

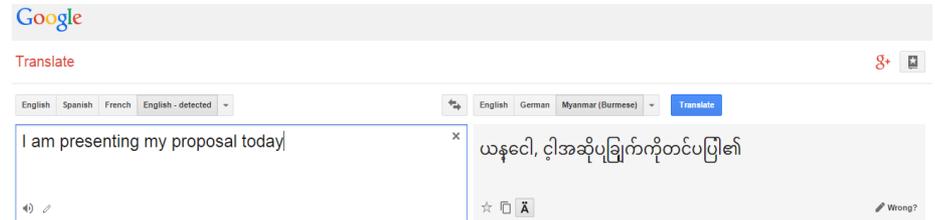
# Machine Learning for Programming

Martin Vechev

ETH Zurich and DeepCode.ai

# Learning and probabilistic models based on Big Data have revolutionized entire fields

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

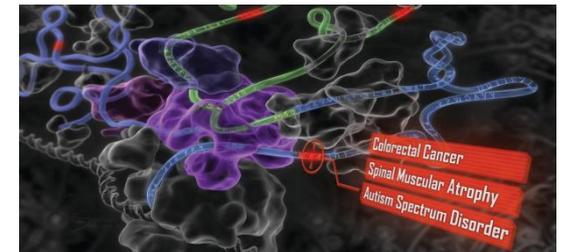
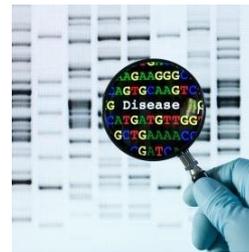


Computer Vision  
(e.g., image captioning)



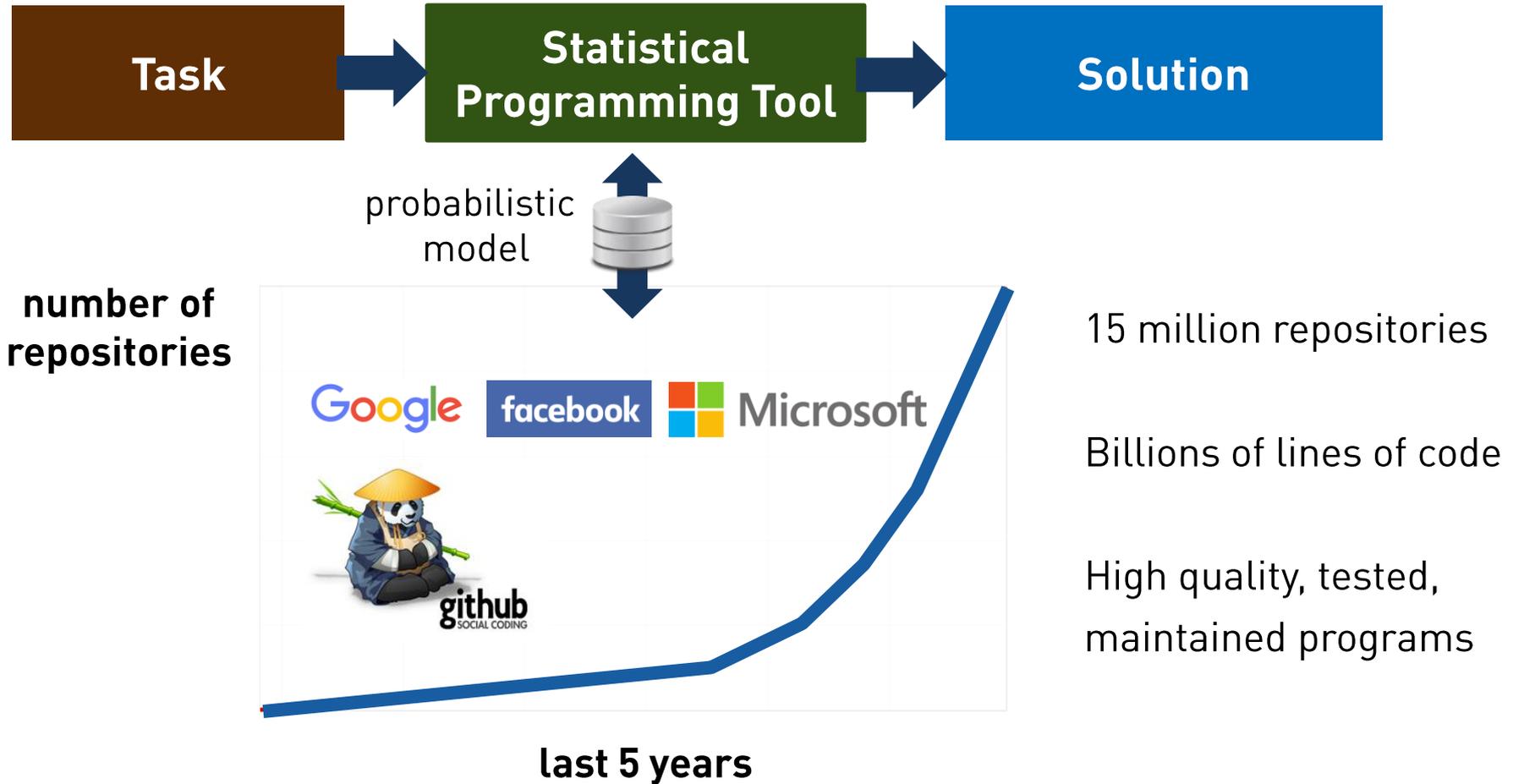
A group of people shopping at an outdoor market. There are many vegetables at the fruit stand.

