

III Научная
конференция ФКН

2025



ФКН

СБОРНИК АННОТАЦИЙ



Данный электронный сборник содержит темы и аннотации докладов основной секции и стендовых докладов третьей научной конференции факультета компьютерных наук НИУ ВШЭ, которая прошла с 26 по 29 октября 2025 года в учебном центре «Вороново», расположенном на территории Новой Москвы.

На конференции исследователи факультета представили свои основные научные результаты и успехи лабораторий. Цель научной конференции ФКН – собрать вместе коллег, которые заняты весьма разными темами. Задача конференции скорее не глубоко погрузиться в конкретную узкую предметную область, но рассказать, исследования по каким темам проводятся на факультете, показать коллегам значение и красоту каждого из направлений.



НКФКН2024

Темы конференции уже традиционно включают следующие направления, но не ограничиваются ими:

- машинное обучение и искусственный интеллект,
- анализ и обработка текстов на естественных языках,
- биоинформатика и медицинские информационные системы,
- программная инженерия и системное программирование,
- человеко-машинное взаимодействие и компьютерные игры,
- облачные и мобильные технологии,
- математическое моделирование,
- теоретическая информатика,
- анализ данных и процессов,
- управление сложными системами.

К участию в работе конференции были приглашены все научные сотрудники и преподаватели факультета, а также аспиранты и студенты, готовые представить результаты своей исследовательской работы. В 2025 году конференция вызвала большой интерес. Для основного трека конференции были отобраны 38 докладов. Также в программе конференции почти 50 стендовых докладов, круглый стол и демонстрация технологических достижений проектных групп.

Официальная страница конференции НКФКН 2025 в интернете:
https://cs.hse.ru/sci_conf2025

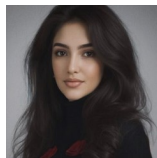
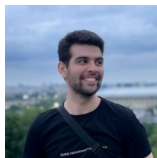
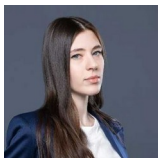
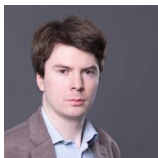
Темы и аннотации докладов, круглых столов, стендовых докладов приводятся в авторской редакции.

Оглавление

Организационный комитет 2025.....	7
Программа в целом.....	8
Приглашённая лекция «Модели долгосрочной социально-экономической динамики».....	12
Круглый стол «Интеллектуальные сервисы с анонимизацией и управлением данными на базе SMARTMLOPS: проблемы и подходы интеграции».....	13
Доклады основных секций.....	15
→ Первая секция.....	16
Генеративная модель для кластерного анализа.....	16
Прощай, объектно-ориентированное программирование?.....	17
Равновесия Нэша в чистых стационарных стратегиях для игр двух игроков с положительными аддитивными штрафами.....	18
Вычислительный юмор: подходы, данные, оценка.....	20
Применение машинно-обучаемых межатомных потенциалов в атомистическом моделировании.....	21
3MDBench: Medical Multimodal Multi-agent Dialogue Benchmark.....	22
Анализ структурной сложности сетей потоков работ для моделирования асинхронного взаимодействия агентов.....	23
Enhancing Claim Fact-Checking Against Wikipedia: A Diagnostic Taxonomy and a Generative Framework.....	24
Робастное глубинное обучение.....	25
→ Вторая секция.....	27
Фазовые переходы для квантовых блужданий.....	27
Методы оценки ошибок генерации диффузионных моделей.....	28
Физически-информированные нейронные сети в инженерных задачах.....	29
Интеграция генеративного ИИ в преподавание баз данных: от инструмента для списывания к ассистенту для глубокого обучения.....	31
Кластеризация графов на основании средних расстояний методом WPGMA.....	33
→ Третья секция.....	34
Решения с открытым кодом для локальных голосовых помощников.....	34
Logistic regression in high dimension.....	36
Выявление коннектома у глухих посредством решения обратной задачи.....	37
Об одной задаче, связанной с моментом первого достижения заданного уровня случайным процессом.....	39
Нейронные стохастические дифференциальные уравнения для генеративного предсказания временных рядов.....	40
Машинное обучение для открытия и создания новых генетических систем.....	41
Tight bounds for Schrodinger potential estimation in unpaired image-to-image translation problems.....	42

Высокоиммерсивная интерактивная коллаборативная среда - новый аспект зловещей долины.....	43
Оптимальный транспорт и диффузионные модели.....	45
→ ЧЕТВЁРТАЯ СЕКЦИЯ.....	47
<null>.....	47
From computability to quasicrystals.....	48
Extreme Events in Action: What Triggers Them in Complex Systems?.....	49
О матричной задаче Прокруста и ее приложениях в глубинном обучении.....	50
Агентно-ориентированная модель обработки логов, основанная на синтезе методов промпт-инжиниринга и цепочки рассуждений.....	51
Mathematical theory of optimal processes in a common path interferometer using a diffraction phase microscope.....	53
→ ПЯТАЯ СЕКЦИЯ.....	54
Алгоритм клиент-серверного взаимодействия по незашифрованному соединению.....	54
О некоторых методологических аспектах обучения генеративных потоковых сетей.....	55
Кооперативная игра «Ханаби» – новый вызов ИИ?.....	56
Приватность данных в эпоху ML: обезличивание, анонимизация и синтетические данные на практике.....	58
Improved Stochastic Optimization of LogSumExp.....	60
Длинные строки составных значений полиномов и базис порядка 2.....	62
Бифуркационные модели нелинейных уравнений в частных производных.....	63
GNN-based neutron reconstruction in the HGND at the BM@N experiment.....	65
Мозг и искусственный интеллект: интерпретируемые модели и новые горизонты анализа нейроданных.....	66
ПОСТЕРЫ (СТЕНДОВЫЕ ДОКЛАДЫ).....	68
→ ПЕРВАЯ ПОСТЕРНАЯ СЕССИЯ.....	68
→ ВТОРАЯ ПОСТЕРНАЯ СЕССИЯ.....	84
→ ЛУЧШИЕ ПОСТЕРЫ.....	97
ДЕМОНСТРАЦИЯ ТЕХНОЛОГИЧЕСКИХ РЕШЕНИЙ.....	98
ПОЛЕЗНАЯ ИНФОРМАЦИЯ.....	99
МЕСТО ПРОВЕДЕНИЯ.....	99
СХЕМА ПРОЕЗДА.....	100
Общественным транспортом.....	100
На личном автомобиле.....	100
РАЗМЕЩЕНИЕ.....	101

Организационный комитет 2025



- Проф. Иван Владимирович Аржанцев², д.ф.-м.н., декан ФКН
- Доц. Алексей Александрович Мицюк², к.комп.н., зам. декана ФКН по науке
- Начальник отдела Ксения Антоновна Кузнецова², отдел сопровождения проектов
- Преподаватель Мишан Хаммад оглы Алиев², лаборатория ТОМИИ ФКН
- Администратор Сабина Истамовна Абдуллаева², отдел сопровождения проектов

Организаторы приложили немало сил для того, чтобы НКФКН 2025 состоялась, и надеются, что участники конференции получат много новой информации и приятные эмоции, выработают новые идеи и лучше узнают друг друга. Желаем хорошей конференции!

Техническая поддержка и фото – Игорь Леонидович Тимошук².

Благодарности

Конференция проводится факультетом компьютерных наук НИУ ВШЭ. Благодарность выражается компании Яндекс за участие в подготовке призов авторам лучших стендовых докладов постерных сессий. Также благодарность выражается Андрею Олеговичу Бондареву² и Ивану Станиславовичу Копылову² и лаборатории робототехники ФКН за помощь в печати уникальных памятных призов. Лучшие доклады выбраны по итогам голосования участников конференции.

Программа в целом

День 1 · 26 октября · Воскресенье

15:30	Трансфер от кампуса на Покровском бульваре до УЦ Вороново (корпус НИУ ВШЭ «Покровский бул., 11», вход №3)
18:00	Ужин в столовой
19:00	Приглашённая лекция « Модели долгосрочной социально-экономической динамики » Доц. Дмитрий Александрович Веселов, к.э.н., зам. декана по научной работе факультета экономических наук НИУ ВШЭ
21:00	Свободное общение · Настольные игры (аудитории 6 этажа)

День 2 · 27 октября · Понедельник

9:00	Завтрак в столовой
9:50	Открытие конференции в большом зале
10:00	Генеративная модель для кластерного анализа ‣ Проф. Борис Григорьевич Миркин, д.т.н.
10:20	Прощай, объектно-ориентированное программирование? ‣ Проф. Александр Иванович Легалов, д.т.н.
10:40	Равновесия Нэша в чистых стационарных стратегиях для игр двух игроков с положительными аддитивными штрафами ‣ Проф. Михаил Николаевич Вялый, к.ф.-м.н.
11:00	Вычислительный юмор: подходы, данные, оценка ‣ Доц. Павел Исаакович Браславский, к.т.н.
11:20	Применение машинно-обучаемых межатомных потенциалов в атомистическом моделировании ‣ Доц. Иван Сергеевич Новиков, к.ф.-м.н.
11:40	3MDBench: Medical Multimodal Multi-agent Dialogue Benchmark ‣ Проф. Андрей Владимирович Савченко, д.т.н.
12:00	Анализ структурной сложности сетей потоков работ для моделирования асинхронного взаимодействия агентов ‣ Доц. Роман Александрович Нестеров, к.комп.н.
12:20	Enhancing Claim Fact-Checking Against Wikipedia: A Diagnostic Taxonomy and a Generative Framework ‣ Доц. Дмитрий Алексеевич Ильвовский, к.т.н.

