

<p>Номер и название направления Программы Корпусная лингвистика. Создание и развитие корпусных ресурсов по языкам народов России</p>	
<p>Название проекта “Создание корпуса текстов республиканских газет на башкирском языке”</p>	
<p>Научный руководитель проекта Хисамитдинова Фирдаус Гильмитдиновна, д.ф.н.</p>	
<p>Е-mail, факс почтовый адрес руководителя проекта</p>	<p>hisamitdinova@list.ru</p>
<p>Полное и краткое название организации — адресат финансирования Федеральное государственное бюджетное учреждение науки Институт истории, языка и литературы Уфимского научного центра Российской академии наук ИИЯЛ УНЦ РАН</p>	<p>ФИО руководителя организации — адресата финансирования Хисамитдинова Фирдаус Гильмитдиновна</p>
	<p>ФИО главного бухгалтера организации — адресата финансирования Иркина Регина Венеровна</p>
	<p>Телефон, факс, Е-mail организации (347)2356077, rihll@anrb.ru</p>
<p>Год начала – год окончания проекта</p>	<p>2012 - 2014</p>
<p>Объем финансирования, полученного в 2012 г.</p>	<p>Объем финансирования, запрашиваемый на 2013 г.</p>
<p>Исполнители</p>	<p>Сиразитдинов Зиннур Амирович, к.ф.н.</p>
	<p>Бускунбаева Лилия Айсовна, к.ф.н.</p>
<p>Дата сдачи отчета 19.11.2012</p>	<p>Подпись руководителя проекта</p>

УТВЕРЖДАЕМ
Координатор Программы

акад. РАН Вяч.Вс. Иванов

Координатор Программы

чл.-корр. РАН В.А.Плунгян

“ _____ ” _____ 2012

1. Название направления Корпусная лингвистика. Создание и развитие корпусных ресурсов по языкам народов России

2. Название проекта Создание корпуса текстов республиканских газет на башкирском языке

3. Руководитель проекта (ФИО, ученая степень, должность, место работы)
Хисамитдинова Фирдаус Гильмитдиновна, доктор филологических наук, директор
Института истории, языка и литературы УНЦ РАН

4. Основные участники проекта

Сиразитдинов Зиннур Амирович, кандидат филологических наук, старший научный сотрудник отдела языкознания Института истории, языка и литературы УНЦ РАН
Бускунбаева Лилия Айсовна, кандидат филологических наук, научный сотрудник отдела языкознания Института истории, языка и литературы УНЦ РАН

5. Содержание фактически проделанной работы, полученные результаты, их новизна, научная и практическая значимость, актуальность, соответствие мировому научному уровню.

I. Определены системы разметок текстов публицистики:

а) Экстралингвистические разметки:

дата издания (*для газетных текстов – год, месяц, число, для журнальных – год, месяц, номер выпуска*);

автор;

название источника;

жанр текста:

- интервью, беседа;
- статья, очерк, репортаж, обозрение;
- рецензия;
- обзор печати;
- советы;
- письма;
- поздравления;
- художественно-публицистические жанры: (эссе, фельетон, рассказ, стихи, эпиграммы);

тематика текста:

- политическая и социальная жизнь (политика, право, философия);
- экономика (бизнес, финансы, коммерция, производство, сельское хозяйство, строительство);
- искусство и культура;
- наука и техника;

- образование;
- природа, путешествие;
- здравоохранение, медицина;
- частная жизнь.

объем текста (число предложений, словоформ);

тип носителя (газета, журнал).

Данный тип разметки дает возможность пользователю, ограничив область поиска и отметив в совокупности или отдельно интересующие его параметры, за короткий промежуток времени получить нужную информацию.

б) Морфологические разметки:

исходная форма слова (лемма),

признак части речи,

признаки грамматических категорий.

Соответственно грамматике башкирского языка морфологическая разметка имеет вид:

