

# Adapting Natural Language Processing to Source Code Processing: Handling Syntactic Structure and Identifiers

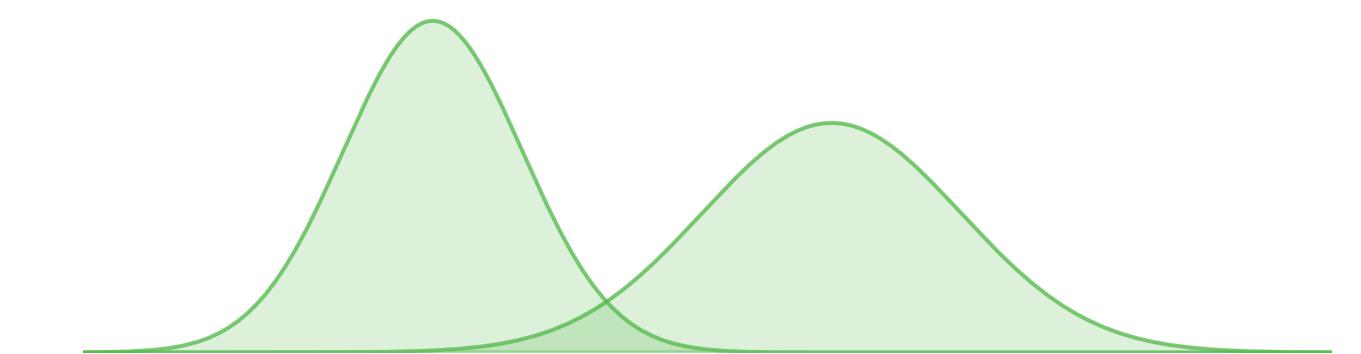
Nadia Chirkova

Higher School of Economics, Samsung-HSE Laboratory  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



$p(\mathbf{B}|\mathbf{A})$ yesgroup.ru

# Deep learning for source code

## Code as input

Predicting bugs  
(code → bug location)  
(code → bug fix)

Code summarization  
(code → comment)  
(code → function name)

Commit message  
generation  
(code changes →  
commit message)

## Code as output

Code completion  
(code → next token(s))

Semantic parsing  
(text → code)

Code retrieval  
(text → code from a  
database)

## Code as both input and output

Code translation  
(Java code →  
Python code)  
(updating code for a  
new library version)

# What type of data is source code?

```
1 def masked_softmax(x, mask=None):
2     if mask is not None:
3         mask = mask.float()
4     if mask is not None:
5         x_masked = x.masked_fill(1-mask, -1e18)
6     else:
7         x_masked = x
8     x_max = x_masked.max(1)[0]
9     x_exp = (x - x_max.unsqueeze(-1)).exp()
10    if mask is not None:
11        x_exp = x_exp * mask.float()
12    return x_exp / x_exp.sum(1).unsqueeze(-1)
```

# What type of data is source code?

Text?

```
1 def masked_softmax(x, mask=None):
2     if mask is not None:
3         mask = mask.float()
4     if mask is not None:
5         x_masked = x.masked_fill(1-mask, -1e18)
6     else:
7         x_masked = x
8     x_max = x_masked.max(1)[0]
9     x_exp = (x - x_max.unsqueeze(-1)).exp()
10    if mask is not None:
11        x_exp = x_exp * mask.float()
12    return x_exp / x_exp.sum(1).unsqueeze(-1)
```

# What type of data is source code?

Text?

✓ Syntactic structure

✓ Invariance to changing  
identifier names

(and other  
specifics)

```
1 def masked_softmax(x, mask=None):
2     if mask is not None:
3         mask = mask.float()
4     if mask is not None:
5         x_masked = x.masked_fill(1-mask, -1e18)
6     else:
7         x_masked = x
8     x_max = x_masked.max(1)[0]
9     x_exp = (x - x_max.unsqueeze(-1)).exp()
10    if mask is not None:
11        x_exp = x_exp * mask.float()
12    return x_exp / x_exp.sum(1).unsqueeze(-1)
```

# Plan

- Empirical study of syntax-based Transformers for source code
- Handling out-of-vocabulary identifiers in Transformers for source code
- Handling anonymized identifiers in RNNs for source code

# Plan

- Empirical study of syntax-based Transformers for source code
- Handling out-of-vocabulary identifiers in Transformers for source code
- Handling anonymized identifiers in RNNs for source code

# Empirical Study of Transformers for Source Code

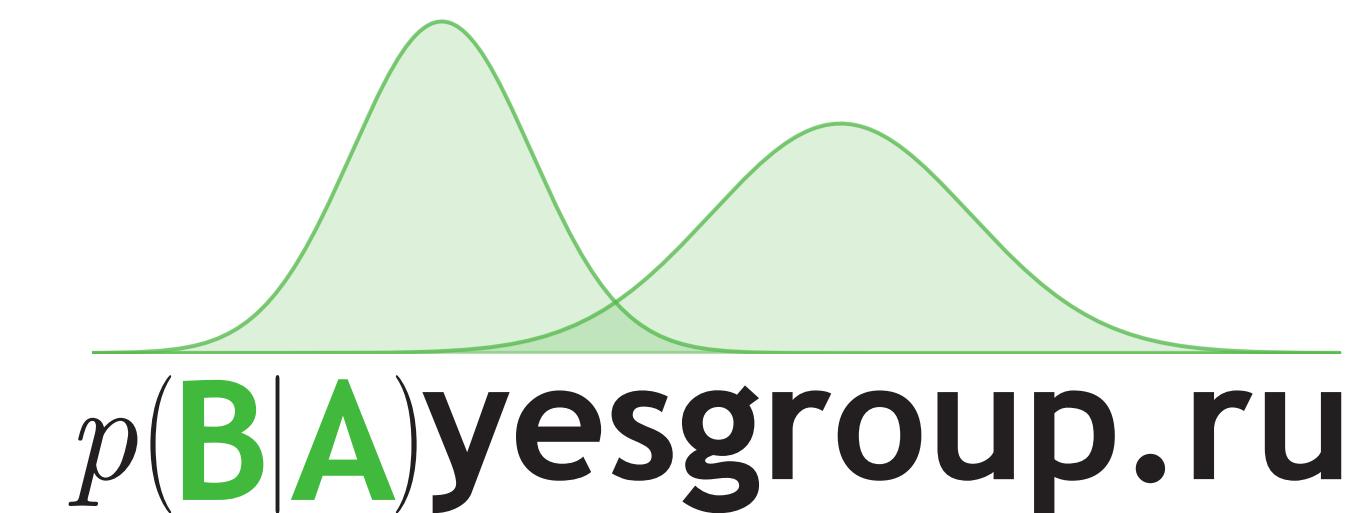
Nadezhda Chirkova, Sergey Troshin

Higher School of Economics, Samsung-HSE Laboratory  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



# Motivation

- A lot of recent works adjust Transformer for handling the syntactic structure of the input:

Which one performs better?

Model	Task
Tree positional encodings (Shiv19)	program translation semantic parsing
Tree relative attention (Kim20)	code competition
GGNN-Sandwich (Hellendoorn20)	bug fixing
...	...

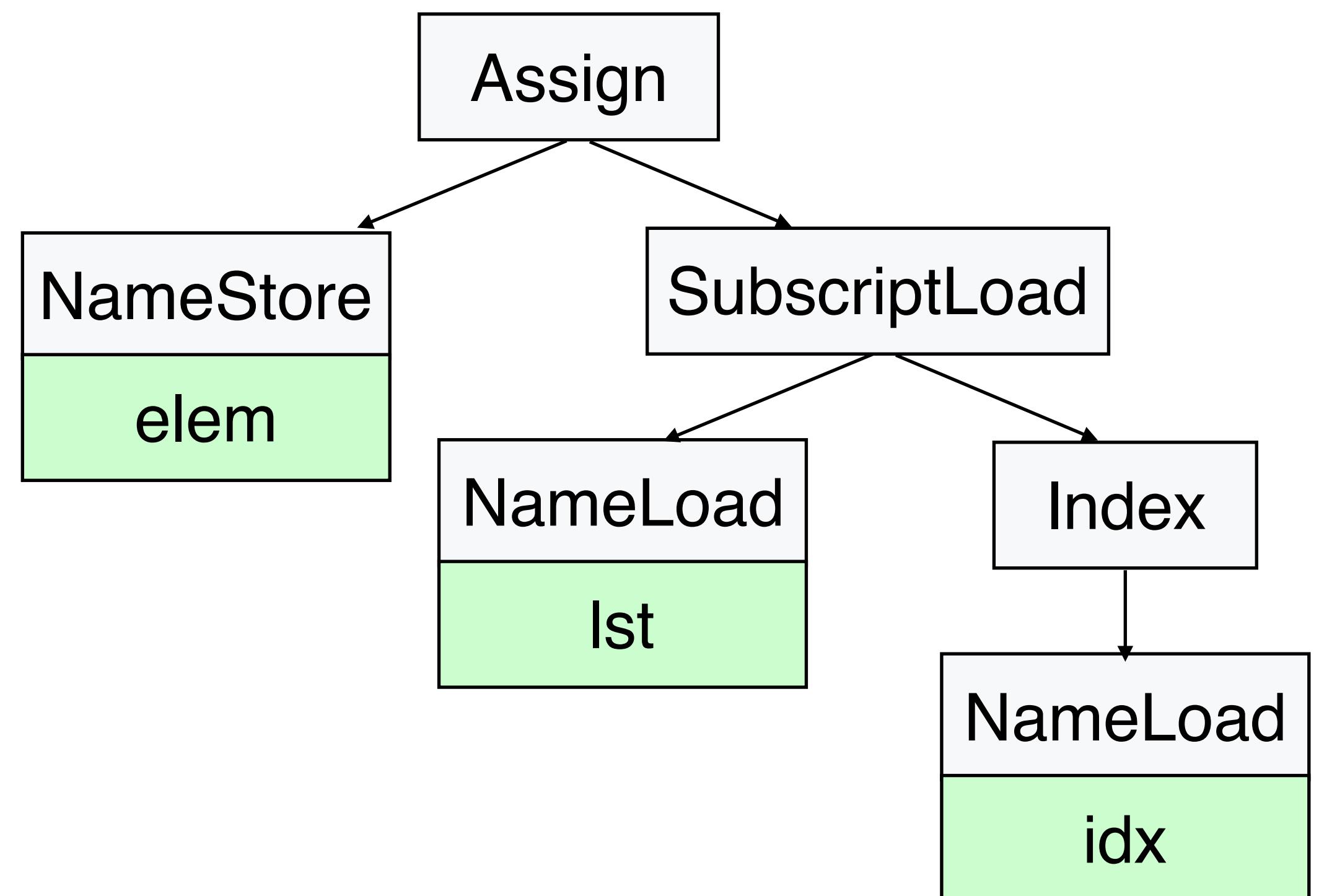
# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
- 5 syntax-capturing Transformer modifications
  - 3 tasks: code completion, bug fixing, function naming
  - 2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework (e. g. same data preprocessing for all techniques)

# Abstract syntax tree (AST)

Code:      `elem = lst[idx]`

Abstract syntax  
tree (AST):

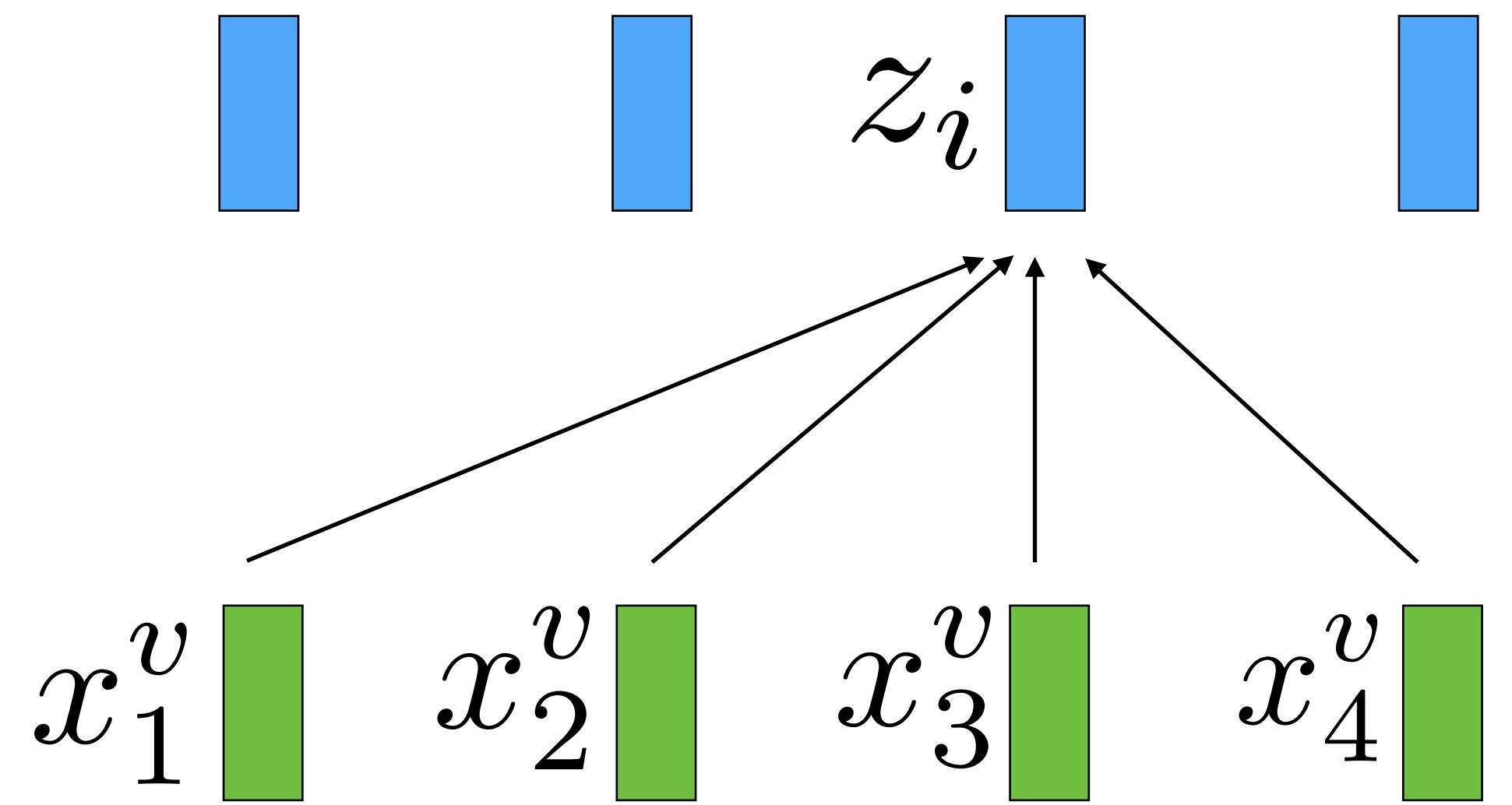


✓ Type in each node

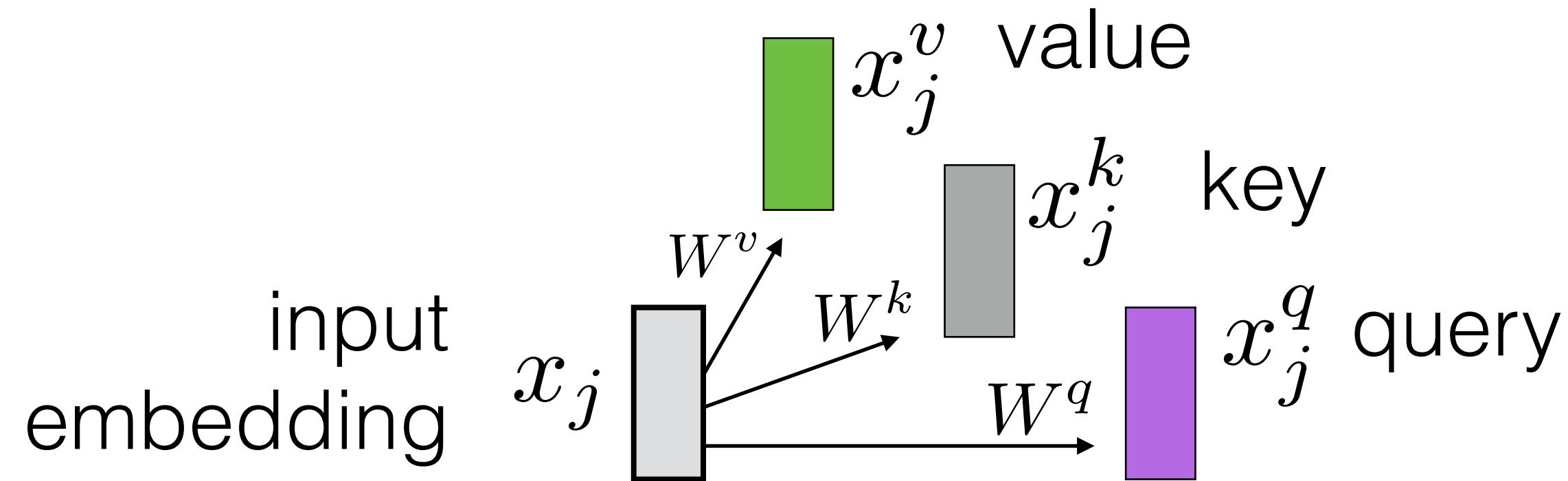
✓ Values in some nodes  
(identifiers, constants,  
function names etc.)

