Bridging the Divide Between Linguistics and NLP: From Vectors to Symbols and Back Again

Tom McCoy June 11, 2025

Yale Department of Linguistics

2

An apparent paradox

Traditional view: Symbolic structure



the doctor by the lawyer saw the artist

An apparent paradox

Traditional view: Symbolic structure



the doctor by the lawyer saw the artist

Neural networks: Vector representations

An apparent paradox

Traditional view: Symbolic structure

the doctor by the lawyer saw the artist

Neural networks: Vector representations

[-1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. 19112, -0.090826, 0.26399]



Vectors

the doctor by the lawyer saw the artist

 $\begin{bmatrix} -1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 \\ 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 \\ 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, \\ 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, \\ -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 \\ , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 \\ 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. \\ 19112, -0.090826, 0.26399]$





the doctor by the lawyer saw the artist

[-1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. 19112, -0.090826, 0.26399]







the doctor by the lawyer saw the artist

[-1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. 19112, -0.090826, 0.26399]







the doctor by the lawyer saw the artist

 $\begin{bmatrix} -1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 \\ 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 \\ 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, \\ 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, \\ -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 \\ , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 \\ 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. \\ 19112, -0.090826, 0.26399]$







the doctor by the lawyer saw the artist

[-1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. 19112, -0.090826, 0.26399]



Girl with a Pearl Earring





the doctor by the lawyer saw the artist

[-1.4792, 1.7206, -0.74454, -0.66112, -0.24197, 0.54995, -0.553 69, -0.35791, -0.27327, 0.56298, 0.48538, 0.1956, 0.76609, 0.00 14633, -0.49319, -0.25365, -0.60558, 1.3433, -0.47967, 1.0888, 0.93805, 1.1932, -0.03101, -0.29201, -0.4451, -1.252, -0.09721, -1.168, -0.37394, 0.24645, 1.5268, 0.12353, -0.98737, -0.82833 , 0.84111, 0.48287, -0.45142, 0.9825, 1.3721, -0.34363, -0.1575 8, -0.34484, 0.0048065, 0.84408, -1.3145, -0.61091, 0.7188, -0. 19112, -0.090826, 0.26399]



Girl with a Pearl Earring



Jackson Pollock

The surprising success of neural nets

	Google	
٩		Ļ
	Google Search I'm Feeling Lucky	

	17:21 - #∷ ♀ Œ
	Welcome to ChatGPT
	This official app is free, syncs your history across devices, and brings you the latest model improvements from OpenAI.
S V	ChatGPT can be inaccurate ChatGPT may provide inaccurate information about people, places, or facts.
\sim	Don't share sensitive info Anonymized chats may be reviewed by our Al trainers to improve our systems.



A: Vectors that act like symbols

A: Vectors that <u>act like</u> <u>symbols</u>







"Vectors that act like symbols"

- Why not just "vectors"?
- Deep learning already excels in:
 - Language
 - Math
 - Board games
 - ...
- Answer: Despite appearances, deep learning systems already have implicit symbolic structure!

Simplified case

Number-to-word mapping:

 $16540 \rightarrow sixteen five forty$

Simplified case

Number-to-word mapping:

 $16540 \rightarrow sixteen five forty$

Requires symbolic structure: Which numbers and where?

• Number-to-word mapping (reversed):

 $16540 \rightarrow \text{forty five sixteen}$

- Two sub-networks:
 - Encoder
 - Decoder





- RNN
- Transformer
- MLP

Hypothesis: Neural networks implicitly implement symbolic representations

Hypothesis: Neural networks implicitly implement symbolic representations

But how?

Hypothesis: Neural networks implicitly implement symbolic representations

But how?

By constructing Tensor Product Representations

Smolensky (1987)

• Example: Representing a sequence of numbers: 3, 6, 7

• Example: Representing a sequence of numbers: 3, 6, 7


































- Hypothesis: Neural network representations are implicitly structured as Tensor Product Representations
 - Even without being designed to have this structure!

Symbolic structures



Classical painting

Vectors



Jackson Pollock

Symbolic structures



Classical painting

<u>Vectors</u>



Jackson Pollock

Tensor Product Representations



Pointillism

- Hypothesis: Neural network representations are implicitly structured as Tensor Product Representations
 - Even without being designed to have this structure!
- To test the hypothesis: train Tensor Product Representations to approximate the neural network representations

(McCoy, Linzen, Dunbar, Smolensky 2019: ICLR)

(McCoy 2022: Dissertation)



Model being analyzed:



Model being analyzed:





Model being analyzed:





Model being analyzed:

















Very strong approximation: Shows that the representations of these neural networks have implicit symbolic structure!

