

Learning from “Big Code”

21th PeWe, April 7, 2017

Pavol Bielik

Software Reliability Lab
Department of Computer Science
ETH Zurich



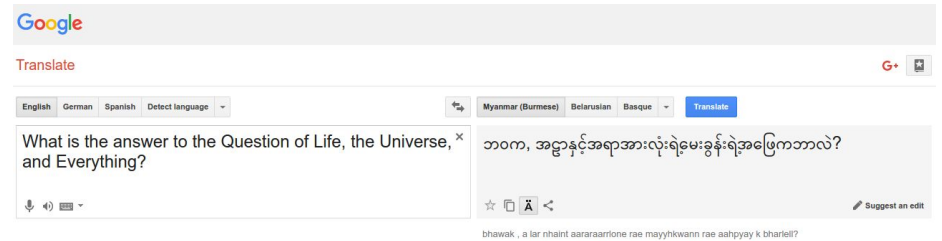
European Research Council

Established by the European Commission

**Supporting top researchers
from anywhere in the world**

Big Data Impact

Natural Language Processing
(e.g., machine translation)

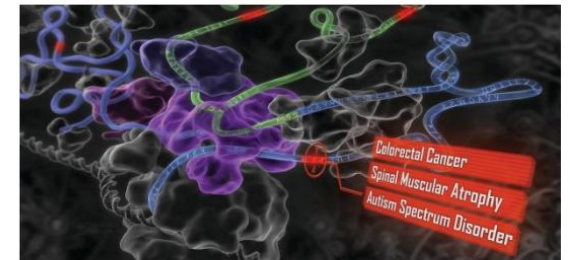


Computer Vision
(e.g., image captioning)



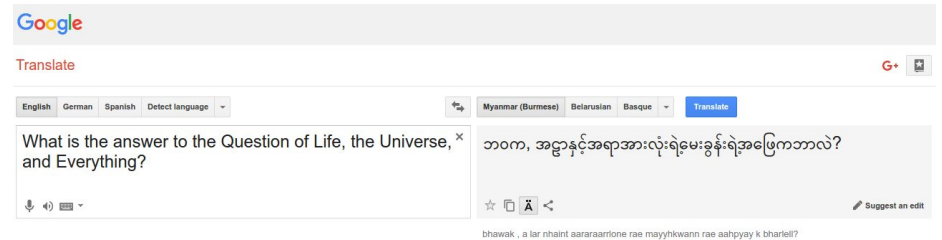
A group of people shopping
at an outdoor market.
There are many vegetables
at the fruit stand.

Medical Computing
(e.g., disease prediction)



Big Data Impact

Natural Language Processing
(e.g., machine translation)

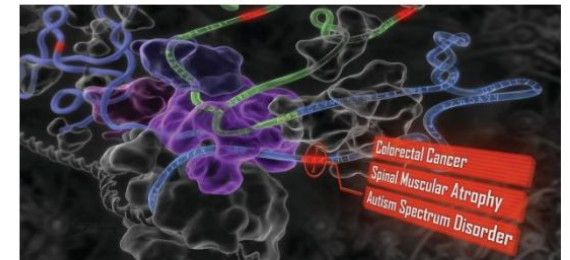


Computer Vision
(e.g., image captioning)



A group of people shopping
at an outdoor market.
There are many vegetables
at the fruit stand.

Medical Computing
(e.g., disease prediction)



Can we bring this revolution to programmers?

Vision

Create new kinds of software tools that leverage **massive codebases** to solve problems **beyond** what is possible with traditional techniques.

number of
repositories



15 million repositories

Billions of lines of code

High quality, tested,
maintained programs

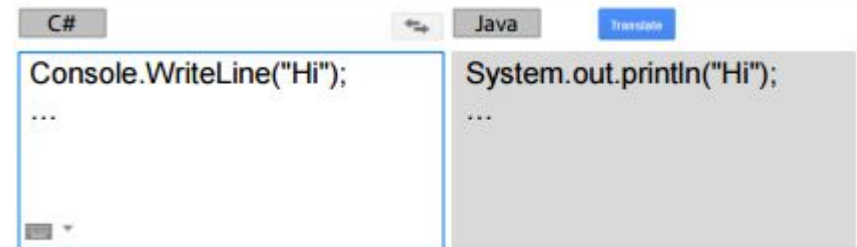
last 8 years

Statistical Software Tools

Write new code:
Code Completion

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
?
```

Port code:
Programming Language Translation

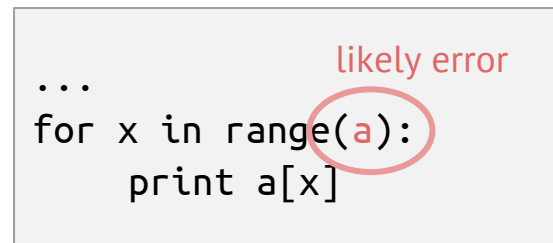


```
C# Console.WriteLine("Hi");  
...  
Java System.out.println("Hi");  
...  
Translate
```

Understand code/security:
JavaScript Deobfuscation
Type Prediction



Debug code:
Statistical Bug Detection



```
...  
for x in range(a):  
    print a[x]
```

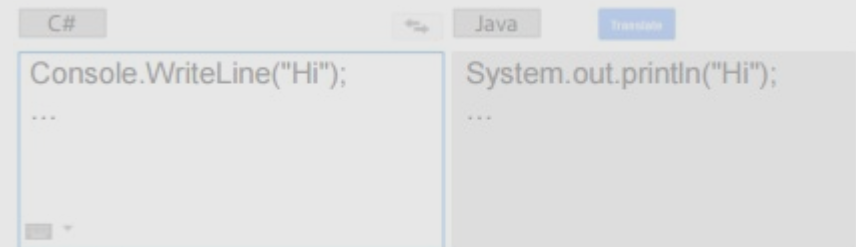
All of these benefit from the “Big Code” and lead to applications not possible with previous techniques

Statistical Software Tools

Write new code:
Code Completion

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
?
```

Port code:
Programming Language Translation



```
C# | Java | Translate  
Console.WriteLine("Hi");  
...  
System.out.println("Hi");  
...
```

Understand code/security:
JavaScript Deobfuscation
Type Prediction



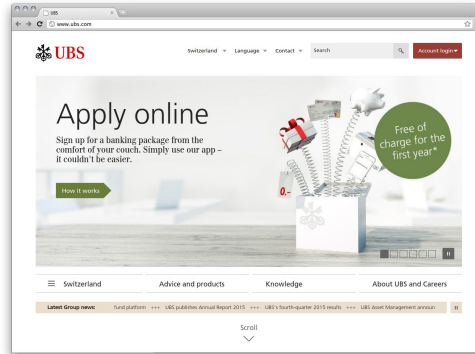
Debug code:
Statistical Bug Detection

```
...  
for x in range(a):  
    print a[x]
```

likely error

All of these benefit from the “Big Code” and lead to applications not possible with previous techniques

Code understanding/security

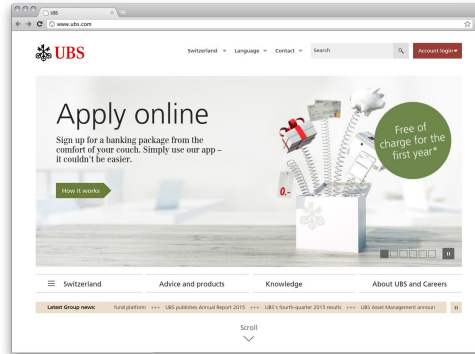


+



JavaScript

Code understanding/security



+



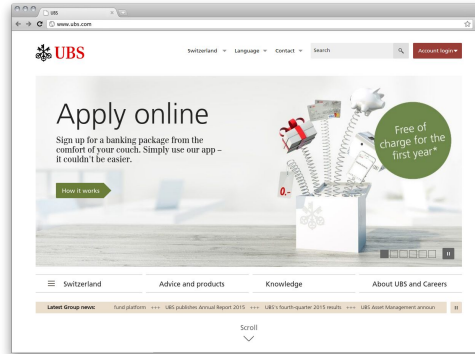
JavaScript

`savePassword(user, password)` →

Obfuscation/
Minification

→ $B(d, c)$

Code understanding/security



+



iPhone



JavaScript



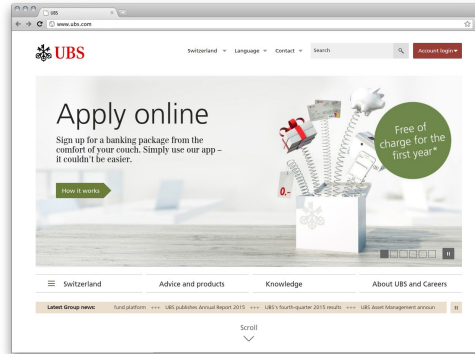
ANDROID

`savePassword(user, password)` →

Obfuscation/
Minification

→ `B(d, c)`

Code understanding/security



Security Analyst

+



iPhone



JavaScript



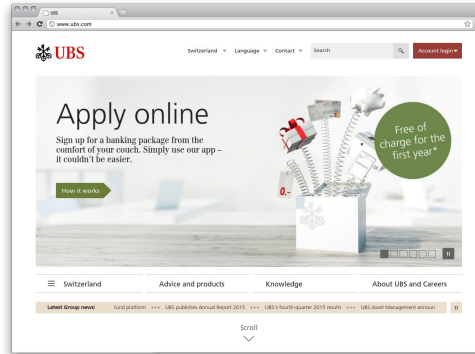
ANDROID

`savePassword(user, password)` →

Obfuscation/
Minification

→ `B(d, c)`

Code understanding/security



Security Analyst

+



iPhone



JavaScript



ANDROID

`savePassword(user, password)` →

Obfuscation/
Minification

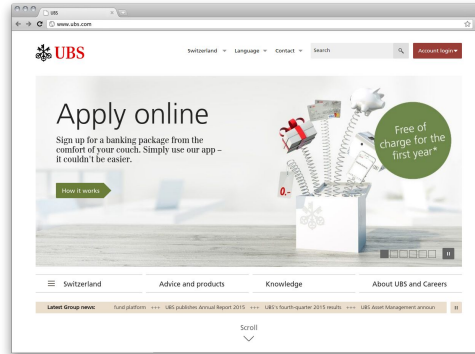
→ `B(d, c)`

`savePassword(user, password)` ←

?

