

# Probabilistic Models of Source Code: An Overview

**Miltos Allamanis**, University of Edinburgh

November, 2016



THE UNIVERSITY of EDINBURGH  
**informatics**

My PhD is supported by

Microsoft®

**Research**

Developers implicitly embed **knowledge** in code that may be useful for the same or other projects.



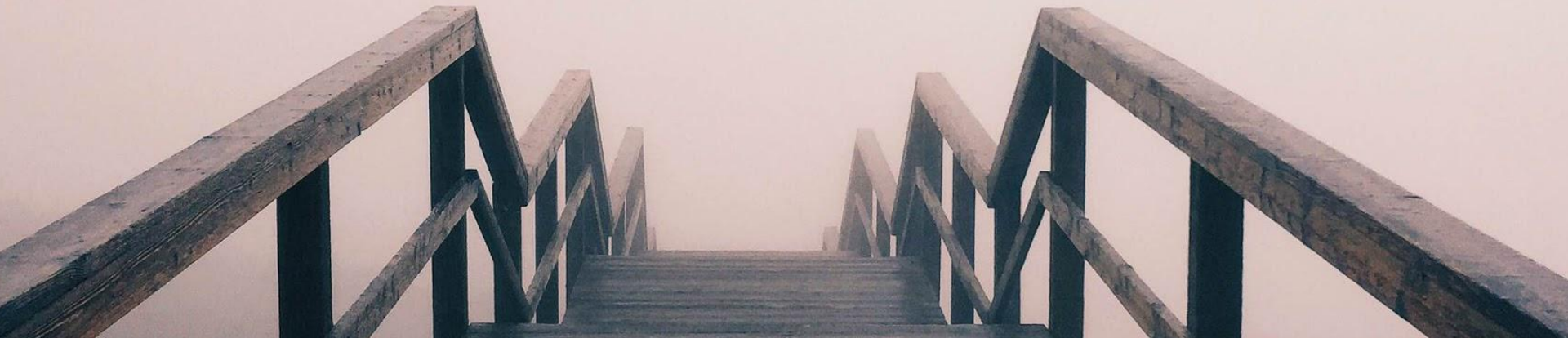
**GitHub**

Atlassian  
**Bitbucket**

internal &  
external  
codebases

**Mine** the hidden knowledge to create **smart** software engineering tools.

Probabilistic reasoning is a **principled** way of handling **ambiguities** and **partial information**.



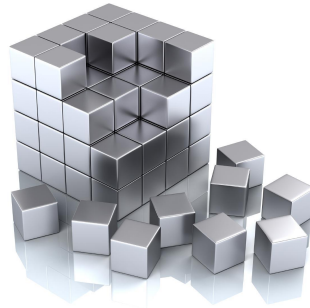
# Natural Language Processing (NLP)



Some **Knowledge**  
of **Linguistics**



**Data:** Corpora of  
Text, Speech etc



**Models of Aspects  
of Natural  
Language**



**Parsing**

**Named Entity  
Recognition**

**Machine  
Translation**

....

➤ **Resolve** language ambiguities with principled probabilistic models of language.

➤ Learn models from annotated corpora.

# Natural Language Processing with Machine Learning



- **Resolve** language ambiguities with principled probabilistic models of language.
- Learn models from annotated corpora.

# Machine Learning Models of Source Code

*“All models are wrong, some are useful” - George Box*



Software Engineers



Codebases



Machine Learning Models  
of Aspects of Source Code



Software Engineering  
Tools

*Just use an n-gram language model.*



*Expert*

# Naturalness of Code

O RLY?

*Miltos*

*Use all the code you can find around.*



*Elementary*

# Big Code

O RLY?

*Miltos*



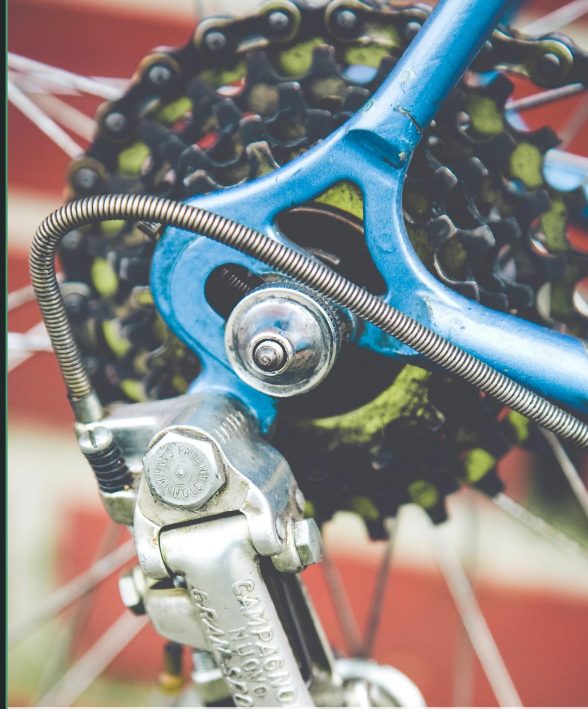
# **Very Brief Intro to Machine Learning**

## **A Taxonomy of Models**

## **Applications**

Outline

# Machine Learning



**Machine Learning Model**

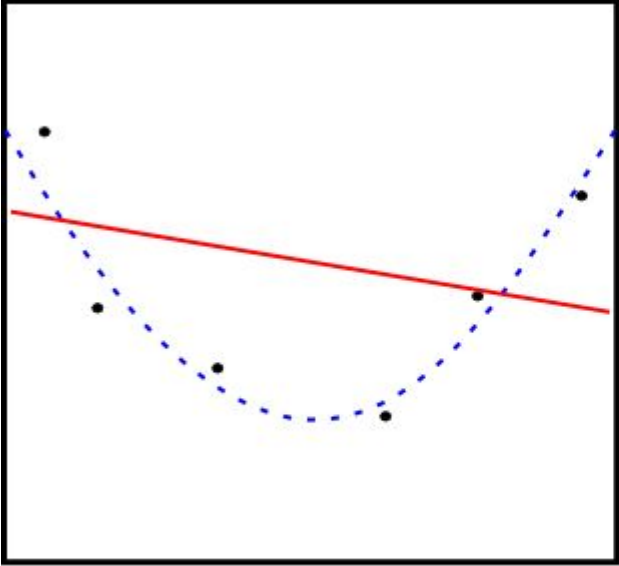
Designed by humans



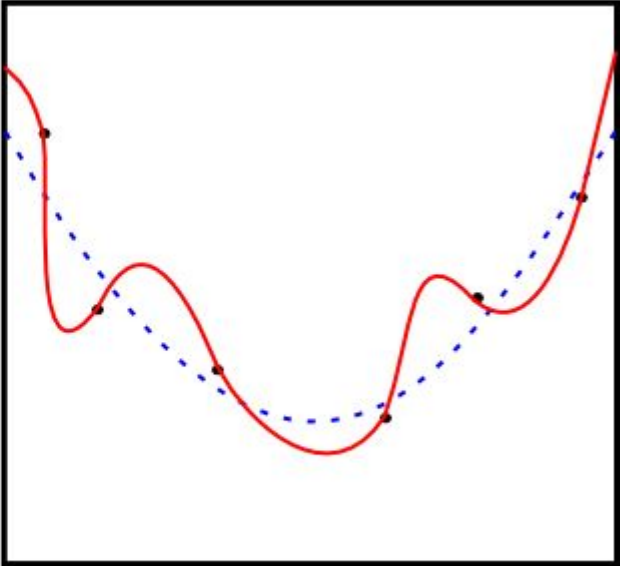
**Model Parameters**

Learned from data

# Finding a good model



Underfitting



Overfitting



# Learning Model Parameters

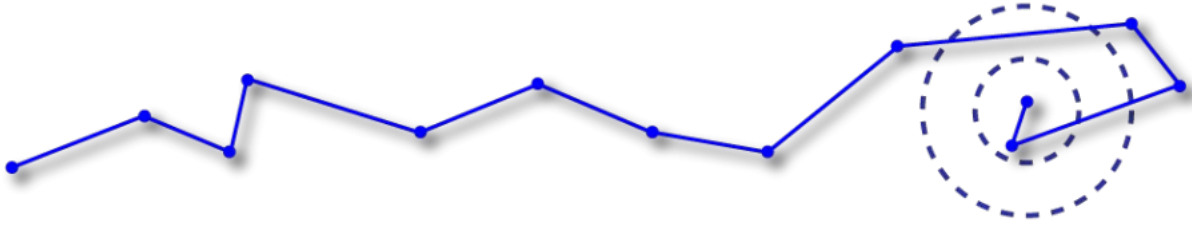


Image from [marple.eeb.uconn.edu](http://marple.eeb.uconn.edu)

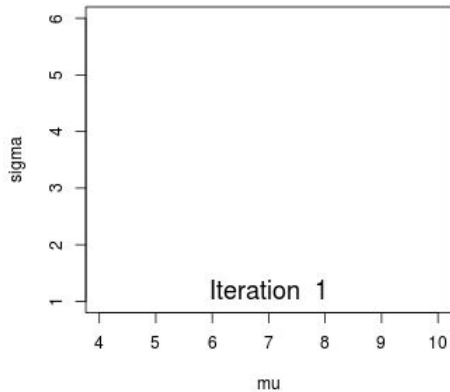
- › **Optimize** objective function in training set
- › Use **computational methods** of optimization



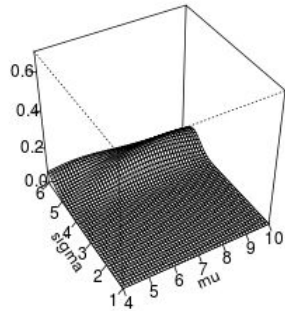
# Learning the Distribution of Model Parameters

$$\underbrace{p(\mu \mid \text{Data})}_{\text{posterior}} = \frac{\underbrace{p(\text{Data} \mid \mu)}_{\text{likelihood}} \cdot \underbrace{p(\mu)}_{\text{prior}}}{\underbrace{p(\text{Data})}_{\text{marginal likelihood}}}$$

Markov chains



Posterior density

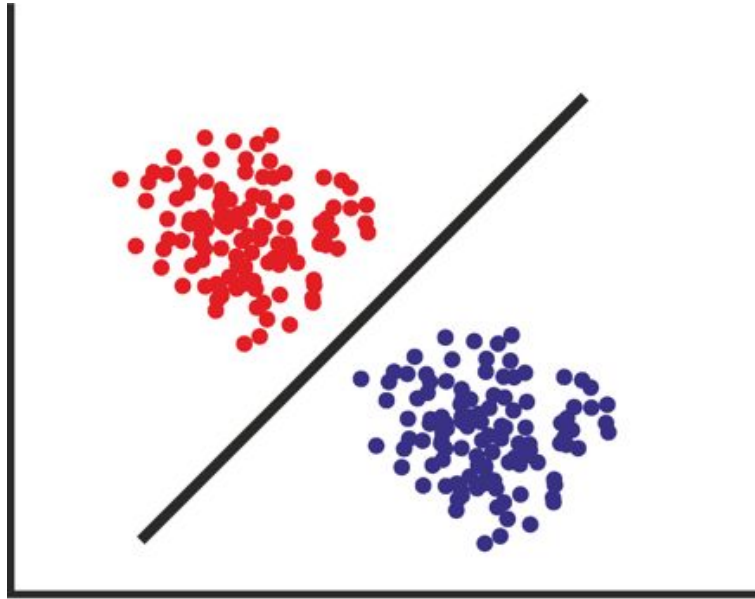


Images from "Bayesian Modelling in Python"

- › **Learn** the distribution of parameters given some data.
- › Use **computational methods**

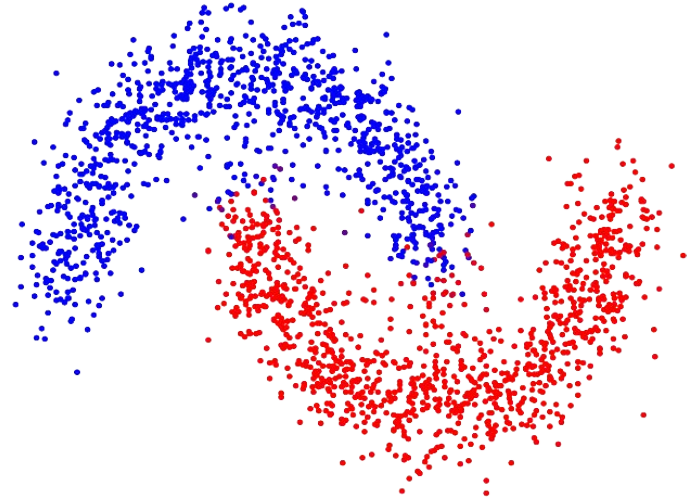


# A Spectrum of Common Machine Learning Tasks



Classification

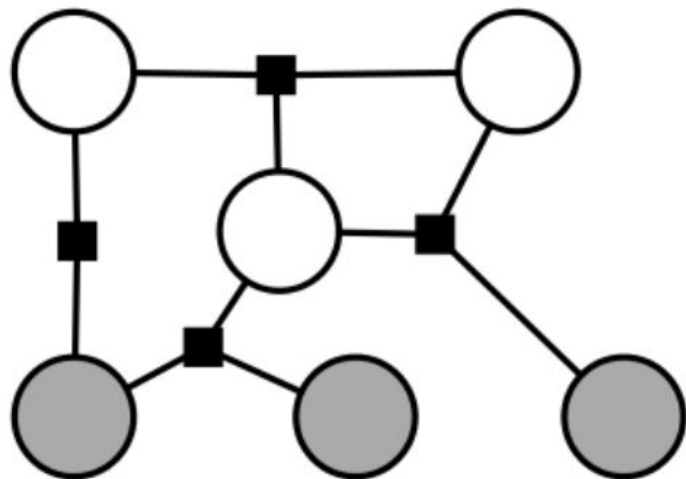
*Supervised*



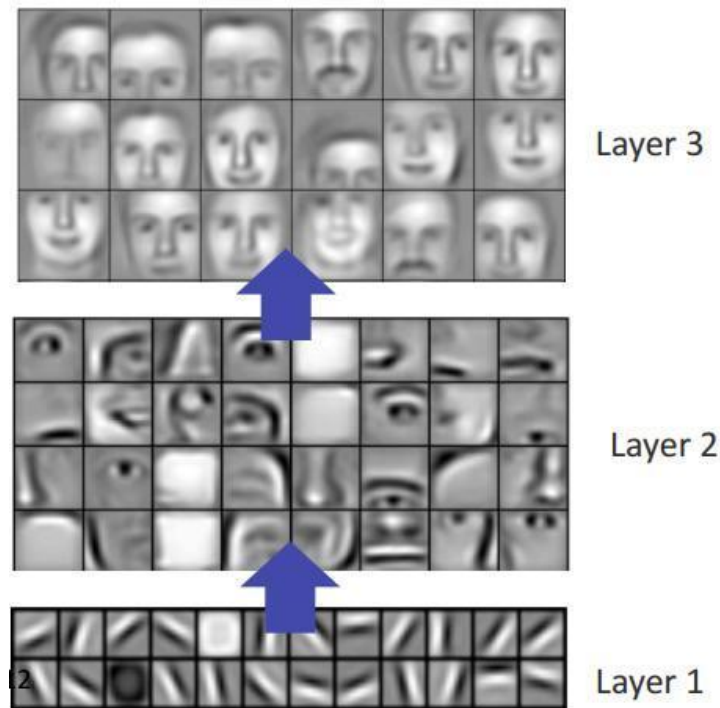
Clustering

*Unsupervised*

# A Spectrum of Common Machine Learning Tasks



Joint Classification



Learning Features

# Automatic Evaluation in Machine Learning

Imperfect measures of performance such as

› Prediction Accuracy

› Model Fit



› **Quantify** performance in a **reproducible** manner

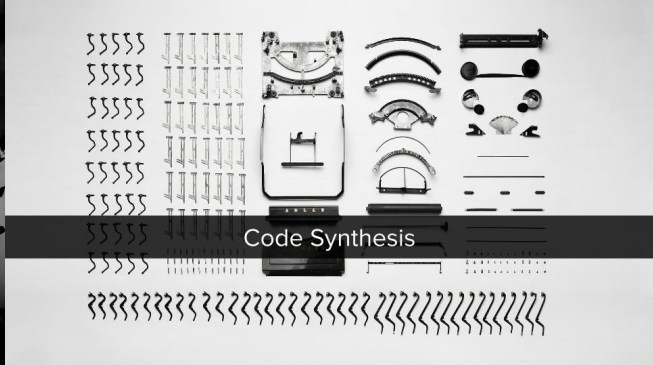
› **Drive** improvement of systems in a **measurable** way



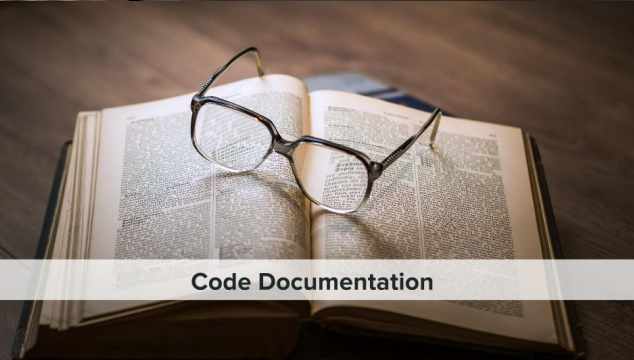
Recommender Systems



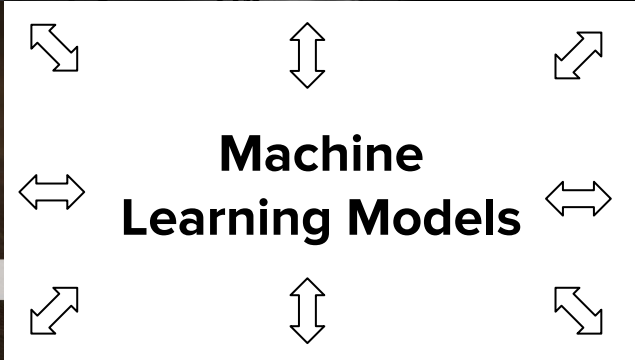
Statistical Code Migration



Code Synthesis



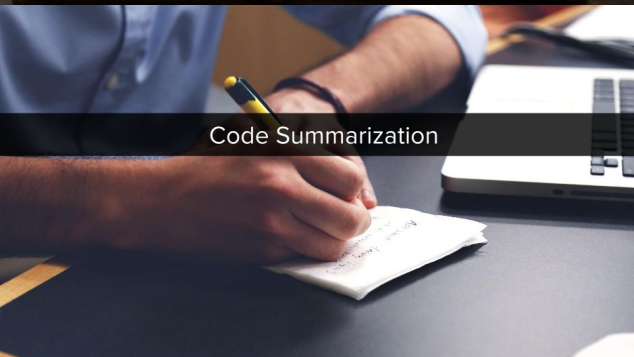
Code Documentation



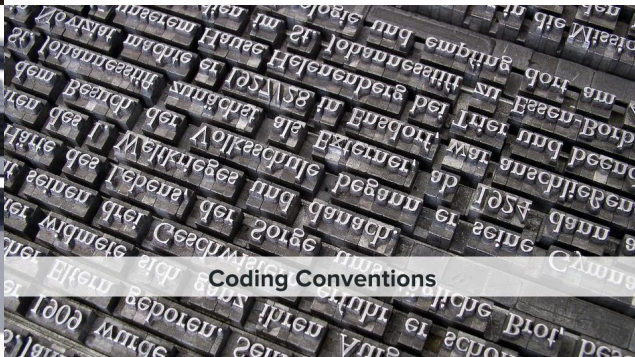
Machine Learning Models



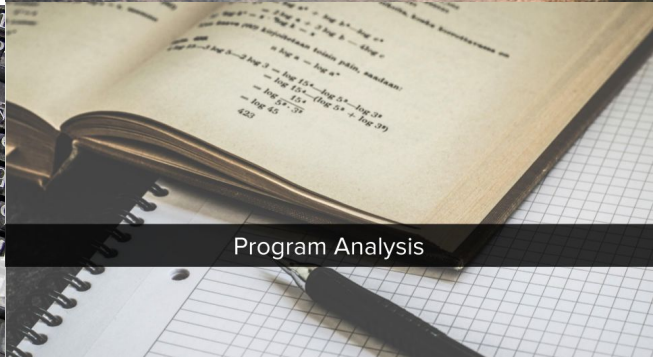
Code Defects



Code Summarization



Coding Conventions



Program Analysis



A close-up photograph of a person's hands holding a folded paper map. The map shows a city street grid with various street names and landmarks. The text 'Part II' is overlaid on the left side of the map.