Medical Computing  
(e.g., disease prediction)



## Can we bring this revolution to programmers?

# Machine Learning for Programs



# Why now?

**Advances in Programming Languages**  
[Automated Reasoning, Synthesis, Constraint Solving]

Confluence of streams

**Advances in Machine Learning**  
[Deep Learning, Graphical Models, Language Models]

**Data**  
[> 15 million public repositories]



**machine learning-based  
programming tools**  
new rules, new ideas, new opportunities

# Machine Learning for Programs

Probabilistically likely solutions to problems hard to solve otherwise

## Joint work with :



Veselin  
Raychev



Andreas  
Krause



Pavol  
Bielik



Christine  
Zeller



Svetoslav  
Karaivanov



Pascal  
Roos



Benjamin  
Bichsel



Timon  
Gehr



Petar  
Tsankov



Mateo  
Panzacchi

## Publications

- Program Synthesis for Character Level Language Modeling, **ICLR'17**
- Learning a Static Analyzer from Data, **CAV'17**
- Statistical Deobfuscation of Android Applications, **ACM CCS'16**
- Probabilistic Mode for Code with Decision Trees, **ACM OOPSLA'16**
- PHOG: Probabilistic Mode for Code, **ACM ICML'16**
- Learning Programs from Noisy Data, **ACM POPL'16**
- Predicting Program Properties from “Big Code”, **ACM POPL'15**
- Code Completion with Statistical Language Models, **ACM PLDI'14**
- Machine Translation for Programming Languages, **ACM Onward'14**

## Statistical Engines

[apk-deguard.com](http://apk-deguard.com)



DEEP3

[jsnice.org](http://jsnice.org)



SLANG

[nice2predict.org](http://nice2predict.org)

NICE 2 Predict

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

# Dimensions:

## Machine Learning for Programming

Applications

Intermediate  
Representation

Analyze Program  
(PL)

Train Model  
(ML)

Query Model  
(ML)

# Dimensions:

## Machine Learning for Programming

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

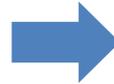
# Statistical Code Synthesis

[Code Completion with Statistical Language Models, ACM PLDI 2014

Veselin Raychev, Eran Yahav, M.V.]

```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);
```

SLANG



```
Camera camera = Camera.open();  
camera.setDisplayOrientation(90);  
  
camera.unlock();  
SurfaceHolder holder = getHolder();  
holder.addCallback(this);  
holder.setType(SurfaceHolder.STP);  
MediaRecorder r = new MediaRecorder();  
r.setCamera(camera);  
r.setAudioSource(MediaRecorder.AS);  
r.setVideoSource(MediaRecorder.VS);  
r.setOutputFormat(MediaRecorder.MPEG4);
```

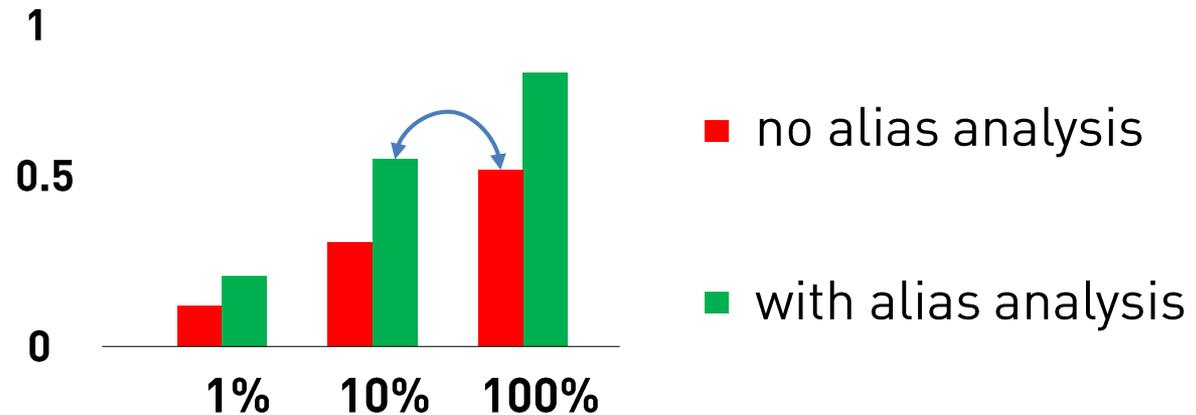
Statistical language models  
Recurrent neural networks

+

Typestate analysis  
Alias analysis

# Importance of Static (Semantic) Analysis

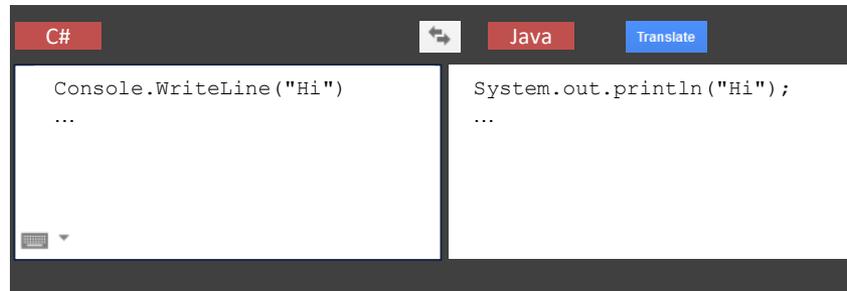
precision vs. % of data used



static analysis benefit = 10x more data

# Statistical Language Translation

[ACM Onward'14, Svetoslav Karaivanov, Veselin Raychev, M.V.]



The screenshot shows a code editor interface with two panes. The left pane is labeled 'C#' and contains the code: `Console.WriteLine("Hi");` followed by an ellipsis. The right pane is labeled 'Java' and contains the code: `System.out.println("Hi");` followed by an ellipsis. A 'Translate' button is visible in the top right corner of the editor.

Phrase-based Statistical  
Machine Translation

+

Prefix Parsing of  
Context-Free Grammars

# Statistical Program de-obfuscation

Predicting program properties from “Big Code”  
Veselin Raychev, Andreas Krause, M.V.