← `B(d, c)`

Code understanding/security



Security Analyst

+



iPhone



JavaScript



ANDROID

`savePassword(user, password)` →

Obfuscation/
Minification

→ `B(d, c)`

`savePassword(user, password)` ←

JS NICE

← `B(d, c)`

Demo

available online:

www.jsnice.org

Impact

30, 000 Users in 1st week

Impact

30, 000 Users in 1st week



Alex Vanston @mrvdot · 7 Jun

I've been looking for this for years. JS NICE buff.ly/1pQ5qfr #javascript #unminify #deobfuscate #makeItReadable



pieter-paulus @vertongen · 23 Jun

Wow, deobfuscate your javascript with machine learning jsnice.org #javascript #jsnice



Antarctic Design @antarcticdesign · 14 Jan

De-obfuscate #JavaScript with JSNice - jsnice.org - so cool and super useful!



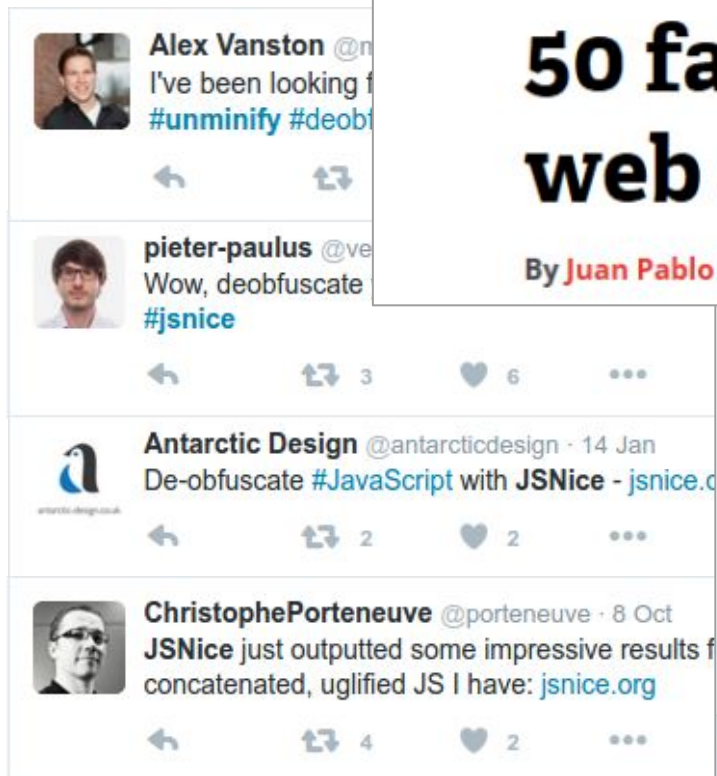
ChristophePorteneuve @porteneuve · 8 Oct

JSNice just outputted some impressive results for production-grade massive concatenated, uglified JS I have: jsnice.org



Impact

30, 000 Users in 1st week



50 fantastic freebies for web designers, July 2014

By **Juan Pablo Sarmiento** |



20 Essential Tools for Coders

14 Essential Tools for Programmers

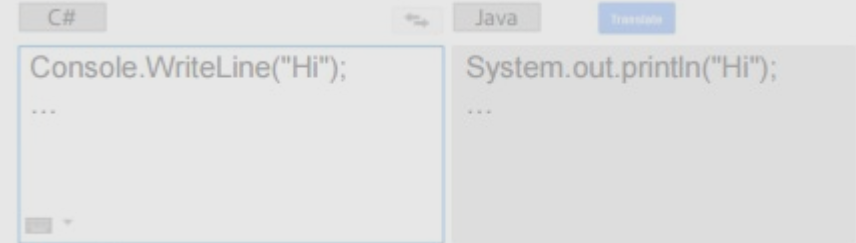
Everyone knows very well, that when it comes to web development than coding is the most important part where every web developer spends it ...

How to build such tools?

Write new code:
Code Completion

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
?
```

Port code:
Programming Language Translation



```
C# | Java | Translate  
Console.WriteLine("Hi");  
...  
System.out.println("Hi");  
...
```

Understand code/security:
JavaScript Deobfuscation
Type Prediction



Debug code:
Statistical Bug Detection

```
...  
for x in range(a):  
    print a[x]
```

likely error

JS **NICE** Statistical Type Inference and Renaming

Applications

```
function get(a, b, c) {  
    b.open("GET", a, false);  
    b.send(c);  
}
```

Predicting types

Predicting names



```
// @param {string} url  
// @param {Object} client  
// @param {string} data  
function get(url, client, data) {  
    client.open("GET", url, false);  
    client.send(data);  
}
```

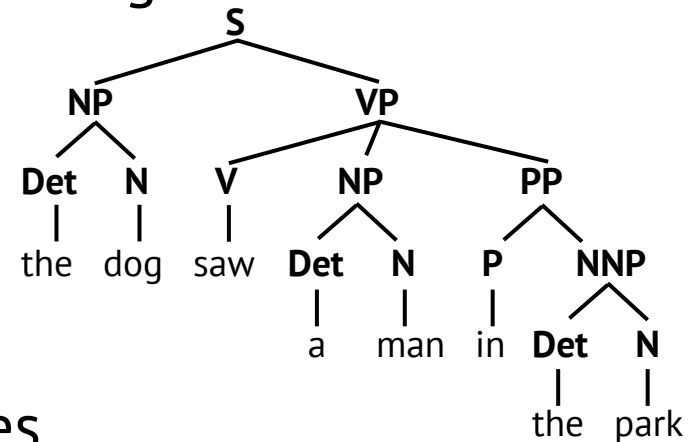
What is a suitable program representation?

Applications

Intermediate
Representation

Natural Language Processing

The dog saw a
man in the park



Programming Languages



?

What is a suitable program representation?

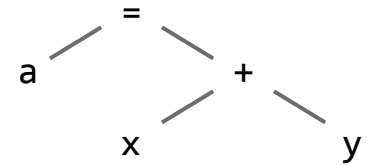
Applications

Intermediate
Representation

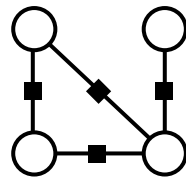
Sequences

req \rightarrow {<open, 0>, <send, 0>}
source \rightarrow {..., <open, 2>}

Trees



Graphical Models



Feature Vectors

req \rightarrow (0,0,1,1,0)
source \rightarrow (1,0,0,0,0)
...

What is a suitable program representation?

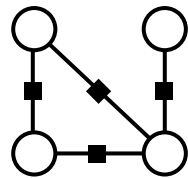
Applications

Intermediate
Representation

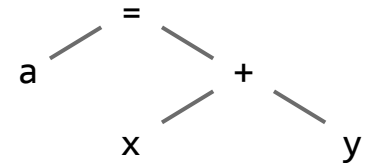
Sequences

req → {<open, 0>, <send, 0>}
source → {..., <open, 2>}

Graphical Models



Trees



Feature Vectors

req → (0, 0, 1, 1, 0)
source → (1, 0, 0, 0, 0)
...

JS Representation

Applications

Intermediate
Representation

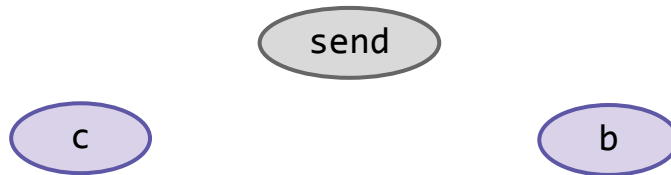
```
function get(a, b, c) {  
    b.open("GET", a, false);  
    b.send(c);  
}
```

JS **NICE** Representation

Applications

Intermediate
Representation

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```

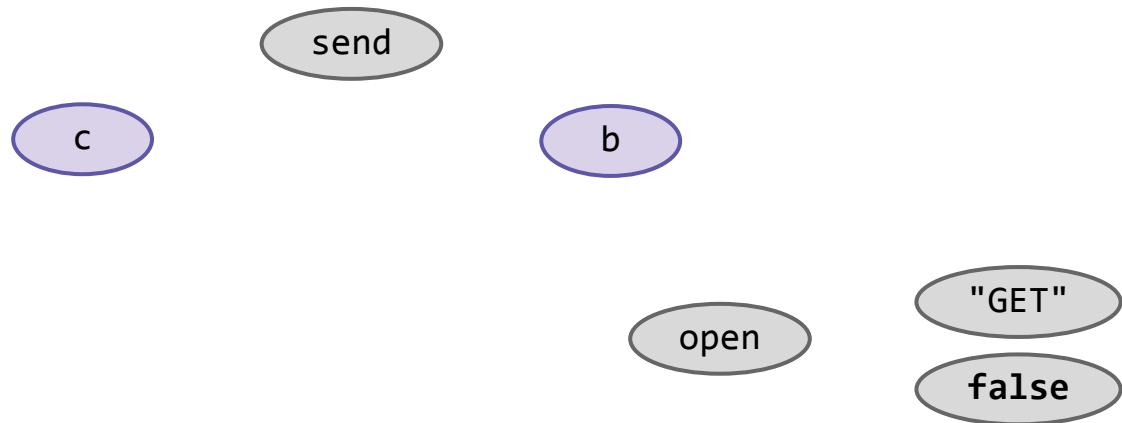


JS **NICE** Representation

Applications

Intermediate
Representation

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```

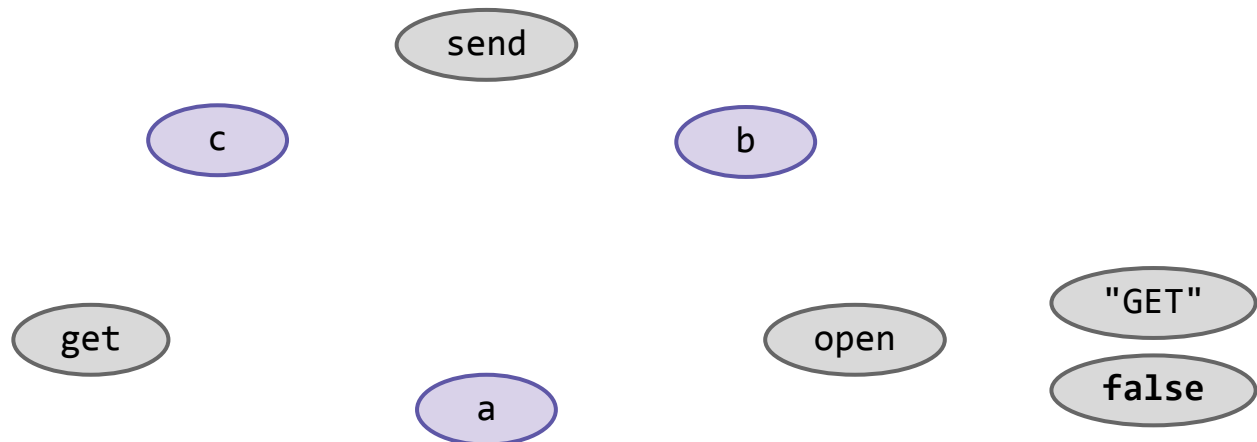


JS **nICE** Representation

Applications

Intermediate
Representation

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```



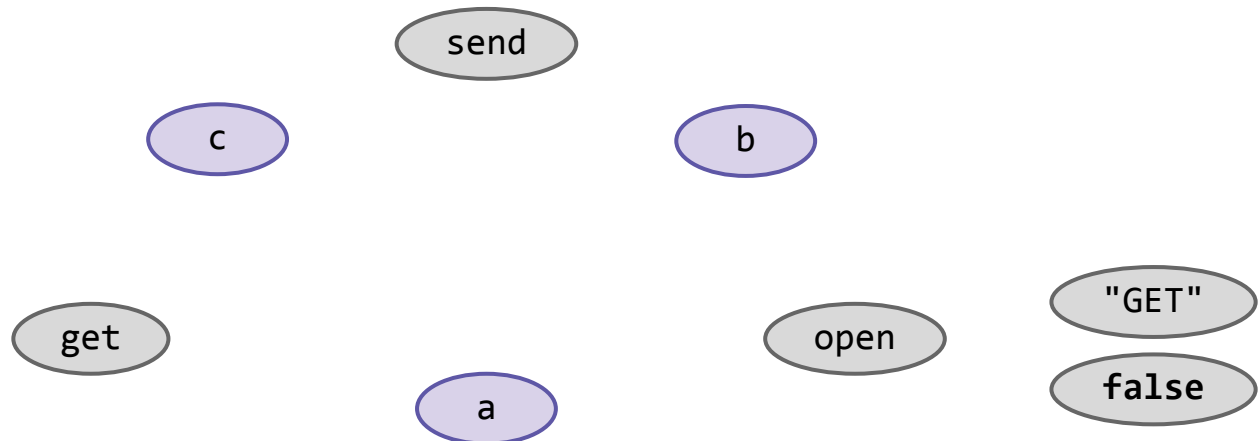
JS **nICE** Representation

Applications

Intermediate
Representation

Analyze Program
(PL)

