

**Об усилении хигманова  
вложения при построении  
регулярных приближений циклов**

**Антонина Н. Непейвода  
Институт программных систем РАН  
г. Переславль-Залесский**

**Объединенный семинар ИЦСА и ИЦПУ  
26 февраля 2013 г., Переславль–Залесский**

*Семантическое дерево программы* — дерево путей исполнения программы на возможных входных данных.

## **Пример**

Рекурсивное определение факториала в арифметике Пеано

$$f(0)=1;$$

$$f(x+1)=m(x+1, f(x));$$

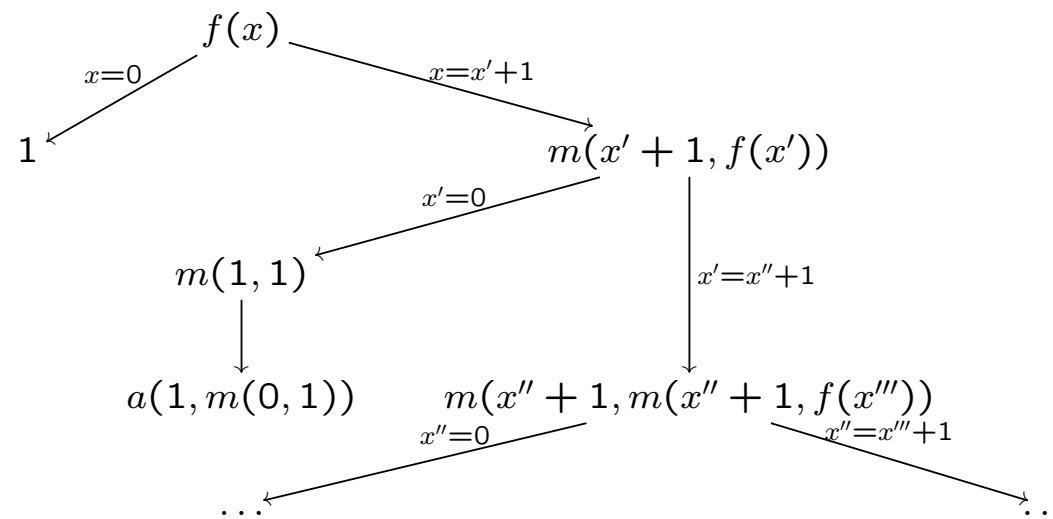
$$m(0, y)=0;$$

$$m(x+1, y)=a(y, m(x, y));$$

$$a(0, y)=y;$$

$$a(x+1, y)=a(x, y)+1;$$

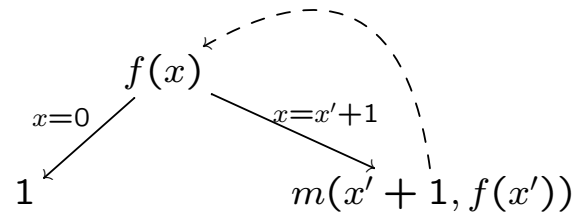
## Начало развертки семантического дерева исполнения $f(x)$



Это дерево может содержать в себе пути сколь угодно большой длины (зависящей от  $x$ ). Такие пути семантически соответствуют циклам.

Имеет смысл задача нахождения критериев для построения приближений циклов в программах. В частности, для семантических деревьев такой критерий должен сравнивать два узла  $Node_1$ ,  $Node_2$ , таких, что  $Node_1$  — предок  $Node_2$ , и строить аппроксимацию цикла, если структура  $Node_2$  некоторым образом повторяет структуру  $Node_1$ .

## Пример



$f(x)$  — подтерм  $m(x'+1, f(x'))$  с точностью до переименования переменных. На основании этого наблюдения ветвь дерева сворачивается в граф.

Отношение  $R$  называется *почти полным (almost well)* на множестве последовательностей строк  $S \subset \Sigma^*$  в некотором алфавите  $\Sigma$ , если всякая бесконечная последовательность строк  $\{a_n\}$  содержит  $a_i, a_j$ , такие, что  $i < j$  и  $R(a_i, a_j)$ .

