

МИНИСТЕРСТВО ТРАНСПОРТА РОССИЙСКОЙ ФЕДЕРАЦИИ
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ
УЧРЕЖДЕНИЕ ВЫСШЕГО ОБРАЗОВАНИЯ
«РОССИЙСКИЙ УНИВЕРСИТЕТ ТРАНСПОРТА»
(РУТ (МИИТ))



Рабочая программа дисциплины (модуля),
как компонент образовательной программы
высшего образования - программа бакалавриата
по направлению подготовки
09.03.01 Информатика и вычислительная техника,
утвержденной первым проректором РУТ (МИИТ)
Тимониным В.С.

РАБОЧАЯ ПРОГРАММА ДИСЦИПЛИНЫ (МОДУЛЯ)

Методы анализа и обработки больших данных

Направление подготовки: 09.03.01 Информатика и вычислительная техника

Направленность (профиль): IT-сервисы и технологии обработки данных на транспорте

Форма обучения: Очная

Рабочая программа дисциплины (модуля) в виде электронного документа выгружена из единой корпоративной информационной системы управления университетом и соответствует оригиналу

Простая электронная подпись, выданная РУТ (МИИТ)
ID подписи: 937226
Подписал: руководитель образовательной программы
Проневич Ольга Борисовна
Дата: 19.05.2025

1. Общие сведения о дисциплине (модуле).

Целью освоение дисциплины (модуля) является приобретения практически навыков и умений, позволяющих создавать высокопроизводительные реализации известных методов вычислительной математики, анализа и обработки больших данных.

Задачами освоения дисциплины (модуля) являются:

освоение базовых знаний в области архитектуры современных многопроцессорных вычислительных систем параллельной обработки информации,

освоений технологий организации параллельных вычислений на многопроцессорных вычислительных комплексах с распределенной или общей оперативной памятью.

2. Планируемые результаты обучения по дисциплине (модулю).

Перечень формируемых результатов освоения образовательной программы (компетенций) в результате обучения по дисциплине (модулю):

ПК-1 - Способен анализировать большие данные с использованием существующей в организации методологической и технологической инфраструктуры;

ПК-7 - Способен к организации процессов разработки программного обеспечения .

Обучение по дисциплине (модулю) предполагает, что по его результатам обучающийся будет:

Знать:

программное обеспечение, необходимое для работы с большими данными,

необходимые элементы инфраструктуры для обработки больших данных.

Уметь:

Обрабатывать данные в Hadoop

Работать с распределенными файловыми системами в Hadoop

Настраивать кластеры в Hadoop

Составлять спецификации оборудования для различных задач

Владеть:

Инструментами настройки сетей

Знаниями о высокопроизводительных вычислениях

Навыками работы со статистическими параметрами вычислений
Технологиями параллельной обработки данных
Навыками работы с vCenter Server
Инструментами работы с Hadoop
Техниками настройки кластеров в Hadoop

3. Объем дисциплины (модуля).

3.1. Общая трудоемкость дисциплины (модуля).

Общая трудоемкость дисциплины (модуля) составляет 3 з.е. (108 академических часа(ов)).

3.2. Объем дисциплины (модуля) в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении учебных занятий:

Тип учебных занятий	Количество часов	
	Всего	Семестр №5
Контактная работа при проведении учебных занятий (всего):	64	64
В том числе:		
Занятия лекционного типа	32	32
Занятия семинарского типа	32	32

3.3. Объем дисциплины (модуля) в форме самостоятельной работы обучающихся, а также в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении промежуточной аттестации составляет 44 академических часа (ов).

3.4. При обучении по индивидуальному учебному плану, в том числе при ускоренном обучении, объем дисциплины (модуля) может быть реализован полностью в форме самостоятельной работы обучающихся, а также в форме контактной работы обучающихся с педагогическими работниками и (или) лицами, привлекаемыми к реализации образовательной программы на иных условиях, при проведении промежуточной аттестации.

4. Содержание дисциплины (модуля).

4.1. Занятия лекционного типа.

