

# Probabilistic Model for Code with Decision Trees



**Veselin Raychev, Pavol Bielik, Martin Vechev**

Software Reliability Lab  
Department of Computer Science  
ETH Zurich

# Big Data Revolution

Research area

Big Data

Application

Computer  
vision



Image labeling



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

Programming  
languages



?

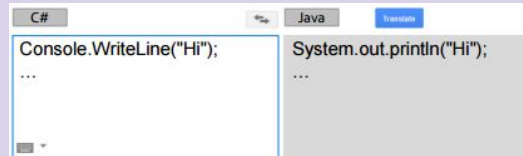
# Statistical programming tools



**Write new code**  
**[PLDI'14, POPL'16]:**  
Code Completion

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

**Port code [ONWARD'14]:**  
Programming Language  
Translation



**Understand code/security [POPL'15]:**  
Deobfuscation  
Type Prediction



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

Other groups

Natural language to code  
Automatic patch generation  
Bug finding  
Code completion  
others...

# 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

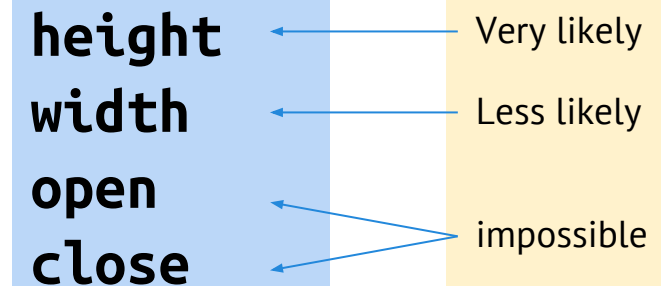
Model is a key part of Statistical Programming Tools

Goal: score programs

Select best among several candidates

Example:

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



# Model usability

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

# Existing works: naive models

## Most common model: n-gram

Training (3-gram model):

```
return a . width * a .  
    + a . open (mode) ;  
+ a . close () ;  
    * a . height ;  
    * a . close () ;
```

*Note: A blue bracket above the first line of code spans from 'a' to 'width' and is labeled '3-gram'.*

Count 3-grams

e.g **a.close** (2 times)

Prediction:

width + a .

a.close - 2 times

a.open - 1 times

a.width - 1 times

a.height - 1 times

Probability

close	0.4
open	0.2
width	0.2
height	0.2

*Note: A blue arrow points from the 'Probability' label to the top of the table. Four blue arrows point from the prediction counts on the left to the corresponding rows in the table.*

# Existing works: naive models

## Most common model: n-gram

Training (3-gram model):

Prediction:

**Main problem:**

**Bad context** leads to bad probability estimates

width +

a .

Probability

a.close - 2 times

a.open - 1 times

a.width - 1 times

a.height - 1 times

close

open

width

height

0.4

0.2

0.2

0.2



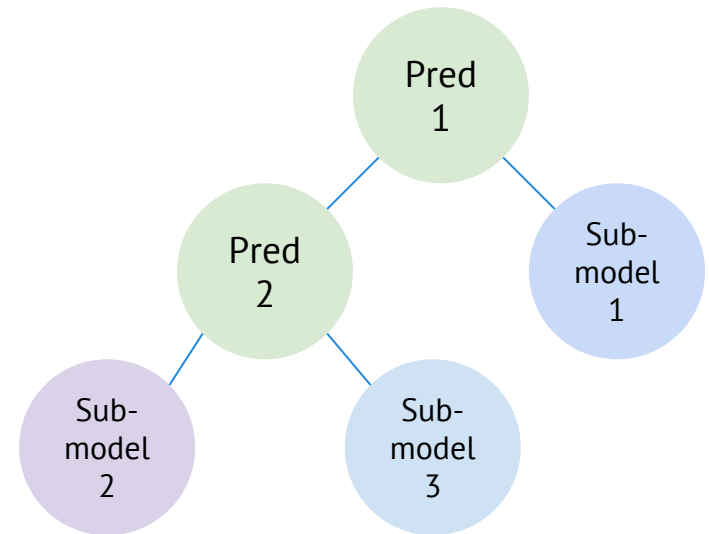
## 1. Richer conditioning context

return a.**width** \* a. ?