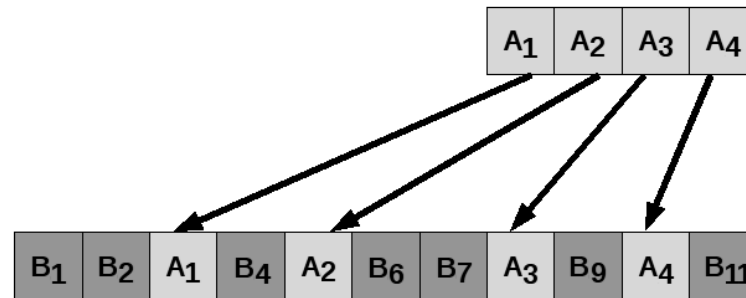
Это свойство гарантирует, что аппроксимация цикла будет построена на всех бесконечных ветвях семантического дерева программы.

Последовательность  $\{a_n\}$ , обладающая свойством  $\forall i, j (i < j \Rightarrow (a_i, a_j) \notin R)$ , называется *плохой последовательностью*.

Если отношение транзитивно и почти полно, оно называется *полным квазипорядком (wqo)*.

## Пример

Пусть дан конечный алфавит  $\Sigma$ .  $A = a_1a_2\dots a_m$ ,  $B = b_1b_2\dots b_n$ ,  $\forall i, j (i \geq 1 \ \& \ j \geq 1)$ . Если  $\forall i (i < m \Rightarrow \exists j, k (b_j = a_i \ \& \ b_k = a_{i+1} \ \& \ k > j))$ , то  $A$  вкладывается в  $B$  в смысле Хигмана ( $A \trianglelefteq B$ ).



$\trianglelefteq$  — почти полное отношение на произвольных последовательностях слов в конечном алфавите (Хигман, 1952).

В 1995 году М.Г. Серенсен предложил использовать хигманово вложение термов в узлах деревьев для построения аппроксимаций циклов в семантическом дереве программы.

Ранее, в 1988 году, В.Ф. Турчин предложил использовать для построения аппроксимаций циклов отношение над стеками функциональных вызовов, сходное с хигмановым вложением.



Набор  $\langle \Sigma, \mathbf{R}, \Gamma_0 \rangle$ , где  $\Sigma$  — алфавит  $\Gamma_0 \in \Sigma^+$  — начальное слово и  $\mathbf{R} \subset \Sigma^+ \rightarrow \Sigma^*$  — правила переписывания, называется *системой переписывания префиксов* (СПП), если правило  $R : R_l \rightarrow R_r$  может применимо лишь к словам вида  $R_l\Phi$  и порождает слова вида  $R_r\Phi$ .

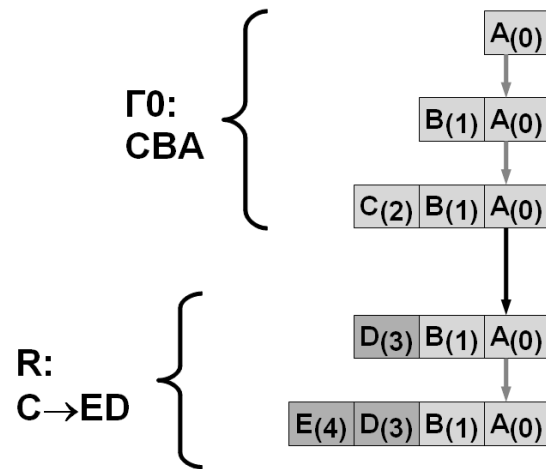
### Пример

$aabab$  после применения правила  $a \rightarrow bb$  превращается в  $bbabab$ , но не в  $abbbab$  или  $aabbbb$ .

Если  $|R_l| = 1$  для всех  $R : R_l \rightarrow R_r$ , то система переписывания префиксов называется *алфавитной* (АСПП).

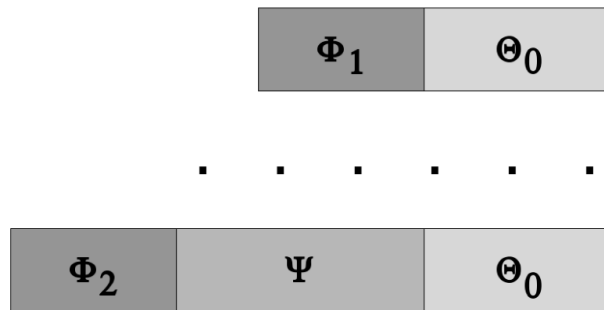
Пусть последовательность слов  $\{\Phi_i\}_{i=1}^n$  порождена СПП  $G = \langle \Sigma, \mathbf{R}, \Gamma_0 \rangle$ . Пометим все  $\Phi_i$  индексами времени:

1.  $\Phi_1[i]$  помечается индексом  $|\Gamma_0| - i$ ;
2. Пусть наибольший индекс времени в отрезке последовательности  $\{\Phi_i\}_{i=1}^k$  ( $k < n$ ) равен  $M$  и  $\Phi_{k+1}$  порождено из  $\Phi_k$  правилом  $R : R_l \rightarrow R_r$ . Тогда  $\Phi_{k+1}[i]$  ( $i \leq |R_r|$ ) получает временной индекс  $M + |R_r| - i + 1$ . Остальные буквы в  $\Phi_{k+1}$  сохраняют прежние временные индексы.



## Определение отношения Турчина

$$\Gamma \preceq \Delta \Leftrightarrow \Gamma = \Phi_1 \Theta_0 \ \& \ \Delta = \Phi_2 \Psi \Theta_0 \ \& \ \Phi_1 \approx \Phi_2$$



## Теорема Турчина

$\preceq$  — почти полное отношение на последовательностях, порожденных АСПП над конечным алфавитом.

## Пример

$G_F$ :

$$R_1 : f \rightarrow \Lambda \quad R_3 : m \rightarrow \Lambda \quad R_5 : a \rightarrow \Lambda$$

$$R_2 : f \rightarrow fm \quad R_4 : m \rightarrow ma \quad R_6 : a \rightarrow a$$

Начальное слово  $\Gamma_0$  —  $fa$ . Начальный отрезок порождаемой последовательности может выглядеть следующим образом

$$\begin{array}{c} \Gamma_0 : \mathbf{f}_{(1)}a_{(0)} \\ \downarrow R_2 \\ \Gamma_1 : \mathbf{f}_{(3)}m_{(2)}a_{(0)} \\ \downarrow R_2 \\ \Gamma_2 : \mathbf{f}_{(5)}m_{(4)}m_{(2)}a_{(0)} \\ \downarrow R_1 \\ \Gamma_3 : \mathbf{m}_{(4)}m_{(2)}a_{(0)} \\ \downarrow R_4 \\ \Gamma_4 : \mathbf{m}_{(7)}a_{(6)}m_{(2)}a_{(0)} \end{array}$$

$$\Gamma_0 \preceq \Gamma_1, \Gamma_0 \preceq \Gamma_2, \Gamma_1 \preceq \Gamma_2, \text{ и } \Gamma_3 \preceq \Gamma_4.$$

Верхняя оценка на длину слова в плохой последовательности, порожденной АСПП

$$|\Gamma_i| \leq |\Gamma_0| + \sum (|(R_i)_r| - 1)$$

$\sum (|(R_i)_r| - 1)$  пробегает все правила с  $|R_r| \neq \Lambda$ .

Обозначим  $\Gamma^-$  слово  $\Gamma$  без первой буквы. Правило  $R$  неукорачивающее, если  $R_r \neq \Lambda$ .

Пусть  $R$  — неукорачивающее правило,  $R(\Theta_0[1])\Theta_0^-$  предшествует  $R(\Theta_1[1])\Theta_1^-$  и  $\exists i(\Theta_1^-[i] = \Theta_0^-[1])$ , тогда  $R(\Theta_0[1])\Theta_0^- \preceq R(\Theta_1[1])\Theta_1^-$ .

Ограниченность длины слова в плохой последовательности доказывает почти полноту отношения  $\preceq$  в конечном алфавите.

Грубая верхняя оценка на длину плохой последовательности относительно  $\preceq$

$$C_{Max} = \text{card}(\Sigma)^{|\Gamma_0| + \sum (|(R_i)_r| - 1)}.$$

АСПП  $G = \langle \Sigma, \mathbf{R} \subset \Sigma \rightarrow \Sigma^*, \Gamma_0 \rangle$  называется *обогащенной* или *ϑ-системой*, если

1. Для всяких двух правил  $R : a \rightarrow R_r, R' : b \rightarrow R'_r \exists i, j (R_r[i] \approx R'_r[j]) \Rightarrow i = j \ \& \ |R| = |R'| \ \& \ \forall i (i < |R| \Rightarrow R_r[i] \approx R'_r[i])$ ;
2. Если  $R^a \in \mathbf{R}, R^a : a \rightarrow R_r$  и  $b \in \Sigma$ , то  $R^b : b \rightarrow R_r \in \mathbf{R}$ .

