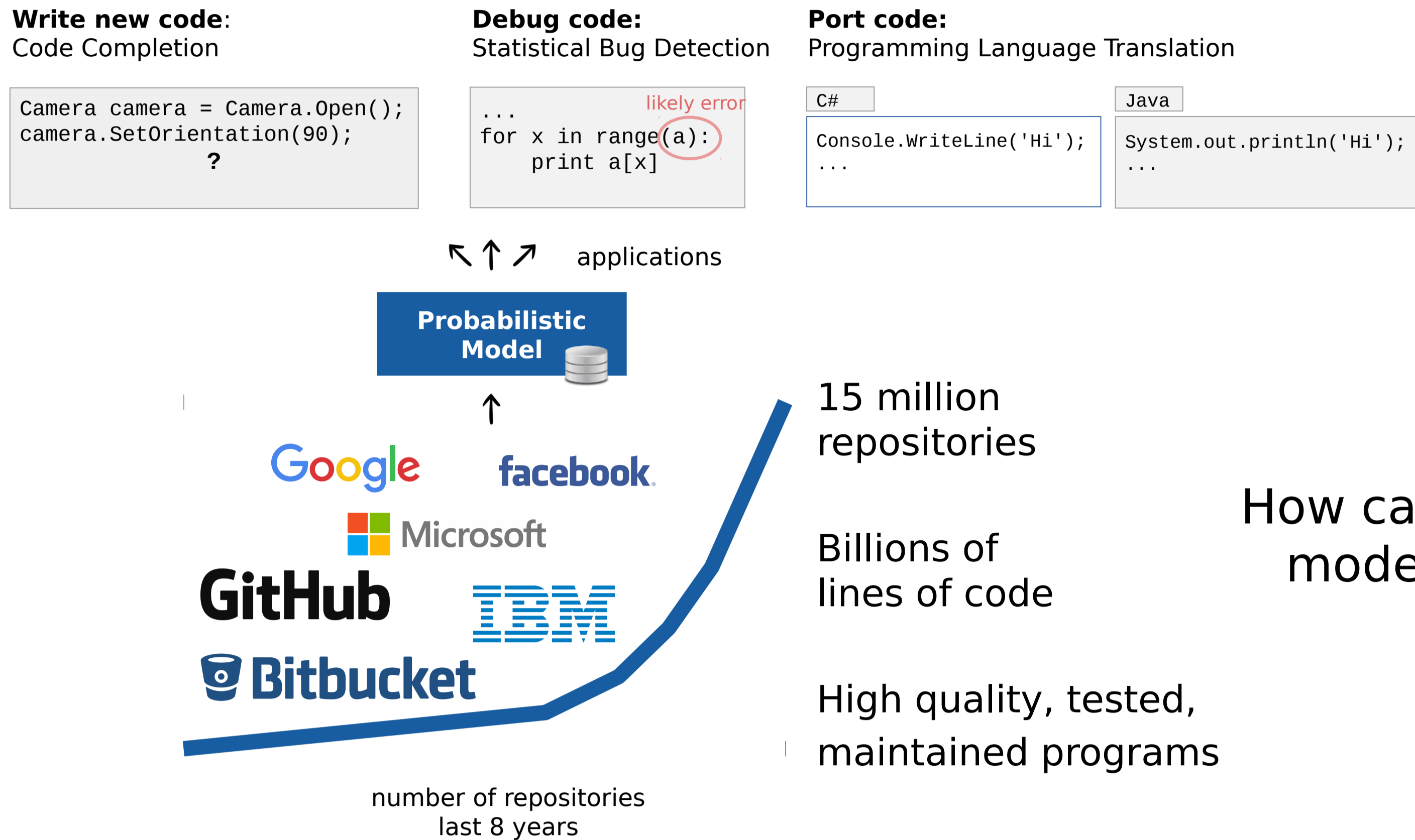


# Probabilistic Higher Order Grammar: Probabilistic Model for Code

Pavol Bielik, Veselin Raychev, Martin Vechev



**Efficient Learning**  
Trained as