

# *Structured Reordering* for Modeling Latent Alignments in Sequence Transduction

---

Bailin Wang<sup>†</sup>, Mirella Lapata<sup>†</sup> and Ivan Titov<sup>†§</sup>

<sup>†</sup> ILCC, University of Edinburgh, <sup>§</sup> ILLC, University of Amsterdam



UNIVERSITY  
OF AMSTERDAM

# Systematic Generalization

## Training Examples

what is the length of the colorado river ?

*len( river( riverid ( 'colorado' ) ) )*

what is the longest river ?

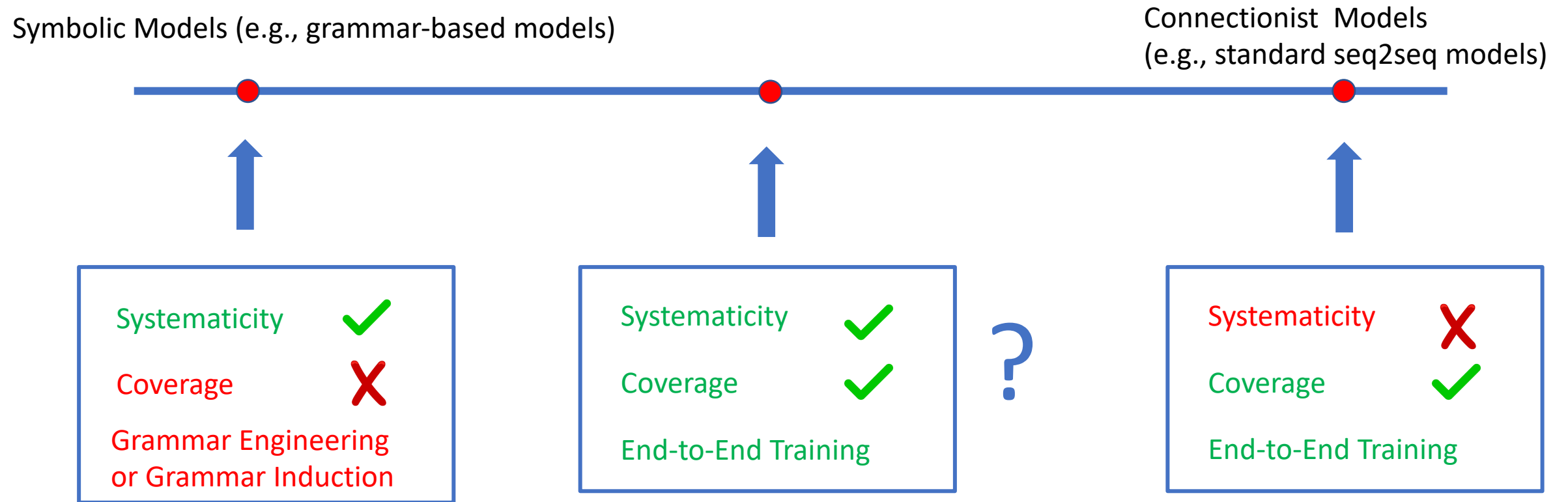
*longest( river( all ) )*

## Test Example

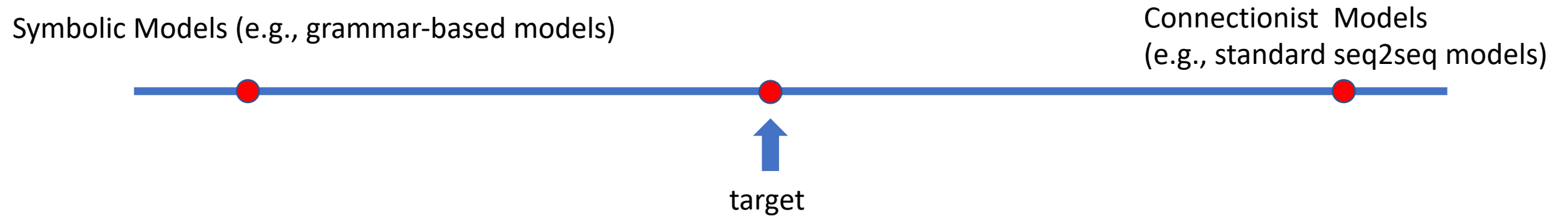
what is the length of the longest river ?

