



软件分析

程序综合：概率

熊英飞
北京大学



目录

- 程序估计问题
 - 和加速的关联
- 扩展枚举的方法解决程序估计问题
- 定义不同的展开方向
- 扩展FlashMeta解决程序估计问题

很多应用需要概率最大的程序



典型应用-自动编写重复程序



```
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)
```

典型应用-缺陷修复



```
/** Compute the maximum of two values
 * @param a first value
 * @param b second value
 * @return b if a is lesser or equal to b, a otherwise
 */
public static int max(final int a, final int b) {
    return (a <= b) ? a : b;
}
```

综合出新的表达式来替换掉旧的



程序估计 Program Estimation

- 输入:
 - 一个程序空间 $Prog$
 - 一条规约 $Spec$
 - 概率模型 P ，用于计算程序的概率
- 输出:
 - 一个程序 $prog$ ，满足
 - $prog = \operatorname{argmax}_{prog \in Prog \wedge prog \vdash spec} P(prog)$
- 如果 P 估计程序满足规约的概率，那么可以用来加速传统程序综合



基本算法：穷举

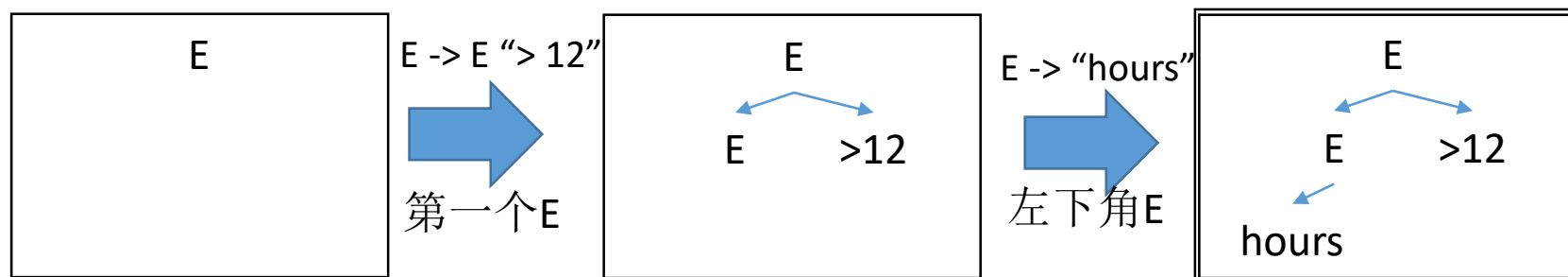
- 用枚举的方法遍历空间中的程序
 - 对每个程序计算概率
 - 返回概率最大的程序
-
- 能否优化这个过程？



扩展枚举算法求解程序估计问题



规则展开概率模型



- $P(\text{rule} \mid \text{prog}, \text{position})$
 - Prog: 当前已经展开的部分程序
 - Position: 准备展开的终结符的位置
 - Rule: 展开该终结符的概率