## 3. Evaluation


## 2. Domain Specific Language encoding Decision Trees



Lang ::= BasicPart | BranchPart

BranchPart ::= if pred(x) then Lang else Lang

# Conditioning context

Training (3-gram model):

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

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

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

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

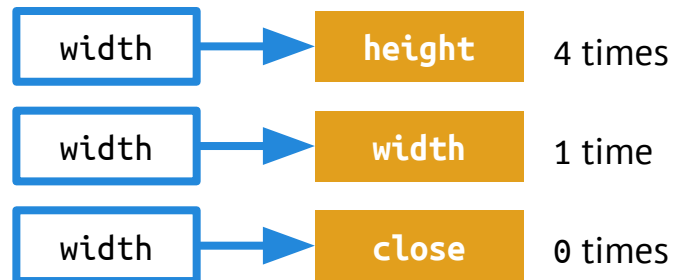
q.depth \* q.**width** \* q.**height**

Prediction:

return a.**width** \* a.**?**

Context: relevant for this prediction

height	0.8
width	0.2



# Indirection

Query:  $x$

```
return a.width * a.?
```

Completion:  $y$

```
? = height
```

Context:  $ctx = f(x)$

```
return a.width * a.?
```



Model:  $P(y|f(x))$



# Example contexts

Query: x

return a.width \* a. ?

ctx = f(x): last two tokens  
3-gram language model

return a.width \* a. ?

ctx = f(x): previous actions  
on the same object [PLDI'14]

return a.width \* a. ?

ctx = f(x): the third token  
before the completion

Question:

Which context is  
the best?

Try many of them  
and evaluate

# Evaluation metric: entropy

Average entropy:  
~0.66667 bits

console. info	$P(\text{info} \mid \text{console.}) = 1$	0
console. info	$P(\text{info} \mid \text{console.}) = 1$	0
a.width * a. height	$P(\text{height} \mid \text{a.}) = 0.5$	1
b.width * b. height	$P(\text{height} \mid \text{b.}) = 0.5$	1
a.right - a. left	$P(\text{left} \mid \text{a.}) = 0.5$	1
b.bottom - b. top	$P(\text{top} \mid \text{b.}) = 0.5$	1

# Evaluation metric: entropy

Average entropy:  
~1.91 bits

console.`.``info`

$$P(\text{info} | .) = 1/3$$

~1.58

console.`.``info`

$$P(\text{info} | .) = 1/3$$

~1.58

a.width \* a.`.``height`

$$P(\text{height} | .) = 1/3$$

~1.58

b.width \* b.`.``height`

$$P(\text{height} | .) = 1/3$$

~1.58

a.right - a.`.``left`

$$P(\text{left} | .) = 1/6$$

~2.58

b.bottom - b.`.``top`

$$P(\text{top} | .) = 1/6$$

~2.58

# Unconditional model

No conditioning in the probability distribution:  
ctx=⊥

Average entropy:  
~1.91 bits

console. info	$P(\text{info}) = 1/3$	~1.58
console. info	$P(\text{info}) = 1/3$	~1.58
a.width * a. height	$P(\text{height}) = 1/3$	~1.58
b.width * b. height	$P(\text{height}) = 1/3$	~1.58
a.right - a. left	$P(\text{left}) = 1/6$	~2.58
b.bottom - b. top	$P(\text{top}) = 1/6$	~2.58

## 1. Richer conditioning context

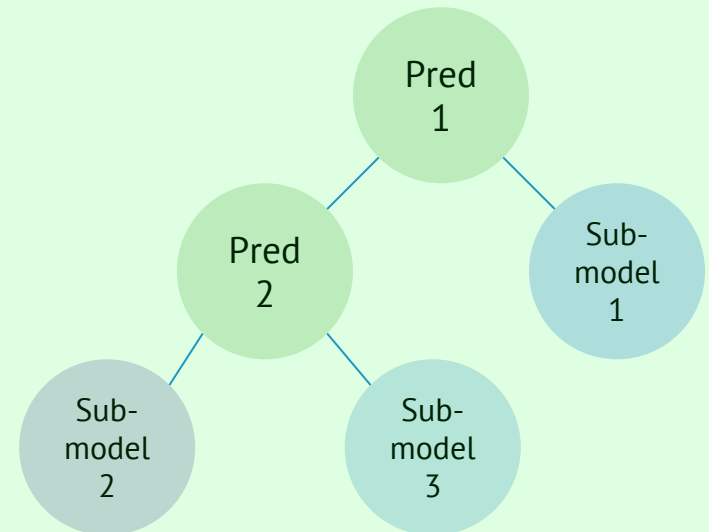
return a.**width** \* a. ?



