



Towards the Future of Automated Programming: On the Use of Deep Learning for Source Code Understanding

Nghi Bui, Ph.D. | Oct 2021

About the Speaker

- Ph.D. in Computer Science from Singapore Management University.
Research Interests:
 - Software Analytics, Programming Languages, Deep Learning.
- Previous Affiliations:
 - Research Scientist: Software Analytics Research Lab, Singapore Management University.
 - Research Scientist: Trustworthy Open Source and Software Engineering Lab, Huawei Research Center in Ireland.
- Current Affiliation:
 - AI Scientist/Mentor: FSoft AI Lab.
- Website: <https://bdqnghi.github.io/>
- Github: <https://github.com/bdqnghi>

Software (1.0) is eating the world, and now AI (Software 2.0) is eating software.

Andrej Karpathy, Director of AI at Tesla

The screenshot shows a section of an O'Reilly Radar / AI & ML article. The header includes links for TEAMS, INDIVIDUALS, FEATURES, BLOG, and CONTENT SPONSORSHIP. The main title is "The road to Software 2.0". Below it, a subtext reads: "It's clear that AI can and will have a big influence on how we develop software."

The screenshot shows another section of the same O'Reilly Radar / AI & ML article. The main title is "What machine learning means for software development". Below it, a subtext reads: "'Human in the loop' software development will be a big part of the future."

The “classical stack” of **Software 1.0** is what we’re all familiar with — it is written in languages such as Python, C++, etc. It consists of explicit instructions to the computer written by a programmer. By writing each line of code, the programmer identifies a specific point in program space with some desirable behavior.

In contrast, **Software 2.0** is written in much more abstract, human unfriendly language, such as the weights of a neural network. No human is involved in writing this code because there are a lot of weights (typical networks might have millions), and coding directly in weights is kind of hard.

Software 2.0

Intelligent Software

Machine Programming

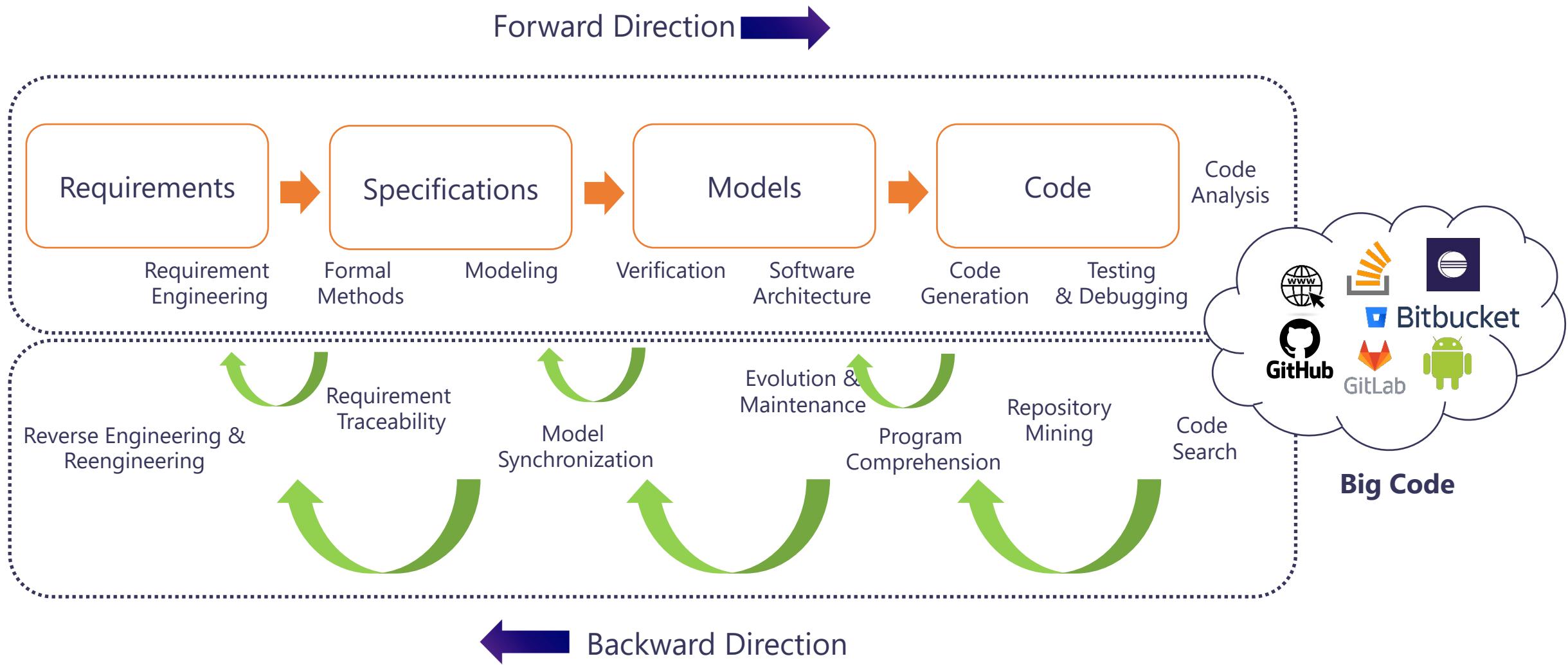
AI4Code

Automated Programming

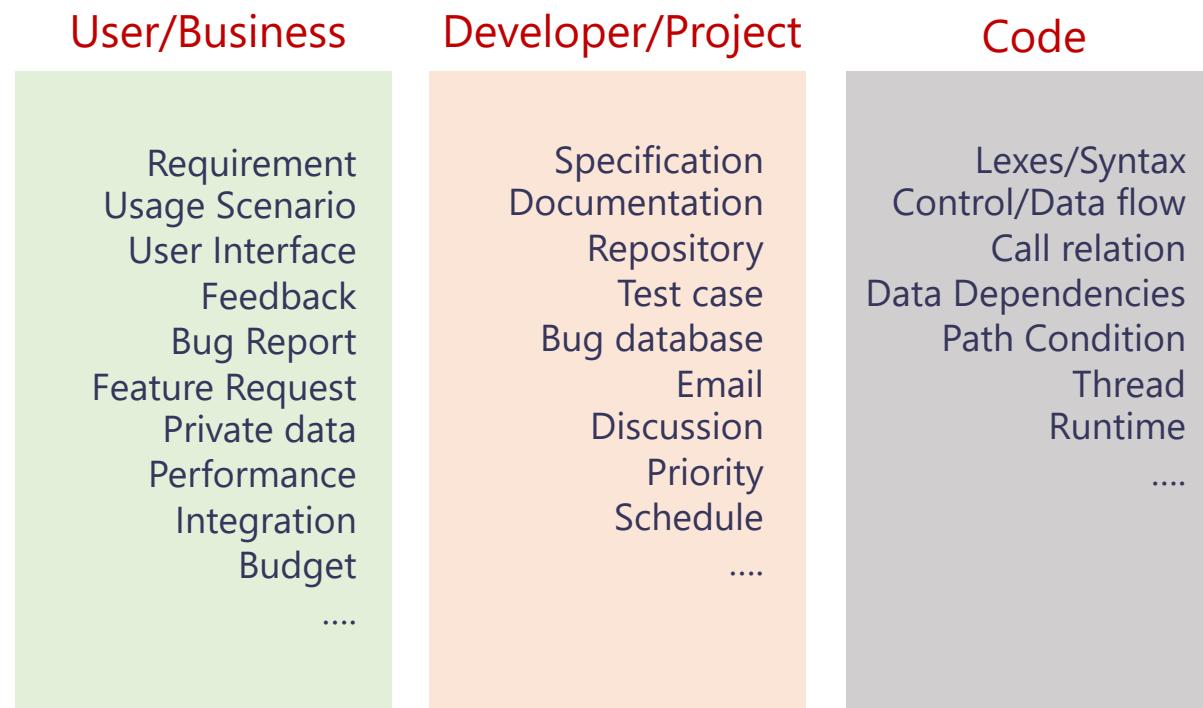
The Automation of Software Development



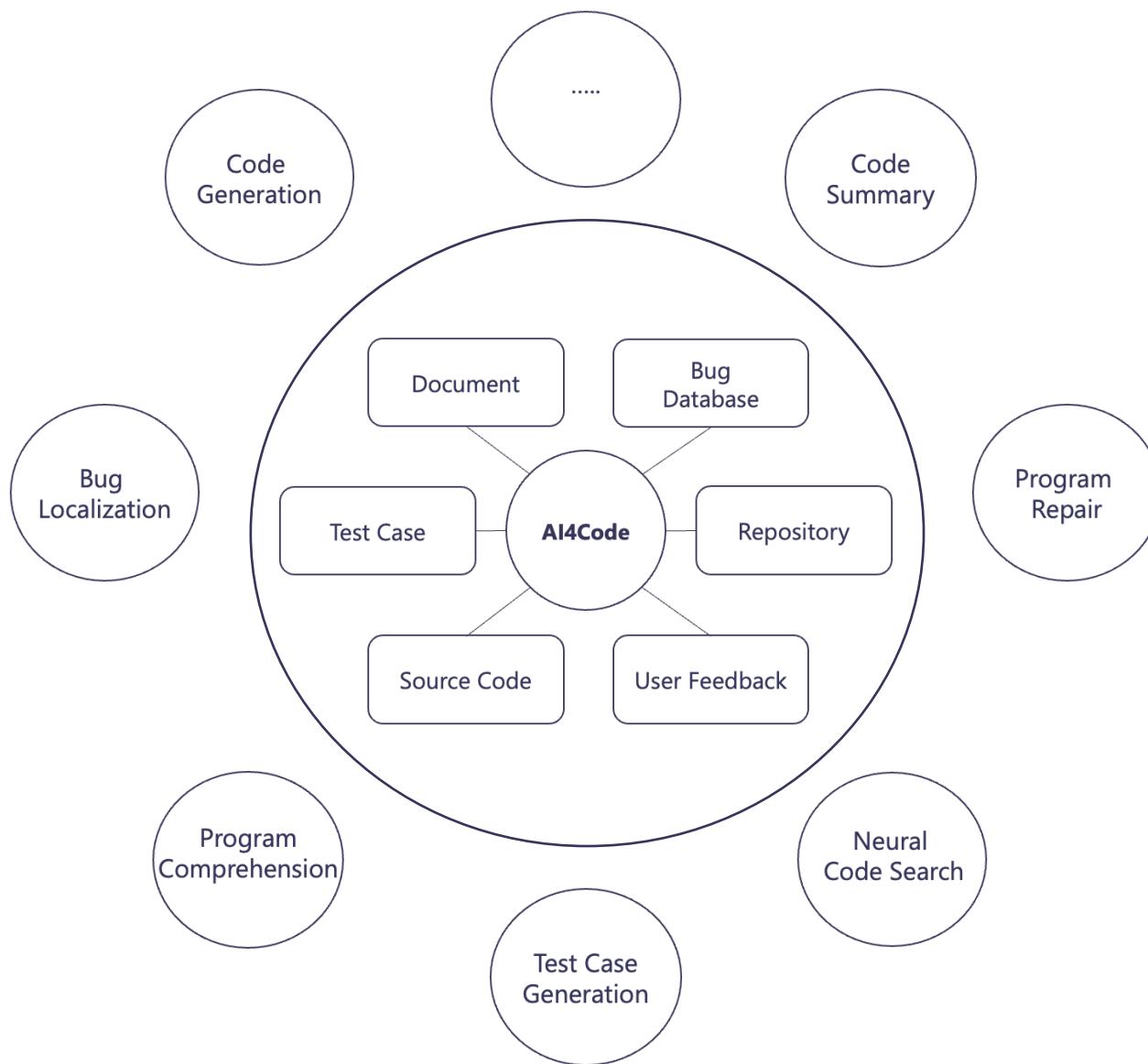
Areas in Software Development



Data Generated from Software Development Process



Building AI Models for Big Code



HOW WOULD THAT IMPACT SOFTWARE DEVELOPMENT?



PRODUCTIVITY



PERFORMANCE



RELIABILITY

Software's Productivity

Jun 20, 2019, 10:22pm EDT | 28,193 views

The Future With Level 5 Autonomous Cars

Well-publicized accidents and missteps show that even at this limited autonomy level, AI systems struggle in application.