N (noun) – существительное: *эсәй, дәрт*

V (verb) – глагол: *уйнау, ярДамлашыу*

NUM (numeral) – числительное: *биш, йөДәр*

ADV (adverb) – наречие: *тиД, алдан*

ADJ (adjective) – имя прилагательное: *матур, Дур*

PRON (pronoun) – местоимение: *мин, ниндәй*

POST (postposition) – послелог: *менән, өсөн*

CONJ (conjunction) – союз: *әммә, йәиһә*

PART (particle) – частица: *ғына, әле*

INTJ (interjection) – междометие: *ай, тфү*

PARENTH (parenthesis) – вводное слово: *әлбиттә, ахыры*

PRAEDIC (predicative) – предикатив: *түгел, мөмкин, кәрәк*

Имена существительные

Категория числа

sg (singular) – единств. число: *китан, Дләм*

pl (plural) – множ. число: *китантар, Дләмдәр*

Категория падежа

nom (nominative) – номинатив, основной падеж (төп): *ДагыД*

gen (genitive) – генитив, родительный (эйәлек): *ДагыДДың*

dat (dative) – датив, дательный (төбәү): *ДагыДка*

acc (accusative) – аккузатив, винительный падеж (төшөм): *ДагыДды*

abl (ablative) аблатив, исходный (сығанаД): *ДагыДдан*

loc (locative) – локатив, местный падеж (урын-ваДыт): *ДагыДда*

Категория принадлежности

poss1, sg – 1 л. ед.ч.: *уДыусым*

poss2, sg – 2 л. ед.ч.: *уДыусың*

poss3, sg – 3 л. ед.ч.: *уДыусыһы*

poss1, pl – 1 л. мн.ч.: *уДыусыбыД*

poss2, pl – 2 л. мн.ч.: *уДыусығыД*

poss3, pl – 3 л. мн.ч.: *уДыусылары*

Категория одушевленности/неодушевленности

anim – одушевленность

inan – неодушевленность

Категория сказуемости

pred, 1p, sg – 1 л. ед.ч.: у^лыусымын

pred, 2p, sg – 2 л. ед.ч.: у^лыусыһың

pred, 1p, pl – 1 л. мн.ч.: у^лыусыбы^л

pred, 1p, pl – 2 л. мн.ч.: у^лыусыһығы^л

Глагол

Категория наклонения

ind (indicative) – изъявительное наклонение: алды^л

imp (imperative) – повелительное наклонение: алыгы^л

int (intentional) наклонение намерения: алма^лсымын

cond (conditional) – условное наклонение: аһа^л

sbjv (subjunctive) – сослагательное наклонение: алыр ине, алган булыр ине, ала ине, аласа^л
ине

opt (optative mood) – желательное наклонение: баргы килә

Неспрягаемые формы глагола

inf (infinitive) – инфинитив: у^лырға

ptcp (participle) – причастие: айткән һү^л китеүсә, һөйләшәһе һү^л

ger (gerund) – деепричастие: ауылға барып, о^лата барыу, хәбәр ишеткәс, туйгансы
ашаны, айткәнсә эшлә, өйгә ^лайтышлай

имя действия: инеу, барыу менән

Категория времени

pst def (Past definite tense) – Прошедшее определенное время: килдем

pst indf (Past indefinite tense) – Прошедшее неопределенное время: килгәнмен

pqpf def – (Plusquamperfect definite tense) – Предпрошедшее определенное время:
килгәйнәм

prs (Present tense) – настоящее время: киләм

fut def (Future definite tense) – будущее определенное время: киләсәкмен

fut indf (Future indefinite tense) – будущее неопределенное время: килермен

Категория отрицания

neg (negation, negative) – отрицание: килмә

Местоимение

n-pron – местоимение-существительное: был, кем

adv-pron – местоимение-наречие: шунлы^лтан, ^лай^л

adj-pron – местоимение-прилагательное: ошондай, бындай, бе^лкең

Прилагательное

Степени сравнения

sup (superlative degree) – ^лып-^лы^лыл, ап-а^л

com (comparative degree) – сравнительная степень: ^лурыра^л, бей

slac (slackening degree) – степень ослабления: а^лһыл, аһыу

Морфологическая разметка в национальном корпусе башкирского языка включает в себя исходную форму слова (лемму), признак части речи, признаки грамматических категорий.

Например, глагол *Пара-ма-ган-һың* ‘ты не смотрел’ содержит в себе признаки следующих грамматических категорий: изъявительное наклонение, основной залог, отрицательная форма, прошедшее неопределенное время, 2 лицо, единственное число.

В национальном корпусе данный пример выглядеть таким образом:

лемма: *Парау*

часть речи: V

Семантическая разметка. Данная разметка в корпусе башкирского языка предназначена для поиска по лексико-семантическим признакам.

