

Finding Syntax with Structural Probes

Jun Li



Representation learned by NN

The chef who ran to the store was out of food.

Representation learned by NN

The chef who ran to the store was out of food.

The	chef	who	ran	to	the	store	was	out	of	food
$\begin{bmatrix} .4 \\ -.2 \\ .3 \end{bmatrix}$	$\begin{bmatrix} .1 \\ .9 \\ -.2 \end{bmatrix}$	$\begin{bmatrix} .3 \\ -.4 \\ .2 \end{bmatrix}$	$\begin{bmatrix} .7 \\ -.4 \\ 0 \end{bmatrix}$	$\begin{bmatrix} .4 \\ 0 \\ -.5 \end{bmatrix}$	$\begin{bmatrix} .1 \\ -.6 \\ .2 \end{bmatrix}$	$\begin{bmatrix} .3 \\ .1 \\ -.6 \end{bmatrix}$	$\begin{bmatrix} .1 \\ .9 \\ -.8 \end{bmatrix}$	$\begin{bmatrix} .3 \\ .1 \\ .8 \end{bmatrix}$	$\begin{bmatrix} -.8 \\ .3 \\ -.6 \end{bmatrix}$	$\begin{bmatrix} 0 \\ .7 \\ -.9 \end{bmatrix}$

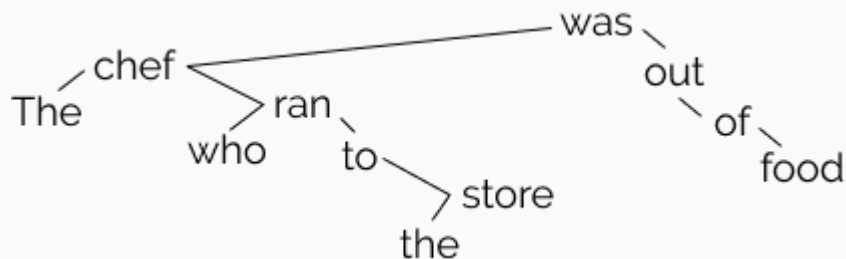
Representation learned by NN



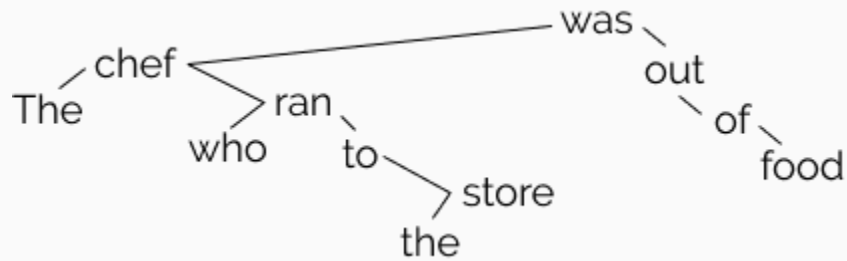
The chef who ran to the store was out of food.

Learned by `human`

The chef who ran to the store was out of food.



Are they *reconcilable*?



Previous works

- We can use these vectors to predict depth of the parse tree
- But can we do better?
- Can we recover a parse tree?

The	chef	who	ran	to	the	store	was	out	of	food
$\begin{bmatrix} .4 \\ -.2 \\ .3 \end{bmatrix}$	$\begin{bmatrix} .1 \\ .9 \\ -.2 \end{bmatrix}$	$\begin{bmatrix} .3 \\ -.4 \\ .2 \end{bmatrix}$	$\begin{bmatrix} .7 \\ -.4 \\ 0 \end{bmatrix}$	$\begin{bmatrix} .4 \\ 0 \\ -.5 \end{bmatrix}$	$\begin{bmatrix} .1 \\ -.6 \\ .2 \end{bmatrix}$	$\begin{bmatrix} .3 \\ .1 \\ -.6 \end{bmatrix}$	$\begin{bmatrix} .1 \\ .9 \\ -.8 \end{bmatrix}$	$\begin{bmatrix} .3 \\ .1 \\ .8 \end{bmatrix}$	$\begin{bmatrix} -.8 \\ .3 \\ -.6 \end{bmatrix}$	$\begin{bmatrix} 0 \\ .7 \\ -.9 \end{bmatrix}$

Question

Q: What do `vector space` and `parse tree space` have in common?

