

Flow2Vec: Value-Flow-Based Precise Code Embedding

Yulei Sui¹, Xiao Cheng², Guanqin Zhang¹, Haoyu Wang²

¹ University of Technology Sydney, Australia

² Beijing University of Posts and Telecommunications, China

November 17 2020

Contribution

A new code embedding approach which marries **precise static analysis** and recent advances in **high-order proximity embedding**, by preserving

- **interprocedural alias-aware program dependence**, and
- **context-free language reachability**,

to better support subsequent code analysis tasks, such as code summarization and semantic labeling.

Code Embedding

Source Code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++) →
        myArray[i] = i;
    return myArray;
}
```

Model



Code Property Prediction



Code Embedding


Source Code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++) →
        myArray[i] = i;
    return myArray;
}
```

Model

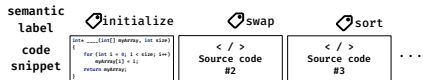


Semantic Label

 initialize

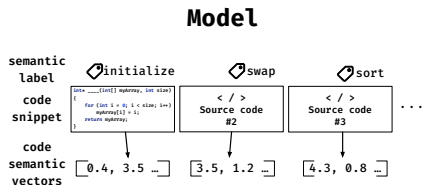
Code Embedding

Model



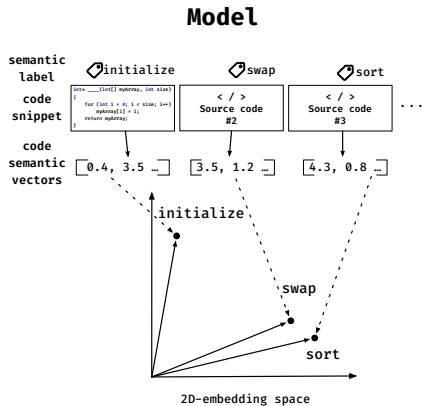
Code Embedding

Code semantic vector in geometric space



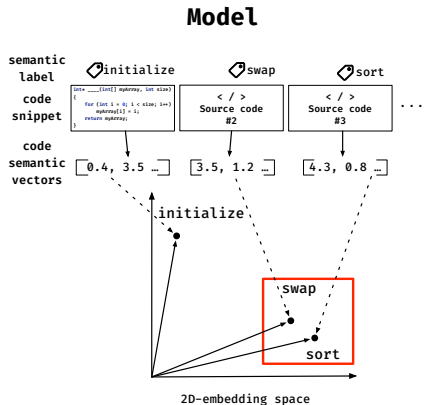
Code Embedding

Code semantic vector in geometric space



Code Embedding

Code semantic vector in geometric space



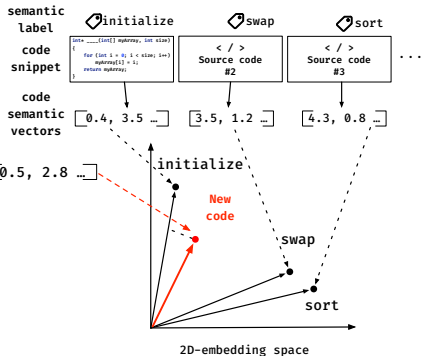
Code Embedding

Code semantic vector in geometric space

New code

```
char* ___(char[] cArray, int size)
{
    int i = 0;
    while ( i < size )
        cArray[i] = i;
        i++;
    return cArray;
}
```

Model



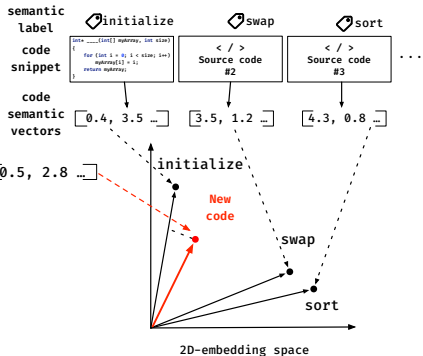
Code Embedding

Code semantic vector in geometric space

New code

```
char* ___(char[] cArray, int size)
{
    int i = 0;
    while ( i < size )
        cArray[i] = i;
        i++;
    return cArray;
}
```

Model



Semantic label

initialize

Existing Embedding Approaches

Structure-oblivious embedding

Source code → A bag of 'sentences'^[1,2]

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

Existing Embedding Approaches

Structure-oblivious embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'sentences'^[1,2]

➔ int*____(int[]myArray,int size)

[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

Existing Embedding Approaches

Structure-oblivious embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'sentences'^[1,2]

→ int * ____ (int [] myArray , int size)

[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

Existing Embedding Approaches

Structure-oblivious embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'sentences'^[1,2]

```
int * ____ ( int [] myArray , int size )
for ( int i = 0 ...
:
:
[ ] [ ] [ ] ...
```

[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

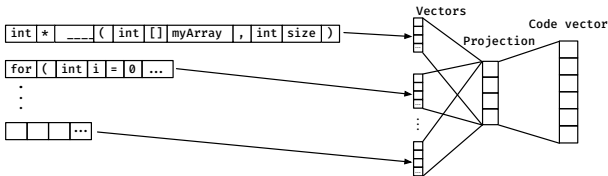
Existing Embedding Approaches

Structure-oblivious embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'sentences'^[1,2]



[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

Existing Embedding Approaches

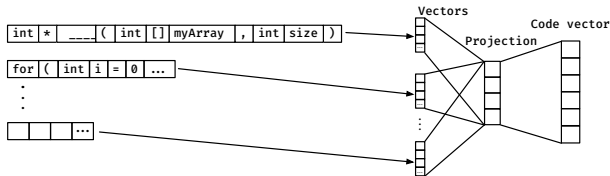
Structure-oblivious embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}

char* ___ (char[] cArray, int size)
{
    int i = 0;
    while ( i < size )
        cArray[i] = i;
        i++;
    return cArray;
}
```

A bag of 'sentences'^[1,2]



Embedding

Textually different but semantically equivalent

[1] Distributed representations of words and phrases and their compositionality. In NeurIPS '13

[2] Distributed representations of sentences and documents. In ICML '14

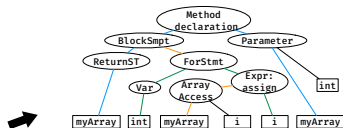
Existing Embedding Approaches

Structure-preserving embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

Abstract Syntax Tree ^[3]



[3] code2vec: Learning distributed representations of code. POPL .2019

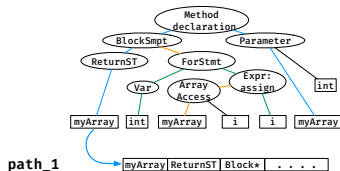
Existing Embedding Approaches

Structure-preserving embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'paths' on AST ^[3]



[3] code2vec: Learning distributed representations of code. POPL .2019

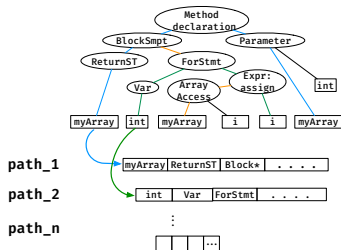
Existing Embedding Approaches

Structure-preserving embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'paths' on AST ^[3]



[3] code2vec: Learning distributed representations of code. POPL .2019

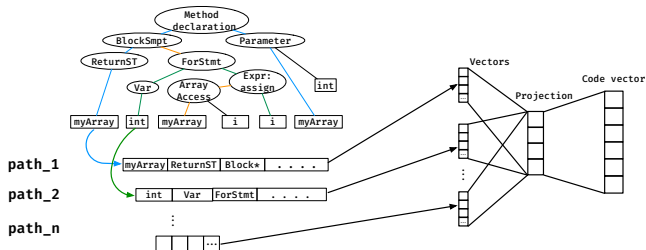
Existing Embedding Approaches

Structure-preserving embedding

Source code