№ п/п	Тематика лекционных занятий / краткое содержание
1	<p>Тема 1. Большие данные. История развития методов анализа и обработки</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - эволюция понятия больших данных - хранилища данных - история развития методов обработки больших данных - отличия методов анализа от методов обработки больших данных
2	<p>Тема 2. Жизненный цикл данных и метаданные</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - элементы жизненного цикла и методы управления жизненным циклом - метаданные
3	<p>Тема 3. Архитектура систем обработки больших данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - прием данных - сбор данных - анализ данных - представления результатов - витрины данных
4	<p>Тема 4. Введение в понятия высокопроизводительных вычислений.</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Важность проблематики параллельных вычислений - Пути достижения параллелизма. Векторная и конвейерная обработка данных. - Многопроцессорная и многомашинная, параллельная обработка данных. - Закон Мура, сдерживающие факторы наращивания количества транзисторов на кристалле и частоты процессоров. Привлекательность подхода параллельной обработки данных
5	<p>Тема 5. Основные направления развития высокопроизводительных компьютеров.</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Сдерживающие факторы повсеместного внедрения параллельных вычислений - Ведомственные, национальные и другие программы, направленные на развитие параллельных вычислений в России. - Необходимость изучения дисциплины параллельного программирования. Перечень критических задач, решение которых без использования параллельных вычислений затруднено или вовсе невозможно.
6	<p>Тема 6. Задачи параллельной обработки данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - программные структуры алгоритмов параллельной обработки данных - инструменты параллельной обработки данных - конвейерно-параллельной обработки интегрированных потоков данных
7	<p>Тема 7. Классификация вычислительных сетей</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Системы с распределенной, общей памятью, примеры систем. - Массивно-параллельные системы (MPP). Симметричные мультипроцессорные системы (SMP). - Параллельные векторные системы (PVP). - Системы с неоднородным доступом к памяти (Numa) - Компьютерные кластеры – специализированные и полнофункциональные. История возникновения компьютерных кластеров – проект Weowulf. Метакомпьютинг. Классификация Флинна, Шора и т.д. Организация межпроцессорных связей – коммуникационные топологии. - Примеры сетевых решений для создания кластерных систем
8	<p>Тема 8. Основные принципы организации параллельной обработки данных: модели, методы и технологии параллельного программирования</p> <p>Рассматриваемые вопросы:</p>

№ п/п	Тематика лекционных занятий / краткое содержание
	<p>-Функциональный параллелизм, параллелизм по данным.</p> <p>-Парадигма master-slave. Парадигма SPMD. Парадигма конвейеризации. Парадигма «разделяй и властвуй». Спекулятивный параллелизм. Важность выбора технологии для реализации алгоритма</p> <p>-Модель обмена сообщениями – MPI.</p> <p>-Модель общей памяти – OpenMP. Концепция виртуальной, разделяемой памяти – Linda.</p> <p>Российские разработки – Т-система, система DVM. Проблемы создания средства автоматического распараллеливания программ</p>
9	<p>Тема 9. Параллельное программирование с использованием интерфейса передачи сообщений MPI</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> -Библиотека MPI. - Модель SIMD. Инициализация и завершение MPI-приложения. <p>Точечные обмены данными между процессами MPI-программы. Режимы буферизации.</p> <p>Проблема deadlock'ов. Коллективные взаимодействия процессов в MPI. Управление группами и коммутаторами в MPI</p>
10	<p>Тема 10. Параллельное программирование на системах с общей памятью (OpenMP)</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> -Введение в OpenMP -Стандарты программирования для систем с разделяемой памятью. -Создание многопоточных приложений. Использование многопоточности при программировании для многоядерных платформ. -Синхронизация данных между ветвями в параллельной программе. Директивы языка OpenMP
11	<p>Тема 11. Параллельное программирование многоядерных GPU. Кластеры из GPU и суперкомпьютеры на гибридной схеме</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> -Существующие многоядерные системы. -GPU. -Использование OpenMP и MPI технологий совместно с CUDA. -Степень параллелизма численного алгоритма. Закон Амдала. Параллельный алгоритм решения СЛАУ
12	<p>Тема 12. Программные платформы и системы для больших данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - системы управления потоками данных - системы хранения больших данных - платформы больших данных - обработка больших данных в реальном времени

4.2. Занятия семинарского типа.