Рассмотрим следующее преобразование АСПП  $G$  в обогащенную  $G'$ .

1. Пусть  $a = R_r[i]$ ,  $a \in \Sigma$ . Заменяем  $a$  на пару  $(a, 2^r * 3^{i-1})$ . Начальное слово считаем правилом номер 0.
2. К каждой паре правил  $a_1 \rightarrow \Phi$ ,  $a_2 \rightarrow \Psi$  добавим правила  $a_1 \rightarrow \Psi$  и  $a_2 \rightarrow \Phi$  ( $a_i \in \Sigma$ ).
3. Пусть  $(a_i, n_i) = R'_r[i]$ . Для каждого  $a_i \rightarrow \Phi$  добавим в  $R$  правила  $(a_i, n_i) \rightarrow \Phi$ .

Если в  $\mathcal{G}$ -системе  $\exists R(R : (a, n) \rightarrow R_r)$ , то для всякого  $(b, m)$   $\exists R(R : (b, m) \rightarrow R_r)$ , поэтому правила переписывания в  $\mathcal{G}$ -системах будут записываться в общем виде:  $R : x \rightarrow R_r$ .



## Пример

Переведем АСПП  $G_F$  в  $\mathfrak{G}$ -форму.

$G'_F$ :

$$R_0 : x \rightarrow f^1 a^3 \quad R_2 : x \rightarrow m^4 a^{12}$$

$$R_1 : x \rightarrow f^2 m^6 \quad R_3 : x \rightarrow a^8$$

$$R_4 : x \rightarrow \Lambda$$

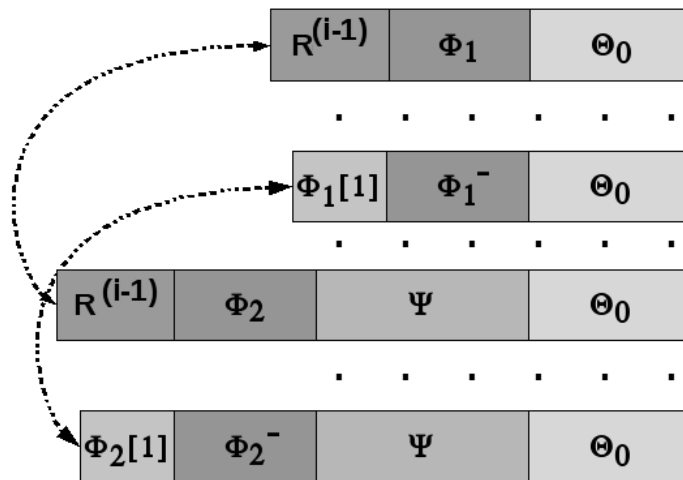
Порожденная последовательность после преобразования:

$$\begin{array}{c} \Gamma_0 : \mathbf{f}_{(1)}^1 a_{(0)}^3 \\ \downarrow R_1 \\ \Gamma_1 : \mathbf{f}_{(3)}^2 m_{(2)}^6 a_{(0)}^3 \\ \downarrow R_1 \\ \Gamma_2 : \mathbf{f}_{(5)}^2 m_{(4)}^6 m_{(2)}^6 a_{(0)}^3 \\ \downarrow R_4 \\ \Gamma_3 : \mathbf{m}_{(4)}^6 m_{(2)}^6 a_{(0)}^3 \\ \downarrow R_2 \\ \Gamma_4 : \mathbf{m}_{(7)}^4 a_{(6)}^{12} m_{(2)}^6 a_{(0)}^3 \end{array}$$

Теперь  $\Gamma_0 \not\leq \Gamma_1$ .

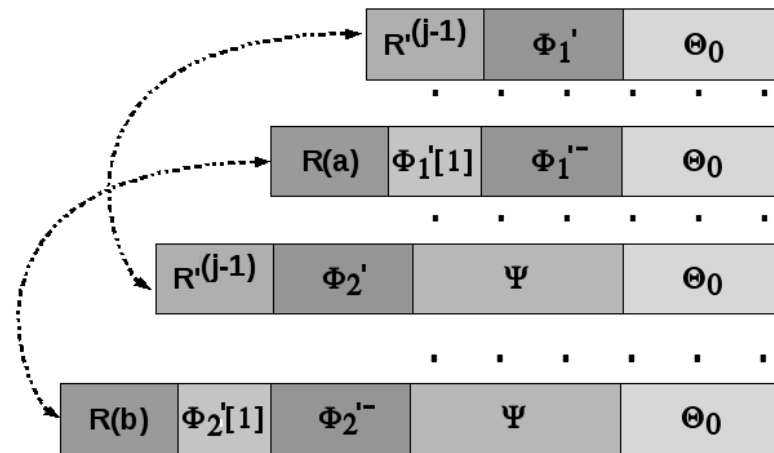
Самая первая пара вложенных по Турчину слов в последовательности, порожденной  $\mathfrak{G}$ -системой, имеет вид  $R(a)\Theta_0, R(b)\Psi\Theta_0$ .

