

# ПРЕДСТАВЛЕНИЕ И ИНТЕРПРЕТАЦИЯ ТЕКСТА В ТЕХНОЛОГИИ СМЕШАННОГО НАБОРА

## TEXT REPRESENTATION AND INTERPRETATION IN A «MIXED SCRIPTS TECHNIQUE»

А. В. Коваленин

Издательство «Книжица», Новосибирск  
[kovalenin@mail.ru](mailto:kovalenin@mail.ru)

Для создания корпуса старинных певческих рукописей потребовалось найти способ работы с текстами, содержащими знаки крюкового письма, нот и церковнославянского языка. Найденный подход оказался применимым и для других письменностей, а его ключевой автомат — для некоторых задач прикладной лингвистики.

**1. Создание электронного фонда** старинных певческих рукописей («Фонда знаменных песнопений»<sup>1</sup>) стало возможным благодаря их набору в формате создаваемой «Технологии смешанного набора», поскольку этот формат позволяет:

– разделить содержание текста и детали для его интерпретации;

– отразить в тексте разные взгляды на то, что является его содержанием, то есть позволить каждой группе исследователей помещать в единицу хранения информацию, не нужную другим;

– использовать широкое понимание интерпретации, включающее многоступенчатые преобразования самого текста, опирающиеся на контекстно-обусловленные описания специального вида.

Последнее свойство не мешает Технологии смешанного набора быть открытой для подключения к процессу интерпретации любых внешних преобразователей. Однако пока удалось обойтись одним автоматом, описав необходимые для нашей практики преобразования в нужном ему виде, что позволило реализовать компактную рабочую версию инструментария (70 К в zip-архиве).

### 2. Недостатки непосредственного набора

Технология смешанного набора первоначально создавалась для работы с текстами древнерусских певческих рукописей, после того как стандартный подход (набор непосредственно шрифтами в привычном текстовом процессоре) оказался неудобным. Это неудобство объективно, сложность крюкового письма только помогла осознать необходимость альтернативы. Перечислю эти трудности.

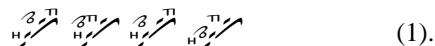
1. *При наборе.* а). Знак в крюковом письме («знамя») — составной объект, в котором на основной элемент накладываются до 5 дополнительных. Реализация в шрифте всех комбинаций отдельными символами громоздка и затрудняет выбор символа, а реализация накладных элементов символами, не сдвигающими курсор, затрудняет редактирование.

Последняя трудность становится критической, так как некоторые элементы традиционно отображаются другим цветом.

б). Многие накладные элементы нужно иметь в шрифте в нескольких вариантах. Это видно уже на простом тексте с ударениями — например, ударение, хорошо стоящее над **ш**, плохо встаёт над **л**. В крюковом письме помета свободно ставится рядом со знаменем; в нашем шрифте, например, помета **т** сочетается с крюком восемью способами:



Когда накладных элементов несколько, подбор взаимно подходящих сочетаний становится головоломным. Например, знамя *стрела мрачная борзая* с двумя пометами **н** и **т** может набираться разными способами, не все из которых удачны:



Возникает парадокс — чем лучше (с точки зрения полиграфической гибкости) продуман шрифт, тем труднее с ним работать.

в). Работа с разными шрифтами требует овладения их кодировками, различие которых даже для одной письменности неизбежно. Практически это означает, что интересные во многих отношениях шрифты остаются невостребованными.

2. *Работа с набранным текстом* попадает в зависимость от кодировки выбранных шрифтов. Сменить шрифт на шрифт другой кодировки невозможно. При использовании нескольких шрифтов даже простой контекстный поиск оказывается затруднён. При использовании в тексте нескольких письменностей все эти трудности умножаются, и добавляются проблемы смены раскладки клавиатуры.

Причина всех этих сложностей — концептуальное смешение процессов ввода, хранения и отображения информации, нарушение принципа «отдельные вещи делаются отдельно». Текст, по сути, неизбежно набирается в расчете на определенный внешний вид, хотя исследователю это может быть вовсе не нужно или нужно в другом виде.

<sup>1</sup> Фонд знаменных песнопений — <http://znamen.ru>.

В итоге, хотя наш первый опыт крюкового набора был также связан с разработкой специального шрифта, мы теперь не набираем и не рекомендуем другим набирать непосредственно шрифтом.

### 3. Примеры набора

Далее *текстом* будем называть просто последовательность знаков, предназначенную для последующей интерпретации читающим (человеком или автоматом). Примеры показывают возможный результат наиболее распространенного вида интерпретации текста — получения внешнего вида.