ACM POPL'15

SIGPLAN and ACM Research Highlight

```
function FZ(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;  
}
```



```
function chunkData(str, step) {  
var colNames = [];  
var len = str.length;  
var i = 0;  
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;  
}
```

# JSNice.org

JS NICE

✓ Every country

✓ 200,000 users

✓ Top ranked tool



 **This Page Amsterdam** @thispage\_ams · Jul 16  
Do you write ugly JavaScript code? Not to worry. JSNice will make it look like you are a **superstar** coder. Yay! - [buff.ly/1HR4JL7](http://buff.ly/1HR4JL7)

 **Ingvar Stepanyan** @RRverser · Aug 6  
JSNice.org became my **must-have** tool for code deobfuscation.  
Expand ↩ Reply ↻ Retweet ★ Favorite ⋮ More

 **Brevity** @seekbrevity · Jul 28  
JSNice is an **amazing** tool for de-minifying #javascript files. JSNice.org, its great for #learning and reverse engineering.  
Expand ↩ Reply ↻ Retweet ★ Favorite ⋮ More

 **Alvaro Sanchez** @alvasavi · Jun 10  
**This is gold.**  
Statistical renaming, Type inference and Deobfuscation.  
[jsnice.org](http://jsnice.org)  
Expand ↩ Reply ↻ Retweet ★ Favorite ⋮ More

 **Alex Vanston** @mrvdot · Jun 7  
I've been looking for this for years: JS NICE [buff.ly/1pQ5qfr](http://buff.ly/1pQ5qfr) #javascript #unminify #deobfuscate #makeitReadable  
Expand ↩ Reply ↻ Retweet ★ Favorite ⋮ More

 **Kamil Tomšik** @cztomsik · Jun 6  
tell me **how this** works!  
de-minify #jquery #javascript incl. args, vars & #jsdoc impressive! [jsnice.org](http://jsnice.org)



COMMENTS  
MOST RECENT

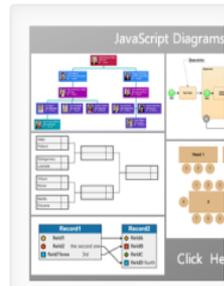
Frici Komez on:  
20 Online Code Editors and Tools for Developers

HOME WORDPRESS WEB DESIGN TYPOGRAPHY RESOURCES ARCHIVES CONTACT ADVERTISE

## 20 Essential Tools for Coders

BY GAVIN IN DEVELOPMENT RESOURCES - 28 JUL, 2014

### 20 Best Freebies for Web Designers of Year 2014



After spending a lot of time in a particular field, we normally get a lot of experience and we suddenly start doing things easily and in short span of time to maximize our efforts towards our goal. This same procedure also holds true in case of web designing also. Because it's all about getting an experience in a particular field or language. Learning web designing is not an easy task, but when you hold experience in the same field, it becomes much easier for you to do most of the work done in very short span of time.

AUGUST 28, 2014

#### JavaScript迷宮救星 — JS NICE



我們網頁前端工程師常常要在網站上用瀏覽器debug開發網站上的JavaScript檔案時，通常那些檔案都已經經過極小化、最佳化，還有被加密了，以確保檔案讀取的安全，並防止他人輕易盜用程式碼。不過一旦要除錯，就苦了前端工程師了。

其實用JS的開發器如Chrome、Safari跟Firefox都支援美化程式碼的功能，讓極小化程式碼變成有條理的狀況，但遇到複雜的檔案時，面對一堆變數名稱如a、b、c、d的程式，還是會覺得很痛苦。

不過網絡上有些神奇的工具可以幫助我們不僅是把變成一團亂碼的變數的，甚至它會替我們選擇出已經過時的變數名稱！

那你是怎麼編譯介紹JS NICE網站。

這個網站的介紹超級乾淨，就像下圖一樣，左邊是複製貼上的JS原始碼，只要按下NICYIFY JAVASCRIPT按鈕，經過除錯與變換變數的結果就會出現在右邊了！

JSnice is tool that help you to make even more obfuscated JavaScript code readable. JSnice is Static Renaming, Type Inference and De-obfuscation.

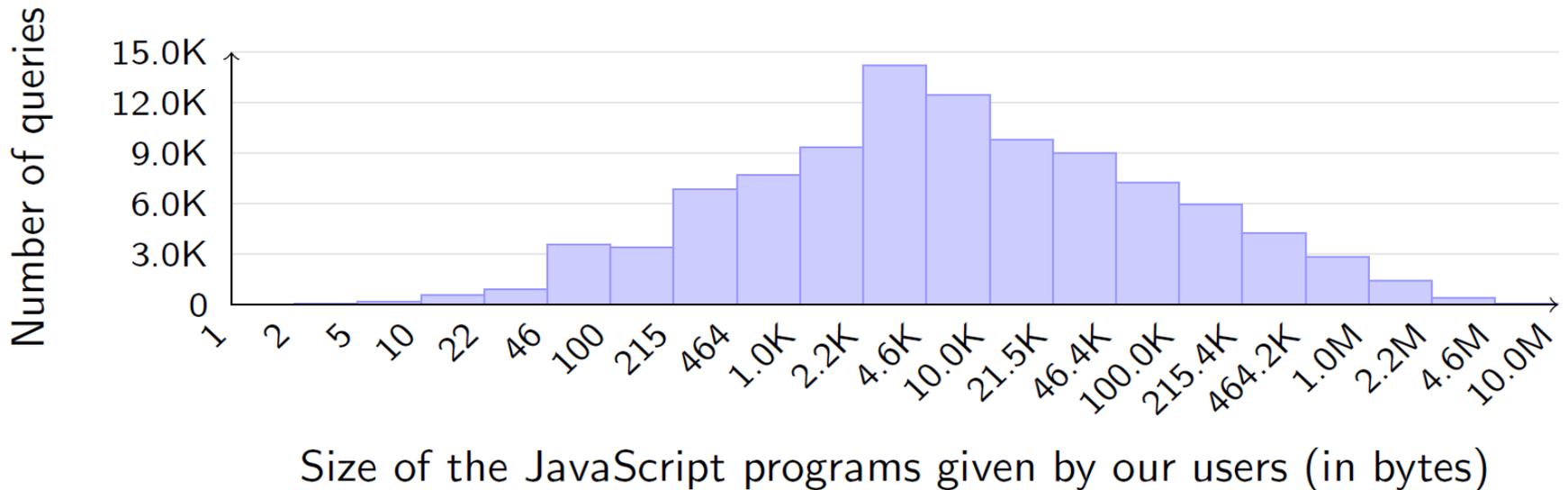
right place, I have gathered 20 Essential JavaScript development tasks. Following dev beautify, store your snippet and h you will find the list handy and acc

