

# Machine Learning Approaches in Programming Language Type Inference

Presentation: Jiancheng Qian

# Overview

- Code Naturalness and Type Inference
- Related Work
  - *Deep Learning Type Inference* (Hellendorn, FSE 2018)
  - *Typilus: Neural Type Hints* (Allamanis, PLDI 2020)
  - *LambdaNet: probabilistic type inference using graph neural networks* (Jiayi Wei, ICLR 2020)
- Limitation, Challenges and My Own Opinions

# Code Naturalness and Type Inference

- Code Naturalness:
  - **The naturalness hypothesis:** *Programming languages, in theory, are complex, flexible and powerful, but, “natural” programs, the ones that real people actually write, are mostly simple and rather repetitive; thus they have usefully predictable statistical properties that can be captured in statistical language models and leveraged for software engineering tasks.*
    - *A. Hindle, On the Naturalness of Software, 2018*
  - **Machine learning ‘Language’ models:**
    - Code generating models
    - Representational models
    - Pattern mining models

# Code Naturalness and Type Inference

- Applications of code naturalness:
  - Recommender systems
  - Code completion
  - Finding defects
  - Code to natural language/natural language to code
  - Pattern mining
  - Code translation
  - Name/type inference
  - Traceability
  - ...

# Code Naturalness and Type Inference

- Type Inference

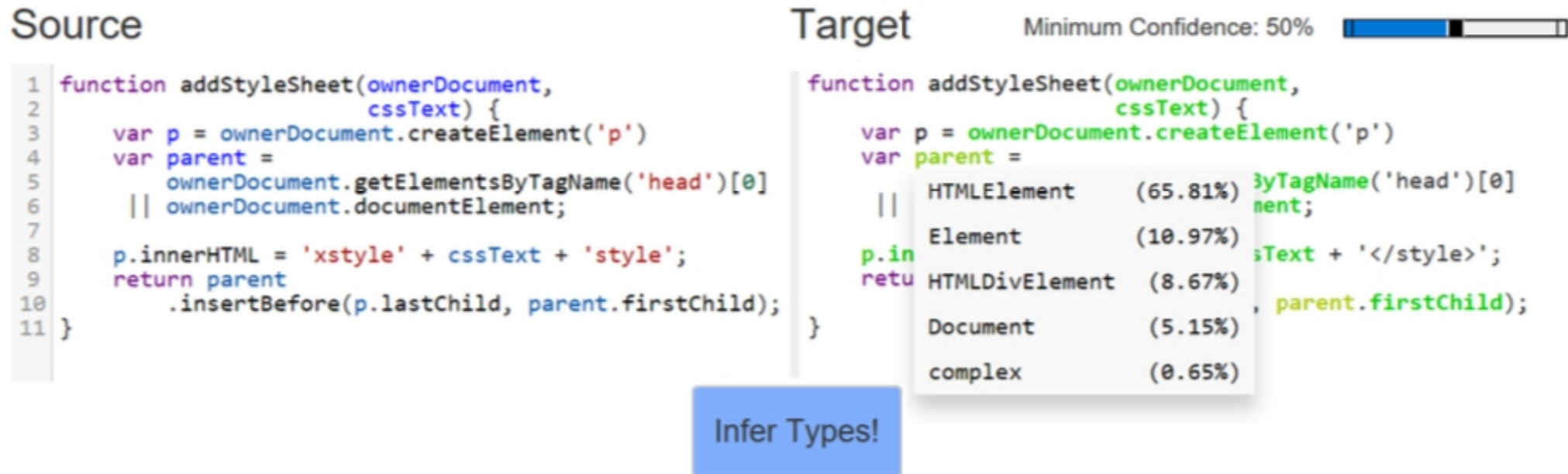


Figure: Deep Learning Type Inference (Hellendorn, FSE 2018)

# Code Naturalness and Type Inference

- Why we need type inference:
  - Dynamic languages like Python/JavaScript
  - Enhance maintainability and readability
  - Avoid 15% of the bugs, lower fault occurrence
  - Automation in code maintenance and testing
- Traditional methods:
  - Type annotation (e.g. *mypy* and *pytype*, *TypeScript*)

```
1 function addStyleSheet(ownerDocument,  
2                        cssText) {  
3     var p = ownerDocument.createElement('p')  
4     var parent =  
5         ownerDocument.getElementsByTagName('head')[0]  
6         || ownerDocument.documentElement;  
7  
8     p.innerHTML = 'x<style>' + cssText + '</style>';  
9     return parent.insertBefore(p.lastChild, parent.firstChild);  
10 }
```

```
1 function addStyleSheet(ownerDocument : Document,  
2                        cssText : string) : any {  
3     var p : HTMLElement = ownerDocument.createElement('p')  
4     var parent : HTMLElement =  
5         ownerDocument.getElementsByTagName('head')[0]  
6         || ownerDocument.documentElement;  
7  
8     p.innerHTML : string = 'x<style>' + cssText + '</style>';  
9     return parent.insertBefore(p.lastChild, parent.firstChild);  
10 }
```

