

Самообучающаяся система машинной транскрипции с использованием нестохастического конечного автомата*

© В.К. Логачева, Э.С. Клышинский

Российский государственный гуманитарный университет, ИПМ им. М.В. Келдыша РАН,
г. Москва
logacheva_vk@mail.ru, klyshinsky@mail.ru

Аннотация

Описан метод транскрипции имен собственных, использующий конечный автомат и правила, записанные в виде продукций. Представлен метод автоматического порождения правил. Сгенерированные правила могут быть преобразованы в конечный автомат, с помощью которого осуществляется транслитерация.

1 Введение

Для электронных архивов большого объема очень важна автоматическая подготовка документов. При поступлении в библиотеку иноязычного документа для помещения его в каталог может понадобиться перевод его названия, аннотации, выходных данных. Системы машинного перевода зачастую оставляют имена собственные (в частности, фамильно-именные группы) без перевода, тогда как фамилия автора – самая важная информация о документе. Перевод имен вручную затруднителен не только потому, что на него тратится много времени, но и из-за различающихся в разных языках правил чтения, которые могут быть неизвестны переводчику, не владеющему данным языком.

Существует несколько подходов к передаче имени собственного средствами другого языка:

- перевод (например, Easter Island – остров Пасхи). Этот способ очень редко можно применить для имен собственных, так как они зачастую не имеют лексического значения;
- транслитерация:
 - строгая – сопоставление каждой букве исходного языка буквы языка перевода; этот способ может исказить звучание слова, так как почти во всех языках существуют диграфы – устойчивые сочетания букв, которые читаются особым образом; даже правил расширенной транслитерации, т. е. правил, допускающих сопоставление одной букве нескольких (sh → ш), не всегда хватает, чтобы описать все зависимости фонетики и графики, существ-

вующие в языке;

- транслитерация с учетом фонетического облика слова, также называемая транскрипцией; в отличие от транскрипции в классическом лингвистическом понимании слова исходного языка записываются не особым фонетическим алфавитом, а символами целевого языка.

В разное время наибольшей популярностью пользовались разные подходы, но с середины XX века господствующей тенденцией стала передача звучания имени при его переводе на иностранный язык. С развитием компьютерной лингвистики и средств автоматизации встал вопрос о создании систем автоматической транслитерации и транскрипции имен собственных. К ранним работам в этой области относятся работы группы Найта. Первоначально был разработан метод, позволяющий проводить транслитерацию с японского языка на английский [1]. Вскоре этот метод был адаптирован для работы с арабским языком [2]. Основной задачей метода являлось обнаружение и восстановление имен собственных, уже переданных с использованием другого алфавита. В качестве основы использовался модифицированный алгоритм Витерби. Однако базовый принцип использования лишь отдельных символов алфавита не позволял получить высокие показатели качества передачи. В связи с этим стали применяться методы, использующие подстроки [3]. Применение подстрок взамен отдельных символов позволило повысить качество транслитерации примерно с 30 до 90%.

Еще одной важной задачей здесь является автоматическое обучение системы переходов конечного автомата, используемого в перечисленных методах. На данный момент развиваются методы как обучения с учителем [1], так и без него [4]. Последние работы позволили перейти к обучению по одноязычному корпусу [5], что привело, правда, к существенному падению качества передачи.

Данная работа основана на программе «Транскриба» [6]. В отличие от приведенных работ, здесь используется не стохастический, а детерминированный подход. Лингвистам предлагается написать правила, которые поступают на вход системы. Метод позволяет получить высокую эффективность передачи, однако требует затрат на создание систе-

Труды 12^й Всероссийской научной конференции «Электронные библиотеки: перспективные методы и технологии, электронные коллекции» – RCDL'2010, Казань, Россия, 2010

мы правил. Следует заметить, что правила для одного языка создаются за срок от нескольких недель до нескольких месяцев, в зависимости от сложности задачи.

