

# code2vec:

## Learning Distributed Representations of Code



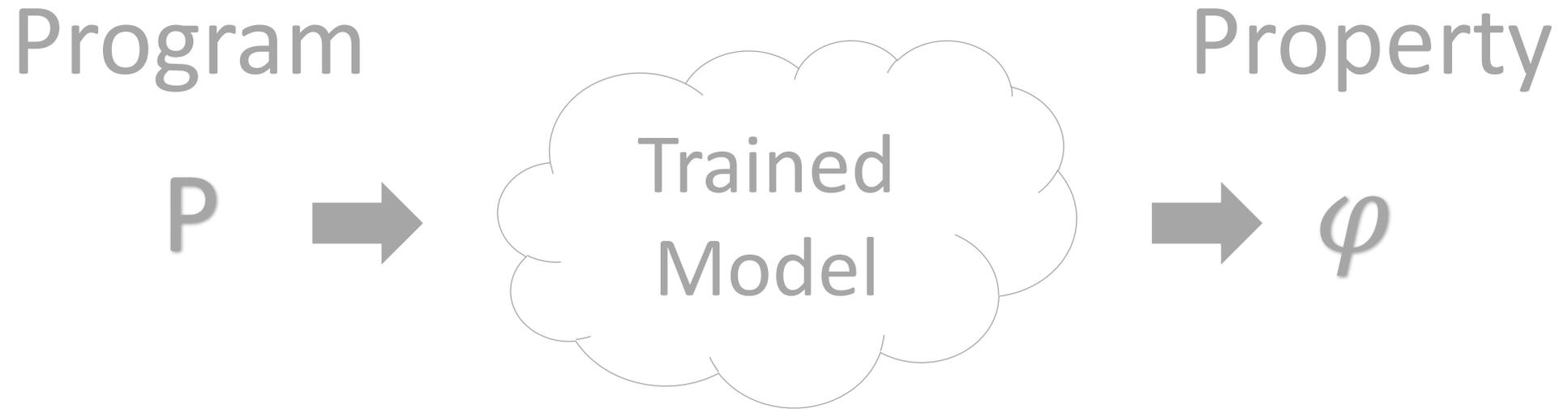
**Uri Alon, Meital Zilberstein, Omer Levy, Eran Yahav**

Facebook AI  
Research



**Technion**

# Predicting Properties of Programs



# A Motivating Example: Semantic Labeling of Code

```
String[] _____(final String[] array) {  
    final String[] newArray = new String[array.length];  
    for (int index = 0; index < array.length; index++) {  
        newArray[array.length - index - 1] = array[index];  
    }  
    return newArray;  
}
```

**reverseArray**



**77.34%**

**reverse**



**18.18%**

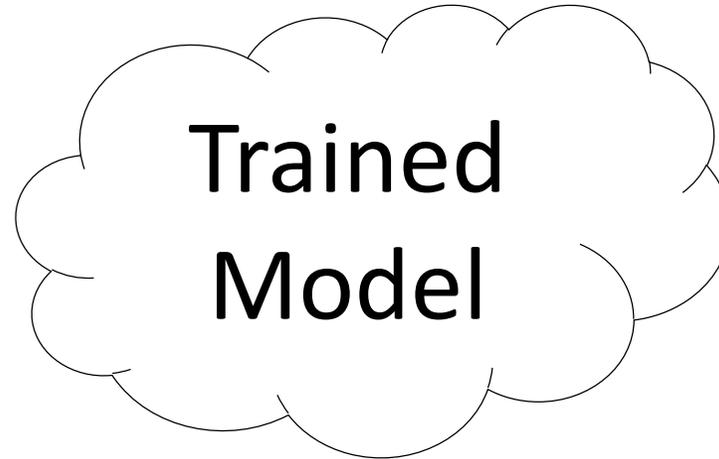
**subArray**



**1.45%**

Program

$P$



Property

$\varphi$

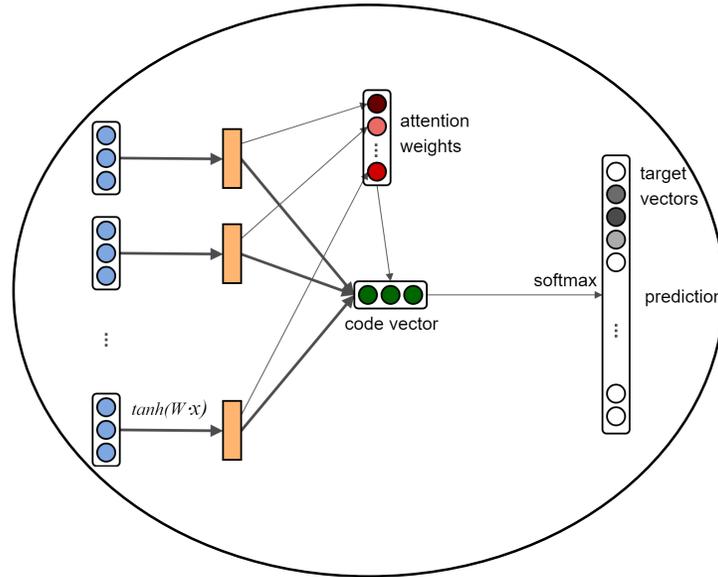
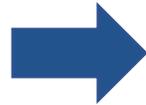
Training data  
(millions of examples):

{  
   $(P_1, \varphi_1)$   
   $(P_2, \varphi_2)$   
  ...  
   $(P_m, \varphi_m)$

Test data:  $(P', \varphi')$

# Program

$P$



# Property

$\varphi$



Training data  
(millions of examples):

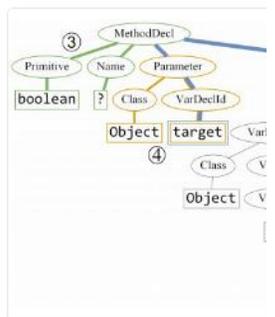
$\left\{ \begin{array}{l} (P_1, \varphi_1) \\ (P_2, \varphi_2) \\ \dots \\ (P_m, \varphi_m) \end{array} \right.$

Test data:  $(P', \varphi')$

# Code2vec: a neural network for predicting properties of code

**samim**  
@samim

code2vec: Learning Representations of  
[arxiv.org/abs/1803](https://arxiv.org/abs/1803.09473)



2:13 AM - 31 Mar 2018

22 Retweets 65 Likes

8:00 AM - 4 Apr 2018

17 Retweets 44 Likes

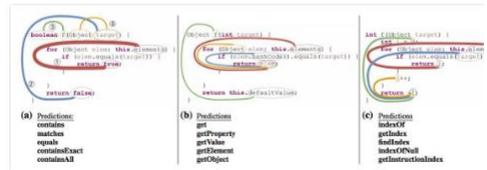
**source(d)**  
@sourcecdtech

Really interesting paper on #MLonCode with an approach based on paths in ASTs

"code2vec: Learning Distributed Representations of Code"  
by Uri Alon, Meital Zilberstein, @omerlevy\_, and @yahave

