УРОВЕНЬ ГЛУБИННОГО СИНТАКСИСА И ВЫДЕЛЕНИЕ СВЕРХФРАЗОВЫХ ЕДИНСТВ В ТЕКСТАХ ПРИ УСТАНОВЛЕНИИ ИХ СЕМАНТИЧЕСКОЙ ЭКВИВАЛЕНТНОСТИ

Цель состоит в разработке и исследовании математической модели процесса распознавания и построения формальных семантических образов сверхфразовых единств в высказываниях на Естественном Языке (ЕЯ) при установлении их семантической эквивалентности.

Задачи:

- 1)Построение концептуальной модели процесса распознавания смысловой взаимной дополняемости фраз анализируемого высказывания;
- 2)Исследование алгоритмической разрешимости и сложности модели;
- 3)Разработка методов и алгоритмов для программной реализации полученной модели;
- 4)Исследование вопросов взаимодействия процессов построения образов сверхфразовых единств в анализируемом тексте и установления его эквивалентности смысловому эталону.

Системообразующие свойства правил расширенной лексико-синтаксической ∆-грамматики

Определение 1. Лексическая Синонимическая Конструкция (ЛСК) — заменяемый лексическим правилом комплекс лексических единиц и связывающих их отношений глубинного синтаксиса. Т₁ и Т₂ — деревья Глубинных Синтаксических

Структур (ГСС) фраз F_1 и F_2 :

$$T_1 = Cig(T_1^0; lpha_0^1 \mid Cig(t_1; lpha_1^1, lpha_2^1, ..., lpha_k^1 \mid T_1^1, T_1^2, ..., T_1^kig)ig)$$
 $T_2 = Cig(T_2^0; lpha_0^2 \mid Cig(t_2; lpha_1^2, lpha_2^2, ..., lpha_l^2 \mid T_2^1, T_2^2, ..., T_2^lig)ig)$,где
 $t_1 = Cig(t_1^0; lpha_0^w \mid Cig(t_1^w(C_0); lpha_1^w, lpha_2^w, ..., lpha_n^w \mid t_1^1, t_1^2, ..., t_1^nig)ig)$
 $t_2 = Cig(t_2^0; eta_0^w \mid Cig(t_2^w(C_0); eta_1^w, eta_2^w, ..., eta_n^w \mid t_2^1, t_2^2, ..., t_2^mig)ig)$

Синтаксическая замена : $t_1 \Rightarrow t_2$;

Лексическая замена : $t_1^w \Rightarrow t_2^w$; (*) С-операция композиции;

 t_1^w и t_2^w - деревья ЛСК, C_0 - ключевое слово ЛСК. Определение. 2. Деревья ГСС ϕ 1 и ϕ 2 удовлетворяют необходимому, но не достаточному условию Л ϕ -синонимии, если их ЛСК относится к

одному и тому же ключевому слову C_0 . Определение 3. Отвечающие необходимому условию ЛФ-синонимии деревья T1 и T2 удовлетворяют необходимому (но не достаточному!) условию взаимной дополняемости (смысловой), если существует последовательность преобразований (*), приводящих T1 и T2 к виду с одинаковой ЛСК.

Функциональные требования к модели

- Анализ применимости каждого преобразования из заданного множества к каждому дереву в смысловом описании высказывания с представлением результата в виде списка : $((npaвилo_i \quad C_0(i))...(npaвилo_k \quad C_0(k)));$
- Построение последовательности преобразований для приведения удовлетворяющих необходимому условию ЛФсинонимии деревьев ГСС к виду с одинаковой ЛСК;
- Определение наличия взаимной дополняемости* деревьев ГСС;
- Суммирование взаимно дополняющих друг друга деревьев с последующей их заменой в смысловом описании анализируемого высказывания на ГСС их суммы.
- *Определение. Приведенные к виду с одинаковой ЛСК деревья глубинного синтаксиса T1 и T2 взаимно дополняют друг друга, если они изоморфны таким образом, что для всякого узла α дерева T1 его образ $f(\alpha)$ в дереве T2:
- либо содержит информацию об одной и той же ненулевой характеризованной обобщенной лексеме;
- либо представляет фиктивную лексему с теми же семантическими словоизменительными характеристиками, что и ненулевая характеризованная обобщенная лексема в узле α;

