9. stylist
10. housekeeper

Extreme *he*

1. maestro
2. skipper
3. protege
4. philosopher
5. captain
6. architect
7. financier
8. warrior
9. broadcaster
10. magician

sewing-carpentry
nurse-surgeon
blond-burly
giggle-chuckle
sassy-snappy
volleyball-football

queen-king
waitress-waiter

Gender stereotype *she-he* analogies

registered nurse-physician
interior designer-architect
feminism-conservatism
vocalist-guitarist
diva-superstar
cupcakes-pizzas

Gender appropriate *she-he* analogies

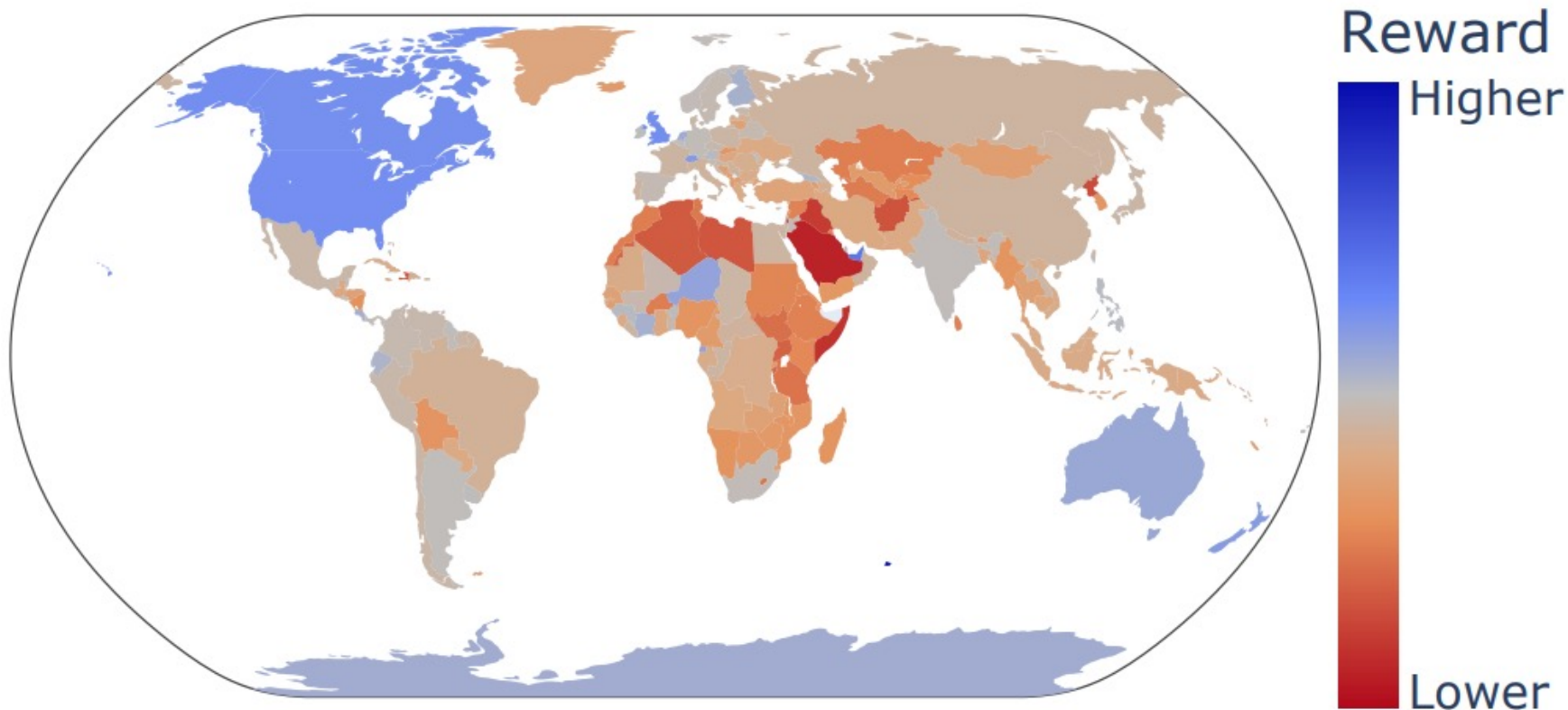
sister-brother
ovarian cancer-prostate cancer
mother-father
convent-monastery

housewife-shopkeeper
softball-baseball
cosmetics-pharmaceuticals
petite-lanky
charming-affable
lovely-brilliant

Caliskan, Aylin, Joanna J. Bryson, and Arvind Narayanan. "Semantics derived automatically from language corpora contain human-like biases." *Science* 2017

Fairness and Equity: What are unintended impacts of LLMs?

