

Взаимно однозначный англо-русский транслитератор

Пономарев А.С.

2023 год

1 Введение

Однажды автор этого текста задался вопросом, можно ли построить однозначную транслитерацию между русским и английским алфавитами. Иначе говоря, возможно ли взаимно однозначно сопоставить друг другу русские и английские слова, чтобы при этом получившееся соответствие было как можно ближе к привычному «транслиту», то есть к переписыванию побуквенно. Под словом, естественно, имеется в виду формальное слово, то есть любая конечная последовательность символов языка.

Ответ на поставленный вопрос – да, можно. Но насколько подобная транслитерация может соответствовать нашим желаниям и насколько она может быть простой – вопросы более сложные и к тому же совершенно расплывчатые. В этой работе предлагается определенная математическая конструкция и на ее основе строится одна из бесчисленного множества взаимно однозначных англо-русских транслитераций, которая, по мнению автора, и проста, и лингвистически удачна.

Зачем нужна такая транслитерация? В наш век высоких технологий большинство компьютерных систем оперируют английским языком. Подавляющее большинство языков программирования созданы на основе английского языка, базовый набор символов в компьютерах – латиница. Человек общается с машиной по-английски. При внедрении русского языка важно поддерживать совместимость имен: чтобы, например, программную функцию с английским названием можно было вызвать по-русски и наоборот. При транслитерации имен взаимная однозначность имеет ключевое значение.

Автор надеется, что работа послужит шагом на пути к созданию полноценного языка программирования на русском языке.

2 Алфавит. Язык

Для начала перечислим основные определения и условные обозначения.

Определение 1. *Алфавит* \mathcal{X} – непустое конечное множество; элементы алфавита называются *символами*, или *буквами*. В дальнейшем будем работать с двумя конкретными алфавитами: русским $\mathcal{R} = \{a, б, в, \dots, я\}$ и английским $\mathcal{A} = \{a, b, c, \dots, z\}$. $|\mathcal{X}|$ – размер алфавита, то есть количество символов в нем. $|\mathcal{R}| = 33$; $|\mathcal{A}| = 26$.

Определение 2. *Словом* над алфавитом \mathcal{X} называется конечная последовательность символов из \mathcal{X} : $\omega = \tilde{x}_1\tilde{x}_2\dots\tilde{x}_n$, $n \in \mathbb{N}_0$, $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n \in \mathcal{X}$. В частности, при $n = 0$ получается пустое слово ε . Слова будем обозначать строчными греческими буквами.

Определение 3. Пусть $\omega_1 = \tilde{x}_1\tilde{x}_2\dots\tilde{x}_n$ и $\omega_2 = \tilde{y}_1\tilde{y}_2\dots\tilde{y}_m$ – слова над алфавитом \mathcal{X} . *Конкатенацией* ω_1 и ω_2 называется слово $\omega = \tilde{x}_1\tilde{x}_2\dots\tilde{x}_n\tilde{y}_1\tilde{y}_2\dots\tilde{y}_m$, то есть полученное записью ω_1 и ω_2 «одного за другим». Пишут $\omega = \omega_1\omega_2$.

Определение 4. Пусть ω_1 и ω_2 – слова над алфавитом \mathcal{X} . Будем говорить, что ω_2 есть *подслово* ω_1 , если существуют слова α, β над \mathcal{X} , такие что $\omega_1 = \alpha\omega_2\beta$.

Определение 5. *Языком* над алфавитом \mathcal{X} называется произвольное множество слов над \mathcal{X} . Множество *всех* слов над алфавитом \mathcal{X} обозначим \mathcal{X}^* . Язык \mathcal{R}^* будем для краткости называть русским языком, \mathcal{A}^* – английским.

Замечание 1. Любой язык \mathcal{X}^* – счетно бесконечное множество.

Взаимно однозначная англо-русская транслитерация есть биективная функция $F : \mathcal{R}^* \longrightarrow \mathcal{A}^*$. Иначе говоря, нам нужно каждому русскому слову ρ поставить в соответствие некоторое английское слово $\alpha = F(\rho)$, так чтобы:

- 1) $\forall \rho_1, \rho_2 \in \mathcal{R}^* \quad (\rho_1 \neq \rho_2 \implies F(\rho_1) \neq F(\rho_2))$ – инъективность;
- 2) $\forall \alpha \in \mathcal{A}^* \quad \exists \rho \in \mathcal{R}^* : \alpha = F(\rho)$ – сюръективность.

3 Разбиение слова. Система правил

Пусть \mathcal{X} и \mathcal{Y} – произвольные алфавиты. Хотелось бы иметь инструмент для удобного описания отображений $F : \mathcal{X}^* \longrightarrow \mathcal{Y}^*$.

Определение 6. *Правил*ом из \mathcal{X}^* в \mathcal{Y}^* будем называть пару (ξ, η) , где ξ и η – непустые слова, $\xi \in \mathcal{X}^*$, $\eta \in \mathcal{Y}^*$. Можно записать более наглядно: $\xi \longrightarrow \eta$. Правило есть «предписание к переводу».

Далее логично рассмотреть множество правил и сделать так, чтобы оно некоторым образом порождало функцию $F : \mathcal{X}^* \longrightarrow \mathcal{Y}^*$.

Определение 7. Пусть ω – слово над алфавитом \mathcal{X} . *Разбиением* слова ω на n частей назовем конечную последовательность $T = (t_0, t_2, \dots, t_n)$, $n \in \mathbb{N}$, где $t_0 = 0$, $t_n = |\omega|$, $t_0 < t_1 < \dots < t_n$. Применяя разбиение к слову, получаем набор из n соответствующих подслов: пусть $\omega = \tilde{x}_1 \tilde{x}_2 \dots \tilde{x}_m$, $\tilde{x}_i \in \mathcal{X}$, $m \in \mathbb{N}$, тогда $T(\omega) = (\omega_1, \omega_2, \dots, \omega_n)$, где $\omega_i = \tilde{x}_{t_{i-1}+1} \tilde{x}_{t_{i-1}+2} \dots \tilde{x}_{t_i}$.

Пример 1. Пусть $\omega = \text{абвг}$, $T = (0, 1, 3, 4)$. Тогда $T(\omega) = (\text{а}, \text{бв}, \text{г})$.

Определение 8. Будем говорить, что из двух разбиений T_1 и T_2 слова ω T_1 больше T_2 , и писать $T_1 > T_2$, если $T_1 \subsetneq T_2$, то есть T_1 строго вложено в T_2 как множество. Если $T_1 > T_2$, то T_2 можно получить из T_1 , добавив один или несколько «разрезов» слова ω .

Замечание 2. Определенный выше порядок на множестве всех разбиений произвольного слова ω не является линейным: в общем случае существуют несравнимые разбиения. Например, $|\omega| = 3$, $T_1 = (0, 1, 3)$, $T_2 = (0, 2, 3)$.

Определение 9. Пусть \mathcal{T} – некоторое множество разбиений слова ω , $T_0 \in \mathcal{T}$. T_0 называется *максимальным элементом* множества \mathcal{T} , если $\forall T \in \mathcal{T} (T \not> T_0)$, то есть не существует элемента больше T_0 . T_0 называется *наибольшим элементом* множества \mathcal{T} , если $\forall T \in \mathcal{T} (T_0 > T)$, то есть если T_0 больше любого элемента.

Замечание 3. Из того, что элемент T_0 наибольший, следует, что он максимальный. Обратное в общем случае неверно, так как порядок на множестве \mathcal{T} не линейный.

Определение 10. Пусть $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$ – конечное или счетно бесконечное множество правил из \mathcal{X}^* в \mathcal{Y}^* , $\omega \in \mathcal{X}^*$. Разбиение T слова ω называется *допустимым* относительно множества правил Γ , если $\exists i_1, i_2, \dots, i_n \in \mathbb{N} : T(\omega) = (\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_n}), (\xi_{i_1}, \eta_{i_1}), \dots, (\xi_{i_n}, \eta_{i_n}) \in \Gamma$.

Определение 11. *Корректной системой правил*, или просто *системой правил*, из \mathcal{X}^* в \mathcal{Y}^* называется конечное или счетно бесконечное множество правил из \mathcal{X}^* в \mathcal{Y}^* $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$, такое что:

- 1) все левые части правил ξ_i различны;

2) для любого слова $\omega \in \mathcal{X}^*$ в множестве его допустимых относительно Γ разбиений существует наибольший элемент.

Замечание 4. Существование наибольшего элемента в множестве допустимых разбиений в том числе означает, что это множество непустое.

Пример 2. Пусть $\mathcal{X} = \{k, c\} \subset \mathcal{R}$, $\mathcal{Y} = \mathcal{A}$. Тогда $\Gamma = \{k \rightarrow k, c \rightarrow s, kc \rightarrow x\}$ – корректная система правил. Действительно, если в слово $\omega \in \mathcal{X}^*$ подслово kc входит m раз, то существует ровно 2^m допустимых разбиений (каждое kc можно оставить слитным либо разбить в (k, c)). Среди них, очевидно, есть наибольшее – то, где все подслова kc оставлены слитно.

