

Synthesis for interpretable machine learning

Veselin Raychev

veselin@deepcode.ai

DEEPCODE



Spinoff

ETH zürich

Example application: code completion

Task: write an android application that records a video

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

?

Completion based on existing programs

But this program is not in the training data

```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
camera.unlock();  
  
SurfaceHolder h = getHolder();  
h.addCallback(this);  
h.setType(SurfaceHolder.STP);  
  
MediaRecorder r = new MediaRecorder();  
r.setCamera(camera);
```

Need to generate and score solutions:

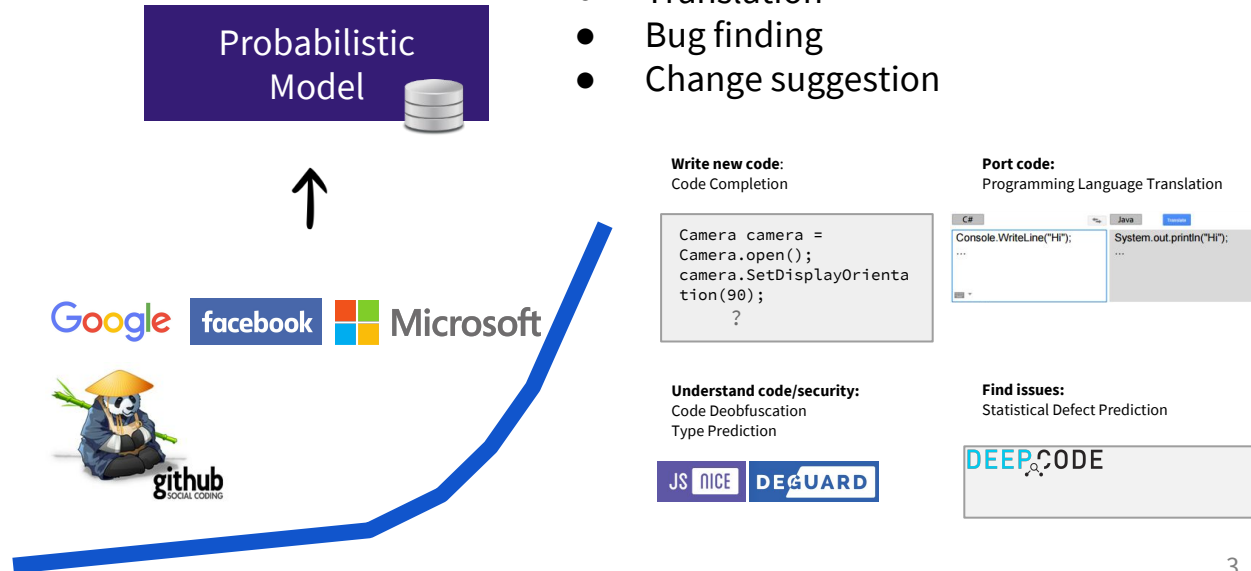
```
Camera camera = Camera.open();  
camera.SetDisplayOrientation(90);  
camera.release();  
camera.release();
```

Probabilistic model for code or text

Can assign probabilities to **programs**, just like a model for **text** assigns probabilities to words and sentences

Can be used as a building block for a variety of systems:

- Code completion
- Translation
- Bug finding
- Change suggestion



Naive approach

N-gram model

Such models are shallow, but
work well with huge datasets

→ English (most likely next word prediction)

Read the file , please !

→ Code (most likely next token prediction)

file . read () ;

Example

query:

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

```
f. ?
```

Training data ***D***

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

```
f.open("f2" , "rw");  
f.write("c");
```

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

query:

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

```
f. ?
```

Training data D

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

```
f.open("f2" , "rw");  
f.write("c");
```

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



3-gram Probabilistic Model P

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

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

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

context γ

Prediction: open

Hindle et. al. [ICSE'12]

query:

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

```
f. open
```

Training data D

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

```
f.open("f2" , "rw");  
f.write("c");
```

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



Prior API Probabilistic Model P

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

context γ

3-gram Probabilistic Model P

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

context γ



Better prediction: read

Raychev et. al. [PLDI'14]

query:

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


Different conditioning **contexts** lead to different probabilistic models

$$P(\text{prediction} \mid \mathbf{context})$$

Manually providing contexts is hard, error-prone and suboptimal !

