

**Разработка правил генерации именных словоформ
для новописьменных вариантов карельского языка**

И. П. Новак

*Институт языка, литературы и истории
Карельского научного центра Российской академии наук,
г. Петрозаводск, Российская Федерация,
novak@krc.karelia.ru*

Н. Б. Крижановская

*Институт прикладных математических исследований
Карельского научного центра Российской академии наук,
г. Петрозаводск, Российская Федерация,
nataly@krc.karelia.ru*

Т. П. Бойко

*Институт языка, литературы и истории
Карельского научного центра Российской академии наук,
г. Петрозаводск, Российская Федерация,
boiko@krc.karelia.ru*

Н. А. Пеллинен

*Институт языка, литературы и истории
Карельского научного центра Российской академии наук,
г. Петрозаводск, Российская Федерация,
nataliapellinen@gmail.com*

АННОТАЦИЯ

Введение. Связывание слов текста (токенов) со значениями лемм в словаре корпуса ВепКар существенно облегчает дальнейшую работу по семантической разметке текстов. Для вепского подкорпуса ВепКар в 2019 г. были разработаны флективные правила, на их основе в корпус была добавлена функция генерации полной парадигмы по базовым словоформам.

При пополнении словарными статьями трёх подкорпусов карельского языка редакторам необходимо вводить большое число словоформ (около 30 для имён и 150 для глаголов), поэтому разработка алгоритма и компьютерной программы генерации словоформ карельского языка оказалась своевременной.

Цель: проиллюстрировать, как с помощью списка основ именных частей речи двух новописьменных наречий карельского языка можно составить правила для автоматической генерации словоформ.

Материалы исследования: леммы и словоформы из Открытого корпуса вепского и карельского языков, Корпуса Приграничной Карелии, электронной версии Словаря карельского языка.

Результаты и научная новизна. На основе изученных по теоретическим источникам и выявленных в ходе многолетних наблюдений грамматических закономерностей, а также проведённых в исследовании экспериментов сформирован список основ и псевдооснов именного словоизменения, разработана система правил генерации словоформ, написана и проверена соответствующая программа. Научная новизна исследования заключается во впервые принимаемой разработке системы единых правил автоматической генерации словоформ для двух наречий карельского языка.

Ключевые слова: карельский язык, новописьменный язык, корпусная лингвистика, морфология, именное словоизменение, генерация словоформ.

Благодарности: Исследование проведено в рамках выполнения государственного задания КарНЦ РАН. Раздел «Разработка программы генерации» подготовлен Н.Б. Крижановской в рамках проекта РФФИ 18-012-00117.

Для цитирования: Новак И. П., Крижановская Н. Б., Бойко Т. П., Пеллинен Н. А. Разработка правил генерации именных словоформ для новописьменных вариантов карельского языка // Вестник угроведения. 2020. Т. 10. № 4. С. 679–691.

**Development of rules of generation of nominal word forms
for new-written variants of the Karelian language**

I. P. Novak

*Institute of Linguistics, Literature and History,
Karelian Research Centre of the Russian Academy of Sciences,
Petrozavodsk, Russian Federation,
novak@krc.karelia.ru*

N. B. Krizhanovskaya

*Institute of Applied Mathematical Research,
Karelian Research Centre of the Russian Academy of Sciences,
Petrozavodsk, Russian Federation,
nataly@krc.karelia.ru*

T. P. Boyko

*Institute of Linguistics, Literature and History,
Karelian Research Centre of the Russian Academy of Sciences,
Petrozavodsk, Russian Federation,
boiko@krc.karelia.ru*

N. A. Pellinen

*Institute of Linguistics, Literature and History,
Karelian Research Centre of the Russian Academy of Sciences,
Petrozavodsk, Russian Federation,
nataliapellinen@gmail.com*

ABSTRACT

Introduction: linking of words of texts (tokens) with meanings of lemmas in the dictionary of VepKar corpus significantly facilitates further work on semantic markup of texts. In 2019, inflectional rules were developed for the Vepsian subcorpora VepKar. To the corpus on the base of these rules a function for generation of a complete paradigm on basic word forms was added.

VepKar editors need to enter a large number of word forms when they create dictionary entries in three Karelian subcorpora (about 30 for names and 150 for verbs). Therefore, the development of an algorithm and a computer program for generation of word forms of the Karelian language turned out to be timely.

Objective: to illustrate how you can use the list of the stems of the nominal parts of speech of two new-written dialects of the Karelian language to create rules for automatic generation of word forms.

Research materials: lemmas and word forms from the Open corpus of the Vepsian and Karelian languages, the Corpus of Border Karelia, and the electronic version of the Dictionary of the Karelian language.

Results and novelty of the research: grammatical patterns were studied over many years from theoretical sources, and they were also discovered through experiments. Thanks to this, the list of stems and pseudo-stems of word forms was formed for the nominal parts of speech, the system of rules for generation of word forms was developed, and the corresponding computer program is written and tested. The scientific novelty of the study lies in the first attempt to develop uniform rules for the automatic generation of word forms for two dialects of the Karelian language.

Key words: Karelian language, new-written language, corpus linguistics, morphology, nominal inflection, generation of word forms.

Acknowledgements: the study is carried out under the state order of the Karelian Research Centre of the Russian Academy of Sciences. The section «Development of the program of generation» was written by N. B. Krizhanovskaya in the framework of the project of the Russian Foundation for Basic Research No. 18-012-00117.

For citation: Novak I. P., Krizhanovskaya N. B., Boiko T. P., Pellinen N. A. Development of rules of generation of nominal word forms for new-written variants of the Karelian language // *Vestnik ugrovedenia* = *Bulletin of Ugric Studies*. 2020; 10 (4): 679–691.

Введение

Открытый корпус вепского и карельского языков (VepKar¹) является продолжением Корпуса вепского языка (рук. Н. Г. Зайцева), созданного в 2009 г. [16, 391]. Карельский язык был включён в него в 2016 г. Корпус представляет собой электронную информационно-справочную систему,

содержащую публицистические, художественные и диалектные тексты, разбитые на предложения и слова-токены. Часть текстов снабжена переводом на русский язык, выравненным с оригиналом по предложениям.

Важнейшей составляющей системы является словарь. Связывание слов текста со значениями

лемм в словаре существенно облегчает дальнейшую работу по их семантической разметке. Однако процесс по заполнению словаря продвигается очень медленно, что объясняется нехваткой специалистов и отсутствием необходимых для этого программ. Например, заполнение вручную именной словоизменительной парадигмы одной леммы (31 словоформа) занимает в среднем около 4–5 минут, а глагольной (149 словоформ) – 15–20 минут.