В примере 1 используются четыре разных письменности — знамена, нотно-линейное письмо, церковнославянские буквы, современные русские буквы; в примере 2 — буквы старославянского, русского (до 1918 г.), латинского и греческого алфавитов.

Пример 1 (фрагмент крюковой азбуки):

7. Чашка полна встречается только в восьмом глазе.

Пример 2 (фрагмент словаря Срезневского):

**ВЫЛЬ** — трава, herba: — земля... РАЖДАЖШТИ ВЪЛЬ ВЛАЖ. *Панѳ. Ант. XI в. 49. СВЕРЪНА ВО ЕСТЬ ПОХОТЬ АВЫ ДИКАА ВЪЛЬ О СОВЪ ВЪЗНИКШИ НА НЕДЪЛАНЪИ ЗЕМЛИ. Ган. Посл. Дм. п. 1078 г. — Ср. ВЪТИ, Гр. φύλλον, Lat. folium.*

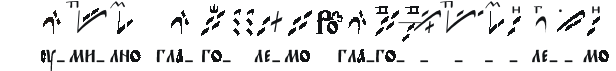
Пример 3 (для современного издания):



Пример 4 (как синодальное издание):



Пример 5 (приблизенно к оригиналу):



Примеры 3–5 — результаты различной интерпретации одного и того же *текста*, набранного из тропаря празднику Покрова. Можно видеть, что по выбору пользователя:

- а) текст песнопения представлен шрифтами *Irmologion Ucs* (пример 3) или *knijitsa* (4, 5), имеющими разные кодировки;
- б) оставлена крюковая запись мелодии (5), нотная (4) или обе (3), во втором случае удалён и потерявший значение фрагмент текста;
- в) формат нотной записи выбран современный (3) или старый (4);
- г-д) раздельное речение и ударения показаны (3, 4) или нет (5);
- е) подтекстовка дана в орфографии оригинала (5) или в нормативной орфографии «позднего извода» (3, 4).

### 4. Технология смешанного набора

§1. Текст набирается в простейшем текстовом редакторе; фрагменты каждой письменности изображаются посредством разработанного для неё представления, основанного на набираемых с клавиатуры символах. Например, в церковнославянском тексте буква ѧ изображается «я'», буква ѡ — «ја=», ѧ — «ъ»); в крюковом письме упомянутая стрела борзая (1) изображается «нпб\_./». Такое представление обычно легко читается, не заставляя исследователя разбираться в шрифтах и их поддержке. Для церковнославянского языка оно сделано на основе стандарта HIP<sup>2</sup>, для нот и крюков разработано самостоятельно.

§2. Каждый фрагмент, требующий особой интерпретации, помечается соответствующей этой особенности меткой. Метка — это только сигнал о наличии в *её фрагменте* особого содержания, а конкретизация этой особенности — способ интерпретации этого фрагмента задается отдельно в *стилевом файле*. Пример 1 — результат интерпретации следующего текста (подчёркнуты метки):

НОМ 7. <в> сч <род>Чашка <н>с1н1 <в>стЧ. <т>по'льная <н>с2н1г1 %объ Чашка полная встречается только в восьмом глазе.

То есть само понятие метки при этом не отличается от привычного понятия *стиля* (*tag*) известных текстовых процессоров. Отличие состоит в двух особенностях интерпретации меток (§§ 3, 4).

§3. Интерпретация по умолчанию (если она не задана в стилевом файле) — не игнорирование метки, как в HTML, а уничтожение следующего за ней фрагмента. Это позволяет вносить в текст любое количество независимых друг от друга разметок. Например, можно перед любой имеющейся в тексте меткой добавить свой фрагмент: «%кстати, к этому месту надо ещё вернуться». Если в стилевых файлах пользователей метка «%кстати» не предусмотрена, для них этот фрагмент будет несуществующим. Так, исключение нот или крюков в примерах 4 и 5 произведены просто удалением из стилевого файла описания интерпретации соответствующих меток.

Таким образом, в тексте может содержаться много слоёв информации, соответствующей различным аспектам исследования. Стилевой файл становится держателем *взгляда на текст* — представления о том, что в нём для пользователя значимо, а что будет игнорироваться.

Поскольку не возникает необходимости ограничить сосуществование различных разметок их взаимной вложенностью, постольку эти взгляды могут быть совершенно независимы.