Results with natural language

- Have also applied this approach to neural networks trained on natural language
 - And have gotten promising results











Modeling language learning

- Representations are one area where vectors and symbols seem mismatched
- Learning is another such area

Modeling language learning

- Neural networks have some attractive properties as models of learning:
 - Trained on naturalistic corpora
 - By the end of training: They capture many aspects of linguistic structure

Problem: Data quantity



Another candidate: Probabilistic models

- Represent hypotheses using symbolic grammars
- E.g., a context-free grammar

 $S \rightarrow NP VP$ $NP \rightarrow Det N$ $VP \rightarrow V NP$ $Det \rightarrow the$ $Det \rightarrow a$

. . .

 Goal: Given a corpus, find the grammar that best describes the corpus

- Goal: Given a corpus, find the grammar that best describes the corpus
- Need some way to decide what it means to "best describe"
 - Answer: Probability

- Goal: Given a corpus, find the grammar that best describes the corpus
- Need some way to decide what it means to "best describe"
 - Answer: Probability

Prominent type of probabilistic model: Bayesian model

- Can learn effectively from small amounts of data
 - Why? Using structured grammars as hypotheses guides the search

- Can learn effectively from small amounts of data
 - Why? Using structured grammars as hypotheses guides the search
- But they are typically intractable in naturalistic settings

- Can learn effectively from small amounts of data
 - Why? Using structured grammars as hypotheses guides the search
- But they are typically intractable in naturalistic settings
 - Hard to find a grammar that captures all the complexities of the data
| Probabilistic models
(e.g., Bayesian) | Neural networks |
|--|-----------------|
| | |
| | |
| | |
| | |

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	
Strong inductive biases	

	Probabilistic models (e.g., Bayesian)	Neural networks
	Strong representational commitments: Symbols	
	Strong inductive biases	
The factors that guide generalization		

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	
Strong inductive biases	
Effective generalization from limited data	

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	
Strong inductive biases	
Effective generalization from limited data	
Struggle to tractably handle natural data	

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	Representational flexibility: Vectors
Strong inductive biases	
Effective generalization from limited data	
Struggle to tractably handle natural data	

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	Representational flexibility: Vectors
Strong inductive biases	
Effective generalization from limited data	
Struggle to tractably handle natural data	Can handle complex, natural data

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	Representational flexibility: Vectors
Strong inductive biases	Weak inductive biases
Effective generalization from limited data	
Struggle to tractably handle natural data	Can handle complex, natural data

Probabilistic models (e.g., Bayesian)	Neural networks
Strong representational commitments: Symbols	Representational flexibility: Vectors
Strong inductive biases	Weak inductive biases
Effective generalization from limited data	Require lots of data
Struggle to tractably handle natural data	Can handle complex, natural data

Proposal: Inductive bias distillation

- Distill Bayesian inductive biases into a neural network
- Neural network → flexible
- Strong inductive bias \rightarrow rapid learning

Sounds nice...but how can we actually do it?

Modeling rapid language learning by distilling Bayesian priors into artificial neural networks

R. Thomas McCoy^{1*} and Thomas L. Griffiths^{2,1}



Universal linguistic inductive biases via meta-learning

R. Thomas McCoy,¹ Erin Grant,² Paul Smolensky,^{3,1} Thomas L. Griffiths,⁴ and Tal Linzen¹











$$p(h|d) = \frac{p(d|h)p(h)}{p(d)}$$
$$h = \text{plus}(D)$$
$$h = \text{concat}($$
or(A, C),
plus(A),
 Σ ,
or(ε , B))

Probabilistic model







One model for the learning of language

Yuan Yang^a and Steven T. Piantadosi^{b,1}

^aCollege of Computing, Georgia Institute of Technology, Atlanta, GA 30332; and ^bDepartment of Psychology, Hele of California, Berkeley, CA 94720

- plus(concatenate(A, B))
 - AB
 - ABAB
 - ABABAB
 - . . .

Inspired by the structure of natural language

- Inspired by the structure of natural language
- Consider plus(concatenate(A, B)): AB, ABAB, ABABAB, …

- Inspired by the structure of natural language
- Consider plus(concatenate(A, B)): AB, ABAB, ABABAB, …
 - A → preposition
 - $B \rightarrow \underline{\text{noun phrase}}$

- Inspired by the structure of natural language
- Consider plus(concatenate(A, B)): AB, ABAB, ABABAB, …
 - A → preposition
 - $B \rightarrow \underline{\text{noun phrase}}$
 - on the table by the door in the kitchen

Probabilistic model

 Probabilistically combine primitives into formal languages









Sampling formal languages

- plus(A)
- D
- concat(A, B, A)
- concat(plus(or(A, concat(C, A, C)), or(plus(D), A))

99

- or(A, A)
- B



100







Learning to learn





• Approach:

- Approach:
 - Show the model many languages
 - Giving it linguistic inductive biases (which types of languages are likely/unlikely?)

