

Cognitive Plausibility of Deep Language Models

Day 4 & 5

Lisa Beinborn
VU Amsterdam

Nora Hollenstein
U Copenhagen

Willem Zuidema
Uv Amsterdam



Interpretability of Neural Language Models

Willem Zuidema
@AmsterdamNLP



Plan

1. LSTM language models

- recap & deep dive - key feature: trainable gates
- interpretability case study I: integrated gradients
- interpretability case study II: diagnostic classifiers (probes)

2. Transformer language models

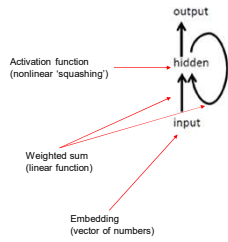
- recap & deep dive - key features: trainable attention & content-addressability
- interpretability case study III: attention tracking
- interpretability case study IV: representational similarity

3. The bigger picture: The Linguistics of Deep Learning

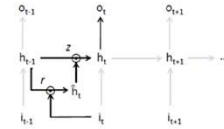


Background I: Gating in Recurrent Networks

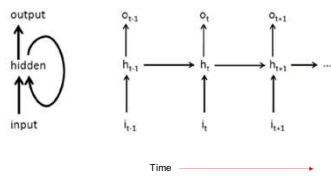
Recurrent network



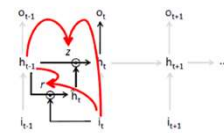
Add reset gate



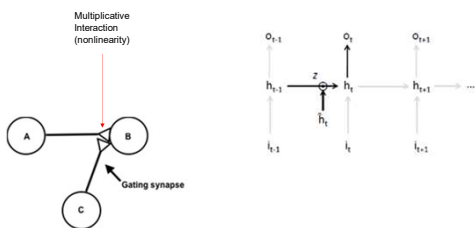
Unfold in time



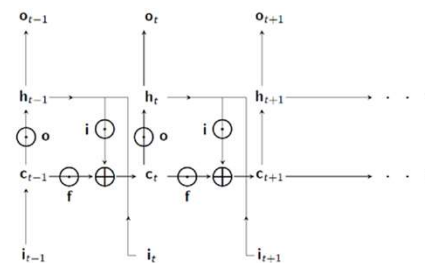
Make gates trainable: GRU (Cho et al., 2015)



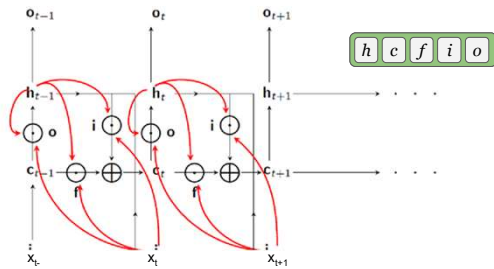
Add update gate



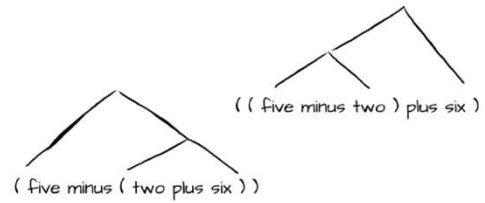
Add memory cell & forget gate



Long-short term memory (Hochreiter & Schmidhuber 1997)



Arithmetic Language



(Veldhoen et al. 2016; Hupkes, Veldhoen & Zuidema, 2018, JAIR)

Can recurrent networks learn languages with combinatorial and hierarchical structure?

Contextfreeness

$A^n B^n, n \geq 1$

AAAABBBB
AABB
AAAAAABBBBBB
...

Negative exs:

AABBB
ABABABAB

Palindromes

$(ww^R), w \in \{A,B,C,D\}^*$

DABBBBAD
BB
ACCCCCCA
...

Negative exs:

AABBCCBBA
DDBCADDBCA

Arithmetics

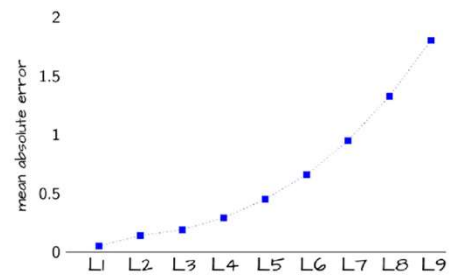
(Expr Op Expr)

$(1 + 3) - (7 - 9)$
 $(2 - 7)$
 $(9 - 3) - ((4 + 2) - (5 + 6))$
...

Targets:

6; -5; 11

Results



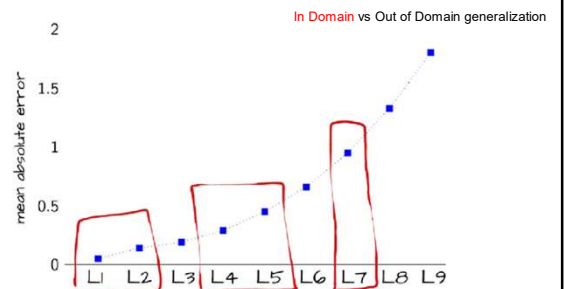
(Veldhoen et al. 2016; Hupkes, Veldhoen & Zuidema, 2018, JAIR)

How do you demonstrate that a neural network has really learned the intended generalization?

