

Документ подписан простой электронной подписью
Информация о владельце:
ФИО: Белгородский Валерий Савельевич
Должность: Ректор
Дата подписания: 25.03.2024 14:28:49
Уникальный программный ключ:
8df276ee93e17c18e7bee9e7cad2d0ed9a082473

Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Российский государственный университет им. А.Н. Косыгина
(Технологии. Дизайн. Искусство)»

Институт Магистратура
Кафедра общего и славянского искусствознания

ОЦЕНОЧНЫЕ МАТЕРИАЛЫ

для проведения текущей и промежуточной аттестации
по учебной дисциплине

Корпусная лингвистика

Уровень образования	магистратура
Направление подготовки	44.04.01 Педагогическое образование
Программа	Дистанционные технологии в гуманитарном образовании
Срок освоения образовательной программы	2 года 6 месяцев
Форма обучения	очно-заочная, заочная

Оценочные материалы учебной дисциплины «Корпусная лингвистика» основной профессиональной образовательной программы высшего образования, рассмотрены и одобрены на заседании кафедры, протокол № 6 от 06.03.2023 г.

Составитель оценочных материалов учебной дисциплины:

профессор Г.В. Варакина

Заведующий кафедрой: Г.В. Варакина

1. ОБЩИЕ СВЕДЕНИЯ

Учебная дисциплина «Корпусная лингвистика» изучается на первом курсе
Курсовая работа/Курсовой проект – не предусмотрены.

Форма промежуточной аттестации:
зачет с оценкой

2. ЦЕЛИ И ЗАДАЧИ ОЦЕНОЧНЫХ СРЕДСТВ, ОБЛАСТЬ ПРИМЕНЕНИЯ

Оценочные средства являются частью рабочей программы учебной дисциплины и предназначены для контроля и оценки образовательных достижений обучающихся, освоивших компетенции, предусмотренные программой.

Целью оценочных средств является установление соответствия фактически достигнутых обучающимся результатов освоения дисциплины, планируемому результату обучения по дисциплине, определение уровня освоения компетенций.

Для достижения поставленной цели решаются следующие задачи:

- оценка уровня освоения общепрофессиональных и профессиональных компетенций, предусмотренных рабочей программой учебной дисциплины;
- обеспечение текущего и промежуточного контроля успеваемости;
- оперативного и регулярного управления учебной, в том числе самостоятельной деятельностью обучающегося;
- Соответствие планируемых результатов обучения задачам будущей профессиональной деятельности через совершенствование традиционных и внедрение инновационных методов обучения в образовательный процесс.

Оценочные материалы по учебной дисциплине включают в себя:

- перечень формируемых компетенций, соотнесённых с планируемыми результатами обучения по учебной дисциплине;
- типовые контрольные задания и иные материалы, необходимые для оценки результатов обучения;

Оценочные материалы сформированы на основе ключевых принципов оценивания:

- валидности: объекты оценки соответствуют поставленным целям обучения;
- надежности: используются единообразные стандарты и критерии для оценивания достижений;
- объективности: разные обучающиеся имеют равные возможности для достижения успеха.

3. ФОРМИРУЕМЫЕ КОМПЕТЕНЦИИ, ИНДИКАТОРЫ ДОСТИЖЕНИЯ КОМПЕТЕНЦИЙ, СООТНЕСЁННЫЕ С ПЛАНИРУЕМЫМИ РЕЗУЛЬТАТАМИ ОБУЧЕНИЯ ПО ДИСЦИПЛИНЕ ИСПОЛЬЗУЕМЫЕ ОЦЕНОЧНЫЕ СРЕДСТВА

Код компетенции, код индикатора достижения компетенции	Планируемые результаты обучения по дисциплине	Наименование оценочного средства	
		текущий контроль (включая контроль самостоятельной работы обучающегося)	промежуточная аттестация
УК-4 ИД-УК-4.1	- редактирует различные академические тексты	Письменная работа: Технология обработки текстового материала (технология критического мышления) Круглый стол 1 Круглый стол 2 Устный опрос Тестирование	зачет с оценкой
УК-4 ИД-УК-4.2	- демонстрирует готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации	Письменная работа: Технология обработки текстового материала (технология критического мышления) Круглый стол 1 Круглый стол 2 Устный опрос Тестирование	

Код компетенции, код индикатора достижения компетенции	Планируемые результаты обучения по дисциплине	Наименование оценочного средства	
		текущий контроль (включая контроль самостоятельной работы обучающегося)	промежуточная аттестация
УК-4 ИД-УК-4.3	- осуществляет межличностное деловое общение, в том числе на иностранных языках с применением профессиональных языковых форм и средств	Круглый стол 1 Круглый стол 2 Устный опрос	

\

4. ТИПОВЫЕ КОНТРОЛЬНЫЕ ЗАДАНИЯ И ДРУГИЕ МАТЕРИАЛЫ, НЕОБХОДИМЫЕ ДЛЯ ОЦЕНКИ ПЛАНИРУЕМЫХ РЕЗУЛЬТАТОВ ОБУЧЕНИЯ И УРОВНЯ СФОРМИРОВАННОСТИ КОМПЕТЕНЦИЙ

4.1. Оценочные материалы **текущего контроля** успеваемости по учебной дисциплине, в том числе самостоятельной работы обучающегося, типовые задания

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.1 Подготовка и редактирование различных академических текстов.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.

Письменная работа: Технология обработки текстового материала (технология критического мышления) (УК-4, ИД-УК-4.1, ИД-УК-4.2)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.1 Подготовка и редактирование различных академических текстов.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

Письменная работа

Время выполнения 90 мин

Форма работы – самостоятельная, индивидуальная.

А) Модель работы с концептуальной таблицей

1. Выделить основания для сопоставления (первый столбик).
2. Провести сравнительный анализ по выделенным основаниям.

Критерии для сопоставления объектов	Национальный корпус русского языка	Фундаментальная электронная библиотека (ФЭБ)
Общее		
1.		
2...		
Различное		
1.		
2. ...		

Материал для сопоставления:

Национальный корпус русского языка и Фундаментальная электронная библиотека (ФЭБ).

ОТВЕТ:

Критерии для сопоставления объектов	Национальный корпус русского языка	Фундаментальная электронная библиотека (ФЭБ)
Общее		

1. Единица хранения	собрание текстов	собрание текстов
2. Формат хранения	электронный	электронный
3. Степень доступности	открытый доступ пользователей к сетевой системе.	открытый доступ пользователей к сетевой системе.
Различное		
1. Цель создания	<i>собрание независимых корпусов текстов для решения определенных лингвистических задач.</i>	<i>создание многофункциональной программно-информационной среды для специалистов по русской филологии и фольклористике.</i>
2. Назначение	<i>изучение текстов, проведение лингвистических исследований по корпусу текстов.</i>	<i>чтение текстов, проведение исследований с использованием полнотекстовой информационной системы, распространение русской словесности и расширение международных культурных связей.</i>
3. Обработка текстов	<i>обработка на техническом языке – разметка (собрание массива аннотированных /размеченных текстов).</i>	<i>обработка – собрание текстов в авторской редакции.</i>
4. Язык объекта сопоставления	<i>русский язык и версия на английском языке</i>	<i>русский язык</i>
4. Характерные признаки сопоставляемых объектов	<i>мультимедийность</i>	<i>энциклопедичность, иерархическая информационная структура</i>
5. Типы собранных текстов	<i>художественная литература, газетные тексты, блогосфера, устная речь, диалектная речь, просторечие, устные тексты разных жанров, научная речь, аудио- и видеофайлы, кино-, теле-, радио- и театральные постановки.</i>	<i>русская литература (XI–XX вв.) и фольклор, история русской филологии и фольклористики, архивные материалы, документы об авторе текстов, мемуары, филологические исследования.</i>
6. Наличие справочной информации	<i>статистика (текстов, слов), частотный словарь, инструменты исследования (портреты корпуса, подкорпуса. слова).</i>	<i>научная литература, словари и энциклопедии, каталог ссылок на филологические и специализированные сайты.</i>
7. Область применения	<i>специальные научные исследования, преподавание русского языка, обучение русскому языку как иностранному.</i>	<i>гуманитарное образование всех уровней.</i>
8. Пользователи	<i>ученые, преподаватели, учащиеся, студенты, аспи-</i>	<i>широкий круг пользователей</i>