How to find the best context?

context γ

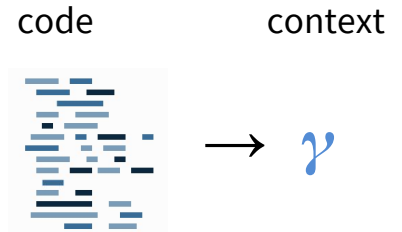
Prior AP

P(read

P(wri

Synthesis approach

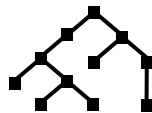
key idea: **synthesize** a **function** f :



Program synthesis: techniques for learning short programs from examples.

Creating a probabilistic model

1. Pick data representation:



2. Define a domain-specific language for expressing functions:

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 contexts, learn and predict:

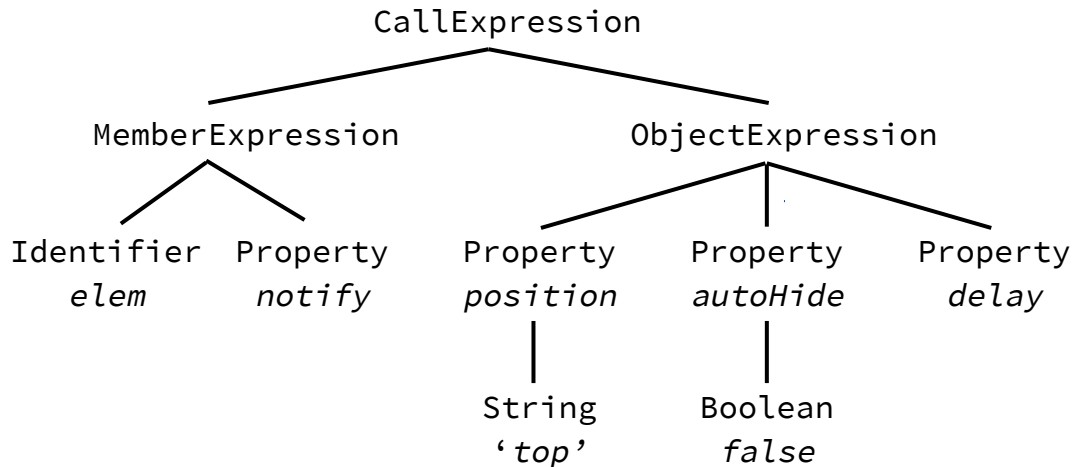
$$f_{best} \left(\begin{array}{c} \text{tree diagram} \end{array} \right) \rightarrow \gamma$$

Step 1: Pick data representation (for code)