Формальная концептуальная модель процесса распознавания смысловой взаимной дополняемости фраз анализируемого высказывания

 Y_S - язык смыслов заданного ЕЯ Y: $Y_S = \langle L_S, \Gamma_S, \Pi, Q, U \rangle$, где

 L_{S} – лексика языка Y_{S} ,

 $\Gamma_{\rm S}$ - синтаксис языка $Y_{\rm S}$,

 П – процедура установления соответствий между фразами языков Y и Y_S,

Q – процедура установления эквивалентности в Y_S , U – процедура преобразования текста на основе учета семантических повторов.

$$U = < Q_U, S_U >$$
, где

 Q_U – процедура приведения связанных по смыслу фраз в Y_S к целевому* представлению, S_U –процедура построения суммарного смысла в языке Y_S .

* Под целевым представлением фраз в языке Y_S понимается представление, допускающее суммирование

Представление системы правил ∆-грамматики ограниченной сетью Петри

Множество $\{T^\pi\}$ входов и выходов правил $\pi \in (\Pi^R \setminus \Pi_U^R)$ произвольных элементарных преобразований в Γ_L^R составляет множество информационных элементов. $\{T^\pi\} = \{T_1^\pi\} \cup \{T_2^\pi\}$, где $\{T_1^\pi\}$ -множество "входов", $\{T_2^\pi\}$ -множество "выходов" правил.

Сеть N_i , моделирующая работу *i*-й системы правил, где $i \in 1$, ..., n_sys:

$$N_i = \{P_i, T_i, F, H, M_{0i}\}$$
, где

множество позиций P_i есть подмножество $\{T^\pi\}$ элементов, составляющих входы и выходы способных образовывать систему правил;

Множество T_i переходов сети составляет множество переходов между состояниями системы: $\pi(r_{12}): T_1^\pi \xrightarrow{\pi(r_{12})} T_2^\pi$

Компонент R в описании правила $\pi \in (\Pi^R \setminus \Pi_U^R)$ отвечает за его применимость и представляется логической функцией :

$$r_j = x^1 \wedge x^2 \wedge \ldots \wedge x^n$$

Правило π может быть применено к дереву T_1^π , если выполняется одно из условий $r_j \in R$: $\vee_{j=1}^m r_j = true$.

Задача приведения деревьев глубинного синтаксиса к виду с одинаковой ЛСК

Применение правила $\pi \in (\Pi^R \setminus \Pi_U^R)$ сводится к выполнению перехода :

$$\pi(r_{12}): T_1^{\pi} \xrightarrow{\pi(r_{12})} T_2^{\pi}$$
, где $r_j = x^1 \wedge x^2 \wedge \ldots \wedge x^n$.

Последовательность применяемых правил моделируется последовательностью $\tau = \left(t_i^1, t_i^2, \dots, t_i^k\right)$ срабатываний переходов :

$$T_1^{\pi} \xrightarrow{\pi_1(r_{12})} T_2^{\pi} \xrightarrow{\pi_2(r_{23})} T_3^{\pi} \to \dots \to T_k^{\pi} \xrightarrow{\pi_k(r_{k+1})} T_{k+1}^{\pi},$$

$$\text{roe } t_i^1 \Leftrightarrow \pi_1(r_{12}), t_i^2 \Leftrightarrow \pi_2(r_{23}), \dots t_i^k \Leftrightarrow \pi_k(r_{k+1})$$

приводящей к последовательной смене разметок :

$$M_{0i} \xrightarrow{t_i^1} M_i^1 \xrightarrow{t_i^2} M_i^2 \dots \xrightarrow{t_i^k} M_i^k$$
 где $M_{0i} \Leftrightarrow T_1^\pi$, $M_i^1 \Leftrightarrow T_2^\pi$, ..., $M_i^k \Leftrightarrow T_{k+1}^\pi$.