### 1. JS Nice



# JSNice: over a year of usage

May 2015 – May 2016



Total: 12GB of code

# JSNice

```
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;
```



```
function chunkData(str, step)
  var colNames = [];
  var len = str.length;
  var i = 0;
  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;
```

## Key Challenges

Facts to be predicted are **dependent**

Prediction should be **fast** (real-time)  
(hard, **millions of possible choices**)

# Key Idea

Phrase the problem of predicting program facts as

**Structured Prediction for Programs**

# MAP inference

```
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;
```

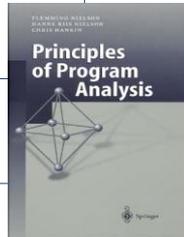
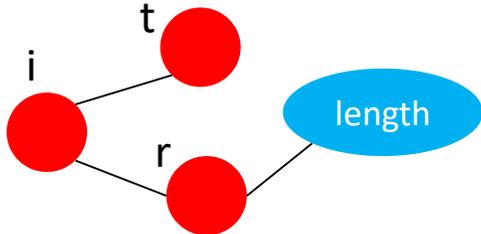
# MAP inference

```
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;
```



Unknown facts: t r i ...

Known facts: length ...



# MAP inference

```

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;

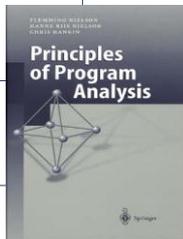
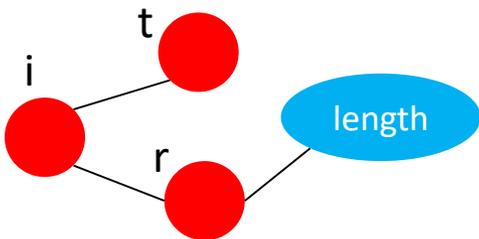
```



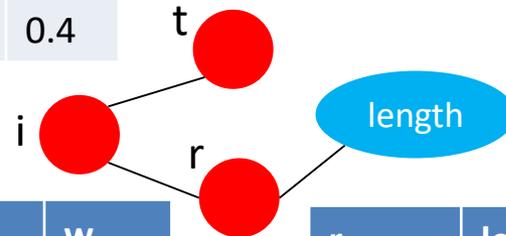
Unknown facts:



Known facts:



i	t	w
i	step	0.5
j	step	0.4



i	r	w
i	len	0.6
i	length	0.3

r	length	w
length	length	0.5
len	length	0.3

# MAP inference

```

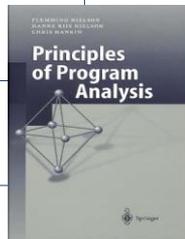
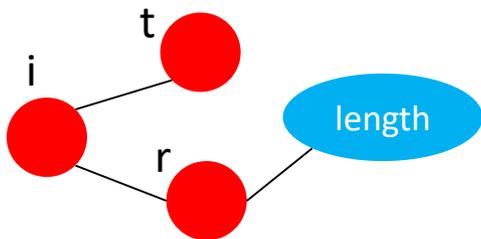
function chunkData(e, t)
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t)
    if (i + t < r)
      n.push(e.substring(i, i + t));
    else
      n.push(e.substring(i, r));
  return n;
  
```



Unknown facts:

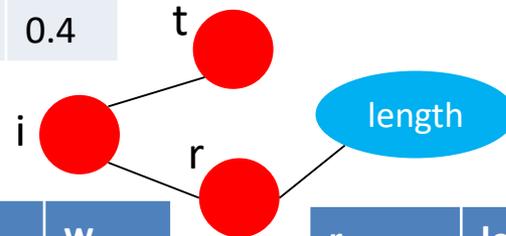


Known facts:



$$\text{argmax } \mathbf{w}^T \mathbf{f}(i, t, r, \text{length})$$

i	t	w
i	step	0.5
j	step	0.4



i	r	w
i	len	0.6
i	length	0.3

r	length	w
length	length	0.5
len	length	0.3

# MAP inference

```

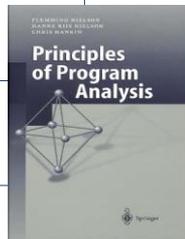
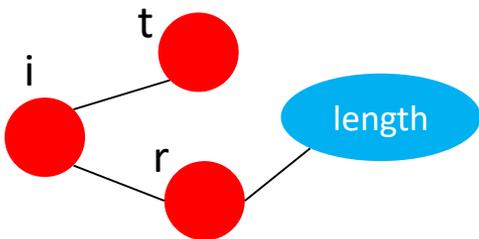
function chunkData(e, t)
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t)
    if (i + t < r)
      n.push(e.substring(i, i + t));
    else
      n.push(e.substring(i, r));
  return n;
  
```



Unknown facts:



Known facts:



$\text{argmax } \mathbf{w}^T \mathbf{f}(i, t, r, \text{length})$

i	t	w
i	step	0.5
j	step	0.4

i	r	w
i	len	0.6
i	length	0.3

r	length	w
length	length	0.5
len	length	0.3

# MAP inference

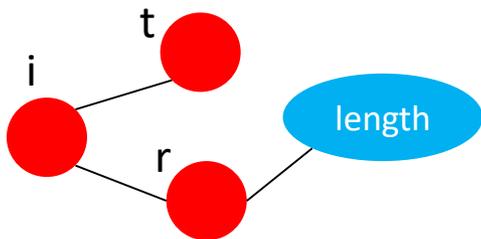
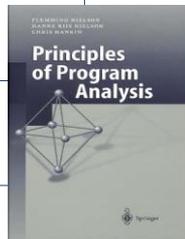
```
function chunkData(e, t)
  var n = [];
  var r = e.length;
  var i = 0;
  for (; i < r; i += t)
    if (i + t < r)
      n.push(e.substring(i, i + t));
    else
      n.push(e.substring(i, r));
  return n;
```

```
function chunkData(str, step)
  var colNames = [];
  var len = str.length;
  var i = 0;
  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;
```

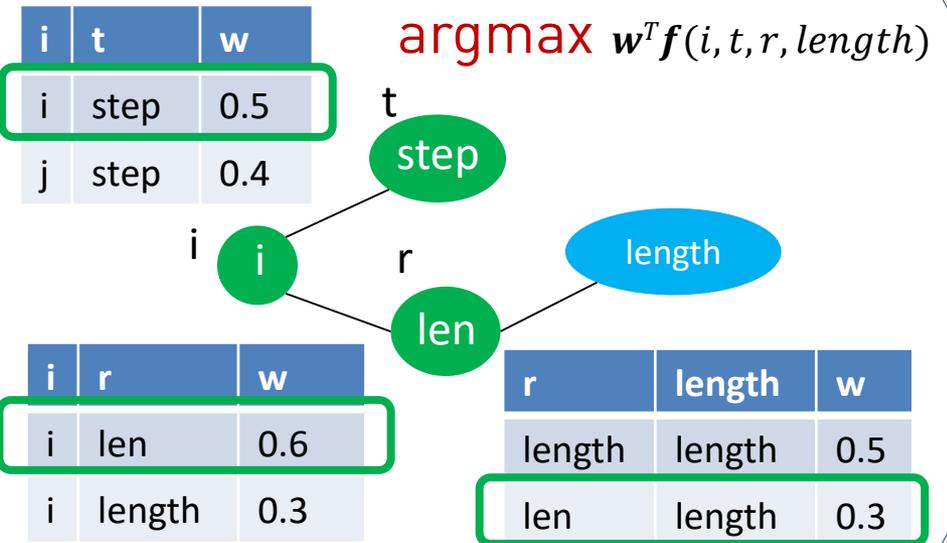
Unknown facts:



Known facts:



$\text{argmax } w^T f(i, t, r, \text{length})$



# A glimpse of learning...

[N. Ratliff, J. Bagnell, M. Zinkevich, AISTATS 2007]

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp(\mathbf{w}^T \mathbf{f}(\mathbf{y}, \mathbf{x}))$$

Given a data set:  $D = \{ \mathbf{x}^j, \mathbf{y}^j \}_{j=1..n}$  learn weights  $\mathbf{w}^T$

Optimization objective (max-margin training):

$$\forall j \forall \mathbf{y} \underbrace{\sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y}^{(j)})}_{\text{Given prediction}} \geq \underbrace{\sum w_i f_i(\mathbf{x}^{(j)}, \mathbf{y})}_{\text{any other prediction}} + \underbrace{\Delta(\mathbf{y}, \mathbf{y}^{(j)})}_{\text{margin}}$$