*len( longest( river( all ) ) )*

# The Spectrum of Sequence Transduction Models



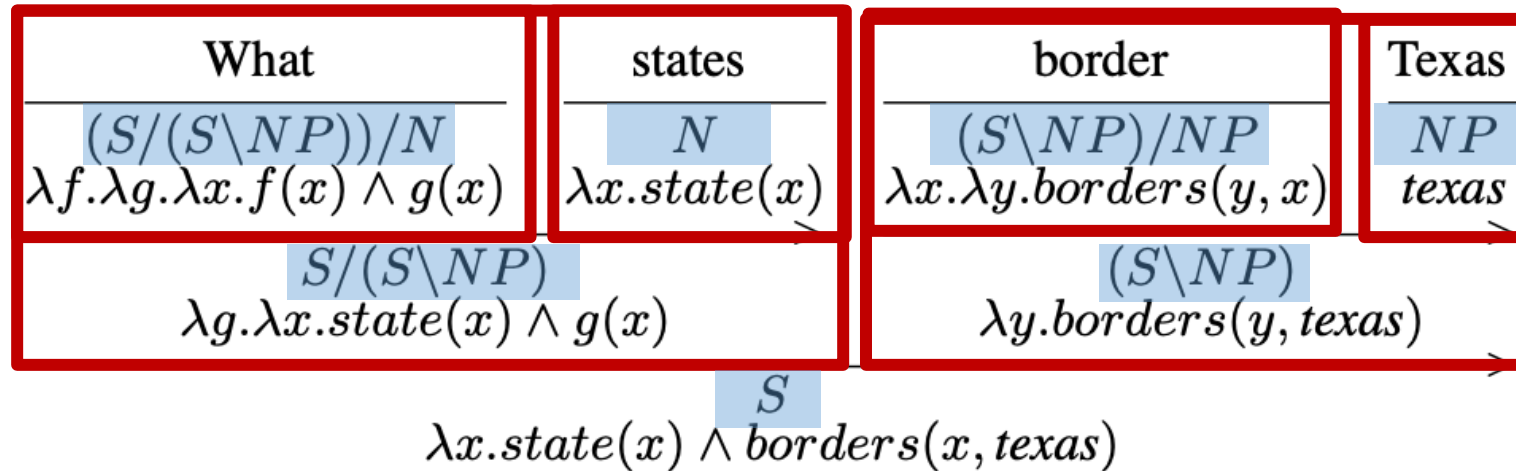
# The Spectrum of Sequence Transduction Models



## Two inter-related questions:

- 1) Why are symbolic models good at systematic generalization?
- 2) What prevents seq2seq models from generalizing systematically?

# Q1: Why are symbolic models good?



**Figure:** A CCG parse of the utterance “what states border Texas”

# Q1: Why are symbolic models good?

$$\frac{\text{What}}{(S/(S \setminus NP))/N}$$
$$\lambda f. \lambda g. \lambda x. f(x) \wedge g(x)$$

**A CCG rule**

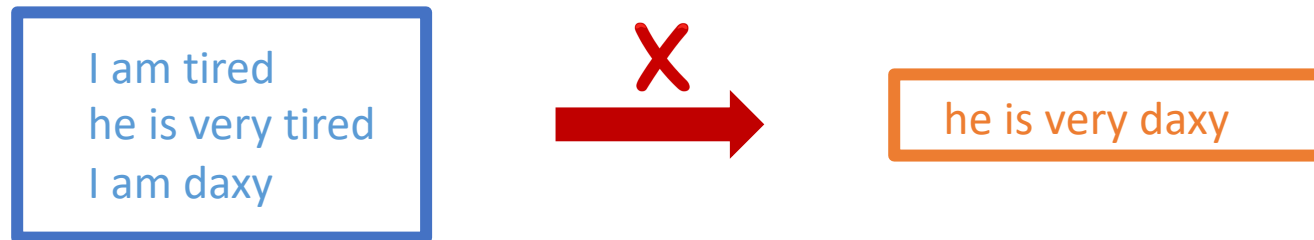
Their grammar rules implicitly encode alignments between *input and output segments*.