2 Метод транскрипции

Предлагаемый в данной работе метод объединяет удобство и скорость разбора входной цепочки с использованием конечных автоматов с удобством представления правил преобразования в привычном для человека виде. Изначально правила преобразования цепочек из одного языка на другой представляются в виде «входная цепочка» → «выходная цепочка». Скорость применения подобных правил будет зависеть от размеров базы и длины входной цепочки. В связи этим в данной статье для преобразования цепочек используется конечный автомат. Скорость разбора с использованием конечного автомата пропорциональна длине цепочки. Здесь мы считаем, что набор правил может быть как создан лингвистом, так и получен в ходе обучения по алгоритму, описанному ниже. Таким образом, при необходимости специалист работает в привычных ему терминах правил. Далее набор правил может быть расширен за счет дообучения. В другом случае набор правил может быть изначально автоматически сформирован с использованием предлагаемого метода. Набор правил преобразуется в конечный автомат, который используется для передачи имен собственных с языка на язык.

2.1 Структура конечного автомата

Используемый нами автомат представляет собой преобразователь. Его можно представить как множество $g = \langle V_I, V_O, Q, q_0, F, \delta \rangle$, где

V_I – входной алфавит (алфавит языка оригинала);

V_O – выходной алфавит (алфавит языка перевода);

Q – множество состояний автомата;

q_0 – начальное состояние автомата;

F – множество конечных состояний;

δ – функция переходов $Q \times V_I \rightarrow Q, a$.

Здесь a – это действие, совершаемое при переходе, причем $a \in A$ – множеству действий, выполняемых конечным автоматом при переходе; $A = \{out(), shift()\}$; $Out(s)$ – функция, подающая на выход строку символов s (в том числе пустую) из алфавита перевода; $Shift(n)$ – функция, выполняющая сдвиг рассматриваемой строки на n символов.

Из каждого конечного состояния существует переход в начальное состояние по пустому символу.

2.2 Описание правил транскрипции

За основу был взят метод, описанный в [6]. Но, поскольку в предлагаемом методе правила используются, прежде всего, для пополнения матрицы переходов конечного автомата и генерируются компьютером, а не человеком, их структура отличается от базового варианта. В нашем случае правило представляется в виде пары $r = \langle p, c \rangle$, где

p – непустое упорядоченное множество символов из алфавита языка оригинала;

c – цепочка символов из алфавита языка перевода, которой передается входная цепочка;

$p = \langle p1, p2, p3 \rangle$, где $p1$ и $p3$ – пред- и постусловия соответственно (левый и правый контекст правила), а $p2$ – преобразуемая подстрока, $p2 = \langle v_1, v_2, \dots, v_{n2} \rangle$, где $v_i \in V_I, i \in [1, n2]$, $p1 = \langle \{v_1\}, \{v_2\}, \dots, \{v_{n1}\} \rangle$, где $v_i \in V_I, i \in [1, n1]$, $p3 = \langle \{v_1\}, \{v_2\}, \dots, \{v_{n3}\} \rangle$, где $v_i \in V_I, i \in [1, n3]$, n_j – количество символов j -й подстроки;

$c = \langle w_1, w_2, \dots, w_{nc} \rangle$, где $w_i \in V_O, i \in [1, nc]$, nc – длина выходной строки.

Исходя из этого, правило применимо с текущей позиции, если на текущей позиции находится подстрока $p2$, перед ней – подстрока, на соответствующих позициях которой находятся символы из $p1$, а после нее – подстрока, перед ней – подстрока, на соответствующих позициях которой находятся символы из $p3$. Считается, что к входной строке последовательно ищутся применяемые правила. При нахождении такого правила текущая позиция сдвигается вправо на $n2$, а на выход подается c .

2.3 Преобразование правил в конечный автомат

Для правила $r = \langle p, c \rangle$ процедура добавления в конечный автомат будет следующей. Здесь, за исключением начального состояния q_0 , все состояния заводятся для данного правила заново. При этом дополнительным алгоритмом обеспечивается уникальность нумерации состояний.