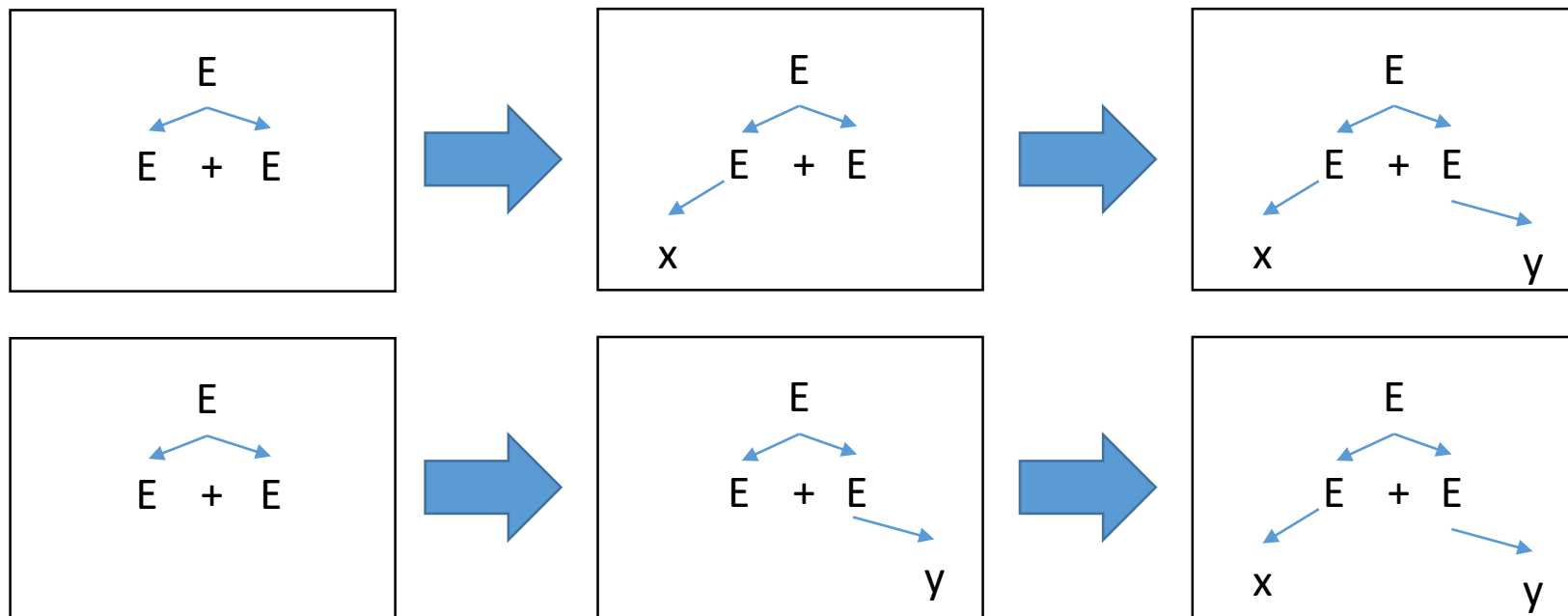


程序概率的计算

- 定理： 给定任意的规则展开序列， 我们有
 - $P(prog) = \prod_i P(rule_i \mid prog_i, position_i)$
 - $prog_i$: 第i步已经生成的程序
 - $position_i$: 第i步准备展开的非终结符的位置
 - rule: 第i步采用的产生式
 - $prog$: 完整程序



程序概率计算的收敛性



- 以上定理表明，任意展开序列都有相同概率



证明

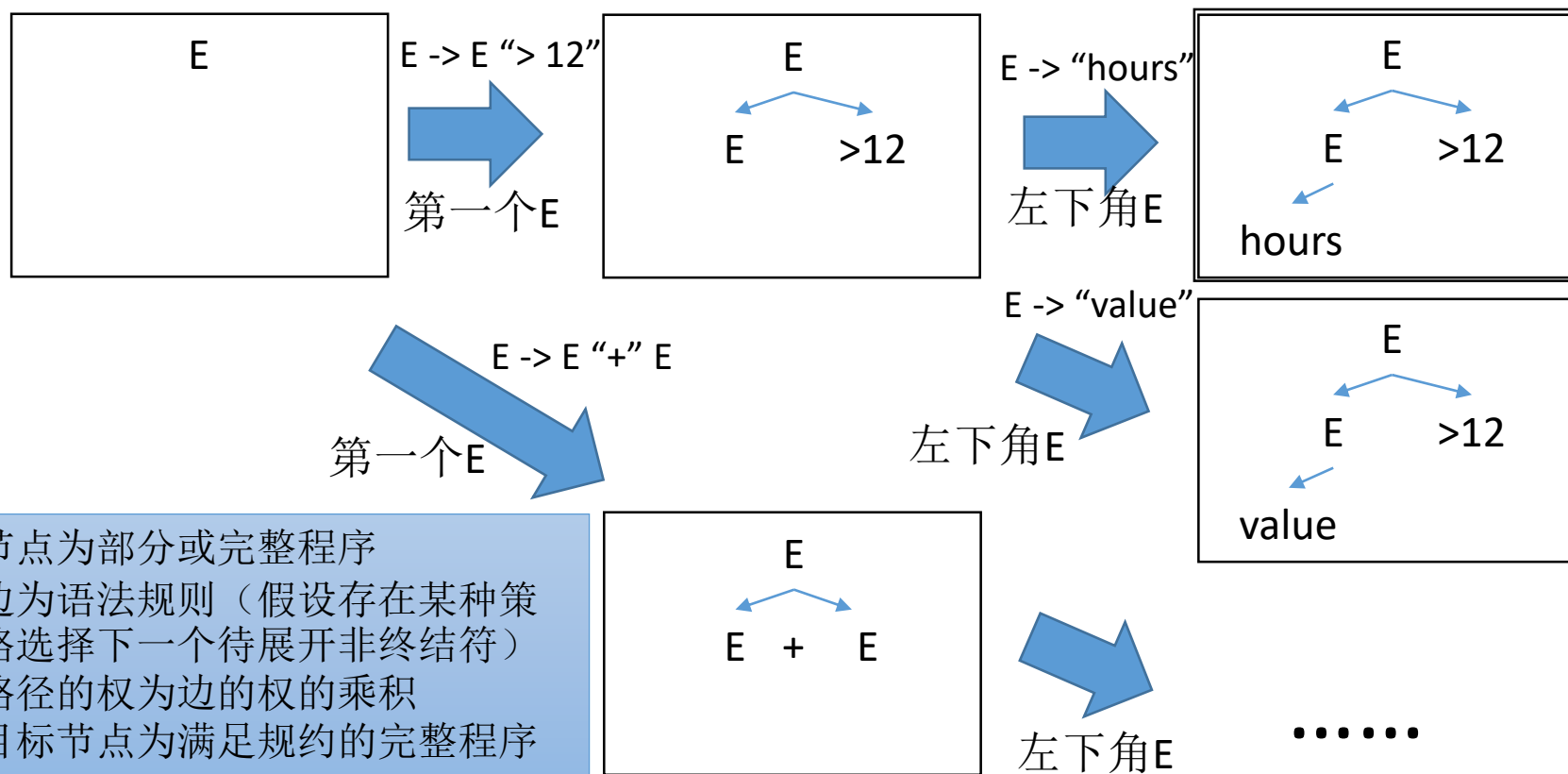
- 假设存在一个policy，决定一个不完整程序中哪个节点先被展开，那么policy的选择和prog的概率是独立的
 - $Pr(prog)$
 - $= Pr(prog \mid policy) // \text{独立性}$
 - $= Pr((\langle prog_i, pos_i, rule_i \rangle)_{i=1}^n \mid policy)$
 - $= Pr(prog_1 \mid policy) Pr(pos_1 \mid policy, prog_1)$
 $Pr(rule_1 \mid policy, prog_1, pos_1)$
 $Pr(eprog_2 \mid policy, prog_1, pos_1, rule_1) \dots$
 $Pr(eprog_{n+1} \mid policy, (eprog_i)_{i=1}^n, (pos_i)_{i=1}^n, (rule_i)_{i=1}^n)$
 - $= \prod_i Pr(rule_i \mid policy, (rule_j)_{j=1}^{i-1}, pos_i) // \text{删除概率为1的项}$
 - $= \prod_i Pr(rule_i \mid policy, prog_i, pos_i)$
 - $= \prod_i Pr(rule_i \mid prog_i, pos_i) // \text{独立性}$