Part II

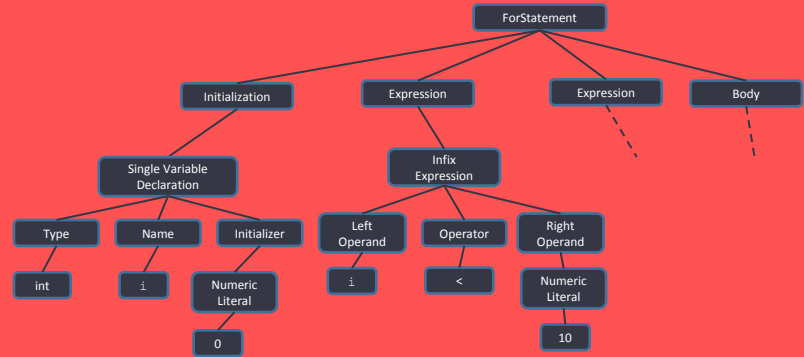
# A Taxonomy of Models

# By Code Representations

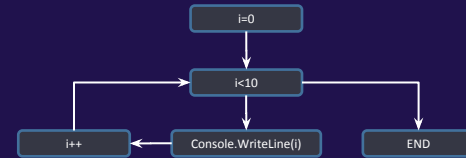
Token

```
for (int i = 0; i < 10; i++){  
    Console.WriteLine(i);  
}
```

Syntax



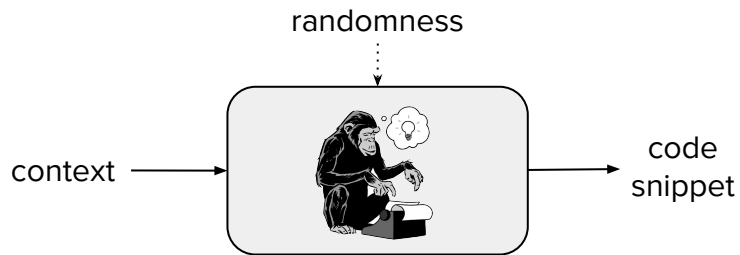
Graph





# Probabilistic Models of Source Code

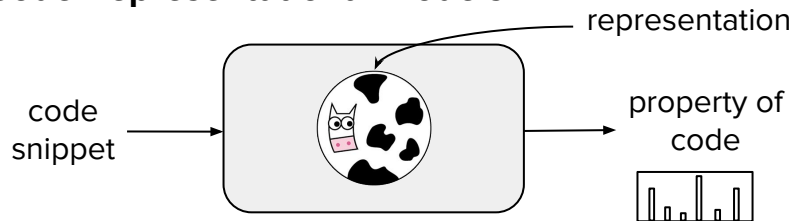
## Code Generating Models



$$P_{\mathcal{D}}(\mathbb{c}|\mathcal{C}(\mathbb{c}))$$

- › Language Models
- › Translation Models
- › Multimodal Models

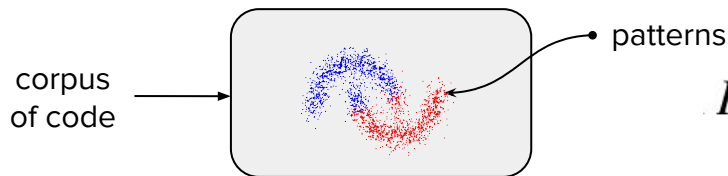
## Code Representational Models



$$P_{\mathcal{D}}(\pi|f(\mathbb{c}))$$

- › Structured Prediction
- › Distributed Representations

## Pattern Mining Models



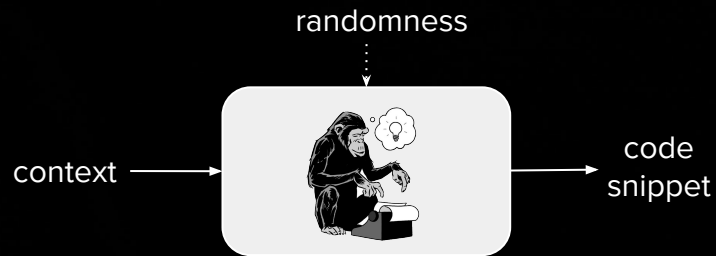
$$P_{\mathcal{D}}(f(\mathbb{c})) = \sum_{\mathbf{I}} P_{\mathcal{D}}(f(\mathbb{c})|\mathbf{I})P(\mathbf{I})$$

- › Latent Variable Models

# Code Generating Models

$$P_{\mathcal{D}}(\mathbb{C} | \underbrace{C(\mathbb{C})}_{\text{context}})$$

↑  
code snippet



Simplifying assumptions about  
*how* code is generated

# Language Models

Assign a non-zero probability to every piece of valid code

$$P(\mathbb{C})$$

```
jQuery(document).ready(function() {  
  $(' .panel-btn').on('click', function(event) {  
    event.preventDefault();  
    $(' .panel').addClass('is-visible');  
  });  
  //close the lateral panel  
  $(' .panel').on('click', function(event) {  
    if( $('event.target').is(' .panel') || $('event.target').is(' .panel-close') ) {  
      $(' .panel').removeClass('is-visible');  
      event.preventDefault();  
    }  
  });  
});
```

```
jQuery(document).ready(function($) {  
  //open the lateral panel  
  $(' .panel-btn').on('click', function(event) {  
    event.preventDefault();  
    $(' .panel').addClass('is-visible');  
  });  
  //close the lateral panel  
  $(' .panel').on('click', function(event) {  
    if( $('event.target').is(' .panel') || $('event.target').is(' .panel-close') ) {  
      {  
        $(' .panel').removeClass('is-visible');  
        event.preventDefault();  
      }  
    }  
  });  
});
```

# n-gram Language Models

```
public void execute(Runnable task) {
    if (task == null)
        throw new NullPointerException();
    ForkJoinTask<?> job;
    if (task instanceof ForkJoinTask<?>) // avoid re-wrap
        job = (ForkJoinTask<?>) task;
    else
        job = new ForkJoinTask.AdaptedRunnableAction(task);
    externalPush(job);
}
```

# n-gram Language Models

```
public void execute(           task) {  
    if (task == null)  
        throw new NullPointerException();  
    ForkJoinTask<?> job;  
    if (task instanceof ForkJoinTask<?>) // avoid re-wrap  
        job = (ForkJoinTask<?>) task;  
    else  
        job = new ForkJoinTask.AdaptedRunnableAction(task);  
    externalPush(job);  
}
```

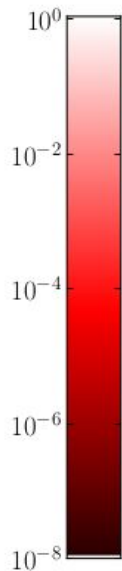
# n-gram Language Models

```
public void execute( task) {  
    if (task == null)  
        throw new NullPointerException();  
    ForkJoinTask<?> job;  
    if (task instanceof ForkJoinTask<?>) // avoid re-wrap  
        job = (ForkJoinTask<?>) task;  
    else  
        job = new ForkJoinTask.AdaptedRunnableAction(task);  
    externalPush(job);  
}
```

$$P(t_0 \dots t_M) = \prod_{m=0}^M P(t_m | t_{m-1} \dots t_{m-n+1})$$



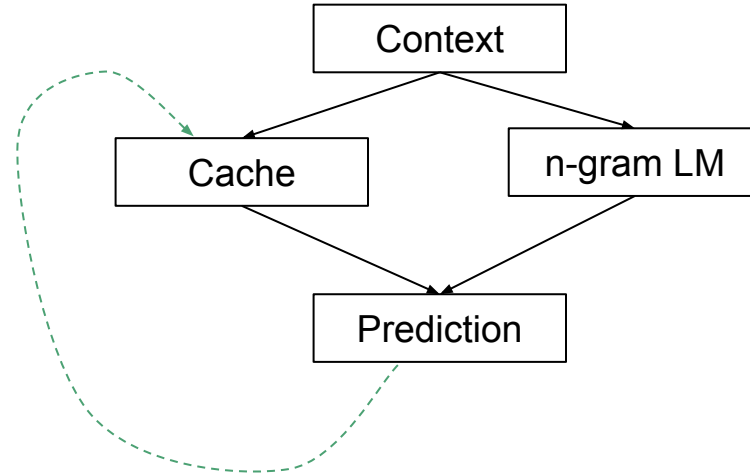
# *n*-gram Language Models



```
1 public void execute(Runnable task) {  
2     if (task == null)  
3         throw new NullPointerException();  
4     ForkJoinTask<?> job;  
5     if (task instanceof ForkJoinTask<?>) // avoid re-wrap  
6         job = (ForkJoinTask<?>) task;  
7     else  
8         job = new ForkJoinTask.AdaptedRunnableAction(task);  
9     doSubmit(job);  
10 }
```

# Cache $n$ -gram Language Models

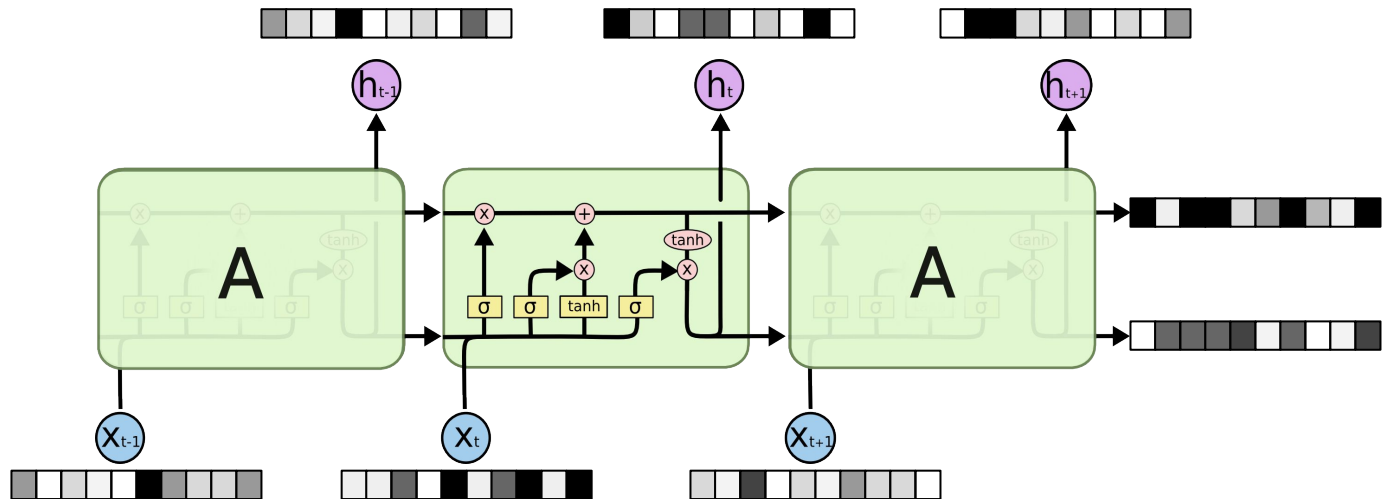
Tu *et al.* 2014



$$P(t_i|h, cache) = \lambda \cdot P_{n\text{-gram}}(t_i|h) + (1 - \lambda) \cdot P_{cache}(t_i|h)$$



