

LSTMs Exploit Linguistic Attributes of Data

Nelson F. Liu^{♠♥}, Omer Levy[♠], Roy Schwartz^{♠♣}, Chenhao Tan[◇], Noah A. Smith^{♠♣}

♠ Paul G. Allen School of Computer Science & Engineering; ♥ Department of Linguistics, University of Washington;

♣ Allen Institute for Artificial Intelligence; ◇ Department of Computer Science, University of Colorado Boulder

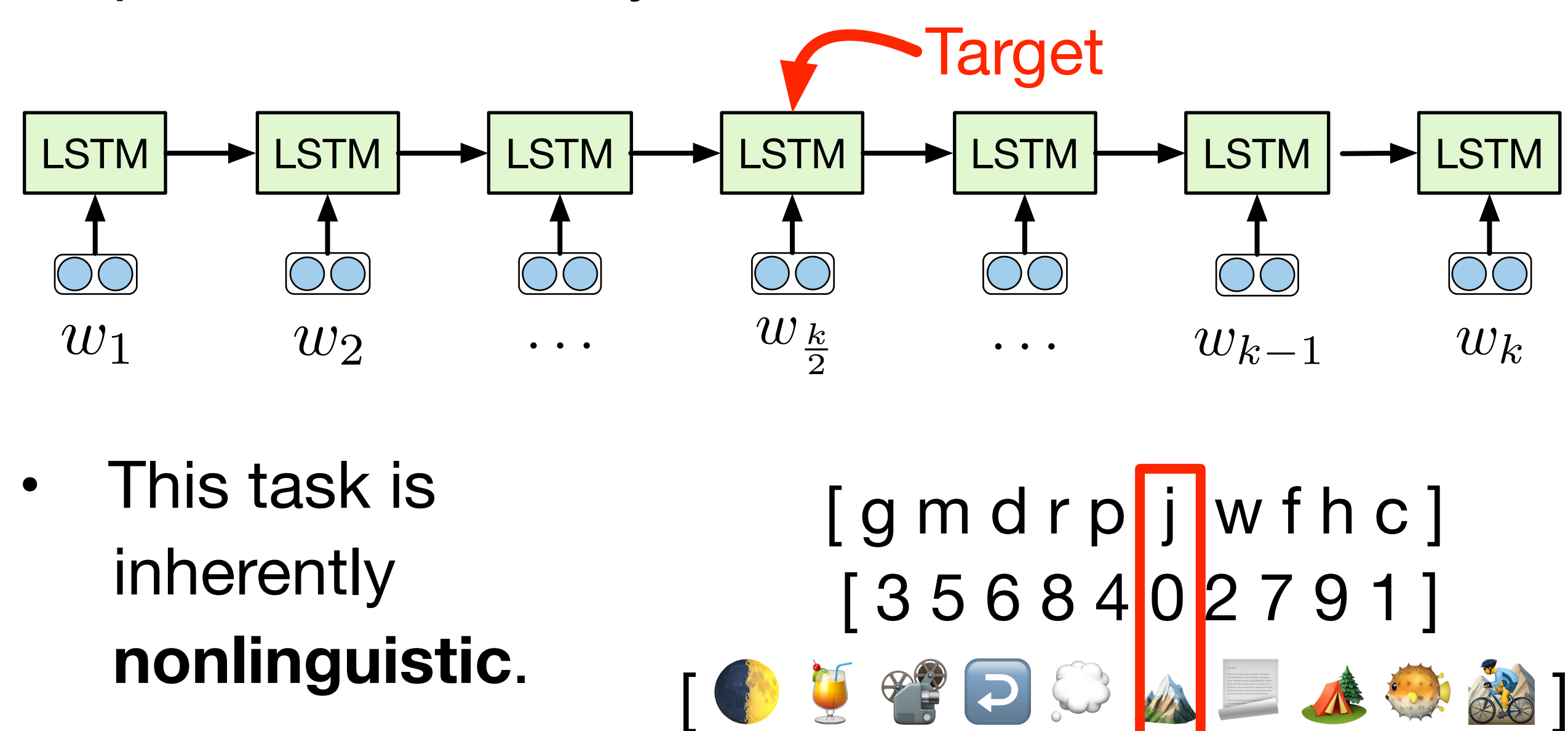


TL;DR

