

# **Reliable and Interpretable Artificial Intelligence**

Martin Vechev

ETH Zurich

Fall 2017

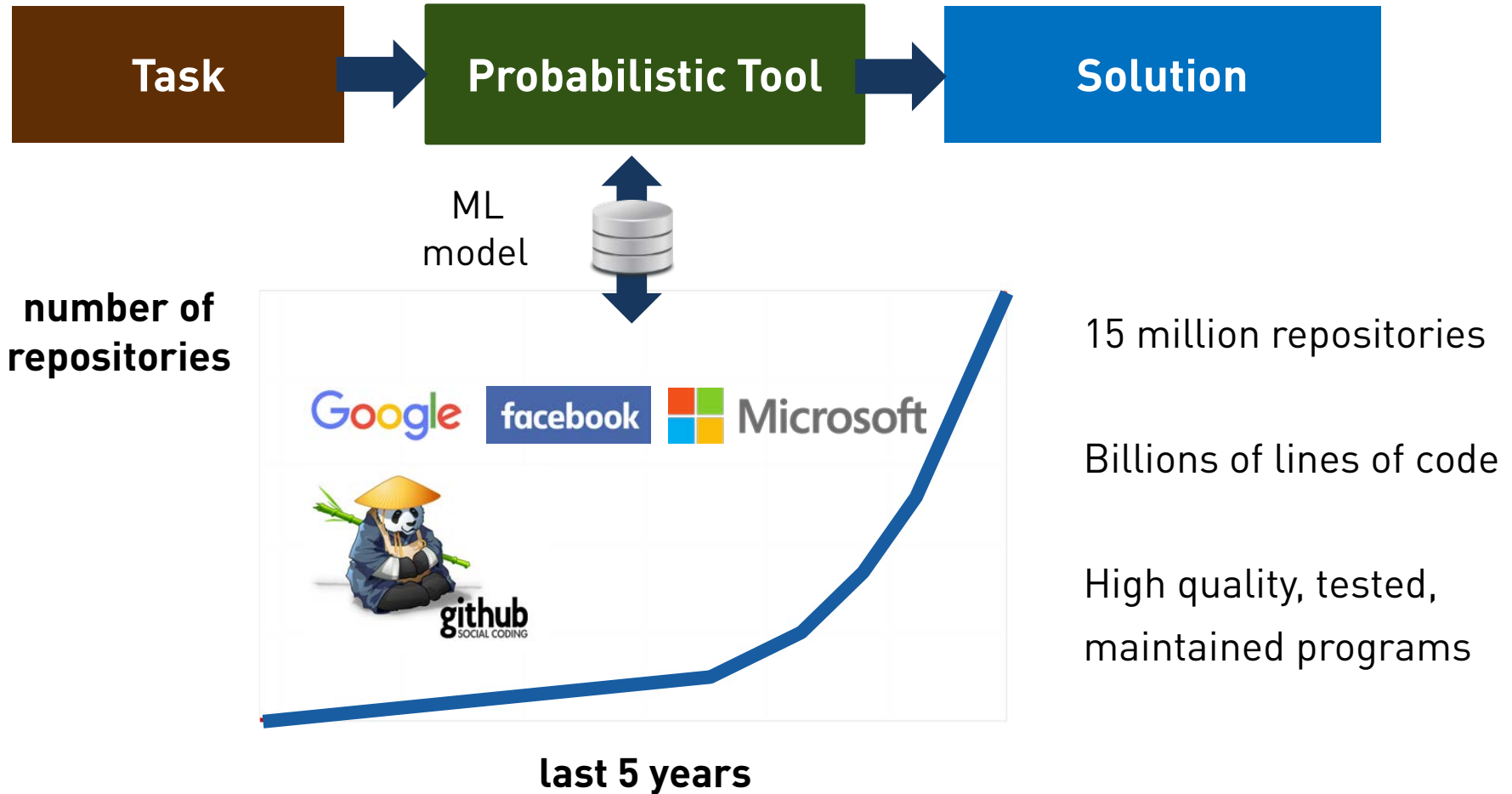
# Machine Learning (ML) + Synthesis

So far we looked at “program synthesis in the small”, meaning we **did not leverage knowledge of prior, similar programs**, when synthesizing a new program. In this lecture, we are going to study both directions:

**Part I:** Leverage **ML models to enable program synthesis**. This will involve learning from a new kind of data sets (e.g., now programs themselves are examples). Motivated by this we will move to **Part II**

**Part II:** Leverage **synthesis to build more interpretable ML models**. An alternative to attention mechanism in neural networks. Applies the iterative synthesis (oracle-based) loop concepts from previous two lectures to learn ML models. The techniques here apply beyond learning from code but natural language and other data.