```
function get(a, b, c) {  
    b.open("GET", a, false);  
    b.send(c);  
}
```



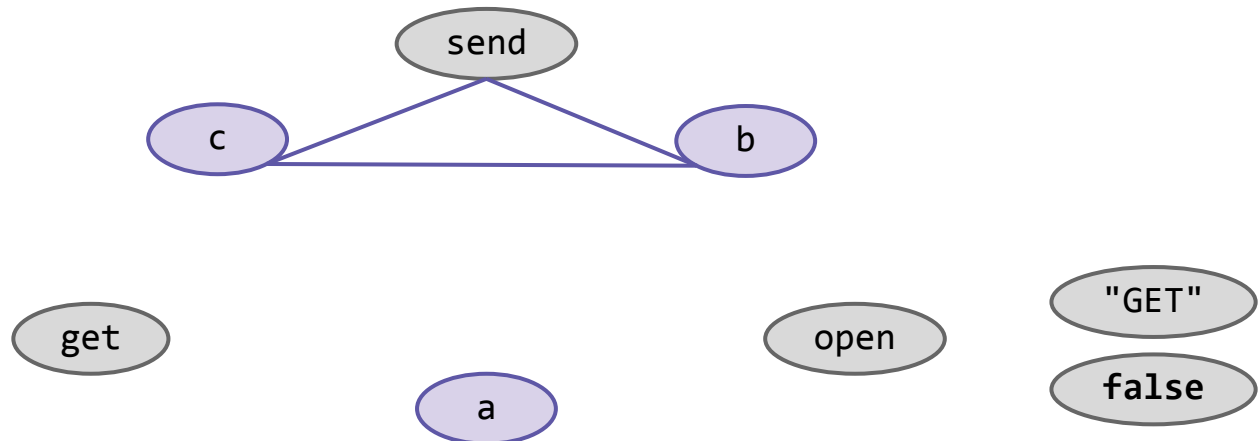
JS **NICE** Representation

Applications

Intermediate
Representation

Analyze Program
(PL)

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```



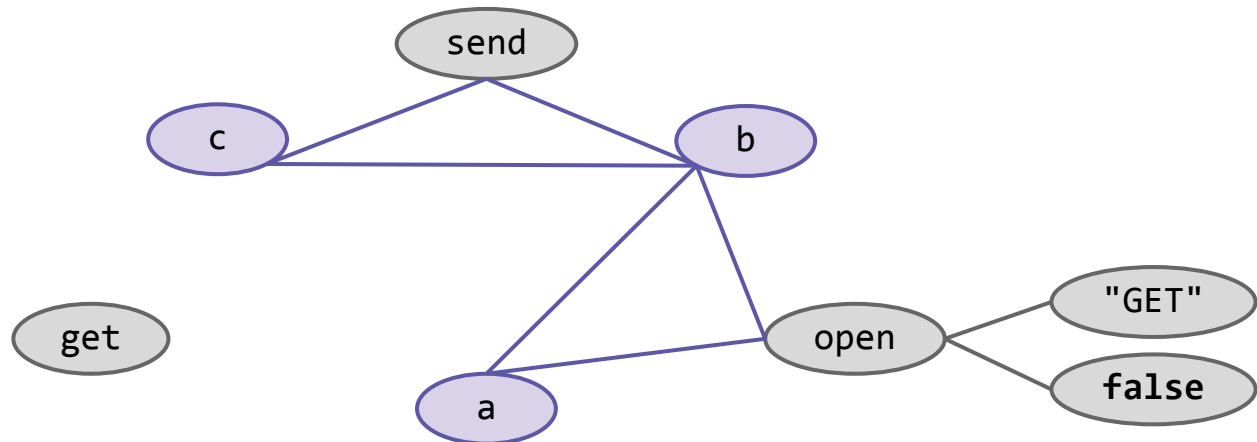
JS **NICE** Representation

Applications

Intermediate
Representation

Analyze Program
(PL)

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```



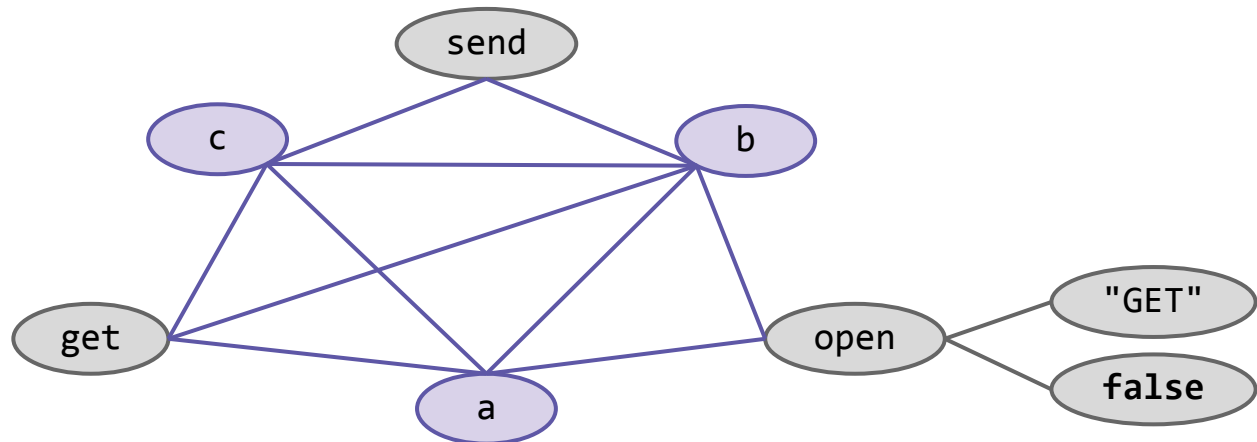
JS **NICE** Representation

Applications

Intermediate
Representation

Analyze Program
(PL)

```
function get(a, b, c) {  
  b.open("GET", a, false);  
  b.send(c);  
}
```



JS **NICE** Representation

Applications

Call graph analysis

Who are the callers of get function?

Alias analysis

b and c point to the same object?

Intermediate Representation