规则展开概率模型的实现

- 通常计算 $P(rule_i | prog_i, position_i, context)$
 - 其中context根据需要可以为程序规约、补全的上下文等
- 可以用任意统计模型或机器学习模型实现

程序估计问题作为路径查找问题



- 节点为部分或完整程序
- 边为语法规则（假设存在某种策略选择下一个待展开非终结符）
- 路径的权为边的权的乘积
- 目标节点为满足规约的完整程序



如何求解概率最大的程序？

- 采用求解路径查找问题的标准算法
- 迪杰斯特拉算法
- 定向搜索（Beam Search）
- A*算法



迪杰斯特拉算法

- 定义节点的权为到达该节点的路径的最大权
- 维护一个可达节点列表，并记录每个节点的权
- 选择(1)权最大的节点, (2)从该节点出发的未探索的权最大的边，将新到达节点加入列表
- 如果某个节点已经没有未探索出边，则从列表中删除
- 反复上一步直到找到目标节点

注：在本问题中只能被一条路径到达，而在一般路径查找问题中，每个节点可以被多条路径达到，所以通用算法还需到达了旧节点时更新最大权。



迪杰斯特拉算法求解的例子

- $\langle E, 1 \rangle$
- $\langle E+E, 0.5 \rangle, \langle E-E, 0.4 \rangle, \langle \cancel{x}, 0.05 \rangle, \langle \cancel{y}, 0.05 \rangle$
- $\langle E-E, 0.4 \rangle, \langle x+E, 0.3 \rangle, \langle (E+E)+E, 0.1 \rangle, \langle y+E, 0.1 \rangle$
- $\langle x+E, 0.3 \rangle, \langle x-E, 0.2 \rangle, \langle y-E, 0.1 \rangle, \langle (E+E)+E, 0.1 \rangle, \langle y+E, 0.1 \rangle, \langle (E+E)-E, 0.05 \rangle, \langle (E-E)-E, 0.05 \rangle$
-



定向搜索 (Beam Search)

- 在迪杰斯特拉算法中不保留所有节点，只保留概率最大的k个
- 近似算法，不保证最优，也不保证找到结果



A*算法

- 节点n的权=到达该节点的权*h(n)
 - $h(n)$ =剩余路径权的上界
- 其他同迪杰斯特拉算法
- 如何知道剩余路径权的上界?
 - 假设存在函数 $\hat{P}(rule)$ ，满足
 - $\forall prog, position: \hat{P}(rule) \leq P(rule \mid prog, position)$
 - 在语法展开式上做静态分析，分析出每个非终结符的概率上界
 - 从 $E \rightarrow E + E \mid x \mid y \mid \dots$
 - 得到方程 $\hat{P}(E) = \max(\hat{P}(E \rightarrow E + E)\hat{P}(E)\hat{P}(E), \hat{P}(E \rightarrow x), \hat{P}(E \rightarrow y), \dots)$
 - 剩余路径权的上界为所有未展开非终结符概率上界的积



剪枝

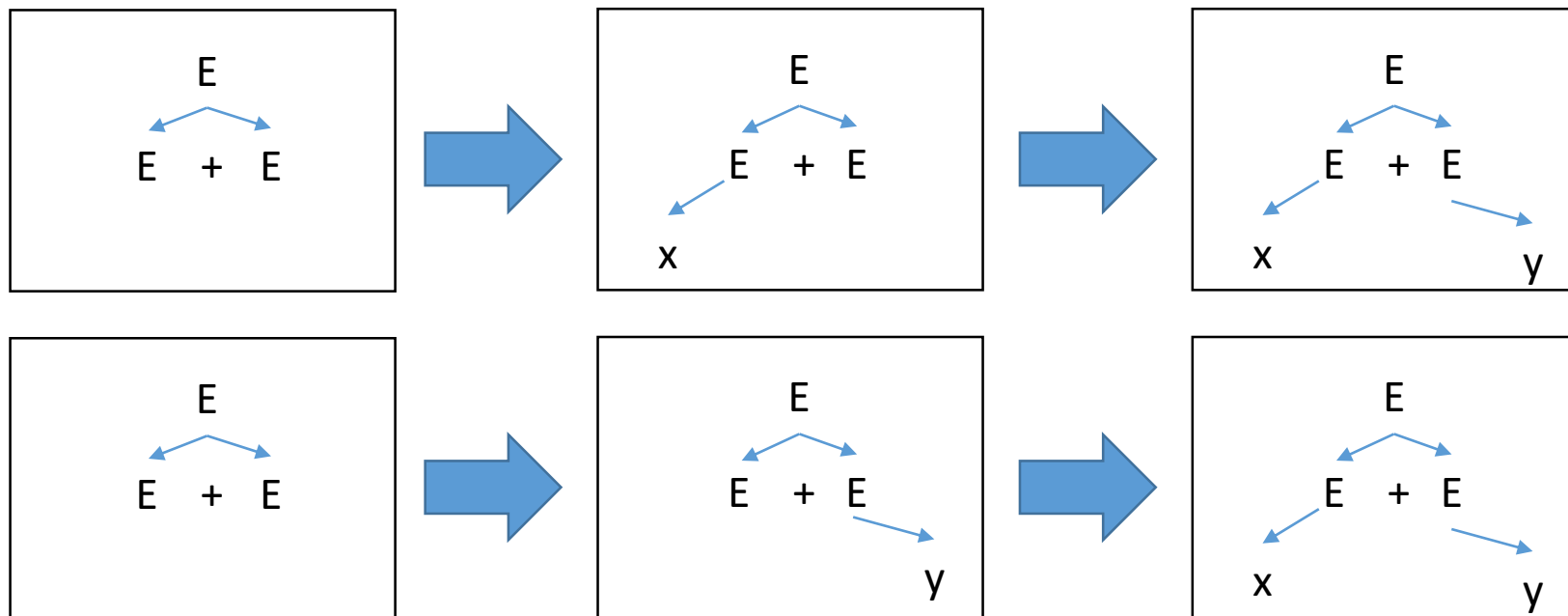
- 之前描述的剪枝过程仍然可以用于求解程序估计问题
- 判断出一个部分程序无法满足规约时，从列表中移除对应节点