A: Geometry between words

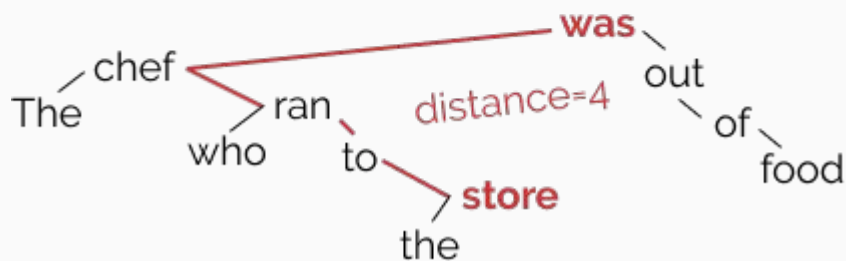
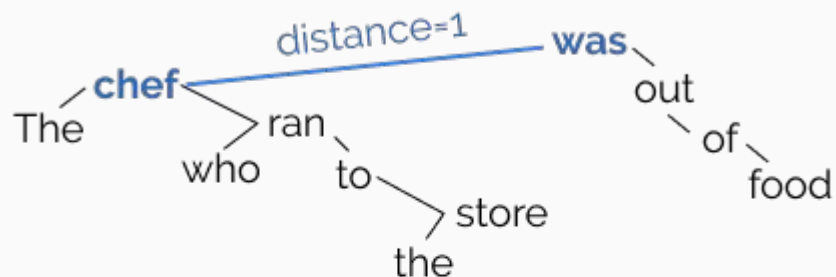
Vector vs. Parse Tree

- Words are represented as vectors
- Words are represented as words

Vector vs. Parse Tree

- Words are represented as vectors
 - Distance between words can be obtained by calculating distance between vectors
- Words are represented as words
 - Distance between words is the length of the tree path

Distance on tree



Vector vs. Parse Tree

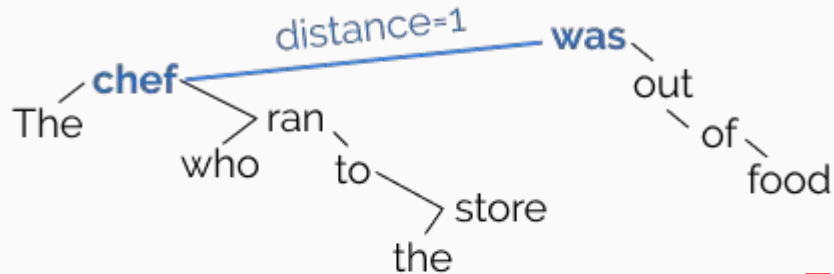
- Words are represented as vectors
- Distance between words can be obtained by calculating distance between vectors
- Norm of a vector represents how far it is from the origin

- Words are represented as words
- Distance between words is the length of the tree path
- Depth of a node represents how far it is from the root node

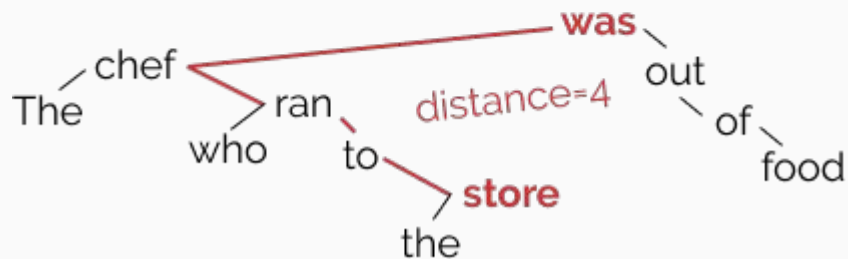
Why do we mention this?