✓ Ordered children

# Self-attention mechanism



Some random short sequence

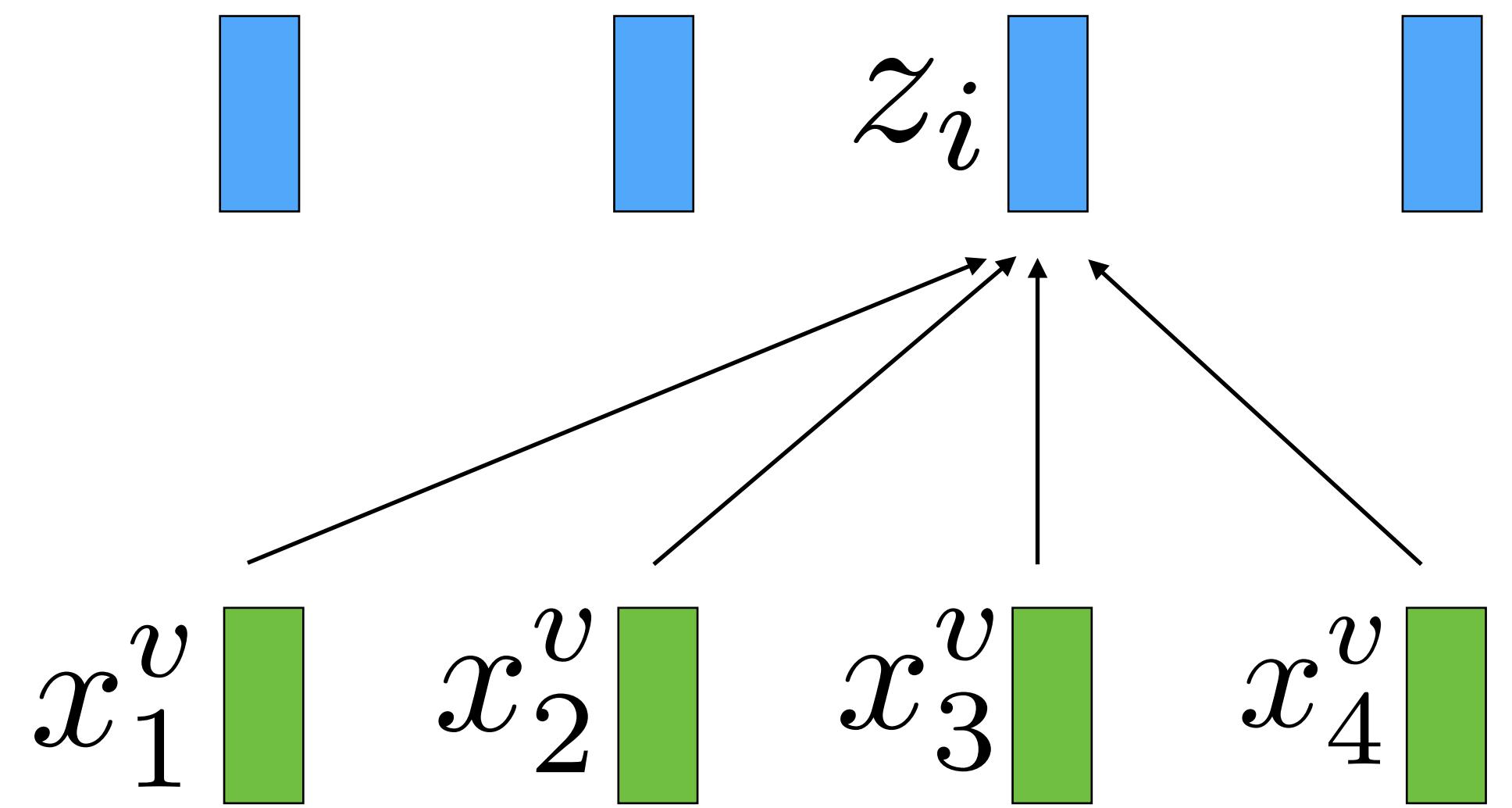


$$z_i = \sum_j \tilde{\alpha}_{ij} x_j^v$$

$$\tilde{\alpha}_{ij} = \frac{\exp(a_{ij})}{\sum_j \exp(a_{ij})}$$

$$a_{ij} = \frac{x_i^q x_j^k {}^T}{\sqrt{d_z}}$$

# Self-attention mechanism



Some random short sequence

+multiple heads

+skip-connection,  
layer normalization,  
fully-connected layer

+ 6 layers

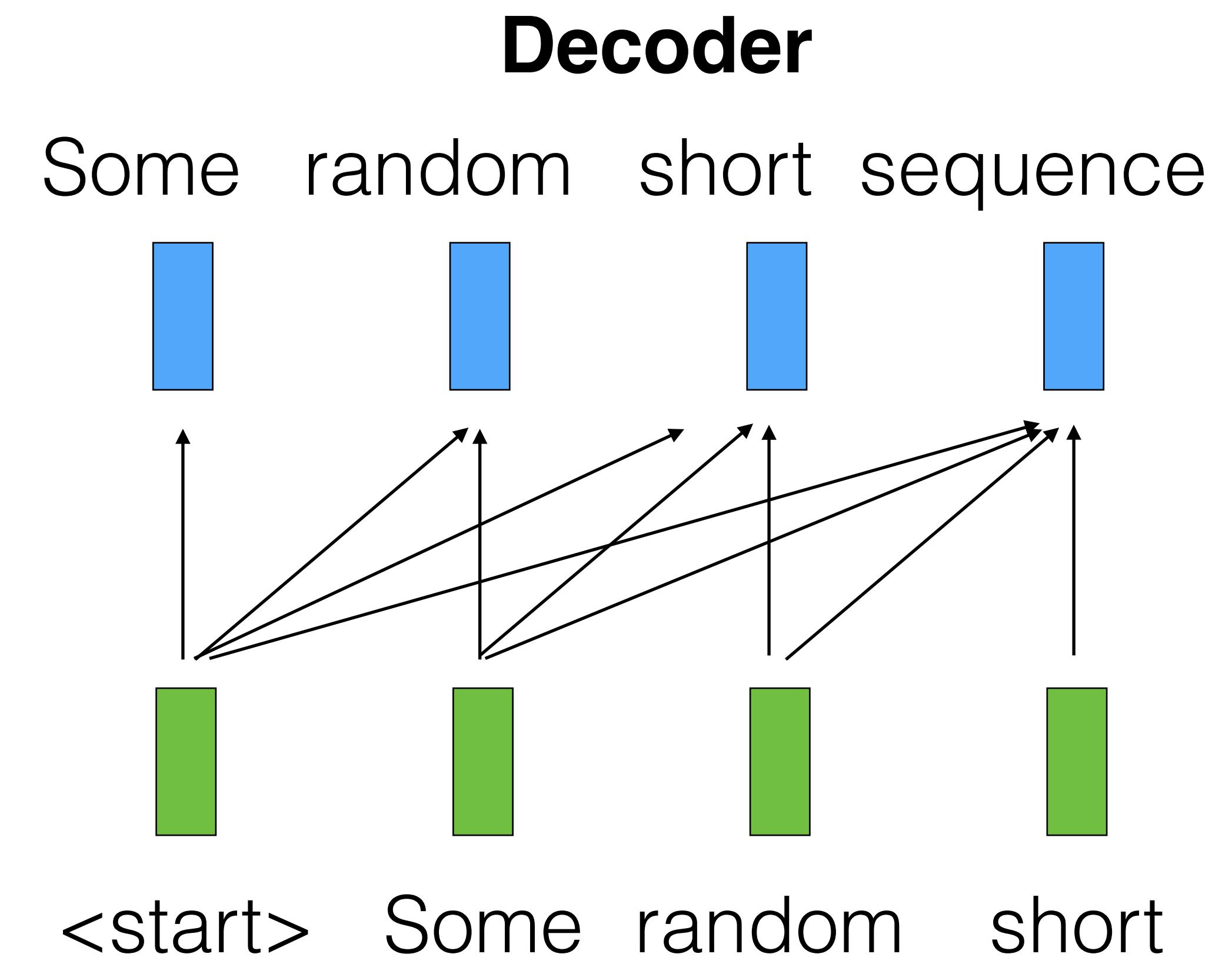
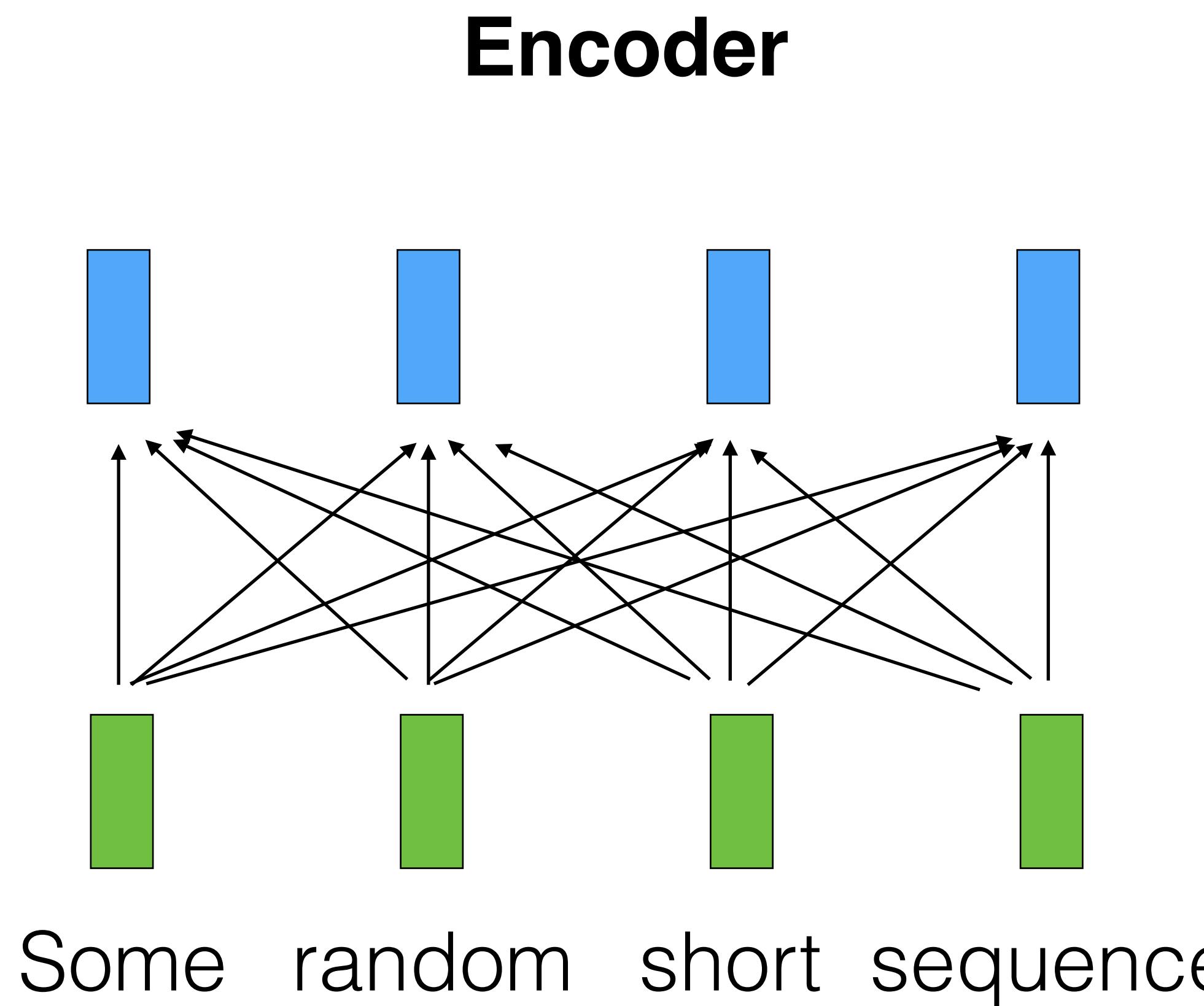
+encoder /  
decoder

$$z_i = \sum_j \tilde{\alpha}_{ij} x_j^v$$

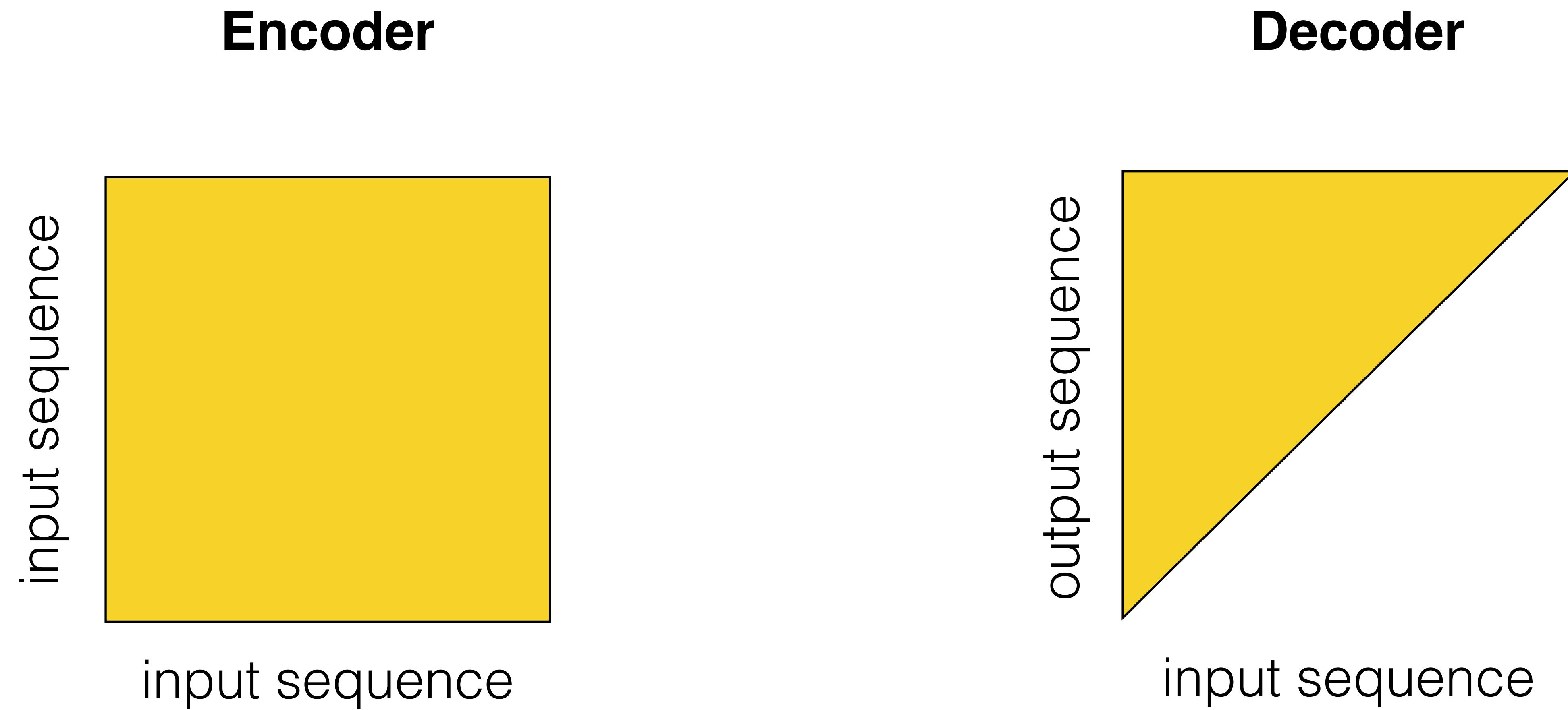
$$\tilde{\alpha}_{ij} = \frac{\exp(a_{ij})}{\sum_j \exp(a_{ij})}$$

$$a_{ij} = \frac{x_i^q x_j^k{}^T}{\sqrt{d_z}}$$

# Encoder and Decoder in Transformer



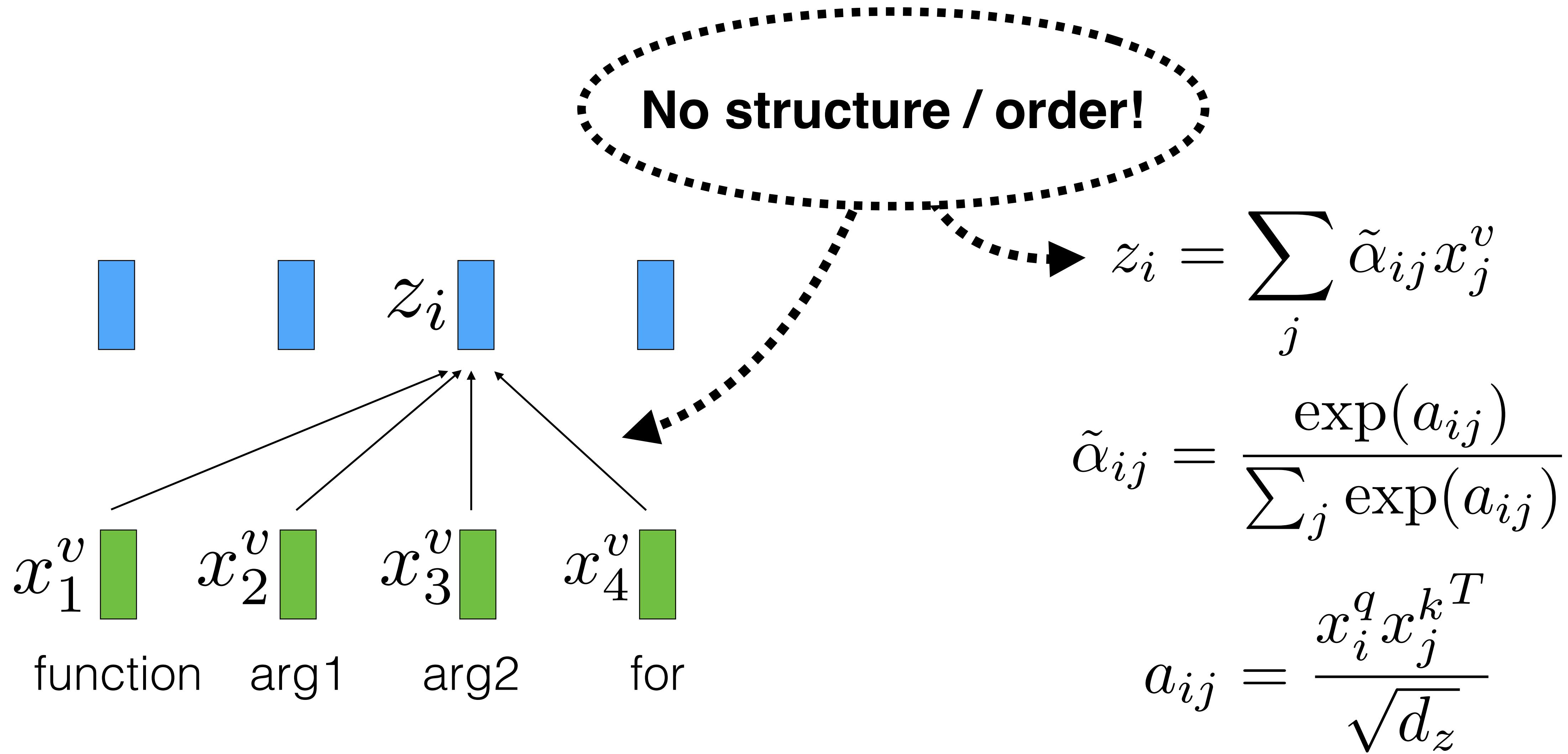
# Encoder and Decoder: attention masks



# Transformer: GIF

<https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>

# Self-attention: “bag” of elements



# AST depth-first traversal

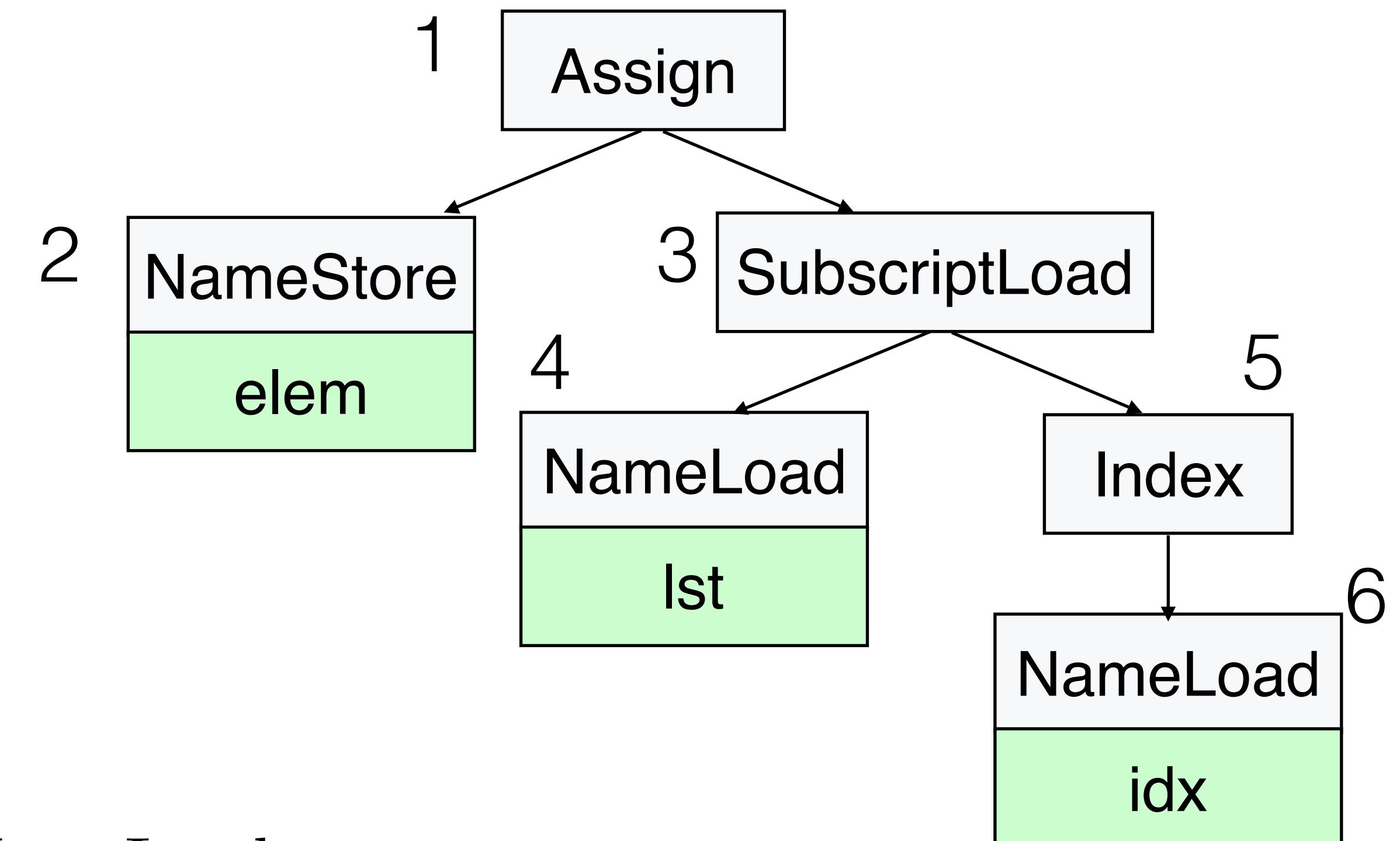
Code: `elem = lst[idx]`



AST depth-first traversal:

Assign	NameStore	SubscriptLoad	...
<code>&lt;empty&gt;</code>	elem	<code>&lt;empty&gt;</code>	...
...	NameLoad	Index	NameLoad
...	lst	<code>&lt;empty&gt;</code>	idx

Abstract syntax tree (AST):



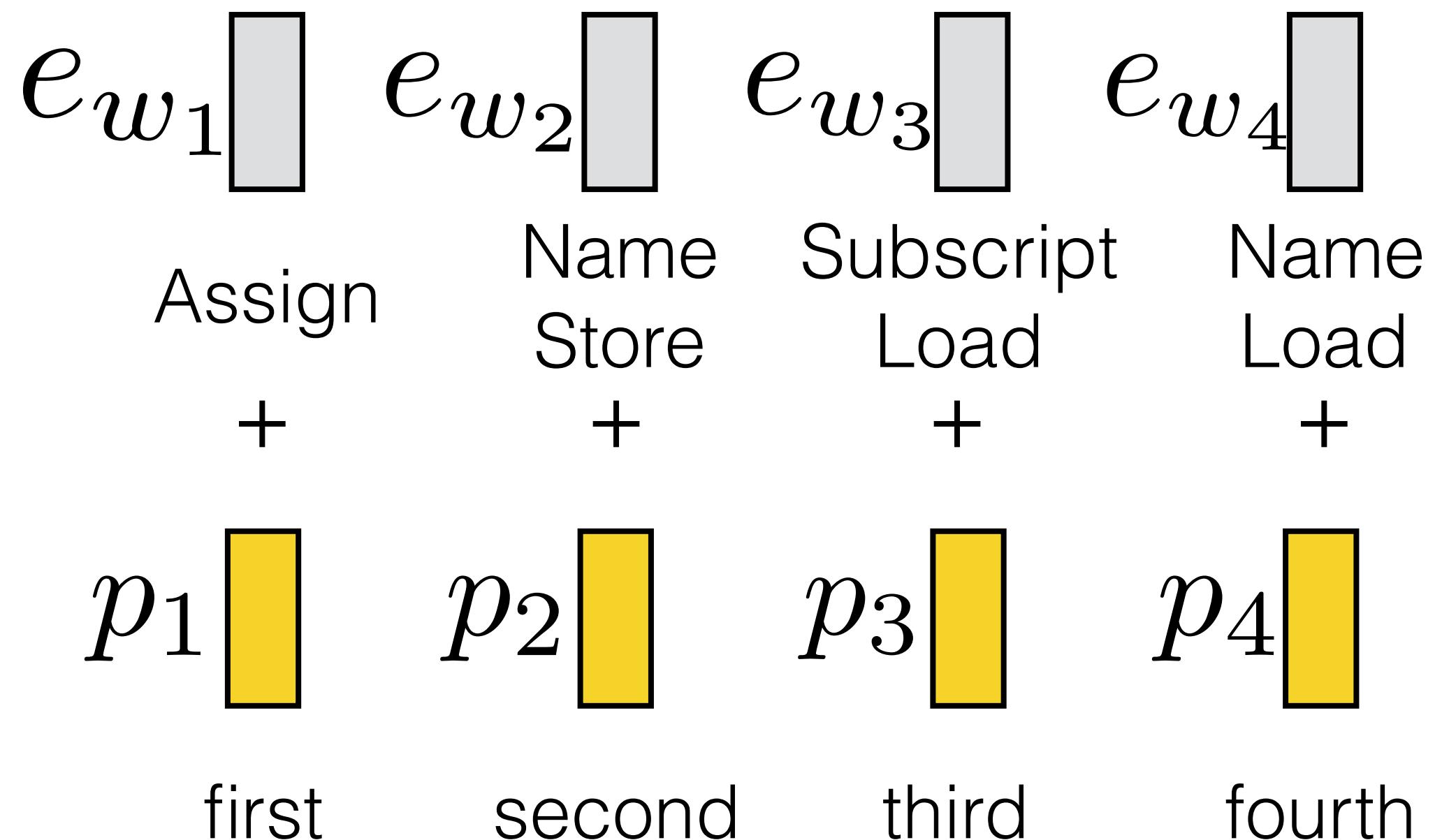
# Sequential positional encodings / embeddings

Input layer:

$$x_j = e_{w_j} + p_j$$

Positional embeddings:

$$p_j = e_j \text{ — learnable embeddings}$$

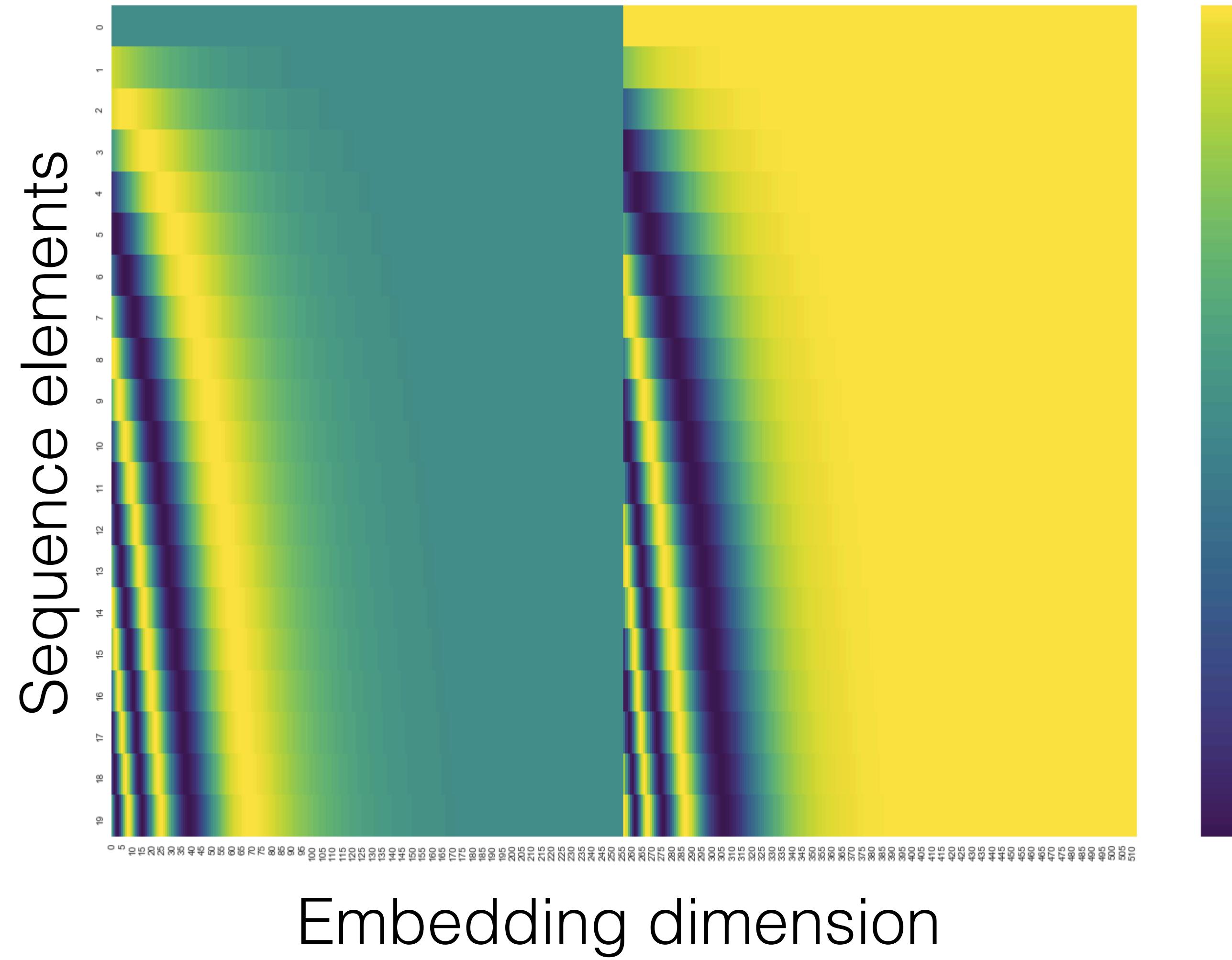


Positional encoding:

$$p_{j,m} = \begin{cases} \sin(\omega_k j) & \text{if } m = 2k \\ \cos(\omega_k j) & \text{if } m = 2k + 1 \end{cases}$$

$$w_k = \frac{1}{10000^{2k/d_x}}$$

# Sequential positional encodings: properties



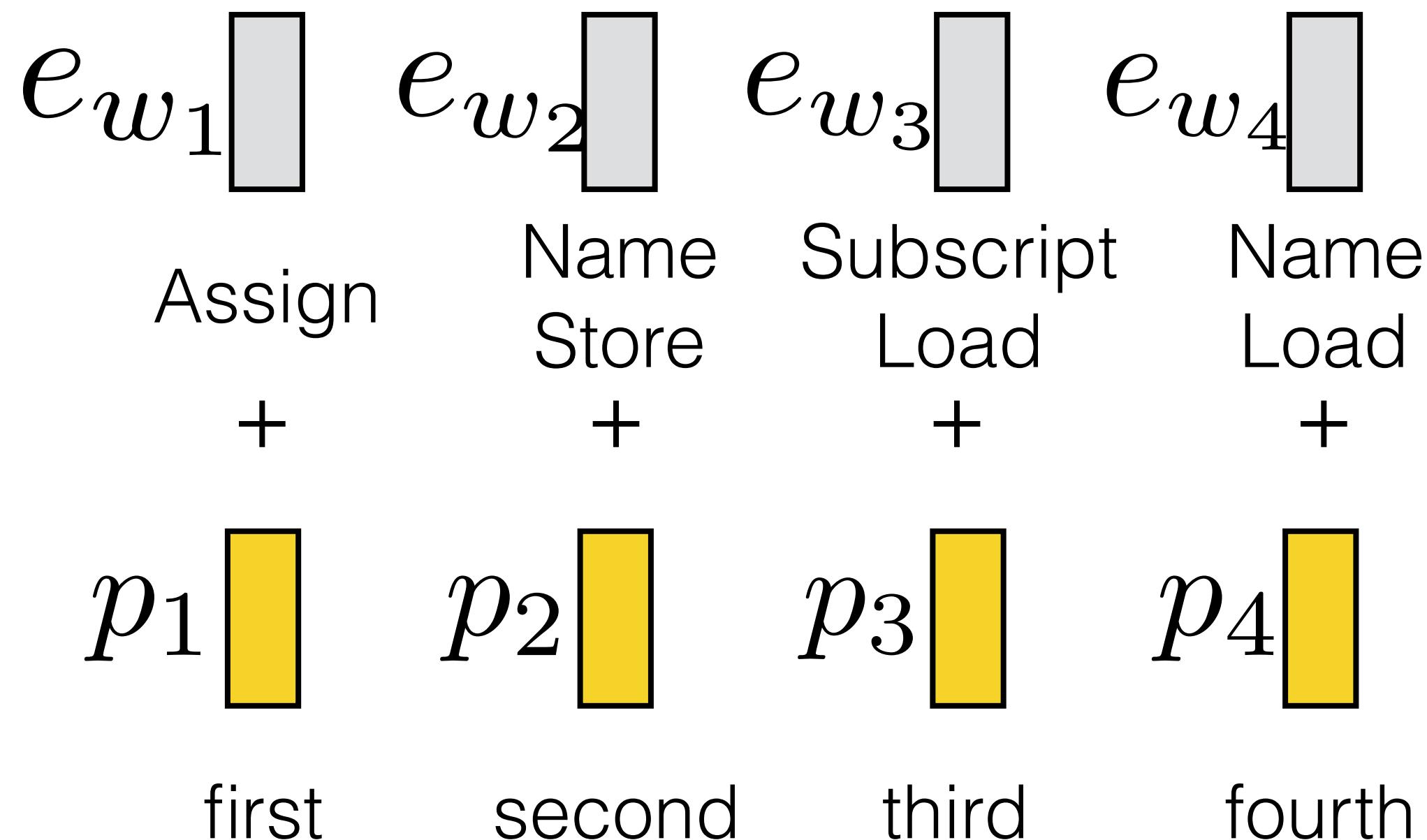
If  $j_1$  and  $j_2$  are two positions  
with distance  $\phi$ ,  
then their encodings satisfy

$$p_{j_1} = A_\phi p_{j_2}$$

# Sequential positional encodings / embeddings

Input layer:

$$x_j = e_{w_j} + p_j$$



Positional embeddings:

Parameters: learnable embeddings

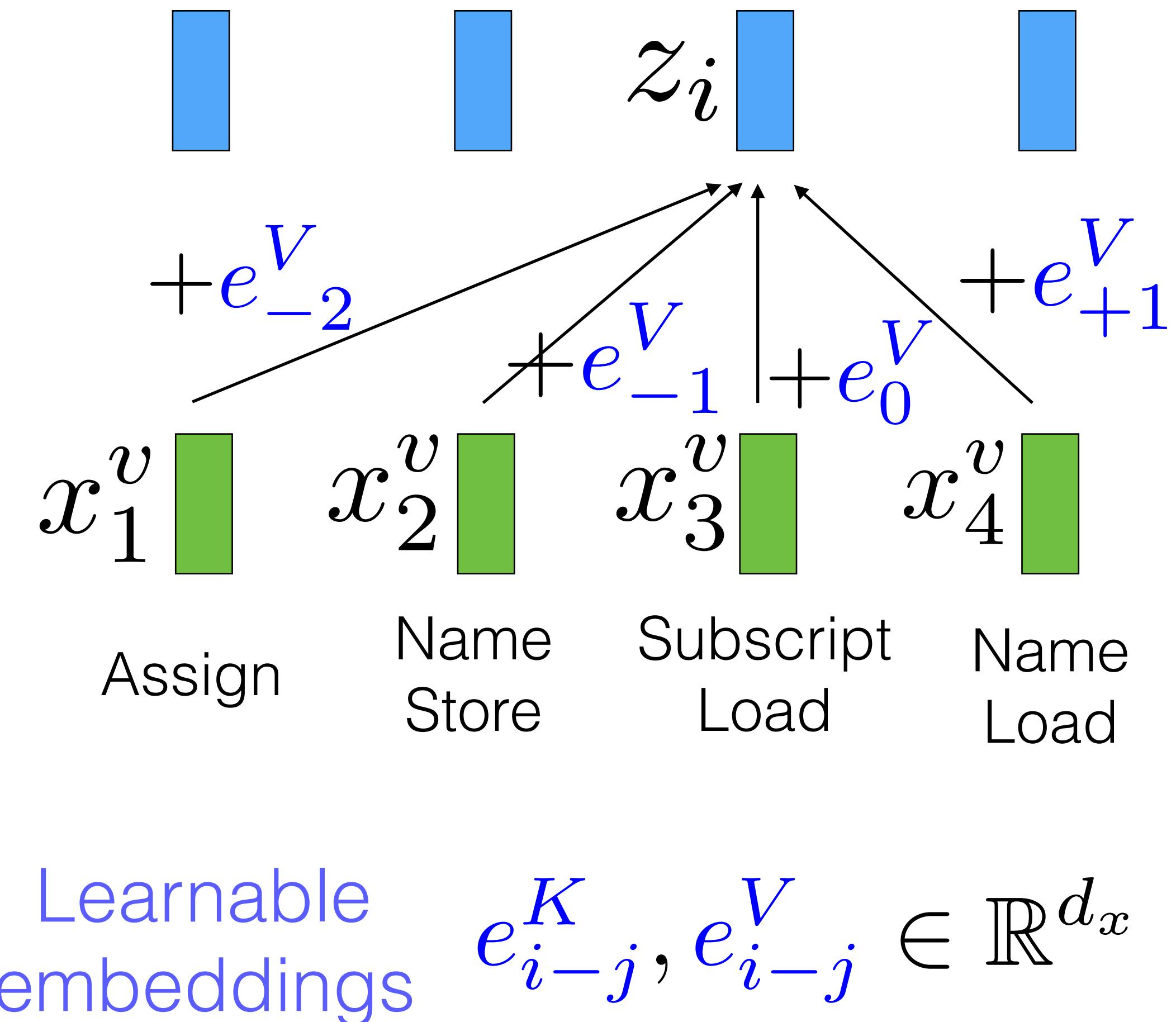
Hyperparameters: none

Positional encoding:

Parameters: none

Hyperparameters: none

# Sequential relative attention

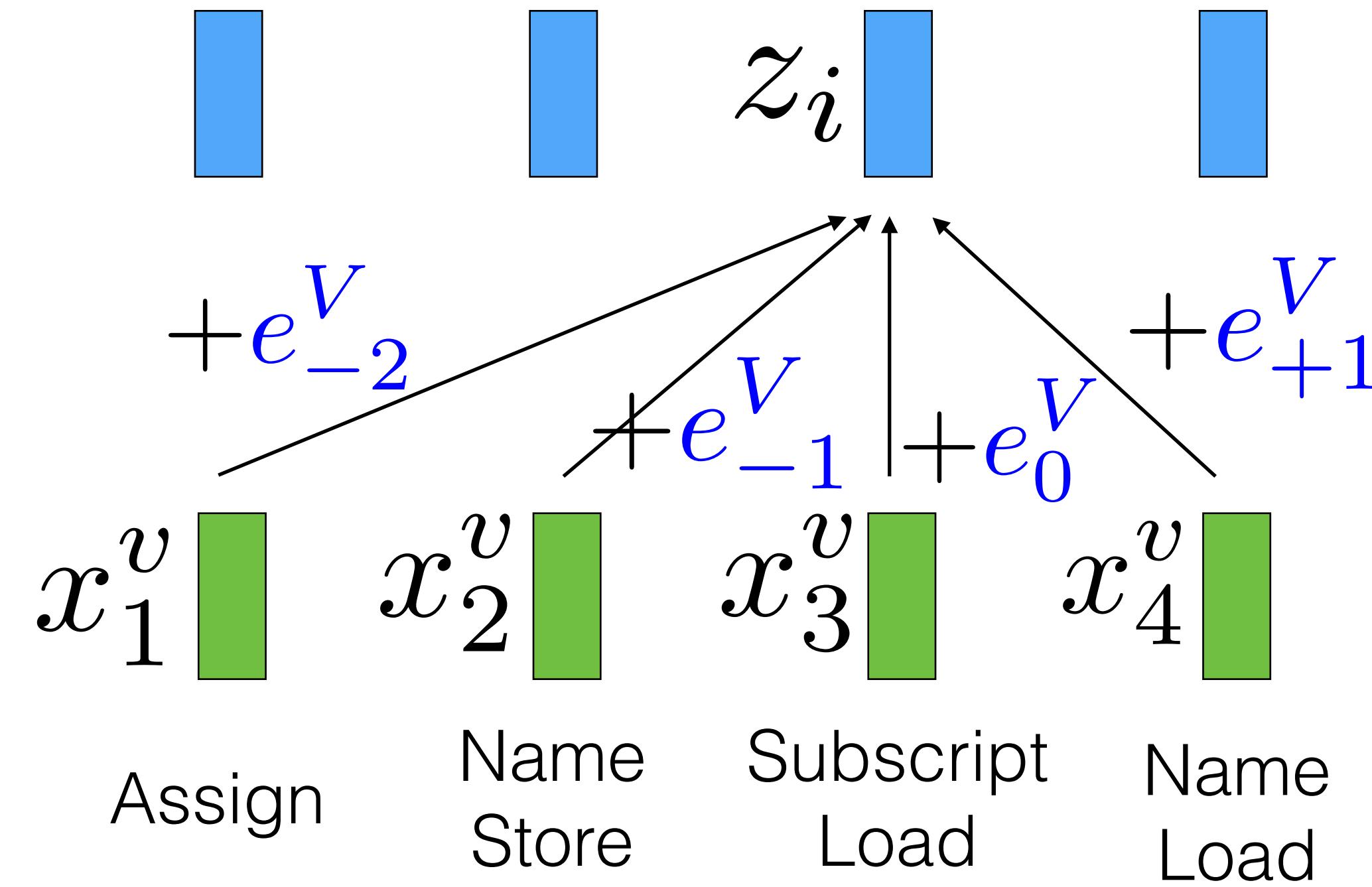


$$z_i = \sum_j \tilde{\alpha}_{ij} (x_j^v + e_{i-j}^V)$$

$$\tilde{\alpha}_{ij} = \frac{\exp(a_{ij})}{\sum_j \exp(a_{ij})}$$

$$a_{ij} = \frac{x_i^q (x_j^k + e_{i-j}^K)^T}{\sqrt{d_z}}$$

# Sequential relative attention



Learnable embeddings  $e_{i-j}^K, e_{i-j}^V \in \mathbb{R}^{d_x}$

Parameters:  
embeddings of “relations”

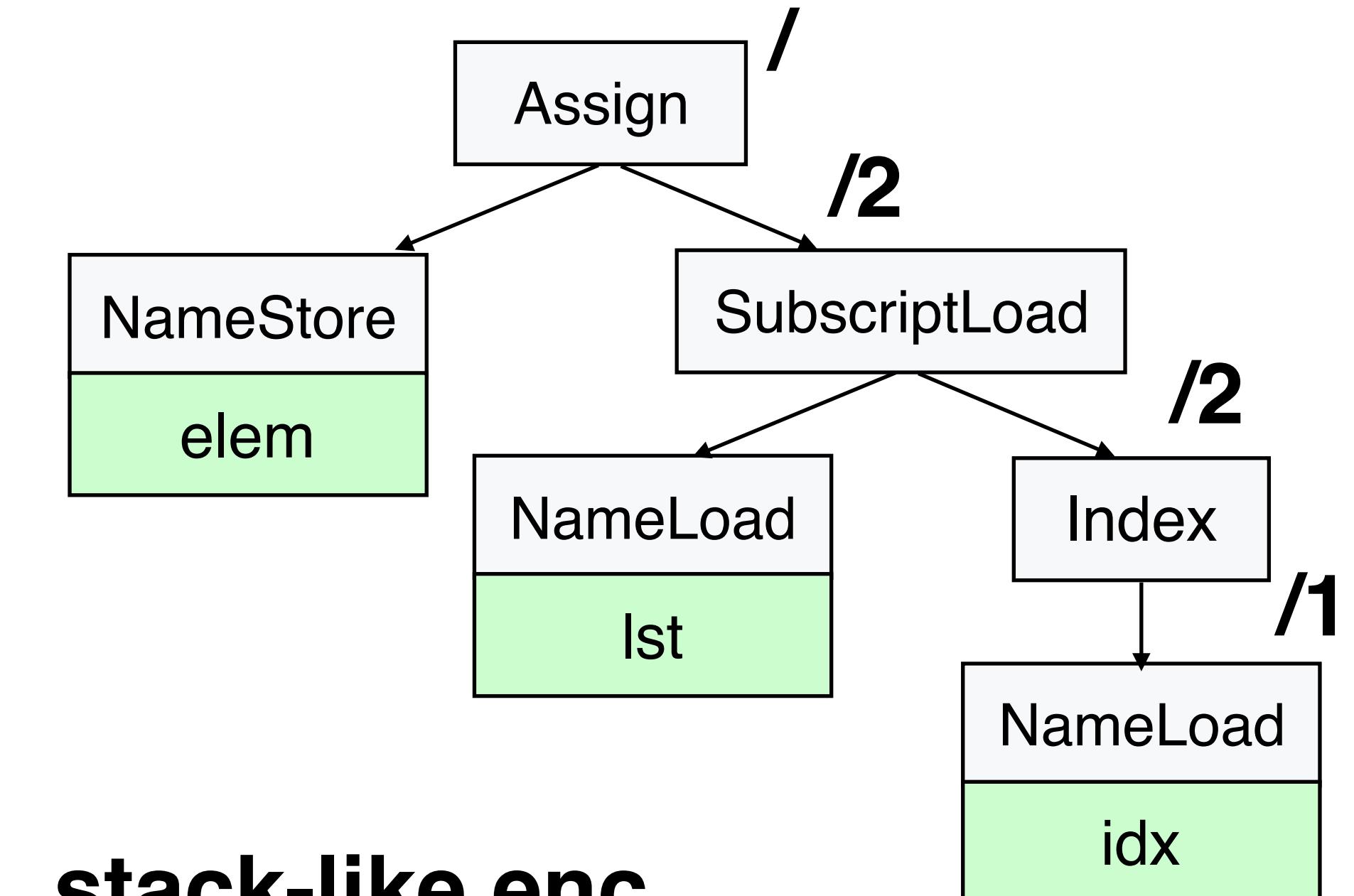
Hyperparameters:  
maximum distance between tokens

# Tree positional encoding

Input layer:

$$x_j = e_{w_j} + p_{\text{path-to-}j}$$

$e_{w_1}$	$e_{w_2}$	$e_{w_3}$	$e_{w_4}$
Assign	Name Store	Subscript Load	Name Load
+	+	+	+
$p_1$	$p_2$	$p_3$	$p_4$
/	/1	/2	/2/1



**stack-like enc.**

/	000 000 000
/1	100 000 000
/2	010 000 000
/2/1	100 010 000
/2/2	010 010 000
/2/2/1	100 010 010

# Tree positional encoding

Input layer:

$$x_j = e_{w_j} + p_{\text{path-to-}j}$$

If  $\text{path}_1$  and  $\text{path}_2$  are two positions with path  $\phi$ , then their encodings satisfy

$$p_{\text{path}_1} = A_\phi p_{\text{path}_2}$$

$e_{w_1}$	$e_{w_2}$	$e_{w_3}$	$e_{w_4}$
Assign	Name Store	Subscript Load	Name Load
+	+	+	+
$p_1$	$p_2$	$p_3$	$p_4$
/	/1	/2	/2/1

	<b>stack-like enc.</b>
/	000 000 000
/1	100 000 000
/2	010 000 000
/2/1	100 010 000
/2/2	010 010 000
/2/2/1	100 010 010

product of  
“up” and “down”  
transformation  
matrices

# Tree positional encoding

Input layer:

$$x_j = e_{w_j} + p_{\text{path-to-}j}$$

Parameters: p's

Hyperparameters:

maximum children count

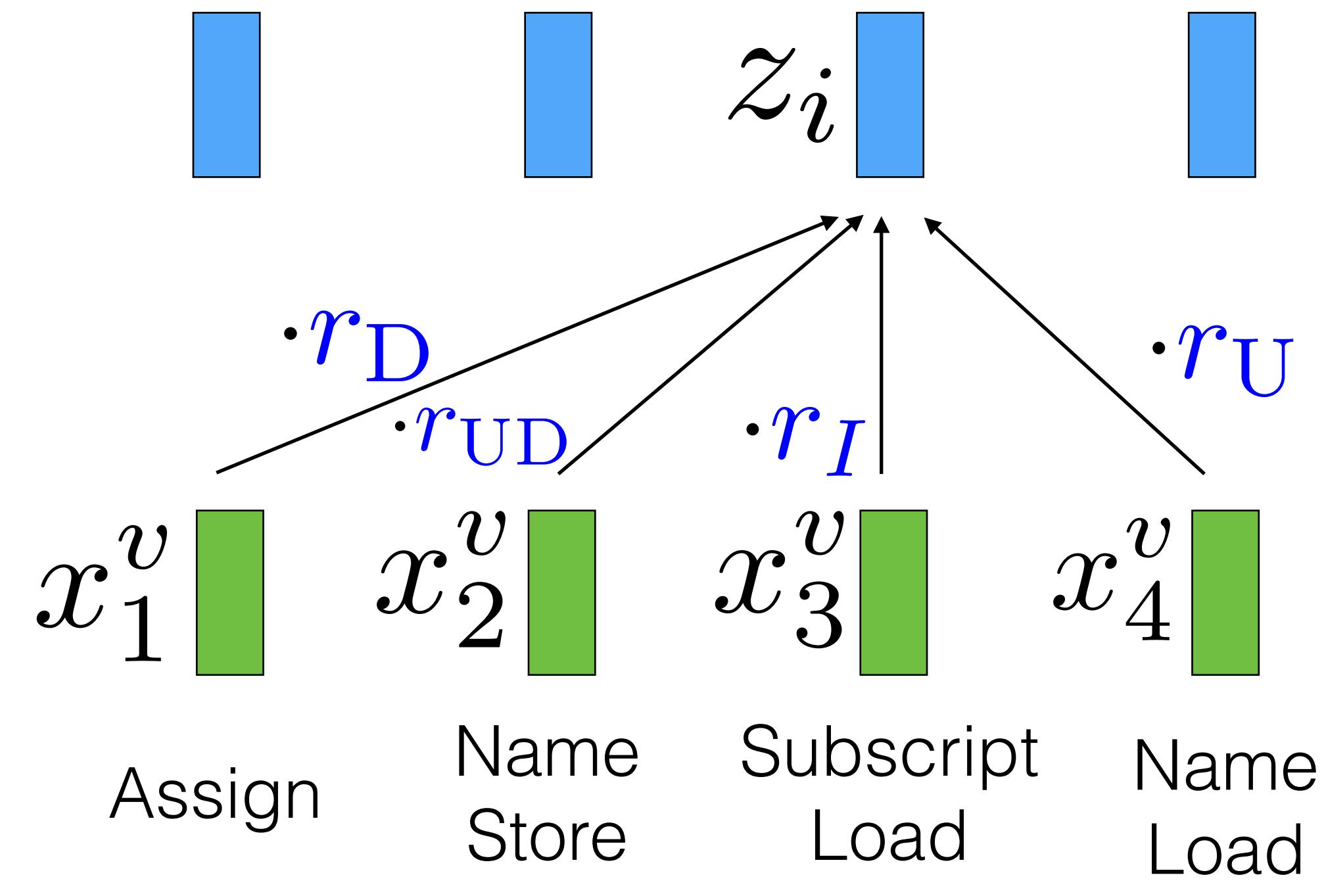
maximum path length

$e_{w_1}$	$e_{w_2}$	$e_{w_3}$	$e_{w_4}$	
Assign	Name Store	Subscript Load	Name Load	
+	+	+	+	
$p_1$	$p_2$	$p_3$	$p_4$	
/	/1	/2	/2/1	

**stack-like enc.**

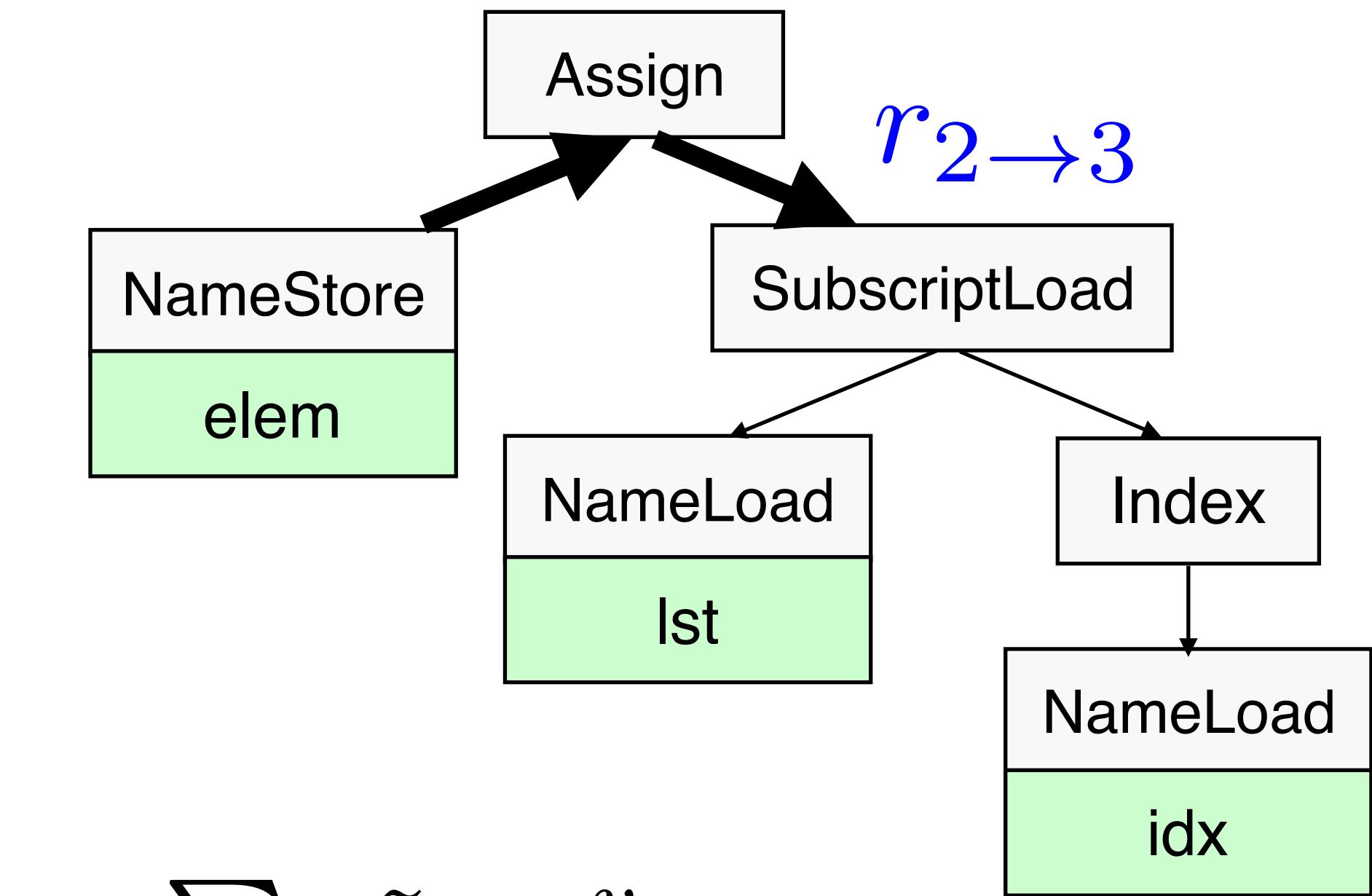
/	000 000 000	$\odot[111 \ p \ p \ p \ p^2 \ p^2 \ p^2]$
/1	100 000 000	
/2	010 000 000	
/2/1	100 010 000	
/2/2	010 010 000	
/2/2/1	100 010 010	$\odot[111 \ p \ p \ p \ p^2 \ p^2 \ p^2]$