12:40	Робастное глубинное обучение ► Доц. Алексей Сергеевич Болдырев, к.ф.-м.н.	
13:00	Обед в столовой	
14:00	Круглый стол « Интеллектуальные сервисы с анонимизацией и управлением данными на базе SmartMLOps: проблемы и подходы интеграции » (модератор: Доц. Сергей Аркадьевич Лебедев, к.э.н.)	Первая постерная сессия (аудитории 6 этажа)
16:30	Перерыв на кофе в столовой	
17:00	Фазовые переходы для квантовых блужданий ► Проф. Алексей Владимирович Устинов, д.ф.-м.н.	
17:20	Методы оценки ошибок генерации диффузионных моделей ► Стажёр Фёдор Константинович Пахуров	
17:40	Физически-информированные нейронные сети в инженерных задачах ► Доц. Александр Александрович Тараканов, Ph.D.	
18:00	Интеграция генеративного ИИ в преподавание баз данных: от инструмента для списывания к ассистенту для глубокого обучения ► Доц. Александр Давидович Брейман, к.т.н.	
18:20	Кластеризация графов на основании средних расстояний методом WPGMA на практике ► Алексей Андреевич Смоленцев, CS RnD лаборатория Т-Банка Т-банка	
18:40	Ужин в столовой	
19:30	Сессия вопросов руководителям ФКН в большом зале (модератор: Преп. Мишан Хаммад оглы Алиев)	
21:00	Свободное общение • Настольные игры (аудитории 6 этажа)	

День 3 · 28 октября · Вторник

9:00	Завтрак в столовой	
10:00	Решения с открытым кодом для локальных голосовых помощников ► Проф. Игорь Рубенович Агамирзян, к.ф.-м.н.	3
10:20	Logistic regression in high dimension ► Проф. Владимир Григорьевич Спокойный, д.ф.-м.н.	
10:40	Выявление коннектома у глухих посредством решения обратной задачи ► Проф. Александр Александрович Харламов, д.т.н.	

11:00	Об одной задаче, связанной с моментом первого достижения заданного уровня случайным процессом ▶ Проф. Сергей Львович Семаков, д.ф.-м.н.	
11:20	Нейронные стохастические дифференциальные уравнения для генеративного предсказания временных рядов ▶ Доц. Максим Львович Каледин, к.комп.н.	
11:40	Машинное обучение для открытия и создания новых генетических систем ▶ С.н.с. Сергей Анатольевич Шмаков, к.б.н.	
12:00	Tight bounds for Schrodinger potential estimation in unpaired image-to-image translation problems ▶ Доц. Пучкин Никита Андреевич, к.комп.н.	
12:20	Высокоиммерсивная интерактивная коллаборативная среда - новый аспект зловещей долины ▶ Доц. Ольга Вениаминовна Максименкова, к.т.н.	
12:40	Оптимальный транспорт и диффузионные модели ▶ Ст.преп. Денис Романович Ракитин	
13:00	Обед в столовой	
13:45	Общее фото участников конференции (атриум учебного центра)	
14:00	Вторая постерная сессия (аудитории 6 этажа)	Площадка технологий (демо лабораторий)
16:00	Перерыв на кофе в столовой	
16:30	<null> ▶ Проф. Аттила Кертес-Фаркаш, д.комп.н. (Prof. Attila Kertesz-Farkas, Dr.Sc.)	
16:50	From computability to quasicrystals ▶ Доц. Тома Ферник, Ph.D. (Assoc.Prof. Thomas Fernique, Ph.D.)	
17:10	Extreme Events in Action: What Triggers Them in Complex Systems? ▶ Доц. Сабаратинам Сринивасан, Ph.D. (Assoc.Prof. Sabarathinam Srinivasan, Ph.D.)	
17:30	О матричной задаче Прокруста и ее приложениях в глубинном обучении ▶ Доц. Максим Владимирович Рахуба, к.ф.-м.н.	
17:50	Агентно-ориентированная модель обработки логов, основанная на синтезе методов промпт-инжиниринга и цепочки рассуждений ▶ Пригл.преп. Кирилл Игоревич Пашигорев, Сбер-Технологии	
18:10	Mathematical theory of optimal processes in a common path interferometer using a diffraction phase microscope ▶ Вед.аналитик Наталья Анатольевна Талайкова	

18:30

Награждение лучших постеров

18:40

Ужин в столовой

19:30

Авторский квиз от **Павла Николаевича Азарова²**, начальника отдела по работе с абитуриентами, студентами и выпускниками, энтузиаста интеллектуальных игр и соревнований.

Свободное общение
(аудитории 6 этажа)



День 4 · 29 октября · Среда

9:00

Завтрак в столовой

10:00

Алгоритм клиент-серверного взаимодействия по незашифрованному соединению

► Проф. Дмитрий Владимирович Александров, д.т.н.

10:20

О некоторых методологических аспектах обучения генеративных потоковых сетей

► Доц. Сергей Владимирович Самсонов, к.мат.н.

10:40

Кооперативная игра «Ханаби» - новый вызов ИИ?

► Доц. Анастасия Александровна Оноприенко, к.ф.-м.н.

11:00

Приватность данных в эпоху ML: обезличивание, анонимизация и синтетические данные на практике

► Ст.преп. Юрий Владимирович Силаев

11:20

Improved Stochastic Optimization of LogSumExp

► Доц. Егор Леонидович Гладин, Ph.D.

11:40

Длинные строки составных значений полиномов и базис порядка 2

► Доц. Артем Олегович Радомский, к.ф.-м.н.

12:00

Бифуркационные модели нелинейных уравнений в частных производных

► Проф. Василий Александрович Громов, д.ф.-м.н.

12:20

GNN-based neutron reconstruction in the HGND at the BM@N experiment

► М.н.с. Владимир Олегович Бочарников

12:40

Мозг и искусственный интеллект: интерпретируемые модели и новые горизонты анализа нейроданных

► Н.с. Илья Владимирович Семенков

13:00

Обед в столовой

14:30

Трансфер из УЦ Вороново до кампуса на Покровском бульваре (атриум УЦ)

5

Приглашённая лекция «Модели долгосрочной социально-экономической динамики»



► Доц. **Дмитрий Александрович Веселов**^{[a](#)}, к.э.н., зам. декана по научной работе факультета экономических наук НИУ ВШЭ

↪ За последние 200 лет человечество достигло невиданного ранее прогресса в уровне доходах на душу населения и других показателях экономического развития. На лекции на основе новейших исследований в области экономики развития мы обсудим следующие вопросы: в чем причины появления феномена роста уровня жизни? Можно ли считать данный тренд устойчивым, и какие условия должны быть выполнены для его поддержания? Каковы закономерности успехов и провалов развития в доиндустриальную и современную эпоху?

Круглый стол «Интеллектуальные сервисы с анонимизацией и управлением данными на базе SmartMLOps: проблемы и подходы интеграции»

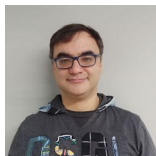
↳ На круглом столе рассматриваются проблемы сбора, обработки, обезличивания и управления данными и подходы к интеграции интеллектуальных сервисов на базе SmartMLOps – системы размещения и управления сервисами искусственного интеллекта, на примере прикладной задачи оценки благополучия студентов НИУ ВШЭ.

Модератор:



► Доц. **Сергей Аркадьевич Лебедев**^{[↗](#)}, к.э.н.,
руководитель департамента программной инженерии

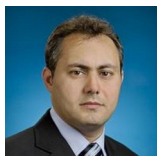
Участники:



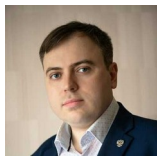
► Доц. **Игорь Борисович Архимандритов**^{[↗](#)},
департамент информатики ШИФТ СПб НИУ ВШЭ



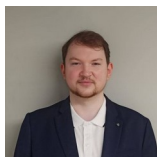
► Ст.преп. **Дмитрий Валерьевич Пантюхин**^{[1](#)},
департамент программной инженерии



► Доц. **Хади Мухаммед Салех**^{[1](#)}, к.т.н., департамент
программной инженерии



► Ст.преп. **Юрий Владимирович Силаев**^{[1](#)}, департамент
программной инженерии



► Ст.преп. **Александр Евгеньевич Трифонов**^{[1](#)},
департамент информатики ШИФТ СПб НИУ ВШЭ



► **Данил Игоревич Швецов**^{[1](#)}, программист Института
искусственного интеллекта и цифровых наук

Доклады основных секций

Далее приводятся аннотации всех докладов пяти основных секций конференции.

→ Первая секция

Генеративная модель для кластерного анализа



► Проф. **Борис Григорьевич Миркин**², д.т.н., департамент анализа данных и искусственного интеллекта, ординарный профессор

↳ Рассматривается простая модель порождения кластеризованных данных и критерий наименьших квадратов для ее идентификации. Упоминается, что популярный метод k -средних для кластерного анализа может рассматриваться как метод подгонки для этой модели. Пифагорово разложение разброса данных на основе этой модели используется как основа для методов: (1) формирования аномальных кластеров в различных приложениях; (2) формирования консенсусного разбиения; (3) формирования древовидных решающих правил для анализа категориальных данных; (4) оценки уровня сложности таблицы данных. В задаче (3) обнаруживается связь между видом решающей функции (типа Джини или Хи-квадрат) и способом учета значимости категории в соответствии с ее частотой. Методы решения задачи (4) в известной нам литературе не встречались.

Прощай, объектно-ориентированное программирование?



▸ Проф. **Александр Иванович Легалов**², д.т.н., департамент программной инженерии

↪ Объектно-ориентированное (ОО) программирование и проектирование находится на пике популярности. Создаются новые ОО языки программирования, предлагаются различные технические и методологические приемы и правила хорошего ОО кодирования. Одной из основ его популярности является поддержка динамического полиморфизма. Вместе с тем, появляются альтернативные инструменты, обеспечивающие поддержку динамического полиморфизма, в языках процедурного и функционального программирования. Примерами могут служить Go и Rust. Однако приняты в них технические решения ничем не превосходят ОО подход. Ситуацию может изменить процедурно-параметрический полиморфизм, обеспечивающий более гибкую разработку эволюционно расширяемых программ. Насколько это реально?

Равновесия Нэша в чистых стационарных стратегиях для игр двух игроков с положительными аддитивными штрафами



► Проф. Михаил Николаевич Вялый², к.ф.-м.н., департамент больших данных и информационного поиска, международная лаборатория теоретической информатики

↳ В игре с аддитивными штрафами каждый ход стоит каждому из игроков некоторого платежа. Общий платеж игрока равен сумме платежей по всем ходам партии. Платежные функции игроков могут различаться. Множество позиций обычно предполагается конечным. Интересует существование равновесия Нэша в чистых стационарных стратегиях. Было известно, что такого равновесия нет для игр, в которых три игрока и все платежи положительные, а также для игр двух игроков, в которых платежи произвольные. В докладе будет рассказано о новом результате в этой области: для двух игроков и положительных платежей равновесие Нэша существует. Существование равновесия Нэша следует из существования дважды кратчайшего пути: такого пути, который является кратчайшим как в графе, ограниченном на стратегию первого игрока с платежной функцией второго, так и в аналогичном графе, где игроки меняются местами. Существование дважды

кратчайшего пути доказано для случая, когда хотя бы один игрок не в состоянии отрезать терминалы от начальной позиции. Если оба игрока могут отрезать терминалы, равновесие Нэша существует по тривиальной причине и равновесная цена бесконечна. Вопрос о существовании бикратчайшего пути в этом случае остаётся открытым.