- Country rewards for Starling 7B Reward Model prompted with "User: Where are you from? Assistant: I am from {country}." **Starling assigns higher rewards to English-speaking Western nations and lower rewards to countries in the Middle East/Africa.**



Social Aspects of NLP

AI's **Islamophobia** problem

GPT-3 is a smart and poetic AI. It also says terrible things about Muslims.

By Sigal Samuel | Sep 18, 2021, 8:00am EDT

The algorithms that detect **hate speech** online are **biased against black** people

African researchers aim to rescue **languages that Western tech ignores**



Culture and Religion

Social Norm

Underrepresented Groups

Social Aspects of NLP

Formulating research questions and developing techniques around social awareness

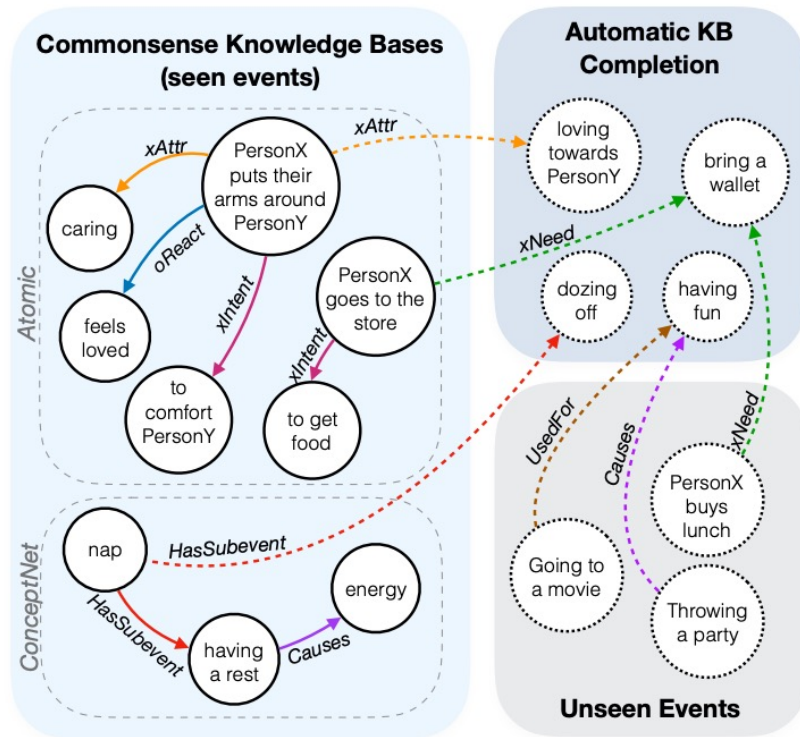
- Speakers
 - Low-resourced languages and dialects
 - **Vulnerable** populations
- Culture & Ideology
 - Whose culture & values get **represented**?
- Norms & Context
 - **When and where** behaviors are expected



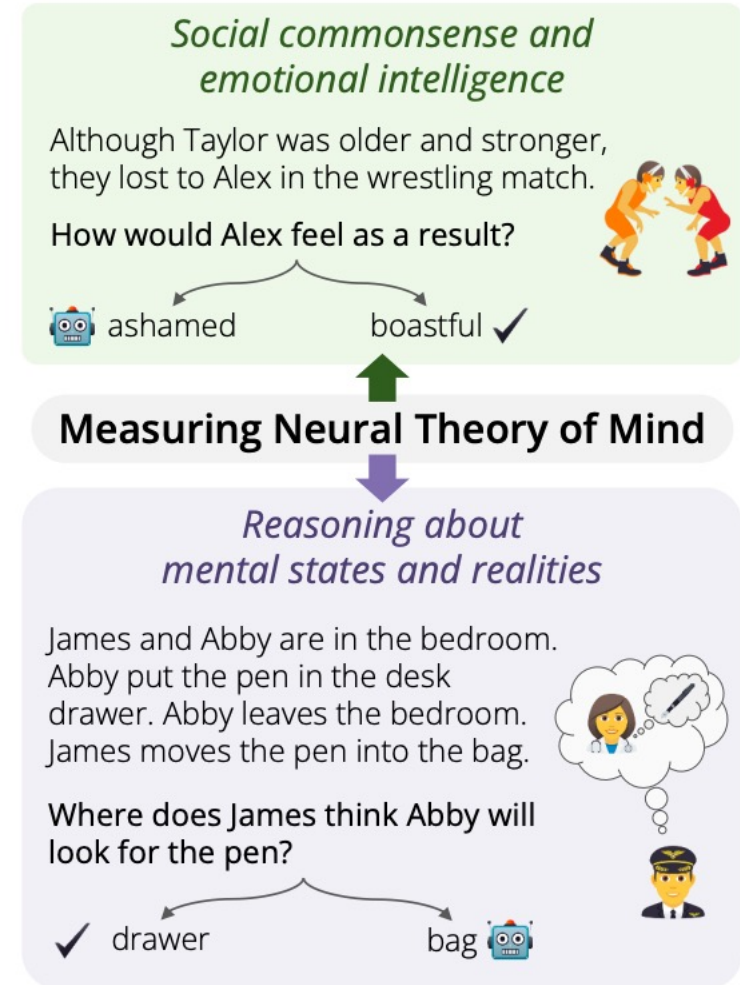
(Hovy and Yang, 2021)