- *Explicit decomposition* of input and output into segments
- *consistent mappings* from input segments to output ones

## Q2: What makes seq2seq fail?

- Primitive units (e.g., words) are inconsistently mapped across different contexts [1].



- Standard seq2seq models tend to memorize large chunks, e.g., **they process 'I am daxy' as a whole** [2].

*They do not exhibit a strategy of decomposition and consistent mapping!*

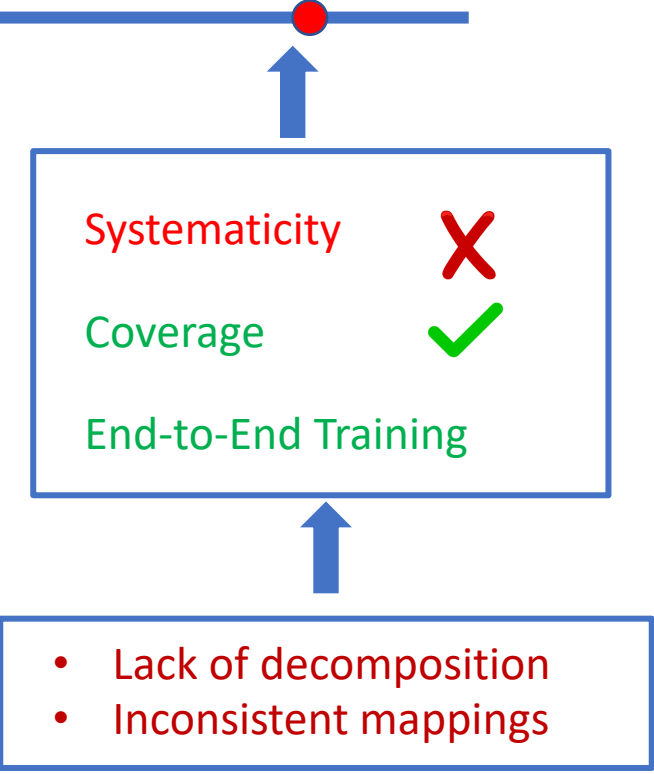
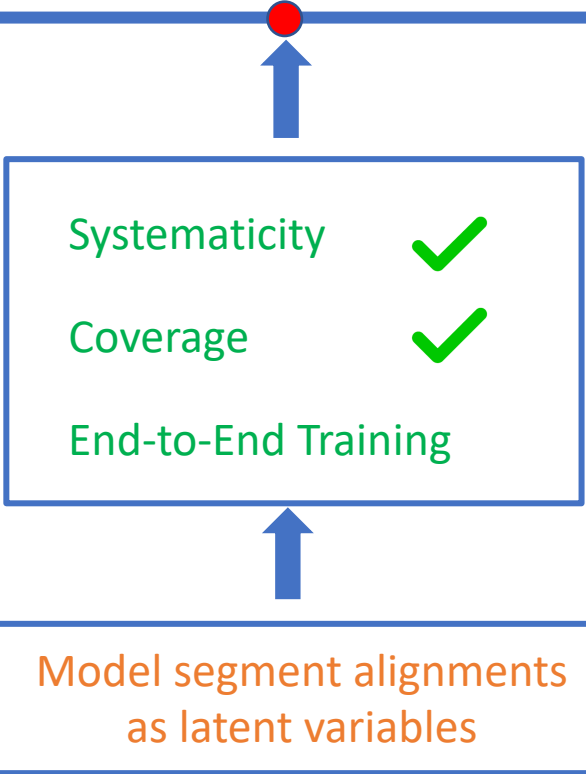
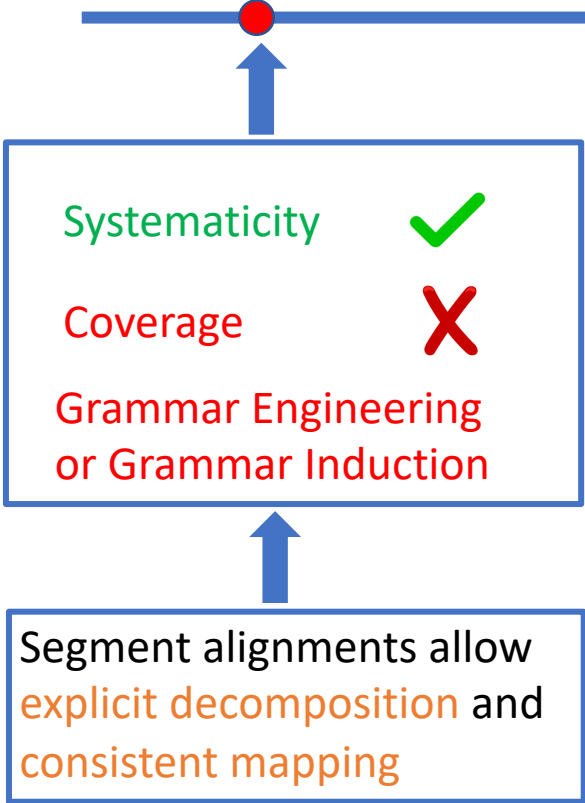
1. [Brenden Lake, Marco Baroni, *Generalization without Systematicity: On the Compositional Skills of Sequence-to-Sequence Recurrent Networks*, 2017]