# Tree relative attention



Learnable 1-dim  
embeddings

$$r_{j \rightarrow i} \in \mathbb{R}$$



$$z_i = \sum_j \tilde{\alpha}_{ij} x_j^v$$

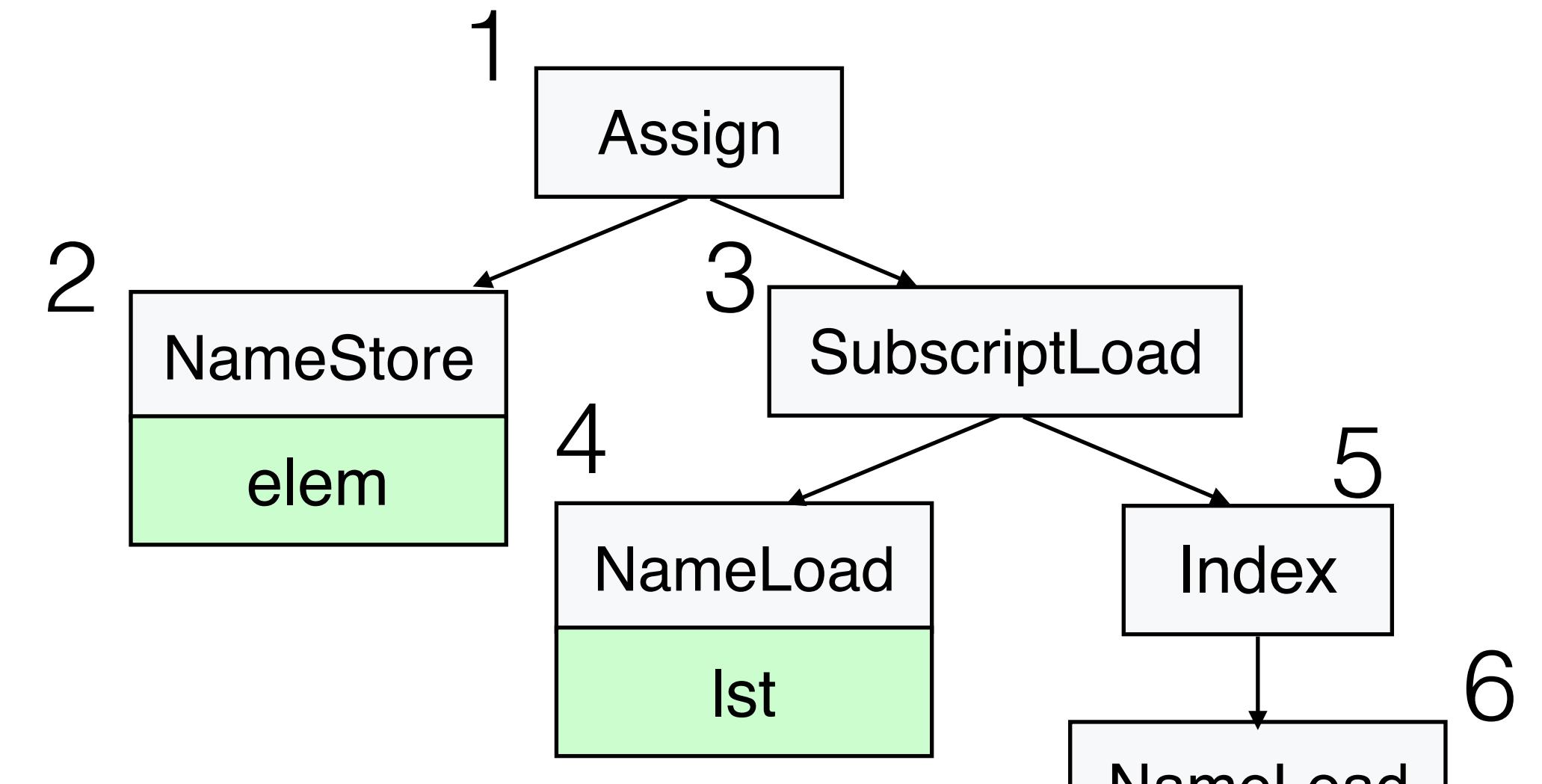
$$\tilde{\alpha}_{ij} = \frac{\exp(a_{ij} \cdot r_{j \rightarrow i})}{\sum_j \exp(a_{ij} \cdot r_{j \rightarrow i})}$$

$$a_{ij} = \frac{x_i^q x_j^k{}^T}{\sqrt{d_z}}$$

# Tree relative attention

Pairwise relations between AST nodes:

	1	2	3	4	5	6
1	I	D	D	DD	DD	DDD
2	U	I	UD	UDD	UDD	UDDD
3	U	UD	I	D	D	DD
4	UU	UUD	U	I	UD	UDD
5	UU	UUD	U	UD	I	D
6	UUU	UUUD	UU	UUD	U	I

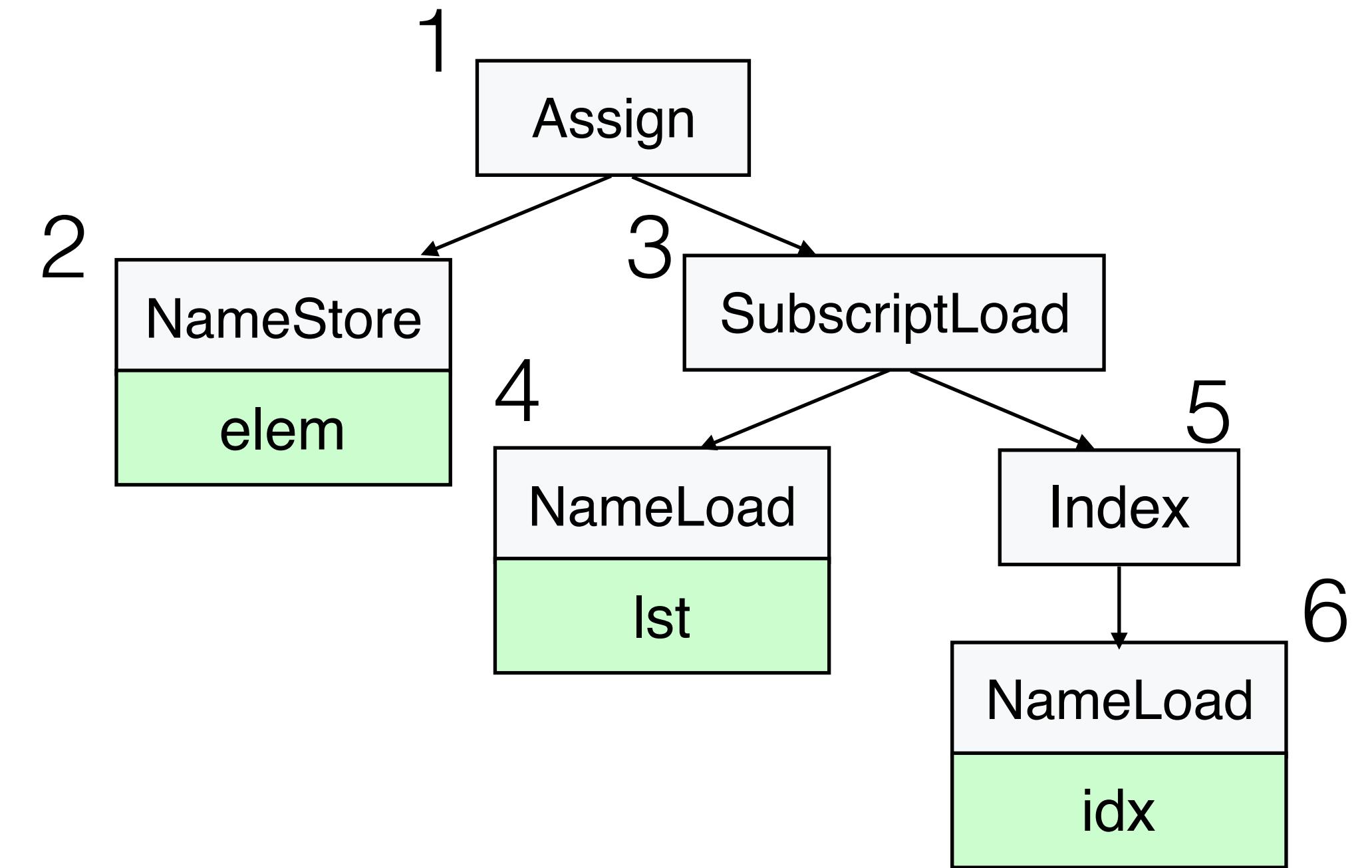


Learn 1-dim  
embedding  
for each relation

# Tree relative attention

Pairwise relations between AST nodes:

	1	2	3	4	5	6
1	I	D	D	DD	DD	DDD
2	U	I	UD	UDD	UDD	UDDD
3	U	UD	I	D	D	DD
4	UU	UUD	U	I	UD	UDD
5	UU	UUD	U	UD	I	D
6	UUU	UUUD	UU	UUD	U	I



Parameters:

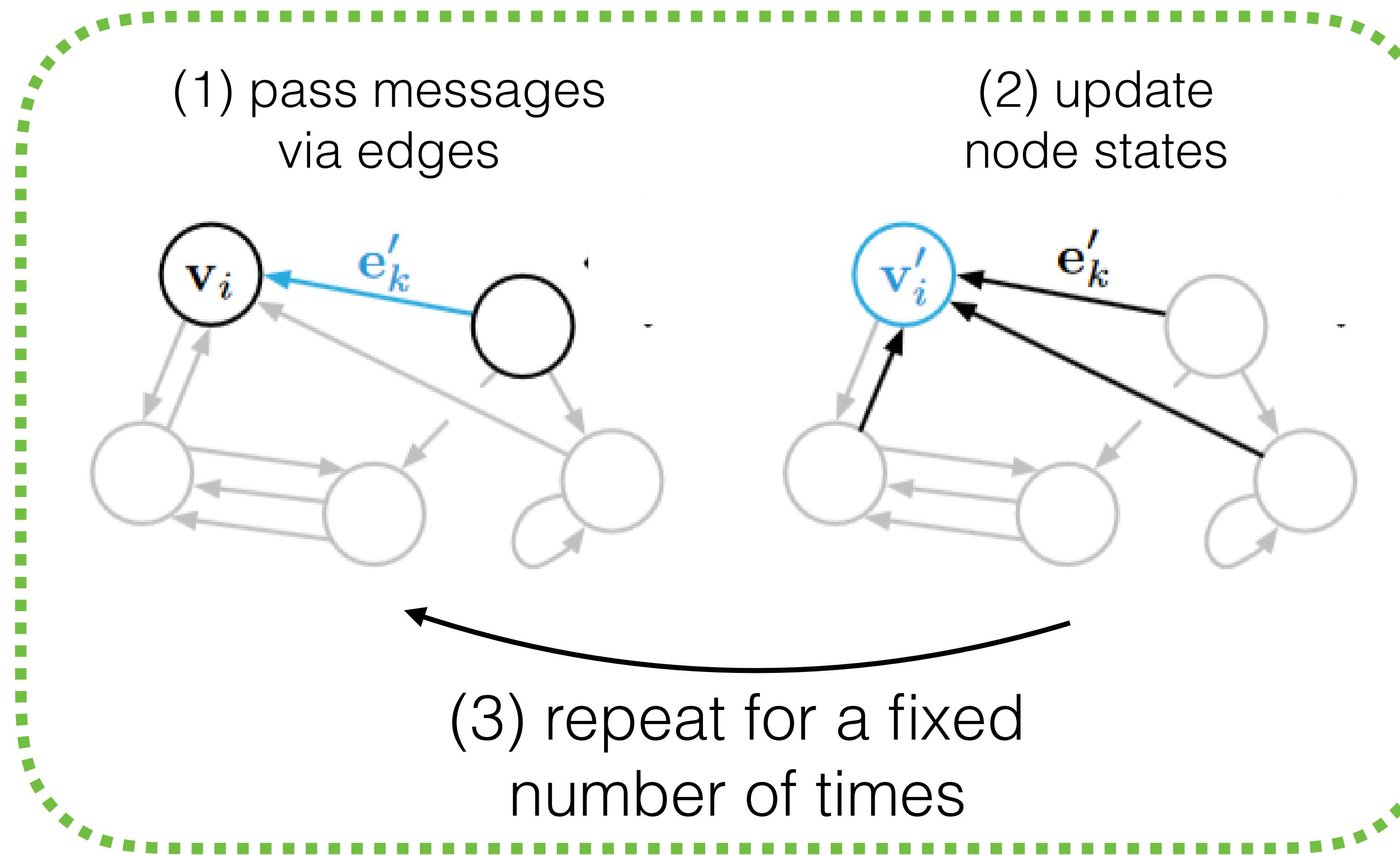
1-dim embeddings  
of “relations”

Hyperparameters: none

(or the maximum number of relations)

# GGNN-Sandwich

Graph Gated Neural Network (GGNN):

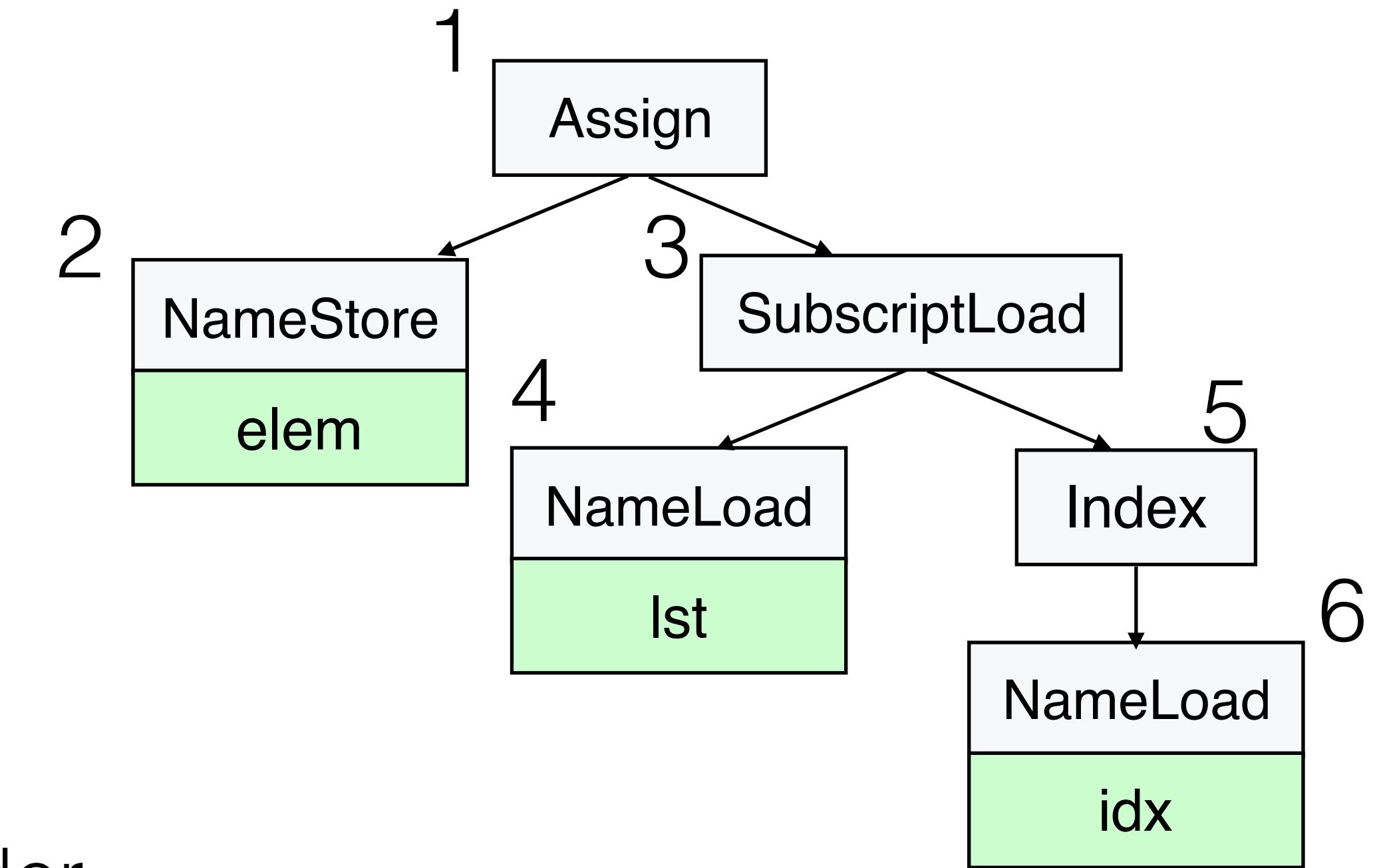
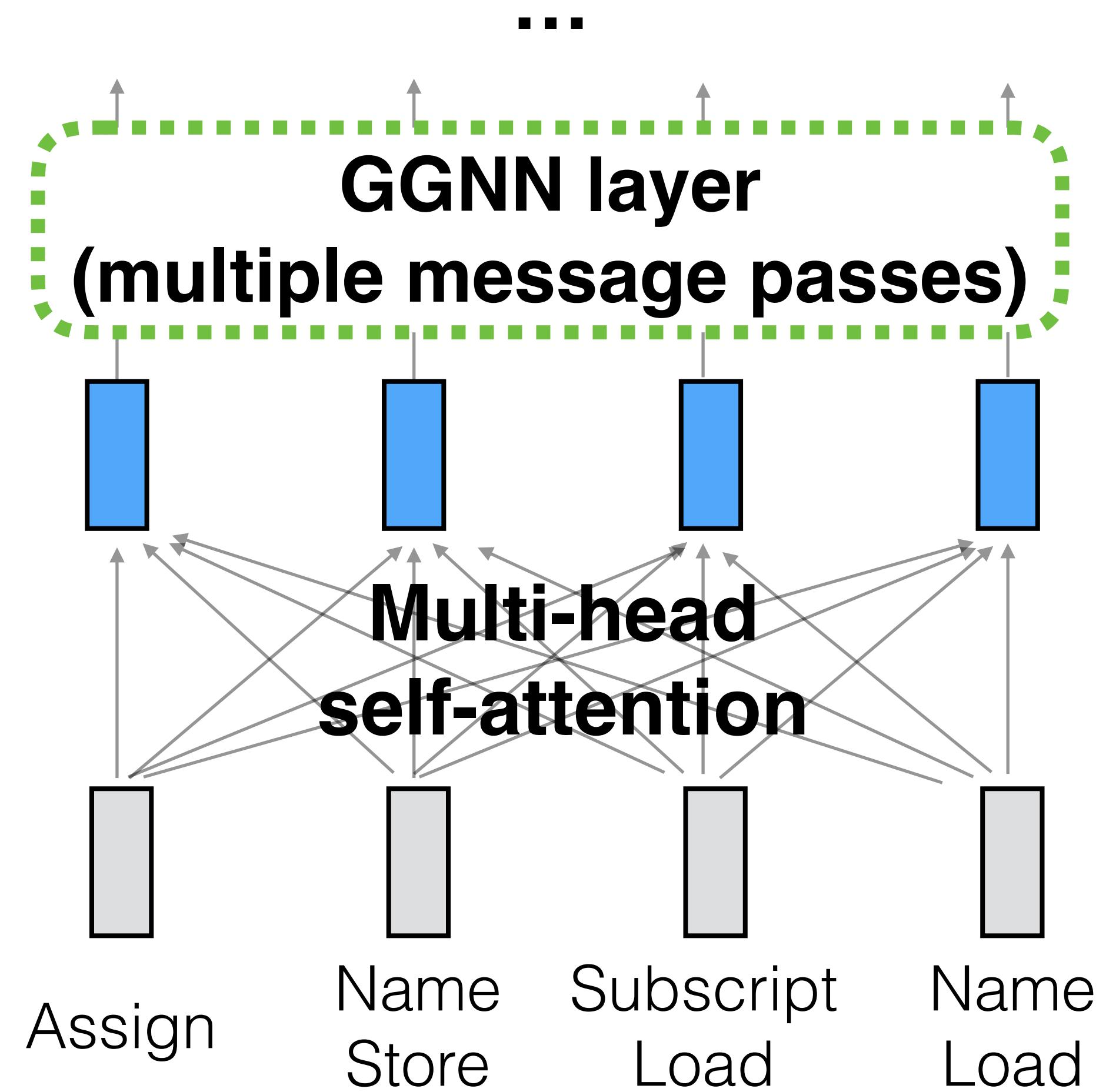


Input:  
states of all nodes

Output:  
updated states  
of all nodes

Arbitrary graph  
structure supported,  
e.g. different edge types

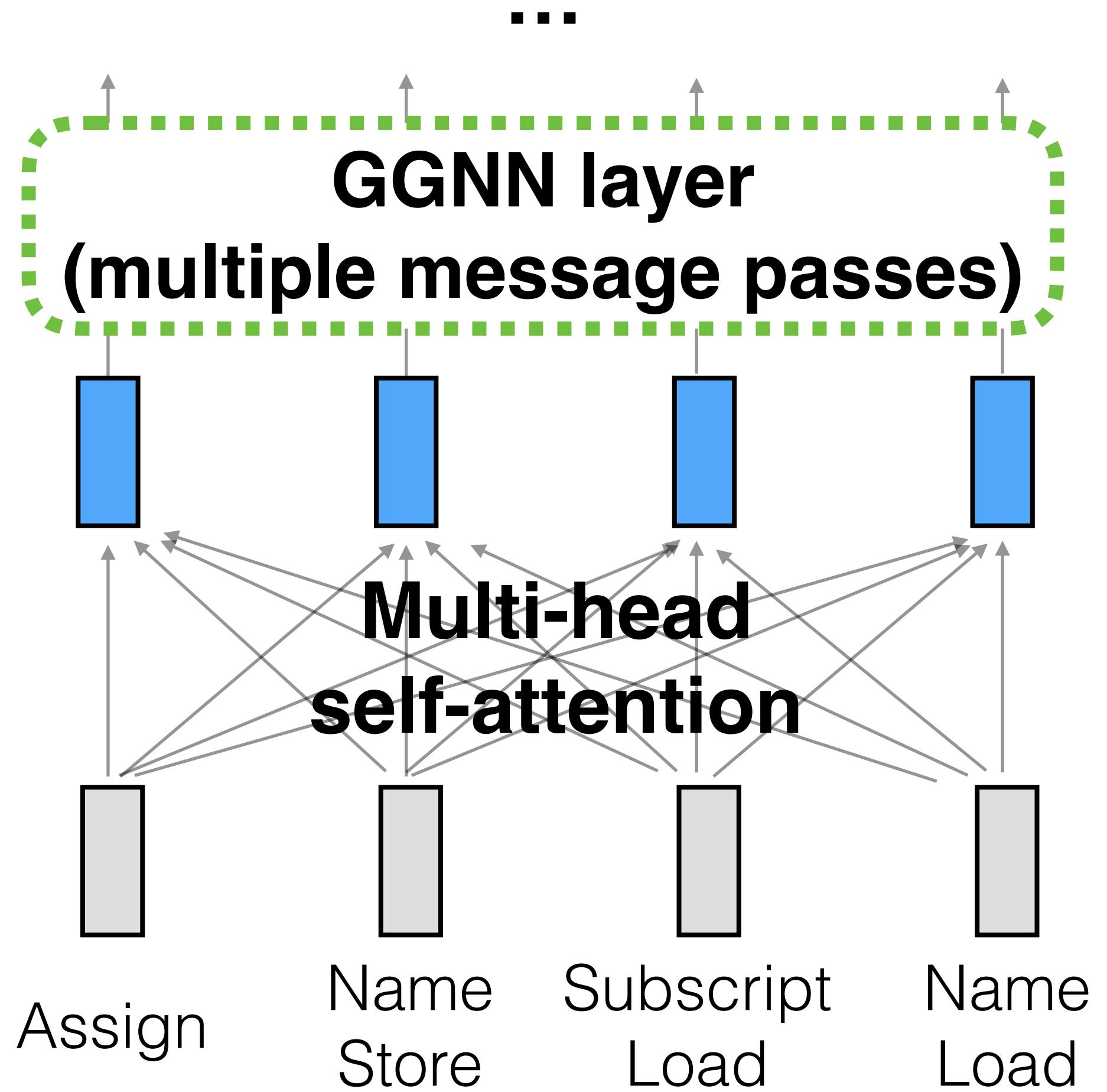
# GGNN-Sandwich



We consider  
4 types  
of edges:  
• Parent (P)  
• Child (C)  
• Left (L)  
• Right (R)

	1	2	3	4	5	6
1		P, L	P			
2	C, R		L			
3	C	R		P, L	P	
4			C, R		L	
5			C	R		P, L
6					C, R	

# GGNN-Sandwich



Parameters:  
same as n GGNN

Hyperparameters:  
same as in GGNN  
(e. g. the number  
of message passes)

