



Understanding Source Code with Deep Learning

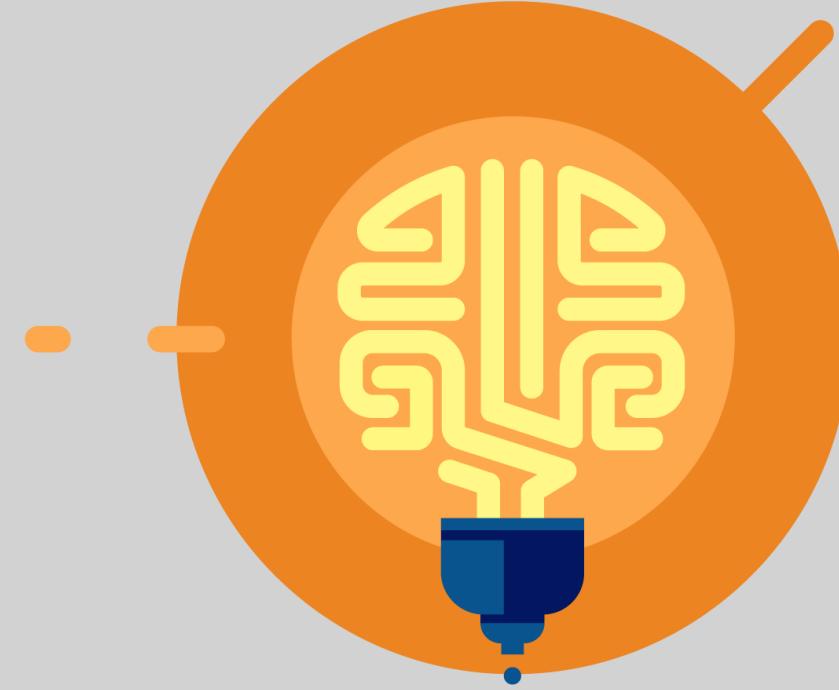
Miltos Allamanis

FOSDEM 2019

 miltos1

 <https://miltos.allamanis.com>

Microsoft Research Cambridge



Source Code is Bimodal

```
        selectedScopes, el
        var previous = $attr.ngScope;
        var selectedTranscludes = [];
        var previousElements = [];
        var previousElements = [];
        selectedScopes = [];

        scope.$watch($watchExpr, function() {
            var i, ii;
            for (ii = 0, ii = previous.length; ii--) {
                previousElements[ii].remove();
            }
            previousElements.length = 0;

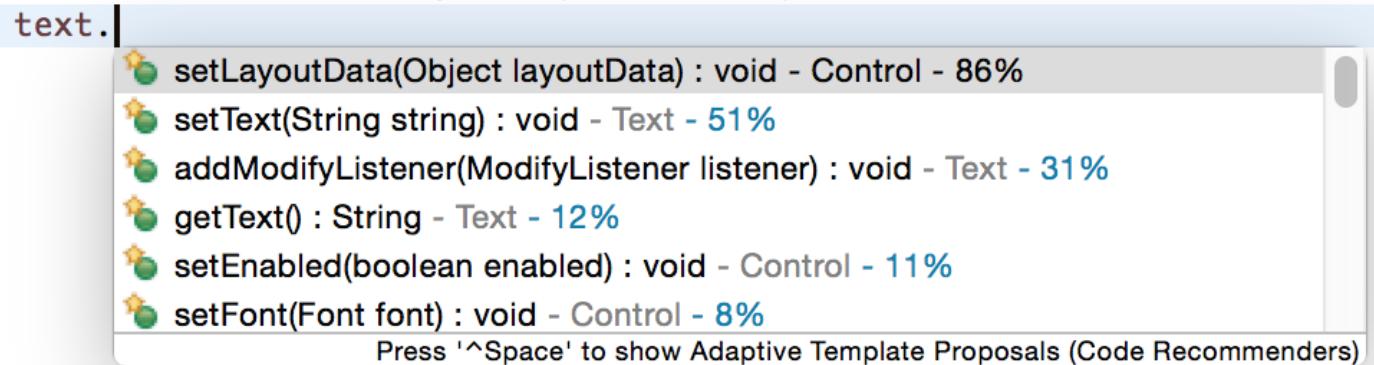
            for (ii = 0, ii = selectedScopes.length; ii++) {
                var selected = selectedElements[ii];
                selectedScopes[ii].$destroy();
                previousElements[ii] = selected;
                $animate.leave(selected, previousElements);
                previousElements.splice(ii, 1);
            }
        });

        selectedElements.length = 0;
        selectedScopes.length = 0;

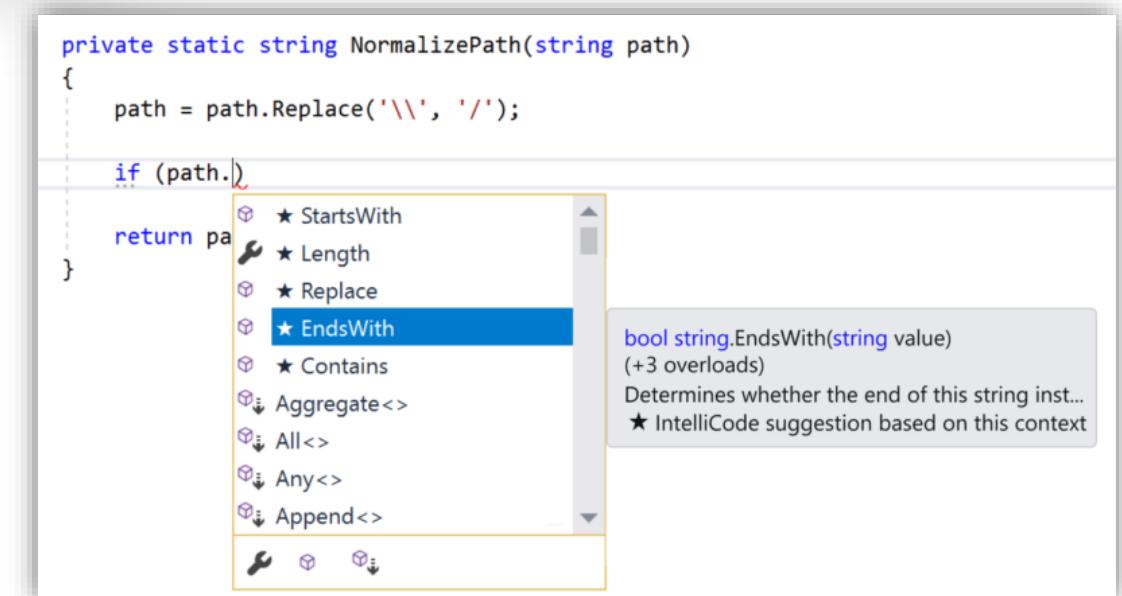
        if ((selectedTranscludes = ngScope.$eval(attr.change));
            forEach(selectedTranscludes, function(transclude) {
                var selectedScope = scope.$new();
                selectedScopes.push(selectedScope);
                selectedScope.$on('$destroy', function() {
                    selectedScopes.pop();
                });
            })
        );
    }
});
```