# Probabilistic Learning from Code



# 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 Programming

## Applications

Code synthesis  
Deobfuscation  
Program repair  
Feedback generation  
Translation

## Intermediate Representation

Sequences (sentences)  
Trees  
Translation Table  
Graphical Models (CRFs)  
Feature Vectors

## Semantic Analysis

typestate analysis  
scope analysis  
control-flow analysis  
alias analysis

## Train Model (ML)

Neural Networks  
N-gram language model  
SVM  
Structured SVM

## Query Model (ML)

$\operatorname{argmax}_{y \in \Omega} P(y | x)$   
Greedy MAP inference

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

ML Systems we have built, such as [jsnice.org](http://jsnice.org) and [apk-deguard.com](http://apk-deguard.com) have > 300,000 users...

# Nice2Predict.org:

scalable structured prediction framework

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

# Machine Learning for Programming

Applications

Code synthesis

Deobfuscation

Program repair

Feedback generation

Translation

Intermediate Representation

Sequences (sentences)

Trees

Translation Table

Graphical Models (CRFs)

Feature Vectors

Semantic Analysis

typestate analysis

scope analysis

control-flow analysis

alias analysis

Train Model (ML)

Neural Networks

SVM

Structured SVM

N-gram language model

Query Model (ML)

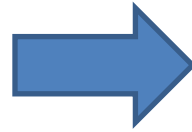
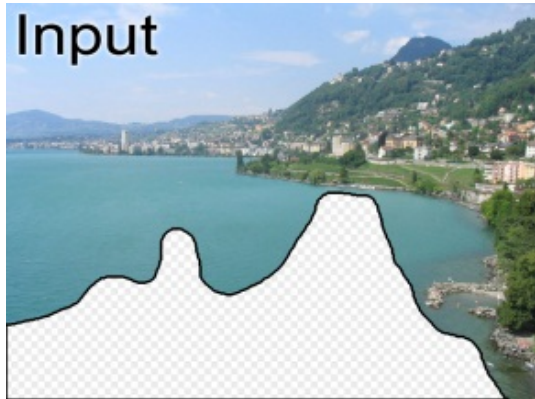
$\operatorname{argmax}_{y \in \Omega} P(y | x)$

Greedy MAP inference

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

ML Systems we have built, such as [jsnice.org](http://jsnice.org) and [apk-deguard.com](http://apk-deguard.com) have > 300,000 users...

# Intuition: Scene Completion





# Statistical Code Synthesis

[Code Completion with Statistical Language Models, ACM PLDI 2014]

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

?

```
MediaRecorder rec = new MediaRecorder();
```

?

```
rec.setAudioSource(MediaRecorder.AudioSource.MIC);  
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

?

```
rec.setOutputFile("file.mp4");
```

...

# Statistical Code Synthesis

[Code Completion with Statistical Language Models, ACM PLDI 2014]

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

```
camera.unlock();
```

```
MediaRecorder rec = new MediaRecorder();
```

```
rec.setCamera(camera);
```

```
rec.setAudioSource(MediaRecorder.AudioSource.MIC);  
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);  
rec.setVideoEncoder(3);
```

```
rec.setOutputFile("file.mp4");
```

```
...
```

Infers sequences  
not in training data

Handles multiple  
objects

Infers multiple  
statements

# Key Insight

Regularities in code are similar to regularities in natural language

# Key Insight

Regularities in code are similar to regularities in natural language

We want to learn that

```
MediaRecorder rec = new MediaRecorder();
```

is before

```
rec.setCamera(camera);
```

# Key Insight

Regularities in code are similar to regularities in natural language

We want to learn that

```
MediaRecorder rec = new MediaRecorder();
```

is before

```
rec.setCamera(camera);
```

like in natural languages

```
Hello
```

is before

```
World !
```

# The SLANG System

partial  
program

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

?

```
MediaRecorder rec = new  
MediaRecorder();
```

?

```
rec.setAudioSource(MediaRecorder.Au  
dioSource.MIC);  
rec.setVideoSource(MediaRecorder.Vid  
eoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.O  
utputFormat.MPEG_4);
```

?

