# Sequence to Sequence models

**Rasmus Berg Palm**

*Technical University of Denmark*
*&*
*Tradeshift*

# Sequence to Sequence models

$$\mathbf{x} \rightarrow \text{Model} \rightarrow \mathbf{y}$$

Sequence to Sequence models

**x**: "the dog ate my homework"

**y**: "El perro se comió mi tarea"

Sequence to Sequence models

$$\mathbf{x}: [x_1, x_2, \ldots, x_n]$$

$$x_i: \text{one-hot encoded word}$$

$$\mathbf{x}.\text{shape} = [n, \text{x-vocab}]$$

Sequence to Sequence models

$$\mathbf{y}: [y_1, y_2, \ldots, y_m]$$

$$y_i: \text{one-hot encoded word}$$

$$\mathbf{y}.\text{shape} = [m, \text{y-vocab}]$$

# Sequence to Sequence models

Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. **"Sequence to sequence learning with neural networks**." Advances in neural information processing systems. 2014.

https://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks.pdf

Sequence to Sequence models

1. Squeeze the entire **x** into a single vector v
2. Generate **y** conditioned on v

Sequence to Sequence models

1. Squeeze the entire **x** into a single vector v

Ideas?

Sequence to Sequence models

# Bag of Words

$$v = \text{sum}(\mathbf{x})$$

nah...

# Bag of Embeddings

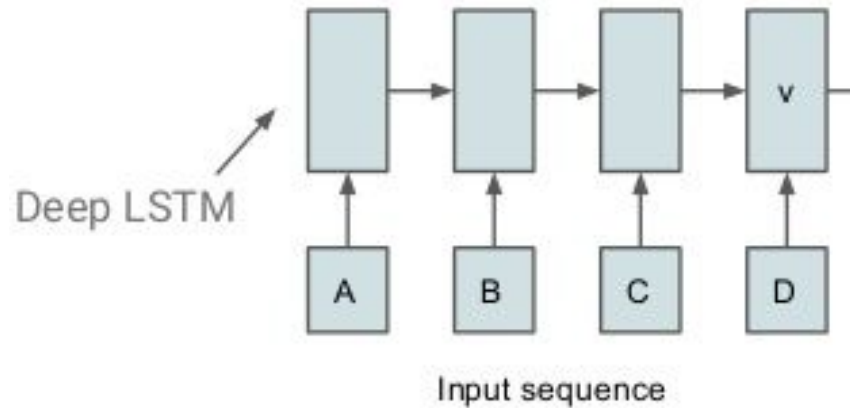$$v = \text{sum}(\text{embed}(\mathbf{x}))$$

nah...

Sequence to Sequence models

# LSTM encoding

$v = LSTM(\mathbf{x})[-1]$

yup!

# Sequence to Sequence models

[Sutskever & Vinyals & Le NIPS 2014]



Deep LSTM

Input sequence

Sequence to Sequence models

1. Squeeze the entire **x** into a single vector v

2. Generate **y** conditioned on v
   Ideas?

Sequence to Sequence models

$$\mathbf{y} = \text{LSTM*}$$

*Set initial hidden state to v

# Sequence to Sequence models



[Sutskever & Vinyals & Le NIPS 2014]

Deep LSTM

Input sequence

Target sequence

Sequence to Sequence models

# Variant that is easier to code (and better)

# Sequence to Sequence models



[Sutskever & Vinyals & Le NIPS 2014]

Deep LSTM

A  B  C  D

Input sequence

Target sequence

X  Y  Z  Q

Sequence to Sequence models

# OK.

# Let's do it!

Sequence to Sequence models

**x**: 12 November 2016

**y**: 2016-11-12

# Sequence to Sequence models

```python
1   # Encoder
2   source = Input(shape=(None,), dtype='int32', name='source')
3   embedded = Embedding(output_dim=128, input_dim=train.source_vocab_size(), mask_zero=True)(source)
4   last_hid = LSTM(output_dim=128)(embedded)
5
6   # Decoder
7   repeated = RepeatVector(train.target.padded.shape[1])(last_hid)
8   decoder = LSTM(output_dim=128, return_sequences=True)(repeated)
9   output = TimeDistributed(Dense(output_dim=train.target_vocab_size(), activation='softmax'))(decoder)
10  model = Model([source], output=[output])
```

Sequence to Sequence models

# http://localhost:5000

Sequence to Sequence models

# Trick

Feed the LSTM the last output it made

Sequence to Sequence models

# Two ways to implement

1. Feed the actual probabilities outputted

Hard, not used very often

2. Feed the target shifted by one

Easy, very used. AKA. "teacher forcing"

# Sequence to Sequence models

[Sutskever & Vinyals & Le NIPS 2014]

Target sequence

Deep LSTM

Input sequence

Sequence to Sequence models

But...

How to generate at test time then?

# Sequence to Sequence models



[Sutskever & Vinyals & Le NIPS 2014]

Deep LSTM

Input sequence

Target sequence

# Sequence to Sequence models



[Sutskever & Vinyals & Le NIPS 2014]

Deep LSTM

Input sequence

Target sequence

# Sequence to Sequence models

[Sutskever & Vinyals & Le NIPS 2014]

Target sequence

Deep LSTM

Input sequence

# Sequence to Sequence models

[Sutskever & Vinyals & Le NIPS 2014]

Target sequence

Deep LSTM

Input sequence

Sequence to Sequence models

# Exercise left for the reader

Implement teacher forcing in the date parser

Sequence to Sequence models

# The great weakness.

# Ideas?

Sequence to Sequence models

# The great weakness.

**x** is n long

v is fixed size

Hard to compress when n grows

# Sequence to Sequence models

Sequence to Sequence models

# The great solution.

## Ideas?

# Sequence to Sequence models

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. **"Neural machine translation by jointly learning to align and translate."** arXiv preprint arXiv:1409.0473 (2014)

Sequence to Sequence models

Let the decoder look at the entire input sequence for every output

AKA. "Attention"

# Sequence to Sequence models

Sequence to Sequence models

Attention is tricky...