Задача приведения деревьев T_1^π и T_{k+1}^π к виду с одинаковой ЛСК включает три задачи :

- Определение достижимости M_i^k из M_{0i} есть определение наличия слова $\tau \in T_i^* | M_{0i} \xrightarrow{\tau} M_i^k$, где T_i^* множество всех слов в алфавите T_i ;
- Задача обратимости слова \mathcal{T} : если $au \in T_i^* \big| M_{0i} \overset{ au}{\longrightarrow} M_i^k$, то существует ли слово $au' = \left(t_i^{k'}, t_i^{(k-1)'}, \ldots, t_i^{2'}, t_i^{1'}\right)$: $M_{0i} \overset{t_i^{l'}}{\longleftarrow} M_i^1 \overset{t_i^{2'}}{\longleftarrow} M_i^2 \ldots M_i^{k-1} \overset{t_i^{k'}}{\longleftarrow} M_i^k$ гле $M_{0i} \Leftrightarrow T_1^{\pi}$, $M_i^1 \Leftrightarrow T_2^{\pi}$, $M_i^k \Leftrightarrow T_{k+1}^{\pi}$:
- Задача определения оптимального слова $\tau \in T_i^*|M_{0i} \xrightarrow{\tau} M_i^k \text{ . Если } \exists \quad \tau_1, \tau_2, ..., \tau_l \ : \ M_{0i} \xrightarrow{\tau_1} M_i^k \text{ , } M_{0i} \xrightarrow{\tau_2} M_i^k \text{ и } M_{0i} \xrightarrow{\tau_3} M_i^k \text{ , то берется обратимое слово минимальной длины.}$

Свойства языка сети, моделирующей систему правил ∆-грамматики

Лемма 1. Все правила $\pi \in (\Pi^R \setminus \Pi_U^R)$ в расширенной лексико-синтаксической грамматике Γ_L^R различны. **Лемма 2.** Проблема достижимости заданной разметки M_i^k из начальной M_{0i} в сети N_i разрешима.

Теорема 2. Все символы-переходы $t_i^j \in T_i$ сети N_i различны.

Теорема 3. Проблема определения обратимости слова $\tau \in T_i^* | M_{0i} \xrightarrow{\tau} M_i^k$ языка $L(N_i)$ разрешима. Теорема 4. Проблема определения оптимального слова $\tau \in T_i^* | M_{0i} \xrightarrow{\tau} M_i^k$ (T_i^* - множество всех слов в алфавите T_i) в языке $L(N_i)$ сети N_i является разрешимой.

Исчисление сценариев на информационном пространстве системы правил Δ-грамматики

$$S_{i}^{j} \in S_{i} : S_{i}^{j} = \left\{T_{k}^{\pi}, T_{l}^{\pi}\right\}$$
, где $T_{k}^{\pi} \in \left\{T^{\pi}\right\}$, $T_{l}^{\pi} \in \left\{T^{\pi}\right\}$.
$$S_{i}^{j} = \left\{ref_{i}^{j}\left(k_{i}^{j}\right), P_{i}^{j}\right\}$$
 $ref_{i}^{j}\left(k_{i}^{j}\right) = \left\{S_{i}^{j1}, \ldots, S_{i}^{jk}\right\}$ - множество сценариев, связанных с S_{i}^{j} через разрешенные в его рамках переходы $t_{i}^{j} \in T_{i}$, $P_{i}^{j} \in \left\{p_{j}^{1}, p_{j}^{2}\right\}$ - множество позиций сети N_{i} , активных в рамках S_{i}^{j} . Целевое состояние системы $S_{i}^{j} = \left\{T_{k}^{\pi}, T_{l}^{\pi}\right\}$: $T_{k}^{\pi} = T_{l}^{\pi} \Leftrightarrow p_{i}^{j} : M_{i}\left(p_{i}^{j}\right) = 2$

