



软件分析

# 程序生成学习

熊英飞  
北京大学



# This Lecture

## Classic Synthesis

- Problem Definition
- Enumerative
- Constraint-based
- Presentation-based

## Program Estimation

- Problem Definition
- Estimating Probabilities
- Locating the most-likely one



# 程序估计 Program Estimation

- 输入:
  - 程序空间（用语法表示）  $G$
  - 规约  $S$
  - 上下文  $C$
  - 包含上下文-程序对的训练集  $T$
- 输出:
  - 程序  $P$  满足
    - $P \in G \wedge P \mapsto S \wedge \text{Pr}(P | C)$
  - $\text{Pr}$  表示从  $T$  学习到的概率
- 如果  $\text{Pr}$  是满足规约的概率，那么可以用来加速传统程序分析



# Example: Condition Completion

- Given a program without a conditional expression, completing the condition

```
public static long fibonacci(int n) {  
    if ( ?? ) return n;  
    else return fibonacci(n-1) + fibonacci(n-2);  
}
```

```
E → E ">12"  
    | E ">0"  
    | E "+" E  
    | "hours"  
    | "value"  
    | ...
```

Space of Conditions  
defined by a grammar

- Specification is a set of tests
- Useful in program repair
  - Many bugs are caused by incorrect conditions
  - Existing work could localize the faulty condition
  - Can we generate a correct condition to replace the incorrect one?