- Data with linguistic attributes helps LSTMs learn a non-linguistic memorization task.
- To solve the task, LSTMs use individual neurons to count timesteps.
- We hypothesize that LSTMs pick up on the patterns and structure in linguistic data and use them as additional noisy training signal.

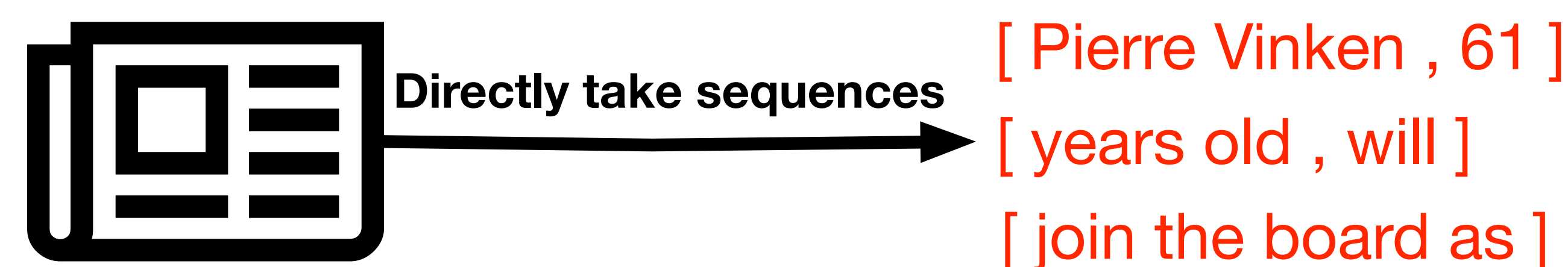
Testbed Memorization Task

- Given a **constant-length** sequence of tokens, predict the identity of the middle token seen.