Практические занятия

№ п/п	Тематика практических занятий/краткое содержание
1	<p>Тема 1. Изучение технологий Hadoop и MapReduce</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> - Hadoop - MapReduce
2	<p>Тема 2. Анализ больших массивов данных инструментами python</p> <p>Рассматриваемые вопросы:</p>

№ п/п	Тематика практических занятий/краткое содержание
	- библиотеки анализа и загрузки больших данных - фреймворки Python с параллельной обработкой данных
3	Тема 3. Высокопроизводительные вычисления с фреймворком Apache Spark Рассматриваемые вопросы: - знакомство фреймворком - знакомство с архитектурой системы - реализации концепции MapReduce - разработки программ с использованием Apache Spark
4	Тема 4. Введение в высокопроизводительные серверы на python Рассматриваемые вопросы: - работа с распределенными системами - создание кода для работы на GRU
5	Тема 5. Параллельное программирование. Параллельная обработка данных Рассматриваемые вопросы: 1. Параллельное программирование 2. Параллельная обработка данных
6	Тема 6. Работа с vCenter Server Рассматриваемые вопросы: 1. Понятие vCenter 2. Установка и развертывание vCenter
7	Тема 7. Развертывание среды для обработки данных при помощи Hadoop Рассматриваемые вопросы: 1. Обработка данных в Hadoop 2. Распределенная файловая система Hadoop 3. Разработка приложений MapReduce 4. Настройка кластера Hadoop
8	Тема 8. Составление спецификаций оборудования для работы с высокопроизводительными вычислениями Рассматриваемые вопросы: 1. Составление спецификаций оборудования для работы с высокопроизводительными вычислениями

4.3. Самостоятельная работа обучающихся.

№ п/п	Вид самостоятельной работы
1	Работа с учебной литературой
2	Участие в онлайн-конференциях и мастер-классах
3	Поиск алгоритмов обработки данных в открытых источниках
4	Подготовка к промежуточной аттестации.
5	Подготовка к текущему контролю.

5. Перечень изданий, которые рекомендуется использовать при освоении дисциплины (модуля).

№ п/п	Библиографическое описание	Место доступа
1	Железнов, М. М. Методы и технологии обработки больших данных : учебно-методическое пособие / М. М. Железнов. — Москва : МИСИ – МГСУ, 2020. — 46 с. — ISBN 978-5-7264-2193-3.	https://e.lanbook.com/book/145102
2	Ланских, Ю. В. Введение в большие данные : учебное пособие / Ю. В. Ланских, В. Г. Ланских, К. В. Родионов. — Киров : ВятГУ, 2023. — 172 с.	https://e.lanbook.com/book/408566

6. Перечень современных профессиональных баз данных и информационных справочных систем, которые могут использоваться при освоении дисциплины (модуля).

<https://habr.com/ru> - база знаний в виде статей, обзоров

<https://journal.tinkoff.ru/short/ai-for-all/> - база данных нейронных сетей

<https://vc.ru/services/916617-luchshie-neyroseti-bolshaya-podborka-iz-top-200-ii-generatorov-po-kategoriyam> - база данных нейронных сетей

<https://github.com/abalmumcu/bert-rest-api> - профессиональная платформа для командой работы над проектов (нейронная сеть bert)

<http://library.miit.ru/> - электронно-библиотечная система Научно-технической библиотеки МИИТ

<https://proglib.io/p/raspoznavanie-obektov-s-pomoshchyu-yolo-v3-na-tensorflow-2-0-2020-11-08> - профессиональная библиотека программистов

https://yandex.cloud/ru/blog/posts/2022/12/andrey-berger-and-yandex-cloud?utm_referrer=https%3A%2F%2Fyandex.ru%2F - библиотека профессиональных статей разработчиков Яндекс

<https://yandex.cloud/ru/blog> - библиотека профессиональных статей разработчиков Яндекс

<https://tproger.ru/translations/opencv-python-guide> - библиотека основных команд OpenCV

7. Перечень лицензионного и свободно распространяемого программного обеспечения, в том числе отечественного производства, необходимого для освоения дисциплины (модуля).

Microsoft Office 2010

VMware Workstation

8. Описание материально-технической базы, необходимой для осуществления образовательного процесса по дисциплине (модулю).

1 учебный класс

Компьютер преподавателя

Компьютеры студентов

Монитор

Проектор

Экран для проектора

Маркерная доска

9. Форма промежуточной аттестации:

Зачет в 5 семестре.

10. Оценочные материалы.

Оценочные материалы, применяемые при проведении промежуточной аттестации, разрабатываются в соответствии с локальным нормативным актом РУТ (МИИТ).

Авторы:

руководитель образовательной
программы

О.Б. Проневич

Согласовано:

Директор

Д.В. Паринов

Руководитель образовательной
программы

О.Б. Проневич

Председатель учебно-методической
комиссии

Д.В. Паринов