

**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 - программы 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: 12.03.2025

## 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 5 з.е. (180 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	16	16
Занятия семинарского типа	48	48

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 116 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"> <li>- эволюция понятия больших данных</li> <li>- хранилища данных</li> <li>- история развития методов обработки больших данных</li> <li>- отличия методов анализа от методов обработки больших данных</li> </ul>
2	<p>Тема 2. Жизненный цикл данных и метаданные</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- элементы жизненного цикла и методы управления жизненным циклом</li> <li>- метаданные</li> </ul>
3	<p>Тема 3. Архитектура систем обработки больших данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- прием данных</li> <li>- сбор данных</li> <li>- анализ данных</li> <li>- представления результатов</li> <li>- витрины данных</li> </ul>
4	<p>Тема 4. Введение в понятия высокопроизводительных вычислений.</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- Важность проблематики параллельных вычислений</li> <li>- Пути достижения параллелизма. Векторная и конвейерная обработка данных.</li> <li>- Многопроцессорная и многомашинная, параллельная обработка данных.</li> <li>- Закон Мура, сдерживающие факторы наращивания количества транзисторов на кристалле и частоты процессоров. Привлекательность подхода параллельной обработки данных</li> </ul>
5	<p>Тема 5. Основные направления развития высокопроизводительных компьютеров.</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- Сдерживающие факторы повсеместного внедрения параллельных вычислений</li> <li>- Ведомственные, национальные и другие программы, направленные на развитие параллельных вычислений в России.</li> <li>- Необходимость изучения дисциплины параллельного программирования. Перечень критических задач, решение которых без использования параллельных вычислений затруднено или вовсе невозможно.</li> </ul>
6	<p>Тема 6. Задачи параллельной обработки данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- программные структуры алгоритмов параллельной обработки данных</li> <li>- инструменты параллельной обработки данных</li> <li>- конвейерно-параллельной обработки интегрированных потоков данных</li> </ul>
7	<p>Тема 7. Классификация вычислительных сетей</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- Системы с распределенной, общей памятью, примеры систем.</li> <li>- Массивно-параллельные системы (MPP). Симметричные мультипроцессорные системы (SMP). - Параллельные векторные системы (PVP).</li> <li>- Системы с неоднородным доступом к памяти (Numa)</li> <li>- Компьютерные кластеры – специализированные и полнофункциональные. История возникновения компьютерных кластеров – проект Beowulf. Метакомпьютинг. Классификация Флинна, Шора и т.д. Организация межпроцессорных связей – коммуникационные топологии.</li> <li>- Примеры сетевых решений для создания кластерных систем</li> </ul>
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"> <li>-Библиотека MPI.</li> <li>- Модель SIMD. Инициализация и завершение MPI-приложения.</li> </ul> <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"> <li>-Введение в OpenMP</li> <li>-Стандарты программирования для систем с разделяемой памятью.</li> <li>-Создание многопоточных приложений. Использование многопоточности при программировании для многоядерных платформ.</li> <li>-Синхронизация данных между ветвями в параллельной программе. Директивы языка OpenMP</li> </ul>
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"> <li>-Существующие многоядерные системы.</li> <li>-GPU.</li> <li>-Использование OpenMP и MPI технологий совместно с CUDA.</li> <li>-Степень параллелизма численного алгоритма. Закон Амдала. Параллельный алгоритм решения СЛАУ</li> </ul>
12	<p>Тема 12. Программные платформы и системы для больших данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- системы управления потоками данных</li> <li>- системы хранения больших данных</li> <li>- платформы больших данных</li> <li>- обработка больших данных в реальном времени</li> </ul>

#### 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"> <li>- Hadoop</li> <li>- MapReduce</li> </ul>
2	<p>Тема 2. Анализ больших массивов данных инструментами python</p> <p>Рассматриваемые вопросы:</p>