Пример 3. Пусть $\mathcal{X} = \{k, c, t\} \subset \mathcal{R}$, $\mathcal{Y} = \mathcal{A}$. Тогда $\Gamma_1 = \{k \rightarrow k, c \rightarrow s\}$ и $\Gamma_2 = \{k \rightarrow k, c \rightarrow s, t \rightarrow t, kc \rightarrow x, ct \rightarrow sht\}$ не являются корректными системами правил. Множество допустимых разбиений слова t по Γ_1 пусто. Множество допустимых разбиений слова kct по Γ_2 равно $\{(0, 1, 2, 3), (0, 2, 3), (0, 1, 3)\}$ – оно не имеет наибольшего элемента.

Определение 12. Пусть $\omega \in \mathcal{X}^*$, Γ – система правил из \mathcal{X}^* в \mathcal{Y}^* . Существует наибольшее допустимое разбиение T слова ω по Γ : $T(\omega) = (\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_n})$, $(\xi_{i_1}, \eta_{i_1}), \dots, (\xi_{i_n}, \eta_{i_n}) \in \Gamma$. *Образом* ω по системе правил Γ называется слово $\Gamma(\omega) = \eta_{i_1}\eta_{i_2}\dots\eta_{i_n}$, то есть полученное конкатенацией правых частей правил.

Пример 4. Рассмотрим систему правил Γ из примера 2. $\Gamma(\varepsilon) = \varepsilon$ (единственное допустимое разбиение – (0)), $\Gamma(k) = k$, $\Gamma(kccsksc) = kxsxx$.

Замечание 5. Система правил Γ из \mathcal{X}^* в \mathcal{Y}^* определяет функцию $F : \mathcal{X}^* \rightarrow \mathcal{Y}^*$, $F(\omega) = \Gamma(\omega)$, притом $F(\varepsilon) = \varepsilon$.

Утверждение 1. Пусть $F : \mathcal{X}^* \rightarrow \mathcal{Y}^*$ – произвольная функция из \mathcal{X}^* в \mathcal{Y}^* , такая что $F(\varepsilon) = \varepsilon$. Тогда существует система правил Γ , определяющая F .

Доказательство. \mathcal{X}^* – счетно бесконечное множество, рассмотрим его индексацию $\mathcal{X}^* = \{\omega_i\}_{i \in \mathbb{N}}$, все ω_i различны. Пусть $\Gamma = \{(\omega_i, F(\omega_i))\}_{i \in \mathbb{N}}$. Γ – корректная система правил (для любого непустого $\omega \in \mathcal{X}^*$ найдется наибольшее допустимое разбиение $(0, |\omega|)$); она задает F . \square

Утверждение 2. Пусть Γ – корректная система правил из \mathcal{X}^* в \mathcal{Y}^* . Тогда $\forall \tilde{x} \in \mathcal{X} \exists \eta \in \mathcal{Y}^* : (\tilde{x}, \eta) \in \Gamma$.

Доказательство. Рассмотрим слово $\omega = \tilde{x}$ (из одной буквы). Существует наибольшее допустимое разбиение ω относительно Γ , но есть всего одно разбиение $\omega - T = (0, 1)$. Значит, оно допустимо $\implies \exists \eta \in \mathcal{Y}^* : (\tilde{x}, \eta) \in \Gamma$. \square

Теорема 1. Пусть $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$ – конечное или счетно бесконечное множество правил из \mathcal{X}^* в \mathcal{Y}^* . Тогда Γ – корректная система правил \iff выполнены условия:

- 1) все левые части правил ξ_i различны;
- 2) $\forall \tilde{x} \in \mathcal{X} \exists \eta \in \mathcal{Y}^* : (\tilde{x}, \eta) \in \Gamma$;
- 3) $\forall i, j \in \mathbf{N} \forall \alpha, \beta, \gamma \in \mathcal{X}^* (\alpha, \beta, \gamma \neq \varepsilon, \xi_i = \alpha\beta, \xi_j = \beta\gamma \implies \exists \eta \in \mathcal{Y}^* : (\alpha\beta\gamma, \eta) \in \Gamma)$.

Третий пункт означает, что если некоторые две левые части правил ξ_i и ξ_j «накладываются» концом одной на начало другой, то существует правило с «объемлющей» левой частью.

Доказательство. (\implies) Пункт 1 по определению; пункт 2 по утверждению 2. Установим пункт 3. Пусть $(\xi_i, \eta_i), (\xi_j, \eta_j) \in \Gamma$, $\alpha, \beta, \gamma \neq \varepsilon$ (непустые слова), $\xi_i = \alpha\beta$, $\xi_j = \beta\gamma$. Пусть $\omega = \alpha\beta\gamma$. $T_1 = (0, |\alpha| + |\beta|, |\alpha| + |\beta| + 1, \dots, |\alpha| + |\beta| + |\gamma|)$ и $T_2 = (0, 1, \dots, |\alpha|, |\alpha| + |\beta| + |\gamma|)$ есть допустимые разбиения слова $\omega = \alpha\beta\gamma$ относительно Γ , так как в Γ лежат правила $(\alpha\beta, \eta_i)$, $(\beta\gamma, \eta_j)$ и, согласно утверждению 2, правила с левыми частями из одной буквы для любой буквы из \mathcal{X} . По определению системы правил для ω существует наибольшее допустимое разбиение T , $T > T_1$, $T > T_2$. В терминах множеств $T \subset T_1$, $T \subset T_2 \implies T \subset T_1 \cap T_2 = (0, |\alpha| + |\beta| + |\gamma|) \implies T = (0, |\alpha| + |\beta| + |\gamma|)$. Значит, в Γ существует правило с левой частью $\alpha\beta\gamma$.

(\impliedby) От противного: пусть Γ не является корректной системой правил из \mathcal{X}^* в \mathcal{Y}^* – значит, $\exists \omega \in \mathcal{X}^*$, для которого не существует наибольшего допустимого разбиения. Пусть \mathcal{T} – множество всех допустимых разбиений слова ω . Оно не пусто (из пункта 1 следует, что $(0, 1, \dots, |\omega|) \in \mathcal{T}$), оно конечно, в нем нет наибольшего элемента \implies , как несложно установить, $\exists T, U \in \mathcal{T} : T \not\geq U, U \not\geq T$ и $\forall V \in \mathcal{T} (V \not\geq T \text{ и } V \not\geq U)$. Элементы T, U максимальны в \mathcal{T} и не сравнимы между собой. Пусть $T = (t_0, t_1, \dots, t_n)$, $U = (u_0, u_1, \dots, u_m)$. Пусть индекс k такой, что $\forall i \in \overline{0, k} (t_i = u_i)$, но $t_{k+1} \neq u_{k+1}$. Без ограничения общности $t_{k+1} < u_{k+1}$. Пусть индекс p такой, что $\forall i \in \overline{0, p-1} (t_i < u_{k+1})$, но $t_p \geq u_{k+1}$. Очевидно, $p > k + 1$. Если $t_p = u_{k+1}$, то разбиение $(t_0, t_1, \dots, t_k (= u_k), t_p (= u_{k+1}), \dots, t_n)$ допустимо и к тому же больше T – противоречие с выбором T . Значит, $t_p > u_{k+1}$. $\omega = \tilde{x}_1 \tilde{x}_2 \dots \tilde{x}_{|\omega|}$, $\tilde{x}_j \in \mathcal{X}$. $\alpha = \tilde{x}_{t_{k+1}} \dots \tilde{x}_{t_p-1}$; $\beta = \tilde{x}_{t_p-1+1} \dots \tilde{x}_{u_{k+1}}$; $\gamma = \tilde{x}_{u_{k+1}+1} \dots \tilde{x}_{t_p}$. $\alpha\beta$ – часть номер $k + 1$ при разбиении U ; $\beta\gamma$ – часть номер p при разбиении T . По пункту 2 существует правило

в левой частью $\alpha\beta\gamma$, следовательно, разбиение $(t_0, t_1, \dots, t_k(= u_k), t_p, \dots, t_n)$ допустимо. Оно больше T – противоречие с выбором T . \square

Замечание 6. Теорема описывает системы правил куда более наглядно, чем исходное определение. С ее помощью по явно записанному множеству правил можно легко сказать, является оно корректной системой правил или нет: достаточно проверить три условия на слова ξ_i .

4 Алгоритм получения образа по системе правил

Замысловатое определение корректной системы правил имело целью дать свободу при переводе слова над одним алфавитом в слово над другим.

Определение 13. Пусть ω – слово над алфавитом \mathcal{X} . *Отрезком* слова ω назовем пару (a, b) , $a, b \in \mathbb{N}_0$, $a \leq b \leq |\omega|$, которую для удобства будем записывать в виде $[a, b]$. Пусть $\omega = \tilde{x}_1\tilde{x}_2\dots\tilde{x}_n$, $\tilde{x}_i \in \mathcal{X}$. Подсловом отрезка $[a, b]$ назовем $\omega[a, b] = \tilde{x}_{a+1}\tilde{x}_{a+2}\dots\tilde{x}_b$. В частности, $\omega[a, a] = \varepsilon$.

Определение 14. Пусть $[a, b]$ и $[c, d]$ – отрезки слова ω . Будем говорить, что $[a, b]$ *больше* $[c, d]$, и писать $[a, b] > [c, d]$, если $a \leq c$, $d \leq b$, причем $a < c$ или $d < b$ (то есть $[a, b] \neq [c, d]$).