Для правил, в которых левый контекст пуст, в КА добавляется переход из состояния q_{i-1} в состояние q_i по символу v_i из $p2, i \in [1, n2]$. При этом никакие действия не совершаются. Состояние q_{n2} помечается как конечное.

Для правил, в которых правый контекст пуст, в КА добавляется переход из q_{n2} в q_0 по пустому символу, сопровождаемый $Shift(n2)$ и $Out(c)$.

Для правил, в которых правый контекст не пуст, в КА добавляется переход из состояния q_{n2+i-1} в состояние q_{n2+i} по всем символам из $\{v_i\}$ из $p3, i \in [1, n3]$. Далее в КА добавляется переход из q_{n2+n3} в q_0 по пустому символу, сопровождаемый действиями $Shift(n2)$ и $Out(c)$.

Для правил, в которых левый контекст не пуст, для всех $q_f \in F$ в КА добавляются переходы из q_f в $q_{n2+n3+1}$ по всем символам $\{v_i\}$ из $p1$ и переходы из состояния $q_{n2+n3+1}$ в состояние $q_{n2+n3+i}$ по всем символам из $\{v_i\}$ из $p1, i \in (1, n1]$. При этом никакие действия не совершаются. Далее в КА добавляется переход из $q_{n1+n2+n3}$ в q_0 по пустому символу, сопровождаемый действиями $Shift(n2)$ и $Out(c)$.

Подобные операции проводятся для всех правил.

3 Порождение правил

Для облегчения поиска соответствий в строках на языке оригинала и перевода каждое слово и его перевод делятся на слоги. Поскольку мы не имеем никакой информации о структуре языка и правилах

слогоделения в этом языке, деление производится по формальным признакам. В порождении правил участвуют только слова, для которых установлено взаимно однозначное соответствие слогов в оригинале и переводе.

3.1 Разделение слова на слоги

Термин «слог» в данном контексте употребляется в смысле, отличном от того, в котором он используется в трудах по фонетике. Поэтому необходимо дать его определение.

Слог – это непустое множество букв, содержащее один или более символов, обозначающих гласные звуки. Слово делится на слоги по следующим правилам:

- деление осуществляется по гласной букве; граница слога – после гласной;
- 2 и более гласных подряд не разделяются;
- 2 и более согласных не разделяются;
- множество элементов, среди которых нет гласной, не выделяется в отдельный слог;
- символы начала и конца слова считаются согласными буквами.

Дополнительное правило для русского языка:

- не отделять «ь» от предыдущей согласной.

Таким образом, слог – это цепочка вида C^*V^+ , где C – буква, обозначающая согласный звук, V – буква, обозначающая гласный звук. Слог может иметь вид $C^*V^+C^+$ только в том случае, если это последний слог в слове, и заключительное множество согласных не может быть выделено в отдельный слог, так как среди них нет слогообразующего символа.

Формально почти все слоги, полученные по изложенным выше правилам, являются открытыми (то есть заканчиваются на гласную). Однако у слога присутствует параметр открытости/закрытости, определяющийся следующим образом: при делении для каждого слога проверяется количество начальных согласных в следующем слоге: если их больше одной, слог считается закрытым. Введение этого параметра связано с тем, что во многих языках (например, в английском и других языках германской группы) он может влиять на правила чтения гласных букв.

3.2 Порождение первичных правил

Первый этап порождения правил – порождение первичных правил. Первичными мы называем правила транслитерации, то есть правила, для которых $|p2| = |c| = 1$.

Первичные правила порождаются на основе слогов вида $CV \rightarrow CV$: i -му символу слога ставится в соответствие i -й символ его перевода при условии, что оба символа обозначают звуки одного и того же типа (гласные или согласные). При таком подходе вероятность порождения некорректного правила очень мала.