```
int* ____ (int[] myArray, int size)
{
    for (int i = 0; i < size; i++)
        myArray[i] = i;
    return myArray;
}
```

A bag of 'paths' on AST^[3]



Embedding

[3] code2vec: Learning distributed representations of code. POPL .2019

Problems and Limitations

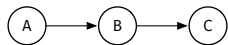
(a) Fail to capture asymmetric transitivity

(b) Alias-unaware

(c) Intraprocedural/context-insensitivity

Problems and Limitations

(a) Fail to capture asymmetric transitivity



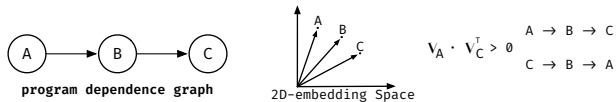
program dependence graph

(b) Alias-unaware

(c) Intraprocedural/context-insensitivity

Problems and Limitations

(a) Fail to capture asymmetric transitivity

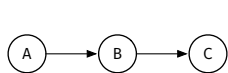


(b) Alias-unaware

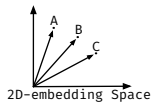
(c) Intraprocedural/context-insensitivity

Problems and Limitations

(a) Fail to capture asymmetric transitivity



program dependence graph



$$v_A \cdot v_C > 0$$

A → B → C ✓ Real reachability and correctly preserved

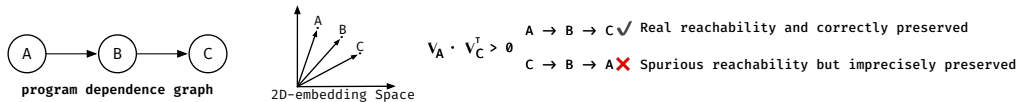
C → B → A ✗ Spurious reachability but imprecisely preserved

(b) Alias-unaware

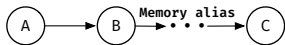
(c) Intraprocedural/context-insensitivity

Problems and Limitations

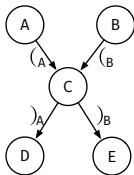
(a) Fail to capture asymmetric transitivity



(b) Alias-unaware

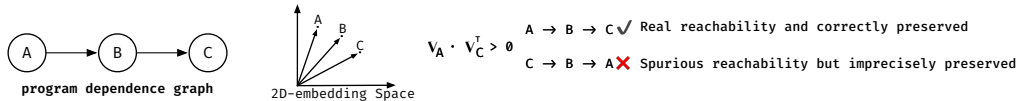


(c) Intraprocedural/context-insensitivity

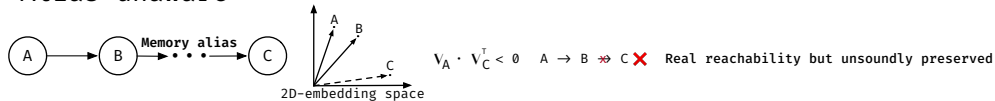


Problems and Limitations

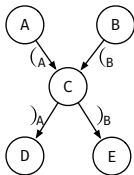