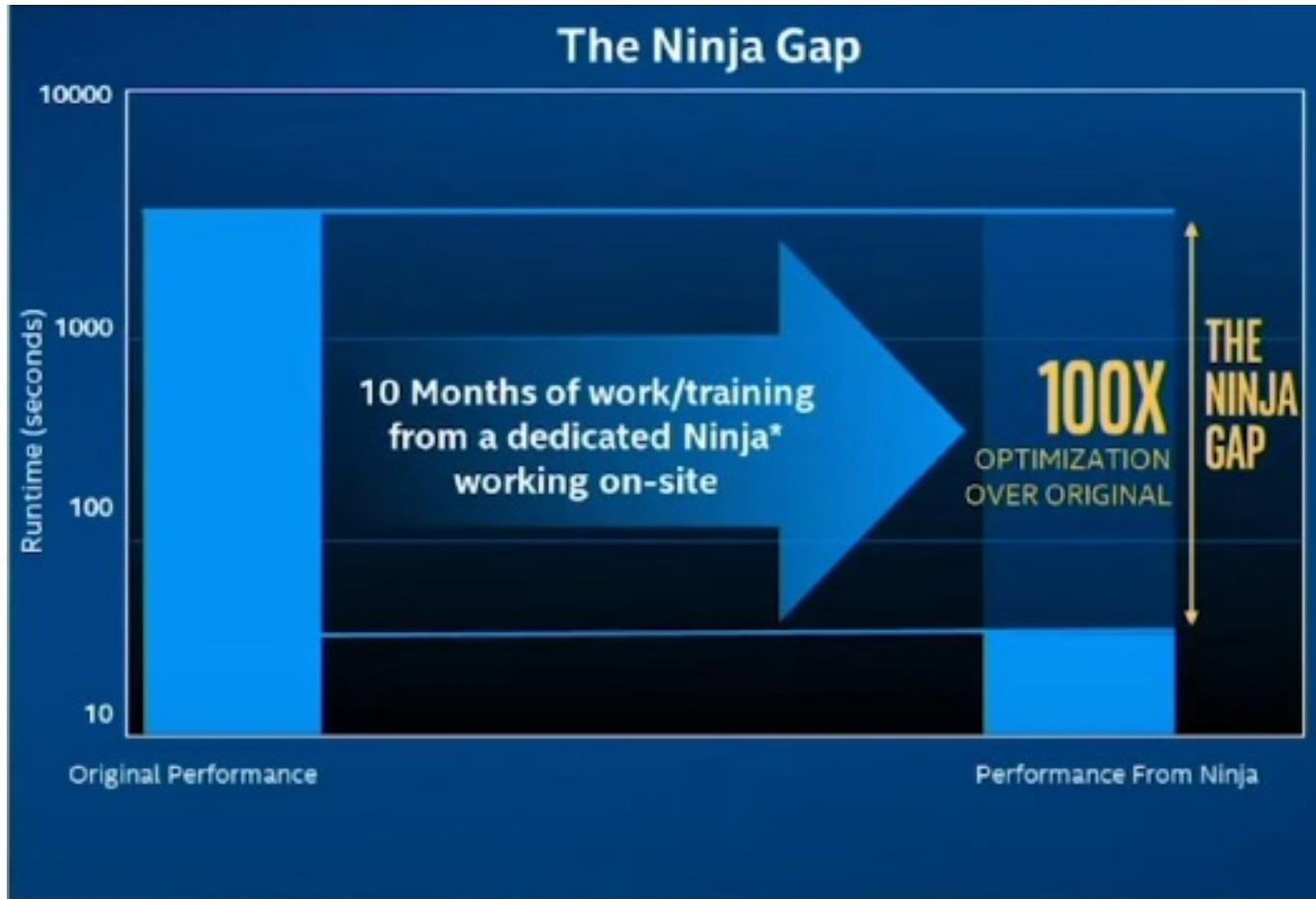


Software

**SOFTWARE IS
LIMITING
AUTONOMOUS
SYSTEMS**



Programmer's Performance



**YOU CODE IS
SUFFERING
FROM THE
NINA GAP**

**Only 10% of programmers
are ninjas**

Ninja Gap: Initiated by Intel, which is the performance difference between an amateur programmer and an expert programmer.

System's Reliability

What Boeing's 737 MAX Has to Do With Cars: Software

Investigators believe faulty software contributed to two fatal crashes. A newly discovered fault will likely keep the 737 MAX grounded until the fall.

CNN BUSINESS Markets Tech Media Success Video



Who's responsible when an autonomous car crashes?

THE VERGE TECH ▾ REVIEWS ▾ SCIENCE ▾ CREATORS ▾ ENTERTAINMENT ▾ VIDEO MORE ▾ f ▾ t ▾ r ▾ u ▾

US & WORLD ▾ TECH ▾ POLITICS ▾

Facebook apologizes after wrong translation sees Palestinian man arrested for posting 'good morning'

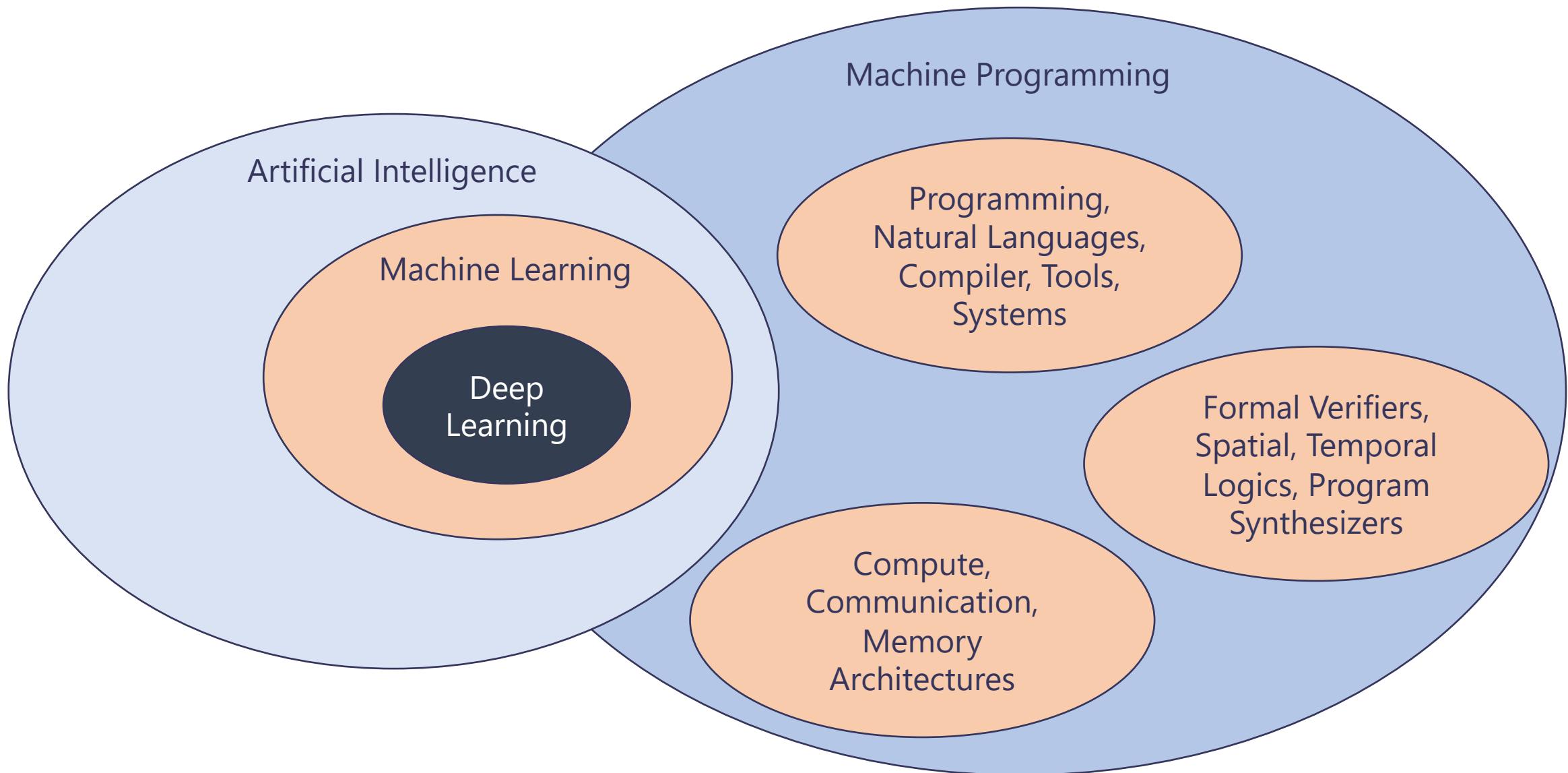
**AUTONOMOUS SYSTEMS
IS SUFFERING FROM
LIFE-THREATENING BUGS**

Machine Programming

- ✓ Humans communicate intention to machine.
- ✓ Machine handle the programming.
Thus, “Machine Programming”

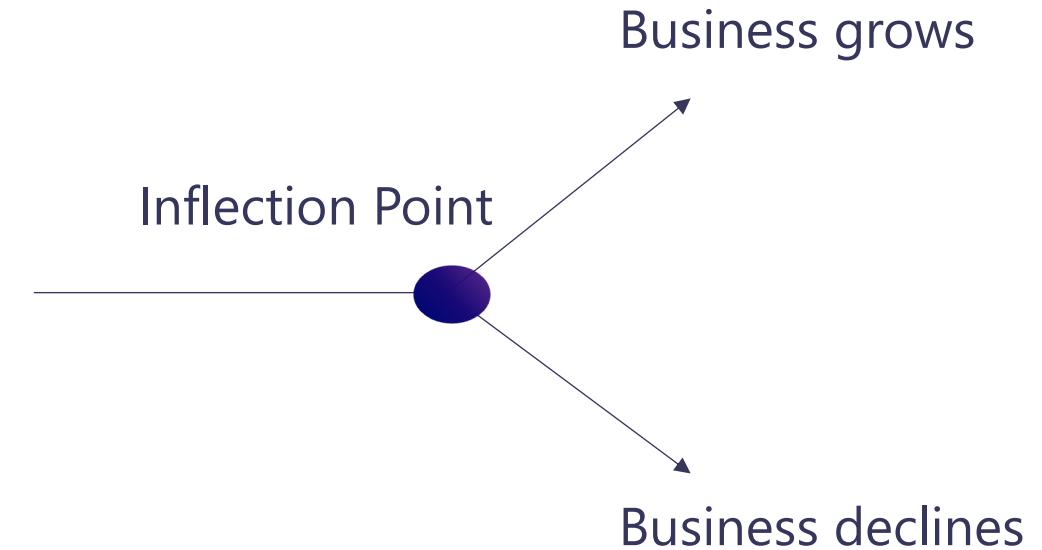
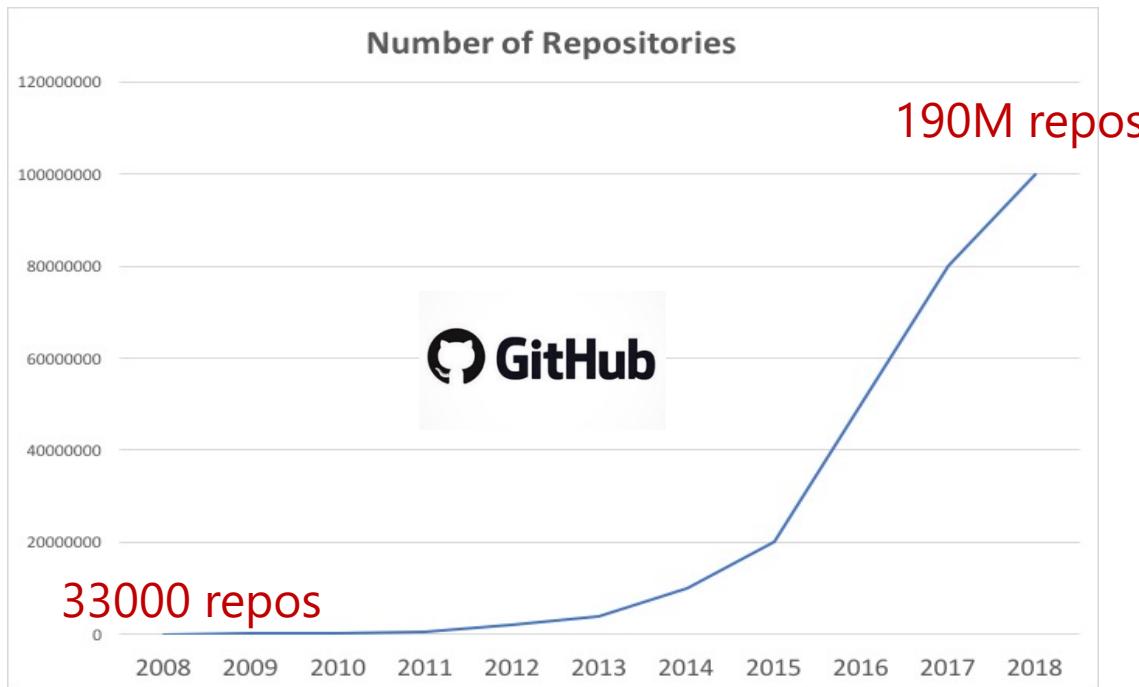
- Increasing the software's **productivity**.
- Closing the **performance** gap between amateur developer and expert developer.
- Enhancing the **reliability** of software system.