2. [Dieuwke Hupkes, Verna Dankers, Mathijs Mul, Elia Bruni, *Compositionality decomposed: how do neural networks generalise?* 2019]

# Alignments for Systematicity

Symbolic Models (e.g., grammar-based models)

Connectionist Models  
(e.g., standard seq2seq models)





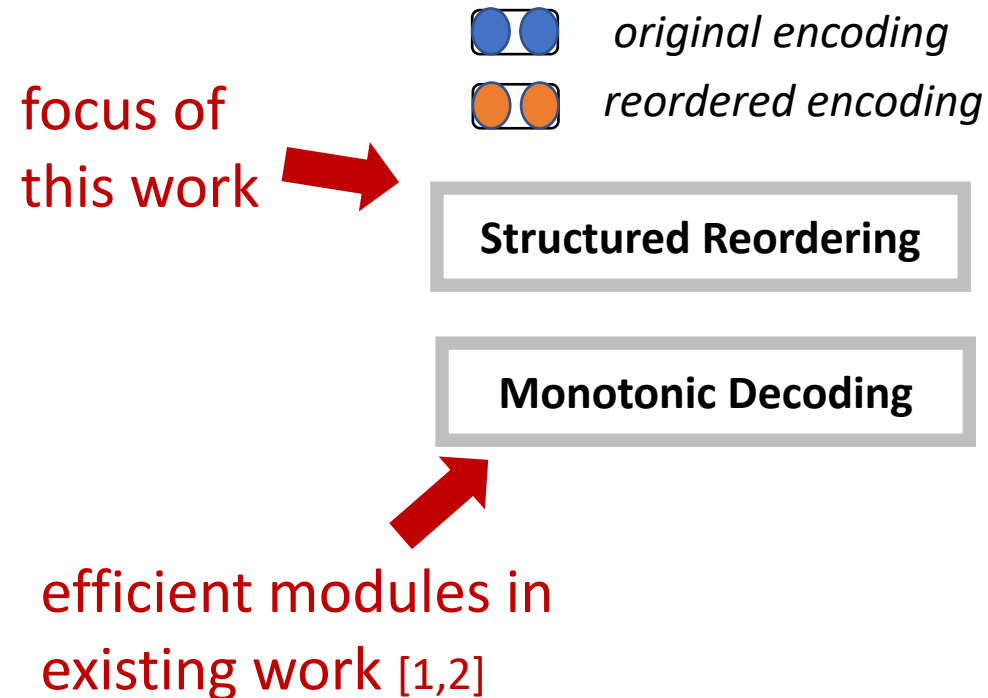
A seq2seq model endowed with *latent  
segment alignments*