§4. Интерпретация фрагмента состоит в задании необходимых *преобразований* помеченного фрагмента. Если, например, предполагается распечатка текста, то необходимы преобразователи пред-

<sup>2</sup> Стандарт HIP, разработка М.Гринчука. — <http://orthlib.ru/hip/>

ставлений используемых в нём письменностей в кодировку какого-нибудь шрифта. В инструменте Фонда знаменных песнопений используются и такие преобразования, как устранение надстрочных знаков, выбор варианта орфографии (между оригиналом и нормой), преобразование арабских чисел в славянские, применение буквицы, выключение двухцветности крюкового фрагмента.

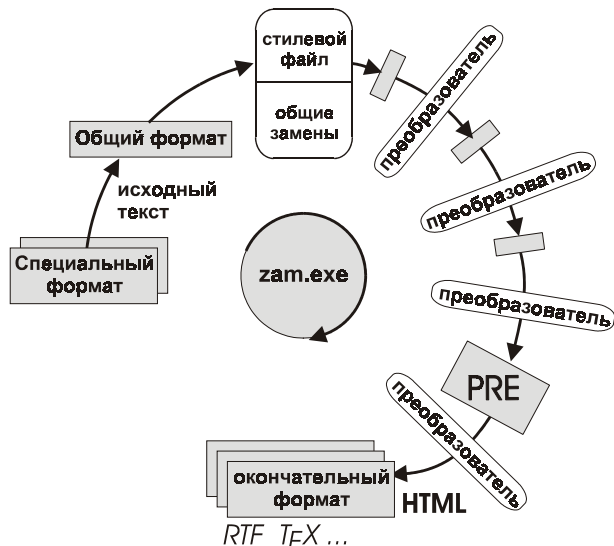


Рис. 1. Схема интерпретации текста

Интерпретация текста состоит в следующем: согласно стилевому файлу метки заменяются символическими цепочками, содержащими сигналы для включения других преобразователей. Преобразователи прочитывают текст и при обнаружении «своего сигнала» обрабатывают начинаемый им фрагмент. В результате этот фрагмент или остаётся в формате представления своей письменности (той же или другой), или превращается в цепочку кодов назначенного шрифта и управляющих символов выходного промежуточного языка. Если интерпретация производится ради получения внешнего вида, то производится преобразование из этого языка в выходной формат. Для действующей системы крюкового набора в качестве выходного формата пока реализован только HTML

### 5. Заменитель

Рассмотрим устройство основного автомата, реализующего цепочку имеющихся преобразований — как системных (начиная со стилевого файла, и кончая получением окончательного вида), так и перечисленных прикладных.

Объяснение его работы удобно начать с такой подзадачи морфологического разбора, как **выделение приставок**. Задача была решена автором в 1988 году в ВЦ СО АН СССР [1] в полной постановке. Здесь упрощенно рассмотрим её как задачу вставки в слово символа морфемной границы «\*». Анализ словаря показал, что для определения, например, приставки «анти» не обязательно держать в памяти

все словоформы на «анти-», а достаточно проверить три возможных продолжения слова. Это можно записать в виде четырёх строк:

$$\begin{array}{l} | \text{анти} | \text{анти}^* | \\ | \text{античн} | \text{античн} | \\ | \text{антиквар} | \text{антиквар} | \\ | \text{антикварк} | \text{анти}^* \text{кварк} | \end{array} \quad (2)$$

Каждая из них представляет собой **замену**, так как свою левую часть («заменяемое») предлагает заменить своей правой. Наш автомат («заменитель»), пробегая по тексту, сравнивает продолжение текста от текущей позиции на совпадение с левыми частями замен, и выполняет ту замену, для которой это совпадение наибольшее.

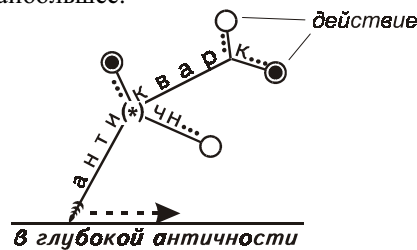


Рис. 2. Дерево замен в поиске места своего применения к тексту

На рисунке 2 заменяемые собраны в древовидную структуру, а возможные вставки на их место правых частей обозначены обобщающим словом «действие».

Эффективность этой операции зависит не от количества заменяемых (объёма полезной части словаря), а от числа ветвлений дерева на пути к «наиболее подходящей» замене, то есть при самом развесистом дереве ограничивается величиной, пропорциональной длине заменяемых. Это выгодно отличает такой механизм замен от привычного последовательного применения каждой отдельной замены ко всему тексту.

