

Interpretability and Analysis

Ling 282/482: Deep Learning for Computational Linguistics

C.M. Downey

Fall 2024

Today's Plan

- NLP's "Clever Hans" Moment: motivating interpretability and analysis
- Survey of several different methods:
 - Neuron-level
 - Psycholinguistic experiments
 - Diagnostic classifiers
 - Attention analysis
 - Adversarial datasets

NLP's “Clever Hans Moment”



ME EDITOR'S NOTE OVERVIEWS PERSPECTIVES ABOUT SUBSCRIBE Q

Clever Hans

BERT



NLP's Clever Hans Moment has Arrived

26.AUG.2019

[link](#)

Clever Hans

- Early 1900s, a horse trained by his owner to do:
 - Addition
 - Division
 - Multiplication
 - Tell time
 - Read German
 - ...
- Wow! Hans is really smart!

Clever Hans Effect

Clever Hans Effect

- Upon closer examination / experimentation...

Clever Hans Effect

- Upon closer examination / experimentation...
- Hans' success:

Clever Hans Effect

- Upon closer examination / experimentation...
- Hans' success:
 - 89% when **questioner knows answer**

Clever Hans Effect

- Upon closer examination / experimentation...
- Hans' success:
 - 89% when **questioner knows answer**
 - 6% when questioner doesn't know answer

Clever Hans Effect

- Upon closer examination / experimentation...
- Hans' success:
 - 89% when **questioner knows answer**
 - 6% when questioner doesn't know answer
- Further experiments: as Hans' taps got closer to correct answer, **facial tension in questioner** increased

Clever Hans Effect

- Upon closer examination / experimentation...
- Hans' success:
 - 89% when **questioner knows answer**
 - 6% when questioner doesn't know answer
- Further experiments: as Hans' taps got closer to correct answer, **facial tension in questioner** increased
- Hans didn't solve the task but **exploited a spuriously correlated cue**

Central question

- Do BERT et al's major successes at solving NLP tasks show that we have achieved robust natural language understanding in machines?
- Or: are we seeing a “Clever BERT” phenomenon?

Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference

R. Thomas McCoy,¹ Ellie Pavlick,² & Tal Linzen¹

¹Department of Cognitive Science, Johns Hopkins University

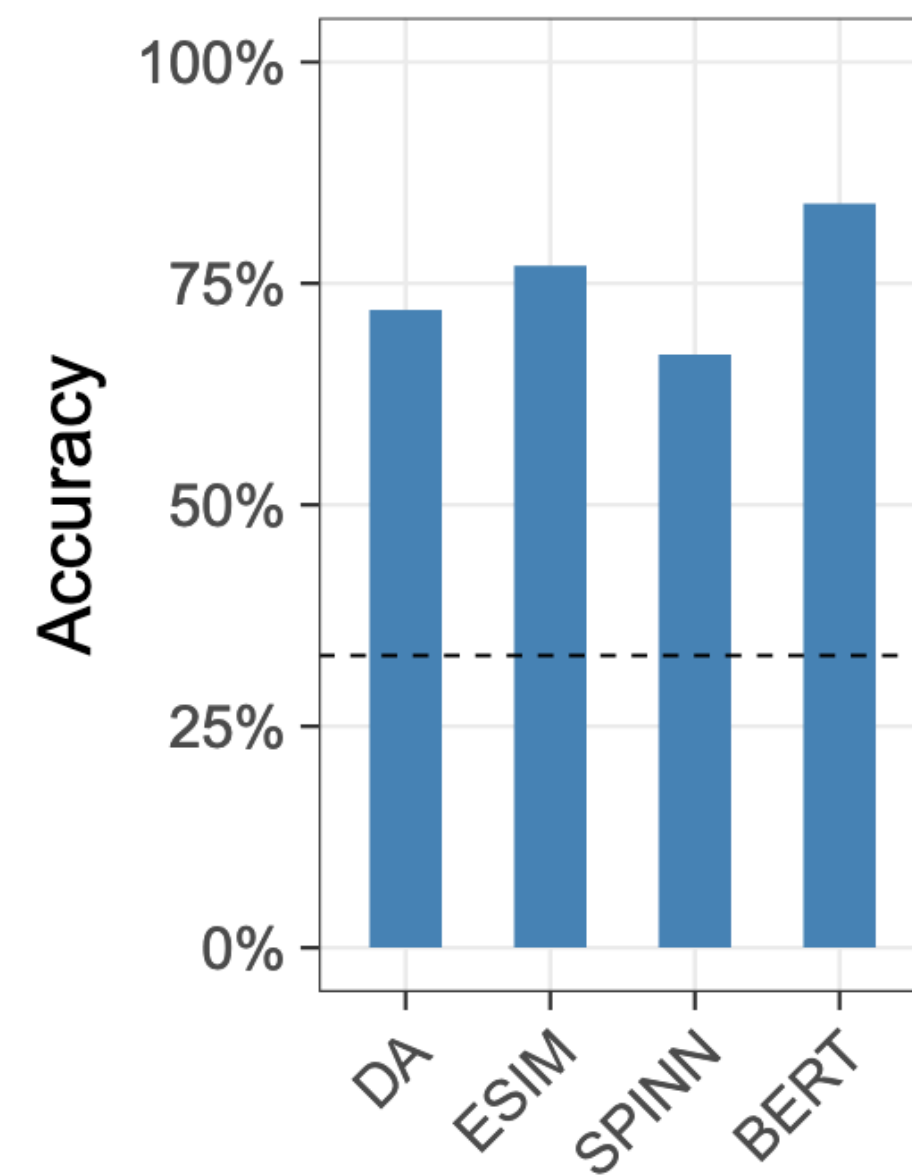
²Department of Computer Science, Brown University

tom.mccoy@jhu.edu, ellie_pavlick@brown.edu, tal.linzen@jhu.edu

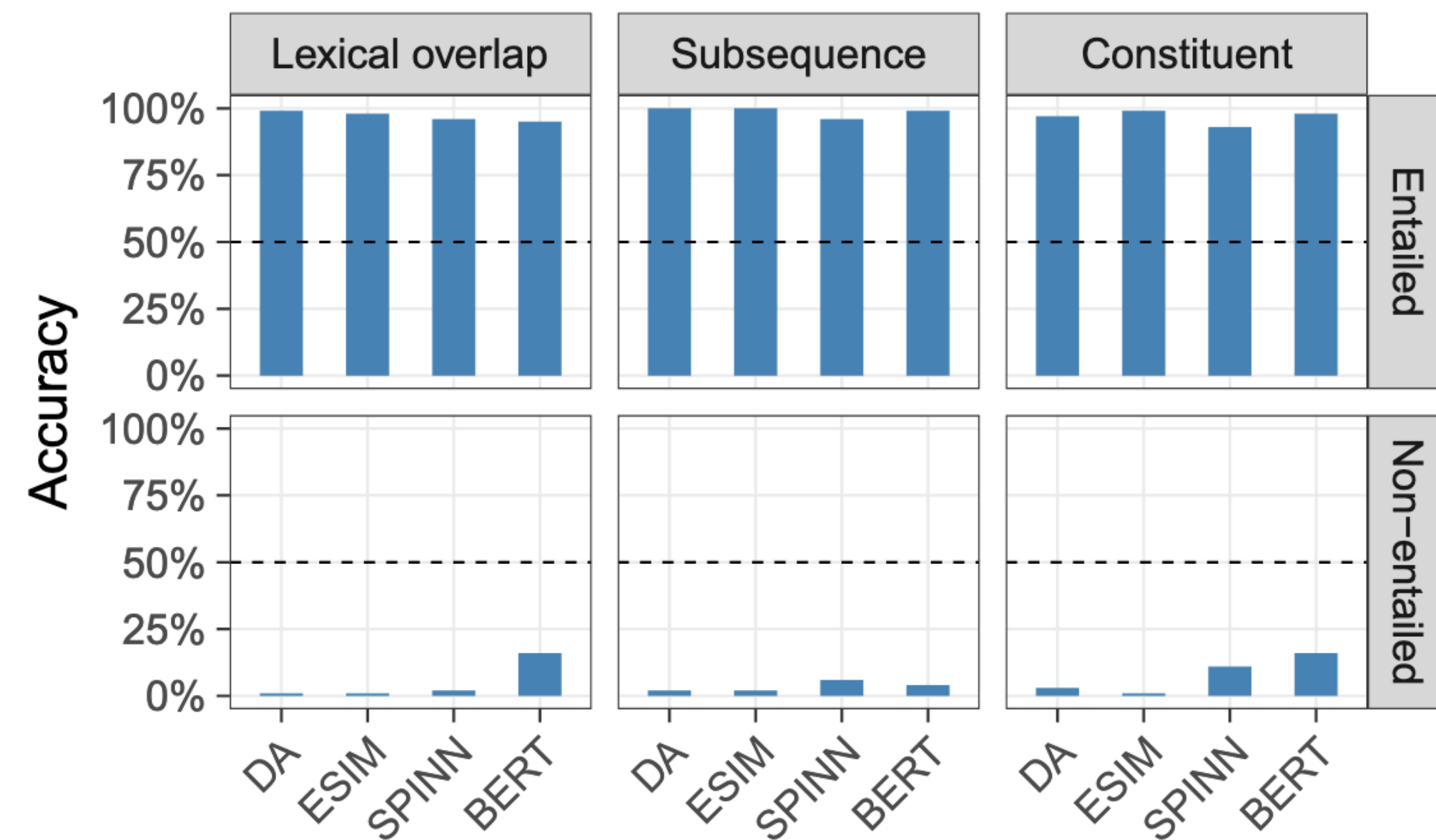
[McCoy et al 2019](#)

Heuristic	Premise	Hypothesis	Label
Lexical overlap heuristic	The banker near the judge saw the actor.	The banker saw the actor.	E
	The lawyer was advised by the actor.	The actor advised the lawyer.	E
	The doctors visited the lawyer.	The lawyer visited the doctors.	N
	The judge by the actor stopped the banker.	The banker stopped the actor.	N
Subsequence heuristic	The artist and the student called the judge.	The student called the judge.	E
	Angry tourists helped the lawyer.	Tourists helped the lawyer.	E
	The judges heard the actors resigned.	The judges heard the actors.	N
	The senator near the lawyer danced.	The lawyer danced.	N
Constituent heuristic	Before the actor slept, the senator ran.	The actor slept.	E
	The lawyer knew that the judges shouted.	The judges shouted.	E
	If the actor slept, the judge saw the artist.	The actor slept.	N
	The lawyers resigned, or the artist slept.	The artist slept.	N

Results



(a)



(b)

(performance improves if fine-tuned on this challenge set)

Why care?

- Engineering: can help build **better language technologies** via improved models, data, training protocols, ...
 - **Trust**, critical applications
- Theoretical: can help us **understand biases** in different architectures (e.g. LSTMs vs Transformers), similarities to **human learning biases**
 - Which linguistic features / properties are *learnable* from raw text alone?
- Ethical: e.g. do some models reflect **harmful social biases** more than others?

Visualization / neuron-level analysis

Main Idea

- Individual neurons in a network have activations that depend on the input
- Check to see whether any of them have **activations that correlate** with (linguistically) interesting features of the input
- Historical discourse of alleged “Jennifer Anniston cells”, aka grandmother cells

Learning to Generate Reviews and Discovering Sentiment

Alec Radford¹ Rafal Jozefowicz¹ Ilya Sutskever¹

Abstract

We explore the properties of byte-level recurrent language models. When given sufficient amounts of capacity, training data, and compute time, the representations learned by these models include disentangled features corresponding to high-level concepts. Specifically, we find a single unit which performs sentiment analysis. These representations, learned in an unsupervised manner, achieve state of the art on the binary subset of the Stanford Sentiment Treebank. They are also very data efficient. When using only a handful of labeled examples, our approach matches the performance of strong baselines trained on full datasets. We also demonstrate the sentiment unit has a direct influence on the generative process of the model. Simply fixing its value to be positive or negative generates samples with the corresponding positive or negative sentiment.

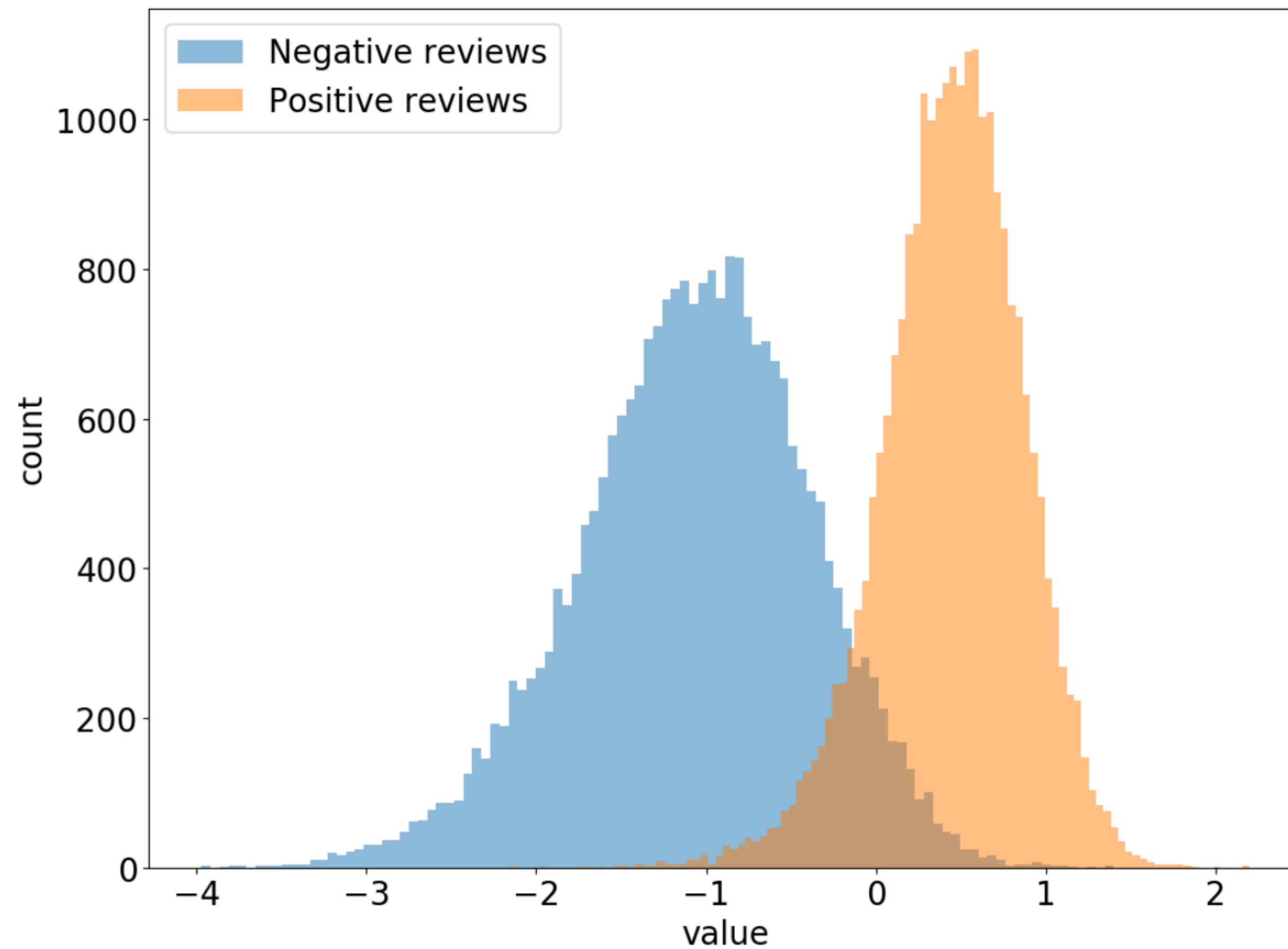
it is now commonplace to reuse these representations on a broad suite of related tasks - one of the most successful examples of transfer learning to date (Oquab et al., 2014).