```
rec.setOutputFile("file.mp4");  
...
```

synthesized  
program

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

```
camera.unlock();
```

```
MediaRecorder rec = new  
MediaRecorder();
```

```
rec.setCamera(camera);
```

```
rec.setAudioSource(MediaRecorder.Au  
dioSource.MIC);  
rec.setVideoSource(MediaRecorder.Vid  
eoSource.DEFAULT);  
rec.setOutputFormat(MediaRecorder.O  
utputFormat.MPEG_4);
```

```
rec.setAudioEncoder(1);  
rec.setVideoEncoder(3);
```

```
rec.setOutputFile("file.mp4");  
...
```

sentences  
with holes

completed  
sentences

Semantic  
Analysis

Query

Combine

Completion phase

Training phase

Language  
model

sentences

Semantic  
Analysis

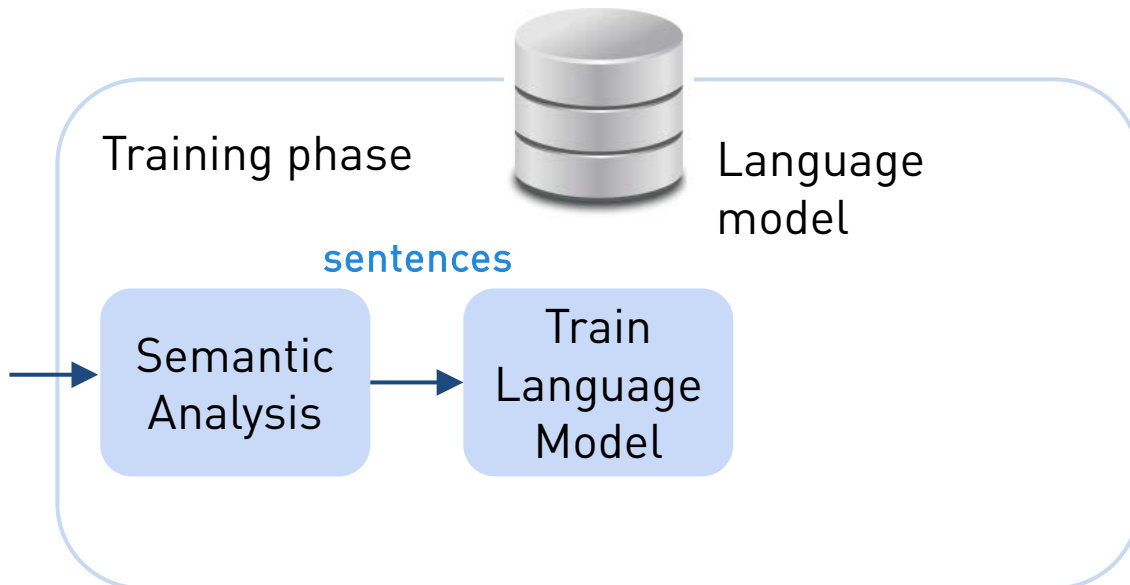
Train  
Language  
Model



github  
SOCIAL CODING

Atlassian  
Bitbucket

# The SLANG System



# Step1: From Programs to Sentences

```
she = new X();  
me = new Y();  
  
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

A note on:

learning from code vs.  
learning from natural language (NLP):

Unlike natural language, code has **formal semantics** that one can benefit from if they have a way to automatically extract these...

This means that with code the actual data can be:

**Choice 1** (*Syntactic*): a sequence of tokens (words) that make up the program. This follows the NLP mindset.

**Choice 2** (*Semantic*): a sequence of words where each word **captures something about the semantics** of the program

Lets see how Choice 2 works...



# From Programs to Sentences

```
she = new X();  
me = new Y();  
  
me.sleep();  
if (random()) {  
    me.eat();  
}  
she.enter();  
me.talk(she);
```

We get 3 sequences:

for abstract object **me**:

$Y_{init}$  sleep talk  
 $Y_{init}$  sleep eat talk

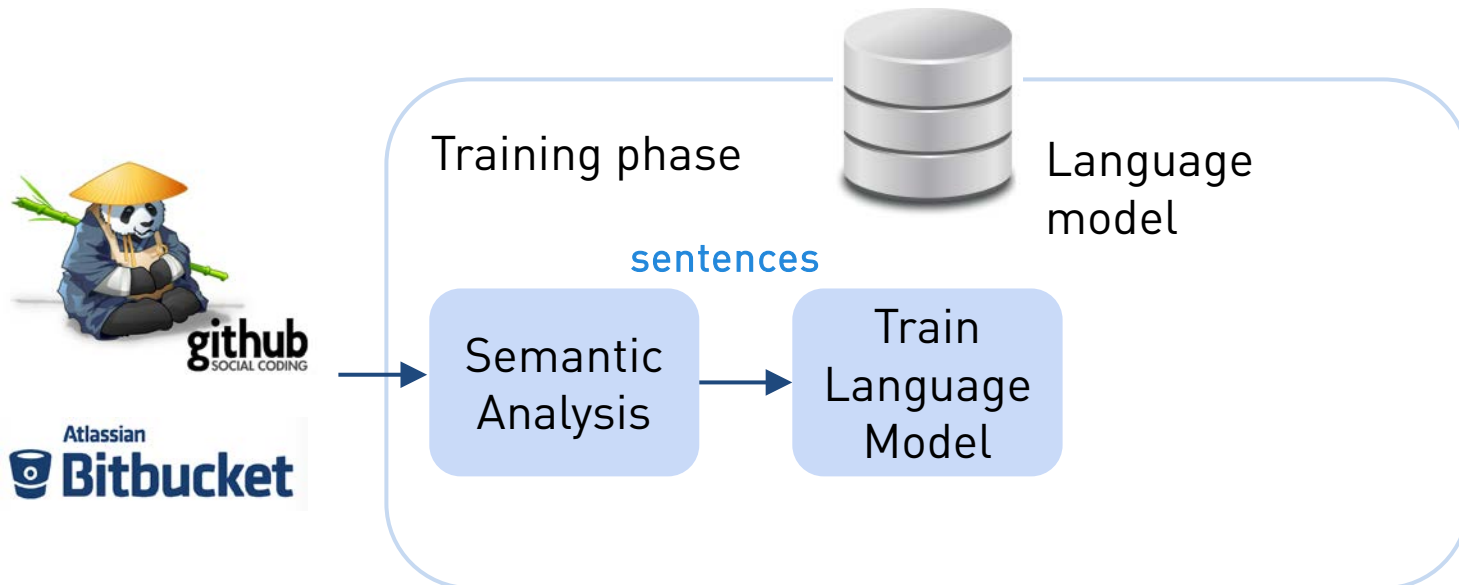
for abstract object **she**:

$X_{init}$  enter talk<sub>param1</sub>

**abstract (semantic) object:**

think of it as the variable denoting all concrete objects assigned to that variable. So many objects assigned to the **she** variable will all be mapped to the same **she** abstract object. In our example we only have 1 concrete object but if we had a loop around **she**=new X(), we would have more concrete objects but still only 1 abstract object.

# Step 2: Learn Regularities in Sentences



# Learn Regularities

Learn regularities in obtained sentences

Regularities in sentences  $\Leftrightarrow$  regularities in code usage

If we see many sequences like:

$X_{\text{init}}$  enter talk<sub>param1</sub>

then we should learn that talk<sub>param1</sub> is often after enter

# Statistical Language Models

Given a sentence  $s = w_1 w_2 w_3 \dots w_n$

estimate  $\mathbf{P}(w_1 w_2 w_3 \dots w_n)$

Decomposed to conditional probabilities

$$\mathbf{P}(w_1 w_2 w_3 \dots w_n) = \prod_{i=1..n} \mathbf{P}(w_i \mid w_1 \dots w_{i-1})$$

# Statistical Language Models

Given a sentence  $s = w_1 w_2 w_3 \dots w_n$

estimate  $\mathbf{P}(w_1 w_2 w_3 \dots w_n)$

Decomposed to conditional probabilities  
(via chain rule):

$$\mathbf{P}(w_1 w_2 w_3 \dots w_n) = \prod_{i=1..n} \mathbf{P}(w_i \mid w_1 \dots w_{i-1})$$

$\mathbf{P}(\text{The quick brown fox jumped}) =$

$$\mathbf{P}(\text{The}) \mathbf{P}(\text{quick} \mid \text{The}) \mathbf{P}(\text{brown} \mid \text{The quick}) \mathbf{P}(\text{fox} \mid \text{The quick brown}) \mathbf{P}(\text{jumped} \mid \text{The quick brown fox})$$

# Choice I: N-gram language model

Conditional probability based only on previous **n-1** words

$$\mathbf{P}(w_i \mid w_1 \dots w_{i-1}) \approx \mathbf{P}(w_i \mid w_{i-n+1} \dots w_{i-1})$$

# N-gram language model

Conditional probability based only on previous **n-1** words

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

# N-gram language model

Conditional probability based only on previous  $n-1$  words

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

Training is achieved by **counting** n-grams, called MLE = maximum likelihood estimation. E.g., with 3-gram language model, we get:

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

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

Time complexity for each word encountered in training is **constant**, so training is usually **fast**.