- Approach:
 - Show the model many languages
 - Giving it linguistic inductive biases (which types of languages are likely/unlikely?)
 - By controlling the languages, we control the model's inductive biases

- Approach:
 - Show the model many languages
 - Giving it linguistic inductive biases (which types of languages are likely/unlikely?)
 - By controlling the languages, we control the model's inductive biases
- Variant that we use: MAML (Finn et al. 2017)

• Result: A prior-trained neural network

Trained to have a particular prior
Meta-learning

- Result: A prior-trained neural network
 - Trained to have a particular prior
- Different from pre-trained:
 - Trained to learn aspects of the intended task learning, not metalearning
 - The prior is an indirect byproduct, not a direct target









Time

- Bayesian model: 1 minute to 7 days
 - Not feasible to train on naturalistic data

Time

- Bayesian model: 1 minute to 7 days
 - Not feasible to train on naturalistic data
- Prior-trained neural network: 10 milliseconds to 3 minutes
 - Can train on naturalistic data!

Training on English

- Child-directed speech
- 8 million words

Training on English: Results



- The prior included the Kleene **plus** primitive
 - plus(B) = {B, BB, BBB, ...}

- The prior included the Kleene **plus** primitive
 - plus(B) = {B, BB, BBB, ...}
- Does the prior-trained model handle recursion well?

- ✓ 1. The book sitting on the table is blue.
- \times 2. The book sits on the table is blue.

- ✓ 1. The book sitting on the table is blue.
- \mathbf{X} 2. The book sits on the table is blue.

- ✓ 1. The book sitting <u>on the table</u> is blue.
- \times 2. The book sits <u>on the table</u> is blue.

- 1. The book sitting <u>on the table</u> in the kitchen is blue.
- \times 2. The book sits <u>on the table</u> in the kitchen is blue.

- 1. The book sitting <u>on the table</u> in the kitchen by <u>the door</u> is blue.
- X 2. The book sits <u>on the table in the kitchen by the</u> <u>door</u> is blue.

Recursion: Results



Recursion: Results



Learning formal languages from few examples		
Learning aspects of English from naturalistic data		

	Bayesian model	
Learning formal languages from few examples		
Learning aspects of English from naturalistic data	×	

	Bayesian model	Standard neural network	
Learning formal languages from few examples		×	
Learning aspects of English from naturalistic data	×		

	Bayesian model	Standard neural network	Prior-trained neural network
Learning formal languages from few examples		×	
Learning aspects of English from naturalistic data	×		

Concept learning

Distilling Symbolic Priors for Concept Learning into Neural Networks

Ioana Marinescu,¹ R. Thomas McCoy,² Thomas L. Griffiths^{1,3}

ioanam@princeton.edu, tom.mccoy@yale.edu, tomg@princeton.edu ¹Department of Computer Science, Princeton University ²Department of Linguistics, Yale University ³Department of Psychology, Princeton University







Concept learning



Conclusion

Q: What type of system should we use to represent language?

A: Vectors that act like symbols

Q: What type of system should we use to represent language?

A: Vectors that <u>act like</u> <u>symbols</u> Q: What type of system should we use to represent language?

A: Vectors that <u>act like</u> <u>symbols</u> <u>Compositional representations</u>

Tensor Product Representations





Symbolic structures



Classical painting

<u>Vectors</u>



Jackson Pollock

<u>Tensor Product Representations</u> and <u>Neural Networks with Meta-Learning</u>





A symbolic victory?

- The framing I've used might seem like a resounding victory for the symbolic perspective
 - I.e., the role of vectors is just to implement symbols!
- But the truth is probably more complex
- The fact that the symbols are implemented in vectors might have important consequences
 - Ability to deviate from pure symbolic structure to handle exceptions, context, etc.
 - Flexibility for learning

How to get there?

- On one hand, neural networks naturally develop compositional structure on their own
 - So maybe we don't need to consciously have symbols in mind when developing systems
- On the other hand: Incorporating soft versions of symbols might be useful as an inductive bias
 - Architectures with compositional structure
 - Training paradigms that encourage symbolic processing

Thank you!

Collaborators:



Ewan Dunbar





Grant

Tom Griffiths







Tal Linzen

loana Marinescu

Paul Smolensky



- Funding: NSF GRFP #1746891, NSF SPRF #2204152
- You!