There is also a long history of unsupervised representation learning (Olshausen & Field, 1997). Much of the early research into modern deep learning was developed and validated via this approach (Hinton & Salakhutdinov, 2006) (Huang et al., 2007) (Vincent et al., 2008) (Coates et al., 2010) (Le, 2013). Unsupervised learning is promising due to its ability to scale beyond only the subsets and domains of data that can be cleaned and labeled given resource, privacy, or other constraints. This advantage is also its difficulty. While supervised approaches have clear objectives that can be directly optimized, unsupervised approaches rely on proxy tasks such as reconstruction, density estimation, or generation, which do not directly encourage useful representations for specific tasks. As a result, much work has gone into designing objectives, priors, and architectures meant to encourage the learning of useful representations.

Approach

- Character-level language model (LSTM variant)
 - One layer; 4096 dim hidden state
 - Training: ~1 month on 4 GPUs
- Data: Amazon product reviews
- Fine-tune: sentiment analysis
 - This data partially overlaps with training data (but a different task)

A sentiment neuron



Samples of the sentiment neuron

I found this to be a charming adaptation, very lively and full of fun. With the exception of a couple of major errors, the cast is wonderful. I have to echo some of the earlier comments -- Chynna Phillips is horribly miscast as a teenager. At 27, she's just too old (and, yes, it DOES show), and lacks the singing "chops" for Broadway-style music. Vanessa Williams is a decent-enough singer and, for a non-dancer, she's adequate. However, she is NOT Latina, and her character definitely is. She's also very STRIDENT throughout, which gets tiresome. The girls of Sweet Apple's Conrad Birdie fan club really sparkle -- with special kudos to Brigitta Dau and Chiara Zanni. I also enjoyed Tyne Daly's performance, though I'm not generally a fan of her work. Finally, the dancing Shriners are a riot, especially the dorky three in the bar. The movie is suitable for the whole family, and I highly recommend it.

Judy Holliday struck gold in 1950 with the George Cukor's film version of "Born Yesterday," and from that point forward, her career consisted of trying to find material good enough to allow her to strike gold again. It never happened. In "It Should Happen to You" (I can't think of a blander title, by the way), Holliday does yet one more variation on the dumb blonde who's maybe not so dumb after all, but everything about this movie feels warmed over and half hearted. Even Jack Lemmon, in what I believe was his first film role, can't muster up enough energy to enliven this recycled comedy. The audience knows how the movie will end virtually from the beginning, so mostly it just sits around waiting for the film to catch up. Maybe if you're enamored of Holliday you'll enjoy this; otherwise I wouldn't bother. Grade: C

Sentiment unit does all the work!

Table 2. IMDB sentiment classification

METHOD	ERROR
FULLUNLABELED BOW (MAAS ET AL., 2011)	11.11%
NB-SVM TRIGRAM (MESNIL ET AL., 2014)	8.13%
SENTIMENT UNIT (OURS)	7.70%
SA-LSTM (DAI & LE, 2015)	7.24%
BYTE MLSTM (OURS)	7.12%
TOPICRNN (DIENG ET AL., 2016)	6.24%
VIRTUAL ADV (MIYATO ET AL., 2016)	5.91%

The Emergence of Number and Syntax Units in LSTM Language Models

Yair Lakretz

Cognitive Neuroimaging Unit
NeuroSpin center
91191, Gif-sur-Yvette, France
yair.lakretz@gmail.com

German Kruszewski

Facebook AI Research
Paris, France
germank@gmail.com

Theo Desbordes

Facebook AI Research
Paris, France
tdesbordes@fb.com

Dieuwke Hupkes

ILLC, University of Amsterdam
Amsterdam, Netherlands
d.hupkes@uva.nl

Stanislas Dehaene

Cognitive Neuroimaging Unit
NeuroSpin center
91191, Gif-sur-Yvette, France
stanislas.dehaene@gmail.com

Marco Baroni

Facebook AI Research
Paris, France
mbaroni@fb.com

Approach

- Evaluating an LSTM LM
- **Number agreement tasks:** as in Linzen et al 2016 (to be discussed shortly!)

- Plus synthetic:

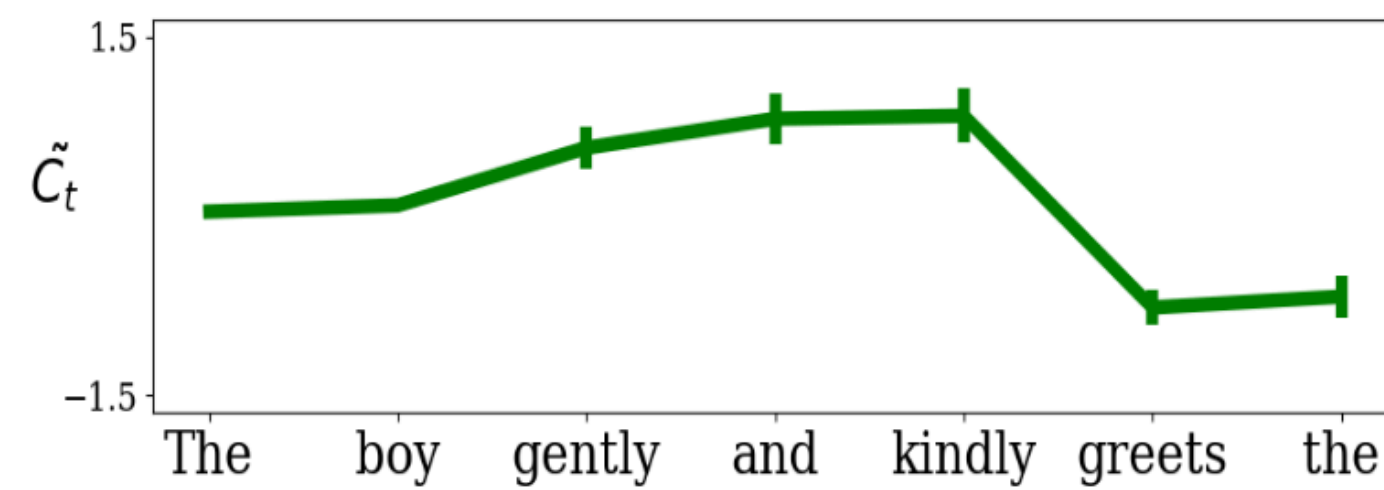
Simple	the boy greet s the guy
Adv	the boy probably greet s the guy
2Adv	the boy most probably greet s the guy
CoAdv	the boy openly and deliberately greet s the guy
NamePP	the boy near Pat greet s the guy
NounPP	the boy near the car greet s the guy
NounPPAdv	the boy near the car kindly greet s the guy

- Find important cells by **ablation**: set activation to 0, **see if performance suffers**.
(Also by regression; more in a minute)

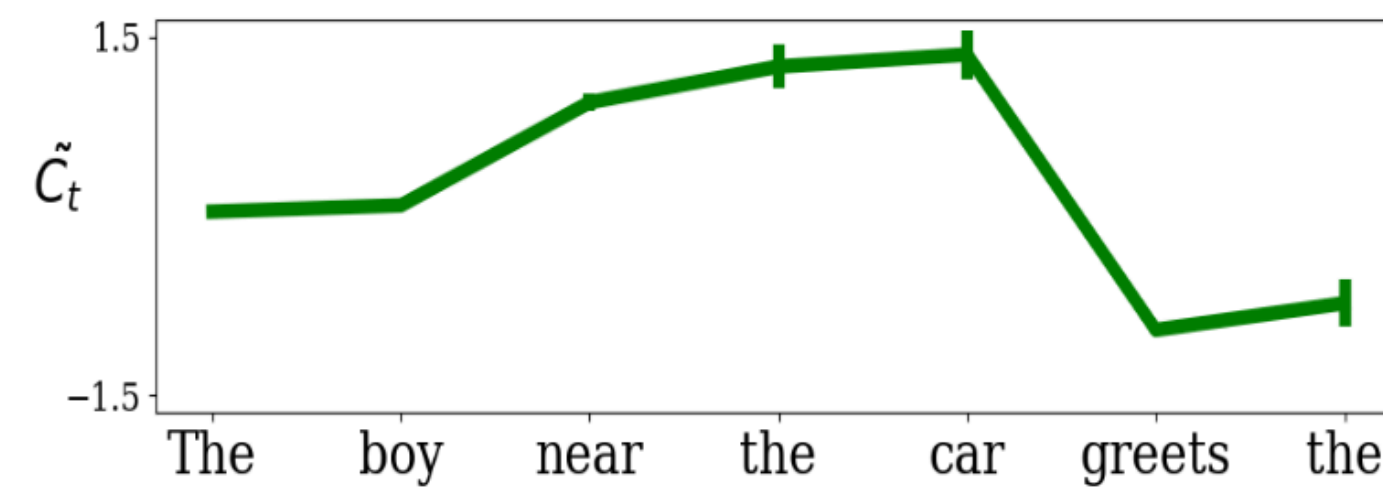
Finding a syntax unit

- Predict, via linear regression, from the cell:
 - Depth of the word in syntactic parse of the sentence
 - (Works pretty well: $R^2 = 0.85$. More on this idea later.)
- Identify cells that are assigned very high weight in the regression

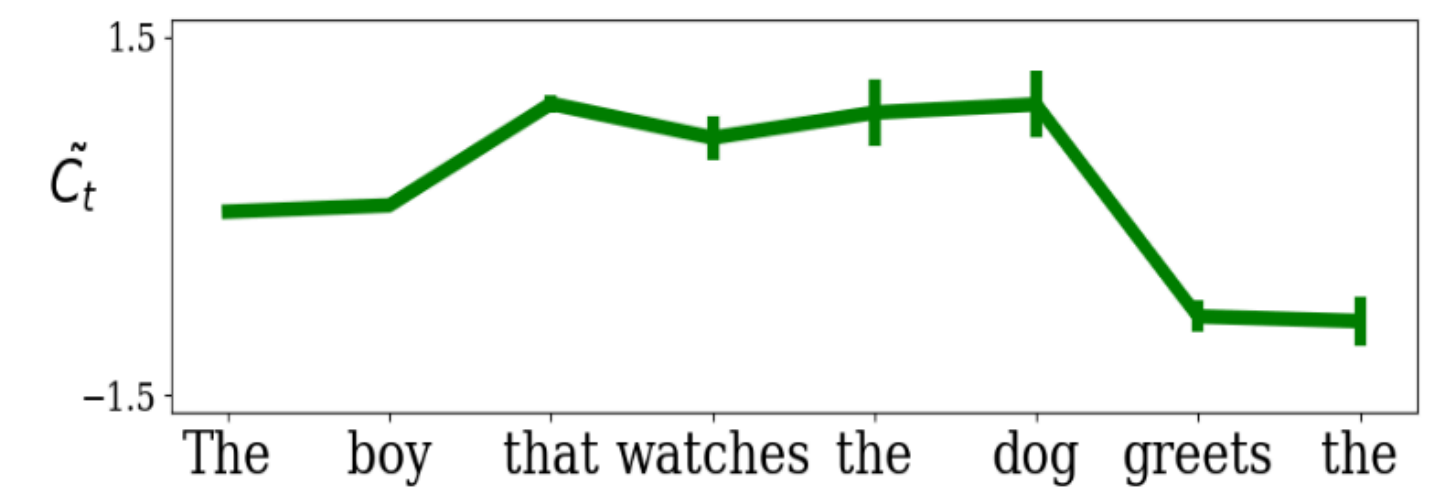
Cell dynamics for a syntax unit



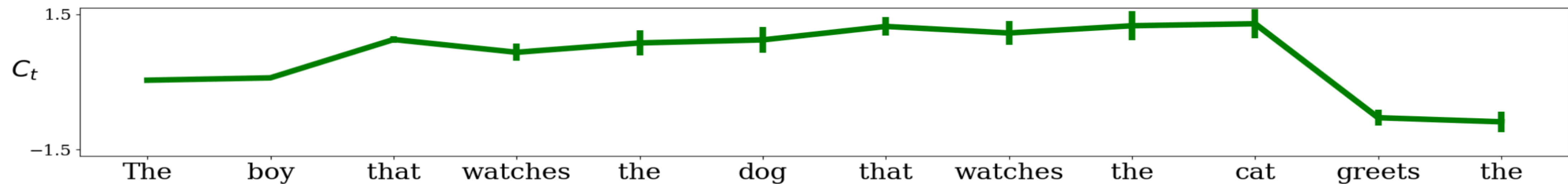
(a) 2Adv



(b) nounPP



(c) subject relative



Neuron-level analysis: summary

Neuron-level analysis: summary

- Very promising and exciting when it does work: a good look “inside the black box”, with very interpretable neural/cell dynamics. BUT