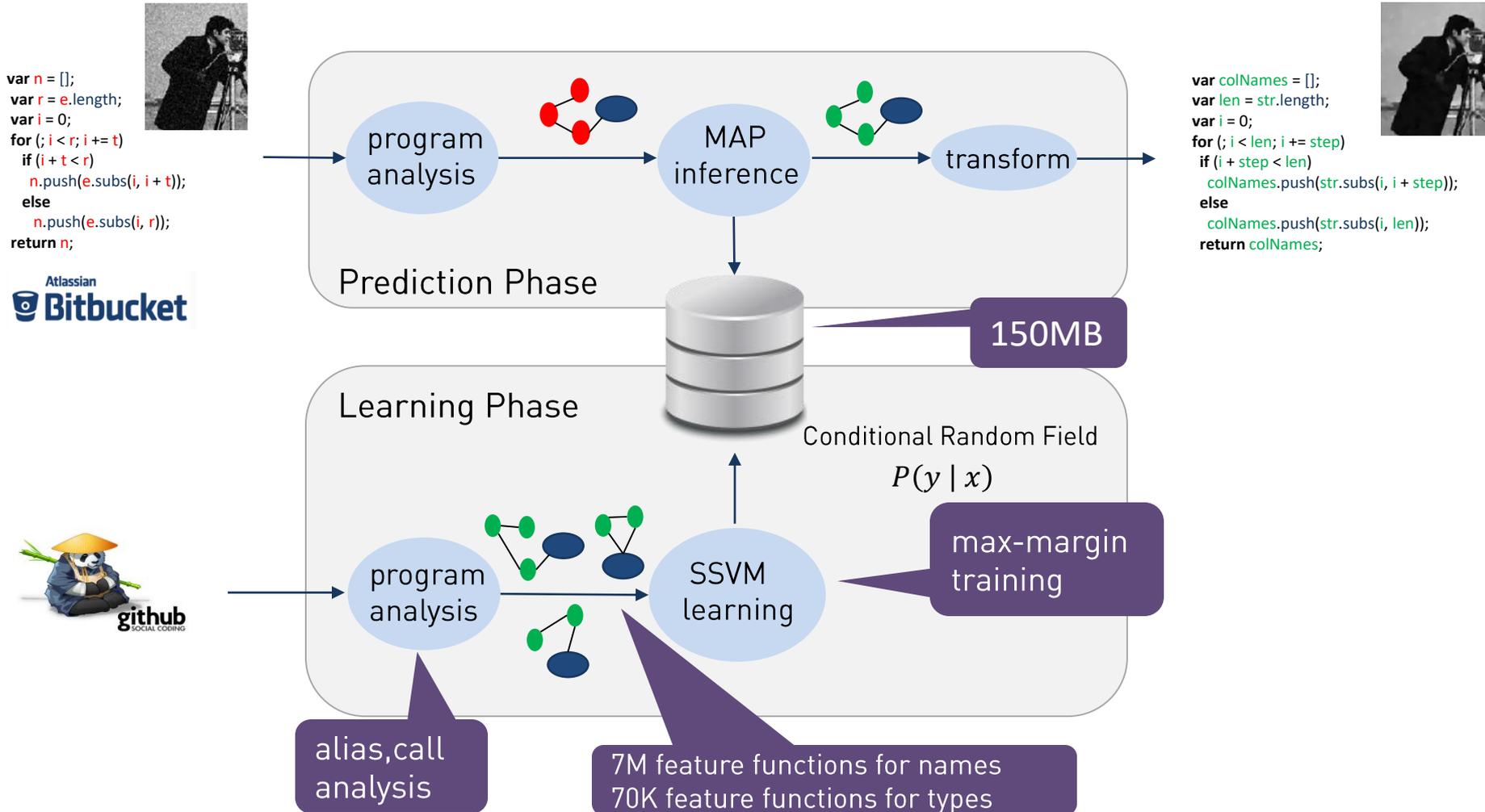
for all samples

Given prediction is better than any other prediction by a margin

Avoids expensive computation of the partition function  $Z(\mathbf{x})$

# JSNice

[Predicting program properties from Big Code, ACM POPL 2015]



# Statistical Program de-obfuscation

[ACM CCS'16, Benjamin Bichsel, Petar Tsankov, Veselin Raychev, M.V.]

```
package a.b.c
class a extends SQLiteHelper
{
    SQLiteDatabase b; public
a(Context ctx) { b =
getWritableDatabase(); }
Cursor c(String str) {
return b.rawQuery(str); }}
```



```
package com.example.dbhelper

class DBHelper extends
SQLiteHelper {
    SQLiteDatabase db;

    public DBHelper(Context ctx) {
        db = getWritableDatabase();
    }

    Cursor execSQL(String str) {
        return db.rawQuery(str);
    }
}
```



## Funny Reddit post/comment



[–] **Tycon712** • 3 points 2 days ago



Can someone tell me **what the point of using Proguard** is if there are tools out there like this?

[permalink](#) [embed](#) [pocket](#)



[–] **theheartbreakpug** • 6 points 2 days ago



As far as I know, this is brand new. I asked the creator of ProGuard a week ago how hard it is to unobfuscate code after it's run through proguard. He said it strips all the names out of the code so it's essentially impossible. **I'm super impressed by what they've done here.**

# Nice2Predict.org:

scalable structured prediction framework

[Pavol Bielik, Veselin Raychev, Matteo Panzacchi, Nick Baumann, M.V.]

fully, **open sourced**,  
Apache license

used by **various**  
**groups worldwide**

JS **NICE**

**DEGUARD**

Your System  
Here

**NICE** 2 Predict

Fast, Approximate  
MAP inference

Fast, Parallel, Structured SVM  
and Pseudo-Likelihood Training

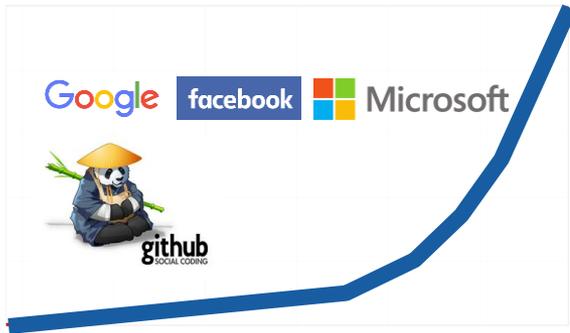
Arbitrary factors and  
indicator functions

# Fundamental Problem

Data

Learning

Model



Probabilistic  
Model



Widely  
Applicable

Efficient  
Learning

High  
Precision

Explainable  
Predictions

## Training dataset ***D***

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

query:

```
f.open("file" | "r");  
f.?
```

## Training dataset $D$

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

query:

```
f.open("file" | "r");  
f. ?
```

### 3-gram model on tokens

Hindle et. al., ACM ICSE'12

P(open		f. )	~ 3/6
P(read		f. )	~ 2/6
P(write		f. )	~ 1/6

  
context  $\gamma$

## Training dataset $D$

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

query:

```
f.open("file" | "r");  
f. open
```

### 3-gram model on tokens

Hindle et. al., ACM ICSE'12

P(open		f. )	~ 3/6
P(read		f. )	~ 2/6
P(write		f. )	~ 1/6

  
context  $\gamma$

## Training dataset $D$

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

## probabilistic model on APIs

Raychev et. al., ACM PLDI'14

```
P(read | open) ~ 2/3  
P(write | open) ~ 1/3
```

context  $\gamma$

query:

```
f.open("file" | "r");  
f. ?
```

## 3-gram model on tokens

Hindle et. al., ACM ICSE'12

```
P(open | f.) ~ 3/6  
P(read | f.) ~ 2/6  
P(write | f.) ~ 1/6
```