JavaScript program

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

AST



Step 2: Define domain-specific language

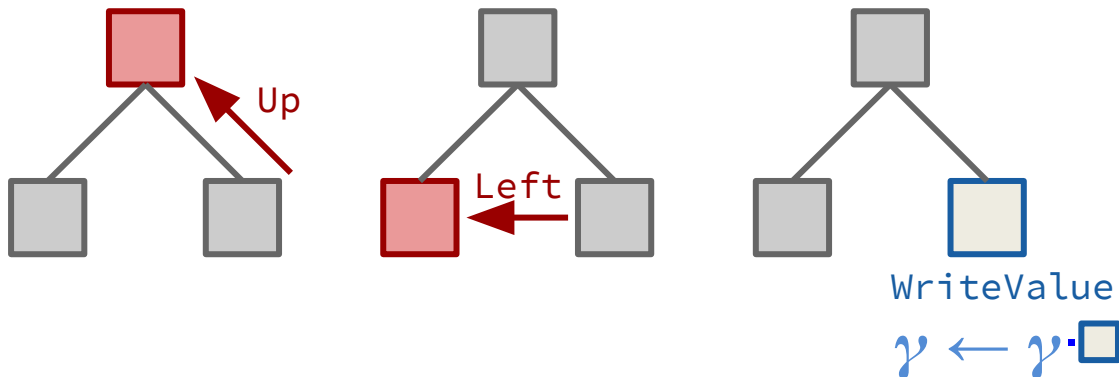
Syntax

TCond ::= ε | WriteOp TCond | MoveOp TCond |
if TCond == const TCond else TCond

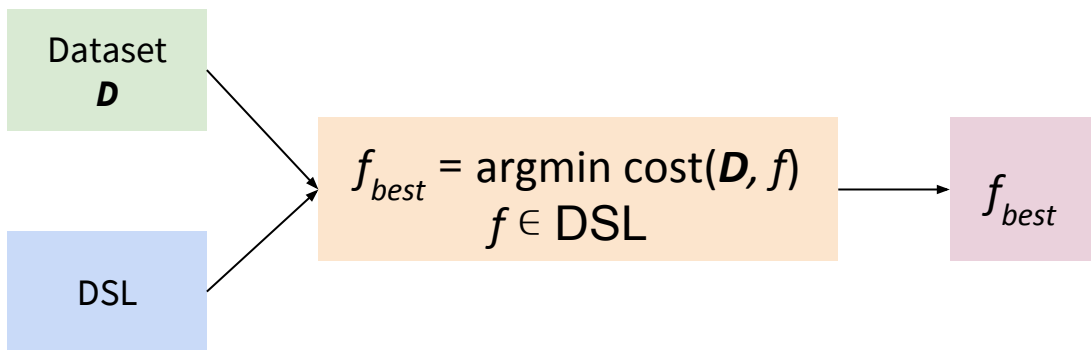
MoveOp ::= Up, Left, Right, DownFirst, DownLast,
NextDFS, PrevDFS, NextLeaf,
PrevLeaf, PrevNodeType, PrevNodeValue, PrevNodeContext

WriteOp ::= WriteValue, WriteType, WritePos

Semantics



Step 3: Synthesize f_{best}



Enabled via advances in synthesis

Discussed later

[1] **Learning Programs from Noisy Data**

Veselin Raychev, Pavol Bielik, Martin Vechev, Andreas Krause
ACM POPL 2016

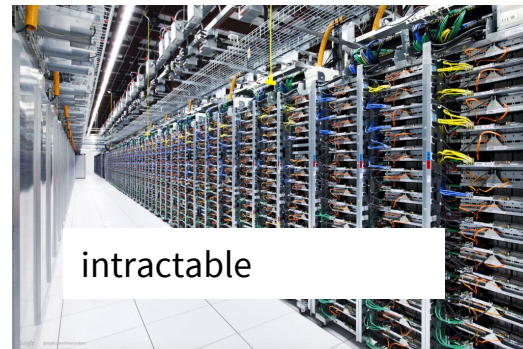
[2] **PHOG: Probabilistic Model for Code**