Neuron-level analysis: summary

- Very promising and exciting when it does work: a good look “inside the black box”, with very interpretable neural/cell dynamics. BUT
- “A needle in a haystack”: how to **find the “good” neurons?**
 - Some principled methods (ablation, regression); **not all of them scale well**
 - Also important differences in **which questions can be tested**
 - Is there a neuron that tracks property P?
 - **Not: what are you tracking?**

Neuron-level analysis: summary

- Very promising and exciting when it does work: a good look “inside the black box”, with very interpretable neural/cell dynamics. BUT
- “A needle in a haystack”: how to **find the “good” neurons?**
 - Some principled methods (ablation, regression); **not all of them scale well**
 - Also important differences in **which questions can be tested**
 - Is there a neuron that tracks property P?
 - **Not: what are you tracking?**
- Deleting interpretable neurons may not effect performance in the original or downstream task ([Morcos et al 2018](#))

Psycholinguistic methods

Animating Idea

- NLMs are a bit of a “black box”. How can we figure out what they’re doing?
- **Humans are also black boxes!**
- Let’s test NLMs **the way we test people** when we try to figure out the nature of their linguistic knowledge
 - In other words: treat NLMs as if they were **participants in psycholinguistic experiments**

Assessing the Ability of LSTMs to Learn Syntax-Sensitive Dependencies

Tal Linzen^{1,2} **Emmanuel Dupoux**¹
LSCP¹ & IJN², CNRS,
EHESS and ENS, PSL Research University
{tal.linzen,
emmanuel.dupoux}@ens.fr

Yoav Goldberg
Computer Science Department
Bar Ilan University
yoav.goldberg@gmail.com

Abstract

The success of long short-term memory (LSTM) neural networks in language processing is typically attributed to their ability to capture long-distance statistical regularities. Linguistic regularities are often sensitive to syntactic structure; can such dependencies be captured by LSTMs, which do not have explicit structural representations? We begin addressing this question using number agreement in English subject-verb dependencies. We probe the architecture’s grammatical competence both using training objectives with an explicit grammatical target (number prediction, grammaticality judgments) and using language models. In the strongly supervised settings,

(Hochreiter and Schmidhuber, 1997) or gated recurrent units (GRU) (Cho et al., 2014), has led to significant gains in language modeling (Mikolov et al., 2010; Sundermeyer et al., 2012), parsing (Vinyals et al., 2015; Kiperwasser and Goldberg, 2016; Dyer et al., 2016), machine translation (Bahdanau et al., 2015) and other tasks.

The effectiveness of RNNs¹ is attributed to their ability to capture statistical contingencies that may span an arbitrary number of words. The word *France*, for example, is more likely to occur somewhere in a sentence that begins with *Paris* than in a sentence that begins with *Penguins*. The fact that an arbitrary number of words can intervene between the mutually predictive words implies that they cannot be captured

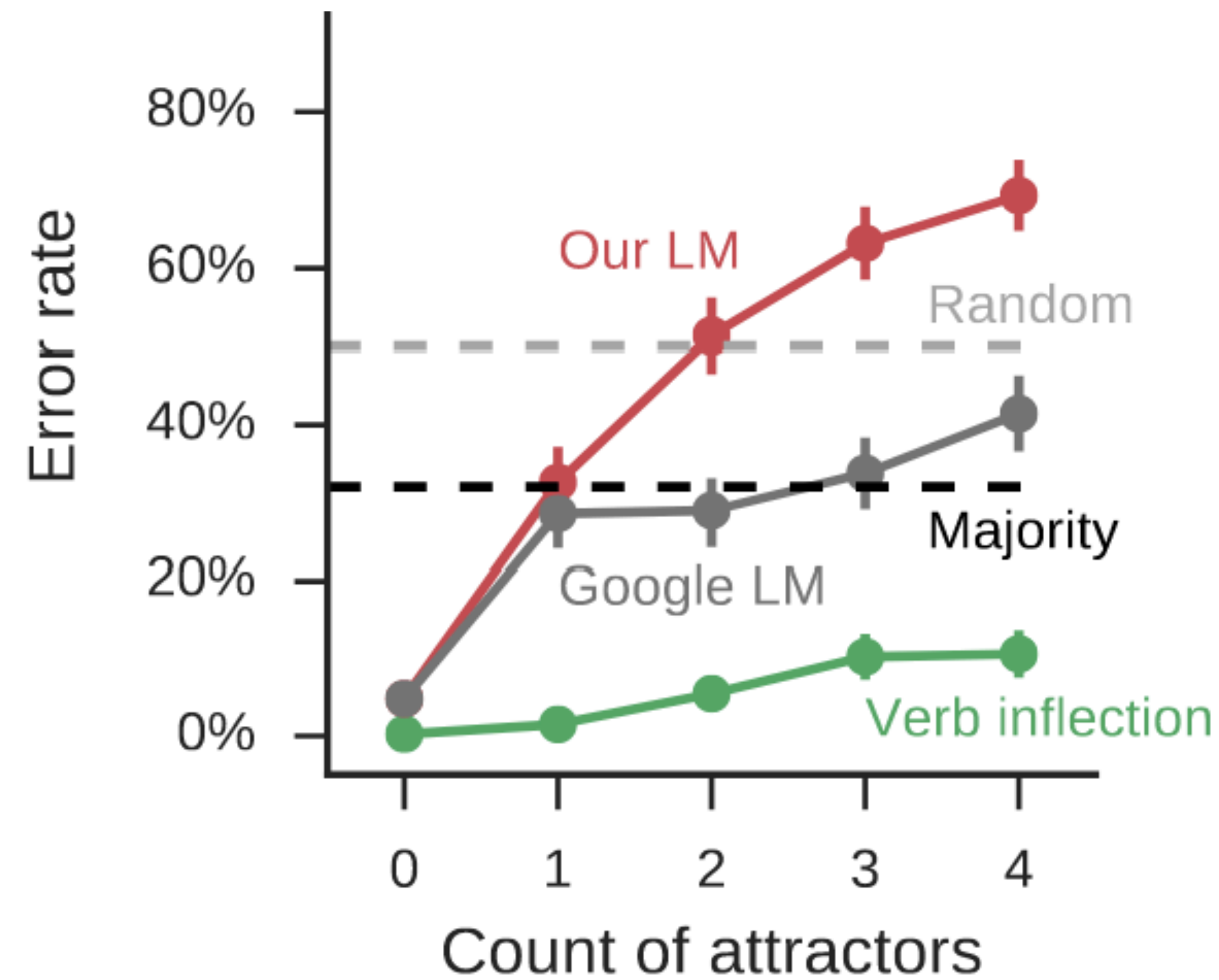
Subject-verb agreement

- Adjacent:
 - The key is on the table [SS]
 - * The key are on the table [SP]
 - * The keys is on the table [PS]
 - The keys are on the table [PP]
- Arbitrarily many **attractors** (nouns w/ different number) in between
 - But even the **city** with several tall buildings and many thriving industries **is** struggling.

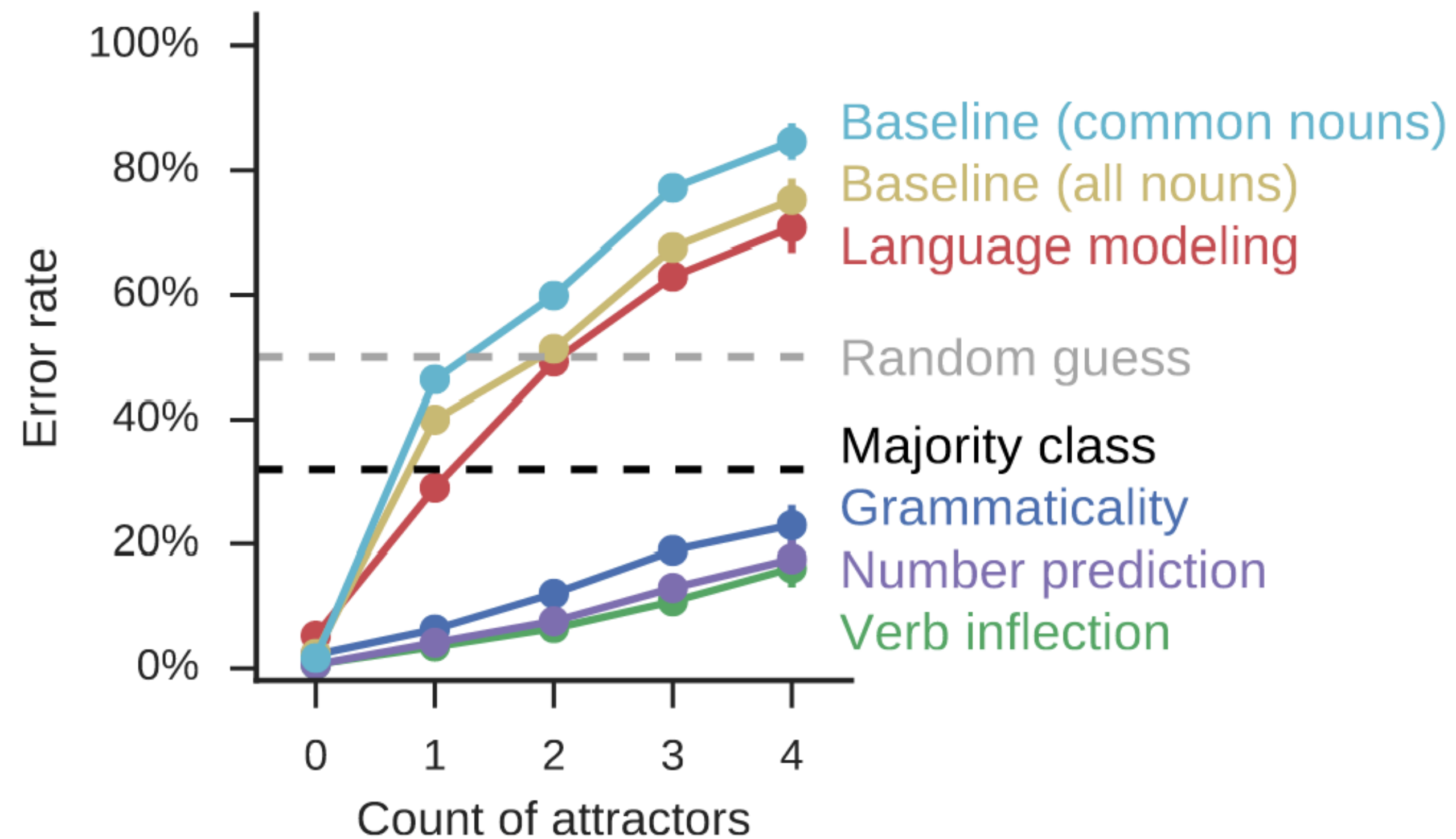
Method

- Does LM predict the right form of the verb?
 - “The keys on the cabinet ...”
 - $P_{LM}(\text{are}) > P_{LM}(\text{is})$?
- Single layer LSTM w/ 50 hidden units
- A lot more in the paper than we’ll talk about here
- Later: other methods for getting LM grammaticality judgments

Accuracy vs. Attractors



Effect of Task



Take Home

- LSTMs can in general learn hierarchical dependencies
- But language modeling **may not provide enough signal** on its own
 - i.e. explicit supervision on the task is required

Take Home

- Language modeling *may* after all provide enough of a signal to learn hierarchical syntactic dependencies
 - But may be very **sensitive to hyper-parameters**, including training data
 - “suggests that the input itself contains enough information to trigger some form of syntactic learning in a system, such as an RNN, that does not contain an explicit prior bias in favour of syntactic structures”
- Good model and data to play with (<https://github.com/facebookresearch/colorlessgreenRNNs>)
- A follow-up, with more constructions than just subject/verb agreement, and artificially generated data: <https://www.aclweb.org/anthology/D18-1151/>

Diagnostic classifiers

Main Idea

- What's in a representation (a vector)? How can we tell?
- For example: does an LSTM's memory encode grammatical number?
 - If we're lucky: a single cell might, as we saw earlier. (Sparse representation)
 - In general: *if we can easily predict the number from the memory*, it's “already in there”.
- Given a representation, train a simple model (usually a linear classifier) to predict a property of interest (usually linguistic) from that representation.