But you're clever :]
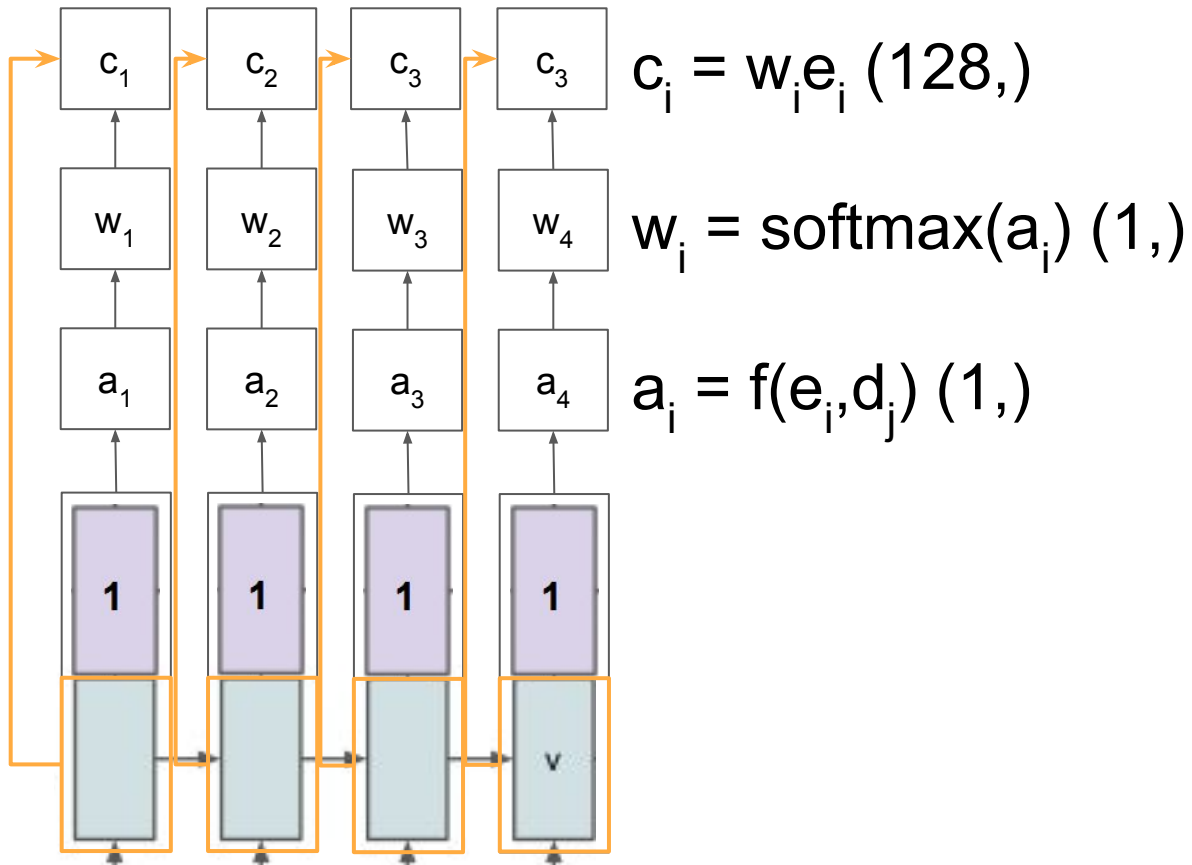
Input sequence
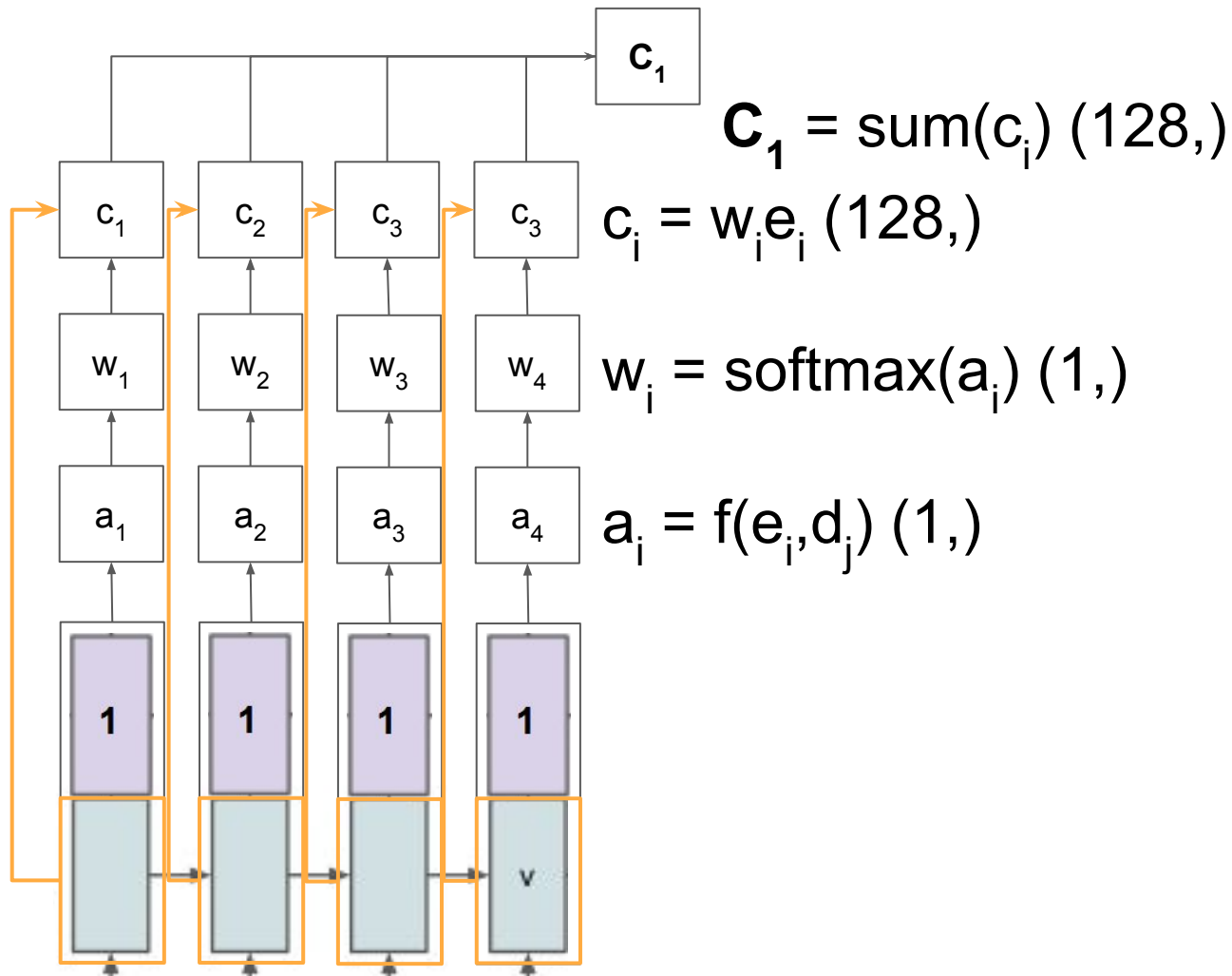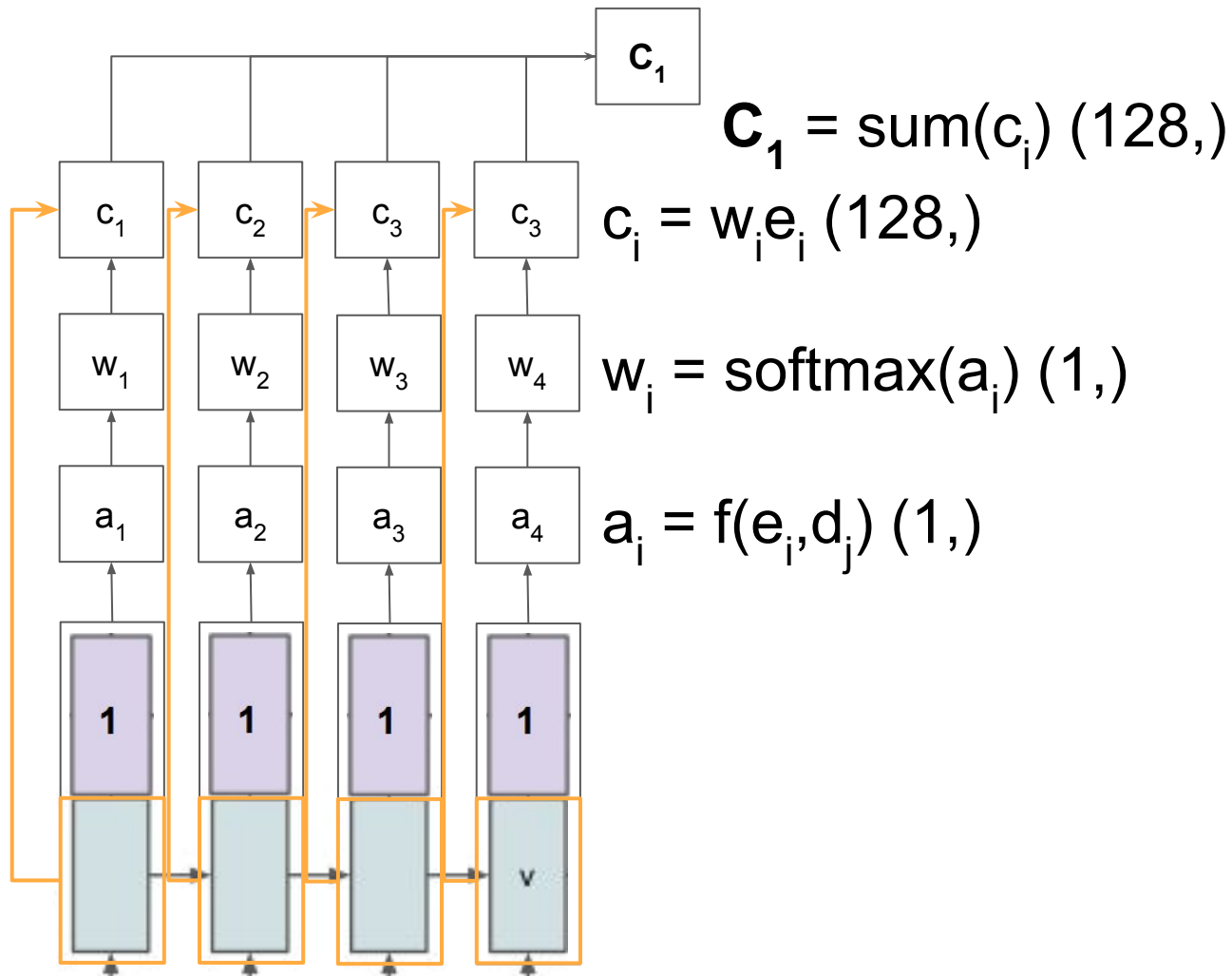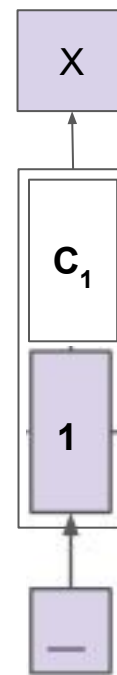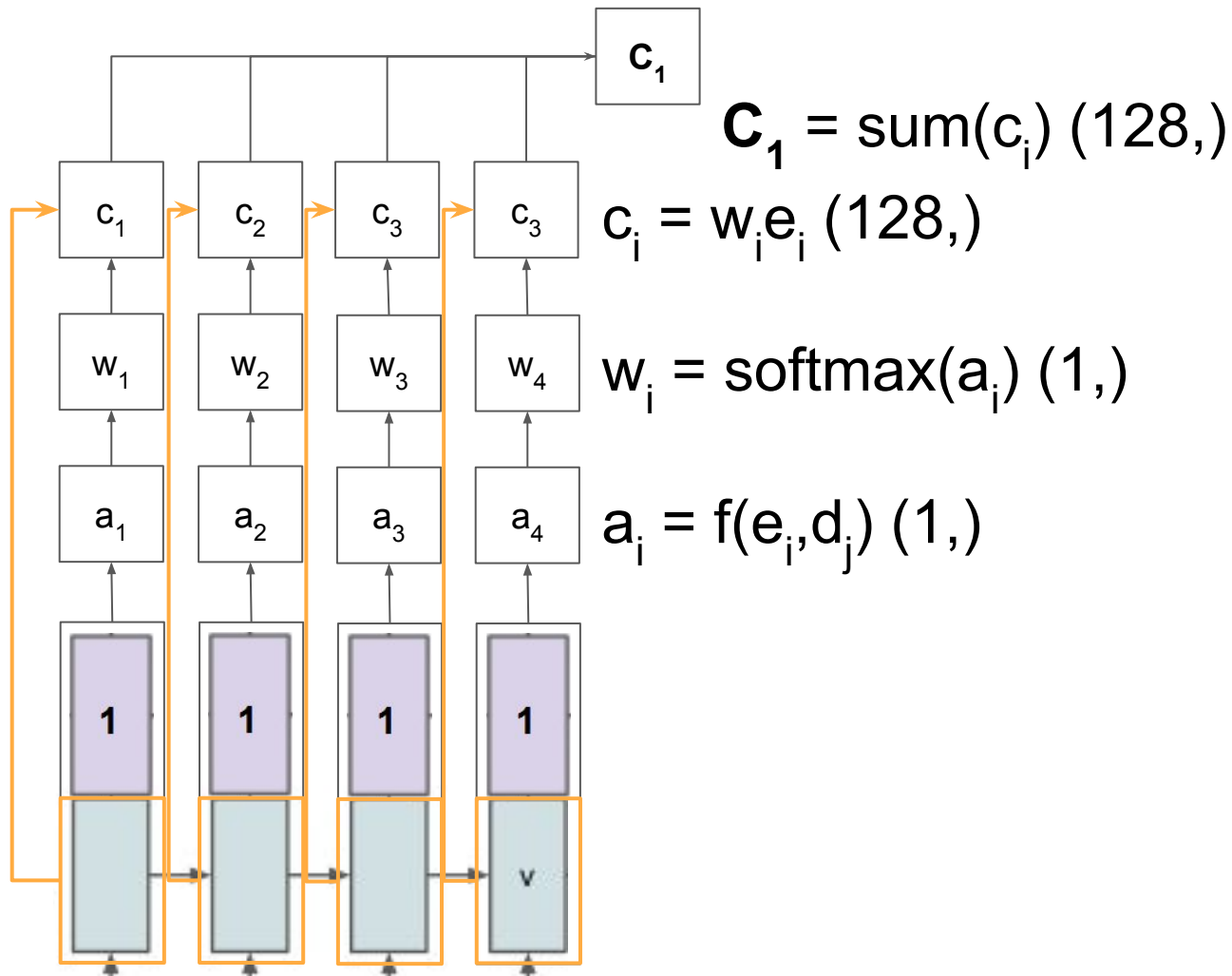
Input sequence