Доказательство основано на существовании дерева кратчайших-длиннейших путей, которое доказано Гурвичем, Жао и Хачияном (2006) с помощью обобщения алгоритма Дейкстры на задачу о кратчайшем пути в графе с запретами, а также следует из существования канонической формы циклической игры, доказанного Гурвичем, Карзановым, Хачияном (1988).

Вычислительный юмор: подходы, данные, оценка



► Доц. Павел Исаакович Браславский²¹, к.т.н., департамент больших данных и информационного поиска, лаборатория моделей и методов вычислительной прагматики

↳ В докладе я сделаю обзор основных методов распознавания и генерации юмора в тексте, существующих датасетов и подходов к оценке. Рассмотрим эволюцию области и ее ближайшие перспективы.

Применение машинно-обучаемых межатомных потенциалов в атомистическом моделировании



► Доц. **Иван Сергеевич Новиков**^{[2](#)}, к.ф.-м.н., департамент больших данных и информационного поиска

↳ В докладе будут представлены некоторые результаты разработки и применения машинно-обучаемых межатомных потенциалов (далее - МОПов). В начале доклада будет дано краткое введение в МОПы и алгоритмы их обучения. Далее речь пойдёт о магнитных МОПах и результатах их применения для исследования магнитных материалов, о МОПах, включающих электростатическое взаимодействие и их приложениях, а также о методах уменьшения числа параметров в МОПах.

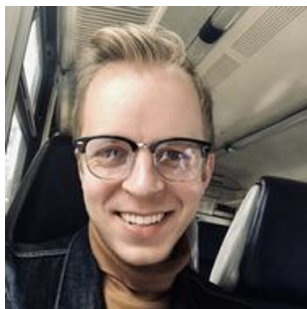
3MDBench: Medical Multimodal Multi-agent Dialogue Benchmark



► Проф. **Андрей Владимирович Савченко**²¹, д.т.н., лаборатория теоретических основ моделей искусственного интеллекта, НИУ ВШЭ в Нижнем Новгороде

↳ Though Large Vision-Language Models (LVLMs) are being actively explored in medicine, their ability to conduct telemedicine consultations combining accurate diagnosis with professional dialogue remains underexplored. In this talk, I present 3MDBench (Medical Multimodal Multi-agent Dialogue Benchmark), an open-source framework for simulating and evaluating LVLM-driven telemedical consultations. 3MDBench simulates patient variability through four temperament-based Patient Agents and an Assessor Agent that jointly evaluate diagnostic accuracy and dialogue quality. It includes 3013 cases across 34 diagnoses drawn from real-world telemedicine interactions, combining textual and image-based data. The experimental study compares diagnostic strategies for popular LVLMs, including GPT-4o-mini, LLaVA-3.2-11B-Vision-Instruct, and Qwen2-VL-7B-Instruct. I will show that injecting predictions from a diagnostic convolutional network into the LVLM's context boosts F1 by up to 20%.

Анализ структурной сложности сетей потоков работ для моделирования асинхронного взаимодействия агентов



► Доц. **Роман Александрович Нестеров**^{[2](#)}, к.комп.н., департамент программной инженерии, лаборатория процессно-ориентированных информационных систем

↪ Структура модели процесса, полученной на основе журнала событий многоагентной системы, часто не отражает архитектуру системы с точки зрения взаимодействий между агентами. Существующие метрики проверки соответствия в основном оценивают степень соответствия поведения обнаруженной модели последовательностям событий, зафиксированным в журнале. Однако такие поведенческие метрики могут оказаться недостаточными для того, чтобы отличить модели процессов, извлеченные из журнала одной и той же многоагентной системы, с учетом степени независимости агентов и сложности их взаимодействий. В данной работе предлагается теоретически обоснованный подход к измерению структурной сложности модели процесса, представляющей многоагентную систему с асинхронным взаимодействием агентов. Также представлены ключевые результаты серии экспериментов, направленных на оценку чувствительности предложенного подхода к структурным изменениям в моделях процессов.

Enhancing Claim Fact-Checking Against Wikipedia: A Diagnostic Taxonomy and a Generative Framework



► Доц. **Дмитрий Алексеевич Ильвовский**², к.т.н., департамент анализа данных и искусственного интеллекта, лаборатория интеллектуальных систем и структурного анализа

↳ Fact-checking is a crucial yet challenging task that continues to gain importance. In an effort to address this issue, the FEVER large-scale dataset was developed to facilitate evidence-based fact-checking using Wikipedia as a reference. Despite numerous proposed approaches and evaluations on this dataset, a comprehensive understanding of the errors made by these approaches is still lacking. Here, we aim to bridge this gap. We introduce a diagnostic taxonomy and a generative framework to enhance FEVER-style fact-checking. We establish a taxonomy of errors and we construct a diagnostic dataset that enables the analysis of the errors made by state-of-the-art models as well as their distribution within the FEVER dataset. Additionally, we provide a set of prompts to generate examples within this taxonomy.

Робастное глубинное обучение



▸ Доц. **Алексей Сергеевич Болдырев**², к.ф.-м.н., департамент больших данных и информационного поиска, лаборатория методов анализа больших данных

↳ Быстрое развитие приложений машинного обучения (ML) и искусственного интеллекта (AI) требует обучения большого количества моделей. Эти растущие требования подчёркивают важность обучения моделей без участия человека при сохранении робастности их предсказаний. В ответ на эту потребность мы предлагаем новый подход к определению робастности модели машинного обучения. Этот подход, дополненный алгоритмом выбора модели, реализованным как мета-алгоритм, является универсальным и применимым к любой модели машинного обучения, при условии её соответствия задаче. В докладе демонстрируется применение нашего подхода для оценки робастности моделей глубокого обучения. Мы рассматриваем небольшие модели, состоящие из нескольких сверточных и полносвязных слоёв, используя распространённые оптимизаторы из-за их простоты интерпретации и вычислительной эффективности. Также исследуется влияние объёма обучающей выборки, инициализации весов модели и априорных предположений на устойчивость моделей глубокого обучения.

Перспективные направления будущих исследований включают применение предложенного подхода для оценки устойчивости глубоких нейронных сетей и современных моделей на различных наборах данных, а также исследование дополнительных критериев выбора, позволяющих гарантировать требуемый уровень робастности моделей.

→ Вторая секция

Фазовые переходы для квантовых блужданий



► Проф. **Алексей Владимирович Устинов**², д.ф.-м.н., департамент больших данных и информационного поиска

↳ Квантовые блуждания, введенные Ричардом Фейнманом, являются архетипической моделью для изучения квантовых алгоритмов. Доклад будет посвящён фазовым переходам, которые демонстрирует одномерное квантовое блуждание.

Методы оценки ошибок генерации диффузионных моделей



► Стажёр **Фёдор Константинович Пахуров**¹, лаборатория
«Искусственный интеллект в математических финансах»

↳ В докладе будут рассмотрены методы оценки галлюцинаций в генерации диффузионных моделей. Основное внимание будет уделено детерминистичному методу, суть которого заключается в отображения изображений в одномерное пространство и кластеризации в плоскости характеристик Complexity–Entropy для полученных из изображений временных рядов, что позволяет оценить степень уверенности в принадлежности генерации к классу галлюцинаций.

Физически-информированные нейронные сети в инженерных задачах



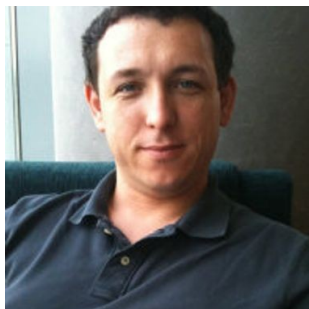
► Доц. **Александр Александрович Тараканов**², Ph.D., департамент больших данных и информационного поиска

↳ Суррогатные модели сложных физических систем широко применяются в инженерных задачах. Например, задачи количественной оценки неопределённости при добыче углеводородов или хранении CO_2 в подземных резервуарах обладают чрезвычайно высокой вычислительной сложностью, которую можно существенно снизить за счёт использования дешёвых аппроксимирующих моделей (прокси-моделей). Для построения таких моделей было применено множество методов машинного обучения — от полиномиальных функций и решающих деревьев до искусственных нейронных сетей.

В данной работе мы предлагаем новый подход к обучению физически обоснованных нейронных сетей (Physics Informed Neural Networks) для аппроксимации течения жидкости в пористых средах. Экспериментально показано, что точную нейронную сеть можно обучить на относительно небольшом наборе данных. Более того, мы демонстрируем, что такие физически обоснованные

нейронные сети оказываются более устойчивыми по сравнению с методами, основанными на полиномиальном хаосе и решающих деревьях, при увеличении сложности моделируемых систем, при условии обучения на одном и том же наборе данных. Другими словами, физически обоснованные нейронные сети способны аппроксимировать системы с высокой степенью нелинейности даже в условиях ограниченного числа обучающих примеров.

Интеграция генеративного ИИ в преподавание баз данных: от инструмента для списывания к ассистенту для глубокого обучения



► Доц. Александр Давидович Брейман², к.т.н.

↪ В докладе рассматривается одна из наиболее острых проблем современного технического образования: использование студентами генеративных моделей ИИ для автоматического решения задач. Стандартный подход, основанный на запретах и попытках контроля, демонстрирует свою неэффективность, поскольку не устраняет первопричину - отсутствие у студентов мотивации к глубокому освоению материала. Вместо борьбы с технологиями предлагается изменить саму педагогическую парадигму, целенаправленно интегрируя ИИ в учебный процесс в качестве интеллектуального ассистента.

В работе представлен конкретный методический инструментарий, позволяющий сместить фокус с получения готового ответа на анализ процесса его создания. Рассматриваются практические форматы заданий: «студент против ИИ», где требуется сравнить и критически оценить собственное и машинное решение; «отладка ИИ», в рамках которого студент должен найти и объяснить ошибки в

сгенерированном коде; а также использование особых задач, которые сложно решить стандартным запросом к нейросети.

Особое внимание уделяется изменению форматов контроля знаний: предлагается внедрение устных защит лабораторных работ в формате код-ревью и экспресс-задания «на бумаге» для проверки фундаментального понимания. Такой подход позволяет не только объективнее оценивать реальные знания, но и развивает у студентов ключевые навыки: критическое мышление, способность к анализу и умение верифицировать информацию. Цель доклада - показать, как трансформировать роль преподавателя от транслятора знаний к наставнику, который учит студентов быть не пассивными потребителями контента, сгенерированного ИИ, а инженерами, способными использовать эти мощные инструменты для решения сложных профессиональных задач.