По карельскому языку вся работа в корпусе автоматически увеличивается в три раза, поскольку в нём представлено три карельских подкорпуса (рис. 1). Такое деление осуществлено в соответствии с тремя наречиями: собственно карельским, ливвиковским и людиковским, обнаруживающими между собой существенные различия на всех языковых уровнях [2, 21–27; 11, 57–58].

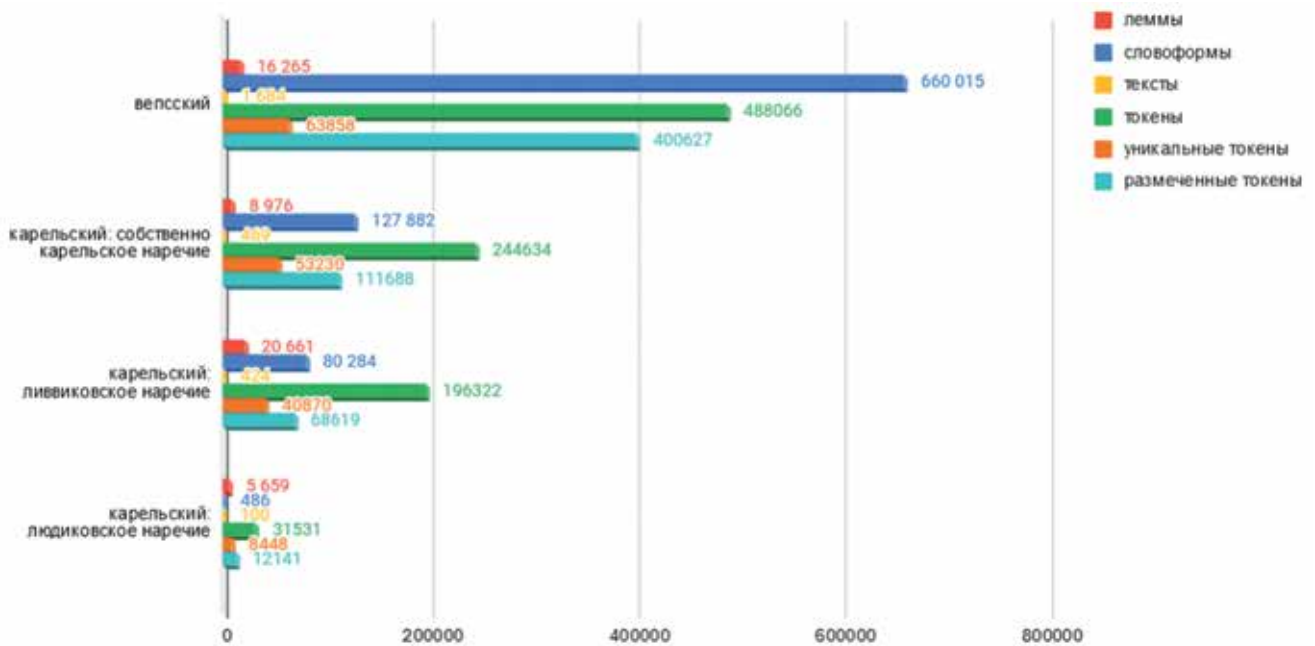


Рис. 1. Соотношение объёмов словарей и подкорпусов ВепКар (данные на 16.05.2020).

Для вепского языка в 2019 г. были разработаны флективные правила и в корпус добавлена функция генерации полной парадигмы по базовым словоформам, отсюда большая доля не только вепских словоформ, но и размеченных токенов (85%). Ускорить работу по разметке текстов и пополнению словарей карельских подкорпусов также можно при помощи программы автоматической генерации словоформ. В настоящем исследовании предполагается проанализировать возможности создания такой программы на примере именной словоизменительной системы новописьменных вариантов карельского языка Республики Карелия, разработанных на базе ливвиковского и собственно карельского наречий [4, 75–77].

Для новописьменных вариантов карельского языка определены орфографические нормы и грамматические правила, закреплённые в грамматиках П. М. Зайкова «Karielan kieliorpi» [17], «Vienankarjalan kieliorpi» [18] (собственно ка-

рельское наречие), Л. Ф. Маркиановой «Karjalan kieliorpi» [12] и Р. Пюёли «Livvin karjalan kieliorpi» [14] (ливвиковское наречие). В них приведены основные словоизменительные типы и парадигмы. Однако в процессе работы был выявлен ряд моментов, оставленных в грамматиках без внимания, что существенно осложнило разработку системы правил генерации и потребовало проведения дополнительных исследований.

В качестве теоретических источников были привлечены также справочник «Карельский язык в грамматиках» [2] и Академическая грамматика близкородственного карельскому финского языка [8].

В ходе исследования были использованы материалы Языкового банка Финляндии [7], в том числе Корпуса Приграничной Карелии (ок. 850 тыс. токенов) [13], основанного на записанных в 1960-х гг. образцах речи карелов Финляндии [10]. В качестве языковой базы для анализа карельской именной словоизменительной систем

¹ См. <http://dictorpus.krc.karelia.ru>

привлечены также данные электронной версии словаря карельского языка (88 тыс. словарных статей на диалектах собственно карельского и ливвиковского наречий) [9].

В процессе работы над программой был использован опыт создания электронных параллельных словарей уральских языков (например, людиковско-русско-финский словарь) [15], а также морфологического анализатора для родственного финскому квенского языка [5]. Кроме того, планируется использовать разработки программ морфологического анализа и морфологической генерации, в том числе для малоресурсных языков, разработанных и обученных на основе тщательно размеченных текстов Библии на английском языке и её переводов на сотни других языков [6]. Наличие переводов Библии на вепсский и карельский языки позволяет надеяться на возможность создания для них морфологических инструментов.

У коллектива редакторов корпуса ВепКар также имеется опыт работы над программой генерации словоформ. В 2019 г. такая программа была разработана для нормированного варианта карельского языка, развивающегося в Тверской области) [1]. Её применение позволило пополнить словарь более чем на 200 тыс. словоформ и увеличить разметку корпуса на 45%. В ходе работы предполагается переработать и распространить программу на собственно карельский и ливвиковский новописьменные варианты.

Материалы и методы

Исследование проведено на базе материалов словаря корпуса ВепКар на собственно карельском (около 9 тыс. лемм и 129 тыс. словоформ) и ливвиковском (около 21 тыс. лемм и 557 тыс. словоформ) новописьменных языках. В 2018 году в ВепКар вручную было добавлено 18 тыс. лемм с основными словоформами (ном., ген., парт. ед. и парт. мн. для имен) из «Большого карельско-русского словаря (ливвиковское наречие)» [1]. Наличие этих материалов позволило провести эксперименты по установлению закономерностей образования основы множественного числа¹.