# Neural Language Models

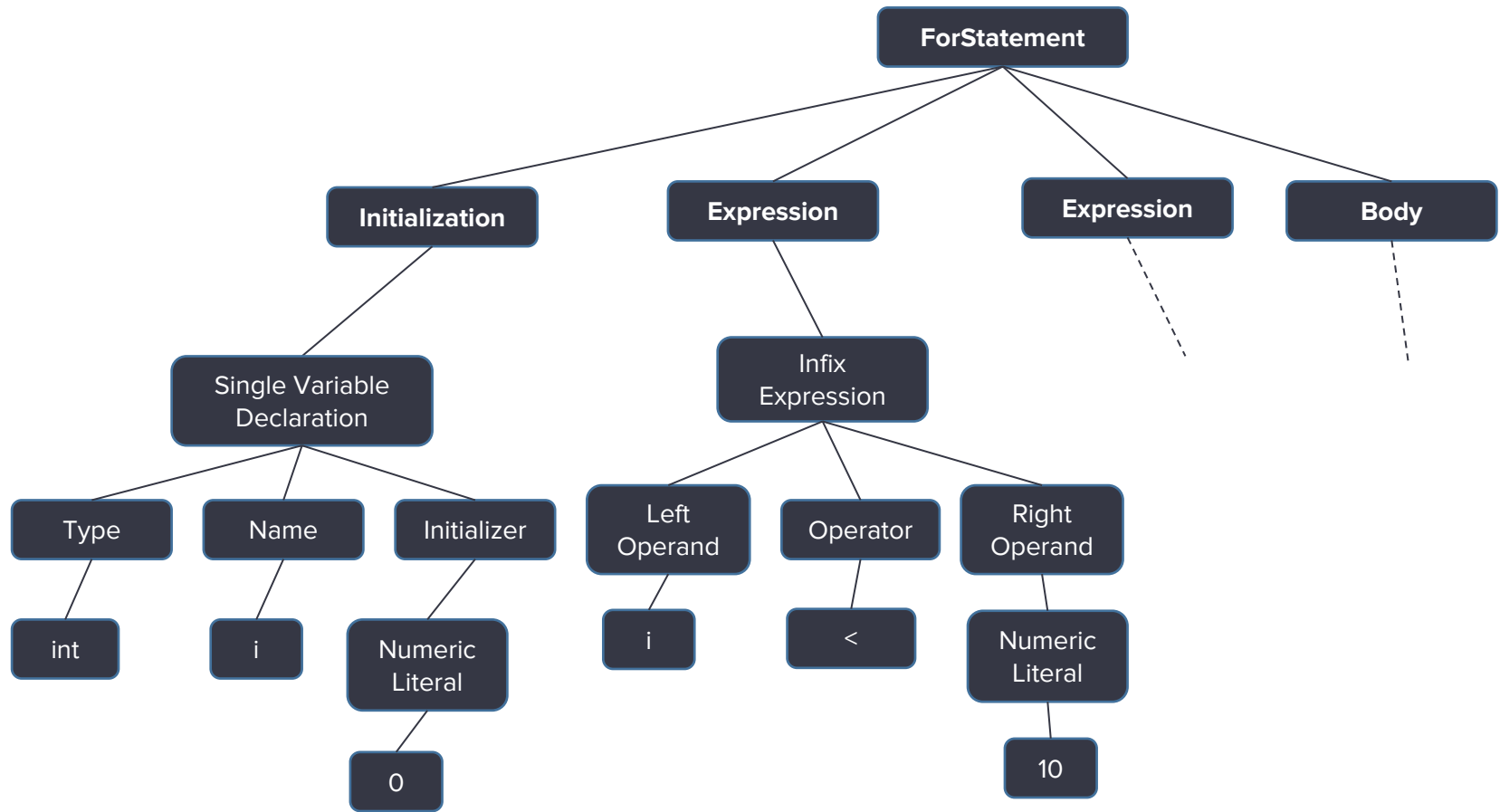


```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

Image from colah.github.io, Karpathy Andrej

Karpathy et al. 2013

White et al. 2015



**Syntactic model** of source code, *i.e.* model how AST is generated

# Tree Generation Model: Probabilistic Context Free Grammars (PCFG)

$E \rightarrow E + E$  (prob 0.7)

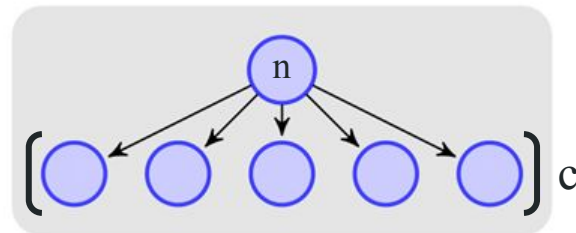
$E \rightarrow T$  (prob 0.3)

$F \rightarrow (E)$  (prob 0.1)

$T \rightarrow F * F$  (prob 0.6)

$T \rightarrow F$  (prob 0.4)

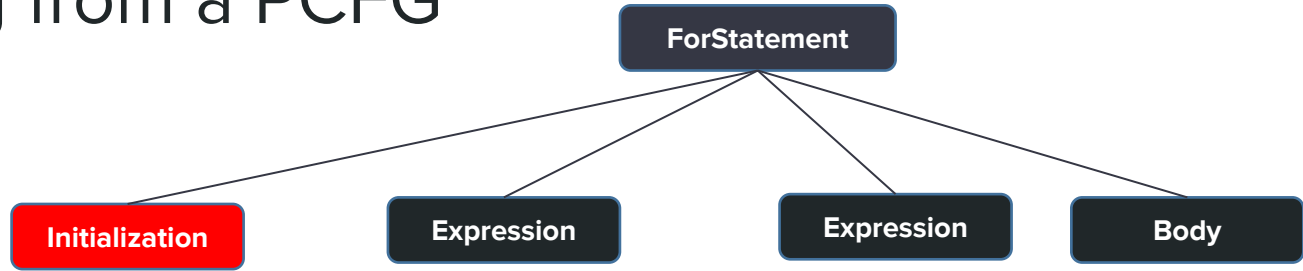
$F \rightarrow id$  (prob 0.9)



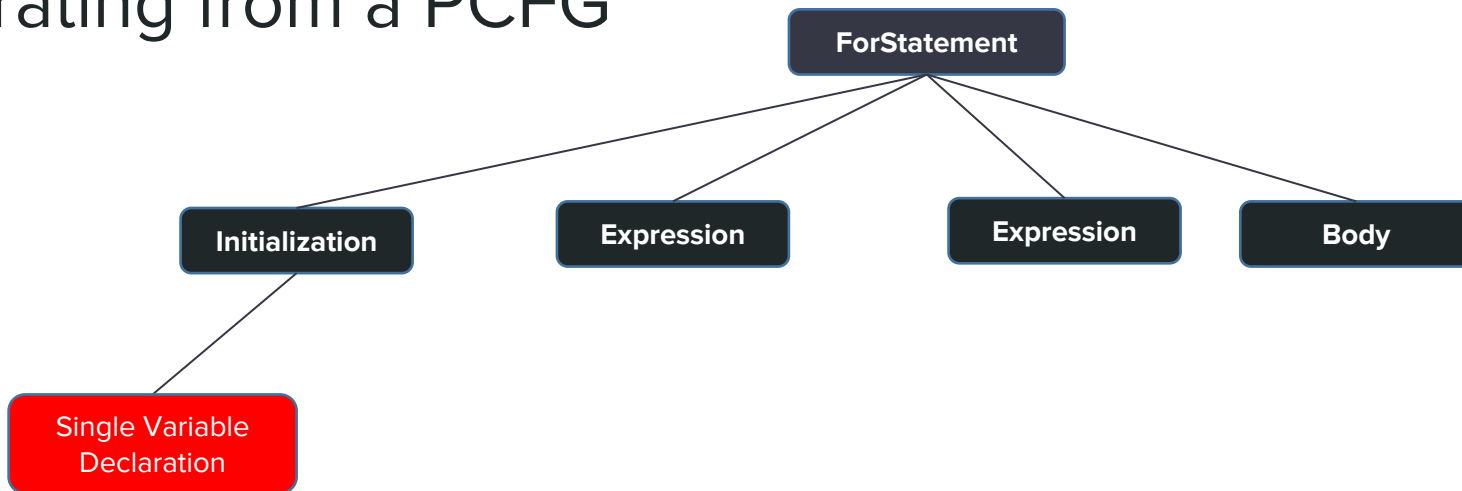
# Generating from a PCFG

ForStatement

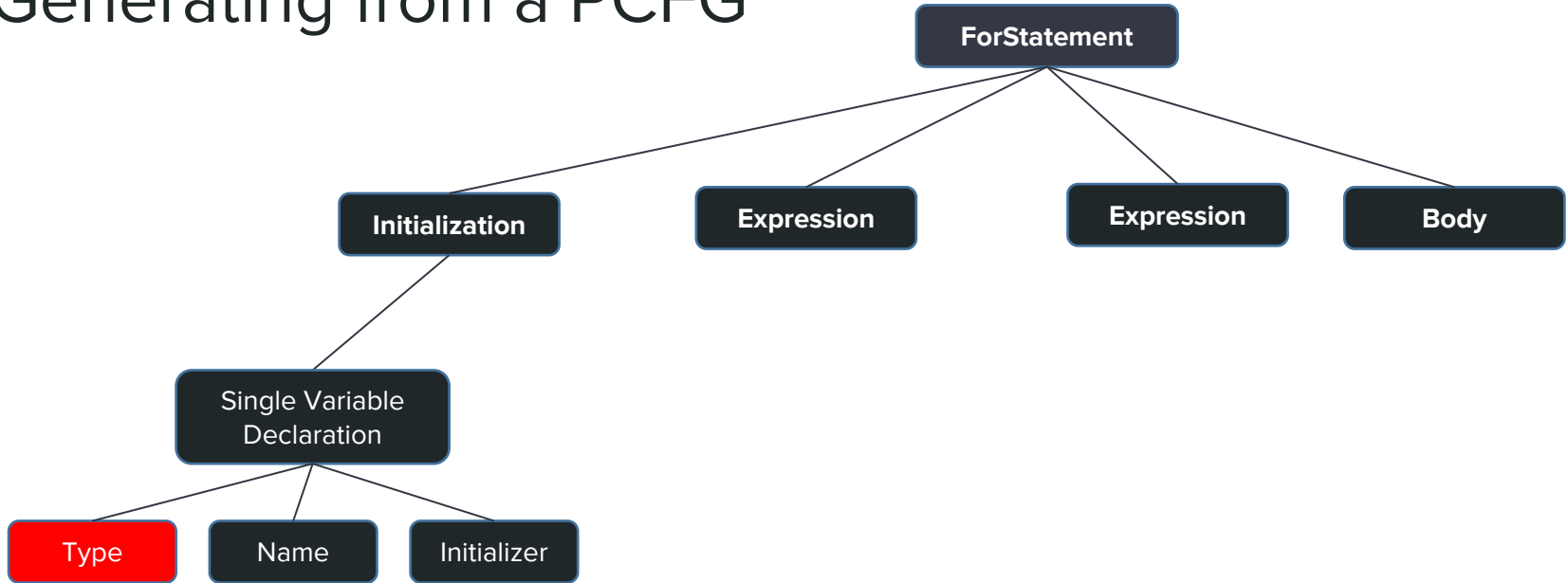
# Generating from a PCFG



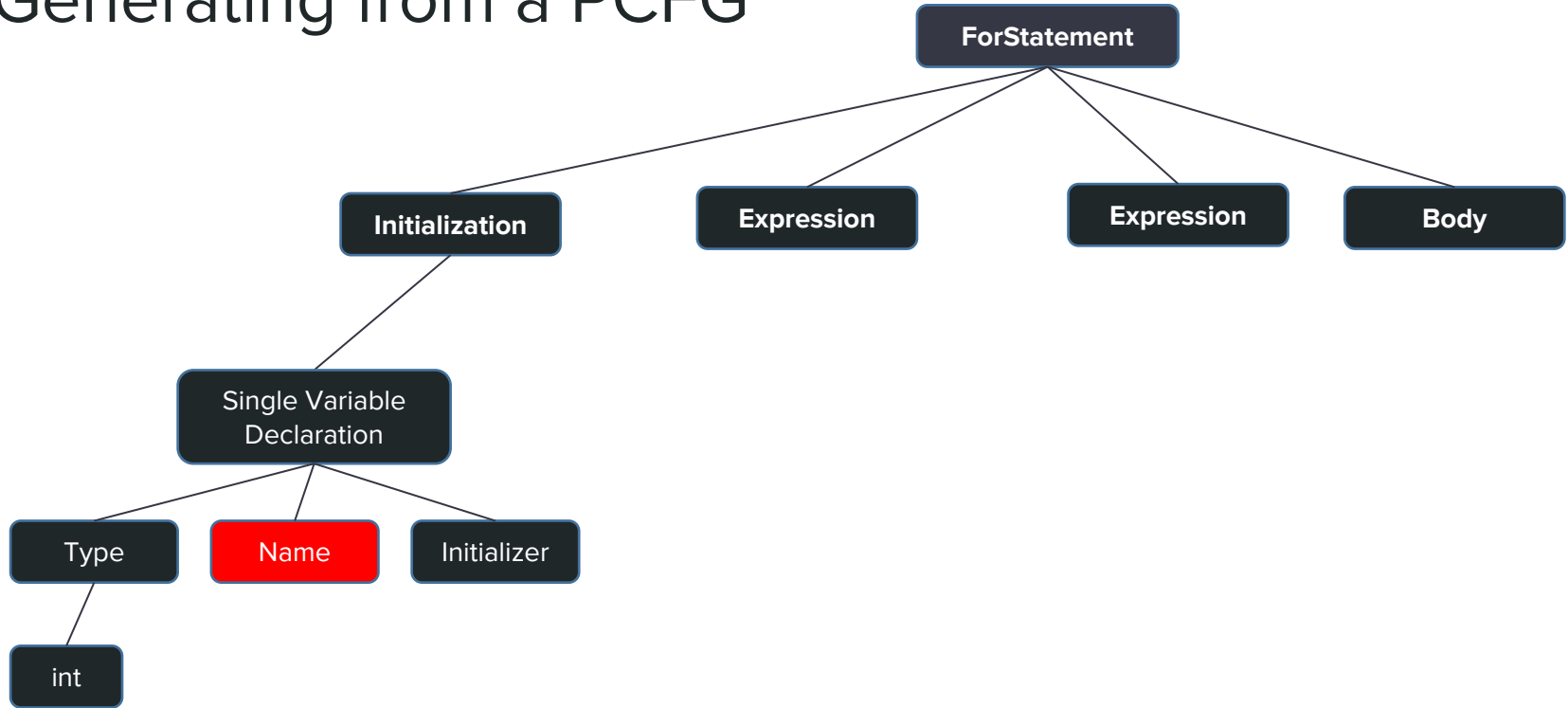
# Generating from a PCFG



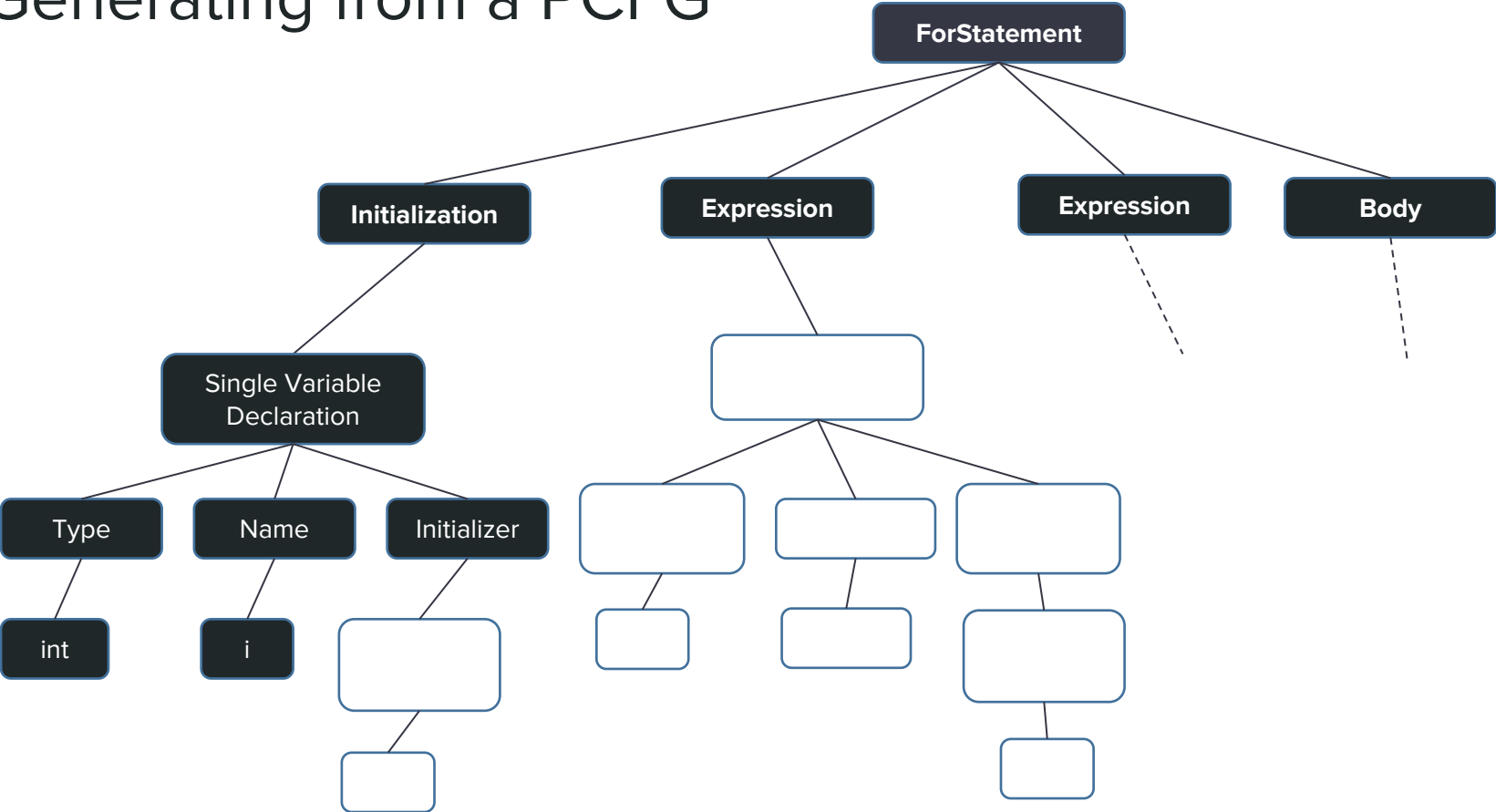
# Generating from a PCFG

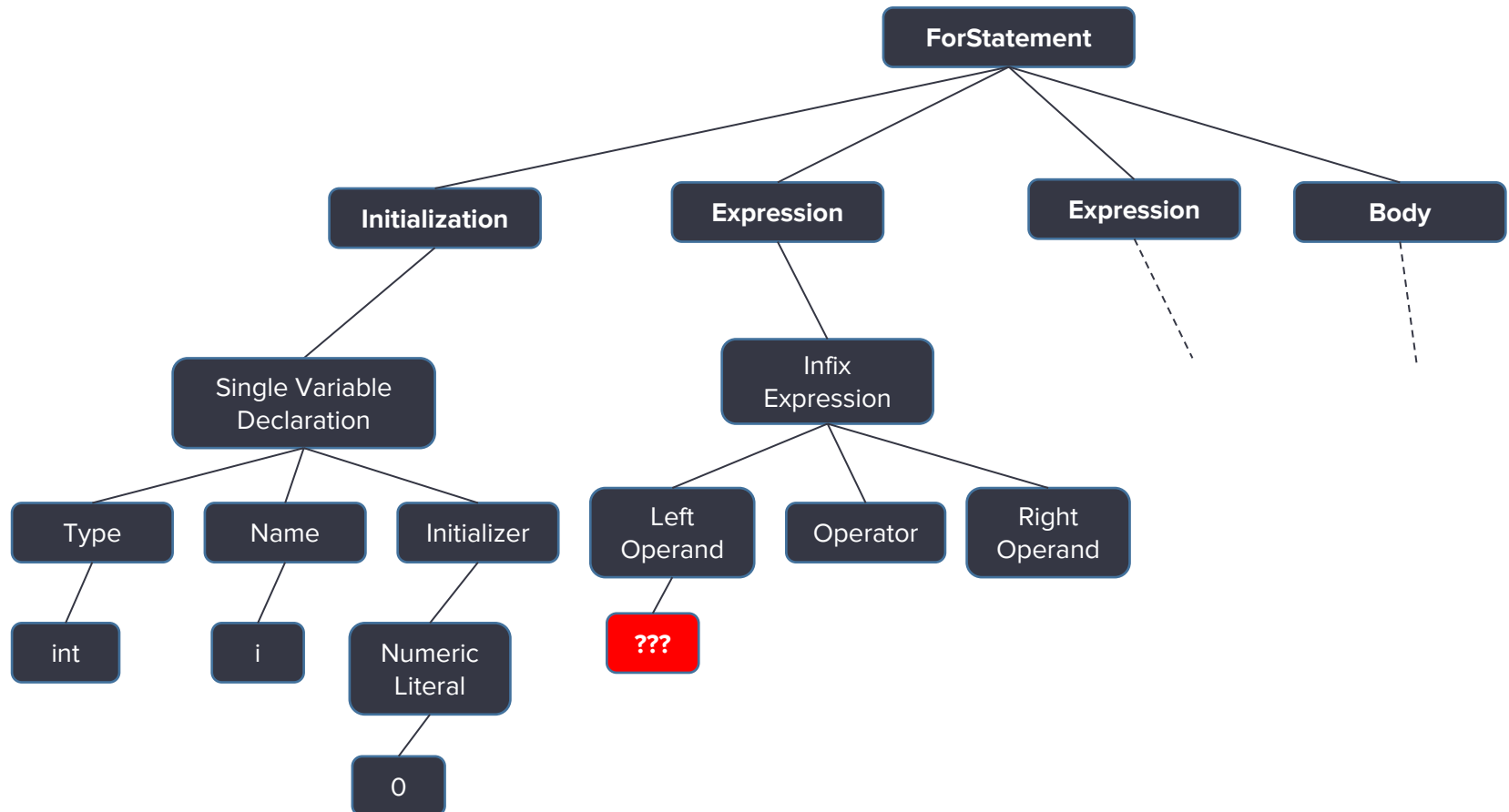


# Generating from a PCFG

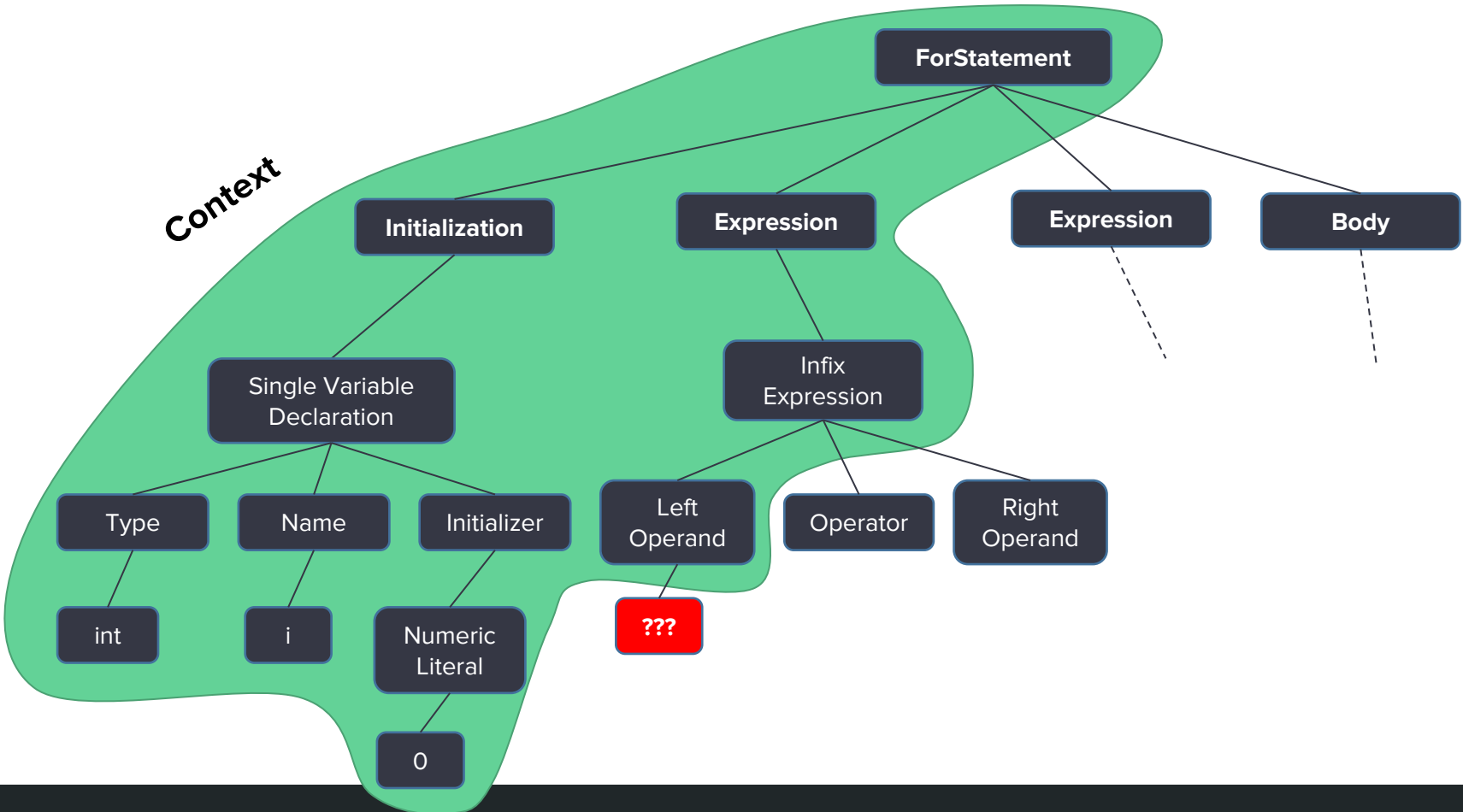


# Generating from a PCFG

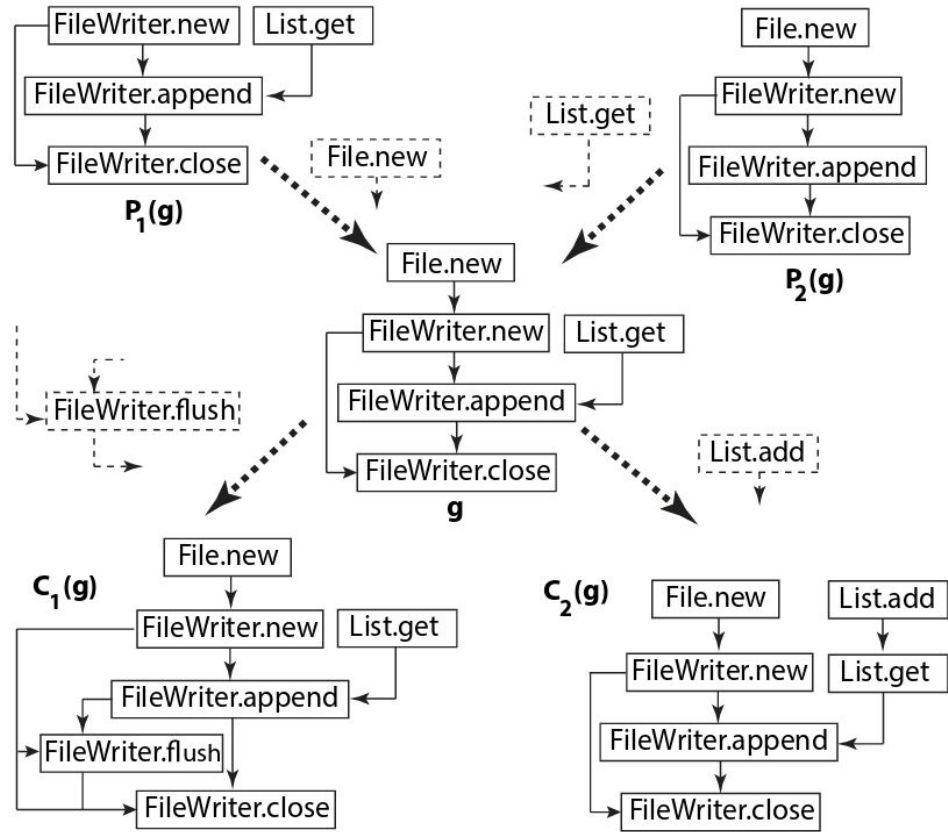




Context



# Graph Language Models



# Evaluation Metrics for Language Models



Log Probability

$$Q(\mathbb{C}, P_{\mathcal{D}}) = -\log_2 (P_{\mathcal{D}}(\mathbb{C}))$$

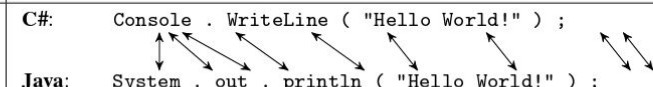
Cross Entropy

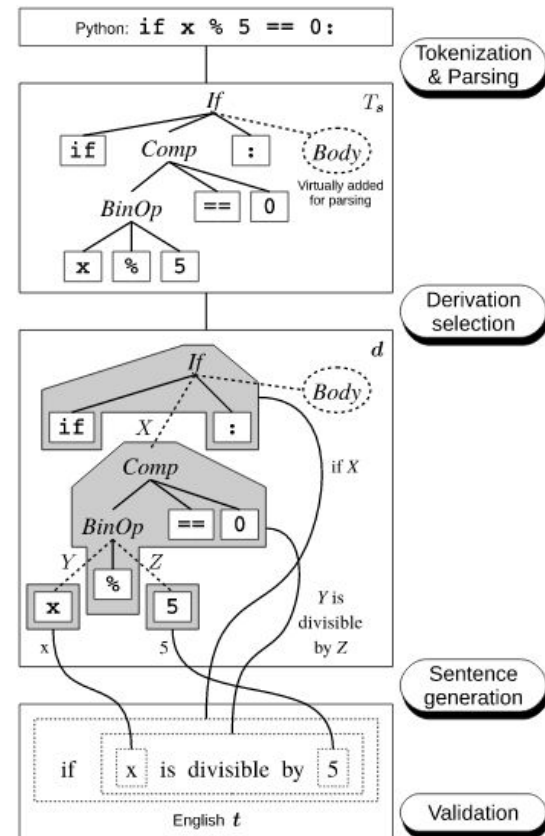
$$H(\mathbb{C}, P_{\mathcal{D}}) = -\frac{1}{M} \log_2 P_{\mathcal{D}}(\mathbb{C})$$

# Translation Models of Source Code

$$\mathfrak{t}^* = \arg \max P_{\mathcal{D}}(\mathfrak{t}|\mathfrak{s}) = \arg \max P_{\mathcal{D}}(\mathfrak{s}|\mathfrak{t}) \underbrace{P_{\mathcal{D}}(\mathfrak{t})}_{\text{Language Model}}$$

Language  
Model

Subphase	Output	Example output																		
Parallel data collection	Method pairs	<b>C#:</b> Console . WriteLine ( "Hello World!" ) ; <b>Java:</b> System . out . println ( "Hello World!" ) ;																		
Word alignment	Aligned method pairs	<b>C#:</b> Console . WriteLine ( "Hello World!" ) ;  <b>Java:</b> System . out . println ( "Hello World!" ) ;																		
Phrase table construction	Phrase table	<table border="1"> <thead> <tr> <th>C# phrase</th> <th>Java phrase</th> <th>Score i.e. <math>Pr(\text{C\# phrase}   \text{Java phrase})</math></th> </tr> </thead> <tbody> <tr> <td>Console</td> <td>System . out</td> <td>0.8</td> </tr> <tr> <td>WriteLine</td> <td>println</td> <td>0.5</td> </tr> <tr> <td>.</td> <td>.</td> <td>0.7</td> </tr> <tr> <td>(</td> <td>(</td> <td>0.9</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> </tr> </tbody> </table>	C# phrase	Java phrase	Score i.e. $Pr(\text{C\# phrase}   \text{Java phrase})$	Console	System . out	0.8	WriteLine	println	0.5	.	.	0.7	(	(	0.9	...	...	...
C# phrase	Java phrase	Score i.e. $Pr(\text{C\# phrase}   \text{Java phrase})$																		
Console	System . out	0.8																		
WriteLine	println	0.5																		
.	.	0.7																		
(	(	0.9																		
...	...	...																		



# Multimodal Models of Source Code

*“get the first letter  
of each word in  
string and  
uppercase”*

→  $P_{\mathcal{D}}(\mathbb{C}|\mathbf{m})$  →