定义程序展开的顺序



展开的顺序

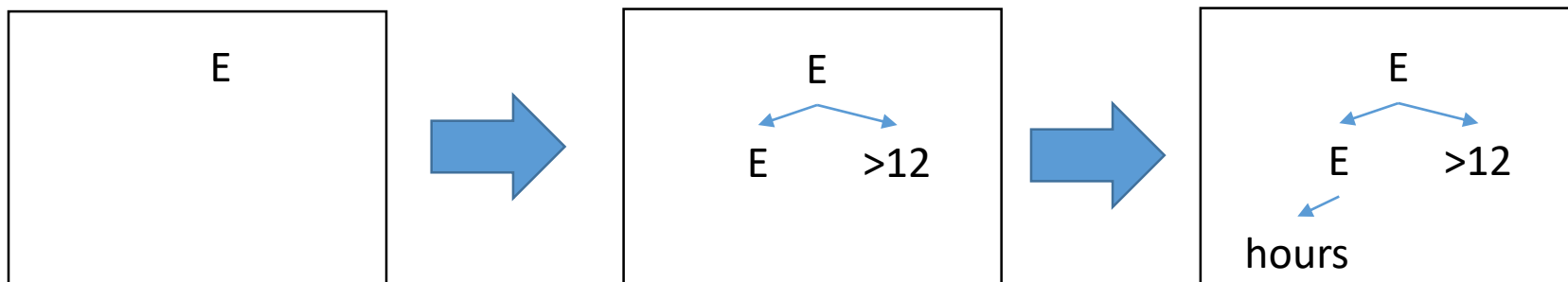


- 如果左下采用 $E \rightarrow x$ 的概率极大，而右下采用 $E \rightarrow y$ 的概率较低，则上面的顺序能显著减少搜索时间
- 需要根据应用特点定义非终结选择策略

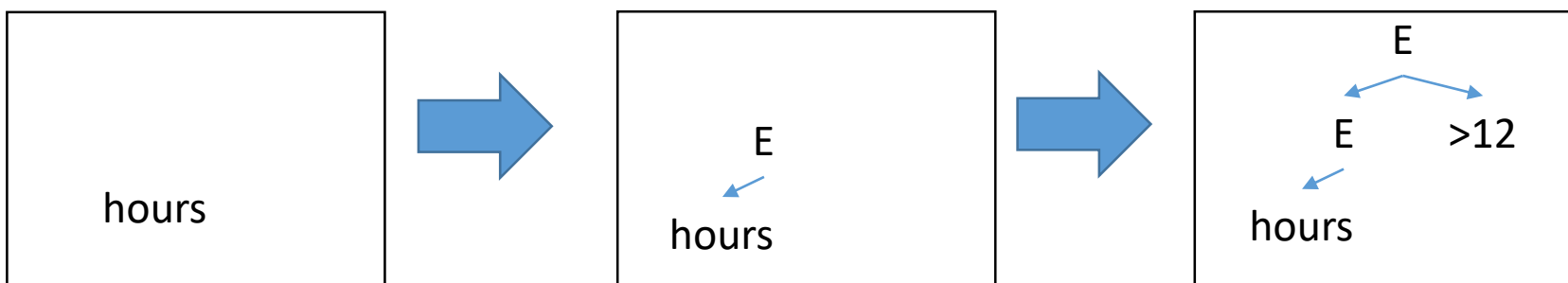


超越上下文无关文法的顺序?