Code Autocompletion

```
Text text = new Text(parent, SWT.NONE);
```



<http://www.eclipse.org/recommenders/>



<https://visualstudio.microsoft.com/services/intellicode/>

Predicting Types

JS NICE STATISTICAL RENAMING, TYPE INFERENCE AND DEOBFUSCATION ABOUT

ENTER JAVASCRIPT NICIFY JAVASCRIPT

```
// Put your JavaScript here that you want to rename, deobfuscate, or infer types for:
function chunkData(e, t) {
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t) {
    if (i + t < r) {
      n.push(e.substring(i, i + t));
    } else {
      n.push(e.substring(i, r));
    }
  }
  return n;
}
// You can also use some ES6 features.
const get = (a,b) => a.getElementById(b);

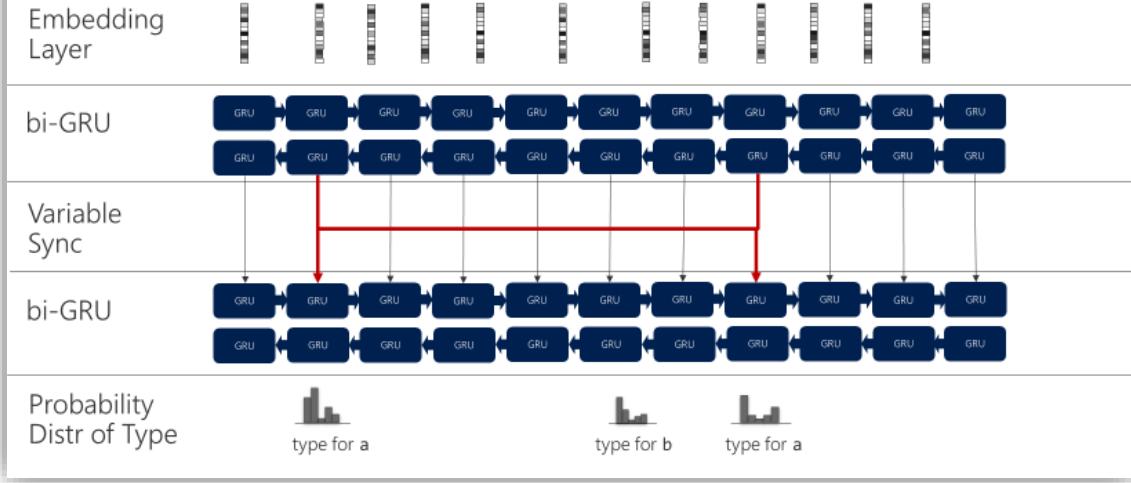
```

RESULT. ARE YOU SATISFIED? YES NO

```
'use strict';
/*
 * @param {string} bin
 * @param {number} size
 * @return {?}
 */
function chunkData(bin, size) {
  /* @type {!Array} */
  var results = [];
  var length = bin.length;
  /* @type {number} */
  var i = 0;
  for (; i < length; i = i + size) {
    if (i + size < length) {
      results.push(bin.substring(i, i + size));
    } else {
      results.push(bin.substring(i, length));
    }
  }
  return results;
}
const get = (doc, key) =>
```

DeepTyper

let **a** = 1 ; let **b** = **a** + 1 ;



Predicting Program Properties from Code

V. Raychev, M. Vechev, A. Krause. 2015

<http://jsnice.org/>

Deep Learning Type Inference

V. Hellendoorn, C. Bird, E.T. Barr, M. Allamanis. 2018

Inferring Type Refinements

Dash *et al.* 2018 "RefiNym: Using Names to Refine Types"

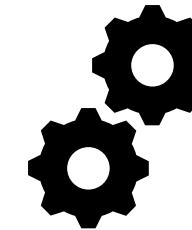


Conceptual Types

"a password"

"a JSON string"

Latent; we don't observe in the *conceptual* types.



Defined Types

`string password;`

`string data = Json.Load();`

Defined explicitly by the programmer.

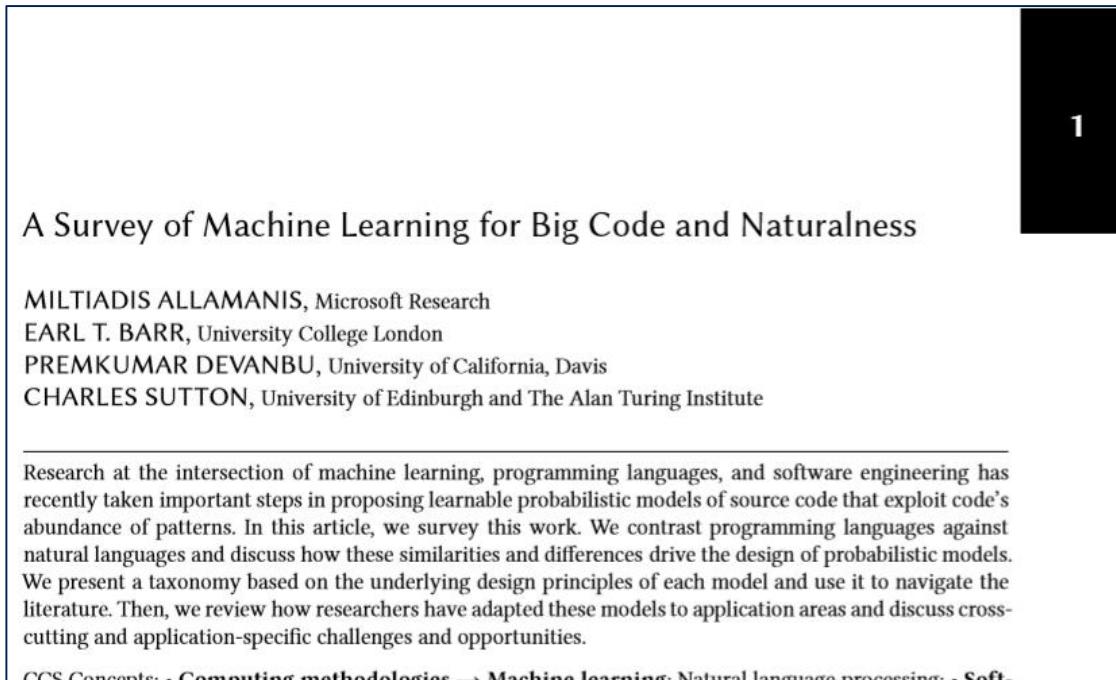
Argument Swapping

Type	Parameter	Original argument	Correct argument
Duration	responseTTLDuration	frequencyCapDuration	responseTTLDuration
Duration	frequencyCapDuration	responseTTLDuration	frequencyCapDuration
List<A>	slotResponse	slotResponse	slotResponse
Builder	builder	builder	builder
boolean	isTransposed	isTransposed	isTransposed
int	startColumnIndex	a.getStartColumnIndex()	0
int	endColumnIndex	a.getEndColumnIndex()	rows.size()
int	startRow	0	a.getStartColumnIndex()
int	endRow	rows.size()	a.getEndColumnIndex()

Research in ML+Code

- Infer latent intent
- Ambiguous information
- Learned heuristics

<https://ml4code.github.io>



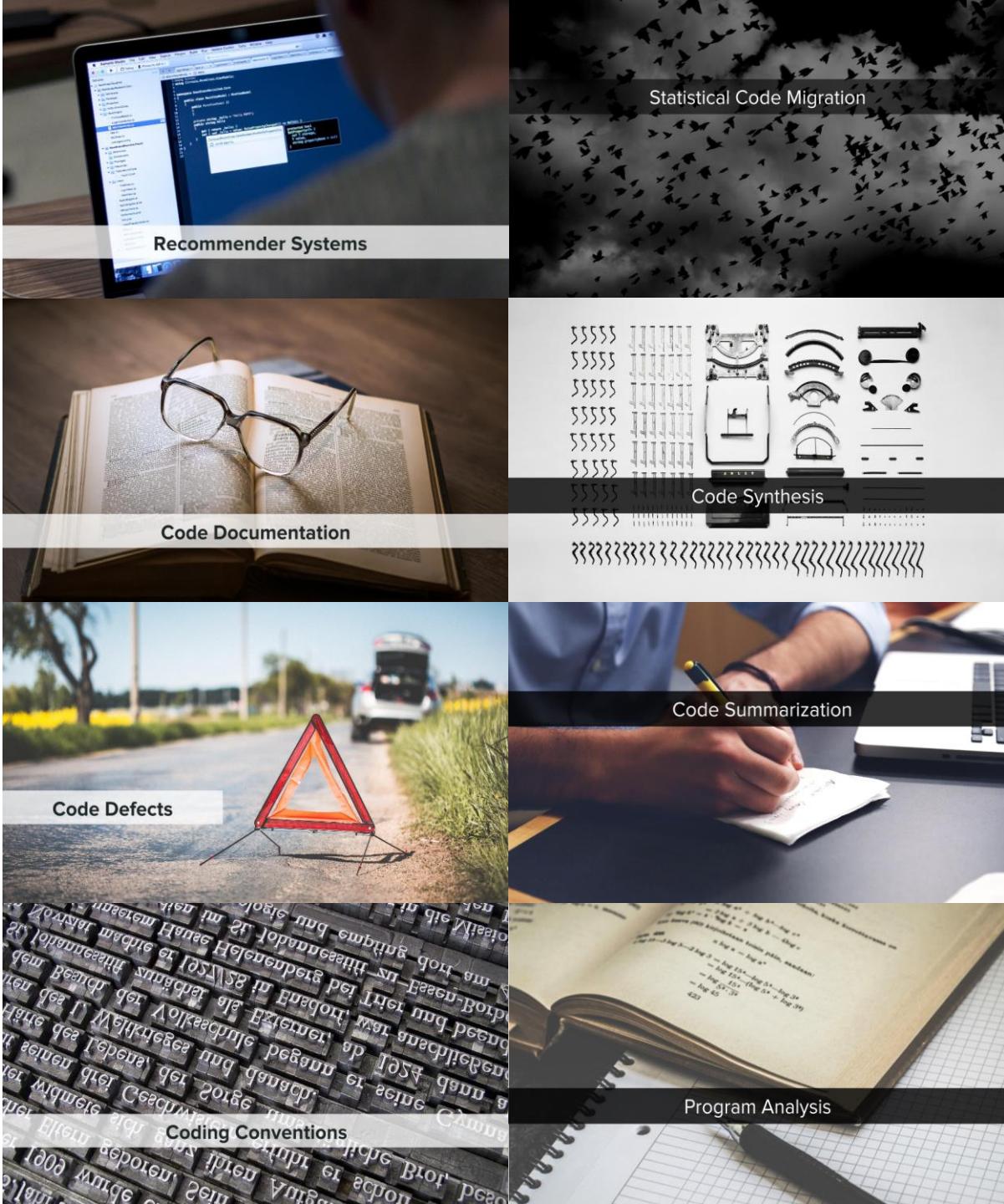
A Survey of Machine Learning for Big Code and Naturalness

