

# Использование конечных автоматов для решения компьютерно-лингвистических задач синтеза речи по тексту

*Гецевич Ю.С.,  
Скопинова Е.Н.,  
Окрут Т.И.*

В данной статье рассматривается общий подход к решению компьютерно-лингвистических задач в приложении к синтезу речи по тексту. Поэтапно описываются процессы решения задачи поиска и классифицирования количественных выражений с единицами измерения, а также задачи автоматизированного оформления диалогов с помощью синтаксических и морфологических грамматик в виде визуальных конечных автоматов, созданных при помощи программы NooJ, для текстового препроцессора синтезатора речи по тексту.

• компьютерно-лингвистическая задача • синтез речи по тексту  
• количественные выражения с единицами измерения • диалоги  
• конечные автоматы • NooJ • белорусский язык • русский язык

In this article there is considered the general approach to solve computer and Linguistic tasks within the appendix to speech synthesis by the text. Gradually it is described the processes of solving task to search and classify quantitative expressions with units of measurement, and also the tasks of automatic decoration the dialogues with synthetic and morphological grammar in a visual finite-state machines, that were made through the program NooJ, for the text pre-processor of speech synthesizer by the text.

*Computer-linguistic task • text-to-speech synthesis • quantitative expressions with measurement units • dialogues • finite automata • NooJ • Belarusian • Russian*

**Общий подход к решению компьютерно-лингвистических задач по электронным текстам**

Под *компьютерно-лингвистической задачей* по электронному тексту будем понимать такую задачу (проблему), которая, во-первых, ставится относительно электронного



текста; во-вторых, затрагивает вопросы конкретного поиска, классификации или переработки последовательностей электронных символов желаемого электронного текста; в-третьих, имеет конечное решение — компьютерную программу для предварительной обработки текста синтезатора речи, работа которой может быть проверена пользователем на неограниченном количестве других электронных текстов. Такое определение является обобщённым подходом к ряду задач, которые формулировались и решались в работах [1–3, 5]. Рассмотрим поэтапно процесс от постановки конкретной задачи к её решению с учётом компьютерных средств и условий. На рис. 1 видно, что данный процесс предусматривает осуществление шести этапов: от определения задачи к созданию продукта (её решения) для пользователя.

Таким образом, сначала перед исследователем ставится проблема (1). После этого экспериментальными путями (постановка и опровержение гипотез) эксперт находит решения данной проблемы относительно небольшого фрагмента текста (2), для этих решений разрабатывается рабочая алгоритмическая модель посредством конечных автоматов, создаваемых с помощью процессора NooJ (3) [4]. Затем разработанные конечные автоматы тестируются на большом количестве текстов (4). Если результаты тестирования согласно оценкам точности и полноты неудовлетворительные, то модель-алгоритм возвращается для доработки на стадию (2), в противном случае он передаётся для создания на его базе экспериментально-программного комплекса (5). При этом тестовые данные, разработанные на этапе (4), используются и на этапах (5) и (6) для точной разработки, тестирования и оценки программного продукта. При положительной оценке (в границах допустимой ошибки) программа попадает в руки к пользователю (7) и приобретает статус окончательного продукта, а при отрицательной — возвращается на стадию (5).

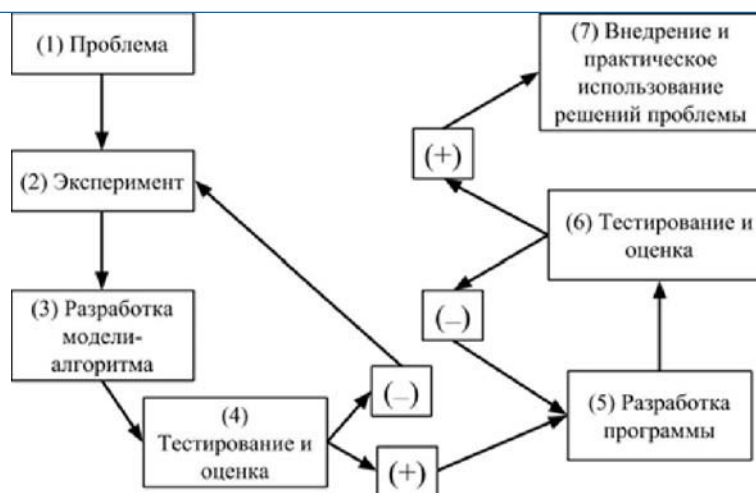


Рис. 1. Обобщенная схема процесса решения компьютерно-лингвистической задачи

Под *пользователем* будем понимать либо экспертов-лингвистов, либо прямых пользователей, которые могут использовать разработанные программные средства для предобработки текста для синтезатора речи. В подтверждение использования рассмотрим решение двух разных компьютерно-лингвистических задач в приложении к синтезу речи по тексту.

## Обработка количественных выражений с единицами измерения

Рассмотрим в качестве компьютерно-лингвистической задачи для синтеза речи проблему поиска и обработки количественных выражений с единицами измерения (КВЕИ) на материале белорусско- и русскоязычных электронных текстовых массивов. В общем виде КВЕИ представляют собой форму особо структурированной информации, которая выражается сочетанием количественного дескриптора (переданного на письме лингвистически с помощью слов либо математически посредством цифр) и обозначения единицы измерения (графического либо буквенного). В качестве примеров приведём следующие буквенно-символьные выражения: *100 МА, 3 тыс. лет назад, +212 °F, 6,022 141 29(27)·10<sup>23</sup> моль<sup>-1</sup>* и т.д. Количественные описания свойственны и общей научной картине мира, и бытовой сфере жизни. Помимо сложноструктурированности, вездесущности, они характеризуются и вариативностью форм содержания и выражения. Поэтому и возникает необходимость специально разработать полноценные алгоритмы и ресурсы для обработки КВЕИ в приложении к синтезу речи по тексту.