- 自顶向下



- 自底向上





扩展规则

- 允许描述不同方向的语法扩展
- 由本课题组提出
- 通过采用合适的扩展规则，求解效率可提高一倍以上



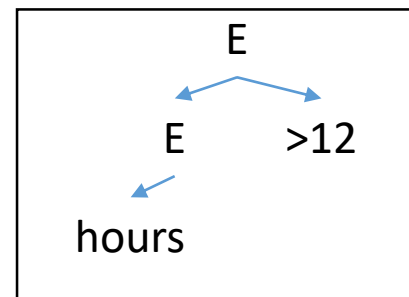
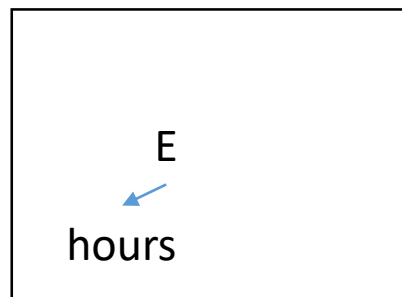
从上下文无关文法到扩展规则

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


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

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

如果第i个子节点已经产生, 产生整棵子树



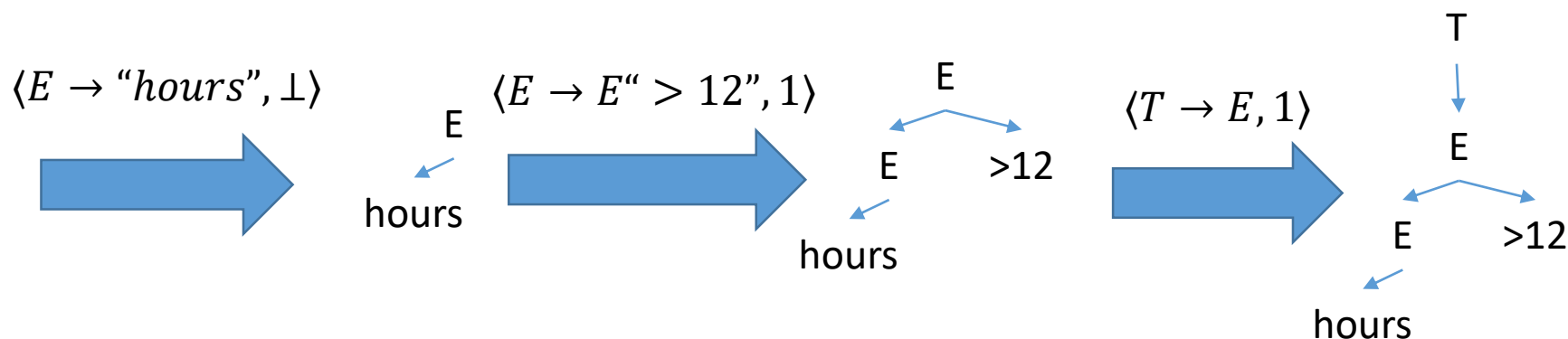
自顶向下规则: $\langle E \rightarrow E " > 12", 0 \rangle$

如果根节点已经产生, 产生整颗子树

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

从零产生一颗子树

基于扩展规则的程序生成过程





扩展规则树Expansion Tree

- 抽象语法树在扩展规则上的对应，记录扩展规则如何被应用的

hours>12	hours+value
$\begin{array}{c} (T \rightarrow E, 1) \\ \uparrow \\ (E \rightarrow E \text{ " > 12" }, 1) \\ \uparrow \\ (E \rightarrow \text{"hours"}, \perp) \end{array}$	$\begin{array}{c} (T \rightarrow E, 1) \\ \uparrow \\ (E \rightarrow E \text{ "+" } E, 1) \\ \swarrow \quad \searrow \\ (E \rightarrow \text{"hours"}, \perp) \quad (E \rightarrow \text{"value"}, 0) \end{array}$



抽象语法树 \rightarrow 扩展规则树

- 扩展规则的性质
 - 完整性: 对任意AST, 至少有一个扩展规则树
 - 唯一性: 对任意AST, 最多有一个扩展规则树
- 是否总是存在完整和唯一的扩展规则集合?



唯一和完整集合的充分条件

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


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

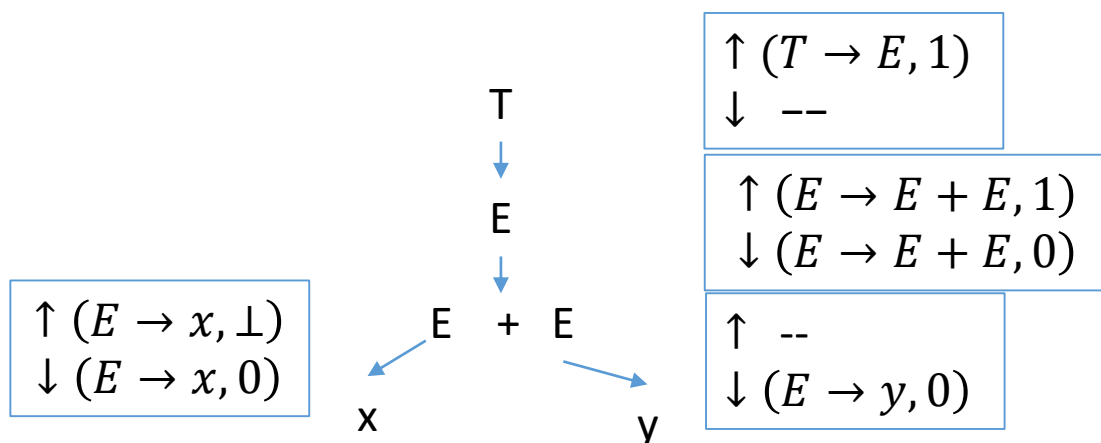
1. 除了初始符号开头的规则，所有语法规则都有对应的自顶向下展开规则
2. 所有语法规则最多只有一条自底向上的展开规则
3. 对于所有从初始符号（延自底向上展开规则）反向可达的非终结符，其所有语法规则都有一条自底向上展开规则或创建规则

