

Generative Code Modeling With Graphs

MarcBrockschmidt, MiltiadisAllamanis, AlexanderGaunt, OleksandrPolozov

From Microsoft

ICLR-2019

Main Tasks

- Typical:
a line of code
- New Task: **more logic**
part of code

```
int ilOffsetIdx =  
    Array.IndexOf(sortedILOffsets, map.ILOffset);  
int nextILOffsetIdx = ilOffsetIdx + 1;  
int nextMapILOffset =  
    nextILOffsetIdx < sortedILOffsets.Length  
    ? sortedILOffsets[nextILOffsetIdx]  
    : int.MaxValue;
```

How to Judge?

- Perplexity

- $H(p) = -\sum p(x) \log p(x)$
- Perplexity = $2^{\{H(p)\}}$

$$\begin{aligned} \text{perplexity}(S) &= p(w_1, w_2, w_3, \dots, w_m)^{-1/m} \\ &= \sqrt[m]{\prod_{i=1}^m \frac{1}{p(w_i | w_1, w_2, \dots, w_{i-1})}} \end{aligned}$$

- Beam Search

- Expand the most promising node
- Reduce memory requirements

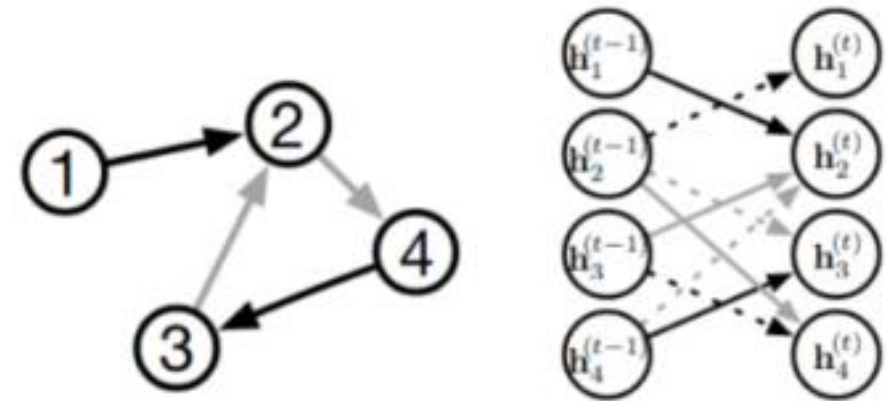
Prerequisite

- Embedding
 - Dense vector of the same length → dimensional reduction
- Attribute Grammar (context-free grammar)
 - Inherited
 - Synthesized
- Graph Neural Network
- Encoder-Decoder
- Attention

```
Expr1 → Expr2 + Term [ Expr1.value = Expr2.value + Term.value ]  
Expr → Term [ Expr.value = Term.value ]  
Term1 → Term2 * Factor [ Term1.value = Term2.value * Factor.value ]  
Term → Factor [ Term.value = Factor.value ]  
Factor → "(" Expr ")" [ Factor.value = Expr.value ]  
Factor → integer [ Factor.value = strToInt(integer.str) ]
```