Figure: Deep Learning Type Inference (Hellendorn, FSE 2018)

# Related Work

- Similar to NLP
- Type inference neural models:
  - 1) Preprocessing: Embed the code/functions into tokens/graph
    - **Word2vec**
    - **AST**
    - **Type Dependency Graph**
  - 2) Train the neural model, obtain the vector representation of identifiers/nodes
    - **RNN**
    - **GNN**
  - 3) Prediction
    - Calculate the probability distribution of type for input tokens/nodes

# Related Work

- Machine learning models:

Title	Input	Machine Learning Model	Prediction
DeepTyper(2018)	Sequential tokens	2 Bi-RNN	
NL2Type(2019)	Vector representation of function	Bi-LSTM	
DTLPy(2019)	Vector representation of function	2 LSTM/GRU	
Typilus(2020)	AST	GNN	KNN
Type4Py(2020)	AST + Sequential tokens	RNN	KNN
LambdaNet(2020)	Type Dependency Graph	GNN	MLP
TypeWriter(2020)	AST based extraction of code + comments	RNN	Feedback directed search



# Deep Learning Type Inference (Hellendorn, FSE 2018)

- *DeepTyper*
  - First work to implement neural network in type inference
  - Training using aligned corpus (TypeScript) and prediction on JS code
  - Explore the probability for open-world type suggestion tasks

# Deep Learning Type Inference (Hellendorn, FSE 2018)

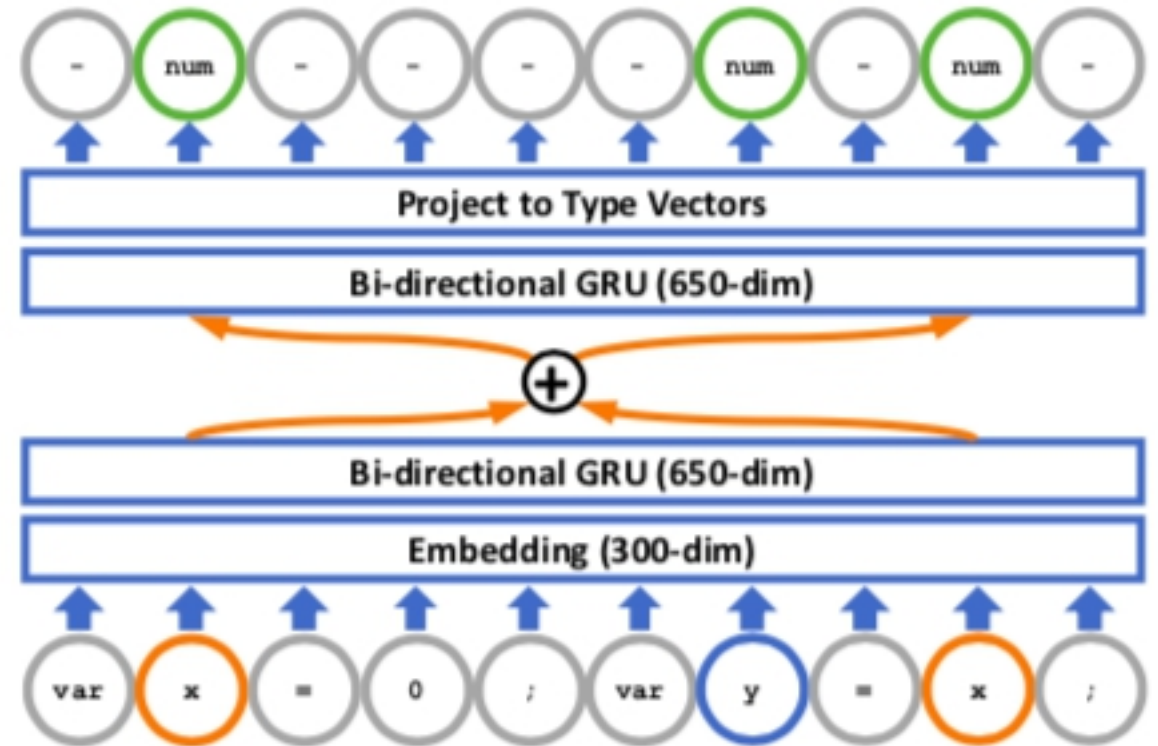
- Input a code sequence  $s_1, \dots, s_N$
- Embedding  $s_t$  to the vector representation  $\mathbf{x}_{s_t}$
- Hidden state of token  $\mathbf{h}_t = \text{RNN}(\mathbf{x}_{s_t}, \mathbf{h}_{t-1})$
- For token  $s_t$ , calculate the probability of type  $\tau$ :