从初始符号开始选择创建/自底向上规则即可



抽象语法树 -> 扩展规则树

- 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





求解程序估计问题

- 给定某种结点选择策略，可以从扩展规则树得到展开序列
- 同样看做路径查找问题求解



扩展FlashMeta求解程序估计问题



FlashMeta vs 程序估计问题

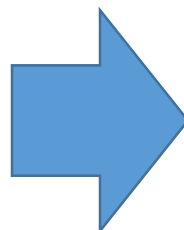
- 能否直接采用FlashMeta求解程序估计问题?
- 之前的方法无法使用
 - 子问题无法支持基于部分模型的概率模型
 - 分解子问题后难以看做路径查找问题
- MaxFlash
 - 2020年由北京大学吉如一等人提出
 - 效率超过FlashMeta达400-2000倍





例：字符串处理程序合成

Start symbol	S	\rightarrow	$N_S \mid N_Z$
String expr	N_S	\rightarrow	Parameters $\mid (+ N_S N_S)$ $\mid (\text{CHARAT } N_S N_Z) \mid \text{'.'}$
Integer value	N_Z	\rightarrow	$0 \mid 1$

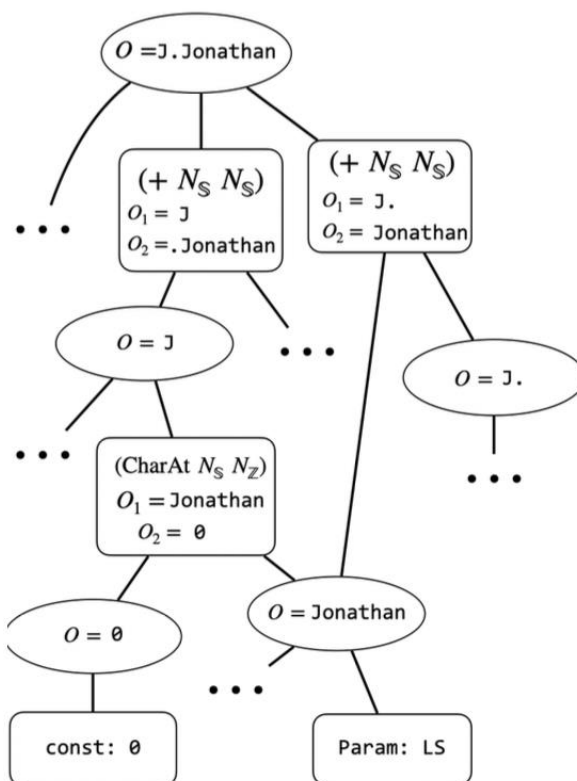


$(+ (\text{CHARAT } FS \ 0) (+ \text{'.' } LS))$

$(\text{'John'}, \text{'Jonathan'}) \rightarrow \text{'J.Jonathan'}$



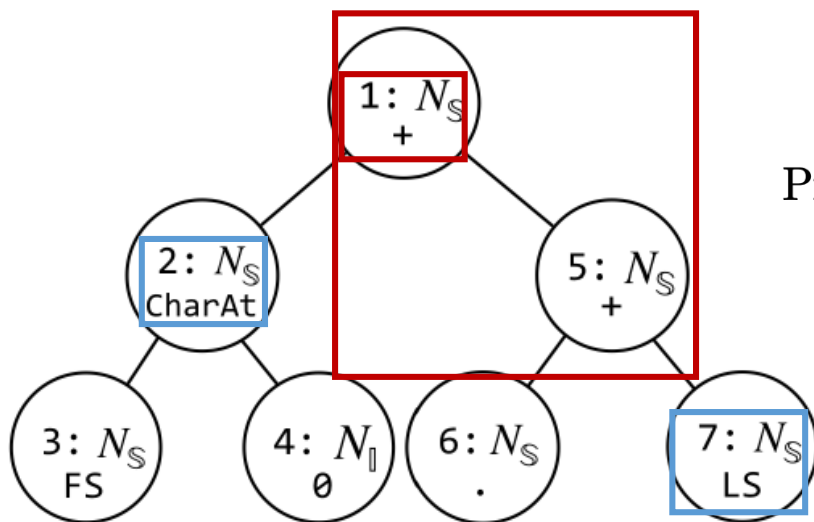
FlashMeta求解过程



概率模型：自顶向下预测模型 Topdown Prediction Model



- **TPM:** 节点展开规则概率只取决于其祖先，即兄弟节点的展开规则相互独立
- 示例:



$$\Pr[\text{CharAt on vertex 2}] = \Pr[\text{CharAt} | (+, 1)]$$

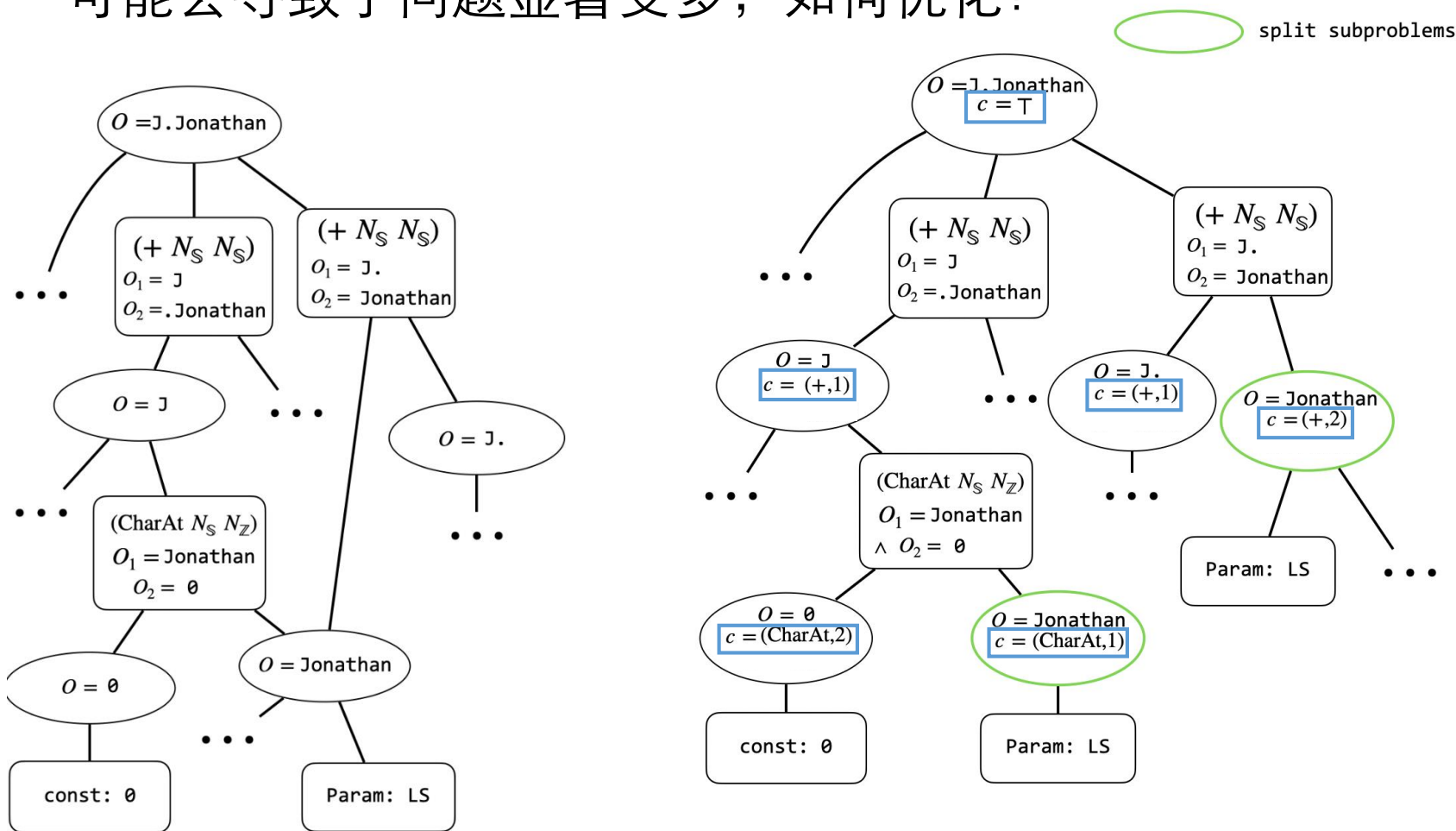
$$\Pr[\text{LS on vertex 7}] = \Pr[\text{LS} | (+, 2), (+, 1)]$$

通常定义为依赖最近k层祖先节点



在FlashMeta中引入概率

- 为子问题添加祖先节点的上下文信息
- 可能会导致子问题显著变多，如何优化？



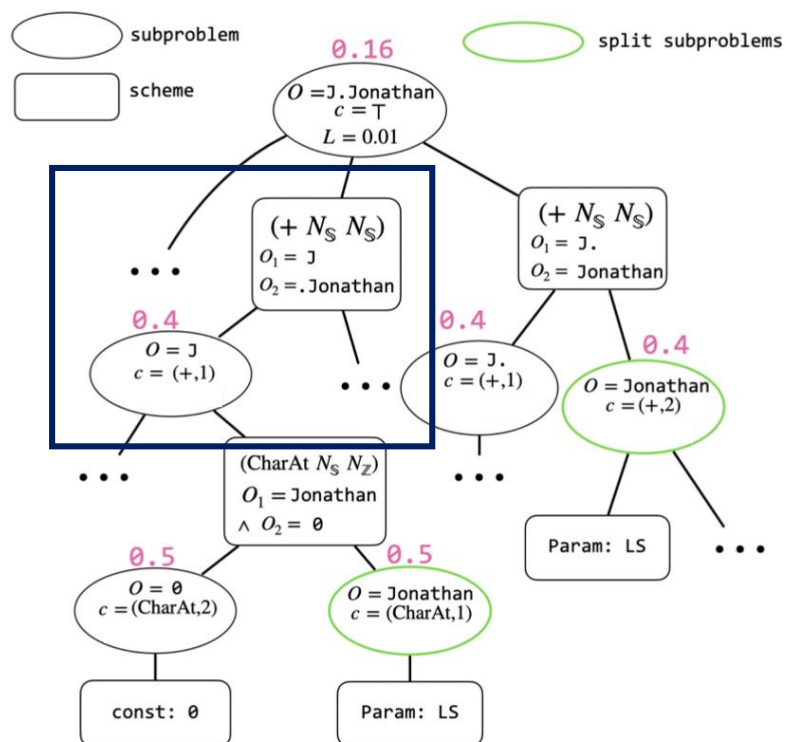


搜索方法

- 分支限界 (Branch and Bound)
 - 在搜索时对每一个子问题带上一个概率下界 l ，表示只有在当前子程序的概率大于 l 的时候，才有可能形成一个可能的程序
- 迭代加深 (Iteratively Deepening)
 - 设置全局的概率下界，如果找不到解则放宽全局的概率下界
- 估价函数 (Heuristic Function)
 - 类似之前A*算法的估价函数
 - 估计每一个子问题的最优解上界
 - 如：每个非终结符的概率上界
 - 如：每个<上下文-非终结符>的概率上界