Behavioral tests:

- Generate examples, divide randomly in train set and test set
- Accuracy of test set := 'in domain' generalization
- Leave specific categories of examples out from train set
- Accuracy on various test sets := 'out of domain' generalization

Results



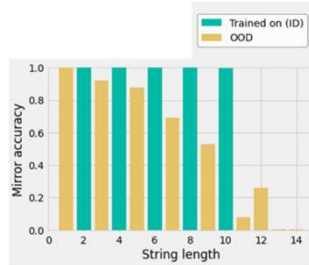
(Veldhoen et al. 2016; Hupkes, Veldhoen & Zuidema, 2018, JAIR)

Palindrome language

Strings of the form: g**b**h**d**d#d**d**h**b**g
 Alphabet: 8 symbols (a-h)
 Lengths train set: 5,9,13,17,21

g**b**h**d**d#d**d**h**b**g
 a**b**b**h**f**f**g**a**#a**g**f**f**h**b**b**a**
 a**a**a**a**a**a**a**a**a#a**a**a**a**a**a**a**a**a

LSTM, 1-layer, Dim=170



(Jumule & Zuidema, in prep)

Case study 1: Attribution methods / Integrated Gradients

Which words are important to generate each prediction?

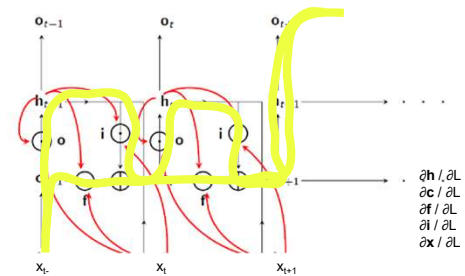
How do you demonstrate that a neural network has really learned the intended generalization?

Behavioral tests:

- Generate examples, divide randomly in train set and test set
- Accuracy of test set := 'in domain' generalization
- Leave specific categories of examples out from train set
- Accuracy on various test sets := 'out of domain' generalization

What level of accuracy do we consider convincing proof?

Tracking gradients



How do you demonstrate that a neural network has really learned the intended generalization?

Behavioral tests:

- Generate examples, divide randomly in train set and test set
- Accuracy of test set := 'in domain' generalization
- Leave specific categories of examples out from train set
- Accuracy on various test sets := 'out of domain' generalization

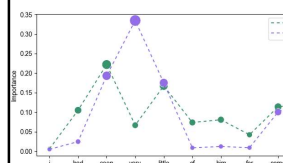
What level of accuracy do we consider convincing proof?

Can we look inside and characterize the learned solution? → **Interpretability**

Saliency 1:= length of the gradient

Beinborn & Hollenstein, "Relative Importance in Sentence Processing", ACL' 21,

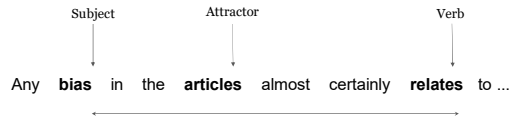
We iterate over each token vector x_i in our input sequence x_1, x_2, \dots, x_n . Let X_i be the input matrix with x_i being masked. The saliency s_{ij} for input token x_j for the prediction of the correct token t_i is then calculated as the Euclidean norm of the gradient of the logit for x_i .



$$s_{ij} = \|\nabla_{x_j} f_{t_i}(X_i)\|_2 \quad (1)$$

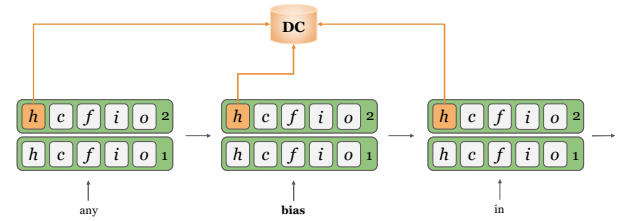
The saliency vector s_i indicates the relevance of each token for the correct prediction of the masked token t_i .¹ The saliency scores are normalized by dividing by the maximum. We determine the rel-

Subject-verb agreement



Linzen et al. 2016
Gulordava et al. 2018

Diagnostic Classification



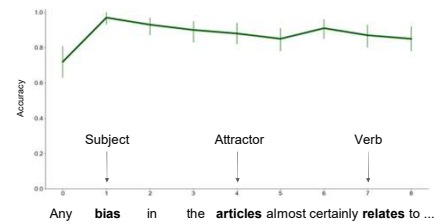
Experimental Setup

- Pretrained Neural Language Model from (Gulordava et al. 2018) with 2 LSTM-layers, with 650 hidden units each
- Wikipedia dependency dataset (Linzen et al. 2016)
- Extract activations for components h_t, c_t, f_t, i_t, o_t during forward pass of the LSTM

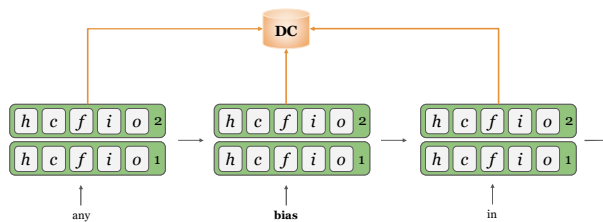
(Giulianelli, Harding, Mohnert, Hupkes & Zuidema, 2018)

Diagnostic Classification to Predict Number

Train: 1000 sentences, context size 5, at least 1 word before subject and at least 1 words after verb.
Test: Two sets of circa 100 sentences with 1 agreement attractor, according to correct/wrong number prediction.