(a) Fail to capture asymmetric transitivity



(b) Alias-unaware

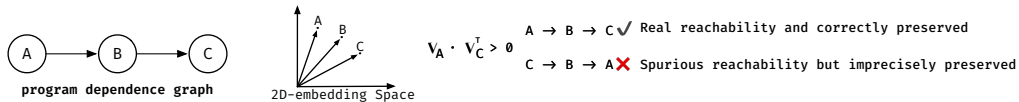


(c) Intraprocedural/context-insensitivity

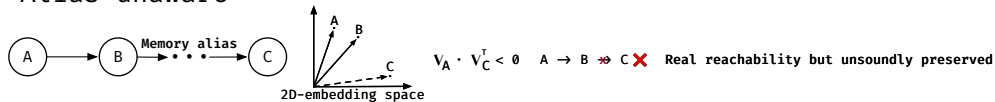


Problems and Limitations

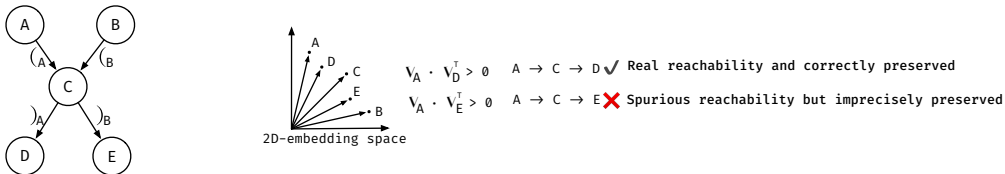
(a) Fail to capture asymmetric transitivity



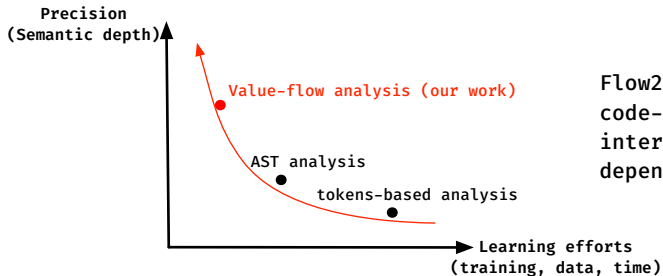
(b) Alias-unaware



(c) Intraprocedural/context-insensitivity

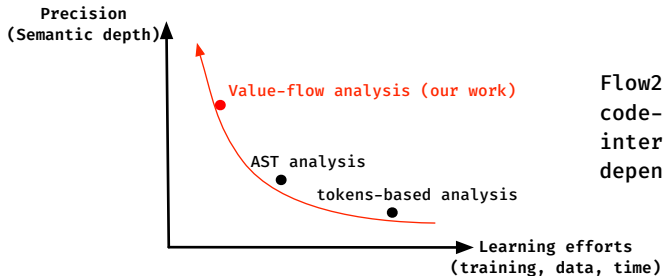


The Aim of This Work

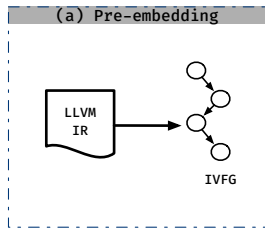


Flow2Vec: a high-order proximity code-embedding approach by preserving interprocedural alias-aware program dependence

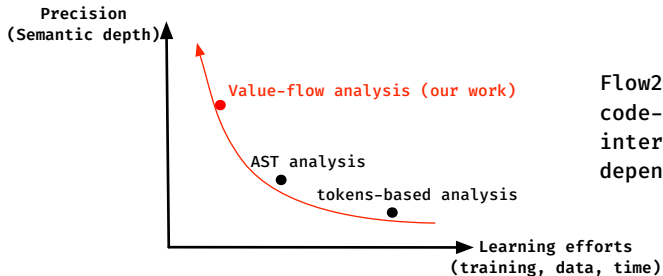
The Aim of This Work



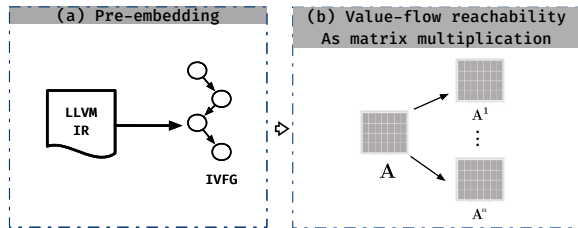
Flow2Vec: a high-order proximity code-embedding approach by preserving interprocedural alias-aware program dependence



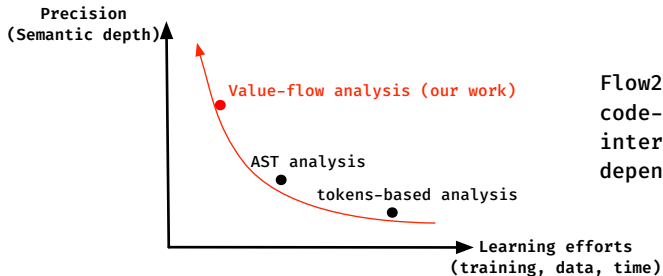
The Aim of This Work



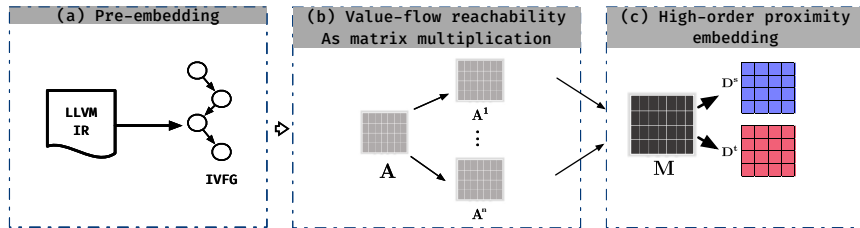
Flow2Vec: a high-order proximity code-embedding approach by preserving interprocedural alias-aware program dependence



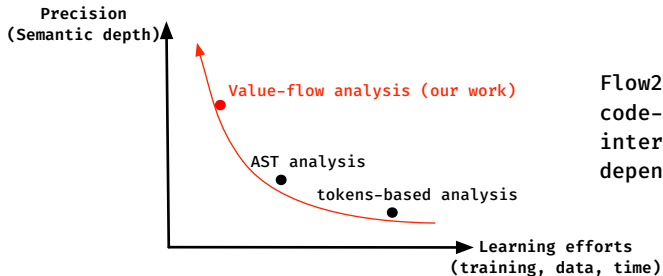
The Aim of This Work



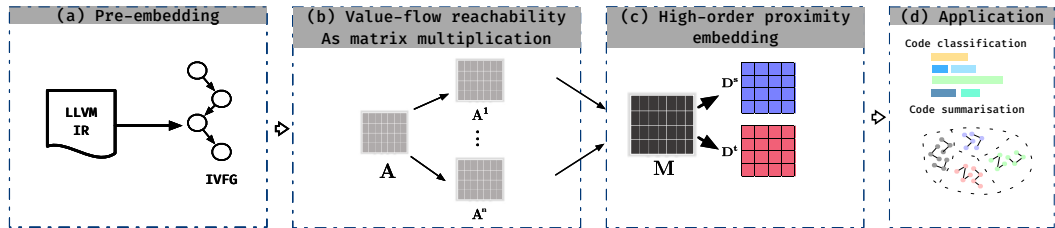
Flow2Vec: a high-order proximity code-embedding approach by preserving interprocedural alias-aware program dependence



The Aim of This Work



Flow2Vec: a high-order proximity code-embedding approach by preserving interprocedural alias-aware program dependence



A Motivating Example

Phase (a) Pre-embedding