MILTIADIS ALLAMANIS, Microsoft Research
EARL T. BARR, University College London
PREMKUMAR DEVANBU, University of California, Davis
CHARLES SUTTON, University of Edinburgh and The Alan Turing Institute

1

Research at the intersection of machine learning, programming languages, and software engineering has recently taken important steps in proposing learnable probabilistic models of source code that exploit code's abundance of patterns. In this article, we survey this work. We contrast programming languages against natural languages and discuss how these similarities and differences drive the design of probabilistic models. We present a taxonomy based on the underlying design principles of each model and use it to navigate the literature. Then, we review how researchers have adapted these models to application areas and discuss cross-cutting and application-specific challenges and opportunities.

CCS Concepts: - Computing methodologies → Machine learning; Natural language processing; - Software engineering → Code analysis and understanding; Code generation; Code transformation; Code evolution; Code mining; Code summarization; Code synthesis; Code translation; Code verification; Code visualization; Code mining; Code summarization; Code synthesis; Code translation; Code verification; Code visualization;





Detecting Variable Misuse Bugs

By representing source code as graphs

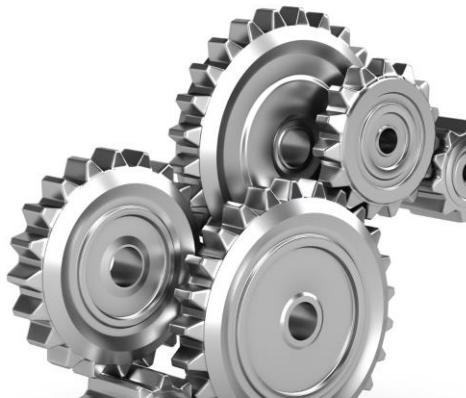
Target Task

```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(first);  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Possible type-correct options: `clazz`, `first`



Not easy to catch with static analysis tools.

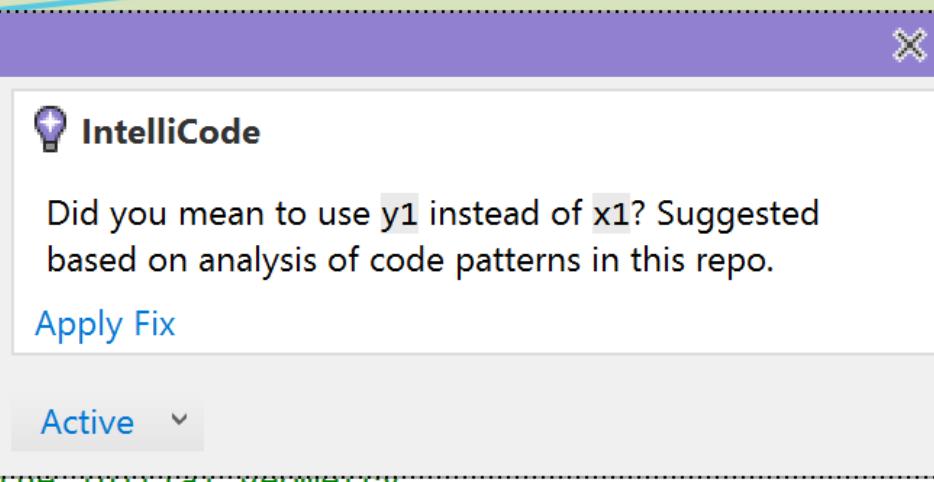


```
        var newBounds = new List<Rect>(bounds.Count);
        foreach (var b in bounds)
        {
            double x1 = b.Left - padding.Left;
            double x2 = b.Right + padding.Right;
            if (x1 < x2)
            {
                double y1 = b.TextTop - padding.Top;
                double y2 = b.TextBottom + padding.Bottom;

                newBounds.Add(new Rect(x1, y1, x2 - x1, y2 - x1));
            }
        }
    }

    return new
}
```

```
public static
{
    if (rectan
        return
    // Set up the initial geometry
}
```



Rect> rectar

Programs as Graphs: Key Idea

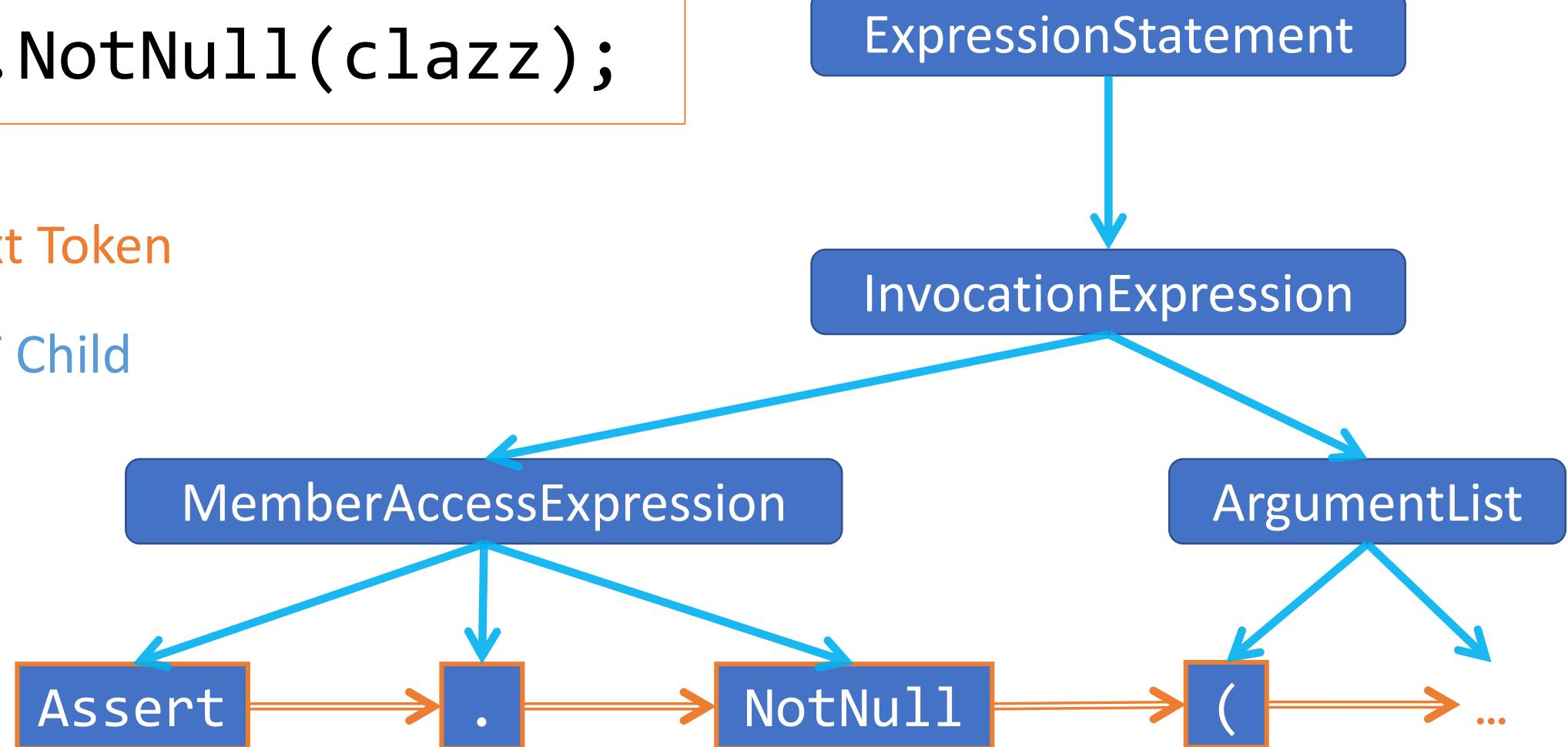
```
int SumPositive(int[] arr, int lim) {  
    int sum = 0;  
    for (int i = 0; i < lim; i++)  
        if (arr[i] > 0)  
            sum += arr[i];  
  