Кластеризация графов на основании средних расстояний методом WPGMA

(на практике)

▸ **Алексей Андреевич Смоленцев**, CS RnD лаборатория Т-банка

↳ WPGMA (Weighted Pair Group Method with Arithmetic Mean) — иерархический алгоритм, который объединяет вершины графа, опираясь на средние попарные расстояния. В докладе рассматриваются:

- принцип работы WPGMA
- оптимизация вычислений с $O(n^2 \log n)$ до $O(n^2)$
- адаптация алгоритма под прикладную задачу

→ Третья секция

Решения с открытым кодом для локальных голосовых помощников



► Проф. Игорь Рубенович Агамирзян², к.ф.-м.н., департамент программной инженерии

↳ Получившие в последние годы массовое распространение голосовые помощники (Siri от Apple, Alexa от Amazon, Google Assistant, Алиса от Яндекса и т. д.) все отличаются общей особенностью – они работают через облачные сервисы и реализуются в Центрах обработки данных соответствующих компаний. Соответственно, при возникновении проблем с доступом в интернет такие помощники оказываются неработоспособными. В то же время развитие аппаратной базы и технологий машинного обучения сделали возможной реализацию голосовых помощников, работающих на краю сети и независимых от доступа в интернет. Сегодня существует ряд решений для STT и TTS, реализованных в модели ПО с открытым кодом и способных работать на относительно небольших и дешёвых компьютерах. К ним относятся Whisper от OpenAI (STT с предобученными сетями для различных

языков), Piper от OpenHome (TTS для различных языков и голосов) и относительно недавно появившийся интеграционный протокол Wyoming (разработка OpenHome Foundation), позволяющий проинтегрировать решения разных производителей в полноценного локального голосового помощника. В докладе проводится обзор таких решений с открытым кодом и рассматривается процесс интеграции и решения возникающих при этом проблем.

Logistic regression in high dimension



► Проф. **Владимир Григорьевич Спокойный**^{[2](#)}, д.ф.-м.н., лаборатория теоретических основ моделей искусственного интеллекта

↳ Logistic regression is widely used e.g. in binary classification in machine learning for binary classification or in binary response models in econometrics. Now issues and phenomena arise when the parameter dimension exceeds the number of observations. This study explains the performance of the classical penalized MLE with a proper ridge penalization for logistic regression with random design in terms of the so called effective dimension p_e . The main results describe a finite sample expansion of pMLE and its quadratic risk under the condition $p_e^2 \ll n$ which is essential and perhaps unavoidable. The proof involve recent results on random matrices and tensors of fourth order.

Выявление коннектома у глухих посредством решения обратной задачи



► Проф. **Александр Александрович Харламов**², д.т.н., департамент программной инженерии

↳ В названии тезисов доклада на сессии РАО «Выявление коннектома у глухих посредством решения обратной задачи» как и в названии выступления на конференции ДПИ ВШЭ в прошлом году прозвучало «у глухих», хотя представленный материал касался анализа записей МЭГ у нормально слышащих испытуемых. И результаты вполне понятные. Тем не менее, решение задачи выявления архитектуры коннектома в разных условиях проведения эксперимента дает возможность понять реальные информационные процессы, которые протекают в мозге человека. Правда эта возможность касается достаточно обобщенного уровня анализа. Полученный коннектом показывает путь прохождения информации в коре, но не показывает подробности этого прохождения. В связи с этим возникают несколько задач, касающихся продолжения этих исследований. (1) Исследование коннектома собственно глухих людей (и, следовательно, попытка выявления особенностей механизма возникновения глухоты). (2) Исследование коннектома нормально слышащих в другой постановке исследовательской

задачи (например, мысленного управления). (3) Исследование коннектома нормально слышащих при использовании другого оборудования (например оборудование EmotivBCI фирмы Emotiv, обеспечивающее аппаратно интерфейс «мозг-компьютер»). (4) Исследование коннектома в другом временном масштабе (выявление элементов произнесенного слова). В рамках решения этих задач выявление коннектома по методике Проля Даса остается неизменным, хотя дополняется (в зависимости от решаемой задачи) своими особенностями. Выявление коннектома является в настоящий момент наиболее актуальной задачей в электроэнцефалографии и магнитоэнцефалографии процессов мозга человека. Одним из путей решения этой задачи является использование свертки несущего (в данном случае – речевого) сигнала с записями ЭЭГ (МЭГ).

Об одной задаче, связанной с моментом первого достижения заданного уровня случайным процессом



► Проф. Сергей Львович Семаков^{[2](#)}, д.ф.-м.н., департамент больших данных и информационного поиска

↳ Рассматривается задача оценки вероятности события, состоящего в том, что первое достижение заданного уровня непрерывным случайным процессом произойдёт в какой-либо момент из заданного промежутка изменения независимой переменной. Ранее полученные результаты общего характера конкретизируются для гауссовского гладкого процесса. Приводятся результаты численных расчетов оценок при различных параметрах процесса.

Нейронные стохастические дифференциальные уравнения для генеративного предсказания временных рядов



► Доц. **Максим Львович Каледин**²¹, к.комп.н., департамент больших данных и информационного поиска, лаборатория методов анализа больших данных

↳ Многие временные ряды по своей сути являются непрерывными во времени процессами, которые наблюдаются в дискретные моменты времени. Примерами могут служить физические параметры атмосферы, необходимые для прогноза погоды, геолокация, а также финансовые ряды, в которых очень высокая частотность наблюдений. Использование моделей временных рядов для подобных приложений приводит с одной стороны к неэффективному использованию априорных предположений о физике модели, которая часто формулируется в непрерывном времени, а с другой – к полному игнорированию или неудобному использованию иррегулярности наблюдений. В нашей работе мы представляем C-SDE-GAN, генеративную модель для предсказания временных рядов на основе нейронных стохастических дифференциальных уравнений, которая способна решить обе проблемы с использованием небольшого количества параметров относительно современных нейросетевых подходов на основе трансформеров.

Машинное обучение для открытия и создания новых генетических систем



► С.н.с. **Сергей Анатольевич Шмаков**², к.б.н., лаборатория статистической и вычислительной геномики

↳ Объемы геномных данных растут экспоненциально, однако наше понимание закодированной в них информации остаётся критически ограниченным. Так, например, в области микробиома кишечника до 40% генов всё ещё имеют неизвестную функцию. Исторически аннотация геномов опиралась на ресурсоёмкие, полуавтоматические методы, которые не справляются с текущим потоком данных.

В докладе будут представлены разработанные вычислительные подходы для аннотации геномов и наши ключевые результаты, включая открытие новых генов CRISPR-Cas в бактериях и их успешную трансформацию в новые инструменты для редактирования геномов. Далее будут обсуждены перспективные стратегии использования машинного обучения для аннотации геномов и поиска генов с заданными свойствами в интересах биотехнологии или создания новых синтетических генов.

Tight bounds for Schrodinger potential estimation in unpaired image-to-image translation problems



► Доц. Пучкин Никита Андреевич², к.комп.н., лаборатория теоретических основ моделей искусственного интеллекта, БК Института проблем передачи информации им. А.А. Харкевича РАН

↳ Modern methods of generative modelling and unpaired image-to-image translation based on Schrodinger bridges and stochastic optimal control theory aim to transform an initial density to a target one in an optimal way. In our work, we assume that we only have access to i.i.d. samples from initial and final distributions. This makes our setup suitable for both generative modelling and unpaired image-to-image translation. Relying on the stochastic optimal control approach, we choose an Ornstein-Uhlenbeck process as the reference one and estimate the corresponding Schrodinger potential. Introducing a risk function as the Kullback-Leibler divergence between couplings, we derive tight bounds on generalization ability of an empirical risk minimizer in a class of Schrodinger.

Высокоиммерсивная интерактивная коллаборативная среда - новый аспект зловещей долины



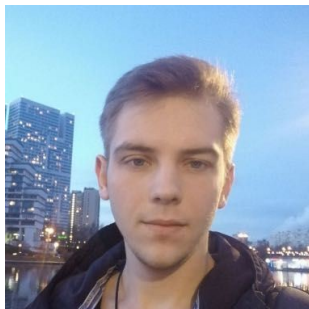
► Доц. **Ольга Вениаминовна Максименкова**^{[2](#)}, к.т.н., департамент программной инженерии

↳ В этом (2025) году опубликованы результаты исследования, впервые показывающие успешное прохождение теста Тьюринга несколькими большими языковыми моделями (БЯМ) вне специально проводимых состязаний (вспомним Turing Award). Не умаляя данного результата, скажем, что тест Тьюринга достаточно старый (в октябре этого года отмечает 75-летний юбилей) и ограниченный (классическое взаимодействие между человеком и машиной посредством текстового терминала). На текущем этапе развития компьютерной техники мы видим куда более развитые варианты подобного взаимодействия в интерактивных иммерсивных средах, например, графические двухмерные и трёхмерные миры компьютерных игр, промышленных симуляторов и т.п.

Опыт взаимодействия человека и машины в высокоиммерсивных средах весьма ограничен (даже для взаимодействия человек-

человек) и для всестороннего изучения и проведения полноценных экспериментов нуждается в накоплении реальных данных и знаний о взаимодействии между человеком и машиной не через терминал, а посредством мультимодальных устройств ввода/вывода всех уровней иммерсивности: от офисного стандарта рабочей станции до шлемов, костюмов и всенаправленных дорожек виртуальной реальности, а также посредством взаимодействия с роботом-андридом в реальном мире. Уже имеющийся фрагментарный опыт заставил по новому взглянуть на «зловещую долину» и высветил неожиданные аспекты человеко-машинного взаимодействия. В докладе речь пойдёт об платформе «ИИГра», как перспективной среде для накопления данных и знаний о совместной деятельности коллективов человеческих и компьютерных агентов в средах разного уровня иммерсивности и сложности моделей описания контекста взаимодействия. С точки зрения системной инженерии это требует взглянуть на задачи развития среды одновременно из области методов искусственного интеллекта, проектирования человеко-машинного взаимодействия, прикладной семиотики и наук о данных. Отметим, что развитие подобных сред упирается не только в технологические и математические проблемы, но и является драйвером значимого философского дискурса.