# 分解为三个问题

- 如何根据训练集计算一个程序的概率？
- 如何找到概率最大的程序？
- 如何保证找到的程序满足规约的要求？

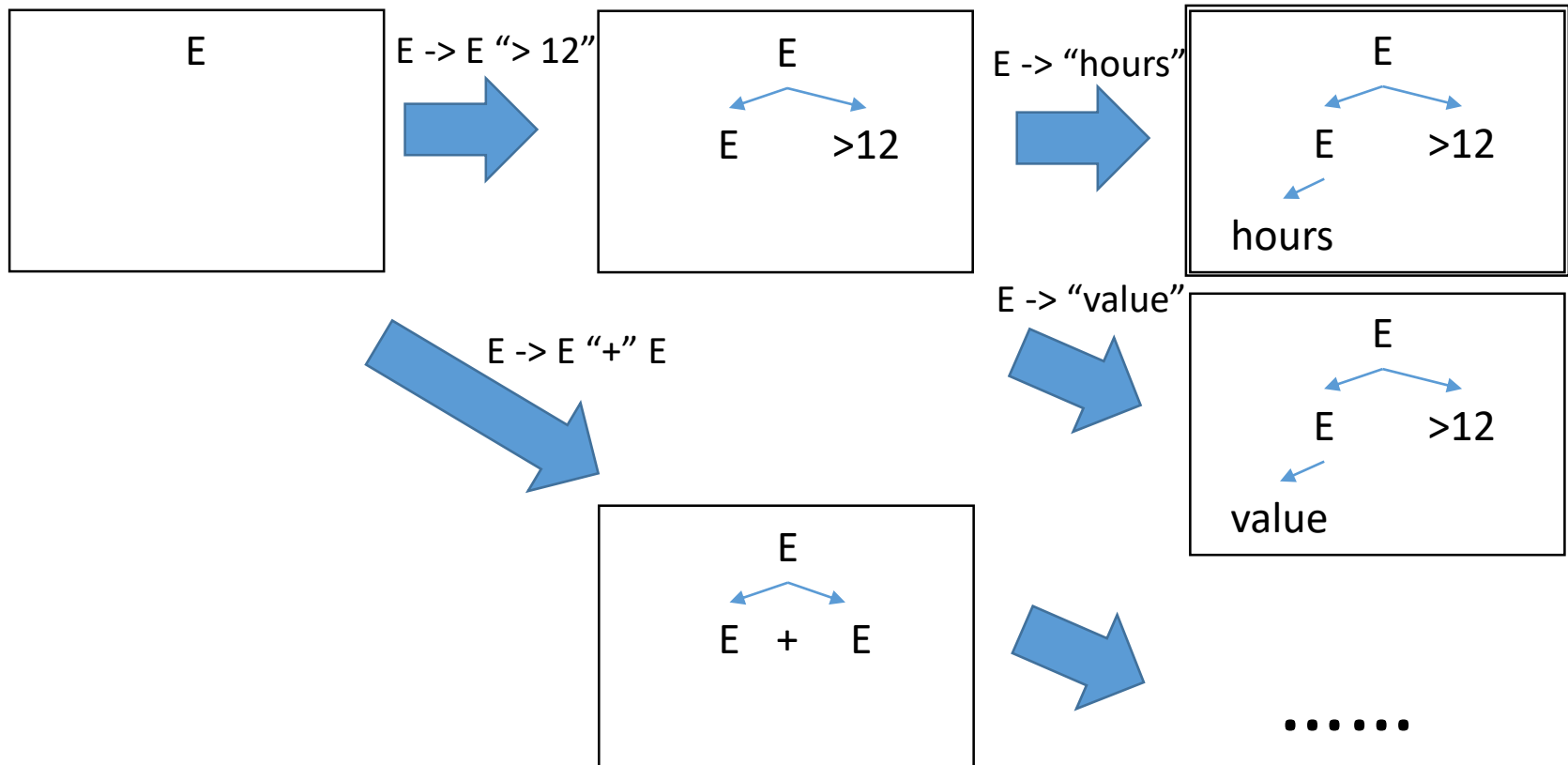


# 程序生成学习框架

## Learning to Synthesize (L2S)

- 本课题组提出的框架
- 泛化了之前的多种具体方法
- 将程序估计问题看做一个图的路径查找问题
  - 起始节点是空程序
  - 目标节点是满足规约的程序
  - 路径的权是终点程序的概率
  - 目标是找到一条从起始节点到任意目标节点的最优路径

# 程序估计问题作为路径查找问题

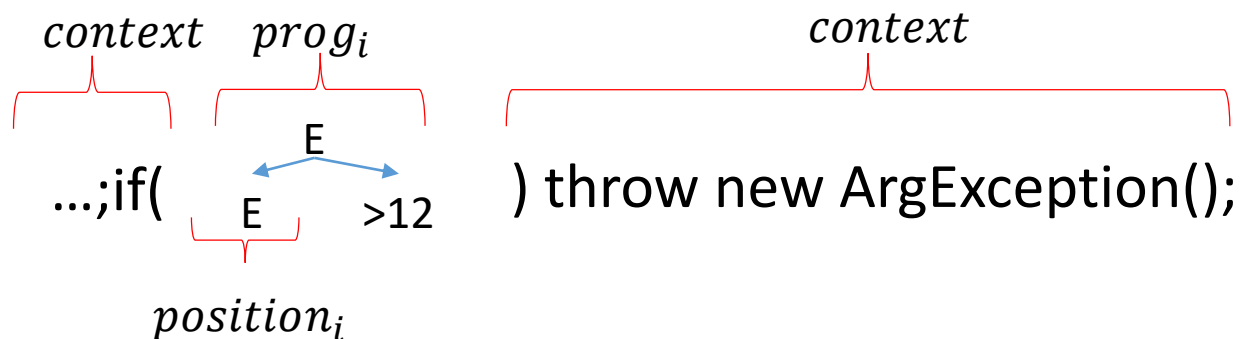




# 如何计算一个程序的概率？

$$\bullet P(\text{prog} \mid \text{context}) = \prod_i P(\text{rule}_i \mid \text{context}, \text{prog}_i, \text{position}_i)$$

- *context*: The context of the program
- *prog<sub>i</sub>*: The AST generated at the *i*th step
- *position<sub>i</sub>*: The non-terminal to be expanded at the *i*th step
- *rule*: the chosen rule at the *i*th step
- *prog*: the complete program







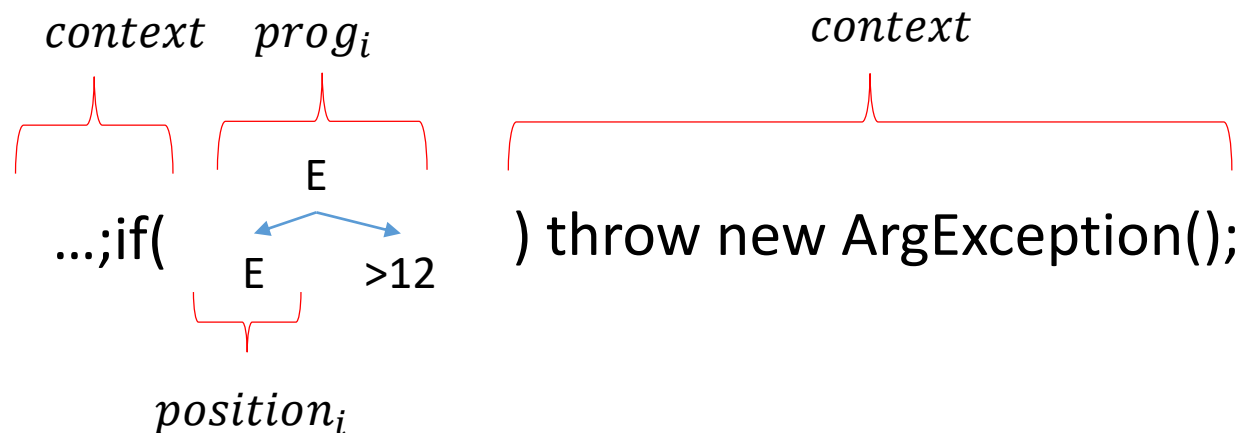
# Training models

- Train a model for each non-terminal
  - to classify rules expanding this non-terminal
- Training set preparation
  - The original training set:
    - A set of programs
    - Their contexts
  - Decomposing the training set:
    - Parse the programs
    - Extract the rules chosen for each non-terminal