Замечание 7. Как и в случае с отношением $>$ на множестве разбиений слова ω , определенное выше отношение порядка не является линейным, то есть в общем случае существуют несравнимые отрезки. Аналогично определяются *наименьший* и *наибольший*, *максимальный* и *минимальный* элементы произвольного множества отрезков слова ω .

Определение 15. Пусть $[a, b]$ – отрезок слова ω ; $T = (t_0, t_1, \dots, t_n)$ – разбиение ω . Будем говорить, что $[a, b]$ *входит* в T , если $\exists i \in \mathbb{N} : t_{i-1} = a, t_i = b$.

Замечание 8. $T(\omega) = (\omega[t_0, t_1], \omega[t_1, t_2], \dots, \omega[t_{n-1}, t_n])$.

Утверждение 3. Пусть $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$ – система правил из \mathcal{X}^* в \mathcal{Y}^* ; ω – слово над \mathcal{X} . Пусть $\exists i, c, d \in \mathbb{N} : \omega[c, d] = \xi_i$. Тогда в наибольшее разбиение T слова ω относительно Γ входит некоторый отрезок $[a, b]$, такой что $[a, b] \geq [c, d]$.

Доказательство. По утверждению 2 $\forall \tilde{x} \in \mathcal{X} \exists \eta \in \mathcal{Y}^* : (\tilde{x}, \eta) \in \Gamma$. Следовательно, разбиение $U = (0, 1, 2, \dots, c, d, d+1, \dots, |\omega|)$ слова ω допустимо относительно Γ . $T > U$; в терминах множеств $T \subset U \implies \forall t \in \mathbb{N} (c < t < d \implies t \notin T) \implies \exists i \in \mathbb{N} : t_{i-1} \leq c, t_i \geq d$. Отрезок $[t_{i-1}, t_i]$ входит в T , $[t_{i-1}, t_i] \geq [c, d]$. \square

Утверждение 4. Пусть $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$ – система правил из \mathcal{X}^* в \mathcal{Y}^* ; ω – слово над \mathcal{X} . Пусть \mathcal{S} – множество всех отрезков $[c, d]$ слова ω , таких что $\exists i \in \mathbb{N} : \omega[c, d] = \xi_i$. Пусть $[a, b]$ – максимальный элемент \mathcal{S} . Тогда $[a, b]$ входит в наибольшее разбиение T слова ω относительно Γ .

Доказательство. По утверждению 3 существует отрезок $[a', b']$ слова ω , входящий в T , такой что $[a', b'] \geq [a, b]$. Но $[a', b'] \in \mathcal{S}$, а в \mathcal{S} не существует элемента, большего $[a, b]$, $\implies [a', b'] = [a, b]$, $[a, b]$ входит в T . \square

Утверждение 5. Пусть в условиях утверждения 4 $[a, b]$ – максимальный элемент \mathcal{S} , $(\xi_i, \eta_i) \in \Gamma$, $\omega[a, b] = \xi_i$. Обозначим $\alpha = \omega[0, a]$, $\beta = \omega[b, |\omega|]$; $\omega = \alpha \omega[a, b] \beta$. Тогда $\Gamma(\omega) = \Gamma(\alpha) \eta_i \Gamma(\beta)$.

Доказательство. Пусть $T = (t_0, t_1, \dots, t_n)$ – наибольшее допустимое разбиение ω относительно Γ , $T(\omega) = (\xi_{j_1}, \xi_{j_2}, \dots, \xi_{j_n})$. По утверждению 4 $[a, b]$ входит в Γ , то есть $\exists k \in \mathbb{N} : a = t_{k-1}, b = t_k, \xi_{j_k} = \xi_i, j_k = i$. Заметим, что $U = (t_0, t_1, \dots, t_{k-1})$ и $V = (0, t_{k+1} - t_k, \dots, t_n - t_k)$ – это наибольшие допустимые разбиения для α и β (иначе можно было бы построить разбиение для ω , большее T). Имеем $\Gamma(\omega) = \eta_{j_1} \dots \eta_{j_{k-1}} \eta_i \eta_{j_{k+1}} \dots \eta_{j_n}$, $\Gamma(\alpha) = \eta_{j_1} \dots \eta_{j_{k-1}}$, $\Gamma(\beta) = \eta_{j_{k+1}} \dots \eta_{j_n}$. \square

Алгоритм 1 (перевода по системе правил). Пусть Γ – система правил из языка \mathcal{X}^* в язык \mathcal{Y}^* ; $\omega \in \mathcal{X}^*$. Получим $\Gamma(\omega)$.

1. Если $\omega = \varepsilon$ (пустое), то результат найден: $\Gamma(\omega) = \varepsilon$.
2. \mathcal{S} – множество отрезков $[c, d]$ слова ω , таких что $\exists i \in \mathbb{N} : \omega[c, d] = \xi_i$, $(\xi_i, \eta_i) \in \Gamma$. Возьмем $[a, b]$ – максимальный элемент \mathcal{S} .
3. Рекурсивно переведем по системе правил Γ слова $\alpha = \omega[0, a]$ и $\beta = \omega[b, |\omega|]$.
4. Получим ответ конкатенацией: $\Gamma(\omega) = \Gamma(\alpha) \eta_i \Gamma(\beta)$.

В шаге 2 максимальный элемент $[a, b]$ найдется, так как \mathcal{S} непустое и конечное. Алгоритм верен в силу утверждения 5.

Замечание 9. Алгоритм 1 приведен в общем виде: не определен способ выбора $[a, b]$; максимальных отрезков может быть несколько. Можно, например, рассматривать в множестве \mathcal{S} только отрезки вида $[0, d]$, тогда $a = 0$, $\alpha = \varepsilon$, $b = \max_{[0, d] \in \mathcal{S}} d$. Такой выбор однозначен; рекурсивный перевод шага 3 в этом случае соответствует движению по слову ω слева направо.

5 Системы правил и биективные отображения

Утверждение 6 (общеизвестное). Пусть A, B – произвольные множества. Функция $F : A \rightarrow B$ биективна $\iff \exists F^{-1} : B \rightarrow A$, такая что $\forall a \in A (F^{-1}(F(a)) = a)$ и $\forall b \in B (F(F^{-1}(b)) = b)$. В этом случае F^{-1} называется обратной функцией.

Определение 16. Пусть $\Gamma = \{(\xi_1, \eta_1), (\xi_2, \eta_2), \dots\}$ – конечное или счетно бесконечное множество правил из \mathcal{X}^* в \mathcal{Y}^* . Обратным множеством правил назовем $\Gamma^{-1} = \{(\eta_1, \xi_1), (\eta_2, \xi_2), \dots\}$ правил из \mathcal{Y}^* в \mathcal{X}^* .

Хотелось бы описать биекцию $F : \mathcal{X}^* \rightarrow \mathcal{Y}^*$ системой правил Γ , так чтобы Γ^{-1} было корректной системой правил и задавало F^{-1} . На пути к этому есть проблемы.

Пример 5. Рассмотрим систему правил Γ из примера 2. Γ^{-1} не является системой правил из $\mathcal{Y}^* = \mathcal{A}^*$ в $\mathcal{X}^* = \{k, c\}^*$.

Пример 6. Рассмотрим Γ из примера 2, но заменим \mathcal{Y} на $\{k, s, x\}$. Теперь $\Gamma^{-1} = \{k \rightarrow k, s \rightarrow c, x \rightarrow ks\}$ – корректная система правил. Однако $F(\xi) = \Gamma(\xi)$, $\xi \in \mathcal{X}^*$, не биекция (инъекция, но не сюръекция: например, нельзя получить слово ks). Значит, $G(\eta) = \Gamma^{-1}(\eta)$, $\eta \in \mathcal{Y}^*$, не может быть обратной к F .

Пример 7. Рассмотрим случай, когда система правил Γ определяет биективную функцию $F : \mathcal{X}^* \rightarrow \mathcal{Y}^*$, но Γ^{-1} не является корректной системой правил. Пусть $\mathcal{X} = \{0, 1\}$, $\mathcal{Y} = \{y\}$. Зададим вспомогательную функцию G так: $\omega \in \mathcal{X}^*$, тогда 1ω – двоичная запись некоторого числа $n \in \mathbf{N}$; пусть $G(\omega)$ будет словом из $n - 1$ буквы y . $G : \mathcal{X}^* \rightarrow \mathcal{Y}^*$ – биекция. Пусть $F(01) = G(00) = ууу$, $F(00) = G(01) = уууу$, а на всех остальных аргументах F совпадает с G . Таким образом, F биекция. Ее определяет система правил $\Gamma = \{0 \rightarrow y, 1 \rightarrow уу, 00 \rightarrow уууу, 10 \rightarrow ууууу, 11 \rightarrow уууууу\} \cup \{(\omega \rightarrow F(\omega)) : |\omega| \geq 3\}$ (содержит все ω , кроме 01 , образ которой есть конкатенация правых частей правил для 0 и 1). Но легко видеть, что Γ^{-1} при этом не является системой правил.

Теорема 2. Пусть Γ – система правил из \mathcal{X}^* в \mathcal{Y}^* , причем Γ^{-1} – система правил из \mathcal{Y}^* в \mathcal{X}^* . Пусть $F(\xi) = \Gamma(\xi)$, $\xi \in \mathcal{X}^*$. Тогда F биекция $\iff \Gamma^{-1}$ задает обратную функцию F^{-1} .