$$P_{s_t}(\tau) = \frac{\exp(\hat{\mathbf{r}}_t^T \mathbf{r}_\tau + b_\tau)}{\sum_{\tau'} \exp(\hat{\mathbf{r}}_t^T \mathbf{r}_{\tau'} + b_{\tau'})}$$

- $\mathbf{r}_\tau$  is the vector representation of well-known type  $\tau$  and  $\hat{\mathbf{r}}_t$  is the projection vector of token  $s_t$ :

$$\hat{\mathbf{r}}_t = \mathbf{h}_t^{bi} + \frac{1}{|V(t)|} \sum_{i \in V(t)} \mathbf{h}_i^{bi}$$

- $V(t)$  is the locations which bound to the same identifier at location  $t$



# Typilus: Neural Type Hints (Allamanis, PLDI 2020)

- *Typilus*:
  - A graph-based deep neural network to the type prediction problem by considering source code syntax and semantics
  - Using GNN to catch the relations between distant identifiers
  - A novel training loss to evaluate the learned type features in the type space, can accurately predict types that were rare, or even unseen
  - Trained on Python codes with *mypy/pytype* corpus

# Typilus: Neural Type Hints (Allamanis, PLDI 2020)

- Overview of *Typilus*

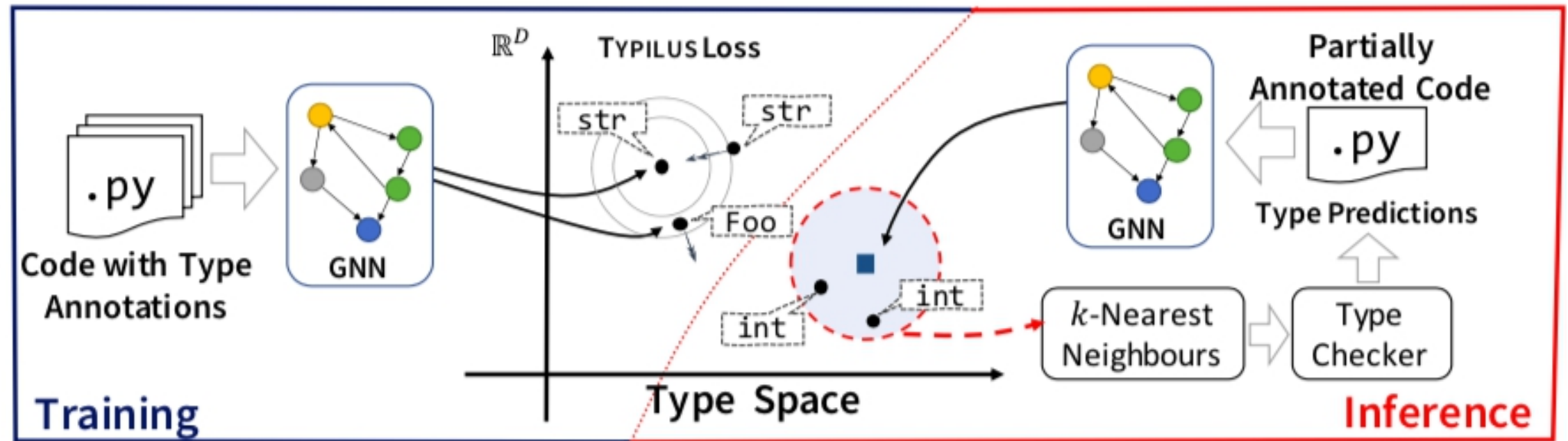


Figure: Allamanis, 2020

# Typilus: Neural Type Hints (Allamanis, PLDI 2020)

- AST(abstract syntax tree) embedding for code
- A sample of `foo = get_foo(i, i+1)`

