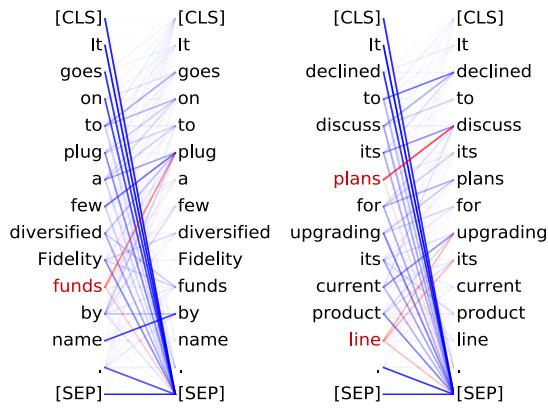


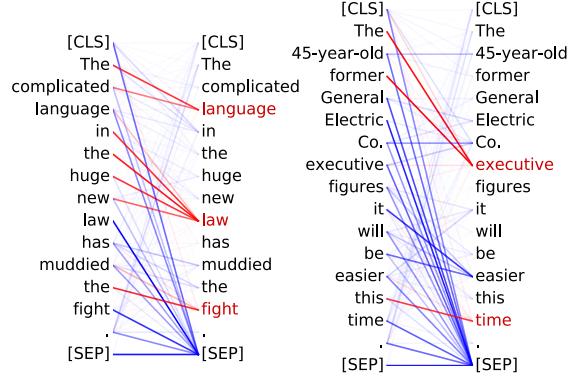
### Head 8-10

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



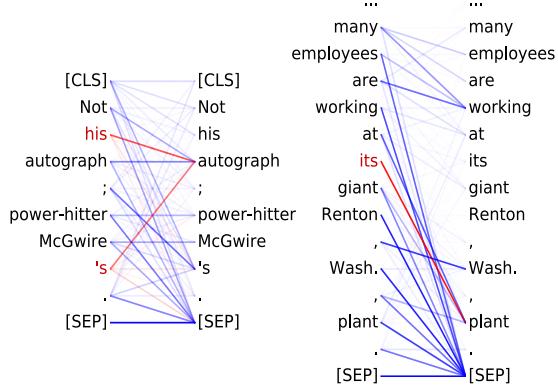
### Head 8-11

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



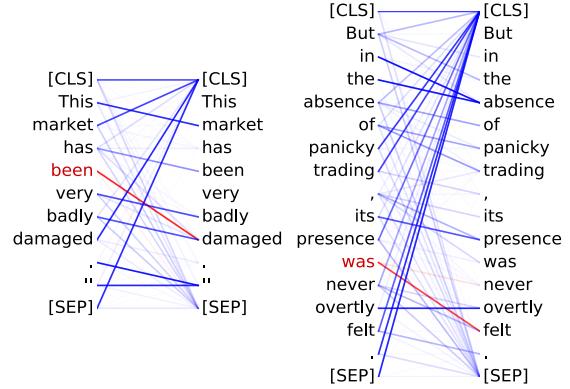
### Head 7-6

- Possessive pronouns and apostrophes attend to the head of the corresponding NP
- 80.5% accuracy at the `poss` relation



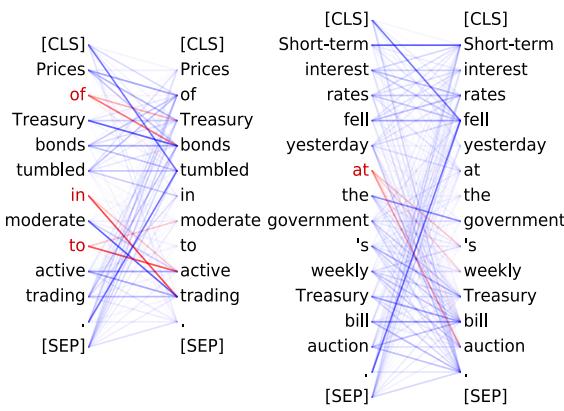
### Head 4-10

- Passive auxiliary verbs attend to the verb they modify
- 82.5% accuracy at the `auxpass` relation



### Head 9-6

- Prepositions attend to their objects
- 76.3% accuracy at the `pobj` 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

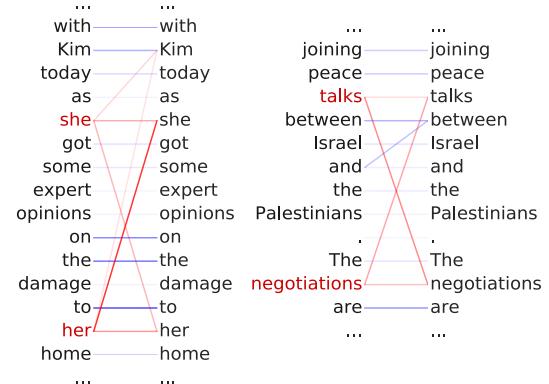


Figure 5: BERT attention heads that correspond to linguistic phenomena. In the example attention maps, the darkness of a line indicates the strength of the attention weight. All attention to/from red words is colored red; these colors are there to highlight certain parts of the attention heads' behaviors. For Head 9-6, we don't show attention to [SEP] for clarity. Despite not being explicitly trained on these tasks, BERT's attention heads perform remarkably well, illustrating how syntax-sensitive behavior can emerge from self-supervised training alone.