

Natural Language Processing

Approaches to sequence labelling

Marco Kuhlmann

Department of Computer and Information Science

Approaches to sequence labelling

Local search

sequence of independent
classification problems

expressive feature models
(RNNs, attention)

no need for specialised
algorithms; fast

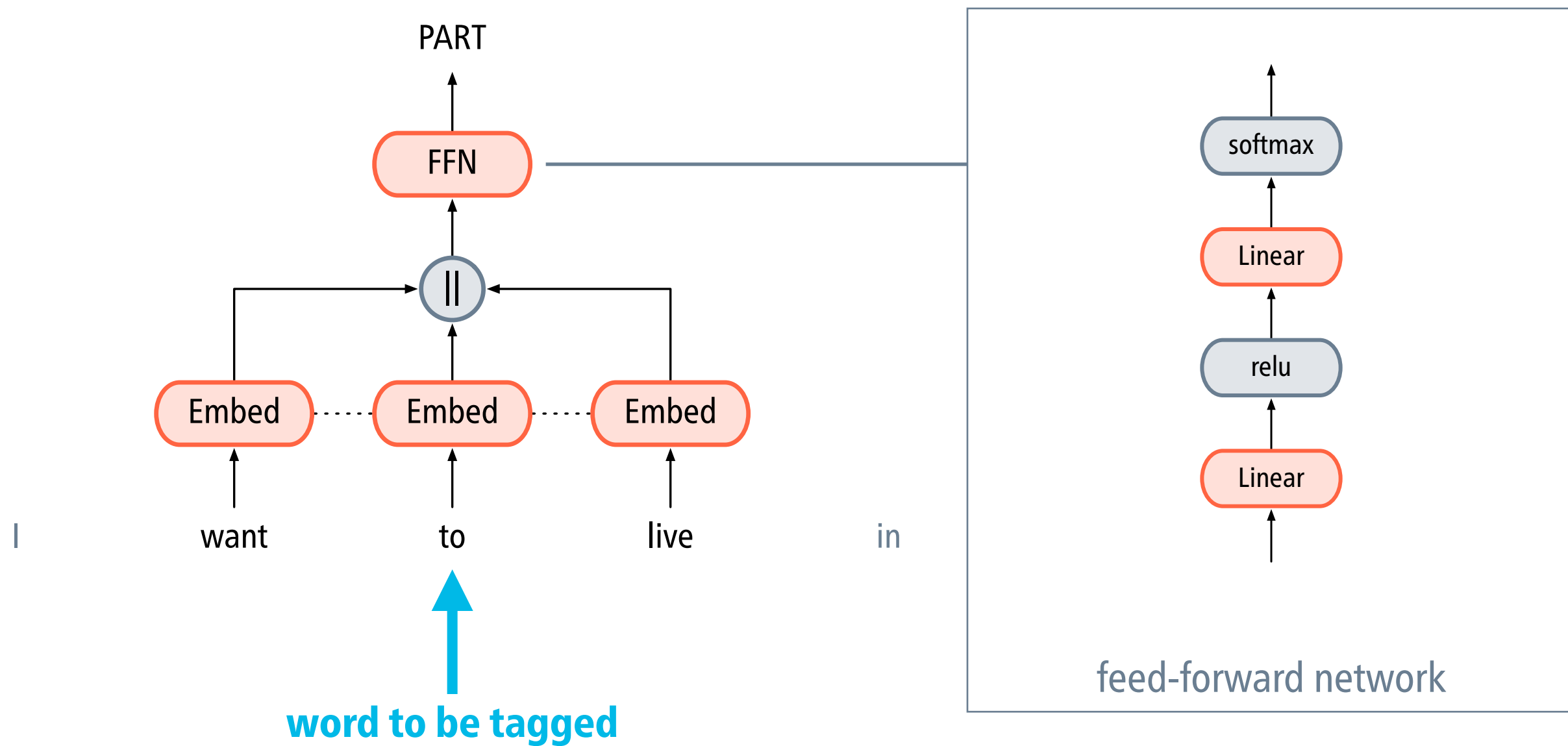
Global search

combinatorial optimisation over
the full search space

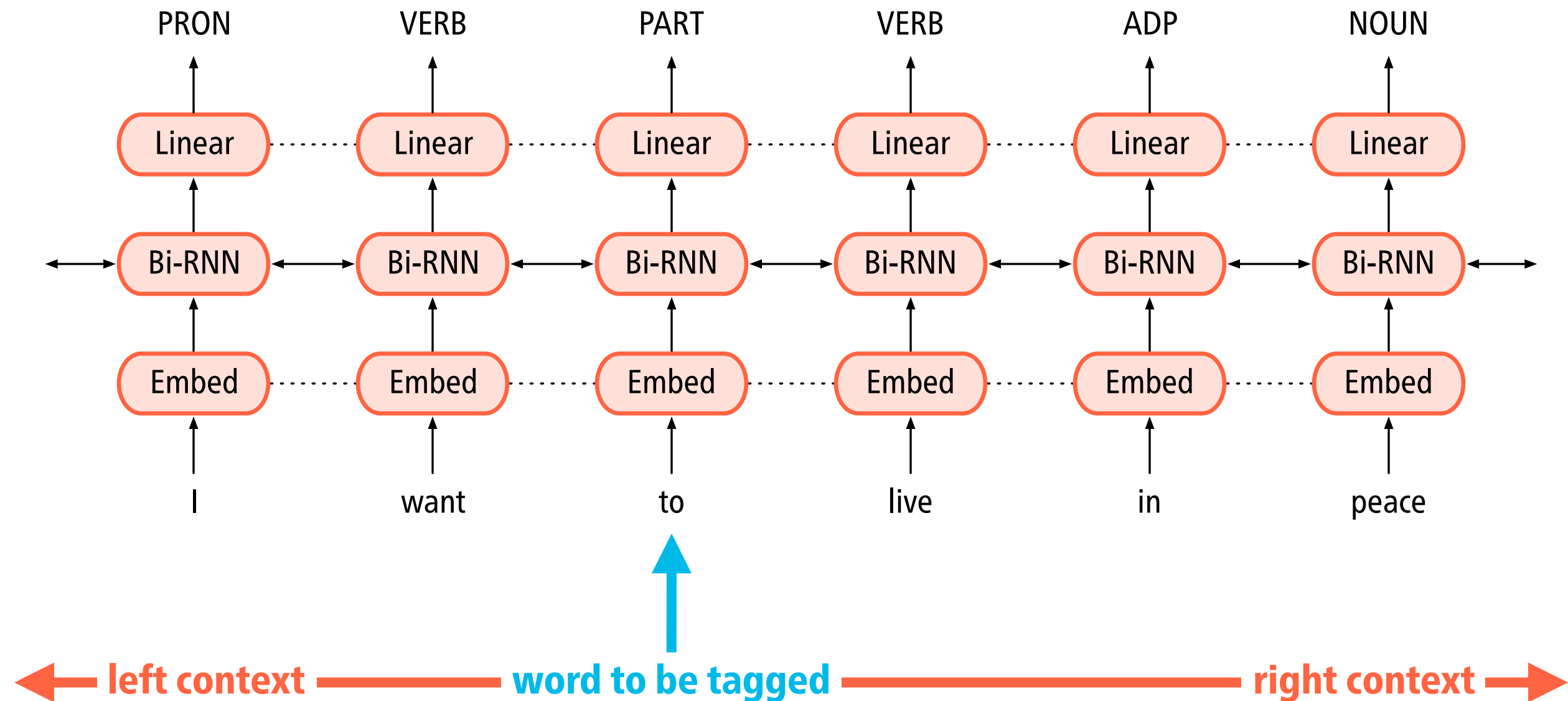
restricted feature models
(factorised feature functions)