На начальном этапе работы над национальным корпусом башкирского языка планируется внедрение следующих лексико-семантических разрядов и разметок:

Имена существительные (N)

Разряды

g:concr — предметные имена (*йөрәк, ултырғыс, айыу*)

g:abstr — непердметные имена (*Параңғылы, мәнфәгәт, дәрт*)

g:propn — имена собственные (*Сибай, АПтырна, Шүлгәнташ*)

Имена собственные

Лексико-семантические пометы

Таксономия:

t:hum | t:hum:supernat — лица (*Вәлиди, Людмила, Черномор*)

t:persn — имена (*Айнур, Зәбилә*)

t:patrn — отчества (*Ғәлиевна*)

t:famn — фамилии (*Баһаутдинов*)

t:topon — топонимы (*Өфө, Ырмы, Палы, Мәскәү, АзиҠел*)

Имена прилагательные (A)

Разряды

g:qual — качественные (*матур, Ыыҫыл, һыуы, тәмле*)

g:rel — относительные (*сәскәле, яҠы, емерек, ташланды*)

Имена числительные (NUM, A-NUM)

Разряды

g:card — количественные (*бер, мең, туһан*)

g:ord — порядковые (*унынсы, йөҠнсө, һикһәнненсе*)

g:distr — разделительные (*берәр, алтышар, етешәр*)

g:coll — собирательные (*берәү, бишәү, йөҠбү*)

g:appr — приблизительные (*бишләп, йөҠләгән*)

g:msr (MEASUREMENT) — меры (*бишле, етеле*)

Местоимения, в том числе:

S-PRO — местоимения-существительные (*мин, кем*)

A-PRO — местоимения-прилагательные (*мине, ниндәй, кемдең*)

ADV-PRO — местоимения-наречия (*ҠайҠа, нисек*)

Разряды

g:pers — личные (*бе, ул*)

g:ref — возвратные (*себя*)

- r:poss — притяжательные (*унылы, бынылы, минеке*)
- r:rel — вопросительные (*нимэ, күпме, дайһы*)
- r:dem — указательные (*бынау, теге*)
- r:indet — неопределенные (*Дасандыр, алла ниндэй, алла нисек*)
- r:neg — отрицательные (*һис кем, бер кем дә*)
- r:spes — кванторные (определительные) (*һәр, һәммә, һәр кем*)

Наречия (ADV)

Лексико-семантические пометы

- t:place — место (*алда, бында*)
- t:time — время (*иртэ, элек, кисэ, бая, башта*)
- t:manner — образ действия (*йэйэу, аҗалата*)
- уподобление (*байрамса, башласа, йэштэрсэ*)
- t:degree — мера и степень (*байта, бөтөнләй, икелэтэ*)
- t:cause — причина и цель (*аңгармадан, күрәләтэ, юрамал*)

Глаголы (V)

- t:move — движение (*барыу ‘идти’, һикерәу ‘прыгать’*);
- t:speech — речь (*әйтеу ‘сказать’, өндәшеу ‘обращаться (к кому)’*);
- t:psych — психическое состояние (*шатланыу ‘радоваться’, аптырау ‘удивляться’*);
- t:ment — мышление (*аңлау ‘понять’, фекерләу ‘мыслить’*);
- t:perc — чувственное восприятие (*Дарау ‘смотреть’, һиләу ‘чувствовать’*);
- t:physiol — физиологическая сфера (*һауығыу ‘выздоровливать’, йоллау ‘спать’*);
- t:action — действие (*ялыу ‘писать’, төлөу ‘строить’*);
- t:weather — природное явление (*буранлау ‘мести (о метели)’, йэшинәу ‘сверкать (о молнии)’*).

Более развернутая система лексико-семантической информации, которая будет включать в себя расширенный тематический класс и словообразовательные характеристики лексемы будет разработана в ходе дальнейших работ.

в) Семантические разметки:

Включает следующие лексико-семантические разряды:

Имена существительные (N)

Разряды

- concr** — предметные имена (*йөрәк, ултырғыс, айыу*)
- abstr** — непердметные имена (*Дараңғылы, мәнфәгәт, дәрт*)
- propn** — имена собственные (*Сибай, Атырна, Шүлгәнташ*)

Имена собственные

Лексико-семантические пометы

- persn** — имена (*Айнур, Зәбилә*)
- patrн** — отчества (*Фәлиевна, Хәбирович*)
- famn** — фамилии (*Баһаутдинов, Ногаманов*)
- topon** — топонимы (*Өфө, Ырмысалы, Мәскәу, Агиҗел*)

Имена прилагательные (A)

Разряды

- qual** — качественные (*матур, ылы, һыуы, тәмле*)
- rel** — относительные (*сәскәле, ялы, емерек, ташланды*)

Имена числительные (NUM)

Разряды

- card** — количественные (*бер, мең, туһан*)
ord — порядковые (*унынсы, йөһөнсө, һикһаненсе*)
distr — разделительные (*берәр, алтышар, етешәр*)
coll — собирательные (*берәү, бишәү, йөһәү*)
appr — приближительные (*бишләп, йөһләгән*)
msr — меры (*бишле, етеле*)