Take this as an open invitation to visit our offices 😊

[buff.ly/2GupEv4](https://buff.ly/2GupEv4)

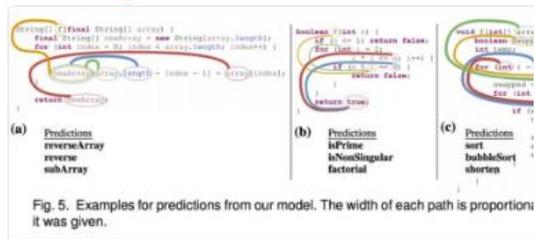


17 Retweets 44 Likes

**ML Review**  
@ml\_review

code2vec: Learning Distributed Representations of Code

[arxiv.org/abs/1803.09473](https://arxiv.org/abs/1803.09473) #Machin



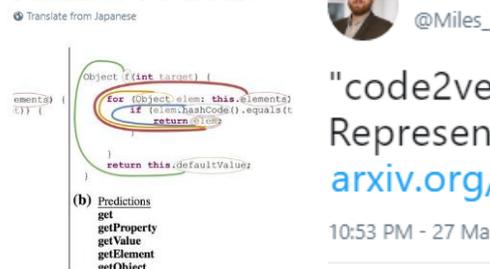
12:40 AM - 30 Mar 2018

19 Retweets 61 Likes

19 Retweets 61 Likes

**Kyosuke Nishida**  
@kyoun

code2vec (Technion)  
[arxiv.org/abs/1803.09473](https://arxiv.org/abs/1803.09473) Javaメソッドの abstract syntax treeのパス集合からその意味を表すタグを出力するアテンションモデルの中でcodeを固定ベクトルに変換, compare + toLower = compareIgnoreCase を理解可 10KのJava GitHubリポジト 14Mのメソッドで学習.



ree methods that have a similar syntactic str differences between them and manages to :hs are proportional to the attention that ea

2:55 PM - 31 Mar 2018

11 Retweets 27 Likes

**Miles Brundage**  
@Miles\_Brundage

"code2vec: Learning Distributed Representations of Code," Alon et al.:  
[arxiv.org/abs/1803.09473](https://arxiv.org/abs/1803.09473)

10:53 PM - 27 Mar 2018

21 Retweets 36 Likes

21 Retweets 36 Likes

# Code2vec: a neural network for predicting properties of code

Example application: predicting method names

```
String[] _____ (...) {  
    final String[] = ...;  
    ...  
    return newArray;  
}
```



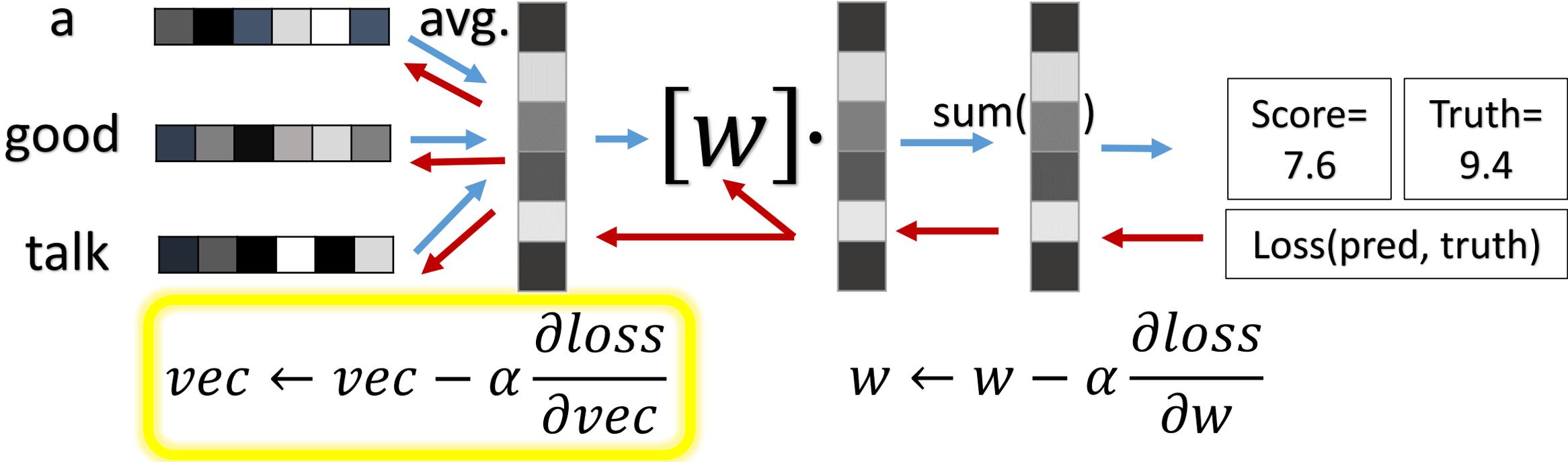
*reverseArray*

- A general approach – has many possible applications:
  - Yes/no malware, required dependencies, keywords / hashtags, clone detection...

*How does code2vec work?*

# Neural Networks

- Sequences of simple algebraic functions over vectors and matrices
- A simple example: Predict how positive is a given sentence (regression)