king - man = queen - woman

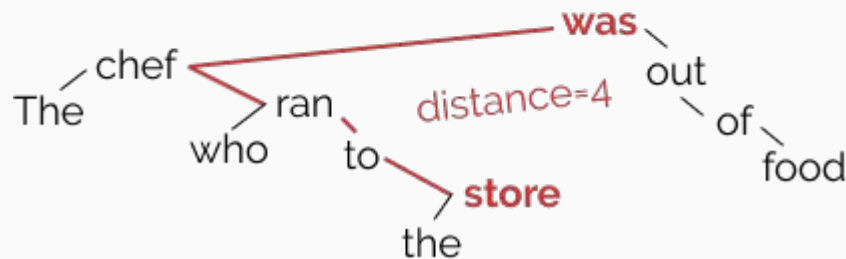
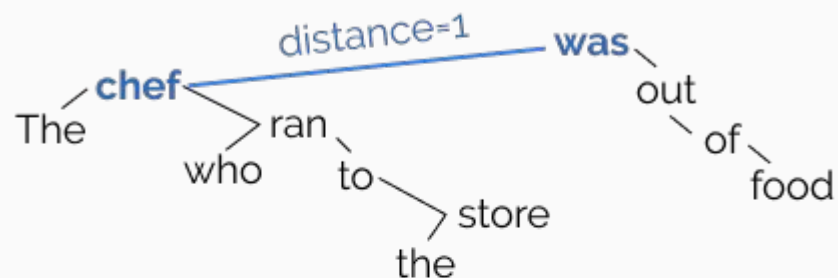
Hypothesis



Tree distance should be encoded in vectors!



Unfortunately...

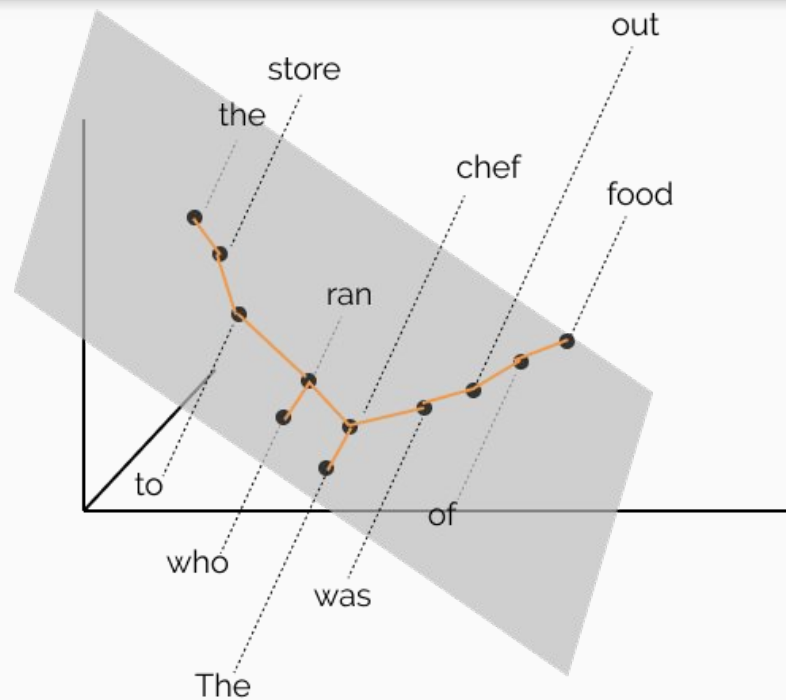


Unfortunately...