Note on Terminology

- Roughly synonyms: diagnostic classifiers, probing classifiers, auxiliary prediction tasks, ...
- (Basically: very simple transfer learning)

Journal of Artificial Intelligence Research 61 (2018) 907-926

Submitted 10/17; published 04/18

**Visualisation and ‘Diagnostic Classifiers’ Reveal
how Recurrent and Recursive Neural Networks
Process Hierarchical Structure**

Dieuwke Hupkes
Sara Veldhoen
Willem Zuidema
ILLC, University of Amsterdam
P.O.Box 94242
1090 CE Amsterdam, Netherlands

D.HUPKES@UVA.NL
SARA.VELDHOEN@GMAIL.COM
ZUIDEMA@UVA.NL

Linguistic Knowledge and Transferability of Contextual Representations

Nelson F. Liu^{♠♡*} Matt Gardner[♣] Yonatan Belinkov[◇]

Matthew E. Peters[♣] Noah A. Smith^{♠♣}

[♠]Paul G. Allen School of Computer Science & Engineering,
University of Washington, Seattle, WA, USA

[♡]Department of Linguistics, University of Washington, Seattle, WA, USA

[♣]Allen Institute for Artificial Intelligence, Seattle, WA, USA

[◇]Harvard John A. Paulson School of Engineering and Applied Sciences and
MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA

`{nfliu, nasmith}@cs.washington.edu`

`{mattg, matthewp}@allenai.org, belinkov@seas.harvard.edu`

Abstract

Contextual word representations derived from large-scale neural language models are successful across a diverse set of NLP tasks, suggesting that they encode useful and transferable features of language. To shed light on the linguistic knowledge they capture, we study the representations produced by several recent pretrained contextualizers (variants of ELMo, the OpenAI transformer language model, and BERT) with a suite of sixteen diverse probing tasks. We find that linear models trained on top of frozen contextual repre-

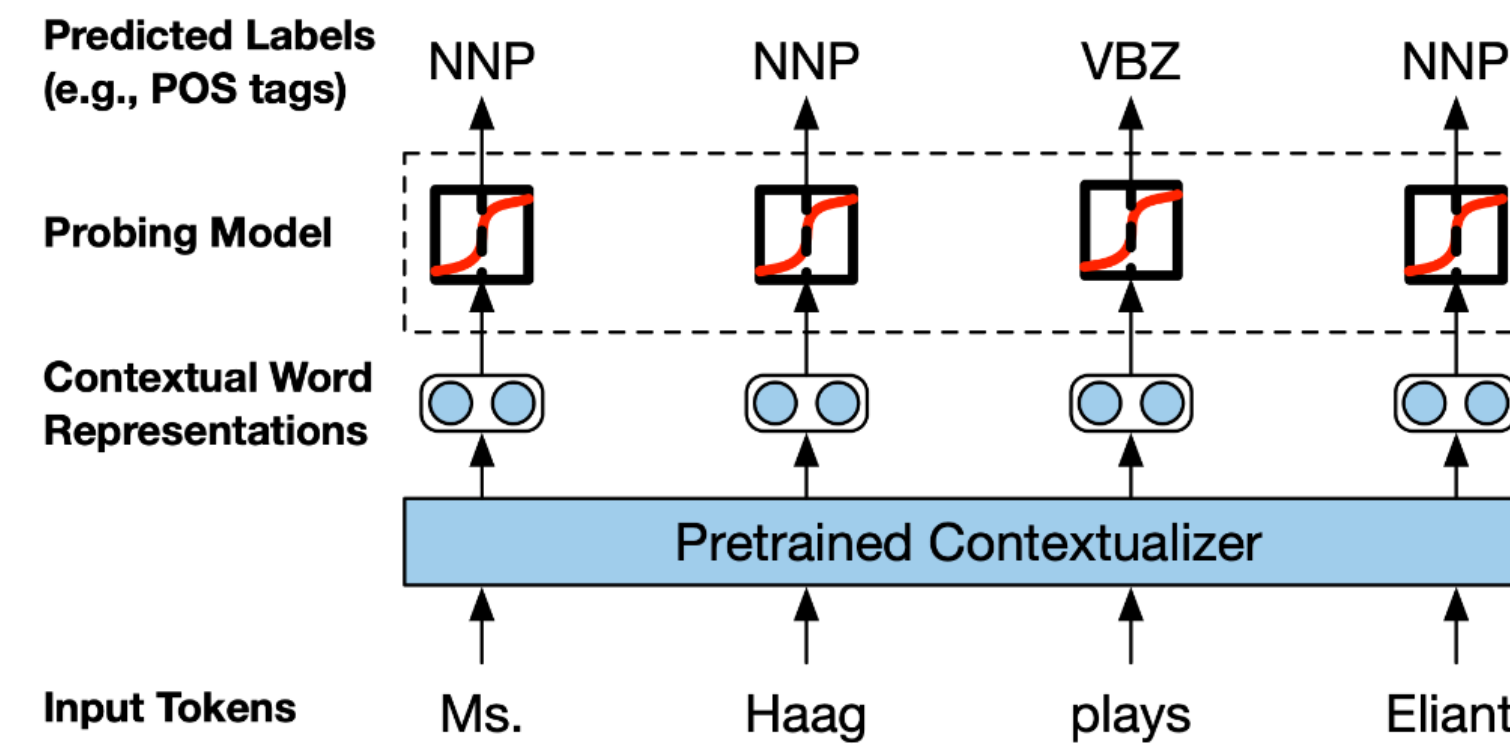


Figure 1: An illustration of the probing model setup used to study the linguistic knowledge within contextual word representations.

Tagging Results

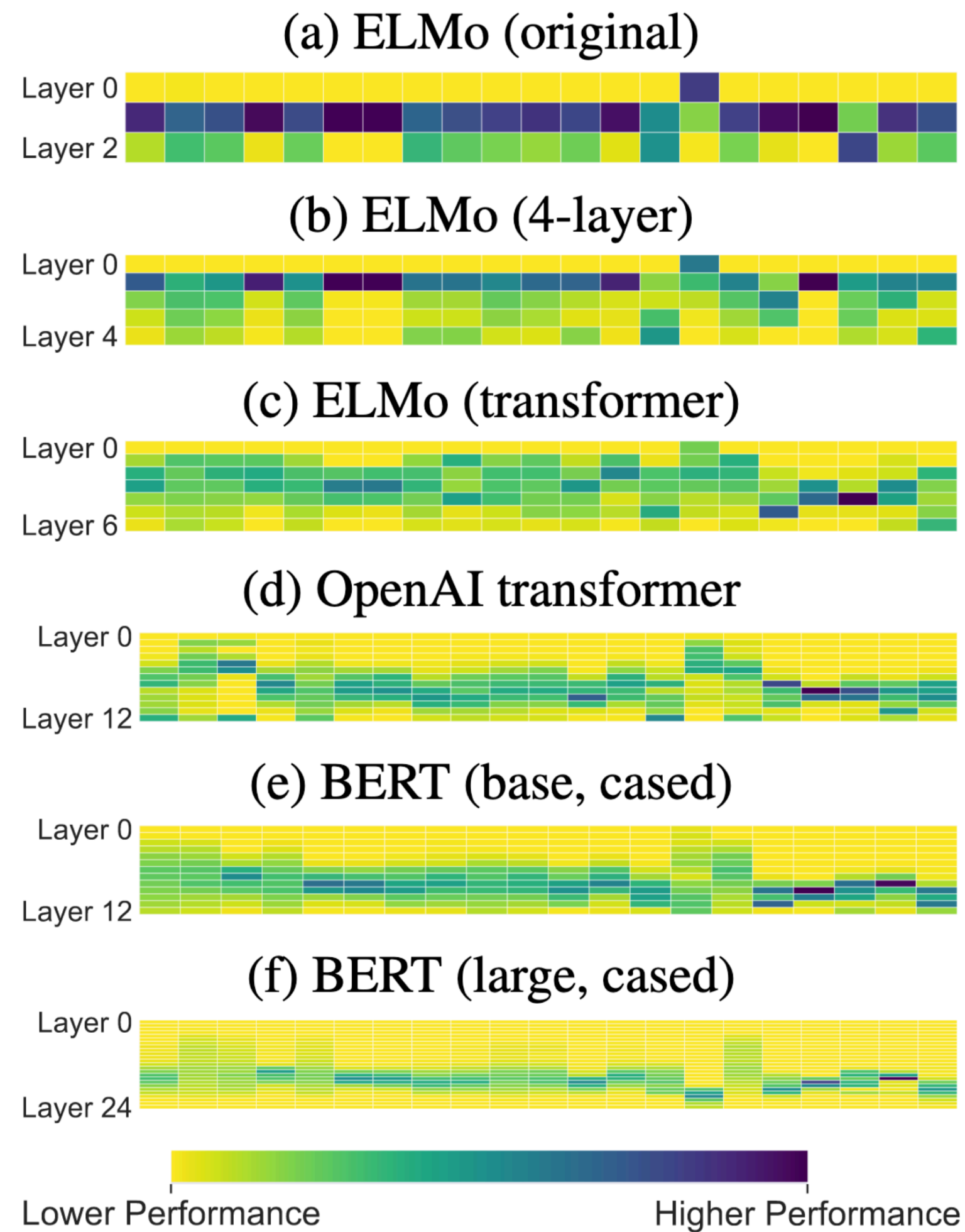
Pretrained Representation	POS			Supersense ID							
	Avg.	CCG	PTB	EWT	Chunk	NER	ST	GED	PS-Role	PS-Fxn	EF
ELMo (original) best layer	81.58	93.31	97.26	95.61	90.04	82.85	93.82	29.37	75.44	84.87	73.20
ELMo (4-layer) best layer	81.58	93.81	97.31	95.60	89.78	82.06	94.18	29.24	74.78	85.96	73.03
ELMo (transformer) best layer	80.97	92.68	97.09	95.13	93.06	81.21	93.78	30.80	72.81	82.24	70.88
OpenAI transformer best layer	75.01	82.69	93.82	91.28	86.06	58.14	87.81	33.10	66.23	76.97	74.03
BERT (base, cased) best layer	84.09	93.67	96.95	95.21	92.64	82.71	93.72	43.30	79.61	87.94	75.11
BERT (large, cased) best layer	85.07	94.28	96.73	95.80	93.64	84.44	93.83	46.46	79.17	90.13	76.25
GloVe (840B.300d)	59.94	71.58	90.49	83.93	62.28	53.22	80.92	14.94	40.79	51.54	49.70
Previous state of the art (without pretraining)	83.44	94.7	97.96	95.82	95.77	91.38	95.15	39.83	66.89	78.29	77.10

Tagging Results

Pretrained Representation	POS			Supersense ID							
	Avg.	CCG	PTB	EWT	Chunk	NER	ST	GED	PS-Role	PS-Fxn	EF
ELMo (original) best layer	81.58	93.31	97.26	95.61	90.04	82.85	93.82	29.37	75.44	84.87	73.20
ELMo (4-layer) best layer	81.58	93.81	97.31	95.60	89.78	82.06	94.18	29.24	74.78	85.96	73.03
ELMo (transformer) best layer	80.97	92.68	97.09	95.13	93.06	81.21	93.78	30.80	72.81	82.24	70.88
OpenAI transformer best layer	75.01	82.69	93.82	91.28	86.06	58.14	87.81	33.10	66.23	76.97	74.03
BERT (base, cased) best layer	84.09	93.67	96.95	95.21	92.64	82.71	93.72	43.30	79.61	87.94	75.11
BERT (large, cased) best layer	85.07	94.28	96.73	95.80	93.64	84.44	93.83	46.46	79.17	90.13	76.25
GloVe (840B.300d)	59.94	71.58	90.49	83.93	62.28	53.22	80.92	14.94	40.79	51.54	49.70
Previous state of the art (without pretraining)	83.44	94.7	97.96	95.82	95.77	91.38	95.15	39.83	66.89	78.29	77.10

Context matters!

Layer-wise Prediction



(each column is
a different task)

Effect of Pretraining Task

- See also:
 - [Zhang and Bowman 2018](#)
 - [Peters et al 2018b](#)
 - [Blevins et al 2018](#)

Pretraining Task	Layer Average Target Task Performance			
	0	1	2	Mix
CCG	56.70	64.45	63.71	66.06
Chunk	54.27	62.69	63.25	63.96
POS	56.21	63.86	64.15	65.13
Parent	54.57	62.46	61.67	64.31
GParent	55.50	62.94	62.91	64.96
GGParent	54.83	61.10	59.84	63.81
Syn. Arc Prediction	53.63	59.94	58.62	62.43
Syn. Arc Classification	56.15	64.41	63.60	66.07
Sem. Arc Prediction	53.19	54.69	53.04	59.84
Sem. Arc Classification	56.28	62.41	61.47	64.67
Conj	50.24	49.93	48.42	56.92
BiLM	66.53	65.91	65.82	66.49
GloVe (840B.300d)	60.55			
Untrained ELMo (original)	52.14	39.26	39.39	54.42
ELMo (original) (BiLM on 1B Benchmark)	64.40	79.05	77.72	78.90

Is it in the probe or the representation?

Designing and Interpreting Probes with Control Tasks

John Hewitt

Stanford University

johnhew@stanford.edu

Percy Liang

Stanford University

pliang@cs.stanford.edu

Is it in the probe or the representation?

Designing and Interpreting Probes with Control Tasks

John Hewitt

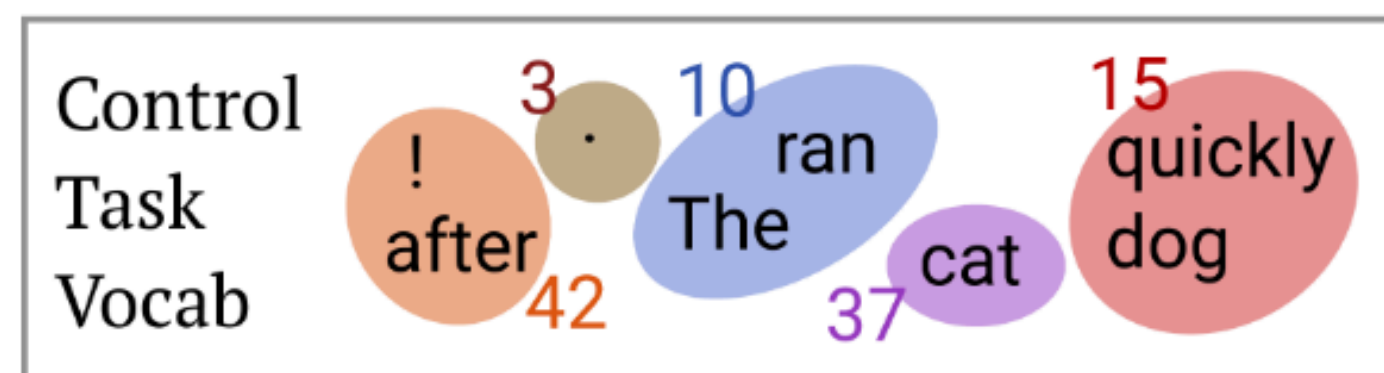
Stanford University

johnhew@stanford.edu

Percy Liang

Stanford University

pliang@cs.stanford.edu



Sentence 1	The	cat	ran	quickly	.
Part-of-speech	DT	NN	VBD	RB	.
Control task	10	37	10	15	3

Sentence 2	The	dog	ran	after	!
Part-of-speech	DT	NN	VBD	IN	.
Control task	10	15	10	42	42

Is it in the probe or the representation?