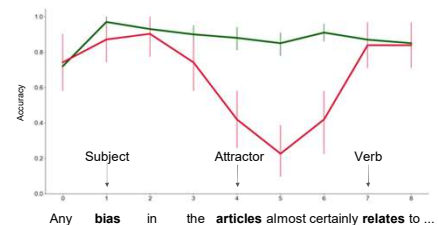


Diagnostic Classification

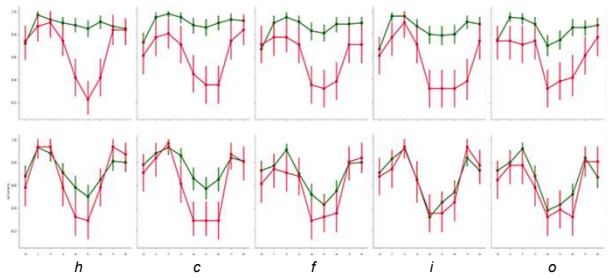


Diagnostic Classification to Predict Number

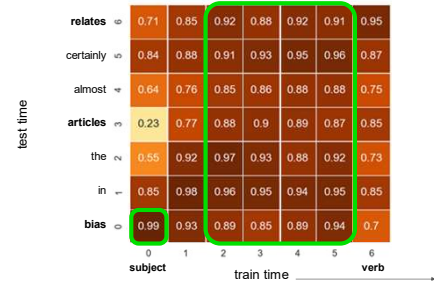
Train: 1000 sentences, context size 5, at least 1 word before subject and at least 1 words after verb.
Test: Two sets of circa 100 sentences with 1 agreement attractor, according to correct/wrong number prediction.



Diagnostic Classification to Predict Number

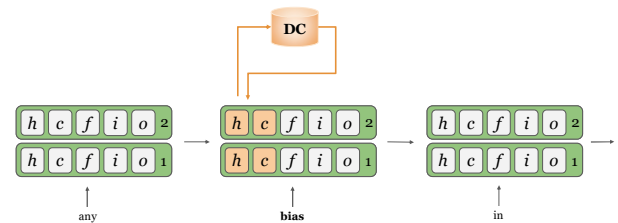


Diagnostic Classification to Predict Number



How is number agreement information processed across timesteps?

Influencing language models with diagnostic classifiers



Characterizing the dynamics of mental representations: the temporal generalization method

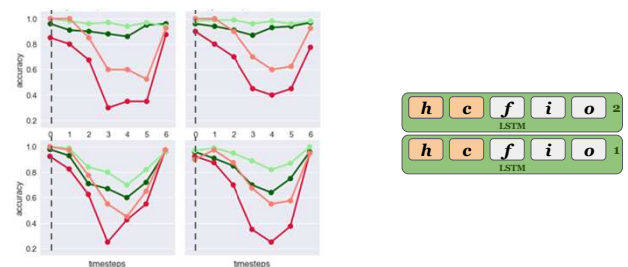
J.-R. King^{1,2,3} and S. Dehaene^{1,2,4,5}

¹Cognitive Neuroimaging Unit, Institut National de la Santé et de la Recherche Médicale, U992, F-91191 Gif/Yvette, France
²NeuroSpin Center, Institute of Biomedicine Commissariat à l'Energie Atomique, F-91191 Gif/Yvette, France
³Institut du Cerveau et de la Moelle Epinière Research Center, Institut National de la Santé et de la Recherche Médicale, U975 Paris, France
⁴Université Paris 11, Orsay, France
⁵Collège de France, F-75005 Paris, France

Parsing a cognitive task into a sequence of operations is a central problem in cognitive neuroscience. We argue that a major advance is now possible owing to the series of steps from those operating as a continuous flow or 'cascade' of overlapping stages [10]. More recently, the advent of brain-imaging techniques

Trends in Cognitive Science, 2014

Diagnostic Classification to Predict Number



Influencing language models with diagnostic classifiers

	An	official	estimate	issued	in	2003	suggests	suggest
Original	-11.05	-8.426	-8.472	-1.243	-3.951	-5.753	-5.6979	-6.4361
Intervention	-11.05	-8.426	-8.472	-1.268	-3.97	-5.691	-6.4361	



	without intervention	with intervention
	78.0	85.4

Plan

1. LSTM language models

- recap & deep dive - key feature: trainable gates
- interpretability case study I: integrated gradients
- interpretability case study II: diagnostic classifiers (probes)

2. Transformer language models

- recap & deep dive - key features: trainable attention & content-addressability
- interpretability case study III: attention tracking
- interpretability case study IV: representational similarity

3. The bigger picture: The Linguistics of Deep Learning

Lessons case study 2

Positives

- Diagnostic Classifiers allow us to track the dynamics of subject-verb agreement in an LSTM-based language model
- Temporal Generalization Method shows the LSTM represents number information in at least two different ways
- An intervention study allows us to go beyond correlation, but shows a causal role for the representations we identified

Negatives

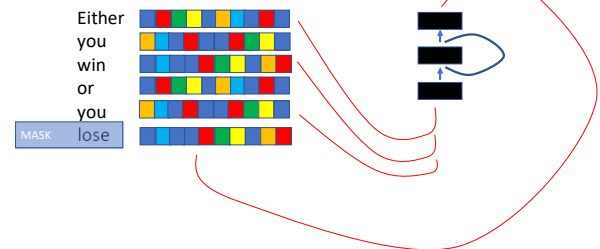
- Diagnostic classification is hypothesis-driven: we need good hypotheses first -- no complete understanding of the learned function reached yet
- Classifiers (especially nonlinear) may yield false positives, and false negatives

Background II: Attention in Transformers

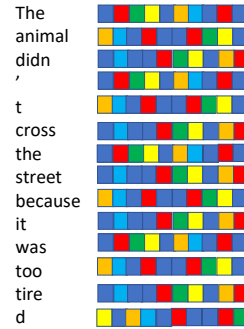
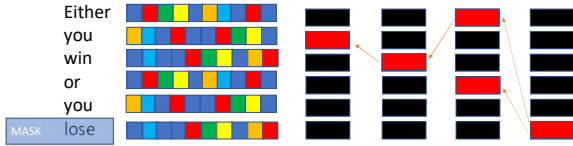
Questions?