```
string s;  
string[] words = s.ToUpper().split(' ');  
string[] firstLetters = new string[words.Length];  
for (int i=0; i < words.Length; i++) {  
    firstLetters[i] = words.Substring(0,1);  
}
```

Non-Code  
Modality

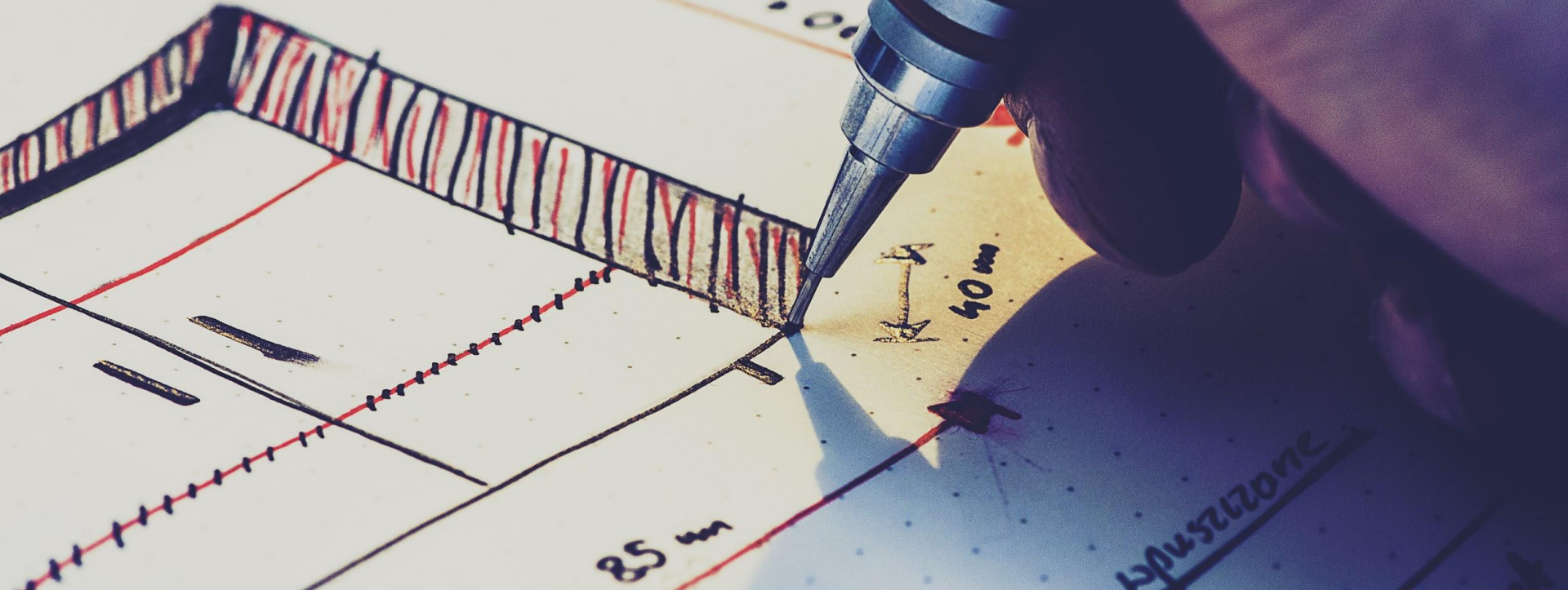
Synthesize/Score Code Snippet



- all elements to small letters for each line and sort
- all elements to small letters for each line and order
- all elements to lowercase for each line and sort
- all elements to lowercase for each line and order
- all elements lowercase for each line and sort
- all elements lowercase for each line and order
- all elements to lower case for each line and sort
- all elements to lower case for each line and order
- all elements to small letters for each new line and sort



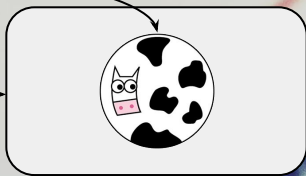
```
string result = String.Join("\n",input_string.Split('\n'))  
    .Select((string x) => x.ToLower()).OrderBy(x => x);
```



# Representational Models of Code

representation

code snippet

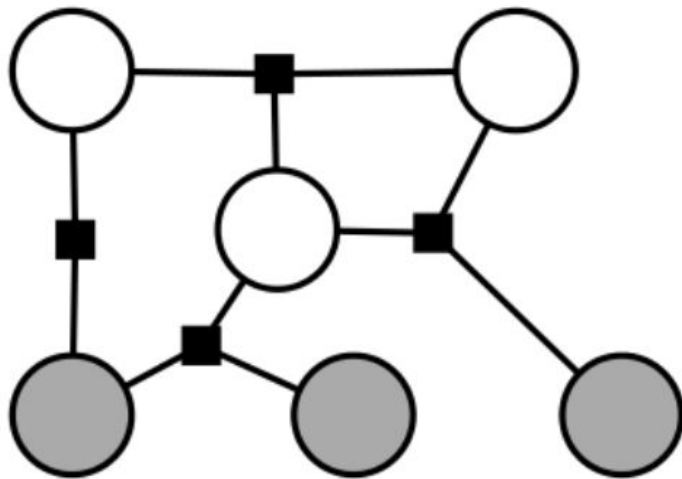


property of code



$$P_{\mathcal{D}}(\pi|f(c))$$

# Structured Prediction



$$P_{\mathcal{D}}(\boldsymbol{\pi} | f(\mathbf{c}))$$

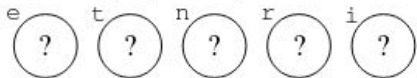
› Output variables are mutually dependent or constrained.