context  $\gamma$

## Training dataset $D$

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

## probabilistic model on APIs

Raychev et. al., ACM PLDI'14

$$P(\text{read} \mid \text{open}) \sim 2/3$$
$$P(\text{write} \mid \text{open}) \sim 1/3$$

context  $\gamma$

query:

```
f.open("file" | "r");  
f.read
```

## 3-gram model on tokens

Hindle et. al., ACM ICSE'12

$$P(\text{open} \mid \text{f.}) \sim 3/6$$
$$P(\text{read} \mid \text{f.}) \sim 2/6$$
$$P(\text{write} \mid \text{f.}) \sim 1/6$$

context  $\gamma$

## Training dataset $D$

```
f.open("f2" | "r");  
f.read();
```

```
f.open("f2" | "w");  
f.write("c");
```

```
f.open("f1" | "r");  
f.read();
```

## probabilistic model on APIs

Raychev et. al., ACM PLDI'14

```
P(read | open) ~ 2/3  
P(write | open) ~ 1/3
```

context  $\gamma$

query:

```
f.open("file" | "r");  
f. ?
```

## 3-gram model on tokens

Hindle et. al., ACM ICSE'12

```
P(open | f.) ~ 3/6  
P(read | f.) ~ 2/6  
P(write | f.) ~ 1/6
```

context  $\gamma$

What should the context be?



“...All problems in computer science can be solved by **another level of indirection...**”

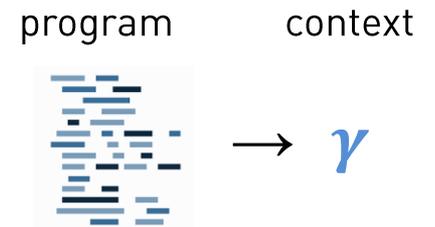
-- David Wheeler



“...All problems in computer science can be solved by **another level of indirection...**”

-- David Wheeler

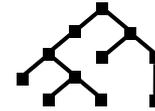
key idea: **synthesize a function**  $f$ :



# Creating probabilistic models: our method

[“Learning Programs from Noisy Data”, ACM POPL’16, “PHOG: Probabilistic Model for Code”, ACM ICML’16,  
“Probabilistic Model for Code with Decision Trees”, ACM OOPSLA’16, “Synthesis for Char. Level Language Modeling, ICLR’17]

1. **Pick** a structure of interest, e.g., ASTs:



2. **Define** a DSL for expressing functions:  
(can be Turing complete)

DSL

3. **Synthesize**  $f_{best} \in \text{DSL}$  from Dataset  $\mathbf{D}$ :

$$f_{best} = \underset{f \in \text{DSL}}{\operatorname{argmin}} \operatorname{cost}(\mathbf{D}, f)$$

4. **Use**  $f_{best}$  to compute context and predict:

$$f_{best} \left( \text{AST} \right) \rightarrow \gamma$$

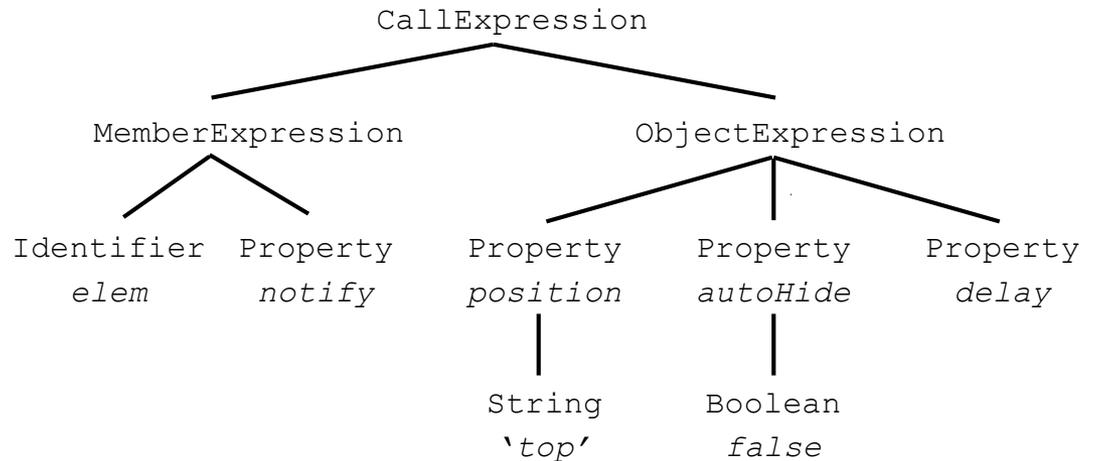
# Step 1: Pick Structure of Interest

Let it be abstract syntax trees (ASTs) of programs

## JavaScript program

```
elem.notify({  
  position: 'top',  
  autoHide: false,  
  delay: 100  
});
```

## AST



# Step 2: Define a DSL over structure

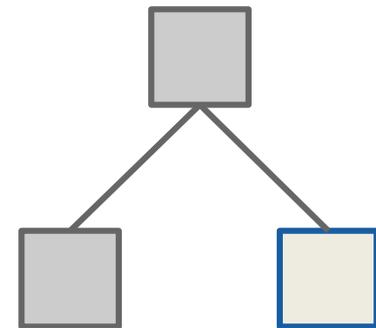
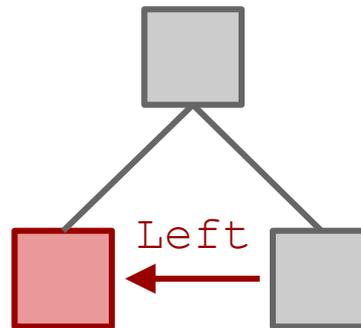
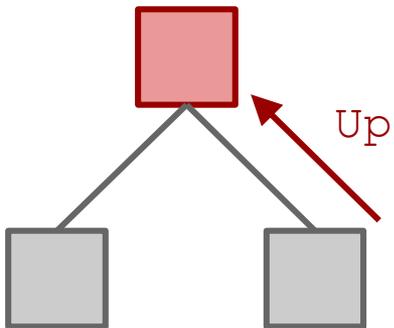
## Syntax