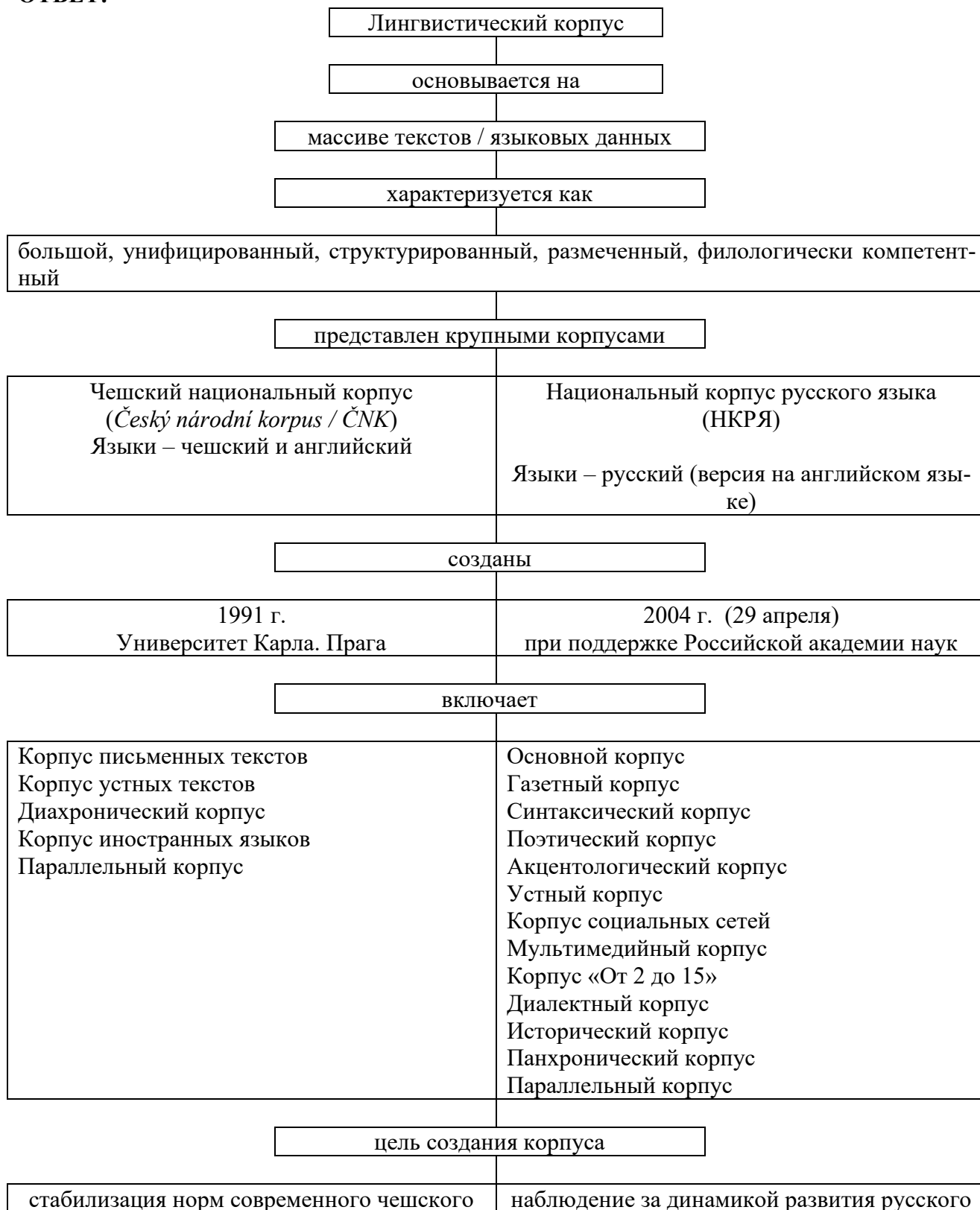
	<i>ранты.</i>	
--	---------------	--

Б) Модель работы с денотатным графом

1. Представьте самые крупные корпуса в формате денотатного графа.

2. Обозначьте основные свойства Чешского национального корпуса (*Český národní korpus / ČNK*) и Национального корпуса русского языка (НКРЯ)

ОТВЕТ:



языка, нормализация политической ситуации (сотрудничество с международным научным сообществом)	языка
---	-------

Круглый стол 1 (УК-4, ИД-УК-4.2, ИД-УК-4.3)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств

Круглый стол на тему: «Корпусная лингвистика: исторический и лингводидактический аспекты»

Время выполнения 90 мин.

Проводится в подгруппах по 5-7 чел.

Форма работы – самостоятельная, групповая.

Проблема: новые технологии в лингвистике

Ожидаемый результат: Систематизация положений, обсуждаемых на семинаре

№	Вопрос	Ответ	Компетенция
		<p>УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.</p> <p>ИД-УК-4.1 Подготовка и редактирование различных академических текстов.</p> <p>ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.</p> <p>ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.</p>	
1	Создание и разметка корпусов текстов	<p>Лингвистическая разметка текстовых корпусов бывает внешней (метаразметкой) и внутренней (собственно лингвистической), при этом синонимами выступают понятия <i>тэггинг</i> и <i>аннотирование</i>. Несмотря на то, что в иностранной литературе зачастую понятия <i>annotation</i>, <i>markup</i> и <i>tagging</i> используются взаимозаменяемо, именно аннотирование относится к внешней разметке. В литературе приводятся различные определения аннотирования. Аннотирование – это процесс составления аннотации, а сама аннотация – краткая характеристика документа, поясняющая его содержание, назначение, форму, другие особенности. Если семантика аннотации выражена явным образом, то такая аннотация – семантическая. Таким образом, процесс формирования списка ключевых слов является видом семантического аннотирования.</p> <p>Разметка текста заключается в приписывании текстам и их компонентам дополнительной информации (метаданных). Метаданные можно поделить на 3 типа: экстралингвистические, относящиеся ко всему тексту; данные о структуре текста; лингвистические</p>	УК-4, ИД-УК-4.2, ИД-УК-4.3