Input sequence

$$a_i = f(e_i, d_j) \ (1,)$$

Input sequence

$w_i = \text{softmax}(a_i)$ (1,)

$a_i = f(e_i, d_j)$ (1,)

Input sequence

$w_i$ = softmax($a_i$) (1,)

$a_i$ = f($e_i$,$d_j$) (1,)

$c_i = w_i e_i \ (128,)$

$w_i = \text{softmax}(a_i) \ (1,)$

$a_i = f(e_i, d_j) \ (1,)$

$$C_1 = \text{sum}(c_i) \; (128,)$$

$$c_i = w_i e_i \; (128,)$$

$$w_i = \text{softmax}(a_i) \; (1,)$$
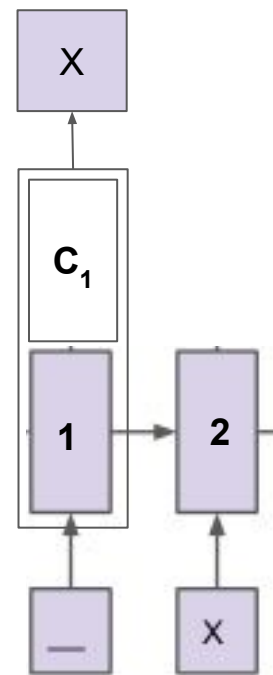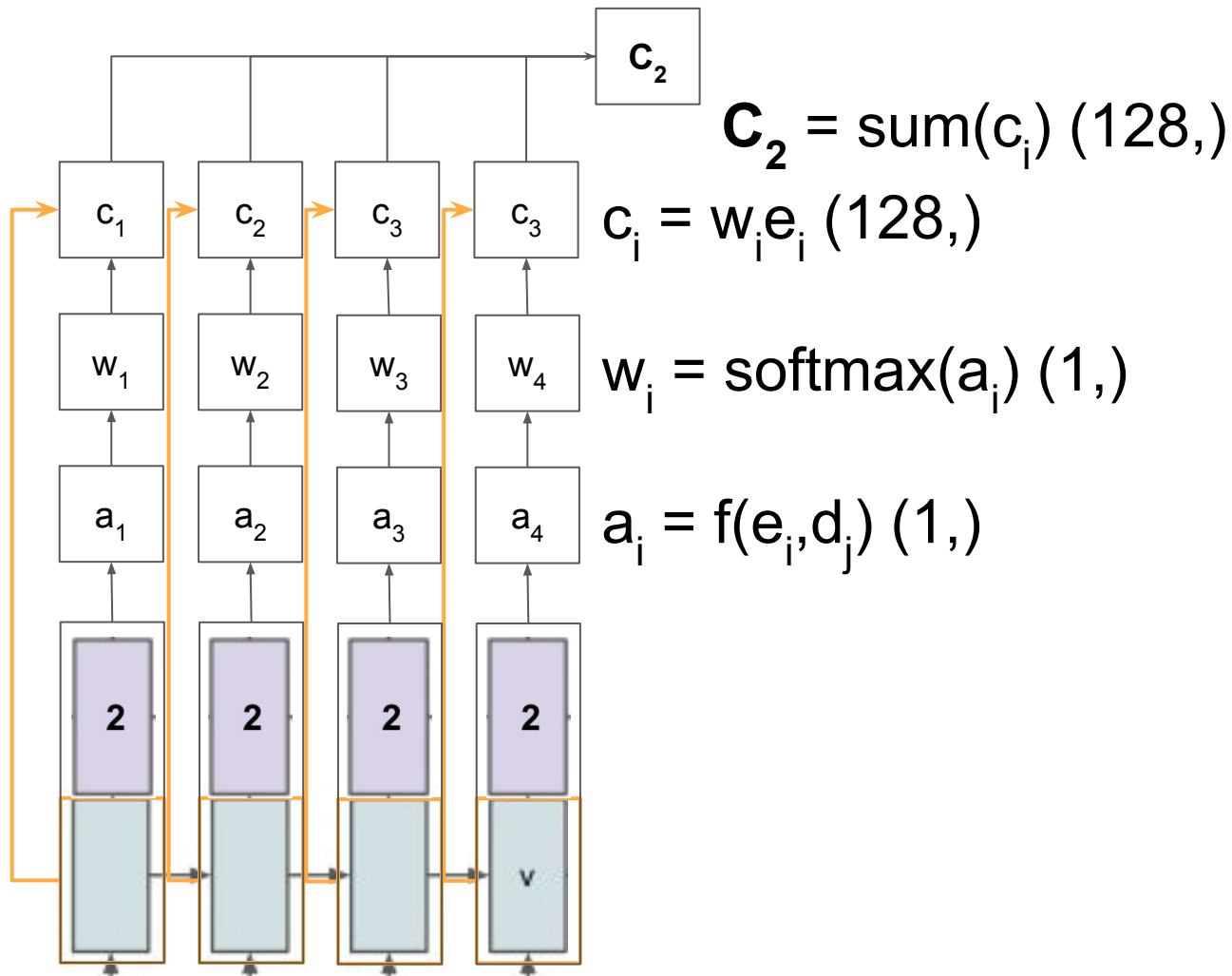
$$a_i = f(e_i, d_j) \; (1,)$$

$$\mathbf{C_1} = \text{sum}(c_i) \ (128,)$$

$$c_i = w_i e_i \ (128,)$$

$$w_i = \text{softmax}(a_i) \ (1,)$$

$$a_i = f(e_i, d_j) \ (1,)$$

$$C_1 = \text{sum}(c_i)\ (128,)$$

$$c_i = w_i e_i\ (128,)$$

$$w_i = \text{softmax}(a_i)\ (1,)$$

$$a_i = f(e_i, d_j)\ (1,)$$

$C_2$ = sum($c_i$) (128,)

$c_i = w_i e_i$ (128,)

$w_i$ = softmax($a_i$) (1,)

$a_i = f(e_i, d_j)$ (1,)

$$C_2 = \text{sum}(c_i) \ (128,)$$

$$c_i = w_i e_i \ (128,)$$

$$w_i = \text{softmax}(a_i) \ (1,)$$

$$a_i = f(e_i, d_j) \ (1,)$$

$$\mathbf{C_2} = \text{sum}(c_i)\ (128,)$$

$$c_i = w_i e_i\ (128,)$$

$$w_i = \text{softmax}(a_i)\ (1,)$$

$$a_i = f(e_i, d_j)\ (1,)$$

# Sequence to Sequence models

Phew!