Рассмотрим первую пару  $\Phi_1\Theta_0, \Phi_2\Psi\Theta_0$ , такую, что  $\Phi_1\Theta_0 \preceq \Phi_2\Psi\Theta_0$  ( $\Phi_1 \approx \Phi_2$ ), обрывающую плохую последовательность.  $\Phi_1[1]$  и  $\Phi_2[1]$  должны порождаться различными применениями  $R$ , и если  $\Phi_1[1] \approx R_r[i]$ , то  $\Phi_2[1] \approx R_r[i]$ . Обозначим слово  $R_r[1]R_r[2]\dots R_r[i-1]$  как  $R^{(i-1)}$  и рассмотрим моменты применения  $R$ . Первый имеет вид  $R_{(k_1)}^{(i-1)}\Phi_1\Theta_0$ , второй —  $R_{(k_2)}^{(i-1)}\Phi_2\Psi\Theta_0$ .



Они образуют турчинскую пару, следовательно, совпадают с  $\Phi_1\Theta_0$  и  $\Phi_2\Psi\Theta_0$ .

Итак,  $\Phi_1 = R(a)\Phi'_1$ ,  $\Phi_2 = R(b)\Phi'_2$  ( $\Phi'_1 \approx \Phi'_2$ ). Пусть  $\Phi'_1 \neq \Lambda$ . Тогда  $\exists R', j (\Phi'_1[1] \approx R'_r[j] \ \& \ \Phi'_2[1] \approx R'_r[j])$  и  $\Phi'_1[1] \neq \Phi'_2[1]$ . Обозначим слово  $R'[1]R'[2]\dots R'[j-1]$  как  $R^{(j-1)}$  и рассмотрим моменты применения  $R'$  к предшественникам  $\Phi_1\Theta_0$  и  $\Phi_2\Psi\Theta_0$ .



Они выглядят как  $R^{(j-1)}_{(l_1)}\Phi'_1\Theta_0$  и  $R^{(j-1)}_{(l_2)}\Phi'_2\Psi\Theta_0$  и образуют турчинскую пару. Следовательно,  $\Phi_1\Theta_0 = R(a)\Theta_0$  и  $\Phi_2\Psi\Theta_0 = R(b)\Psi\Theta_0$ .

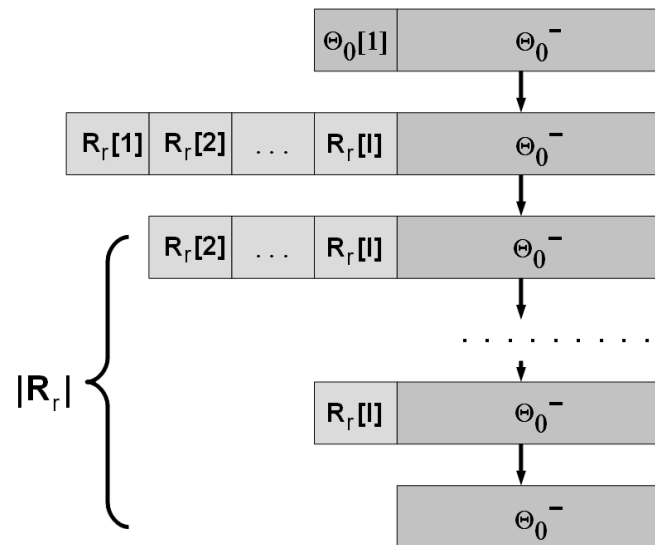
Применение правила  $R(\Delta[1])$  отменено в  $\Delta_2 \Leftrightarrow \Delta_2[1] = \Delta_1^- [1] \vee \forall i(\Delta_2[i] \neq \Delta_1^- [1])$ .