\* technicality: to actually get a probability distribution across all sentences we typically add N-1 start symbols and N-1 end symbols </s> to each sentence.



# 3-gram language model: example

$$P(w_1 \cdot w_2 \cdot w_3 \cdot \dots \cdot w_n) \approx P(w_1) * P(w_2 | w_1) * P(w_3 | w_1 \cdot w_2) * \dots * P(w_n | w_{n-2} \cdot w_{n-1})$$

3-grams	# of occurrences
brown fox jumped	125
brown fox walked	45
brown fox snapped	30


$$P(\text{jumped} | \text{brown fox}) \approx 125/200 \approx 0.625$$

$$P(\text{brown fox jumped}) \approx$$

$$\begin{array}{c} P(\text{brown}) \\ 200 / 600 \end{array} * \begin{array}{c} P(\text{fox} | \text{brown}) \\ 200 / 200 \end{array} * \begin{array}{c} P(\text{jumped} | \text{brown fox}) \\ 125 / 200 \end{array} \approx 0.208$$

# Key Problem: Sparsity of Data

What if this number is 0?


$$P(\text{jumped} \mid \text{brown fox}) \approx \frac{\#(\text{brown fox jumped})}{\#(\text{brown fox})}$$

The problem of sparsity **gets worse** as the size of the n-gram becomes **larger**.

We need to handle n-grams with 0 or few occurrences in the training data.

Techniques that can do that are: **smoothing, discounting**

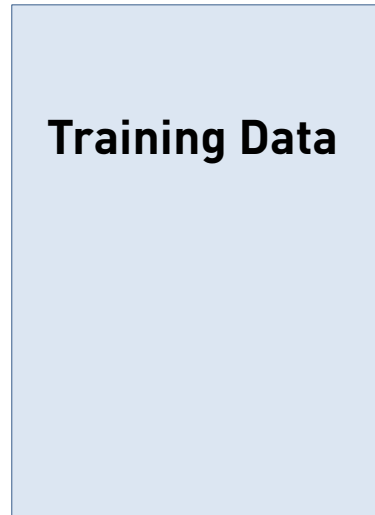
# Solution: Smoothing Techniques

- Smoothing techniques give **non-zero probability** to n-grams **not in the training data**.
- Essentially, they try to **estimate how likely** it is that the n-gram is missing due to the limited size of the training data.
- They work by taking the probability mass of the existing n-grams and redistributing that mass over n-grams that occur zero times.
- Typically, probability of such n-grams is estimated by looking at the probabilities of n-1 grams , n-2 grams, unigrams, etc.

**Example** smoothing techniques are: Witten-Bell Interpolated (WBI), Witten-Bell Backoff (WBB), Natural Discounting (ND), Stupid-Backoff (SB)