Designing and Interpreting Probes with Control Tasks

John Hewitt

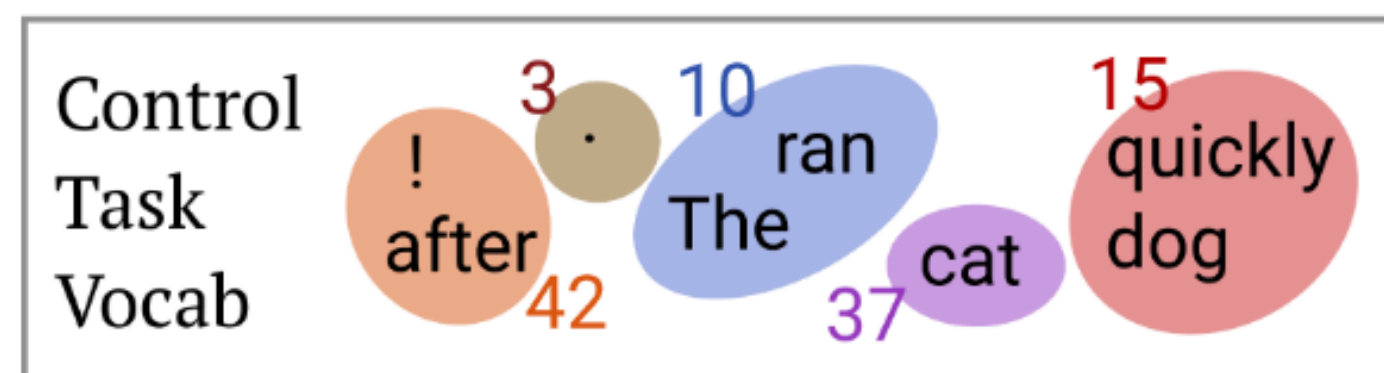
Stanford University

johnhew@stanford.edu

Percy Liang

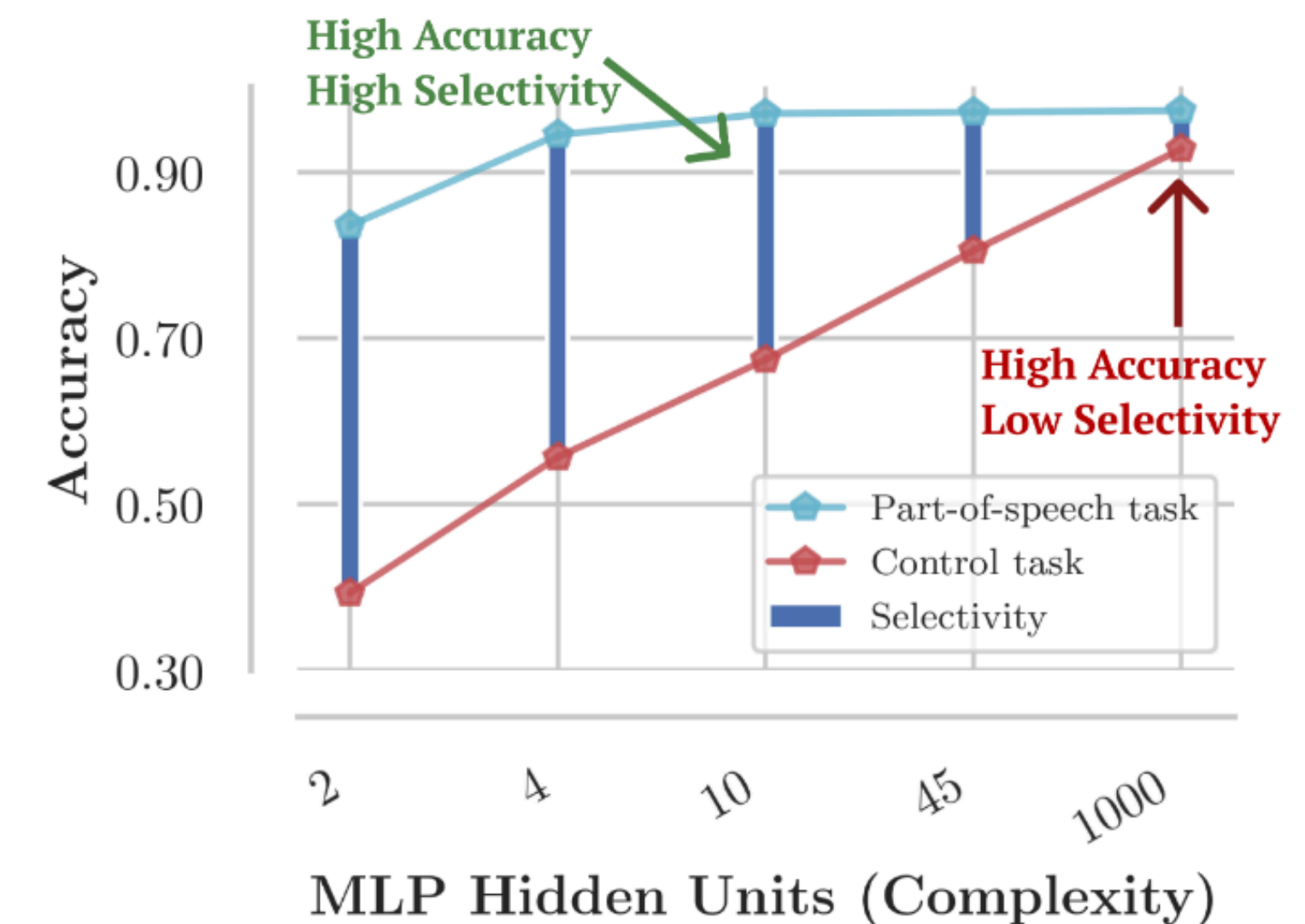
Stanford University

pliang@cs.stanford.edu



Sentence 1	The	cat	ran	quickly	.
Part-of-speech	DT	NN	VBD	RB	.
Control task	10	37	10	15	3

Sentence 2	The	dog	ran	after	!
Part-of-speech	DT	NN	VBD	IN	.
Control task	10	15	10	42	42



Summary

- Use simple classifiers to see what can be extracted from a model's representations.
- Some clear trends:
 - Contextualized representations have more info than global ones (GloVe e.g.)
 - Especially for syntax
 - Layer-wise: early recurrent layers are more transferrable, less clear on Transformers
 - Language modeling a very good task for building transferrable representations

Summary, cont.

- Promises:
 - Lets us learn what's encoded in a model's opaque representation
- Shortcomings:
 - Comparison/control
 - Correlation vs causation: encoding \neq used by the model
 - New methods try to overcome this

Probing Classifiers: Promises, Shortcomings, and Advances

Yonatan Belinkov*
Technion – Israel Institute of Technology
belinkov@technion.ac.il

Probing classifiers have emerged as one of the prominent methodologies for interpreting and analyzing deep neural network models of natural language processing. The basic idea is simple—a classifier is trained to predict some linguistic property from a model's representations—and has been used to examine a wide variety of models and properties. However, recent studies have demonstrated various methodological limitations of this approach. This squib critically reviews the probing classifiers framework, highlighting their promises, shortcomings, and advances.

Attention-based

What does BERT look at?

An Analysis of BERT's Attention

Kevin Clark[†] Urvashi Khandelwal[†] Omer Levy[‡] Christopher D. Manning[†]

[†]Computer Science Department, Stanford University

[‡]Facebook AI Research

{kevclark, urvashik, manning}@cs.stanford.edu

omerlevy@fb.com

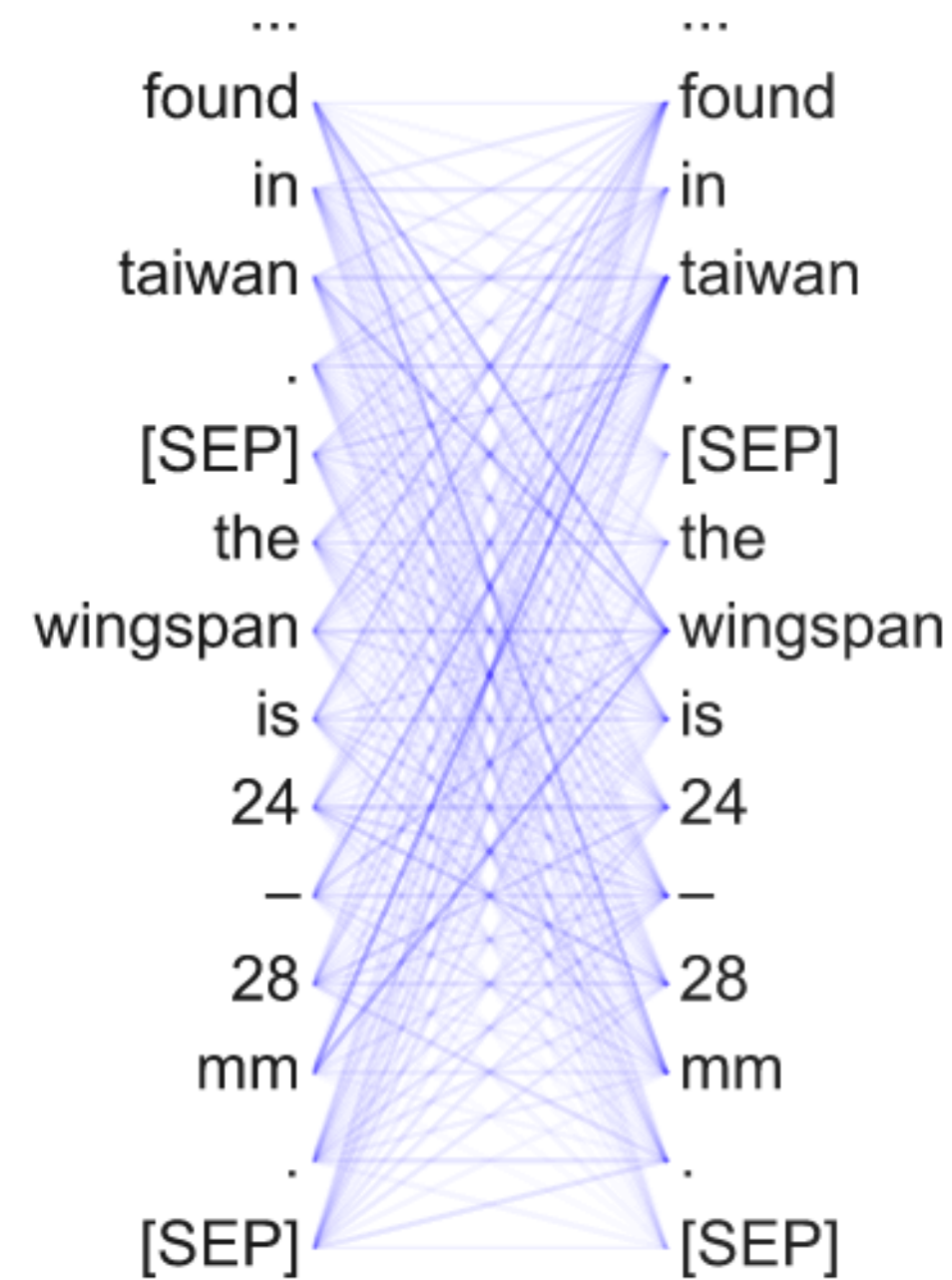
Abstract

Large pre-trained neural networks such as BERT have had great recent success in NLP, motivating a growing body of research investigating what aspects of language they are able to learn from unlabeled data. Most recent analysis has focused on model outputs (e.g., lan-

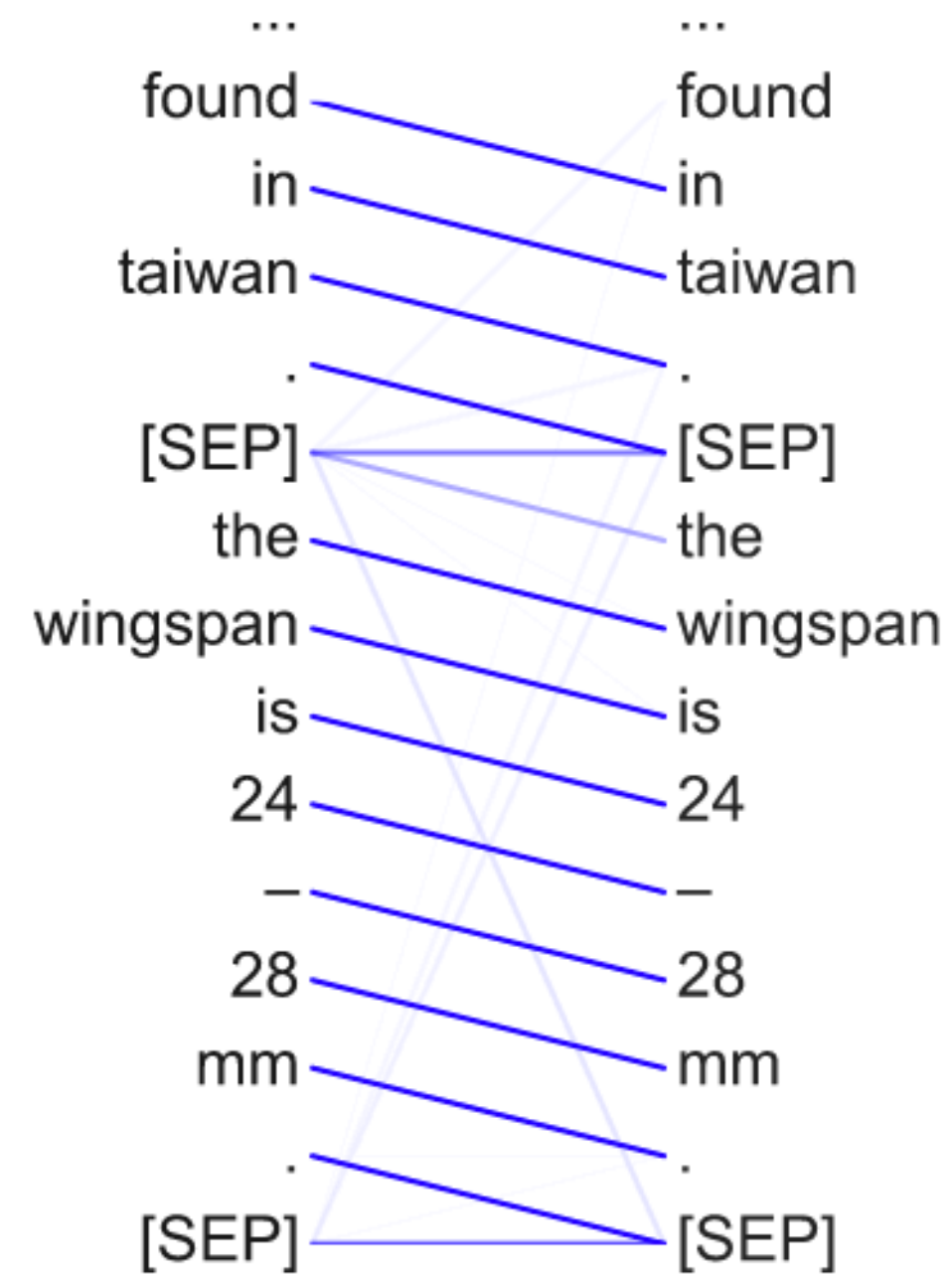
study¹ the *attention maps* of a pre-trained model. Attention (Bahdanau et al., 2015) has been a highly successful neural network component. It is naturally interpretable because an attention weight has a clear meaning: how much a particular word will be weighted when computing the next representation for the current word. Our analysis fo-