# Model Architecture

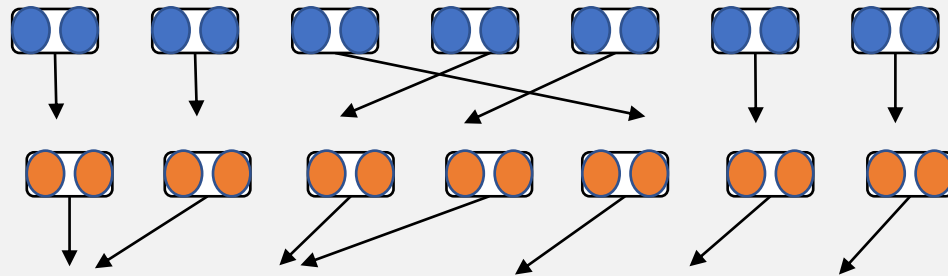
## ReMoto

Reordered-then-Monotone alignments

**Discrete latent variables rather than attention**



how many states do not have rivers



count exclude state\_all loc 1 river\_all

1. [Lei Yu, Jan Buys, and Phil Blunsom. Online segment to segment neural transduction, 2016]

2. [Chong Wang, Yining Wang, Po-Sen Huang, Abdelrahman Mohamed, Dengyong Zhou, and Li Deng. Sequence modeling via segmentations. 2017]

# Structured reordering via separable permutations

**Permutation:**

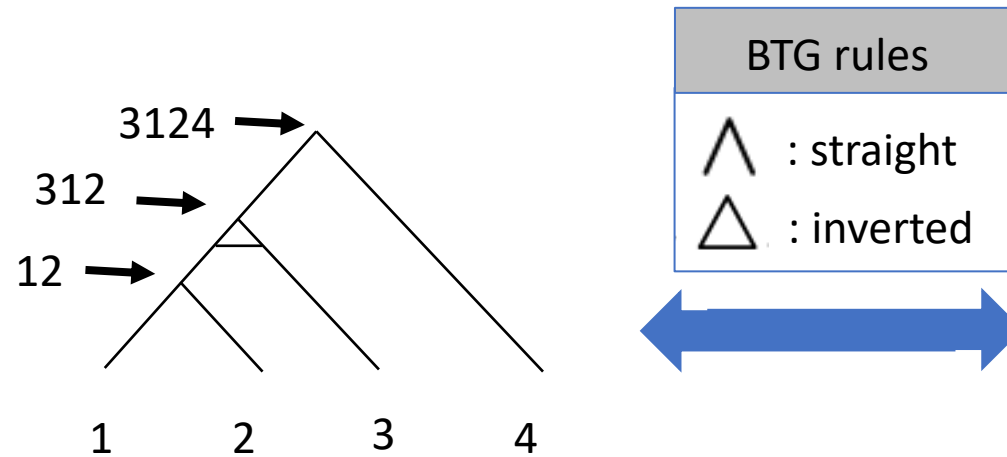


*these five young lads*



*young these five lads*

**Separable Permutation:**

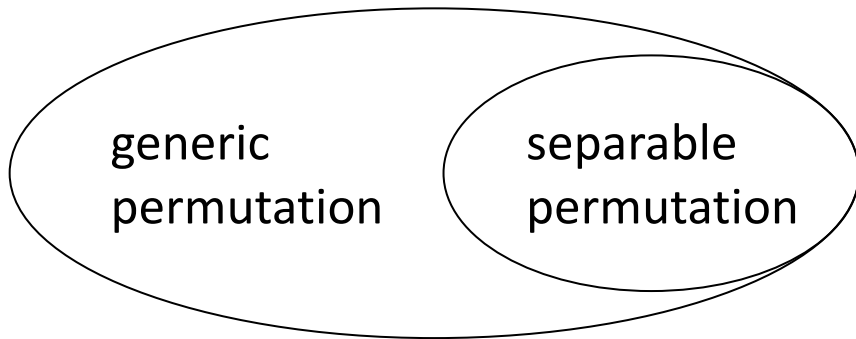


Forbidden patterns!