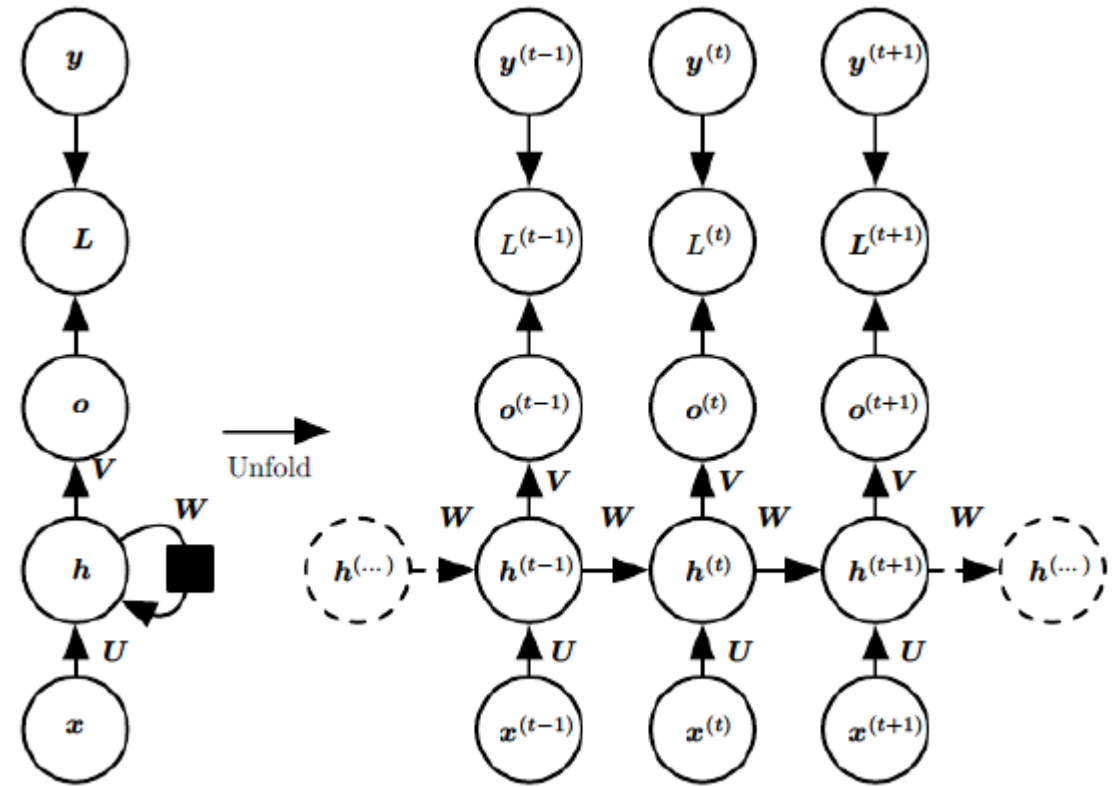
Prerequisite: GNN

- semi-supervised
- message passing
- neighborhood aggregation
- Deep walk – embedding (KDD' 14)
 - Random sequence
 - Skip-gram (stride: 1)



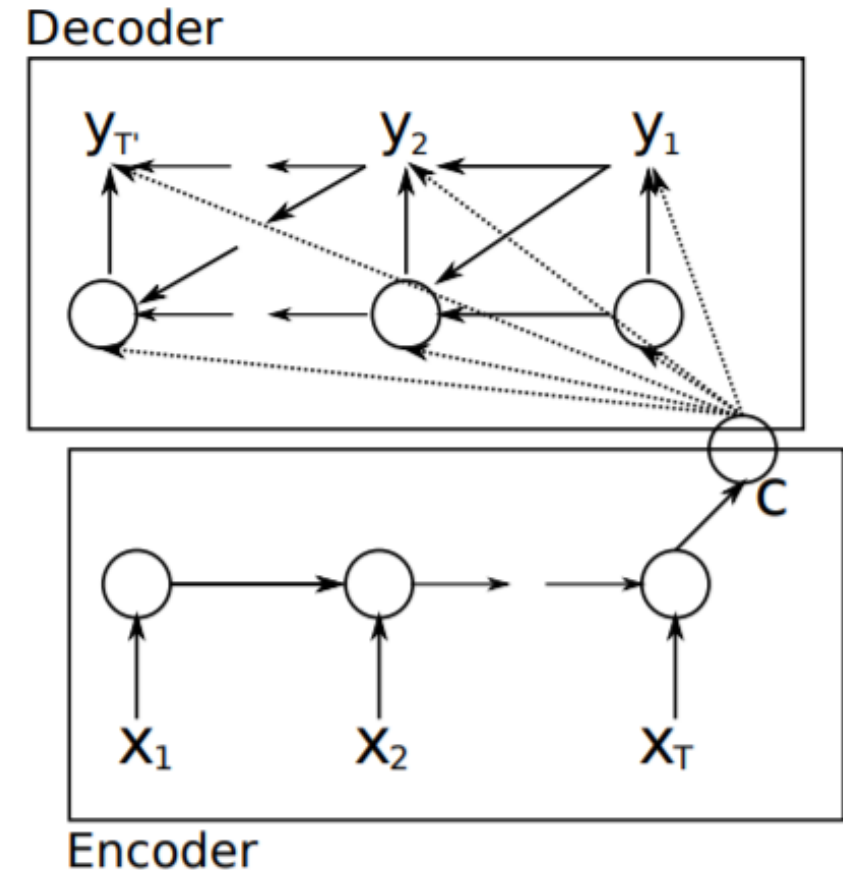
Prerequisite: RNN

- Influenced by what it has learnt from the past
- Learnt from prior input while generating output
- Bidirectional RNNs
 - Know about the past and the future