    return sum;  
}
```

Programs as Graphs: Syntax

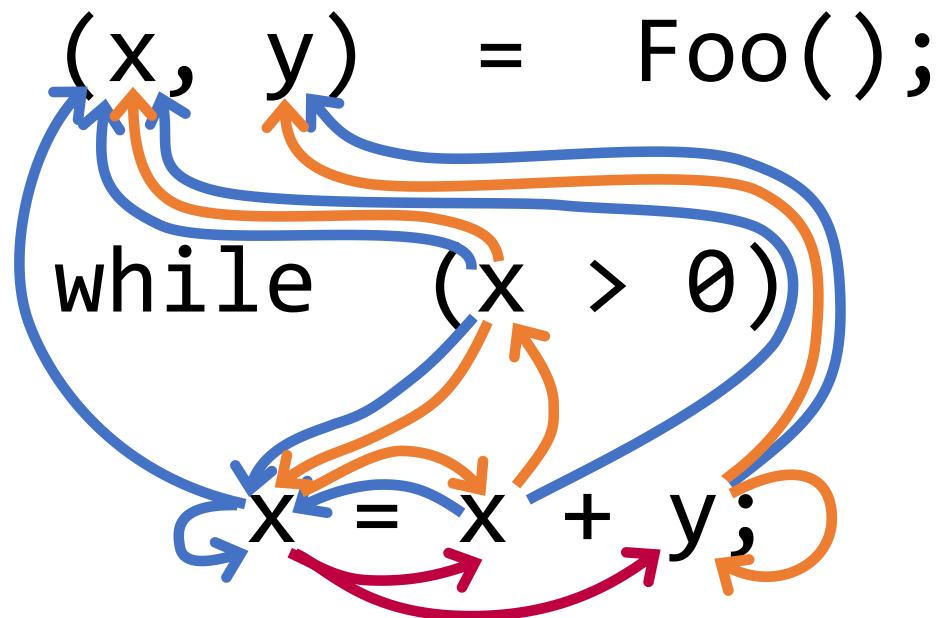
Assert.NotNull(clazz);

→ Next Token

→ AST Child



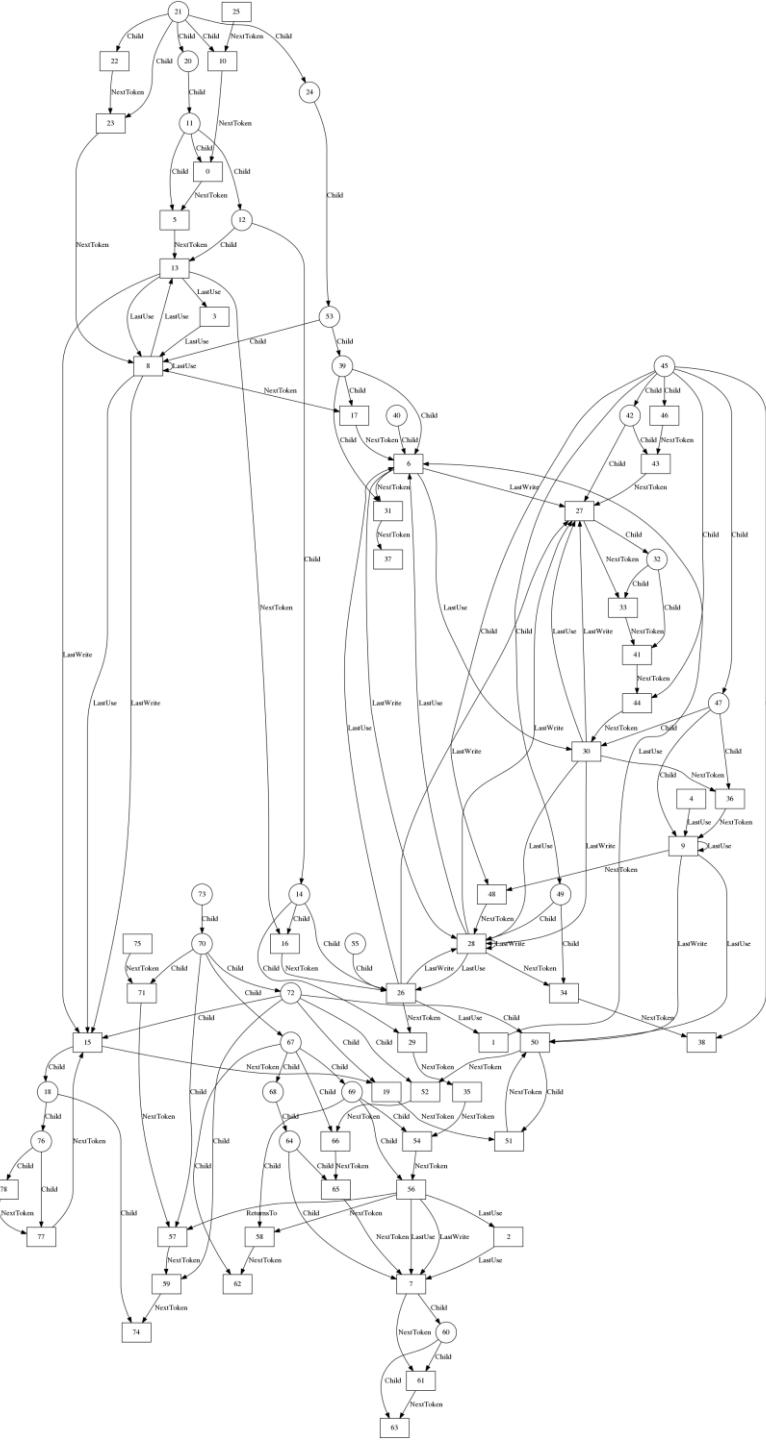
Programs as Graphs: Data Flow



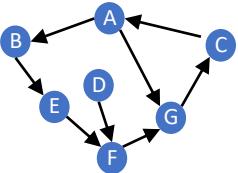
- Last Write
- Last Use
- Computed From

Programs as Graphs

```
int SumPositive(int[] arr, int lim) {  
    int sum = 0;  
    for (int i=0; i < lim; ++)  
        if (arr[i] > 0)  
            sum += arr[i];  
  
    return sum;  
}
```



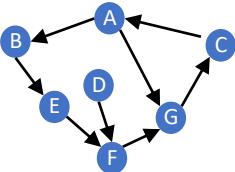
~900 nodes/graph ~8k edges/graph



Graph Representation for Variable Misuse

```
var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(clazz);  
  