Pavol Bielik, Veselin Raychev, Martin Vechev
ICML 2016

[3] **Probabilistic Model for Code with Decision Trees**

Veselin Raychev, Pavol Bielik, Martin Vechev
ACM OOPSLA 2016

Enumerative search



Step 4: Use f_{best} to train and predict

Query

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

f_{best}

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

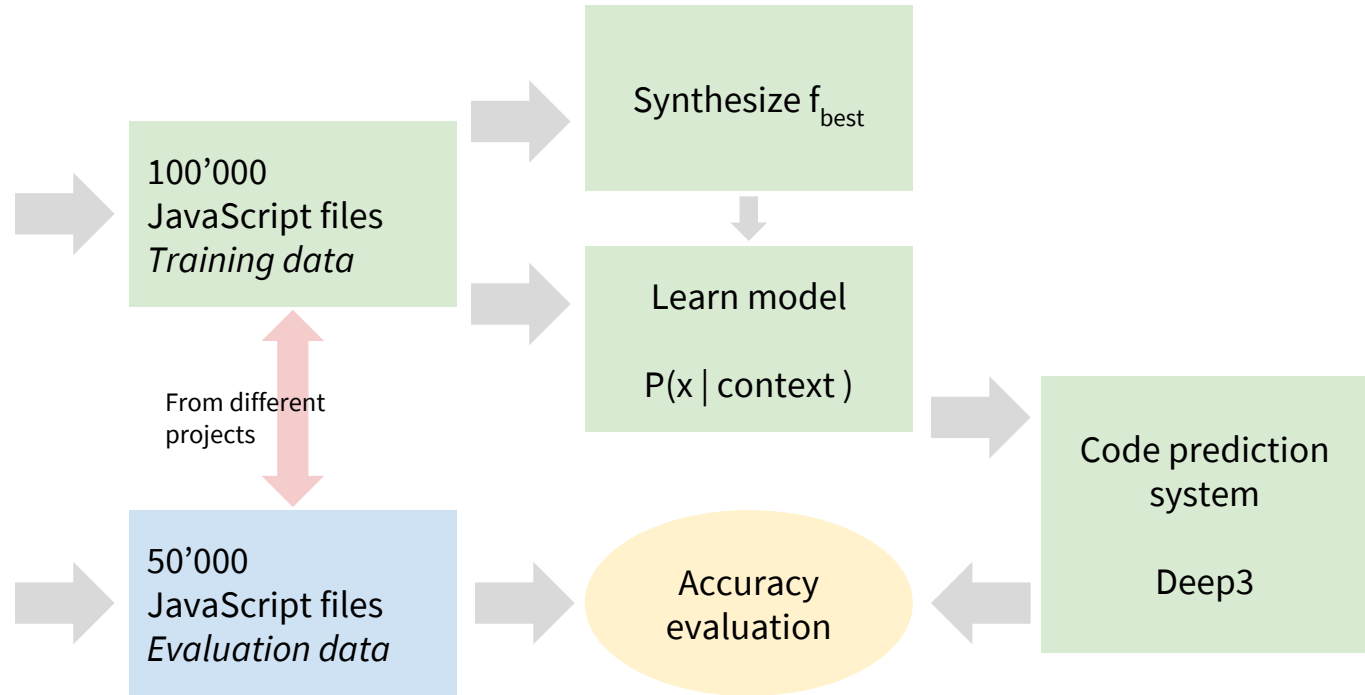
Context

γ

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

{Previous Property, Parameter Position, API name}

Evaluation of a machine learning system



Experiment: JavaScript API prediction

Synthesis time: ~**100** hours

f_{best} ~ **2000** instructions

Probabilistic Model	Accuracy
Last two tokens, Hindle et. al. [ICSE'12]	22.2%
Last two APIs, Raychev et. al. [PLDI'14]	30.4%
<i>Deep3 (Synthesized Model)</i>	66.6%

There is only one context and even it is suboptimal

Beyond code... to Natural Languages

Pavol Bielik, Veselin Raychev, Martin Vechev
ICLR 2017

Writing in natural lan

g

Predict next character

Immediately useful for virtual keyboards (phones)

DSL for natural language: **see exercise.**

Beyond code... to Natural Languages

Dataset **D**:

Synthesis time: ~8 hours

f_{best} ~ 9000 instructions

Hutter Prize Wikipedia Dataset

Probabilistic Model

Entropy (bits per char)

Stacked neural network LSTM (Graves 2013)