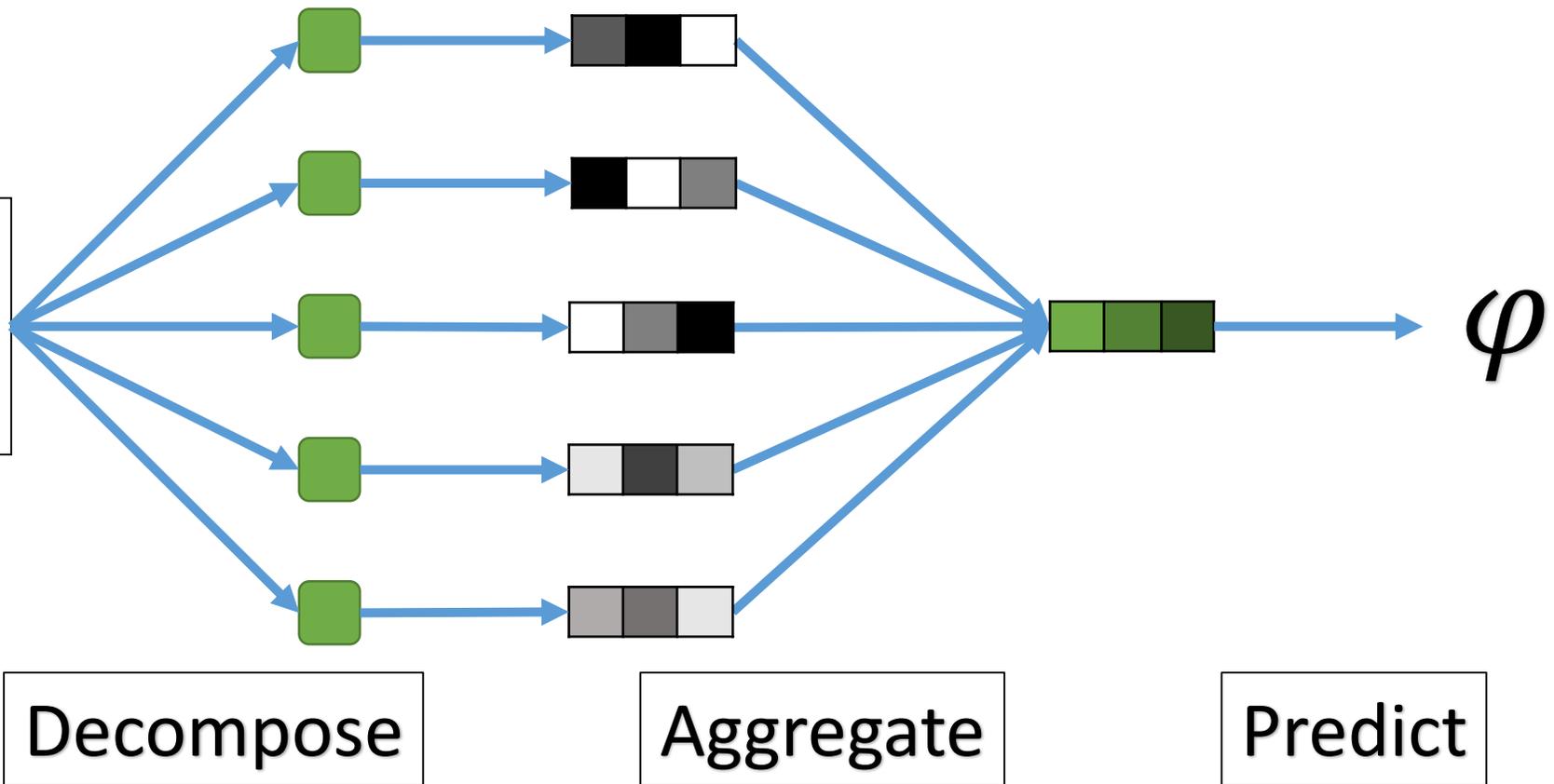
# Back to our problem

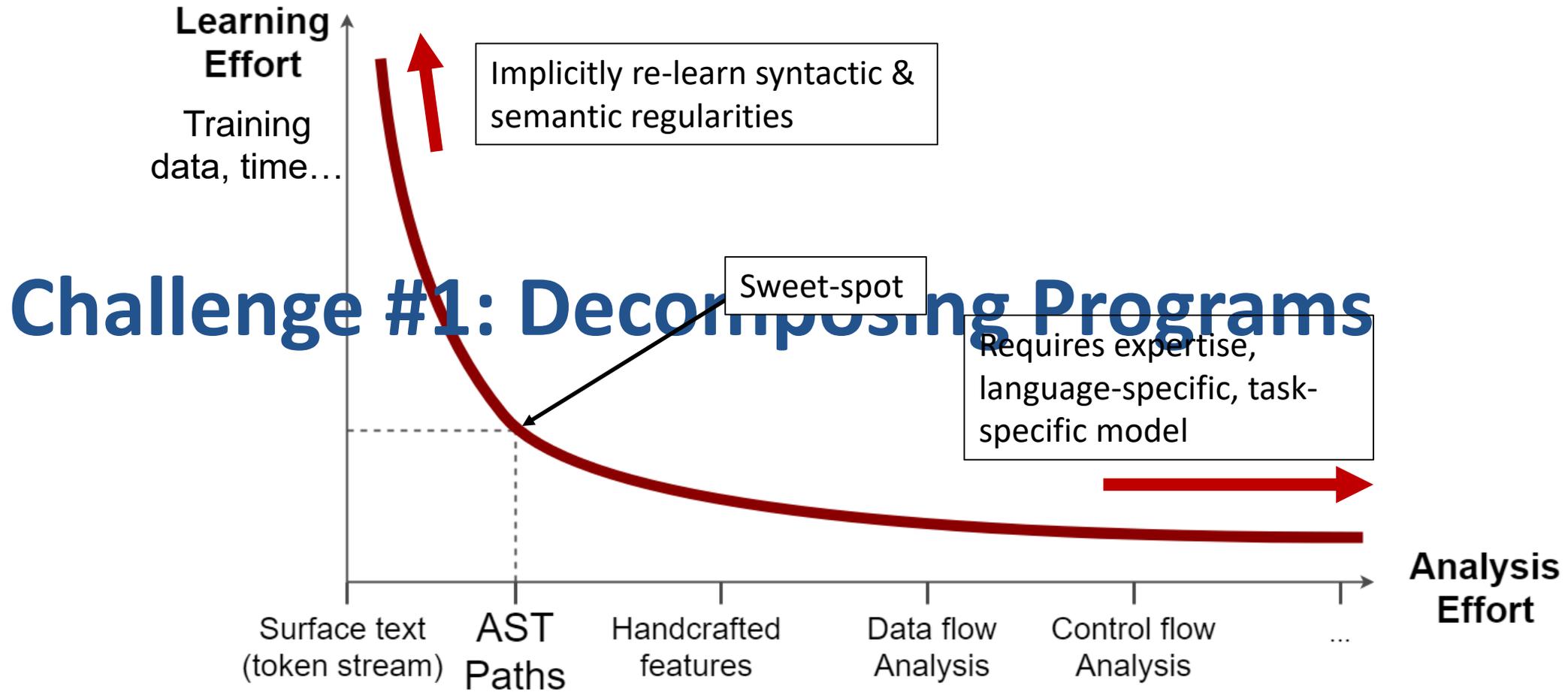
Two main challenges in encoding programs:

1. How to decompose programs to smaller building blocks?
  - Small enough to repeat across programs
  - Large enough to be meaningful the “bias-variance tradeoff”
2. How do we aggregate a set of these building blocks?