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;  
Assert.NotNull(first);  
  
Assert.Equal("string", first.Properties["Name"].Name);  
Assert.False(clazz.Properties["Name"].IsArray);
```

Possible type-correct options: `clazz`, `first`



Graph Representation for Variable Misuse

```

var clazz=classTypes["Root"].Single() as JsonCodeGenerator.ClassType;
Assert.NotNull(clazz);

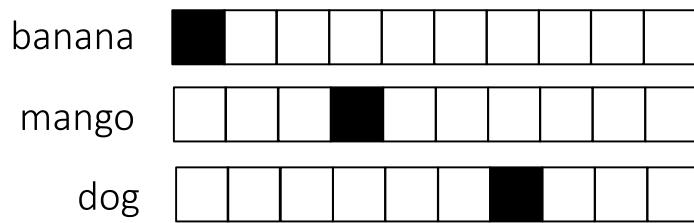
var first=classTypes["RecClass"].Single() as JsonCodeGenerator.ClassType;
Assert.NotNull( SLOT );           first      clazz
                                ↓          ↓
Assert.Equal("string", first.Properties["Name"].Name);
Assert.False(clazz.Properties["Name"].IsArray);

```

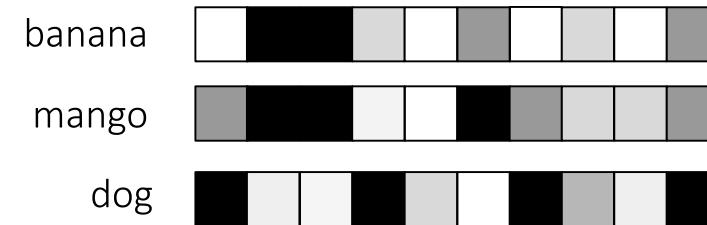
Goal: make the representation of SLOT as close as possible to the representation of the correct candidate node

$$(\mathbf{h}_T^{SLOT})^T \mathbf{h}_T^{first} \gg (\mathbf{h}_T^{SLOT})^T \mathbf{h}_T^{clazz}$$