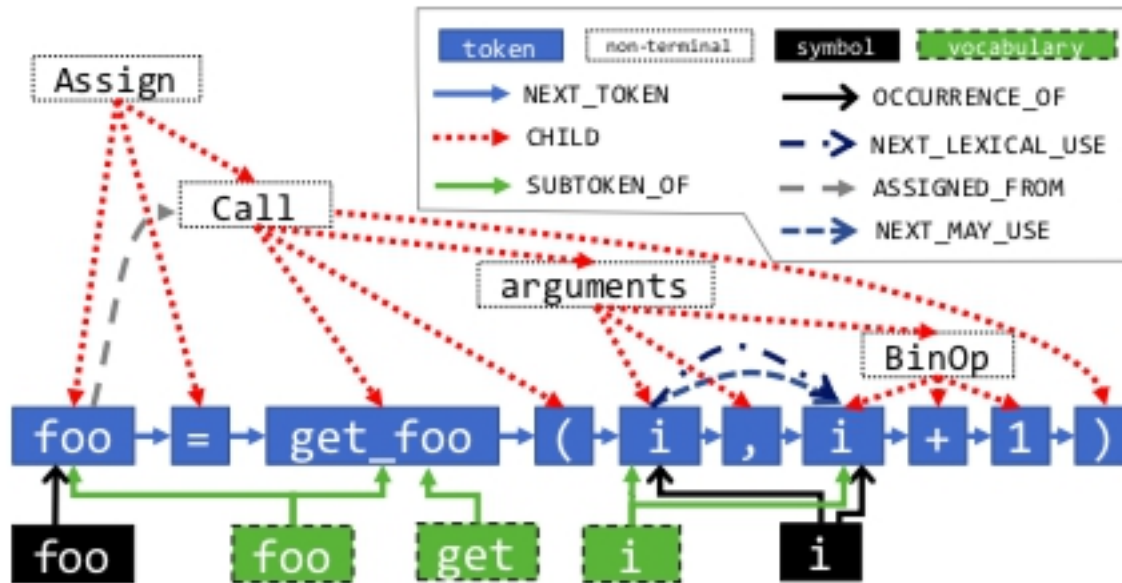


Figure: Allamanis, 2020

# Typilus: Neural Type Hints (Allamanis, PLDI 2020)

- Input code sequences  $S$  with each identifier  $s \in S$ , learn the vector representation (type embedding):  $e(S)[s] = \mathbf{r}_s \in \mathbb{R}^D$ , where  $e$  is the GNN network serves as a map
- Define the well-known types  $\mathcal{T} = \{\tau_i\}$ , learn their representation  $\tilde{\mathbf{r}}_{\tau_i}$
- Calculate the classification loss for  $s$ :

$$\mathcal{L}_{\text{CLASS}}(\mathbf{r}_s, \tau) = -\log \underbrace{\frac{\exp(\mathbf{r}_s \tilde{\mathbf{r}}_{\tau}^T + b_{\tau})}{\sum_{\tau_j \in \mathcal{T}} \exp(\mathbf{r}_s \tilde{\mathbf{r}}_{\tau_j}^T + b_{\tau_j})}}_{-\log P(s:\tau)}$$

- Similarity learning: calculate the similarity loss for  $s$ ,  $\mathbf{r}_{s^+}$  is the representation vector of same type as  $s$ , while  $\mathbf{r}_{s^-}$  is a different type and :

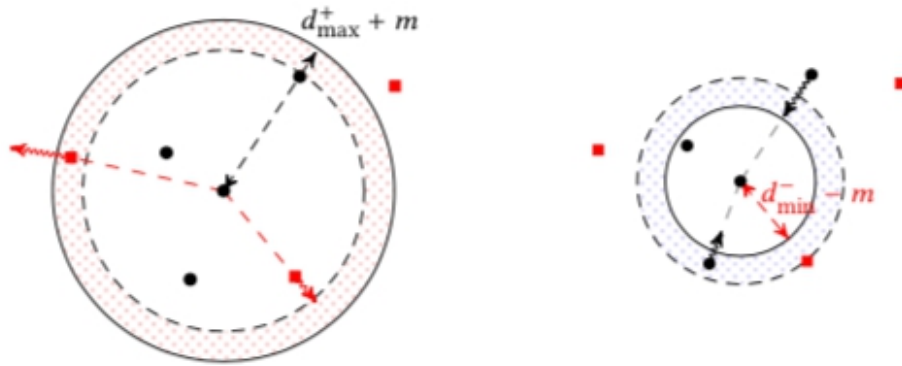
$$\mathcal{L}_{\text{SPACE}}(s) = \sum_{s_i^+ \in P_+(s)} \frac{\|\mathbf{r}_{s_i^+} - \mathbf{r}_s\|}{|P_+(s)|} - \sum_{s_i^- \in P_-(s)} \frac{\|\mathbf{r}_{s_i^-} - \mathbf{r}_s\|}{|P_-(s)|}, \quad \begin{aligned} P_+(s) &= \{x_i^+ : \|\mathbf{r}_{s_i^+} - \mathbf{r}_s\| > d_{\min}^-(s) - m\} \\ P_-(s) &= \{x_i^- : \|\mathbf{r}_{s_i^-} - \mathbf{r}_s\| < d_{\max}^+(s) + m\} \end{aligned}$$