В порожденном множестве первичных правил могут содержаться неоднозначности, то есть пары правил $r1 = \langle p, c \rangle$ и $r2 = \langle p, c \rangle$, такие, что

$r1(p2) = r2(p2)$ и $r1(c) \neq r2(c)$. Такие ситуации в принципе возможны в языке, но на данном этапе неоднозначность скорее всего обозначает влияние на букву окружающих ее букв. Для избавления от неоднозначностей вводится новый тип правил – временные правила.

Временными правилами мы называем правила, используемые при обучении системы правил. Временное правило представляет собой тройку $r = \langle p, c, s \rangle$, где

p, c – те же, что в определенном выше формате правила (см. п. 2.2);

s – множество слогов, удовлетворяющих правилу.

При генерации правила в его множество s добавляется слог, на основе которого оно было сгенерировано. При обучении системы правил слог, к которому удалось применить правило, добавляется в множество s . Хранение всех слогов позволяет при обнаружении неоднозначностей в правилах составлять более полные контексты, а также выявлять частотность употребления правила (например, если правилу удовлетворяет только один слог, можно с большой вероятностью утверждать, что это исключение).

Избавление от неоднозначностей производится следующим образом.

Для каждого правила r_i , для которого существует r_j , такое, что $r_i(p2) = r_j(p2)$, $r_i(c) \neq r_j(c)$, для каждого s_k из $r_i(s)$ составляется новое правило $r_{ik} = \langle p, c, s \rangle$, такое, что $r_{ik}(p2) = r_i(p2)$, $r_{ik}(c) = r_i(c)$, $r_{ik}(p1) = \langle \{v_1\} \rangle$, где $v_1 \in V_1$ – символ, предшествующий $p2$ в слоге s_k , $r_{ik}(p3) = \langle \{v_2\} \rangle$, $v_2 \in V_1$ – символ, следующий за $p2$ в слоге s_k . В случае, если $p2$ – начальная или заключительная подстрока в слоге, v_1 и v_2 берутся из предыдущего или следующего слога соответственно.

После порождения системы первичных правил производится расширение их контекстов на основе слогов вида $CnV \rightarrow CnV$, где $n > 1$. Каждый слог s_i , который не может быть целиком разобран с помощью существующей системы правил, можно представить как $\langle p_{i1}, \dots, p_{ik}, p_x, p_{ik+1}, \dots, p_{in} \rangle \rightarrow \langle c_{i1}, \dots, c_{ik}, c_x, c_{ik+1}, \dots, c_{in} \rangle$, где $p_x \rightarrow c_x$ – подстрока, не удовлетворяющая ни одному из существующих правил. Можно выделить три случая несоответствия p_x правилам:

- существует r_i , такое, что $p_x = r_i(p2)$, $c_x = r_i(c)$, но $p_{ik} \notin \{v_{11}\}$ или $p_{ik+1} \notin \{v_{31}\}$ (т. к. $r_i(p1) = \langle \{v_{11}\} \rangle$, $r_i(p3) = \langle \{v_{31}\} \rangle$); в этом случае контекст правила r_i расширяется: $r_i(p1) = \langle \{v_{11}, p_{ik}\} \rangle$, $r_i(p3) = \langle \{v_{31}, p_{ik+1}\} \rangle$;

- не существует r_i , такого, что $p_x = r_i(p2)$ и $c_x = r_i(c)$; в этом случае составляется новое правило r_j , такое, что $r_j(p1) = \langle \{p_{ik}\} \rangle$, $r_j(p2) = p_x$, $r_j(p3) = \langle \{p_{ik+1}\} \rangle$, $r_j(c) = c_x$, $r_j(s) = \{s_i\}$;

- существует r_i , такое, что $p_x = r_i(p2)$, $p_{ik} \in \{v_{11}\}$ и $p_{ik+1} \in \{v_{31}\}$ (т. к. $r_i(p1) = \langle \{v_{11}\} \rangle$, $r_i(p3) = \langle \{v_{31}\} \rangle$, $c_x \neq r_i(c)$); это может объясняться одной из следующих причин:

○ контекст правила r_i недостаточен для правильной интерпретации подстроки, и надо рассматривать не один, а несколько символов, предшествующих ей или следующих за ней. Например, для имени Marin \rightarrow Марен (французский язык) было порождено правило $i \rightarrow e$, которое при избавлении от неоднозначностей приобрело вид $i\{n\} \rightarrow e$. Но имя Marine \rightarrow Марин не удовлетворяет этому правилу, потому что i переходит в e (э), только если после него стоит n и слог является закрытым (заканчивается на согласную). В этом случае нужно проверять не один, а два следующих за i символа;

○ подстрока может читаться двумя различными способами в силу причин, не зависящих от контекста. В шведском языке \grave{a} читается как «о», однако в машиночитаемых текстах без диакритик \grave{a} заменяется на aa . Между тем сочетание двух шведских букв a будет передаваться на русский язык как «аа». В тестовой выборке для шведского языка можно встретить примеры неоднозначностей Baad \rightarrow Баад, Vaang \rightarrow Бонг, Naapanen \rightarrow Хаапанен, Naafman \rightarrow Хофман, которые нельзя предугадать, имея в качестве обучающей выборки машиночитаемый файл. Значит, для каждого шведского имени, содержащего подстроку «аа», будет сгенерировано два варианта перевода.

В настоящее время нет алгоритма выбора одной из двух описанных причин и избавления от подобных неоднозначностей.

3.3 Порождение сложных правил

Более сложные правила порождаются при анализе слогов более сложной структуры, то есть слогов вида $C^pV^n \rightarrow C^qV^m$, где $p \neq q$ и/или $n \neq m$. В этом случае, как и в вышеописанном алгоритме пополнения контекстов с помощью слогов вида $C^nV \rightarrow C^mV$, каждый слог можно представить как $\langle p_{i1}, \dots, p_{ik}, p_x, p_{ik+1}, \dots, p_m \rangle \rightarrow \langle c_{i1}, \dots, c_{ik}, c_x, c_{ik+1}, \dots, c_{im} \rangle$, где $p_x \rightarrow c_x$ – подстрока, не удовлетворяющая ни одному из существующих правил. Можно выделить три случая несоответствия p_x правилам:

- $p_x = \emptyset, c_x \neq \emptyset$ – эта ситуация обозначает появление в транскрипции букв, которых не было в оригинале (при составлении правил транскрипции с произвольного языка на русский чаще всего в такой позиции встречается буква «ь», которая обозначает палатализацию предыдущего согласного, которая никак не отражается на письме в языке оригинала). В этом случае составляется правило r_i , такое, что $r_i(p1) = p_{ik-1}, r_i(p2) = p_{ik}, r_i(p3) = p_{ik+1}, r_i(c) = c_{ik} + c_x$. – иными словами, правило для предшествующего p_x символа, который в определенном контексте передается несколькими символами;

- $p_x \neq \emptyset, c_x = \emptyset$ – эта ситуация обозначает наличие букв, которые не читаются, а значит, не записываются в транскрипции (например, буквы t, d, s и др. во французском языке на конце слова). Для таких случаев составляется правило r_i , такое, что $r_i(p1) = p_{ik}, r_i(p2) = p_x, r_i(p3) = p_{ik+1}, r_i(c) = \emptyset$, то есть определяется контекст, в котором подстрока p_x не читается;

- $p_x \neq \emptyset, c_x \neq \emptyset$ – в случае, когда $|c_x| = 1$, такая ситуация с большой вероятностью говорит о том, что p_x – устойчивое сочетание букв, обозначающее один звук – ди- или триграф, если же $|p_x| = 1$, это означает появление буквы, обозначающей звук, которого нет в русском языке и для записи которого используется 2 или более русских букв. В обоих случаях составляется новое правило r_i , такое, что $r_i(p1) = r_i(p3) = \emptyset, r_i(p2) = p_x, r_i(c) = c_x$.