В качестве базовых, наряду с экспериментами, выступили метод словообразовательного анализа, а также сравнительный и статистический методы, применение которых в совокупности позволило определить алгоритм разработки системы правил автоматической

генерации словоформ: 1) формирование списка грамматических форм; 2) выделение словоизменительных типов; 3) определение количества необходимых для генерации основ; 4) разработка системы правил генерации; 5) написание программы генерации и её проверка. В рамках статьи предлагается рассмотреть все этапы этого процесса на примере именной словоизменительной системы.

Результаты

1. Формирование парадигмы именного словоизменения

Карельский язык, как все прибалтийско-финские языки, является агглютинативным. Словоизменение в нём происходит путём присоединения аффиксов, содержащих грамматическое значение, к лексической основе слова [2, 163; 8, 53]. Последовательность аффиксов строго регламентирована: лексическая основа – показатель множественного числа – падежное окончание [14, 22; 18, 59]. Напр., *pien-i-h* от *pieni* ‘маленький’, где *piene-* – основа, *-i-* – суффикс множественного числа, *-h* – падежное окончание иллатива.

Склонение имён карельского языка происходит по падежам и числам. Категории притяжательности и сравнения в прибалтийско-финском языкознании рассматриваются как пограничные случаи словоизменения и словообразования [8, 62, 95, 296], в связи с чем они были изъяты из анализа.

Для карельского языка характерна бинарная структура категории числа. Форма единственного числа не маркирована (*hyvä-h tyttö-h* илл. ед. от *hyvä tyttö* ‘хорошая девочка’). Множественное число номинатива образуется при помощи окончания *-t*. Для образования форм множественного числа косвенных падежей употребляются показатели *-i-*, *-loi-* / *-löi-*, *-zi-* (*hyvi-h tyttö-löi-h* илл. мн.) [12, 34–36; 14, 26–29; 17, 48–49; 18, 71–72].

В диалектах карельского языка выделяют от 9 до 16 продуктивных падежей, поскольку система находится в процессе непрерывного становления. Для собственно карельского новописьменного варианта определено 14 продуктивных падежей [17, 55; 18, 75], а для ливвиковского – 16 [12, 37–38] / 17 [14, 30]. Большая часть падежей имеет единое окончание как в единственном, так и во множественном числе. Исключением является форма номинатива, особое внимание следует уделить также образованию форм партитива и генитива.

¹ См. http://dictorpus.krc.karelia.ru/ru/experiments/vowel_gradation

Парадигма именного словоизменения новописьменных вариантов карельского языка¹

Падеж	Одноосновное имя		Двуосновное имя	
	Ед. ч.	Мн. ч.	Ед. ч.	Мн. ч.
Ном.	tyttö 'девочка'	tytö-t	<i>lammaš/s</i> 'овца'	lamp/baha-t
Ген.	tytö-n	tyttö -j/löi-en/n	lamp/baha-n	lamp/bahj-en
Парт.	tytty -ö	tyttö -jä/löi	<i>lammaš/s-ta/tu</i>	lamp/bahj-e/i
Эсс.	tyttö -nä / tytö-nny	tyttö -löi-nä/nny	lamp/baha-na/nnu	lamp/bahj-na/nnu
Транс.	tytö-ksi/kse	tyttö -löi-ksi/kse	lamp/baha-kši/kse	lamp/bahj-ksi/kse
Инесс.	tytö-ssä/s	tyttö -löi-ssä/s	lamp/baha-šša/s	lamp/bahj-ssa/s
Элат.	tytö-štä/s(späi)	tyttö -löi-štä/s(späi)	lamp/baha-šta/s(späi)	lamp/bahj-sta/s(späi)
Илл.	tyttö -h	tyttö -löi-h	lamp/baha-h	lamp/bahj-h
Адесс.	tytö-llä/l	tyttö -löi-llä/l	lamp/baha-lla/l	lamp/bahj-lla/l
Абл.	tytö-ltä/l(päi)	tyttö -löi-ltä/l(päi)	lamp/baha-lta/l(päi)	lamp/bahj-lta/l(päi)
Алл.	- / tytö-le	- / tyttö -löi-le	lamp/baha-lla/le	lamp/bahj-lla/le
Абесс.	tytö-ttä/ttäh	tyttö -löi-ttä/ttäh	lamp/baha-tta/ttah	lamp/bahj-tta/ttah
Ком.	- / tytö-nke	tyttö -löi-neh/nneh(nke)	- / lamp/baha-nke	lamp/bahj-eh/nneh(nke)
Прол.	- / tytö-či	tyttö -löi-neh	- / lamp/baha-či	- / lamp/bahj-či
Инстр.	-	tyttö -löi-n	-	lamp/bahj-n

Представленные в таблице 1 словоизменительные парадигмы демонстрируют, что карельский язык далёк от агглютинативного идеала, поскольку в нём развились такие флективные способы словоизменения, как чередования конечных гласных основы (*tyttö-h* илл. ед. – *tytty-ö* парт. ед., *lamp/baha-h* илл. ед. – *lamp/bahj-h* илл. мн.) и консонантная альтернатива, на основании которой производится разделение основ на слабые и сильные (*tytö-t* ном. мн. – *tyttö-h* илл. ед.). Для языка характерно наличие как одноосновных, обладающих гласной основой, и двуосновных, обладающих гласной и согласной основами, имён (*tytö-n* ном. ед. – *tytty-ö* парт. ед., *lamp/baha-n* ген. ед. – *lammaš/s-ta/tu* парт. ед.). Следовательно, для выявления словоизменительных типов имён необходимо наличие информации о следующих словоформах: **словарная**, формы **генитива** (сл. гл. осн. одноосн. имён и сильн. гл. осн. двуосн.) и **пар-**

титива (сильн. гл. осн. с чередованием гласных одноосн. имён, согл. осн. двуосн.) **единственного числа**, форма **генитива множественного числа** (осн. мн. ч.).

2. Определение количества типов словоизменения имён

Словоизменительные типы имён карельского языка восходят к прибалтийско-финскому языку-основе, поэтому практически не имеют отличий на уровне наречий. Основной критерий, позволяющий отнести имя к тому или иному словоизменительному типу, – наличие или отсутствие у него согласной основы. Количество именных словоизменительных типов в нормативных грамматиках варьируется от 6 до 18 [12, 31–33; 14, 22–25; 17, 51–55; 18, 60–69]. В таблице 2 приведены примеры основных возможных типов имён, выявленных в материалах словарей корпуса ВепКар.