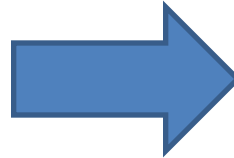
# Smoothing: Intuitively

Distribution of probability mass before smoothing

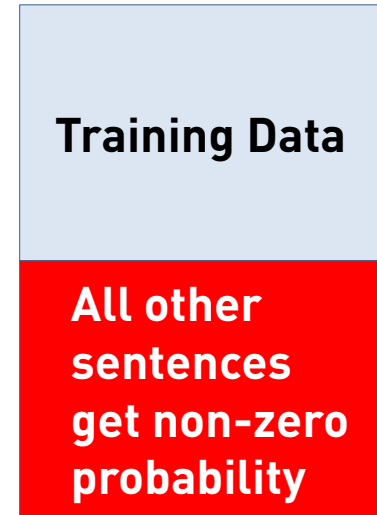


**All other sentences get 0 probability**

**Smoothing**



Distribution of probability mass after smoothing



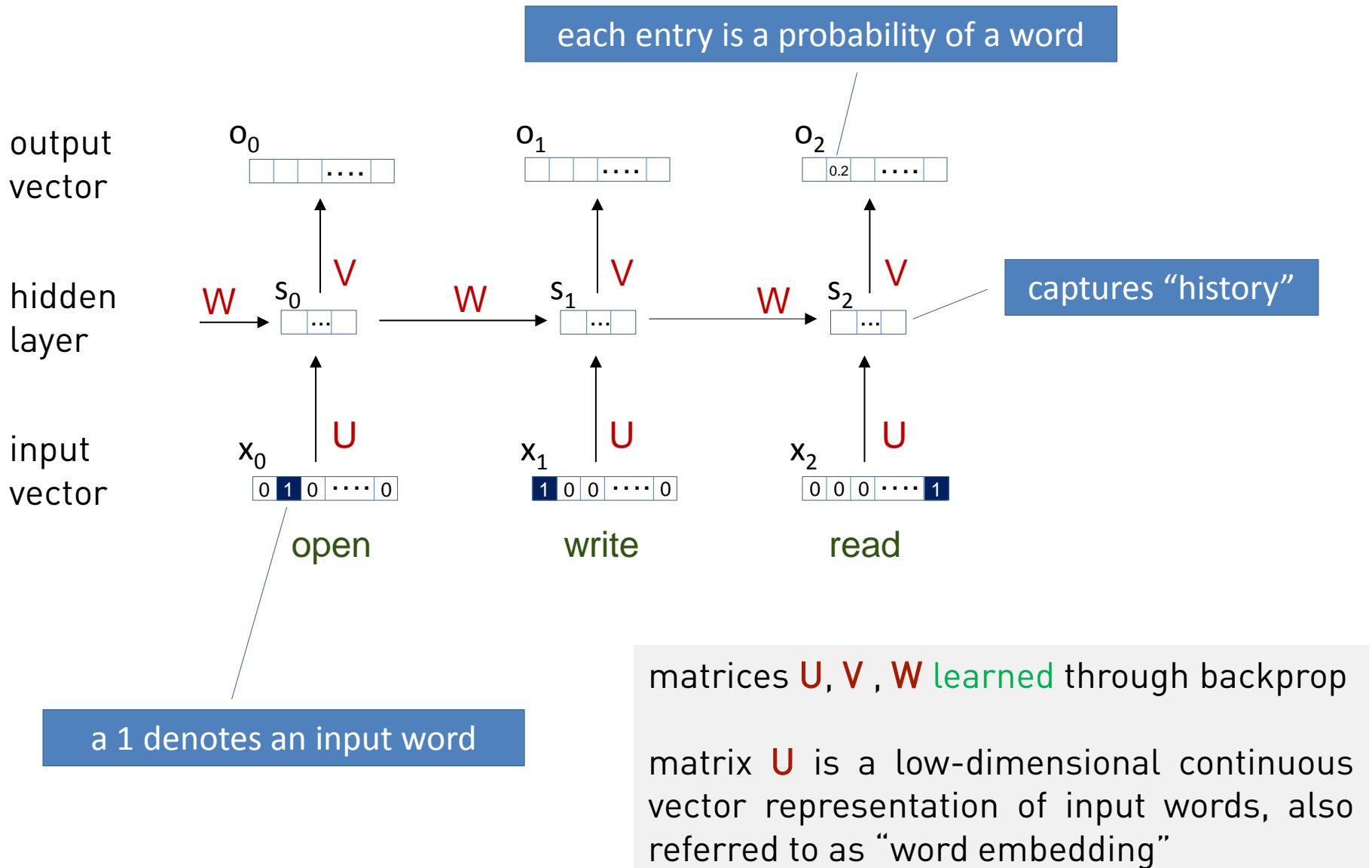
# N-gram language model

Conditional probability only on previous **n-1** words

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

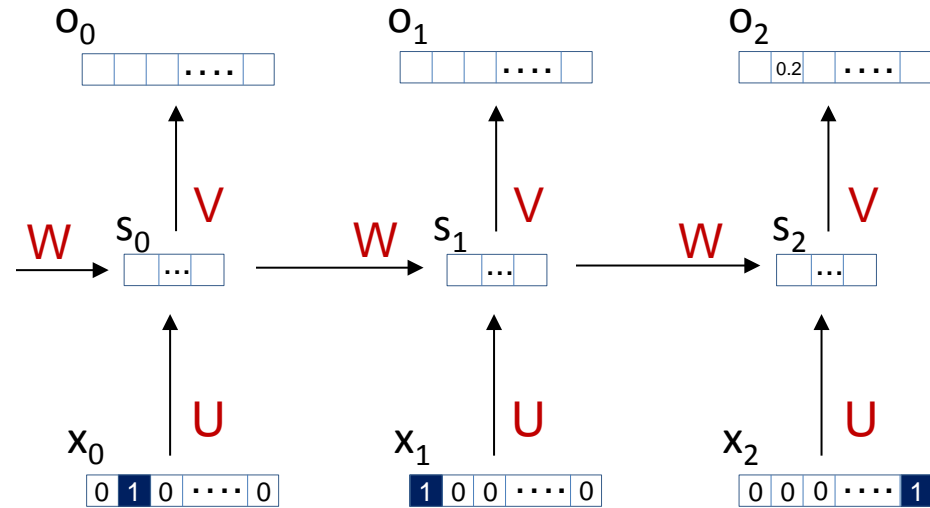
Existing library implementing language models is SRILM:  
<http://www.speech.sri.com/projects/srilm/>

# Choice II: Recurrent Neural Networks (RNN)



RNNs: more details on board

# Choice II: Recurrent Neural Networks (RNN)



- + can learn dependencies beyond the prior several words (e.g., LSTM)
- + learns continuous representation of vocabulary (helps avoid sparsity)
- slow and tricky to train
- predictions hard to explain

RNNLM-40 Download: <http://rnnlm.org/>



<https://www.tensorflow.org/>



# ML-based Code Synthesis

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    ? {smsMgr, list} // (Hole H1)
} else {
    ? {smsMgr, message} // (Hole H2)
}
```

# ML-based Code Synthesis

```
smsMgr = SmsManager.getDefault();  
int length = message.length();  
if (length > MAX_SMS_MESSAGE_LENGTH) {  
    list = smsMgr.divideMessage(message);  
    ? {smsMgr, list} // (Hole H1)  
} else {  
    ? {smsMgr, message} // (Hole H2)  
}
```

Abstract object **smsMgr**:

getDefault<sub>result</sub> divideMessage H1  
getDefault<sub>result</sub> H2

Abstract object **list**:

divideMessage<sub>result</sub> H1

Abstract object **message**:

length divideMessage<sub>param1</sub>  
length H2

# ML-based Code Synthesis

getDefault <sub>result</sub> divideMessage <b>H1</b>	
--	--

getDefault <sub>result</sub> <b>H2</b>	
--	--

divideMessage <sub>result</sub> <b>H1</b>	
---	--

length <b>H2</b>	
------------------	--

# ML-based Code Synthesis

Completion probability

getDefault <sub>result</sub> divideMessage <b>sendMultipartTextMessage</b>	0.0033
getDefault <sub>result</sub> divideMessage <b>sendTextMessage</b>	0.0016
getDefault <sub>result</sub> <b>sendTextMessage</b>	0.0073
getDefault <sub>result</sub> <b>sendMultipartTextMessage</b>	0.0010