Оптимальный транспорт и диффузионные модели



► Ст.преп. **Денис Романович Ракитин**², Центр глубинного обучения и байесовских методов, Департамент больших данных и информационного поиска

↳ В задачах генеративного моделирования непрерывных данных (изображения, видео) последние несколько лет лучше всего себя показывают диффузионные модели. Их идея состоит в рассмотрении процесса постепенного зашумления данных и в обучении обратного к нему, переводящего шум в целевое распределение. Помимо text-to-image генерации, в которой диффузионные модели прекрасно изучены и отлично работают, остается большое количество задач, на которые напрямую конструкция диффузионных моделей не переносится. Среди них, например, перевод между доменами (в частности, редактирование изображений) и задача семплирования из распределений без доступа к выборке. В обоих случаях обобщением диффузионных моделей является так называемый Мост Шрёдингера (МШ). Задача МШ состоит в построении процесса, стартующего из входного распределения и приходящего в целевое, при этом

минимизирующего некоторый функционал вдоль пути. В докладе мы обсудим задачу моста Шрёдингера, ее связь с оптимальным транспортом, потенциальные приложения и ценность на практике, а также возможные подходы к эффективному решению.

→ Четвёртая секция

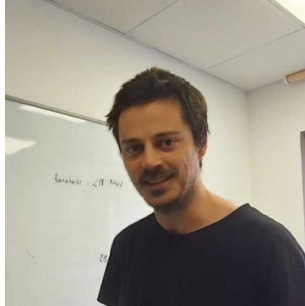
<null>



► Prof. **Attila Kertesz-Farkas**², Dr.Sc., департамент анализа данных и искусственного интеллекта, лаборатория искусственного интеллекта для вычислительной биологии

↳ Yes, the title is correct. We will be talking about how null distribution can be used to control risk (False Discovery Rate, Accuracy, etc.) in classification problems.

From computability to quasicrystals



► Assoc.Prof. **Thomas Fernique**^{[2](#)}, Ph.D., департамент анализа данных и искусственного интеллекта

↳ In the 1960s, computer scientists introduced a model of computation called 'tilings' — essentially infinite puzzles. Independently, in the 1980s, physicists discovered strange crystals, soon named quasicrystals. We will present the connection between the two, with beautiful illustrations (and some results or open questions).

Extreme Events in Action: What Triggers Them in Complex Systems?



► Assoc.Prof. **Srinivasan Sabarathinam**², Ph.D., департамент анализа данных и искусственного интеллекта, лаборатория моделирования и управления сложными системами

↳ Have you ever wondered what causes sudden, extreme events -like rogue waves in the ocean, unexpected power outages, or unpredictable spikes in biological systems? In this presentation, I explore how such rare but powerful bursts can emerge in complex systems, even when everything seems stable. Using mathematical models of chemical reactions, gene networks, KdV equations, we uncover the hidden triggers (mechanism) behind these extreme events and self organised criticalities. By studying these patterns across different systems from physics to biology, we aim to predict, and possibly prevent, extreme events in real-world applications. This work brings us closer to understanding how small changes can lead to big, sometimes dangerous, consequences in nature and technology.

О матричной задаче Прокруста и ее приложениях в глубинном обучении



► Доц. **Максим Владимирович Рахуба**²¹, к.ф.-м.н., базовая кафедра Института вычислительной математики им. Г.И. Марчука РАН, лаборатория матричных и тензорных методов в машинном обучении

↳ В рамках доклада будет рассмотрена матричная задача Прокруста о поиске ближайшей ортогональной матрицы. Мы обсудим применения этой задачи в глубинном обучении, включая новые методы обучения и сжатия нейронных сетей. Также будут представлены новые эффективные итерационные методы решения, наиболее подходящие при запуске на GPU.

Агентно-ориентированная модель обработки логов, основанная на синтезе методов промпт-инжиниринга и цепочки рассуждений



► Пригл. преп. Кирилл Игоревич Пашигорев², Сбер

↳ Рассмотрены гипотезы о возможности сокращения количества токенов промпта, полученного большой языковой моделью, с целью исключить малозначимые данные. Также проверяется способность большой языковой модели рассуждать на тему удаления токенов посредством итеративного взаимодействия приложения с большой языковой моделью, насколько результат такого рассуждения соответствует человеческому суждению. Такие очищенные промпты, которыми являются технические логи, требуются для генерации текстов, при этом важной особенностью работы является способность сгенерировать человекочитаемый текст из технических логов без применения методик дообучения модели и генерации с дополненным поиском. Для постановки задачи приведен принцип работы токенизатора как алгоритма. В качестве составляющих эксперимента представлены спроектированная модель архитектуры приложения в микросервисном стиле, системные промпты, метаданные датасетов, приведено описание самих

экспериментов. Эксперимент проводился на технических логах подсистем большой программной платформы. Приведен анализ полученных результатов, которые представлены графически, и сделаны выводы о проделанной работе. В заключении представлено подтверждение рассматриваемых гипотез, что большие языковые модели с помощью ограниченного набора методик на примере промпт-инжиниринга и цепочек рассуждений могут самостоятельно выполнять специфические задачи без участия человека, такие как самостоятельно сократить избыточные малозначимые данные в технических логах, при этом корректно определить полезные данные, достаточные для генерации текста о том, что написано в переданных логах. Также сформулированы вопросы, которые подлежат отдельным исследованиям, такие как фактчекинг ответов большой языковой моделью, и преимущества либо недостатки цепочки размышлений агентами, что является «нормальным» результатом, вариативность или постоянство.

Mathematical theory of optimal processes in a common path interferometer using a diffraction phase microscope



► Вед. аналитик **Наталья Анатольевна Талайкова**², лаборатория методов анализа больших данных

↪ The diffraction phase microscope is a common path interferometer that is installed in the image plane of an optical microscope. The paper presents a mathematical description of the processes of image formation in an optical microscope, which is the basis for the diffraction phase microscope. Also, a mathematical description of the processes of image formation in the interferometer for various parameters of the optical scheme is presented. We will discuss the choice of optimal parameters for the optical scheme of the interferometer, as well as its adjustment for observing objects of both biological and technical nature.

→ Пятая секция

Алгоритм клиент-серверного взаимодействия по незашифрованному соединению



► Проф. Дмитрий Владимирович Александров²¹, д.т.н., департамент программной инженерии, лаборатория облачных и мобильных технологий

↳ Рассматриваются различные известные сценарии клиент-серверного взаимодействия, результаты их сравнительного анализа, включая основные достоинства и недостатки. Кроме того, будут представлены сценарии, обеспечивающие определенную безопасность при обмене данными, включая аутентификацию, между клиентскими и серверными приложениями в незащищенных сетях.

О некоторых методологических аспектах обучения генеративных потоковых сетей



► Доц. **Сергей Владимирович Самсонов**²¹, к.мат.н., лаборатория стохастических алгоритмов и анализа многомерных данных

↳ Генеративные потоковые сети (generative flow networks, GFlowNets) - это семейство генеративных моделей, которые учатся генерировать объекты из заданного распределения вероятностей, известного, вообще говоря, лишь с точностью до нормирующей константы. Ключевая идея GFlowNets заключается в использовании двух политик: прямой, поэтапно конструирующей составные объекты, и обратной, последовательно их декомпозирующей. На этапе генерации работа GFlowNet'a сводится к генерации траекторий в соответствующим образом сконструированной среде в виде ориентированного ациклического графа.

В данном докладе мы рассмотрим проблемы и ограничения обучения генеративных потоковых сетей. В частности, мы обсудим целесообразность оптимизации обратной политики в GFlowNets, и изложим составные элементы теории, ослабляющей требование ацикличности порождающего графа. Если позволит время, мы также обсудим обобщения нашего подхода на непрерывное время в контексте диффузионных сэмплеров.

Кооперативная игра «Ханаби» – новый вызов ИИ?



► Доц. **Анастасия Александровна Оноприенко**²¹, к.ф.-м.н.,
департамент больших данных и информационного поиска

↳ Игры представляют собой вид человеческой деятельности, где условия задачи совершенно ясны и легко формализуются. В некоторых видах игр, таких как шахматы и го, успешная игра рассматривается как высшее достижение человеческого, «естественного» интеллекта. С середины XX столетия игры рассматриваются в качестве полигона для тестирования возможностей компьютера.

Игры часто представляют собой примеры многоагентного взаимодействия участников с противоположными интересами. Однако «Ханаби» является примером игры сотрудничества, в которой участники совместно достигают общей цели. На данный момент успехи ИИ в игре «Ханаби» довольно скромные: компьютер уступает даже командам из игроков-новичков.

Очевидное препятствие для «лобового» решения задачи автоматизации игры – «экспоненциальный взрыв». С одной

стороны, такой «взрыв» очевиден на практике при попытке запрограммировать игру, а с другой стороны, математически это выражается в виде утверждения об NP-трудности соответствующих вычислительных задач.

NP-полнота игры «Ханаби» была установлена даже для простейшего варианта игры в случае одного игрока, который видит всю колоду и пытается «разложить пасьянс»: выложить на столе карточки всех цветов. При этом карточки каждого цвета должны выкладываться по возрастанию (на каждой карточке написано число), и в любой момент времени у игрока в руке может быть лишь небольшое (заранее фиксированное) количество карт. Нами установлена точная граница параметров игры «Ханаби», при которой она всё ещё остаётся NP-полной, а при уменьшении любого из этих чисел игра «Ханаби» перестаёт быть NP-трудной (разумеется, если P не равно NP). Найденные нами значения параметров оказываются очень маленькими, что демонстрирует практическую невозможность точного анализа «Ханаби» даже при небольших параметрах игры.

Приватность данных в эпоху ML: обезличивание, анонимизация и синтетические данные на практике



► Ст.преп. **Юрий Владимирович Силаев**², департамент программной инженерии

↳ Доклад предлагает целостный, практический взгляд на приватность данных в эпоху активной аналитики с использованием ML. Цель доклада – перевести размытые требования к обезличиванию и анонимизации в воспроизводимый инженерный процесс с проверяемыми метриками риска и полезности. Сначала аккуратно разводим термины и режимы обработки, чтобы не смешивать псевдонимизацию, обезличивание и полную анонимизацию. Затем строим простую модель угроз и карту мест, где приватность может быть утрачена: на сыром источнике, при объединении наборов, в трансформациях подготовки данных, на витринах и в признаках (feature store), а также в артефактах обучения и инференса.