		<p>метаданные, описывающие элементы текста. Метаописание текстов корпуса включает как содержательные элементы данных (библиографические данные, признаки, характеризующие жанровые и стилевые особенности текста, сведения об авторе), так и формальные (имя файла, параметры кодирования, версия языка разметки, исполнители этапов работ). Эти данные обычно вводятся вручную. Структурная разметка документа (выделение абзацев, предложений, слов) и собственно лингвистическая разметка обычно осуществляются автоматически. Чем богаче и разнообразнее разметка, тем выше научная и учебная ценность корпуса. В Национальном корпусе русского языка в настоящее время используется пять типов разметки: метатекстовая, морфологическая (словоизменяемая), синтаксическая, акцентная и семантическая.</p>	
2	<p>Словарь-конкорданс и его применение в рамках корпусной лингвистики</p>	<p>Поиск в корпусе данных позволяет по выбранному слову построить конкорданс. Конкорданс – это список всех употреблений определенного языкового выражения в контексте и со ссылкой на источник при необходимости. В данной коннотации этот термин часто применяется в корпусной лингвистике. В исходном своем толковании этот термин также используется для обозначения списка ключевых слов книги или работы, градированных в алфавитном порядке, с их контекстами.</p> <p>Конкордансы часто используются в прикладной лингвистике: при переводе, в лексикографии, при обучении и изучении языка, при анализе текста.</p> <p>Конкордансы применяются для – решения следующих лингвистических задач:</p> <ul style="list-style-type: none"> -создавать списки слов; -сопоставлять разные использования одного слова; -искать и исследовать фразы и идиомы; -анализировать ключевые слова; -искать перевод различных словосочетаний или слов; -анализировать частоту употребления какого-либо слова и словосочетания. <p>Словарь-конкорданс является особым типом словаря, в котором каждое слово приводится с минимальным контекстом.</p> <p>При изучение текста можно выделить поисковую, эвристическую, аналитическую функции конкорданса, а так же функции индексации и сравнения. <i>Поисковая функция</i> позволяет быстро находить нужный фрагмент текста, используя заданное слово или словосочетание. <i>Эвристическая функция</i> конкорданса отличает его от именных указателей наличием контекста. Зачастую, контексты позволяют лингвисту или простому читателю увидеть новую трактовку текста. <i>Аналитическая функция</i> позволяет проводить анализ различных языковых показателей, таких как лексемы, ключевые слова, частота их употребления в тексте и т.д. <i>Функция индексации</i> дает возможность при подготовке текста к публикации создавать индексы и списки слов. <i>Функция сравнения</i> применяется при сравнении в тексте всевозможных коннотаций и употреблений слова.</p> <p>В отличие от словаря, который опирается на словарную статью, конкорданс приводит примеры – контексты словоупотреблений. Конкорданс – фун-</p>	<p>УК-4, ИД-УК-4.2, ИД-УК-4.3</p>

		дамент для созданий разнообразных дифференциальных, в том числе и толковых, словарей.	
3	Метод автоматической кластеризации текстов и его применение.	<p>Основные процедуры обработки естественного языка: токенизация, лемматизация, стемминг, парсинг.</p> <p><i>Токенизация</i> – разбиение потока символов в естественном языке на отдельные значимые единицы (токены, словоформы), является необходимым условием для дальнейшей обработки естественного языка. Если бы языки обладали совершенной пунктуацией, токенизация не представляла бы сложности – даже самая простая программа могла бы разделить текст на слова, руководствуясь пробелами и знаками препинания. Но в действительности языки подобной пунктуацией не обладают, что усложняет задачу токенизации. Например, в английском языке встречаются случаи, которые не могут быть однозначно токенизированы.</p> <p>Другая специфическая задача морфологического анализа – это <i>лемматизация</i>, процесс образования первоначальной формы слова, исходя из других его словоформ. Во многих языках слово может встречаться в нескольких формах с различными флексиями. Базовая форма, зафиксированная в словаре, называется <i>леммой</i> слова. Лемматизация представляет собой процесс группировки различных флексивных форм одного слова таким образом, чтобы при анализе они обрабатывались как одно слово. Лемматизация связана с идентификацией частей речи и включает в себя сокращение слов из корпуса до соответствующих им лексем. Именно лемматизация позволяет исследователю выделять и изучать все варианты отдельной лексемы без необходимости введения всех возможных вариантов.</p> <p>Процесс, несколько отличный от лемматизации, называется <i>стеммингом</i>, он состоит в нахождении стема (основы) слова. Разница заключается в том, что стеммер обрабатывает отдельное слово без знания контекста, и, таким образом, не может дифференцировать слова, которые имеют разные значения в силу отнесенности к разным частям речи. Стеммеры обычно более просты для реализации и быстрее обрабатывают данные, а более низкая точность их работы может не иметь решающего значения для многих приложений. Лемма является базовой формой для токена, и это соответствие будет обнаружено как при стемминге, так и при лемматизации.</p> <p><i>Парсинг</i> – это процесс сопоставления линейной последовательности лексем (слов, токенов) языка с его формальной грамматикой. Результатом обычно является дерево зависимостей (синтаксическое дерево). Построение автоматических синтаксических анализаторов (парсеров) для больших корпусов является одной из самых важных областей компьютерной лингвистики. Большинство подходов объединяет качественные и количественные измерения. Наряду с разными статистическими подходами, которые тренируются на снабженных вручную пометами синтаксических деревьях, многие синтаксические анализаторы используют основанные на правилах</p>	УК-4, ИД-УК-4.2, ИД-УК-4.3

		или основанные на ограничениях подходы, которые прямо моделируют специфические лингвистические теории. Разработка этих синтаксических анализаторов тесно переплетается с развитием этих теорий. Поскольку большинство предложений неоднозначны в любой теории, на основе правил (или перечня ограничений) должна быть разработана стратегия снятия неоднозначности. Многие стратегии снятия неоднозначности полагаются на количественные данные – частоту данной структуры в данном корпусе (тип), ограничения на выборку для данных лексических единиц, которые были получены или выделены из корпусных данных, и т.д.	
--	--	---	--

Круглый стол 2 (УК-4, ИД-УК-4.2, ИД-УК-4.3)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.

Круглый стол на тему: «Семантические исследования на материале корпусов текстов».

Время выполнения 90 мин.

Проводится в подгруппах по 5-7 чел.

Форма работы – самостоятельная, групповая.

Проблема: новые технологии в семантике

Ожидаемый результат: Систематизация положений, обсуждаемых на семинаре

№	Вопрос	Ответ	Компетенция
		УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия. ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации. ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.	
1	Синтаксически и семантически аннотированный корпус русского языка: современное состояние и перспективы.	Ценность корпуса определяется глубиной разметки: каждому слову приписывается исчерпывающая морфологическая информация, а для каждого предложения строится полное синтаксическое дерево зависимостей. Разметка производится автоматически, однако работа над корпусом включает ручную правку всех предложений человеком. Работа по обогащению корпуса семантической информацией включает в себя четыре этапа: (1) разработку инвентаря семантических дескрипторов, (2) создание семантического словаря с приписанными лексемам семантическими дескрипторами и согласование этого словаря с комбинаторным словарем (КС) системы ЭТАП, (3) внедрение семантической информации в уже размеченный морфологически и синтаксически корпус текстов; (4) создание инструментария для	УК-4, ИД-УК-4.2, ИД-УК-4.3