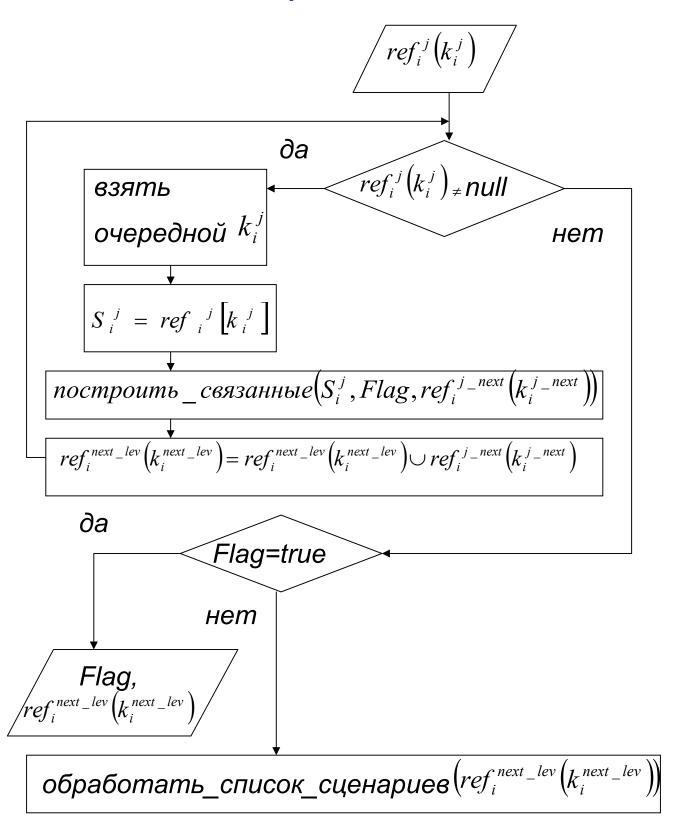
Теорема 5. Для каждого задаваемого над сетью N_i сценария $S_i^j \in S_i$ можно указать максимум два перехода $t_i^j \in T_i$ и $t_i^k \in T_i$: $t_i^j \neq t_i^k$ и $\exists S_i^{j1} \in S_i$. $S_i^{j2} \in S_i$. $S_i^{j1} \neq S_i^{j2} : S_i^k \xrightarrow{t_i^j} S_i^{j1} : S_i^{j1} \xrightarrow{t_i^k} S_i^j : S_i^j \xrightarrow{t_i^k} S_i^{j2} : S_i^{j2} \xrightarrow{t_i^j} S_i^j$ Для формирования пути к найденному решению в сценарий, с которым S_i^J связан посредством некоторого перехода $t_i^j \in T_i$:

$$S_{i}^{j} = \left\{ ref_back_{i}^{j}, P_{i}^{j} \right\}$$

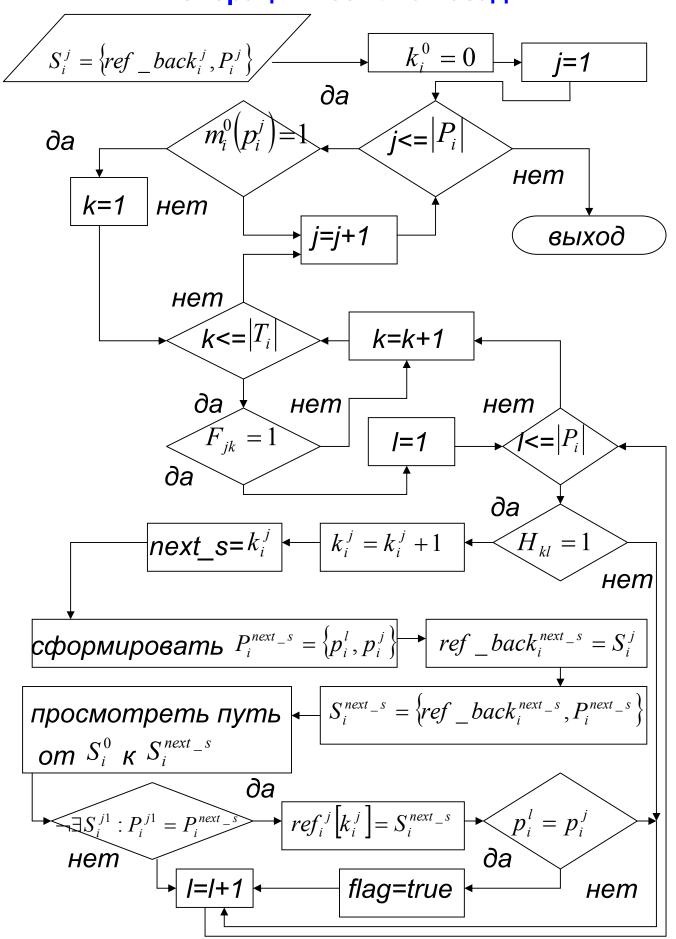
Генерация $S_i^j \in S_i$ происходит путем обработки матрицы F с использованием массива ссылок на описания входов/выходов правил $\Sigma_{dbfi} = P_i = \left\{ p_i^1, p_i^2, ..., p_i^{|P_i|} \right\}$ и описания условий