Доказательство. (\Leftarrow) Очевидно по утверждению 6.

(\Rightarrow) От противного: пусть $\exists \eta \in \mathcal{Y}^* : F^{-1}(\eta) \neq \Gamma^{-1}(\eta)$. Обозначим

$\xi = F^{-1}(\eta)$, $\omega = \Gamma^{-1}(\eta)$. $F(\xi) = \Gamma(\xi) = \eta$: по определению системы правил среди допустимых разбиений слова ξ относительно Γ существует наибольшее T , $\exists i_1, i_2, \dots, i_n \in \mathbb{N} : T(\xi) = (\xi_{i_1}, \xi_{i_2}, \dots, \xi_{i_n}), (\xi_{i_1}, \eta_{i_1}), \dots, (\xi_{i_n}, \eta_{i_n}) \in \Gamma$; $\eta = \eta_{i_1} \eta_{i_2} \dots \eta_{i_n}$. Отсюда получаем допустимое разбиение U слова η относительно Γ^{-1} : $U = (0, |\eta_{i_1}|, |\eta_{i_1}| + |\eta_{i_2}|, \dots, |\eta|)$; $U(\eta) = (\eta_{i_1}, \eta_{i_2}, \dots, \eta_{i_n})$. Существует наибольшее допустимое разбиение V слова η относительно Γ^{-1} . $V \neq U$, так как в таком случае было бы $\Gamma^{-1}(\eta) = \xi$. Значит, $V > U$. Пусть $U = (u_0, u_1, \dots, u_n)$, $V = (v_0, v_1, \dots, v_m)$. В терминах множеств $V \subset U \implies \exists j_1, j_2, k \in \mathbb{N}_0 : u_{j_1} = v_{k-1}, u_{j_2} = v_k, j_2 > j_1 + 1$, то есть в V существуют два подряд идущих элемента v_{k-1}, v_k , равные двум элементам u_{j_1}, u_{j_2} из U , не идущим подряд. $\eta = \tilde{y}_1 \tilde{y}_2 \dots \tilde{y}_{|\eta|}$, $\tilde{y}_j \in \mathcal{Y}$; рассмотрим $\beta = \tilde{y}_{u_{j_1}+1} \tilde{y}_{u_{j_1}+2} \dots \tilde{y}_{u_{j_2}}$. $\xi = \tilde{x}_1 \tilde{x}_2 \dots \tilde{x}_{|\xi|}$, $\tilde{x}_j \in \mathcal{X}$; $T = (t_1, t_2, \dots, t_n)$; рассмотрим $\alpha = \tilde{x}_{t_{j_1}+1} \tilde{x}_{t_{j_1}+2} \dots \tilde{x}_{t_{j_2}}$. $T' = (0, t_{j_1+1} - t_{j_1}, t_{j_1+2} - t_{j_1}, \dots, t_{j_2} - t_{j_1})$ есть разбиение слова α , допустимое относительно Γ . Более того, оно наибольшее среди допустимых для α , так как в противном случае можно было бы построить допустимое разбиение слова ξ , большее T . $T'(\alpha) = (\xi_{i_{t_1}}, \xi_{i_{t_1+1}}, \dots, \xi_{i_{t_2}}) \implies \Gamma(\alpha) = \eta_{i_{t_1}} \eta_{i_{t_1+1}} \dots \eta_{i_{t_2}} = \beta$. Притом β – это часть номер k в разбиении V слова $\eta \implies \exists \gamma \in \mathcal{X}^* : (\beta, \gamma) \in \Gamma^{-1}$. Значит, $(\gamma, \beta) \in \Gamma$, $F(\gamma) = \Gamma(\gamma) = \beta$. $\alpha \neq \gamma$, так как в таком случае для α существовало бы допустимое разбиение $(0, |\alpha|)$, а по построению α в наибольшем допустимом разбиении T' хотя бы три числа ($j_2 > j_1 + 1$). Получается, что $\alpha \neq \gamma$, $F(\alpha) = F(\gamma) = \beta$ – противоречие с биективностью F . \square

6 Построение биективной англо-русской транслитерации

Необходимо построить биективную функцию $F : \mathcal{R}^* \longrightarrow \mathcal{A}^*$. Опишем ее системой правил Γ , такой что Γ^{-1} тоже система правил. Тогда по теореме 2 Γ^{-1} задаст обратную функцию F^{-1} .

Сократим рассматриваемые алфавиты: перейдем к $\mathcal{R}' \subset \mathcal{R}$, $\mathcal{A}' \subset \mathcal{A}$. Будем шаг за шагом добавлять новые буквы и пополнять Γ , следя за тем, чтобы на рассматриваемых алфавитах Γ и Γ^{-1} были системами правил и чтобы функция, определяемая Γ , обладала инъективностью (никакие $\omega_1 \neq \omega_2 \in \mathcal{R}^*$ не дают один и тот же образ) и сюръективностью (любое $\eta \in \mathcal{A}^*$ имеет прообраз). \mathcal{R}' и \mathcal{A}' по ходу рассуждений указывать явно не станем.

Выбор образов для тех или иных букв не имеет строгого логического обоснования. При их выборе автор руководствовался двумя принци-

нами. Во-первых, система правил должна быть простой, насколько это возможно, для удобства практического применения. Во-вторых, автор следовал собственной лингвистической интуиции, которую способен объяснить словами лишь отчасти.

Так как мы строим систему правил Γ , такую что Γ^{-1} тоже система правил, то для наглядности в записи правил вместо \longrightarrow будем использовать \longleftrightarrow .

Шаг 1 (очевидные правила буква-буква).

$a \longleftrightarrow a$	$и \longleftrightarrow i$	$o \longleftrightarrow o$	$\phi \longleftrightarrow f$
$б \longleftrightarrow b$	$к \longleftrightarrow k$	$п \longleftrightarrow p$	$x \longleftrightarrow h$
$в \longleftrightarrow v$	$л \longleftrightarrow l$	$р \longleftrightarrow r$	
$д \longleftrightarrow d$	$м \longleftrightarrow m$	$с \longleftrightarrow s$	
$з \longleftrightarrow z$	$н \longleftrightarrow n$	$т \longleftrightarrow t$	

Шаг 2 (правила буква-буква, требующие внимания).

Выбор следующих правил может быть не очевиден.

$г \longleftrightarrow g$	$ж \longleftrightarrow j$	$й \longleftrightarrow y$	$ц \longleftrightarrow c$
---------------------------	---------------------------	---------------------------	---------------------------

Шаг 3 (ч, ш, щ; бесконечная цепочка правил).

Рассмотрим следующий частный случай. Хочется добавить правило $ш \longleftrightarrow sh$. Но в этом случае Γ перестанет быть биекцией: $\Gamma(ш) = \Gamma(сх) = sh$. Мы только что «забрали» образ у любого русского слова, содержащего $сх$ как подслово – придется добавить для $сх$ новое правило. Пусть это будет $сх \longleftrightarrow skh$. Возникает та же проблема: $\Gamma(сх) = \Gamma(скх) = skh$. Ее можно снова решить тем же образом и так далее до бесконечности.

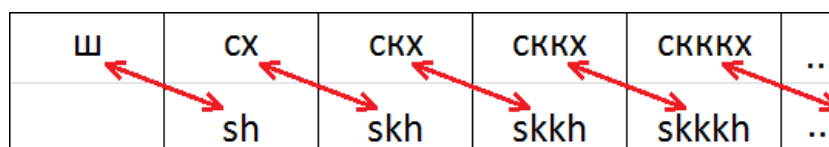


Рис. 1: цепочка исправлений для ш

Если проделать любое конечное число описанных исправлений, проблема, очевидно, останется. Но если переопределить образы для всей счетно бесконечной цепочки слов (рис. 1), она разрешится! Переопределение достигается добавлением счетно бесконечной цепочки правил: $ш \longleftrightarrow sh$, $сх \longleftrightarrow skh$, $скх \longleftrightarrow skkh$, $сккх \longleftrightarrow skkkh$,

Для лаконичного описания подобных цепочек введем следующий способ записи.