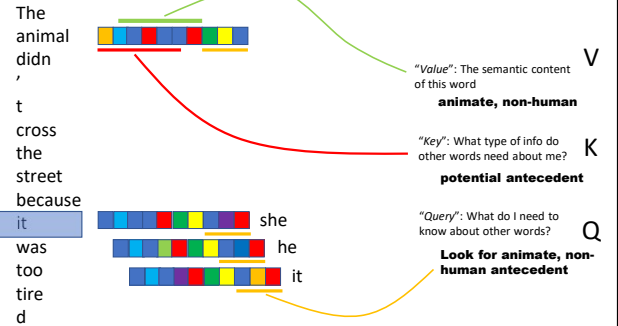
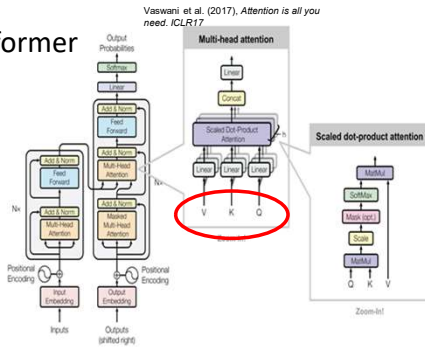
Until 2018: Recurrent networks



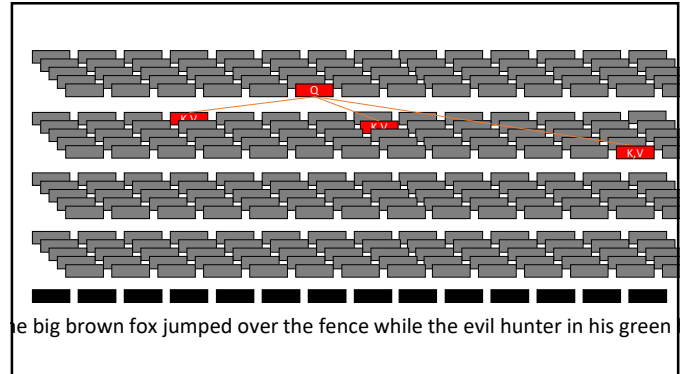
Since 2018: Transformers



The Transformer



The
animal
didn't
cross
the
street
because
it
was
too
tired
d

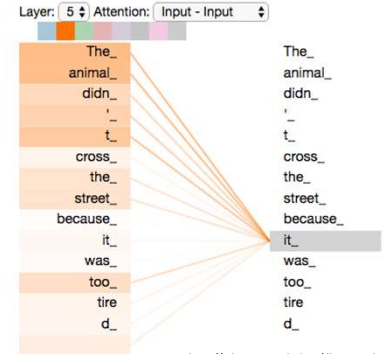


Transformer: "Attention is all you need"

- When considering the next word (w) to predict, each attention head can access information (i) of thousands of previous processing steps
 - w determines the *query*
 - i determines the *key*
 - If key and query 'match', the *value* extracted from i is used to compute the new state of the head
- Model with many layers, and many 'attention heads' per layer
 - ideal for parallelization on GPU's

Vaswani et al. (2017), *Attention is all you need*.
ICLR17
<https://arxiv.org/abs/1706.03762>

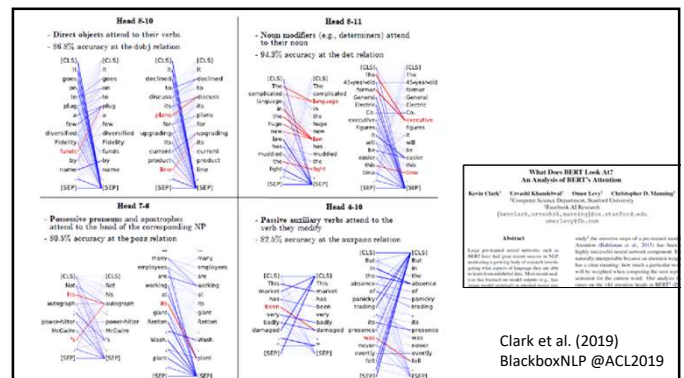
Attention tracking



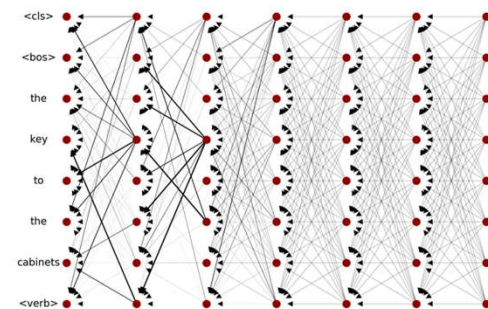
<http://jalammar.github.io/illustrated-transformer/>

Extremely large language models

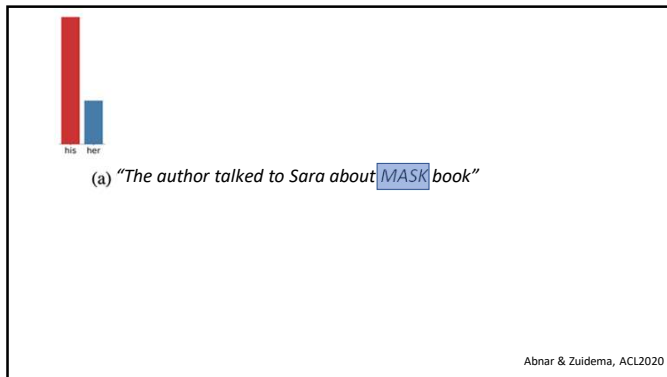
- Bert, GPT3: "Transformer" architecture (Vaswani et al. 2017; 3d most cited paper across academic fields in 2020)
- Extremely large deep learning model (GPT3: 178B parameters)
- Trained on enormous dataset (GPT3: 300M words, extracted from 1B word CommonCrawl + a number of custom datasets)
- Trained with enormous amount of compute (GPT3: ~\$12M), using a generalization of backpropagation of error.