```
function get(a, b, c) {  
    b.open("GET", a, false);  
    b.send(c);  
}
```

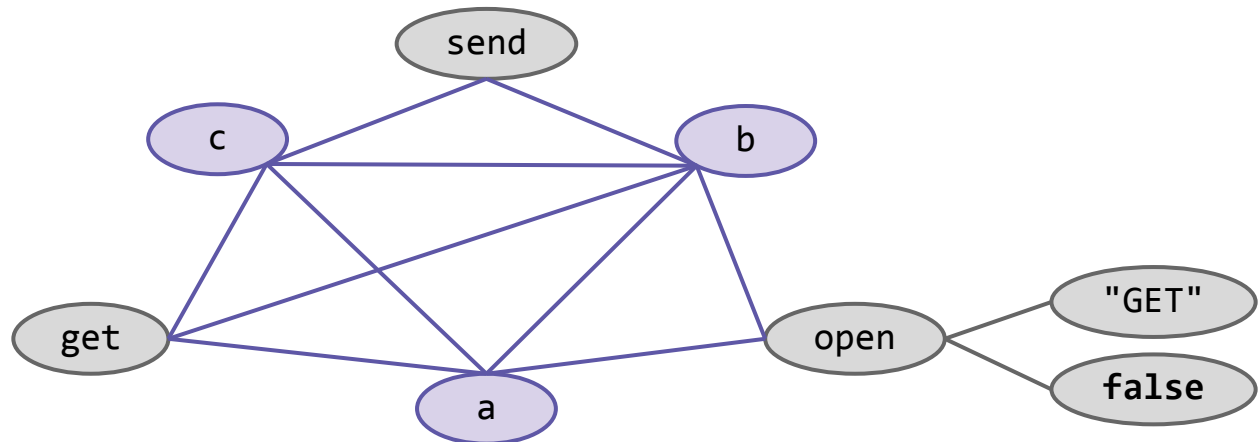
Analyze Program (PL)

Scope analysis

Program location where b was defined?

Typestate analysis

State of the object b?



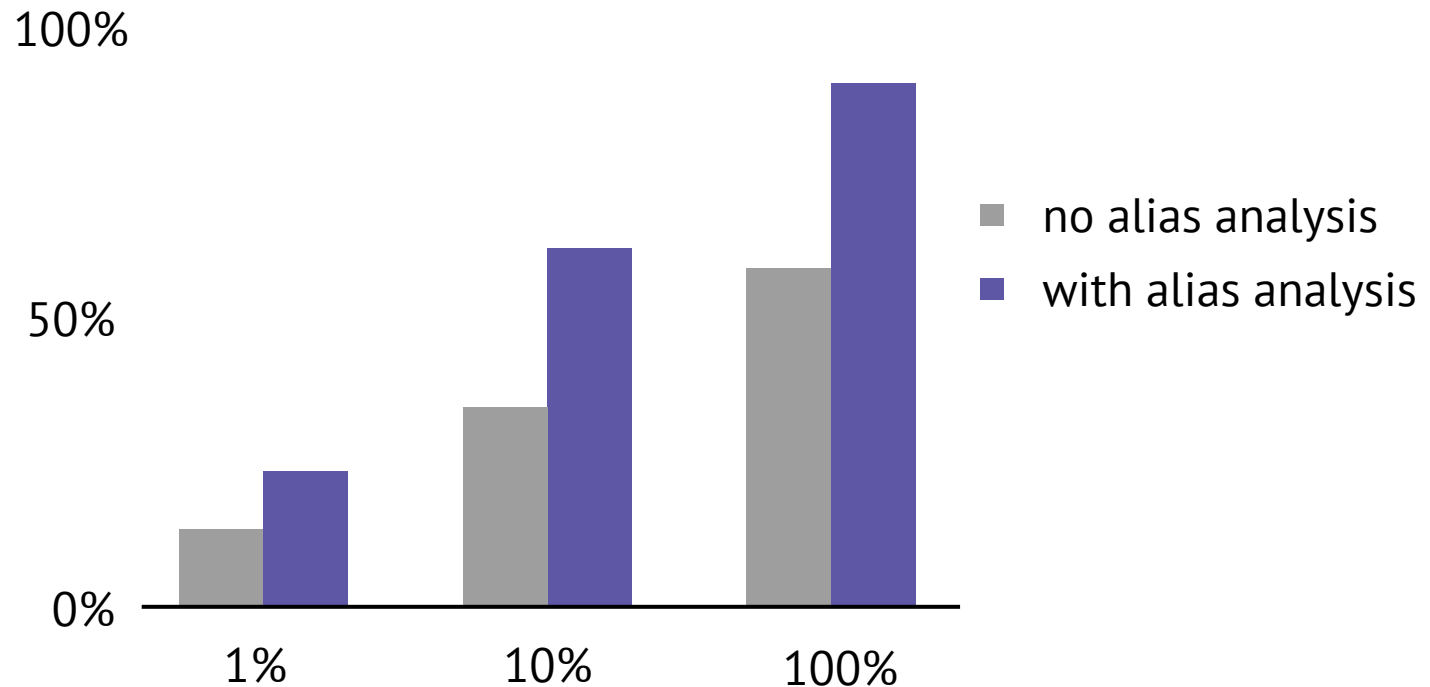
Importance of Program Analysis

Applications

Intermediate
Representation

Analyze Program
(PL)

[Precision vs % of data used]



Markov Network

Applications

Undirected graphical model

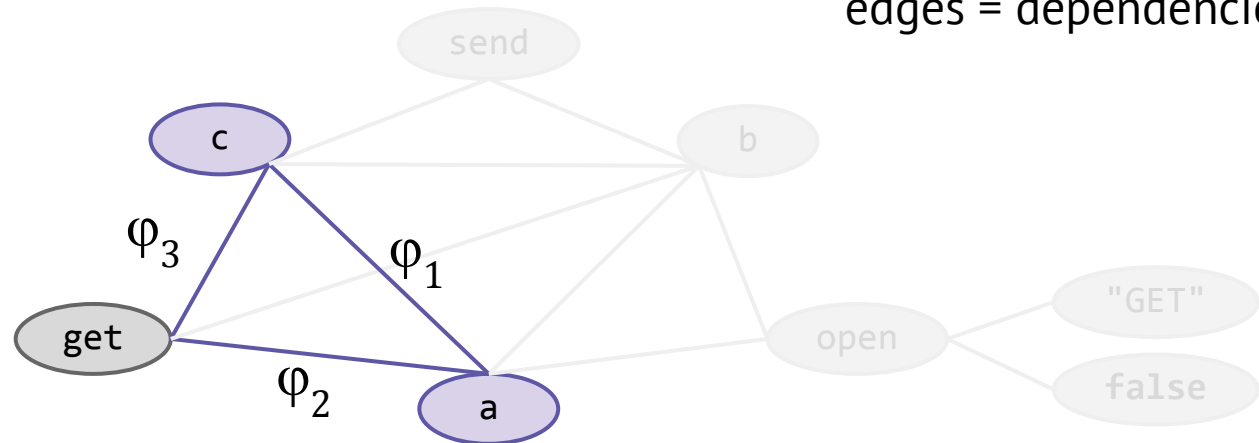
Captures **dependencies** between facts to be predicted

Intermediate Representation

nodes = random variables

edges = dependencies

Analyze Program (PL)



Train Model (ML)

$$P(\text{get}, c, a) = \varphi_1(a, c) * \varphi_2(a, \text{get}) * \varphi_2(c, \text{get}) * \dots / Z(\text{get}, c, a)$$

makes **P** a valid probability distribution
very expensive to compute

Conditional Random Fields

(J. Lafferty, A. McCallum, F. Pereira, ICML 2001)

Applications

Undirected graphical model

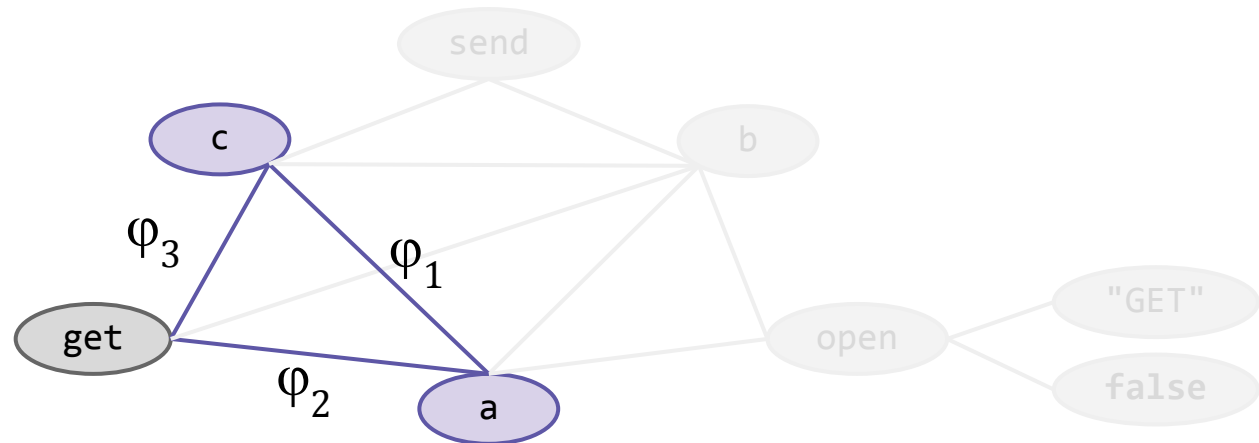
Captures **dependencies** between facts to be predicted

Captures **conditional distribution** on known facts

Intermediate Representation

Analyze Program (PL)

Train Model (ML)

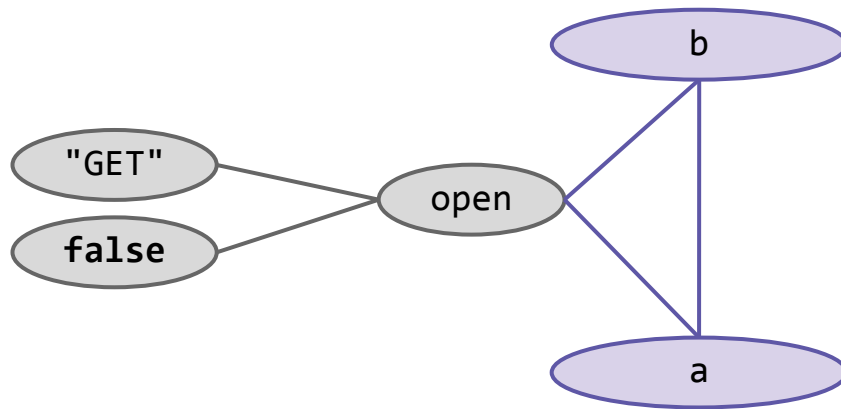


$$P(\mathbf{Y} | \mathbf{X}) = \varphi_1(a, c) * \varphi_2(a, get) * \varphi_2(c, get) * \dots / Z(\mathbf{Y})$$

makes \mathbf{P} a valid probability distribution
very expensive to compute

MAP inference example

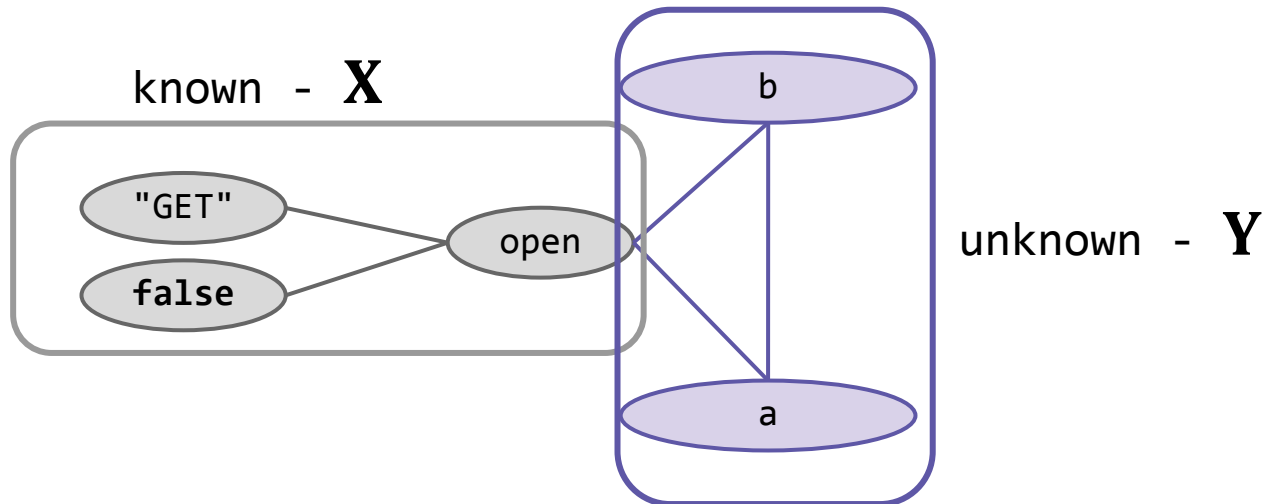
$$\operatorname{argmax}_Y \mathbf{P}(Y|\mathbf{X})$$