Вторым, решающим шагом в развитии заменителя явилось небольшое усложнение действия, выполняемого после замены — введение *сдвига*. Рассмотренной системе замен эквивалентна такая:

$$\begin{array}{l} | \text{анти} | \text{анти} (*) | - 3 \\ (*) | * | \\ (*) \text{чн} | \text{чн} | \\ (*) \text{квар} | \text{квар} | \\ (*) \text{кварк} | * \text{кварк} | \end{array} \quad (3)$$

Число  $-3$  («сдвиг») указывает, что после выполнения первой замены (если заменитель находился в позиции перед «анти») автомат сдвинется на три позиции назад, то есть встанет не после появившегося в тексте «(\*)», как если бы сдвига не было, а перед ним. Можно сказать, что автомат находится в состоянии «разбираем приставку *анти*», в котором уже будут учитываться только необходимые в этой ситуации замены.

Более простое применение сдвига — замена «| '20 '20 | '20 | -1», заменяющая любое количество стоящих подряд пробелов ('20) на один пробел.

Введение сдвига дало возможность такими заменами описывать сложные преобразования с движением по тексту не только вперёд, но и назад.

<b>Ёщик</b>	/EK}/[/-1
<b>3101</b>	/EKa/ЁEKa/[-1
<b>нормальная</b>	/EKa}/*a}/[-1
<b>форма</b>	/EKе/*е/[-1
	/EKи/ЁEKи/[-1
	/EKи}/*и}/[-1
/ {весел/ {весЁл*EK/-2	/EKил/*ил/[-1
/ {зелен/ {зелЁн*EK/-2	/EKит/*ит/[-1
/ {темл/ {тЁмлEK/-2	/EKо/*о/[-1
/ {далек/ {далЁкEK/-2	/EKо}/-o}/[-1
/ {тяжел/ {тяжЁлEK/-2	/EKо}/*о}/[-1
/ {желт/ {жЁлтEK/-2	/EKое/ЁEKое/[-1
/ {тепл/ {тЁплEK/-2	/EKой/ЁEKой/[-1
/ {легк/ {лЁгкEK/-2	/EKом/ЁEKом/[-1
/ {мертв/ {мЁртвEK/-2	/EKу/ЁEKу/[-1
/ {черн/ {чЁрнEK/-2	/EKы/ЁEKы/[-1
	/EKы}/*ы}/[-1
/*EK/[-1	/EKь/*ь/[-1
/ЁEK/[-1	/EKя/*я/[-1

Рис.3. Фрагмент системы замен для расстановки точек над «е». Фигурные скобки обозначают начало и конец слова, а «[-1» — переход после замены назад, на одну позицию раньше прежнего положения автомата.

В качестве примера можно рассмотреть систему замен для расстановки в тексте «ё». Целый ряд основ имеет общую, так сказать, «ё-парадигму» — набор продолжений слова, говорящих о наличии или отсутствии в этой основе «ё». Согласно заменам на рисунке 3 слева, заменитель вставляет в эти основы временную (прописную) Ё и сочетание ЕК, с которым он затем работает согласно заменам справа. Эти замены, в свою очередь, вставляют в слово знак «&», подтверждающий наличие «ё», или знак «\*», отрицающий его. После этого ещё одна группа замен перемещает этот знак назад к ранее вставленной Ё, чтобы сделать окончательный выбор.

### 6. Задачи о заменителе

Поиск по словарю эффективно реализуется через его представление в виде изображенного на рисунке 1 дерева. Тем самым словарь с уже расставленными нужными признаками (морфемными границами, диакритическими знаками, кодами грамматических значений...) даёт тривиальное решение задачи их расстановки. Но, как видно из систем замен (2), на 118 словарных слов [2] достаточно иметь только 4 замены, и это даёт не только сокращение перебора, но и решение задачи для большого множества несловарных слов. В то же время, на некоторых из них это решение может быть неправильным, что оставляет актуальным поддержание эталонного словаря и ставит задачу автоматического построения минимально достаточного дерева замен.

Например, для корпуса текстов на церковнославянском языке<sup>3</sup> интернет-сообщество «Славянская типографика»<sup>4</sup> созданы словники, которые по-

полняются по мере пополнения корпуса. Решение этой задачи позволило бы справиться с проблемой расстановки надстрочных знаков. При этом проблема разноударяемых слов решалась бы так же, только вместо перечня слов, использовался бы перечень более длинных контекстов, составление которого по готовому корпусу текстов не представляет технической сложности.