Во избежание появления недетерминированного вывода в конечном автомате стоит проверять каждое вновь добавляемое правило на однозначность. Если для добавляемого правила $r1$ среди уже принятых правил существует такое $r2$, что $r1(p)$ и $r2(p)$ имеют общий префикс, формируется правый контекст для правила с более короткой входной строкой.

После каждого этапа (порождение новых правил, расширение контекста существующих правил) осуществляется нормализация системы правил.

3.4 Нормализация системы правил

Нормализация системы правил состоит в избавлении от правил, полностью или частично дублирующих друг друга.

Для двух правил $r_1 = \langle p, c, s \rangle$ и $r_2 = \langle p, c, s \rangle$, где $r_1(p1) = r_2(p1) = r_1(p3) = r_2(p3) = \emptyset$, справедливо утверждение $r_1 = r_2$, если $r_1(p2) = r_2(p2)$ и $r_1(c) = r_2(c)$. В этом случае правила r_1 и r_2 заменяются правилом $r_3 = r_1 \cup r_2 = \langle p, c, s \rangle$, где $r_3(p1) = r_3(p3) = \emptyset, r_3(p2) = r_1(p2) = r_2(p2), r_3(c) = r_1(c) = r_2(c)$ и $r_3(s) = r_1(s) \cup r_2(s)$.

Два правила $r_1 = \langle p, c, s \rangle$ и $r_2 = \langle p, c, s \rangle$, для которых хотя бы одно из множеств $r_1(p1), r_2(p1), r_1(p3), r_2(p3)$ не пусто, считаются эквивалентными, если $r_1(p2) = r_2(p2), r_1(c) = r_2(c), |r_1(p1)| = |r_2(p1)|, |r_1(p3)| = |r_2(p3)|$, и выполняется хотя бы одно из неравенств $r_1(p1) \cap r_2(p1) \neq \emptyset, r_1(p3) \cap r_2(p3) \neq \emptyset$. В этом случае правила r_1 и r_2 заменяются правилом $r_3 = r_1 \cup r_2 = \langle p, c, s \rangle$, где $r_3(p1) = r_1(p1) \cup r_2(p1), r_3(p3) = r_1(p3) \cup r_2(p3), r_3(p2) = r_1(p2) = r_2(p2), r_3(c) = r_1(c) = r_2(c)$ и $r_3(s) = r_1(s) \cup r_2(s)$.

В конце этапа нормализации проводится проверка на предмет порождения «фантомных» правил, для которых имеются противоречащие им.

4 Реализация метода. Результаты экспериментов

В настоящее время программно реализованы лишь некоторые этапы описанного метода: деление на слоги, отбор слов, подходящих для участия в порождении правил, составление первичной системы правил, избавление ее от неоднозначностей путем добавления контекстов, проверка слогов более сложной структуры с помощью полученной системы правил и пополнение контекстов на основе этих слогов. Каждый этап завершается нормализацией системы правил. Примеры ниже демонстрируют различные типы полученных правил с контекстами:

(1) $\{aeuy\}c\{i\} \rightarrow c$

$\{<r\}c\{e\} \rightarrow c$
 $\{<eins\}c\{aklo\} \rightarrow k$
 (2) $\{bcdgijklmnrstvz\}e\{bcdgijklmnrstvz\} \rightarrow e$
 $\{<\}e\{dlmntx\} \rightarrow e$

Из примера (1) видно, что буква *c* передается на русский язык как «с», если после нее стоят буквы *i* или *e*, и как «к» в остальных случаях. Левый контекст здесь не имеет значения, но на первых этапах анализа невозможно определить, какая из букв – следующая или предыдущая – влияет на данную. Понять это можно, только собрав достаточное количество примеров. Если контекст включает все символы алфавита, его можно не учитывать. Если контекст включает большую часть символов алфавита, не объединенных каким-либо общим свойством (например, принадлежность к согласным буквам), его также можно не учитывать. Такая ситуация показана в примере (2): буква *e* передается на русский язык как «е» во всех позициях, кроме начальной позиции в слове.