efficiently as PCFGs and n-gram models

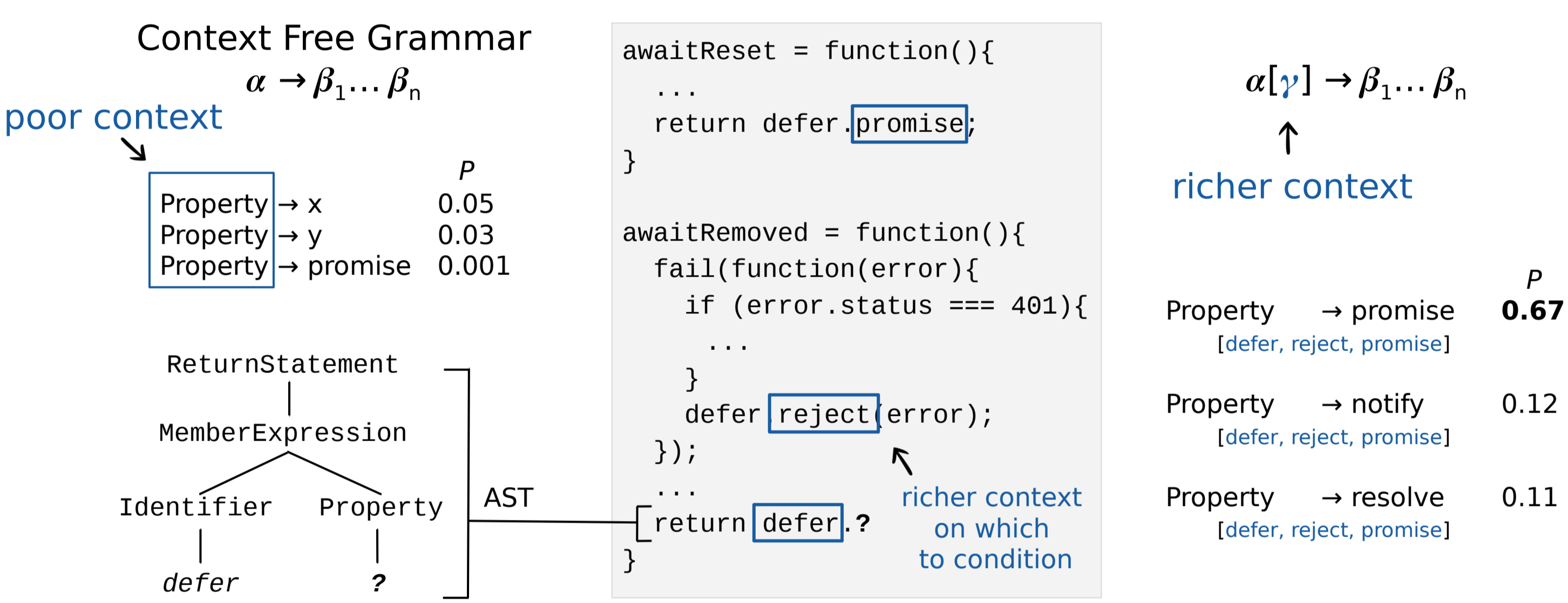
**Widely Applicable**  
Agnostic to programming language

How can we learn probabilistic models directly from data?