# Sequence to Sequence models

# Sequence to Sequence models

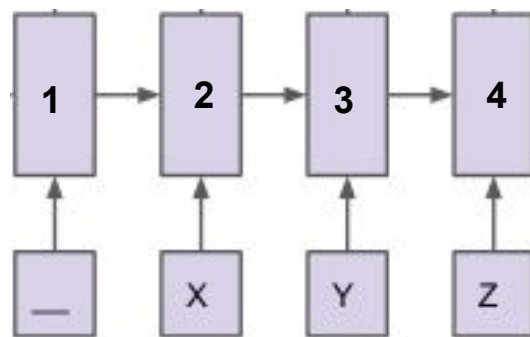Sequence to Sequence models

# Exercise left for the reader
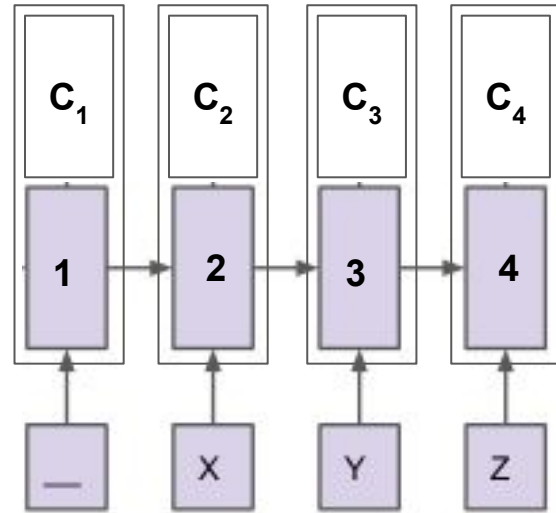
## Implement attention for the date parser
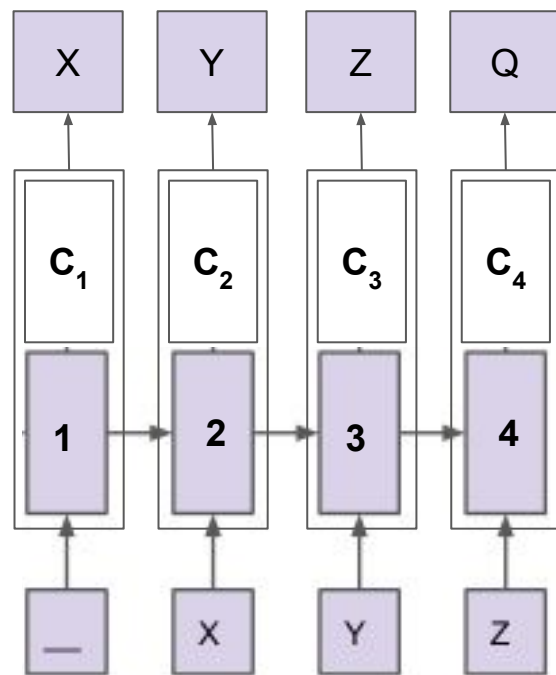
Sequence to Sequence models

# Trick

Teacher forcing makes attention easier to implement

# My own work

# Cognitive Systems
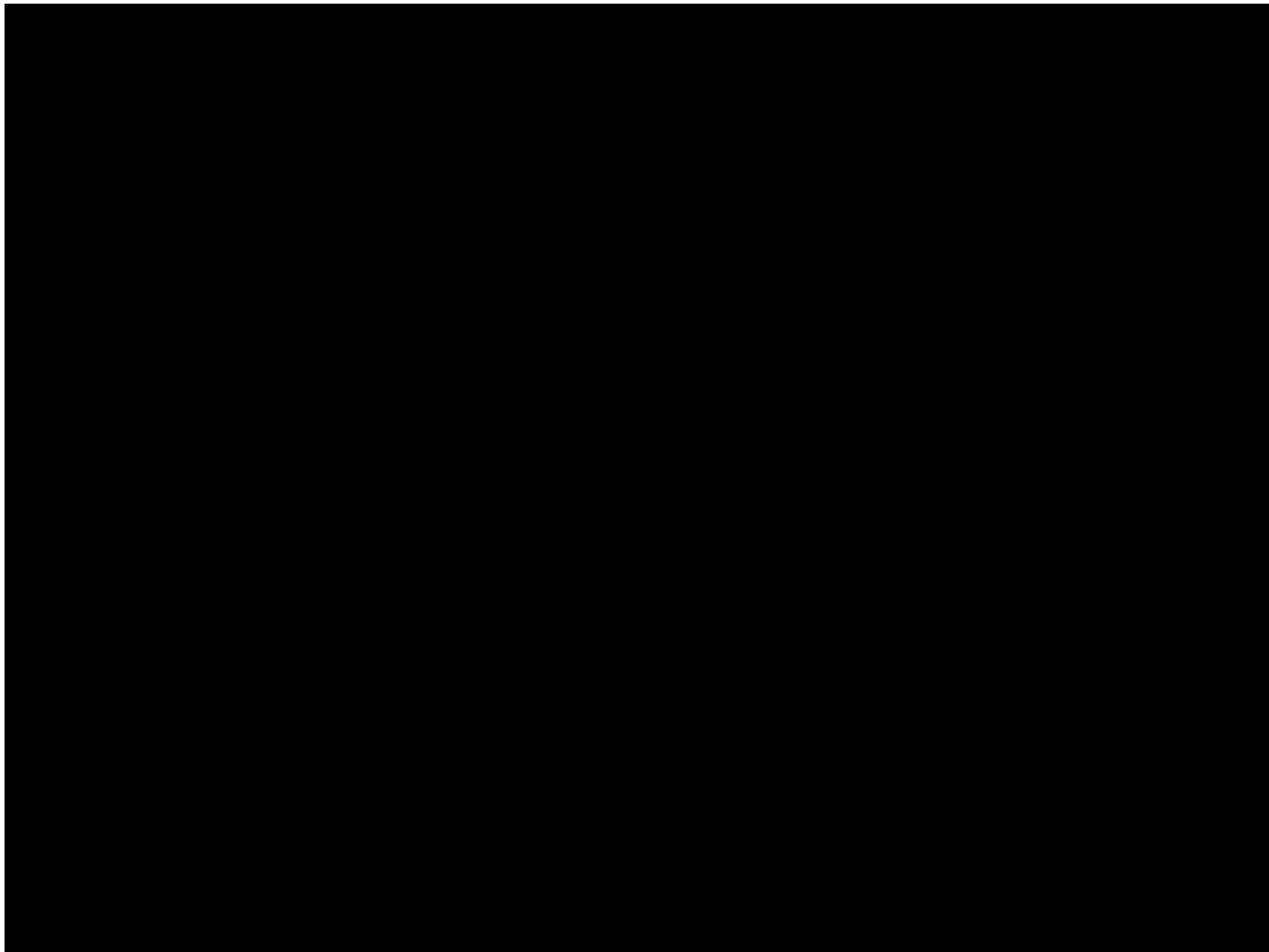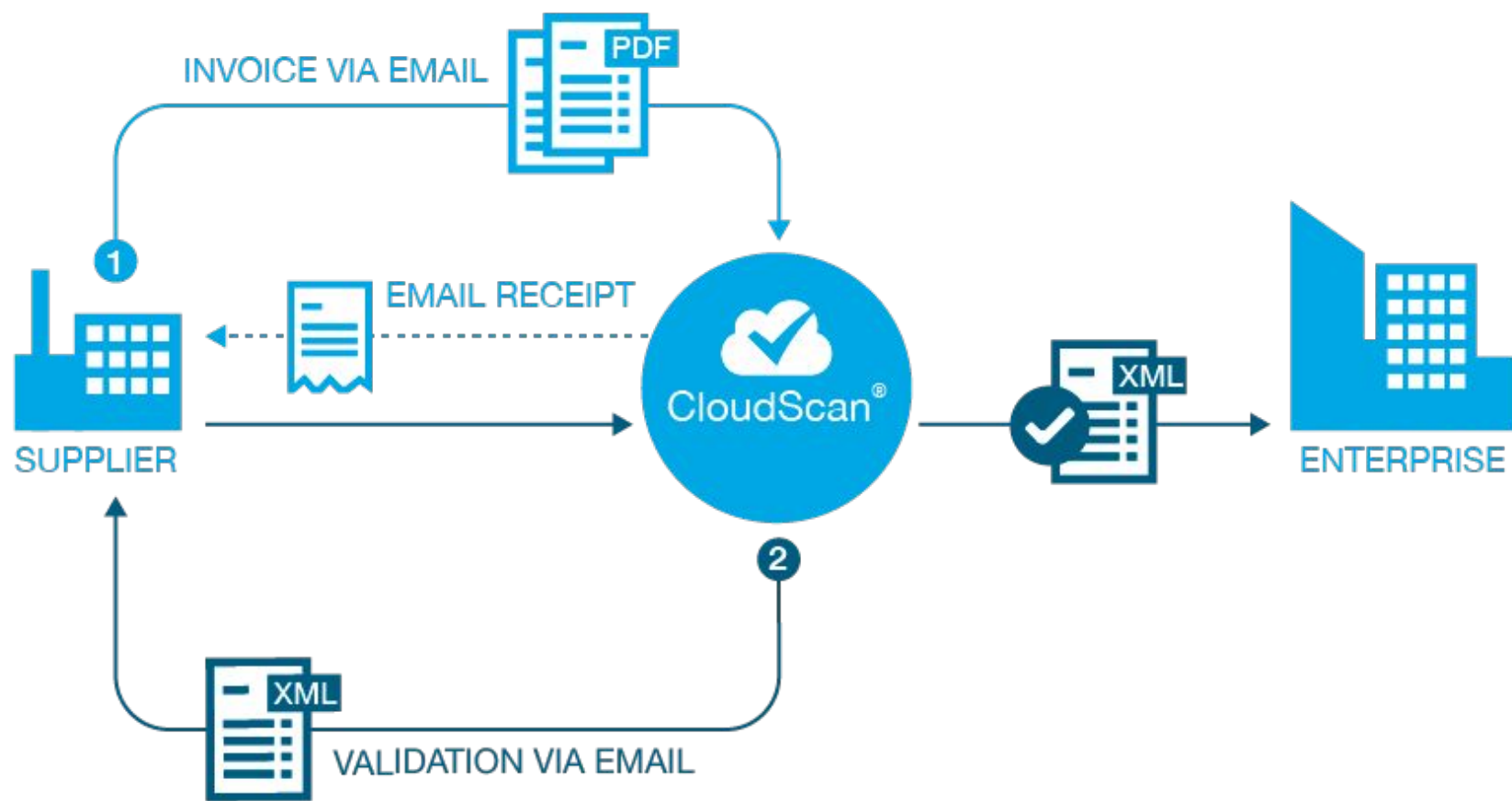*Technical University of Denmark*