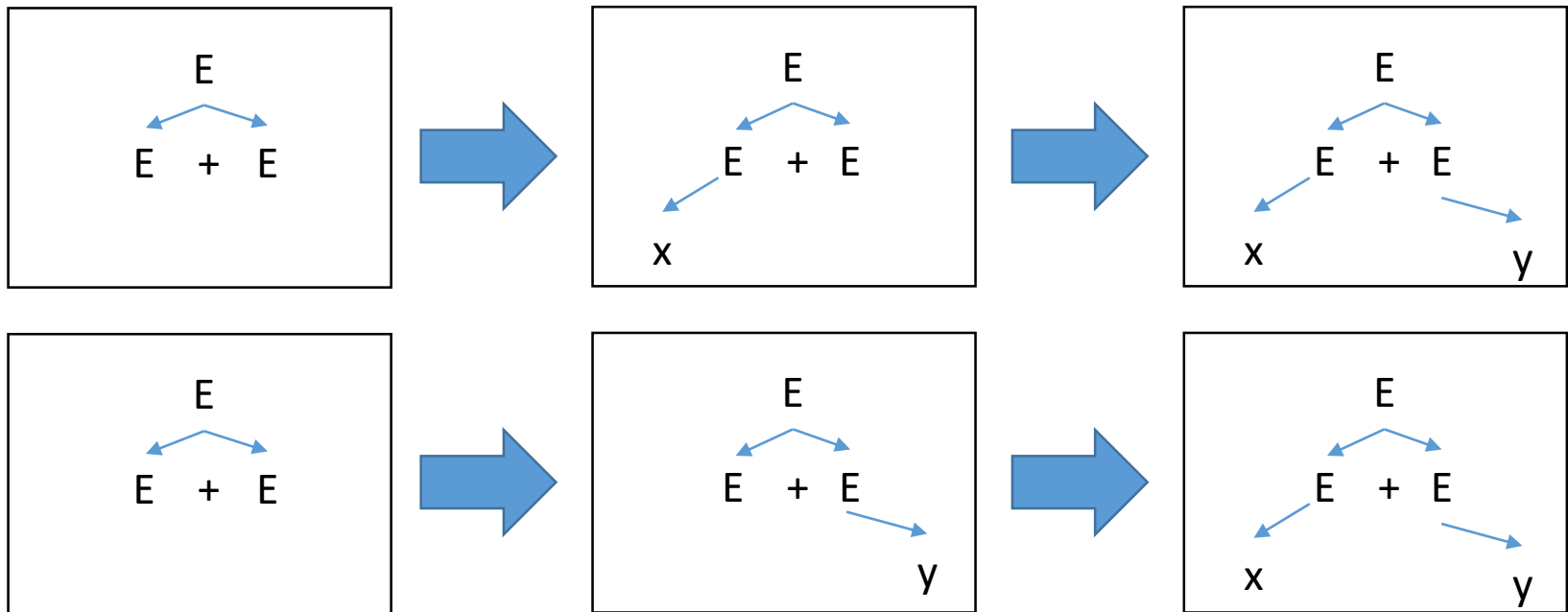
# Feature Engineering

- Extract features from
  - *context* : The context
  - *prog<sub>i</sub>* : The generated partial AST
  - *position<sub>i</sub>* : The position of the node to be expanded





# The order of expansion



- Does different order make a difference?



# The order of expansion

- 如果所有概率都是精确的，两者的结果没有差别
  - $P(\text{prog} \mid \text{context}) = \prod_i P(\text{rule}_i \mid \text{context}, \text{prog}_i, \text{position}_i)$
- 证明：假设存在一个 **policy**，决定一个不完整程序中哪个节点先被展开，那么 **policy** 的选择和 **prog** 的概率是独立的
  - $\Pr(\text{prog})$
  - $= \Pr(\text{prog} \mid \text{policy})$  // 独立性
  - $= \Pr(\langle \text{prog}_i, \text{pos}_i, \text{rule}_i \rangle)_{i=1}^n \mid \text{policy})$
  - $= \Pr(\text{prog}_1 \mid \text{policy}) \Pr(\text{pos}_1 \mid \text{policy}, \text{prog}_1)$   
 $\Pr(\text{rule}_1 \mid \text{policy}, \text{prog}_1, \text{pos}_1)$   
 $\Pr(\text{eprog}_2 \mid \text{policy}, \text{prog}_1, \text{pos}_1, \text{rule}_1) \dots$   
 $\Pr(\text{eprog}_{n+1} \mid \text{policy}, (\text{eprog}_i)_{i=1}^n, (\text{pos}_i)_{i=1}^n, (\text{rule}_i)_{i=1}^n)$
  - $= \prod_i \Pr(\text{rule}_i \mid \text{policy}, (\text{rule}_j)_{j=1}^{i-1}, \text{pos}_i)$  // 删除概率为1的项
  - $= \prod_i \Pr(\text{rule}_i \mid \text{policy}, \text{prog}_i, \text{pos}_i)$
  - $= \prod_i \Pr(\text{rule}_i \mid \text{prog}_i, \text{pos}_i)$  // 独立性