Vector Space Representations

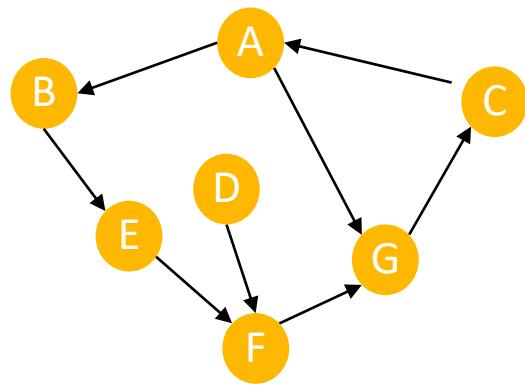


Local representation

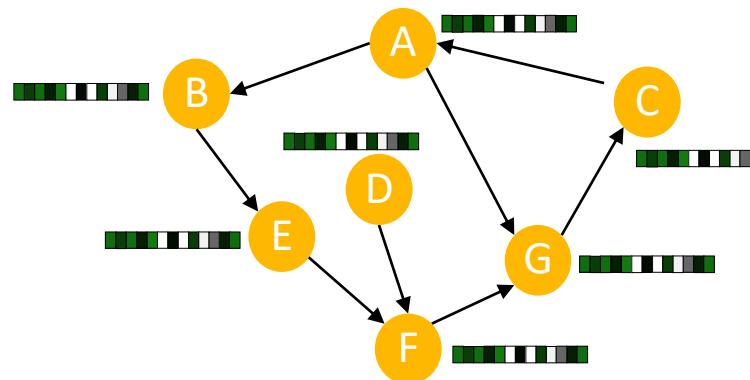


Distributed representation

Graph Neural Networks



Graph Representation
of Problem

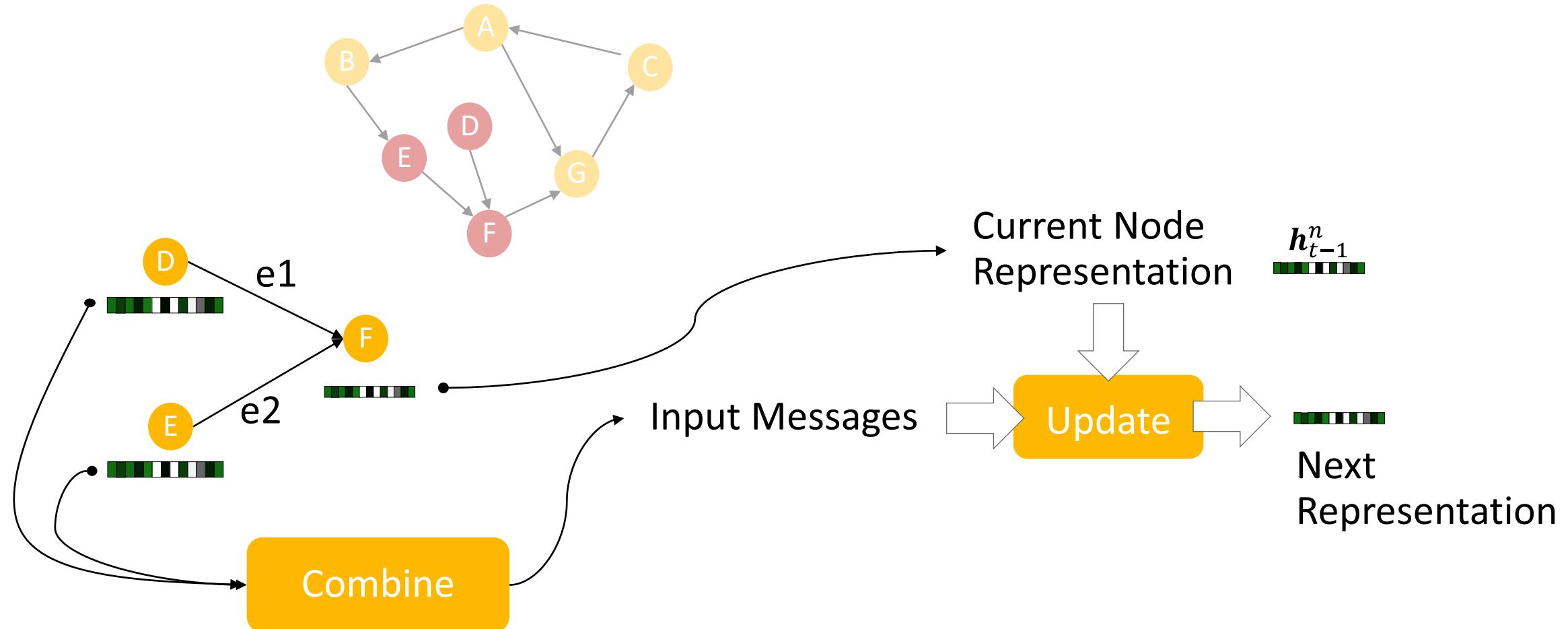


Initial Representation
of each node

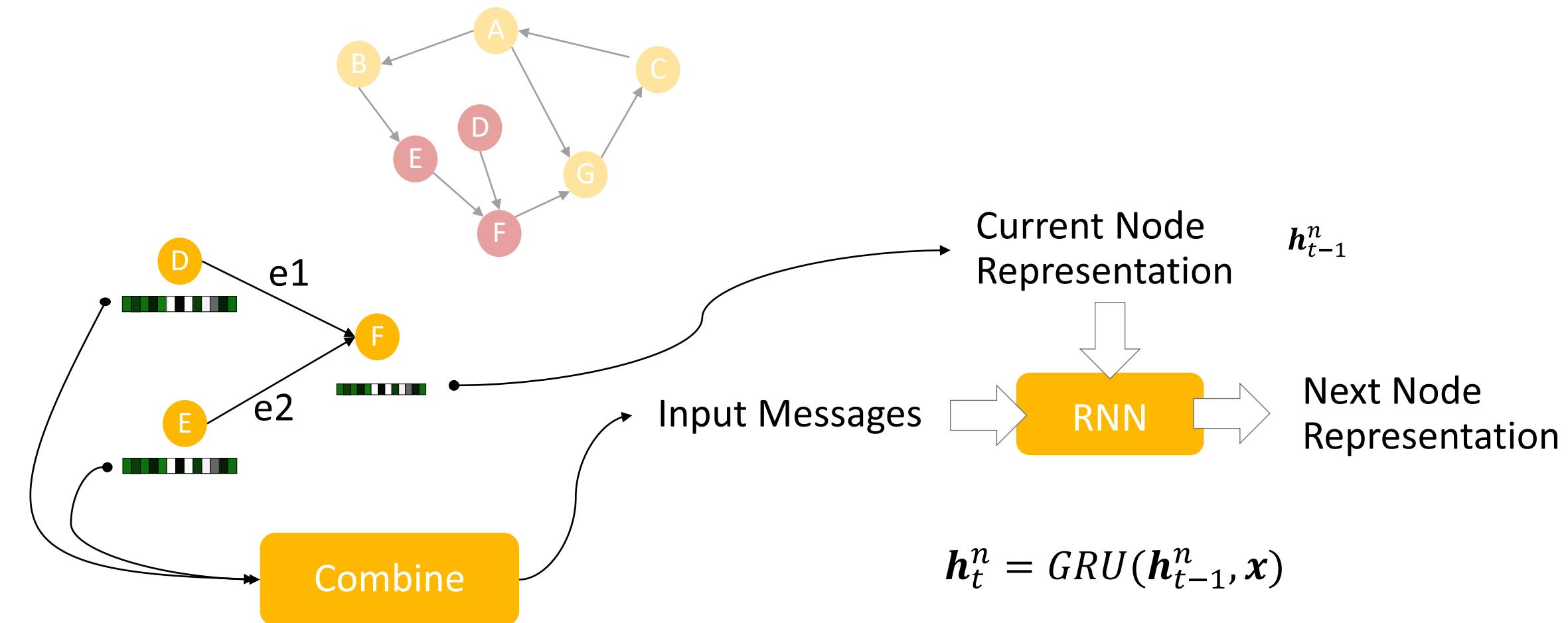
Li et al (2015). Gated Graph Sequence Neural Networks.

Gilmer et al (2017). Neural Message Passing for Quantum Chemistry.

Graph Neural Networks: Message Passing

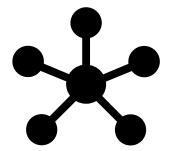


Graph Neural Networks: Message Passing

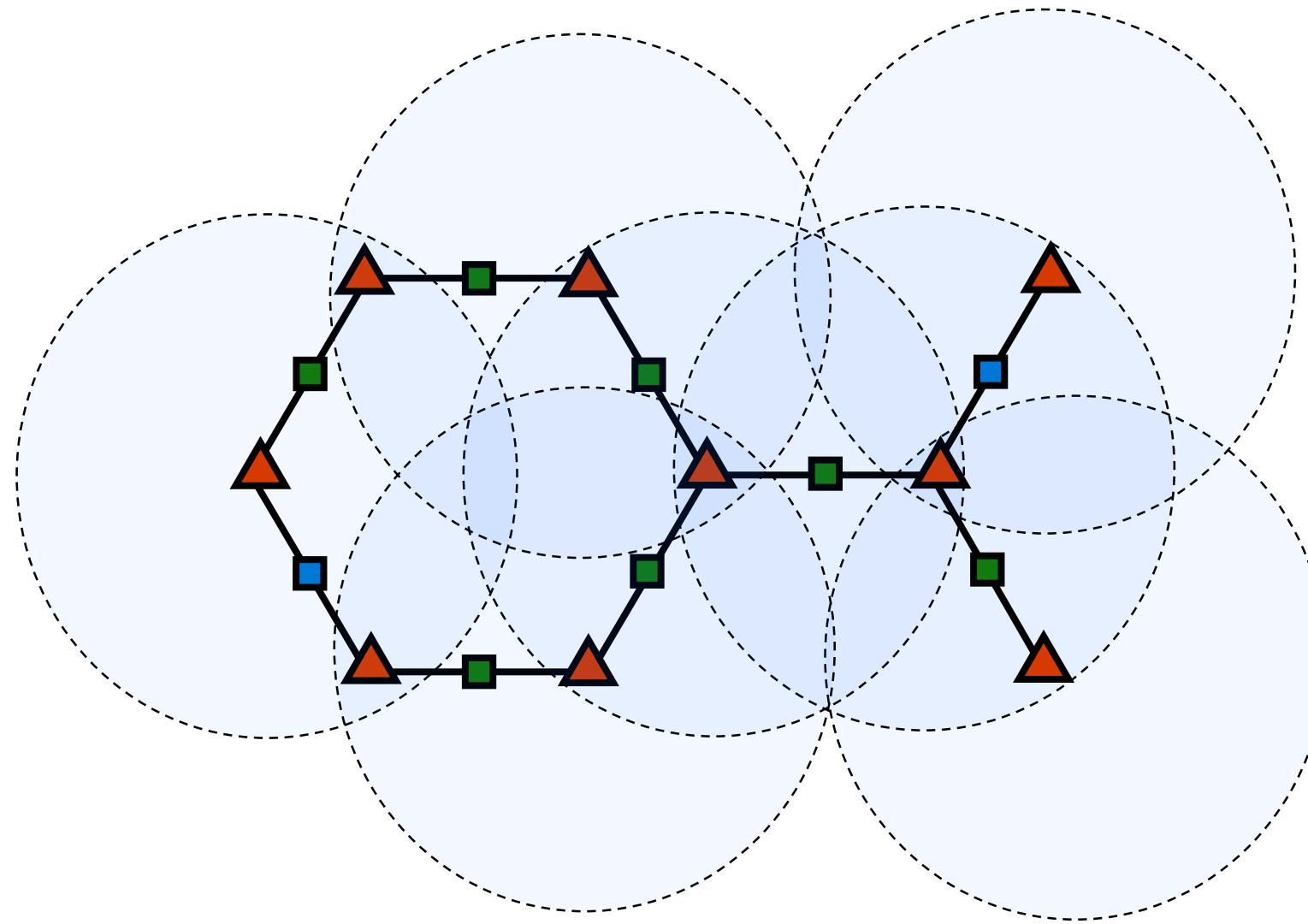


$$x = \sum_{n' \in \text{neig}(n)} E_{\tau(n' \rightarrow n)} h_{t-1}^{n'} + b_{\tau(n' \rightarrow n)}$$

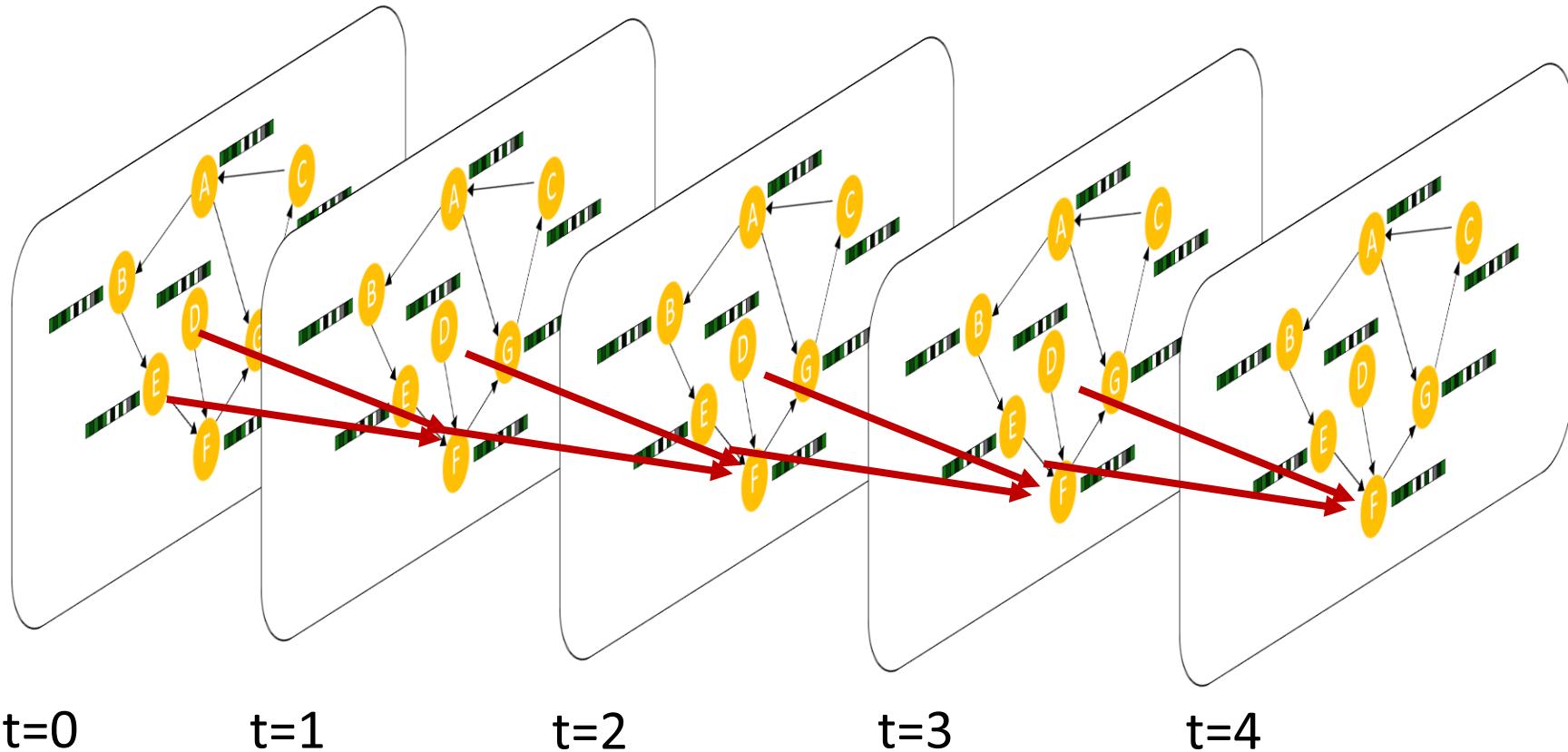
Li et al (2015). Gated graph sequence neural networks.



Graph Neural Networks: Unrolling

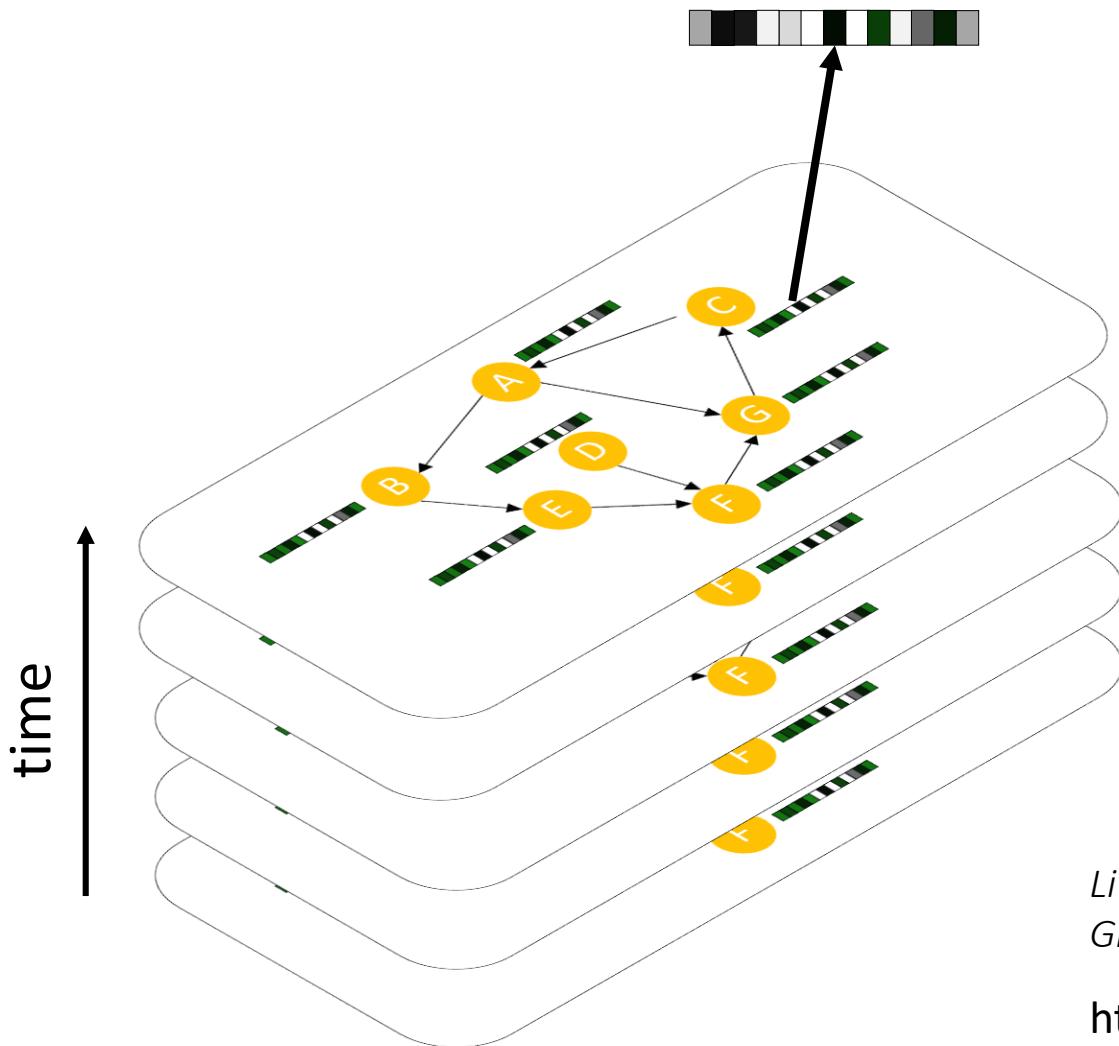


Graph Neural Networks: Unrolling



Li et al (2015). Gated graph sequence neural networks.

Graph Neural Networks: Unrolling



- node selection
- node classification
- graph classification

*Li et al (2015). Gated Graph Sequence Neural Networks.
Gilmer et al (2017). Neural Message Passing for Quantum Chemistry.*

<https://github.com/Microsoft/gated-graph-neural-network-samples>

Quantitative Results – Variable Misuse

Accuracy (%)	BiGRU	BiGRU+Dataflow	GGNN
Seen Projects	50.0	73.7	85.5

Seen Projects: 24 F/OSS C# projects (2060 kLOC): Used for train and test

3.8 type-correct alternative variables per slot (median 3, $\sigma=2.6$)

Quantitative Results – Variable Misuse

Accuracy (%)	BiGRU	BiGRU+Dataflow	GGNN
Seen Projects	50.0	73.7	85.5
Unseen Projects	28.9	60.2	78.2

Seen Projects: 24 F/OSS C# projects (2060 kLOC): Used for train and test

Unseen Projects: 3 F/OSS C# projects (228 kLOC): Used only for test

3.8 type-correct alternative variables per slot (median 3, $\sigma=2.6$)