Probably the vectors encode too many information, not just syntax

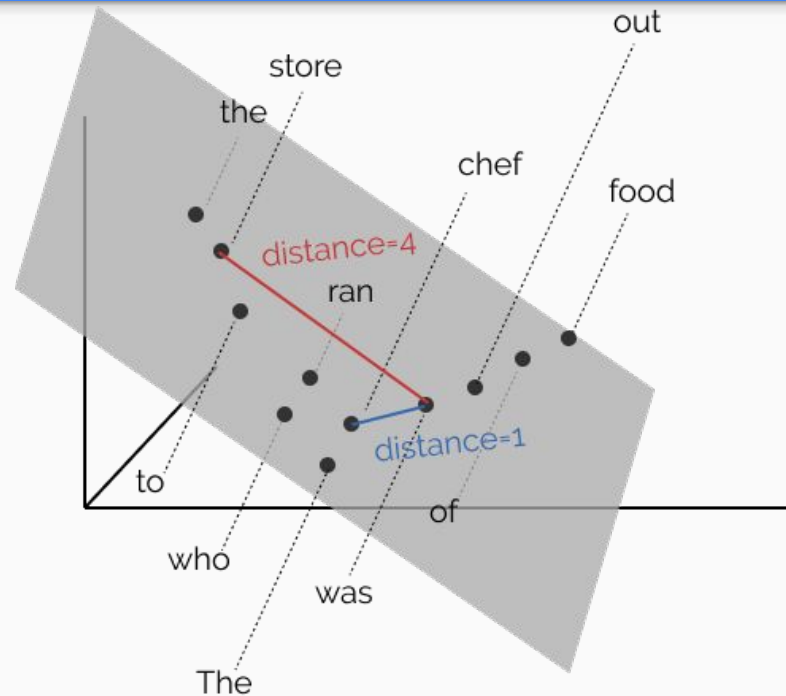
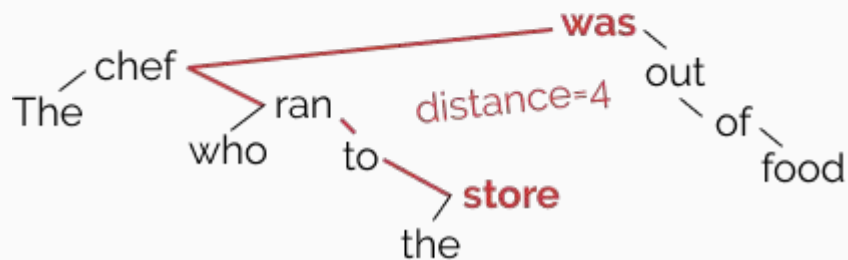
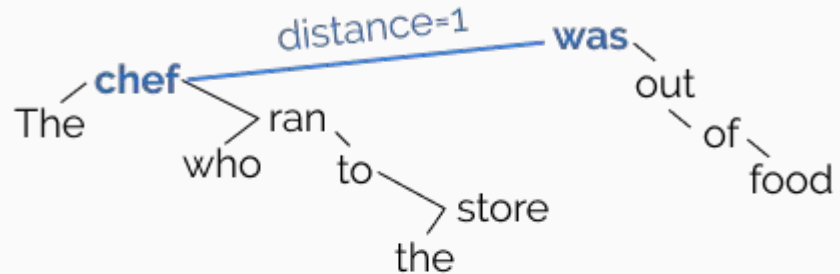
Maybe...from a different perspective?



Relaxed version

Tree distance should be encoded in vectors(after a linear transformation)!

We hope



Formal definition

We can define a new distance:

$$d_B(\mathbf{h}_i^\ell, \mathbf{h}_j^\ell)^2 = (B(\mathbf{h}_i^\ell - \mathbf{h}_j^\ell))^T (B(\mathbf{h}_i^\ell - \mathbf{h}_j^\ell)) \quad (1)$$

$$d(\mathbf{h}_i, \mathbf{h}_j) = (\mathbf{h}_i - \mathbf{h}_j)^\top (\mathbf{h}_i - \mathbf{h}_j)$$

And the objective is:

$$\min_B \sum_{\ell} \frac{1}{|s^\ell|^2} \sum_{i,j} |d_{T^\ell}(w_i^\ell, w_j^\ell) - d_B(\mathbf{h}_i^\ell, \mathbf{h}_j^\ell)^2|$$

Finding syntax with structural probes

What is a probe by the way?