Преимуществом системы NooJ касательно задачи поиска и обработки КВЕИ является то, что разрабатываемые с помощью встроенного визуального редактора и представляющие собой конечные автоматы алгоритмы обладают наглядностью, благодаря чему их можно относительно легко корректировать и пополнять, что исключительно важно в силу вышеперечисленных свойств КВЕИ.

Следуя вышеописанной в разделе 1 схеме решения компьютерно-лингвистической проблемы, для начала чётко сформулируем задачу (1): найти и классифицировать в электронных текстах количественные выражения с единицами измерения. На этапе (2) экспертом проводились наблюдения и анализ КВЕИ относительно их строения и использования на материале электронных текстов научно-технического и правового тематического домена на белорусском и русском языках. К данному моменту на этапе (3) авторами смоделировано три взаимодополняющих алгоритмических комплекса для поиска количественных выражений с единицами измерения в больших корпусах текстов, которые позволяют:

- находить КВЕИ и классифицировать их по трём типам согласно международной системе единиц СИ (*основные, производные, внесистемные*) [2] (рис. 2);
- находить КВЕИ с метрологическими приставками (*кратными или дольными, сокращёнными или в полной форме*) и классифицировать их согласно словообразовательным особенностям [3] (рис. 3);
- преобразовывать КВЕИ в орфографические слова [1] (рис. 4).

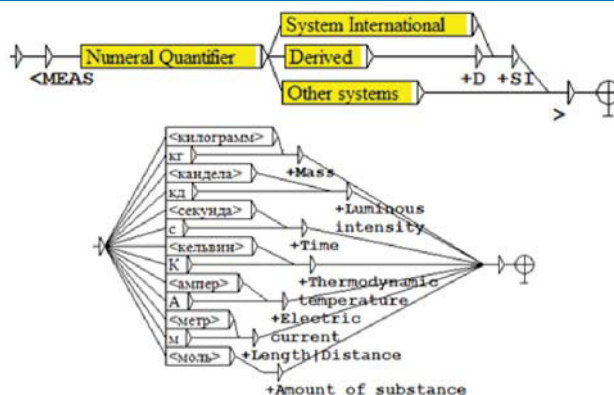


Рис. 2. Конечный автомат поиска и классифицирования КВЕИ согласно международной системе единиц СИ для белорусского и русского языков (NooJ)

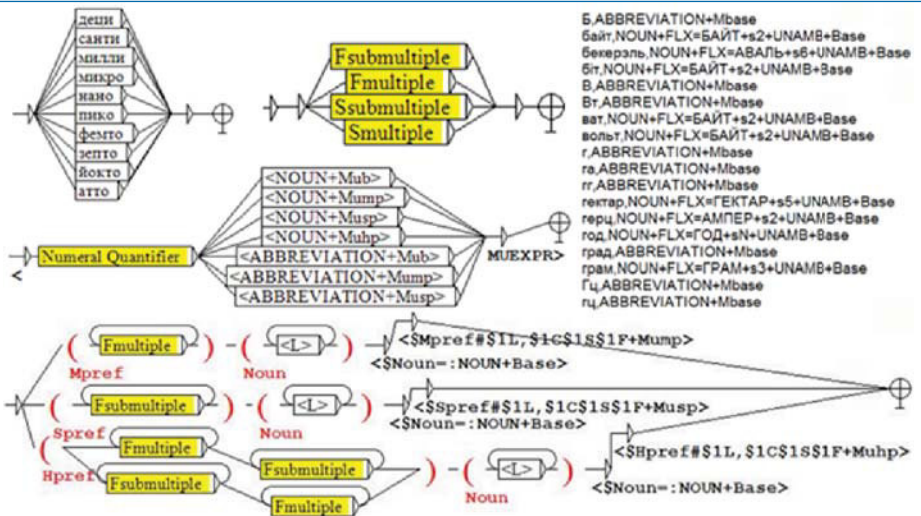


Рис. 3. Конечный автомат поиска и классифицирования КВЕИ согласно словообразовательным особенностям для белорусского и русского языков (NooJ)