```
bool TryFindGlobalDirectivesFile(string baseDirectory, string fullPath, out string path) {
    baseDirectory = baseDirectory.TrimEnd(Path.DirectorySeparatorChar);
    var directivesDirectory = Path.GetDirectoryName(████████)
                                .TrimEnd(Path.DirectorySeparatorChar);
    while (directivesDirectory != null && directivesDirectory.Length >= baseDirectory.Length)
    {
        path = Path.Combine(directivesDirectory, GlobalDirectivesFileName);
        if (File.Exists(path)) return true;

        directivesDirectory = Path.GetDirectoryName(directivesDirectory)
                                .TrimEnd(Path.DirectorySeparatorChar);
    }
    path = null;
    return false;
}
```

What the model sees...



UI/UX



ML Capabilities



Metrics

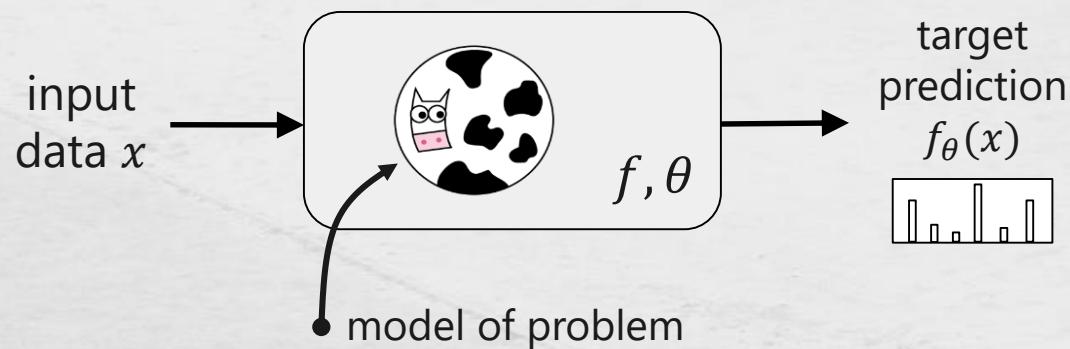


Low resources

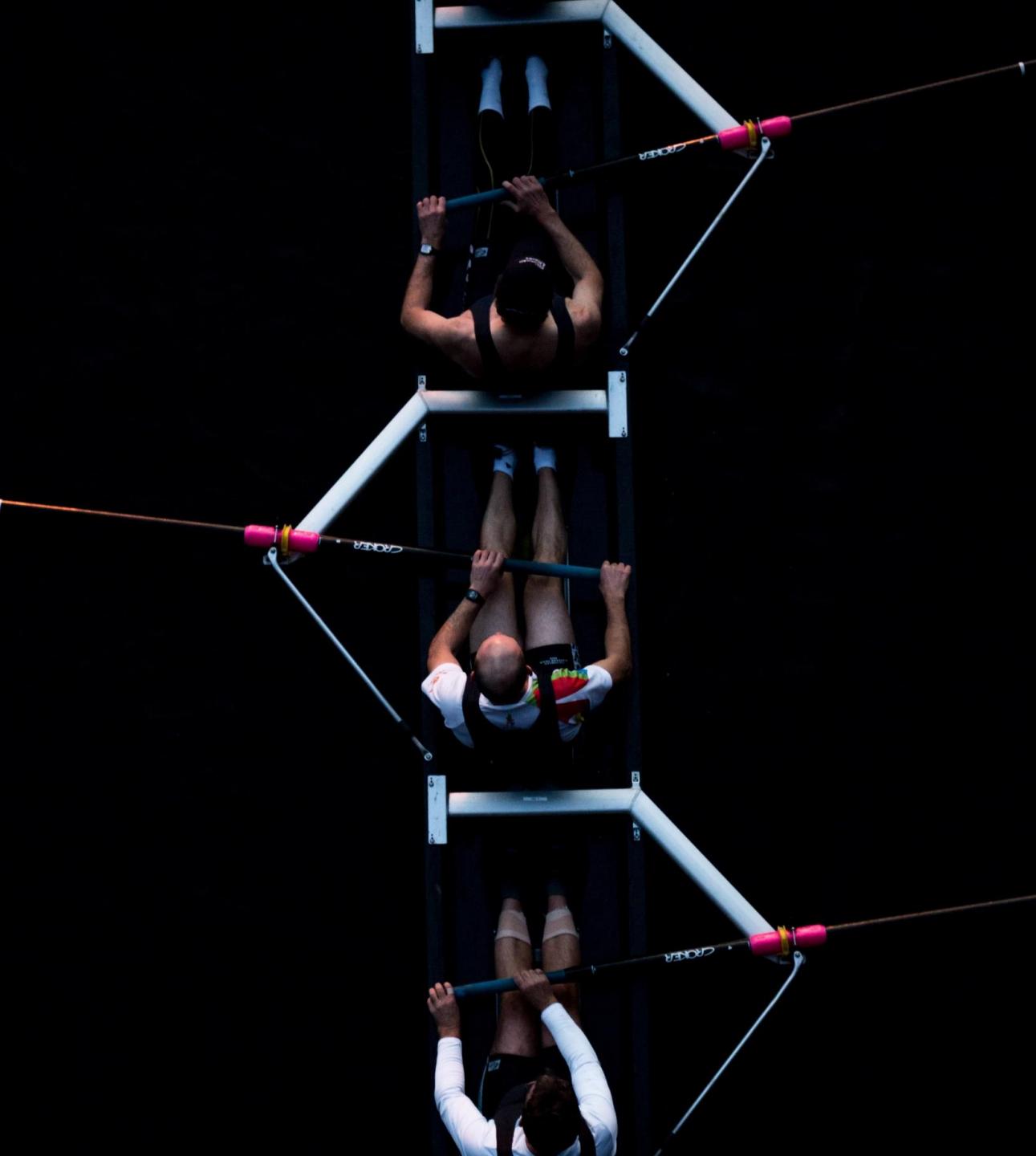




Learning Signals



- Given dataset $\{(x_1, y_0), \dots, (x_N, y_N)\}$
- Minimize Loss $\mathcal{L}(\theta) = \frac{1}{N} \sum_i L(f_\theta(x_i), y_i)$





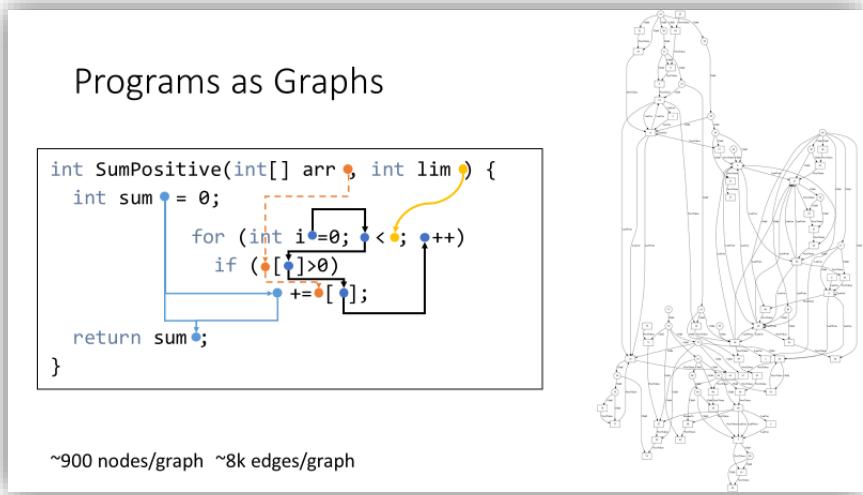
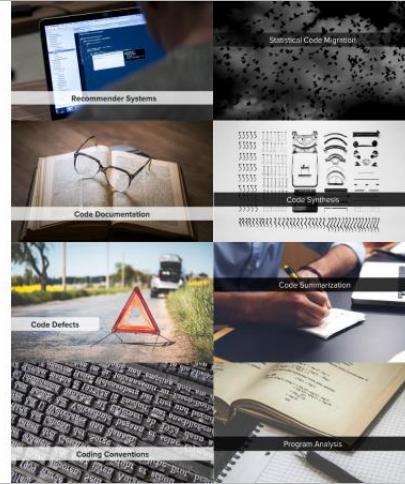
Research in ML+Code

- Infer latent intent
- Ambiguous information
- Learned heuristics

<https://ml4code.github.io>

A Survey of Machine Learning for Big Code and Naturalness
MILTIADIS ALLAMANIS, Microsoft Research
EARL T. BARR, University College London
PREMKUMAR DEVANBU, University of California, Davis
CHARLES SUTTON, University of Edinburgh and The Alan Turing Institute

Research at the intersection of machine learning, programming languages, and software engineering has recently taken important steps in proposing learnable probabilistic models of source code that exploit code's abundance of patterns. In this survey, we survey this work. We contrast programming languages against natural languages and then discuss how they differ in terms of the types of patterns they exhibit and how they can be learned. We present a taxonomy based on the underlying design principles of each model and use it to navigate the literature. Then, we review how researchers have adapted these models to application areas and discuss cross-cutting and application-specific challenges and opportunities.



Graph Neural Networks: Message Passing