A supervised model for finding information in a representation

Observation Evidence: Whether a given desired behaviour is observed(S-V)

Constructive Evidence: The model may encode the phenomenon of interest, and we train a probe supervisely to recover it

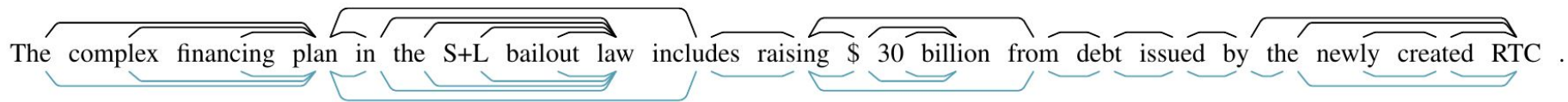
Experiments

Method	Distance		Depth	
	UUAS	DSpr.	Root%	NSpr.
LINEAR	48.9	0.58	2.9	0.27
ELMo0	26.8	0.44	54.3	0.56
DECAY0	51.7	0.61	54.3	0.56
PROJ0	59.8	0.73	64.4	0.75
ELMo1	77.0	0.83	86.5	0.87
BERTBASE7	79.8	0.85	88.0	0.87
BERTLARGE15	82.5	0.86	89.4	0.88
BERTLARGE16	81.7	0.87	90.1	0.89

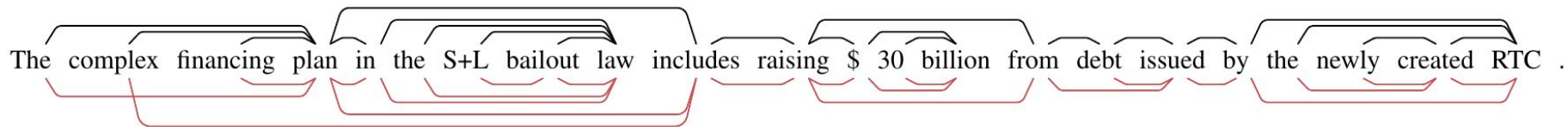
Table 1: Results of structural probes on the PTB WSJ test set; baselines in the top half, models hypothesized to encode syntax in the bottom half. For the distance probes, we show the Undirected Unlabeled Attachment Score (UUAS) as well as the average Spearman correlation of true to predicted distances, DSpr. For the norm probes, we show the root prediction accuracy and the average Spearman correlation of true to predicted norms, NSpr.