**Flexible Representation**  
Conditioning for the predictions is determined dynamically

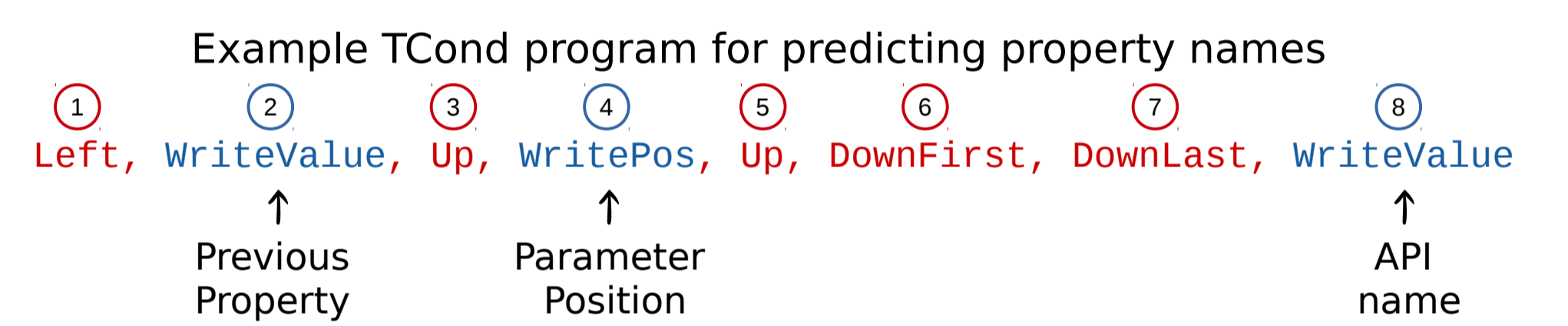
## The Need for Better Conditioning

## Using Programs to Explain Data



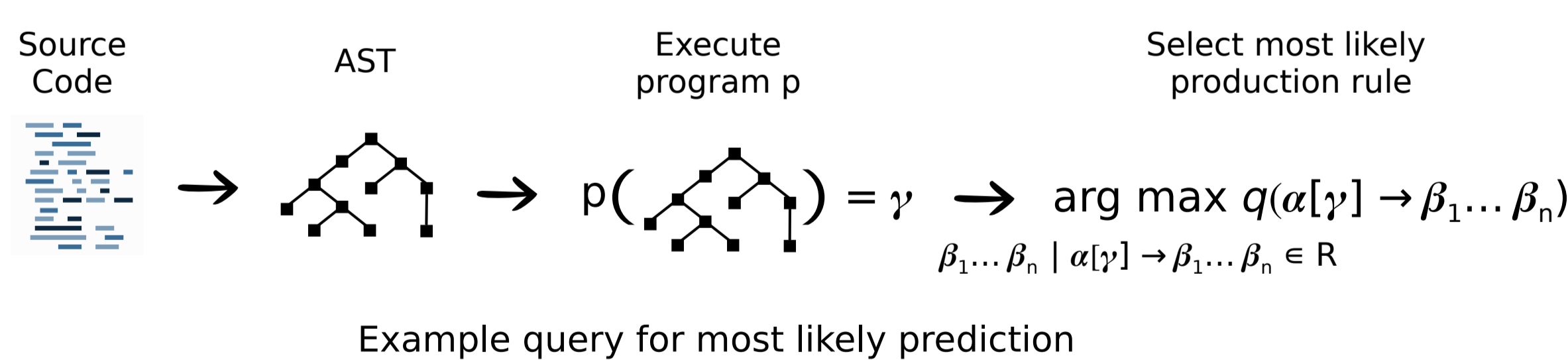
In general: Unrestricted programs (Turing complete) vs. Our Work: TCond Language for navigating over trees and accumulating context

TCond ::=  $\epsilon$  | WriteOp TCond | MoveOp TCond  
 MoveOp ::= Up, Left, Right, DownFirst, DownLast, NextDFS, PrevDFS, NextLeaf, PrevLeaf, PrevNodeType, PrevNodeValue, PrevNodeContext  
 WriteOp ::= WriteValue, WriteType, WritePos