requires specialised
algorithms; slow

Fixed-window model



Bidirectional RNN model

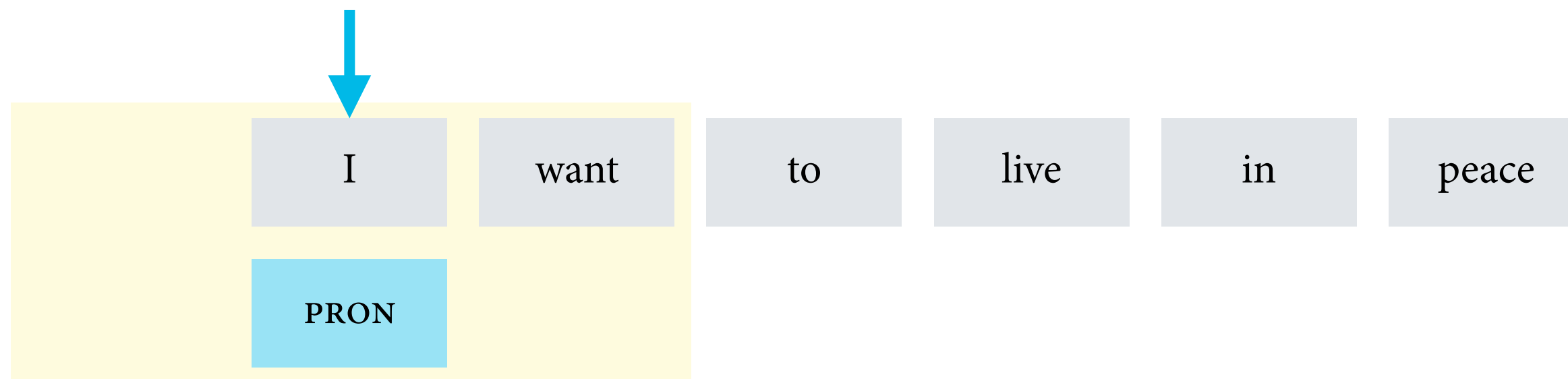


Labels are interdependent

I	want	to	live	in	peace
PRON	VERB	PART	VERB	ADP	NOUN
PRON	VERB	ADP	VERB	ADP	NOUN

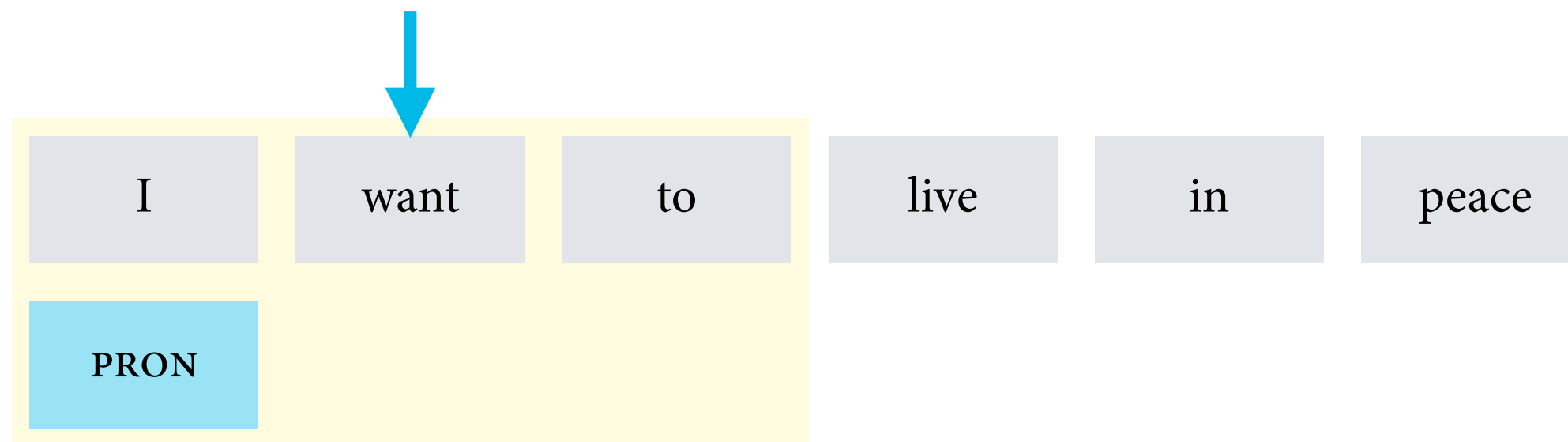
Some combinations of part-of-speech tags
are more likely than others.