Если некоторое неукорачивающее правило  $R$  никогда не применялось при порождении слова  $\Gamma$  либо все его применения были отменены (назовем такие правила *свободными*), то применение  $R$  позволит построить более длинную плохую последовательность, чем стирание первой буквы  $\Gamma$ .

Если неукорачивающие правила применяются к начальному слову длины 1 в порядке  $R_0, R_1, \dots, R_N$  (то есть  $R_i$  всегда приоритетнее, чем  $R_{i+j}$ , если оба эти правила применимы без образования турчинской пары), то наибольшая длина плохой последовательности, которая может быть построена этими правилами

$$C'_{Max} = 1 + |(R_0)_r| * (1 + |(R_1)_r| * (\dots * (1 + |(R_N)_r|) \dots))$$

Пусть в момент  $\Theta_0$  осталось только одно свободное неукорачивающее правило  $R_0$ . Длина наибольшего сегмента плохой последовательности, построенном на  $\Theta_0$  и заканчивающегося  $\Theta_0^-$ , равна  $|(R_0)_r| + 1$ .



Пусть осталось  $N + 1$  свободных неукорачивающих правил.

Применим самое приоритетное правило  $R_k$ , получим  $R_k(\Theta_0[1])\Theta_0^-$ . На нем построим самую длинную плохую последовательность, сохраняющую  $R_k^-(\Theta_0[1])\Theta_0^-$  (длину этого сегмента будем считать известной и обозначим за  $C(N)$ ). Сегмент завершится словом  $R_k^-(\Theta_0[1])\Theta_0^-$ . На нем построим самую длинную плохую последовательность, сохраняющую  $R_k^{--}(\Theta_0[1])\Theta_0^-$ , и т.д. до момента  $\Theta_0^-$ . Каждая буква в  $(R_k)_r$  породит  $C(N)$  элементов плохой последовательности. Итоговое увеличение длины равно  $1 + |(R_k)_r| * C(N)$ .

## Пример

Построим самую длинную плохую последовательность в  $\mathcal{G}$ -системе  $G'_F$ . Ее длина равна  $2 * (1 + 2 * (1 + 2 * (1 + 1))) = 22$ .

$$\begin{aligned} G'_F: \\ R_0 : x \rightarrow ab \quad R_2 : x \rightarrow ef \\ R_1 : x \rightarrow cd \quad R_3 : x \rightarrow g \\ R_4 : x \rightarrow \Lambda \end{aligned}$$

$$\begin{array}{ll} \Gamma_0 : & a_{(1)}b_{(0)} \\ \Gamma_1 : & c_{(3)}d_{(2)}b_{(0)} \\ \Gamma_2 : & e_{(5)}f_{(4)}d_{(2)}b_{(0)} \\ \Gamma_3 : & g_{(6)}f_{(4)}d_{(2)}b_{(0)} \\ \Gamma_4 : & f_{(4)}d_{(2)}b_{(0)} \\ \Gamma_5 : & g_{(7)}d_{(2)}b_{(0)} \\ \Gamma_6 : & d_{(2)}b_{(0)} \\ \Gamma_7 : & e_{(9)}f_{(8)}b_{(0)} \\ \Gamma_8 : & g_{(10)}f_{(8)}b_{(0)} \\ \Gamma_9 : & f_{(8)}b_{(0)} \\ \Gamma_{10} : & g_{(11)}b_{(0)} \\ \Gamma_{11} : & b_{(0)} \\ \Gamma_{12} : & c_{(13)}d_{(12)} \\ \Gamma_{13} : & e_{15}f_{(14)}d_{(12)} \\ \Gamma_{14} : & g_{16}f_{(14)}d_{(12)} \\ \Gamma_{15} : & f_{(14)}d_{(12)} \\ \Gamma_{16} : & g_{(17)}d_{(12)} \\ \Gamma_{17} : & d_{(12)} \\ \Gamma_{18} : & e_{(19)}f_{(18)} \\ \Gamma_{19} : & g_{(20)}f_{(18)} \\ \Gamma_{20} : & f_{(18)} \\ \Gamma_{21} : & g_{(21)} \end{array}$$

Применение любого из  $R_1$ – $R_4$  к  $\Gamma_{21}$  породит турчинскую пару.



$\preceq$  не является транзитивным, но существует транзитивное отношение  $T$ ,  $T \subset \preceq$ , являющееся почти полным на всех последовательностях, порожденных АСПП.