## 3. Evaluation

Red	Orange	Yellow	Green
Red	Orange	Yellow	Green
Red	Orange	Green	Yellow

## 2. Domain Specific Language encoding Decision Trees



Lang ::= BasicPart | BranchPart

BranchPart ::= if pred(x) then Lang else Lang



# Synthesize the best model

Query:  $x$

```
return a.width * a. ?
```

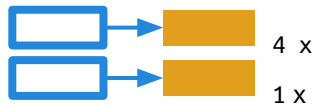
Completion:  $y$

```
? = height
```

Context:  $ctx = f(x)$

```
return a. width * a. ?
```

Model:  $P(y|f(x))$



Find function  $f$

From a domain  
specific language

# Main DSL requirement

**Program** ::= BasicPart | **BranchPart**

**BranchPart** ::= if **pred(x)** then **Program** else **Program**

## Basic part

Includes models from a simple DSL

```
return a.width * a. ?
```

ctx = f(x): last two tokens  
3-gram language model

```
return a.width * a. ?
```

ctx = f(x): previous actions on the same  
object [PLDI'14]

```
return a.width * a. ?
```

ctx = f(x): empty  
Unconditional model

## Branch part

Combines programs based on predicates. Language of predicates **pred**

e.g. are we completing a field name?  
Is there another operation on the same object?

Special empty program  $\epsilon$

# Best program

If there is no prior field access, then use 3-gram model,  
otherwise prior field access

Entropy (bits)

<code>console.</code> <code>info</code>	$P(\text{info} \mid \text{console.}) = 1$	0
<code>console.</code> <code>info</code>	$P(\text{info} \mid \text{console.}) = 1$	0
<code>a.</code> <code>width</code> * <code>a.</code> <code>height</code>	$P(\text{height} \mid \text{a.}) = 1$	0
<code>b.</code> <code>width</code> * <code>b.</code> <code>height</code>	$P(\text{height} \mid \text{b.}) = 1$	0
<code>a.</code> <code>right</code> - <code>a.</code> <code>left</code>	$P(\text{left} \mid \text{a.}) = 1$	0
<code>b.</code> <code>bottom</code> - <code>b.</code> <code>top</code>	$P(\text{top} \mid \text{b.}) = 1$	0

# Synthesis using basic part of DSL

Involved  
See POPL'16

At high level:  
Search through  
**thousands** of  
candidate  
programs that  
describe  
conditioning.

## Basic part

Includes models from a simple DSL

```
return a.width * a. ?
```

ctx = f(x): last two tokens  
3-gram language model

```
return a.width * a. ?
```

ctx = f(x): previous actions on the same  
object [PLDI'14]

```
return a.width * a. ?
```

ctx = f(x): empty  
Unconditional model

# Synthesis of branch part

Space of possible programs:

**BranchPart** ::= if  $\text{pred}(x)$  then **Program** else **Program**

Search  
thousands

Search  
thousands

Search  
thousands

Billions even without nesting

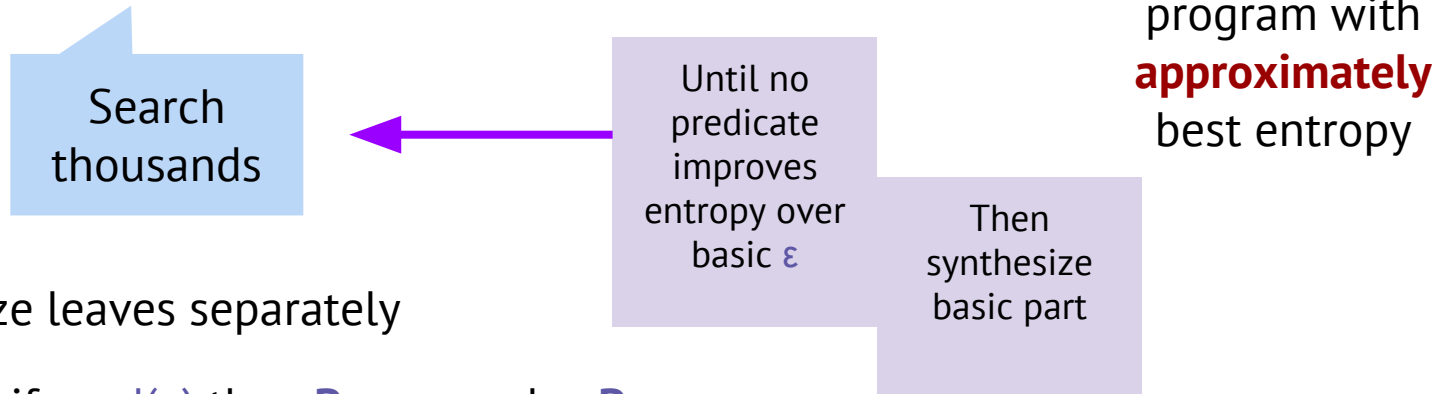
Infinite with nesting

Goal: find  
program with  
best entropy

# Idea 1: synthesis in parts

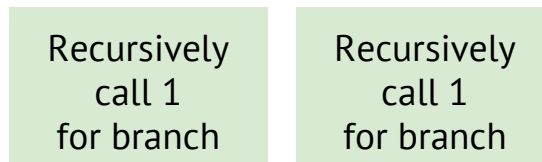
1. Synthesize branch with empty programs in leaves