Autoregressive tagging with a fixed-window model



The model predicts the tag for the first word in the sentence.
Features are extracted from a context window.

Autoregressive tagging with a fixed-window model



The prediction at the second position can use features defined over the tags that have already been predicted.

Training autoregressive models

- At test time, we run the model incrementally, and feed it with its own predicted labels.
- At training time, we feed the model with the gold-standard label. This regime is called **teacher forcing**.
- Teacher forcing can be problematic, because the model does not learn to deal with its own prediction errors.

difference between training time and prediction time (exposure bias)

Approaches to sequence labelling

Local search

sequence of independent
classification problems

expressive feature models
(RNNs, attention)

no need for specialised
algorithms; fast

Global search

combinatorial optimisation over
the full search space

restricted feature models
(factorised feature functions)

requires specialised
algorithms; slow

Dealing with combinatorial explosion

- The number of candidate sequences is exponential in the length of the input sequence, so naive optimisation is doomed to fail.
- To make the search problem tractable, we will restrict ourselves to **factorised scoring functions**.
- Factorised scoring functions will allow us to use specialised algorithms to solve the optimisation problem in polynomial time.

Example: Viterbi algorithm

Factorised scoring function

candidate
output sequence

sequence
length

$$\text{score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = \sum_{i=1}^{|\mathbf{x}|} \text{score}_1(\mathbf{x}, i, y_i; \boldsymbol{\theta}) + \sum_{i=1}^{|\mathbf{x}|} \text{score}_2(\mathbf{x}, i, y_{i-1}, y_i; \boldsymbol{\theta})$$

input
sequence

score for a
single label

score for a
pair of labels

Maximum Entropy Markov Model (MEMM)

$$\text{score}(\mathbf{x}, i, y_{i-1}, y_i) = \text{score}_1(\mathbf{x}, i, y_i) + \text{score}_2(\mathbf{x}, i, y_{i-1}, y_i)$$

$$P(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^{|\mathbf{x}|} \frac{\exp(\text{score}(\mathbf{x}, i, y_{i-1}, y_i))}{\sum_{y'} \exp(\text{score}(\mathbf{x}, i, y_{i-1}, y'))}$$

candidate
label

Algorithmic problems for MEMMs

- **Training:** To train a model, we want to minimise the negative log likelihood on gold-standard examples $((\mathbf{x}, y_{i-1}), y_i)$.

standard softmax regression problem

- **Decoding:** For a trained model, we want to find the most probable label sequence \mathbf{y} , given the input sequence \mathbf{x} .

involves the search over exponentially many candidate sequences \rightarrow Viterbi