Определение 17. Введем символы $+$ и $*$. Пусть \mathcal{X} и \mathcal{Y} – алфавиты; $n \in \mathbb{N}$; $\alpha_0, \alpha_1, \dots, \alpha_n \in \mathcal{X}^*$, $\beta_0, \beta_1, \dots, \beta_n \in \mathcal{Y}^*$ – некоторые слова; $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n \in \mathcal{X}$; $\tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_n \in \mathcal{Y}$ – некоторые буквы; $a_1, a_2, \dots, a_n, b_1, b_2, \dots, b_n \in \{+, *\}$. Выражение $\alpha_0 \tilde{x}_1^{a_1} \alpha_1 \tilde{x}_2^{a_2} \dots \tilde{x}_n^{a_n} \alpha_n \longleftrightarrow \beta_0 \tilde{y}_1^{b_1} \beta_1 \tilde{y}_2^{b_2} \dots \tilde{y}_n^{b_n} \beta_n$ обозначает множество тех и только тех правил, которые имеют вид

$$\begin{aligned} \alpha_0 \tilde{x}_1 \dots \tilde{x}_1(i_1 \text{ раз}) \alpha_1 \tilde{x}_2 \dots \tilde{x}_2(i_2 \text{ раз}) \dots \tilde{x}_n \dots \tilde{x}_n(i_n \text{ раз}) \alpha_n &\longleftrightarrow \\ &\longleftrightarrow \beta_0 \tilde{y}_1 \dots \tilde{y}_1(j_1 \text{ раз}) \beta_1 \tilde{y}_2 \dots \tilde{y}_2(j_2 \text{ раз}) \dots \tilde{y}_n \dots \tilde{y}_n(j_n \text{ раз}) \beta_n, \end{aligned}$$

где числа i_k, j_k соотносятся внутри каждой пары следующим образом:

- 1) если $a_k = *$ и $b_k = *$, то $0 \leq i_k = j_k$;
- 2) если $a_k = *$ и $b_k = +$, то $0 \leq i_k = j_k - 1$;
- 3) если $a_k = +$ и $b_k = *$, то $0 \leq i_k - 1 = j_k$;
- 4) если $a_k = +$ и $b_k = +$, то $0 \leq i_k - 1 = j_k - 1$.

Описанный способ записи будем называть $(*, +)$ -нотацией.

Используя $(*, +)$ -нотацию, можно записать цепочку правил для ш как $ш \longleftrightarrow sh$, $ск^*х \longleftrightarrow sk^+h$. Средние символы (из обоих алфавитов) в таких цепочках будем в дальнейшем называть *разделителями*.

Аналогичную цепочку правил построим для ч: $ч \longleftrightarrow ch$, $цк^*х \longleftrightarrow ck^+h$.

Для щ хочется добавить правило $щ \longleftrightarrow shch$. Тогда, аналогично рассмотренному ранее случаю, образ «потеряют» все русские слова, содержащие шч как подслово (обратим внимание, что именно шч, а не, например, схцх, ведь мы уже ввели правила для ш и ч). Проблему решим по известному образцу, используя в качестве разделителя букву т: $шт^*ч \longleftrightarrow sht^+ch$.

Итак, были добавлены следующие правила:

$$\begin{array}{lll} ч \longleftrightarrow ch & ш \longleftrightarrow sh & щ \longleftrightarrow shch \\ цк^*х \longleftrightarrow ck^+h & ск^*х \longleftrightarrow sk^+h & шт^*ч \longleftrightarrow sht^+ch \end{array}$$

Шаг 4 (th и ph).

Пока что $\Gamma^{-1}(th) = tx$, что плохо отражает суть этого английского буквосочетания. Автору кажется хорошим решением ввести правило $th \longleftrightarrow th$.

Заметим, что теперь $\Gamma(tx) = \Gamma(th) = th$, а в th не переходит никакое русское слово: образ th «передали» th . Самым простым решением будет «отдать» th старый образ th , а именно добавить $th \longleftrightarrow th$.

Неудачен также перевод $\Gamma^{-1}(ph) = px$. Гораздо лучше $\Gamma^{-1}(ph) = pf$. Поступим аналогично!

Итак, добавлены следующие правила:

$$\text{тф} \longleftrightarrow \text{th} \quad \text{тх} \longleftrightarrow \text{tf} \quad \text{пф} \longleftrightarrow \text{ph} \quad \text{пх} \longleftrightarrow \text{pf}$$

Факт того, что решения двух, казалось бы, совершенно не связанных проблем образуют стройную систему, несомненно, свидетельствует в пользу избранной автором транслитерации.

Шаг 5 (у, э и йотированные гласные).

Добавим правила $у \longleftrightarrow оу$, $ю \longleftrightarrow и$. Такое сопоставление автору подсказывает лингвистическая интуиция. Перевод $и$ в $ю$ во многих английских словах смотрится гармонично, а английская буква $и$ даже называется «ю». Записывать $у$ как «оу» вообще в духе и русского, и английского языка: в древнерусском $у$ записывалась диграфом $оу$, в английском буквосочетание ou произносится как «у».

Введем правила $э \longleftrightarrow ое$, $е \longleftrightarrow е$. Соответствие $е \longleftrightarrow е$ крайне удобно, а правило для $э$ сконструируем аналогично правилу для $у$.

Для $ё$ и $я$ добавим правила $ё \longleftrightarrow уо$, $я \longleftrightarrow уа$.

Как и на шаге 3, чтобы сохранить биективность Γ , нужно добавить бесконечные цепочки правил с символами-разделителями. Для $у$ и $э$ разделителем пусть будет $й$, для $ё$ и $я$ – $и$. В результате имеем:

$$\begin{array}{llll} ю \longleftrightarrow и & & е \longleftrightarrow е & \\ у \longleftrightarrow оу & э \longleftrightarrow ое & ё \longleftrightarrow уо & я \longleftrightarrow уа \\ ой^*ю \longleftrightarrow оу^+и & ой^*е \longleftrightarrow оу^+е & йи^*о \longleftrightarrow уи^+о & йи^*а \longleftrightarrow уи^+а \end{array}$$

Замечание 10. На текущий момент Γ не является системой правил: например, есть правила с левыми частями $йо$ и $ою$, но нет – с $йою$ (не выполнено условие критерия 1). Будем говорить, что цепочки для $ё$ и $у$ *конфликтуют*. Эту ситуацию разрешим позже.

Шаг 6 (ъ, ы, ь).

Собрав в кулак имеющуюся лингвистическую интуицию и представления о практичности и эстетике, автор сочинил следующие правила: $ъ \longleftrightarrow оа$, $ы \longleftrightarrow еа$, $ь \longleftrightarrow ie$.

Добавив цепочки с символами-разделителями $й$ и $х$ (перед $а$), получаем:

$$\begin{array}{lll} ъ \longleftrightarrow оа & ы \longleftrightarrow еа & ь \longleftrightarrow ie \\ ох^*а \longleftrightarrow ох^+а & ех^*а \longleftrightarrow ех^+а & ий^*е \longleftrightarrow ий^+е \end{array}$$

Шаг 7 (разрешение конфликтов гласных: $ё$ и $у$).

Рассмотрим бесконечные цепочки правил $ё \longleftrightarrow уо$, $йи^*о \longleftrightarrow уи^+о$ и $у \longleftrightarrow оу$, $ой^*ю \longleftrightarrow оу^+и$.

ё	йо	йио	йиио	йииио	...
	уо	уіо	уііо	уіііо	...

у	ою	ойю	оййю	ойййю	...
	оу	оуу	оууу	оуууу	...

Рис. 2: цепочки правил для ё и у

Как видно на рис. 2, цепочки правил получаются при «сдвиге» соответствий слов на один относительно порожденных побуквенно. В начало добавляется правило для нового, ранее не встречавшегося, символа. Этот символ будем называть *головой* цепочки.

Если взять любое (кроме головы) слово из первой цепочки и любое (кроме головы) слово на том же языке из второй цепочки, то их можно «наложить» друг на друга одним символом (например, йио и ою, уо и оу). Значит, ни Γ , ни Γ^{-1} не станут системами правил, пока не будут добавлены правила для всех «объемлющих» слов (таких как, например, йиою и оуу). Это и есть *конфликт*. В рассматриваемом случае обе головы русские, принадлежат к одному языку – возникающий конфликт назовем *односторонним*.

Итак, для всех слов, в которые слова из цепочек входят как подслова «без наложений», образы уже корректно определены. Выпишем все возможные вхождения «с наложениями» в бесконечную таблицу (рис. 3). Номер клетки по горизонтали – количество первого по порядку символа-разделителя, по вертикали – второго.

	йу	йиу	...
ёю	йою	йиою	...
ёйю	йойю	йиойю	...
⋮	⋮	⋮	

Рис. 3: таблица конфликтующих «наложений»

йу	йиу	йиуу	...
ёю	йою	йиою	...
ёйю	йойю	йиойю	...
⋮	⋮	⋮	

Рис. 4: разрешение конфликта

Русские слова заполняют всю таблицу, кроме верхней левой клетки, английские – кроме первого ряда и первого столбца. Нам нужно сгруппировать все слова по парам. Поступим следующим образом: сдвинем верхний ряд на одну клетку влево и каждому русскому слову поставим в пару английское, стоящее от него снизу справа (рис. 4). Все полученные правила добавим в Γ . Таким образом, вся счетная бесконечность конфликтующих наложений разрешается.

Естественная структура множества добавленных правил не одномерная (цепочка), а двумерная. Его можно представить в $(*,+)$ -нотации, однако такая запись мало помогает пониманию:

$$\begin{aligned} \text{йй}^* \text{у} &\longleftrightarrow \text{уі}^* \text{оу} \\ \text{ёй}^* \text{ю} &\longleftrightarrow \text{уоу}^+ \text{у} \\ \text{йй}^* \text{ой}^* \text{ю} &\longleftrightarrow \text{уі}^+ \text{оу}^+ \text{у} \end{aligned}$$

Гораздо легче осмыслить его в виде следующей закономерности: в описанной ситуации перед использованием правила с головным символом разделитель не требуется. Примеры:

- 1) $\Gamma(\text{йу}) = \text{уоу}$ (что не может не радовать), $\Gamma(\text{йиу}) = \text{уіоу}$, $\Gamma(\text{йиииу}) = \text{уіііоу}$: использовано правило с головой у, перед ним разделитель не требуется, а именно, букв и столько же, сколько і;
- 2) $\Gamma(\text{ёу}) = \text{уоуу}$, $\Gamma(\text{ёйю}) = \text{уоууу}$, $\Gamma(\text{ёйййю}) = \text{уоууууу}$: использовано правило с головой ё, но после разделитель требуется, а именно, букв у на одну больше, чем й;
- 3) $\Gamma(\text{йюю}) = \text{уіоуу}$, $\Gamma(\text{йиоййю}) = \text{уііоуууу}$ – правил с головным символом не используется, разделители требуются.