Prerequisite: Encoder – Decoder

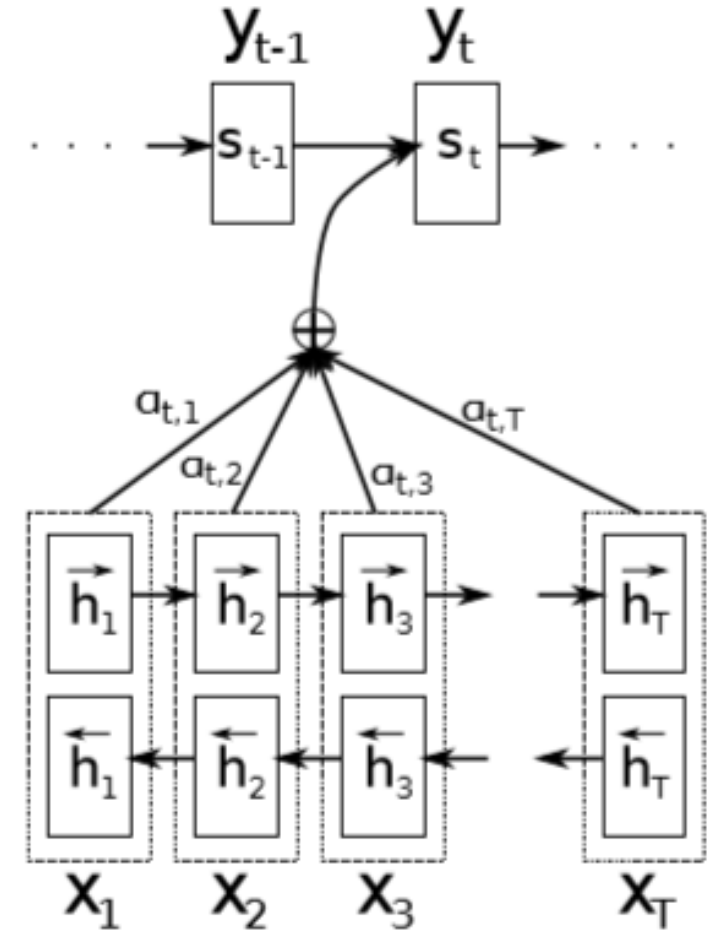
- Encoder
 - A stack of several recurrent units
- Intermediate vector
 - Initial hidden state of the decoder part of the model
- Decoder
 - accepts a hidden state from the previous unit
 - produces and output as well as its own hidden state



$$p(a|c) = \prod_t p(a_t|c, a_{<t})$$

Prerequisite: Attention

- Why attention is needed?
 - fixed-length context vector
 - forget the earlier parts of the sequence
- How it work?
 - Compute a score each encoder state
 - Compute the attention weights
 - Compute the context vector
 - Concatenate context vector with output of previous time step
 - Decoder Output