```
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;
}
```

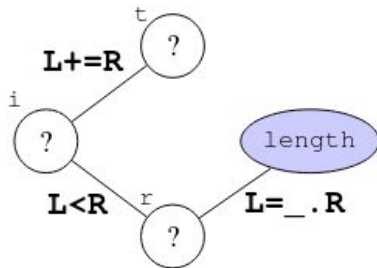
(a) JavaScript program with minified identifier names



Unknown properties (variable names):

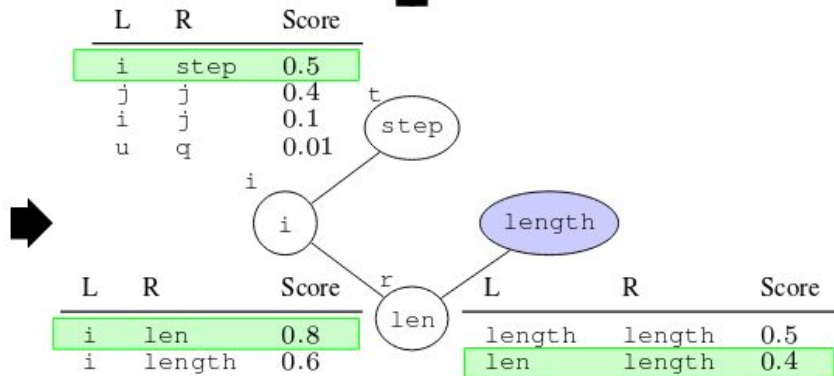


Known properties (constants, APIs):

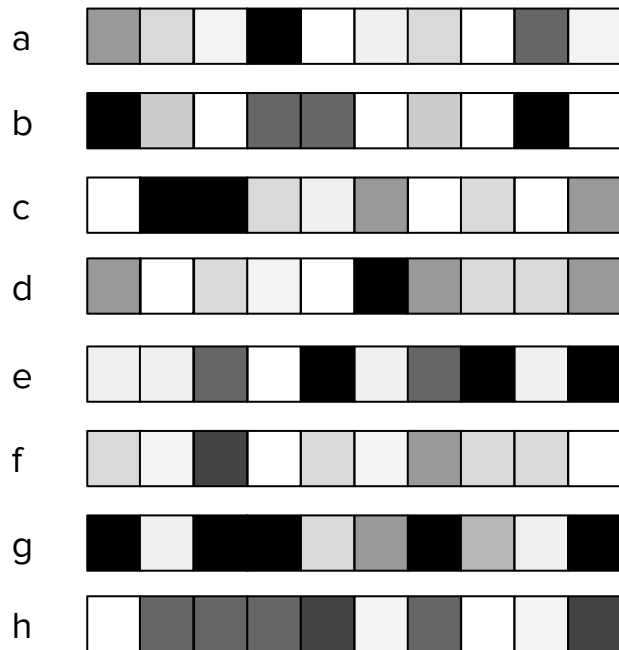
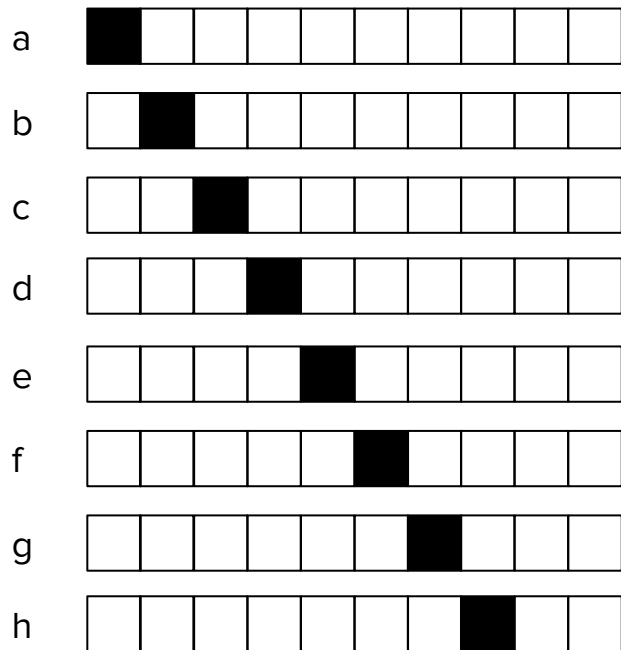


```
/* str: string, step: number, return: Array */
function chunkData(str, step) {
  var colNames = []; /* colNames: Array */
  var len = str.length;
  var i = 0; /* i: number */
  for (; i < len; i += step) {
    if (i + step < len) {
      colNames.push(str.substring(i, i + step));
    } else {
      colNames.push(str.substring(i, len));
    }
  }
  return colNames;
}
```

(e) JavaScript program with new identifier names and types



# Localized vs Distributed (Vector) Representations



$$P_{\mathcal{D}}(\boldsymbol{\pi} | f(\mathbb{c}))$$

$\downarrow$   
 $\mathbb{R}^D$



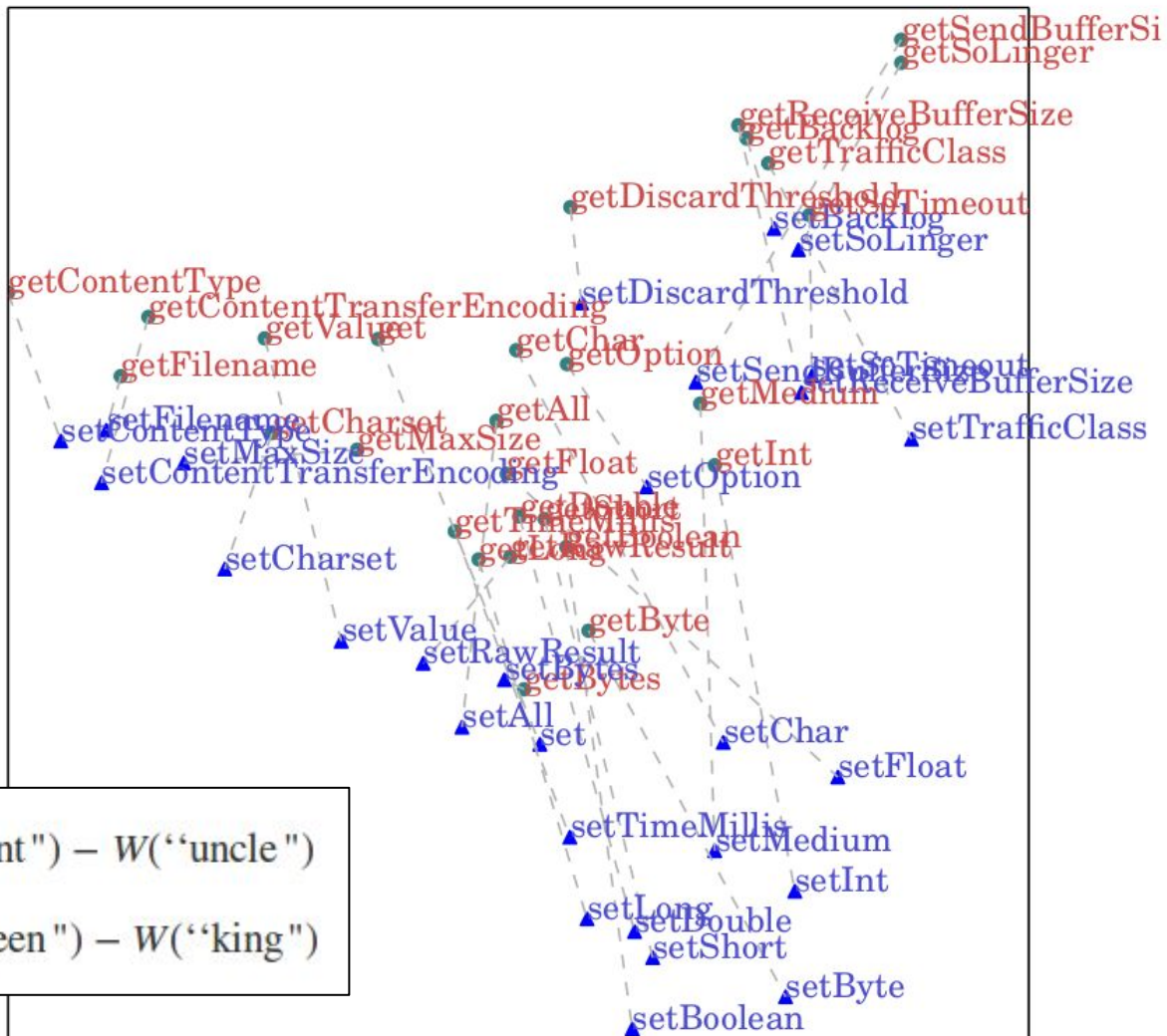


# Distributed Representations



Allamanis *et al.*, “Suggesting accurate method and class names”, 2015

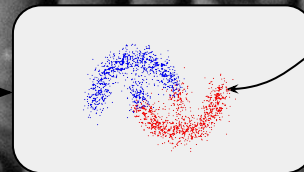
# Setters vs Getters in bigbluebutton



$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"aunt"}) - W(\text{"uncle"})$

$W(\text{"woman"}) - W(\text{"man"}) \simeq W(\text{"queen"}) - W(\text{"king"})$

corpus  
of code

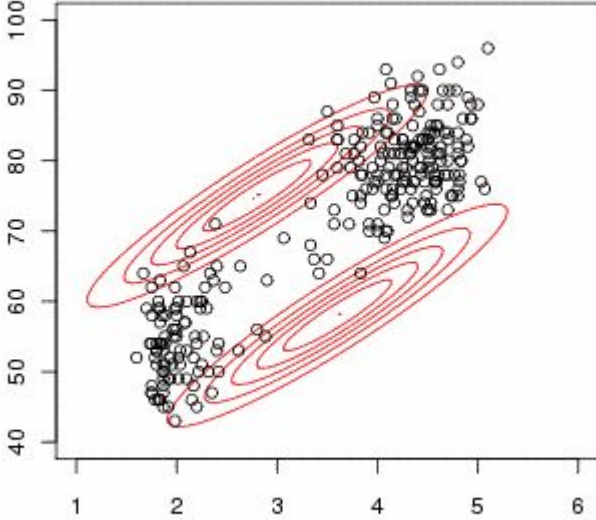
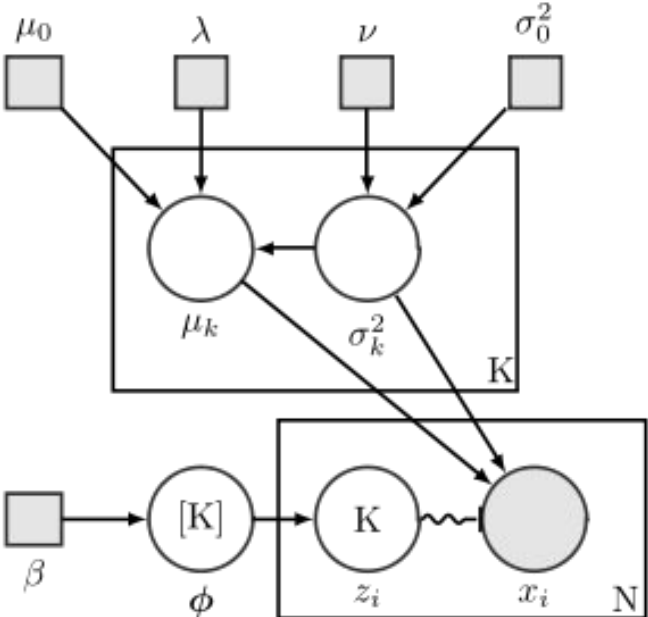


patterns

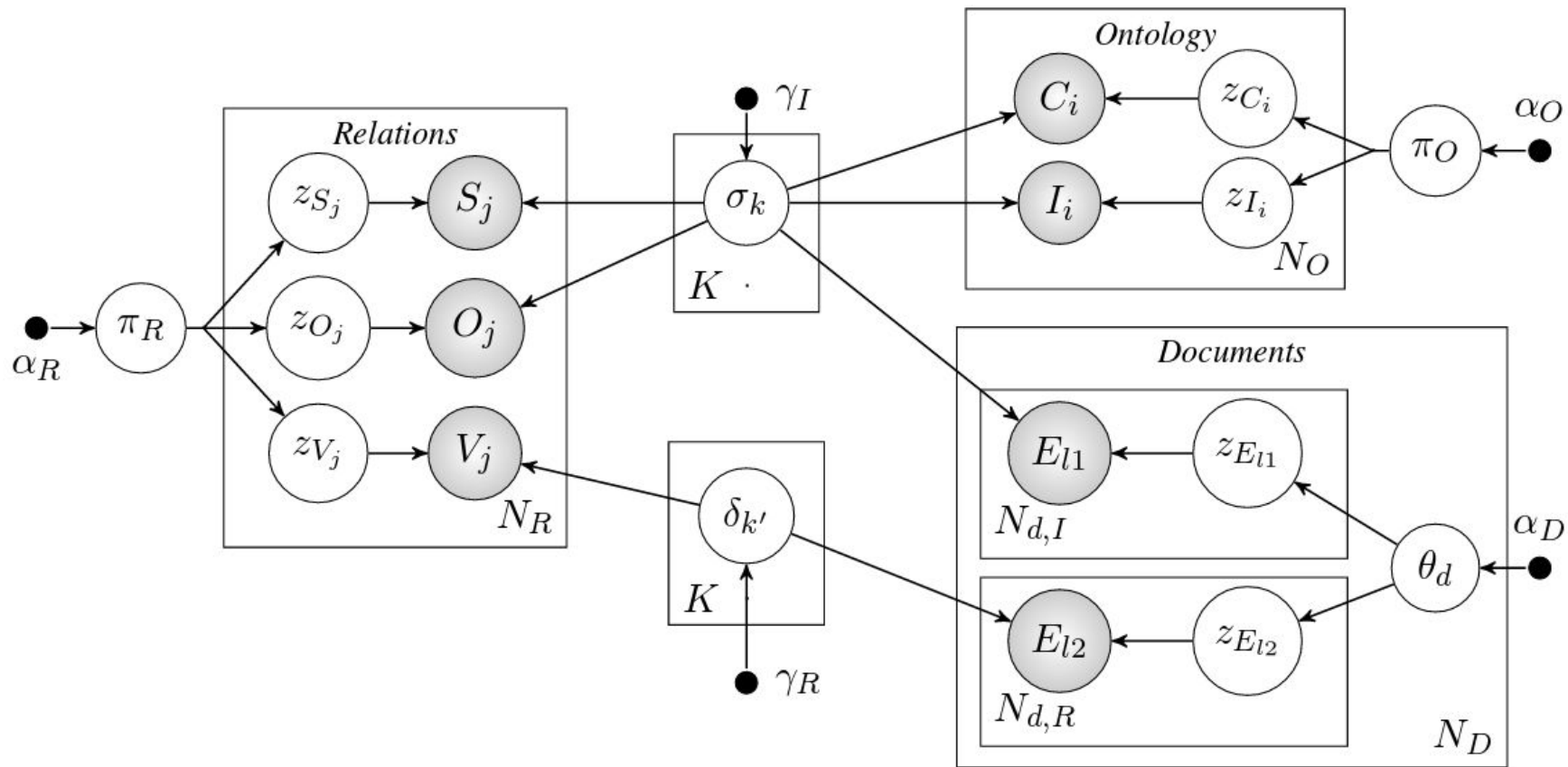
## Pattern Mining Models of Code

$$P_{\mathcal{D}}(f(\mathbb{C})) = \sum_{\mathbf{l}} P_{\mathcal{D}}(f(\mathbb{C})|\mathbf{l})P(\mathbf{l})$$

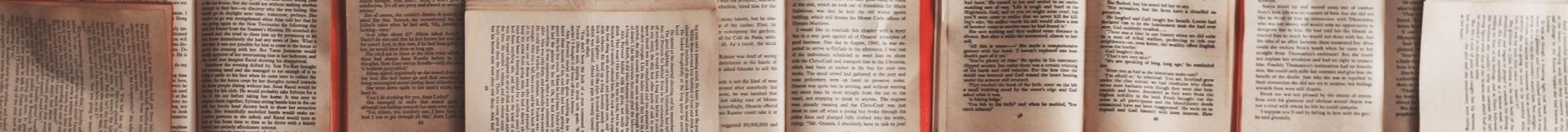
# Latent Variable Models



Images from Wikipedia



Movshovitz-Attias & Cohen, "KB-LDA: Jointly learning a knowledge base of hierarchy, relations, and facts", 2015



# Part III Applications of Models of Source Code



**Machine learning** helps each application area by **learning** to resolve (some) **ambiguities** from **data**.

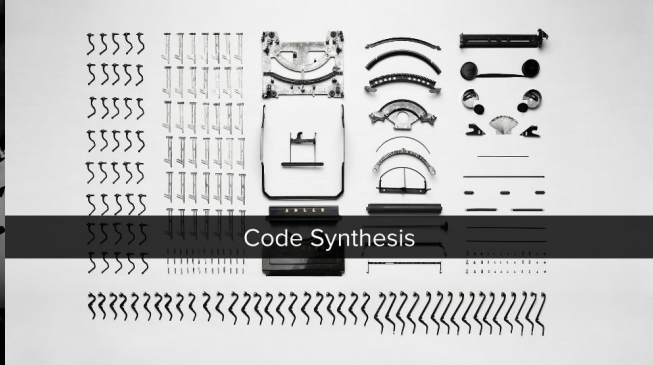




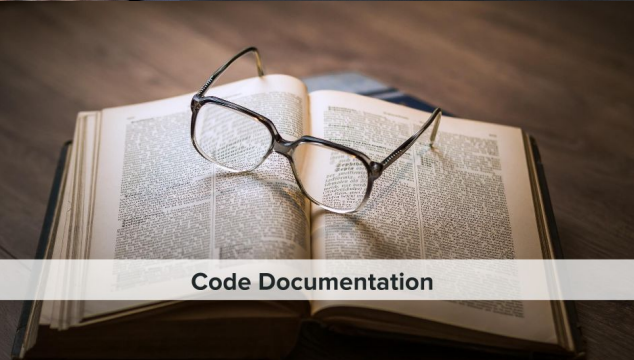
Recommender Systems



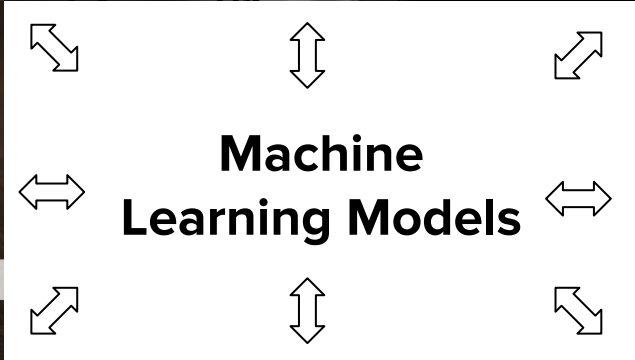
Statistical Code Migration



Code Synthesis



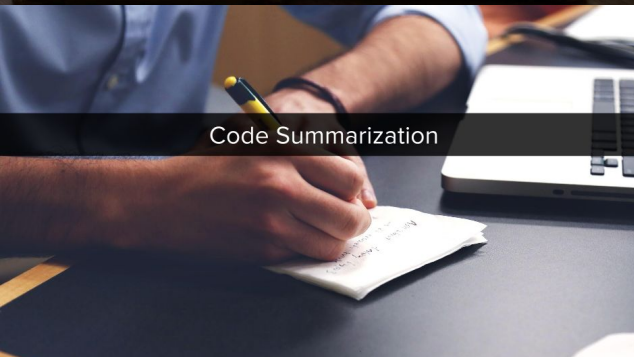
Code Documentation



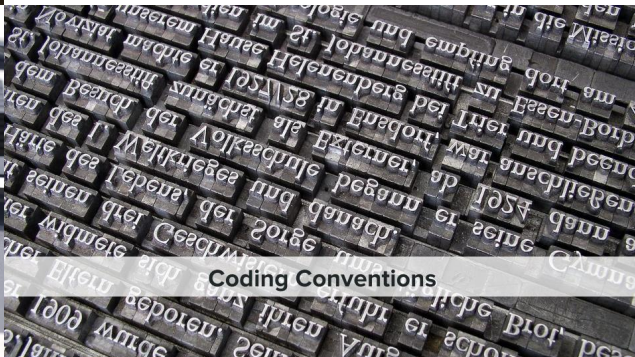
Machine Learning Models



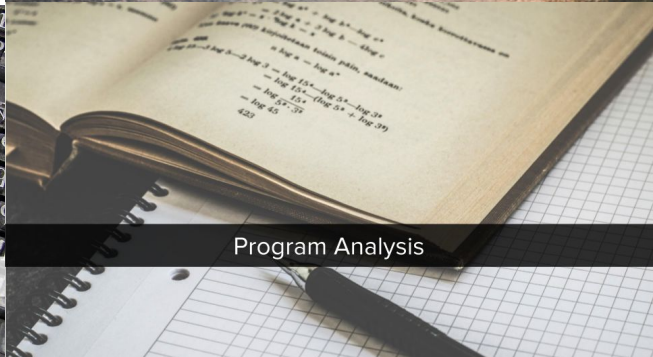
Code Defects



Code Summarization

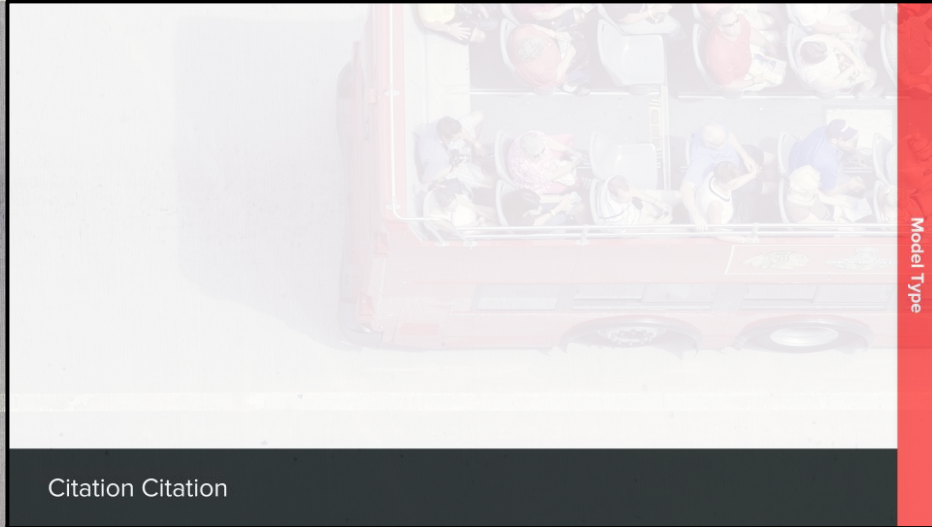


Coding Conventions



Program Analysis

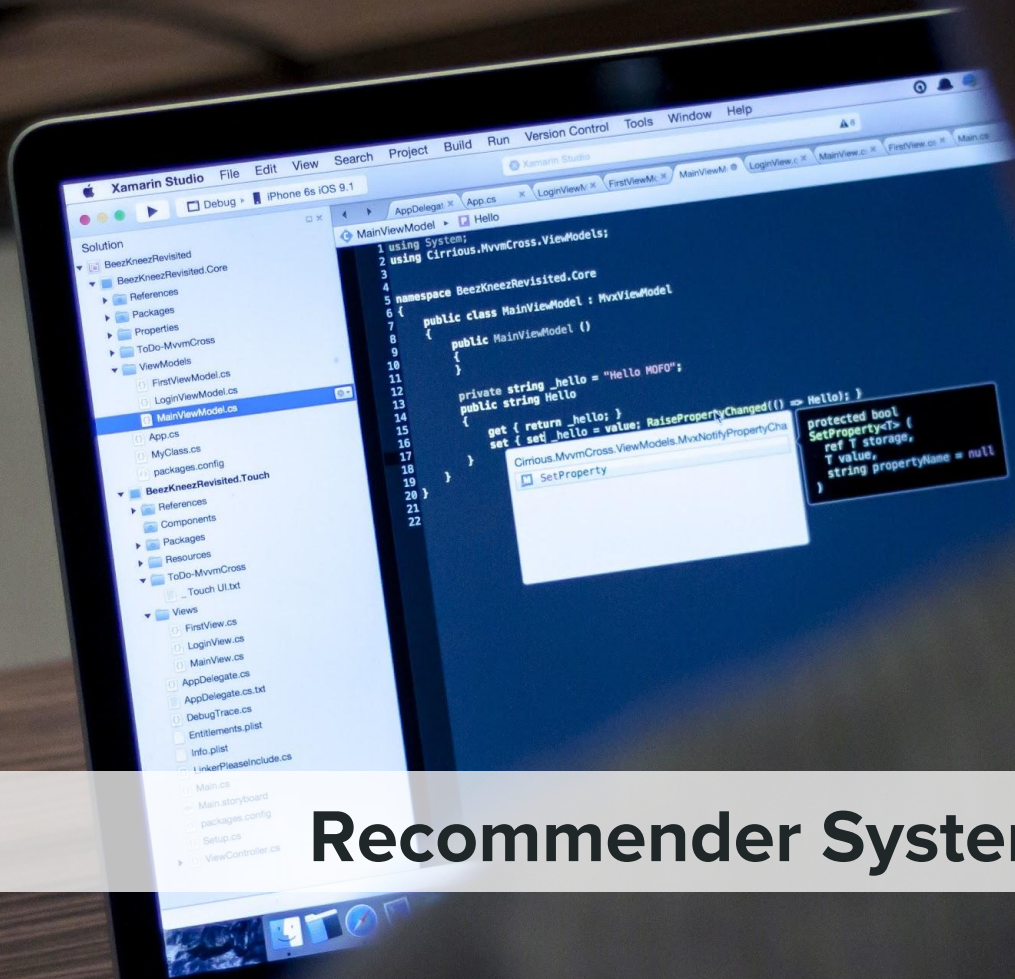
# Applications: A (Very) Brief Tour



Model Type



Citation Citation



# Recommender Systems