示例



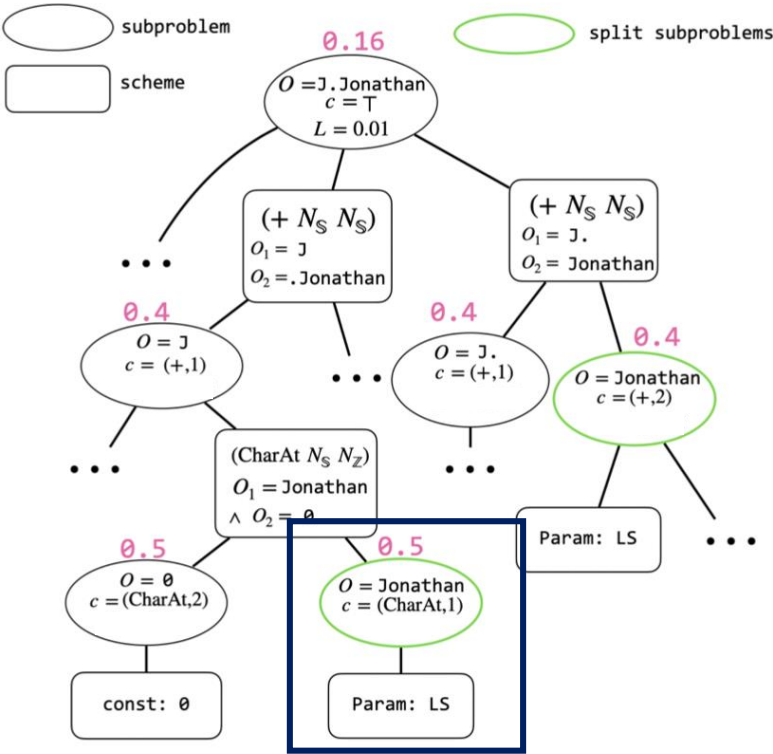
$$\Pr[+ | \text{empty}] = 1$$

$$\begin{array}{c} \boxed{\begin{array}{ccc} (+ & P_1 & P_2) \\ 1 & \nearrow & \leq 0.4 \end{array}} \geq 0.01 \\ \geq 0.025 \end{array}$$

	CHARAT	+	FS	LS	'.'	0	1
Start	0	1	0	0	0	0	0
(+, 1)	0.5	0.01	0.09	0	0.4	0	0
(+, 2)	0.05	0.5	0.05	0.4	0	0	0
(CHARAT, 1)	0.05	0.05	0.5	0.4	0	0	0
(CHARAT, 2)	0	0	0	0	0	0.5	0.5



示例



	CHARAT	+	FS	LS	'.'	0	1
(CHARAT, 1)	0.05	0.05	0.5	0.4	0	0	0

跳过

	CHARAT	+	FS	LS	'.'	0	1
Start	0	1	0	0	0	0	0
(+, 1)	0.5	0.01	0.09	0	0.4	0	0
(+, 2)	0.05	0.5	0.05	0.4	0	0	0
(CHARAT, 1)	0.05	0.05	0.5	0.4	0	0	0
(CHARAT, 2)	0	0	0	0	0	0.5	0.5



如何重用子问题

- 重复子问题的结果复用是动态规划的核心
 - FlashMeta 的子问题：当前的非终结符 S ，输入输出用例 A
 - 目前的子问题：非终结符 S ，输入输出用例 A ，**概率下界** L ，上下文 c
- 需要复用概率下界不同的子问题
 - 考虑两个除了概率下界不同以外，其他都一样的子问题 $(P, 0.2)$, $(P, 0.1)$
 - Case 1: $(P, 0.2)$ 先于 $(P, 0.1)$
 - 有解，则同样是 $(P, 0.1)$ 的解；
 - 无解，则可以更新 P 的估价函数
 - Case 2: $(P, 0.1)$ 先于 $(P, 0.2)$
 - 有解，则同样是 $(P, 0.2)$ 的解（因为总是搜索概率最大的结果）
 - 无解， $(P, 0.2)$ 同样无解



MaxFlash总结

- 分支限界和启发式函数：给定概率下界的时候，只搜索部分满足下界的程序空间
- 迭代加深：逐步放宽概率下界，避免一次探索过大空间
- 子问题重用：概率下界不同的子问题也能复用，使得不同迭代之间的计算不会被浪费



参考文献

- Yingfei Xiong, Bo Wang, Guirong Fu, Linfei Zang. Learning to Synthesize. GI'18: Genetic Improvement Workshop, May 2018.
- Ruyi Ji, Yican Sun, Yingfei Xiong, Zhenjiang Hu. Guiding Dynamic Programing via Structural Probability for Accelerating Programming by Example. OOPSLA'20: Object-Oriented Programming, Systems, Languages, and Applications 2020, November 2020.