К этой закономерности мы еще вернемся.

Шаг 8 (разрешение конфликтов гласных: двойные конфликты).

Только что был разрешен односторонний конфликт цепочек для букв ё и у. После добавления правил для гласных, ъ и ь возникло пять подобных конфликтов: ё-у, ё-ъ, ё-э, э-ы, ь-ы. Наглядно их можно представить на графе (рис. 5), где вершины – английские гласные, а направленные ребра – двухбуквенные образы русских гласных, ъ и ь.

Конфликту двух букв соответствует путь длиной в два ребра, их в графе как раз пять штук. Легко видеть, что в графе к тому же есть один путь ё-э-ы длиной в три ребра: он соответствует тройному одностороннему конфликту. Никаких более путей (как и циклов) в графе гласных нет.

Оставшиеся четыре двойных односторонних конфликта разрешим полностью аналогично конфликту ё-у, рассмотренному на предыдущем

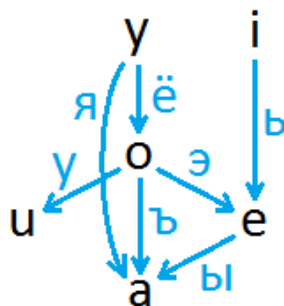


Рис. 5: граф гласных

шаге. Будет верен тот же принцип: перед использованием правила с головным символом разделитель не требуется. Результат можно описать и в $(*,+)$ -нотации:

йи*ъ \longleftrightarrow уі*оа
 ёх*а \longleftrightarrow уоh⁺а
 йи*ох*а \longleftrightarrow уі⁺оh⁺а

ой*ы \longleftrightarrow оу*еа
 эх*а \longleftrightarrow оeh⁺а
 ой*ех*а \longleftrightarrow оу⁺eh⁺а

йи*э \longleftrightarrow уі*ое
 ёй*е \longleftrightarrow уоу⁺е
 йи*ой*е \longleftrightarrow уі⁺оу⁺е

ий*ы \longleftrightarrow іу*еа
 ьх*а \longleftrightarrow іeh⁺а
 ий*ех*а \longleftrightarrow іу⁺eh⁺а

Шаг 9 (разрешение конфликтов гласных: тройной конфликт).

Осталось разрешить односторонний конфликт ё-э-ы. Как и раньше, построим бесконечную таблицу конфликтующих «наложений», но она будет уже не двумерной, а трехмерной (приведены первые три слоя – рис. 6, 7, 8). Номер клетки по горизонтали – количество первого символа-разделителя, по вертикали – второго, в глубину – третьего.

Вид слоев, начиная со второго, повторяет вид двумерной таблицы для двойного конфликта. На первом слое английских слов нет вовсе, русские занимают все клетки, кроме первого ряда. Подобная сложная форма множества рассматриваемых русских слов объясняется тем, что правила с первой и второй, а также со второй и третьей головами не могут быть применены одновременно, а с первой и третьей – могут (первый столбец первого слоя).

Чтобы разрешить конфликт, поступим следующим образом. Весь первый слой сдвинем на клетку вверх; во всех последующих слоях сдвинем верхний ряд на клетку влево. Теперь объединим каждое русское слово в пару с английским, стоящим от него снизу справа на следующем слое (рис. 9, 10, 11). Все полученные правила добавим в Г.

			...
ёы	йоы	йиоы	...
ёйы	йойы	йиойы	...
⋮	⋮	⋮	

Рис. 6: конфликтующие «наложения», слой 1

	йэа	йиэа	...
ёеа	йоеа	йиоеа	...
	уоеа	уіоеа	...
ёйеа	йойеа	йиойеа	...
	уоуеа	уіоуеа	...
⋮	⋮	⋮	

Рис. 7: конфликтующие «наложения», слой 2

	йэха	йиэха	...
ёеха	йоеха	йиоеха	...
	уоеха	уіоеха	...
ёйеха	йойеха	йиойеха	...
	уоуеха	уіоуеха	...
⋮	⋮	⋮	

Рис. 8: конфликтующие «наложения», слой 3

ёы	йоы	йиоы	...
ёйы	йойы	йиойы	...
ёййы	йоййы	йиоййы	...
⋮	⋮	⋮	

Рис. 9: разрешение конфликта, слой 1

йэа	йиэа	йииэа	...
ёеа	йоеа	йиоеа	...
	уоеа		...
ёйеа	йойеа	йиойеа	...
	уоуеа	уіоуеа	...
⋮	⋮	⋮	

Рис. 10: разрешение конфликта, слой 2

йэха	йиэха	йииэха	...
ёеха	йоеха	йиоеха	...
	уоеха	уіоеха	...
ёйеха	йойеха	йиойеха	...
	уоуеха	уіоуеха	...
⋮	⋮	⋮	

Рис. 11: разрешение конфликта, слой 3

Естественная структура множества добавленных правил трехмерная. Его можно представить в $(*, +)$ -нотации:

$$\begin{aligned}\ddot{y}y^* &\longleftrightarrow yoy^*ea \\ \ddot{y}i^*oy^* &\longleftrightarrow yi^+oy^*ea \\ \ddot{y}i^*ex^*a &\longleftrightarrow yi^*oeh^+a \\ \ddot{y}y^*ex^*a &\longleftrightarrow yoy^+eh^+a \\ \ddot{y}i^*oy^*ex^*a &\longleftrightarrow yi^+oy^+eh^+a\end{aligned}$$

Заметим, что добавленное множество правил описывается практически той же закономерностью, что и разрешение двойного одностороннего конфликта: непосредственно перед использованием правила с головным символом разделитель не требуется. В качестве уточнения добавлено слово непосредственно. Примеры:

- 1) $\Gamma(\ddot{y}y) = yoea$, $\Gamma(\ddot{y}y) = yoyea$: одновременно использованы правила с первой (\ddot{e}) и третьей (y) головами, разделитель \ddot{y}/y перед третьей не требуется;
- 2) $\Gamma(\ddot{y}oy) = yioea$, $\Gamma(\ddot{y}ioy) = yiiioea$, $\Gamma(\ddot{y}ioy) = yioyea$, $\Gamma(\ddot{y}ioy) = yiiioyea$: использовано правило с головой y , второй разделитель (\ddot{y}/y , непосредственно перед) не нужен, но первый (i/i) нужен;
- 3) $\Gamma(\ddot{y}эa) = yoeha$, $\Gamma(\ddot{y}иэa) = yioeha$, $\Gamma(\ddot{y}эха) = yoehhha$, $\Gamma(\ddot{y}иэхa) = yioehhha$: использовано правило с головой $э$, первый разделитель не нужен, третий нужен;
- 4) $\Gamma(\ddot{y}ea) = yoyeha$, $\Gamma(\ddot{y}йea) = yoyueha$, $\Gamma(\ddot{y}eha) = yoyehhha$, $\Gamma(\ddot{y}йeha) = yoyuehhha$: использовано правило с головой \ddot{e} , второй и третий разделители нужны;
- 5) $\Gamma(\ddot{y}oea) = yioyeha$, $\Gamma(\ddot{y}ioйeha) = yiiioyuehhha$: правила с головами не встречаются, все разделители нужны.

Итак, верен следующий закон.

Закон 1. Чтобы выполнить перевод при одностороннем конфликте ($\ddot{e}-y$, $\ddot{e}-\ddot{y}$, $\ddot{e}-э$, $э-ы$, $ь-ы$ либо $\ddot{e}-э-ы$), нужно применить обычные правила из конфликтующих цепочек с тем лишь отличием, что непосредственно перед использованием правила с головным символом разделитель не требуется.

Шаг 10 (q).

Для всех русских букв были добавлены правила, содержащие их в качестве левой части. Однако осталось три английских буквы, для которых еще нет правил: q , w и x .

Добавим правило $k\ddot{y} \longleftrightarrow q$. Как и раньше, для биективности Γ необходимо добавить бесконечную цепочку правил. Ситуация симметричная привычной нам: символов-разделителей на один больше не в правых частях, а в левых. Пусть разделителем будет g/g .

$$\begin{aligned} \text{кѣ} &\longleftrightarrow \text{q} \\ \text{кг}^+\text{ѣ} &\longleftrightarrow \text{kg}^*\text{oa} \end{aligned}$$

Шаг 11 (х, разрешение конфликтов).

Для буквы х добавим следующее правило, вместе с ним – бесконечную цепочку:

$$\begin{aligned} \text{кс} &\longleftrightarrow \text{х} \\ \text{кг}^+\text{с} &\longleftrightarrow \text{kg}^*\text{s} \end{aligned}$$

Обратим внимание, что между цепочками для х и ш и для х и щ возникли конфликты. Рассмотрим для начала конфликт х-ш (рис. 12).