```
public static void main(String[] args) {
```

```
    " "
}
```

- Code Recommenders is enabled
- charAt(int index) : char - String
- chars() : IntStream - CharSequence
- codePointAt(int index) : int - String
- codePointBefore(int index) : int - String
- codePointCount(int beginIndex, int endIndex)
- codePoints() : IntStream - CharSequence
- compareTo(String anotherString) : int - String
- compareToIgnoreCase(String str) : int - String
- concat(String str) : String - String

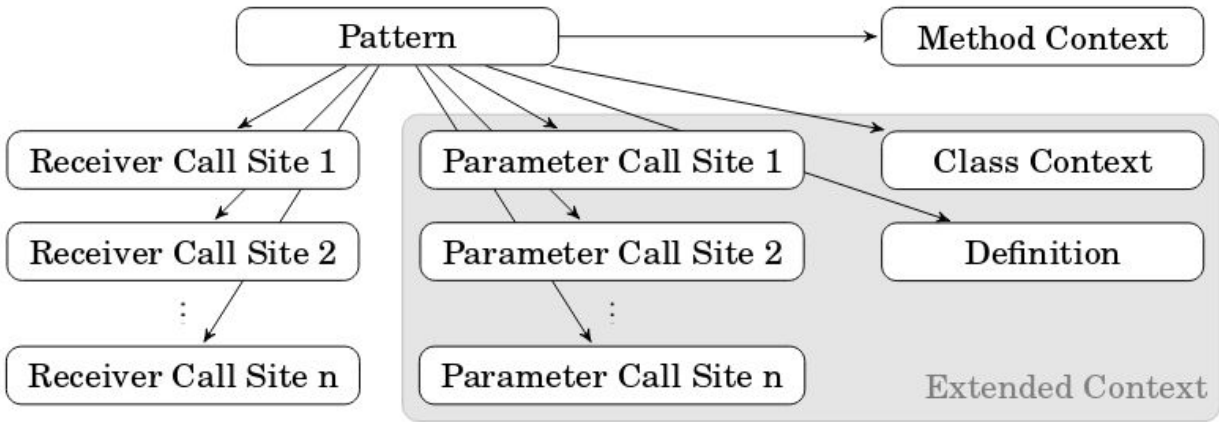
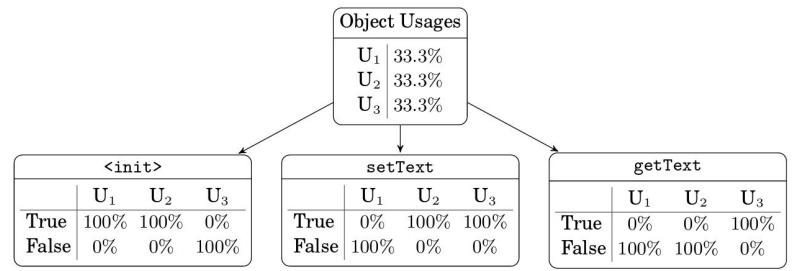
Code Recommenders improves Eclipse's built-in code completion by showing you how others have used an API before. Additionally, subwords completion allows you to search for subword matches rather than just prefix matches in Content Assist.

By default, both Code Recommenders and subwords completion are enabled but you may customize their behavior at [Code Recommenders > Completions](#).

To learn more about Code Recommenders and Subwords, please visit the [project homepage](#) or consult the [manual](#).

Press '^Space' to show Template Proposals

Eclipse Code Recommenders



Token Role	Construction Rule
Data type $T$	TYPE[ $T$ ]
Variable $x$	VAR[typeof( $x$ )]
Literal $v$	LIT[typeof( $v$ )]
Function decl $m$	FUNC[type( $m$ ),lexeme( $m$ ),paralist( $m$ ),rettype( $m$ )]
Function call $m$	CALL[type( $m$ ),lexeme( $m$ ),paracount( $m$ ),rettype( $m$ )]
Parameter $x$	PARA[typeof( $x$ )]
Field $f$	FIELD[type( $f$ ), lexeme( $f$ )]
Operator $o$	OP[name( $o$ )]
Cast ( $T$ )	CAST[ $T$ ]
Keyword	To corresponding reserved token
Block open & close	To corresponding reserved token
Special literal	To corresponding reserved token
Unknown	To special lexical token LEX

```

offset -= replacementLength;
char c = ctx.getDocument().getChar(offset);
boolean prevCharWasWhitespace = false;
while (numberOfSeparators < 10 && offset >= 0) {
    if (!(Character.isWhitespace(ctx.getDocument().getChar(
        prev
        b.ap
    )
    )
    else {
        if (p
        prev
    )
    if (c ==
        numb
    )
    offset--
    if (offs
        brea
    )
    c = ctx.
}

```

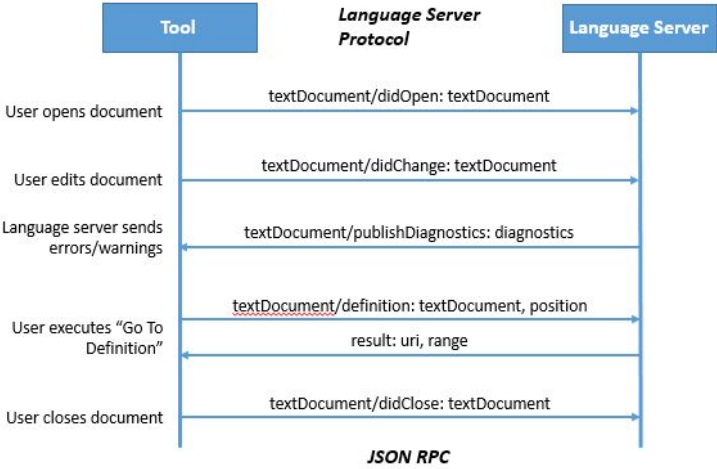
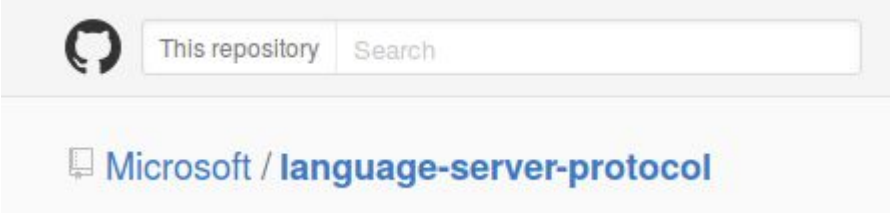
The screenshot shows a list of method suggestions for the call `ctx.getDocument()`. The suggestions include:

- `getDocument() : IDocument - ContentAssistInvocationContext`
- `lookup`
- `computeIdentifierPrefix() : CharSequence - ContentAssistInvocationContext`
- `equals(Object obj) : boolean - ContentAssistInvocationContext`
- `getClass() : Class<?> - Object`
- `getCompilationUnit() : ICompilationUnit - JavaContentAssistInvocationContext`
- `getCoreContext() : CompletionContext - JavaContentAssistInvocationContext`
- `getExpectedType() : IType - JavaContentAssistInvocationContext`
- `getHistoryRelevance(String qualifiedTypeName) : float - JavaContentAssistInvocationContext`
- `getInvocationOffset() : int - ContentAssistInvocationContext`
- `getKeywordProposals() : IJavaCompletionProposal[] - JavaContentAssistInvocationContext`

At the bottom of the list, it says "Press 'Ctrl+Space' to show CACHECA Proposals".

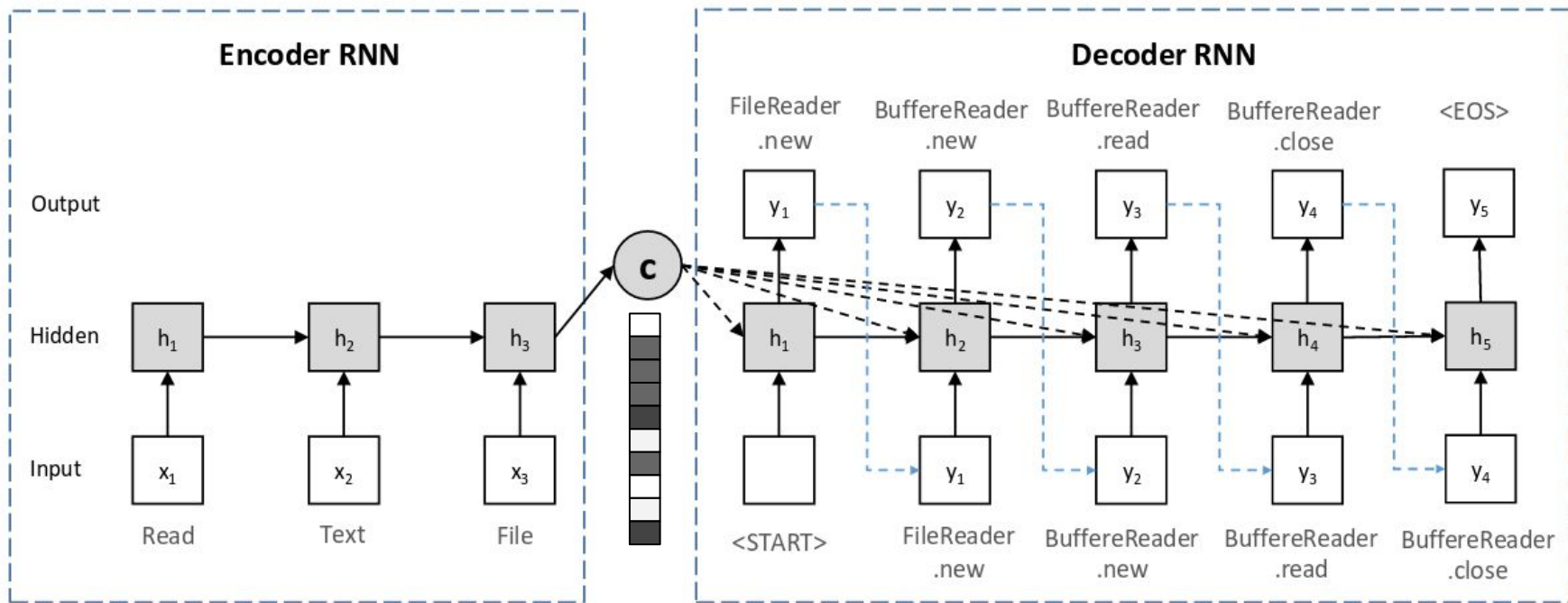


# Language Server Protocol



An open book with a pair of glasses resting on it, symbolizing reading and documentation. The book is open to a page with dense text, and the glasses are positioned over the text. The background is a soft, out-of-focus light color.

# Code Documentation





```
TwitterFactory.<init>  
TwitterFactory.getInstance
```

```
Status.getUser  
Status.getText
```

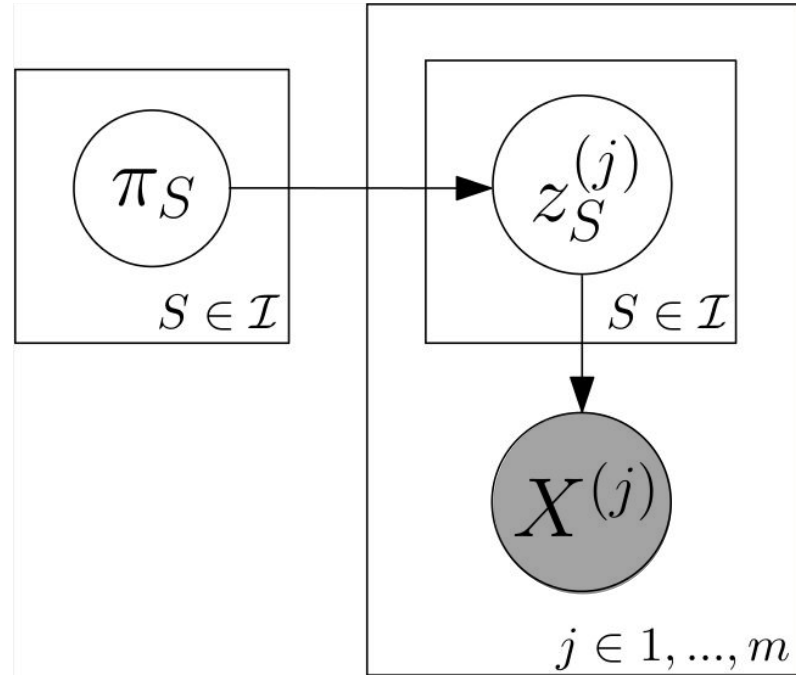
```
ConfigurationBuilder.<init>  
ConfigurationBuilder.build
```