divideMessage <sub>result</sub> <b>sendMultipartTextMessage</b> <sub>param3</sub>	0.0821
length <b>length</b>	0.0132
length <b>split</b>	0.0080
length <b>sendTextMessage</b> <sub>param3</sub>	0.0017

# ML-based Code Synthesis

Completion probability

getDefault <sub>result</sub> divideMessage <b>sendMultipartTextMessage</b>	0.0033
getDefault <sub>result</sub> divideMessage sendTextMessage	0.0016
getDefault <sub>result</sub> <b>sendTextMessage</b>	0.0073
getDefault <sub>result</sub> sendMultipartTextMessage	0.0010
divideMessage <sub>result</sub> <b>sendMultipartTextMessage</b> <sub>param3</sub>	0.0821
length <b>length</b>	0.0132
length split	0.0080
length sendTextMessage <sub>param3</sub>	0.0017

# ML-based Code Synthesis

Completion probability

getDefault <sub>result</sub> divideMessage <b>sendMultipartTextMessage</b>	0.0033
getDefault <sub>result</sub> divideMessage sendTextMessage	0.0016
getDefault <sub>result</sub> <b>sendTextMessage</b>	0.0073
getDefault <sub>result</sub> sendMultipartTextMessage	
divideMessage <sub>result</sub> <b>sendMultipartTextMessage</b>	
length <b>length</b>	
length split	0.0080
length sendTextMessage <sub>param3</sub>	0.0017

Not a feasible solution:  
completions disagree on  
selected method

The solution must satisfy  
program constraints

# ML-based Code Synthesis

Completion probability

getDefault <sub>result</sub> divideMessage <b>sendMultipartTextMessage</b>	0.0033
getDefault <sub>result</sub> divideMessage sendTextMessage	0.0016
getDefault <sub>result</sub> <b>sendTextMessage</b>	0.0073
getDefault <sub>result</sub> sendMultipartTextMessage	0.0010
divideMessage <sub>result</sub> <b>sendMultipartTextMessage</b> <sub>param3</sub>	0.0821
length <b>length</b>	0.0132
length split	0.0080
length sendTextMessage <sub>param3</sub>	0.0017

# ML-based Code Synthesis

Completion probability

getDefault <sub>result</sub> divideMessage <b>sendMultipartTextMessage</b>	0.0033
getDefault <sub>result</sub> divideMessage sendTextMessage	0.0016
getDefault <sub>result</sub> <b>sendTextMessage</b>	0.0073
getDefault <sub>result</sub> sendMultipartTextMessage	0.0010
divideMessage <sub>result</sub> <b>sendMultipartTextMessage</b> <sub>param3</sub>	0.0821
length length	0.0132
length split	0.0080
length <b>sendTextMessage</b> <sub>param3</sub>	0.0017



# ML-based Code Synthesis