A View of Machine Programming



The Inflection Point of Machine Programming

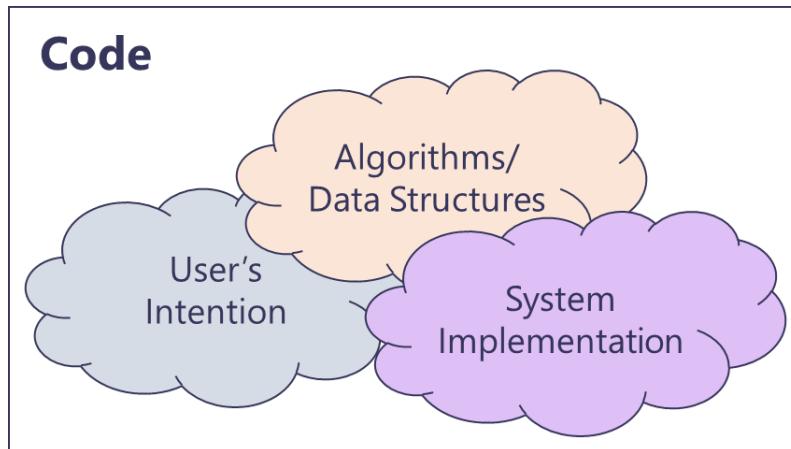
- ML and formal methods
 - Transformer, DeepCoder, BERT, NetSyn, ...
- New & improved hardware
 - CPU, GPU, NNP, TPU, FPGAs,
- Large codebase dataset (**BIG CODE !!**)



**Machine Programming
needs all of this to work**

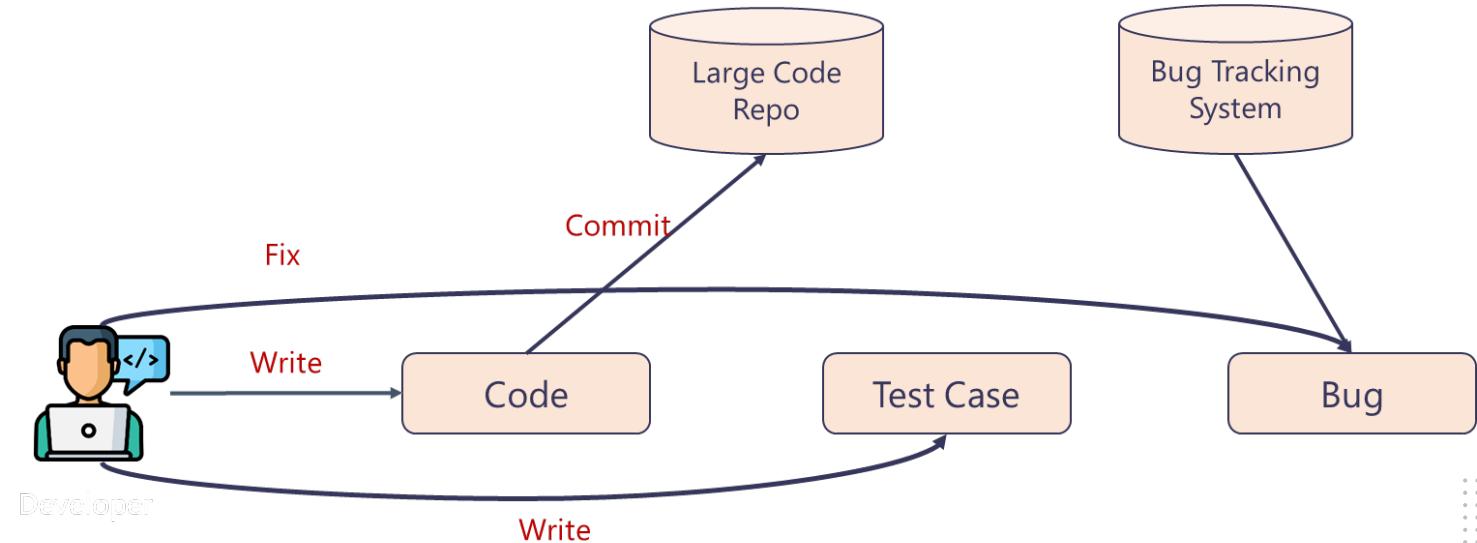
Today's Programming

- Programming is slow and error-prone.
- Programmers struggle to ensure correctness, security, and performance.
- Code constantly evolves to keep pace with an external environment.

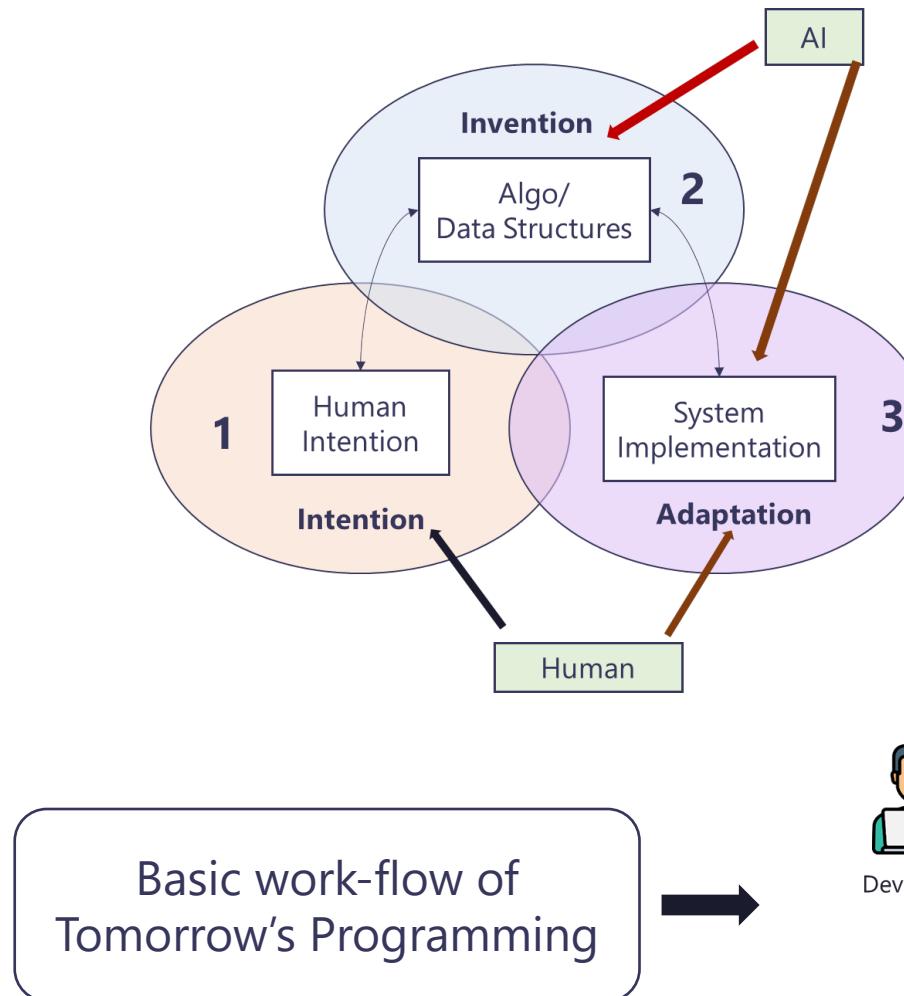


- A program usually contains:
 - The developer's intent
 - Code to express that intent.
 - Specifications for software/hardware adaptation.
- Human developers must handle all of these steps.

Basic work-flow of
Today's Programming

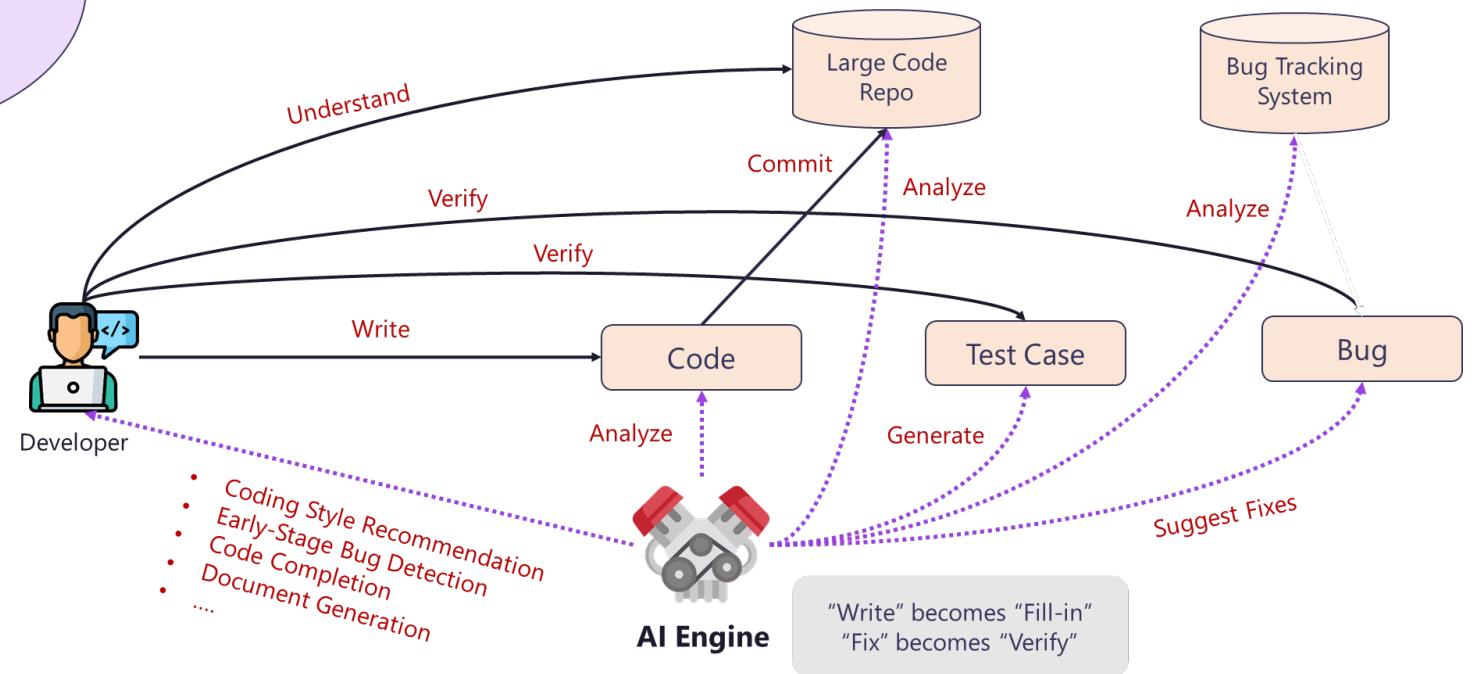


Tomorrow's Programming



- Human only specifies intention:
"Computer, gives me a program X that does Y".

- Next, the system determines what it needs to construct that program.
- May have dialogue interface between human and machine.



How is Machine Programming being implemented in the industry?



copilot.github.com

GitHub Copilot

Technical Preview

Your AI pair programmer

With GitHub Copilot, get suggestions for whole lines or entire functions right inside your editor.

Sign up >

sentiment.ts write.sql.go parse.expenses.py addresses.rb

```
1 #!/usr/bin/env ts-node
2
3 import { fetch } from "fetch-h2";
4
5 const response = await fetch("https://www.example.com");
6
7 console.log(response);
8
9 // Output: Response with status 200 OK and headers {"Content-Type": "text/html; charset=UTF-8"}.
```

Amazon CodeGuru

Automate code reviews and optimize application performance with ML-powered recommendations

[Get started with Amazon CodeGuru](#)

FEATURED
AWS Certification
Explore the resources available to help you prepare for your AWS Certification.

Find your most expensive lines of code and improve code quality

Jun 4, 2021, 10:25am EDT | 26,279 views

IBM CodeNet: Artificial Intelligence That Can Program Computers And Solve A \$100 Billion Legacy Code Problem

Paul Smith-Goodson Contributor
Moor Insights and Strategy Contributor Group Ⓛ
Cloud
Analyst-in-residence, Quantum Computing

Intel Machine Programming Tool Detects Bugs in Code



Today, Intel unveiled ControlFlag – a machine programming research system that can autonomously detect errors in code.

News

- December 3, 2020
- Contact Intel PR

[More New Technologies News →](#)

» Watch video: "Intel Labs Day 2020: Justin Gottschlich Second Session from Intel PR on Vimeo"

What's New: Today, Intel unveiled ControlFlag – a machine programming research system that can autonomously detect errors in code. Even in its infancy, this novel, self-supervised system shows promise as a powerful productivity tool to assist software developers with the labor-intensive task of debugging. In preliminary tests, ControlFlag trained and learned novel defects on over 1 billion unlabeled lines of production-quality code.

intel newsroom Top News Sections ▾ News By Category ▾ All News ▾ Search News ➔ ai.facebook.com/blog/aroma-ml-for-code-recommendation/ FACEBOOK AI Research Publications

News Byte July 29, 2020 Contact Intel PR

Intel, MIT and Georgia Tech Deliver Improved Machine-Programming Code Similarity System

What's New: Today, Intel unveiled a new machine programming (MP) system – in conjunction with Massachusetts Institute of Technology (MIT) and Georgia Institute of Technology (Georgia Tech). The system, machine inferred code similarity (MISIM), is an automated engine designed to learn what a piece of software intends to do by studying the structure of the code and analyzing syntactic differences of other code with similar behavior.

AI gives software development tools a boost

GitHub Copilot, DeepDev, IntelliCode, and other code-focused applications of machine learning can help us deliver better code, faster.

[f](#) [t](#) [in](#) [s](#) [e](#) [d](#)

POSTED ON NOVEMBER 6, 2018 TO DEVINFRA, ML APPLICATIONS

Getafix: How Facebook tools learn to fix bugs automatically

Microsoft | Visual Studio Products ▾ Downloads Buy ▾ Support

Visual Studio IntelliCode

AI-assisted development

[Sign up for news and updates >](#)

trusted-programming.github.io

Trustworthy Programming lab, Huawei Research Ireland

We are leading the smooth transition towards more trustworthy software engineering through innovative R&D in programming theory, e.g., lambda calculus-based functional programming, deep code learning-based intelligent software engineering, programming technology and methods, e.g., model-driven software development (MDSD), domain-specific languages (DSL), and programming tools, e.g., code transpiler, deep code learner, bug localizer, etc.

Blog posts [RSS](#)

- Rustconf Reflections
- Explainer of Deep Code Learning
- Our Rust Mission at Huawei

Machine Programming Principles and Applications



Programming Language vs Natural Language

Executability: Programming language is executable, natural language is not

Formality: programming languages are designed top-down by a few designers for many users. Natural languages, in contrast, emerge, bottom up, "through social dynamics".

Cross-Channel Interaction: Code's two channels, the algorithmic and the explanatory channels, interact through their semantic units.



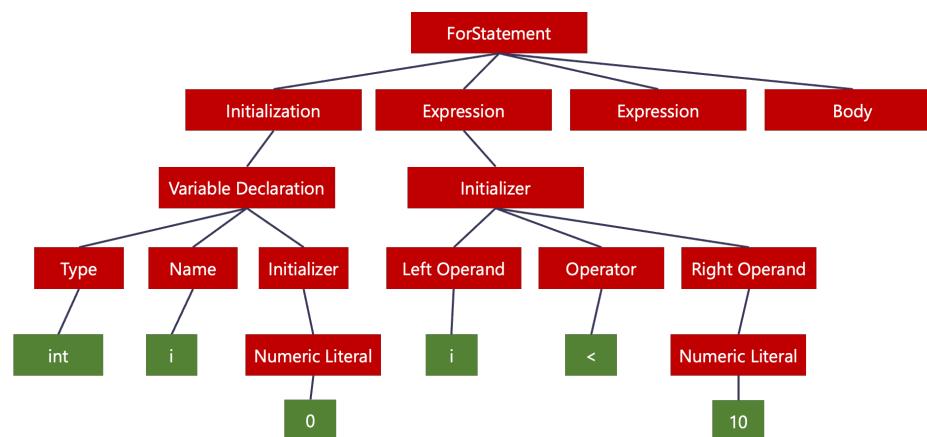
NLP techniques are applicable for programming language, but needs additional contexts from Compiler/Formal Language theory.