1.67

MI-LSTM, HM-LSTM

1.40-1.44

Synthesis

Two orders of magnitude
faster to query than LSTM

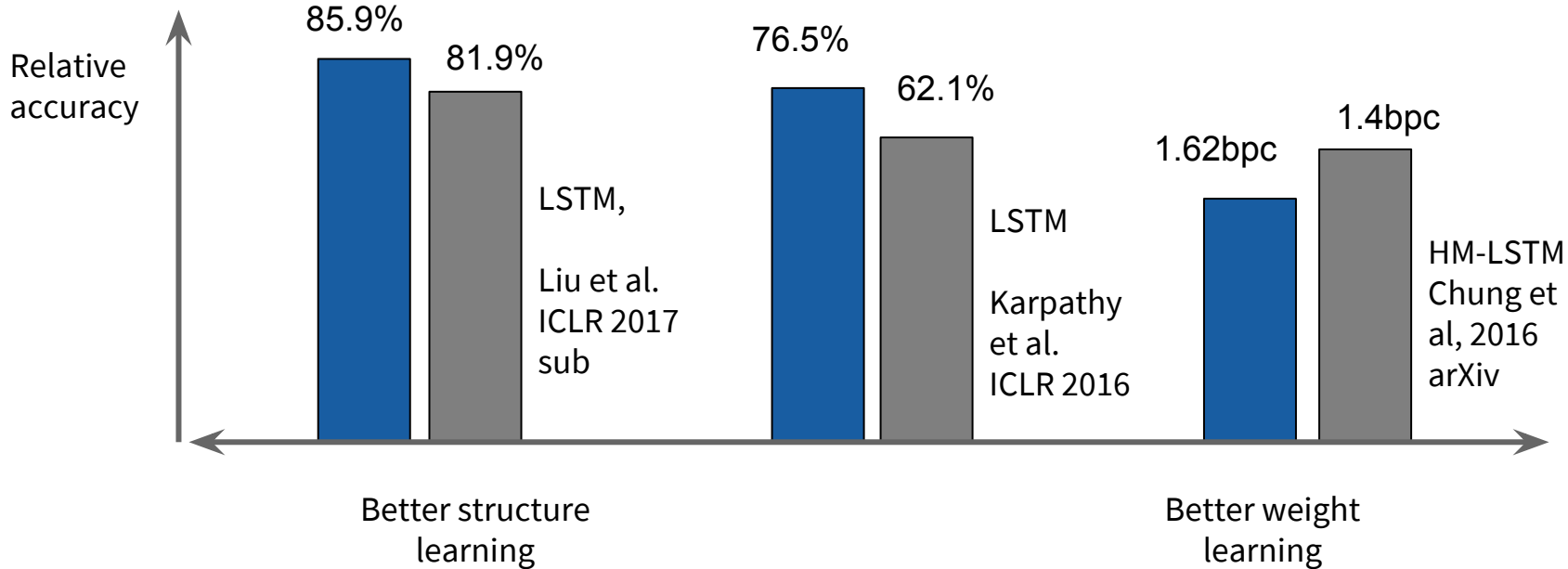
1.62

Synthesis vs Neural networks

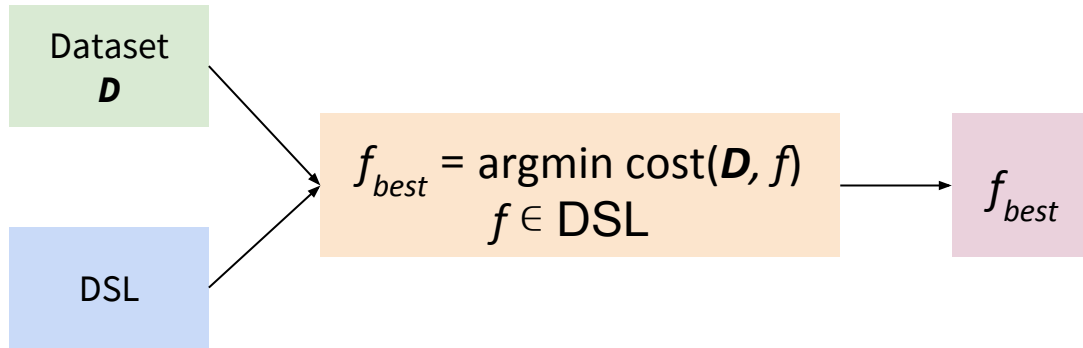
JavaScript ASTs
Node value prediction

Program **text/comments**
Next character prediction

Wikipedia (**English** text)
Next character prediction



How does it work?



Cost: entropy

Example dataset:
abcbabc

Entropy = Number of bits to encode the dataset

$f_1 = \text{Left WriteChar}$

$$H(\mathbf{D}) = -\sum_{x \in \mathbf{D}} 1/n \log_2 Pr(x)$$

$x = (\text{prefix}, \text{char})$

$$Pr(x) = P(\text{char} | f(\text{prefix}))$$

$f_1(\text{prefix})$	char	
$f_1(\text{abcbab}_-) = \mathbf{b}$	c	$P(\mathbf{c} \mathbf{b}) = 2/3$
$f_1(\text{abcba}_-) = \mathbf{a}$	b	$P(\mathbf{a} \mathbf{b}) = 1/3$
$f_1(\text{abcb}_-) = \mathbf{b}$	a	
$f_1(\text{abc}_-) = \mathbf{c}$	b	$P(\mathbf{b} \mathbf{a}) = 1$
$f_1(\text{ab}_-) = \mathbf{b}$	c	$P(\mathbf{b} \mathbf{c}) = 1$
$f_1(\mathbf{a}_-) = \mathbf{a}$	b	
$f_1(\mathbf{-}) = \varepsilon$	a	$P(\mathbf{a} \varepsilon) = 1$