Case study 3: Attention Tracking & Attention Flow



Abnar & Zuidema, ACL2020

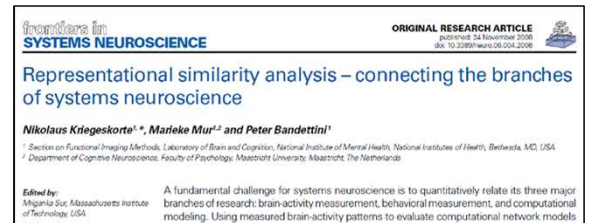
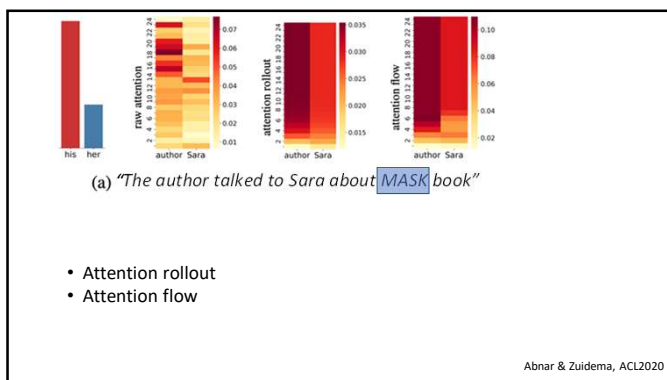


Case study 4: Representational Similarity & Stability Analysis

How similar are representations learned by different models, and how similar are they to representations in the brain?

(Abnar, Beinborn, Choenni & Zuidema, 2019)

BlackboxNLP @ACL2019



Lessons case study 3

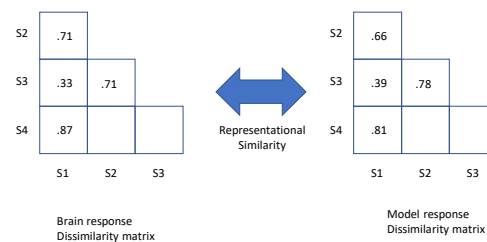
Positives

- Attention is the central mechanism in Transformer, and often turns out to correspond to linguistic functions
- Attention Roll-out & Flow take into account the structure of the attention network and approximate the effective attention at each point
- Attention Flow corresponds to a well-known attribution method: Shapley values

Negatives

- Attention is a local method; characterizing the globally learned function requires much extra work
- Attention is only one component; may yield false positives, and false negatives

Representational Similarity Analysis



How similar are representations learned by different deep learning models of language, and how similar are they to the brain?

(Abnar, Beinborn, Choenni & Zuidema, 2019)

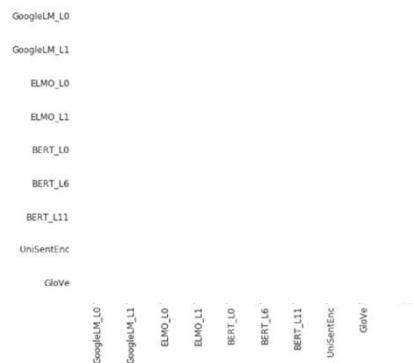
BlackboxNLP @ACL2019

Representational Stability Analysis

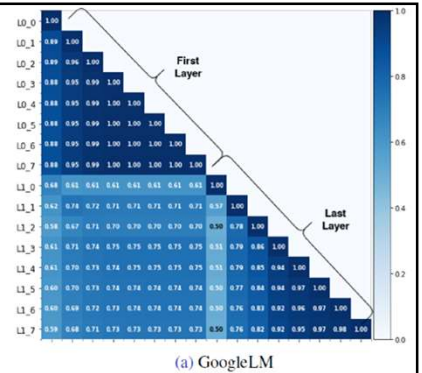
U0_0
U0_1
U0_2
U0_3
U0_4
U0_5
U0_6
U0_7
U1_0
U1_1
U1_2
U1_3
U1_4
U1_5
U1_6
U1_7