	кс	кгс	кггс	кгггс	...
х	ks	kgs	kggs	kgggs	...

ш	сх	скх	сккх	скккх	...
	sh	skh	skkh	skkkh	...

Рис. 12: цепочки правил для х и ш

Ситуация похожа на уже известную нам (например, конфликт ё-у), но теперь один головной символ английский, другой русский. Такой конфликт назовем *двусторонним*.

Как и раньше, выпишем все конфликтующие «наложения» слов из цепочек в бесконечную таблицу (рис. 13); номер клетки по горизонтали – количество первых символов-разделителей, по вертикали – количество вторых. Русские слова заполняют всю таблицу кроме первого столбца, английские – кроме первой строки. Чтобы разрешить конфликт, каждому русскому слову поставим в пару английское, стоящее от него снизу слева (рис. 14). Полученные правила добавим в Г.

Множество добавленных правил можно представить в $(*,+)$ -нотации:

$$\begin{aligned} \text{кш} &\longleftrightarrow \text{xh} \\ \text{кг}^+\text{ш} &\longleftrightarrow \text{kg}^*\text{sh} \\ \text{кск}^*\text{x} &\longleftrightarrow \text{xk}^+\text{h} \\ \text{кг}^+\text{ск}^*\text{x} &\longleftrightarrow \text{kg}^*\text{sk}^+\text{h} \end{aligned}$$

	кш	кгш	
			...
	кcx	кгcx	
xh	ksh	kgsh	...
	кскx	кгскx	
xkh	kskh	kgskh	...
⋮	⋮	⋮	

Рис. 13: таблица конфликтующих «наложений»

	кш	кгш	
			...
	кcx	кгcx	
xh	ksh	kgsh	...
	кскx	кгскx	
xkh	kskh	kgskh	...
⋮	⋮	⋮	

Рис. 14: разрешение конфликта

Проще осмыслить его через следующую закономерность. Заметим, что для перевода любого слова из таблицы нужно сначала провести замену на одно слово вперед по цепочке с головой на исходном языке, затем перевод побуквенно, и, наконец, замену на одно слово назад по цепочке с головой на целевом языке. Примеры:

- 1) $\Gamma(\text{кш}) = \text{xh}$: $\text{кш} \rightarrow (\text{замена вперед по цепочке для ш}) \text{кcx} \rightarrow \text{ksh} \rightarrow (\text{замена назад по цепочке для х}) \text{xh}$; обратный перевод: $\text{xh} \rightarrow (\text{замена вперед по цепочке для х}) \text{ksh} \rightarrow \text{кcx} \rightarrow (\text{замена назад по цепочке для ш}) \text{кш}$;
- 2) $\Gamma(\text{кгш}) = \text{ksh}$: $\text{кгш} \rightarrow \text{кгcx} \rightarrow \text{kgsh} \rightarrow \text{ksh}$; обратно: $\text{ksh} \rightarrow \text{kgsh} \rightarrow \text{кгcx} \rightarrow \text{кгш}$;
- 3) $\Gamma(\text{кcx}) = \text{xkh}$: $\text{кcx} \rightarrow \text{кскx} \rightarrow \text{kskh} \rightarrow \text{xkh}$; обратно: $\text{xkh} \rightarrow \text{kskh} \rightarrow \text{кскx} \rightarrow \text{кcx}$;
- 4) $\Gamma(\text{кхcx}) = \text{kskh}$: $\text{кхcx} \rightarrow \text{кхскx} \rightarrow \text{khskh} \rightarrow \text{kskh}$; обратно: $\text{kskh} \rightarrow \text{khskh} \rightarrow \text{кхскx} \rightarrow \text{кхcx}$.

Конфликт цепочек х и ш имеет аналогичный вид (рис. 15). Разрешается он тем же способом, верна та же закономерность.

	кc	кгc	кггc	кгггc	...
x	ks	kgs	kggs	kgggs	...

щ	шч	штч	шттч	штттч	...
	shch	shtch	shttch	shtttch	...

Рис. 15: цепочки правил для х и ш

Добавленные правила можно описать в $(*, +)$ -нотации:

$$\begin{aligned} \text{кщ} &\longleftrightarrow \text{xhsh} \\ \text{кг}^+\text{щ} &\longleftrightarrow \text{kg}^*\text{shch} \\ \text{кшт}^*\text{ч} &\longleftrightarrow \text{xht}^+\text{ch} \\ \text{кг}^+\text{шт}^*\text{ч} &\longleftrightarrow \text{kg}^*\text{sht}^+\text{ch} \end{aligned}$$

Обратим внимание на следующую особенность. При «наложении» русских слов из цепочек «на стыке» стоит ш (в него «сворачивается» sh). Так, например, кс и щч могут образовать конфликтующее «наложение» кшч, хотя на первый взгляд это может показаться противостественным. Также нужно использовать уже введенные правила для ш и ч при побуквенном переводе. Примеры:

- 1) $\Gamma(\text{кщ}) = \text{xhch}$: $\text{кщ} \longrightarrow (\text{замена вперед по цепочке для щ}) \text{кшч} \longrightarrow \text{kshch} \longrightarrow (\text{замена назад по цепочке для х}) \text{xhch}$; обратный перевод: $\text{xhch} \longrightarrow (\text{замена вперед по цепочке для х}) \text{kshch} \longrightarrow \text{кшч} \longrightarrow (\text{замена назад по цепочке для ш}) \text{кщ}$;
- 2) $\Gamma(\text{кгщ}) = \text{kshch}$: $\text{кгщ} \longrightarrow \text{кгшч} \longrightarrow \text{kgshch} \longrightarrow \text{kshch}$; обратно: $\text{kshch} \longrightarrow \text{kgshch} \longrightarrow \text{кгшч} \longrightarrow \text{кгщ}$;
- 3) $\Gamma(\text{кшч}) = \text{xhtch}$: $\text{кшч} \longrightarrow \text{кштч} \longrightarrow \text{kshtch} \longrightarrow \text{xhtch}$; обратно: $\text{xhtch} \longrightarrow \text{kshtch} \longrightarrow \text{кштч} \longrightarrow \text{кшч}$;
- 4) $\Gamma(\text{кгшч}) = \text{kshtch}$: $\text{кгшч} \longrightarrow \text{кгштч} \longrightarrow \text{kgshtch} \longrightarrow \text{kshtch}$; обратно: $\text{kshtch} \longrightarrow \text{kgshtch} \longrightarrow \text{кгштч} \longrightarrow \text{кгшч}$.

Итак, верен следующий закон.

Закон 2. Чтобы выполнить перевод при двустороннем конфликте (х-ш либо х-щ), нужно провести замену на одно слово вперед по цепочке с головой на исходном языке, затем перевод побуквенно, и, наконец, замену на одно слово назад по цепочке с головой на целевом языке.

Шаг 12 (w, разрешение конфликта).

Добавим следующее правило, вместе с ним – бесконечную цепочку:

$$\begin{aligned} \text{ув} &\longleftrightarrow \text{w} \\ \text{уф}^+\text{в} &\longleftrightarrow \text{ouf}^*\text{v} \end{aligned}$$

Возникает конфликт (рис. 16) между этой цепочкой и $\text{йи}^*\text{у} \longleftrightarrow \text{yi}^*\text{ou}$ (правилами верхней строки таблицы на рис. 4).

В цепочке $\text{йи}^*\text{у} \longleftrightarrow \text{yi}^*\text{ou}$ слова в некотором смысле соответствуют побуквенно. Она не является привычной нам цепочкой правил, получающейся сдвигом сопоставлений слов на один относительно побуквенных и прибавлением головы. У конфликта только одна голова – обозначим его йу-w.

	йу	йиу	йииу	йиииу	...
	↕	↕	↕	↕	
	you	yіou	yііou	yіііou	...

	уѳ	уѳѳ	уѳѳѳ	уѳѳѳѳ	...
	↗	↗	↗	↗	
w	ouѳ	ouѳѳ	ouѳѳѳ	ouѳѳѳѳ	...

Рис. 16: цепочки правил йу и w

Построим бесконечную таблицу конфликтующих «наложений» слов из цепочек (рис. 17). Русские слова заполняют всю таблицу кроме первого ряда, английские – всю целиком. Чтобы разрешить конфликт, сдвинем все русские слова на строку вверх и поставим в пару слова, стоящие в одной клетке (рис. 18)

			...
уw	yіw	yііw	...
йуѳ	йиуѳ	йииуѳ	...
youѳ	yіouѳ	yііouѳ	...
йуѳѳ	йиуѳѳ	йииуѳѳ	...
youѳѳ	yіouѳѳ	yііouѳѳ	...
⋮	⋮	⋮	

Рис. 17: таблица конфликтующих «наложений»

йуѳ	йиуѳ	йииуѳ	...
↕	↕	↕	
уw	yіw	yііw	...
йуѳѳ	йиуѳѳ	йииуѳѳ	...
↕	↕	↕	
youѳ	yіouѳ	yііouѳ	...
йуѳѳѳ	йиуѳѳѳ	йииуѳѳѳ	...
↕	↕	↕	
youѳѳѳ	yіouѳѳѳ	yііouѳѳѳ	...
⋮	⋮	⋮	

Рис. 18: разрешение конфликта

В $(*,+)$ -нотации множество добавленных правил описывается так:

$$\begin{aligned} \text{йи}^* \text{уѳ} &\longleftrightarrow \text{yі}^* \text{w} \\ \text{йи}^* \text{уѳ}^+ \text{ѳ} &\longleftrightarrow \text{yі}^* \text{ouѳ}^* \text{ѳ} \end{aligned}$$

Как видно, выполнена несложная закономерность: перевод происходит так, будто правила из цепочки йу игнорируются.