```
b.open("GET", a, false);
```

MAP inference example

$$\operatorname{argmax}_{\mathbf{Y}} \mathbf{P}(\mathbf{Y}|\mathbf{X})$$



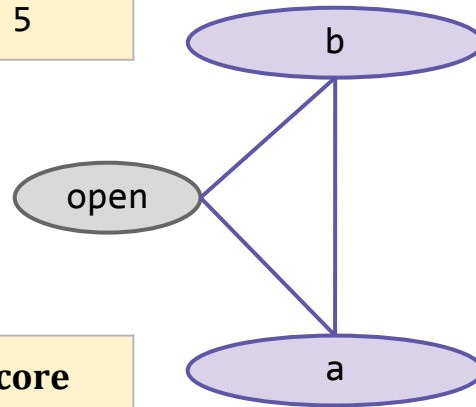
```
b.open("GET", a, false);
```

MAP inference example

$$\operatorname{argmax}_{c,t} \mathbf{P}(c,t|v=\text{open})$$

ϕ_1

v	c	Score
open	req	6
open	client	5



t	c	Score
client	url	8
client	link	5
req	link	2

ϕ_2

v	t	Score
open	link	7
open	url	5
open	address	2

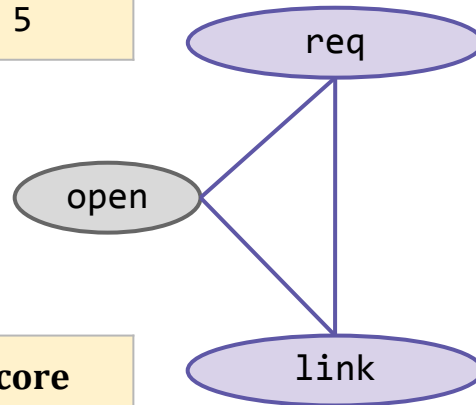
`b.open("GET", a, false);`

MAP inference example

$$\operatorname{argmax}_{c,t} \mathbf{P}(c,t|v=\text{open})$$

ϕ_1

v	c	Score
open	req	6
open	client	5



Maximize product of scores:

$$6 * 7 * 2 = 84$$

t	c	Score
client	url	8
client	link	5
req	link	2

ϕ_2

v	t	Score
open	link	7
open	url	5
open	address	2

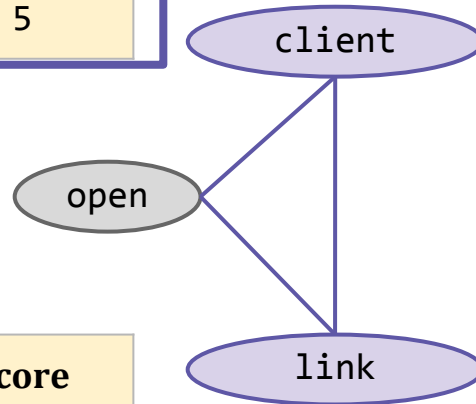
```
b.open("GET", a, false);
```

MAP inference example

$$\operatorname{argmax}_{c,t} \mathbf{P}(c,t|v=\text{open})$$

ϕ_1

v	c	Score
open	req	6
open	client	5



Maximize product of scores:

$$5 * 7 * 5 = 175$$

t	c	Score
client	url	8
client	link	5
req	link	2

ϕ_2

v	t	Score
open	link	7
open	url	5
open	address	2

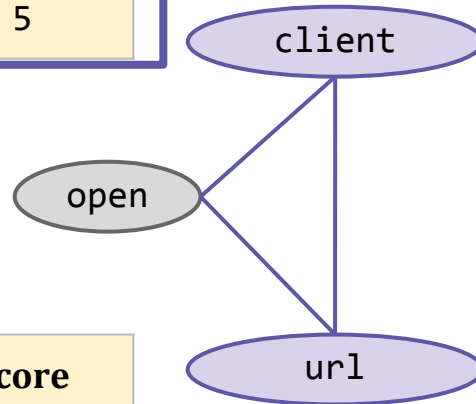
```
b.open("GET", a, false);
```

MAP inference example

$$\operatorname{argmax}_{c,t} \mathbf{P}(c,t|v=\text{open})$$

ϕ_1

v	c	Score
open	req	6
open	client	5



Maximize product of scores:

$$5 * 5 * 8 = 200$$

t	c	Score
client	url	8
client	link	5
req	link	2

ϕ_2

v	t	Score
open	link	7
open	url	5
open	address	2

```
b.open("GET", a, false);
```

Query

Our goal is to find the most likely assignment of \mathbf{y} that satisfies the constraints, also known as **MAP inference**:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} \mathbf{P}(\mathbf{y}' | \mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} 1/\mathbf{Z} \prod \phi_i(\mathbf{x}, \mathbf{y})$$

Good news:

the expensive partition function $\mathbf{Z}(\mathbf{x})$ is unnecessary

Query

Our goal is to find the most likely assignment of \mathbf{y} that satisfies the constraints, also known as **MAP inference**:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} \mathbf{P}(\mathbf{y}' | \mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} 1/Z \prod \phi_i(\mathbf{x}, \mathbf{y})$$

Good news:

the expensive partition function $\mathbf{Z}(\mathbf{x})$ is unnecessary

Bad news:

computing the argmax is still NP-hard (Max-SAT)

Query

Our goal is to find the most likely assignment of \mathbf{y} that satisfies the constraints, also known as **MAP inference**:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} \mathbf{P}(\mathbf{y}'|\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} 1/Z \prod \phi_i(\mathbf{x},\mathbf{y})$$

Good news:

many approximate algorithm exists (Variational Methods, EM, Gibbs sampling, Elimination Algorithm, Junction-Tree algorithm)

Query

Our goal is to find the most likely assignment of \mathbf{y} that satisfies the constraints, also known as **MAP inference**:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} \mathbf{P}(\mathbf{y}' | \mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \frac{1}{Z} \prod \phi_i(\mathbf{x}, \mathbf{y})$$

Good news:

many approximate algorithm exists (Variational Methods, EM, Gibbs sampling, Elimination Algorithm, Junction-Tree algorithm)

Bad news:

still too slow for learning

Query

Our goal is to find the most likely assignment of \mathbf{y} that satisfies the constraints, also known as **MAP inference**:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} \mathbf{P}(\mathbf{y}' | \mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} 1/Z \prod \phi_i(\mathbf{x}, \mathbf{y})$$

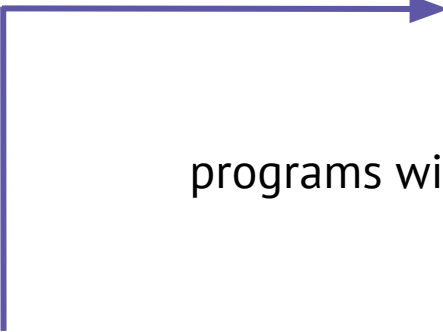
Good news:

approximate algorithms designed to fit our setting

Learning

$$P(\mathbf{y}|\mathbf{x}) = 1/Z \exp \sum \lambda_i f_i(\mathbf{x}, \mathbf{y})$$

Learning finds weights λ_i from training data


$$D = \{ \mathbf{x}^{(j)}, \mathbf{y}^{(j)} \}_{j=1..n}$$

programs with facts of interest already manually annotated

Big codebase to learn from

Programmers have spent countless hours to develop, maintain and annotate

Structured SVM

Generalizes SVM, learns weights such that:

$$\forall j \forall \mathbf{y} \sum \lambda_i f_i(\mathbf{x}^{(j)}, \mathbf{y}^{(j)}) \geq \sum \lambda_i f_i(\mathbf{x}^{(j)}, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}^{(j)})$$

for all training data samples

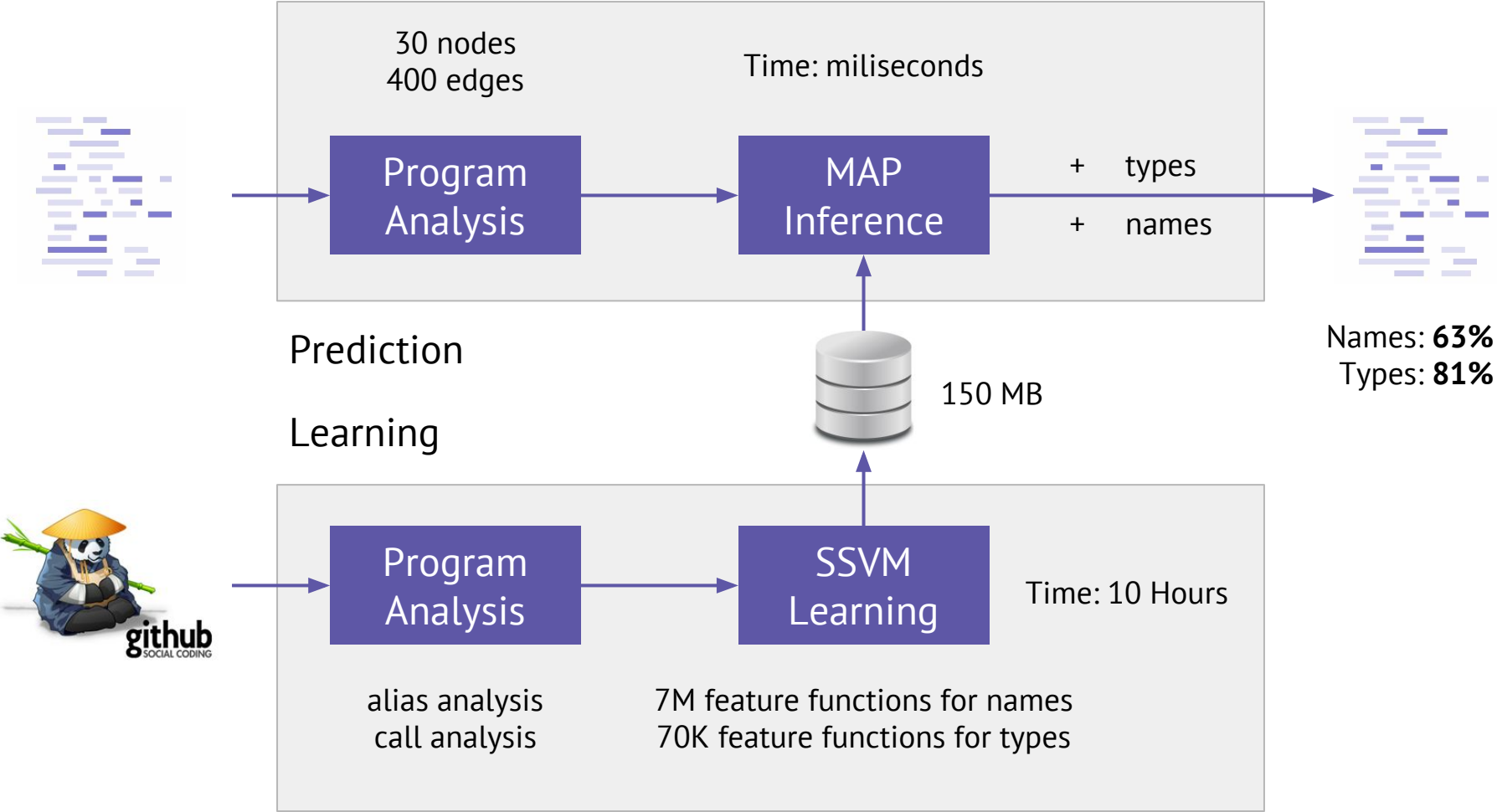
the **given prediction** is better than **any other prediction** by at least a **margin**

Training procedure:

N.Ratliff, J. Bagnell, M. Zinkevich: (Online) Subgradient Methods for Structured Prediction, AISTATS'07

Memory efficient
Fast and scalable

Structured Prediction for Programs



Programming with “Big Code”

Applications	Code completion Deobfuscation	Program synthesis	Translation Feedback generation
Intermediate Representation	Sequences (sentences) Trees	Translation Table	Graphical Models Feature Vectors
Analyze Program (PL)	alias analysis scope analysis	control-flow analysis typestate analysis	
Train Model (ML)	Neural Networks N-gram language model	SVM	Structured SVM
Query Model	$\operatorname{argmax}_{y \in \Omega} P(y x)$		Greedy MAP Inference

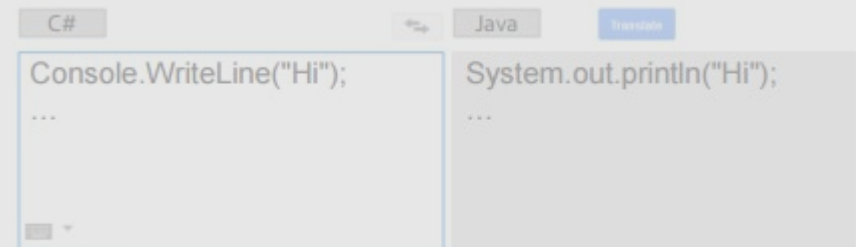
More information and tutorials at: <http://www.nice2predict.org/>
<http://plml.ethz.ch/>

Statistical Software Tools

Write new code:
Code Completion

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
?
```