TCond ::=  $\varepsilon$  | WriteOp TCond | MoveOp TCond

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

WriteOp ::= WriteValue, WriteType, WritePos

## Semantics



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

Step 3: synthesize  $f_{best}$

$$f_{best} = \underset{f \in \text{DSL}}{\operatorname{argmin}} \operatorname{cost}(\mathbf{D}, f)$$

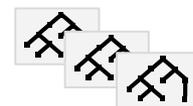
# Step 3: synthesize $f_{best}$

DSL

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

generate candidate  $f$

dataset  $D$



millions ( $\approx 10^8$ )

Synthesizer

$$f_{best} = \operatorname{argmin}_{f \in \text{DSL}} \operatorname{cost}(D, f)$$

Build Probabilistic Model  $P$

use  $D$  and  $f$  to compute

$$P(\text{element} \mid f(\text{tree}))$$

to scale: iterative synthesis on fraction of examples

$$\operatorname{cost}(D, f) = \operatorname{entropy}(P)$$

$$O(|D|)$$

# Step 4: use $f_{best}$ to predict

program

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

$f_{best}$

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



{Previous Property, Parameter Position, API name}

# Deep3: Experimental Results

[Probabilistic Model of JavaScript]

Dataset **D**: 150,000 files    Training Time: ~**100** hours     $f_{best} \sim$  **50,000** instr.

## Probabilistic Model

---

## Accuracy (APIs)

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

22.2%

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

30.4%

*Deep3*

**66.6%**

Details in: "Probabilistic Model for Code with Decision Trees", ACM OOPSLA'16

# Applying the Concept to Natural Language

[Program Synthesis for Character Level Language Modeling, **ICLR'17**]

Dataset **D**:

Hutter Prize Wikipedia Dataset

Training Time: ~ **8** hours

$f_{best} \sim$  **9,000** instr

uses a char-level DSL with state

Interpretable model, browse here:

<http://www.srl.inf.ethz.ch/charmmodel.html>

## Probabilistic Model

## Bits-per-Character

7-gram (best)	1.94
Stacked LSTM (Graves 2013)	1.67
Char-based DSL synthesis	1.62
MRNN (Sutskever 2011)	1.60
MI-LSTM (Wu et al. 2016)	1.44
HM-LSTM* (Chung et al. 2016)	1.40

# Learning a Static Analyzer from Data

[Pavol Bielik, Veselin Raychev, M.V., CAV'17]

```
function isBig(v) {  
  return v < this.length  
}  
[12, 5].filter(isBig);
```

```
VarPtsTo("global", h)  
checkIfInsideMethodCall  
checkMethodCallName  
checkReceiverType  
checkNumberOfArguments ...
```

---

```
VarPtsTo(this, h)
```

Can be understood by experts

Found issues in Facebook's Flow



# More Resources

**Learning from Large Codebases,**

**PhD Thesis, ETH Zurich, 2016**

**ACM Doctoral Dissertation, Honorable Mention Award**



Veselin  
Raychev

- <http://plml.ethz.ch>
- Dagstuhl Seminar on Big Code Analytics, Nov 2015
- Data sets, tools, challenges: <http://learningfrombigcode.org>

# Summary

## OPPORTUNITY

**Advances in Programming Languages**  
[Automated Reasoning, Synthesis, Constraint Solving]

**Advances in Machine Learning**

[Deep Learning, Graphical Models, Language Models]

**Data**   
[> 15 million public repositories]

**Learning-based programming tools**  
new rules, new ideas, new opportunities

## RICH PROBLEM SPACE

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

## NEW PROBABILISTIC MODELS

1. **Pick** a structure of interest, e.g., trees:



2. **Define** a DSL for expressing functions:  
(can be Turing complete)

```
TCond ::= ε | WriteOp TCond | MoveOp TCond
MoveOp ::= Up, Left, Right, DownFirst, DownLast,
NextDFS, PrevDFS, NextLeaf, PrevLeaf,
PrevNodeType, PrevNodeValue,
WriteOp ::= WriteValue, WriteType, WritePos
```

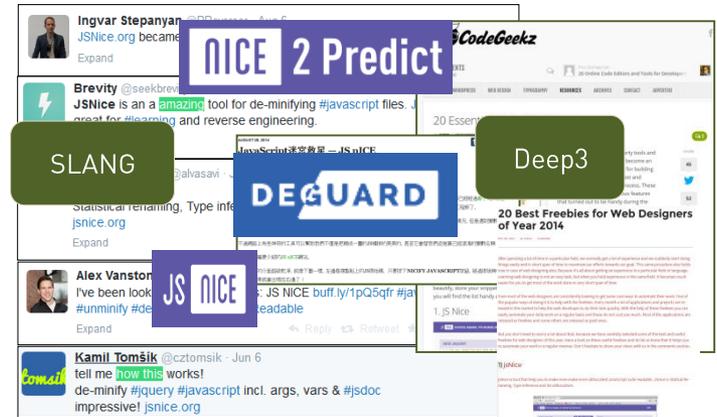
3. **Synthesize**  $f_{best} \in \text{DSL}$  from Dataset  $D$ :

$$f_{best} = \operatorname{argmin}_{f \in \text{DSL}} \operatorname{cost}(D, f)$$

4. **Use**  $f_{best}$  on new structures:

$$f_{best} (\text{tree}) \rightarrow \gamma$$

## PRACTICAL IMPACT



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