$G_{ABC}$ :

$$\begin{aligned} R_1 &: a \rightarrow \Lambda & R_3 &: d \rightarrow \Lambda \\ R_2 &: a \rightarrow ad & R_4 &: b \rightarrow fbe \\ R_5 &: f \rightarrow ad \end{aligned}$$

$$\begin{array}{c} \Gamma_0 : \mathbf{a}_{(2)}\mathbf{b}_{(1)}\mathbf{c}_{(0)} \\ \downarrow R_2 \\ \Gamma_1 : \mathbf{a}_{(4)}\mathbf{d}_{(3)}\mathbf{b}_{(1)}\mathbf{c}_{(0)} \\ \downarrow R_1 \\ \Gamma_2 : \mathbf{d}_{(3)}\mathbf{b}_{(1)}\mathbf{c}_{(0)} \\ \downarrow R_3 \\ \Gamma_3 : \mathbf{b}_{(1)}\mathbf{c}_{(0)} \\ \downarrow R_4 \\ \Gamma_4 : \mathbf{f}_{(7)}\mathbf{b}_{(6)}\mathbf{e}_{(5)}\mathbf{c}_{(0)} \\ \downarrow R_5 \\ \Gamma_5 : \mathbf{a}_{(9)}\mathbf{d}_{(8)}\mathbf{b}_{(6)}\mathbf{e}_{(5)}\mathbf{c}_{(0)} \end{array}$$

$\Gamma_0 \preceq \Gamma_1$  и  $\Gamma_1 \preceq \Gamma_5$ , но  $\Gamma_0 \not\preceq \Gamma_5$ .

Рассмотрим все порожденные последовательности  $\{\Phi_i\}_{i=1}^{\infty}$ , такие, что  $\exists N \forall i \exists j (i < j \ \& \ |\Phi_j| \leq N)$ .

Для каждой такой последовательности выберем наименьшее  $N$ , удовлетворяющее этому свойству. Существует конечное число классов эквивалентности  $Q'$  слов длины не больше  $N$ , таких, что  $\forall i, j (|\Phi_i| \leq N \Rightarrow (\Phi_i \in Q' \ \& \ \Phi_j \in Q' \Leftrightarrow \Phi_i \approx \Phi_j))$ . Хотя бы один из них содержит бесконечное число элементов — обозначим его  $\{\Phi'_i\}$ . Если для некоторых  $i, j, k$   $\Phi'_i[k] \neq \Phi'_j[k]$ , то моменты, в которые  $\Phi'_i[k]$  и  $\Phi'_j[k]$  были порождены — результаты применения неукорачивающих правил. Значит, какое-то из них применяется бесконечное число раз.

Все прочие порожденные последовательности удовлетворяют условию  $\forall N \exists i_N \forall j (j > i_N \Rightarrow |\Phi_j| > N)$ .

Для каждого  $N$  выберем такое первое  $i_N$ , что  $\forall j (j < i_N \Rightarrow \exists k (k \geq j \ \& \ |\Phi_k| < N))$ . Итак,  $|\Phi_{i_N-1}| < N$  и  $|\Phi_{i_N}| \geq N$ , и  $\Phi_{i_N}$  порождено из предыдущего слова применением неукорачивающего правила  $R$ ,  $|R| \geq 2$ :  $\Phi_{i_N} = R(\Phi_{i_N-1}[1])\Phi_{i_N-1}^-$ .  $\Phi_{i_N-1}^-$  неизменно, поскольку  $|\Phi_{i_N-1}| < N$ . Все элементы последовательности  $\{\Phi_{i_N}\}_{i=1}^\infty$  начинаются с правой части некоторого неукорачивающего правила, значит, существует бесконечная подпоследовательность  $\{\Phi_{i_N}\}_{i=1}^\infty$ , такая, что все ее элементы начинаются с правой части одного и того же правила.

## Следствие

Пусть  $R$  — квазипорядок, полный на произвольных последовательностях.  $R \cap \preceq$  — почти полное отношение.

Все турчинские пары в  $T$  являются результатами применения одного и того же неукорачивающего правила без его отмены. Поэтому турчинское отношение можно сузить до отношения  $\preceq' = \{(\Gamma, \Delta) \mid \Gamma = R(a)\Theta_0 \ \& \ \Delta = R(b)\Psi\Theta_0\}$  без потери почти полноты.