```
smsMgr = SmsManager.getDefault();
int length = message.length();
if (length > MAX_SMS_MESSAGE_LENGTH) {
    list = smsMgr.divideMessage(message);
    smsMgr.sendMultipartTextMessage(...list...);
} else {
    smsMgr.sendMessage(...message...);
}
```

# The SLANG System

partial program

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

?

MediaRecorder rec = new
MediaRecorder();

?

AudioSource.MIC);
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);

?

rec.set
```

84 testing samples

sentences with holes

Semantic Analysis

Query

completed sentences

Combine

Completion phase

synthesized program

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

camera.unlock();

MediaRecorder rec = new
MediaRecorder();

rec.setCamera(camera);

rec.setAudioSource(MediaRecorder.AudioSource.MIC);
rec.setVideoSource(MediaRecorder.VideoSource.DEFAULT);
rec.setOutputFormat(MediaRecorder.OutputFormat.MPEG_4);
```

Correct completion in top 3 for ~90% of cases

~100MB

Training phase



Language model

sentences

Semantic Analysis

Train Language Model

~700MB

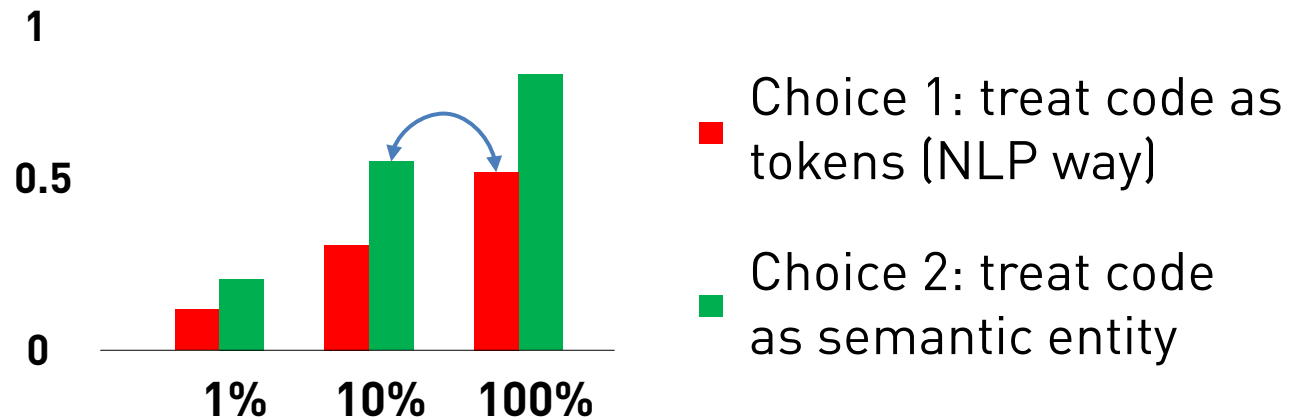
github

1M Java methods



# Do semantics help?

precision vs. % of data used



\*need to decide what to ``compile'' into a word

semantic analysis benefit = 10x more data

Next: Synthesis for Interpretable ML

# Fundamental Question

Data

Learning

Model



Probabilistic  
Model



Widely  
Applicable

Efficient  
Learning

High  
Precision

Interpretable  
Model