(a) GoogleLM

RSA across models

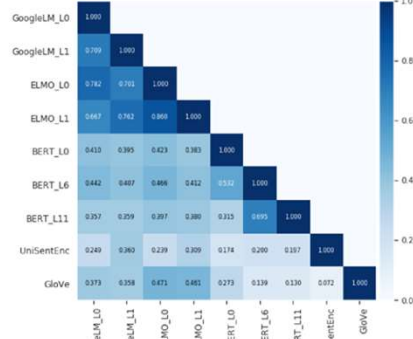


Representational Stability Analysis

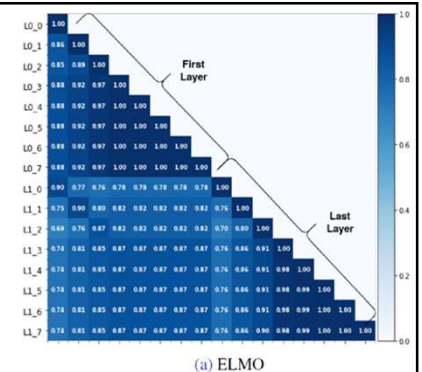


(a) GoogleLM

RSA across models

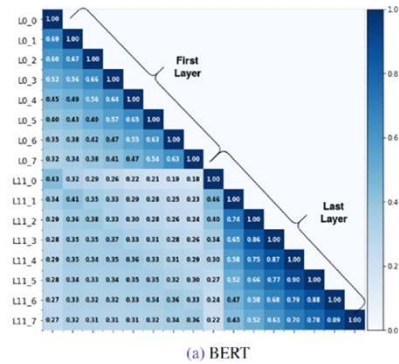


Representational Stability Analysis



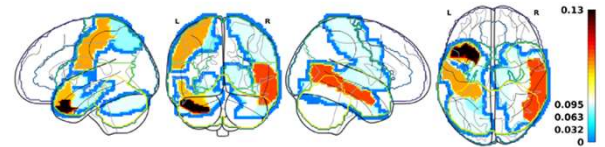
(a) ELMO

Representational Stability Analysis



Representational Similarity Analysis

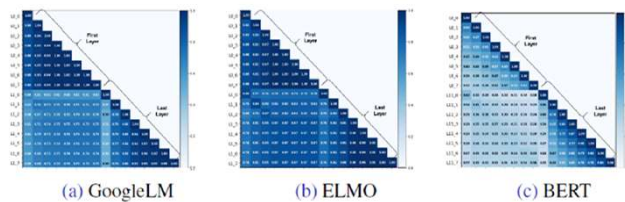
→ great differences between brain areas in their similarity to the models



<http://projects.illc.uva.nl/LaCo/clclab/>

Representational Stability Analysis

→ great differences between current deep language models in their dependence on context



Lessons case study 4

Positives

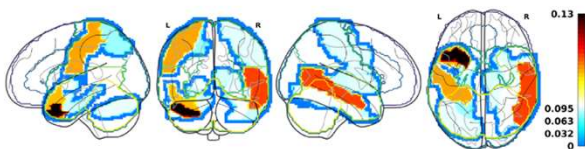
- Representational Similarity Analysis allows us to compare the representational spaces built up by models/brains with radically different architectures
- RSA is a (fairly) global interpretability method, that allows us to characterize the globally learned function
- Representation Stability Analysis allows us to characterize differences in sensitivity to manipulations of different components of the model

Negatives

- If we only compare uninterpretable models with uninterpretable brains, what have we really learned? → need to apply interpretability methods to the models

Representational Similarity Analysis

→ great differences between brain areas in their similarity to the models



<http://projects.illc.uva.nl/LaCo/clclab/>

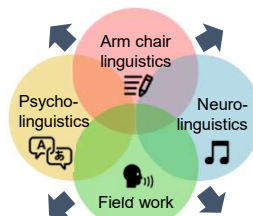
Lessons from all 4 case studies

- All state-of-the-art models in NLP are based on deep learning
- Presents us with the blackbox problem, making it difficult to:
 - Generate *explanations* to users and *justify* decisions based on the systems
 - Allow users to *interact* with the learned solutions and *adapt* them to their needs
 - Use prior knowledge to *augment* machine learned solutions
- Diagnostic classification is a way to test specific hypotheses on what information is represented; should be applied with as much rigor as model testing in (cognitive) neuroscience

Lessons from all 4 case studies (ctd.)

- Representational Similarity Analysis is a way to compare models across paradigms, and test the sensitivity of the learned representations to parameter choices
- Data- & hypothesis-driven methods, local & global methods are complementary, and can often be used together
- There is no silver bullet: the excellent performance of current models is found away from the easily interpretable points in hypothesis space
- We need to systematically apply the ever increasing toolbox of interpretability tools and see how far we get!

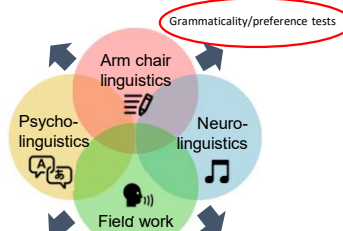
The linguistics of deep learning



Homework

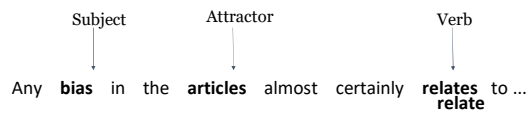
- Practice with Representation Similarity Analysis using this notebook: <https://colab.research.google.com/drive/1LtEwlaqXVEQXCvVTNaTG10ZlatkRL3M?usp=sharing>
- Practice with Word Embeddings and Dissimilarity Matrices here: https://colab.research.google.com/drive/1JeHTIZ9qi0LRHZurBerwHXOD06nrQb_2?usp=sharing

The linguistics of deep learning



The bigger picture:
The virtuous interaction between deep
learning, cognitive neuroscience &
linguistics

Subject-verb agreement



Linzen et al. 2016
Gulordava et al. 2018
Guilianelli, Harding, Mohnert, Hupkes & Zuidema, 2019

Bias

The author talked to Sara about MASK book



Mary convinced John of MASK love

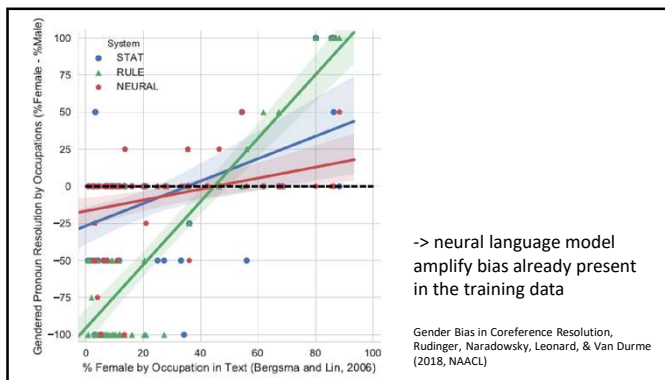


Abnar & Zuidema, ACL2020

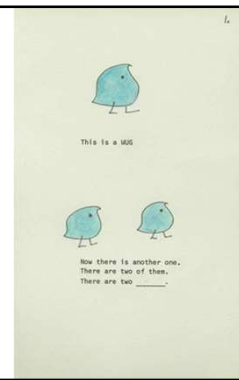
Compositionality

Principle of compositionality: the meaning of whole is a function of the meaning of the parts and the way they are put together