```
ConfigurationBuilder.<init>  
TwitterFactory.<init>
```

```
ConfigurationBuilder.<init>  
ConfigurationBuilder.setOAuthConsumerKey
```


```
ConfigurationBuilder.build  
TwitterFactory.<init>
```

```
ConfigurationBuilder.<init>  
ConfigurationBuilder.build  
TwitterFactory.<init>
```





# Statistical Code Migration

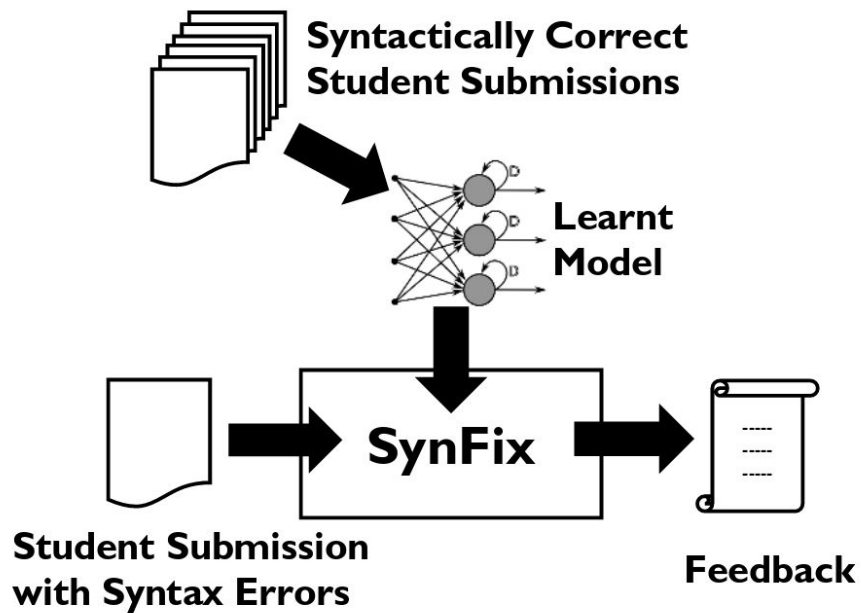
Subphase	Output	Example output																		
Parallel data collection	Method pairs	<b>C#:</b> Console . WriteLine ( "Hello World!" ) ; <b>Java:</b> System . out . println ( "Hello World!" ) ;																		
Word alignment	Aligned method pairs	<b>C#:</b> Console . WriteLine ( "Hello World!" ) ;  <b>Java:</b> System . out . println ( "Hello World!" ) ;																		
Phrase table construction	Phrase table	<table border="1"> <thead> <tr> <th>C# phrase</th> <th>Java phrase</th> <th>Score i.e. <math>Pr(\text{C\# phrase}   \text{Java phrase})</math></th> </tr> </thead> <tbody> <tr> <td>Console</td> <td>System . out</td> <td>0.8</td> </tr> <tr> <td>WriteLine</td> <td>println</td> <td>0.5</td> </tr> <tr> <td>.</td> <td>.</td> <td>0.7</td> </tr> <tr> <td>(</td> <td>(</td> <td>0.9</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> </tr> </tbody> </table>	C# phrase	Java phrase	Score i.e. $Pr(\text{C\# phrase}   \text{Java phrase})$	Console	System . out	0.8	WriteLine	println	0.5	.	.	0.7	(	(	0.9	...	...	...
C# phrase	Java phrase	Score i.e. $Pr(\text{C\# phrase}   \text{Java phrase})$																		
Console	System . out	0.8																		
WriteLine	println	0.5																		
.	.	0.7																		
(	(	0.9																		
...	...	...																		

$$\mathbb{t}^* = \arg \max P_{\mathcal{D}}(\mathbb{t} | \mathbb{s}) = \arg \max P_{\mathcal{D}}(\mathbb{s} | \mathbb{t}) P_{\mathcal{D}}(\mathbb{t})$$



## Code Defects



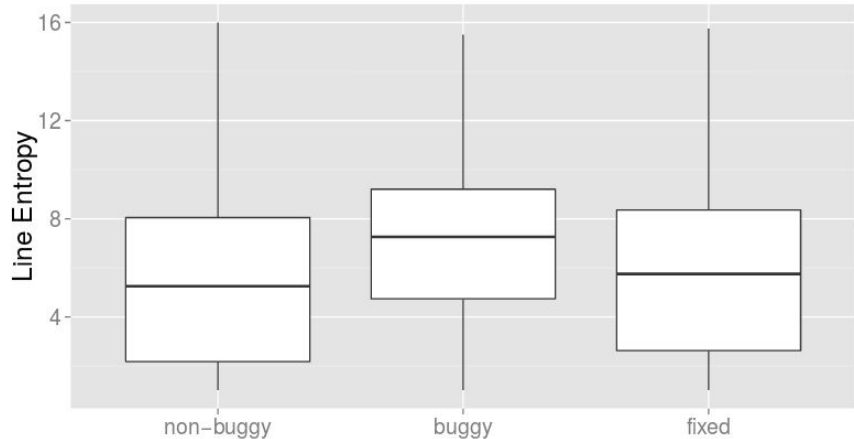


```

1 for (int i = 0; i < scorers.length; i++) {
2     if (scorers[i].nextDoc() == NO_MORE_DOCS)
5         lastDoc = NO_MORE_DOCS;
6         return;
7     }
8 }

```





Netty (2013-08-20)

File: ThreadPerChannelEventLoopGroup.java

Entropy dropped after bugfix : **4.6257**

```
if (isTerminated()) {  
    // Before (entropy = 5.96485):  
-   terminationFuture.setSuccess(null);  
    // After (entropy = 1.33915):  
+   terminationFuture.trySuccess(null);  
}
```

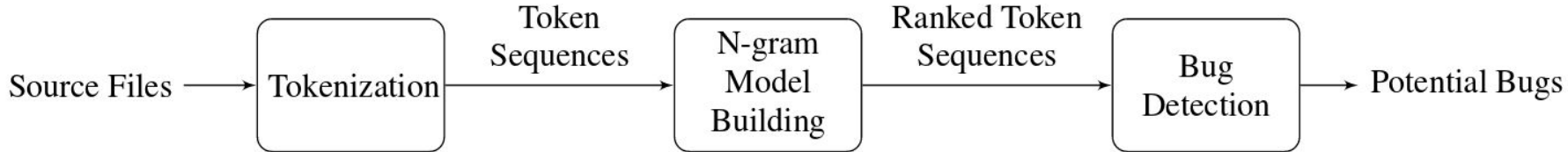
Ray et al. "On the Naturalness of Buggy Code", 2016



```
5
6
7 # CodexLint Example:
8 "a string".split("\n").to_s
9
10
```



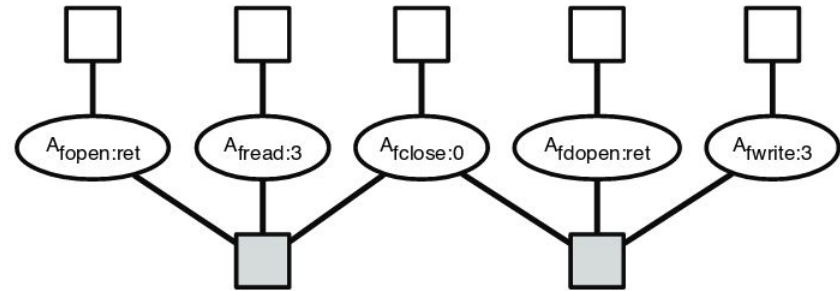
Function to\_s has appeared 12 times and split has appeared 29 times, and they've appeared 0 times together.



---

1. FILE \* fp1 = fopen( "myfile.txt", "r" );
2. FILE \* fp2 = fdopen( fd, "w" );
3. fread( buffer, n, 1, fp1 );
4. fwrite( buffer, n, 1, fp2 );
5. fclose( fp1 );
6. fclose( fp2 );

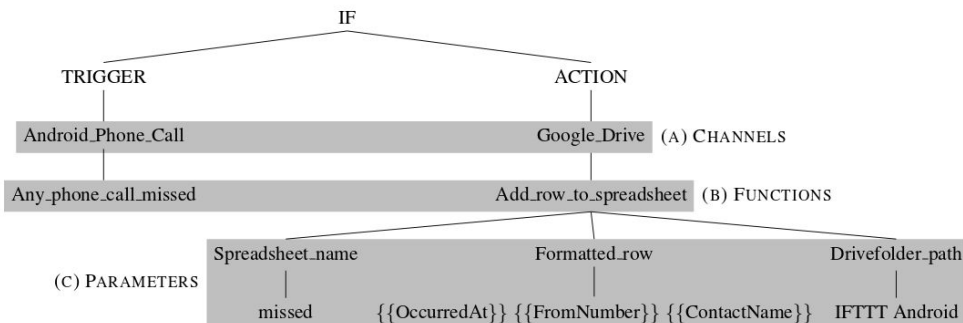
---





# Code Synthesis





Archive your missed calls from Android to Google Drive

sum the totalpay for the capitol hill baristas

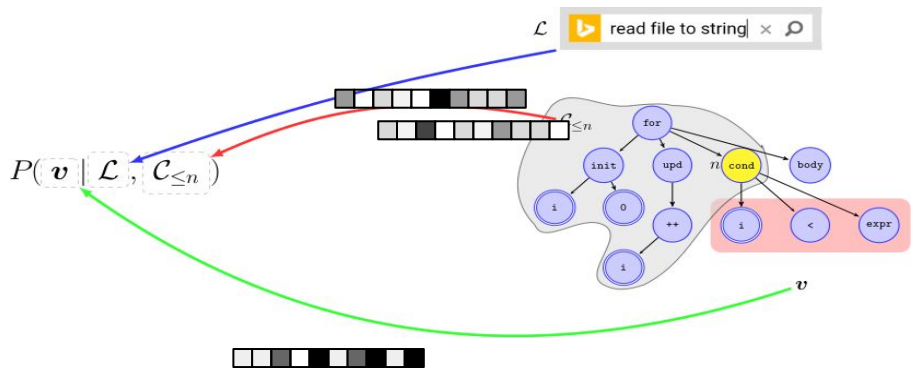
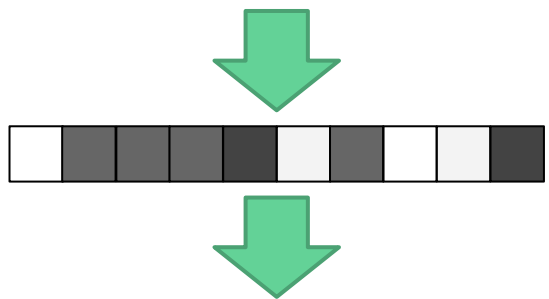
sum the totalpay for the capitol hill baristas =SUMIFS(WeeklyHours[totalpay], WeeklyHours[location], "capitol hill", WeeklyHours[title], "barista", WeeklyHours[totalpay])  
 sum the totalpay for the capitol hill baristas =SUMIF(WeeklyHours[title], "barista", WeeklyHours[totalpay])  
 sum the totalpay for the capitol hill baristas =SUMIF(WeeklyHours[location], "capitol hill", WeeklyHours[totalpay])

I2 : x ✓ ✕ =SUMIFS(WeeklyHours[totalpay], WeeklyHours[location], "capitol hill", WeeklyHours[title], "barista")

	A	B	C	D	E	F	G	H	I
1	location	name	title	hours	othours	basepay	otpay	totalpay	
2	capitol hill	aaron	chef	18	0	\$243.00	\$0.00	\$243.00	\$888.80
3	capitol hill	blanca	cashier	25	0	\$193.75	\$0.00	\$193.75	
4	capitol hill	chris	manager	40	10	\$990.00	\$371.25	\$1,361.25	
5	capitol hill	deeraj	barista	40	0	\$352.00	\$0.00	\$352.00	
6	capitol hill	grace	barista	21	0	\$184.80	\$0.00	\$184.80	
7	capitol hill	hannah	cashier	40	3	\$310.00	\$34.88	\$344.88	
8	capitol hill	irene	barista	40	0	\$352.00	\$0.00	\$352.00	
9	queen anne	tom	chef	16	0	\$216.00	\$0.00	\$216.00	
10	queen anne	susan	barista	26	0	\$228.80	\$0.00	\$228.80	
11	queen anne	steve	cashier	22	0	\$170.50	\$0.00	\$170.50	

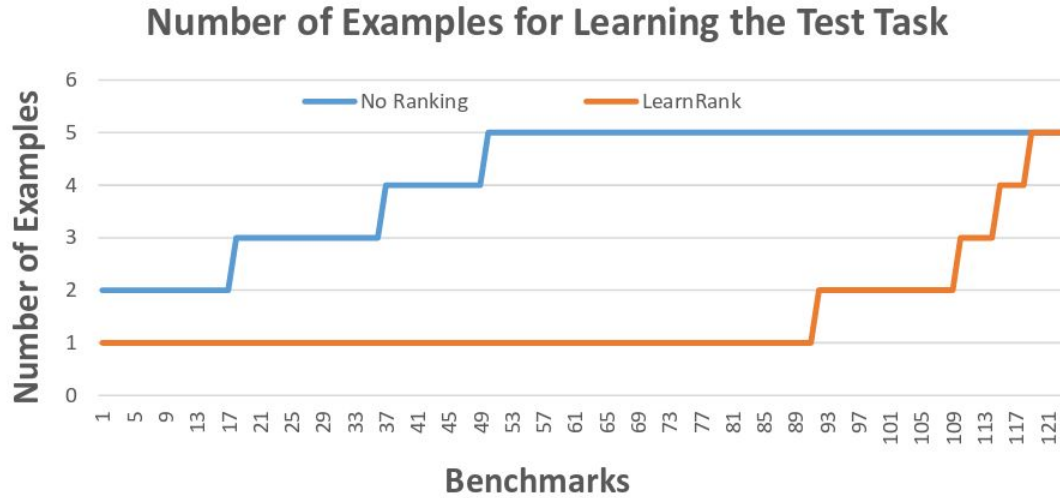


- all elements to small letters for each line and sort
- all elements to small letters for each line and order
- all elements to lowercase for each line and sort
- all elements to lowercase for each line and order
- all elements lowercase for each line and sort
- all elements lowercase for each line and order
- all elements to lower case for each line and sort
- all elements to lower case for each line and order
- all elements to small letters for each new line and sort



```
string result = String.Join("\n",input_string.Split('\n'))
    .Select((string x) => x.ToLower()).OrderBy(x => x);
```





A close-up photograph of a person's hands writing on a white notepad with a yellow and black pen. The person is wearing a light blue shirt and a black wristband. A silver laptop is visible in the background on a dark desk. A semi-transparent black banner is overlaid across the middle of the image, containing the text "Code Summarization" in white.

# Code Summarization

```

{
protected int code;
protected String reason;

/** Construct StatusLine */
public StatusLine(int c, String r)
{ code=c;
  reason=r;
}

/** Create a new copy of the request-line*/
public Object clone()
{ }

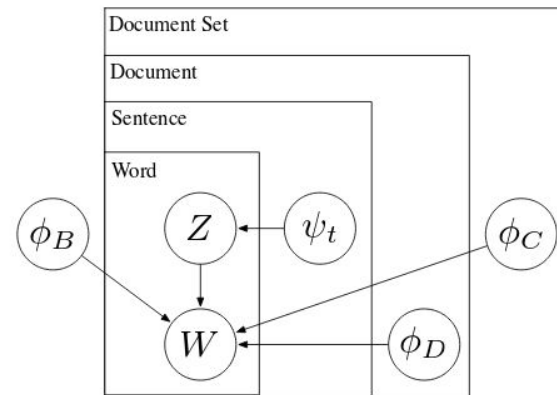
/** Indicates whether some other Object is "equal to" this StatusLine
public boolean equals(Object obj)
{ }

public String toString()
{ }

public int getCode()
{ }

public String getReason()
{ }
}

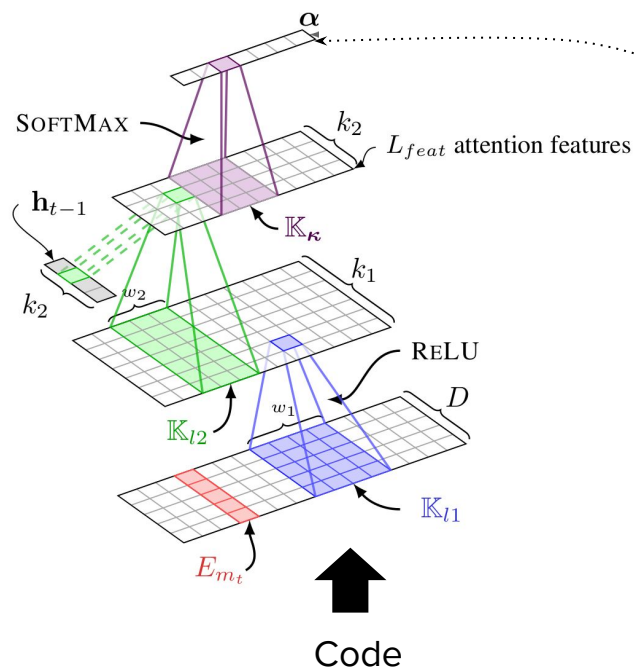
```



Haghighi & Vanderwende, 2009



# Summarization via Method Name Prediction



**conv\_attention** (code  $c$ , previous state  $\mathbf{h}_{t-1}$ )  
 $\alpha \leftarrow \text{attention\_weights} (L_{feat}, \mathbb{K}_{att})$



Target		Attention Vectors		$\lambda$
$m_1$	set	$\alpha =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	0.012
		$\kappa =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	
$m_2$	use	$\alpha =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	0.974
		$\kappa =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	
$m_3$	browser	$\alpha =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	0.969
		$\kappa =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	
$m_4$	cache	$\alpha =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	0.583
		$\kappa =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	
$m_5$	END	$\alpha =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	0.066
		$\kappa =$	<code>&lt;s&gt;{ this . use <u>Browser</u> Cache = use <u>Browser</u> Cache ; } &lt;/s&gt;</code>	

## Attention Visualization

Allamanis *et al*, “A Convolutional Attention Network for Extreme Summarization of Source Code”, ICML 2016



# Code Summarization to Natural Language

## 1. Source Code (C#):