Port code:
Programming Language Translation



The screenshot shows a web-based code translation interface. On the left, under a 'C#' tab, the code is: `Console.WriteLine("Hi");` followed by an ellipsis. On the right, under a 'Java' tab, the code is: `System.out.println("Hi");` followed by an ellipsis. A 'Translate' button is visible in the top right corner.

Understand code/security:
JavaScript Deobfuscation
Type Prediction



Debug code:
Statistical Bug Detection

```
...  
for x in range(a):  
    print a[x]
```

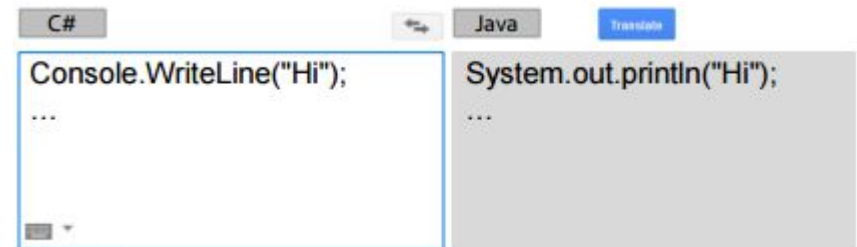
likely error

Statistical Software Tools

Write new code:
Code Completion

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
?
```

Port code:
Programming Language Translation



The screenshot shows a web-based code translation interface. On the left, under a 'C#' tab, the code is: `Console.WriteLine("Hi");` followed by an ellipsis. On the right, under a 'Java' tab, the code is: `System.out.println("Hi");` followed by an ellipsis. A 'Translate' button is visible in the top right corner.

Understand code/security:
JavaScript Deobfuscation
Type Prediction



Debug code:
Statistical Bug Detection

```
...  
for x in range(a):  
    print a[x]
```

likely error

Probabilistic Model for Code

Model is a key part of the Statistical Programming Tools

Goal: score programs

Select best among
several candidates

Probabilistic Model for Code

Model is a key part of the Statistical Programming Tools

Goal: score programs

Select best among several candidates

Example: Which function is more likely?

```
(a) function area(a) {  
      return a.width * a.height  
}
```

```
(b) function area(a) {  
      return a.width * a.close()  
}
```

Statistical Code Completion

Model is a key part of Statistical Programming Tools

Goal: score programs

Select best among several candidates

Example:

```
function area(a) {  
    return a.width * a.  
}
```

A diagram illustrating the likelihood of different code completion candidates. A blue box on the left contains the candidates: **height**, **width**, **open**, and **close**. A yellow box on the right contains the likelihood categories: **Very likely**, **Less likely**, and **impossible**. Blue arrows point from the likelihood categories to the candidates: **Very likely** points to **height**, **Less likely** points to **width**, and **impossible** points to both **open** and **close**.

height	← Very likely
width	← Less likely
open	← impossible
close	← impossible

Probabilistic Model for Code

Directly applicable to **code completion**, but is a ***key statistical component*** for many others tasks:
e.g. **natural language to code**,
statistical bug localization

Model Requirements

Existing Programs

Learning

Model



Probabilistic
Model



Widely
Applicable

Efficient
Learning

High
Precision

Explainable
Predictions

Observation

Regularities in code are similar
to regularities in natural language

The quick brown fox jumps over the lazy ?

Observation

Regularities in code are similar
to regularities in natural language

The quick brown fox jumps over the lazy dog

Observation

Regularities in code are similar
to regularities in natural language

The quick brown fox jumps over the lazy dog

```
file = open(filename, "r")  
file.??
```

Observation

Regularities in code are similar
to regularities in natural language

The quick brown fox jumps over the lazy dog

```
file = open(filename, "r")  
file.read()
```

N-gram Language Model

Conditional probability only on previous n-1 words

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{n-1 \text{ words}}) \approx \frac{\#(w_{i-n+1} \dots w_{i-1} w_i)}{\#(w_{i-n+1} \dots w_{i-1})}$$

#(n-gram) - number of occurrences of n-gram in training data

N-gram Language Model

Conditional probability only on previous n-1 words

$$P(w_i | w_1 \dots w_{i-1}) \approx P(w_i | \underbrace{w_{i-n+1} \dots w_{i-1}}_{n-1 \text{ words}}) \approx \frac{\#(w_{i-n+1} \dots w_{i-1} w_i)}{\#(w_{i-n+1} \dots w_{i-1})}$$

Training is achieved by counting n-grams (3-gram)

$$P(\text{jumped} | \text{The quick brown fox}) \approx P(\text{jumped} | \text{brown fox}) \approx \frac{\#(\text{brown fox jumped})}{\#(\text{brown fox})}$$

$\#(n\text{-gram})$ - number of occurrences of n-gram in training data

N-gram Language Model

Training

```
return f . height * scale;
```

```
f . open (mode) ;
```

```
f . close () ;
```

```
2 * f . width;
```

```
f . close ()
```

Prediction

```
f.width + f. ?
```

N-gram Language Model

Training (3-gram model)

```
return f . height * scale;
```

```
f . open (mode) ;
```

```
f . close () ;
```

```
2 * f . width;
```

```
f . close ()
```

3-gram

Prediction

```
f.width + f. ?
```

N-gram Language Model

Training (3-gram model)

```
return f . height * scale;
```

```
f . open (mode) ;
```

```
f . close () ;
```

```
2 * f . width;
```

```
f . close ()
```



3-gram

Prediction

```
f.width + f. ? P
```

close	0.4
open	0.2
width	0.2
height	0.2

N-gram Language Model

Conditioning

Accuracy

Last two tokens, Hindle et. al. [ICSE'12]

22.2%

Main Problem:

f.width +

f.

?

P

Bad context leads to bad probability estimates

close	0.4
open	0.2
width	0.2
height	0.2

Better Context

Training

```
return x.width * x.height
```

```
return y.width * y.height
```

```
area = s.width * s.height
```

```
s.width = s.width + 10
```

```
q.depth * q.width * q.height
```

Prediction

```
f.width + f.?
```

Better Context

Training

return x.**width** * x.**height**

return y.**width** * y.**height**

area = s.**width** * s.**height**

s.**width** = s.**width** + 10

q.depth * q.**width** * q.**height**

Prediction

f.**width** + f.**?**

Context: relevant for
this prediction

height	0.8
width	0.2
open	0.0
close	0.0

Better Context

Conditioning	Accuracy
<i>Last two tokens, Hindle et. al. [ICSE'12]</i>	22.2%
<i>Last two APIs, Raychev et. al. [PLDI'14]</i>	30.4%

JavaScript APIs

Conditioning	Accuracy
---------------------	-----------------

<i>Last two tokens, Hindle et. al. [ICSE'12]</i>	22.2%
--	-------

<i>Last two APIs, Raychev et. al. [PLDI'14]</i>	30.4%
---	-------

is this the best we can do?

JavaScript APIs

Conditioning	Accuracy
---------------------	-----------------

<i>Last two tokens, Hindle et. al. [ICSE'12]</i>	22.2%
--	-------

<i>Last two APIs, Raychev et. al. [PLDI'14]</i>	30.4%
---	-------

is this the best we can do?

<i>Program synthesis</i>	66.4%
--------------------------	-------

JavaScript APIs

Conditioning

Accuracy

Last two tokens, Hindle et. al. [ICSE'12]

22.2%

Last two APIs, Raychev et. al. [PLDI'14]

30.4%

Last three APIs

Declaration Site + Last two APIs

Variable Name + Method Name + Last API

...

***How do we know
that which is the
best context?***

JavaScript APIs

Conditioning

Accuracy

Last two tokens, Hindle et. al. [ICSE'12]

22.2%

Last two APIs, Raychev et. al. [PLDI'14]

30.4%

Last three APIs

Declaration Site + Last two APIs

Variable Name + Method Name + Last API

...

***How do we know
that which is the
best context?***

JavaScript APIs

Identifiers

Strings

Numbers

Arguments

Properties

Statements

RegExp

Structure

Solution: Synthesise the Best Model

$$\{\text{width}\} = f(x.\text{width} + x.\text{?})$$

Synthesise a function f from a domain specific language that explains the data

Function Examples

$f(p_1) = \{ \}$