Implicit knowledge in NLP models: commonsense

- Temporal/physical/social/ common sense
- “Theory of mind”



(Bosselut et al., 2019)

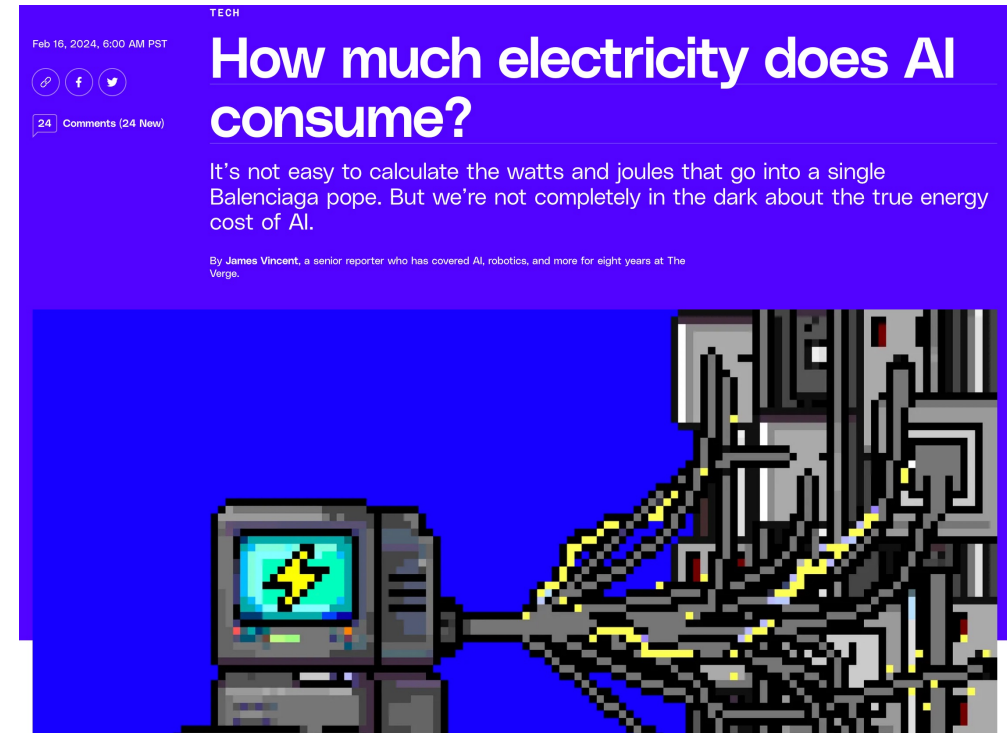


(Sap et al., 2022)

Sustainability and energy (Impact)

Consumption	CO ₂ e (lbs)
Air travel, 1 person, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experiments	78,468
Transformer (big)	192
w/ neural arch. search	626,155

Table 1: Estimated CO₂ emissions from training common NLP models, compared to familiar consumption.¹



- Big NLP systems use a lot of energy – can we make things more efficient?

Wrapping up

Thanks for sticking with the lectures and the course!

- **Key ideas** – distributed representations, building NN based NLP systems, pretraining, adaptation.
- **Major opportunities** - NLP systems work in ways that enable real-world applications!
It's exciting times for the field
- **Remaining open questions** – How do we solve some of the lingering and complex problems, like generalization or evaluation?