```
foo(){
l1:  stack = malloc(...);
l2:  queue = malloc(...);
l3:  p = initialize(stack); // cs1
l4:  q = initialize(queue); // cs2

.....
}
l5: initialize(x){
    // initialization for
    // objects that 'x' points to
    .....
l6:  return x;
}
```

A Motivating Example

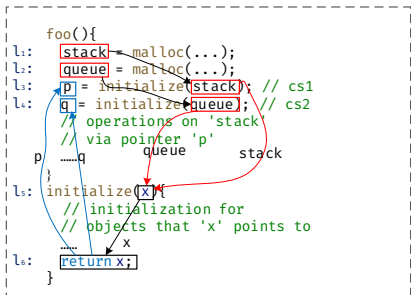
Phase (a) Pre-embedding

```
foo(){
l1:  stack = malloc(...);
l2:  queue = malloc(...);
l3:  p = initialize(stack); // cs1
l4:  q = initialize(queue); // cs2
      // operations on 'stack'
      // via pointer 'p'
      .....
}
l5: initialize(x){
      // initialization for
      // objects that 'x' points to
      .....
l6:  return x;
}
```

Call site 1 p refers to stack
Call site 2 q refers to queue

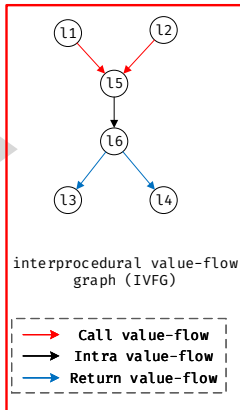
A Motivating Example

Phase (a) Pre-embedding



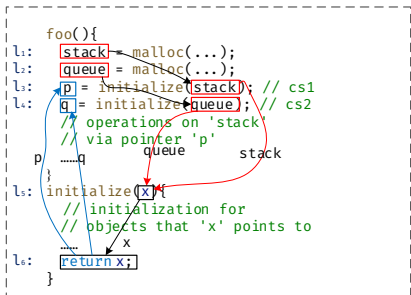
Call site 1 `p` refers to `stack`

Call site 2 `q` refers to `queue`



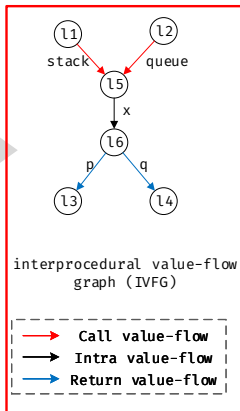
A Motivating Example

Phase (a) Pre-embedding



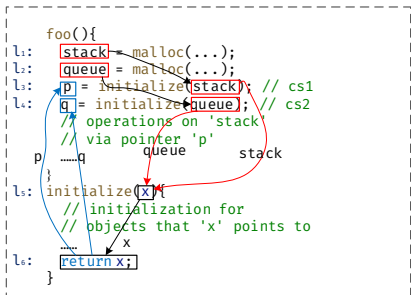
Call site 1 `p refer to stack`

Call site 2 `q refer to queue`



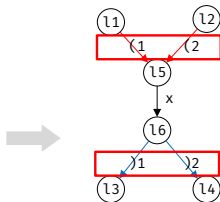
A Motivating Example

Phase (a) Pre-embedding

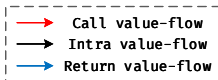


Call site 1 `p` refers to `stack`

Call site 2 `q` refers to `queue`

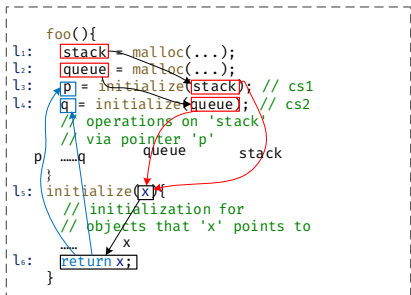


interprocedural value-flow graph (IVFG)



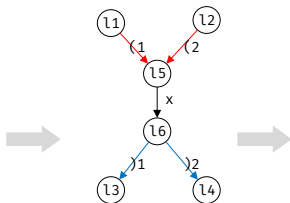
A Motivating Example

Phase (a) Pre-embedding

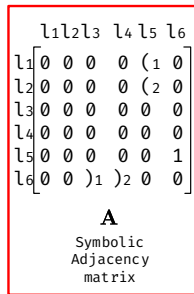
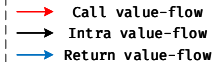


Call site 1 p refers to stack

Call site 2 q refers to queue

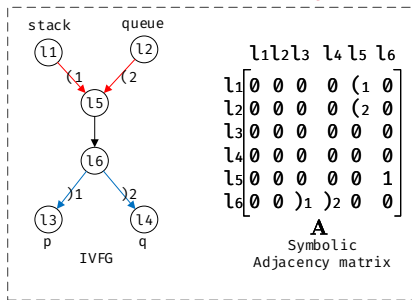


interprocedural value-flow graph (IVFG)



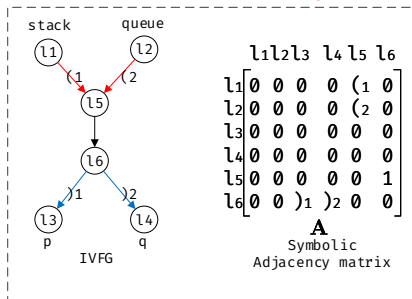
A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication



A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

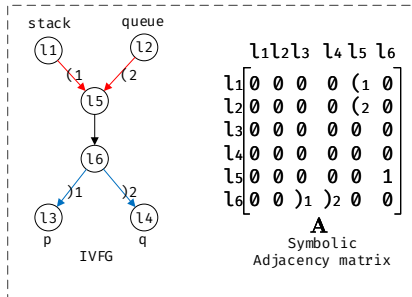


$$\mathbf{A}^h = \underbrace{\mathbf{A} \cdot \dots \cdot \mathbf{A}}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication



$$A^2 = A \cdot A$$

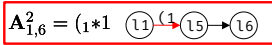
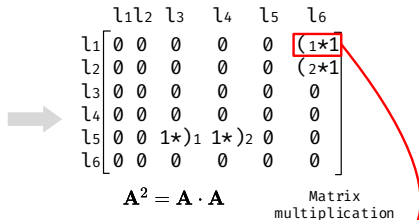
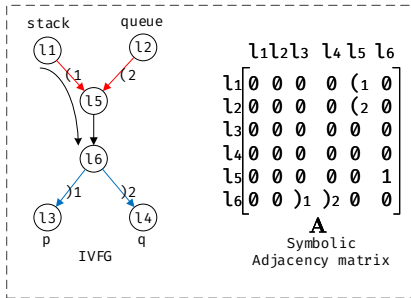
Matrix multiplication

$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

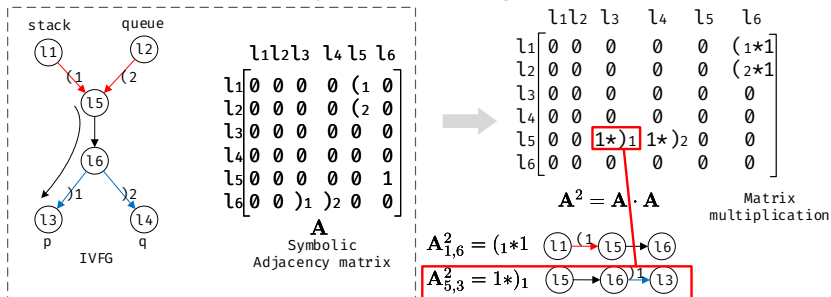


$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

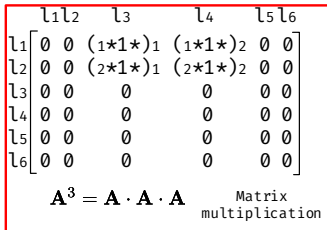
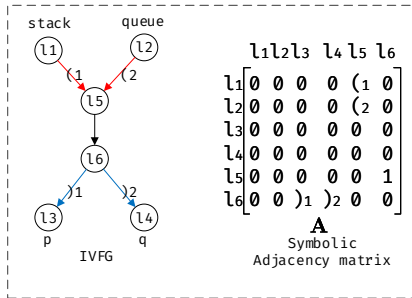


$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

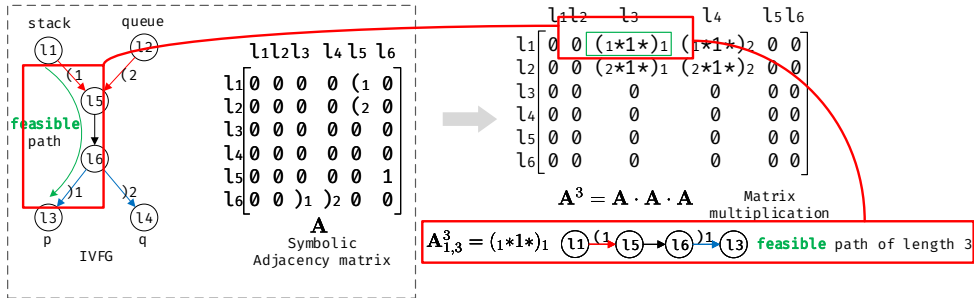


$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

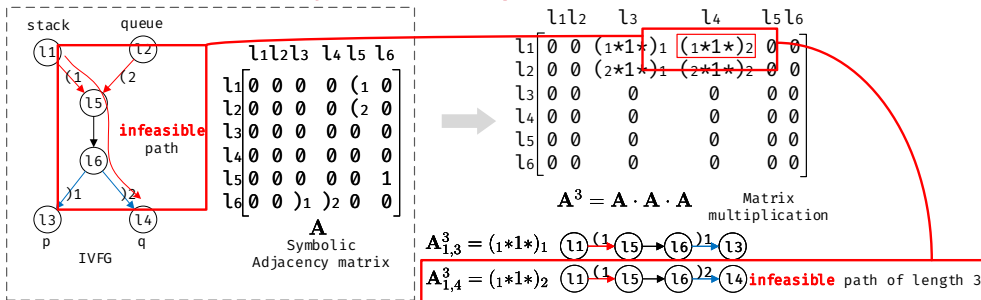


$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication

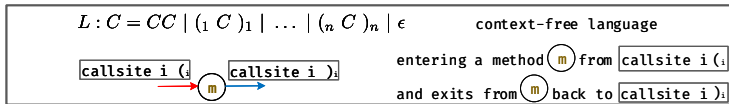
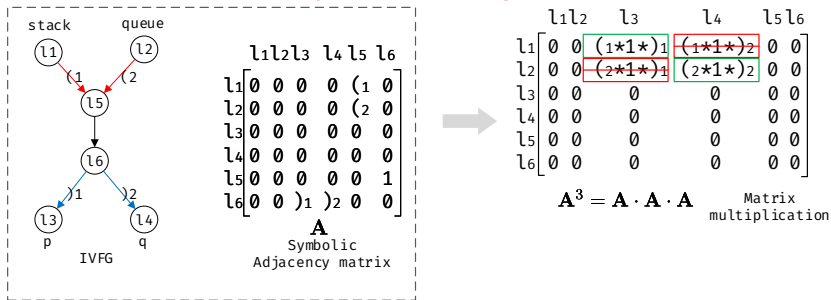


$$A^h = \underbrace{A \cdot \dots \cdot A}_{\times h}$$

The power of matrix
h-th order reachability

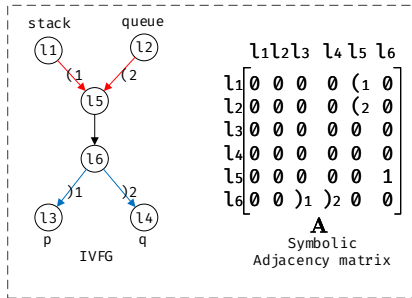
A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication



A Motivating Example

Phase (b) Value-flow reachability as matrix multiplication



$$A^3 = A \cdot A \cdot A$$

Matrix multiplication

	l1	l2	l3	l4	l5	l6
l1	0	0	1	0	0	0
l2	0	0	0	1	0	0
l3	0	0	0	0	0	0
l4	0	0	0	0	0	0
l5	0	0	0	0	0	0
l6	0	0	0	0	0	0

A Motivating Example

Phase (c) High-order proximity embedding

$$\begin{array}{c} \begin{array}{cccccc} l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \end{array} \\ \begin{array}{l} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{array} \end{array} \begin{bmatrix} 0 & 0 & 0 & 0 & (1 & 0) \\ 0 & 0 & 0 & 0 & (2 & 0) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 &)_1 &)_2 & 0 & 0 \end{bmatrix}$$

\mathbf{A}

$$\begin{array}{c} \begin{array}{cccccc} l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \end{array} \\ \begin{array}{l} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{array} \end{array} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & (1*1) \\ 0 & 0 & 0 & 0 & 0 & (2*1) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1*)_1 & 1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$\mathbf{A}^2 = \mathbf{A} \cdot \mathbf{A}$

$$\begin{array}{c} \begin{array}{cccccc} l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \end{array} \\ \begin{array}{l} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{array} \end{array} \begin{bmatrix} 0 & 0 & (1*1*)_1 & (1*1*)_2 & 0 & 0 \\ 0 & 0 & (2*1*)_1 & (2*1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$\mathbf{A}^3 = \mathbf{A} \cdot \mathbf{A} \cdot \mathbf{A}$

Katz Index

$$\mathbf{M} = \sum_{h=1}^H (\beta \cdot \mathbf{A})^h = 0.8 \cdot \mathbf{A} + 0.8^2 \cdot \mathbf{A}^2 + 0.8^3 \cdot \mathbf{A}^3$$

A Motivating Example

Phase (c) High-order proximity embedding

$$\begin{array}{c}
 \begin{array}{cccccc}
 l_1 & l_2 & l_3 & l_4 & l_5 & l_6
 \end{array} \\
 \begin{array}{c}
 l_1 \\
 l_2 \\
 l_3 \\
 l_4 \\
 l_5 \\
 l_6
 \end{array}
 \begin{bmatrix}
 0 & 0 & 0 & 0 & (1 & 0 \\
 0 & 0 & 0 & 0 & (2 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 1 \\
 0 & 0 &)_1 &)_2 & 0 & 0
 \end{bmatrix}
 \end{array}$$

A

$$\begin{array}{c}
 \begin{array}{cccccc}
 l_1 & l_2 & l_3 & l_4 & l_5 & l_6
 \end{array} \\
 \begin{array}{c}
 l_1 \\
 l_2 \\
 l_3 \\
 l_4 \\
 l_5 \\
 l_6
 \end{array}
 \begin{bmatrix}
 0 & 0 & 0 & 0 & 0 & (1*1 \\
 0 & 0 & 0 & 0 & 0 & (2*1 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 1*)_1 & 1*)_2 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}
 \end{array}$$

$A^2 = A \cdot A$

$$\begin{array}{c}
 \begin{array}{cccccc}
 l_1 & l_2 & l_3 & l_4 & l_5 & l_6
 \end{array} \\
 \begin{array}{c}
 l_1 \\
 l_2 \\
 l_3 \\
 l_4 \\
 l_5 \\
 l_6
 \end{array}
 \begin{bmatrix}
 0 & 0 & (1*1*)_1 & (1*1*)_2 & 0 & 0 \\
 0 & 0 & (2*1*)_1 & (2*1*)_2 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0
 \end{bmatrix}
 \end{array}$$

$A^3 = A \cdot A \cdot A$

Katz Index

$$M = \sum_{h=1}^H (\beta \cdot A)^h = 0.8 \cdot A + 0.8^2 \cdot A^2 + 0.8^3 \cdot A^3$$

Context-sensitive