		<p>работы с семантической информацией.</p> <p>Предлагаемый набор семантических дескрипторов (семантический метаязык) должен в конечном счете решать две задачи: 1. обеспечивать лингвистически содержательную классификацию всей лексики – и предметной, и предикатной; 2., в соединении с морфологической и синтаксической разметкой текстов предоставлять исследователю существенную информацию о закономерностях поведения элементов различных лексико-семантических классов в текстах. В качестве дескрипторов везде, где это возможно, используются слова естественного языка. При разработке инвентаря семантических дескрипторов важно разделить лексемы на два основных типа – предметные (названия животных, птиц, рыб, овощей, фруктов, камней, гор, планет, светил и т. п.) и предикатные. Предметные и предикатные дескрипторы делятся на две подгруппы – родовые и видовые. Родовые дескрипторы обозначаются существительными (например, «животное», «совокупность», «состояние», «действие»), тогда как видовые – прилагательными («домашний», «природный», «речевой», «ментальный», «физический»). Предикатным словам, кроме родовых и видовых дескрипторов, приписываются семантические роли по каждой из валентностей. Например, глаголу вязать в значении «плести спицами или крючком» приписываются семантические роли «агенс» (<i>Маша вяжет</i>), «результат» (<i>вязать шарф</i>), «пациенс» (<i>вязать из шерсти</i>) и «инструмент» (<i>вязать крючком</i>). С учетом семантических ролей общий объем дескрипторов составляет 250–300 единиц. Предметной и предикатной лексике соответствуют две разные семантические классификации языковых единиц – таксономическая и фундаментальная. Предметные дескрипторы членят словарь не с научной, а с наивно-энциклопедической точки зрения.</p> <p>Синтаксическая разметка является результатом <i>парсинга</i>, выполняемого на основе данных морфологического анализа. Этот вид разметки описывает синтаксические связи между лексическими единицами и различные синтаксические конструкции (например, придаточное предложение, глагольное словосочетание и т.д.).</p> <p>В отличие от морфологии, способы представления синтаксической структуры и синтаксических отношений не столь унифицированы. Наблюдается разнообразие синтаксических теорий и формализмов:</p> <ul style="list-style-type: none"> – грамматика зависимостей; – грамматика непосредственно составляющих; – грамматика структурных схем; – традиционные синтаксические учения о членах предложения; – функциональная грамматика; – семантический синтаксис и др. <p>Синтаксический анализ для русского языка чаще всего представлен структурами зависимостей.</p>	
2	Анализ семантических помет в Национальном корпусе русского языка.	<p>Национальный корпус русского языка снабжен разными видами разметки, в том числе – морфологической и семантической. Семантическая разметка включает пометы таксономического класса, мереологию, оценку и др. типы информации. В разных видах</p>	УК-4, ИД-УК-4.2, ИД-УК-4.3

		<p>разметки есть разные виды неоднозначности. Что касается семантической разметки, то связанная с ней проблема неоднозначности состоит в следующем. В словаре у каждого значения многозначного слова есть своя собственная семантическая помета. Однако когда программа автоматически расставляет пометы в тексте, то она каждому вхождению слова приписывает все пометы, которые есть в словаре, потому что программа не знает, в каком значении выступает слово в данном тексте. Тогда нужно различить эти значения. Одним из эффективных путей решения этой проблемы являются семантические фильтры, т.е. семантические правила, которые позволяют оставлять при каждом вхождении слова только одну помету. Таким образом, многозначность снимается с точностью до семантического класса (т.е. с точностью до семантической пометы). Для снятия многозначности с помощью фильтров используется принцип контекстной однозначности. В словаре у слова может быть несколько значений, а в тексте (кроме специальных случаев языковой игры) – одно значение. Поскольку слово обычно не выступает в тексте изолированно, а включено в определенный контекст (в другой терминологии – в конструкцию), то этот контекст (конструкция) и является, в самом грубом и общем виде, фильтром. Фильтр работает следующим образом: формулируются семантические, морфологические и синтаксические условия, в которых реализуется некоторое значение слова; эти условия записываются в виде контекста; осуществляется поиск слова в заданном контексте; результатом этого поиска будет корпус примеров, где слово выступает в заданном значении, соответствующем определенной семантической помете. Остальные пометы стираются. Далее процедура повторяется для остальных семантически размеченных значений слова.</p>	
3	<p>Статистические исследования на материале корпусов текстов.</p>	<p>Статистические методы играют важную роль в анализе текстов в рамках корпусных исследований. Они позволяют исследователям получить количественные данные о различных языковых явлениях и провести объективный анализ текстов.</p> <p><i>Частотный анализ.</i> Один из основных статистических методов, используемых в анализе текстов, – это частотный анализ. Он позволяет определить частотность и распределение определенных слов, выражений или грамматических конструкций в тексте или в корпусе текстов.</p> <p>Частотный анализ может быть полезен для выявления наиболее употребляемых слов или выражений в определенном жанре текстов, для сравнения частотности слов в разных языках или для изучения изменений в употреблении слов в течение времени.</p> <p><i>Коллокационный анализ.</i> Коллокационный анализ – это метод, который позволяет исследователям выявить наиболее типичные сочетания слов в тексте или в корпусе текстов. Он основан на статистическом анализе совместной встречаемости слов и позволяет выявить семантические и грамматические связи между словами. Коллокационный анализ может быть полезен для изучения лексических особенностей языка, выявления идиоматических выражений или для определения наиболее типичных словосоче-</p>	<p>УК-4, ИД-УК-4.2, ИД-УК-4.3</p>

		<p>таний в определенном жанре текстов.</p> <p><i>Статистический анализ стилей и жанров.</i> Статистические методы также могут быть использованы для анализа стилей и жанров текстов. Исследователи могут проводить статистический анализ различных лингвистических признаков, таких как длина предложений, использование определенных частей речи или грамматических конструкций, чтобы выявить характерные особенности стиля или жанра. Статистический анализ стилей и жанров может быть полезен для автоматической классификации текстов по жанрам или для выявления стилистических особенностей в текстах разных авторов.</p> <p><i>Машинное обучение и анализ текстов.</i> Статистические методы также широко применяются в машинном обучении и анализе текстов. Используя статистические модели и алгоритмы, исследователи могут разрабатывать системы автоматического распознавания и классификации текстов, а также системы автоматического извлечения информации из текстов. Машинное обучение и анализ текстов на основе статистических методов могут быть полезными для обработки больших объемов текстов, автоматического анализа семантики текстов или для разработки систем машинного перевода.</p>	
--	--	--	--

Устный опрос (УК-4, ИД-УК-4.1, ИД-УК-4.2, ИД-УК-4.3)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.1 Подготовка и редактирование различных академических текстов.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.

Устный опрос по вопросам:

Форма работы – самостоятельная, индивидуальная

№	Вопрос	Ответ	Компетенция
<p>УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.</p> <p>ИД-УК-4.1 Подготовка и редактирование различных академических текстов.</p>			
1	Что может являться единицей корпуса?	Корпус текстов - это вид корпуса данных, единицами которого являются тексты или их достаточно значительные фрагменты, включающие, например, какие-то полные фрагменты макроструктуры текстов данной проблемной области. Корпус текстов характеризуется четырьмя основными параметрами: 1. он должен быть достаточно большого объема; 2. корпус должен быть структурированным или размеченным; 3. тексты, составляющие определенного кор-	УК -4 ИД-УК-4.1