# Code2vec: High-level Overview

```
String[] ____ (... ) {  
  final String[] = ...;  
  ...  
  return newArray;  
}
```

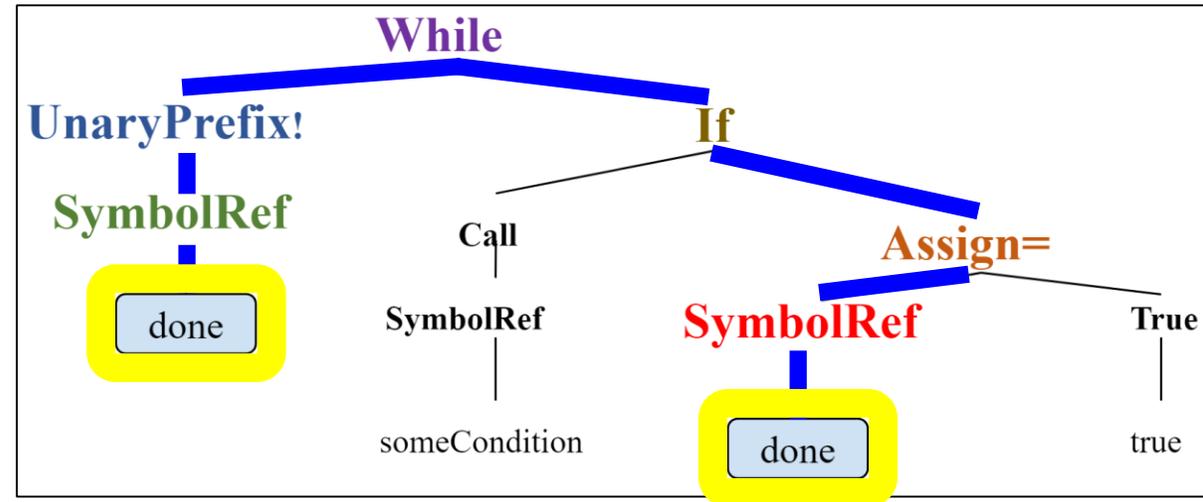
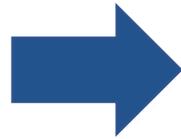




[“A General Path-based Representation for Predicting Program Properties”, PLDI’2018]

# A Program as a Set of AST Paths

```
while (!done) {  
    if (someCondition()) {  
        done = true;  
    }  
}
```



(done, **SymbolRef** ↑

- AST paths capture some of the semantics, by using only the syntax.
- We represent a program as the set of all its paths.

# Representing AST-Paths as Vectors

Two sets of learned vectors:

- Token vectors
- Path vectors

(done, **SymbolRef**↑**UnaryPrefix!**↑**While**↓**If**↓**Assign=**↓**SymbolRef**, done)

lookup  
(token vectors)

lookup  
(path vectors)

lookup  
(token vectors)

$$\underbrace{\tanh([w] \cdot [\text{lookup}(\text{done}), \text{lookup}(\text{SymbolRef}), \text{lookup}(\text{UnaryPrefix!}), \text{lookup}(\text{While}), \text{lookup}(\text{If}), \text{lookup}(\text{Assign=}), \text{lookup}(\text{SymbolRef}), \text{lookup}(\text{done})])}_{\text{fully-connected layer}} = \text{path\_context}$$

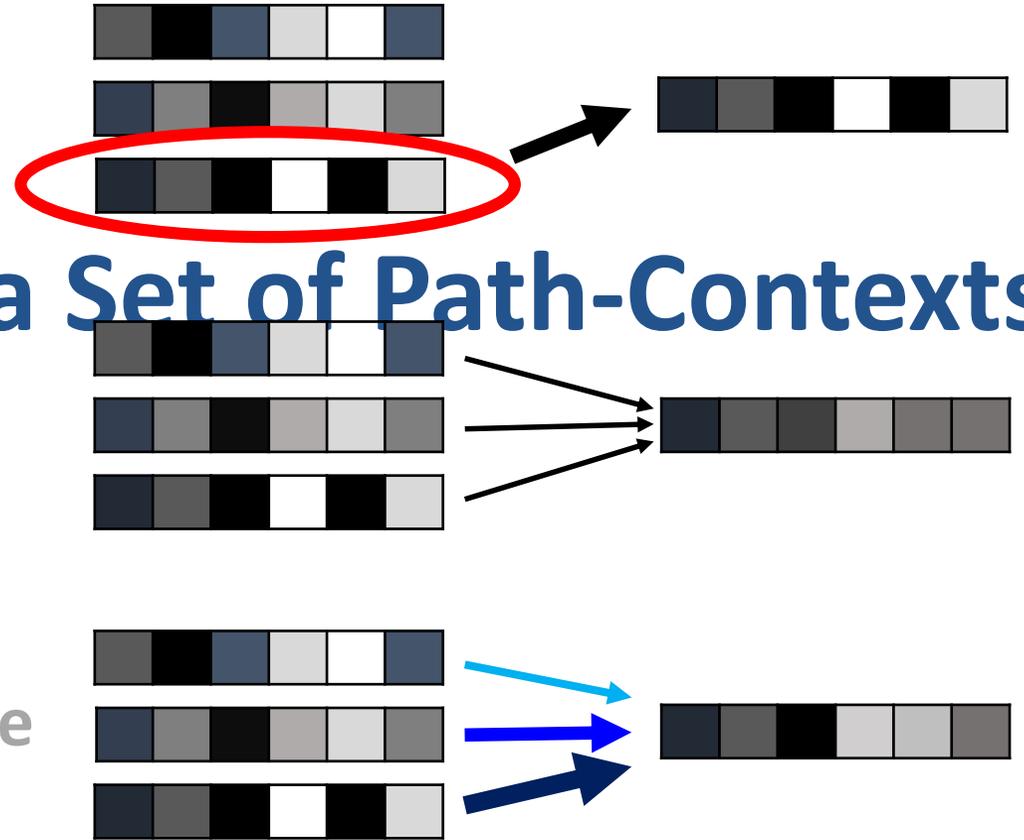
- Input: an arbitrary-sized set of vectors representing AST paths

- Select the “most important vector”