Related Works

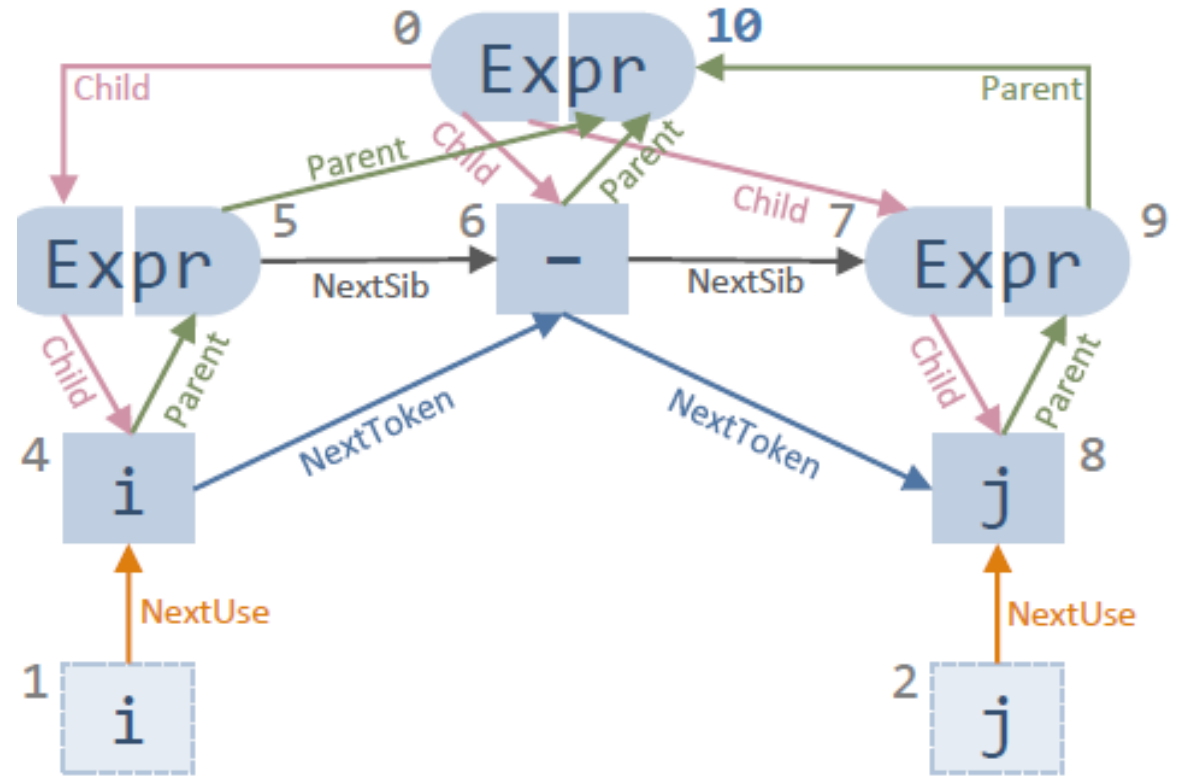
- Natural Language Processing
 - Can not distinguish unlikely from likely
 - Fail to be semantically relevant
 - Wrong on human inspection
- Abstract Syntax Tree
 - Costly on calculation resources
 - Modify model to fit the AST

→ New method: sequence

```
expression: BinaryExpression {  
  type: "BinaryExpression"  
  start: 0  
  end: 7  
  - left: BinaryExpression {  
    type: "BinaryExpression"  
    start: 1  
    end: 4  
    + left: Literal {type, start, end, value, raw}  
    operator: "+"  
    + right: Literal {type, start, end, value, raw}  
  }  
  operator: "*"   
  - right: Literal {  
    type: "Literal"  
    start: 6  
    end: 7  
    value: 3  
    raw: "3"  
  }  
}
```

Overview

- Build basic AST and graph
- Get context information & representation of root node
- Keep expanding node
- Update (MLE)



Expand Node

- Common generation strategy
- Expand AST & graph

$$\mathbf{h}_v = g(\text{emb}(\ell_v), \sum_{\langle u_i, t_i, v \rangle \in \mathcal{E}_v} f_{t_i}(\mathbf{h}_{u_i}))$$

Algorithm 1 Pseudocode for Expand

Input: Context c , partial AST a , node v to expand

```
1:  $\mathbf{h}_v \leftarrow \text{getRepresentation}(c, a, v)$ 
2:  $rhs \leftarrow \text{pickProduction}(v, \mathbf{h}_v)$ 
3: for child node type  $\ell \in rhs$  do
4:    $(a, u) \leftarrow \text{insertChild}(a, \ell)$ 
5:   if  $\ell$  is nonterminal type then
6:      $a \leftarrow \text{Expand}(c, a, u)$ 
7: return  $a$ 
```