Hard to use  
in Transformer decoder

# Techniques for capturing AST structure in Transformer

- Sequential positional embeddings (+AST depth-first traversal)
- Sequential relative attention (+AST depth-first traversal)
- Tree positional encoding
- Tree relative attention
- GGNN-Sandwich

# Techniques out-of-scope

- Recursive Transformer [1] — not a “global” model, slow
- Transformer with hierarchical accumulation [2] — applied only to texts, future work
- Inferring a tree for the input sequence [3] — other problem statement

[1] Mahtab Ahmed, Muhammad Rifayat Samee, Robert E. Mercer. You Only Need Attention to Traverse Trees, ACL19

[2] Xuan-Phi Nguyen, Shafiq Joty, Steven Hoi, Richard Socher. Tree-Structured Attention with Hierarchical Accumulation. ICLR20

[3] Yau-Shian Wang, Hung-Yi Lee, Yun-Nung Chen. Tree Transformer: Integrating Tree Structures into Self-Attention. EMNLP19

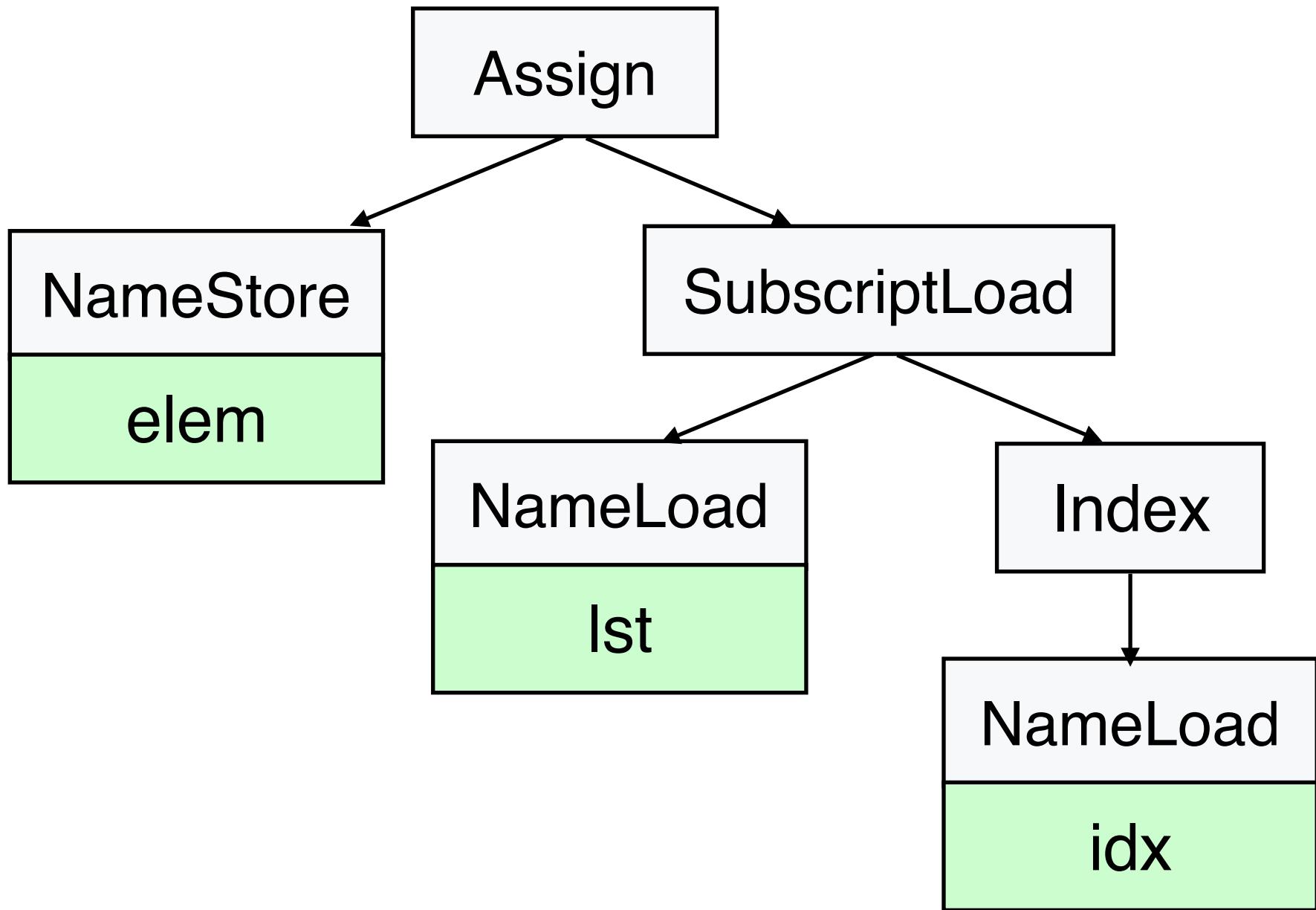
# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
- 5 syntax-capturing Transformer modifications
  - 3 tasks: code completion, bug fixing, function naming
  - 2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework  
(e. g. same data preprocessing for all techniques)

# Some technical details

Each AST node contains type,  
but not all nodes contain values

How to define type and value vocabularies?

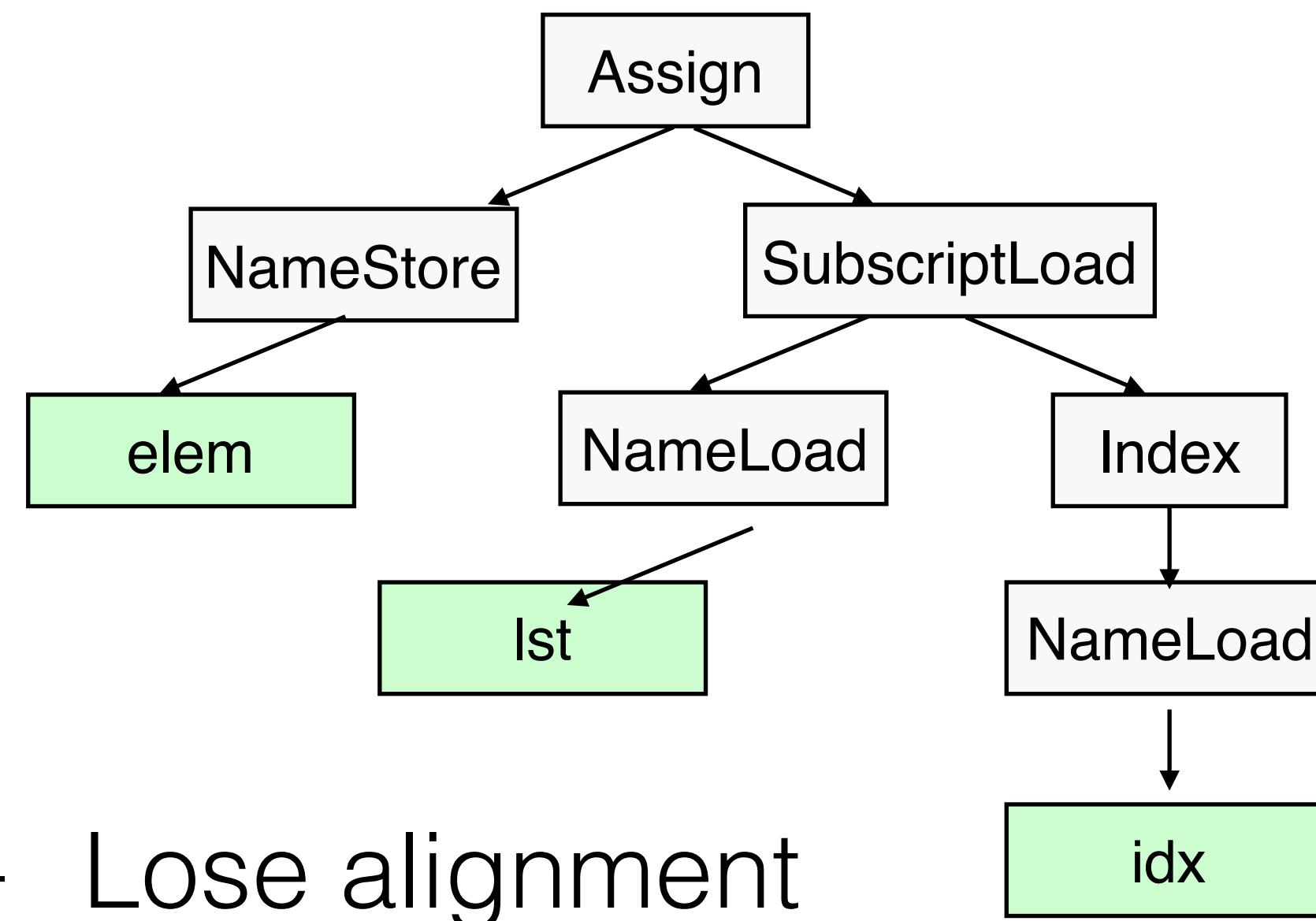


# Some technical details

Each AST node contains type,  
but not all nodes contain values

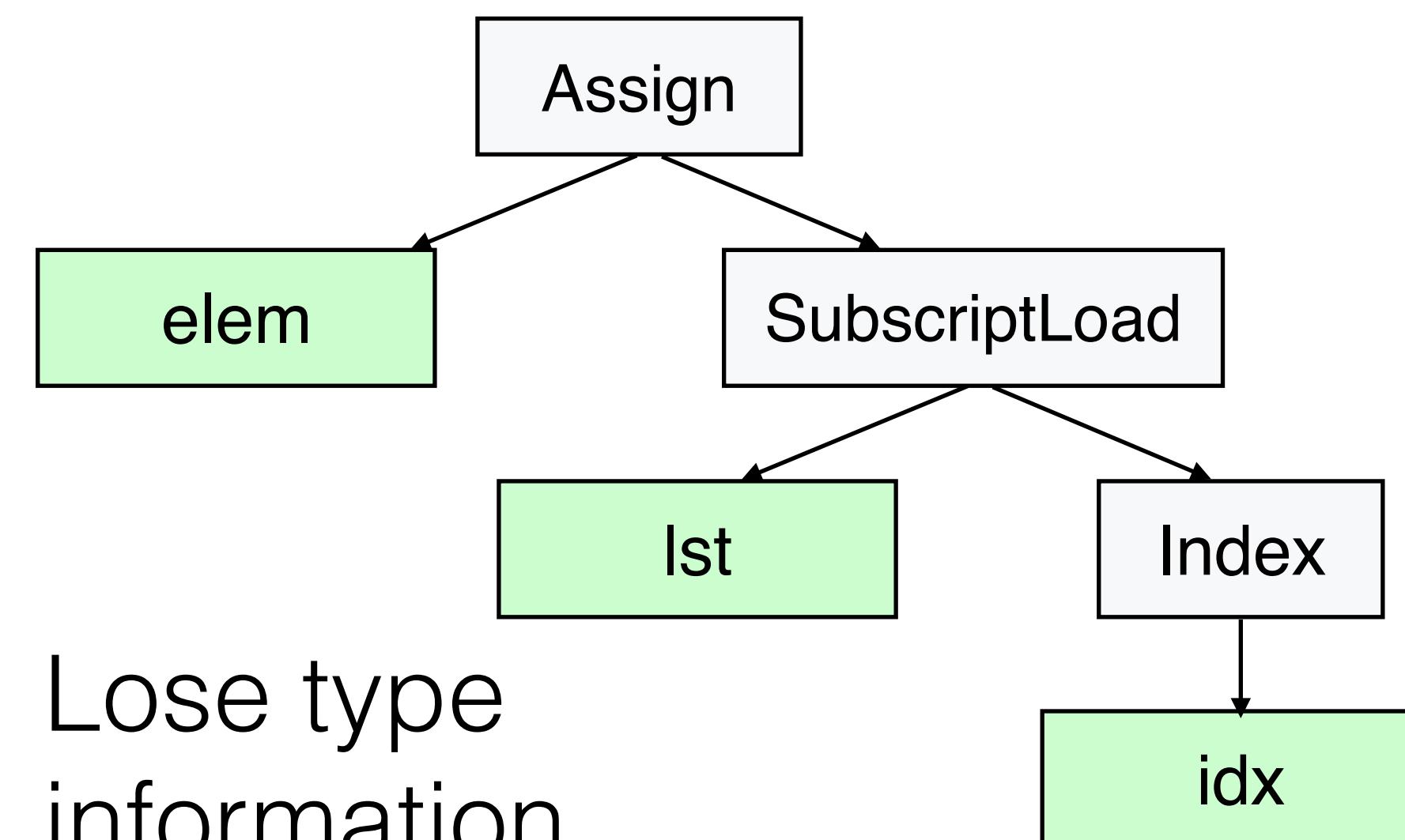
1 node — 1 token strategies:

(Kim20) Re-assign values:



- Lose alignment
- Longer sequences

(Hellendoorn20)  
Omit types:



- Lose type information

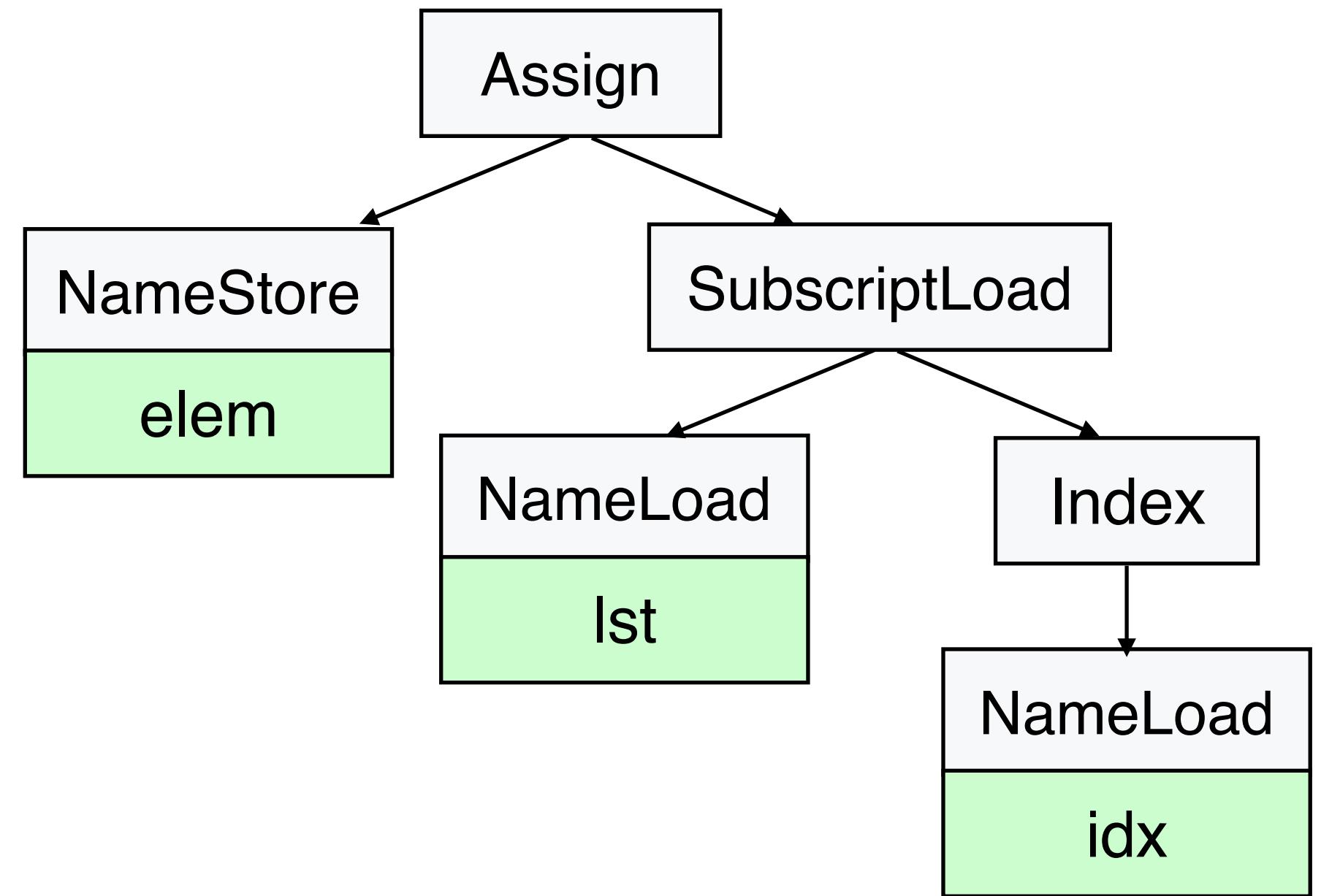
# Some technical details

Each AST node contains type,  
but not all nodes contain values

Our strategy (borrowed from (Li18)):

Use <empty> values, sum type embedding and value embedding:

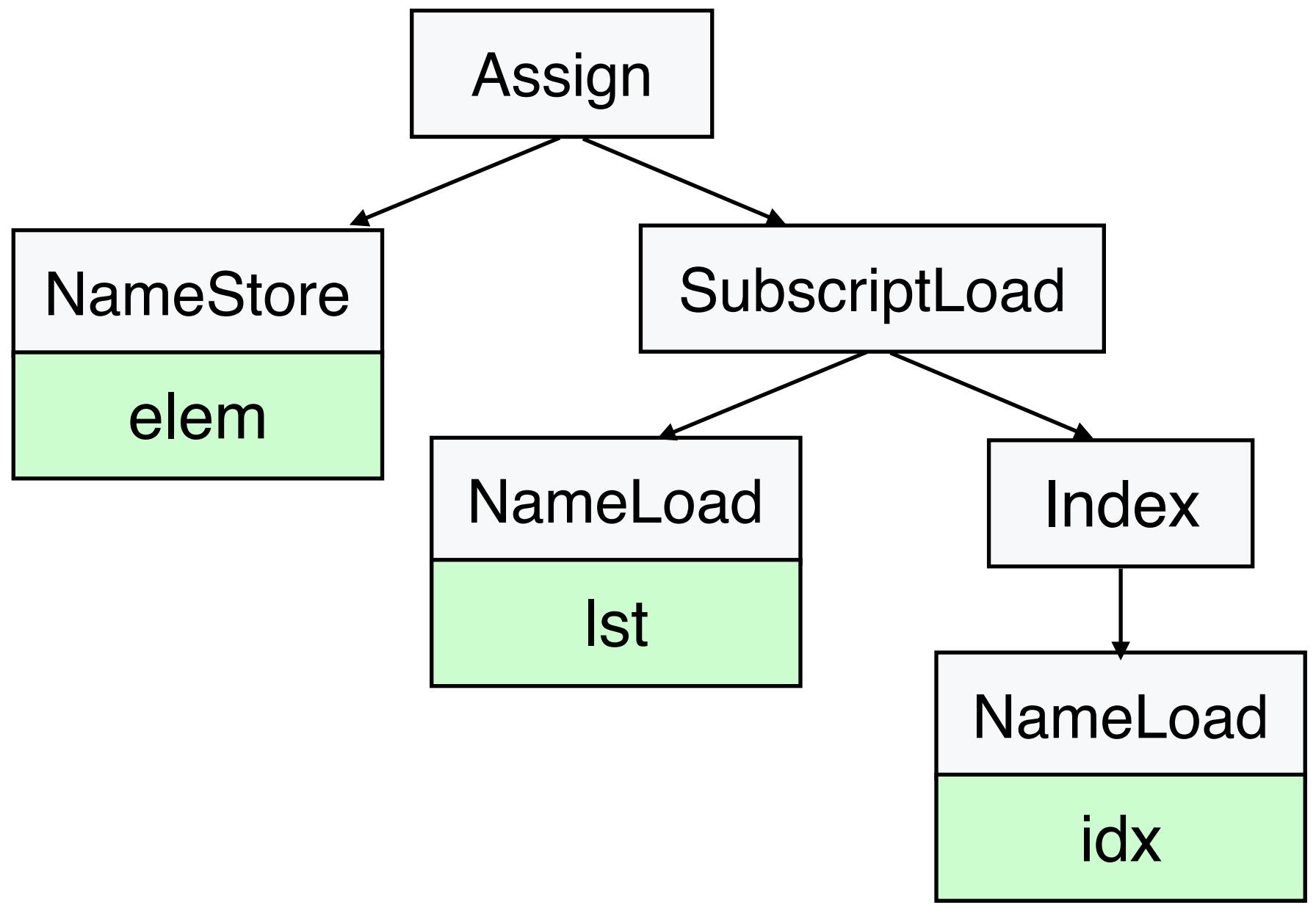
Assign	NameStore	SubscriptLoad	...
<empty>	elem	<empty>	...
...	NameLoad	Index	NameLoad
...	lst	<empty>	idx



- + No information lost
- + Preserve input length
- Small redundancy

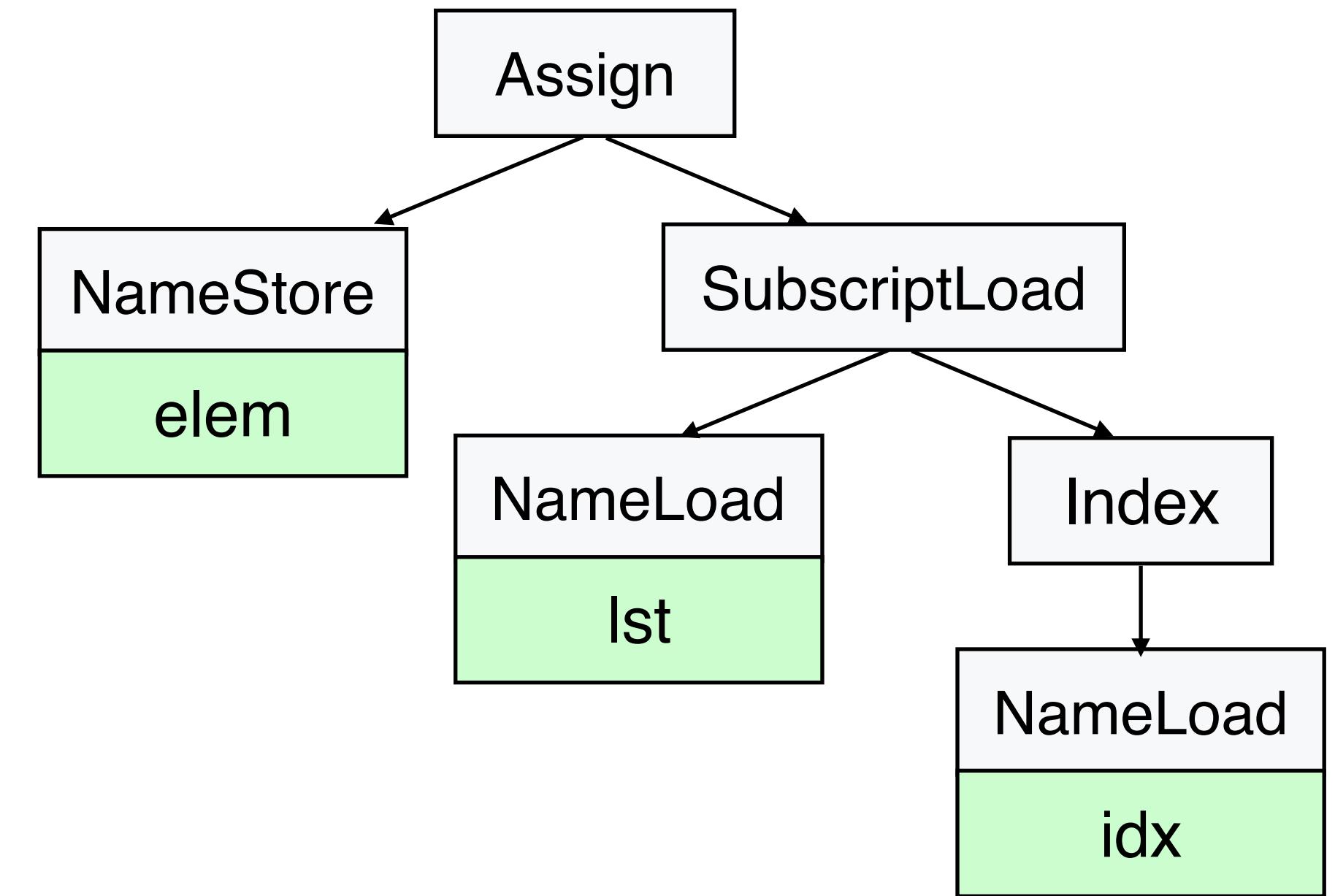
# Some technical details

- Use <empty> values, sum type embedding and value embedding
- Vocabulary preprocessing:
  - for code: no token splitting  
**very\_long\_identifier** used “as is”
  - for function names: split by tokens  
**very long function name**



# Some technical details

- Use <empty> values, sum type embedding and value embedding
- Vocabulary preprocessing:
  - for code: no token splitting  
**very\_long\_identifier** used “as is”
  - for function names: split by tokens  
**very long function name**
- Hyperparameters for tree pos. encoding, seq. rel. attention and GGNN-Sandwich were chosen using grid search



# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
- 5 syntax-capturing Transformer modifications
  - 3 tasks: code completion, bug fixing, function naming
  - 2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework (e. g. same data preprocessing for all techniques)

# Tasks

Code completion:

```
import sys
from utils.parsing import parse
if __name__ == '__main__':
    arguments = parse(
        sys.argv[1:])
)
print( next token?
```

Function naming:

```
def <fun_name>(filename):
    with open(filename) as fin:
        return len(
            fin.read()
            .split("\n"))
)
```

**Function name?**

Variable misuse

```
def count_lines(filename):
    with open(filename) as fin:
        return len(
            filename.read()
            .split("\n"))
)
```

**A bug? Fix?**

# Tasks

Code completion:

```
import sys
from utils.parsing import parse
if __name__ == '__main__':
    arguments = parse(
        sys.argv[1:])
)
print( argmnents
```

Kim, S., Zhao, J., Tian, Y., and  
Chandra, S. Code Prediction  
by Feeding Trees to  
Transformers.  
2020

Function naming:

```
def <fun_name>(filename):
    with open(filename) as fin:
        return len(
            fin.read()
            .split("\n"))
)
```

**count\_lines**

Variable misuse

```
def count_lines(filename):
    with open(filename) as fin:
        return len(
bug → filename.read()
            .split("\n"))
)
```

**fix**

# Tasks

Code completion:

```
import sys
from utils.parsing import parse
if __name__ == '__main__':
    arguments = parse(
        sys.argv[1:])
)
print( argmnents
```

**Decoder only**

**Encoder only**

Function naming:

```
def <fun_name>(filename):
    with open(filename) as fin:
        return len(
            fin.read()
            .split("\n"))
)
```

**count\_lines**

Variable misuse

```
def count_lines(filename):
    with open(filename) as fin:
        return len(
bug → filename.read()
            .split("\n"))
)
```

**fix**

**Code encoder +  
Text decoder**

# Tasks: metrics

Code completion:

```
import sys
from utils.parsing import parse
if __name__ == '__main__':
    arguments = parse(
        sys.argv[1:])
)
print( argmnents
```

MRR

Joint loc. and  
repair accuracy

Function naming:

```
def <fun_name>(filename):
    with open(filename) as fin:
        return len(
            fin.read()
            .split("\n"))
)
```

**count\_lines**

F-measure

Variable misuse

```
def count_lines(filename):
    with open(filename) as fin:
        return len(
bug → filename.read()
            .split("\n"))
)
```

**fix**

# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
  - 5 syntax-capturing Transformer modifications  
3 tasks: code completion, bug fixing, function naming  
2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework  
(e. g. same data preprocessing for all techniques)

# Datasets

- Python150k / JavaScript150k (100K train files + 50K test files)  
( $\sim 10^8$  train chars,  $\sim 5 \cdot 10^7$  test chars)
- Custom train / test split:
  - repository-based
  - deduplication

# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
- 5 syntax-capturing Transformer modifications
  - 3 tasks: code completion, bug fixing, function naming
  - 2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework (e. g. same data preprocessing for all techniques)

# What does this code snippet do?

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

# What does this code snippet do?

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

```
filtered_seqs = []
for seq in seqs:
    new_seq = []
    for word in seq:
        if word in vocab:
            new_seq += [word]
    filtered_seqs += [new_seq]
```

# Anonymized setting vs Full data setting

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

Only one source of information  
— syntactic structure

var78, var5, var12 etc  
— anonymized identifiers

```
filtered_seqs = []
for seq in seqs:
    new_seq = []
    for word in seq:
        if word in vocab:
            new_seq += [word]
    filtered_seqs.append(new_seq)
```

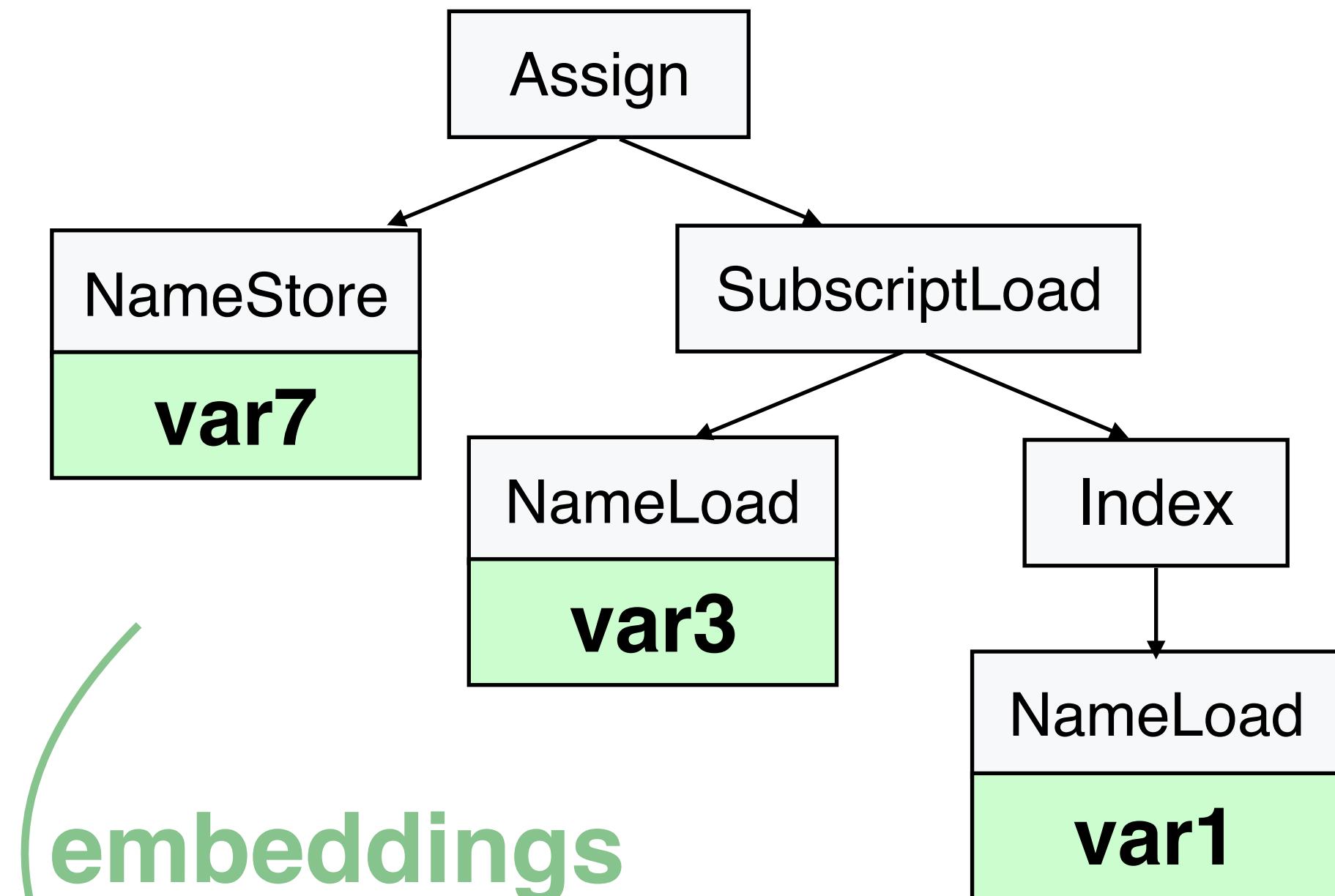
Two sources of information:

- Syntactic structure
- Identifier names (seq, word etc)

# Anonymization procedure

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

parse AST



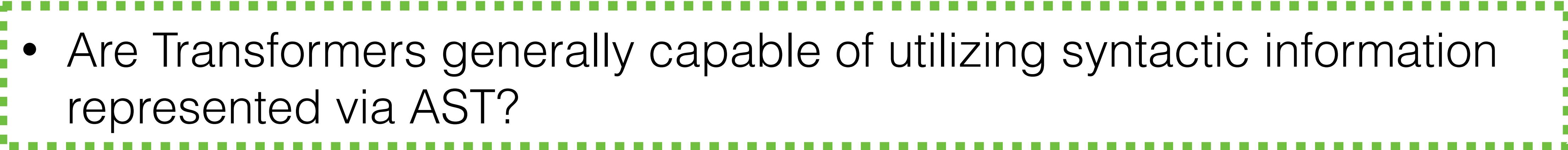
- Random subset of **var1** ... **var1000**
- Save equality/inequality information
- An identifier may have different anonymizations in different input snippets

Assign	<empty>
NameStore	var7
SubscriptLoad	<empty>
NameLoad	var3
Index	<empty>
NameLoad	var1

# Overview

- Goal: understanding whether Transformers are able to utilize syntactic structure and what is the best approach to do it
- 5 syntax-capturing Transformer modifications
  - 3 tasks: code completion, bug fixing, function naming
  - 2 datasets: Python150k and Javascript150k
- Thoughtful experimental setup:
  - repository-based train / test split
  - 2 settings: full data and anonymized data
  - re-implementing all techniques in a unified framework (e. g. same data preprocessing for all techniques)

# Research questions

- 
- Are Transformers generally capable of utilizing syntactic information represented via AST?
  - What components of AST (structure, node types and values) do Transformers use in different tasks?
  - What is the most effective approach for utilizing AST *structure* in Transformer?

# Capability of Transformers to utilize syntactic information

*Syntax<sup>\*</sup>+Text*

1	Assign	<empty>
2	NameStore	elem
3	SubscriptLoad	<empty>
4	NameLoad	lst
5	Index	<empty>
6	NameLoad	idx

vs *Text only*

1	elem
2	lst
3	idx

*Syntax<sup>\*</sup> only*

1	Assign	<empty>
2	NameStore	var7
3	SubscriptLoad	<empty>
4	NameLoad	var3
5	Index	<empty>
6	NameLoad	var1

vs *Dummy*

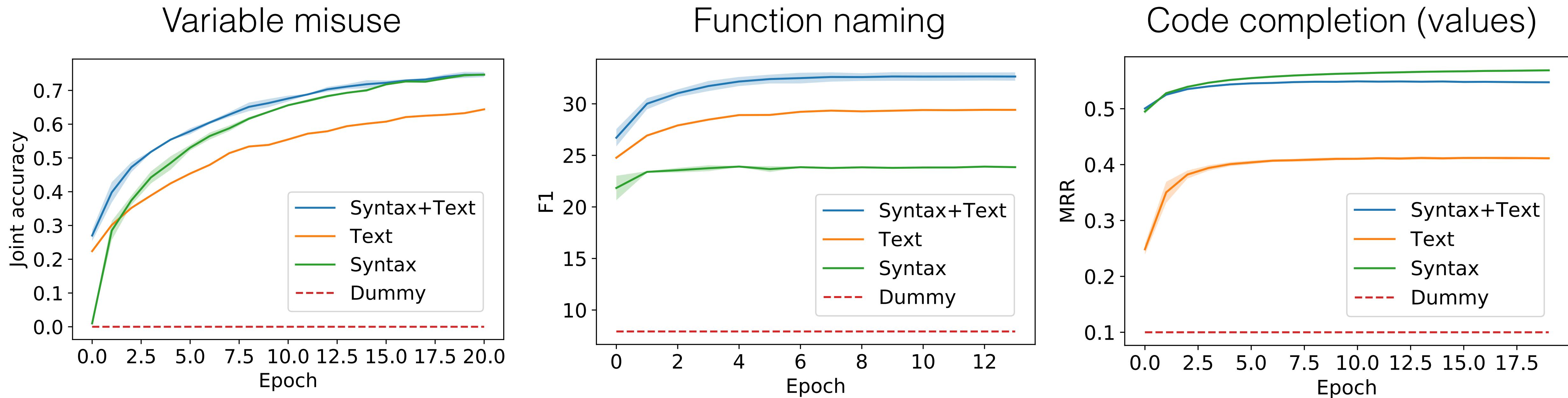
Most frequent constant prediction

Does using AST  
in addition to identifier names  
improve the performance?

Does using only AST  
perform any better  
than a Dummy baseline?

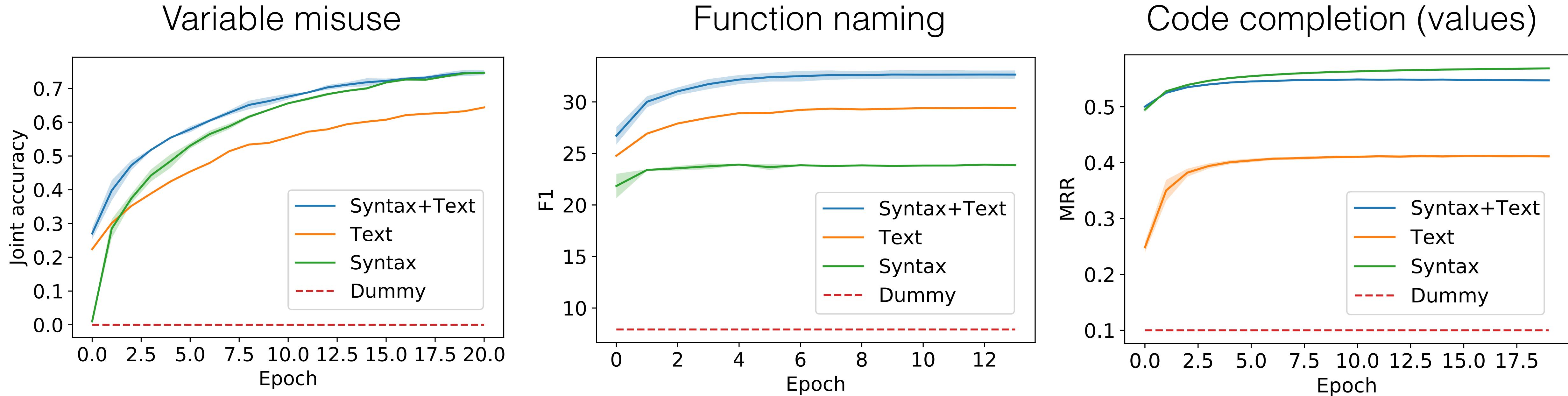
\* AST depth-first traversal + standard sequential positional embeddings

# Capability of Transformers to utilize syntactic information: results



- *Syntax + Text > Text only*
- *Syntax only > Dummy*

# Capability of Transformers to utilize syntactic information: results



- *Syntax + Text > Text only*
- *Syntax only > Dummy*
- Interestingly, sometimes *Syntax only ≥ Syntax + Text*

# Code completion: demo

```
def get_or_create_bucket(s3_connection):
    bucket = s3_connection.get_bucket(
        settings.S3_BUCKET_NAME
    )
    if bucket is None:
        bucket = s3_connection.
```

```
import sys
from utils.parsing import parse
if __name__ == '__main__':
    arguments = parse(
        sys.argv[1:]
    )
    print(
```

(a) **S+T:** [create\_bucket get\_bucket get\_key]  
**S:** [get\_bucket S3\_BUCKET\_NAME var107]  
**T:** [bucket get\_bucket settings]  
**Dummy:** [self 0 None]  
**Gold:** create\_bucket

(b) **S+T:** [Usage: done. \n]  
**S:** [argmnents var440 var289]  
**T:** [sys outfn len]  
**Dummy:** [self 0 None]  
**Gold:** argmnents

# Function naming: demo

```
def <fun_name>(seqs):
    dt = {}
    for seq in seqs:
        for word in seq.split():
            if not word in dt:
                dt[word] = 1
            else:
                dt[word] += 1
    return dt
```

- (c) **S+T**: get dict
  - S**: get all
  - T**: get split
  - Dummy**: init
  - Gold**: get dictionary

```
def <fun_name>(filename):
    with open(filename) as fin:
        return len(
            fin.read()
            .split("\n")
        )
```

- (d) **S+T**: read file
  - S**: read
  - T**: read file
  - Dummy**: init
  - Gold**: count lines

# Research questions

- Are Transformers generally capable of utilizing syntactic information represented via AST?
- What components of AST (structure, node types and values) do Transformers use in different tasks?
- What is the most effective approach for utilizing AST *structure* in Transformer?

# Ablation study of mechanisms for utilizing syntactic information in Transformer

3 syntax-capturing mechanisms:

- types
- syntactic structure  
(AST edges)
- information about value equality/inequality

**What mechanisms are essential for achieving the benefit in quality?**

# Ablation study of mechanisms for utilizing syntactic information in Transformer

3 syntax-capturing mechanisms:

- types

(already  
ablated  
in the previous  
experiment)

- syntactic structure  
(AST edges)

a bag of (type, (an) value) pairs  
— turn off  
seq. pos. embeddings

Assign	<empty>
NameStore	var7 / elem
SubscriptLoad	<empty>
NameLoad	var3 / 1st
Index	<empty>
NameLoad	var1 / idx

- information about value equality/inequality

a sequence of types

1	Assign
2	NameStore
3	SubscriptLoad
4	NameLoad
5	Index
6	NameLoad

# Ablation study of mechanisms for utilizing syntactic information in Transformer: results

	Full data			Anonymized data		
	Var. misuse	Fun. naming	Comp. (val.)	Var. misuse	Fun. naming	Comp. (val.)
Full AST	<b>74.59 %</b>	<b>32.83±0.35%</b>	<b>54.89 ±0.07%</b>	<b>74.71%</b>	<b>23.92 ±0.004%</b>	<b>56.87 ±0.09%</b>
AST w/o struct.	26.1%	<b>33.0 ±0.08%</b>	54.2 ±0.1%	12.12%	23.28 ±0.02%	55.07 ±0.05%
AST w/o an.val.	N/A	N/A	N/A	32.1%	<b>24.00 ±0.15%</b>	N/A

- Var. misuse, Code completion: all mechanisms important
- Function naming: mostly *types* contribute

Full AST

1	Assign	<empty>
2	NameStore	var7 / elem
3	SubscriptLoad	<empty>
4	NameLoad	var3 / 1st
5	. Index	<empty>
6	NameLoad	var1 / idx

AST w/o struct.

Assign	<empty>
NameStore	var7 / elem
SubscriptLoad	<empty>
NameLoad	var3 / 1st
Index	<empty>
NameLoad	var1 / idx

AST w/o an. val.

1	Assign
2	NameStore
3	SubscriptLoad
4	NameLoad
5	Index
6	NameLoad

# Research questions

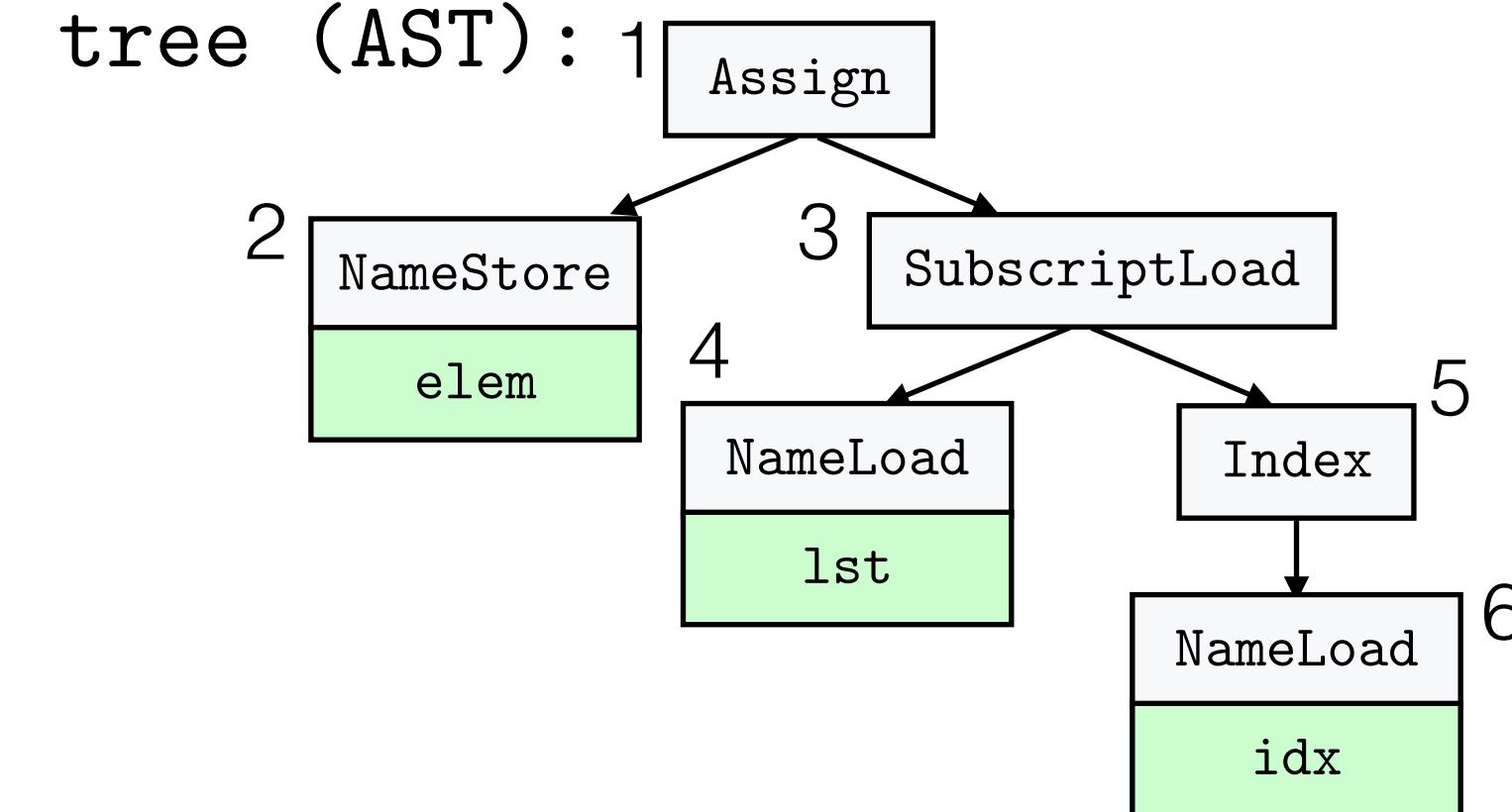
- Are Transformers generally capable of utilizing syntactic information represented via AST?
- What components of AST (structure, node types and values) do Transformers use in different tasks?

- 
- What is the most effective approach for utilizing AST *structure* in Transformer?

# Comparison of approaches for utilizing syntactic structure in Transformer

(a) Code: `elem = lst[idx]`

(b) Abstract syntax

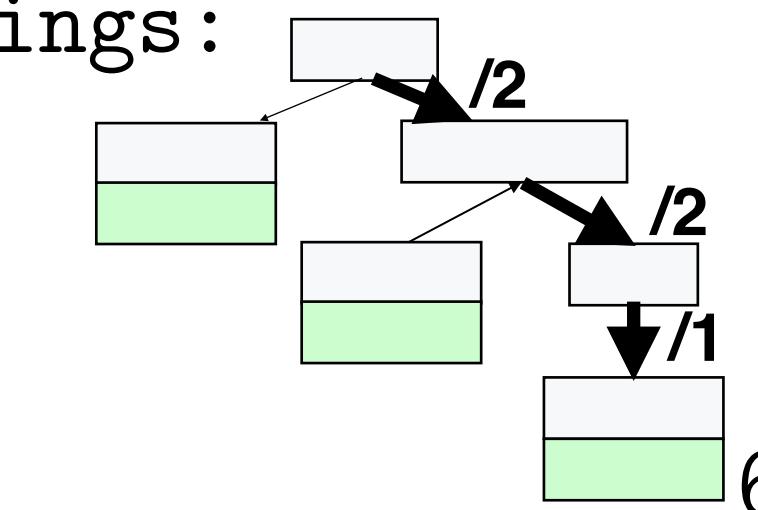


(c) AST depth-first traversal:

Assign	NameStore	SubscriptLoad	...
$\langle \text{empty} \rangle$	elem	$\langle \text{empty} \rangle$	...
...	NameLoad	Index	NameLoad
...	lst	$\langle \text{empty} \rangle$	idx

(d) Tree positional encodings:

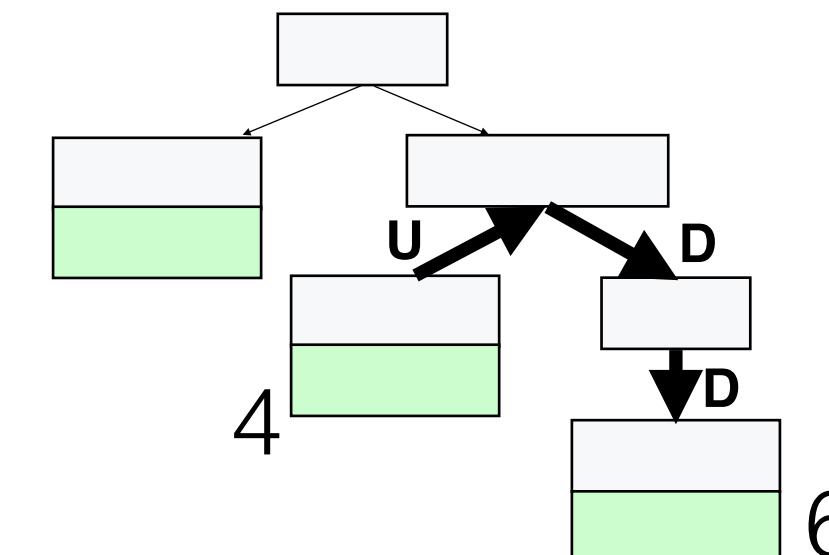
1	2	3	4	5	6
/	/1	/2	/2/1	/2/2	/2/2/1



stack-like enc.