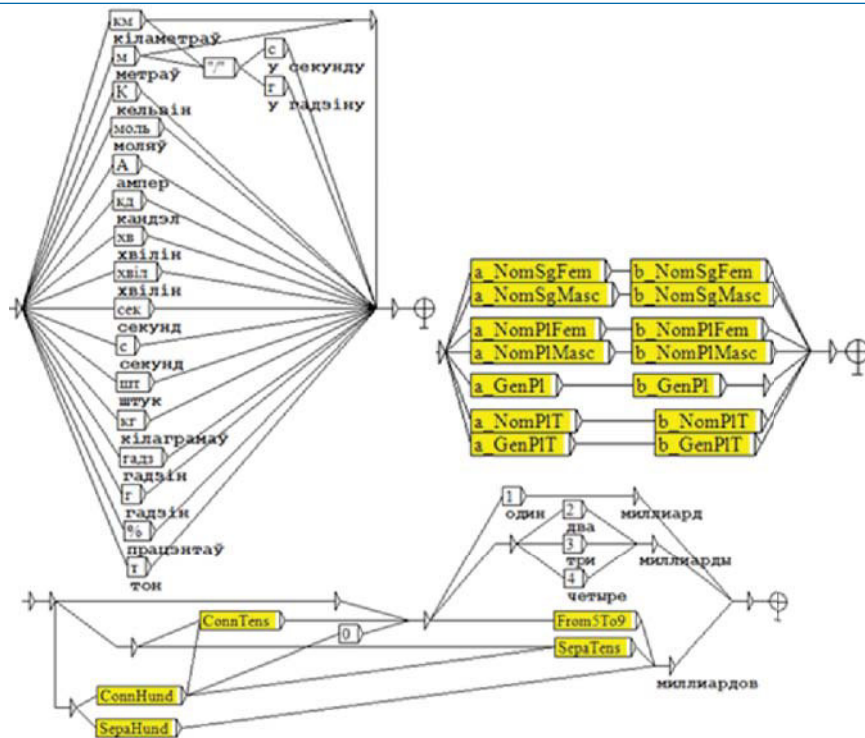


Рис. 4. Конечный автомат преобразования КВЕИ в орфографические слова (NooJ)

Этап тестирования (4) показал, что, например, первый алгоритмический комплекс даёт поисковые результаты с точностью в 72%. Этот относительно высокий показатель дал основание для разработки экспериментального программного комплекса (5), который бы находил КВЕИ и классифицировал их по трём типам согласно международной системе единиц СИ (основные, производные, внесистемные) [6]. Программа получила название «QEMU Identifier», т.е. «Quantitative Expressions with Measurement Units Identifier» или «Идентификатор количественных выражений с единицами измерения». Панель управления включает в себя вкладки:

меню (открытие файла, выход из программы), обработки (анализ текстов), настроек (выбор белорусско-, русско- либо англоязычного интерфейса) и справки (контактные данные разработчиков). Рассмотрим работу QEMU Identifier на примере анализа фрагмента русского-язычного текста из научно-технического текстового корпуса (рис. 5). Тексты для анализа можно вводить произвольно в верхнее поле либо открывать готовые тексты в формате txt через соответствующую команду во вкладке меню.

Команда «Преобразовать» (Transform) осуществляет поиск КВЕИ и присваивает им определённые маркеры. К примеру, на рисунке 6 было идентифицировано 40°C как <MEAS+Temperature in Celsius scale+D>. Это означает, что данный буквенно-символьный набор является количественным выражением с единицей измерения температуры по шкале Цельсия, причём по стандартам системы СИ данная единица является производной.

Команда «Найти» (Find) на выходе выдаёт список всех найденных КВЕИ и предоставляет количественные сведения (рис. 7). В данном тестовом тексте нашлось 31 количественное выражение с единицами измерения, из которых 8 выражений включают базовые мерные единицы СИ; 8 выражений имеют в своём составе единицы, производные от СИ; в оставшихся 15 выражениях присутствуют внесистемные единицы.

Осуществление команды «Фильтрация» (Filter) даёт возможность совершать разнообразные поисковые запросы через формальный язык регулярных выражений. Например, после ввода *Voltage/Frequency* список найденных количественных выражений ограничивается лишь теми, которые вмещают единицы измерения либо электрического напряжения, либо частоты тока (рис. 8).

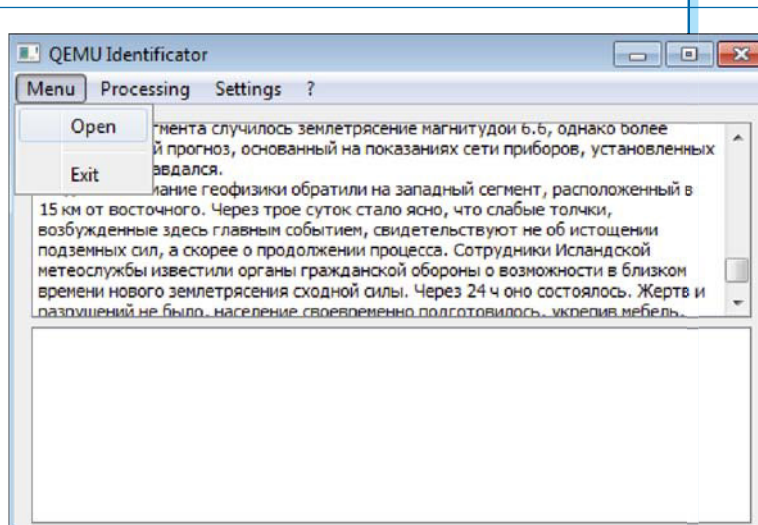


Рис. 5. Интерфейс экспериментальной программы QEMU Identifier

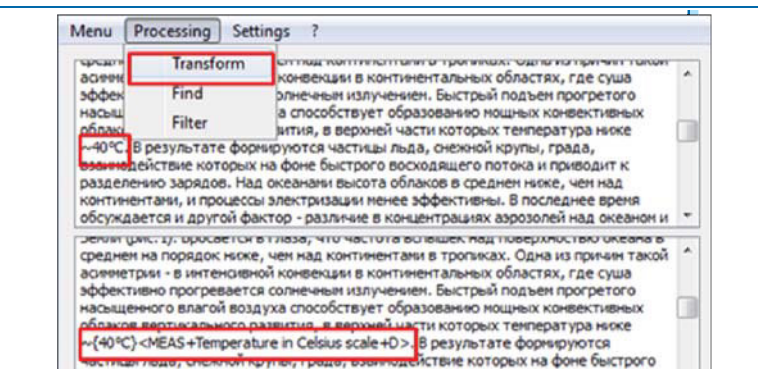


Рис. 6. Исполнение команды преобразования текста в программе QEMU Identifier

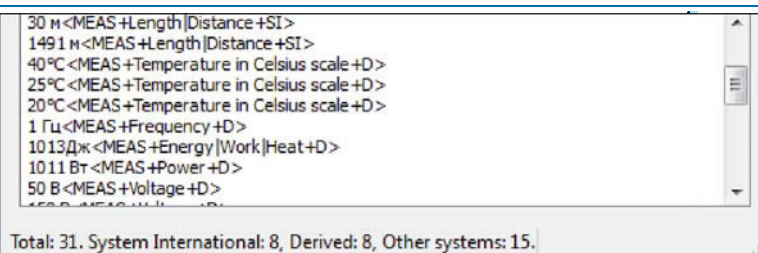


Рис. 7. Результаты поиска количественных выражений с единицами измерения в электронных текстах с помощью экспериментальной программы QEMU Identifier