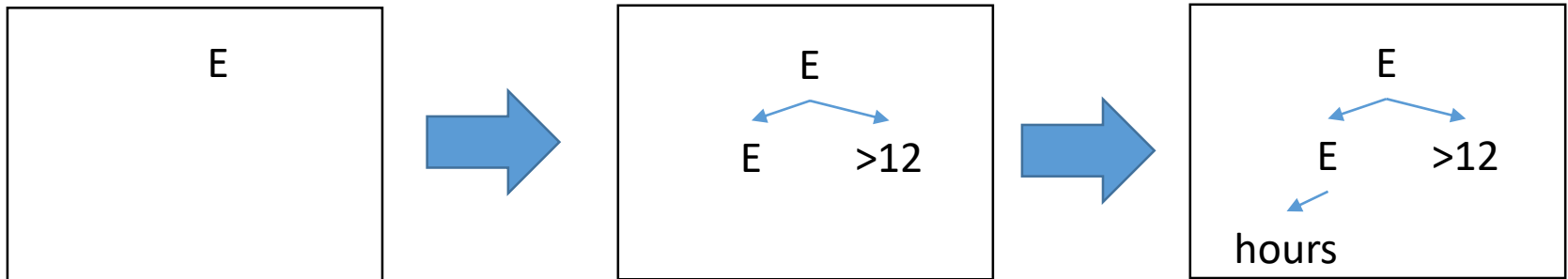
# The order of expansion

- 虽然精确的概率乘积是相同的，但
  - 统计模型/机器学习的预测精度可能不同
  - 对路径查找问题求解算法的影响不同
- 根据经验，采用不同的顺序可能对结果产生很大影响

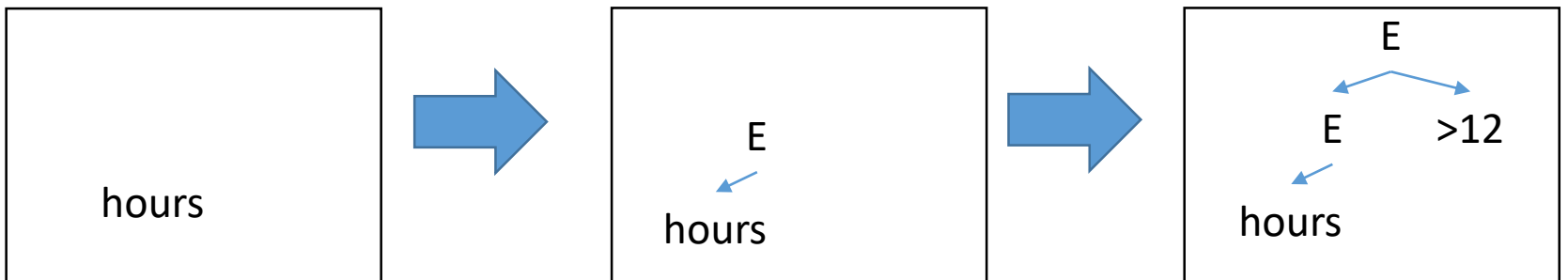


# Order beyond CFG?

- Top-down



- Bottom-up





# CFG to Expansion Rules

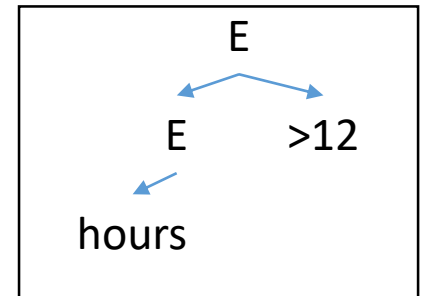
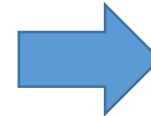
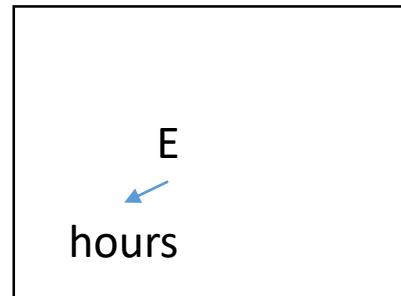
$T \rightarrow E$   
 $E \rightarrow E \text{ " > 12" } \mid E \text{ " > 0" } \mid E \text{ " + " } E \mid \text{ "hours" } \mid \text{ "value" } \mid \dots$