value-flow
reachability
matrix

$$\begin{array}{c}
 \begin{array}{cccccc}
 l_1 & l_2 & l_3 & l_4 & l_5 & l_6
 \end{array} \\
 \begin{array}{c}
 l_1 \\
 l_2 \\
 l_3 \\
 l_4 \\
 l_5 \\
 l_6
 \end{array}
 \begin{bmatrix}
 0 & 0 & 0.512*(1*1*)_1 & 0.512*(1*1*)_2 & 0.8*(1 & 0.64*(1*1 \\
 0 & 0 & 0.512*(2*1*)_1 & 0.512*(2*1*)_2 & 0.8*(2 & 0.64*(2*1 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0.64*1*)_1 & 0.64*1*)_2 & 0 & 0.8 \\
 0 & 0 & 0.8*)_1 & 0.8*)_2 & 0 & 0
 \end{bmatrix}
 \end{array}$$

A Motivating Example

Phase (c) High-order proximity embedding

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & (1 & 0) \\ 0 & 0 & 0 & 0 & (2 & 0) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 &)_1 &)_2 & 0 & 0 \end{bmatrix} \end{matrix}$$

A

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & (1*1) \\ 0 & 0 & 0 & 0 & 0 & (2*1) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1*)_1 & 1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

A² = A · A

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & (1*1*)_1 & (1*1*)_2 & 0 & 0 \\ 0 & 0 & (2*1*)_1 & (2*1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

A³ = A · A · A

Katz Index

$$\mathbf{M} = \sum_{h=1}^H (\beta \cdot \mathbf{A})^h = 0.8 \cdot \mathbf{A} + 0.8^2 \cdot \mathbf{A}^2 + 0.8^3 \cdot \mathbf{A}^3$$

Context-sensitive
value-flow
reachability
matrix

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0.512 & 0 & 0.8 & 0.64 \\ 0 & 0 & 0 & 0.512 & 0.8 & 0.64 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.64 & 0.64 & 0 & 0.8 \\ 0 & 0 & 0.8 & 0.8 & 0 & 0 \end{bmatrix} \end{matrix}$$

A Motivating Example

Phase (c) High-order proximity embedding

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & (1 & 0) \\ 0 & 0 & 0 & 0 & (2 & 0) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 &)_1 &)_2 & 0 & 0 \end{bmatrix}
 \end{matrix}$$

\mathbf{A}

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & (1*1) \\ 0 & 0 & 0 & 0 & 0 & (2*1) \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1*)_1 & 1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}
 \end{matrix}$$

$\mathbf{A}^2 = \mathbf{A} \cdot \mathbf{A}$

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & (1*1*)_1 & (1*1*)_2 & 0 & 0 \\ 0 & 0 & (2*1*)_1 & (2*1*)_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}
 \end{matrix}$$

$\mathbf{A}^3 = \mathbf{A} \cdot \mathbf{A} \cdot \mathbf{A}$

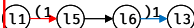
Katz Index

$$\mathbf{M} = \sum_{h=1}^H (\beta \cdot \mathbf{A})^h = 0.8 \cdot \mathbf{A} + 0.8^2 \cdot \mathbf{A}^2 + 0.8^3 \cdot \mathbf{A}^3$$

Context-sensitive
value-flow
reachability
matrix

$$\begin{matrix} & l_1 & l_2 & l_3 & l_4 & l_5 & l_6 \\ \begin{matrix} l_1 \\ l_2 \\ l_3 \\ l_4 \\ l_5 \\ l_6 \end{matrix} & \begin{bmatrix} 0 & 0 & 0.512 & 0 & 0.8 & 0.64 \\ 0 & 0 & 0 & 0.512 & 0.8 & 0.64 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0.64 & 0.64 & 0 & 0.8 \\ 0 & 0 & 0.8 & 0.8 & 0 & 0 \end{bmatrix}
 \end{matrix}$$

Example: l_1 to l_3