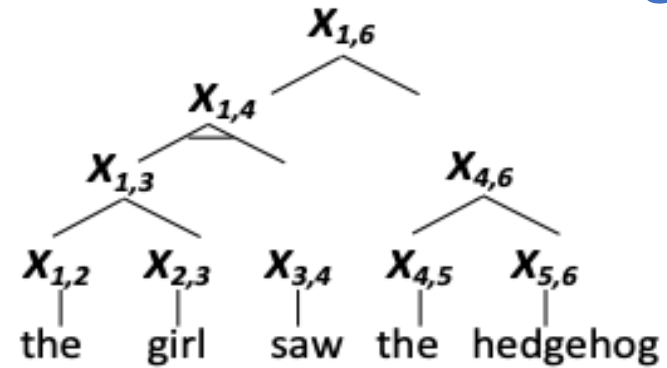


# Why separable permutation?

## Computational efficiency



## Hierarchical modeling



Linguistic inductive bias [1,2]

1. [Mark Steedman. *A formal universal of natural language grammar*, 2020]
2. [Miloš Stanojevic and Mark Steedman. *Formal basis of a language universal*, 2021]

# Training Objective:

$$-\log \mathbb{E}_{\underbrace{p_\phi(D|x)}_{\text{parser for reordering}}} \underbrace{p_\theta(y|M^D X)}_{\text{monotonic decoding}}$$

*reordered representations*

|        |   |
|--------|---|
| $x, y$ | input and output sequence                   |
| $X$    | encodings of input $x$                      |
| $D$    | permutation tree                            |
| $M^D$  | the permutation matrix corresponding to $D$ |

# Surrogate Objective:

$$-\log p_\theta(y|M' X)$$

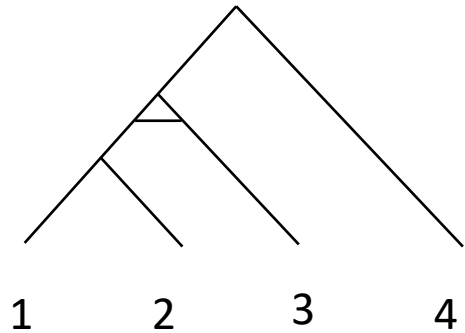
where

$$\underbrace{M'}_{\text{expectation of permutation matrices}} = \mathbb{E}_{p_\phi(D|x)} M^D \quad \text{or} \quad \underbrace{M'}_{\text{an approximated sample of permutation matrix}} = \text{Perturb-and-MAP}[p_\phi(D|x)]$$

Marginal Inference

MAP Inference

# Constructing Permutation Matrices



Permutation Tree

3 1 2 4



|   |   |   |   |
|---|---|---|---|
| 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 |

Permutation Matrix

$$((M_1^2 \oplus M_2^3) \ominus M_3^4 \oplus M_4^5)$$

$$M_i^{i+1} = \mathbf{1} \quad \mathbf{A} \oplus \mathbf{B} = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \quad \mathbf{A} \ominus \mathbf{B} = \begin{bmatrix} \mathbf{0} & \mathbf{A} \\ \mathbf{B} & \mathbf{0} \end{bmatrix}$$

# Recursion underlying Marginal Inference

*expected permutation matrix  
for the segment from i to k*

$$\underline{E_i^k} = \sum_{i < j < k} \underline{p(\text{straight})(E_i^j \oplus E_j^k) + p(\text{inverted})(E_i^j \ominus E_j^k)}$$

*enumerate all middle points*

*for every middle point, there are two ways to combine  
the permutation matrices from the left and right segments*

**Surrogate Objective:**

$$-\log p_\theta(y|M'X)$$

$$\underline{M' = E_1^n}$$

*soft-ReMoto*

$$\underline{M' = \tilde{E}_1^n}$$

*hard-ReMoto*

*a perturbed variant based on  
Straight-through Gumbel*

# Experiments



# Diagnostic Task: Infix-to-Postfix Conversion

- **Train** (nesting depth  $< 7$ )