$$H(\mathbf{D}) = -2/7 \log_2(2/3) - 1/7 \log_2(1/3)$$

0.39 bits per char

Cost: entropy

Example dataset:
abcbabc

Entropy = Number of bits to encode the dataset

$f_2 =$ Left Left WriteChar

$$H(\mathbf{D}) = -\sum_{x \in \mathbf{D}} 1/n \log_2 Pr(x)$$

$x = (\text{prefix}, \text{char})$

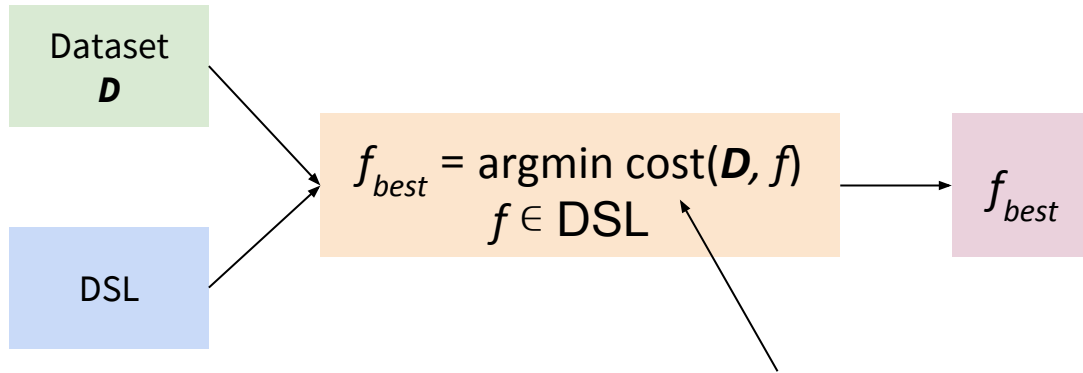
$$Pr(x) = P(\text{char} | f(\text{prefix}))$$

$f_2(\text{prefix})$	char	
$f_2(\text{abcbab}_-) = a$	c	$P(c a) = 1$
$f_2(\text{abcba}_-) = b$	b	$P(b b) = 1$
$f_2(\text{abcb}_-) = c$	a	$P(a c) = 1$
$f_2(\text{abc}_-) = b$	b	
$f_2(\text{ab}_-) = a$	c	
$f_2(a_-) = \varepsilon$	b	$P(b \varepsilon) = 1/2$
$f_2(_ -) = \varepsilon$	a	$P(a \varepsilon) = 1/2$

$$H(\mathbf{D}) = -2/7 \log_2(1/2)$$

0.29 bit per character

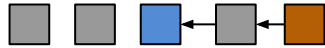
How does it work?



Evaluating **COST** is slow for large dataset

Idea: use counter-example-based synthesis

Synthesis with large datasets



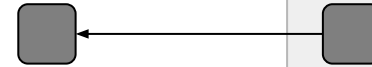
Model described by a short program

Left Left Write

2) search candidate programs on small dataset

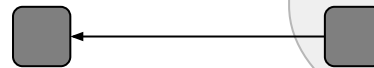
- Left Left Write Left Left Write
- Left Left Write Left Write
- Left Write Left Write
- ϵ

1) Pick a random subset



Large dataset

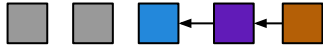
3) evaluate best program on large dataset



4) Pick a next dataset

repeat ...