Compositional generalization: generalizing to new examples by reusing parts of earlier experiences in novel combinations

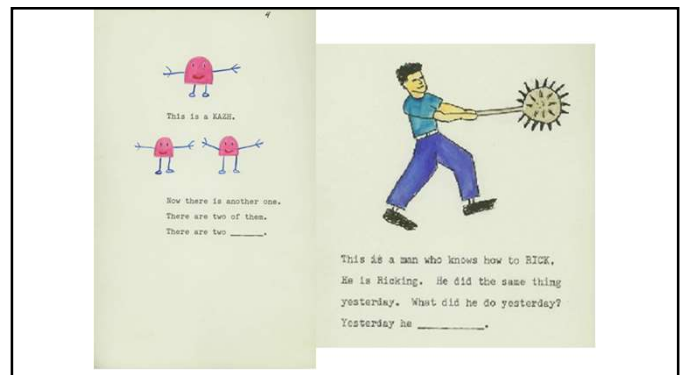
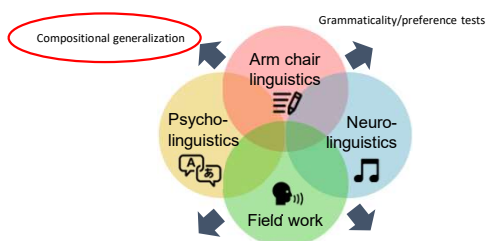


Wug test



Jean Berko Gleason (1958)

The linguistics of deep learning



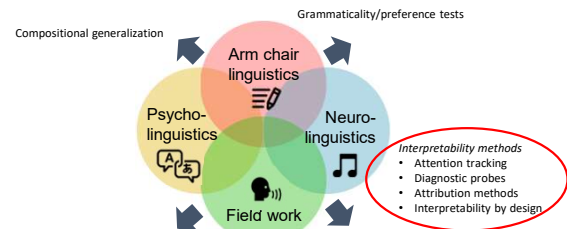
Zero-shot, one-shot, few-shot generalization

- Zero-shot generalization: generalizing to a new pattern without having seen single example of the target pattern
- One-shot generalization: generalizing to new examples based only on a single example of the target pattern
- Few-shot generalization: generalizing to new examples based only on a handful of examples of the target pattern

Language Models are Few-Shot Learners

Tom R. Brown*	Shayne Neel*	Nik Rishi*	Melanie Snickel*
Jared Kaplan*	Prithvi Rajaram*	Arvind Neelakantan	Prafulla Mishra
Shyam Prasad	Samuel Alpert	Alec Hefner	Constance Keller
Benjamin Lester	Jeffrey Berman	Ronald M. Davis	Jeffrey Wei
Cheng-Yu Hsieh			

The linguistics of deep learning

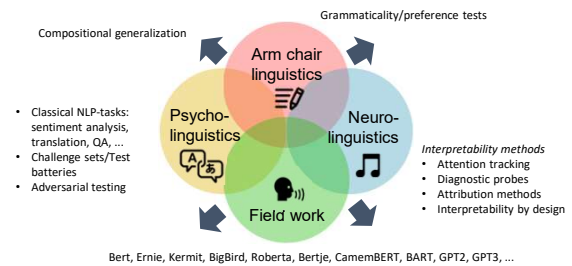


GPT-3

[Human prompt] To do a “farduddle” means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

[GPT-3 continuation] One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduddles.

The linguistics of deep learning



Compositional generalization

all Germans love all Italians
implies
some Germans love some Romans

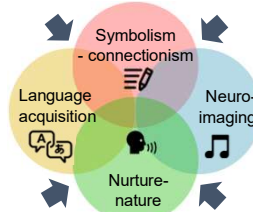
Zero-shot generalization to:

all Germans love all Italians
implies
some Germans love some Venetians

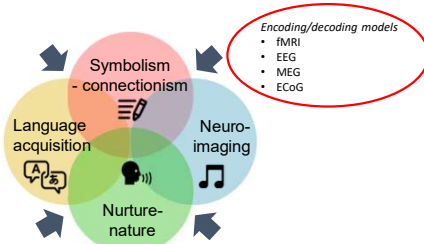
(but not from “all French hate Parisians” to “all French detest Parisians”)

Mul & Zuidema (2019): Siamese recurrent networks learn first-order logic reasoning and exhibit zero-shot compositional generalization

Deep learning contributing to linguistics



Deep learning contributing to linguistics

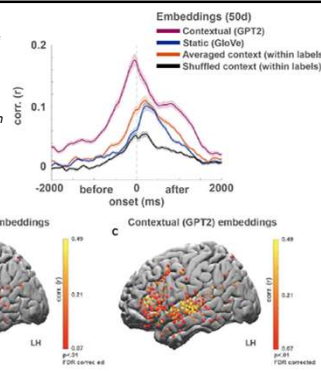


Compositionality revisited

- The principle of compositionality
 - the meaning of whole is a function of the meaning of the parts and the way they are put together*
- Do these deep learning models operate according to the principle of compositionality"?
 - They do not generalize perfectly to all novel combination or arbitrary levels of embedding;
 - They reach a level of performance incompatible with a memorization strategy;
 - Generalization is noisy -- but the networks approximate a truly compositional strategy

Ariel Goldstein et al (2021, BioRxiv), Thinking ahead: spontaneous next word predictions in context as a keystone of language in humans and machines

using spoken narrative and electrocorticographic recordings... we demonstrate that the brain constantly and spontaneously predicts the identity of the next word in natural speech, hundreds of milliseconds before they are perceived.



The principle & the approximation

- Is it disappointing that the networks (g) only approximate true compositionality (f)?
 - Not at all. If g approximates f , then f also approximates g ;
- Why did we adopt f in the first place?
 - Accumulation of evidence that humans perform combinatorial, recursive generalization;
 - But all that evidence was noisy -- humans too might closely approximate true compositionality.