Таблица 2

Словоизменительные типы имён новописьменных вариантов карельского языка²

С.ф. на	Ном. ед.	Ген. ед.	Парт. ед.	Ген. мн.
Одноосновные имена				
-a, -ä / + -u, -y	ranta / randu 'берег'	ranna-n	rant/du-a	ranto-j-en / ranno-in
-a, -ä / + -u, -y	leipä / leiby 'хлеб'	leivä-n	leipy/leibi-ä	leip/bi-en
-u, -y	tukku 'куча'	tuku-n	tukku-o	tukku-j/loi-en/n
-o, -ö	pelt/do 'поле'	pello-n	pelt/du-o	pelt/do-j/loi-en/n

¹ Принятые обозначения: / – отделение собственно карельского от ливвиковского варианта (фонема, окончание, словоформа), обычный шрифт – слабая гласная основа, полужирный – сильная гласная основа, курсив – согласная основа, подчёркивание – изменения гласного компонента основы.

² Принятые обозначения: / – отделение собственно карельского от ливвиковского варианта, курсив – отличающиеся в наречиях фонемы, + – наличие второго возможного варианта, наряду с первым, в ливвиковском наречии, полужирный шрифт – чередования гласных и согласных (слабоступенный вариант).

-i	pappi 'священник'	papi-n	pappi-e/i	pappi-en / -loi-n
-i	lehti 'лист'	leh/hte-n	lehti-e	lehti-en / +-löi-n
VV	mua 'земля'	mua-n	mua-ta/du	mai-j-en / mua-loi-n
Двуосновные имена				
-hi, -li, -mi, -ni, -ri,	tuohi 'береста'	tuoh-e-n	tuoh-ta/tu	tuohi-en / +-loi-n
-si, -ši	nuori 'молодой'	nuore-n	nuor-ta/du	nuori-en
-si, -ši / -zi, -ži	ves/zi 'вода'	veje/ie-n	vet-tä/ty	ves/zi-en / +-löi-n
-ni / -ne	naini/e 'женщина'	nais/ze-n	nais-ta/tu	nais/zi-en
-Ce	vuate 'одежда'	vuattie-n	vuatet-ta/tu	vuattei-j/loi-en/n
-C	pereh 'семья'	perehe-n	pereh-tä/ty	perehi-en
-š / -s	armaš/s 'любимый'	armaha-n	armaš/s-ta/tu	armahi-en
-š, -s	kirveš/s 'топор'	kirvehe-n	kirveš/s-ta/ty	kirvehi-en
-š / -s	vereš/s 'свежий'	verekš/se-n	vereš/s-tä/ty	verekš/si-en
-š / -s, -t	vähyš/s 'недостаток'	vähyö-n	vähyt-tä/ty	vähyi-jen / vähyzi-en
	lyhyt 'короткий'	lyhyö-n	lyhyt-tä/ty	lyhyi-jen / lyhyzi-en

3. Определение основ для генерации словоформ

Важнейшим условием этого этапа (таблица 3) является стремление к использованию минимального набора исходных грамматических данных о лексеме. Поскольку словоформа далеко не всегда

содержит в себе однозначное указание на внешний вид гласной (слабой гласной) основы (*kirveš/s – kirvehe-*, *vereš/s – verekš/se-*), а также на отсутствие или наличие у лексемы согласной основы (*lehti – leh/hte-*; *pieni – piene-*, *pien-*), автоматическое их заполнение невозможно.

Таблица 3

Выделение основ для генерации словоформ¹

Осн.	Правила формирования основы	Пример
с.ф.	– и из словарного слова. С.с., в котором изменяемая часть отделяется от неизменяемой символом , а между составными частями сложного слова ставится .	ranta/du 'берег' pelt/do 'поле' mua 'земля' nuori 'молодой' lyhyt 'короткий' ves/zi 'вода'
о.1 (сл. / гл. осн.)	А. 1) часть с.с. до + первое в скобках; 2) если нет , то с.с. + первое в скобках. Б. Если в скобках два варианта п.о., то формируется о.1.1 (с левым вариантом п.о.) и о.1.2 (с правым вариантом п.о.). П.о. вносится в скобки: <i>mua []</i> , <i>ran ta [na]</i> / <i>ran du [na]</i> , <i>pel to [lo]</i> / <i>pel do [lo]</i> , <i>nuor i [e]</i> , <i>lyhy t [ö]</i> , <i>ve si [je/te]</i> / <i>v ezi [ie/ede]</i> .	А. ran+na= ranna pel+lo= pello mua nuor+e= nuore lyhy+ö= lyhyö Б. ve+je= veje , ve+te= vete / v+ie= vie , v+ede= vede
о.2 (сильн. / согл. осн.)	А. Если в скобках одно (нуль) п.о., то с.ф. – конечные V, VV → + конечные V, VV из о.1. Б. Если в скобках два п.о., то: 1) часть с.с. до + второе в скобках; 2) если нет , то с.с. + второе в скобках. Второе п.о. для двуосн. имен вносится в скобки через запятую от первого п.о.: <i>nuor i [e, -]</i> , <i>lyhy t [ö, t]</i> , <i>ve si [je/te, t]</i> / <i>v ezi [ie/ede, et]</i> .	А. rant/d+a= rant/da pelt/d+o= pelt/do mua Б. nuor+0= nuor lyhy+t= lyhyt ve+t/v+et= vet

¹ Принятые для таблиц 3 и 4 обозначения: || – отделение составных частей сложного слова, | – отделение неизменяемой части слова от изменяемой, – – убрать, + – добавить, = – равно, > – переход, / – отделение с.к. варианта от ливв., ~ – отделение заднерядного варианта показателя от переднерядного, ^ – отделение вариантов показателей, содержащих шипящие или свистящие согласные компоненты, V – гласный, C – согласный, → – переход к следующему действию, курсив в примерах – отличающиеся в наречиях фонемы, курсив в правилах – алгоритм выведения шаблона леммы.