применимости
$$\Sigma_{Ri} = \{t_i^1, t_i^2, ..., t_i^{|T_i|}\}$$

Построение системы целевых выводов на информационном пространстве системы правил ∆-грамматики



Построение множества сценариев с применением операции "ссылка назад"



Функциональное описание информационного наполнения дерева глубинного синтаксиса

Определение. W/V-дерево есть помеченное дерево с пометками в узлах из множества W и с пометками на ветвях из множества V.

Информационное наполнение узла $w_\chi \in W_\chi^{W/V}$: $w_\chi = (lex_in_\chi, gram_in_\chi, arrow_label., composition_label)$, где $W_\chi^{W/V}$ - множество узлов дерева ГСС $T_\chi^{W/V}$ фразы χ ; lex_in_χ - лексическая часть w_χ ; $gram_in_\chi$ - грамматическая часть w_χ ; $arrow_label$ - пометка входящей в узел ветви; $composition_label$ - композиционная метка узла. $lex_in_\chi = (C_0, fun_n, ..., fun_1)$, где

 C_0 - ключевое слово ЛСК,

 $fun_n,...,fun_1$ - символы Лексических Функций (ЛФ). $gram_in_\chi = (part_of_speech, list_semant_categ)$, где $part_of_speech$ - символьный атом, обозначающий часть речи;

S – существительное Сопј – союз

V – глагол Num – числительное

А – прилагательное Р – причастие

Adv – наречие Prep – предлог

 $list_semant_categ$ - список семантически обусловленных словоизменительных категорий. Аналогично представляется информация узла $w_{\pi} \in W_{\pi}^{W/V}$ входного/выходного дерева $T_{\pi}^{W/V}$ правила π

Функционально-логическая модель входа/выхода правила ∆-грамматики

Порождаемые входом/выходом правила процессы моделируются сетью действий :

 P_{π} - множество состояний входа/выхода.

Каждому $t_{\pi}^i \in T_{\pi}$ соответствуют вычисления функций lex_in_{π} , $gram_in_{\pi}$, $arrow_label_{\pi}$ для очередного w_{π} .

Цветам маркера $color i \in C$ соответствуют способы использования информационного элемента; $Color i \in C$ соответствуют способы использования информационного элемента; $Color i \in C$ соответствуют способы $Color i \in C$ соответствуют способы использования $Color i \in C$ соответствуют способы $Color i \in C$ соответствуют

$$egin{align*} & |T_\pi^{'}||T_\pi^{'}|+1|P_\pi^{'}| \\ & \wedge & \vee & \wedge \\ j=1 & \underset{i
eq j}{\wedge} & k=1 \end{pmatrix} (f_\pi^{'}[k,j] \wedge t_\pi^{j} o h_\pi^{*'}[i,k] \wedge t_\pi^{i}) = true \ ,$$
 где