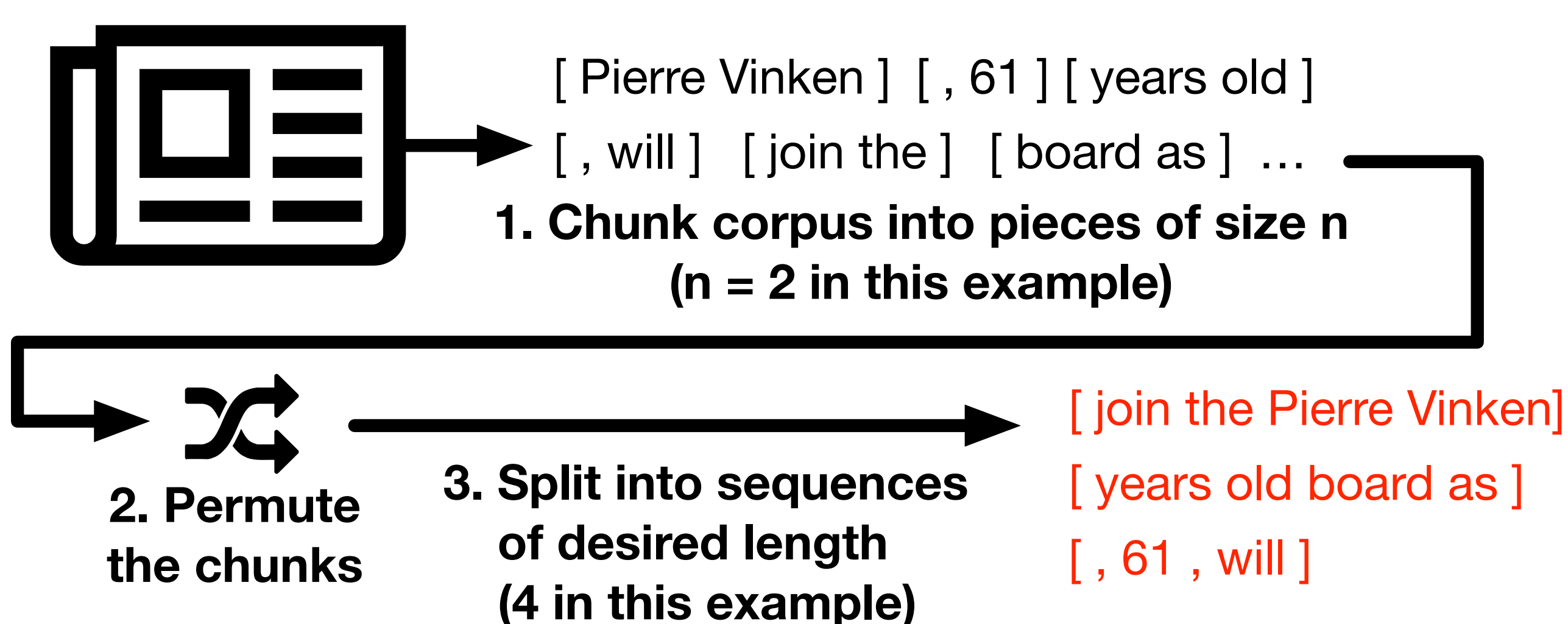


Training Datasets with Various Linguistic Attributes

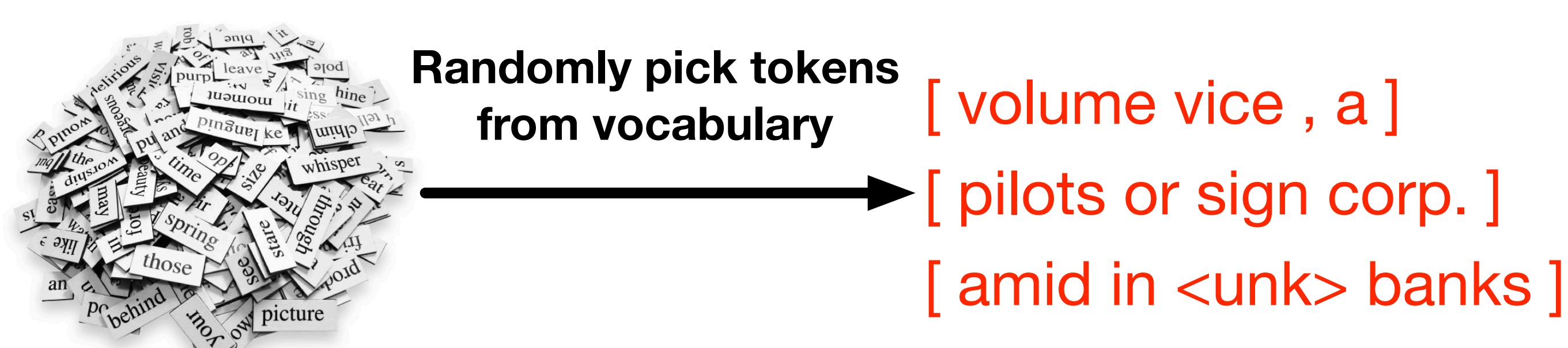
1. Language setting



2. n-gram setting



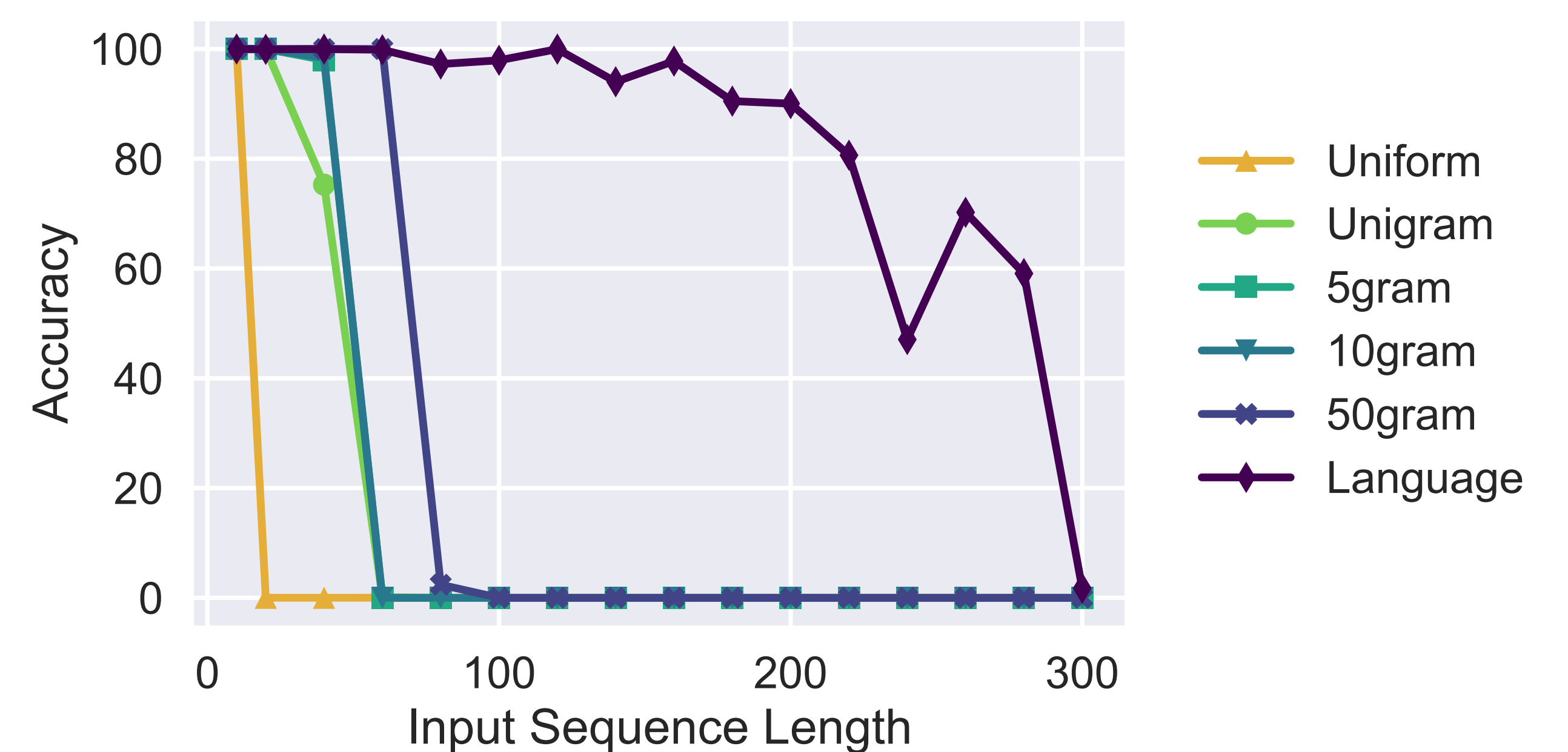
3. Uniform setting



Experiments

- Test data: uniform distribution over the 100 rarest words in the PTB.
- Ensures that models truly generalize and are not just using training data-specific features.

Models trained on data with linguistic features generalize better



What happens if we add more hidden units?