<p>0.3 (вспом. осн. мн., сильн.)</p>	<p>А. Если о.2 заканч. на 1) Ci, Cu, Cy, а также Co, Cö (в <i>ливв.</i> всегда, в <i>с.к.</i> если в о.2 2 слога), то о.2+loi~löi; 2) с.к.: Co, Cö и в о.2 3 слога, то о.2+i (<i>karpalo+i</i>); 3) Ce, то e>i; 4) Ca, то a>i или a>oi¹; 5) Cä, то ä>i или ä>öi²; 6) Vi, то о.2+loi~löi (täi+löi); 7) VV (не Vi), то: с.к.: – первый V → +i; <i>ливв.:</i> о.2+loi~löi. Б. Если о.2 заканч. на C, то о.2>о.1 → если о.1 заканч. на 1) e и с.ф. заканч. на <i>si, ši, zi, ži</i>, то о.3=с.ф.; 2) Ce и с.ф. не заканч. на <i>si, ši, zi, ži</i>, то e>i; 3) Ca, Cä, то a, ä>i; 4) VV, то с.к.: <i>uo>ui, yö>yi, ie>ei</i> (<i>vuattie > vuattei</i>); <i>ливв.:</i> о.2+loi~löi (<i>vuattie+loi</i>), если о.1 на <i>uo, yö</i>, то – <i>o, ö</i> → +zi.</p>	<p>A. <i>pelt/do+loi=pelt/doloi</i> <i>rant/da>rant/doi</i> <i>mua-u=ma+i=mai /</i> <i>mua+loi=mualoi</i> Б. <i>nuore>nuori</i> <i>lyhyö>lyhyi /</i> <i>lyhy-ö= lyhy+zi=lyhyzi</i> <i>veje/vie>ves/zi</i></p>
<p>0.4 (вспом. осн. мн., сл.)</p>	<p>А. Если о.3 заканч. на <i>Coi~Cöi</i>, о.2 на <i>Ca, Cä</i> и кол-во слогов в о.3=о.2, то о.3>о.1 → если о.1 заканч. на 1) Ca, Cä, то a>oi, ä>öi; 2) ua, iä, то ua>avoi, iä>ävöi (<i>padoi > puu > pavo</i>). Б. Если о.3 заканч. на <i>loi~löi</i> и кол-во слогов в о.3 больше чем в о.2, то о.4=о.3. В. Если о.3 заканч. на <i>Vi</i> и о.1 на <i>VV</i>, то о.4=о.3. Г. Если о.3 заканч. на <i>Ci</i> и перед конечным <i>i</i> согл. <i>čč, šš, ss, k, t, p, g, d, b</i> 1) нет, то о.4=о.3; 2) есть, то о.3>о.1 → если о.1 заканч. на <i>CV</i>, то V>i; если о.1 заканч. на <i>ie>ei, ua>ai, iä>äi</i> (<i>regi>rie>rei</i>).</p>	<p>A. <i>rant/doi>ranna</i> >rannoi Б. pelt/doloi mualoi (ливв.) В. mai (с.к.) lyhyi (с.к.) Г. nuori lyhyzi (ливв.) ves/zi</p>

4. Разработка системы правил генерации словоформ

Прежде чем перейти к рассмотрению правил (таблица 4), важно обратить внимание на такую особенность фонетической системы карельского языка, как сингармонизм. Многие аффиксы представлены как в переднерядном, так и заднерядном вокалическом оформлении (*na~nä, nnu~nny, tta~ttä, ttah~ttäh*). В связи с этим для выбора варианта грамматического показателя вводится следующее дополнительное правило: 1) если в с.ф. есть гласные *a, o, u*, используется левый от ~ вариант показателя; 2) если в с.ф.

нет гласных *a, o, u*, используется правый от ~ вариант показателя. Для сложных слов программа должна ориентироваться на правую от || часть.

Кроме того, в собственно карельском новописьменном варианте важную роль играет дистрибуция свистящих и шипящих согласных, от которой также зависит выбор одного из парных вариантов падежных окончаний (*kši^ksi, šša~ššä^ssa~ssä*). В таком случае действует правило: 1) если о.1 не заканчивается на *i*, используется левый от ^ вариант; 2) если о.1 заканчивается на *i*, используется правый от ^ вариант.

Таблица 4

Система правил генерации словоформ именного словоизменения новописьменных вариантов карельского языка

Правило	Пример	Правило	Пример
Единственное число		Множественное число	
Ном.: с.ф.	<i>ranta/du</i> 'берег' <i>pelt/do</i> 'поле' <i>mua</i> 'земля' <i>nuogi</i> 'молодой' <i>lyhyt</i> 'короткий' <i>ves/zi</i> 'вода'	Ном.: о.1(о.1.1)+t	<i>ranna+t</i> <i>pello+t</i> <i>mua+t</i> <i>nuore+t</i> <i>lyhyö+t</i> <i>veje+t / vie+t</i>

¹ А. а > i если в о.2: 1) 2 слога и в первом слоге есть ц, о (в том числе в составе VV); 2) более 2-х слогов и о.2 заканч. на *mra / mba*, *ma* или является прил. с о.2. на *va*.

Б. а > oi если в о.2: 1) 2 слога и в первом слоге есть гласные а, е, i (в том числе в составе VV); 2) более 2-х слогов (кроме случаев А.2).

² А. ä > i если в о.2: 1) 2 слога; 2) более 2-х слогов и о.2 заканч. на *mä, vä, zä, sä, jä, pä, bä*.

Б. ä > öi если в о.2 более 2-х слогов (кроме случаев А.2).