Ole Winther
Professor, PhD

# Machine Learning Team

*Tradeshift*



Florian Laws
Team Lead, PhD

INVOICE VIA EMAIL

PDF

EMAIL RECEIPT

CloudScan®

SUPPLIER

XML

ENTERPRISE

1

2

XML

VALIDATION VIA EMAIL

# The actual data and problem

Interesting challenges

1.  **Training data is PDF and XML pairs. No word annotations!**
2.  Large handcrafted post-processing stage
3.  Structured output (totals have to add up, etc.)
4.  Modelling word context
5.  Using image features

# Missing word annotations

# Addressing the lack of word annotations

## 'End-to-End Information Extraction without Token-Level Supervision'

*Rasmus Berg Palm, Dirk Hovy, Ole Winther, Florian Laws*

https://arxiv.org/abs/1707.04913

Let's create a travel concierge app

Takes natural language input.

Proposes flights.



GO
Hi Freya, what can I do for you?

Freya
I need a flight from Oslo to JFK, landing tomorrow at 8pm.

GO
I found this flight for you:

Aug 9                                    KLM
Oslo                                New York
OSL                                     JFK
12:55 PM                             7:20 PM

Behind the covers..

We use a flight search engine

The search engine accepts a fixed set of fields, e.g. "`from`", "`to`", "`day`", etc.

We return the top-1 hit

GO
Hi Freya, what can I do for you?

Freya
I need a flight from Oslo to JFK, landing tomorrow at 8pm.

GO
I found this flight for you:

Aug 9                          KLM

Oslo                        New York
OSL          ✈           JFK
12:55 PM                   7:20 PM

Behind the covers..

We hire human operators to extract the values for the search engine fields.
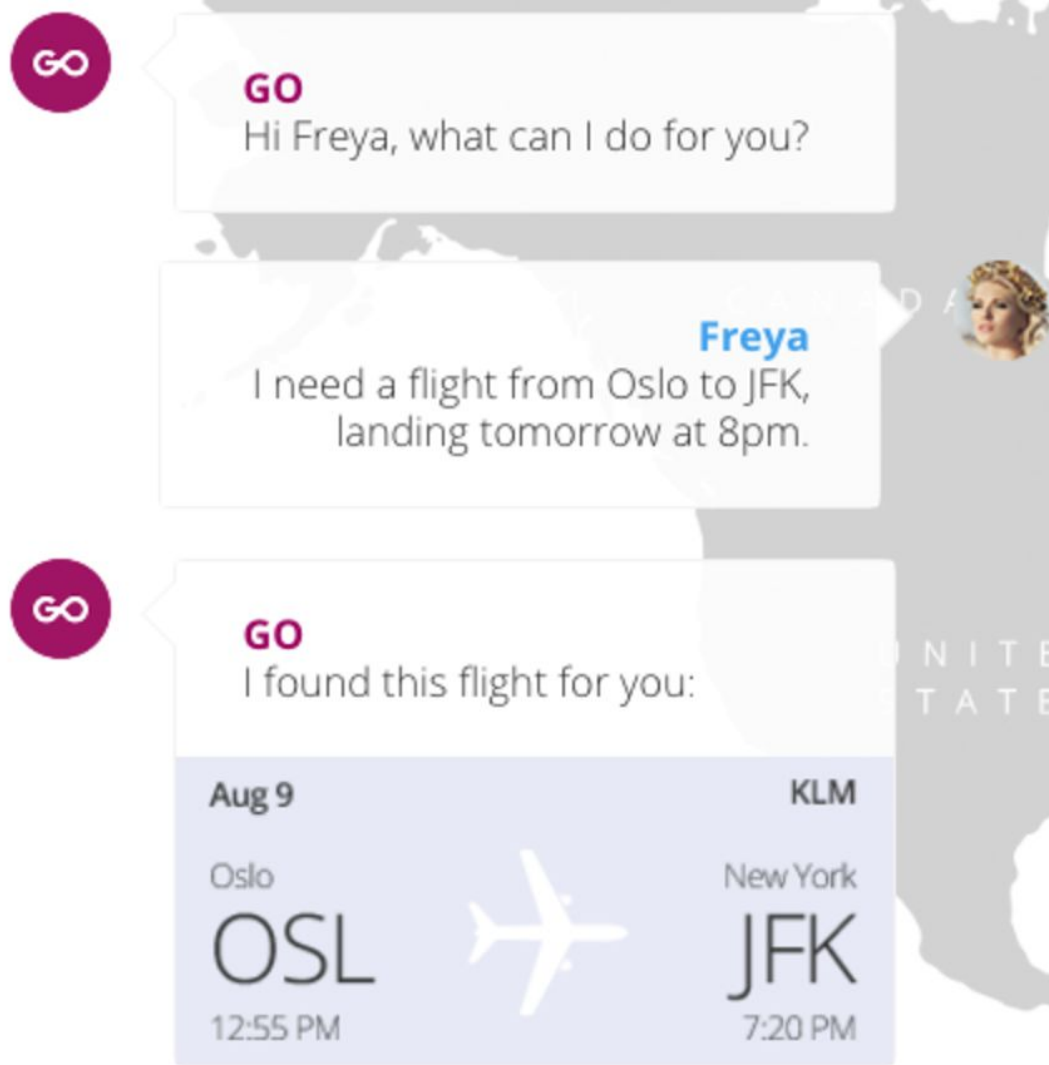
This labor is tedious for the operators and costly for us



GO
Hi Freya, what can I do for you?

Freya
I need a flight from Oslo to JFK, landing tomorrow at 8pm.

GO
I found this flight for you:

Aug 9                                          KLM

Oslo                                        New York
OSL                          ✈           JFK
12:55 PM                                   7:20 PM

So let's automate this extraction of information

## Token level

| Token | Label |
|---|---|
| cheapest | B-PRICE |
| airfare | O |
| from | O |
| tacoma | B-FROM |
| to | O |
| st. | B-TO |
| louis | I-TO |

## Token level

| Token | Label |
|---|---|
| cheapest | B-PRICE |
| airfare | O |
| from | O |
| tacoma | B-FROM |
| to | O |
| st. | B-TO |
| louis | I-TO |

| Token level | | | Field level | |
|---|---|---|---|---|
| **Token** | **Label** | | **Field** | **Value** |
| cheapest | B-PRICE | | FROM | tacoma |
| airfare | O | | TO | st. louis |
| from | O | | PRICE | cheapest |
| tacoma | B-FROM | | DAY | - |
| to | O | | MONTH | - |
| st. | B-TO | | YEAR | - |
| louis | I-TO | | AIRLINE | - |

## Token level

| Token | Label |
|---|---|
| cheapest | B-PRICE |
| airfare | O |
| from | O |
| tacoma | B-FROM |
| to | O |
| st. | B-TO |
| louis | I-TO |

Chunking →

## Field level

| Field | Value |
|---|---|
| FROM | tacoma |
| TO | st. louis |
| PRICE | cheapest |
| DAY | - |
| MONTH | - |
| YEAR | - |
| AIRLINE | - |