Reconstructed parse trees

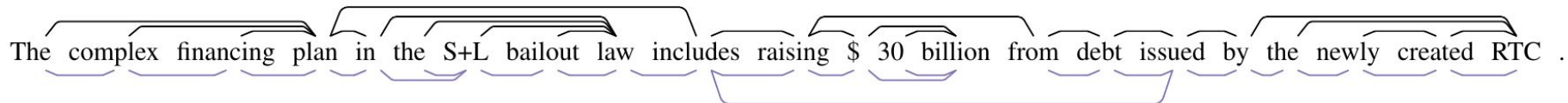
BERTlarge16



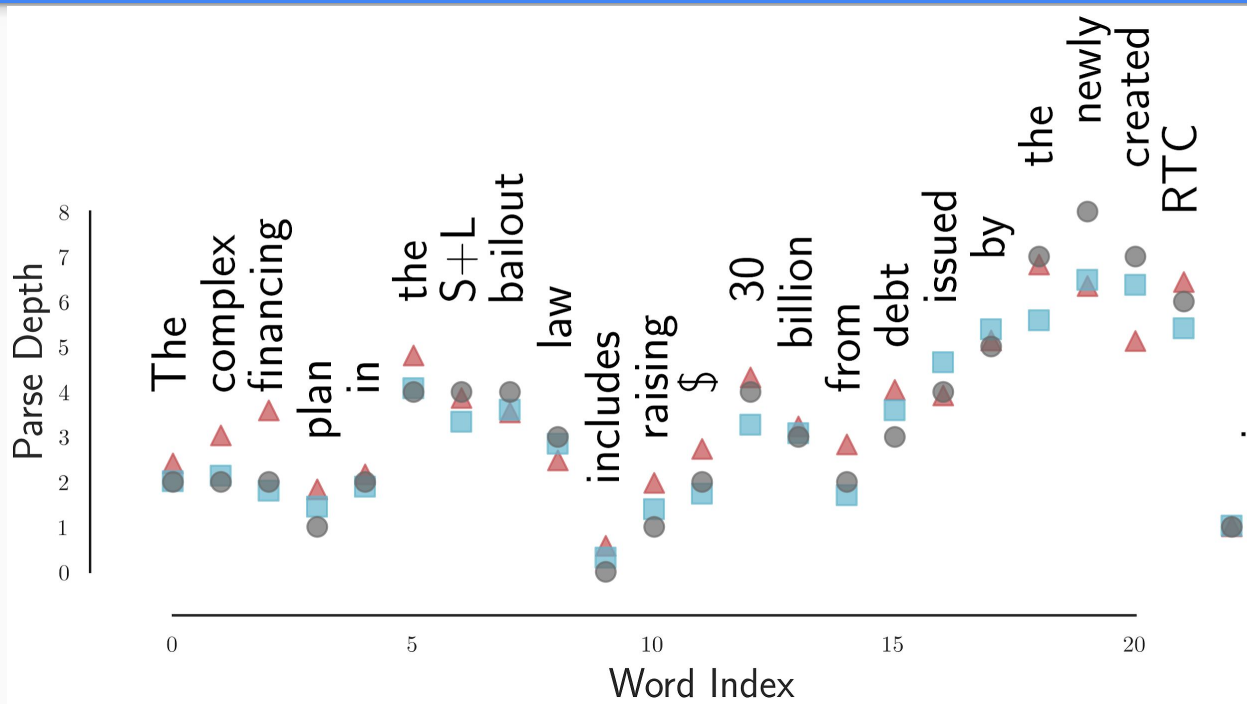
ELMo1