## Challenge #2: Aggregating a Set of Path-Contexts

- Use all vectors, e.g., by averaging them

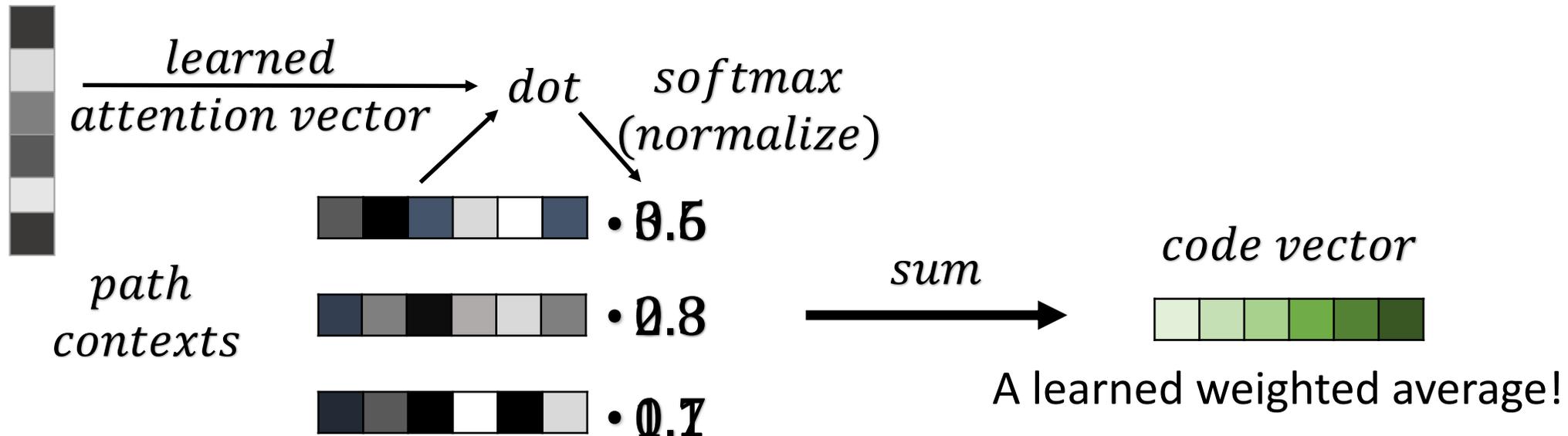
- Attention – a learned weighted average



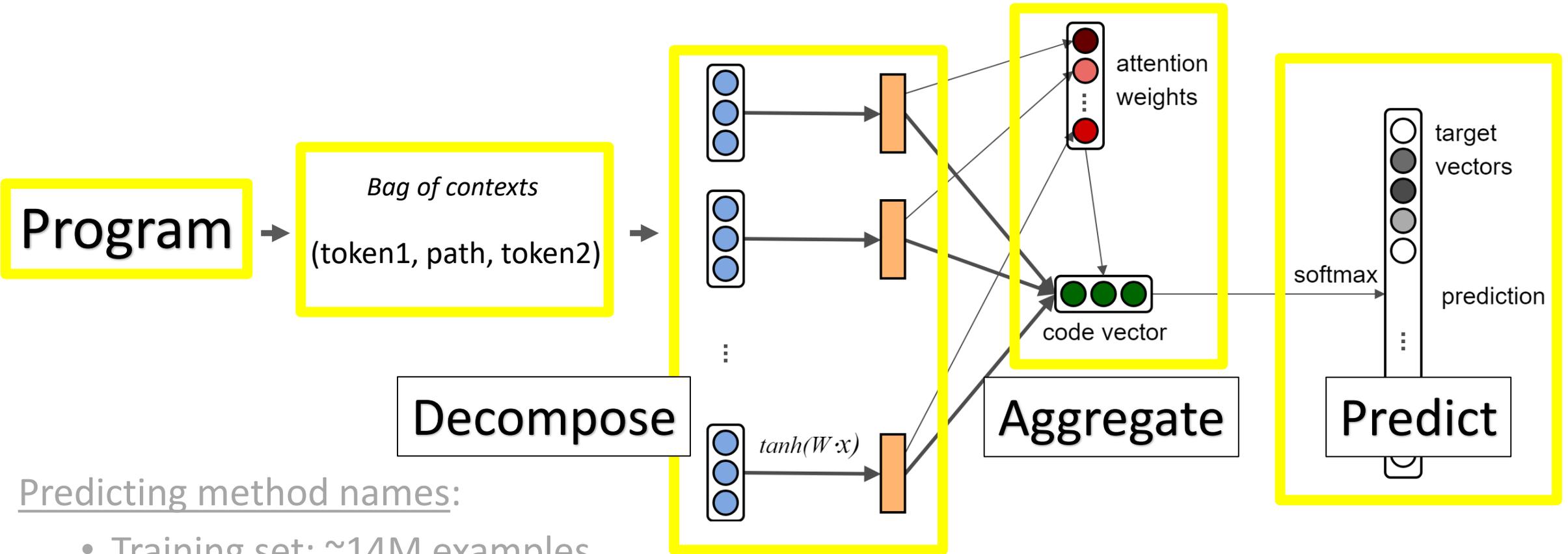
# Attention

Core idea - the values of the vectors learn two distinct goals:

1. The **semantic meaning** of the path-context
2. The amount of **attention** this path-context should get



# Code2vec Architecture



Predicting method names:

- Training set: ~14M examples
- Training time: <1 day (very fast) thanks to its simplicity
- End-to-end: the entire network is trained simultaneously

```

boolean f(Object target) {
    for (Object elem: this.elements) {
        if (elem.equals(target)) {
            return true;
        }
    }
    return false;
}

```

```

Object f(int target) {
    for (Object elem: this.elements) {
        if (elem.hashCode().equals(target)) {
            return elem;
        }
    }
    return this.defaultValue;
}

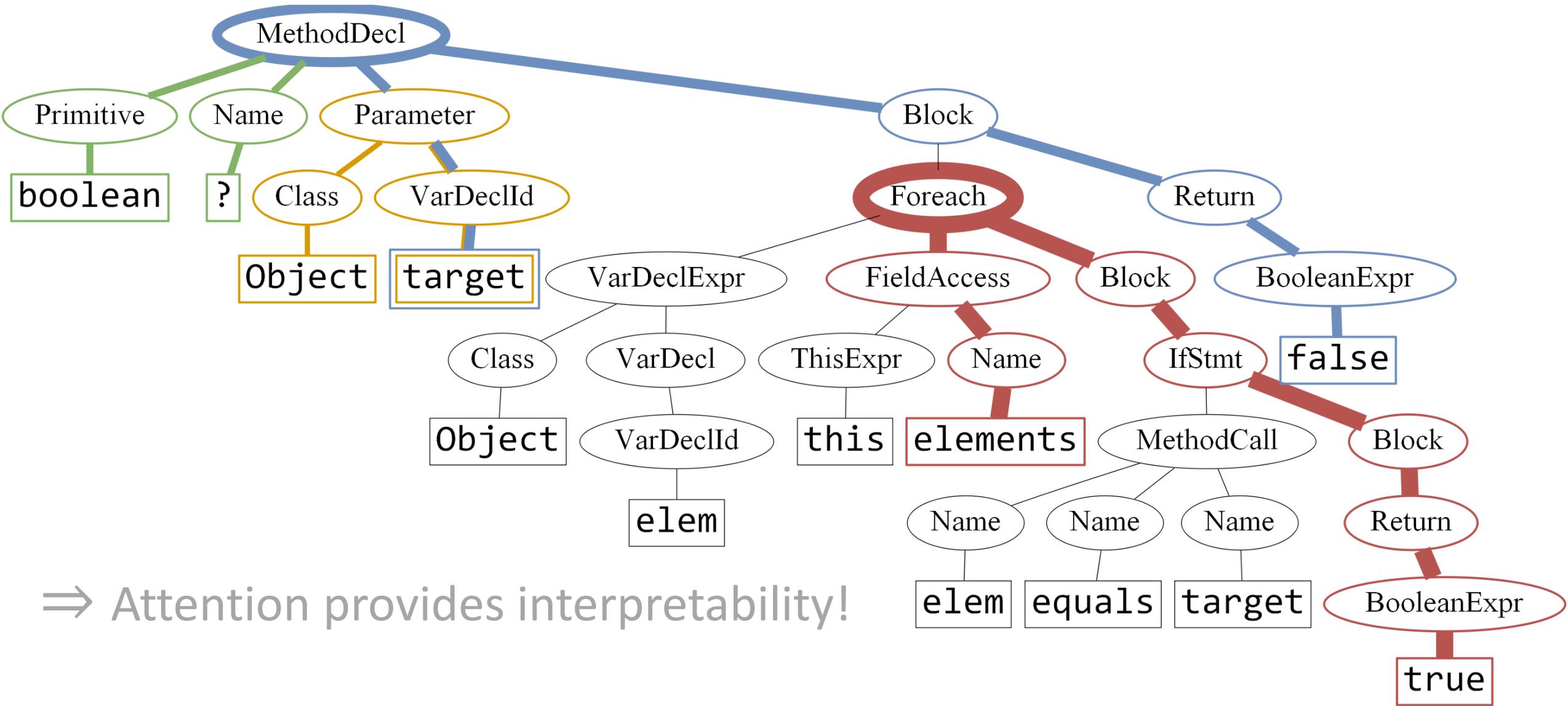
```