*Input:*  $((1 + 9) * ((7 + 8) / 4))$

*Output:*  $((1 9 +)((7 8 +) 4 /) *)$

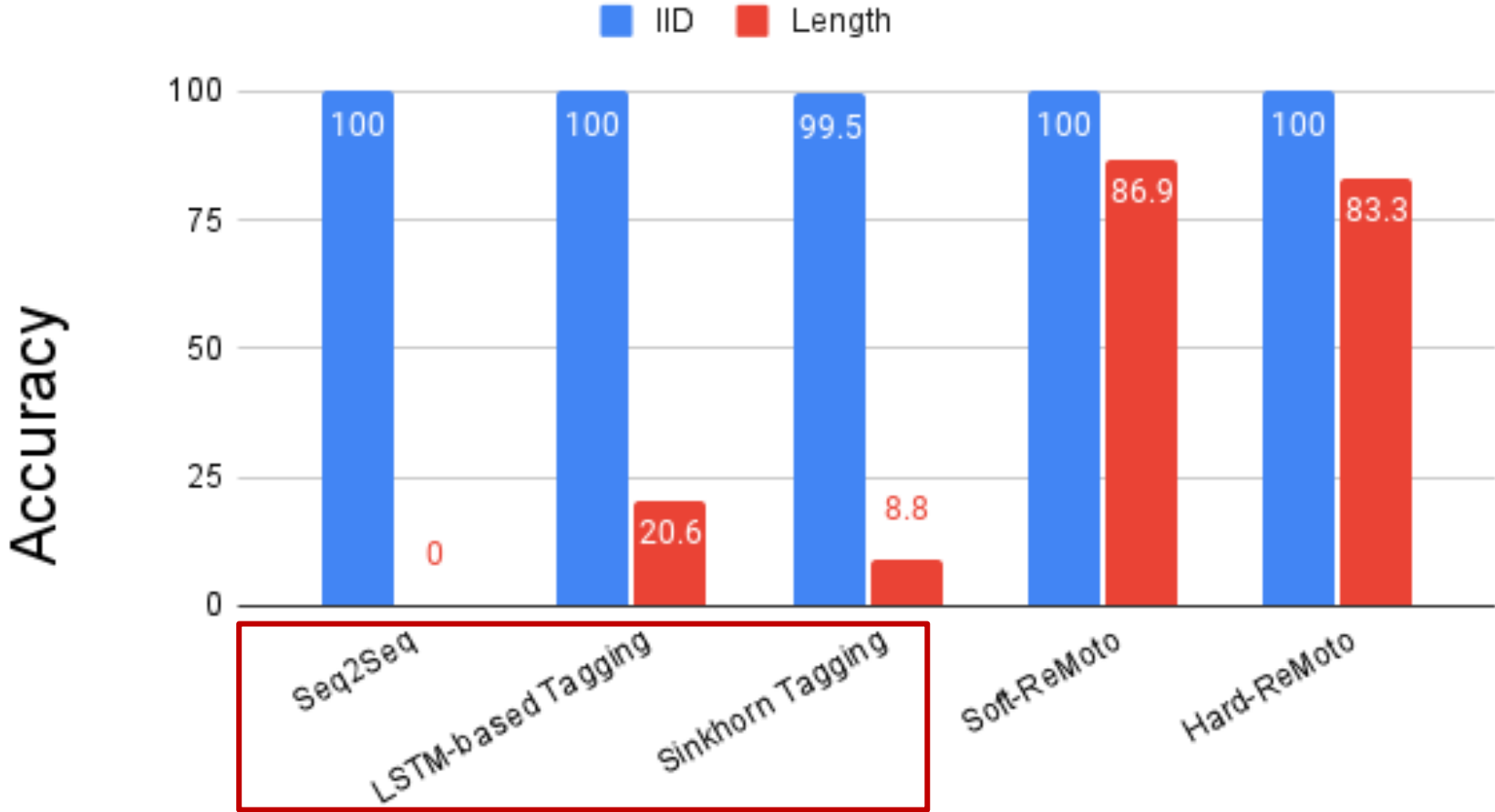
*Input:*  $((6 + 5) * (3 + 2))$

*Output:*  $((6 5 +)(3 2 +) *)$

- **IID Evaluation** (nesting depth  $< 7$ )
- **Length Evaluation** (nesting depth  $= 7$ )

# Results

## Infix-to-Postfix Conversion



# Experiment: Semantic Parsing

- Mapping natural language utterances to executable programs

*Input:* how many states do not have rivers ?