<p>Ген.: о.1(о.1.1)+n</p>	<p>ranna+n pello+n mua+n nuore+n lyhyö+n veje+n / vie+n</p>	<p>Ген.: если о.3 заканч. на 1) Ci, то о.3+en; ливв.: 2) oi~öi, то о.3+n и о.4+n; с.к.: 3) loi~löi и кол-во слогов в о.3 больше чем в о.2, то loi~löi>jen; 4) Vi и кол-во слогов в о.3=о.2, а о.2 заканч. на – CV, то i>jen; – VV, то о.3+jen</p>	<p>rantoi>rantojen / randoi+n и rannoi+n peltoloi>peltojen / peldoloi+n mai+jen / mualoi+n nuori+en lyhyi/zi+jen/en ves/zi+en</p>
<p>Парт.: если о.2 заканч. на 1) CV, то a>ua, ä>yä/iä, u>uo, y>yö, o>uo, ö>yö, e>ie, i>ie/i; 2) VV, то о.2+ta~tä/ du~dy; 3) C, то о.2+ta~tä/du~dy (о.2 на l, n, r) или ty~ty (о.2 на h, s, t)</p>	<p>rant/dua pelt/duo mua+ta/du nuor+ta/du lyhyt+tä/ty vet+tä/ty</p>	<p>Парт.: если о.3 заканч. на 1) Ci, то о.3+e/i; ливв.: 2) oi~öi, то=о.3; с.к.: 3) loi~löi и кол-во слогов в о.3 больше чем в о.2, то loi~löi>ja~jä; 4) Vi и кол-во слогов в о.3=о.2 о, а о.2 заканч. на – CV, то i>ja~jä; – VV, то о.3+ta~tä</p>	<p>rantoi>rantoja / randoi peltoloi>peltoja / peldoloi mai+ta / mualoi nuori+e/i lyhyi/zi+tä/i ves/zi+e/i</p>
<p>Эсс.: с.к.: если о.2 заканч. на 1) V, то о.2+na~nä; 2) C, то о.1(о.1.2)+na~nä ливв.: о.1(о.1.1)+nnu~nny</p>	<p>ranta+na / ranna+nnu pelto+na / pello+nnu mua+na/nnu nuore+na/nnu perehe+nä/nny lyhyö+nä/nny vete+nä / vie+nny</p>	<p>Эсс.: о.3+na~nä / о.4+nnu~nny От о.3/о.4 образуются: <i>ком.</i> (с.к.): о.3+neh <i>ком.</i> (ливв.): о.4+nke(nneh) (о.4. на oi, öi) или enke(nneh) (о.4. на Ci)</p>	<p>rantoi+na / rannoi+nnu leipi+nä / leivi-nny pelt/doloi+na/nnu mai+na / mualoi+nnu nuori+na/nnu lyhyi/zi+nä/nny ves/zi+nä/nny</p>
<p>Транс.: о.1(о.1.1)+kši^ksi/ kse От о.1(о.1.1)+ образуются: <i>инесс.:</i> šša~ššä^ssa~ssä/s <i>элат.:</i> šta~štä^sta~stä/späi <i>адесс.:</i> lla~llä/l <i>абл.:</i> lta~ltä/lpäi <i>абесс.:</i> tta~ttä/ttah~ttäh <i>алл.</i> (ливв.): le <i>прол.</i> (ливв.): či <i>ком.</i> (ливв.): nke</p>	<p>ranna+t pello+t mua+t nuore+t lyhyö+t veje+t / vie+t</p>	<p>Транс.: о.4+ksi/kse От о.4+ образуются: <i>инесс.:</i> ssa~ssä/s <i>элат.:</i> sta~stä/späi <i>адесс.:</i> lla~llä/l <i>абл.:</i> lta~ltä/lpäi <i>абесс.:</i> tta~ttä/ttah~ttäh <i>инстр.:</i> n <i>алл.</i> (ливв.): le <i>прол.</i> (ливв.): či</p>	<p>rannoi+ksi/kse pelt/doloi+ksi/kse mai+ksi / mualoi+kse nuori+ksi/kse lyhyi/zi+ksi/kse ves/zi+ksi/kse</p>
<p>Илл.: если о.2 заканч. на 1) V, то о.2+h; 2) C, то о.1(о.1.2)+h</p>	<p>rant/da+h pelt/do+h mua+h nuore+h lyhyö+h vet/de+h</p>	<p>Илл.: о.3+h</p>	<p>rant/doi+h pelt/doloi+h mai/mualoi+h nuori+h lyhyi/zi+h ves/zi+h</p>

5. Разработка программы генерации

На основе полученных правил в комплексе программ Dictopus реализованы функции автоматической генерации словоформ по минимизированному шаблону для двух наречий карельского языка. Ре-

дактор словаря ВепКар в поле «лемма» вводит шаблон (mua [], ran|ta [na] / ran|du [na], nuori|i [e,]) и система автоматически создаёт 24 словоформы для собственно карельского или 31 – для ливвиковского наречия (рис. 2).

The screenshot shows the dictionary entry for 'aika' in the Veikko Karjalainen dictionary. The entry includes the language (Karelian), part of speech (noun), and various phonetic variants. A modal window titled 'Редактирование леммы: aika' is open, showing a form to add a new word form. The form includes fields for language, part of speech, lemma, dialect, and phonetic variants. A table on the right lists 24 grammatical forms of the word 'aika'.

No	грамматические признаки	Младописьменный севернокарельский (24)
Единственное число		
1.	номинатив	aika
2.	генитив	aikan
3.	партитив	aikua
4.	эссив	aikana
5.	транслатив	aikaksi
6.	абессив	aijatta
7.	инессив	aijašša
8.	аллатив	aijašta
9.	иллатив	aikah
10.	адессив-аллатив	aijaalla
11.	аблатив	aijaalta
12.	комитатив	aijaat
13.	инструктив	aikoien
14.	инфинитив	aikoja
15.	эссив	aikoina
16.	транслатив	aijoiksi
17.	абессив	aijoitta
18.	инессив	aijoissa
19.	аллатив	aijoista
20.	иллатив	aikoih
21.	адессив-аллатив	aijoilla
22.	аблатив	aijoilta
23.	комитатив	aikoineh
24.	инструктив	aijoin

Рис. 2. Добавление словоформ через шаблон леммы в словаре ВепКар.

Благодаря новым правилам в словарь ВепКар добавлено 600 тыс. ливвиковских словоформ, что увеличило количество размеченных токенов в текстах подкорпуса на 35%. В ближайшее время предстоит заполнение собственно карельской части словаря.

Программный код написан на языке PHP с использованием платформы Laravel и размещён в открытом доступе на сайте GitHub¹.

Обсуждение и заключение

В процессе исследования именной словоизменительной системы новописьменных вариантов карельского языка с целью создания программы автоматической генерации словоформ удалось выделить и унифицировать правила образования всех грамматических форм для имён основных словоизменительных типов, совпадающих в ливвиковском и собственно карельском наречиях. Количество минимальной грамматической информации об имени, которую необходимо вручную вносить в шаблоны для дальнейшего автоматического запол-

нения полной парадигмы слова, удалось сократить до нуля или одного псевдоокончания для одноосновных имён (*mua* [], *koivu* [], *uk|ko* [o]) и до двух или трёх – для двуосновных (*arma|š* [ha, š] / *arma|s* [ha, s], *ve|si* [jel|te, t] / *v|ezi* [ie|ede, et]). Достигнуть таких результатов помогли эксперименты, на основе которых для одноосновных имён с основой на *a-*, *ä-* были разработаны правила образования форм множественного числа.

В процессе проверки программы было выявлено отсутствие унификации в образовании формы множественного числа генитива в ливвиковском новописьменном варианте: в литературных текстах обнаружено два варианта образования этой формы (от сильно- и от слабоступенной основы). Для корпуса ВепКар было принято решение о включении обоих вариантов с целью увеличения процента разметки текстов. При этом отмечена крайняя необходимость доработки правила для нормированного варианта языка.

Предложенные в статье правила применимы не ко всем именам, а только к существительным,

¹ См. <https://github.com/componavt/dictorpus>

прилагательным, числительным и второстепенным местоимениям (пришедшим из других частей речи). Первостепенные по происхождению местоимения (личные, указательные, вопросительные) и их словоформы в процессе развития языка претерпели значительные изменения и приобрели целый ряд особенностей, что делает невозможным применение к ним программы генерации. Поскольку количество таких местоимений невелико, ручное заполнение их словоизменительных парадигм не представляет труда.