# Our model



cheapest airfare from tacoma to st. louis <EOS>
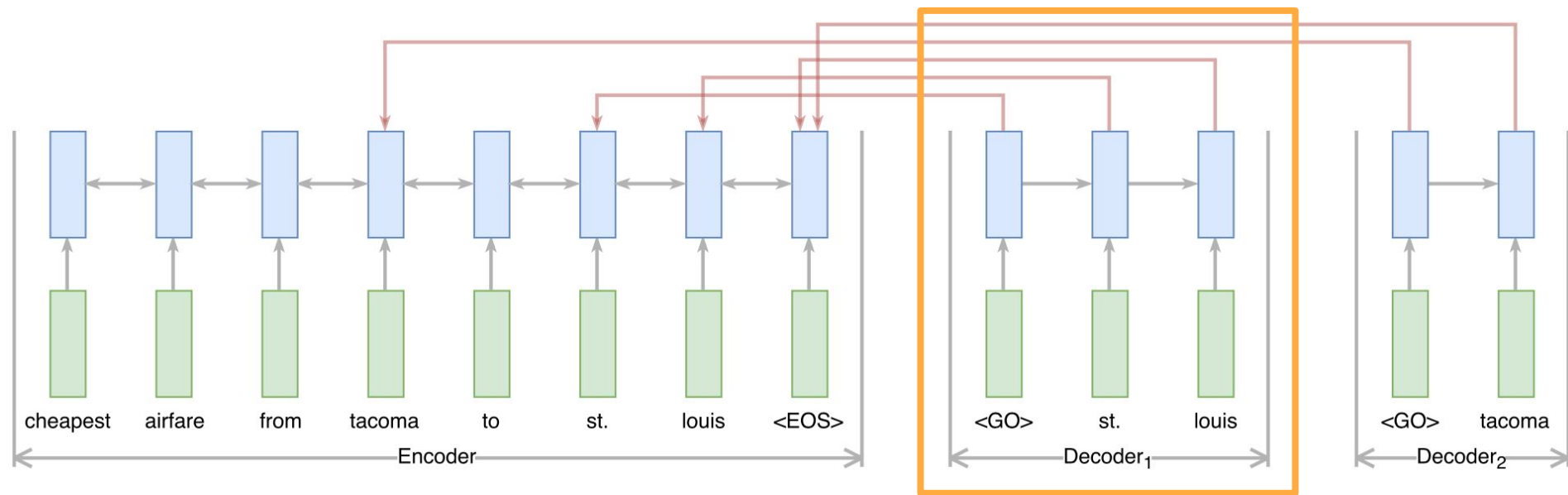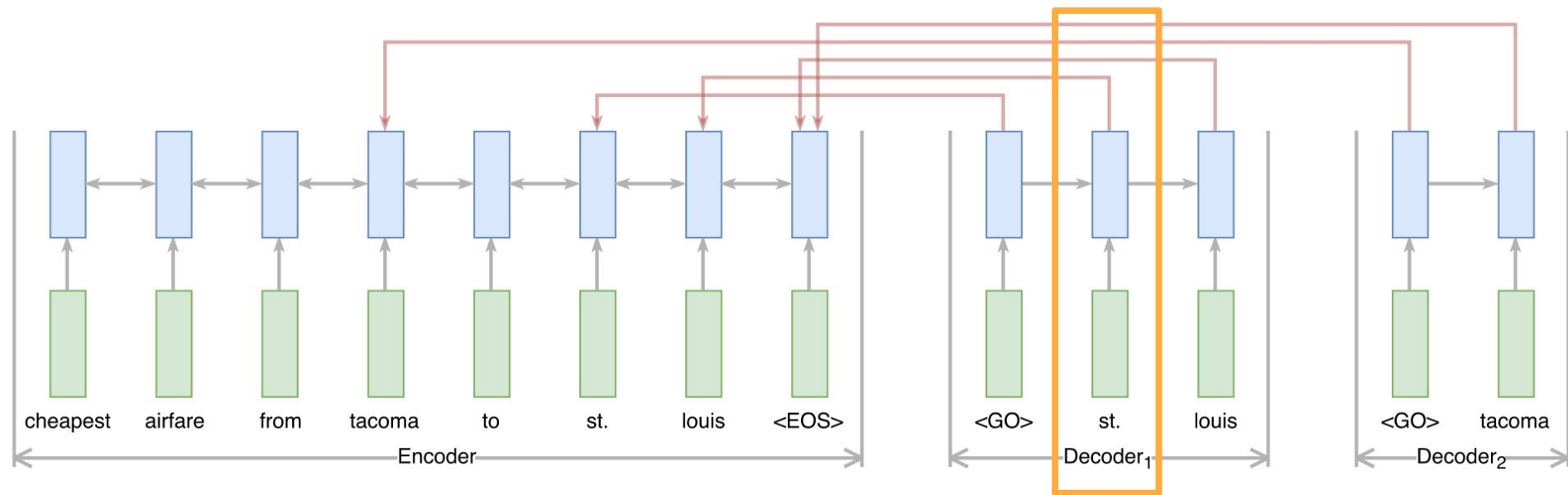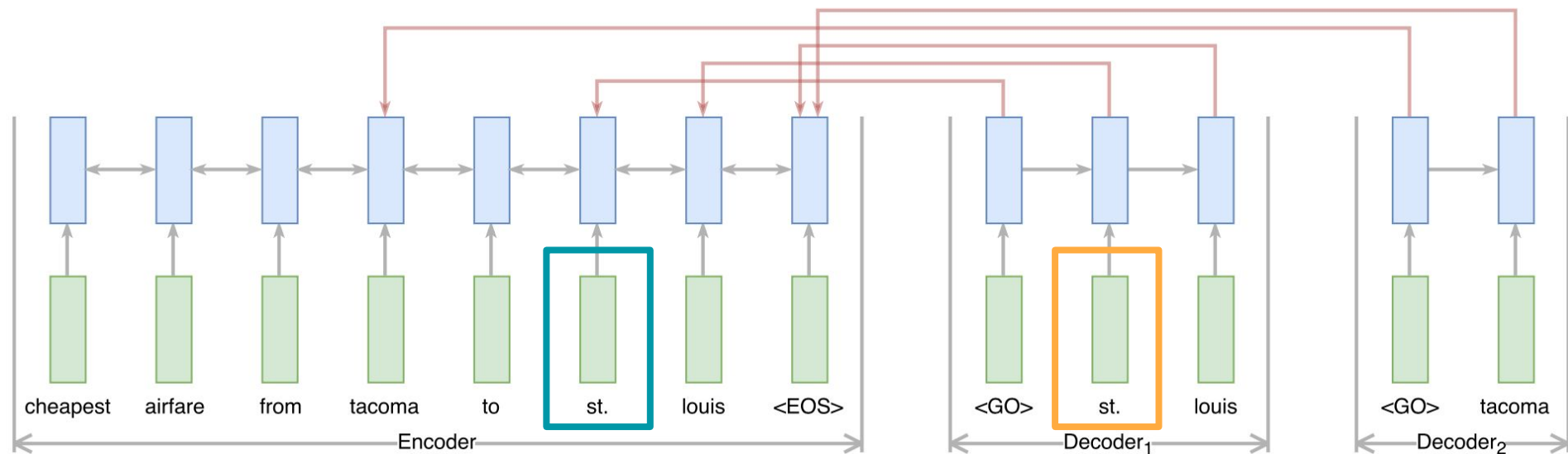