000 000 000
100 000 000
010 000 000
100 010 000
010 010 000
100 010 010

(e) Tree relative attention:

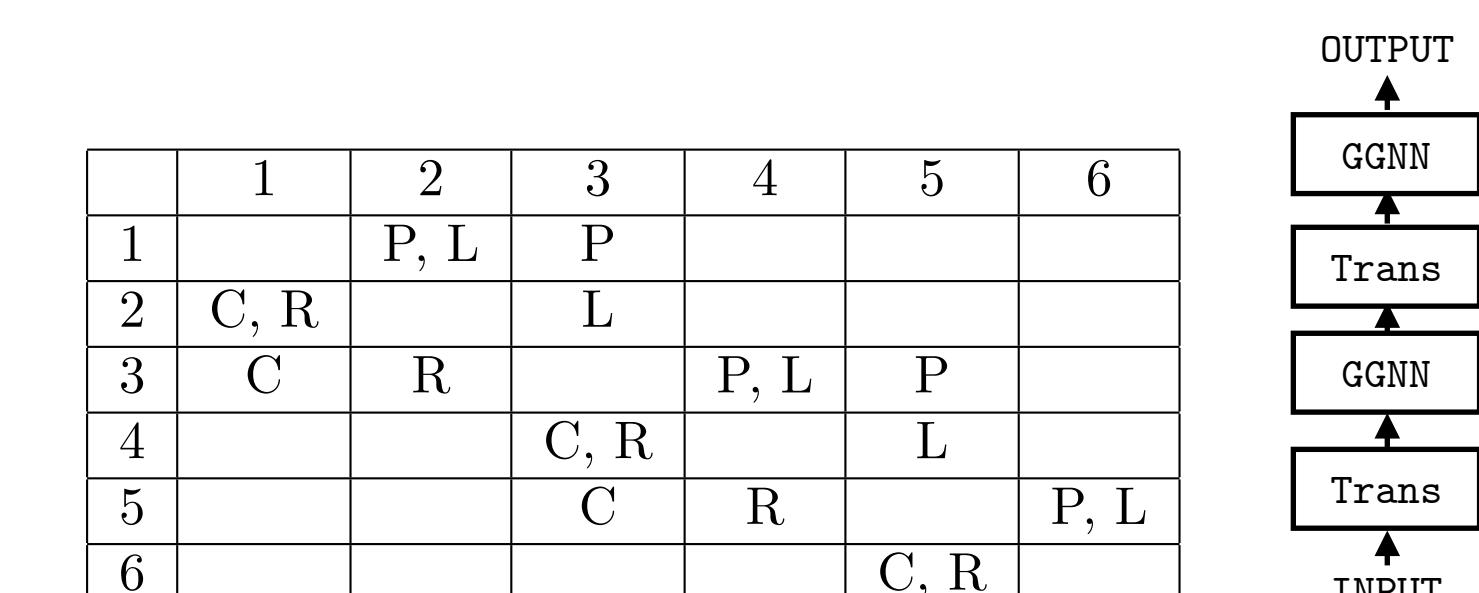


	1	2	3	4	5	6
1	I	D	D	DD	DD	DDD
2	U	I	UD	UDD	UDD	UDDD
3	U	UD	I	D	D	DD
4	UU	UUD	U	I	UD	UDD
5	UU	UUD	U	UD	I	D
6	UUU	UUUD	UU	UUD	U	I

(f) GGNN Sandwich:

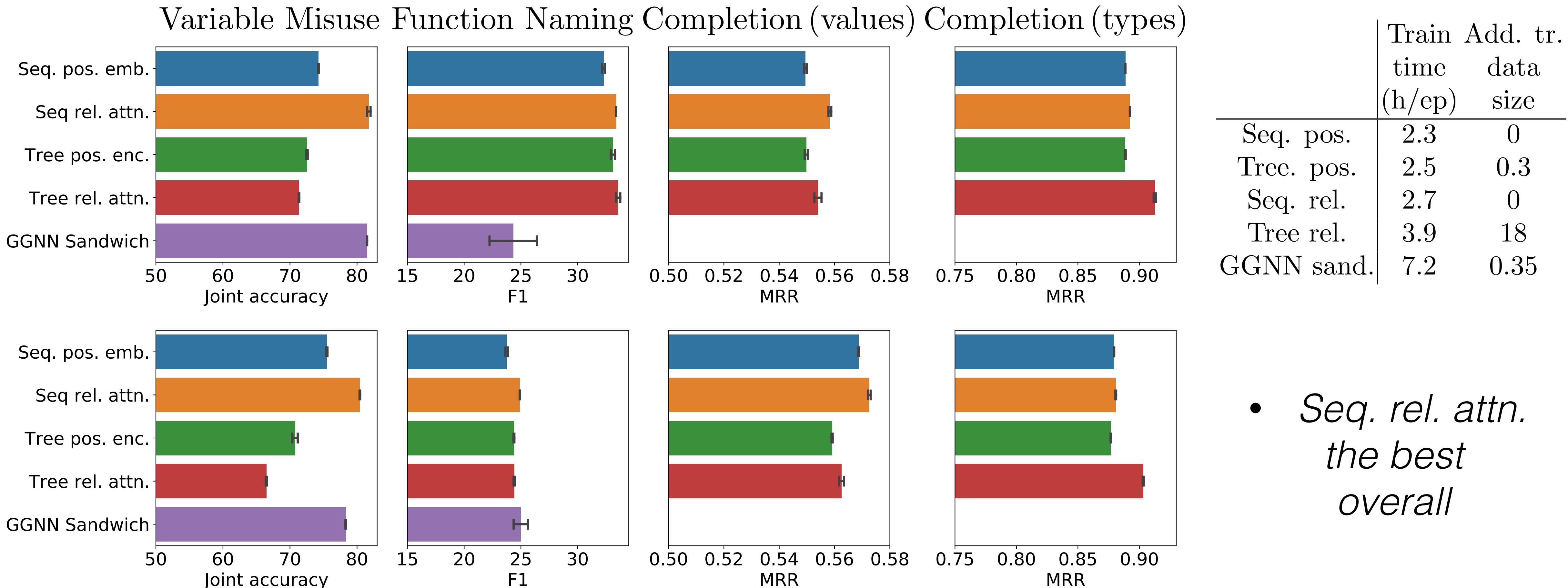
4 types of edges:

- Parent (P)
- Child (C)
- Left (L)
- Right (R)



- Sequential positional embeddings
- Sequential relative attention

# Comparison of approaches for utilizing syntactic structure in Transformer: results



- *GGNN-Sandwich fine for VM*

- *Tree.rel. attn. fine for CC-types*

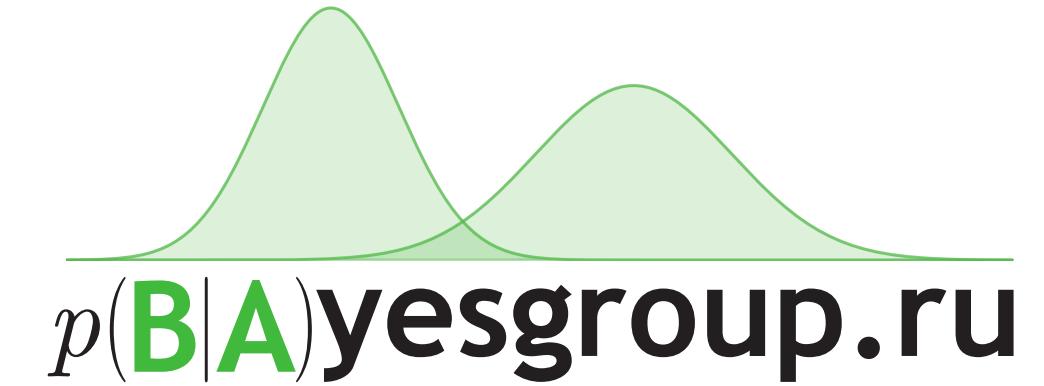
- *Seq. rel. attn. the best overall*

# Combining AST-capturing techniques

Model	VM	FN	CC (val.)
SRA	81.54 %	33.38 %	55.79
SRA + Seq. pos. emb.	80.31 %	31.78 %	55.73
SRA + Tree rel. attn.	80.59 %	33.01 %	<b>56.26</b>
SRA + Tree pos. enc.	81.74 %	32.52 %	55.98
SRA + GGNN sand.	<b>84.14</b> %	27.67 %	N/A

- Other techniques may be beneficial in combination with sequential relative attention

# Ensembling of syntax-based models



Models	VM	FN	CC (typ.)	CC (val.)
ST	81.54 %	<b>33.38%</b>	89.1 %	55.79 %
ST & ST	83.15 %	32.57 %	<b>89.4</b> %	56.86 %
ST & S	<b>85.73%</b>	29.42%	88.51 %	<b>59.87</b> %

- Ensembling *Syntax+Text* and *Syntax* outperforms ensembling two *Syntax+Text* (standard) models in several tasks

ST: *Syntax + Text*  
S : *Syntax* only  
& denotes ensembling

# Summary

- Using syntactic structure in Transformer increases the quality in all three tasks. However, in the function naming task, mostly type information contributes to the increase, while syntax itself does not.
- Sequential relative attention outperforms other techniques in almost all cases, with less computational overhead.  
This baseline is missing in a lot of works.
- Anonymization increases the quality of the ensemble and sometimes increases even the quality of a single model.

# A Simple Approach for Handling Out-of-Vocabulary Identifiers in Deep Learning for Source Code

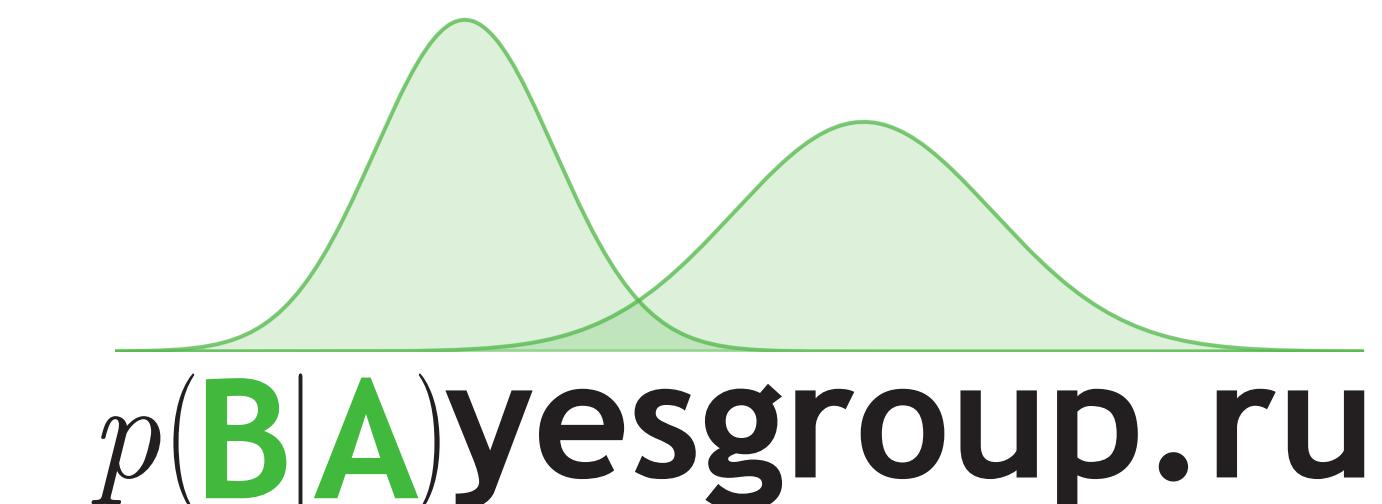
Nadezhda Chirkova\*, Sergey Troshin\*

Higher School of Economics, Samsung-HSE Laboratory  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



# Motivation

Advantages of identifier anonymization:

- Uniform token frequencies — dealing with rare or misspelled identifiers
- Other form of same data — beneficial for ensembling, (possibly) regularization
- Disadvantage: losing semantic information stored in variable names

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

# Motivation

Advantages of identifier anonymization:

- Uniform token frequencies — dealing with rare or misspelled identifiers
- Other form of same data — beneficial for ensembling, (possibly) regularization
- Disadvantage: losing semantic information stored in variable names

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

# Proposed method

Preprocessing steps:

1. Select the vocabulary of frequent identifiers, e. g. 50K
2. Anonymize out-of-vocabulary identifiers

Vocabulary: {**np, sin**}

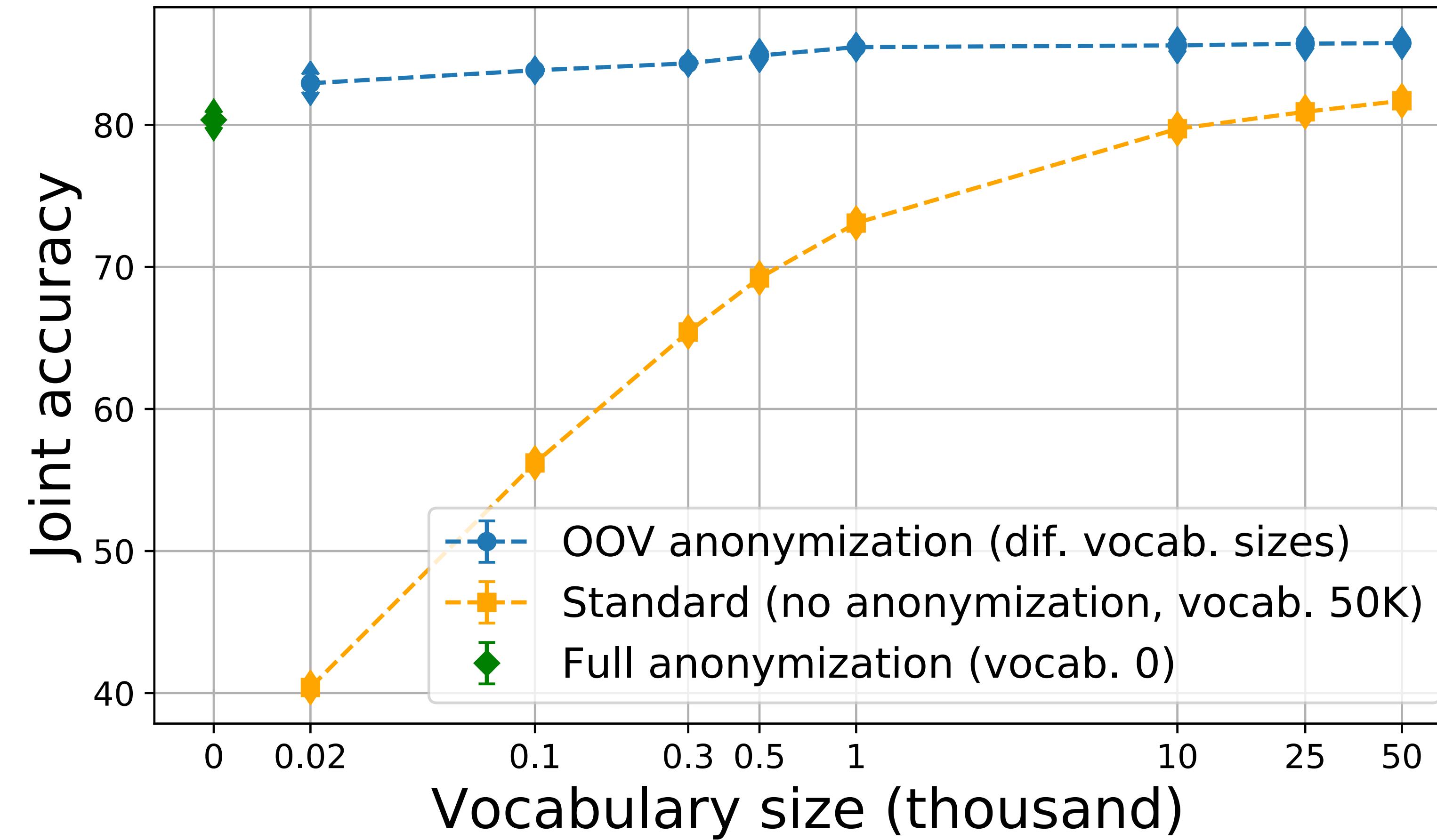
Input: **my\_y = np.sin(my\_x) + my\_x**

Standard OOV  
processing procedure: **UNK = np.sin(UNK) + UNK**

Proposed OOV  
anonymization procedure: **VAR1 = np.sin(VAR2) + VAR2**

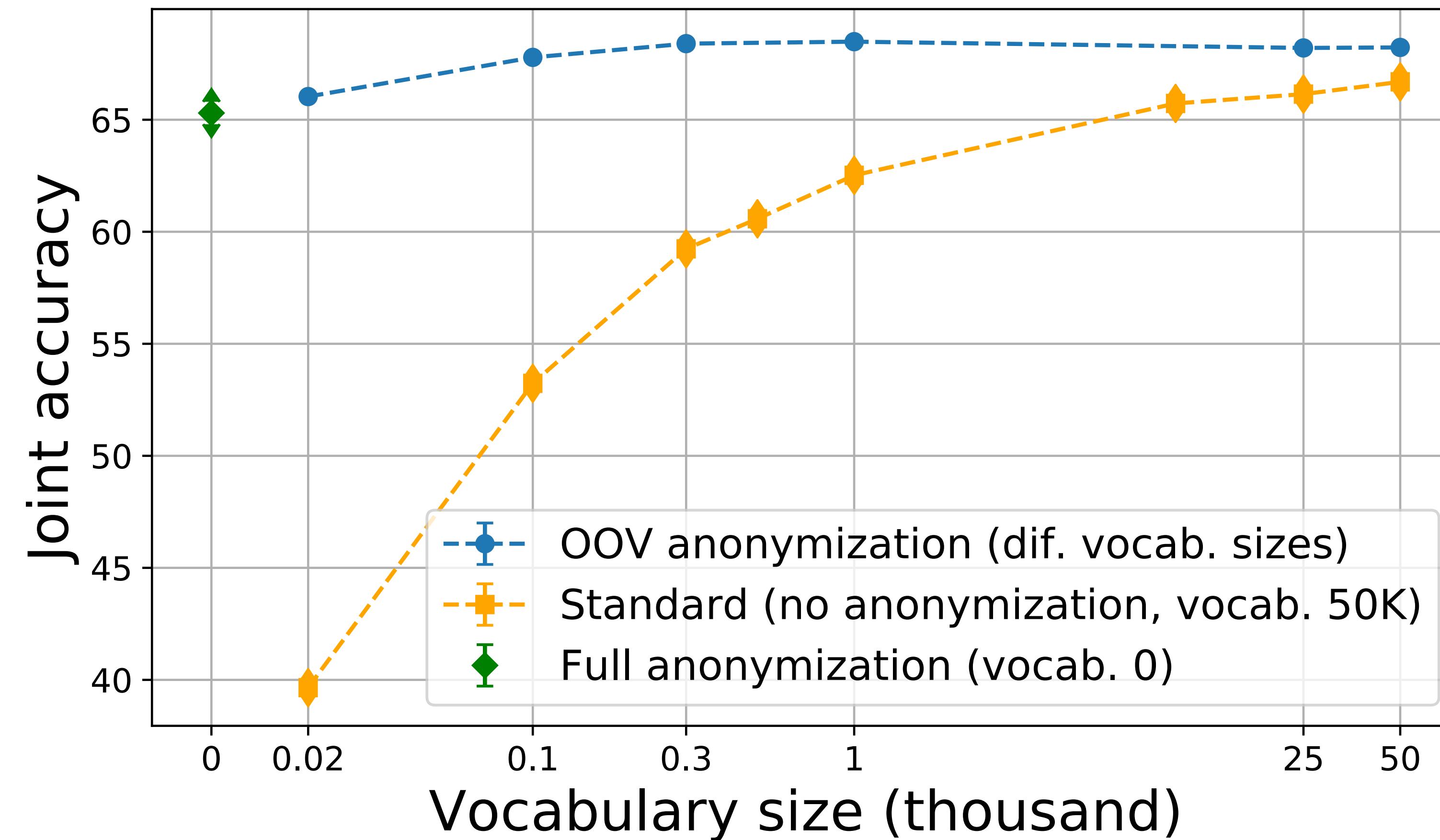
# Results: variable misuse task, Python150k

Same setting as below, sequential relative attention

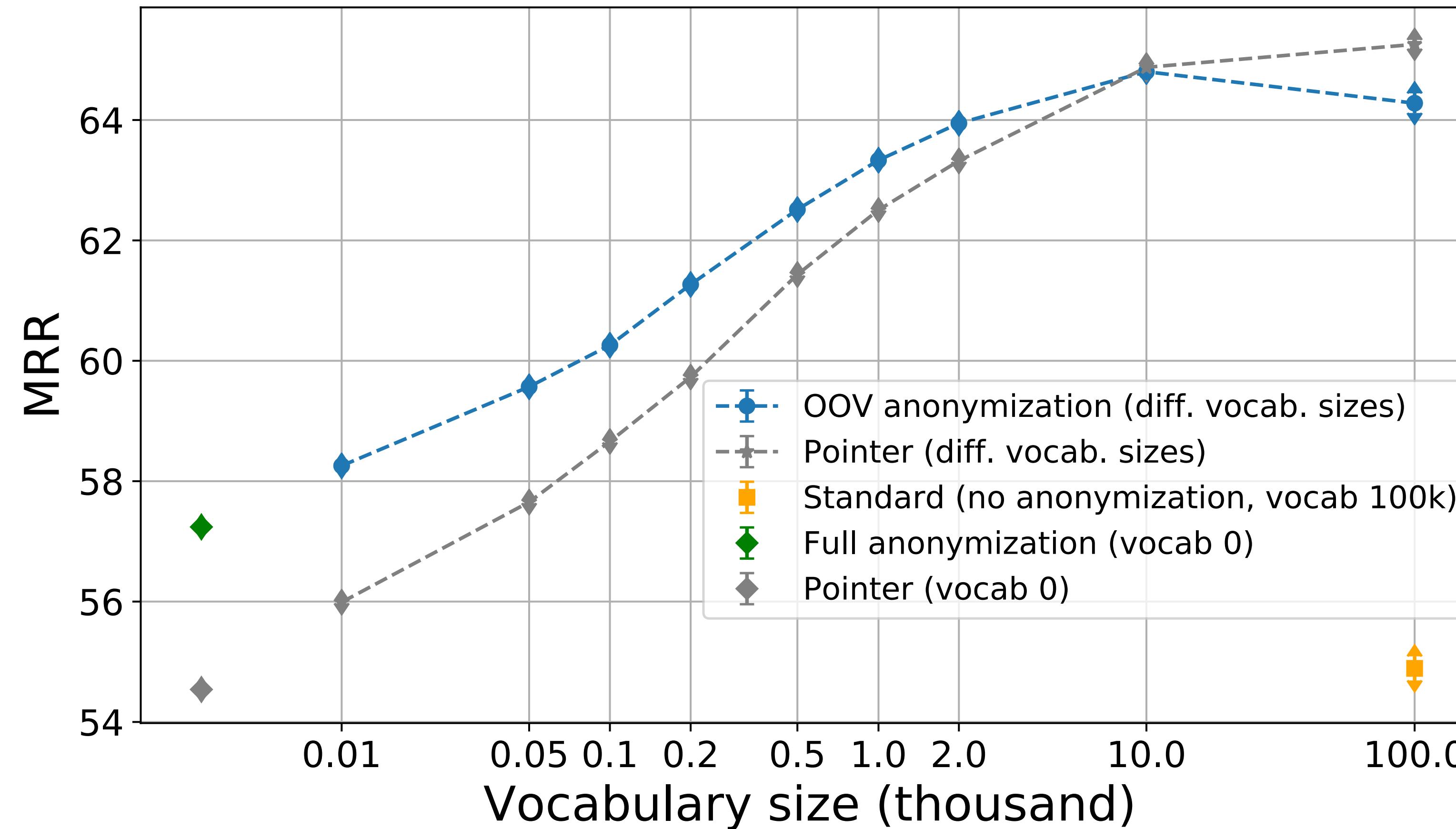


# Results: variable misuse task, Javascript150k

Same setting as below, sequential relative attention

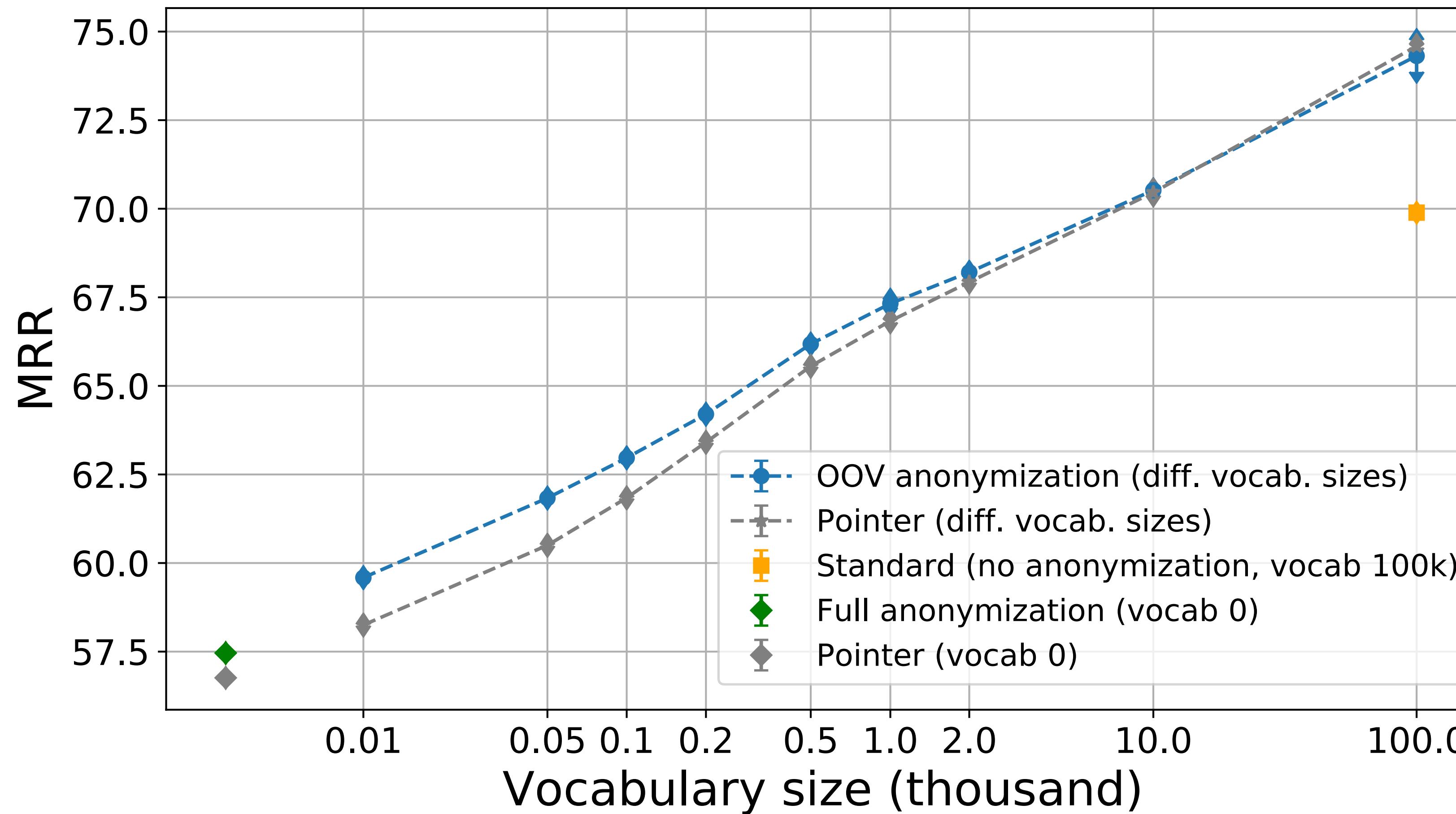


# Results: code completion task, Python150k



The proposed OOV anonymization is much simpler in implementation than the pointer mechanism