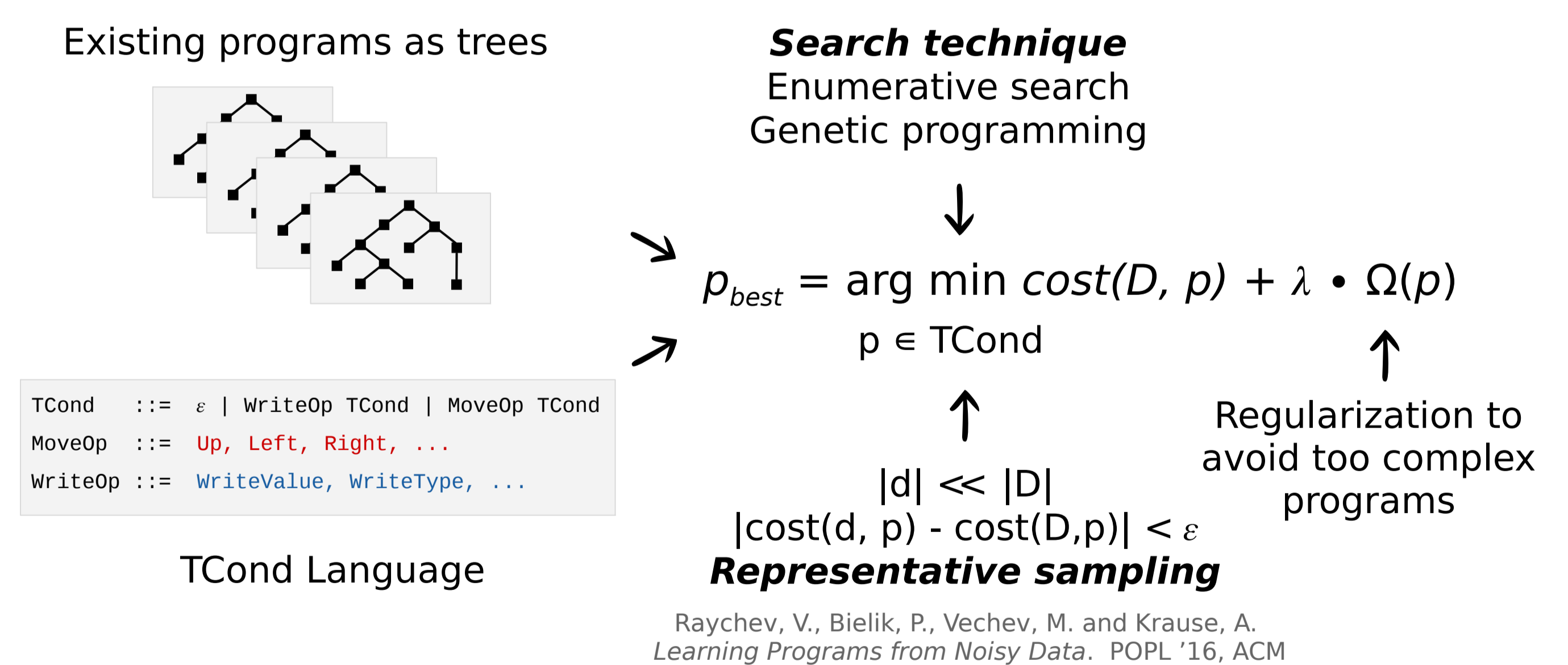


## Probabilistic Higher Order Grammar (PHOG)

N: set of non-terminal symbols  
 $\Sigma$ : set of terminal symbols  
 s: start symbol  
 Condition the predictions on richer context  $\rightarrow$  R: set of rules in form:  $\alpha[\gamma] \rightarrow \beta_1 \dots \beta_n$   
 program that dynamically obtains the context for given query  $\rightarrow$  p: AST  $\rightarrow$  C  
 q: R  $\rightarrow \mathbb{R}^+$   $\leftarrow$  valid probability distribution  
 $\sum q(\alpha[\gamma] \rightarrow \beta_1 \dots \beta_n) = 1$   
 $\alpha[\gamma] \rightarrow \beta_1 \dots \beta_n \in R$



## Learning



## Evaluation

**GitHub** 150k JavaScript Programs  $\rightarrow$  100k: Training Set (1.07  $\cdot$  10<sup>8</sup> completion queries)  
<http://www.srl.inf.ethz.ch/js150.php> 50k: Evaluation Set (5.3  $\cdot$  10<sup>7</sup> completion queries)

	Code Completion Error Rate		Training Time	Queries per Second
	Non-Terminals	Terminals		
PCFG	48.5%	49.9%	1 min	71 000
n-gram	30.8%	28.7%	4 min	15 000
Naive Bayes	41.6%	45.8%	3 min	10 000
SVM	32.5%	29.5%	36 hours	12 500
<b>PHOG</b>	<b>25.9%</b>	<b>18.5%</b>	<b>162 + 3 min</b>	<b>50 000</b>

Code Completion Examples		
Completion Kind	Error Rate	Completion Example
Identifier	38%	contains = <b>jQuery</b> ...
Property	35%	start = list. <b>length</b> ;
String	48%	'[ ' + attrs + <b>' ]</b>
Number	36%	canvas(xy[0], xy[ <b>1</b> ], ...)
RegExp	34%	line.replace(/( <b>&amp;nbsp;</b> ; <b> </b> )+/, ...)
UnaryExpr	3%	if (!events    <b>!</b> ...)
BinaryExpr	26%	while (++index <b><math>\leq</math></b> ...)
LogicalExpr	8%	frame = frame    <b>!</b> ...