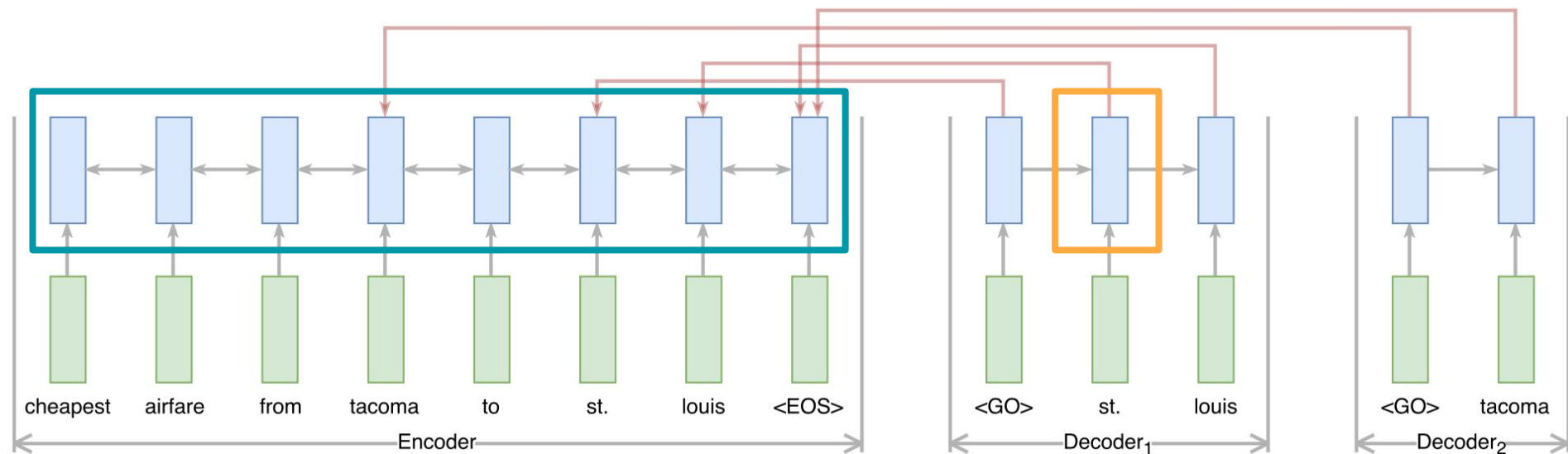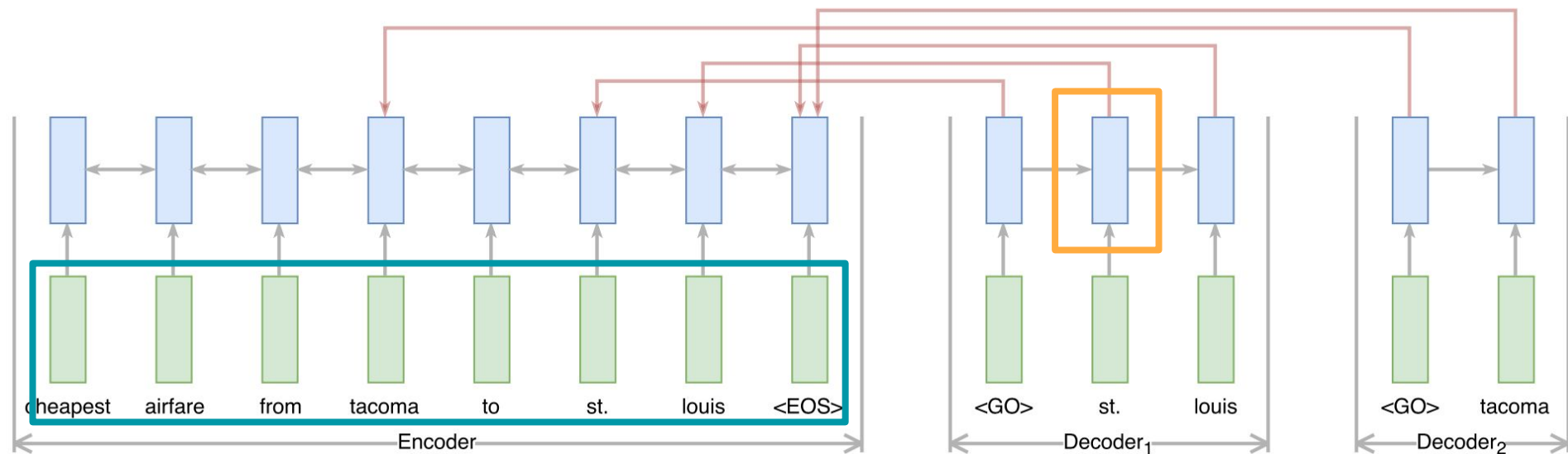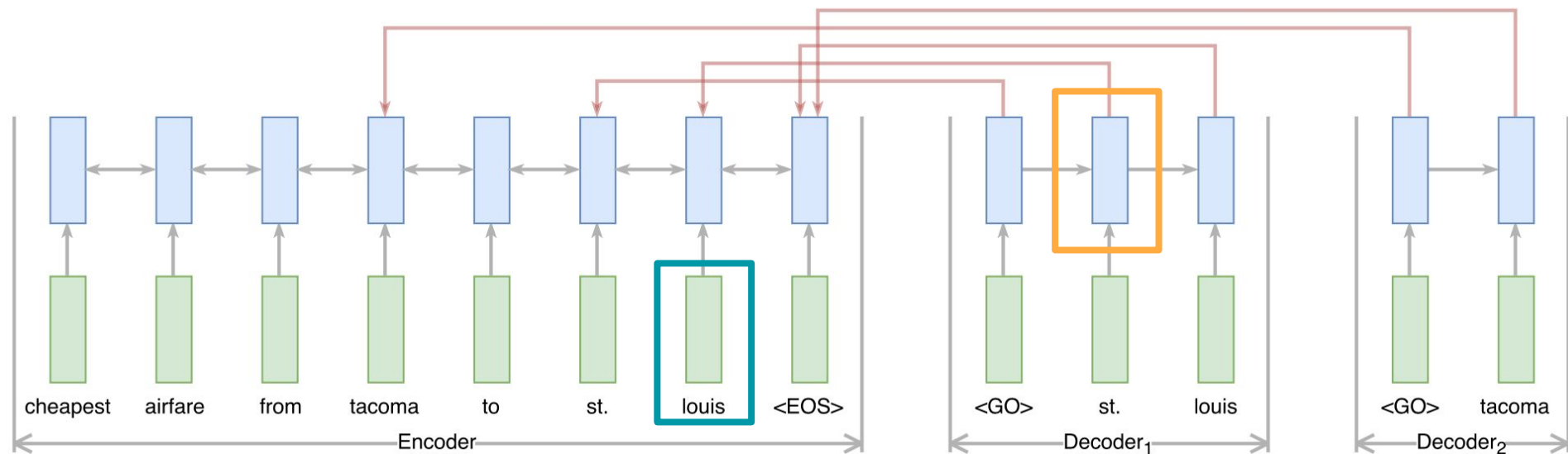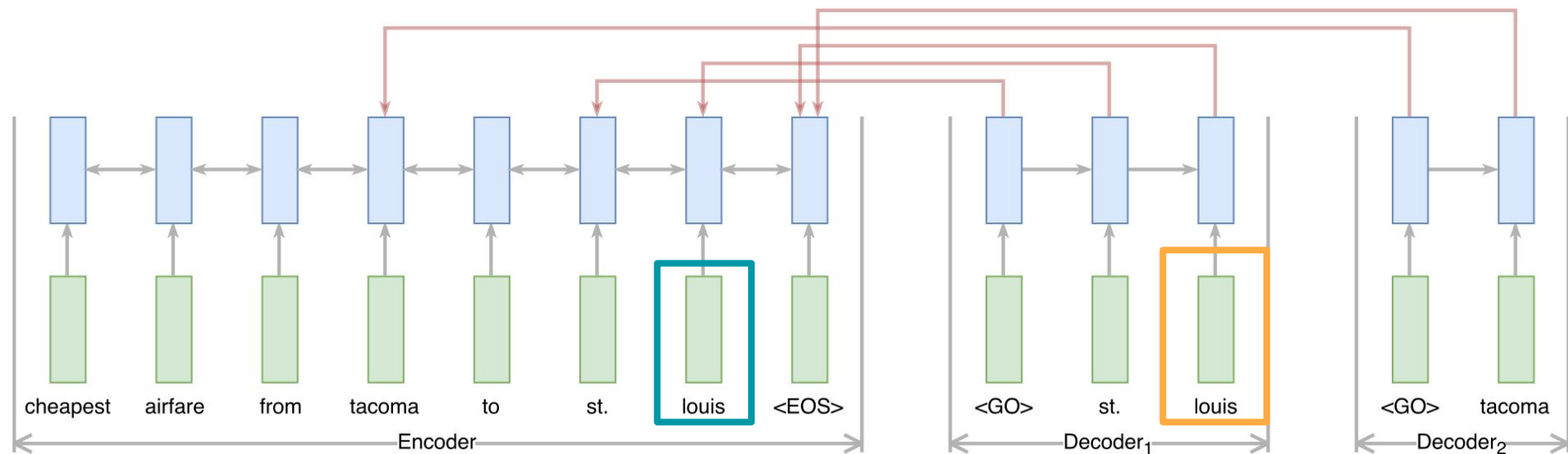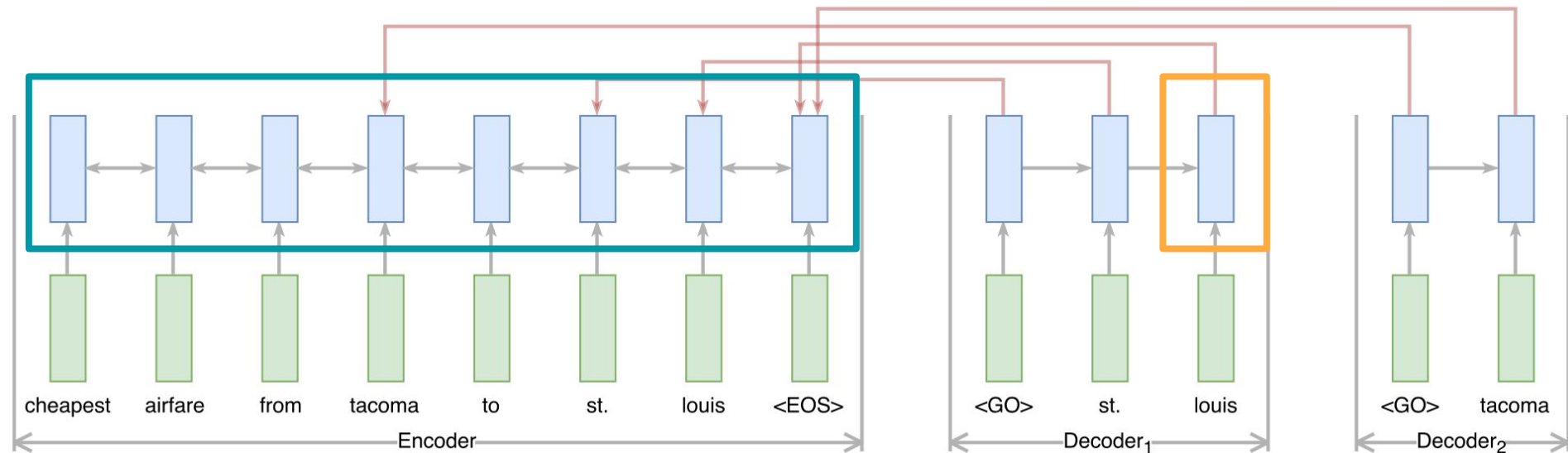
Encoder