Вторая задача — обнаружение в построенном дереве одинаковых поддеревьев и перестройка дерева, аналогичная переходу от (2) к (3). Вставляемую при этом метку «(\*)» можно считать идентификатором некоторого словоизменительного класса: основы для изменения выписаны в строчках, ею заканчиваемых, а варианты изменения — в строчках, ею начинаемых.

За этими классами могут обнаруживаться реальные языковые явления. Но поскольку, как правило, конкретные основы представлены в словнике неполным набором возможных продолжений, интересно рассмотреть более сложный вариант второй задачи — когда для обнаруживаемых поддеревьев требуется не полное совпадение, но вложение друг в друга, а ещё лучше — когда поддерево, в которое могут вкладываться другие, реконструируется (восстановление полной парадигмы по неполным). Постановка такой задачи ещё нуждается в уточнении.

### 7. Добавление новой письменности

Показанный в системе замен (3) приём введения в текст «метки ситуации» позволяет так организовать систему замен, чтобы она работала только с фрагментами текста, начинающимися особым сигналом. Так создаются преобразователи для отдельных письменностей.

Допустим, мы хотим набрать подтекстовку песнопения японской слоговой азбукой хираганой. Покажем, как включить поддержку этой азбуки в систему. Сам японский текст напишем в удобной для нас транслитерации, предварив меткой-стилем:

<подтекст>шу авареме ё

В стилевом файле этой метке сопоставим выбор шрифта и сигнал для преобразователя хираганы:

```
|<подтекст>|<кг 16><гарн "Code2000">
<::япон hiragana-uni>|
```

Для составления преобразователя (системы замен) воспользуемся таблицей кодов хираганы в стандарте Unicode как показано на рисунке 4.

### 8. Состояние и перспективы проекта

Изложенная технология реализована в качестве рабочего инструмента для набора Фонда знаменых песнопений. Разработаны форматы для представления церковнославянских текстов, нот, знамен столбового и демественного распево, реализованы преобразователи из них в отдельные шрифты, энтузиастами ведётся набор единиц хранения фонда.

<sup>3</sup> «Библиотека святоотеческой литературы». — <http://orthlib.ru>

<sup>4</sup> «Славянская типографика». — <http://fonts.improvement.ru/>

Для этого был сформулирован специальный формат песнопения (рис.1), свободный от оформительских деталей.

В качестве заменителя в инструменте Фонда знаменных песнопений используется программа zam.exe, придуманная автором и написанная в 1994 году А.Такмаковым в виде DOS-приложения. Это обстоятельство позволяет пока говорить только о «скелетной» версии инструмента, в которой средства набора и интерпретации разделены.

Технология смешанного набора позволяет преодолеть эту разделённость и даже вернуться к удобствам режима немедленной визуализации набираемого (wysiwyg), но без описанных в п.2 недостатков. Для этого структура текста внутри текстового процессора может оставаться в формате этой технологии, а сама технология — использоваться для визуализации текста на экране. При этом отражение движения по тексту курсора будет производиться через описанный механизм замен.

--- типичное для всех регистров a


```
<::японhiragana-uni>|<::><jpon>|-6
<jpon>|<?><err>|[3 {! ::япон: непредусмотренный символ }
<jpon>|<jpon>|[3
<jpon><::>|<::>|-3
<jpon>'0D'0A'0D'0A|'0D'0A'0D'0A
<jpon>'0D'0A'0D'0A<jpon>|-6
```

--- собственно хирагана

```
<jpon>a|&#12354;<jpon>|-6
<jpon>me|&#12417;<jpon>|-6
<jpon>re|&#12428;<jpon>|-6
<jpon>su|&#12375;<jpon>|-6
<jpon>shu|&#12375;&#12421;<jpon>|-6
<jpon>wa|&#12431;<jpon>|-6
<jpon>yo|&#12424;<jpon>|-6
```