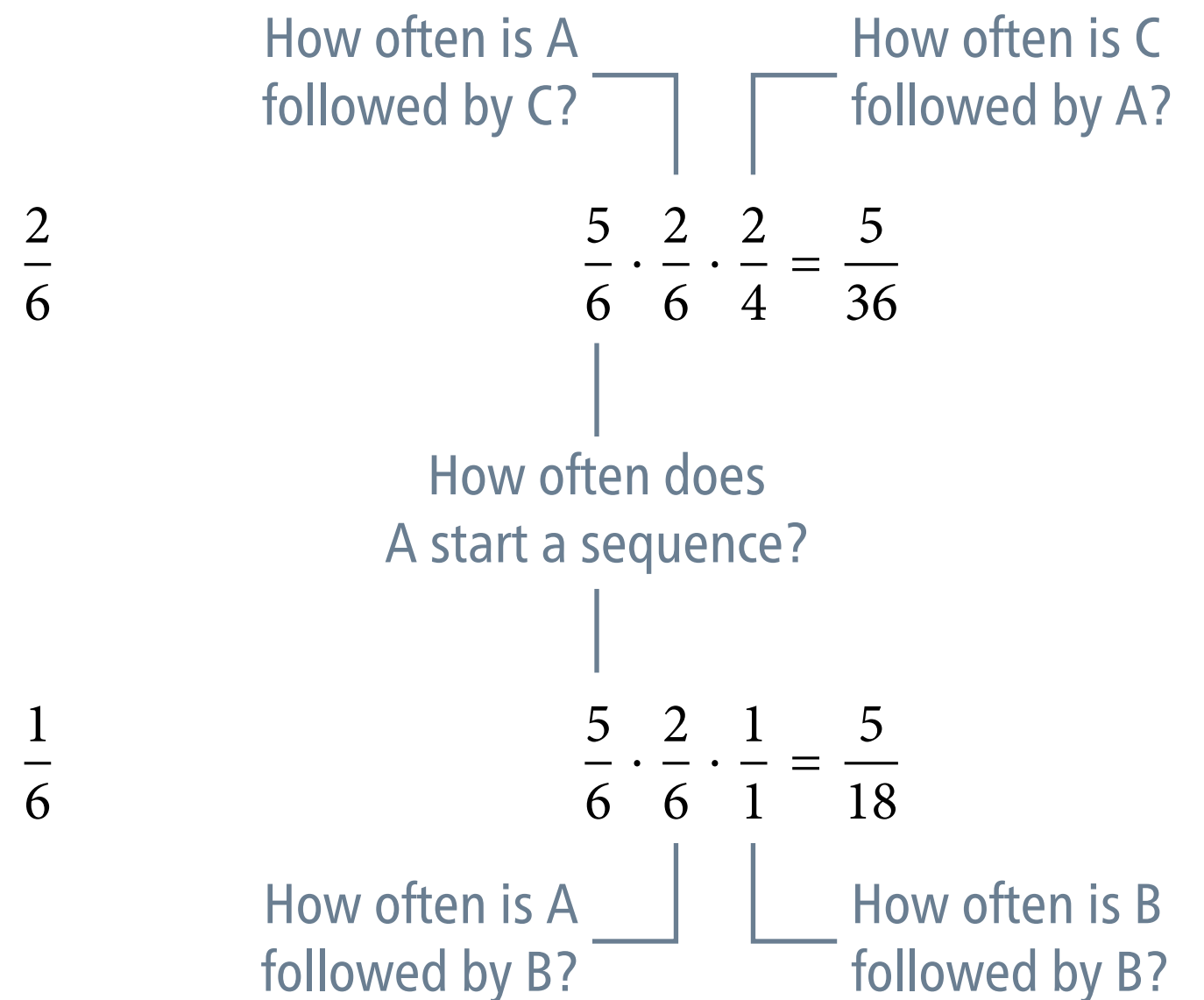
The label bias problem

Zhang and Teng (2021)

Sequence id	Sequence
1	ACA
2	ACA
3	AA
4	AAB
5	ABB
6	CCC

Global probability

Factorised probability



Conditional random field (CRF)

Lafferty et al. (2001)

$$\text{score}(\mathbf{x}, \mathbf{y}) = \sum_i^{|\mathbf{x}|} \text{score}(\mathbf{x}, i, y_{i-1}, y_i)$$

$$P(\mathbf{y} | \mathbf{x}) = \frac{\exp(\text{score}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\text{score}(\mathbf{x}, \mathbf{y}'))}$$

candidate
label sequence