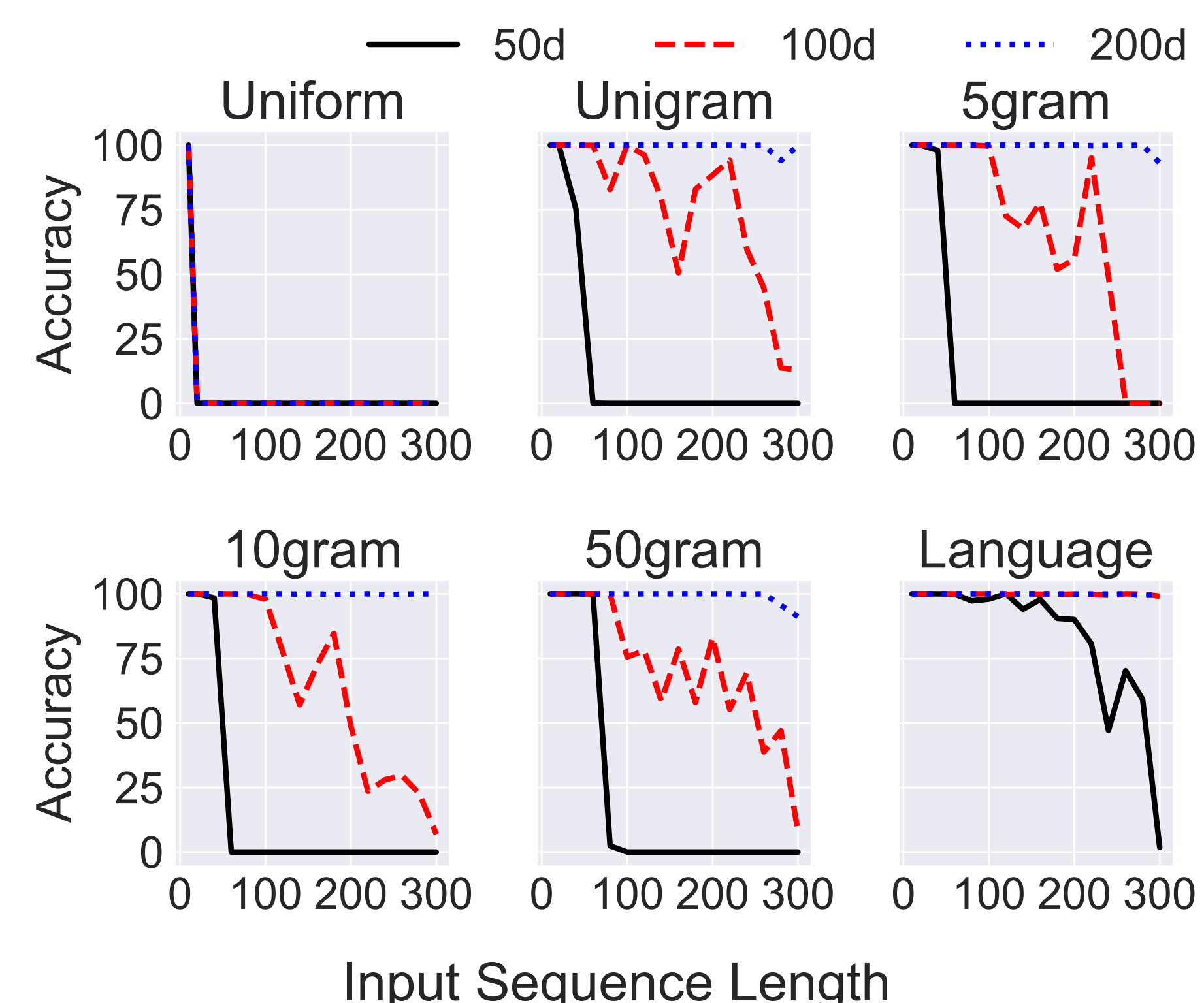