		пуса, должны быть в электронном варианте; 4. в понятие «Электронный корпус» входит, как правило, специальное программное обеспечение для работы с этим корпусом.																	
<p>УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия. ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.</p>																			
2	Определите виды корпуса по указанному признаку	<table border="1"> <thead> <tr> <th>Признак</th> <th>Виды корпуса (ответ)</th> </tr> </thead> <tbody> <tr> <td>Форма хранения</td> <td>звуковые, письменные, смешанные</td> </tr> <tr> <td>Язык текстов</td> <td>русский, английский</td> </tr> <tr> <td>Параллельность</td> <td>одноязычные, двуязычные, многоязычные</td> </tr> <tr> <td>Стиль</td> <td>литературные, диалектные, разговорные, публицистические, терминологические, смешанные</td> </tr> <tr> <td>Способ доступа</td> <td>свободно доступные, коммерческие, закрытые</td> </tr> <tr> <td>Разметка</td> <td>размеченные, неразмеченные</td> </tr> <tr> <td>Характер разметки</td> <td>морфологические, синтаксические, семантические, просодические</td> </tr> </tbody> </table>	Признак	Виды корпуса (ответ)	Форма хранения	звуковые, письменные, смешанные	Язык текстов	русский, английский	Параллельность	одноязычные, двуязычные, многоязычные	Стиль	литературные, диалектные, разговорные, публицистические, терминологические, смешанные	Способ доступа	свободно доступные, коммерческие, закрытые	Разметка	размеченные, неразмеченные	Характер разметки	морфологические, синтаксические, семантические, просодические	ИД-УК-4.1 ИД-УК-4.2
Признак	Виды корпуса (ответ)																		
Форма хранения	звуковые, письменные, смешанные																		
Язык текстов	русский, английский																		
Параллельность	одноязычные, двуязычные, многоязычные																		
Стиль	литературные, диалектные, разговорные, публицистические, терминологические, смешанные																		
Способ доступа	свободно доступные, коммерческие, закрытые																		
Разметка	размеченные, неразмеченные																		
Характер разметки	морфологические, синтаксические, семантические, просодические																		
<p>УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия. ИД-УК-4.1 Подготовка и редактирование различных академических текстов. ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации. ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.</p>																			
3	Поясните, что означают следующие понятия: «исследовательский корпус»	Исследовательский корпус – это корпус, предназначенный для изучения различных аспектов функционирования языковой системы. Они строятся не post factum — после проведения какого-либо исследования, а до его проведения. Этот тип корпусов данных ориентирован на широкий класс лингвистических задач.	УК -4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3																
4	Поясните, что означают следующие понятия: «статический корпус»	Статический корпус Первоначально корпусы текстов создавались как статические образования, отражающие определенное временное состояние языковой системы. Типичными представителями этого вида корпусов – авторские	УК -4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3																

		корпусы, коллекции текстов писателей.	
5	Поясните, что означают следующие понятия: «параллельный корпус»	<p>Параллельный корпус</p> <p>1) корпусы, представляющие множество текстов-оригиналов, написанных на каком-либо исходном языке, и текстов-переводов этих исходных текстов на один или несколько других языков;</p> <p>2) корпусы, объединяющие тексты из одной и той же тематической области, независимо написанные на двух или нескольких языках.</p>	УК -4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3
6	Как можно использовать рассмотренные корпусы в лингвистическом исследовании / в практике перевода?	<p><i>Корпусы параллельных текстов</i> формируются для научных и практических целей (в частности, для преподавания иностранных языков). По своей структуре это подмножество текстов на языке-источнике и одно или несколько подмножеств текстов, которые являются переводами текстов языка-источника на языки-цели. Например, английский текст «Alice in Wonderland» и его переводы на немецкий, французский и русский языки могут формировать такой корпус или быть частью большего корпуса параллельных текстов. Корпусные методы исследования перевода привнесены в область аудиовизуального перевода, применяются возможности корпусных менеджеров в дескриптивном переводоведении, в преподавании языка в иностранной аудитории. Учебные тексты позволяют классифицировать типы ошибок и учитывать их в процессе преподавания. Сведения такого рода учитываются в англоязычных учебных словарях.</p> <p><i>Исследовательский корпус</i> может быть использован как инструмент научного исследования, в грамматических и лексикологических исследованиях, в лингвистической, судебно-лингвистической экспертизе, исследование частотных характеристик лексем (лемм), исследование коллокаций (сочетаний лексем).</p> <p><i>Статический корпус</i> необходим для исследования нормы / узуса, в социолингвистических исследованиях, для изучения устной речи. Наличие электронных текстов, принадлежащих одному автору, дает возможность расширить круг задач, традиционно решаемых стилистикой и авторской стилиметрией.</p>	УК -4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3

Тестирование (УК-4, ИД-УК-4.1, ИД-УК-4.2)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.1 Подготовка и редактирование различных академических текстов.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

Тест

Время выполнения 15 мин.

Количество вопросов 10.

Форма работы – самостоятельная, индивидуальная.

Способ проведения теста: бланковый

Инструкция для тестируемых:

Внимательно читать задания к тестовым вопросам и четко отвечать на них в зависимости от типа задания. Тест выполняется самостоятельно в течение 15 минут.

Номинальная шкала предполагает, что за правильный ответ к каждому заданию выставляется один балл, за не правильный – ноль. В соответствии с номинальной шкалой, оценивается всё задание в целом, а не какая-либо из его частей.

Процентное соотношение баллов и оценок по пятибалльной системе:

«2» - равно или менее 54%

«3» - 55% - 69%

«4» - 70% - 84%

«5» - 85% - 100%.

Инструкция для проверяющих:

Правильные ответы тестовых заданий выделены. Номинальная шкала предполагает, что за правильный ответ к каждому заданию выставляется один балл, за не правильный — ноль. В соответствии с номинальной шкалой, оценивается всё задание в целом, а не какая-либо из его частей.

Процентное соотношение баллов и оценок по пятибалльной системе:

«2» - равно или менее 54%

«3» - 55% - 69%

«4» - 70% - 84%

«5» - 85% - 100%.

Тест с ответами

№	Вопрос	Ответ	Компетенция
	<p>УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.</p> <p>ИД-УК-4.1 Подготовка и редактирование различных академических текстов.</p> <p>ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.</p>		
1	<p>1. Корпусная лингвистика занимается</p> <p>а) объяснением фактов языка;</p> <p>б) объяснением фактов речи;</p> <p>в) изучением системности и структурности языка;</p>		<p>УК-4</p> <p>ИД-УК-4.1</p> <p>ИД-УК-4.2</p>