Алгоритмическая сложность суммирования деревьев глубинного синтаксиса

Определение. W/V деревья $t_1^{'}$ и t_1 считаются изоморфными с точностью до функционального соответствия, если между множествами их узлов существует взаимно-однозначное соответствие так, что в дереве $t_1^{'}$ из узла A' в узел B' идет ветвь с некоторой пометкой тогда и только тогда, когда в дереве t_1 из узла A в узел B идет ветвь с той же пометкой и узел А' удовлетворяет требованиям в узле А, и узел В' удовлетворяет требованиям в узле В. **Теорема 12.** Задача установления функционального соответствия W/V-деревьев $T_{\chi 1}^{W/V}$ и $T_{\chi 2}^{W/V}$ принадлежит классу Р комбинаторных задач с временной оценкой n^{D} , где $n=\max\Bigl(\!\!ig|W_{\chi 1}^{W/V}\Bigr|,\!\!ig|W_{\chi 2}^{W/V}\Bigr|\Bigr)$, $D=\sum_{i=1}^{k}arphi(a_{i})$, где $\varphi = \begin{pmatrix} a_1, a_2, \cdots, a_k \\ n_1, n_2, \cdots, n_k \end{pmatrix}$ - матрица ограничений на характер ветвления и на размещение пометок на ветвях из V, где $\{n_1, n_2, \cdots, n_k\} \subset N$ - подмножество натуральных

Теорема 13. Задача поиска для каждого дерева $T_{\chi l}^{W/V}$ из множества W/V-деревьев forest1 дерева $T_{\chi m}^{W/V}$ из множества forest2, функционально соответствующего $T_{\chi l}^{W/V}$, принадлежит классу P комбинаторных задач с временной оценкой n^D , где множества forest1 и forest2 представимы в виде лесов (forest1= $\left\langle V^{f1},E^{f1}\right\rangle$,

$$forest2=\langle V^{f2},E^{f2}\rangle$$
) $n=\max(V^{f1}|,|V^{f2}|)$, a $D=\sum_{i=1}^k \varphi(a_i)$.

чисел.

Служебная информация правил и относительность синонимических замен

Формат списка применимости правил к анализируемому дереву :

$$egin{pmatrix} \left(npaвило\ (i)\ count\ (i)\ C_0(i,count\ (i)))... \\ \left(npaвилo\ (k)\ count\ (k)\ C_0(k,count\ (k))) \right), \ \text{где} \end{cases}$$

count(i) и count(k) — счетчики вхождений в анализируемое дерево заменяемых i-м и k-м правилами поддеревьев;

 $C_0(i,count\ (i))$, $C_0(k,count\ (k))$ - ключевые слова ЛСК.

Описание поля композиционных меток узла ГСС:

$$ig(ig(composition_label\ (i,j,count\ (i))\ count\ (i)\ npaвuлo\ (i))... ig), \ (composition_label\ (k,l,count\ (k))\ count\ (k)\ npaвuлo\ (k)) ig),$$

где j,l — номера композиционных меток, соответственно, count(i) и count(k)-го вхождений заменяемых i-m и k-m правилами поддеревьев.

Формирование ЛФ-синонимических множеств

 $\Pi^{^R} \setminus \Pi^{^R}_U$ - множество правил произвольных элементарных преобразований Δ -грамматики :

$$\left(\Pi^{R}\setminus\Pi_{U}^{R}\right)=\Pi_{LS}^{R}\cup\Pi_{S}^{R}$$
 , $\Pi_{LS}^{R}\cap\Pi_{S}^{R}=\varnothing$, где

 Π_{LS}^{R} - множество лексико-синтаксических преобразований,

причем $\Pi_L^R \subset \Pi_{LS}^R$ - множество лексических преобразований без синтаксических замен;

 Π_S^R - множество синтаксических преобразований без вспомогательных лексических замен.

$$\Pi_{\mathit{LS}}^{\mathit{R}} = \Pi_{\mathit{LS_\mathit{JICK}}}^{\mathit{R}} \cup \Pi_{\mathit{LS_\mathit{CONV}}}^{\mathit{R}} \cup \Pi_{\mathit{LS_\mathit{IMP}}}^{\mathit{R}}$$
 , где

 $\Pi_{LS_{JICK}}^{R}$ - множество двусторонних правил с сохранением валентных мест ЛСК;

 $\Pi_{LS_CONV}^{R}$ - множество конверсивных замен с утратой валентности;