- Combine the two loss functions for the learning objective:

$$\mathcal{L}_{\text{TYPIIUS}}(s, \tau) = \mathcal{L}_{\text{SPACE}}(s) + \lambda \mathcal{L}_{\text{CLASS}}(W \mathbf{r}_s, \text{ER}(\tau))$$

# Typilus: Neural Type Hints (Allamanis, PLDI 2020)

- The implementation of type space enables the model to predict unseen types without well-known annotation



**Figure 2.** Graphic depiction of the two terms of the similarity objective in Eq. 3. *Left:* all dissimilar points (red squares), *i.e.*  $P_-$  within distance  $d_{\max}^+ + m$  of the query point are pushed away. *Right:* all similar points (black circles) that are further than  $d_{\min}^- - m$  from the query point, *i.e.*  $P_+$ , are pulled towards it. The margin distance  $m$  is shaded.

*Figure: Allamanis, 2020*

# LambdaNet: probabilistic type inference using graph neural networks (Jiayi Wei, ICLR 2020)

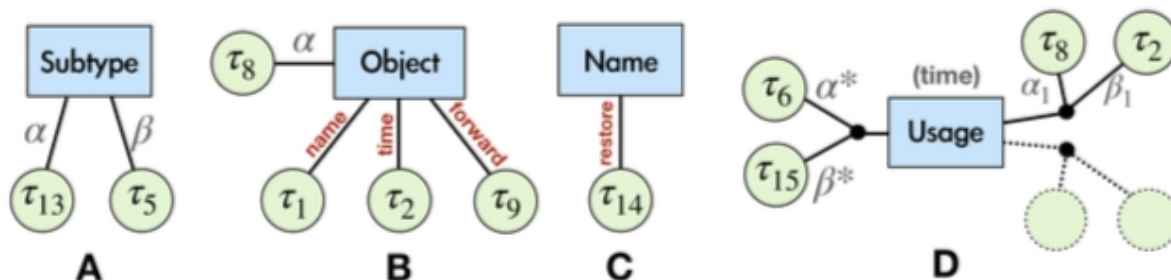
- *LambdaNet*:
  - Using a novel graph representation – *type dependency graph* for the JavaScript code
  - Using GNN to learn the features of input nodes
  - An MLP prediction layer to compute a compatibility score with well-known types for each output vector



# LambdaNet: probabilistic type inference using graph neural networks (Jiayi Wei, ICLR 2020)

- Type dependency graph

```
1  var c1:  $\tau_8$  = class MyNetwork {
2    name:  $\tau_1$ ; time:  $\tau_2$ ;
3    var m1:  $\tau_9$  = function forward(x:  $\tau_3$ , y:  $\tau_4$ ): $\tau_5$  {
4      var v1:  $\tau_{10}$  = x.concat; var v2:  $\tau_{11}$  = v1(y);
5      var v3:  $\tau_{12}$  = v2.TIMES_OP; var v4:  $\tau_{13}$  = v3(NUMBER);
6      return v4;
7    }
8  } // more classes...
9  var f1:  $\tau_{14}$  = function restore (network:  $\tau_6$ ):  $\tau_7$  {
10    var v3:  $\tau_{15}$  = network.time;
11    var v4:  $\tau_{16}$  = readNumber(STRING);
12    network.time = v4; // more code...
13  }
```



Types of hyperedges: bool, subtype, function, call, object, name, usage,...

Figure: LambdaNet: probabilistic type inference using graph neural networks, Jiayi, 2020

# Challenge in Type Inference

- Depends on existing corpus (like mypy/pytype/TypeScript)
- Predicting unknown/complicated types
  - Some types are user-defined and specialized in some projects
    - Need to know the stability among different projects/datasets
  - For some nested types like `List[List[List[int]]]`, `Optional[Map[str, int]]` are hard to predict
- Cross programming language use
  - E.g. Typilus performs bad on JS comparing to Python code
  - Need to find more universal methods since PL is developing fast
    - Focus on the dynamic type features for different language and build dependency graphs
    - The impact of network structure is still unknown

# My Own Opinions

- Recent work (since 2020) are focusing on the graph representation, while not enough research on the network structure (most are using RNN/single layer GNN)
- Working on a NAS(Neural architecture search) attempting to find the network structures implications
- We can using a unsupervised or self-supervised learning method to enhance the ability for training without an existing type annotation corpus