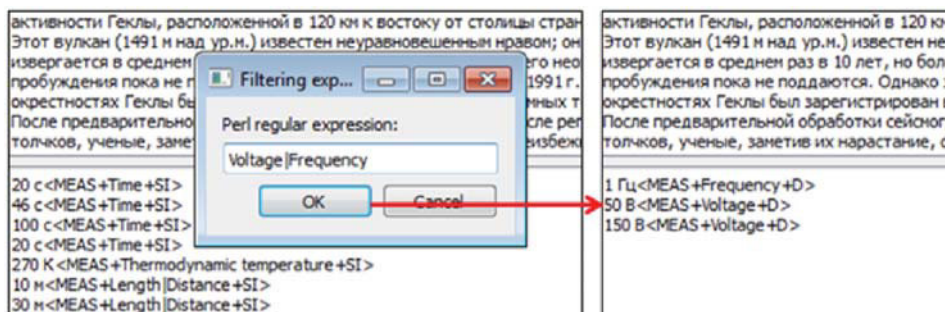
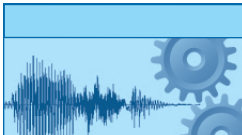


Рис. 8. Осуществление команды фильтрации списка найденных количественных выражений с единицами измерения на примере запроса единиц измерения либо электрического напряжения, либо частоты тока с помощью экспериментальной программы QEMU Identificator

### Идентификация диалогов и родовой принадлежности их участников

Благодаря встроенной в NooJ возможности разрабатывать конечные автоматы для синтаксических грамматик и аннотировать с их помощью электронный текст, стало возможным решение ещё одной компьютерно-лингвистической задачи — поиск и разметка диалогического текста в произведениях на белорусском языке для последующего применения в озвучивании различных аудиокниг. Обычно в существующих для этой цели программах используется одnogолосое озвучивание, однако зачастую электронные тексты состоят не только из слов автора, но и реплик различных персонажей. Таким образом, мы задались целью сделать синтез речи по тексту многоголосым.

Первичной задачей для решения данной проблемы была разработка алгоритмов для автоматического определения прямой речи в тексте. Для этого экспертами был выбран и проанализирован текстовый материал на предмет выявления структур оформления прямой речи (первый этап описанной ранее схемы процесса решения компьютерно-лингвистических задач). Сначала были обработаны первые четыре раздела произведения «Каласы пад сярпом тваім» В.С. Короткевича. Выявленные структуры были заложены в синтаксическую грамматику (посредством программы NooJ) и затем, по ходу последующего тщательного тестирования, были пополнены и скорректированы (второй этап схемы, более детально он описан в работе [7]).

Ниже приведён последний вариант набора структур оформления диалогического текста со следующими обозначениями: *M* — слова участника диалога (персонажа); *A* — слова автора; круглые скобки *()* — начало и завершение набора вариаций знаков препинания; вертикальная черта *|* — символ *или* разделение между знаками препинания в наборе их вариаций.

1. Слова персонажа без слов автора:  
— *M (! !!! !!!! !? !?! | ... |.)*.
2. Слова персонажа со словами автора в конце:  
— *M (, !! !!! !!!! !? !?! | ... |.) — A (... |.)*.
3. Слова персонажа с одной или несколькими авторскими вставками:

— M (, !! !!! !!!! !? !?! ! ... !.) — A (, / ... / : / ) —  
 M (, !! !!! !!!! !? !?! ! ... !.)  
 ( — A (, / ... / : / ) — M (, !! !!! !!!! !? !?! ! ... !.)).

Таким образом, разработанная грамматика DS\_All (рис. 9) срабатывает тогда, когда прямая речь начинается с тире, за которым идут слова персонажа; переход от слов автора и наоборот обозначается через комбинацию тире с запятой, точкой, восклицательным знаком, вопросительным знаком или их сочетаниями. Если после слов персонажа идут слова автора, грамматика продолжает поиск завершения фразы диалога.

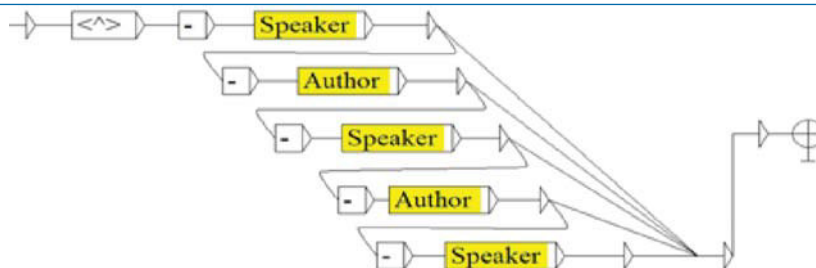


Рис. 9. Общий вид синтаксической грамматики DS\_All

Внутри подграфов Speaker (рис. 10 а) и Author (рис. 10 б) грамматики DS\_All отдельно рассматриваются внутренние знаки препинания. В диалогическом тексте может встретиться любая словоформа, любые числа, любые знаки препинания, кроме комбинаций с тире. Авторский текст имеет меньшую вариативность внутренних знаков препинания — запятая, точка, многоточие и также скобки.