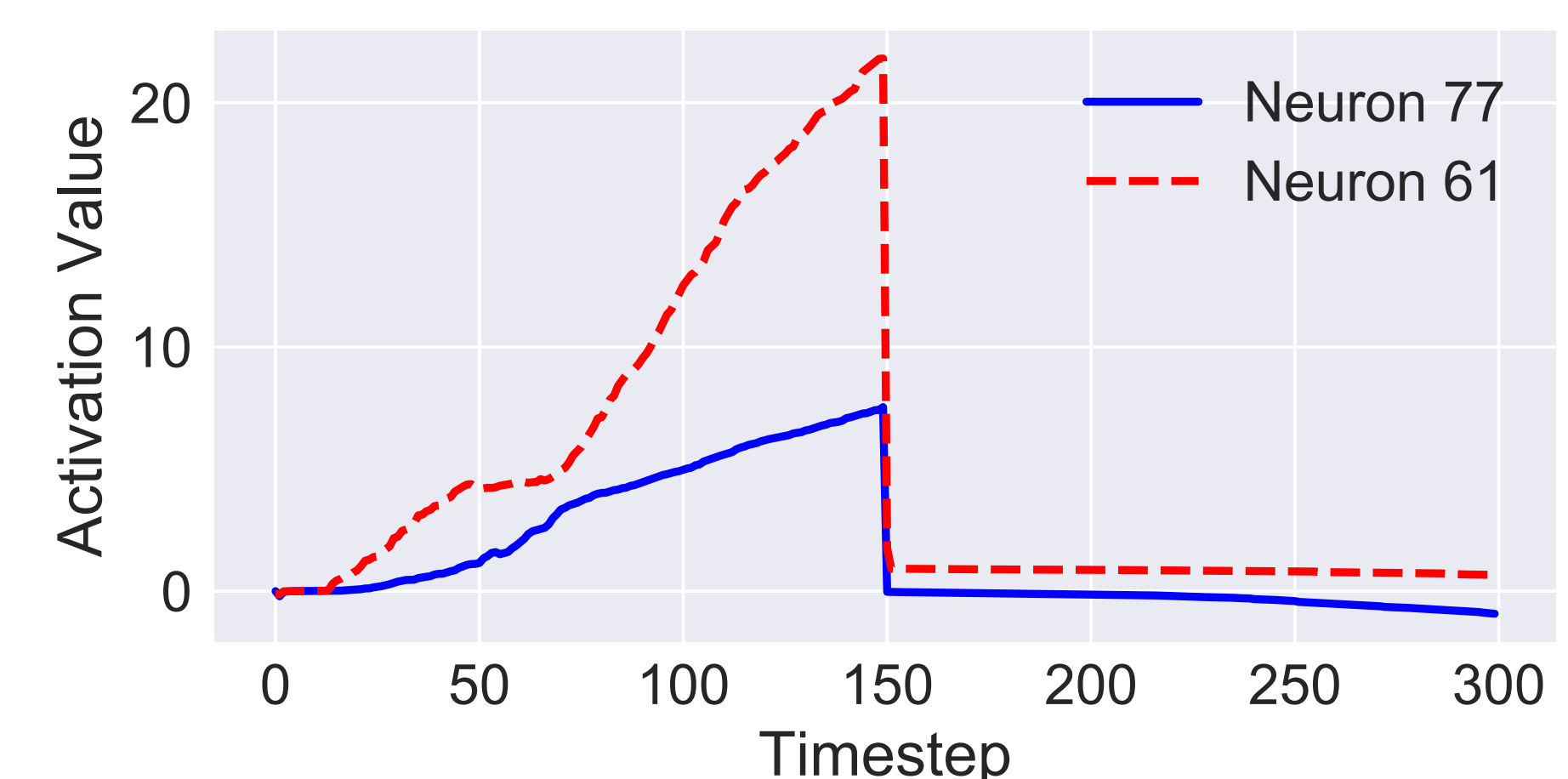