Автоматизировать процесс генерации именной словоизменительной парадигмы полностью, т. е. обучить программу распознавать словоизменительный тип имени и учитывать особенности

чередования ступеней согласных основы, не представляется возможным. Однако использование программы генерации именных словоформ в том виде, в котором она существует на настоящий момент, позволило сократить время заполнения одной именной парадигмы до 10–20 секунд или в среднем в 20 раз.

Проведённое исследование позволило существенно приблизиться к созданию приложения для проверки правописания карельского языка, а также подойти к разработке морфоанализатора и автоматического перевода. Правила образования форм множественного числа, которые удалось дополнить в ходе работы, найдут применение в практике преподавания карельского языка.

Список сокращений

абесс. – абессив	одноосн. – одноосновное имя
абл. – аблатив	осн. – основа
адесс. – адессив	п.о. – псевдоокончание
алл. – аллатив	парт. – партитив
вспом. осн. – вспомогательная основа	прил. – прилагательное
ед. – единственное число	прол. – пролатив
ген. – генитив	сильн. осн. – сильная основа
гл. осн. – гласная основа	с.к. – собственно карельское наречие
двуосн. – двуосновное имя	с.с. – словарное слово
илл. – иллатив	с.ф. – словарная форма
инесс. – инессив	сл. осн. – слабая основа
инстр. – инструктив	согл. осн. – согласная основа
ком. – комитатив	транс. – транслатив
ливв. – ливвиковское наречие	ч. – число
мн. – множественное	элат. – элатив
ном. – номинатив	эсс. – эссив

Список источников и литературы

1. Бойко Т. П. Большой карельско-русский словарь (ливвиковское наречие). Петрозаводск: Периодика. 2016. 352 с.
2. Карельский язык в грамматиках / Новак И., Пенттонен М., Руусканен А., Сиилин Л. Петрозаводск: КарНЦ РАН, 2019. 479 с.
3. Крижановская Н. Б., Новак И. П. New written Tver Karelian dialects word form generator. Свидетельство о гос. регистрации программы для ЭВМ Федеральной службы по интеллектуальной собственности № 2019665163 от 20 ноября 2019 г. URL: <https://new.fips.ru/publication-web/publications/document?type=doc&tab=PrEVM&id=6E578F3E-989F-44C6-A973-54863D199455> (дата обращения: 16.05.2020).
4. Нагурная С. В. Карельская письменность // Народы Карелии: историко-этнографические очерки. Петрозаводск: Периодика, 2019. С. 65–77.
5. A morphological analyzer for Kven / Trosterud S. R., Trosterud T., Räisänen A. K., Niiranen L., Haavisto M., Maliniemi K. // Proceedings of the Third Workshop on Computational Linguistics for Uralic Languages, 2017. Pp. 76–88.
6. Fine-grained Morphosyntactic Analysis and Generation Tools for More Than One Thousand Languages / Garrett N., Dylan L., McCarthy A. D., Mueller A., Wu W., Yarowsky D. // Proceedings of The 12th Language Resources and Evaluation Conference. 2020. Pp. 3963–3972.
7. Institute for the Languages of Finland. The Helsinki Korp Version of Samples of Spoken Finnish [text corpus]. Kielipankki. URL: <http://urn.fi/urn:nbn:fi:lb-2016042702> (дата обращения: 15.05.2020).

8. Iso suomen kielioppi / Hakulinen A., Vilkkuna M., Korhonen R., Koivisto V., Heinonen T.R., Alho I. Helsinki: Suomalaisen Kirjallisuuden Seura. 2008. Verkkoversio. URL: <http://scripta.kotus.fi/visk> (дата обращения: 28.05.2020).
9. Karjalan kielen verkkosanakirja / Toim. M. Torikka. URL: <http://kaino.kotus.fi/cgi-bin/kks/karjala.cgi> (дата обращения: 29.05.2020).
10. Kielikorpuksia Suomen itärajalta / Palander M., Riionheimo H., Kemppanen H., Mäkisalo J. // Multi lingual Finnic. Helsinki: Suomalais-Ugrilainen Seura, 2019. Pp. 425–438.
11. Koivisto V. Border Karelian dialects – a diffuse variety of Karelian // On the border of language and dialect. Helsinki: Suomalaisen Kirjallisuuden Seura, 2018. Pp. 56–84.
12. Markianova L. Karjalan kielioppi. Petroskoi: Periodika, 2002. 296 s.
13. Palander M., Riionheimo H., Koivisto V. Kotimaisten kielten keskus, Institute for the Languages of Finland. Raja-Karjalan korpus. 2017. URL: <http://urn.fi/urn:nbn:fi:lb-2014073033> (дата обращения: 28.05.2020).
14. Pyöli R. Livvinkarjalan kielioppi. Helsinki: Karjalan Kielen Seura, 2011. 200 p.
15. Rueter J., Hämäläinen M. Synchronized Mediawiki based analyzer dictionary development // Proceedings of the Third Workshop on Computational Linguistics for Uralic Languages. 2017. Pp. 1–7.
16. Zaiceva N. Veps language heritage in Karelia // Multi lingual Finnic. Helsinki: Suomalais-Ugrilainen Seura, 2019. Pp. 379–400.
17. Zaikov P. Karjalan kielioppi. Petroskoi: Periodika, 2002. 208 p.
18. Zaikov P. Vienankarjalan kielioppi. Helsinki: Karjalan Sivistysseura, 2013. 284 p.