Рассмотрим правило  $a_1a_2 \dots a_n \rightarrow b_1b_2 \dots b_m$  неалфавитной СПП. Его действие эквивалентно композиции действий:

$$a_1 \rightarrow \Lambda$$

...

$$a_{n-1} \rightarrow \Lambda$$

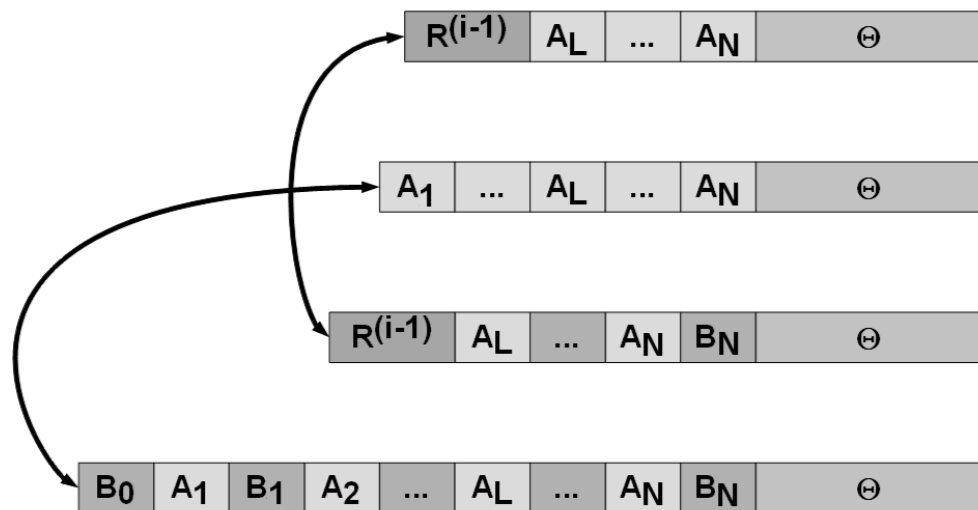
$$a_n \rightarrow b_1b_2 \dots b_m.$$

Отношение  $\preceq'$  почти полно на последовательностях, порожденных АСПП. Значит,  $\preceq$  почти полно на последовательностях, порожденных произвольными СПП с конечным набором правил в конечном алфавите.

Первая хигманова пара в последовательности, порожденной  $\mathfrak{G}$ -системой, является турчинской парой.

Рассмотрим такие  $\Phi_1$  и  $\Phi_2$ , что  $\Phi_1 \trianglelefteq \Phi_2$ , причем последовательность вплоть до  $\Phi_2$  представляет собой плохую последовательность в смысле  $\trianglelefteq$ .

$\Phi_1 = A_1 A_2 \dots A_n \Theta$  и  $\Phi_2 = B_0 A'_1 B_2 A'_2 \dots B_{n-1} A'_n B_n \Theta$ , и  $\forall i (i \geq 1 \ \& \ i \leq n \Rightarrow A_i \approx A'_i)$ . Вернемся к моментам порождения  $A_l[1]$  и  $A'_l[1]$  правилом  $R : x \rightarrow R_r$  (пусть  $A_l[1] \approx R_r[i]$ ).



В эти моменты имеем  $R_{(k_1)}^{(i-1)} A_l A_{l+1} \dots A_n \Theta$  и  $R_{(k_2)}^{(i-1)} A'_l B_l \dots B_{n-1} A_n B_n \Theta$  ( $R_{k_j}^{(i-1)}$  обозначает префикс  $R_r[1]_{(k_j+i-2)} \dots R_r[i-1]_{(k_j)}$ ). Они образуют хигманову пару, следовательно,  $l = 1 = n$ . Тогда  $\Phi_1 = A_1 \Theta$  и  $\Phi_2 = A'_1 B_n \Theta$ , и  $\Phi_1 \preceq \Phi_2$ .

## Выводы

1.  $\preceq$  является почти полным над последовательностями, порожденными конечными СПП над конечным алфавитом. Верхняя оценка длины плохой последовательности по  $\preceq$  экспоненциальна от размера исходной грамматики.
2. Плохие последовательности по  $\trianglelefteq$  и по  $\preceq$  совпадают на множестве последовательностей, порожденных обогащенными конечными СПП.

## Направление развития

Поиск аналога  $\preceq$  для деревьев функциональных вызовов.



**Спасибо за внимание!**