Predictions:

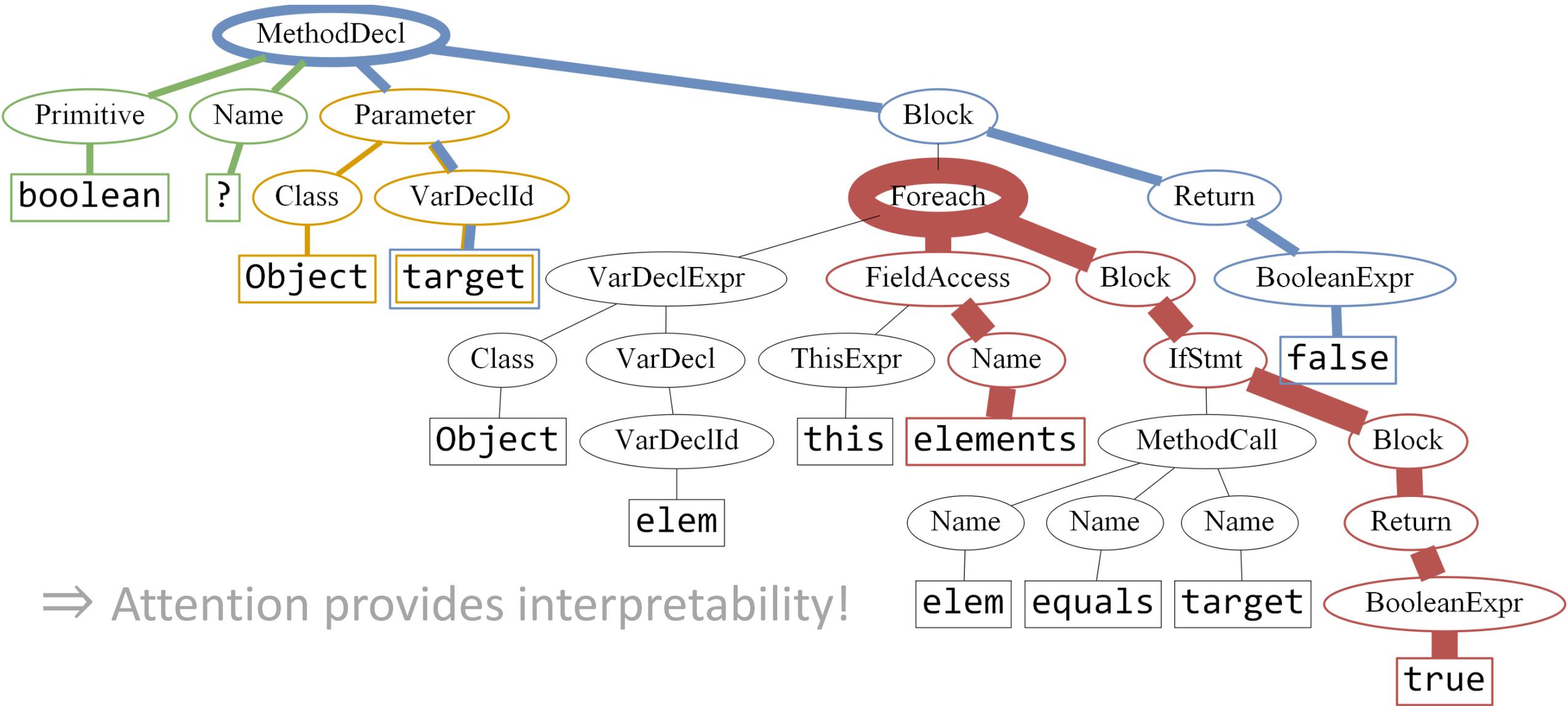
<b>contains</b>		90.93%
matches		3.54%
canHandle		1.15%
equals		0.87%
containsExact		0.77%