Qualitative Patterns

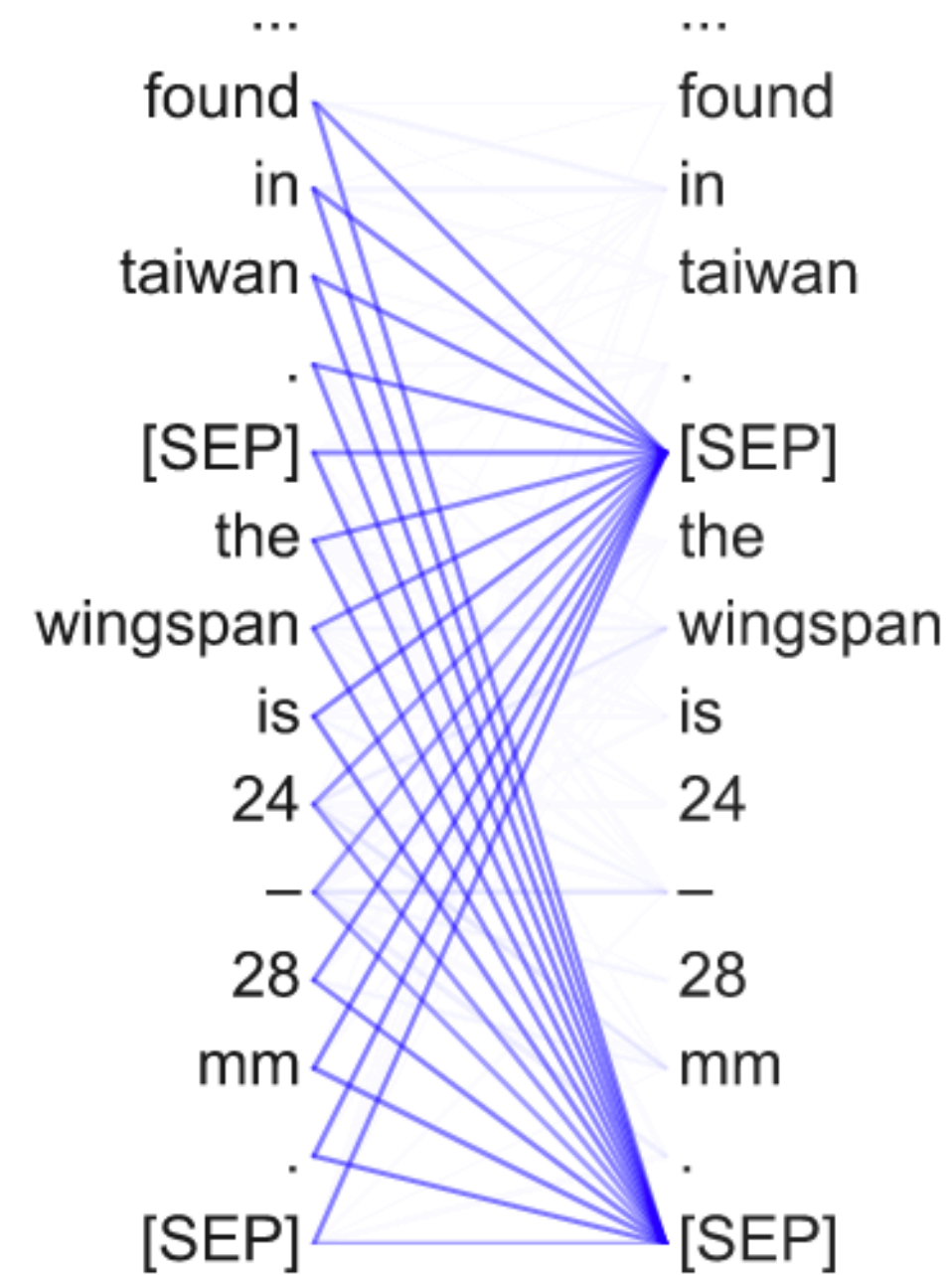
Head 1-1
Attends broadly



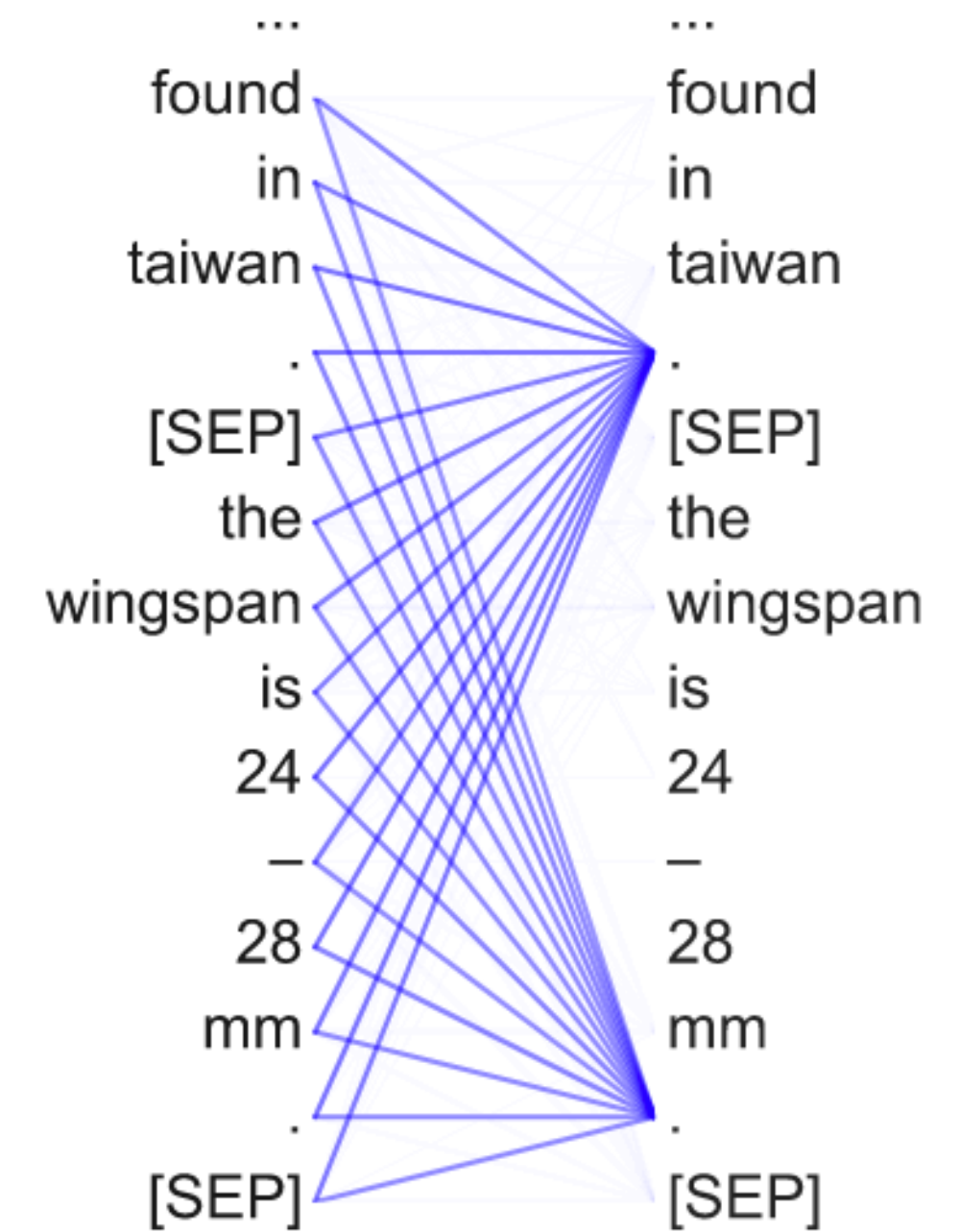
Head 3-1
Attends to next token



Head 8-7
Attends to [SEP]



Head 11-6
Attends to periods



Attention Head as Classifier

- No new training required
- Do any of these work for pairwise classification tasks “off-the-shelf”?

Attention Head as Classifier

- No new training required
- Do any of these work for pairwise classification tasks “off-the-shelf”?

$$\alpha_j = q \cdot k_j$$

$$e_j = e^{\alpha_j} / \sum_j e^{\alpha_j}$$

$$c = \sum_j e_j v_j$$

$$\text{class}(q) = \arg \max_j \alpha_j$$

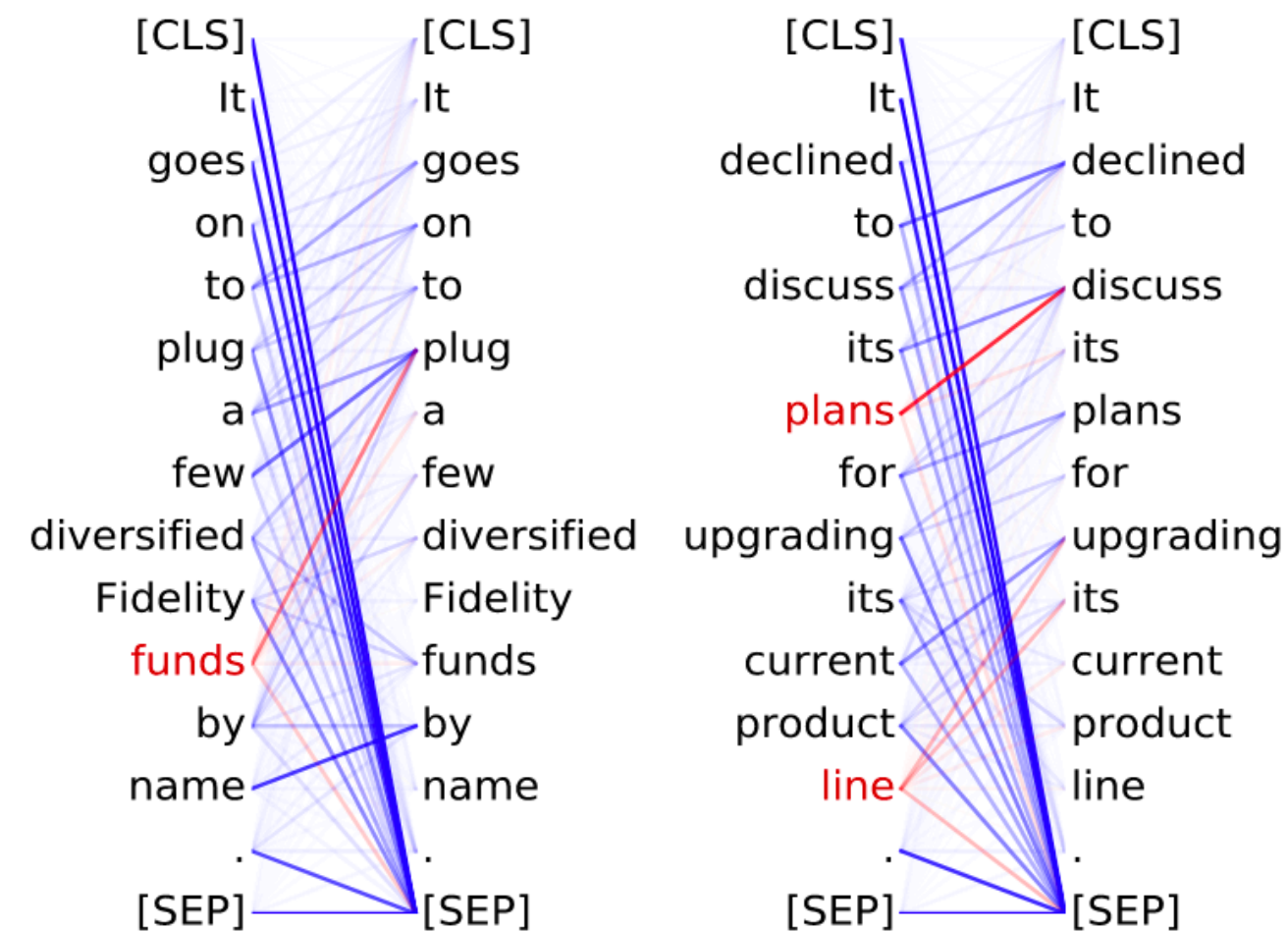
Dependency Parsing

Relation	Head	Accuracy	Baseline
All	7-6	34.5	26.3 (1)
prep	7-4	66.7	61.8 (-1)
pobj	9-6	76.3	34.6 (-2)
det	8-11	94.3	51.7 (1)
nn	4-10	70.4	70.2 (1)
nsubj	8-2	58.5	45.5 (1)
amod	4-10	75.6	68.3 (1)
dobj	8-10	86.8	40.0 (-2)
advmod	7-6	48.8	40.2 (1)
aux	4-10	81.1	71.5 (1)
poss	7-6	80.5	47.7 (1)
auxpass	4-10	82.5	40.5 (1)
ccomp	8-1	48.8	12.4 (-2)
mark	8-2	50.7	14.5 (2)
prt	6-7	99.1	91.4 (-1)