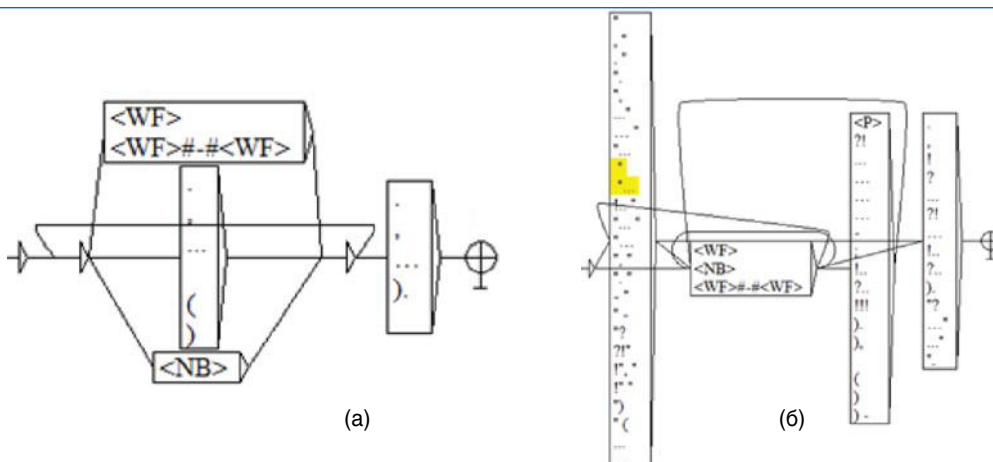


Рис. 10. Подграфы Speaker и Author грамматики DS\_All, где WF – любая словоформа, NB – любая последовательность чисел

Следующей задачей стояла идентификация рода персонажей. Наиболее полезными для её решения оказались вставки слов автора в прямую речь и непосредственные комментарии к словам участников диалога (это также обусловлено тем, что NooJ работает на уровне абзацев и не позволяет использовать связь между словами персонажей и словами автора из левого или правого контекста). Они включают такие слова-индикаторы рода, как глаголы прошедшего времени (казаў, казалa), имена собственные (Алесь, Майка) и существительные-индикаторы лиц (бацька, дзяўчынка). Так, из тренировочного текста были извлечены реплики с авторскими вставками, после чего размечены по роду (таблица 1).



Таблица 1

**Фрагмент разметки текста для выявления индикаторов рода персонажей с помощью анализа вставок слов автора**

Предложение	Индикатор рода	Род
Можа, не твая правіна? — сумеўся Загорскі. Можа, нехта другі?..	сумеў <u>ся</u> Загорскі	М
Накладайце сабе, пан Адам, — сказала маці... (...скарочана)	сказала <u>а</u> маці.	Ж

Индикаторы рода были выписаны отдельно и использованы для создания грамматик идентификации рода (рис. 11). Для этого в подграф Author из грамматики DS\_All были добавлены родозависимые подграфы VERBSfeminine (грамматика DS\_F) и VERBSmasculine (грамматика DS\_M) для определения женского и мужского рода соответственно (рис.12).

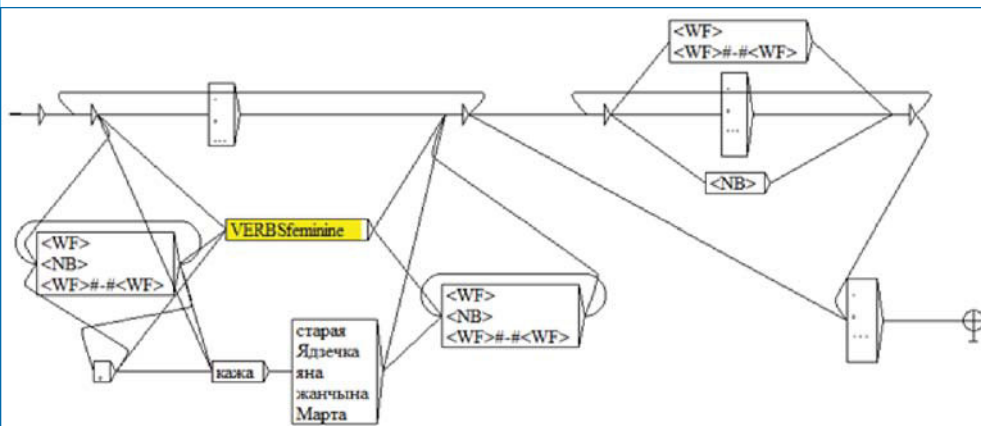


Рис. 11. Подграф Author грамматики DS\_F для определения женского рода участников диалога

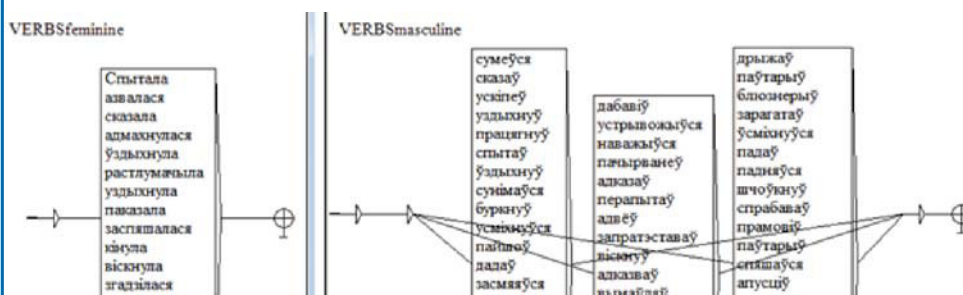


Рис. 12. Фрагменты подграфов VERBSfeminine и VERBSmasculine конечных автоматов грамматик DS\_F и DS\_M