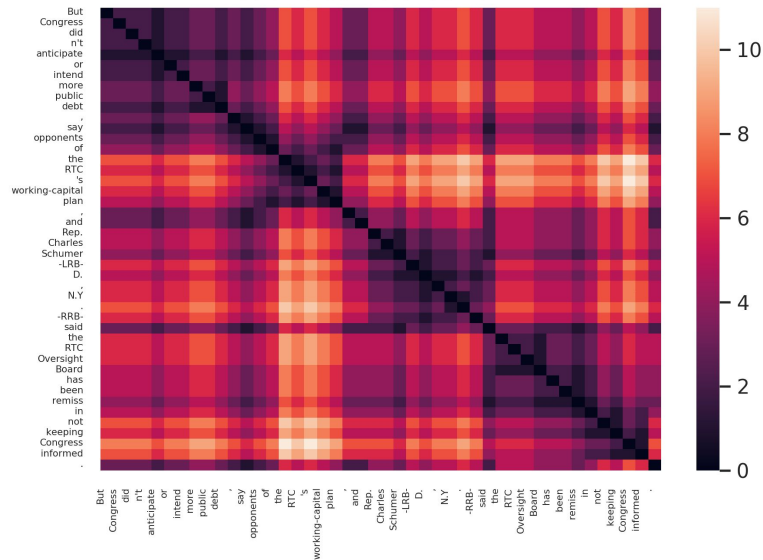
Proj0



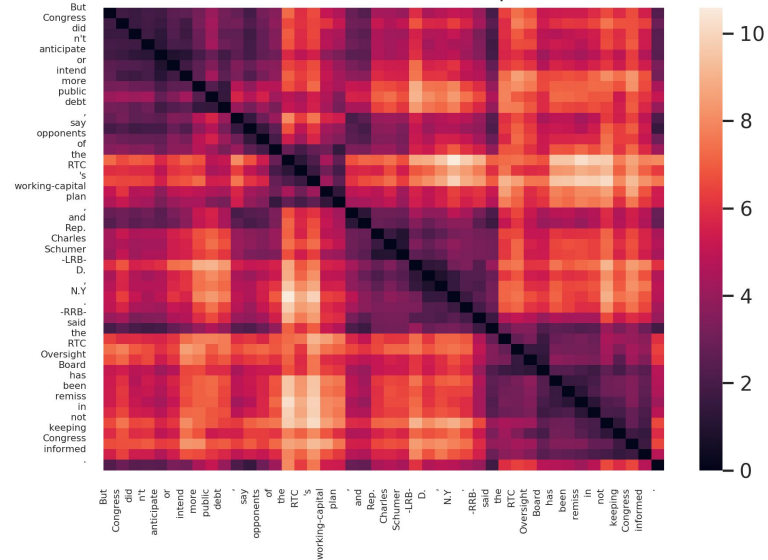
Parse diff



Gold Parse Distance Matrix



Predicted Parse Distance (squared)