```
public int TextWidth(string text) {  
    TextBlock t = new TextBlock();  
    t.Text = text;  
    return  
        (int)Math.Ceiling(t.ActualWidth);  
}
```

### Descriptions:

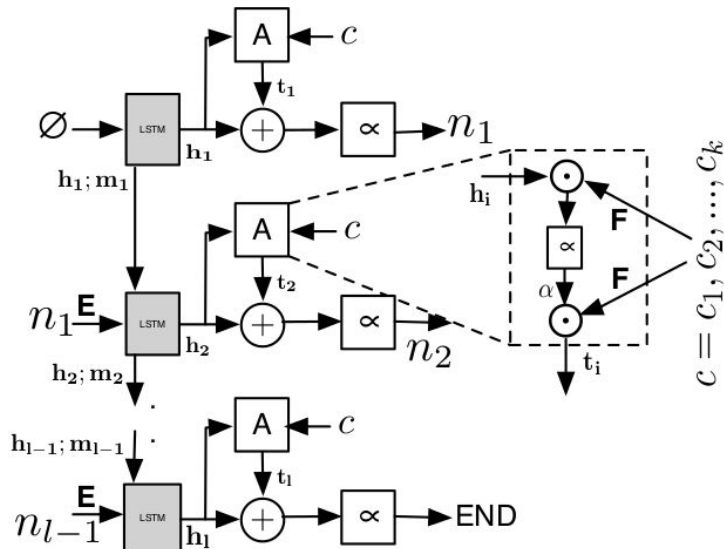
- Get rendered width of string rounded up to the nearest integer
- Compute the actual textwidth inside a textblock

## 2. Source Code (C#):

```
var input = "Hello";  
var regex = new Regex("World");  
return !regex.IsMatch(input);
```

### Descriptions:

- Return if the input doesn't contain a particular word in it
- Lookup a substring in a string using regex





# Coding Conventions

# Naming Conventions



```
public void testRunReturnsResult() {
    PrintStream oldOut = System.out;
    System.setOut(new PrintStream(
        new OutputStream() {
            @Override
            public void write(int arg0)
            }
        ));
    try {
        TestResult result = junit.textui.
        assertTrue(result.wasSuccessful());
    }
}
```

Code for Review



Proposers  
(rename identifiers,  
add formatting)



```
public void testRunReturnsResult() {
    PrintStream oldOut = System.out;
    System.setOut(new PrintStream(
        new OutputStream() {
            @Override
            public void write(int arg0)
            }
        ));
    try {
        TestResult result = junit.textui.
        assertTrue(result.wasSuccessful());
    }
}
```

Candidates



Scoring Function  
(ngram language model, SVM)



```
public void testRunReturnsResult() {
    PrintStream oldOut = System.out;
    System.setOut(new PrintStream(
        new OutputStream() {
            @Override
            public void write(int arg0)
            }
        ));
    try {
        TestResult result = junit.textui.
        assertTrue(result.wasSuccessful());
    }
}
```

Top Suggestions

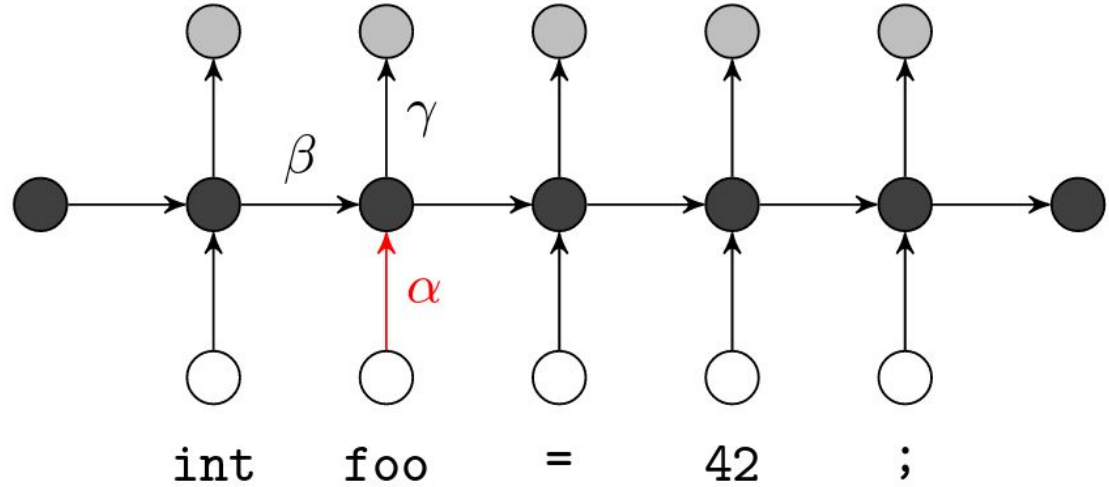
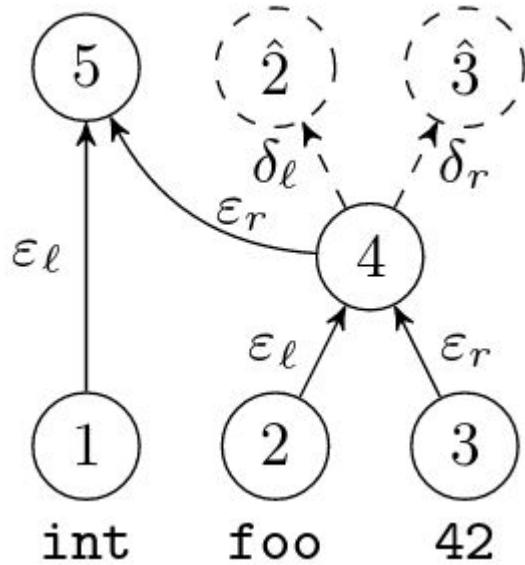
Training Corpus  
(rest of project code)



```
ForkJoinTask<?> ■■■;
if (task instanceof ForkJoinTask<?>) // avoid re-wrap
    ■■■ = (ForkJoinTask<?>) task;
else
    ■■■ = new
    ForkJoinTask.AdaptedRunnableAction(task);
externalPush(■■■);
```

- 1. job (30%)
- 2. task (20%)
- 3. tsk (15%)







```
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;
}
```

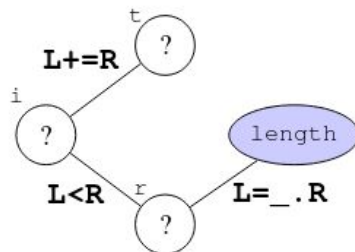
(a) JavaScript program with minified identifier names



Unknown properties (variable names):



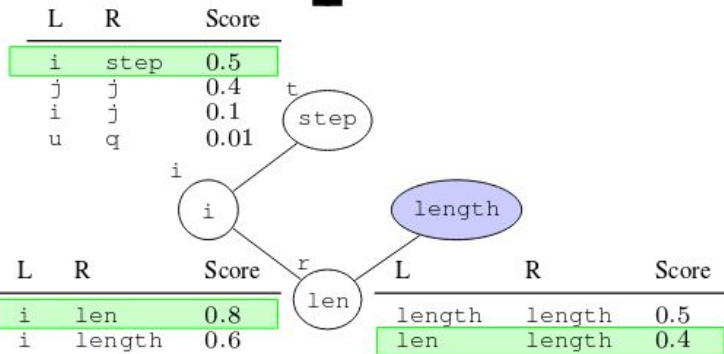
Known properties (constants, APIs):



(c) Dependency network

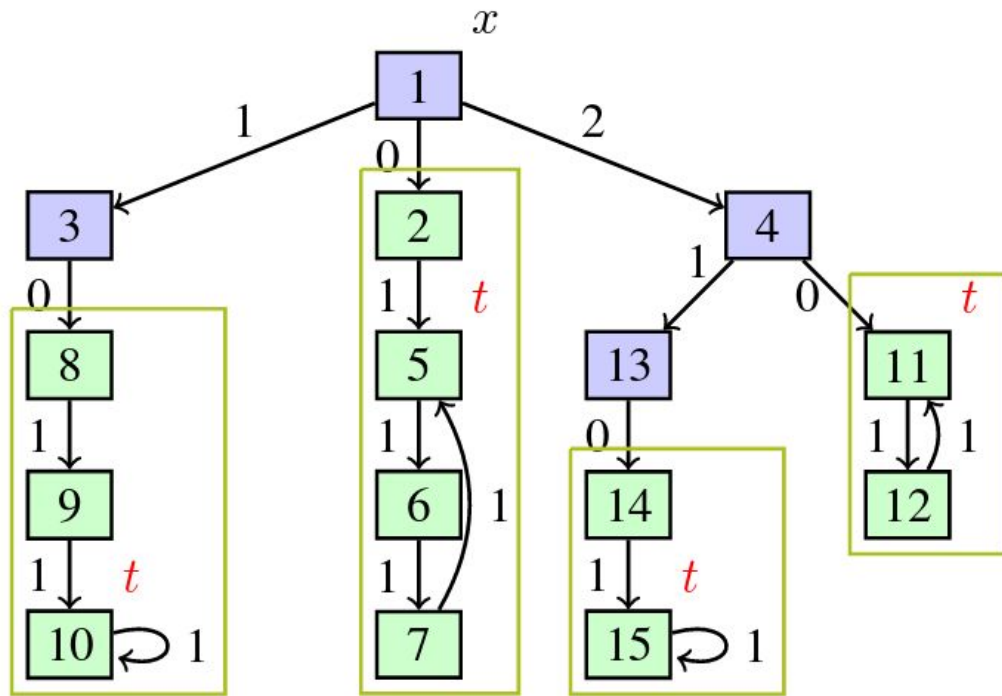
```
/* str: string, step: number, return: Array */
function chunkData(str, step) {
  var colNames = []; /* colNames: Array */
  var len = str.length;
  var i = 0; /* i: number */
  for (; i < len; i += step) {
    if (i + step < len) {
      colNames.push(str.substring(i, i + step));
    } else {
      colNames.push(str.substring(i, len));
    }
  }
  return colNames;
}
```

(e) JavaScript program with new identifier names and types



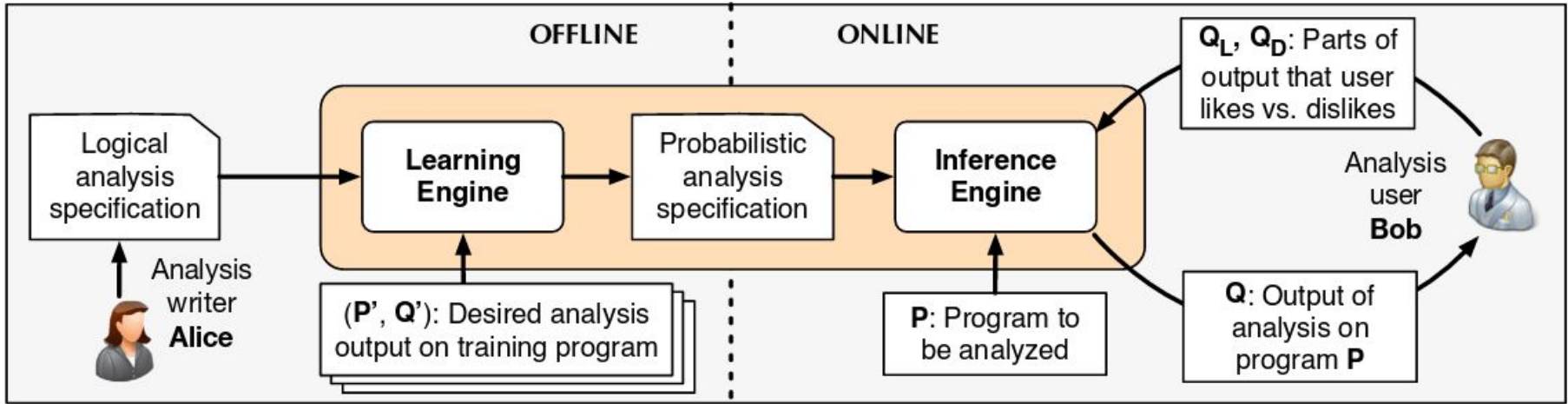
(d) Result of MAP inference





$$\psi = \text{tree}(x, \lambda i_1, i_2, i_3, i_4 \rightarrow \exists t. \text{ls}(i_2, t, \lambda i_5, i_6, i_7, i_8 \rightarrow \top) * \text{ls}(t, t, \lambda i_9, i_{10}, i_{11}, i_{12} \rightarrow \top)).$$







## **Probabilistic Models of Source Code**

# What was *not* discussed

- Non-*Probabilistic* Models
- Models with *no* Learning Component
- Pure Information Retrieval Models
- Off-the-shelf models: topic models, *etc.*  
with hand-crafted features
- Bag-of-Word Representations
- Things that don't Involve Code



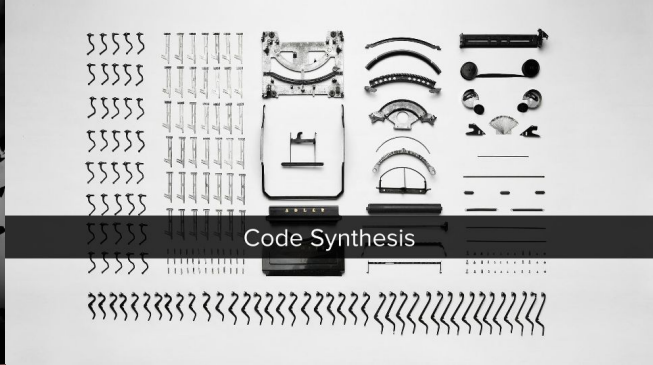
**Challenges Ahead...**



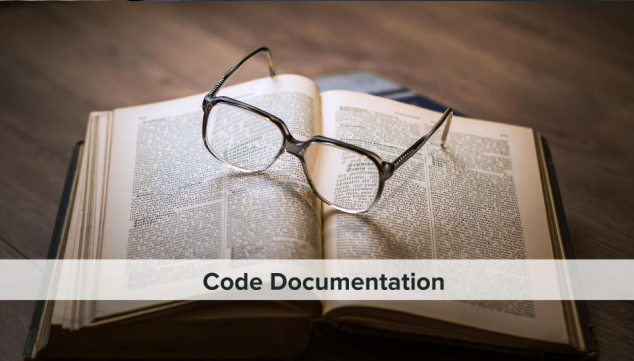
Recommender Systems



Statistical Code Migration



Code Synthesis



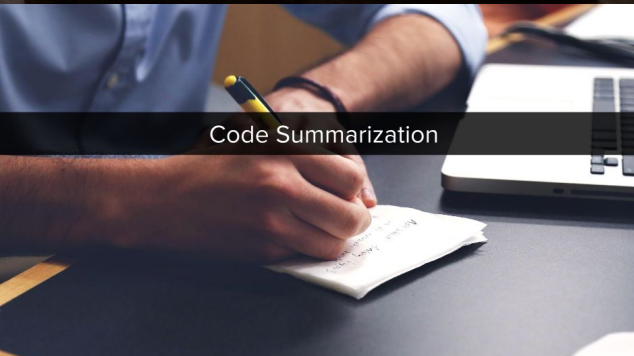
Code Documentation



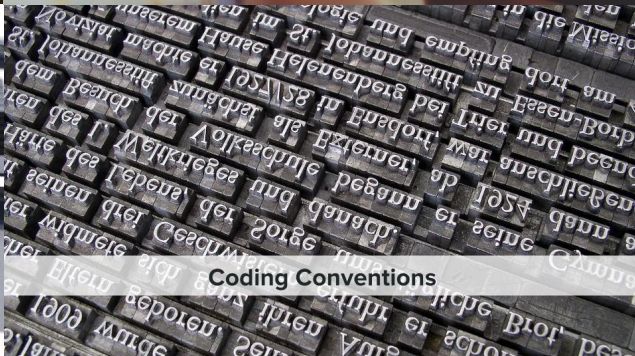
A Taxonomy of Models



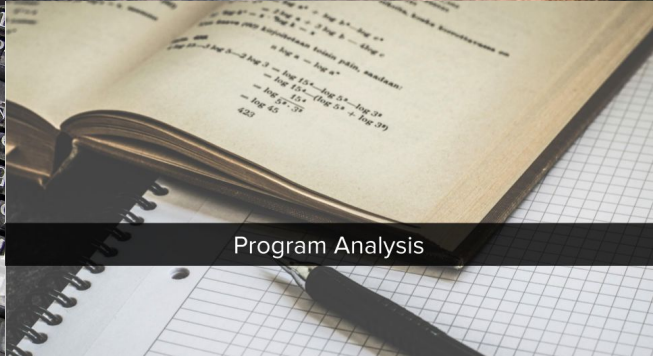
Code Defects



Code Summarization



Coding Conventions



Program Analysis

```
previousElements = [];  
selectedScopes = [];  
  
scope.$watch(watchExpr, function ngSwitchWatchAction(value) {  
  var i, ii;  
  for (i = 0, ii = previousElements.length; i < ii; ++i) {  
    previousElements[i].remove();  
  }  
  previousElements.length = 0;  
  
  for (i = 0, ii = selectedScopes.length; i < ii; ++i) {  
    var selected = selectedElements[i];  
    selectedScopes[i].$destroy();  
    previousElements[i] (the end);  
    $animate.leave(selected, function() {  
      previousElements.splice(i, 1);  
    });  
  }  
});
```

```
selectedElements.length = 0;  
selectedScopes.length = 0;
```

```
if ((selectedTranscludes = ngSwitchController.cases['!' + value] || ngSwitchController.defaultCase)) {  
  scope.$eval(attr.change);  
  forEach(selectedTranscludes, function(selectedTransclude) {  
    var selectedScope = scope.$new();  
    selectedScopes.push(selectedScope);  
    selectedTransclude.$scope = selectedScope;  
    selectedTransclude.$element = selectedTransclude.element.clone(true);  
    selectedTransclude.$element.appendTo(selectedTransclude.container);  
    selectedScope.$parent = scope;  
    selectedScope.$root = scope.$root;  
    selectedScope.$on('$destroy', function() {  
      selectedTransclude.$element.remove();  
    });  
  });  
}
```