	г) типологией языковых ситуаций.	
2	2. Машинный перевод является предметом описания а) структурной лингвистики; б) компьютерной лингвистики; в) корпусной лингвистики; г) социолингвистики.	УК-4 ИД-УК-4.1 ИД-УК-4.2
3	Какой из представленных корпусов реально не существует? а) Мангеймский корпус немецкого языка; б) Национальный корпус русского языка; в) Упсальский корпус русского языка; г) Национальный корпус французского языка; д) Корпус современного американского английского языка.	УК-4 ИД-УК-4.1 ИД-УК-4.2
4	Лингвистический корпус – это а) собрание текстов, сборники, которые носят довольно случайный характер; б) вид издания, которое концентрирует в однотипно оформленных единицах (томах) все или основные произведения одного автора в качестве его научного, литературно-художественного или публицистического наследия; в) тексты, собранные и размеченные по определённому стандарту и обеспеченные специализированной поисковой системой; г) собрание письменных сообщений, объективированных в виде письменных документов, имеющих прагматическую установку и соответственно литературно обработанных; д) информационно-справочная система, основанная на собрании текстов на некотором языке в электронной форме, представляет данный язык на определенном этапе (или этапах) его существования и во всём многообразии жанров, стилей, территориальных и социальных вариантов.	УК-4 ИД-УК-4.1 ИД-УК-4.2
5	Выберите критерии, по которым могут выделяться типы корпусов текстов: а) по языку представления текстов б) по статистическому критерию; в) по жанровой принадлежности; г) по форме хранения; д) по удобству чтения; е) по возможностям читательского восприятия.	УК-4 ИД-УК-4.1 ИД-УК-4.2
6	Цель создания корпуса – а) для отражения художественного содержания текстов; б) для работы в области машинного перевода; в) для научных исследований и обучения языку; г) для создания образовательной среды, объединяющей литературные произведения в единую систему.	УК-4 ИД-УК-4.1 ИД-УК-4.2
7	Найдите среди представленных параметров типы разметок корпуса: а) метаязыковой; б) семантический; в) гносеологический; г) морфологический.	УК-4 ИД-УК-4.1 ИД-УК-4.2
8	Какие из представленных единиц являются подкорпусами Национального корпуса русского языка: а) параллельный корпус; б) обучающий корпус русского языка; в) стратификационный корпус; г) газетный корпус; д) спряжение глаголов; е) динамический корпус; ж) речевой корпус.	УК-4 ИД-УК-4.1 ИД-УК-4.2
9	Разбиение на орфографические слова в корпусе – это а) конкорданс;	УК-4 ИД-УК-4.1

	б) токенизация; в) лемматизация; г) кластеризация.	ИД-УК-4.2															
10	Найдите ошибку в сопоставлении традиционной и корпусной лингвистики	УК-4 ИД-УК-4.1 ИД-УК-4.2															
	<table border="1"> <thead> <tr> <th></th> <th>Корпусная лингвистика</th> <th>Традиционная лингвистика</th> </tr> </thead> <tbody> <tr> <td>а)</td> <td>Текст рассматривается как некоторая физическая сущность.</td> <td>Текст рассматривается как некоторая абстракция.</td> </tr> <tr> <td>б)</td> <td>Составление грамматики конкретных языков.</td> <td>Изучает языковые универсалии</td> </tr> <tr> <td>в)</td> <td>Основное внимание уделяется форме.</td> <td>Основное внимание не только форме, но и содержанию.</td> </tr> <tr> <td>г)</td> <td>Предпочитаются искусственные примеры, из изолированных от текста словоупотреблений.</td> <td>Проводится работа с лингвистическими данными (словоупотреблениями) в том виде, в каком они встречались в контексте.</td> </tr> </tbody> </table>		Корпусная лингвистика	Традиционная лингвистика	а)	Текст рассматривается как некоторая физическая сущность.	Текст рассматривается как некоторая абстракция.	б)	Составление грамматики конкретных языков.	Изучает языковые универсалии	в)	Основное внимание уделяется форме.	Основное внимание не только форме, но и содержанию.	г)	Предпочитаются искусственные примеры, из изолированных от текста словоупотреблений.	Проводится работа с лингвистическими данными (словоупотреблениями) в том виде, в каком они встречались в контексте.	
	Корпусная лингвистика	Традиционная лингвистика															
а)	Текст рассматривается как некоторая физическая сущность.	Текст рассматривается как некоторая абстракция.															
б)	Составление грамматики конкретных языков.	Изучает языковые универсалии															
в)	Основное внимание уделяется форме.	Основное внимание не только форме, но и содержанию.															
г)	Предпочитаются искусственные примеры, из изолированных от текста словоупотреблений.	Проводится работа с лингвистическими данными (словоупотреблениями) в том виде, в каком они встречались в контексте.															

4.2. Оценочные материалы промежуточного контроля успеваемости по учебной дисциплине, в том числе самостоятельной работы обучающегося, типовые задания

Зачет с оценкой (УК-4, ИД-УК-4.1, ИД-УК-4.2, ИД-УК-4.3)

УК-4 Способен применять современные коммуникативные технологии, в том числе на иностранном (ых) языке (ах), для академического и профессионального взаимодействия.

ИД-УК-4.1 Подготовка и редактирование различных академических текстов.

ИД-УК-4.2 Готовность к участию в профессиональных дискуссиях и грамотное использование деловой, устной и письменной коммуникации.

ИД-УК-4.3 Навыки межличностного делового общения, в том числе на иностранных языках с применением профессиональных языковых форм и средств.

Устный опрос по вопросам:

Время на подготовку 20 мин

Перечень теоретических вопросов к зачету:

Вопрос	Ответ	Компетенции
1. Компьютерная лингвистика как направление научной деятельности, ее цели и задачи.	<i>Компьютерная лингвистика – направление в прикладной лингвистике, ориентированное на использование программ, компьютерных технологий организации и обработки данных для моделирования функционирования языка. создание и использование электронных корпусов текстов; создание электронных словарей, тезаурусов; автоматический перевод текстов; автореферирование; автоматическое распознавание речи и автоматический синтез речи.</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3
2. Корпусная лингвистика как новое	<i>Корпусная лингвистика – это раздел прикладной лингвистики, занимающийся</i>	УК-4 ИД-УК-4.1

направление в языкознании.	<i>разработкой общих принципов построения и использованием лингвистических корпусов. Корпус – это информационно-справочная система (специализированная поисковая система), собрание снабженных стандартной разметкой текстов на определенном языке в электронной форме.</i>	ИД-УК-4.2 ИД-УК-4.3
3. Критерии создания корпусов	<i>Стратегия построения корпуса. Репрезентативность (представительность) – способность корпуса текстов отражать все свойства проблемной области. Репрезентативность корпуса указывает на то, что единицы проблемной области отражаются пропорционально в корпусе данных. Полнота требует учета релевантных явлений, даже если это не соответствует идее пропорционального сужения. Экономичность корпуса текстов – экономия усилий исследователя при изучении проблемной области. Чем более «экономичен» корпус, тем выше порог отбора. Структуризация материала – описание данных корпуса, в которой единицы хранения характеризуются по тем параметрам, которые могут оказаться важными для пользователя. Компьютерная поддержка – комплекс программ по обработке данных, обеспечивающих функции составления конкордансов, статистической инвентаризации, автоматической словарной обработки, лемматизации.</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3
4. Корпусная лингвистика и проблемы перевода	<i>Создание автоматизированных лексикографических систем в помощь переводчику. Создание, ведение и издание переводных словарей. Задачи – выбор принципов лексикографирования и создания специального «словаря для переводчика» (определение макроструктуры словаря); выбор формата представления информации в словаре (определение микроструктуры словаря); выбор принципов и методов формирования и использования корпуса текстов. Тексты оригинала и перевода являются параллельными или битекстами. Параллельный корпус текстов – совокупность документов, переведенных на два или более языков, выровненных по предложениям и размеченных, написанных на одну тему и на одном языке авторами с разными родными языками. Решение проблем перевода с использованием корпуса: интерпретация оригинала и текстопостроения перевода; определение лексической сочетаемости; сохранение стилистического рисунка исходного текста, правильность выбранных грамматических и синтаксических конструкций; сокращение затрачиваемого времени на перевод.</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3
5. Британские корпуса и словари.	<i>Британский национальный корпус (British National Corpus). Время создания – 1990-е гг. Объем этого знаменитого корпуса составляет 100 миллионов словоупотреблений. BNC – эталон лингвистического корпуса. Содержит образцы письменного и разговорного британского английского языка из широко-</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3