Местоимения (PRON)

Разряды

- pers** — личные (*беһ, ул*)
poss — притяжательные (*уныһы, быныһы, минеке*)
rel — вопросительные (*нимә, күпме, һайһы*)
dem — указательные (*бынау, теге*)
indet — неопределенные (*һасандыр, аллә һиндәй, аллә һисек*)
neg — отрицательные (*һис кем, бер кем дә*)
spec — кванторные (определительные) (*һәр, һәммә, һәр кем*)

Глаголы (V)

Лексико-семантические пометы

- move** — движение (*һуһеал, тәгәрә, бар, һикер*)
speech — речь (*әйт, өндәш, саһыр, һабатла*)
psych — психическое состояние (*оял, шатлан, аптыра*)
ment — мышление (*аңла, фекерлә, төшөн*)
perc — чувственное восприятие (*һара, һиһ, ишет*)
physiol — физиологическая сфера (*һауыһ, йөһә, ят*)
action — действие (*яһ, төһө, һуһеат*)
weather — природное явление (*буранла, йәшенлә, яу*)

Наречия (ADV)

Лексико-семантические пометы

- place** — место (*алда, бында*)
time — время (*иртә, элек, кисә, бая, башта*)
manner — образ действия (*йәйәү, аһсалата*)
likening — уподобление (*байрамса, баһһаса, йәштәрсә*)
degree — мера и степень (*байтаһ, бөтөнләй, икеләтә*)
cause — причина и цель (*аңгармаһтан, күрәләтә, юрамал*)

Частицы (PART)

Разряды

- lim** — ограничительные (*һына/һенә, һына/һенә*)
intens — усилительные (*уһ, үк; та/тә, да/дә, һа/һә, ла/лә*)
conf — подтверждения (*дабаһа/һабаһа/лабаһа/табаһа*)
indef — неопределенности (*әле*)

Модальные слова (MOD)

Разряды

- affir/neg** — утверждение/отрицание (*әйе, юһ, түгел, әлбиттә*)
necess — необходимость (*кәрәк, тейеш*)
hypoth — вероятность, гипотетические (*ахыры, бәлки, күрәһең*)
possib — возможность (*мөмкин, ярай*)
induc — побуждение (*һинһар, рәһим итегеһ, һуй*)

Подражательные слова (IMT)

echoic – звукоподражательный (*шалтыр, кетер, былт*)

image – образоподражательный (*был□, йым-йым*)

Союзы (CONJ)

Разряды

coord – сочинительный (*һәм, ләкин*)

copul – соединительный (*һәм, да/дә, тағы, йәнә, шуға күрә*)

advers – противительный (*ләкин, юғиһа, фә□вт*)

disjun – разделительный (*йә, хатта, берсә*)

subord – подчинительный (*әгәр, сөнки, тимәк*)

exppound – изъяснительный (*тимәк, йәғни*)

causal – причинные (*сөнки, шуға күрә*)

cond – условный (*әгәр, әгәр □ә*)

comp – сравнительный (*әйтерһең, әйтерһең дә*)

Послелого (POST)

Разряды

nom – управляющий основным падежом (*менән, өсөн, кеүек*)

dat – управляющий дательным падежом (*табан, □арай, са□лы*)

abl – управляющий исходным падежом (*һуң, баш□а, бирле*)

II. Разработан алгоритм и программа программа автоматического анализа башкирской словоформы.

Для агглютинативного башкирского языка предложен подход в многомерном (тензорном) представлении. Любая основа башкирского языка представляется в виде элемента многомерного объекта:

${}^{pt}{}_s a_i^f$ – i-я основа, где p, t – определяют фонетическое строение слова.

f – определяет тип изменений в конечном звуке основы при присоединении аффиксов.

s – определяет принадлежность слова (s=1 — исконно башкирские и ранние заимствованные слова, s=2 — поздние заимствования из русского языка, имеющие значительные отличия при словоизменении (изменения в корне и в типах аффиксов словоизменения, например: *лагерь: лагер□ан, лагерга; волость: волос□а, волостан, волосы*); s=3 — заимствования из арабского и персидского языков, имеющие отличия при словоизменении.

Вводится оператор сцепления текстовых переменных \otimes . Данный оператор сцепляет только те элементы псевдотензора, которые имеют одинаковые значения индексов p, t. Текстовыми переменными для данного оператора являются корень и аффикс, либо корень + аффикс + аффикс и т.д.