На дальнейшем этапе по мере пополнения списка глаголов-индикаторов рода был создан отдельный словарь. В нём парами представлены глаголы прошедшего времени в формах женского и мужского рода. Для глаголов, которые начинаются на [y], была создана специальная парадигма ŸVERB1, которая учитывает переход на [j] после гласных (рис. 13).



```

трымаў, VERB+SpeechAct+Masculine
трымала, VERB+SpeechAct+Feminine
ударыў, VERB+SpeechAct+Masculine+FLX=ўVERB1
ударыла, VERB+SpeechAct+Feminine+FLX=ўVERB1
уздыгнуў, VERB+SpeechAct+Masculine+FLX=ўVERB1
уздыгнула, VERB+SpeechAct+Feminine+FLX=ўVERB1
    
```

Рис. 13. Фрагмент словаря глаголов-индикаторов рода (452 вхождения)

Рисунок 14 демонстрирует, как данный словарь подключён к грамматике: вместо подграфа VERBSfeminine теперь применяются специальные теги (категории) SpeechAct (семантическая помета для глаголов-комментариев) и Feminine, а для подграфа VERBSmasculine используются теги SpeechAct и Masculine.

В ходе разработки алгоритмов также была исследована проблема грамматических омоформ. Так, глагол *кажа* (форма настоящего времени глагола *казаць*) может относиться и к мужскому, и к женскому роду. Очевидно, подобного рода глагольные формы сами по себе не могут быть использованы для конкретизации родовой принадлежности участников диалога. Чтобы решить данную задачу, была создана дополнительная связка графов «глагол-существительное», где первый граф включает грамматические омоформы речевых глаголов, второй — список существительных-индикаторов рода персонажа (графы изображены на рис. 14).

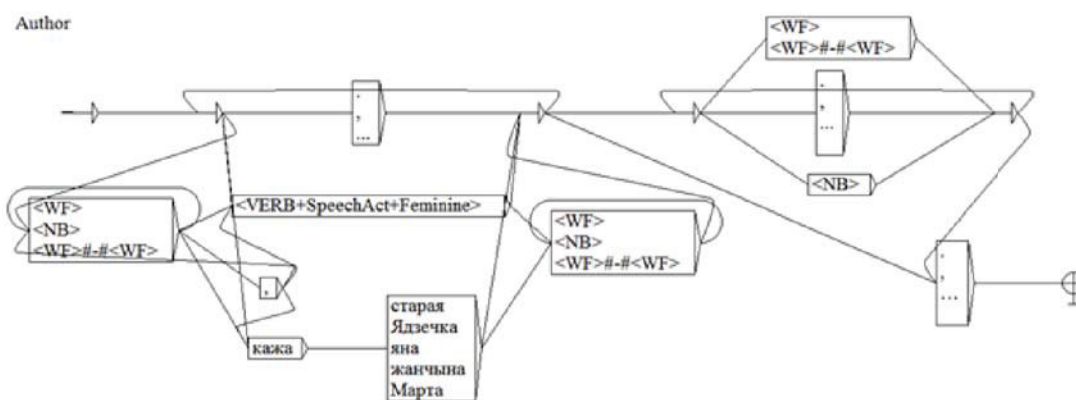


Рис. 14. Подграф Author грамматики DS\_F с подключённым словарём глаголов

Полученные грамматики могут применяться последовательно к желаемому электронному тексту посредством программы NooJ, причём пометы первой грамматики сохраняются при работе второй грамматики. Таким образом, на выходе текст получает разметку по репликам для мужского и женского голосов.

В качестве тренировочного материала для тестирования разработанных алгоритмов были использованы 32 раздела произведения «Каласы пад сярпом тваім» (свыше 100 000 словоупотреблений). Когда точность грамматик превысила 95%, был собран тестовый корпус из отрывков следующих произведений художественной литературы: «На ростанях» Я. Колоса, «У гарах дажджы» И.П. Мележа, «Жалезная кнопка» Л.И. Рублевской, «Асеннія лісты» Тётки. Тексты подбирались методом случайной выборки и в целом насчитывают около 24 000 словоупотреблений. По подсчётам эксперта полученный корпус содержит 481 реплику (реплики с тире), 165 реплик персонажей мужского рода и 68 реплик персонажей женского рода. Для оценки работы грамматик были вычислены значения точности (P), полноты (R) и их средней гармонической величины (таблица 2). При этом M — только те реплики, которые были правильно найдены и проанализированы грамматикой; L — реплики, которые были и правильно, и неправильно определены и проанализированы грамматикой; N — реплики, найденные и проанализированные экспертом (четвёртый этап схемы).