	<p>го круга источников. Корпус охватывает британский английский конца XX века, представленный широким разнообразием жанров, и задуман как образец типичного разговорного и письменного британского английского языка того времени. BNC является одноязычным корпусом, так как он содержит образцы только британского английского языка, хотя иногда в текстах встречаются слова и фразы из других языков. Это синхронический корпус, так как в нём содержатся примеры использования языка только одного временного периода – конец XX века. По этой причине BNC не может служить источником данных об истории развития британского варианта английского языка. BNC сбалансированный корпус, включает данные из различных источников. Корпус BNC содержит частеречную разметку. Включает письменный и разговорный подкорпусы. Недостаток корпуса – слишком общая классификация текстов.</p>	
<p>6. Корпуса американского варианта английского языка.</p>	<p>Корпус современного американского английского (<i>Corpus of Contemporary American English, COCA</i>) – электронный корпус текстов, созданный профессором корпусной лингвистики Марком Дэвисом из Университета Бригама Янга в 2000–2003 годах на основе текстов журнала <i>Time</i>, написанных с 1923 года. Это наибольший (450 млн слов) корпус текстов американского варианта английского языка. Свободно доступный корпус, включающий большое разнообразие текстов различных жанров. Он составлен из более чем 160 тыс. текстов, включая по 20 млн слов за каждый год с 1990 по 2011. Это наиболее широко используемый структурированный корпус текстов, ежемесячно его используют примерно 10 000 человек. Архитектура корпуса – принцип других проектов Марка Дэвиса. Существует центральная база n-граммов, которая содержит информацию о каждом из ста миллионов слов корпуса. Они связаны с таблицами, позволяющими анализировать регистр, а также с отдельными таблицами для синонимов, лемм и форм, появляющихся у слова с течением времени. Характерная особенность корпуса – быстрый поиск (занимает менее секунды даже для самых сложных запросов, содержащих словоформу, часть речи, частоту и регистр). Корпус позволяет исследовать: изменение частоты и контекста использования слов и фраз, связанных с переменами в культурной и социальной жизни общества; языковые перемены в морфологии и грамматических конструкциях, колебания частоты использования тех или иных групп слов со временем, семантические изменения слов на протяжении XX века.</p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>
<p>7. НКРЯ как инструмент семантико-грамматического исследования лек-</p>	<p>Национальный корпус русского языка (НКРЯ, Корпус) – собрание независимых корпусов,</p>	<p>УК-4 ИД-УК-4.1</p>

СИКИ.	<p><i>основной инструмент поиска при лингвистических исследованиях. Корпусы большие по объёму и представительные – ценный материал для количественных и качественных исследований. Структура корпуса – основной (тексты 18-21 вв.), газетный, поэтический, акцентологический, устный корпус, корпус социальных сетей, мультимедийный, синтаксическом, обучающий, исторический, панхронический, диалектный, параллельный корпуса, корпус «От 2 до 15», корпус региональной прессы. Объём основного корпуса на декабрь 2023 года – 375 млн словоупотреблений, а общий объём корпусов – больше 2 млрд словоупотреблений. Тексты снабжены метаразметкой (по дате создания, автору, жанру и тому подобному); словоформы в текстах снабжены автоматической морфологической и семантической разметкой; параллельные тексты выровнены; тексты поэтического корпуса снабжены также особой метрической разметкой.</i></p>	ИД-УК-4.2 ИД-УК-4.3
8. Статистическое описание лексики в писательской лексикографии.	<p><i>Приход больших данных в лингвистику, развитие нейронных сетей позволяют по-новому исследовать материал по авторской лексикографии. Компьютеризация лексикографической деятельности заключается прежде всего в создании специализированных машинных банков данных и в разработке методов формирования этих банков, представления информации в банках и её использовании. На этой основе формируется целое новое направление лингвистики и лексикографии – корпусная лингвистика и лексикография. Ряд специализированных проектов ориентирован на исследование сочетаемости ключевых слов на материале русского языка с привлечением статистического аппарата. Таким ресурсом является система CoCoCo (“Collocations, Colligations, Constructions”), которая разрабатывается под руководством М. Копотева в Хельсинкском университете, представляет информацию о многословных выражениях на основе Национального корпуса русского языка. Многословные выражения понимаются предельно широко: к ним относятся идиомы, составные единицы, коллокации и коллигации. В системе предусмотрен поиск по леммам или токенам (словоформам), который можно ограничить морфологическими параметрами. При выборе части речи есть возможность задать значения грамматических категорий. Результаты поиска демонстрируют не все возможные словосочетания, а только наиболее значимые. Подобная информация основывается на применении статистических метрик. Для решения этой задачи используется расстояние Кульбака–Лейблера (мера KLD), которое позволяет оценить, какие грамматические признаки оказываются наиболее употребительными для данной лексической единицы.</i></p>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3

<p>9. Автоматический семантический анализ текста на русском языке.</p>	<p><i>Семантический анализ текста – наиболее сложная проблема в области искусственного интеллекта и компьютерной лингвистика. Главная проблема заключается в том, как «научить» компьютер однозначно верно трактовать образы, которые пытался передать автор текста. Цель семантического (смыслового) анализа – оценивание смысла передаваемой информации, соотношения ее с информацией, которая хранилась до появления обрабатываемой информации. Семантические связи между словами или другими единицами языка отражаются в семантических словарях. Задачами семантического анализа являются: 1) построение семантической интерпретации слов и конструкций; 2) установление семантических отношений между различными элементами текста. При семантическом анализе предложений используют надежные грамматики и семантические валентности, а семантика предложения задается через связи главного слова (глагола) с его семантическими актантами. Основой семантического анализа – утверждение, что конкретное значение слова не является элементарной семантической единицей. Оно делится на более мелкие единицы – единицы словаря семантического языка, являющиеся своеобразными атомами, комбинации которых складываются в «молекулы» – значения слов естественного языка. Именно семантический анализ дает возможность решить проблемы многозначности (омонимии), которая часто возникает при автоматическом анализе на разных языковых уровнях. В заключении стоит отметить, что ценность автоматического анализа текста на данный момент особенно высока, поскольку человек уже не в состоянии самостоятельно обработать современные объемы информации. Автоматический анализ текста находит применение в филологии (определение авторства произведений, авторского стиля), в экспертных системах.</i></p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>
<p>10. Исходные понятия корпусной лингвистики: проблемная область, корпус данных, корпус текстов.</p>	<p><i>Проблемная область – область реализаций языковой системы, содержащая феномены, подлежащие лингвистическому описанию, она имеет два измерения – языковое и речевое. В корпусной лингвистике языковой аспект фактически игнорируется, поскольку изначально фиксируется область привлекаемых данных – реализаций языковой системы. Проблемная область для разработчика – корпус как множество данных, обработка которых затруднена из-за того, что языковых реализаций слишком много. Корпус данных – сформированная по определенным правилам выборка данных из проблемной области, результат отображения из проблемной области. В отличие от проблемной области, корпус данных имеет только одно измерение – речевое. Лингвисту приходится по отдель-</i></p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>