Code Representations

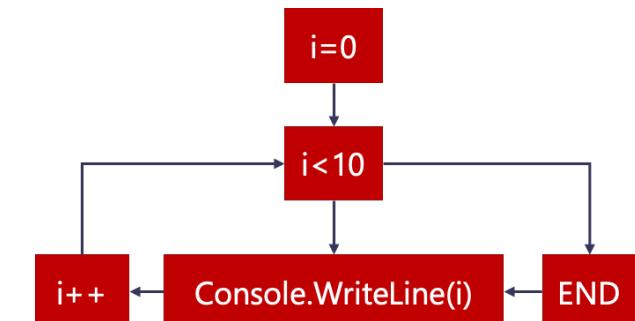
Token-based

```
for (int i = 0; i < 10; i++){\n    Console.WriteLine(i);\n}
```

Syntax-based (Tree)



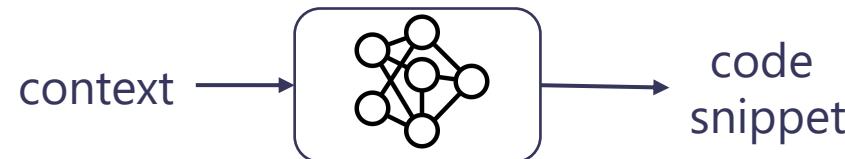
Graph-based



Probabilistic Models of Code

Code Generating Models

- Probability distribution over code, e.g. tokens or AST nodes.
- Model how the code is written

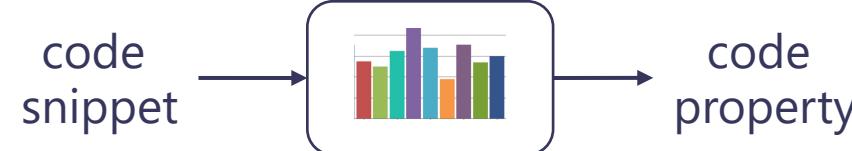


$$P_{\mathcal{D}}(\mathbb{C}|\mathcal{C}(\mathbb{C}))$$

- Context \mathcal{C} (\mathbb{C})
- Code representation \mathbb{C}

Representational Models of Code

- Yield conditional probability distribution over code element property given the code input
- Be able to predict the property

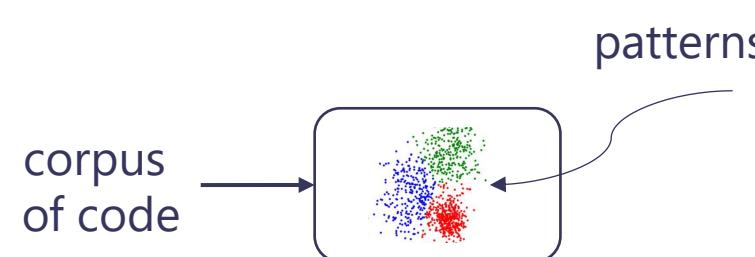


$$P_{\mathcal{D}}(\pi|f(\mathbb{C}))$$

- Code property π
- f to convert code snippet \mathbb{C} into code representation

Pattern Mining of Code

- Infer, without supervision, a likely latent structure within code corpus.
- Find reusable and human-interpretable patterns.



$$P_{\mathcal{D}}(f(\mathbb{C})) = \sum_l P_{\mathcal{D}}(g(\mathbb{C})|l)P(l)$$

- View of code $g(\mathbb{C})$
- l : set of latent variables to infer

Distributed Representation of Code



Suggesting Accurate Method and Class Names

Miltiadis Allamanis[†]

Earl T. Barr[‡]

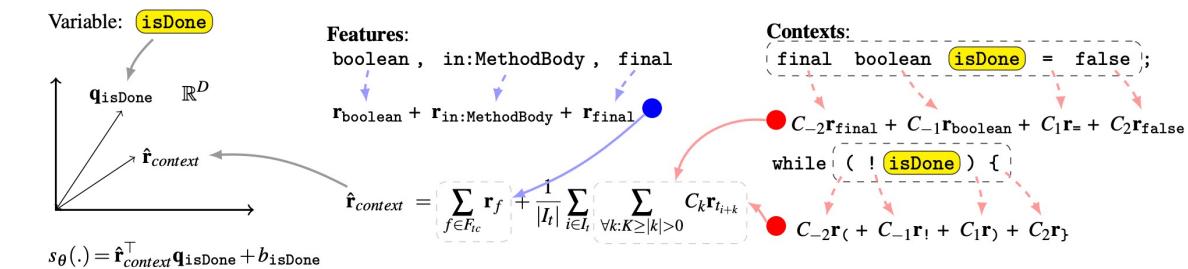
Christian Bird^{*}

Charles Sutton[†]

[†]School of Informatics
University of Edinburgh
Edinburgh, EH8 9AB, UK
{m.allamanis, csutton}@ed.ac.uk

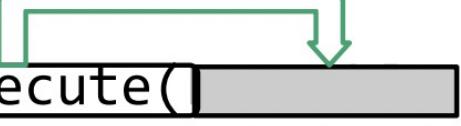
[‡]Dept. of Computer Science
University College London
London, UK
e.barr@ucl.ac.uk

^{*}Microsoft Research
Microsoft
Redmond, WA, USA
cbird@microsoft.com



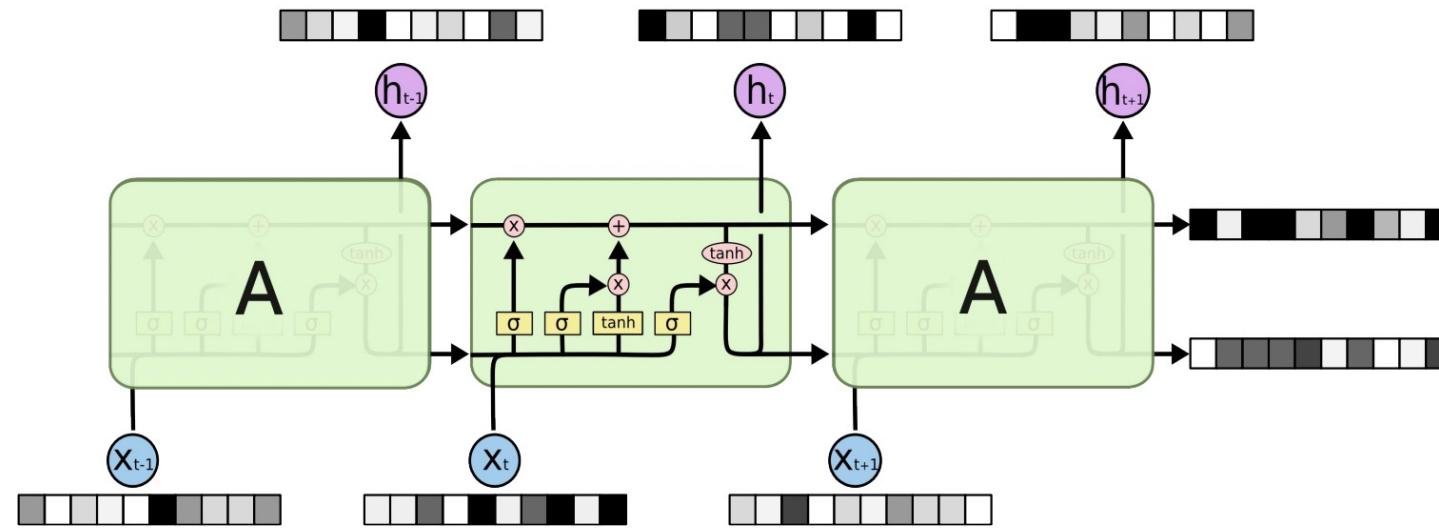
Log-Bilinear model to learn distributed representation of code

N-gram Language Models of Code

```
public void execute( task) {  
    if (task == null)  
        throw new NullPointerException();  
    ForkJoinTask<?> job;  
    if (task instanceof ForkJoinTask<?>) // avoid re-wrap  
        job = (ForkJoinTask<?>) task;  
    else  
        job = new ForkJoinTask.AdaptedRunnableAction(task);  
    externalPush(job);  
}
```

$$P(t_0 \dots t_M) = \prod_{m=0}^M P(t_m | t_{m-1} \dots t_{m-n+1})$$

Neural Language Models of Code



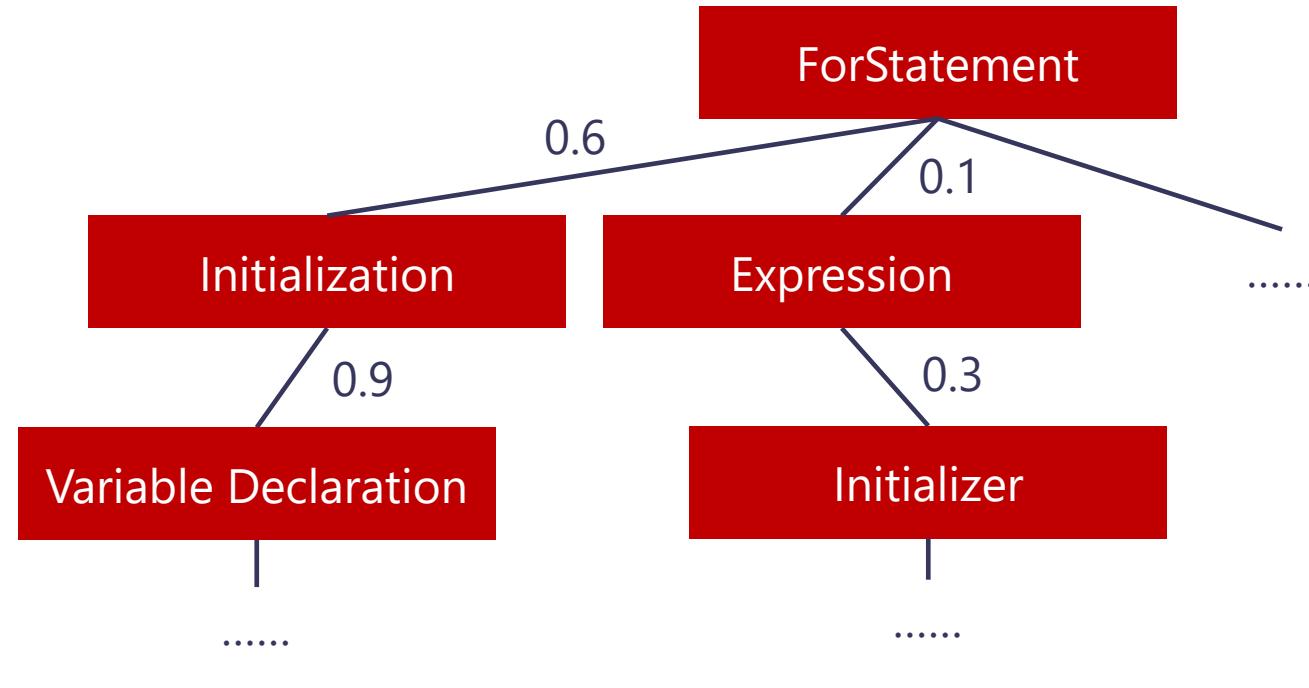
```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

Probabilistic Context Free Grammar for Code Generation

Assign Probability to Rules in the Context Free Grammar

- $E \rightarrow E + E$ (prob=0.7)
- $E \rightarrow T$ (prob=0.3)
- $F \rightarrow (E)$ (prob=0.1)
- $T \rightarrow F * F$ (prob=0.6)
- $F \rightarrow F$ (prob=0.1)
- $F \rightarrow id$ (prob=0.9)
-

Structured Generative Models of Natural Source Code



Structural Language Models of Code

Chris J. Maddison[†]

University of Toronto

Daniel Tarlow

Microsoft Research

CMADDIS@CS.TORONTO.EDU

DTARLOW@MICROSOFT.COM

Uri Alon¹ Roy Sadaka¹ Omer Levy^{2,3} Eran Yahav¹

¹Technion, Israel ²Tel Aviv University

³Facebook AI Research. Correspondence to: Uri Alon <urialon@cs.technion.ac.il>, Roy Sadaka <roysadaka@gmail.com>, Omer Levy <omerlevy@gmail.com>, Eran Yahav <yahave@cs.technion.ac.il>.

Proceedings of the 37th International Conference on Machine Learning, Online, PMLR 119, 2020. Copyright 2020 by the author(s).

Learning to Fix Bug

Bug detection tools are dominated by Program Analysis

sonarqube.org/developer-edition/?gads_campaign=Asia-Code&gads_ad_group=Analysis

WEBINAR Next week: Join The SonarSource Team for the US 2021



Product ▾ What's New Documentation Community

Developer EDITION

Enterprise EDITION

**Built for Developers
By Developers**



Innovative features to systematically track and improve
Code Quality and Code Security in your applications

← → ⌛ 🔒 fbinfer.com

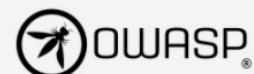
Infer

Docs Support Blog Twitter Facebook GitHub

Search

A tool to detect bugs in Java and C/C++/Objective-C code before it ships

Infer is a static analysis tool - if you give Infer some Java or C/C++/Objective-C code it produces a list of potential bugs. Anyone can use Infer to intercept critical bugs before they have shipped to users, and



PROJECTS CHAPTERS EVENTS ABOUT

Search OWAS

Source Code Analysis Tools

Contributor(s): Dave Wichers, itamarlavender, will-obrien, Eitan Worcel, Prabhu Subramanian, kingthorin, coadaflorin, hblankenship, GovorovViva64, pfhorman, GouveaHeitor, Clint Gibler, DSotnikov, Ajin Abraham, Noam Rathaus, Mike Jang

[Source code analysis](#) tools, also known as Static Application Security Testing (SAST) Tools, can help analyze source code or compiled versions of code to help find security flaws.

SAST tools can be added into your IDE. Such tools can help you detect issues during software development. SAST tool feedback can save time and effort, especially when compared to finding vulnerabilities later in the development cycle.

Learning to Fix Bug (cont.)

Some semantic bugs cannot be fixed by program analysis

```
module.exports = function (grunt) {
  grunt.initConfig({
    execute: {...}, copy: {...}, checktextdomain: {...}
    wp_readme_to_markdown: {...}, makepot: {...})
    ...
  grunt.registerTask('default', ['wp_readme_to_markdown',
    'makepot', 'execute', 'checktextdomain'])
};
```

Missing "copy"

```
function clearEmployeeListOnLinkClick(){
  document.querySelector("a").addEventListener("click",
    function(event){
      document.querySelector("ul").innerHTML = "";
    }
  );
}
```

innerHTML

Method 1.