- $\langle E \rightarrow \text{"hours"}, \perp \rangle$
- $\langle E \rightarrow \text{"value"}, \perp \rangle$
- $\langle E \rightarrow E \text{ " > 12"}, 1 \rangle$
- $\langle E \rightarrow E \text{ " + " } E, 1 \rangle$
- $\langle T \rightarrow E, 1 \rangle$
- $\langle E \rightarrow E \text{ " > 12"}, 0 \rangle$
- $\langle E \rightarrow E \text{ " + " } E, 0 \rangle$
- $\langle E \rightarrow \text{"hours"}, 0 \rangle$
- $\langle E \rightarrow \text{"value"}, 0 \rangle$

自底向上规则:  $\langle E \rightarrow E \text{ " > 12"}, 1 \rangle$

Generate the tree if the ith child is ready



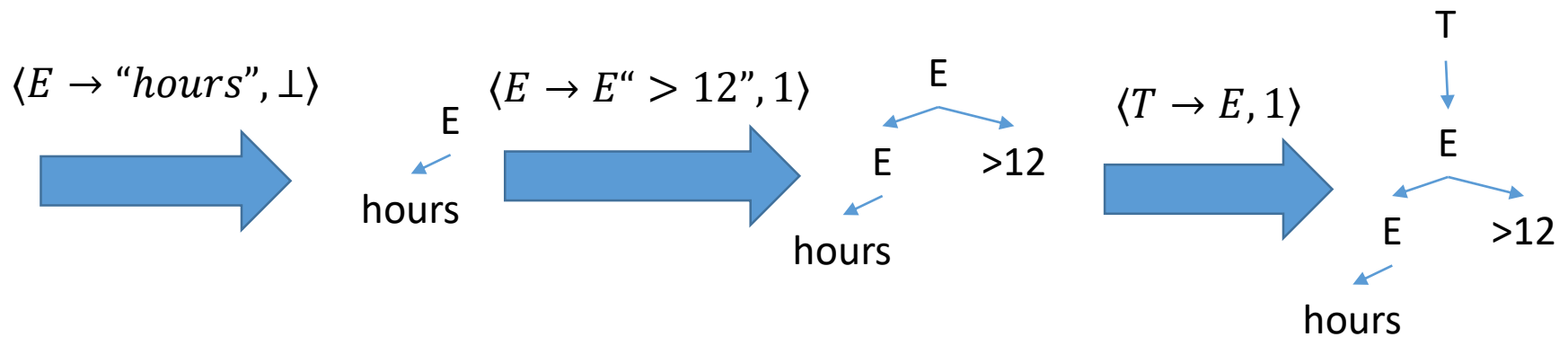
自顶向下规则:  $\langle E \rightarrow E \text{ " > 12"}, 0 \rangle$

Generate the tree is the root is generated

创建规则:  $\langle E \rightarrow \text{"hours"}, \perp \rangle$

Generate the tree at the beginning

# Generation by Expansion Rules







# Expansion Trees

- Expansion trees capture how expansion rules are applied

hours>12	hours+value
$(T \rightarrow E, 1)$ ↑ $(E \rightarrow E \text{ " > 12" }, 1)$ ↑ $(E \rightarrow \text{"hours"}, \perp)$	$(T \rightarrow E, 1)$ ↑ $(E \rightarrow E \text{ "+" } E, 1)$ ↙ ↘ $(E \rightarrow \text{"hours"}, \perp)$ $(E \rightarrow \text{"value"}, 0)$



# AST $\rightarrow$ Expansion Tree

- **完整性Completeness**: for each AST, there is at least one expansion tree
- **唯一性Uniqueness**: for each AST, there is at most one expansion tree
- 是否总是存在完整和唯一的Expansion Rule集合?



# 唯一和完整集合的充分条件

$$T \rightarrow E$$
$$E \rightarrow E \text{ " > 12" } \mid E \text{ " > 0" } \mid E \text{ " + " } E \mid \text{"hours"} \mid \text{"value"} \mid \dots$$


$\langle E \rightarrow \text{"hours"}, \perp \rangle$
$\langle E \rightarrow \text{"value"}, \perp \rangle$
$\langle E \rightarrow E \text{ " > 12"}, 1 \rangle$
$\langle E \rightarrow E \text{ " + " } E, 1 \rangle$
$\langle T \rightarrow E, 1 \rangle$
$\langle E \rightarrow E \text{ " > 12"}, 0 \rangle$
$\langle E \rightarrow E \text{ " + " } E, 0 \rangle$
$\langle E \rightarrow \text{"hours"}, 0 \rangle$
$\langle E \rightarrow \text{"value"}, 0 \rangle$

1. 除了初始符号开头的规则，所有语法规则都有对应的自顶向下展开规则
2. 所有语法规则最多只有一条自底向上的展开规则
3. 对于所有从初始符号（延自底向上展开规则）反向可达的非终结符，其所有语法规则都有一条自底向上展开规则或创建规则