Predictions

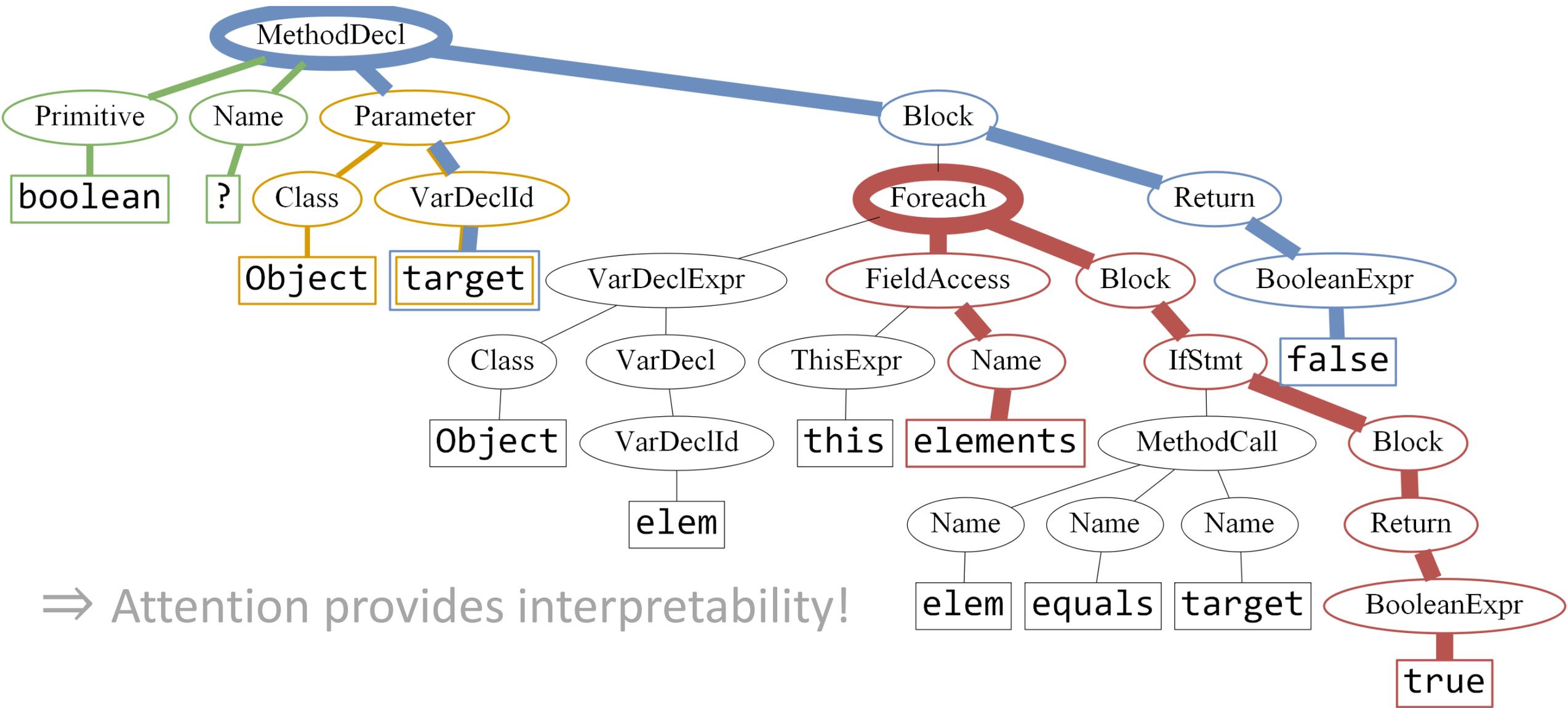
<b>get</b>		31.09%
getProperty		20.25%
getValue		14.34%
getElement		14.00%
getObject		6.05%



⇒ Attention provides interpretability!



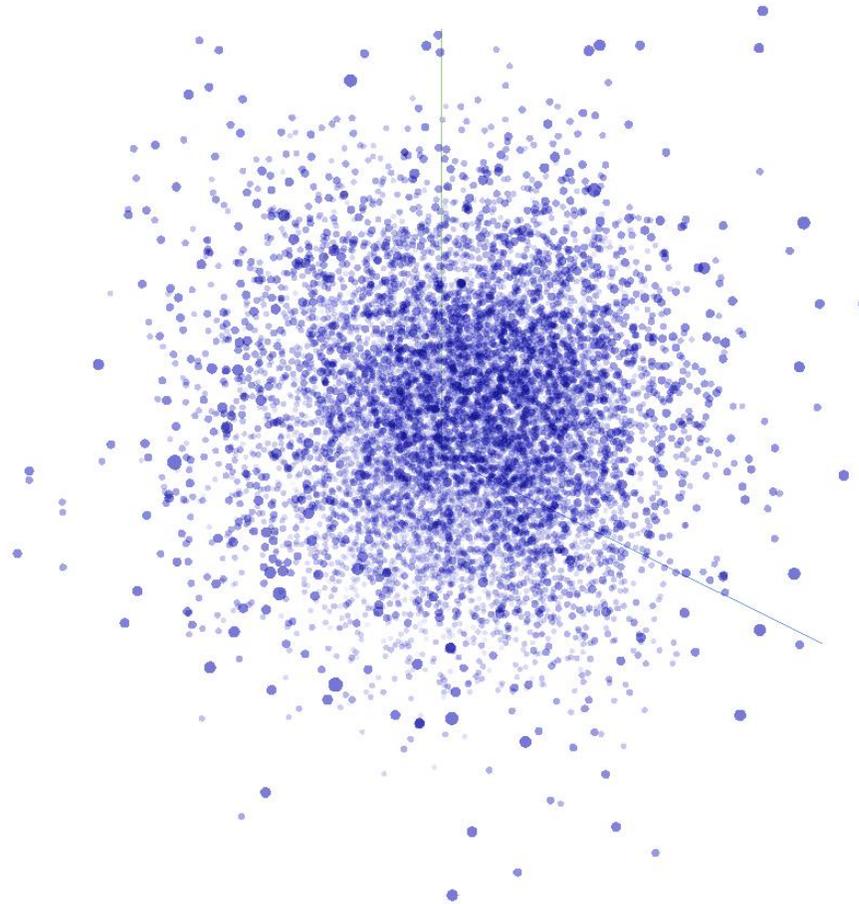
⇒ Attention provides interpretability!



⇒ Attention provides interpretability!

# The Vector Space of Target Labels

Cosine-similar vectors are learned for **semantically similar** labels.





# <http://code2vec.org>

MOST SIMILAR ?

**count**

...is similar to:

PREDICT



getCount		70.02%
size		69.64%
index		64.99%

# <http://code2vec.org>

## COMBINATIONS ?

equals

and

toLowerCase

...combined, are similar to:

PREDICT



equalsIgnoreCase | 78.75%

isUpperCase | 75.82%

equiv | 75.72%

## ANALOGIES ?

receive

is to

download

as...

send

...is to:

PREDICT



upload

| 76.38%

delete

| 71.53%

connect

| 70.51%

# CODE2VEC

Demonstration of principles shown in the paper

## CODE2VEC: LEARNING DISTRIBUTED REPRESENTATIONS OF CODE



Source



Paper



Examples

?

```
boolean f(Object target) {  
    for (Object elem: this.elements) {  
        if (elem.equals(target)) {  
            return true;  
        }  
    }  
    return false;  
}
```