```
1   public List<String> getSkewedColumnNames(String alias) {
2     ...
3     else {
4       ...
5       TypeInfo typeInfo = TypeInfoUtils.getTypeInfoFromObjectInspector(this.metaData.
6       getTableForAlias(tabAlias).getDeserializer().getObjectInspector());
7       desc = new ExprNodeConstantDesc(typeInfo.getStructFieldTypeInfo(colName), null);
8     }
9     ...
10 }
```

Method 2.

```
1   public TypeInfo getStructFieldTypeInfo(String field) {
2     String fieldLowerCase = field.toLowerCase();
3     for(int i=0; i<allStructFieldNames.size(); i++) {
4       if (field.equals(allStructFieldNames.get(i))) {
5         return allStructFieldTypeInfos.get(i);
6       }
7     }
8     throw new RuntimeException("cannot_find_field_" + field + "(lowercase_form:_"
9     + fieldLowerCase + ")_in_" + allStructFieldNames);
10 }
```

Method 3.

```
1   public Table getTableForAlias(String alias) {
2   -   return this.aliasToTable.get(alias);
3   +   return this.aliasToTable.get(alias.toLowerCase());
4   }
```

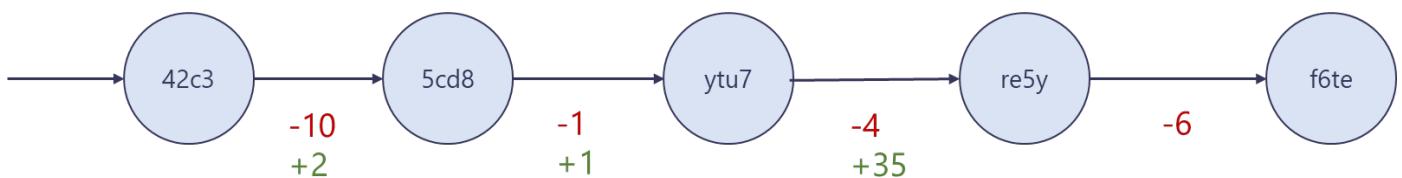
Convert "alias" to lower case

Semantic Bugs in Javascript

Semantic Bugs in Java

Learning to Fix Bug from Commit History

AI can help to fix semantic bugs by learning from commit history



mysql / mysql-server Public

Code Pull requests Actions Projects Security Insights

Commits on Aug 20, 2018

- Bug#27788907 SOME FILE OPERATIONS IN MF_IOCACHE2.C ARE NOT INSTRUMENTED ... marcalff committed on Aug 20, 2018

Commits on Aug 19, 2018

- Bug #26791931: INCORRECT BEHAVIOR IN ALTER TABLE REORGANIZE ... Sreeharsha Ramanavarapu committed on Aug 19, 2018

Commits on Aug 3, 2018

- BUG#28144933 - MYSQL-SERVER RPM DOES NOT INSTALL PERL-DATA-DUMPER AS... bkandas committed on Aug 3, 2018

Commits on Jul 27, 2018

- Merge branch 'mysql-5.5.61-release' into mysql-5.5 hramilson committed on Jul 27, 2018

Code diff between 195 and 196:

```

@@ -195,12 +195,12 @@ public long longValue() {
195   195     */
196   196     @Override
197   197     public float floatValue() {
198 -       @SuppressWarnings("cast")
199 -       float fValue = (float) (value & UNSIGNED_M
200 -       if (value < 0) {
201 -           fValue += 0x1.0p63f;
202 +       if (value >= 0) {
203 +           return (float) value;
204     }
205     -       return fValue;
206     +       // The top bit is set, which means that th
207     +       // So we can ignore the bottom 8, except fi
208     +       return (float) ((value >> 1) | (value & 1
209   }

```

Code diff between 105 and 106:

```

@@ -105,8 +105,8 @@
105   105
106   106     private transient BiEntry<K, V>[] hashTableKToV;
107   107     private transient BiEntry<K, V>[] hashTableVToK;
108 -   -     private transient @Nullable BiEntry<K, V> firstInKeyInsertionOrder;
109 -   -     private transient @Nullable BiEntry<K, V> lastInKeyInsertionOrder;
108 +   +     @Weak private transient @Nullable BiEntry<K, V> firstInKeyInsertionOrder
109 +   +     @Weak private transient @Nullable BiEntry<K, V> lastInKeyInsertionOrder
110   110     private transient int size;
111   111     private transient int mask;
112   112     private transient int modCount;