Different hidden layers

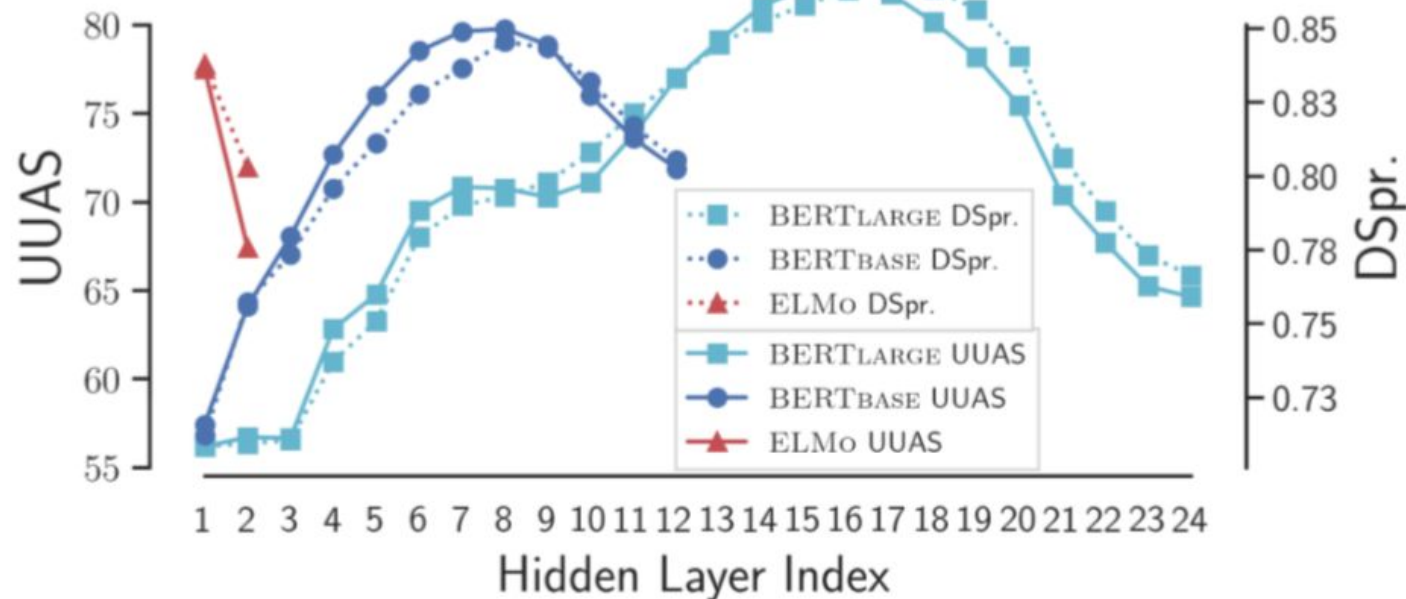


Figure 1: Parse distance UUAS and distance Spearman correlation across the BERT and ELMo model layers.

Rank of the linear transformation

Intuitively, larger k means a more expressive probing model, and a larger fraction of the representational capacity of the model being devoted to syntax.

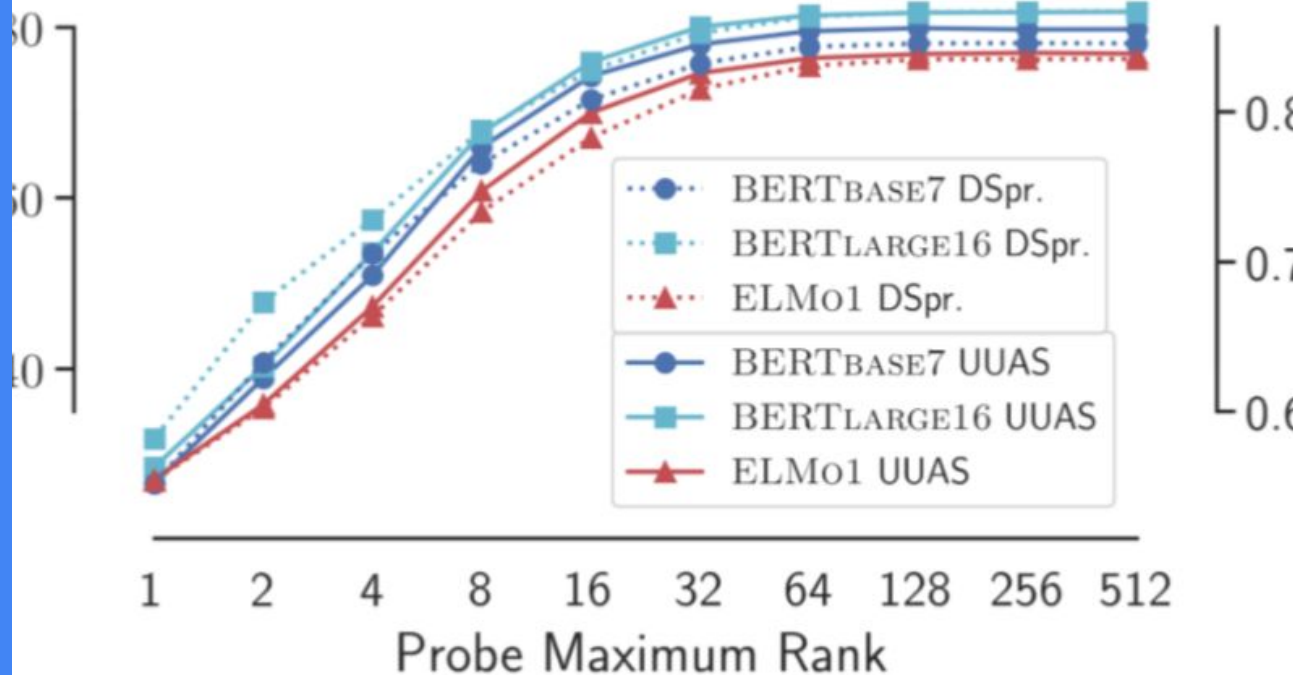


Figure 5: Parse distance tree reconstruction accuracy with linear transformation is constrained to varying maximum dimensionality.

Thanks