Examples

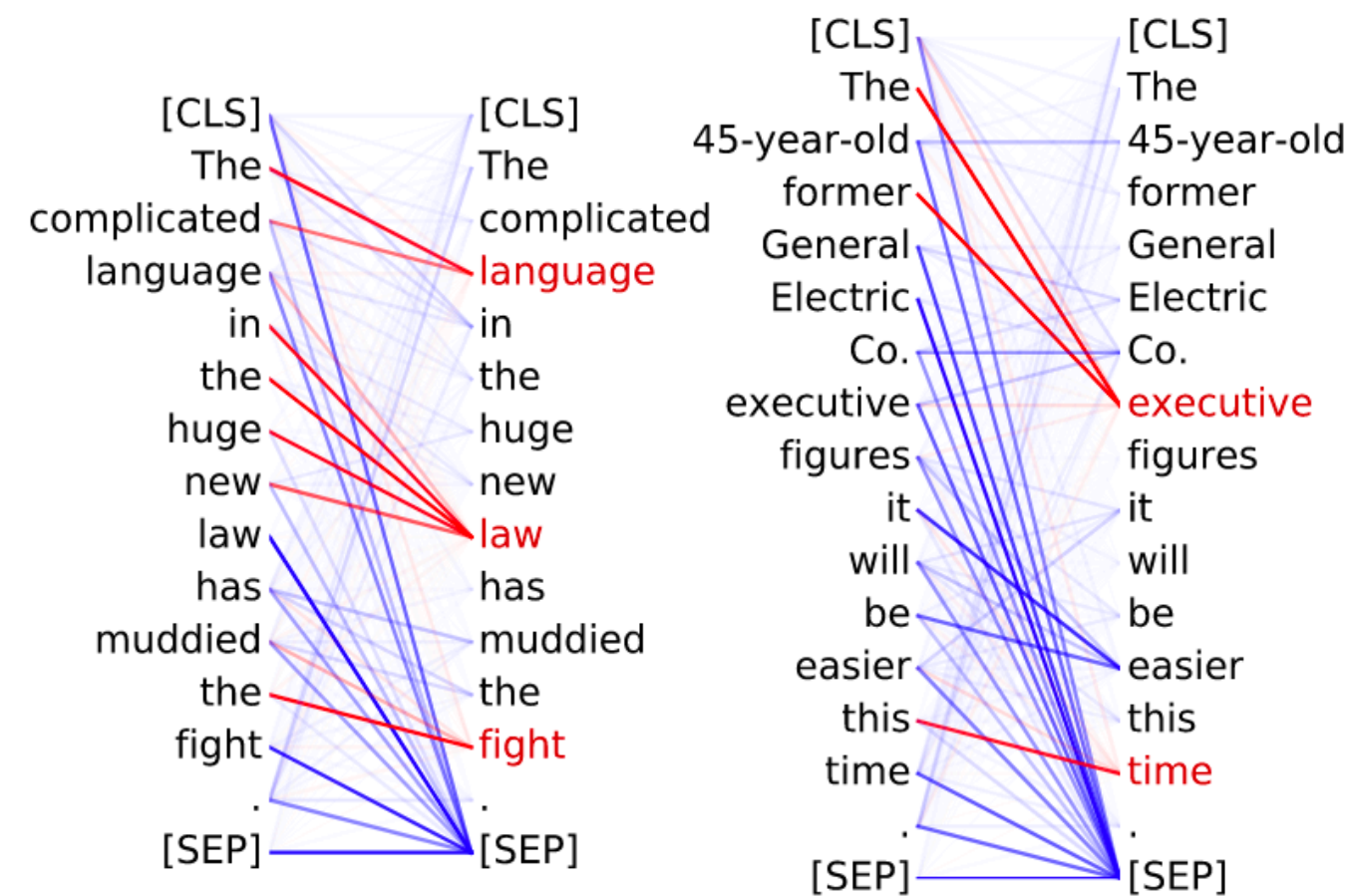
Head 8-10

- **Direct objects** attend to their verbs
- 86.8% accuracy at the dobj relation



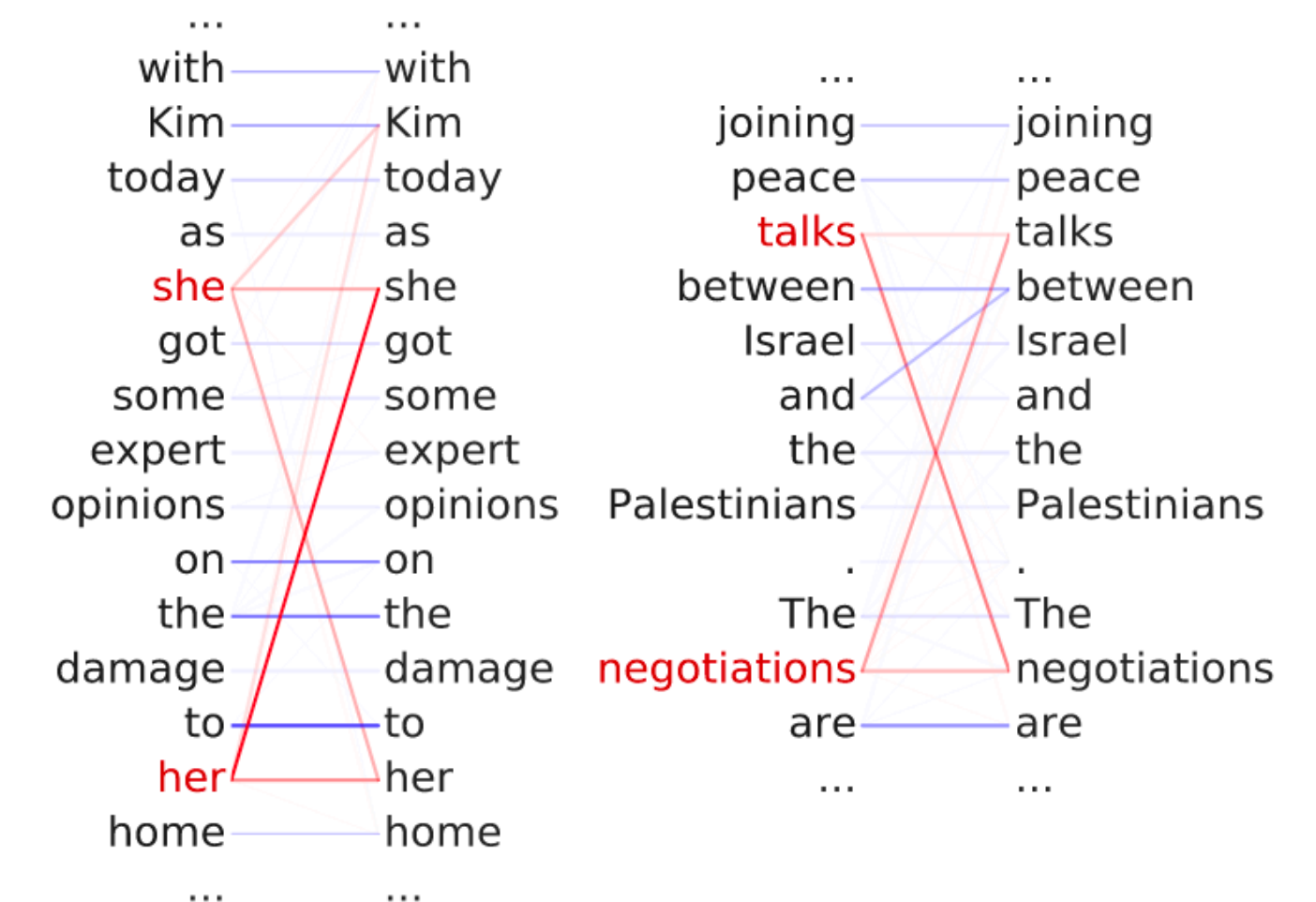
Head 8-11

- **Noun modifiers** (e.g., determiners) attend to their noun
- 94.3% accuracy at the det relation



Head 5-4

- **Coreferent** mentions attend to their antecedents
- 65.1% accuracy at linking the head of a coreferent mention to the head of an antecedent



Revealing the Dark Secrets of BERT

Olga Kovaleva, Alexey Romanov, Anna Rogers, Anna Rumshisky

Department of Computer Science
University of Massachusetts Lowell
Lowell, MA 01854

`{okovalev, arum, aromanov}@cs.uml.edu`

Abstract

BERT-based architectures currently give state-of-the-art performance on many NLP tasks, but little is known about the exact mechanisms that contribute to its success. In the current work, we focus on the interpretation of self-attention, which is one of the fundamental underlying components of BERT. Using a subset of GLUE tasks and a set of handcrafted features-of-interest, we propose the methodology and carry out a qualitative and quantita-

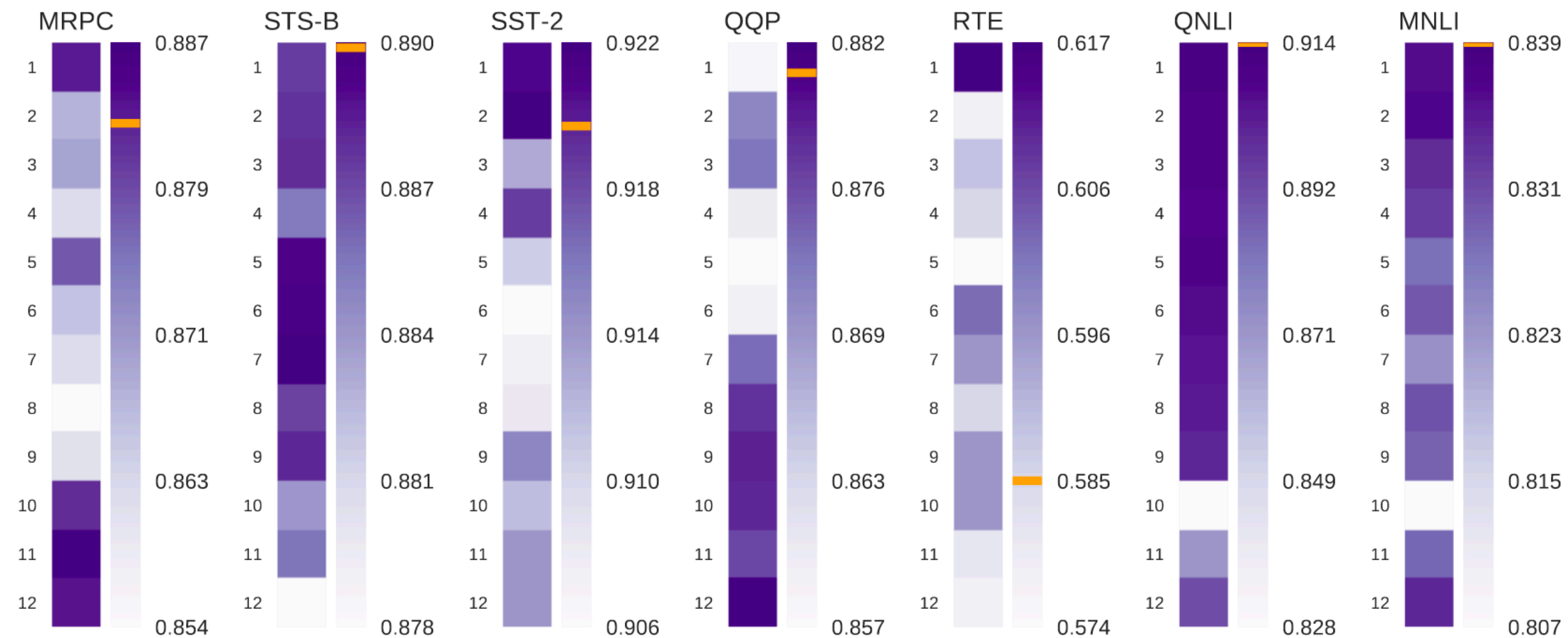
inference. State-of-the-art performance is usually obtained by fine-tuning the pre-trained model on the specific task. In particular, BERT-based models are currently dominating the leaderboards for SQuAD¹ ([Rajpurkar et al., 2016](#)) and GLUE benchmarks² ([Wang et al., 2018](#)).

However, the exact mechanisms that contribute to the BERT’s outstanding performance still remain unclear. We address this problem through selecting a set of linguistic features of interest and

Overall

- Same observation as previous: many heads only pay attention to [SEP] and [CLS] tokens
- Changes in attention before and after fine-tuning
- Pruning some heads can actually improve performance
 - (see also [Voita et al](#) on the original Transformer)

Pruning all attention in a layer



← pay attention to the scales

Summary

- Sometimes, attention heads seem to encode some linguistically interesting properties
 - But there appears to be lots of redundancy
 - And there's much more terrain to explore here
- As before: we can ask if property P can be found in attention, but not what role (independently of a hypothesis) a head is playing
- For the curious: ongoing debate about the connection between attention and model predictions (not as applied to LMs yet): Attention is not explanation;
Attention is not not explanation

Adversarial Datasets

Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference

R. Thomas McCoy,¹ Ellie Pavlick,² & Tal Linzen¹

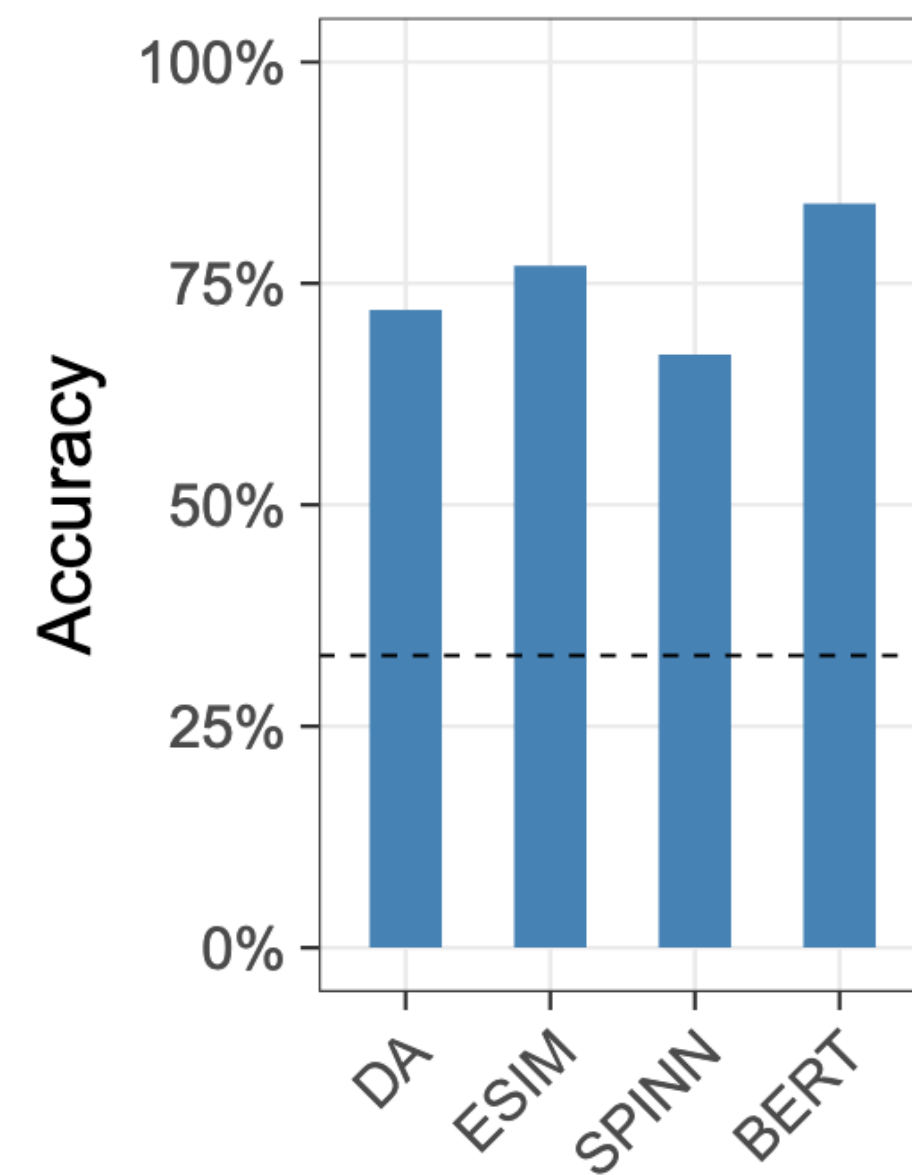
¹Department of Cognitive Science, Johns Hopkins University