Центральная идея доклада – «приватность как код». Рецепты обезличивания выражаются декларативными конфигурациями, проходят ревью, тестируются, исполняются в конвейерах данных,

версионироваться и оставляют трассировку (lineage). Разбираем реестр детекторов персональных данных как основу покрытия: типы детекторов (регулярные выражения, справочники, NER), метрики точности и полноты, приоритизацию, управление версиями и разрешение конфликтов срабатываний. Интегрируем блокирующие проверки в CI/CD: негативные тесты, регрессионные тесты приватности, измерение риска и полезности на эталонных задачах.

Инженерные методы: токенизация и хранение соответствий в защищённом хранилище, формат-сохраняющее шифрование для устойчивых идентификаторов, маскирование и генерализация. Даём практический минимум по k-анонимности, l-разнообразию, t-близости и дифференциальной приватности (включая выбор бюджета ϵ), связывая выбор метода с измеримым снижением риска и с контролем деградации полезности на целевых метриках. Отдельный блок – синтетические данные. Обсуждаем критерии уместности (разработка, тестирование, обмен), способы генерации, метрики похожести и приватности, простые тесты на утечки, включая сценарии membership inference. Показываем, как оформлять отчёты для владельцев данных и аудиторов: гейты допуска, необходимые артефакты, протоколы обновления. Коротко обсуждаем роли и ответственность: владелец данных, инженер по приватности, платформа, служба информационной безопасности. Все приёмы иллюстрируются на реальном примере: создание прототипа озера данных с доменом «Студент», где объединяются неоднородные источники, нет явных ключей, требуется семантическая связка, регулярное обновление витрин и запрет на вывоз исходных данных.

Improved Stochastic Optimization of LogSumExp



► Доц. Егор Леонидович Гладин², Ph.D., департамент больших данных и информационного поиска, лаборатория теоретических основ моделей искусственного интеллекта

↳ Optimization problems across diverse fields often involve the LogSumExp function, or more generally, the log partition function. Examples include multiclass classification with softmax probabilities, entropy-regularized optimal transport (eOT), reinforcement learning (RL) with entropy penalty, and distributionally robust optimization (DRO). First-order methods are commonly used to solve such problems. However, challenges arise when the number of exponential terms inside the logarithm is large, since computing the gradient requires differentiating every term. We propose an approximation to the log partition function (and, in particular, LogSumExp) that can be efficiently optimized using stochastic gradient methods. The resulting optimization problem introduces a single additional scalar variable, and the new objective takes the form of an average (or expectation) of Softplus functions. Importantly, our approximation retains the convexity property of LogSumExp, thereby ensuring existence of global minima and the applicability of efficient first-order methods. The accuracy of approximation is controlled by a tunable parameter and can be made arbitrarily small. Existing approaches often

approximate the objective by computing a LogSumExp over a batch, which requires extremely large batch sizes to keep the bias sufficiently small. In contrast, our method is effective across both large and small batch regimes. We demonstrate the usefulness of the proposed approximation through numerical experiments on distributionally robust optimization with Kullback-Leibler (KL) divergence and continuous entropy-regularized optimal transport problems.

Длинные строки составных значений полиномов и базис порядка 2



► Доц. **Артём Олегович Радомский**², к.ф.-м.н., департамент больших данных и информационного поиска, лаборатория теоретической информатики

↪ Доказано, что для любого неприводимого (над полем рациональных чисел) полинома с целыми коэффициентами для любого достаточно большого числа N найдутся две строки последовательных натуральных чисел $I_1 = \{n_1 - m, \dots, n_1 + m\}$ и $I_2 = \{n_2 - m, \dots, n_2 + m\}$ такие, что $n_1 + n_2 = N$, $m > (\log N) (\log \log N)^{(1/325525)}$ и $f(n)$ – составное число для любого n , принадлежащему объединению строк I_1 и I_2 . Это обобщает предыдущий результат, который был получен только для полинома $f(n) = n$. Также данная теорема дает новый и нетривиальный пример аддитивного базиса порядка 2.

Бифуркационные модели нелинейных уравнений в частных производных



► Проф. **Василий Александрович Громов**², д.ф.-м.н., департамент анализа данных и искусственного интеллекта, лаборатория анализа семантики, лаборатория интеллектуальных систем и структурного анализа

↳ Цель настоящего исследования — разработка численного метода бифуркационного анализа для нелинейных уравнений в частных производных, основанного на методе сведения уравнений в частных производных к обыкновенным с использованием теоремы Колмогорова-Арнольда.

Методы. В данной работе описывается метод сведения уравнений в частных производных к обыкновенным с использованием теоремы Колмогорова-Арнольда, а также метод бифуркационного анализа нелинейных краевых задач для обыкновенных дифференциальных уравнений.

Результаты. В работе представлен новый метод решения и бифуркационного анализа нелинейных краевых задач для уравнений в частных производных, допускающих вариационную

постановку. Метод был применён к нелинейной двумерной задаче Брату с граничными условиями типа Дирихле.

Заключение. Разработан новый метод бифуркационного анализа для нелинейных уравнений в частных производных, а именно был предложен метод сведения уравнений в частных производных к обыкновенным, который позволяет применять разработанный аппарат бифуркационного анализа для краевых задач обыкновенных дифференциальных уравнений. Метод позволяет строить бифуркационные картины для нелинейных уравнений в частных производных произвольного вида.

GNN-based neutron reconstruction in the HGND at the BM@N experiment



► М.н.с. **Владимир Олегович Бочарников**², лаборатория методов анализа больших данных

↳ The Highly Granular Neutron Detector (HGND) is designed for the BM@N experiment to study neutron emission in heavy ion collisions at beam energies up to 4A GeV. This detector allows the identification of neutrons and the reconstruction of their energies using time-of-flight method, facilitating the assessment of neutron yields and azimuthal flow. The challenging neutron energy range of \$0.5-4\$ GeV and large background contributions require sophisticated reconstruction algorithms. In this contribution, we present a Graph Neural Network-based approach to the neutron reconstruction problem and discuss preliminary results of the proposed algorithm.

Мозг и искусственный интеллект: интерпретируемые модели и новые горизонты анализа нейроданных



► **Н.с. Илья Владимирович Семенов**, департамент анализа данных и искусственного интеллекта

↳ В последние годы методы машинного обучения всё активнее применяются для анализа нейрофизиологических данных. Эти подходы позволяют не только решать прикладные задачи в клинической практике, но и формировать новые исследовательские направления на стыке нейронаук и искусственного интеллекта. В докладе рассматриваются несколько проектов, демонстрирующих, как классические и современные алгоритмы могут быть адаптированы для работы с различными типами мозговых сигналов. Одна из ключевых тем – функциональное картирование речевых зон на основе инвазивных электрофизиологических записей. Применение компактных моделей машинного обучения позволяет повысить надёжность определения функционально значимых областей и дополнить существующие клинические методы. Ещё одно направление связано с созданием специализированных датасетов для исследований языковых процессов. Такие данные открывают возможности для изучения реакции мозга на различные

лингвистические стимулы и служат основой для дальнейших сравнительных исследований об обработке языковой информации человеком и большими языковыми моделями. Отдельное внимание уделено применению нейросетевых архитектур к неинвазивным данным ЭЭГ, МЭГ и фМРТ. Рассматриваются примеры разработки и модификации моделей для анализа сигналов, декодирования стимулов и восстановления визуальной информации по активности мозга. Эти проекты объединяет стремление к сочетанию точности, интерпретируемости и практической применимости методов. Таким образом, доклад охватывает широкий спектр задач – от функционального картирования до генеративных подходов, и подчёркивает потенциал машинного обучения как инструмента, позволяющего глубже понять нейрофизиологические процессы и расширить возможности анализа нейроданных.

Постеры (стендовые доклады)

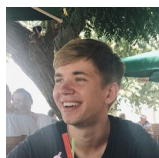
В этом разделе приводятся темы и аннотации постеров, принятых на конференцию.

→ Первая постерная сессия



Анализ паттернов в процессах онлайн-оценивания с использованием иерархических моделей

► Ст.преп. Антонина Константиновна Бегичева^{[2](#)},
департамент программной инженерии, лаборатория
процессно-ориентированных информационных
систем

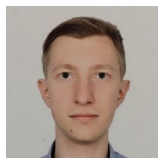


**Генеративное моделирование в нетривиальной
топологии репараметризацией через накрытия – и
байесовское обучение**

► Н.с. Максим Евгеньевич Бекетов^{[2](#)}, департамент
больших данных и информационного поиска,
лаборатория стохастических алгоритмов и анализа
многомерных данных

↳ Представьте, что вы сделали выборку фотографий одного объекта, повернутого на различные углы. Ясно, что за полученными данными скрыто "латентное" многообразие группы вращений. Обучение модели генерации таких топологически нетривиальных данных – область, которой посвящен данный доклад. Давно известно, как строить необходимые для этого модели в случае, когда латентное многообразие данных обладает структурой группы (Ли), на котором ясно как репараметризовать аналог нормального распределения – часто выбираемого вариационным постериором в базовом примере генеративных моделей – вариационном автокодировщике (VAE). В совместной работе с Павлом Сноповым мы наконец показали, что в этой задаче требование наличия у многообразия структуры группы можно ослабить – достаточно уметь накрывать его другим, на

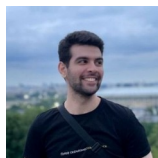
котором репараметризация тем или иным способом возможна. Таким образом класс многообразий, на которых возможно подобное генеративное моделирование – расширяется. Мы демонстрируем приложение этой техники для построения VAE-модели, латентным пространством которой является бутылка Клейна (который мы назвали KleinVAE) – многообразие, которое удивительным образом возникает в задачах компьютерного зрения, и учет структуры которого (путем обучения KleinVAE как априорного распределения весов сверточных фильтров) в теории позволяет реализовать байесовское обучение сверточных нейросетей (CNN), сулящее более быструю сходимость и лучшую обобщающую способность. На данный момент работа находится на ревью конференции AAAI-26.