Deep learning contributing to linguistics

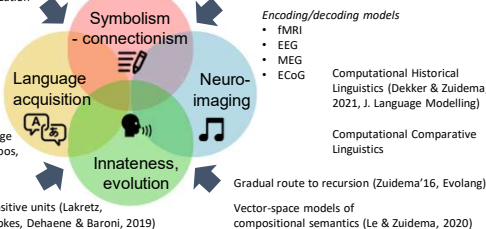
Marr's levels: Symbolic grammars as 'computational level' approximations of an underlying continuous neural reality

Compositional generalization

Existence proof of the learnability of key properties of natural language from data

- Emergence of phonology in grounded models of language learning (Chrupala, Gelderloos, & Alishahi, 2017)

- Emergence of structure sensitive units (Lakretz, Kruszewski, Desbordes, Hupkes, Dehaene & Baroni, 2019)



The principle & the approximation 2/2

- The "principle of compositionality" would then still be a scientific law,
 - more like Gay-Lussac ($P \sim T$) than like the Principle of Conservation of Energy
- Is that disappointing?
 - Yes -- if you are nostalgic for the good old days when formal semantics had the monopoly on modelling sentence meaning
 - No -- if you are satisfied with formal semantics providing an explainable, computational level characterization of the asymptote that neural systems approximate.

Linguistic relevance

- Aren't neural networks just a lower level of description, relevant for neuroscientists but uninterpretable and irrelevant for linguists?
- Much theory and methodology in linguistics is based on the idea that the symbols, rules, relations, and hierarchical are cognitive and neural primitives.
- We tend to think about neural networks as, at best, approximating the symbolic 'truth'. But what if it is the other way around?
- Symbolic models remain extremely useful for characterizing structure in language.
- But if our symbolic models are really approximations of the underlying neural reality, maybe it not so surprising that 'all grammar leak'.
- Multiple models can coexist as they fulfill different roles!

Suggested Readings

- (1) Afra Alishahi, Grzegorz Chrupala, Tal Linzen (2019), Analyzing and Interpreting Neural Networks for NLP: A Report on the First BlackboxNLP Workshop, <https://arxiv.org/abs/1904.04063>
- (2) Yonatan Belinkov, James Glass (2019), Analysis Methods in Neural Language Processing: A Survey, <https://arxiv.org/abs/1812.08951>
- (3) Mario Giulianelli, Jack Harding, Florian Mohnert, Dieuwke Hupkes, Willem Zuidema (2018), Under the Hood: Using Diagnostic Classifiers to Investigate and Improve how Language Models Track Agreement Information, <https://arxiv.org/abs/1808.08079>
- (4) Samira Abnar, Lisa Beinborn, Rochelle Choenni, Willem Zuidema (2019), Blackbox meets blackbox: Representational Similarity and Stability Analysis of Neural Language Models and Brains, <https://arxiv.org/abs/1906.01539>
- (5) The AllenAI NLP Guide has a chapter on Interpretability, with a useful section on Attribution Methods ("Saliency methods") including Integrated Gradients: <https://guide.allennlp.org/interpret>

Conclusions

- Linguists can and should engage with the spectacular progress in deep learning for Natural Language Processing
- Linguistics has much to contribute in trying to open the blackbox of deep learning system: study them as linguistic agents
- Deep learning has much to contribute to linguistics: proofs of concepts (to avoid misunderstandings in classic debates) and concrete tools (to make specific predictions)

zuidema@uva.nl, @wzuidema, <http://projects.illc.uva.nl/LaCo/clclab/>

Collaborators & Contact



Lisa Beinborn



Samira Abnar

<http://projects.illc.uva.nl/LaCo/clclab/>



Mario Giulianelli



Jack Harding



Florian Mohnert



Dieuwke Hupkes

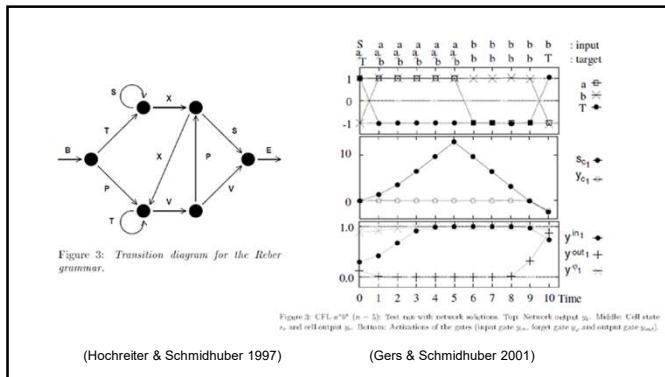


Willem Zuidema
zuidema@uva.nl

Questions?



Extra slides



Sources of error

- True function: $Y = \text{eval}(X)$,
 - with Y the target,
 - X an expression of the form (Expr Op Expr)
 - Expr $\in \{-9, -8, \dots, 8, 9\} \cup \{(\text{Expr Op Expr})\}$
 - Op $\in \{+, -\}$
- Structural error type 1: Not discovering the recursive nature of the problem
 - Because not all lengths are seen at training, OOD generalization will fail catastrophically
- Structural error type 2: errors in learning meaning of brackets
 - Because scope only matters for the minus operation, error will depend on number and depth of minuses, and increase significantly with length of the expression
- Approximation error type 1: errors in learning exact meaning of Expr's
 - Because all operations are linear, total error will grow linearly with length of expression
- Approximation error type 2: errors in learning meaning of Op's
 - Deviations from linearity will yield exponentially increasing error with length of expression

Gradient = partial derivative in a specific point

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}.$$