Тогда алгоритм анализа башкирской словоформы состоит из следующих процедур:

1. Задается линейный список возможных в языке моделей словоизменений ($A^{ijklm\nu vxyz}$): 611 для именных и 665 для глагольных форм.

2. Для удобства индексации и последующего анализа словоформы вводится двухмерный массив. Столбцы массива определяют грамматические категории, а подкатегории определяются строками данного массива (i, j, k и т.д.). Фонетическая строение конкретных аффиксов описывается как элемент массива в виде матрицы 4x4 (p и t).

3. Программа посимвольно выделяет слева возможный корень словоупотребления, сравнивая со словарем основ, определяет его фонетическое строение.

4. В процессе вычленения аффиксов от остатка строится дерево возможных сочетаний аффиксов словоизменения для данной основы. При добавлении нового аффикса к узлу дерева проверяется его конечность.

5. Если словоформа полностью разложена, то разложение сравнивается с заданной моделью словоизменения.

6. При наличии нескольких возможных основ, рассматриваются все возникающие варианты.

Реализованная программа анализа башкирской словоформы показывает удовлетворительную работу и используется нами в для автоматической морфологической разметки корпуса

III. Подготовлены электронные тексты из газет и журналов, издаваемых на башкирском языке.

Для корпуса башкирской публицистики взяты тексты из следующих источников:

а) Газеты

- «Башкортостан» (с 2002 по 2011 гг.);
- «Киске Өфө» (с 2003 по 2012 гг.);
- «Йәшлек» (с 2003 по 2010 гг.).

б) Журналы:

- «Башкортостан кызы» (с 2002 по 2011 гг.);
- «Ағизел» (с 2003 по 2012 гг.);
- «Шоңкар» (с 2003 по 2010 гг.).

Объем подготовленных текстов для корпуса составляет на текущий момент 4 миллиона словоупотреблений. Ведется работа по выставлению корпуса в сети Интернет по модели корпуса прозаических текстов <http://mfbl.ru/bashkorp/korpus>.

6. Общее число опубликованных в 2012 г. по проекту работ

6.1. количество монографий

6.2. количество сборников статей

6.3. количество статей 3.

7. Список опубликованных монографий и сборников статей, с полным указанием выходных данных, а также их краткие аннотации.

8.Список опубликованных по проекту статей

Сиразитдинов Б.З., Сиразитдинов З.А. О моделировании словоизменительной системы и разработке программ автоматического морфологического анализа башкирского языка//Современное казахское языкознание: актуальны вопросы прикладной лингвистики. Алматы, 2012,с. 103-107.

Сиразитдинов З. А., Бускунбаева Л. А., Ишмухаметова А. Ш., Ибрагимова А. Д., Мигранова Л. Г. Корпус текстов периодической печати на башкирском языке//Актуальные проблемы диалектологии народов России: Материалы XII региональной конференции. Уфа, 2012, с. 139-141.

9. Список книг, сданных в печать или поданных на издательские гранты

10. Экспедиции, организованные в рамках проекта

11. Конференции, организованные в рамках проекта.

12. Важнейшие научные результаты работы по проекту. Определена структура и функциональные возможности корпуса башкирских публицистических текстов,разработаны системы лингвистических и экстралингвистических разметок для, создана программа автоматического анализа и разметки текстов, подготовлены сами материалы для корпуса

13. Наиболее значимый научный результат проекта

На базе многомерного представления для агглютинативного башкирского языка разработана система автоматического анализа словоизменительных форм и программа автоматическоюй разметки текстов.

14. Краткий финансовый отчет за 2012 г.

15. Запрашиваемый на 2013 г. объем финансирования, с кратким обоснованием расходов.

16. Краткое обоснование научных работ на 2013 г., ожидаемые результаты – до 1 стр. и заполненная форма²

Подпись руководителя

Форма 2
Планируемое содержание работ на 2013 г.

№	Название проекта	Организация-исполнитель и учреждения-соисполнители	Руководитель проекта (+ кол-во исполнителей)	Запрашиваемый объем финансирования на 2013 г. (тыс. руб.)	Ожидаемые в 2013 г. результаты
	“Создание корпуса текстов республиканских газет на башкирском языке”	Федеральное государственное бюджетное учреждение науки Институт истории, языка и литературы Уфимского научного центра Российской академии наук ИИЯЛ УНЦ РАН	Хисамитдинова Ф.Г. руководитель проекта 5 исполнителей		1. Увеличение объема корпуса на 6 миллионов словупотреблений 2. Увеличение объема базового словаря морфоанализатора на 10 тыс. единиц. 3. Функционирование корпуса в сети Интернет.