Статистический анализ генеративных диффузионных моделей

► М.н.с. Константин Дмитриевич Яковлев [7](#),
департамент больших данных и информационного
поиска, лаборатория теоретических основ моделей
искусственного интеллекта

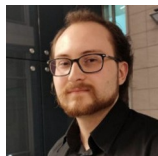
↳ Диффузионные модели являются одним из передовых методов генеративного моделирования, способным создавать изображения высокого разрешения. В их основе лежит идея искажения исходных данных с помощью гауссовского шума и восстановления обратного преобразования. Как правило, для этого используется метод сопоставления градиентов логарифмов плотностей или скор-функций (denoising score matching). В качестве класса всевозможных оценок истинной скор-функции выступает класс нейронных сетей с функцией активации $\text{ReLU}(x) = \max(x, 0)$ специальной архитектуры. Формулируется результат об аппроксимационных свойствах рассматриваемых нейронных сетей в терминах дивергенции Фишера. В качестве оценки скор-функции рассматривается минимизатор эмпирического риска. Установлена скорость сходимости оценки метода сопоставления скор-функций к градиенту логарифма плотности распределения элементов выборки. Скорость сходимости в терминах объема выборки зависит исключительно от эффективной размерности. Кроме того, в оценках отслеживается зависимость от размерности пространства и других параметров диффузионной модели.



Modern Generative Models in 3D Computer Vision

► Преп. Мишан Хаммад оглы Алиев², департамент больших данных и информационного поиска, лаборатория теоретических основ моделей искусственного интеллекта

↳ Доклад посвящён обзору и анализу актуальных достижений в применении генеративных моделей для 3D-компьютерного зрения. Особое внимание будет уделено методам генерации и реконструкции трёхмерных структур на основе диффузионных моделей.



Погружение уточняющих типов данных в системы с зависимыми типами

► Преп. Павел Павлович Соколов², департамент больших данных и информационного поиска, лаборатория теоретической информатики

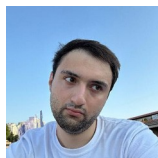
↳ Зависимые типы позволяют создавать программы, корректные по построению, и формально верифицировать математические теоремы. Однако порог вхождения, который необходимо преодолеть для уверенного программирования с использованием зависимых типов, зачастую оказывается слишком высок, так что интересным промежуточным решением являются уточняющие типы - это обычные типы данных, к которым привыкли программисты, снабжённые разрешимым предикатом, с помощью чего можно получать автоматические доказательства несложных утверждений о корректности функциональных программ. Применимость таких систем зачастую ограничена мощностью солвера, доказывающего допустимые предикаты; интересным направлением развития является погружение уточняющих типов в системы с зависимыми типами для, во-первых, предоставления возможности ручного доказательства теорем, с которыми не справляется солвер; во-вторых, для предоставления возможности пользовательского расширения используемого солвера. В настоящем докладе будут рассмотрены мотивирующие примеры программ, существующие языковые решения и будет предложена авторская система типов, подходящая для реализации этой техники.



Генеративные модели для быстрой симуляции черенковского детектора ФАРИЧ

► Преп. Фома Александрович Шипилов², департамент больших данных и информационного поиска, лаборатория методов анализа больших данных

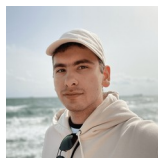
↳ In the end-cap region of the SPD detector complex, particle identification will be provided by a Focusing Aerogel RICH detector (FARICH). FARICH will primarily aid with pion / kaon separation in final open charmonia states (momenta below 5 GeV/c). A free-running (triggerless) data acquisition pipeline to be employed in the SPD results in a high data rate necessitating new approaches to event generation and simulation of detector responses. Several machine learning based approaches are described here, generating high-level reconstruction observables as well as full Cherenkov rings using a generative neural network. The fast simulation is trained using Monte-Carlo simulated data samples. We compare different approaches and demonstrate that they produce high-fidelity samples.



GAS: улучшение генерации диффузионной модели при малом числе шагов

► Преп. Александр Артурович Оганов², департамент больших данных и информационного поиска

↳ В области генеративного моделирования диффузионные модели показали, что генерация изображений может быть качественной и разнообразной. Главным недостатком современных моделей является долгое сэмплирование, так как для качественной генерации требуется большое число шагов (вызовов нейронной сети). Одним из способов ускорения генерации является разработка специальных солверов обыкновенных дифференциальных уравнений (ОДУ). В недавних работах был предложен способ дистилляции солвера учителя с большим числом шагов в солвера студента, который использует меньшее число шагов. Мы предлагаем Generalized Solver (GS) – новую и эффективную параметризацию солвера, которую обучаем с учетом данных и обученной диффузионной модели. Без специальных техник в обучении GS достигает лучшего качества, чем существующие методы. Кроме того, мы объединили стандартную задачу дистилляции с состязательным обучением, что позволило дополнительно улучшить качество. Generalized Adversarial Solver (GAS) использует эффективную параметризацию, состязательное обучение и показывает лучшее качество при генерации за малое число шагов.



Автоматическая классификация 2-значных групп

► Стажёр Александр Михайлович Левин²,
лаборатория алгебраической топологии и ее
приложений

↳ В докладе я расскажу о том, как автоматически классифицировать v_2 -значные группы небольших размеров. 2-значная группа задаётся так: дано множество X и «умножение», которое вместо одного результата выдаёт два (с учётом кратности), при этом остальные свойства группы почти такие же, как у обычной. Такие структуры возникают из обычных групп через «склейку» по инволюции, а ещё полезны как дискретные модели для задач в топологии и алгебре.

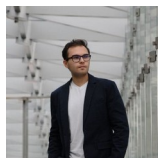


Flow matching models for DNA sequence generation

► Intern Muhammad Hashaam (Стажёр Мухаммад Хасхаам)², лаборатория биоинформатики

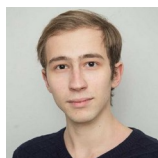
↳ Generating regulatory DNA sequences with defined activity levels in specific cell types has important applications in medicine and synthetic biology. However, the high dimensionality and complex distribution of DNA sequences make this task challenging for existing generative models. Autoregressive approaches, while successful for discrete data such as text, impose a sequential inductive bias that is not biologically meaningful for DNA, and they offer limited flexibility for conditional generation. Diffusion-based methods allow powerful conditioning techniques such as classifier guidance and classifier-free guidance but are naturally designed for continuous rather than discrete data. Here, we introduce discrete flow matching as a new generative paradigm for biological sequence design. Discrete flow matching model training is similar to diffusion models where the model learns to denoise the data. Generation, however shares some similarity with autoregressive models as it unmask the tokens gradually however the unmasking is not sequential. We train models to generate cell-type-specific enhancer sequences and promoter sequences with defined transcription initiation signal profiles. To evaluate our approach, we employ metrics such as Fréchet Biological Distance (FBD) and perplexity. FBD measures the distance between the distributions of the generated DNA sequences and real DNA sequences. These metrics

demonstrate that discrete flow matching produces realistic and controllable DNA sequences. Our results indicate that this method is a promising step toward state-of-the-art generative modeling for regulatory genomics.



Современные методы контроля генеративных моделей для видео и изображений

► Стажёр Ильгар Габитович Мамедов², лаборатория «Искусственный интеллект в математических финансах»



Verifying Soundness of Integrated Process Models with Data and Resources

► Асп. Николай Михайлович Суворов², лаборатория процессно-ориентированных информационных систем

↳ Для формального представления классических моделей процессов, как правило, используются сети Петри. Для моделей, представленных в виде сетей Петри, ранее было предложено множество свойств бездефектности и, соответственно, множество алгоритмов для их верификации. Сложность возникает с моделями процессов с данными – обогащение потока управления данными приводит к значительному увеличению размеров сети и к еще большему увеличению пространства состояний, что делает практически неприменимой проверку корректности моделей процессов с данными на сетях Петри. В случае бесконечных областей значений данных в процессе, например, в случае наличия вещественных переменных, процесс в общем случае не может быть представлен в виде сети Петри. Для решения данных проблем были предложены различные формализмы, позволяющие лаконично представить взаимодействие потоков управления и данных, одним из которых являются сети Петри с данными. В сетях Петри с данными каждому переходу сопоставлено ограничение, включающее входные и выходные условия на ограниченный набор переменных. Входные условия определяют, при каких значениях переменных переход может сработать. Выходные условия декларативно задают трансформацию значений переменных при срабатывании перехода. Проверка свойств бездефектности для таких моделей является непростой задачей, так как срабатывание каждого перехода такой сети изменяет не только маркировку, но и состояние переменных. В данном докладе мы

покажем, как возможно проверить бездефектность классических моделей с данными и ресурсно-ориентированных моделей с данными, представленных в виде сетей Петри с данными, а также продемонстрируем инструмент, реализующий данные алгоритмы и визуализирующий результаты верификации. Сам прототип размещен в открытом доступе и доступен для скачивания. Результаты проведенных экспериментов показывают практическую применимость предложенных алгоритмов для моделей малых и средних размеров, а также более высокую скорость верификации бездефектности по сравнению с существующими решениями.



The EG-TD3 Machine Learning Architecture: Evolutionary-Guided Twin Delayed Deep Deterministic Policy Gradient

► Ph.D. student Djambong Tenkeu Hank-Debain (Асп. Джамбонг Тенке Ханк-Дебэн)², департамент программной инженерии

↳ This research introduces a novel, theoretical, and hybrid machine learning architecture - the Evolutionary-Guided Twin Delayed Deep Deterministic Policy Gradient (EG-TD3) — designed for adaptive control in networking systems such as packet buffering, congestion avoidance, and resource allocation. The model combines the strengths of Reinforcement Learning and evolutionary optimization to address the limitations of pure Reinforcement Learning methods in noisy and non-stationary environments. Specifically, the architecture integrates a Twin Delayed Deep Deterministic Policy Gradient (TD3) actor-critic framework for stable and efficient policy learning in continuous action spaces, a Differential Evolution (DE) layer for global optimization of policy network parameters and hyper-parameters to reduce sensitivity to initial conditions and local optima, and an online transfer learning mechanism that allows pre-trained policies to adapt quickly to new network conditions with minimal fine-tuning. The proposed EG-TD3 architecture is generic and can be applied to various network control problems, although this work focuses on its theoretical application to dynamic packet buffering, a core problem in service meshes, IoT networks, and edge systems.



A Multimodal Deep Learning Framework for Protein–Protein Interaction Prediction Using Surface and Structural Features

► Ph.D. student David Moreano Brandon Arteaga (Асп. Давид Мореано Брандон Артеага)^{[2](#)}, лаборатория биоинформатики

↳ A comprehensive understanding of protein-protein interactions (PPIs) is imperative for elucidating the complex biological mechanisms that are governed by the multifaceted relationships between proteins and other biomolecules. Employing deep learning (DL) algorithms instead of experimental methods to automate PPI prediction is a contemporary solution to this challenge, as it is less resource-intensive and time-consuming. However, the majority of research conducted on PPIs has focused predominantly on sequence and three-dimensional (3D) structural information. There has been comparatively less emphasis on protein surface information, despite the crucial role that this plays in PPIs. In addition, the majority of current approaches that incorporate multimodal information sources typically entail the simple concatenation of features from distinct modalities. However, this approach fails to acknowledge the interrelationships and context that are imperative for optimising the complementarity between the various protein features, thereby compromising the potential for comprehensive analysis and interpretation. To tackle these challenges, the present study proposes a novel approach that incorporates protein surface, structural, and graph based features, as well as the implementation of linear projectors and the transformer attention mechanism for the prediction of PPIs. The experimental results on the test set, in conjunction with a comparison with analogous studies, demonstrate the efficacy of the proposed approach. Furthermore, we conduct ablation studies that demonstrate the relevance of protein surface information for predicting PPIs.