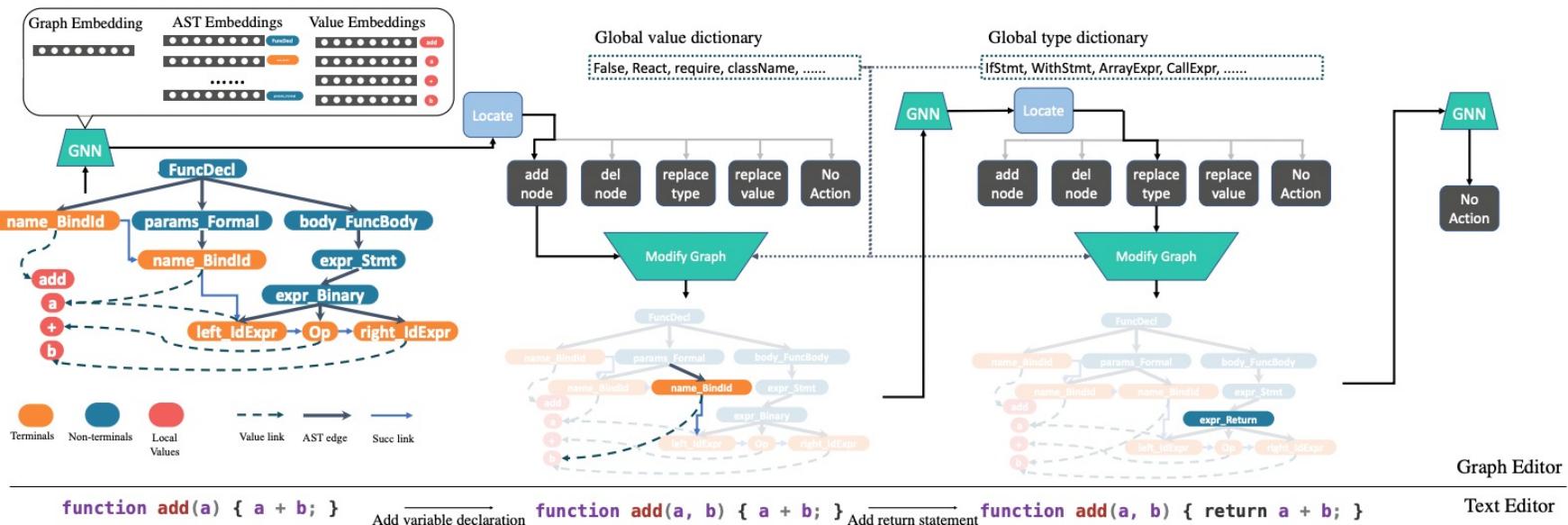
```

The program was first piloted on October 23 in Orava and Bardejov two regions with the highest number of Covid-19 cases in the country. Nearly 141,000 people, 91% of those who were eligible, got tested in the two regions over the three days of the pilot. In the rest of the country, the testing took place simultaneously on Saturday and Sunday. The campaign encouraged everyone older than 10 to take part in the test. People older than 65 years who spend most of their time at home, cancer patients, immunocompromised people, and disabled groups were exempt. Europe suffers record case numbers as France & Macron warns critics may last until summer. Europe suffers record case numbers as France & Macron warns critics may last until summer. Europe suffers record case numbers as France & Macron warns critics may last until summer.



Imagine this is similar to using AI for grammar error correction in Natural Language Processing

Learning to Fix Bug from Commit History (cont.)



Published as a conference paper at ICLR 2020

HOPPITY: LEARNING GRAPH TRANSFORMATIONS TO DETECT AND FIX BUGS IN PROGRAMS

Elizabeth Dinella*
University of Pennsylvania

Hanjun Dai*
Google Brain

Ziyang Li
University of Pennsylvania

Mayur Naik
University of Pennsylvania

Le Song
Georgia Tech

Ke Wang
Visa Research

ABSTRACT

We present a learning-based approach to detect and fix a broad range of bugs in Javascript programs. We frame the problem in terms of learning a sequence of graph transformations: given a buggy program modeled by a graph structure, our model makes a sequence of predictions including the position of bug nodes and

Getafix: Learning to Fix Bugs Automatically

Johannes Bader
Facebook
jobader@fb.com

Andrew Scott
Facebook
andrewscott@fb.com

Michael Pradel
Facebook
michael@binaervarianz.de

Satish Chandra
Facebook
satch@fb.com

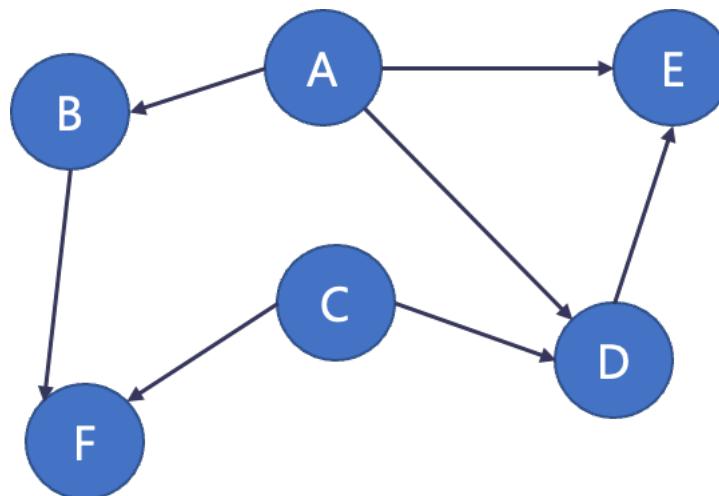
Perform steps of program transformation to convert buggy program to non-buggy program

$$p(g_{fix}|g_{bug}; \theta) = p(g_1|g_{bug}; \theta)p(g_2|g_1; \theta) \dots p(g_{fix}|g_{T-1}; \theta)$$

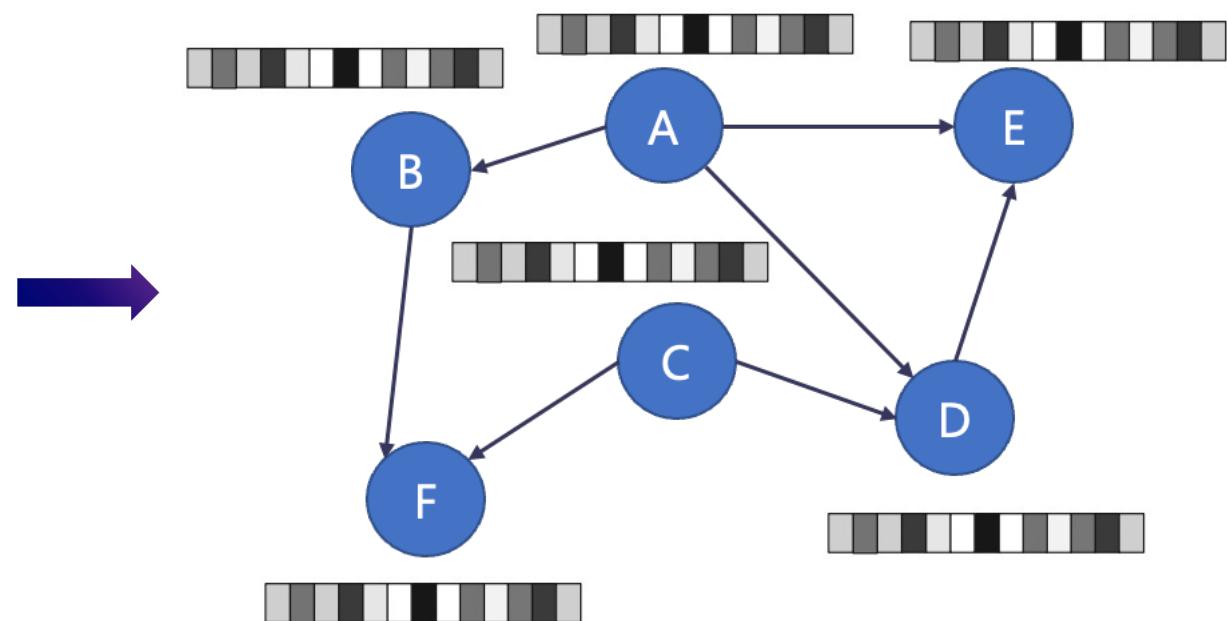
Graph Neural Network for Code Understanding

Graph Notation: $\mathbf{G} = (\mathbf{V}, \mathbf{E})$

- Nodes/Vertices \mathbf{V}
- Edges/Links \mathbf{E}

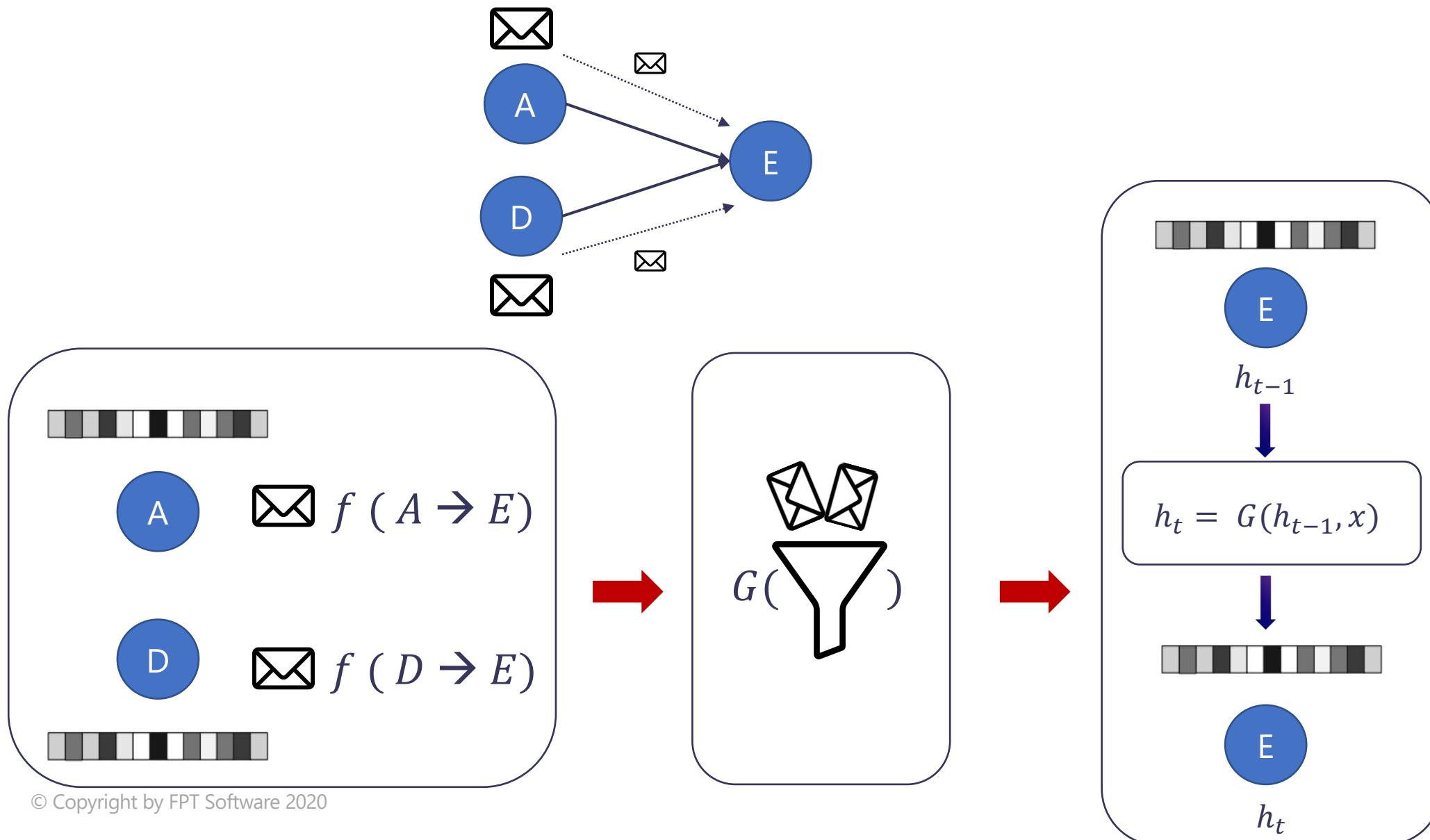


Graph-based Problem

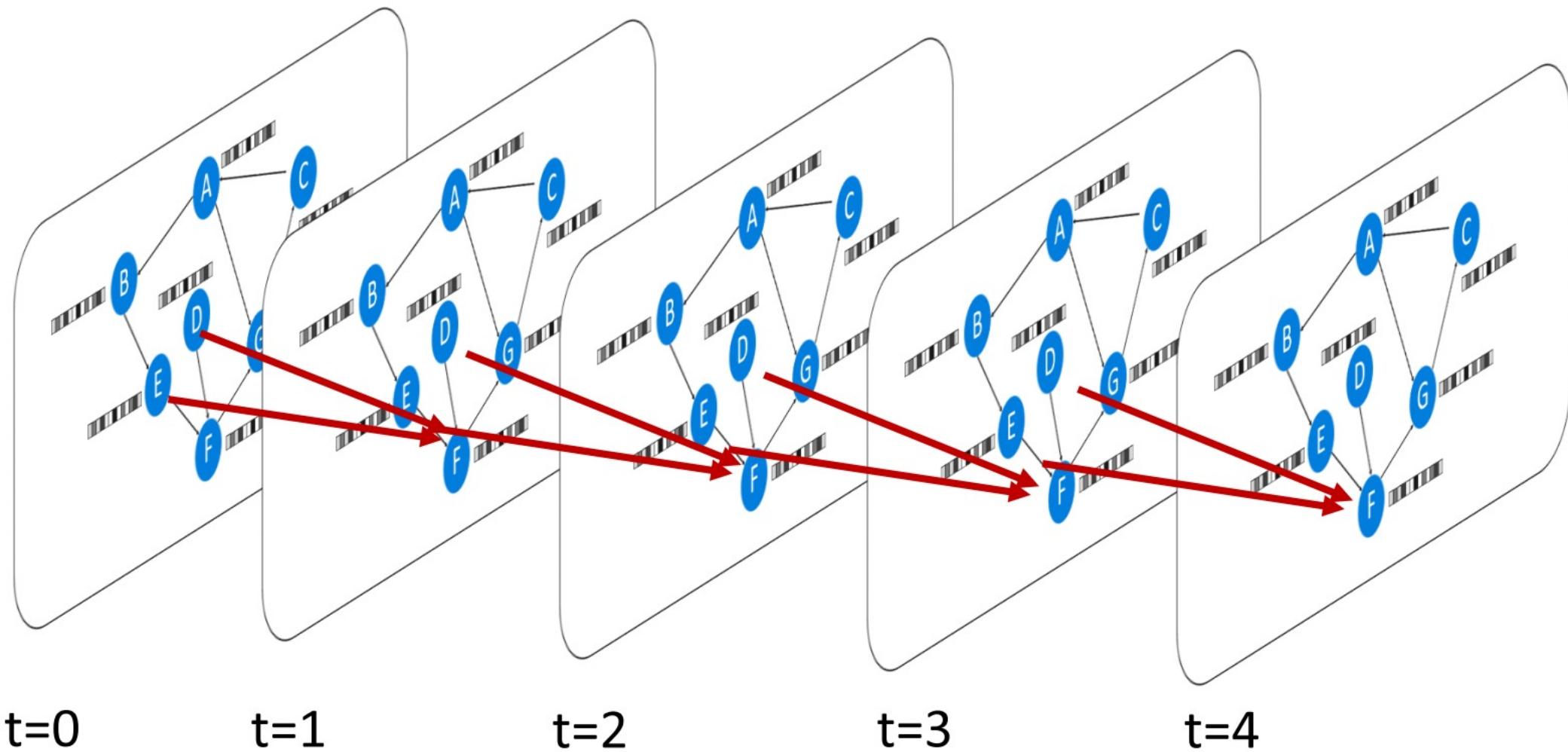


Initial Representation of each node

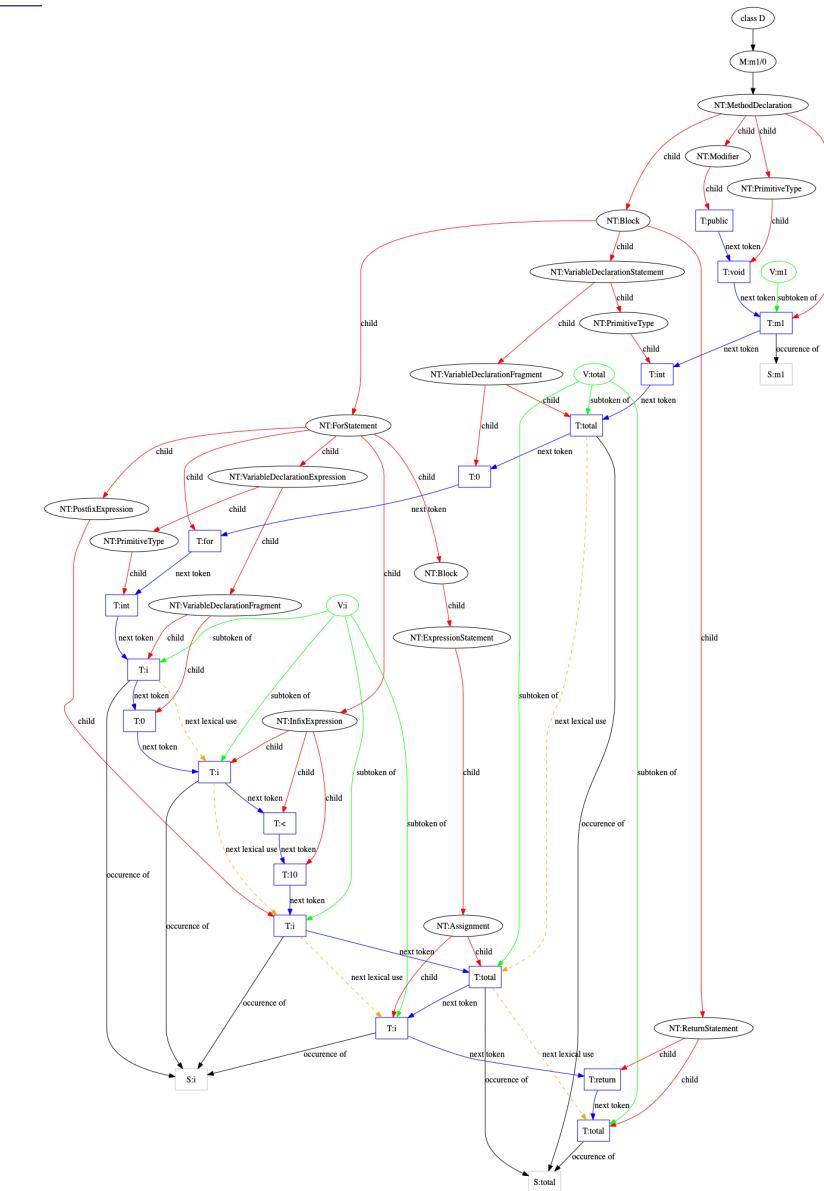
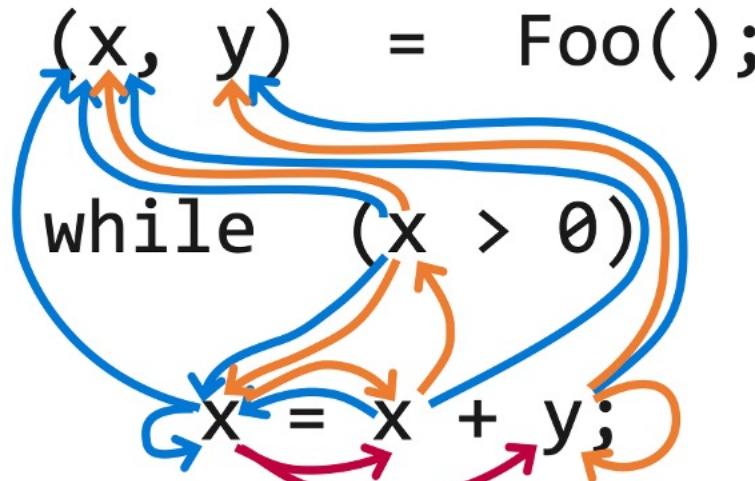
Graph Neural Network: Message Passing



Graph Neural Network: Message Passing (cont.)



Represent a Program as Graph with Data Flow



Graph Neural Network to Fix Variable Mis-used Bugs

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

Real variable mis-used bug found in RavenDB

How to know the variable “**clazz**” is mis-used?

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

Graph Neural Network to Fix Variable Mis-used Bugs

Published as a conference paper at ICLR 2018

```
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);
```



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

Fill-in-the-blank task on program graph

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

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

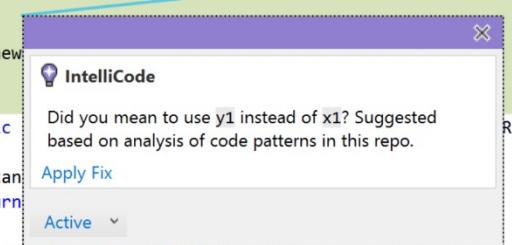
LEARNING TO REPRESENT PROGRAMS WITH GRAPHS

Miltiadis Allamanis
Microsoft Research
Cambridge, UK
miallama@microsoft.com

Marc Brockschmidt
Microsoft Research
Cambridge, UK
mabrocks@microsoft.com

Mahmoud Khademi*
Simon Fraser University