б СЪКТИНЪА ВЕЛІКАА

Источник: Обедница Калашникова, 1909 г.  
Демеством



	Hiragana Range: 3040–309F																																										
<pre>&lt;jpon&gt;a &amp;#12354;&lt;jpon&gt; -6 &lt;jpon&gt;wa &amp;#12431;&lt;jpon&gt; -6 &lt;jpon&gt;me &amp;#12417;&lt;jpon&gt; -6 &lt;jpon&gt;pe &amp;#12417;&lt;jpon&gt; -6 &lt;jpon&gt;cy &amp;#12417;&lt;jpon&gt; -6 &lt;jpon&gt;shu &amp;#12375;&amp;#12421;&lt;jpon&gt; -6 &lt;jpon&gt;e &amp;#12417;&lt;jpon&gt; -6</pre>	<p><b>3041</b> <span style="font-size: 2em;">6</span> <b>Hiragana</b> <span style="float: right;"><b>308C</b></span></p> <p><b>Based on JIS X 0208</b></p> <table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 10px;">3041</td> <td style="width: 10px;">ぁ</td> <td style="width: 10px;">HIRAGANA LETTER SMALL A</td> <td style="width: 10px;">3065</td> <td style="width: 10px;">づ</td> <td style="width: 10px;">HIRAGANA LETTER DU</td> <td style="width: 10px;">= ZU (not unique)</td> </tr> <tr> <td>3042</td> <td>ぁ</td> <td>HIRAGANA LETTER A</td> <td>3066</td> <td>て</td> <td>HIRAGANA LETTER TE</td> <td>= 3064 つ 3069 ち</td> </tr> <tr> <td>3043</td> <td>い</td> <td>HIRAGANA LETTER SMALL I</td> <td>3067</td> <td>で</td> <td>HIRAGANA LETTER DE</td> <td>= 3068 て 3069 ち</td> </tr> <tr> <td>3044</td> <td>い</td> <td>HIRAGANA LETTER I</td> <td>3068</td> <td>と</td> <td>HIRAGANA LETTER TO</td> <td>= 3068 と 3069 ち</td> </tr> <tr> <td>3045</td> <td>う</td> <td>HIRAGANA LETTER SMALL U</td> <td>3069</td> <td>ど</td> <td>HIRAGANA LETTER DO</td> <td>= 3068 と 3069 ち</td> </tr> <tr> <td>3046</td> <td>う</td> <td>HIRAGANA LETTER U</td> <td>306A</td> <td>な</td> <td>HIRAGANA LETTER NA</td> <td></td> </tr> </table>	3041	ぁ	HIRAGANA LETTER SMALL A	3065	づ	HIRAGANA LETTER DU	= ZU (not unique)	3042	ぁ	HIRAGANA LETTER A	3066	て	HIRAGANA LETTER TE	= 3064 つ 3069 ち	3043	い	HIRAGANA LETTER SMALL I	3067	で	HIRAGANA LETTER DE	= 3068 て 3069 ち	3044	い	HIRAGANA LETTER I	3068	と	HIRAGANA LETTER TO	= 3068 と 3069 ち	3045	う	HIRAGANA LETTER SMALL U	3069	ど	HIRAGANA LETTER DO	= 3068 と 3069 ち	3046	う	HIRAGANA LETTER U	306A	な	HIRAGANA LETTER NA	
3041	ぁ	HIRAGANA LETTER SMALL A	3065	づ	HIRAGANA LETTER DU	= ZU (not unique)																																					
3042	ぁ	HIRAGANA LETTER A	3066	て	HIRAGANA LETTER TE	= 3064 つ 3069 ち																																					
3043	い	HIRAGANA LETTER SMALL I	3067	で	HIRAGANA LETTER DE	= 3068 て 3069 ち																																					
3044	い	HIRAGANA LETTER I	3068	と	HIRAGANA LETTER TO	= 3068 と 3069 ち																																					
3045	う	HIRAGANA LETTER SMALL U	3069	ど	HIRAGANA LETTER DO	= 3068 と 3069 ち																																					
3046	う	HIRAGANA LETTER U	306A	な	HIRAGANA LETTER NA																																						

Рис.4. Подключение японского языка.

a — фрагмент преобразователя русской транслитерации японского текста в знаки хираганы (числа после «&#» — десятичные коды символов согласно Unicode); б — внешний вид песнопения с крюками, нотами и японской подтекстовкой; в — таблица хираганы в стандарте Unicode.

#### СПИСОК ЛИТЕРАТУРЫ

1. Коваленин А.В., Несговорова Г.П.. Несловарное описание исключений и автоматизация переноса слов // Третья всесоюзная конференция по созданию Машинного фонда русского языка. Тезисы докладов. Москва, 1989. Ч. I, С.164–166.
2. Зализняк А.А. Грамматический словарь русского языка. М.: Русский язык, 1977.