More expressive DSL with branches



Model described by a short program

But different program for different kinds of predictions - first character of a word is different than other characters

```
if Left Write == space
  • Left PrevChar Right Write
else
  • Left Write
```

Synthesized programs are not short anymore

Search space is intractable

Synthesis: DSL and decision trees

Interesting DSL program:

ϵ

Describes unconditional model

$P(\text{prediction} \mid \text{empty context})$

Can use it to synthesize in pieces

Synthesis: DSL and decision trees

Interesting DSL program:

ϵ

Describes unconditional model

$P(\text{prediction} \mid \text{empty context})$

Can use it to synthesize in pieces

Synthesize condition expression that splits the dataset and gives it to empty programs in branches

1.

```
if Left Write == space
```

• ϵ_1

```
else
```

• ϵ_2

Synthesis: DSL and decision trees

Interesting DSL program:

ϵ

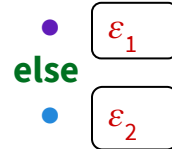
Describes unconditional model

$P(\text{prediction} \mid \text{empty context})$

Can use it to synthesize in pieces

2.

if Left Write == space



Recursively repeat for both branches

Synthesis: DSL and decision trees

Interesting DSL program:

ϵ

Describes unconditional model

$P(\text{prediction} \mid \text{empty context})$

Can use it to synthesize in pieces

3.

if Left Write == space

- Left PrevChar Right Write

else

- Left Write

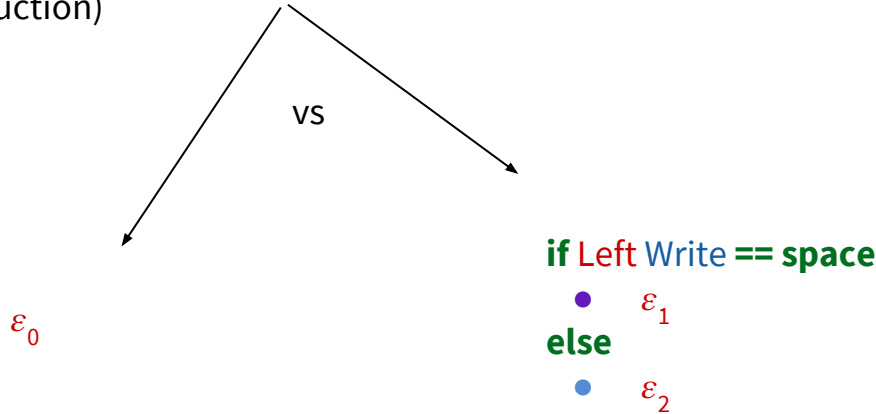
Synthesize branch-free programs if recursive branches do not improve score

Synthesis: DSL and decision trees

Main result 1: When optimization goal is entropy, then this is a formulation of the **ID3** decision tree learning algorithm.

ID3 is the most commonly used decision tree algorithm, but it is not optimal

ID3 maximizes information gain (i.e. entropy reduction)



Synthesis: DSL and decision trees

Main result 2: DSL formulation allows us to instantiate other previously unexplored decision tree algorithms

New decision tree algorithm: E13

1.

Left Write

Synthesize one program for the entire dataset

2.

```
if Left Write == space
  • Left Write
else
  • Left Write
```

Synthesize condition expression for 1st program in both branches

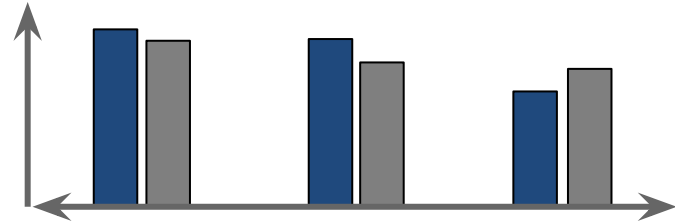
3.

```
if Up WriteType == MethodCallExpression
  • Left PrevChar Right Write
else
  • Left Write
```

Synthesize recursively both both branches

Probabilistic model with synthesis

Accuracy **competitive** with neural networks



Explainable predictions:
Synthesized DSL programs are interpretable by humans

```
if Left Write == space
  • Left Write
else
  • Left Write
```

Application in programming **tools** and **beyond**

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

Natural language processing

DeepCode.AI



ETH spin-off

MSc theses offered

If you are interested to work on interpretable machine learning, write me veselin@deepcode.ai and martin.vechev@inf.ethz.ch

Demo