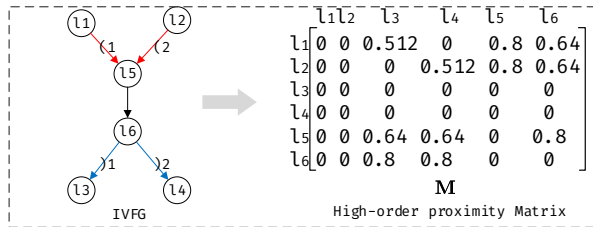


Reachability from l_1 to l_3 :

$$\begin{aligned}
 \mathbf{M}_{1,3} &= 0.8 \cdot \mathbf{A}_{1,3} + 0.8^2 \cdot \mathbf{A}_{1,3}^2 + 0.8^3 \cdot \mathbf{A}_{1,3}^3 \\
 &= 0.8 * 0 + 0.8^2 * 0 + 0.8^3 * \boxed{1} = 0.512
 \end{aligned}$$

A Motivating Example

Phase (c) High-order proximity embedding



A Motivating Example

Phase (c) High-order proximity embedding



Singular Value Decompose

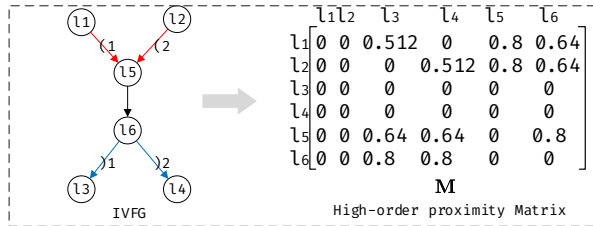
$$M \approx U \Sigma V^T = \sum_{i=1}^N \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \mathbf{D}^s \cdot \mathbf{D}^t$$

$$\mathbf{D}^s = [\sqrt{\sigma_1} \cdot \mathbf{u}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{u}_K] = [\mathbf{d}_1^s, \dots, \mathbf{d}_N^s]$$

$$\mathbf{D}^t = [\sqrt{\sigma_1} \cdot \mathbf{v}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{v}_K] = [\mathbf{d}_1^t, \dots, \mathbf{d}_N^t]$$

A Motivating Example

Phase (c) High-order proximity embedding



Singular Value Decompose

$$\mathbf{M} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{i=1}^N \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \mathbf{D}^s \cdot \mathbf{D}^t$$

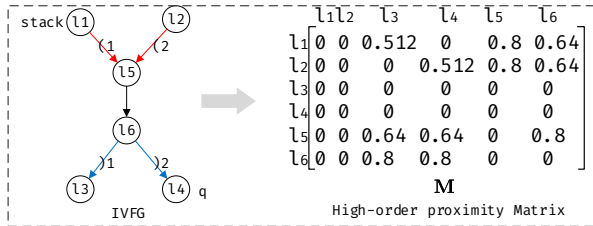
$$\mathbf{D}^s = [\sqrt{\sigma_1} \cdot \mathbf{u}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{u}_K] = [\mathbf{d}_1^s, \dots, \mathbf{d}_N^s]$$

$$\mathbf{D}^t = [\sqrt{\sigma_1} \cdot \mathbf{v}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{v}_K] = [\mathbf{d}_1^t, \dots, \mathbf{d}_N^t]$$

$$\begin{bmatrix} -0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0.33 & -0.79 & | & 0 & 0 & 0 \\ 0 & 0.71 & -0.58 & | & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{D}^s$$
$$\begin{bmatrix} 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ -0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ 0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ 0 & -0.71 & -0.58 & | & 0 & 0 & 0 \\ 0 & -0.33 & -0.79 & | & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{D}^t$$

A Motivating Example

Phase (c) High-order proximity embedding



Singular Value Decompose

$$\mathbf{M} \approx \mathbf{U}\Sigma\mathbf{V}^T = \sum_{i=1}^N \sigma_i \mathbf{u}_i \mathbf{v}_i^T = \mathbf{D}^s \cdot \mathbf{D}^{tT}$$

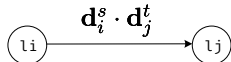
$$\mathbf{D}^s = [\sqrt{\sigma_1} \cdot \mathbf{u}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{u}_K] = [\mathbf{d}_1^s, \dots, \mathbf{d}_N^s]$$

$$\mathbf{D}^t = [\sqrt{\sigma_1} \cdot \mathbf{v}_1, \dots, \sqrt{\sigma_K} \cdot \mathbf{v}_K] = [\mathbf{d}_1^t, \dots, \mathbf{d}_N^t]$$

$$\begin{bmatrix} -0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0.33 & -0.79 & | & 0 & 0 & 0 \\ 0 & 0.71 & -0.58 & | & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{D}^s$$

$$\begin{bmatrix} 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ -0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ 0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ 0 & -0.71 & -0.58 & | & 0 & 0 & 0 \\ 0 & -0.33 & -0.79 & | & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{D}^t$$