Закон 3. При конфликте йу-w перевод происходит без учета правил из цепочки йу.

7 Итоговые формулировки

Итак, было построено следующее множество правил Γ из \mathcal{R}^* в \mathcal{A}^* . Для удобства представим его в виде $\Gamma = \Gamma_1 \cup \Gamma_2$: в Γ_1 лежат ключевые правила, в Γ_2 – добавленные для разрешения конфликтов.

Γ_1 :

$a \longleftrightarrow a$	$ж \longleftrightarrow j$	$м \longleftrightarrow m$	$т \longleftrightarrow t$
$б \longleftrightarrow b$	$з \longleftrightarrow z$	$н \longleftrightarrow n$	$ф \longleftrightarrow f$
$в \longleftrightarrow v$	$и \longleftrightarrow i$	$о \longleftrightarrow o$	$х \longleftrightarrow h$
$г \longleftrightarrow g$	$й \longleftrightarrow y$	$п \longleftrightarrow p$	$ц \longleftrightarrow c$
$д \longleftrightarrow d$	$к \longleftrightarrow k$	$р \longleftrightarrow r$	$ю \longleftrightarrow u$
$е \longleftrightarrow e$	$л \longleftrightarrow l$	$с \longleftrightarrow s$	
$тф \longleftrightarrow th$	$тх \longleftrightarrow tf$	$пф \longleftrightarrow ph$	$пх \longleftrightarrow pf$
$ч \longleftrightarrow ch$	$у \longleftrightarrow ou$	$ъ \longleftrightarrow oa$	$къ \longleftrightarrow q$
$цк^*x \longleftrightarrow ck^+h$	$ой^*ю \longleftrightarrow oy^+u$	$ох^*a \longleftrightarrow oh^+a$	$кг^*ъ \longleftrightarrow kg^*oa$
$ш \longleftrightarrow sh$	$э \longleftrightarrow oe$	$ы \longleftrightarrow ea$	$кс \longleftrightarrow x$
$ск^*x \longleftrightarrow sk^+h$	$ой^*е \longleftrightarrow oy^+e$	$ех^*a \longleftrightarrow eh^+a$	$кг^*с \longleftrightarrow kg^*s$
$щ \longleftrightarrow shch$	$ё \longleftrightarrow yo$	$ь \longleftrightarrow ie$	$ув \longleftrightarrow w$
$шт^*ч \longleftrightarrow sht^+ch$	$йи^*о \longleftrightarrow yi^+o$	$ий^*е \longleftrightarrow iy^+e$	$уф^*в \longleftrightarrow ouf^*v$
	$я \longleftrightarrow ya$		
	$йи^*a \longleftrightarrow yi^+a$		

Γ_2 :

$йи^*у \longleftrightarrow yi^*ou$	$ой^*ы \longleftrightarrow oy^*ea$
$ёй^*ю \longleftrightarrow yo y^+u$	$эх^*a \longleftrightarrow oeh^+a$
$йи^*ой^*ю \longleftrightarrow yi^+oy^+u$	$ой^*ех^*a \longleftrightarrow oy^+eh^+a$
$йи^*ъ \longleftrightarrow yi^*oa$	$ий^*ы \longleftrightarrow iy^*ea$
$ёх^*a \longleftrightarrow yoh^+a$	$ьх^*a \longleftrightarrow ieh^+a$
$йи^*ох^*a \longleftrightarrow yi^+oh^+a$	$ий^*ех^*a \longleftrightarrow iy^+eh^+a$
$йи^*э \longleftrightarrow yi^*oe$	$ёй^*ы \longleftrightarrow yo y^*ea$
$ёй^*е \longleftrightarrow yo y^+e$	$йи^*ой^*ы \longleftrightarrow yi^+oy^*ea$
$йи^*ой^*е \longleftrightarrow yi^+oy^+e$	$йи^*эх^*a \longleftrightarrow yi^*oeh^+a$
	$ёй^*ех^*a \longleftrightarrow yo y^+eh^+a$
	$йи^*ой^*ех^*a \longleftrightarrow yi^+oy^+eh^+a$

$$\begin{array}{ll}
\text{кш} \longleftrightarrow \text{xh} & \text{кщ} \longleftrightarrow \text{xhsh} \\
\text{кг}^+ \text{ш} \longleftrightarrow \text{kg}^* \text{sh} & \text{кг}^+ \text{ш} \longleftrightarrow \text{kg}^* \text{shch} \\
\text{кск}^* \text{x} \longleftrightarrow \text{xk}^+ \text{h} & \text{кшт}^* \text{ч} \longleftrightarrow \text{xht}^+ \text{ch} \\
\text{кг}^+ \text{ск}^* \text{x} \longleftrightarrow \text{kg}^* \text{sk}^+ \text{h} & \text{кг}^+ \text{шт}^* \text{ч} \longleftrightarrow \text{kg}^* \text{sht}^+ \text{ch} \\
\\
\text{йи}^* \text{ув} \longleftrightarrow \text{yi}^* \text{w} & \\
\text{йи}^* \text{уф}^+ \text{в} \longleftrightarrow \text{yi}^* \text{ouf}^* \text{v} &
\end{array}$$

Теорема 3. *Описанное выше множество Γ , а также Γ^{-1} – корректные системы правил, Γ задает англо-русскую биекцию $F : \mathcal{R}^* \longrightarrow \mathcal{A}^*$, Γ^{-1} задает обратную функцию F^{-1} .*

Доказательство. По теореме 1 Γ и Γ^{-1} – корректные системы правил. $F(\omega) = \Gamma(\omega)$ есть биективная функция из \mathcal{R}^* в \mathcal{A}^* (построена биективной в предыдущем разделе). По теореме 2 система правил Γ^{-1} задает обратную функцию F^{-1} . \square

Для получения $\Gamma(\xi)$ по данному $\xi \in \mathcal{R}^*$ и $\Gamma^{-1}(\eta)$ по данному $\eta \in \mathcal{A}^*$ можно использовать алгоритм 1. При переводе на компьютере программа-транслитератор может оперировать правилами из Γ , представленными в $(*, +)$ -нотации либо как-то иначе.

Для транслитерации вручную гораздо проще использовать множество ключевых правил Γ_1 и сформулированные в предыдущем разделе законы 1, 2 и 3. В этом случае в алгоритме 1 выбирается максимальный отрезок с учетом конфликтующих «наложений». В предыдущем разделе показано, что использование трех законов равносильно применению правил из Γ_1 .

8 Обобщенная транслитерация текста

Кратко рассмотрим прикладную задачу транслитерации. Пусть \mathcal{R}^c – алфавит заглавных русских букв, \mathcal{A}^c – заглавных английских; \mathcal{Z} – алфавит посторонних символов.

Определение 18. Алфавит $\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z}$ назовем *расширенным русским алфавитом*, язык $(\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z})^*$ всех слов над ним – *расширенным русским языком*. Аналогично $\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z}$ – *расширенный английский алфавит*, $(\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z})^*$ – *расширенный английский язык*. Слово расширенного языка будем называть *текстом*.

Требуется на основе англо-русской биекции $F : \mathcal{R}^* \longrightarrow \mathcal{A}^*$ построить биекцию $G : (\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z})^* \longrightarrow (\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z})^*$. Это можно сделать просто и естественно следующим образом.

Определение 19. *Словом в узком смысле над расширенным алфавитом $\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z}$ (или $\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z}$) назовем конечную последовательность букв из $\mathcal{R} \cup \mathcal{R}^c$ (соответственно, из $\mathcal{A} \cup \mathcal{A}^c$), в которой все буквы строчные, кроме, быть может, первой, которая допускается заглавной.*

Утверждение 7. *Любой текст над $\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z}$ (или $\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z}$) можно представить в виде конкатенации наименьшего возможного числа слов в узком смысле и последовательностей символов из \mathcal{Z} единственным образом.*

Алгоритм 2 (транслитерации текста). Пусть дано $\tau \in (\mathcal{R} \cup \mathcal{R}^c \cup \mathcal{Z})^*$, получим $G(\tau)$.

1. Прделаем разбиение τ , описанное в утверждении 7.
2. Пусть ω – слово в узком смысле, содержащееся в разбиении. Тогда заменим первую букву на строчную, получим перевод $F(\omega)$, а затем заменим первую букву на заглавную, если у исходного слова она была заглавной.
3. Сконкатенируем транслитерированные слова в узком смысле и оставленные без изменения последовательности символов над \mathcal{Z} .
Получен образ $G(\tau)$.

Полностью аналогично можно получить $G^{-1}(v)$ по данному $v \in (\mathcal{A} \cup \mathcal{A}^c \cup \mathcal{Z})^*$. Несложно установить, что построенная функция G биективна.