Edges

- Typed directed edges
- Embed different edge type into different space

Algorithm 2 Pseudocode for ComputeEdge

Input: Partial AST a , node v

```
1: Edge set  $\mathcal{E} \leftarrow \emptyset$ 
2: if  $v$  is inherited then
3:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle \text{parent}(a, v), \text{Child}, v \rangle\}$ 
4:   if  $v$  is terminal node then
5:      $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle \text{lastToken}(a, v), \text{NextToken}, v \rangle\}$ 
6:     if  $v$  is variable then
7:        $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle \text{lastUse}(a, v), \text{NextUse}, v \rangle\}$ 
8:   if  $v$  is not first child then
9:      $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle \text{lastSibling}(a, v), \text{NextSib}, v \rangle\}$ 
10: else
11:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle u, \text{Parent}, v \rangle \mid u \in \text{children}(a, v)\}$ 
12:    $\mathcal{E} \leftarrow \mathcal{E} \cup \{\langle \text{inheritedAttr}(v), \text{InhToSyn}, v \rangle\}$ 
13: return  $\mathcal{E}$ 
```

pickProduction

- choose production \rightarrow classification
- $\text{pickProduction} = \operatorname{argmax} P(\text{rule} | l_v, h_v)$

pickVariable

- Pointer Network
 - sequence-to-sequence & attention
 - output of pointer networks is discrete and correspond to positions in the input sequence
 - the number of target classes in each step of the output depends on the length of the input
- $\text{pickVariable}(\mathbb{I}, h_v) = \operatorname{argmax} P(\text{var} | h_v)$

pickLiteral

- Copy one of the tokens
 - For each token compute a score
 - Only way to generate a UNK token
-
- $\text{pickLiteral}(\text{token}, h_v) = \operatorname{argmax} p(\text{lit} | h_v)$

Results

Model	Test (from seen projects)				Test-only (from unseen projects)			
	Perplexity	Well-Typed	Acc@1	Acc@5	Perplexity	Well-Typed	Acc@1	Acc@5
<i>PHOG</i> [†]	–	–	34.8%	42.9%	–	–	28.0%	37.3%
<i>Seq</i> → <i>Seq</i>	87.48	32.4%	21.8%	28.1%	130.46	23.4%	10.8%	16.8%
<i>Seq</i> → <i>NAG</i>	6.81	53.2%	17.7%	33.7%	8.38	40.4%	8.4%	15.8%
<i>G</i> → <i>Seq</i>	93.31	40.9%	27.1%	34.8%	28.48	36.3%	17.2%	25.6%
<i>G</i> → <i>Tree</i>	4.37	49.3%	26.8%	48.9%	5.37	41.2%	19.9%	36.8%
<i>G</i> → <i>ASN</i>	2.62	78.7%	45.7%	62.0%	3.03	74.7%	32.4%	48.1%
<i>G</i> → <i>Syn</i>	2.71	84.9%	50.5%	66.8%	3.48	84.5%	36.0%	52.7%
<i>G</i> → <i>NAG</i>	2.56	86.4%	52.3%	69.2%	3.07	84.5%	38.8%	57.0%

```
int methParamCount = 0;
if (paramCount > 0) {
    IParameterTypeInfo[] moduleParamArr =
        GetParamTypeInformations(Dummy.Signature, paramCount);
    methParamCount = moduleParamArr.Length;
}
if (paramCount > methParamCount) {
    IParameterTypeInfo[] moduleParamArr =
        GetParamTypeInformations(Dummy.Signature,
            paramCount - methParamCount);
}
```

G → *NAG*:

paramCount > methParamCount (34.4%)
 paramCount == methParamCount (11.4%)
 paramCount < methParamCount (10.0%)

G → *ASN*:

paramCount == 0 (12.7%)
 paramCount < 0 (11.5%)
 paramCount > 0 (8.0%)

```
public static String URIToPath(String uri) {
    if (System.Text.RegularExpressions
        .Regex.IsMatch(uri, "^file:\\\\[a-z,A-Z]:")) {
        return uri.Substring(6);
    }
    if (uri.StartsWith(@"file:")) {
        return uri.Substring(5);
    }
    return uri;
}
```

G → *NAG*:

uri.Contains(UNK_STRING_LITERAL) (32.4%)
 uri.StartsWith(UNK_STRING_LITERAL) (29.2%)
 uri.HasValue() (7.7%)

G → *Syn*:

uri == UNK_STRING_LITERAL (26.4%)
 uri == "" (8.5%)
 uri.StartsWith(UNK_STRING_LITERAL) (6.7%)

Limitations

- Limited types
- Training data noise
- Bad performance on new project

Prospect

- Unified evaluation index
- Qualified dataset

Thanks