Reachability between
i and j



A Motivating Example

Phase (c) High-order proximity embedding

$$\mathbf{d}_1^s \begin{bmatrix} -0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0.51 & -0.49 & -0.69 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0.33 & -0.79 & | & 0 & 0 & 0 \\ 0 & 0.71 & -0.58 & | & 0 & 0 & 0 \end{bmatrix}$$

\mathbf{D}^s

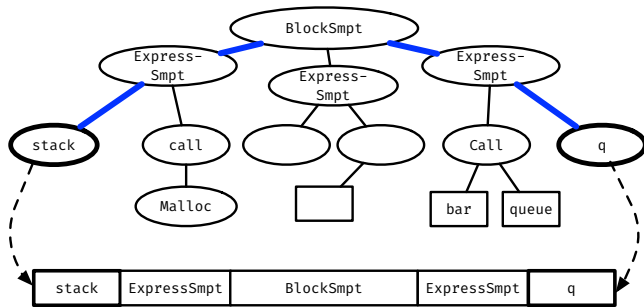
$$\begin{bmatrix} 0 & 0 & 0 & | & 0 & 0 & 0 \\ 0 & 0 & 0 & | & 0 & 0 & 0 \\ \mathbf{d}_3^t & -0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ \mathbf{d}_4^t & 0.51 & 0.49 & -0.69 & | & 0 & 0 & 0 \\ \mathbf{d}_5^t & 0 & -0.71 & -0.58 & | & 0 & 0 & 0 \\ \mathbf{d}_6^t & 0 & -0.33 & -0.79 & | & 0 & 0 & 0 \end{bmatrix}$$

\mathbf{D}^t

Reachability	Path length
$\mathbf{d}_1^s \cdot \mathbf{d}_5^t{}^\top = 0.75$	1
$\mathbf{d}_1^s \cdot \mathbf{d}_6^t{}^\top = 0.71$	2
$\mathbf{d}_1^s \cdot \mathbf{d}_3^t{}^\top = 0.5$	3
$\mathbf{d}_1^s \cdot \mathbf{d}_4^t{}^\top = -0.02$	infeasible

A Motivating Example

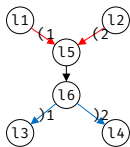
Phase (c) High-order proximity embedding



infeasible dependence relation between `stack` to `q`

A Motivating Example

Phase (d) Value-Flow Vector and Applications



$$\left[\begin{array}{ccc|ccc} -0.51 & -0.49 & -0.69 & 0 & 0 & 0 \\ 0.51 & -0.49 & -0.69 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0.33 & -0.79 & 0 & 0 & 0 \\ 0 & 0.71 & -0.58 & 0 & 0 & 0 \end{array} \right]$$

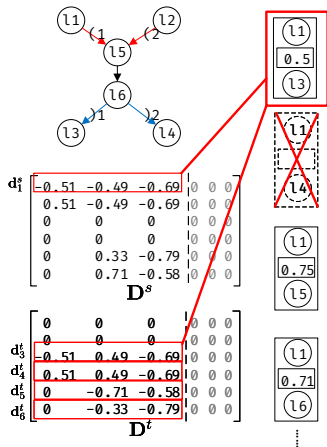
D^s

$$\left[\begin{array}{ccc|ccc} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -0.51 & 0.49 & -0.69 & 0 & 0 & 0 \\ 0.51 & 0.49 & -0.69 & 0 & 0 & 0 \\ 0 & -0.71 & -0.58 & 0 & 0 & 0 \\ 0 & -0.33 & -0.79 & 0 & 0 & 0 \end{array} \right]$$

D^t

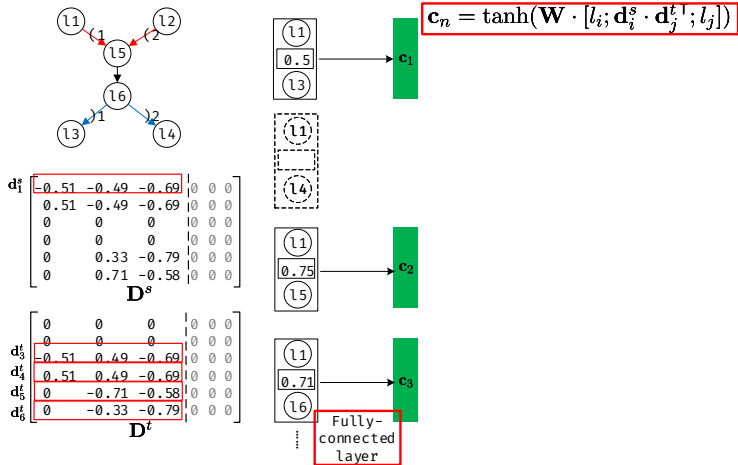
A Motivating Example

Phase (d) Value-flow Vectors and Applications



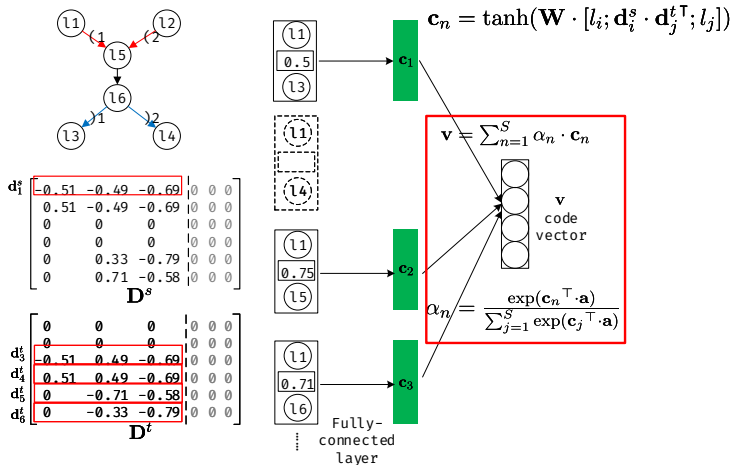
A Motivating Example

Phase (d) Value-flow Vectors and Applications



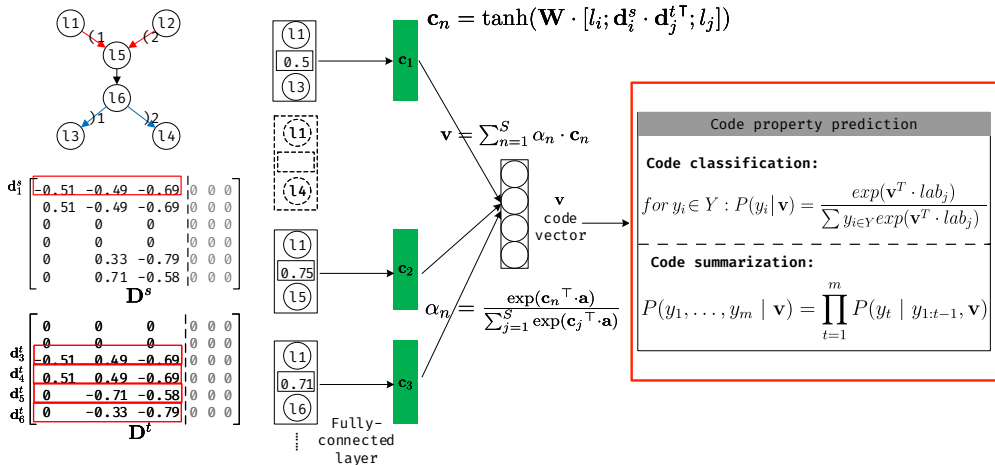
A Motivating Example

Phase (d) Value-flow Vectors and Applications



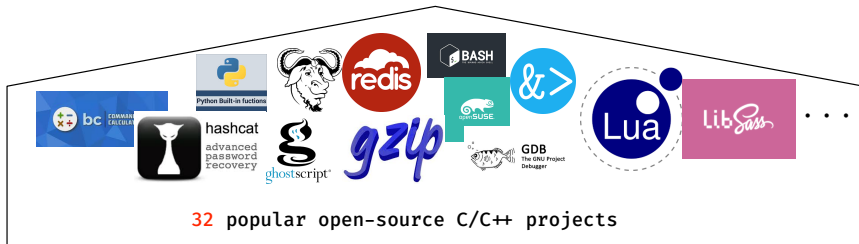
A Motivating Example

Phase (d) Value-flow Vectors and Applications



Experimental Evaluation

Benchmarks



Total Line of Instructions:4,922,162

Total Methods:17,529

Total Pointers: 2,913,748

Total Objects: 190,157

Total Number of Calls:536,033

Total IVFGNodes: 4,637,301

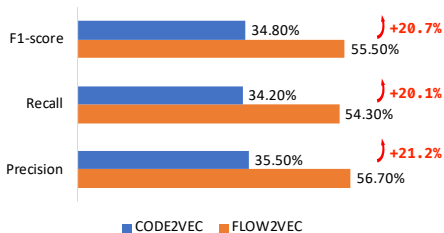
Total IVFGEEdges: 6,531,578

*Conducted machine: Intel Xeon Gold 6132 @ 2.60GHz CPUs and 128GB of RAM (All finish analyzing in 272.5mins)

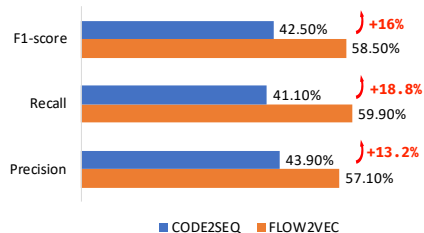
Experimental Evaluation

Comparison with baselines

FLOW2VEC *VS* CODE2VEC

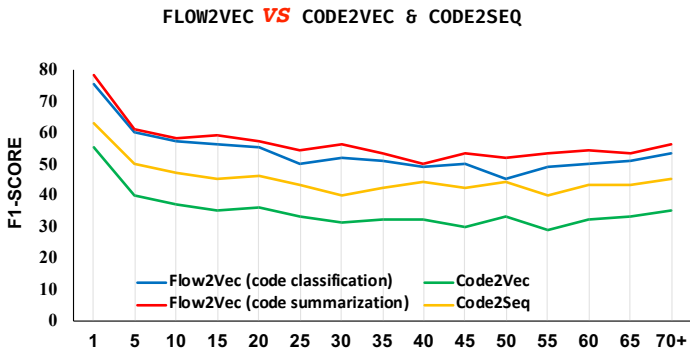


FLOW2VEC *VS* CODE2SEQ



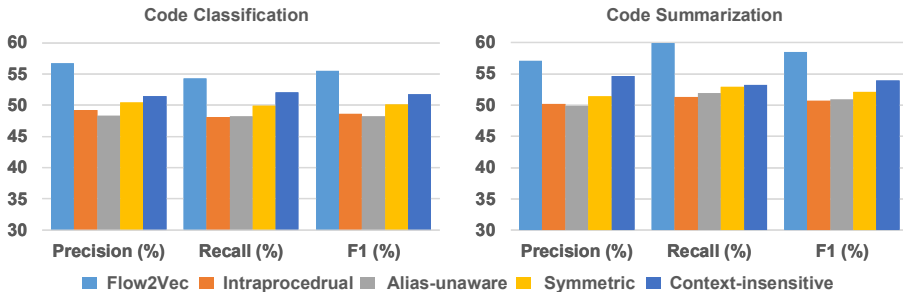
Experimental Evaluation

F1-score under different lengths of code



Experimental Evaluation

Ablation analysis



Thanks!

Q & A