```
for (j = 0; j < groups.length; j++) {  
  idsInGroup = groups[j].filter(  
    function(id) { return id >= 42; }  
  );  
  if (idsInGroup.length == 0) {  
    ?  
  }  
}
```

Function Examples

$f(p_1) = \{ \begin{array}{l} \text{for,} \\ \text{if,} \\ \text{length} == 0 \end{array} \}$

```
for (j = 0; j < groups.length; j++) {  
  idsInGroup = groups[j].filter(  
    function(id) { return id >= 42; }  
  );  
  if (idsInGroup.length == 0) {  
    ?  
  }  
}
```

Function Examples

$f(p_1) = \{$
 for,
 if,
 length==0
 $\}$

$f(p_2) = \{$
 $\}$

```
for (j = 0; j < groups.length; j++) {  
    idsInGroup = groups[j].filter(  
        function(id) { return id >= 42; }  
    );  
    if (idsInGroup.length == 0) {  
        ?  
    }  
}
```

```
elem.notify( ..., {  
    position: 'top',  
    hide: false,  
    ?  
} );
```

Function Examples

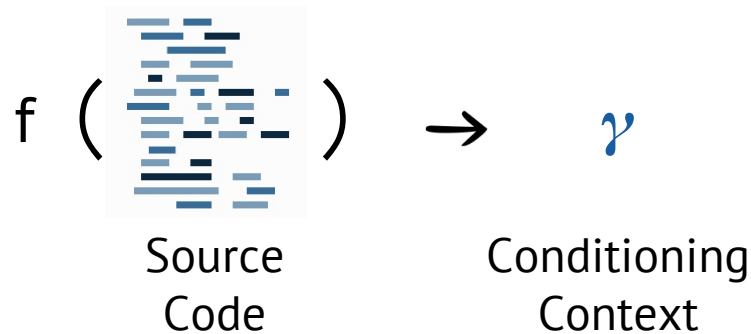
$f(p_1) = \{ \text{for, if, length} == 0 \}$

$f(p_2) = \{ \text{notify, position, hide} \}$

```
for (j = 0; j < groups.length; j++) {  
  idsInGroup = groups[j].filter(  
    function(id) { return id >= 42; }  
  );  
  if (idsInGroup.length == 0) {  
    ?  
  }  
}
```

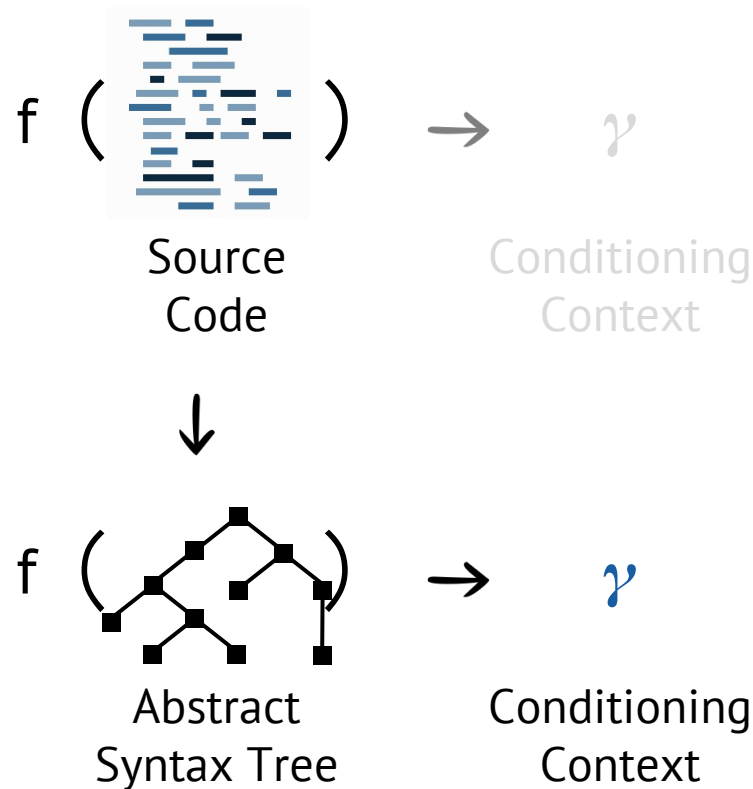
```
elem.notify( ..., {  
  position: 'top',  
  hide: false,  
  ?  
} );
```

Overview



Synthesise a function f from a domain specific language that explains the data

Overview



Synthesise a function f from a domain specific language that explains the data

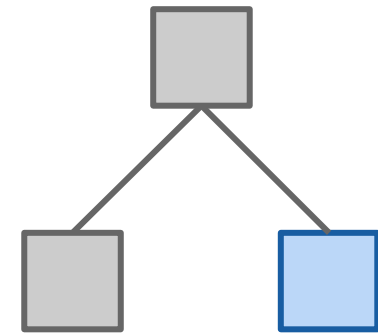
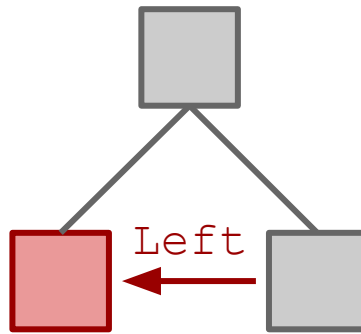
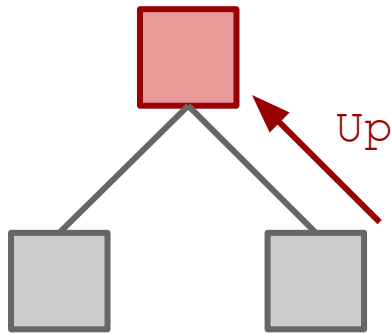
Function Representation

In general:
Unrestricted programs (Turing complete)

Our Work:
TCond Language for navigating over trees
and accumulating context

```
TCond      ::=  $\varepsilon$  | WriteOp TCond | MoveOp TCond | BranchProg
BranchProg ::= if pred(x) then TCond else TCond
MoveOp     ::= Up, Left, Right, DownFirst, DownLast,
              NextDFS, PrevDFS, NextLeaf, PrevLeaf,
              PrevNodeType, PrevNodeValue, PrevNodeContext
WriteOp    ::= WriteValue, WriteType, WritePos
```

Expressing functions: TCond Language



WriteValue

$\gamma \leftarrow \gamma \cdot \square$

TCond ::= ε | WriteOp TCond | MoveOp TCond | BranchProg

BranchProg ::= **if** pred(x) **then** TCond **else** TCond

MoveOp ::= Up, Left, Right, DownFirst, DownLast, NextDFS, PrevDFS, NextLeaf, PrevLeaf, PrevNodeType, PrevNodeValue, PrevNodeContext

WriteOp ::= WriteValue, WriteType, WritePos

Example

Query

TCond

γ

```
elem.notify(  
  ... ,  
  ... ,  
  {  
    position: 'top',  
    hide: false,  
    ?  
  }  
);
```

Example

Query

```
elem.notify(  
  ... ,  
  ... ,  
  {  
    position: 'top',  
    hide: false,  
    ?  
  }  
);
```

TCond

Left
WriteValue

γ

{ }
{hide}

Example

Query

```
elem.notify(  
  ... ,  
  ... ,  
  {  
    position: 'top',  
    hide: false,  
    ?  
  }  
);
```

TCond

Left
WriteValue
Up
WritePos

γ

{ }
{hide}
{hide}
{hide, 3}

Example

Query

```
elem.notify(  
  ... ,  
  ... ,  
  {  
    position: 'top',  
    hide: false,  
    ?  
  }  
);
```

TCond

```
Left  
WriteValue  
Up  
WritePos  
Up  
DownFirst  
DownLast  
WriteValue
```

γ

```
{ }  
{hide}  
{hide}  
{hide, 3}  
{hide, 3}  
{hide, 3}  
{hide, 3}  
{hide, 3, notify}
```

Example

Query

```
elem.notify(  
  ... ,  
  ... ,  
  {  
    position: 'top',  
    hide: false,  
    ?  
  }  
);
```

TCond

```
Left  
WriteValue  
Up  
WritePos  
Up  
DownFirst  
DownLast  
WriteValue {hide, 3, notify}
```

γ