Контексты правил в приведенных примерах, особенно в примере (1), не полны, но это не является серьезной проблемой, так как они могут быть пополнены при анализе слогов (и успешно пополняются в текущей версии системы). Гораздо более серьезной проблемой являются пересекающиеся контексты в правилах с одинаковой входной строкой и разными выходными. Это может обозначать как недостаточность существующего контекста, так и неоднозначность в правилах чтения для данного языка. При анализе только фонетического уровня эта проблема не может быть решена.

Процедура нормализации контекста еще не формализована, поэтому она не включена в настоящую версию алгоритма.

5 Обсуждение метода

Метод в том виде, в котором он представлен в данной статье, имеет некоторые недостатки.

Неизвестно, универсален ли метод, так как еще не все этапы реализованы программно, а генерация правил по этому алгоритму вручную проводилась на примере французского и немецкого языков. Для некоторых языков, например, английского, метод может оказаться неприменимым. Впрочем, в случае английского языка даже составленные экспертом-лингвистом правила не всегда являются однозначными.

Надо отметить также низкую устойчивость метода к ошибкам в тестовых данных и исключениям из правил. Для каждого слога, не объясненного уже существующими правилами, порождается новое правило, все слоги имеют одинаковый вес. С другой стороны, это свойство позволяет генерировать адекватные системы правил с помощью сколь угодно малой тестовой выборки, при условии, что в ней отражены основные правила орфографии данного языка.

Механизм использования контекста в настоящее время также несовершенен. В правилах используются только «простые» контексты: составляются

множества символов, которые могут стоять на *n*-й позиции слева и справа от рассматриваемой цепочки. Однако иногда необходимо учитывать более сложные условия: закрытый слог, последний (первый) слог в слове и др.

Данный метод, в отличие от многих предшествующих, использует информацию о буквах, но весь анализ на данном этапе происходит исключительно на фонетическом уровне – то есть в слове не выделяются морфемы и не учитывается принадлежность буквы к той или иной морфеме. Отчасти это объясняется тем, что систему планируется использовать для машинной транскрипции фамильно именных групп произвольных языков, и зачастую единственной доступной информацией об этих языках будет список тестовых примеров (то есть слов с переводом). Но в некоторых случаях принадлежностью буквы к префиксу или суффиксу определяется, как она будет читаться, и невозможно сформировать соответствующее правило, рассматривая слово исключительно на фонетическом уровне.

Литература

- [1] Al-Onaizan Y., Knight K. Machine transliteration of names in arabic text// Proc. of the ACL Workshop on Computational Approaches to Semitic Languages, 2002.
- [2] Klyshinsky E., Maximov V., Yolkeen S. Cross-language transcription of proper names// Language Forum. – 2008. – V. 34, No 1. – P. 137-152.
- [3] Knight K., Graehl J. Machine transliteration// Computational Linguistics. – 1998. – V. 24, No 4. – P. 599-612.
- [4] Ravi S., Knight K. Learning phoneme mappings for transliteration without parallel data// Proc. of Human Language Technology Conf. The 2009 Annual Conf. of the North American Chapter of the Association for Computational Linguistics, 2009.
- [5] Sherif T., Kondrak G. Substring-based transliteration//Proc. of the ACL Workshop on Computational Approaches to Semitic Languages, 2007.
- [6] Tao T., Yoon S., Fister A., Sproat R., Zhai C. Unsupervised named entity transliteration using temporal and phonetic correlation// Proc. of EMNLP, 2006.

Self-learning system of machine transliteration using non-stochastic finite state automaton

V. Logacheva, E. Klyshinsky

The article describes a method of proper names transliteration that uses a combination of finite automaton and translation rules in the form of production. We introduce a method of automatic rules generation. The method uses a set of proper names and their transliterations. Then rules can be translated into finite automaton which conducts the transliteration method.

* Работа выполнена при частичной финансовой поддержке РФФИ (проект 10-01-00800)