№ п/п	Тематика практических занятий/краткое содержание
	<ul style="list-style-type: none"> <li>- библиотеки анализа и загрузки больших данных</li> <li>- фреймворки Python с параллельной обработкой данных</li> </ul>
3	<p>Тема 3. Высокопроизводительные вычисления с фреймворком Apache Spark</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- знакомство фреймворком</li> <li>- знакомство с архитектурой системы</li> <li>- реализации концепции MapReduce</li> <li>- разработки программ с использованием Apache Spark</li> </ul>
4	<p>Тема 4. Введение в высокопроизводительные серверы на python</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- работа с распределенными системами</li> <li>- создание кода для работы на GRU</li> </ul>
5	<p>Тема 5. Параллельное программирование. Параллельная обработка данных</p> <p>Рассматриваемые вопросы:</p> <ol style="list-style-type: none"> <li>1. Параллельное программирование</li> <li>2. Параллельная обработка данных</li> </ol>
6	<p>Тема 6. Работа с vCenter Server</p> <p>Рассматриваемые вопросы:</p> <ol style="list-style-type: none"> <li>1. Понятие vCenter</li> <li>2. Установка и развертывание vCenter</li> </ol>
7	<p>Тема 7. Развертывание среды для обработки данных при помощи Hadoop</p> <p>Рассматриваемые вопросы:</p> <ol style="list-style-type: none"> <li>1. Обработка данных в Hadoop</li> <li>2. Распределенная файловая система Hadoop</li> <li>3. Разработка приложений MapReduce</li> <li>4. Настройка кластера Hadoop</li> </ol>
8	<p>Тема 8. Составление спецификаций оборудования для работы с высокопроизводительными вычислениями</p> <p>Рассматриваемые вопросы:</p> <ol style="list-style-type: none"> <li>1. Составление спецификаций оборудования для работы с высокопроизводительными вычислениями</li> </ol>
9	<p>Тема 9. Работа с базами данных NoSQL</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- знакомство с базами данных NoSQL</li> <li>- сравнение SQL и NoSQL</li> <li>- работа с MongoDB и Cassandra</li> </ul>
10	<p>Тема 10. Машинное обучение на больших данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- основы машинного обучения</li> <li>- применение машинного обучения для анализа больших данных</li> <li>- использование библиотек Scikit-learn и TensorFlow</li> </ul>
11	<p>Тема 11. Работа с облачными платформами для обработки данных</p> <p>Рассматриваемые вопросы:</p> <ul style="list-style-type: none"> <li>- знакомство с облачными платформами (AWS, Google Cloud, Azure)</li> <li>- развертывание кластеров в облаке</li> <li>- обработка данных в облачных средах</li> </ul>

№ п/п	Тематика практических занятий/краткое содержание
12	Тема 12. Оптимизация производительности в распределенных Рассматриваемые вопросы: - методы оптимизации производительности - анализ узких мест в распределенных системах - инструменты мониторинга и анализа производительности
13	Тема 13. Работа с потоковыми данными Рассматриваемые вопросы: - обработка потоковых данных в реальном времени - использование Apache Kafka и Apache Flink - разработка приложений для обработки потоковых данных
14	Тема 14. Визуализация больших данных Рассматриваемые вопросы: - инструменты визуализации данных (Tableau, Power BI, Matplotlib) - создание интерактивных дашбордов - визуализация данных в реальном времени
15	Тема 15. Работа с графовыми базами данных Рассматриваемые вопросы: - знакомство с графовыми базами данных (Neo4j, ArangoDB) - анализ графовых данных - применение графовых баз данных в транспортных системах
16	Тема 16. Разработка и тестирование приложений для обработки Рассматриваемые вопросы: - разработка приложений для обработки больших данных - тестирование производительности и отказоустойчивости - оптимизация кода для работы с большими объемами данных

#### 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	Габдуллин, Н. М. Развитие человеческого капитала и цифровой экономики в регионах России: факторный и кластерный анализ :	<a href="https://e.lanbook.com/book/173018">https://e.lanbook.com/book/173018</a>

	монография / Н. М. Габдуллин. — Казань : КФУ, 2019. — 268 с. — ISBN 978-5-00130-291-9	
2	Гласснер, Э. Глубокое обучение без математики. Том 2. Практика : руководство / Э. Гласснер ; перевод с английского В. А. Яроцкого. — Москва : ДМК Пресс, 2020. — 610 с. — ISBN 978-5-97060-767-1	<a href="https://e.lanbook.com/book/131710">https://e.lanbook.com/book/131710</a>
3	Кук, Д. Машинное обучение с использованием библиотеки H2O / Д. Кук ; перевод с английского А. Б. Огурцова. — Москва : ДМК Пресс, 2018. — 250 с. — ISBN 978-5-97060-508-0	<a href="https://e.lanbook.com/book/97353">https://e.lanbook.com/book/97353</a>
4	Флах, П. Машинное обучение. Наука и искусство построения алгоритмов, которые извлекают знания из данных / П. Флах. — Москва : ДМК Пресс, 2015. — 400 с. — ISBN 978-5-97060-273-7	<a href="https://e.lanbook.com/book/69955">https://e.lanbook.com/book/69955</a>
5	Шалев-Шварц, Ш. Идеи машинного обучения : учебное пособие / Ш. Шалев-Шварц, Бен-Давид Ш. ; перевод с английского А. А. Слинкина. — Москва : ДМК Пресс, 2019. — 436 с. — ISBN 978-5-97060-673-5	<a href="https://e.lanbook.com/book/131686">https://e.lanbook.com/book/131686</a>
6	Гудфеллоу, Я. Глубокое обучение / Я. Гудфеллоу, И. Бенджио, А. Курвилль ; перевод с английского А. А. Слинкина. — 2-е изд. — Москва : ДМК Пресс, 2018. — 652 с. — ISBN 978-5-97060-618-6	<a href="https://e.lanbook.com/book/107901">https://e.lanbook.com/book/107901</a>

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](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:

старший преподаватель кафедры  
«Цифровые технологии управления  
транспортными процессами»

И.В. Зенковский

Согласовано:

Директор

Б.В. Игольников

Руководитель образовательной  
программы

О.Б. Проневич

Председатель учебно-методической  
комиссии

Д.В. Паринов