**Synthesize** : if  $\text{pred}(x)$  then  $\epsilon$  else  $\epsilon$



2. Synthesize leaves separately

**Synthesize** : if  $\text{pred}(x)$  then **Program** else **Program**



# Synthesis of branch part

Space of possible programs:

**BranchPart** ::= if  $\text{pred}(x)$  then **Program** else **Program**

Search  
thousands

Search  
thousands

Search  
thousands

Search for each component separately

Performance at expense of possible  
non-optimality

Goal: find  
program with  
**approximately**  
best entropy

# Idea 2: synthesis in parts

1. Synthesize basic program

$f \in \text{BasicPart}$

Search  
thousands

Goal: find  
program with  
**approximately**  
best entropy

2. Synthesize branch

**Synthesize** : if  $\text{pred}(x)$  then  $f$  else  $f$

Search  
thousands

Until no  
predicate  
improves  
entropy over  
simply  $f$

If no, return  $f$

3. Synthesize leaves separately

**Synthesize** : if  $\text{pred}(x)$  then **Program** else **Program**

Recursively  
call 1  
for branch

Recursively  
call 1  
for branch



# Main result

Synthesis procedure 1 is a new formulation of a known and popular algorithm for decision tree learning: ID3

In fact, we extended ID3 to support programs in decision tree leaves from the `BasicPart` DSL fragment

Synthesis procedure 2 is new  
And also applicable to decision trees

New name:  
**E13**

## 1. Richer conditioning context

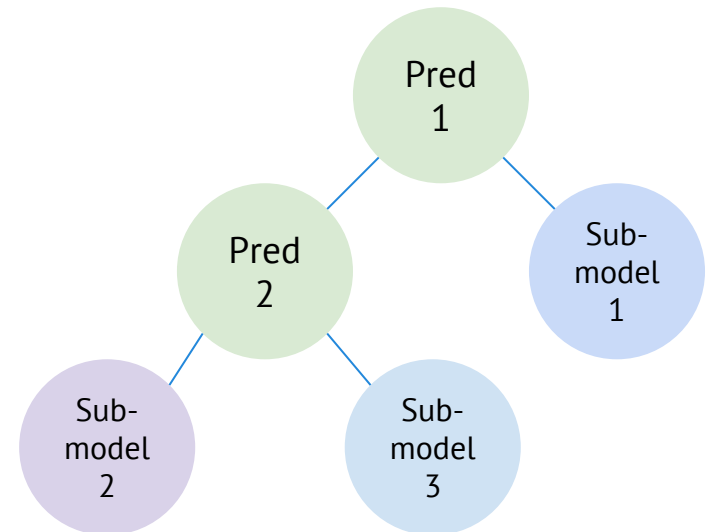
return a.**width** \* a. ?



## 3. Evaluation

tan	light green	yellow	green
tan	light green	yellow	green
tan	light green	green	yellow

## 2. Domain Specific Language encoding Decision Trees



Lang ::= BasicPart | BranchPart

BranchPart ::= if pred(x) then Lang else Lang

# Evaluation

150'000 JavaScript files, from GitHub.com,  
parsed into ASTs.                      Public datasets.

**100'000** files  
Training data  
for synthesis

**50'000** files  
Evaluation  
data

Evaluation files are not on the  
same projects as training

Synthesis time: ~100 hours

Question: How well can we predict program elements?

# Accuracy: JavaScript

Query time for all models is basically the same  
(>10K queries per second)

Task	PCFG	3-gram	ID3+	E13
API completion	0.04%	30.0%	54%	<b>66.6%</b>
Field access completion	3.2%	32.9%	52.5%	<b>67.0%</b>
Predicting loops	0%	37.5%	<b>65.0%</b>	28.3%

Many more evaluation results in the paper.... Also for Python. Easily applicable to all languages.

Big improvement over prior methods

Both ID3+ and E13 learn useful probabilistic models