```
{ }  
{hide}  
{hide}  
{hide, 3}  
{hide, 3}  
{hide, 3}  
{hide, 3}  
{hide, 3, notify}
```



{ Previous Property, Parameter Position, API name }

Results

Probabilistic Model of JavaScript Language

20k Learning

100k Training

50k Blind Set

GitHub

JavaScript APIs

Conditioning

Accuracy

Last two tokens, Hindle et. al. [ICSE'12]

22.2%

Last two APIs, Raychev et. al. [PLDI'14]

30.4%

Program synthesis

66.4%

JavaScript Structure

Model	Accuracy
<i>PCFG</i>	51.1%
<i>N-gram</i>	71.3%
<i>Naïve Bayes</i>	44.2%
<i>SVM</i>	70.5%
<i>Program synthesis</i>	81.5%

JavaScript Values

	Accuracy	Example
<i>Identifier</i>	62%	contains = <u>jQuery</u> ...
<i>Property</i>	65%	start = list. <u>length</u> ;
<i>String</i>	52%	`[` + attrs + <u> `]</u> `
<i>Number</i>	64%	canvas(xy[0], xy[<u>1</u>], ...)
<i>RegExp</i>	66%	line.replace(<u>/(&nbsp;)+/</u> , ...)
<i>UnaryExpr</i>	97%	if (!events <u>!</u> ...)
<i>BinaryExpr</i>	74%	while (++index <u>≤</u> ...)
<i>LogicalExpr</i>	92%	frame = frame <u> </u> ...

Model Requirements

Existing Programs

Learning

Model



Probabilistic
Model



Widely
Applicable

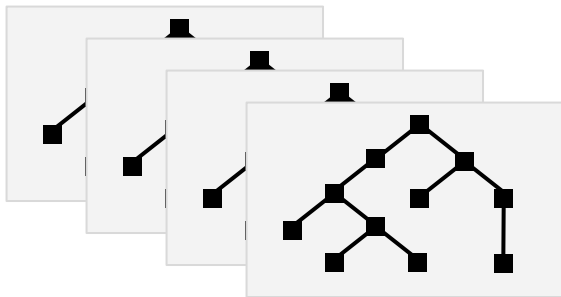
Efficient
Learning

High
Precision

Explainable
Predictions

Learning

Dataset



Program Synthesis

- Enumerative search
- Genetic programming
- Decision tree learning
- MCMC



$$f_{best} = \arg \min_{f \in \text{TCond}} \text{cost}(D, f)$$



$$|d| \ll |D|$$

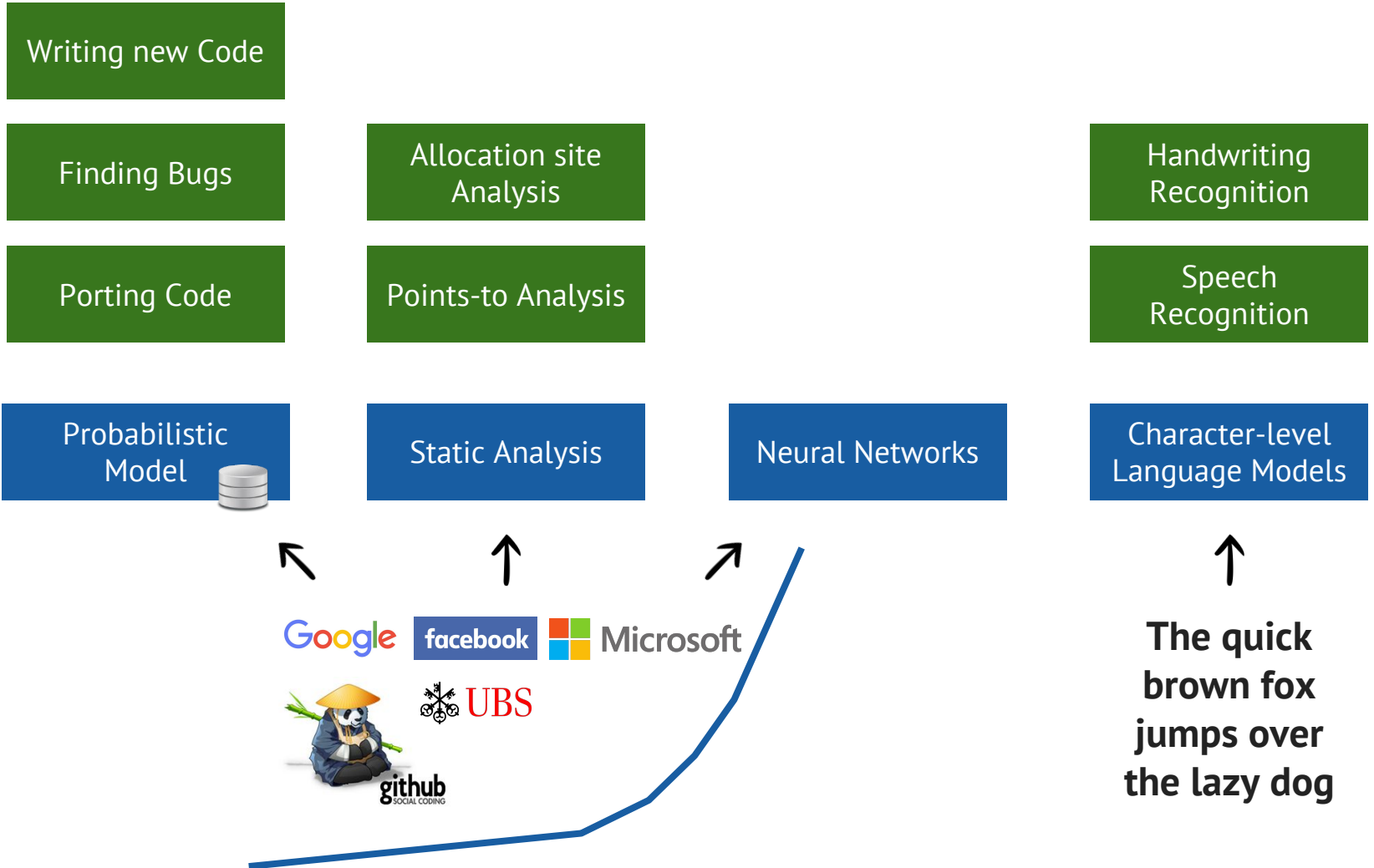
$$|\text{cost}(d, f) - \text{cost}(D, f)| < \varepsilon$$

Representative sampling

```
TCond ::= ε | WriteOp TCond | MoveOp TCond
MoveOp ::= Up, Left, Right, ...
WriteOp ::= WriteValue, WriteType, ...
```

TCond Language

Applications



Work @ ETH Zurich



Prof.
Martin
Vechev



Veselin
Raychev



Pavol
Bielik



Christine
Zeller



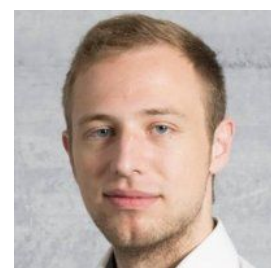
Pascal
Roos



Benjamin
Bischel

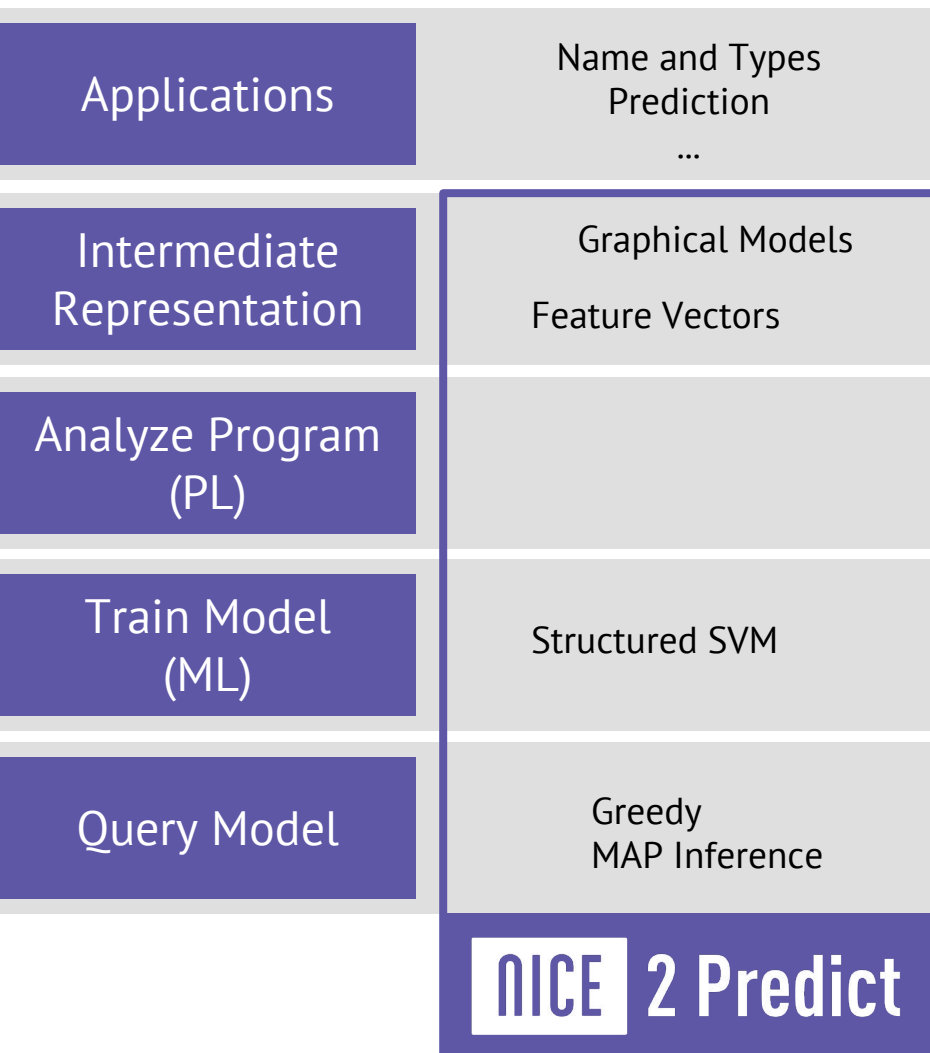


Svetoslav
Karaivanov



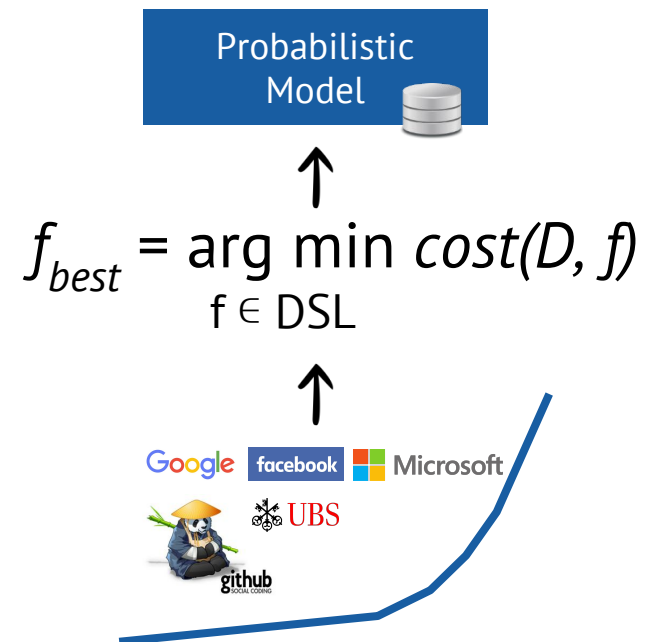
Benjamin
Mularczyk

Learning from “Big Code”



Key Idea:

Learn a function f that explains the data. The function dynamically obtains the best conditioning context for a given query.



<http://plml.ethz.ch/>