从初始符号开始选择创建/自底向上规则即可



# AST $\rightarrow$ Expansion Tree

- 利用一个动态规划算法，AST可以在 $O(n)$ 时间内转成Expansion Tree
  - 后根次序依次判断每个AST结点是否可以被自底向上和自顶向下的方式生成，如果可以，记录下采用的规则
  - 先根次序恢复出Expansion Tree



# 如何找到概率最大的程序？

- 采用求解路径查找问题的标准算法
- 精确算法
  - 迪杰斯特拉算法
  - A\*算法
- 近似算法
  - Beam Search

# 如何保证找到的程序满足规约的要求？



- 基本方法：
  - 寻找概率最大的程序
  - 判断是否满足规约
  - 如不满足，回到第一步
- 能否在搜索过程中就进行剪枝？



# 剪枝

- 搜索过程中剪枝
  - 语义：假设输入变量的取值仅为2，要求输出为3，且语法中只有加号，那么 $E+E$ 肯定无法满足
  - 类型： $E+E \ \&\& \ E$ 肯定无法满足
  - 大小：假设AST树的大小（节点数）限定为4，那么 $E+E$ 肯定无法满足
- 剪枝的条件
  - 所有可展开的程序都无法满足约束



# 语法上的静态分析

- 假设所有约束都是  $\text{Pred}(\text{Prop}(N))$  的形式
  - $N$ : 非终结符
  - $\text{Prop}$ : 以  $N$  为根节点的子树所具有的属性值
  - $\text{Pred}$ : 该属性值所应该满足的谓词
- 如:
  - 语义约束:  $\text{Prop}$  为表达式取值
  - 类型约束:  $\text{Prop}$  为表达式的可能类型
  - 大小约束:  $\text{Prop}$  为表达式的大小
- 通过静态分析获得  $\text{Prop}$  的所有可能取值
  - 要求上近似
- 如果所有可能取值都不能满足  $\text{Pred}$ , 则该部分程序可以减掉





# 语法上静态分析示例：语义

- 抽象域：由1, 2, 3, 4, 5, >5, <1, true, false构成的集合
- 容易定义出抽象域上的计算
- 从语法规则产生方程
- $E \rightarrow E + E \mid \text{"x"} \mid E \text{">5"} \mid \dots$ 
  - $V[E] = (V[E] + V[E]) \cup \{2\} \cup (V[E] > 5) \cup \dots$
- 求解方程得到每一个非终结符可能的取值（在开始时做一次）
- 根据当前的部分程序产生计算式





# 语法上静态分析示例： 类型和大小

- 抽象域：由Int, Float, Boolean构成的集合

- 从语法规则产生方程

- $E \rightarrow E + E \mid \text{"hours"} \mid E > 5 \mid \dots$

- $T[E] = (T[E] + T[E]) \cup \{\text{Int}\} \cup (V[E] > 5) \cup \dots$

- 其中

- $t_1 + t_2 = \begin{cases} \{\text{Int}, \text{Float}\}, & \text{同时满足下面两个条件} \\ \{\text{Int}\}, & \text{Int} \in t_1 \wedge \text{Int} \in t_2 \\ \{\text{Float}\}, & \text{Float} \in t_1 \wedge \text{Float} \in t_2 \\ \emptyset, & \text{否则} \end{cases}$

- $t > 5 = \begin{cases} \{\text{Boolean}\}, & \text{Int} \in t \vee \text{Float} \in t \\ \emptyset, & \text{否则} \end{cases}$

- 用类似方法可以计算非终结符展开的最小大小



# 从AST到Expansion Tree

- 类似的规则可以对不完整Expansion Tree定义
- 需要计算出每个非终结符向上/向下展开的所有可能取值



# Summary

- L2S Combines four tools
  - **Expansion rules**: a novel extension to CFG for defining a path finding problem
  - **Static Analysis on Rules**: pruning off invalid choices in each step
  - **Statistical models**: estimating the probabilities of choices in each step
  - **Search algorithms**: solving the path finding problem



# Evaluation

- Evaluation 1:
  - Repairing Conditional Expressions
- Evaluation 2:
  - Generating Code from Natural Language Expression



# Repairing Conditional Expressions

- Condition bugs are common

```
hours = convert(value);  
+ if (hours > 12)  
+   throw new ArithmeticException();
```

Missing boundary checks

```
- if (hours >= 24)  
+ if (hours > 24)  
    withinOneDay=true;
```

Conditions too weak or too strong

- Steps:
  1. Localize a buggy if condition with SBFL and predicate switching
  2. Synthesize an if condition to replace the buggy one
  3. Validate the new program with tests