# Results: code completion task, Javascript150k



The proposed  
OOV anonymization  
is much simpler in  
implementation than  
the pointer  
mechanism

# Neural Code Completion with Anonymized Variable Names

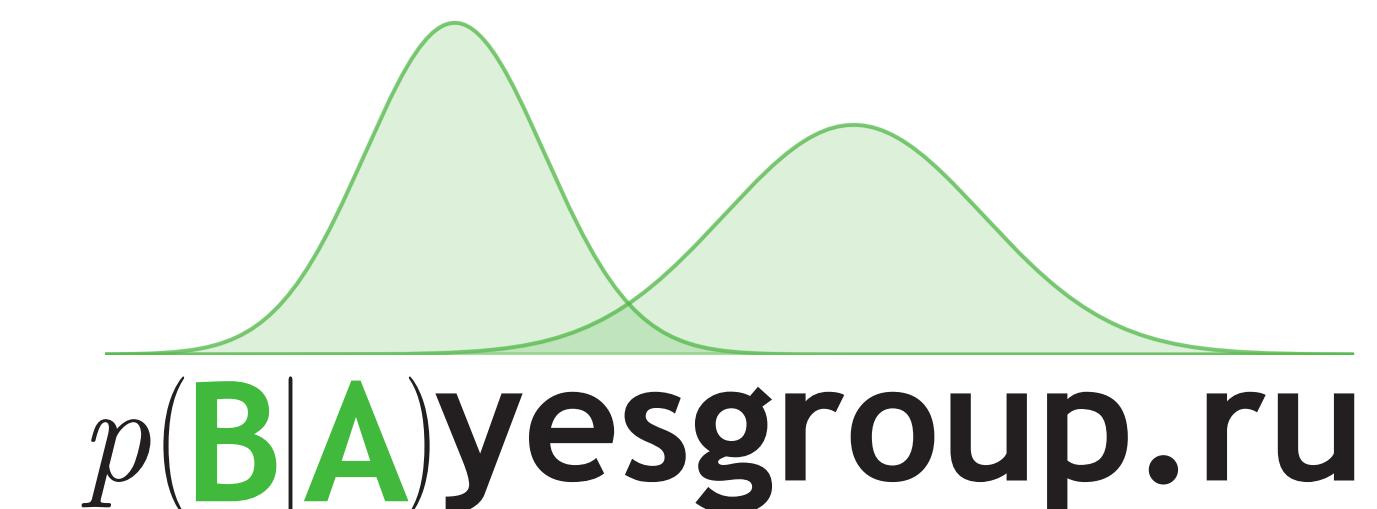
Nadezhda Chirkova

Higher School of Economics, Samsung-HSE Laboratory  
Moscow, Russia



NATIONAL RESEARCH  
UNIVERSITY

**SAMSUNG**  
**Research**



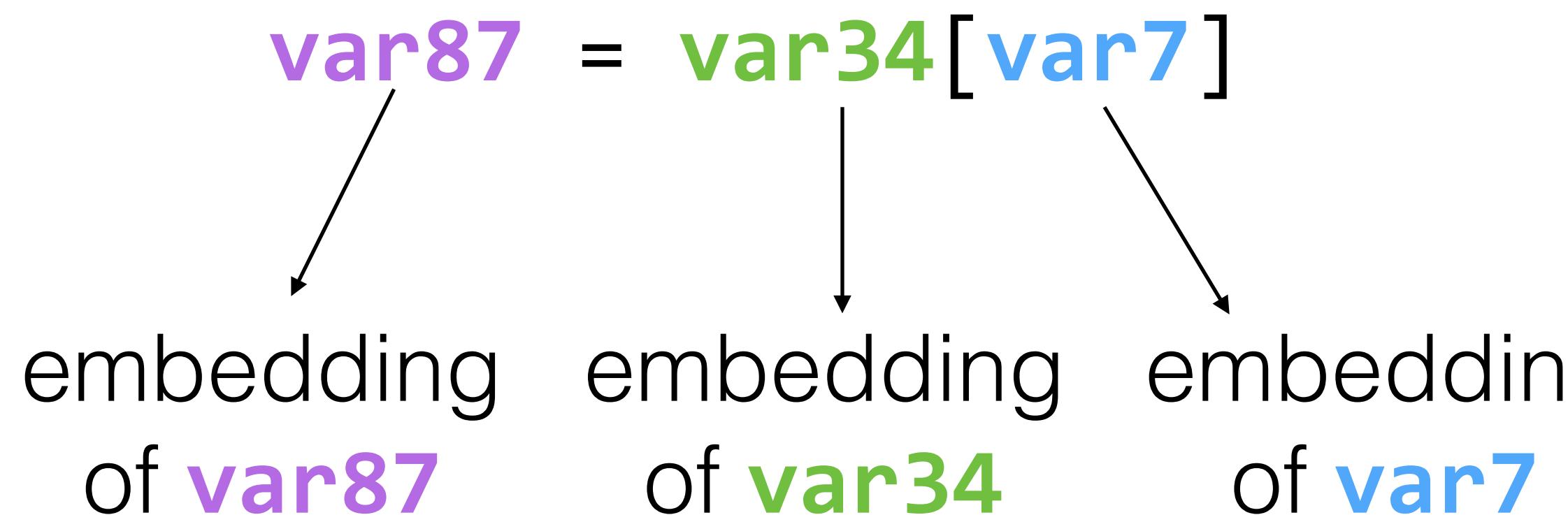
# Motivation

- Two different anonymizations result in (formally) different predictions:

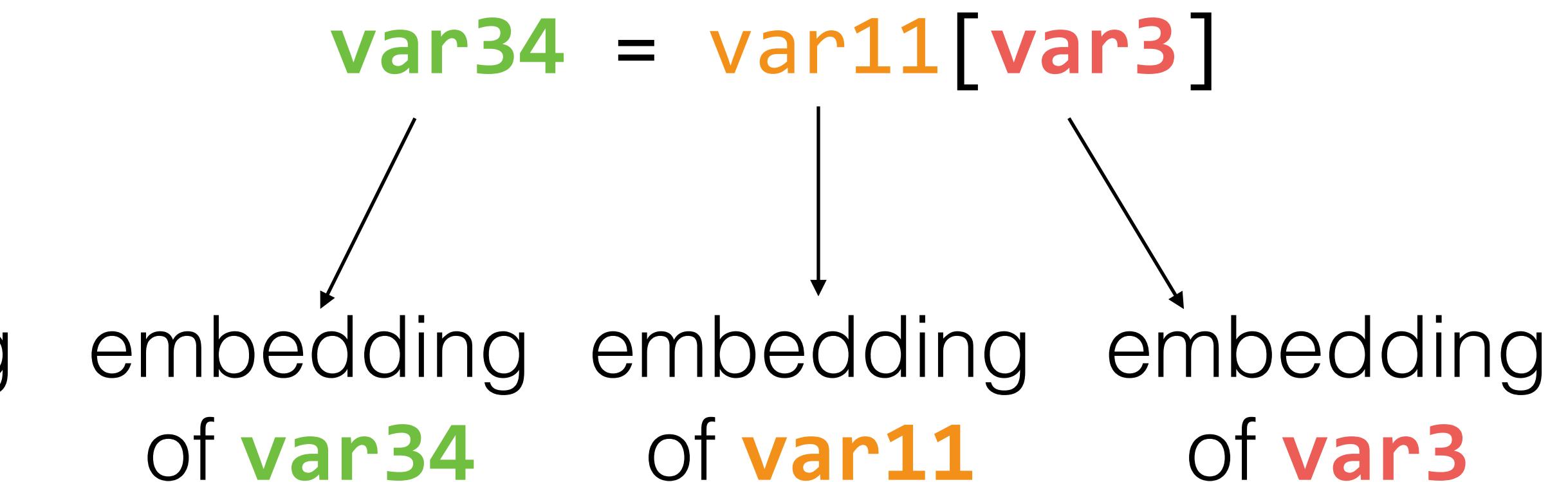
Source:

```
elem = lst[idx]
```

Anonymization 1:



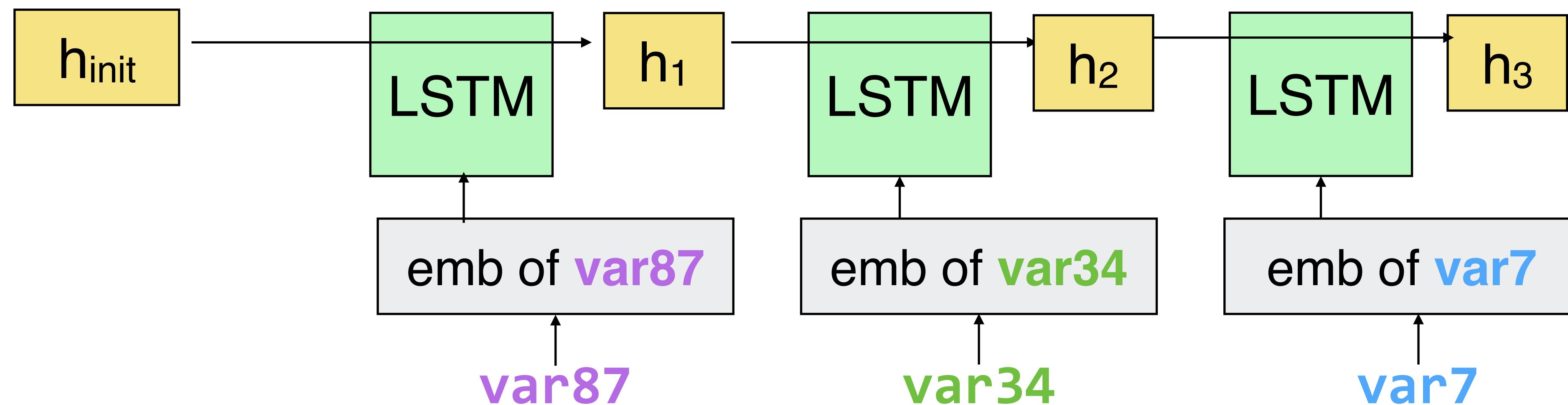
Anonymization 2:



- OK for Transformers, not OK for RNNs

# Goal

- Goal: to develop RNN model that is invariant to identifier renaming,  
i. e. renaming all the identifiers does not change the model's prediction
- Baseline: static embeddings (conventional embedding layer)
  - does not fulfill the invariance property



# Goal

- Goal: to develop RNN model that is invariant to identifier renaming,  
i. e. renaming all the identifiers does not change the model's prediction
- Baseline: static embeddings (conventional embedding layer)
  - does not fulfill the invariance property
- Our model: dynamic embeddings to replace the conventional embedding layer

# Dynamic embeddings: general idea

Initialize the embedding of each variable with  $DE_{initial}$ , update the embedding after each use of the variable

## Input    Dynamic embeddings

var5     $DE_{var5} = DE_{initial}$

var12     $DE_{var12} = DE_{initial}$

var31     $DE_{var31} = DE_{initial}$

var1     $DE_{var1} = DE_{initial}$

var5     $DE_{var5} = UpdDE(H_4, DE_{var5})$

var1     $DE_{var1} = UpdDE(H_5, DE_{var1})$

var25     $DE_{var25} = DE_{initial}$

...    ...

## Hidden state

$$H_0 = H_{init}$$

$$H_1 = UpdH(DE_{var5}, H_0)$$

$$H_2 = UpdH(DE_{var12}, H_1)$$

$$H_3 = UpdH(DE_{var31}, H_2)$$

$$H_4 = UpdH(DE_{var1}, H_3)$$

$$H_5 = UpdH(DE_{var5}, H_4)$$

$$H_6 = UpdH(DE_{var1}, H_5)$$

$$H_7 = UpdH(DE_{var25}, H_6)$$

...

```
var78 = []
for var5 in var12:
    var31 = []
    for var1 in var5:
        if var1 in var25:
            var31 += [var1]
    var78 += [var31]
```

# Dynamic embeddings: algorithm

Input data:

$[(t_1, \tilde{v}_1), \dots, (t_L, \tilde{v}_L)]$

— sequence of  
(type, anon. value) pairs

1	Assign	<empty>
2	NameStore	var7
3	SubscriptLoad	<empty>
4	NameLoad	var3
5	Index	<empty>
6	NameLoad	var1

# Dynamic embeddings: algorithm

Input data:

$[(t_1, \tilde{v}_1), \dots, (t_L, \tilde{v}_L)]$   
— sequence of  
(type, anon. value) pairs

Notation:

$e_{\tilde{v}_i, j}$  — dyn. emb. of an.value  $\tilde{v}_i$   
at timestep  $j$   
 $e_{t_i}$  — (static) emb. of type  $t_i$   
 $h_i$  — hidden state as timestep  $i$

Initialization

Dyn. embs.:

$$e_{\tilde{v}_i, 0} = e_{\text{init}}, \quad i = 1, \dots, L$$

RNN hidden state:

$$h_0 = h_{\text{init}}$$

# Dynamic embeddings: algorithm

Input data:

$[(t_1, \tilde{v}_1), \dots, (t_L, \tilde{v}_L)]$   
— sequence of  
(type, anon. value) pairs

One timestep  $i-1 \rightarrow i$ :

Update current dyn. emb.:  $e_{\tilde{v}_i, i} = \text{LSTM}_{\text{dyn}}(h_{i-1}, e_{t_i}; e_{\tilde{v}_i, i-1})$

Other dyn. embs. stay the same:  $e_{\tilde{v}, i} = e_{\tilde{v}, i-1}, \quad \tilde{v} \neq \tilde{v}_i$

Update RNN state:  $h_i = \text{LSTM}_{\text{main}}(e_{\tilde{v}_i, i-1}, e_{t_i}; h_{i-1})$

Notation:

$e_{\tilde{v}_i, j}$  — dyn. emb. of an.value  $\tilde{v}_i$  at timestep  $j$   
 $e_{t_i}$  — (static) emb. of type  $t_i$   
 $h_i$  — hidden state as timestep  $i$

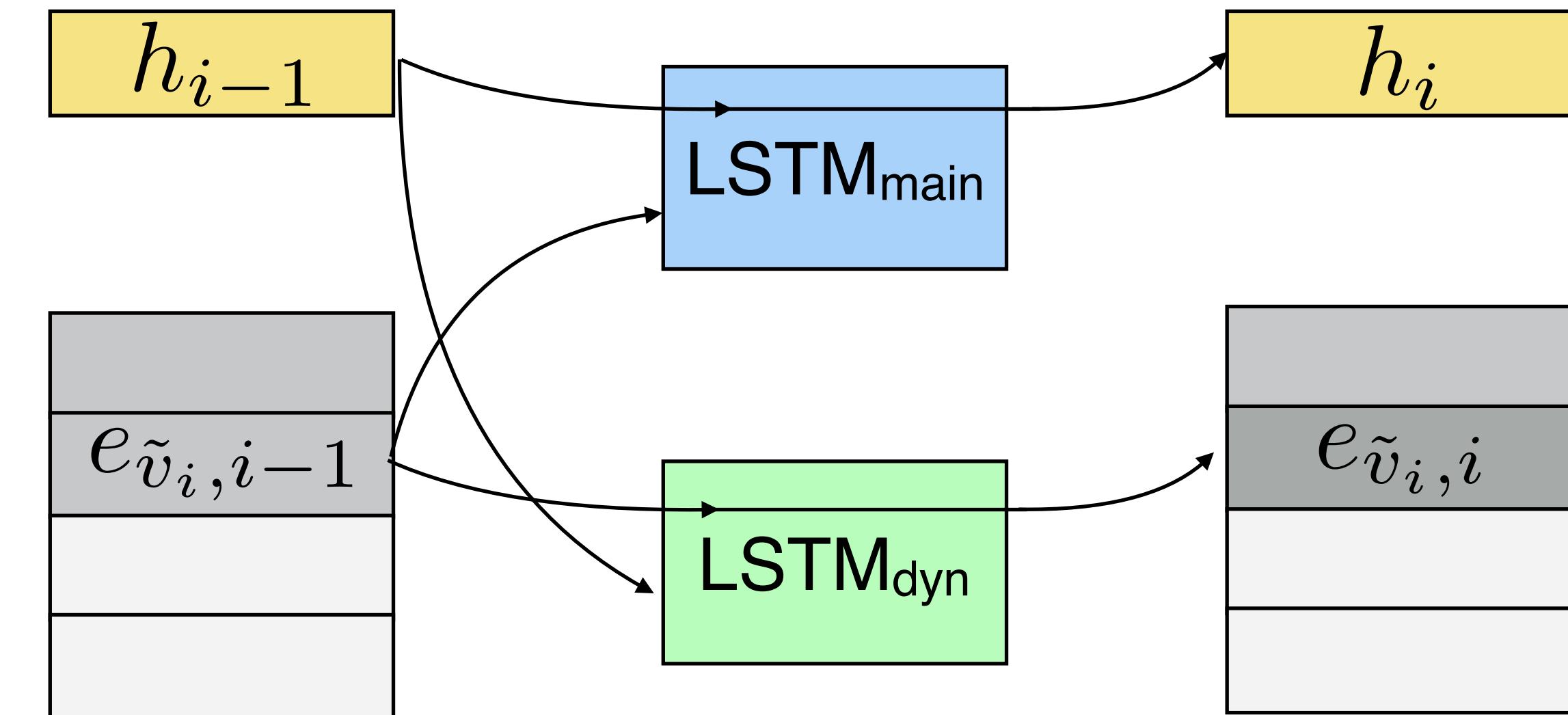
# Dynamic embeddings: algorithm

Input data:

$$[(t_1, \tilde{v}_1), \dots, (t_L, \tilde{v}_L)]$$

— sequence of  
(type, anon. value) pairs

One timestep  $i-1 \rightarrow i$ :



Update current dyn. emb.:  $e_{\tilde{v}_i, i} = \text{LSTM}_{\text{dyn}}(h_{i-1}, e_{t_i} ; e_{\tilde{v}_i, i-1})$

Other dyn. embs. stay the same:  $e_{\tilde{v}, i} = e_{\tilde{v}, i-1}, \quad \tilde{v} \neq \tilde{v}_i$

Update RNN state:  $h_i = \text{LSTM}_{\text{main}}(e_{\tilde{v}_i, i-1}, e_{t_i} ; h_{i-1})$

# Results: code completion, anonymized setting

3 variants of base architecture: plain LSTM, LSTM+attention, LSTM+pointer

2 datasets: Python150k (Py) and Javascript150k (JS)

Metric: accuracy

Setup of (Li18) + custom train / test split

Type embedding: 300

Value embedding: 1200 (standard, stat. emb)  
500 (dyn. emb.)

Model	Py:	LSTM	LSTM+at	LSTM+pt	JS:	LSTM	LSTM+at	LSTM+pt
W/o vars.		47.34	47.33	60.65		45.22	45.21	58.80
Stat. emb.		55.76	59.74	60.28		51.80	56.26	57.67
Dyn. emb.		<b>66.35</b>	<b>66.79</b>	<b>66.90</b>		<b>61.69</b>	<b>62.86</b>	<b>62.85</b>
Standard		61.62	63.73	64.69		62.03	64.28	65.05

# Results: code completion, both settings

3 variants of base architecture: plain LSTM, LSTM+attention, LSTM+pointer

2 datasets: Python150k (Py) and Javascript150k (JS)

Metric: accuracy

Setup of (Li18) + custom train / test split

Type embedding: 300

Value embedding: 1200 (standard, stat. emb.)

500 (dyn. emb., mixed emb.)

Model	Py:	LSTM	LSTM+at	LSTM+pt	JS:	LSTM	LSTM+at	LSTM+pt
W/o vars.		47.34	47.33	60.65		45.22	45.21	58.80
Stat. emb.		55.76	59.74	60.28		51.80	56.26	57.67
Dyn. emb.		<b>66.35</b>	<b>66.79</b>	<b>66.90</b>		<b>61.69</b>	<b>62.86</b>	<b>62.85</b>
Standard		61.62	63.73	64.69		62.03	64.28	65.05
Mixed. emb.		<b>67.18</b>	N/A	<b>68.61</b>		<b>63.96</b>	N/A	<b>65.56</b>

# Results: code completion, OOV anonymization

Standard embeddings for in-vocabulary identifiers  
and dynamic embeddings for OOV identifiers

Type embedding: 300

Value embedding: 1200 (standard)  
500 (dyn. emb. — OOV)

Model	Py	JS
Standard	64.69	65.05
Proposed	<b>66.01</b>	<b>65.26</b>

Vocabulary: {np, sin}

Input:

my\_y = np.sin(my\_x) + my\_x

OOV anonymization:

VAR1 = np.sin(VAR2) + VAR2

# Results: variable misuse, anonymized setting

Base architecture: bidirectional LSTM (dyn. emb. are also bidirectional)

2 datasets: Python150k (Py) and Javascript150k (JS)

# Metric: joint loc. and repair accuracy

Type embedding: 300      Value embedding: 1200 (standard, stat. emb)  
                                  500 (dyn. emb)

Model	Py	JS
Stat. emb.	24.20	23.89
Dyn. emb.	<b>54.39</b>	<b>46.38</b>
Standard	49.38	41.59

# Summary

- Dynamic embeddings: invariant to any renaming identifiers
- Dynamic embeddings significantly outperform static embeddings in the anonymized setting
- Dynamic embeddings in the full data setting: mixed embeddings outperform standard model; dynamic embeddings help to process OOV identifiers more effectively