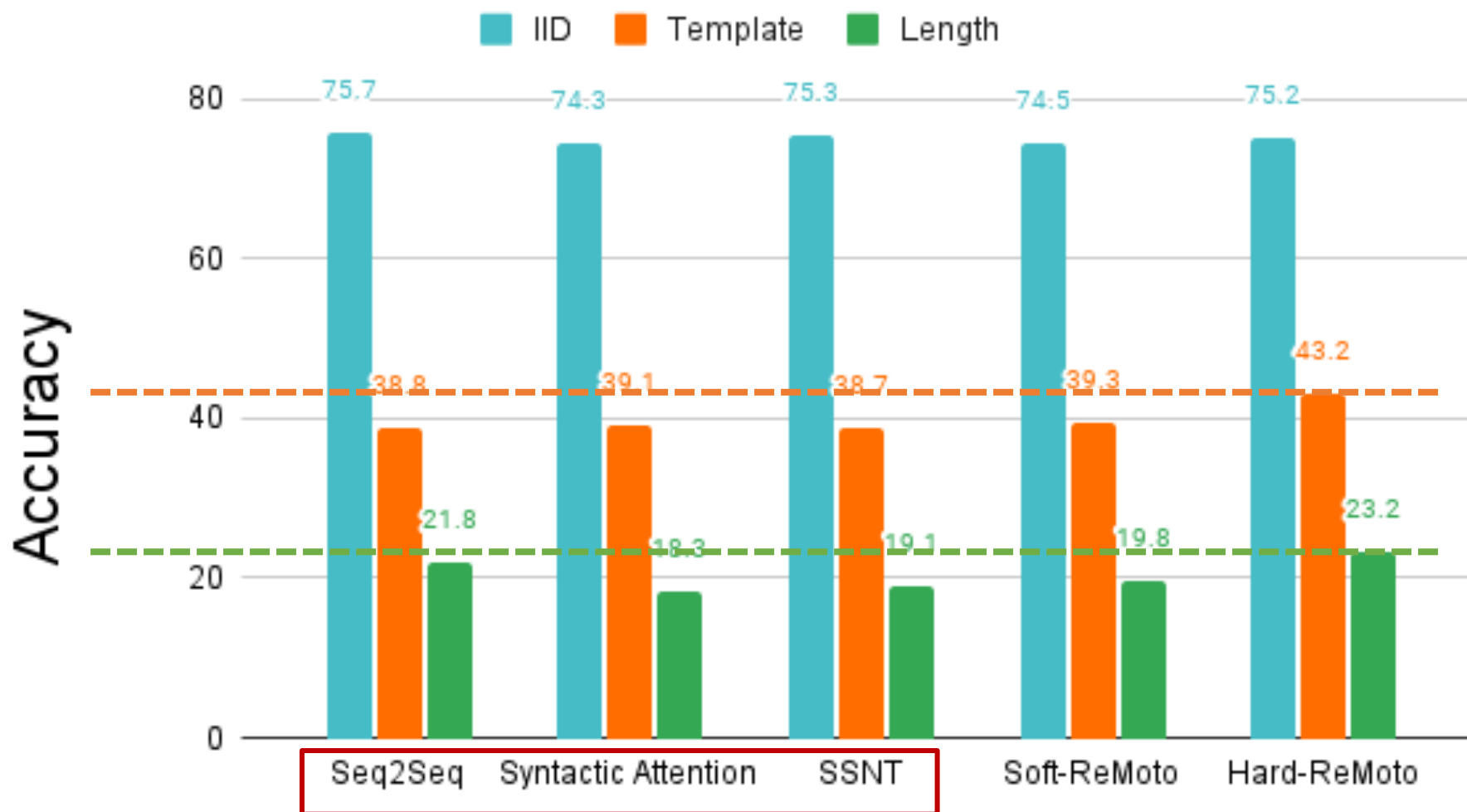
*Output:* `count(exclude(state(all), loc_1(river(all))))`

- Splits

- IID split: a standard split
- Template split: training and test examples have *disjoint templates*
- Length split: test examples are *longer* than training examples

# Results

English



# Summary

- A seq2seq model for NLP tasks that accounts for *latent non-monotonic segment-level alignments*.
- Efficient algorithms for **exact marginal and MAP inference** with separable permutations, allowing for **end-to-end training**
- Better **systematic generalization** on both synthetic and real NLP tasks.
- Code and data are available at <https://github.com/berlino/tensor2struct-public>