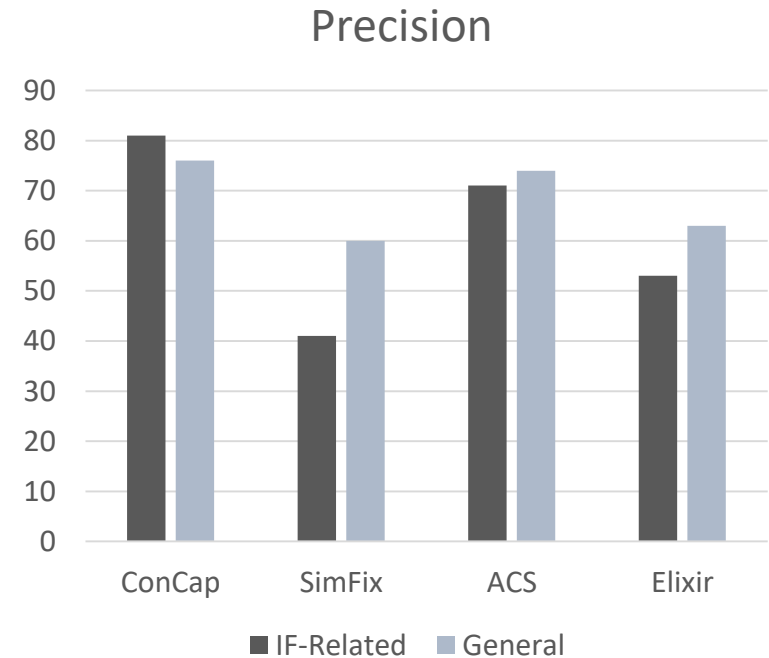
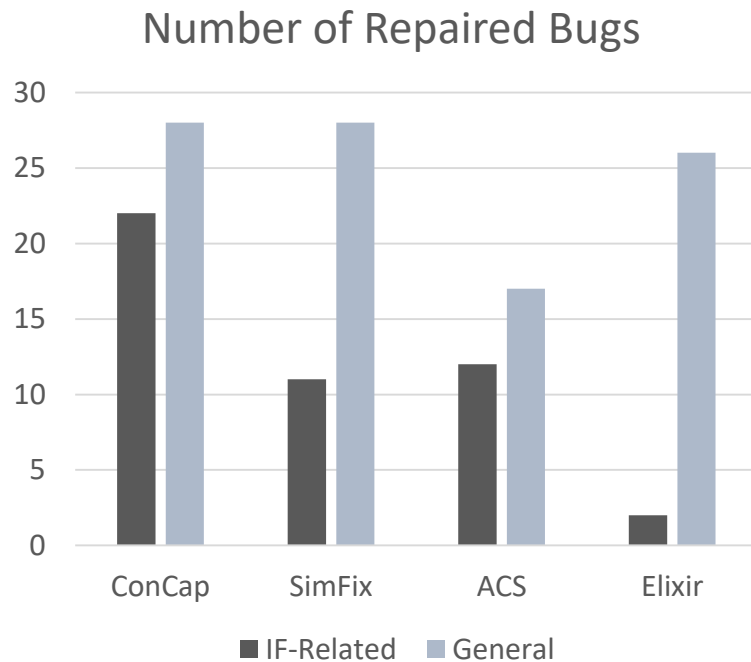
# L2S Configuration

- Expansion rules
  - Bottom-up
  - Estimate the leftmost variable first
- Machine learning
  - Xgboost
  - Manually designed features
- Constraints
  - Type constraints & size constraints
- Search algorithm
  - Beam search



# Results

Benchmark: Defects4J



Also repaired 8 unique bugs that have never been repaired by any approach.





# Generating Code from Natural Language Expression

- Can we generate code automatically to avoid repetitive coding?
- Existing approaches use RNN to translate natural language descriptions to programs
  - **Long dependency problem:** work poorly on long programs



```
[NAME]  
Acidic Swamp Ooze  
[ATK] 3  
[DEF] 2  
[COST] 2  
[DUR] -1  
[TYPE] Minion  
[CLASS] Neutral  
[RACE] NIL  
[RARITY] Common  
[DESCRIPTION]  
"Battlecry: Destroy Your Opponent's Weapon"
```