Figure 2: Models trained on **Uniform** data do poorly, even with more hidden units.

Further Analysis

- We further study an LSTM with 100 hidden units trained on **Language**, where train and test sequences are of length 300.

To solve the task, RNNs learn to count



The RNN exploits linguistic features to bootstrap itself early in training and learns to generalize later.

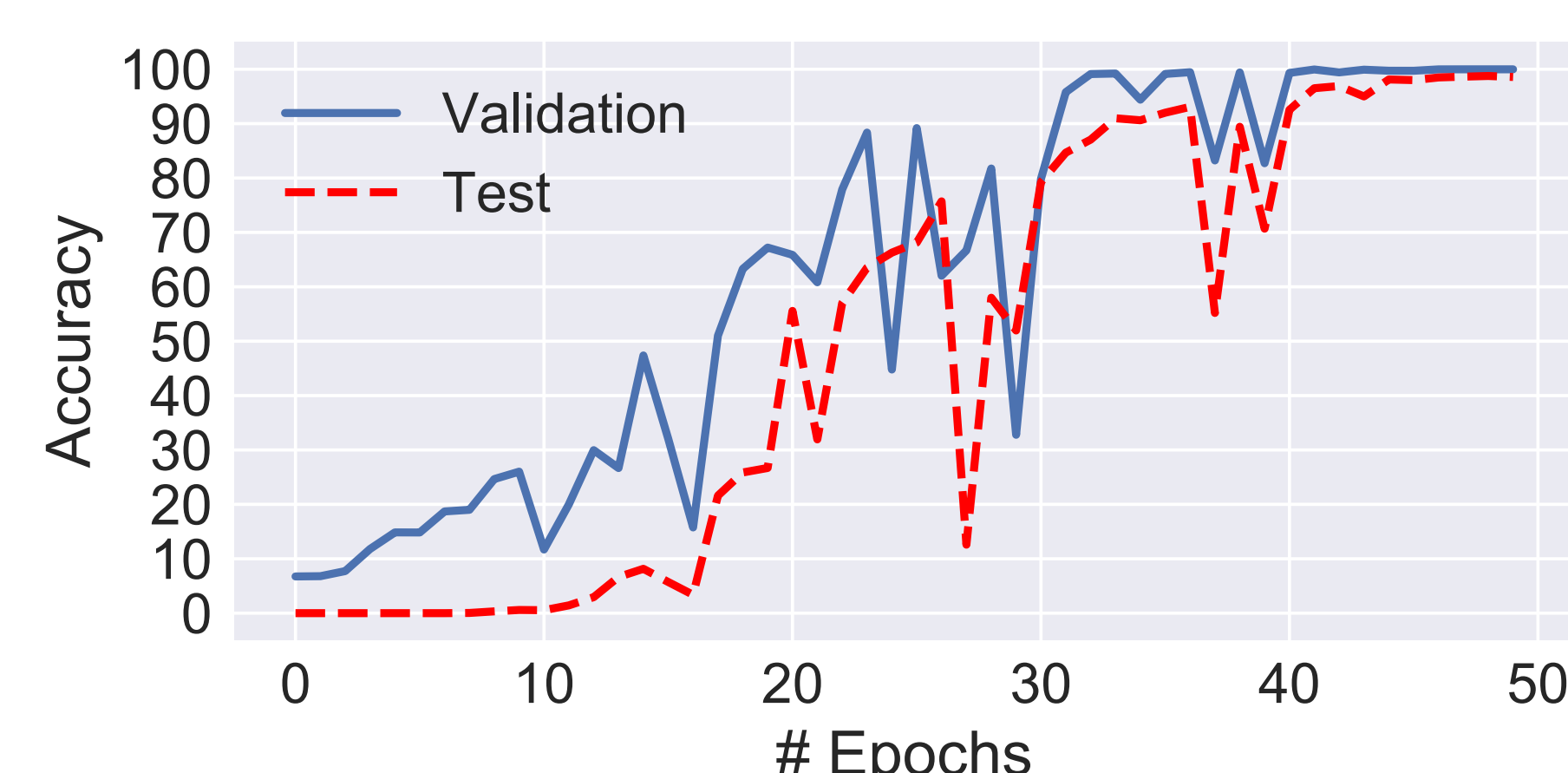


Figure 4: Validation set has the same distribution as train. We see validation accuracy improves faster than test at early epochs (exploiting linguistic features), but the two curves move in unison in later epochs (true memorization).