```
public readonly static Thickness MultiLinePadding = new Thickness(0.0, 1  
  
public static IList<Rect> GetRectanglesFromBounds(IList<TextBounds> bounds)  
{  
    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 - y1));  
        }  
    }  
    return newBounds;  
}  
  
public static Rect GetRectFromTextBounds(TextBounds bounds)  
{  
    if (rectangles == null)  
        return null;  
    Rect rect = rectangles[bounds];  
    if (rect != null)  
        return rect;  
    else  
        return GetRectanglesFromBounds(bounds).First();  
}
```



IntelliCode (Microsoft) can recommend mis-used variable using AI

Refine Variable's Name with Probabilistic Graphical Model

Probabilistic Model for Code with Decision Trees

Veselin Raychev
Department of Computer Science
ETH Zürich, Switzerland
veselin.raychev@inf.ethz.ch

Pavol Bielik
Department of Computer Science
ETH Zürich, Switzerland
pavol.bielik@inf.ethz.ch

Martin Vechev
Department of Computer Science
ETH Zürich, Switzerland
martin.vechev@inf.ethz.ch

Veselin Raychev
Department of Computer Science
ETH Zürich
veselin.raychev@inf.ethz.ch

Martin Vechev
Department of Computer Science
ETH Zürich
martin.vechev@inf.ethz.ch

Andreas Krause
Department of Computer Science
ETH Zürich
andreas.krause@inf.ethz.ch

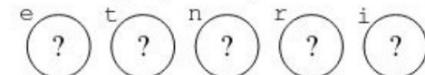


Rename meaningless variable name to meaningful variable name using conditional random field

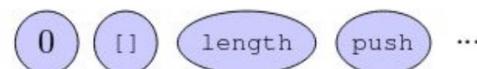
```
function chunkData(e, t) {
    var n = [];
    var r = e.length;
    var i = 0;
    for (; i < r; i += t) {
        if (i + t < r) {
            n.push(e.substring(i, i + t));
        } else {
            n.push(e.substring(i, r));
        }
    }
    return n;
}
```

(a) JavaScript program with minified identifier names

Unknown properties (variable names):



Known properties (constants, APIs):



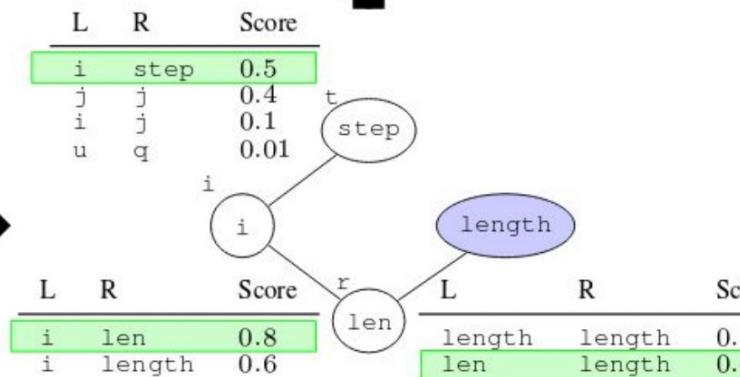
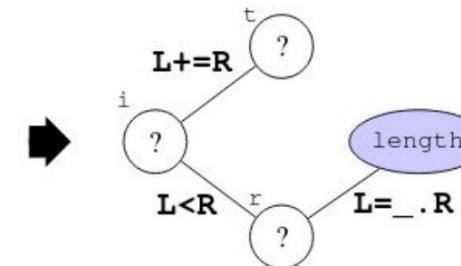
Predicting Program Properties from “Big Code”

Veselin Raychev
Department of Computer Science
ETH Zürich
veselin.raychev@inf.ethz.ch

Martin Vechev
Department of Computer Science
ETH Zürich
martin.vechev@inf.ethz.ch

```
/* str: string, step: number, return: Array */
function chunkData(str, step) {
    var colNames = []; /* colNames: Array */
    var len = str.length;
    var i = 0; /* i: number */
    for (; i < len; i += step) {
        if (i + step < len) {
            colNames.push(str.substring(i, i + step));
        } else {
            colNames.push(str.substring(i, len));
        }
    }
    return colNames;
}
```

(e) JavaScript program with new identifier names and types



Code Summarization

```
/*
A function to | [REDACTED]
*/
void| [REDACTED] () [REDACTED]
    TreeView myTreeView = new TreeView()
    myTreeView.Nodes.Clear();
    foreach (string parentText in xml.parent)
    {
        TreeNode parent = new TreeNode();
        parent.Text = parentText;
        myTreeView.Nodes.Add(treeNodeDivisions);
        foreach (string childText in xml.child)
        {
            TreeNode child = new TreeNode();
            child.Text = childText;
            parent.Nodes.Add(child);
        }
    }
}
```

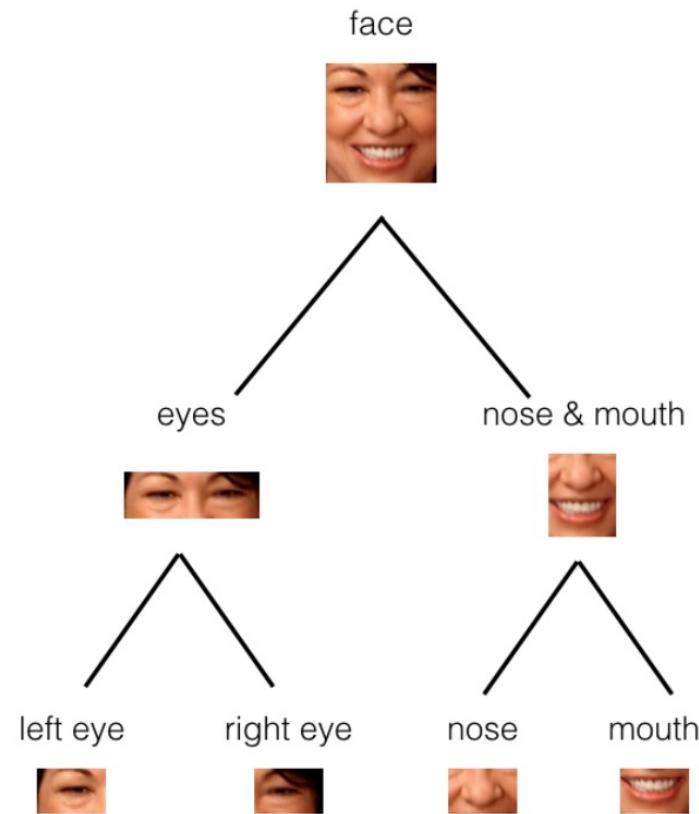
Generate Docstring: **add a child node to a TreeView**

Predict Method's Name: **addChildNode**

Code Summarization includes:

- **Method name prediction**
- **Comment generation**
- **Docstring generation**

Learning to Summarize Code with Capsule Theory



Dynamic Routing Between Capsules

Sara Sabour

Nicholas Frosst

Geoffrey E. Hinton

Google Brain

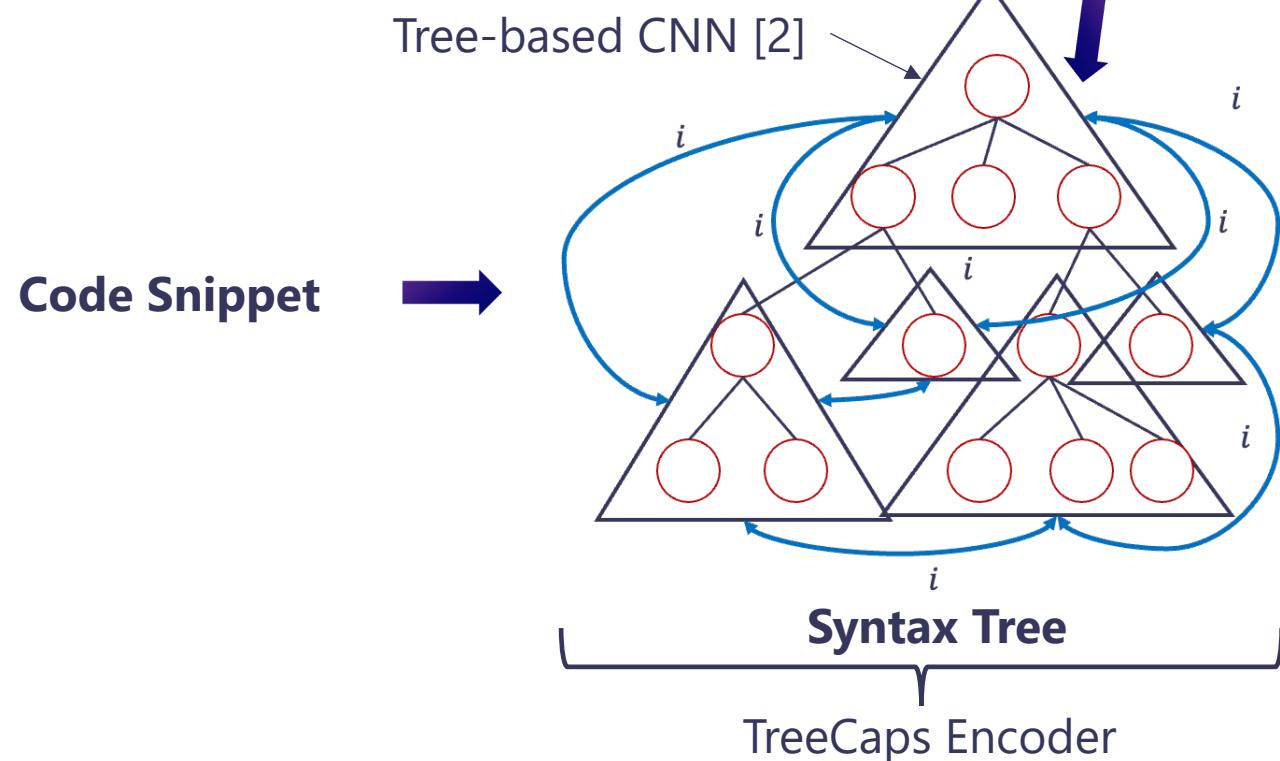
Toronto

{sasabour, frosst, geoffhinton}@google.com

- Combine observations from smaller parts to build up a more complex part of a picture.
- Different parts of the pictures have dependency on each other

Learning to Summarize Code with Capsule Theory (cont.)

Modeling Syntax Tree Using Capsule Theory [1]



TreeCaps: Tree-Based Capsule Networks for Source Code Processing

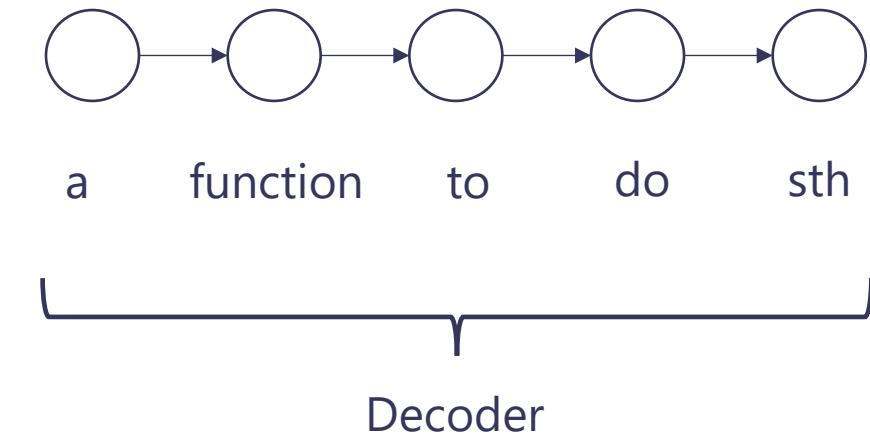
Nghi D. Q. Bui ^{1,3} Yijun Yu ^{1,2} Lingxiao Jiang ³

¹ Trustworthy Open-Source Software Engineering Lab, Huawei Research Centre, Ireland

² School of Computing & Communications, The Open University, UK

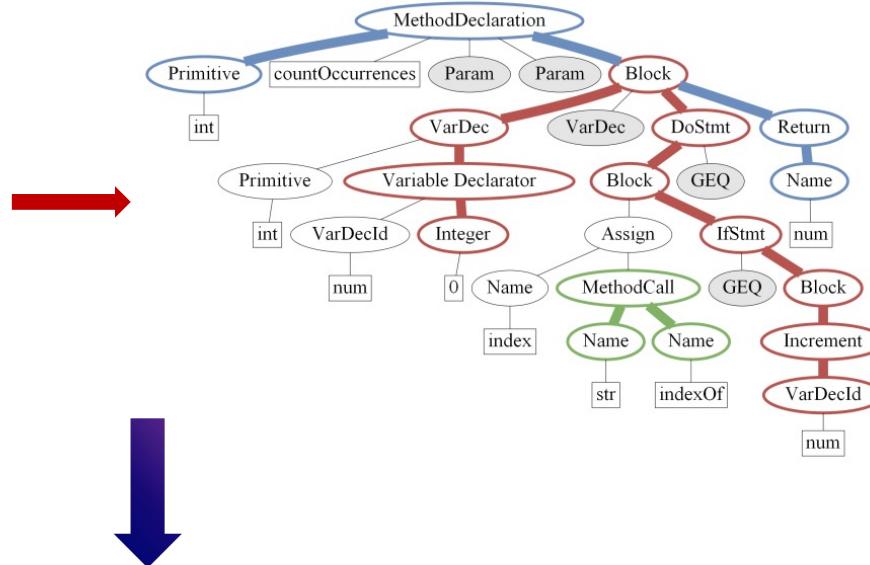
³ School of Computing & Information Systems, Singapore Management University

nghi.bui@huawei.com , y.yu@open.ac.uk, lxjiang@smu.edu.sg



Learning to Summarize Code by Path Abstraction