<GO> st. louis

Decoder₁

<GO> tacoma

Decoder₂

# Our model

# Our model

# Our model

# Our model

# Our model



Encoder: cheapest airfare from tacoma to st. louis <EOS>

Decoder₁: <GO> st. louis
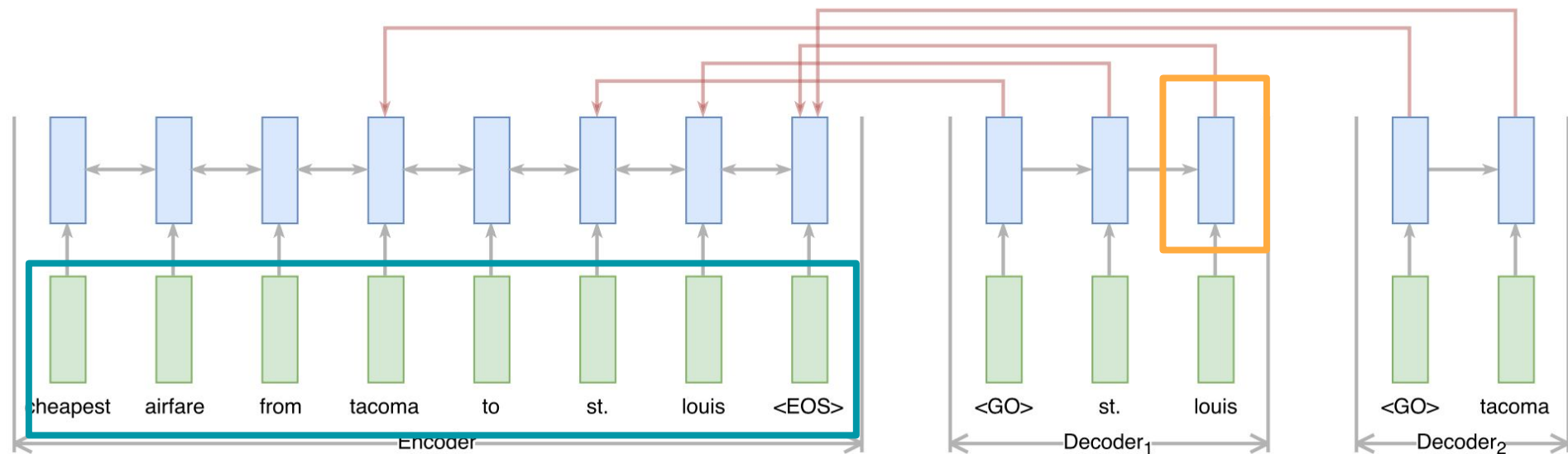
Decoder₂: <GO> tacoma

# Our model

# Our model

# Our model

# Our model



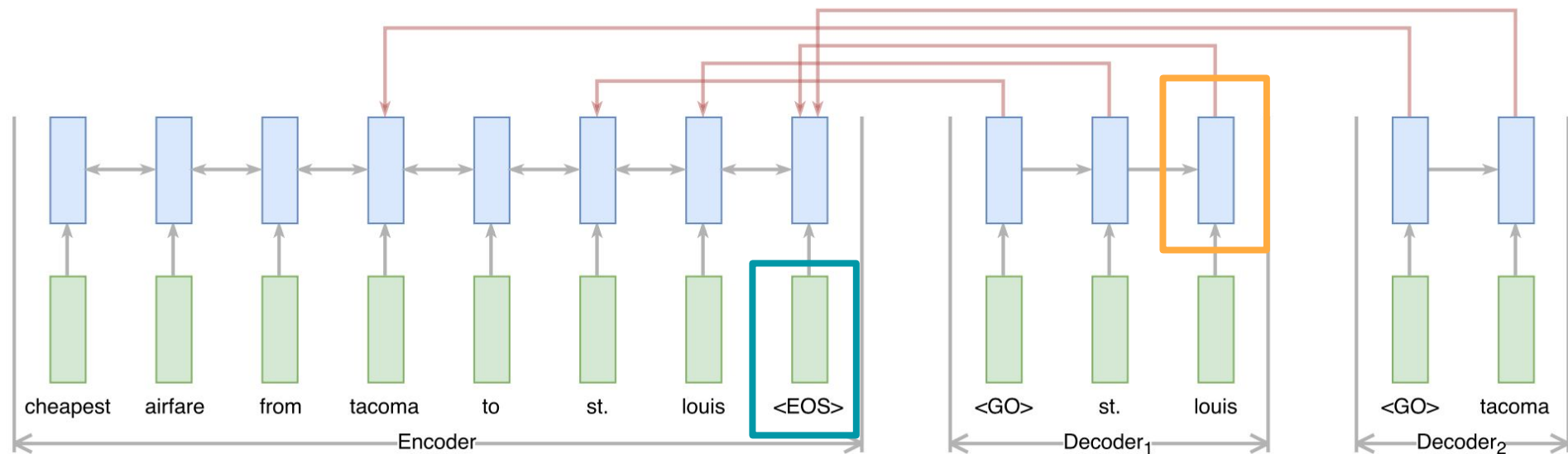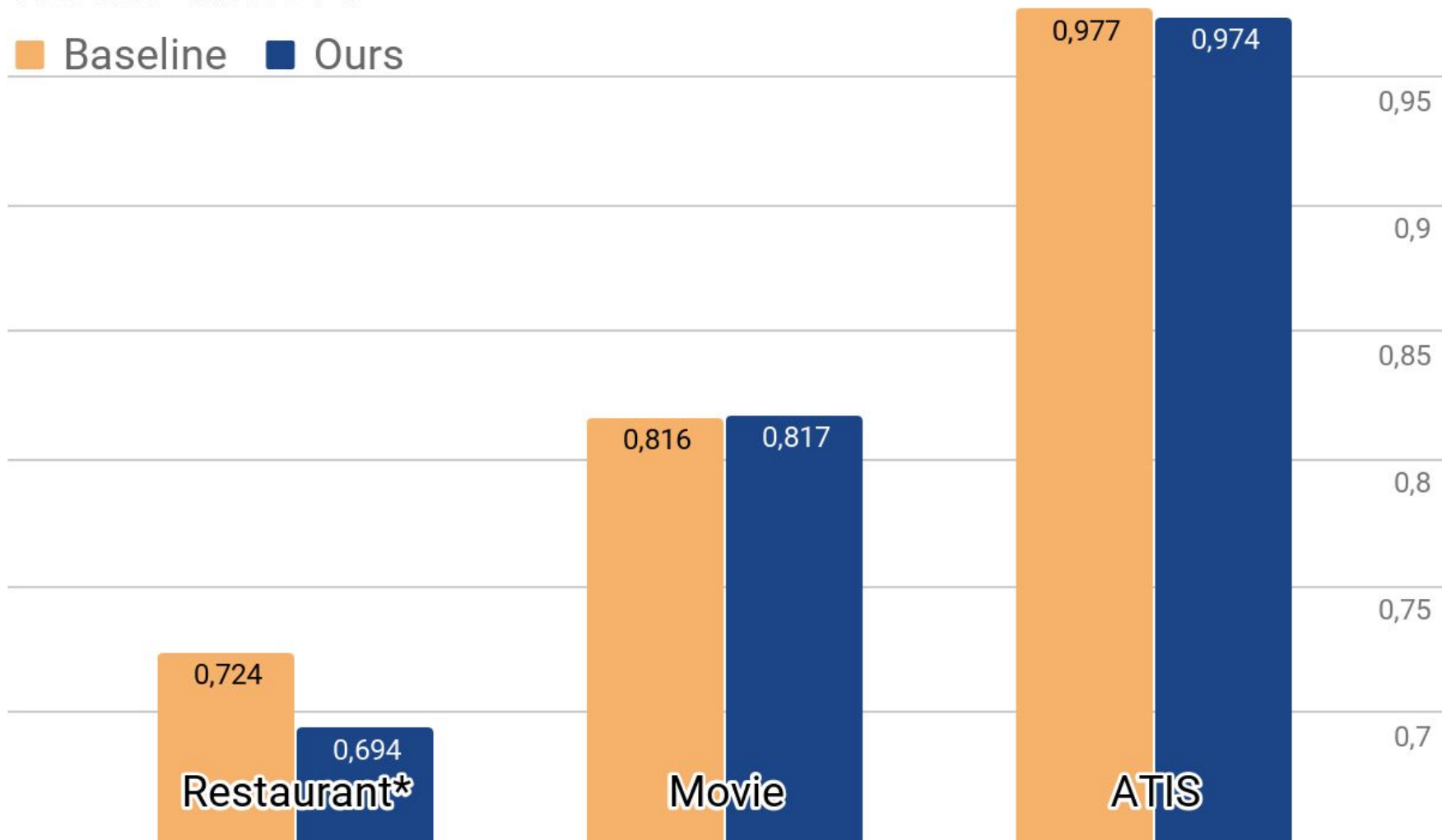Encoder: cheapest airfare from tacoma to st. louis <EOS>

Decoder₁: <GO> st. louis

Decoder₂: <GO> tacoma

# Our model

# Our model

Code is available at

github.com/rasmusbergpalm/e2e-ie-release

# Results - Micro F1

Legend: ■ Baseline ■ Ours

| Dataset | Baseline | Ours |
|---|---|---|
| Restaurant* | 0,724 | 0,694 |
| Movie | 0,816 | 0,817 |
| ATIS | 0,977 | 0,974 |

Axis values: 0,7 · 0,75 · 0,8 · 0,85 · 0,9 · 0,95

There's one major limitation

Normalization

'17 Jan 2012' → '2012-01-17'