	<p>ным результатам деятельности языка делать выводы о функционировании языка как целого, как системы. Корпус текстов – это вид корпуса данных, единицами которого являются тексты или их достаточно значительные фрагменты, полные фрагменты макроструктуры текстов данной проблемной области. Типы корпусов данных – исследовательские корпусы (для изучения различных аспектов функционирования языковой системы), иллюстративные корпусы (для подтверждения и обоснования уже полученные результаты), динамические (функционирования языковых феноменов на временной шкале) и статические корпусы текстов (временное состояние языковой системы – авторские корпусы).</p>	
<p>11. Параллельный многоязычный корпус текстов, его структура и сфера применения.</p>	<p>Параллельный корпус (<i>Parallel Corpora</i>) – электронный аналог параллельных переводных текстов, состоящий из множества блоков «текст-оригинал и один/несколько его переводов». Электронные тексты в корпусе могут представлять собой целое оригинальное словесное произведение или какую-либо его часть.</p> <p>В современной корпусной лингвистике выделяется два вида параллельных корпусов: 1) многоязычный, или <i>Comparable (Multilingual) Corpora</i>, 2) переводной, или <i>Translation Corpora</i>.</p> <p>Структурная организация корпуса: в виде традиционного текста со ссылкой на перевод/ы, в табличной «зеркальной» форме, что более удобно для восприятия и сравнения, в виде базы данных. Параллельные корпуса – прикладной потенциал в методике обучения иностранным языкам и переводу, а также в компьютерной лингвистике. Переводчику необходимы ресурсы-эталон перевода и оценка перевода в «стандартных» условиях. Электронные параллельные корпусы и лингвистические компьютерные технологии позволяют значительно сократить временные затраты и предоставляют образцы профессионального перевода при изучении приемов и способов перевода.</p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>
<p>12. Аннотированные корпусы текстов, автоматизация их создания и коррекции.</p>	<p>Большие текстовые корпуса давно и плодотворно используются в компьютерной лингвистике. Уровни аннотации текста: 1) лемматизированные тексты, в которых для каждого слова указывается его основная форма и часть речи; 2) тексты с морфологической информацией, в которых для каждого слова указываются его основная форма, часть речи и полный набор морфологических характеристик; 3) тексты с синтаксической информацией, в которых для каждого слова указываются его основная форма, часть речи и морфологические характеристики, и для каждого предложения указывается его синтаксическая структура. Построение аннотированного корпуса осуществляется в полу-</p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>

	<p>автоматическом режиме: аннотация вначале порождается системой автоматического морфологического и синтаксического анализа, а затем корректируется специалистом-лингвистом. Степень участия лингвиста в процессе аннотации определяется самим лингвистом в зависимости от сложности структуры текста. Типы лингвистической информации на каждом уровне: морфологическая информация – морфологические характеристики слова (часть речи, одушевленность, род, падеж, число, падеж, степень сравнения, краткость, репрезентация, вид, время, лицо, залог); синтаксическая информация – число узлов в синтаксической структуре равно числу слов в предложении. Инвентарь синтаксических отношений – деление на 6 больших групп: 1) актантные; 2) атрибутивные; 3) количественные; 4) обстоятельственные; 5) сочинительные; 6) служебные. Поиск в корпусе другой информации: распределение по годам, частота употребления слова по годам; статистика (метаатрибуты: автор, пол автора, сфера функционирования, тип текста, тематика текста, жанр).</p>	
<p>13. Опыт разработки корпусов текстов в России и за рубежом.</p>	<p>руководством У. Фрэнсиса существует в компьютерном варианте и на микрофишах. Объем корпуса около 1 млн словоупотреблений. Корпус состоит из 500 текстов, каждый из которых включает 2 000 словоупотреблений). Ланкастерско-Осло-Бергенский корпус (LOB), британский вариант английского языка. Аннотированная версия корпуса LOB появилась в 1985 г. Лондонско-Лундский корпус (фиксация особенности грамматической системы английского языка в речи взрослого образованного носителя; 1960-е г. под руководством Рэндола Квирка в Лондонском университетском колледже; объем корпуса – 1 млн словоупотреблений, письменные текст и тексты устной речи). Бирмингемский корпус (мониторный / динамический корпус, автор проекта – Дж. Синклер, 7,3 млн словоупотреблений, 6 млн письменные тексты, 1,3 млн – устные тексты). Разработка Германии проекта LIMAS-корпуса (система немецко-английского машинного перевода). Корпусы текстов немецкой разговорной речи. «Корпус базового немецкого» (1961 г. в Стэнфорде). Корпусная лингвистика во Франции («Сокровищница французского языка»). Корпусы текстов по русскому языку (начало 1970-х годов). «Уппсальский машинный фонд русского языка», создававшийся с 1987 г. в Уппсальском университете. Собственно российский опыт составления корпусов: Корпус по дискурсивным словам русского языка, Корпус словаря языка Достоевского, Корпус текстов словаря языка Достоевского, Динамический корпус текстов по современной публицистике (90-е гг.), Динами-</p>	<p>УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3</p>

	<i>ческий корпус текстов как новая технология прикладной лингвистики.</i>	
14. Практическое использование аннотированных корпусов текстов в системах автоматической обработки текстов.	<i>Аннотирование текстов – это процесс добавления метаданных и лингвистической информации к текстам в корпусе. Для эффективного поиска и анализа текстов в корпусе необходимо создать поисковые индексы. Аннотирования корпусов требует тщательной работы и экспертизы в области лингвистики, разработки интерфейса для взаимодействия пользователей с корпусом. Примеры практического применения корпусов в лингвистике. Компьютерная обработка текста – преобразование текста на искусственном или естественном языке с помощью компьютера. Системное программирование – создание программного обеспечения функционирования компьютера и работы пользователей. Издательское дело – одно из направлений автоматизации редакционно-издательских процессов. Автоматизированное редактирование – внесение в текст, находящийся в памяти компьютера, исправлений и дополнений. Эти тенденции прогнозируются и прослеживаются на примере развития АОТ-систем (АОТ – автоматическая обработка текста), представляющих коммерческий интерес. Решение следующих прикладных задач – машинного перевода, генерации текста, локализации и интернационализации, работа на ограниченном языке, создание текстовых документов (ввод, редактирование, исправление ошибок), информационный поиск и связанные с ним задачи.</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3
15. Корпус и междисциплинарные исследования	<i>Междисциплинарное исследование – исследование одного объекта методами различных дисциплин: методология нескольких дисциплин – базовая, а другие – новые, инновационными по отношению к исследуемому объекту. Практическая целесообразность междисциплинарного исследования в корпусной лингвистике: корпусы текстов – продукт синергии а) лингвистических дисциплин, б) разработок в области машинной обработки языка, основанных на статистических и кибернетических методах, в) собственно информационных технологий. Исследования в области корпусной лингвистики – это способы и результаты усовершенствования подборки и управления корпусным материалом и усовершенствование имеющихся и разработка новых способов синтаксического, семантического и прочих видов анализа. Важным компонентом исследования является прагматика корпуса текстов – взаимодействие с пользователем. Два направления междисциплинарных исследований на основе корпусной лингвистики – фундаментальный (понимание), тяготеющие к психологии, и прикладной (внешнее моделирование полезных свойств), близкие к когнитивной науке.</i>	УК-4 ИД-УК-4.1 ИД-УК-4.2 ИД-УК-4.3

ЛИСТ УЧЕТА ОБНОВЛЕНИЙ ОЦЕНОЧНЫХ СРЕДСТВ УЧЕБНОЙ ДИСЦИПЛИНЫ

В оценочные средства учебной дисциплины внесены *изменения/обновления*, утверждены на заседании кафедры:

№ пп	год обновления оценочных средств	номер протокола и дата заседания кафедры