```
int countOccurrences(String str, char ch) {
    int num = 0;
    int index = -1;
    do {
        index = str.indexOf(ch, index + 1);
        if (index >= 0) {
            num++;
        }
    } while (index >= 0);
    return num;
}
```



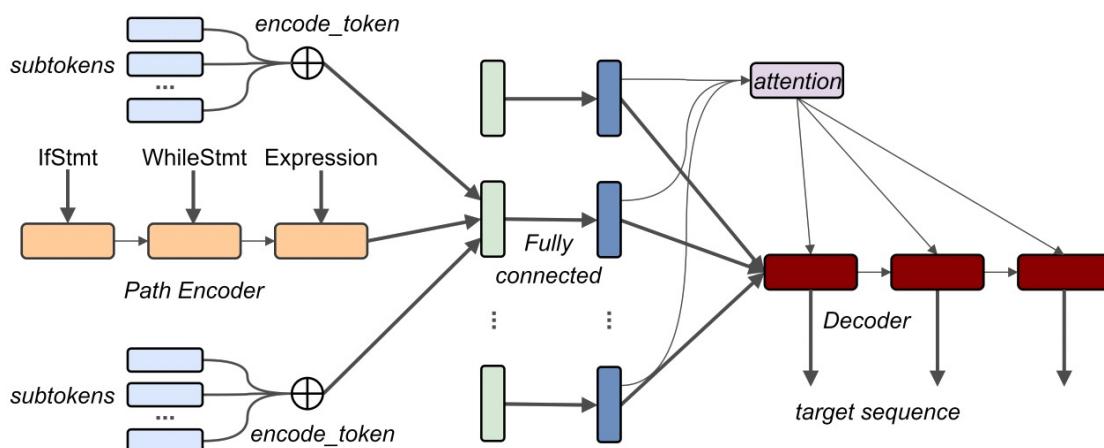
Modeling Paths of the Syntax Tree

code2vec: Learning Distributed Representations of Code

URI ALON, Technion
MEITAL ZILBERSTEIN, Technion
OMER LEVY, Facebook AI Research
ERAN YAHAV, Technion

We present a neural model for representing snippets of code as continuous distributed vectors ("code embeddings").

Published as a conference paper at ICLR 2019



CODE2SEQ: GENERATING SEQUENCES FROM STRUCTURED REPRESENTATIONS OF CODE

Uri Alon
Technion
urialon@cs.technion.ac.il

Shaked Brody
Technion
shakedbr@cs.technion.ac.il

Omer Levy
Facebook AI Research
omerlevy@gmail.com

Eran Yahav
Technion
yahave@cs.technion.ac.il

Program Synthesis: Large Language Models for Code

Program Synthesis with Large Language Models

Jacob Austin* Augustus Odena*

Maxwell Nye† Maarten Bosma Henryk Michalewski David Dohan Ellen Jiang Carrie Cai

Michael Terry Quoc Le Charles Sutton

Google Research

* denotes equal contribution

jaaustin@google.com, augustusodena@google.com

Evaluating Large Language Models Trained on Code

Mark Chen^{*1} Jerry Tworek^{*1} Heewoo Jun^{*1} Qiming Yuan^{*1} Henrique Ponde de Oliveira Pinto^{*1}
Jared Kaplan^{*2} Harri Edwards¹ Yuri Burda¹ Nicholas Joseph² Greg Brockman¹ Alex Ray¹ Raul Puri¹
Gretchen Krueger¹ Michael Petrov¹ Heidy Khlaaf³ Girish Sastry¹ Pamela Mishkin¹ Brooke Chan¹
Scott Gray¹ Nick Ryder¹ Mikhail Pavlov¹ Alethea Power¹ Lukasz Kaiser¹ Mohammad Bavarian¹
Clemens Winter¹ Philippe Tillet¹ Felipe Petroski Such¹ Dave Cummings¹ Matthias Plappert¹
Fotios Chantzis¹ Elizabeth Barnes¹ Ariel Herbert-Voss¹ William Hebgen Guss¹ Alex Nichol¹ Alex Paino¹
Nikolas Tezak¹ Jie Tang¹ Igor Babuschkin¹ Suchir Balaji¹ Shantanu Jain¹ William Saunders¹
Christopher Hesse¹ Andrew N. Carr¹ Jan Leike¹ Josh Achiam¹ Vedant Misra¹ Evan Morikawa¹
Alec Radford¹ Matthew Knight¹ Miles Brundage¹ Mira Murati¹ Katie Mayer¹ Peter Welinder¹
Bob McGrew¹ Dario Amodei² Sam McCandlish² Ilya Sutskever¹ Wojciech Zaremba¹

DEEPCODER: LEARNING TO WRITE PROGRAMS

Matej Balog*
Department of Engineering
University of Cambridge

Alexander L. Gaunt, Marc Brockschmidt,
Sebastian Nowozin, Daniel Tarlow
Microsoft Research

ABSTRACT

We develop a first line of attack for solving programming competition-style problems from input-output examples using deep learning. The approach is to train a neural network to predict properties of the program that generated the outputs from the inputs. We use the neural network's predictions to augment search techniques from the programming languages community, including enumerative search and an SMT-based solver. Empirically, we show that our approach leads to an order of magnitude speedup over the strong non-augmented baselines and a Recurrent Neural Network approach, and that we are able to solve problems of difficulty comparable to the simplest problems on programming competition websites.

© Copyright by FPT Software 2020

prompt

Write a python function to check if a given number is one less than twice its reverse. Your code should satisfy these tests:

```
assert check(70) == False
assert check(23) == False
assert check(73) == True
```

model

```
def check(n):
    if n == 2*int(str(n)[::-1])-1:
        return True
    else:
        return False
```

prompt

Write a function to find the smallest missing element in a sorted array. Your code should satisfy these tests:

```
assert smallest_missing([0, 1, 2, 3, 4, 5, 6], 0, 6) == 7
assert smallest_missing([0, 1, 2, 6, 9, 11, 15], 0, 6) == 3
assert smallest_missing([1, 2, 3, 4, 6, 9, 11, 15], 0, 7) == 0
```

model

```
def smallest_missing(arr, n, m):
    smallest = min(n, m)
    for i in range(n, m + 1):
        if arr[i] <= smallest:
            smallest += 1
    return smallest
```

prompt

Please, solve the mathematical problem: a and b start walking towards each other at 4pm at a speed of 2 kmph and 3 kmph. They were initially 15 km apart. At what time do they meet? $n_0 = 4.0$, $n_1 = 2.0$, $n_3 = 15.0$.

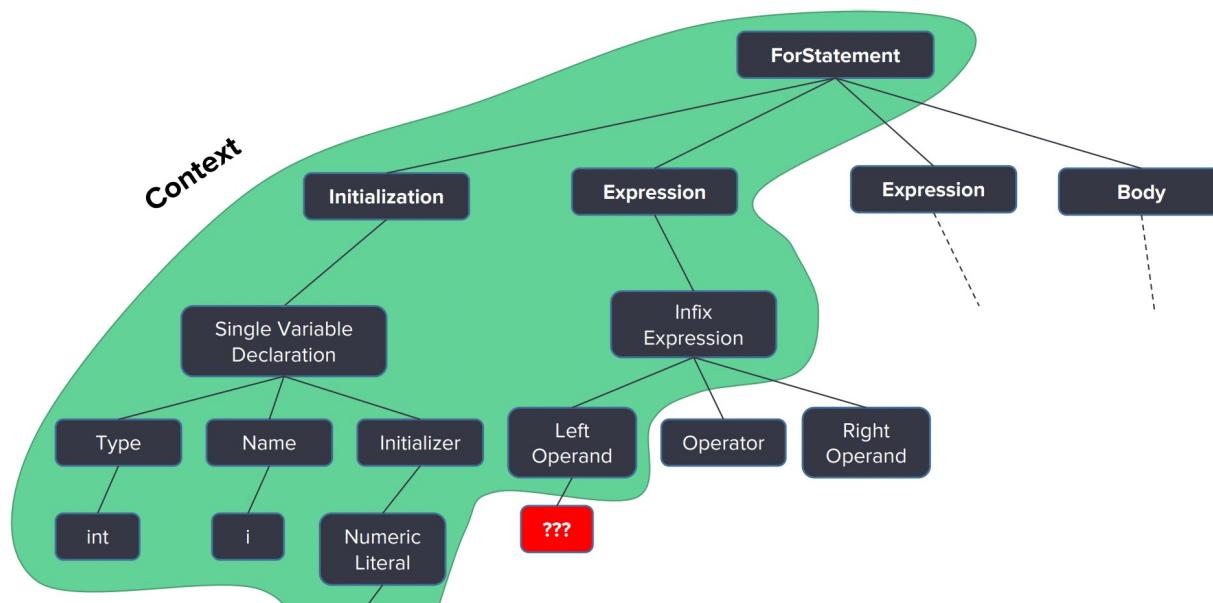
model

```
n0 = 4.0
n1 = 2.0
n2 = 3.0
n3 = 15.0
t0 = n1 + n2
t1 = n3 / t0
answer = n0 + t1
```

How Machine Programming Can Bring Benefits to Fsoft?



Code Generation



```
private static void log(String value) {  
    if (value != null  
        && ???)  
        value = value.substring(0, 55)+"...";  
    LOG.info(value);  
}
```

True ref:	value.length() > 55	
SLM top-5:	(9.6%) value.length() > 0	Ճ
	(7.3%) value.length() > 55	✓
	(1.8%) value.startsWith("...")	
	(1.5%) !value.startsWith("...")	
	(0.9%) value.charAt(0) == '.'	

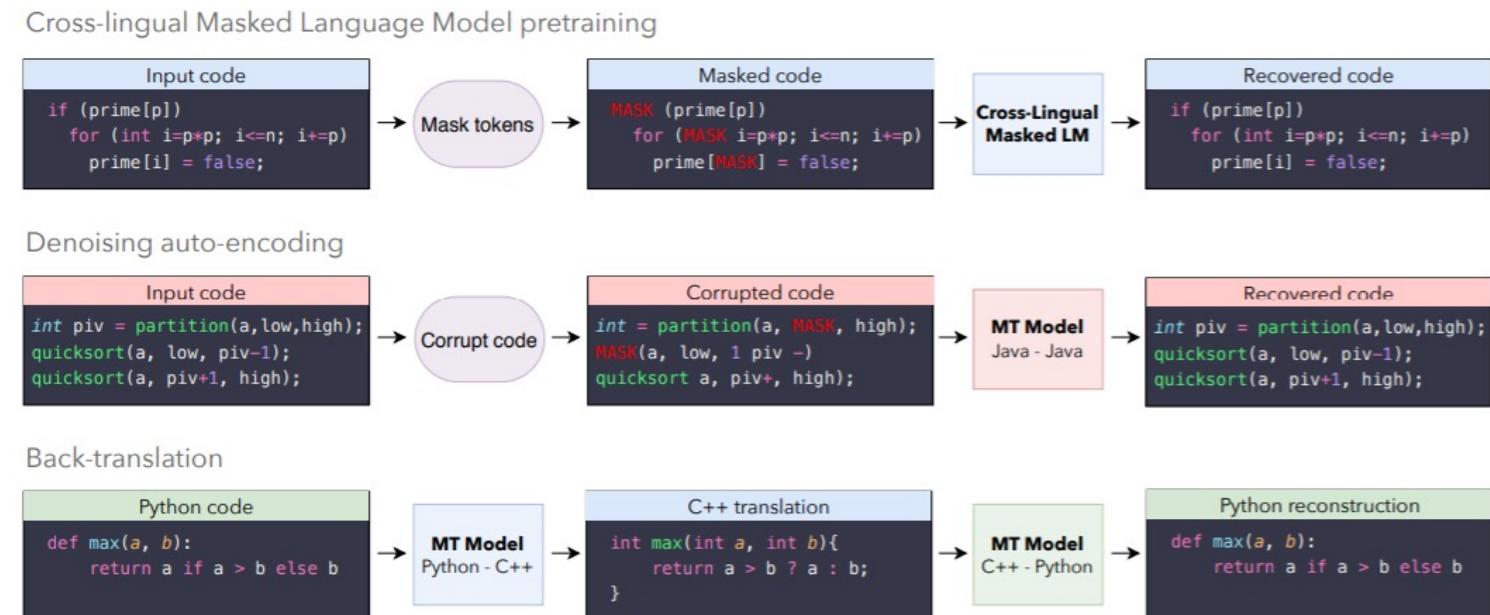
Fill in the missing code given the surrounding context

Code Migration

There is a huge need for language migration

- Language Migration
 - Cobol → Python, Java, etc.
 - C → Rust
 - ...
- Library Version Migration
 - Python2 → Python3
 - Tensorflow 1.x → Tensorflow 2.x
 -
- Platform Migration
 - X86 → ARM
 - Android → OSX
 - Intel 8080 → Intel 8086
 -

- Build neural-based models to translate code (inspired from Natural language processing).
- Use compiler-based techniques to validate the translated programs.



Excerpt from “**Unsupervised Programming Languages Translation**”, Facebook AI Research

AI-Assisted Automation Testing

- **Test Case Generation:** To generate unit test given a function.
- **Test Case Prioritization:** Prioritize which test cases to run first to enhance the fault detection rate as early as possible.
- **Image Recognition:** automatically detect changes on a website and update element to DOM (code generation) without manual effort.
- **Predictive Self-Healing:** update test cases when the code changed.
- **Bug Prediction:** running test cases to predict bugs.
-

Practical Consideration



Complement, Not Replace



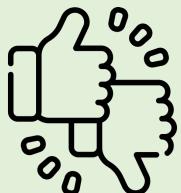
Low-Resource Environment



Explainability



Privacy-Preserving



Feedback Loop



Metrics



THANK YOU

Contact: nghibdq@fsoft.com.vn

Website: <https://bdqnghi.github.io/>