References

1. Boyko T. P. *Bol'shoj karel'sko-russkij slovar'* [Large Russian-Karelian dictionary]. Petrozavodsk: Periodika Publ., 2016. 352 p. (In Russian)
2. Novak I., Penttonen M., Ruuskanen A., Siilin L. *Karel'skij yazyk v grammatikakh* [Karelian language in grammars]. Petrozavodsk: KarRC RAS Publ., 2019. 479 p. (In Russian)
3. Krizhanovskaya N. B., Novak I. P. New written Tver Karelian dialects word form generator. *Svidetel'stvo o gos. registratsii programmy dlya EVM Federal'noy sluzhby po intellektual'noy sobstvennosti № 2019665163 ot 20 noyabrya 2019 g.* [Certificate of state registration of the program, № 2019665163, November 20, 2019]. Available at: <https://new.fips.ru/publication-web/publications/document?type=doc&tab=PrEVM&id=6E578F3E-989F-44C6-A973-54863D199455> (accessed May 16, 2020). (In Russian)
4. Nagurnaya S. V. *Karel'skaya pis'mennost'* [Karelian writing]. *Narody Karelii: istoriko-etnograficheskie ocherki* [Peoples of Karelia: historical-ethnographic essays]. Petrozavodsk: Periodika Publ., 2019. pp. 65–77. (In Russian)
5. A morphological analyzer for Kven. Trosterud S. R., Trosterud T., Räisänen A. K., Niiranen L., Haavisto M., Maliniemi K. *Proceedings of the Third Workshop on Computational Linguistics for Uralic Languages*, 2017, pp. 76–88. (In English)
6. Fine-grained Morphosyntactic Analysis and Generation Tools for More Than One Thousand Languages. Garrett N., Dylan L., McCarthy A. D., Mueller A., Wu W., Yarowsky D. *Proceedings of The 12th Language Resources and Evaluation Conference*, 2020, pp. 3963–3972. (In English)
7. Institute for the Languages of Finland. The Helsinki Korp Version of Samples of Spoken Finnish [text corpus]. Kielipankki. Available at: <http://urn.fi/urn:nbn:fi:lb-2016042702> (accessed May 15, 2020). (In English)
8. Hakulinen A., Vilkkuna M., Korhonen R., Koivisto V., Heinonen T.R., Alho I. *Iso suomen kielioppi*. Helsinki: Suomalaisen Kirjallisuuden Seura. 2008. Available at: <http://scripta.kotus.fi/visk> (accessed May 28, 2020). (In Finnish)
9. *Karjalan kielen verkkosanakirja*. Ed. by M. Torikka. Available at: <http://kaino.kotus.fi/cgi-bin/kks/karjala.cgi> (accessed May 29, 2020). (In Finnish)
10. *Kielikorpuksia Suomen itärajalta*. Palander M., Riionheimo H., Kemppanen H., Mäkisalo J. Multi lingual Finnic. Helsinki: Suomalais-Ugrilainen Seura, 2019. pp. 425–438. (In Finnish)
11. Palander M., Riionheimo H., Kemppanen H., Mäkisalo J. Multi lingual Finnic *Kielikorpuksia Suomen itärajalta*. Helsinki: Suomalais-Ugrilainen Seura, 2019. pp. 425–438. (In Finnish)
12. Markianova L. Karjalan kielioppi. Petrozavodsk: Periodika, 2002. 296 p. (In Karelian)
13. Palander M., Riionheimo H., Koivisto V. *Kotimaisten kielten keskus, Institute for the Languages of Finland. Raja-Karjalan korpus*. 2017. Available at: <http://urn.fi/urn:nbn:fi:lb-2014073033> (accessed May 28, 2020). (In Finnish)
14. Pyöli R. *Livvinkarjalan kielioppi*. Helsinki: Karjalan Kielen Seura, 2011. 200 p. (In Finnish)

15. Rueter J., Hämmäläinen M. Synchronized Mediawiki based analyzer dictionary development. *Proceedings of the Third Workshop on Computational Linguistics for Uralic Languages*, 2017. pp. 1–7. (In English)
16. Zaiceva N. *Veps language heritage in Karelia*. Multi lingual Finnic. Helsinki: Suomalais-Ugrilainen Seura, 2019. pp. 379–400. (In English)
17. Zaikov P. *Karjalan kieliooppi*. Petrozavodsk: Periodika, 2002. 208 p. (In Karelian)
18. Zaikov P. *Vienankarjalan kieliooppi*. Helsinki: Karjalan Sivistysseura, 2013. 284 p. (In Finnish)

ИНФОРМАЦИЯ ОБ АВТОРАХ

Новак Ирина Петровна, научный сотрудник сектора языкознания Института языка, литературы и истории КарНЦ РАН, ФИЦ «Карельский научный центр» (185910, Российская Федерация, Республика Карелия, г. Петрозаводск, ул. Пушкинская, д. 11), кандидат филологических наук.

novak@krc.karelia.ru

ORCID.ID: 0000-0002-9436-9460

Крижановская Наталья Борисовна, ведущий инженер-исследователь лаборатории информационных компьютерных технологий Института прикладных математических исследований КарНЦ РАН, ФИЦ «Карельский научный центр» (185910, Российская Федерация, Республика Карелия, г. Петрозаводск, ул. Пушкинская, д. 11).

nataly@krc.karelia.ru

ORCID ID: 0000-0002-9948-1910

Бойко Татьяна Петровна, научный сотрудник сектора языкознания Института языка, литературы и истории КарНЦ РАН, ФИЦ «Карельский научный центр» (185910, Российская Федерация, Республика Карелия, г. Петрозаводск, ул. Пушкинская, д. 11).

boiko@krc.karelia.ru

ORCID ID: 0000-0001-5095-2921

Пеллинен Наталия Александровна, младший научный сотрудник сектора языкознания Института языка, литературы и истории КарНЦ РАН, ФИЦ «Карельский научный центр» (185910, Российская Федерация, Республика Карелия, г. Петрозаводск, ул. Пушкинская, д. 11), кандидат филологических наук.

nataliapellinen@gmail.com

ORCID ID: 0000-0002-5648-6877

ABOUT THE AUTHORS

Novak Irina Petrovna, Researcher, Department of Linguistics, Institute of Linguistics, Literature and History, Karelian Research Centre of the RAS (185910, Russian Federation, Republic of Karelia, Petrozavodsk, Pushkinskaya st., 11), Candidate of Philological Sciences.

novak@krc.karelia.ru

ORCID ID: 0000-0002-9436-9460

Krizhanovskaya Natalya Borisovna, Leading Engineer Research, Laboratory for Information Computer Technologies, Institute of Applied Mathematical Research, Karelian Research Center of the RAS (185910, Russian Federation, Republic of Karelia, Petrozavodsk, Pushkinskaya st., 11).

nataly@krc.karelia.ru

ORCID ID: 0000-0002-9948-1910

Boiko Tatyana Petrovna, Researcher, Department of Linguistics, Institute of Linguistics, Literature and History, Karelian Research Centre of the RAS (185910, Russian Federation, Republic of Karelia, Petrozavodsk, Pushkinskaya st., 11).

boiko@krc.karelia.ru

ORCID ID: 0000-0001-5095-2921

Pellinen Nataliya Aleksandrovna, Junior Researcher, Department of Linguistics, Institute of Linguistics, Literature and History, Karelian Research Centre of the RAS (185910, Russian Federation, Republic of Karelia, Petrozavodsk, Pushkinskaya st., 11), Candidate of Philological Sciences.

naliapellinen@gmail.com

ORCID ID: 0000-0002-5648-6877