²Department of Computer Science, Brown University

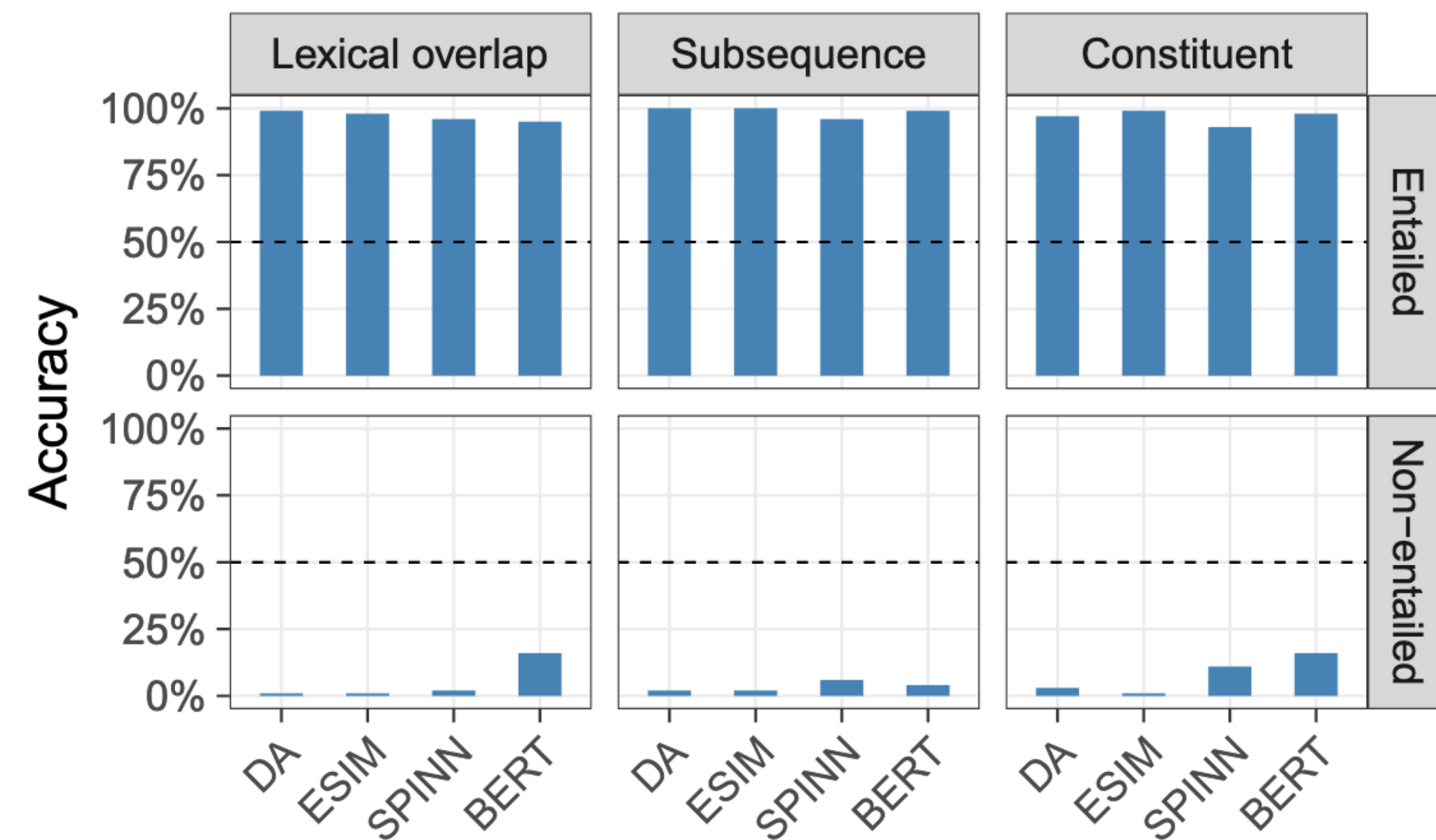
`tom.mccoy@jhu.edu, ellie_pavlick@brown.edu, tal.linzen@jhu.edu`

Heuristic	Premise	Hypothesis	Label
Lexical overlap heuristic	The banker near the judge saw the actor.	The banker saw the actor.	E
	The lawyer was advised by the actor.	The actor advised the lawyer.	E
	The doctors visited the lawyer.	The lawyer visited the doctors.	N
	The judge by the actor stopped the banker.	The banker stopped the actor.	N
Subsequence heuristic	The artist and the student called the judge.	The student called the judge.	E
	Angry tourists helped the lawyer.	Tourists helped the lawyer.	E
	The judges heard the actors resigned.	The judges heard the actors.	N
	The senator near the lawyer danced.	The lawyer danced.	N
Constituent heuristic	Before the actor slept, the senator ran.	The actor slept.	E
	The lawyer knew that the judges shouted.	The judges shouted.	E
	If the actor slept, the judge saw the artist.	The actor slept.	N
	The lawyers resigned, or the artist slept.	The artist slept.	N

Results



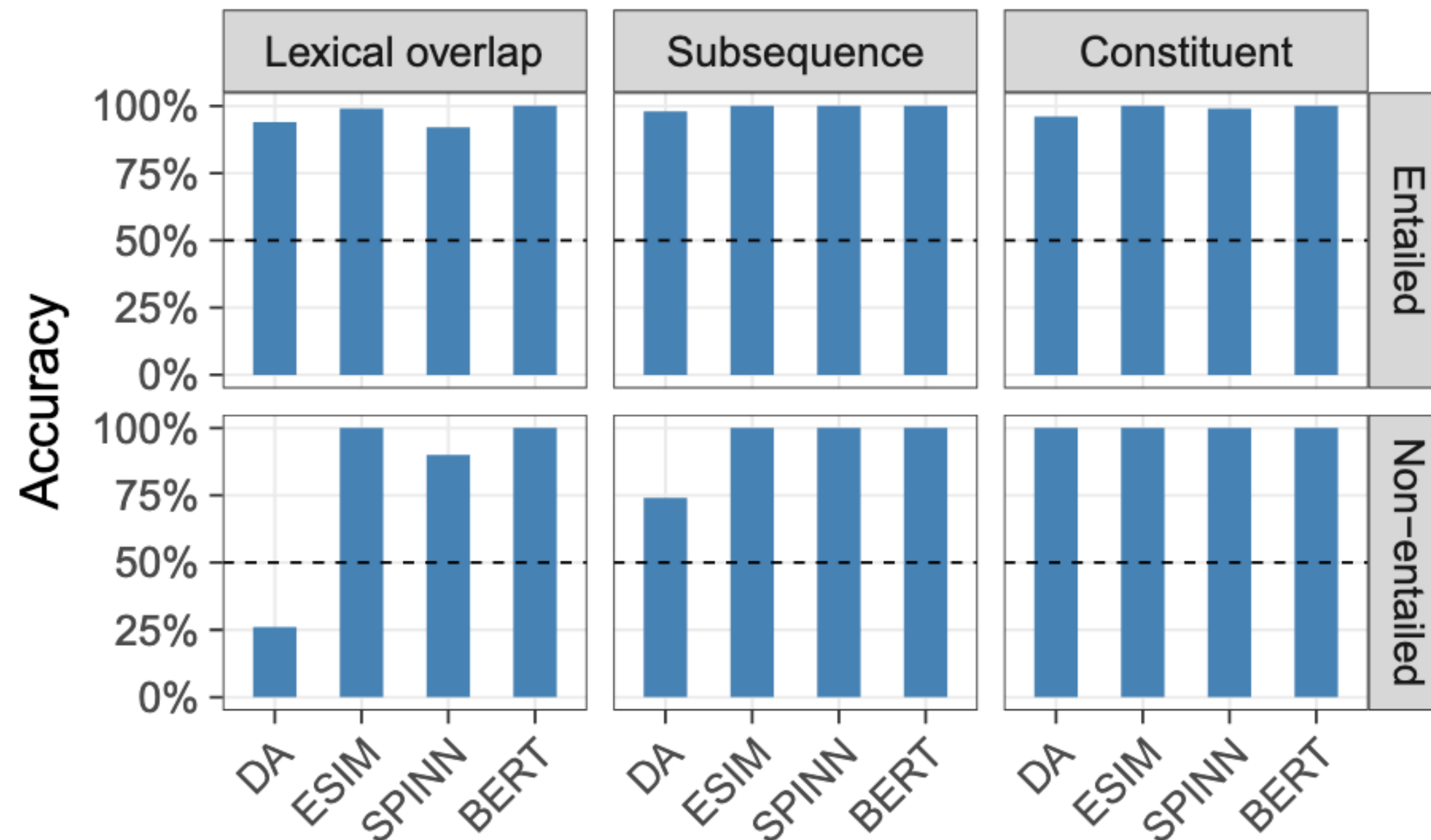
(a)



(b)

(performance improves if fine-tuned on this challenge set)

Fine-tuning augmented with examples



Conclusion

- Solving a dataset \neq solving a task
 - Models are very powerful, can be very “clever”
 - Adopt heuristics that exploit spurious cues in the data
- Careful design of “adversarial” data can both expose the heuristics being relied on and hopefully improve the representations learned

Problem with Probing

- Recall the issue with diagnostic classifiers / probing:
 - We can learn that property X is encoded in representation R
 - But not: does the model use property X in making its decisions
- Main idea here: *causally intervene* on the model and/or data to figure out which properties the model is relying on
 - Somewhat analogous to individual neuron ablation
 - E.g. if we “remove all number information” from R , does the model’s performance on a given task suffer

One last meta-point

Investigating BERT’s Knowledge of Language: Five Analysis Methods with NPIs

Alex Warstadt,^{†,1,2} Yu Cao,^{†,3} Ioana Grosu,^{†,2} Wei Peng,^{†,3} Hagen Blix,^{†,1}
Yining Nie,^{†,1,2} Anna Alsop,^{†,2} Shikha Bordia,^{†,3} Haokun Liu,^{†,3} Alicia Parrish,^{†,2,3}
Sheng-Fu Wang,^{†,3} Jason Phang,^{†,1,3} Anhad Mohananey,^{†,1,3} Phu Mon Htut,^{†,3}
Paloma Jeretić,^{†,1,2} and Samuel R. Bowman
New York University

[†]Equal contribution with roles given below; order assigned randomly. Correspondence: bowman@nyu.edu

¹Framing and organizing the paper ²Creating diagnostic data ³Constructing and running experiments

Abstract

Though state-of-the-art sentence representation models can perform tasks requiring significant knowledge of grammar, it is an open question how best to evaluate their grammatical knowledge. We explore five experimental methods inspired by prior work evaluating pretrained sentence representation models. We use a single linguistic phenomenon, negative polarity item (NPI) licensing in English, as a case study for our experiments. NPIs like *any* are grammatical only if they appear in a *licensing environment* like negation (*Sue doesn’t have any cats* vs. **Sue has any cats*).

acceptability. [Linzen et al. \(2016\)](#), [Warstadt et al. \(2018\)](#), and [Kann et al. \(2019\)](#) use Boolean acceptability judgments inspired by methodologies in generative linguistics. However, we have not yet seen any substantial direct comparison between these methods, and it is not yet clear whether they tend to yield similar conclusions about what a given model knows.

We aim to better understand the trade-offs in task choice by comparing different methods inspired by previous work to evaluate sentence understanding models in a single empirical domain. We choose as our case study negative polarity

Negative polarity items

- NPIs are expressions like *any*, *ever* that are only grammatical in “negative” environments:
 - * Shaan has done *any* of the reading.
 - Shaan hasn’t done *any* of the reading.
- Question: does BERT “understand” NPIs?
- See also [Marvin and Linzen 2018](#); [Jumelet and Hupkes 2018](#); Jumelet et al 2021 Findings of ACL

Does BERT “understand” NPIs?

- It depends!
- “We find that BERT has significant knowledge of these features, but its success varies widely across different experimental methods. We conclude that a variety of methods is necessary to reveal all relevant aspects of a model’s grammatical knowledge in a given domain.”

Wrapping Up

Interpretability and Analysis

- Current NLP models are often a “black box”, trained on huge amounts of data, which makes it very unclear what they are learning from their data
 - Engineering: build better models for the future [though caveat emptor]
 - Theoretical: what kinds of linguistic information are learnable (and not) from what kinds of data
 - Ethical: what harmful effects are learned from the data, and how can these be mitigated
- Methods briefly surveyed: neuron-level, psycholinguistic, diagnostic classifiers (+ causal variants), attention analysis, adversarial data
- A huge and growing area!