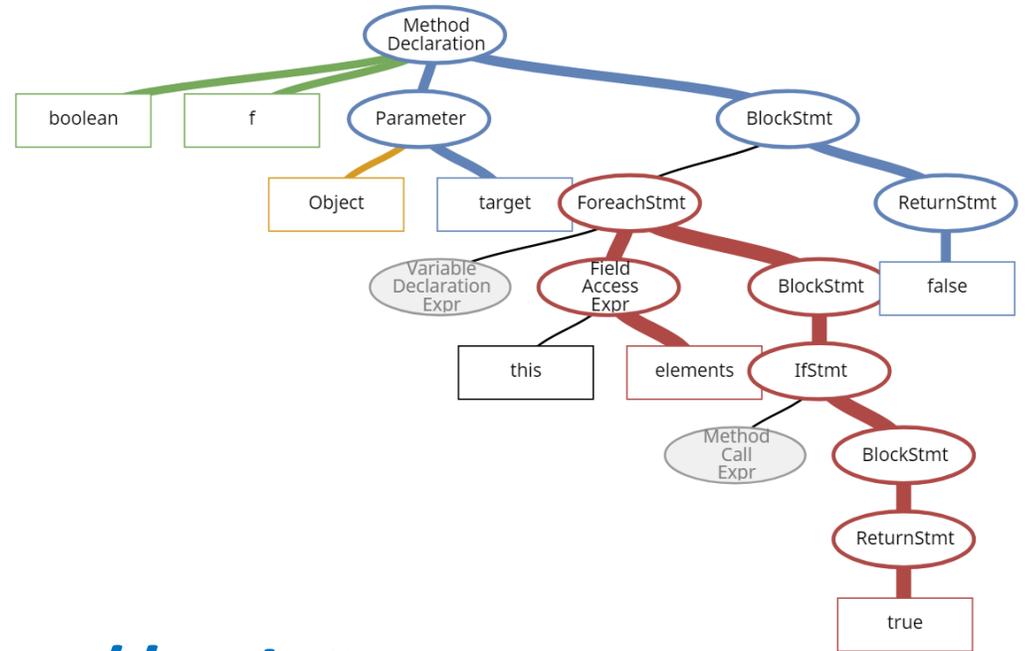
JAVA



contains		90.93%
matches		3.54%
canHandle		1.15%
equals		0.87%
containsExact		0.77%



AST



<http://code2vec.org>

tech-srl / code2vec

<> Code 0 Issues 0 Pull requests 0 Projects 0 Wiki

★ Star 256

TensorFlow code for the neural network presented in the paper: "code2vec: Learning Distributed Representations of Code", POPL'2019  
<https://code2vec.org>

code2vec learning distributed representations of code technion Manage topics

18 commits 1 branch 0 releases 2 contributors MIT

Branch: master New pull request Create new file Upload files Find file Clone or download

Commit	Message	Time
urialon	Update README - adding link to the final POPL PDF	Latest commit dc76680 16 days ago
CSharpExtractor	Adding CSharpExtractor	23 days ago
JavaExtractor	update old version of jackson	23 days ago
images	initial commit	2 months ago
Input.java	initial commit	2 months ago
LICENSE	initial commit	2 months ago
PathContextReader.py	fix reference to config.NUM_EXAMPLES in PathContextReader, which does...	2 months ago
README.md	Update README - adding link to the final POPL PDF	16 days ago
__init__.py	initial commit	2 months ago
build_extractor.sh	initial commit	2 months ago
code2vec.py	initial commit	2 months ago
common.py	initial commit	2 months ago
extractor.py	initial commit	2 months ago
interactive_predict.py	initial commit	2 months ago
model.py	Renaming variables to make it clear where the code vectors are	20 days ago
preprocess.py	initial commit	2 months ago
preprocess.sh	initial commit	2 months ago
preprocess_csharp.sh	Adding CSharpExtractor	23 days ago
train.sh	initial commit	2 months ago

README.md

## Code2vec

A neural network for learning distributed representations of code. This is an official implementation of the model described in:

[Uri Alon, Meital Zilberstein, Omer Levy and Eran Yahav, "code2vec: Learning Distributed Representations of Code", POPL'2019 \[PDF\]](#)

# Summary

- Core ideas in learning code snippets:
  1. Representing a code snippet as a set of syntactic paths
  2. Aggregate all paths using neural attention
- A simple and fast to train architecture
- Interpretable thanks to the attention mechanism
- The learned vectors capture interesting phenomena

**<http://code2vec.org>**

# How many paths do you take from each code snippet? Taking all paths is quadratic!

- An unlimited number!
  - Since attention is simply a weighted average, it can handle an arbitrary number of path-contexts.
  - Empirically, we found that sampling 200 from each code example is sufficient. ~200 is also the average number of paths per example.
  - This number (200) can be easily increased if the dataset contained especially large pieces of code.
  - Paths that are missed due to sampling are “covered” by other paths.

# Why not performing additional control flow or data flow analyses?

- These might help, but we are not sure they are necessary here. Most of the important signals are expressed in the syntax.
- Our pure-syntactic approach has the advantage of generality – the same approach can be easily applied to other languages.
- Semantic analysis is probably necessary in other tasks (for example, when the programs are binaries).

# How robust are the results for variable renaming?

- As any machine learning model, confusing or adversarial examples can mislead our model.
- Since the network was trained on “well-named” examples from top-starred GitHub projects, it does perform worse without names.
- We are exploring similar approaches for obfuscated code as part of ongoing research.

# Do you keep vectors for *all* paths and tokens?

- Almost all!
  - Limiting to the most occurring 1M tokens, 1M paths, and 300k target labels.
- Each token and path vectors has 128 elements of 4 bytes (float32)
- Each target vector has 384 elements of 4 bytes
- Attention vector has 384 elements
- Fully connected layer is a matrix of size  $384 \times 384$

- Total size:  $\underbrace{128 \cdot 4}_{\text{vector}} \cdot \left( \underbrace{1M + 1M}_{\substack{\text{token+path} \\ \text{vocab} \\ \text{sizes}}} \right) + \underbrace{384 \cdot 4}_{\text{vector}} \cdot \underbrace{300k}_{\substack{\text{target} \\ \text{vocab size}}} + \underbrace{384}_{\text{attention}} + \underbrace{384^2}_{\substack{\text{fully-} \\ \text{connected}}} \approx \mathbf{1.5 GB}$

- Standard GPU memory size: **12 GB**

# Did you try Gated Graph NNs (Allamanis et al., ICLR'2018)?

- GGNNs were applied to a simpler task of Var-Misuse.
  - Their code is not fully available.
- Two conceptual advantages of **code2vec** over **GGNNs**:
  1. **Much faster to train** - thus practically easier to leverage huge corpora (our dataset is orders of magnitude larger).
  2. **Our model is purely syntactic** - the same algorithm can work for every programming language. In GGNNs, the edges in the graph are analyses like “ComputedFrom” and “LastWrite”, that need to be re-implemented for different languages.

# Can a non-neural model solve the same task?

- Yes, and pretty well (PLDI'2018).
  - But not as good as a neural model.
- Main advantages of using a neural network:
  1. Much **better generalization** (Section 5 in the paper)
  2. Our neural network can **produce a vector**, which can be fed to a variety of other (neural and non-neural) ML models and tasks.