Таблица 2

**Оценка синтаксических грамматик, выявляющих предложения с прямой речью и определяющих род персонажей путём анализа вставок слов автора (на материале тестового корпуса текстов художественной литературы)**

Название грамматик	Точность (P)	Полнота (R)	Средняя гармоническая величина (F1-measure), %
	(M/L)	(M/N)	$2 \cdot P \cdot R / (P + R)$
DS_All	461/462 = 0,995	461/481 = 0,958	97,6
DS_M	143/145 = 0,986	143/165 = 0,866	92,2

Таким образом, разработанные алгоритмы идентифицируют прямую речь в тексте и определяют род персонажей по вставкам слов автора с точностью выше 90%, что позволяет перейти к следующему этапу — внедрение в систему синтеза речи по тексту (пятый этап схемы). Для того чтобы использование грамматик под стандарт SAPI 5.1 в синтезаторе речи по тексту стало возможным, необходимо привести тексты к виду SAPI TTS XML [8]. Для выбора необходимого голоса синтаксическая аннотация, которую генерирует грамматика, должна быть адаптирована под следующий код:

```
<VOICE Required="name=[Название голоса в системе]">
```

...Тексты для озвучивания...

```
</VOICE>
```

Для этого в DS\_M и DS\_F были вставлены маркеры обозначения путей, по которым срабатывают грамматики (рис. 15). Маркеры настроены так, что неразмеченный текст и слова автора озвучиваются голосом AlesiaBel, мужские реплики — синтезатором BorisBel, а женские — ElenaBel.

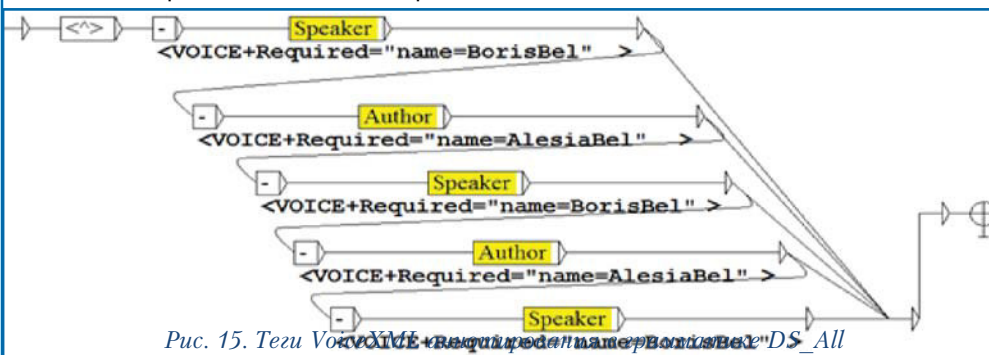


Рис. 15. Tezu VoiceXML-выходной код для DS\_All

Например, для предложений на рис. 16 оформление текста после обработки грамматиками будет иметь вид, демонстрируемый рис. 17.

Яна вагалася. Нават уздыхнула ў цемры. І тут ён пачуў нешта такое...  
 - Мож, і ёсца тут праўда... - амаль з вясковым прыдыханнем сказала яна. Зусім чыста па-мужыцку...  
 - Ты што ж... І гаварыць можаш? - спытаў ён. - Чаго ж прыкідвалася?  
 - Бацька са мною, калі не пры гасцях, заўсёды так гаворыць, - уздыхнула яна. - А прыкідвалася... так проста.  
 - Ну і брыда, - са злосцю сказаў ён. - Ідзі адсюль. Ну, чаго стаіш? Ідзі, кажу.  
 - Я ніколі больш не буду спяваць песню пра медзвездзяня, - сказала яна.

Рис. 16. Пример текста для аннотирования через VoiceXML теги

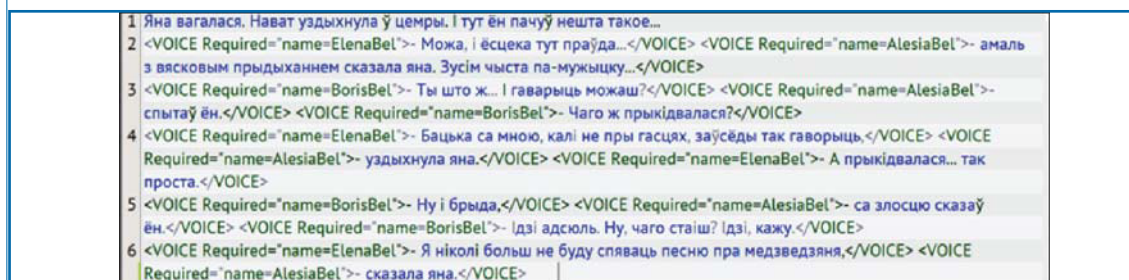


Рис. 17. Аннотирование через VoiceXML теги для автоматического переключения синтезаторов в

Затем размеченный текст можно подавать на вход системы синтеза речи по тексту. На рис. 18 изображено, как программа SAPI5 TTSAPP автоматически переключает поставленные в системе голоса AlesiaBel, BorisBel, ElenaBel при выбранной опции Process XML.

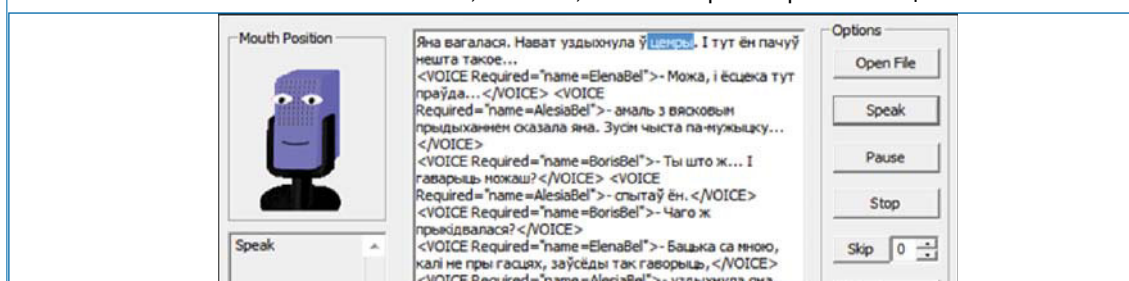


Рис. 18. Озвучивание автоматически размеченного грамматиками DS\_M и DS\_F текста тремя голосами синтезатором речи по тексту

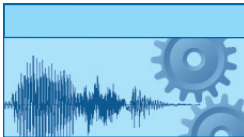
Разработанные модели-алгоритмы показали положительные результаты и в дальнейшем могут быть использованы для создания дополнительного блока автоматического выбора мужского или женского голоса системой синтеза речи по тексту, что предоставляет возможность быстрого многоголосого озвучивания электронных книг с сохранением индивидуальных особенностей персонажей того или иного произведения.