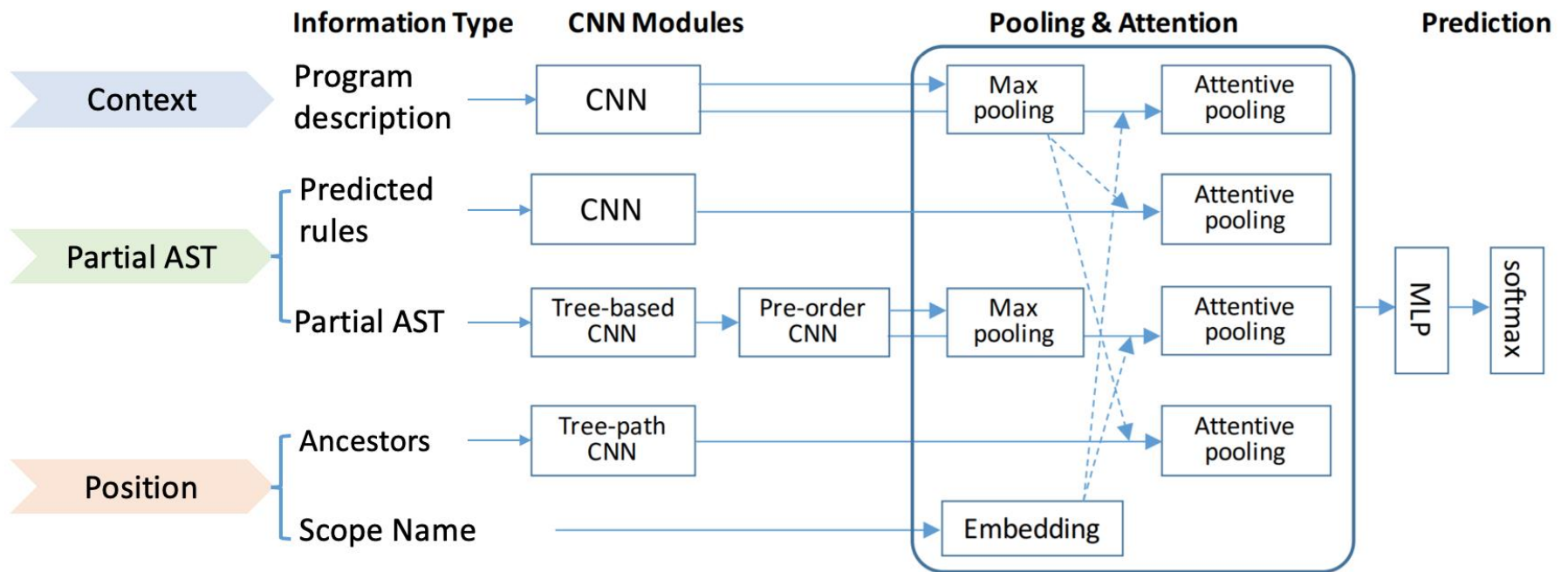
```
class AcidicSwampOoze(MinionCard):  
    def __init__(self):  
        super().__init__("Acidic Swamp Ooze", 2,  
            CHARACTER_CLASS.ALL, CARD_RARITY.COMMON,  
            battlecry=Battlecry(Destroy(), WeaponSelector(EnemyPlayer())))  
  
    def create_minion(self, player):  
        return Minion(3, 2)
```



# L2S Configuration

- Expansion rules
  - Top-down
- Machine learning
  - A CNN-based network
- Constraints
  - Size constraints
- Search algorithm
  - Beam search

# A CNN-based Network Architecture



# Results



## Benchmark: HearthStone

<b>Model</b>	<b>StrAcc</b>	<b>Acc+</b>	<b>BLEU</b>
LPN (Ling et al. 2016)	6.1	–	67.1
SEQ2TREE (Dong and Lapata 2016)	1.5	–	53.4
SNM (Yin and Neubig 2017)	16.2	~18.2	75.8
ASN (Rabinovich, Stern, and Klein 2017)	18.2	–	77.6
ASN+SUPATT (Rabinovich, Stern, and Klein 2017)	22.7	–	79.2
<b>Our system</b>	<b>27.3</b>	<b>30.3</b>	<b>79.6</b>



# Newest Results

- Replacing the CNN with a Transformer
  - Transformer: a new neural architecture at 2017
  - The flexibility of L2S allows to easily utilize new models

	Model	StrAcc	Acc+	BLEU
Plain	LPN (Ling et al., 2016)	6.1	–	67.1
	SEQ2TREE (Dong and Lapata, 2016)	1.5	–	53.4
	YN17 (Yin and Neubig, 2017)	16.2	~18.2	75.8
	ASN (Rabinovich et al., 2017)	18.2	–	77.6
	ReCode (Hayati et al., 2018)	19.6	–	78.4
	<b>CodeTrans-A</b>	<b>25.8</b>	<b>25.8</b>	<b>79.3</b>
Structured	ASN+SUPATT (Rabinovich et al., 2017)	22.7	–	79.2
	SZM19 (Sun et al., 2019)	27.3	30.3	79.6
	<b>CodeTrans-B</b>	<b>31.8</b>	<b>33.3</b>	<b>80.8</b>



# Conclusion

- Program Estimation: to find the most probable program under a context
- L2S: combining four tools to solve program estimation
- Why worked?
  - Machine learning to estimate probability
  - Expansion rules and constraints to confine the space
  - Search algorithms to locate the best program
- Better to combine the tools we have



# 参考资料

- Yingfei Xiong, Bo Wang, Guirong Fu, Linfei Zang. Learning to Synthesize. GI'18: Genetic Improvement Workshop, May 2018.