 $\Pi^{\it R}_{\it LS_\it{IMP}}$ - множество смысловых импликаций;

$$\Pi_{LS-JCK}^R \cap \Pi_{LS-CONV}^R \cap \Pi_{LS-IMP}^R = \emptyset$$

 $\left. \Phi_{_{i}} \right|_{i \; = \; 1}^{n}$ - последовательность ГСС :

$$\left\{ \! oldsymbol{\Phi}_{i}
ight. \left\} \! = \! \left\{ \! oldsymbol{\Phi}_{i}^{\mathcal{N}CK}
ight. \! \right\} \! \! \cdot \left\{ \! oldsymbol{\Phi}_{i}^{\mathcal{N}O_{-}\mathcal{I}CK}
ight. \! \right\}$$
, где

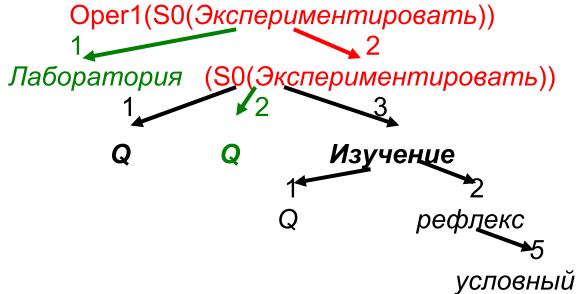
$$\left\{ \Phi_{i}^{1-JICK} \right\} \cap \left\{ \Phi_{i}^{1-NO-JICK} \right\} = \varnothing$$

Таким образом, $\{\Phi\}_i = \{\Phi\}_i^{\mathit{JCK}} \cup \{\Phi\}_i^{\mathit{CONv}} \cup \{\Phi\}_i^{\mathit{IMP}} \cup \{\Phi\}_i^{\mathit{SYNT}}$, причем $\{\Phi\}_i^{\mathit{JCK}} \cap \{\Phi\}_i^{\mathit{CONV}} \cap \{\Phi\}_i^{\mathit{IMP}} \cap \{\Phi\}_i^{\mathit{SYNT}} = \emptyset$

Пример построения образа суммарного смысла (этап 1, исходные деревья первого и второго предложений)

Первое предложение: Лаборатория провела эксперименты по изучению условных рефлексов.

ГСС первого предложения:



Список применимости: ((17 1 Экспериментировать)) Второе предложение: Подопытными животными были собаки.

ГСС второго предложения:



Список применимости: ((16 1 Экспериментировать))

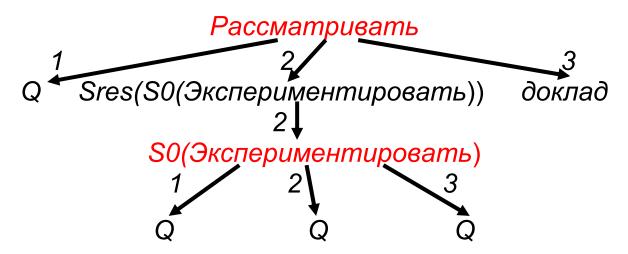
Обозначения:

Поддерево, заменяемое лексическим правиломПоддерево, заменяемое синтаксическим правилом

Пример построения образа суммарного смысла (этап 2, исходные деревья третьего и четвертого предложений)

Третье предложение: Результаты экспериментов рассматривались в докладе на конференции.

ГСС третьего предложения:

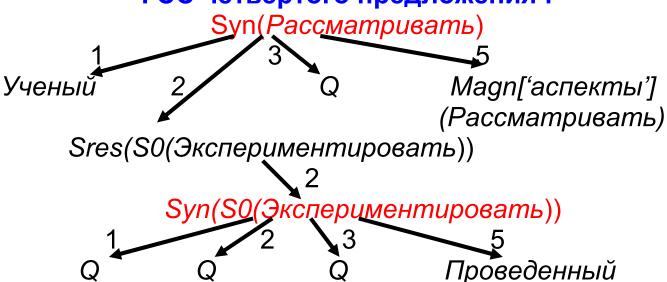


Список применимости:

((1 1 Рассматривать)(1 2 Экспериментировать))

Четвертое предложение: Ученый детально анализировал результаты проведенных опытов.