## Заключение

Таким образом, был предложен общий подход к решению компьютерно-лингвистических задач в приложении к синтезу речи по тексту. Для моделирования алгоритмов-решений использовался настраиваемый лингвистический процессор NooJ. В качестве примеров были приведены конечные автоматы, которые позволили эмулировать работу морфологических и синтаксических грамматик для решения двух разных компьютерно-лингвистических задач. В будущем авторами планируется повышение точности работы алгоритмов идентификации и генерирования орфографического текста по количественным выражениям с единицами измерения, а также разработка алгоритмов идентификации рода персонажей по репликам без вставок слов автора.

## Литература

1. Skopinava A.M., Hetsevlch Y.S., Lobanov B.M. Processing of quantitative expressions with units of measurement in scientific texts as applied to Belarusian and Russian text-to-speech synthesis // Компьютерная лингвистика и интеллектуальные технологии: материалы Междунар. конф. «Диалог», Московская обл., г. Бекасово, 29 мая — 2 июня 2013 г. Вып. 12 (19). В 2 т. Т.1. М.: Изд-во РГГУ, 2013. С. 634–651.



2. Гецевіч Ю.С., Скопінава А.М. Ідэнтыфікацыя выказаў з адзінкамі вымярэння ў навукова-тэхнічных і прававых тэкстах на беларускай і рускай мовах // Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2012): доклады XI Междунар. конф., Минск, 15 нояб. 2012 г. Минск: ОИПИ НАН Беларуси, 2012. С. 260–265.
3. Гецевіч Ю.С., Скопінава А.М. Кампаненты ідэнтыфікацыі колькасных выказаў з адзінкамі вымярэння ў тэкстах на беларускай і рускай мовах // Открытые семантические технологии проектирования интеллектуальных систем = Open Semantic Technologies for Intelligent Systems (OSTIS — 2013): материалы III Междунар. науч.-техн. конф., Минск, 21–23 февр. 2013 г. Минск: БГУИР, 2013 г. С. 319–328.
4. Лінгвістычны працэсар NooJ [Электронны рэсурс]. 2002. Рэжым доступу: <http://www.nooj4nlp.net/pages/nooj.html> — Дата доступу: 01.07.2012.
5. Гецевіч Ю.С. Аўтаматызаваная апрацоўка сімвальных выказаў у тэкстах для сістэмы сінтэзу беларускага маўлення // Информатика. 2011. № 4. С. 82–93.
6. Гецевіч Ю.С., Скопінава А.М., Есіс А.Ф. Мадэляванне і распрацоўка сістэм пошуку колькасных выказаў з адзінкамі вымярэння ў электронных тэкстах на беларускай і рускай мовах // Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2013): доклады XII Международной конференции (Минск, 20 ноября 2013 г.). Минск: ОИПИ НАН Беларуси, 2013. С. 282–287.
7. Гецевіч Ю.С., Окрут Т.И., Лабанаў Б.М. Аўтаматызацыя шматгаласавога стварэння аўдыёкніг на беларускай мове з дапамогай сінтэзатараў маўлення па тэксце // Развитие информатизации и государственной системы научно-технической информации (РИНТИ-2013): доклады XII Международной конференции (Минск, 20 ноября 2013 г.). Минск: ОИПИ НАН Беларуси, 2013. С. 269–276.
8. XML TTS Tutorial (SAPI 5.3) // Microsoft Developer Network [Electronic resource]. 2013. Mode of access: <http://msdn.microsoft.com/en-us/library/ms717077%28v=vs.85%29.aspx> — Date of access: 29.07.2013.

## Сведения об авторах

### **Гецевич Юрий Станиславович,**

кандидат технических наук, заведующий лаборатории распознавания и синтеза речи ОИПИ НАН Беларуси. Область научных интересов: речевые технологии, компьютерная лингвистика, естественно-языковые интерфейсы, робототехника. E-mail: [yury.hetsevich@gmail.com](mailto:yury.hetsevich@gmail.com)

### **Окрут Татьяна Ивановна,**

бакалавр гуманитарных наук, стажёр младшего научного сотрудника лаборатории распознавания и синтеза речи Объединенного института проблем информатики Национальной академии наук Беларуси, г. Минск. Область научных интересов: синтез речи по тексту, автоматическая обработка естественного языка, диалоговые системы, системы автоматического создания аудиокниг. E-mail: [tatberrie@gmail.com](mailto:tatberrie@gmail.com)

### **Скопинова Елена Николаевна,**

магистр филологических наук, младший научный сотрудник лаборатории распознавания и синтеза речи ОИПИ НАН Беларуси, г. Минск. Область научных интересов: синтез речи по тексту, автоматическая обработка естественного языка, обработка количественных выражений с единицами измерения. E-mail: [skelena777@gmail.com](mailto:skelena777@gmail.com)