ГСС четвертого предложения:

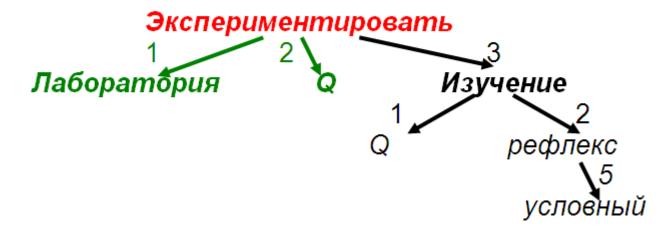


Список применимости:

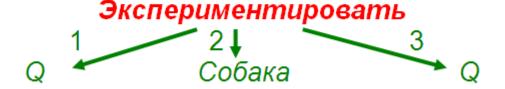
((1 1 Рассматривать)(1 2 Экспериментировать))

Пример построения образа суммарного смысла (этап 3, преобразованные деревья первого и второго предложений)

Преобразованная ГСС первого предложения:

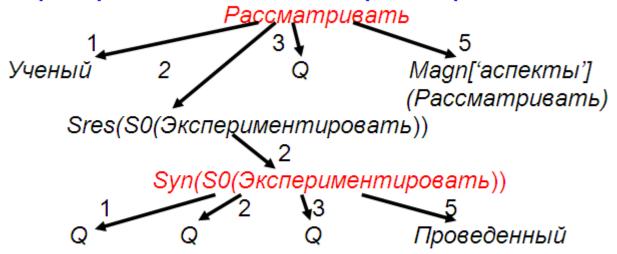


Преобразованная ГСС второго предложения:

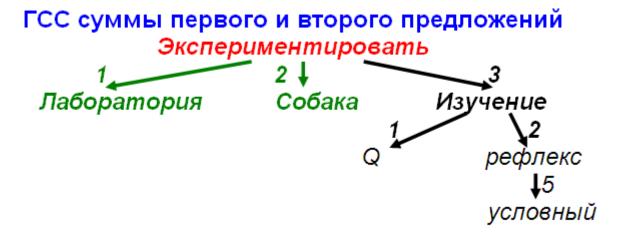


Пример построения образа суммарного смысла (этап 4, преобразованное дерево четвертого предложения и суммарная ГСС для первого и второго предложений)

Преобразованная ГСС четвертого предложения:



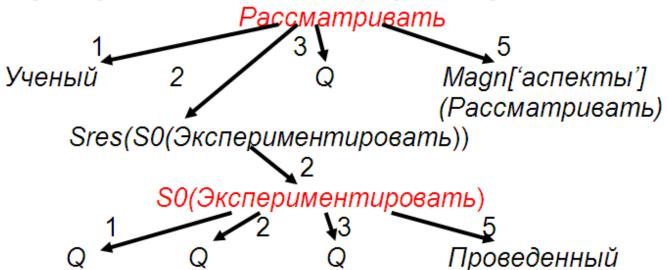
Список применимости: ((1 2 Экспериментировать))



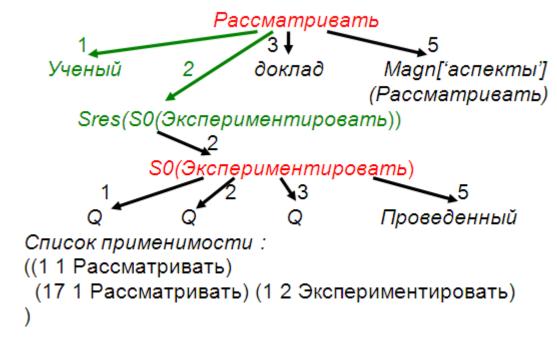
Список применимости: ((17 1 Экспериментировать))

Пример построения образа суммарного смысла (этап 5, преобразованное дерево четвертого предложения и суммарная ГСС для третьего и четвертого предложений)

Преобразованная ГСС четвертого предложения:



ГСС суммы третьего и четвёртого предложений



Преимущества подхода:

- 1. Выделение сверхфразовых единств без существенного ограничения жанра анализируемого текста;
- 2. Единые с задачей установления семантической эквивалентности механизмы оперирования лингвистическими знаниями.