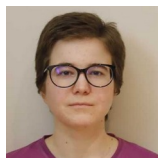


Применение размытых магнитудных гомотопий для классификации ориентированных сетей мозга

► Асп. Александр Сергеевич Качура^{[2](#)}, департамент больших данных и информационного поиска, лаборатория методов анализа больших данных

↳ Анализ сетей является распространённым методом диагностики некоторых неврологических заболеваний и психических расстройств. Одним из типов сетей, изучаемых в этой области, являются функциональные коннектомы. Они отражают

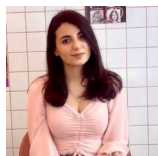
функциональные связи между областями мозга. Уровень взаимосвязи обычно измеряется при помощи ненаправленных корреляций, в то время как функциональные связи между областями мозга по своей сути являются направленными. Устойчивые гомологии являются набирающим популярность подходом в анализе данных, который можно применять в том числе и к анализу графов. Суть этого подхода заключается в отслеживании моментов появления и исчезновения топологических структур в исследуемом пространстве при изменении масштаба его рассмотрения при помощи алгебраических инвариантов, называемых гомологиями. В применении к взвешенным графам, к которым относятся и функциональные коннектомы, метод устойчивых гомологий позволяет изучать, как меняется топологическая структура графа при изменении порога фильтрации рёбер. Это даёт возможность выявить устойчивые к шуму топологические признаки, что особенно важно в случае, когда структура графа не задана изначально, а оценивается по данным. Существует несколько типов гомологий. Симплициальные гомологии, которые исторически первыми стали применяться для вычисления устойчивых гомологий, не позволяют учесть направления рёбер графа. Наиболее прямым обобщением данного типа гомологий на случай орграфов являются гомологии ориентированного флагового комплекса. Однако при их использовании некоторая информация о направлении рёбер теряется. Одной из теорий гомологий, дающих возможность сохранить большую долю информации о направлениях рёбер, являются размытые магнитудные гомологии. В работе предложен метод классификации ориентированных функциональных коннектомов при помощи кривых Бетти размытых магнитудных гомологий, одного из способов численного описания устойчивых гомологий. Изучена применимость представленного алгоритма к задаче диагностирования расстройств аутистического спектра (РАС) по изображению головного мозга, полученного при помощи функциональной магнитно-резонансной томографии. Проведено экспериментальное тестирование предложенного подхода на наборе реальных данных ABIDE. Его результаты позволяют судить о том, что разработанный алгоритм может быть полезен для применения в качестве практического инструмента.



Поиск мутаций в геноме возбудителя туберкулеза, определяющих устойчивость к лекарственным препаратам, с помощью языковых моделей

► Стажёр Наталия Александровна Дьячкова²,
лаборатория статистической и вычислительной
геномики

↳ Туберкулез является ведущей причиной смерти от инфекционных заболеваний. Лечение туберкулеза затрудняется тем, что бактерия-возбудитель может относительно легко развивать устойчивость к антибиотикам, и из-за этого туберкулез лечится одновременно несколькими препаратами. Однако, если бактерия уже обладает устойчивостью к одному из используемых препаратов, то она может достаточно быстро выработать устойчивость и к остальным. По этой причине для эффективного лечения необходимо заранее знать о наличии устойчивости к тому или иному препарату. Для предсказания лекарственной устойчивости существует множество методов, в том числе с использованием машинного обучения, однако, существующие модели не учитывают корреляцию между устойчивостью к разным препаратам, возникающую вследствие их одновременного применения. Из-за этого мутации, ассоциированные с устойчивостью к одному препарату, часто интерпретируются моделями как ассоциированные с устойчивостью к препарату, скоррелированному с ним. Такие модели биологически неправдоподобны, и потому их трудно достоверно интерпретировать. Их нельзя использовать для поиска новых признаков, ассоциированных с устойчивостью к конкретному препарату, или для достоверной классификации монорезистентных штаммов - штаммов, устойчивых только к одному препарату. В данном исследовании разрабатывается нейросетевая модель многометковой классификации для определения устойчивости одновременно к нескольким препаратам, способная выявлять меткоспецифичные признаки.

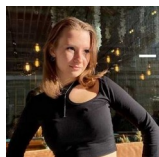


Optimization of Beamlines in particle accelerators using RL: a simulation based approach

► Ph.D. student Anwar Ibrahim (Асп. Анвар Ибрахим)²,
департамент больших данных и информационного
поиска, лаборатория методов анализа больших
данных

↳ Optimizing particle accelerator beamlines is a challenging task due to the high-dimensional parameter space, nonlinear beam dynamics, and strong interdependencies

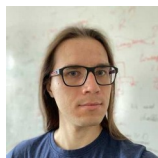
between controllable elements such as quadrupole and dipole magnets. Manual tuning requires extensive expertise and time, often leading to suboptimal performance, while traditional optimization methods struggle to generalize across diverse accelerator configurations. In this work, we address this problem by introducing RLACC, a reinforcement learning (RL) framework for beamline optimization based on Elegant simulations. RL excels in such optimization tasks by enabling agents to learn through trial and error in complex environments, as demonstrated in robotics and control systems. RLACC provides a flexible, custom RL environment that encodes beam states (particle distributions at watch points), continuous 4D actions (magnet strengths and steering), and rewards based on transmission efficiency, compatible with libraries like Stable-Baselines3. We implemented agents using Deep Deterministic Policy Gradient (DDPG) and Soft Actor-Critic (SAC) from scratch. Results show 100% particle transmission on a simple 7-quadrupole beamline, 68% on a higher-emittance variant, and 91.3% on a complex 10-quadrupole, 3-dipole layout—consistent with human experts. These outcomes highlight RL's potential to automate tuning, minimize manual effort, and scale to intricate setups, with ongoing work exploring advanced algorithms for real-world integration.



Реализация алгоритма для рекомендательных систем на основе масштабируемых контекстуальных бандитов

► Стажёр Евгения Константиновна Шустова²,
лаборатория стохастических алгоритмов и анализа
многомерных данных

↳ В работе рассматривается подход к обновлению на основе динамического низкорангового приближения с использованием метода Projector-Splitting Integrator. Данный метод обладает адаптивной природой: величина модификации разложения напрямую зависит от масштаба изменений в данных, PSI реализует инкрементальную схему, при которой обновление выполняется только на основе вновь поступившей информации. Дополнительно учитывается, что распределения контекста для каждого действия обычно имеют низкий ранг, поэтому поддержание и использование этой низкоранговой структуры позволяет эффективно вычислять обратные матрицы в онлайн-сценарии. С помощью интегратора был модифицирован алгоритм для рекомендательных систем на основе контекстуальных бандитов (LinUCB), который сохраняет качество базового алгоритма, но работает значительно быстрее и требует существенно меньше памяти.



Aspartik b3 --- быстрый и эффективный филогенетический анализ

► Стажёр Андрей Андреевич Колчин², лаборатория
статистической и вычислительной геномики

↳ Филогенетика – это наука о восстановлении дерева жизни. Благодаря цифровой революции и экспоненциальному росту генетических данных за последние 20 лет на первый план вышла вычислительная филогенетика. Эта дисциплина использует статистические методы и настраиваемые биологические модели для анализа генетических данных: восстановления филогенетических деревьев, миграций древних популяций, географического распространения видов, скорости мутации, и многого другого.

В докладе будет представлен Aspartik b3 – быстрый и эффективный по памяти программный пакет для филогенетического анализа. Будет пояснено, как b3 добивается большей скорости по сравнению с существующими программами. Также будет показано, как b3 позволяет повысить точность передовых исследований современной филогенетики.



Автоматизация формирования текстовых описаний глазного дна с использованием языковых моделей в информационной системе поддержки врачей- офтальмологов

► Максим Александрович Ставцев, департамент
программной инженерии

↳ В условиях возрастающей нагрузки на систему здравоохранения и нехватки квалифицированных специалистов, автоматизация рутинных задач врачей-офтальмологов становится критически важной. Значительную часть времени приема, до 35%, занимает составление текстовых описаний параметров глазного дна.

Данный доклад посвящен описанию разработки и интеграции модуля автоматизированного текстового описания в информационную систему помощи офтальмологам «EYAS». Целью работы являлось создание инструмента, способного

генерировать связные, структурированные и соответствующие медицинским стандартам текстовые заключения на основе параметров, полученных из анализа изображений глазного дна.

В ходе работы был проведен анализ открытых языковых моделей на собственном наборе данных. В итоге для решения этой задачи была выбрана и дообучена языковая модель YandexGPT-5-lite. Применение технологии LoRA позволило улучшить качество генерируемых текстов, что подтверждено экспертной оценкой врача-офтальмолога.

Разработанный модуль интегрирован в микросервисную архитектуру системы «EYAS» и включает в себя функционал по созданию и детализации текстовых описаний. В докладе будут представлены архитектурные решения, этапы программной реализации, а также результаты тестирования и демонстрация работы.

Диффузионные модели для физических процессов

► Александр Андреевич Дорош

↳ Моделирование физических экспериментов является ключевым инструментом в физике высоких энергий (НЕР), позволяя эффективно анализировать огромные объемы данных, получаемых на таких установках, как Большой адронный коллайдер (LHC). Традиционные методы, в частности Монте-Карло, сталкиваются с растущими вычислительными затратами, что стимулирует поиск альтернатив. В данной работе исследуется применение глубоких генеративных моделей, а именно диффузионных моделей, для решения задачи моделирования развития ливней частиц в калориметре. Актуальность работы обусловлена необходимостью создания точных и эффективных инструментов для связи теоретических предсказаний с экспериментальными данными.

Целью исследования является разработка диффузионной модели, превосходящей по своим характеристикам существующее решение на основе Generative Adversarial Networks (GAN) – модель CaloGAN. В рамках работы планируется выбор оптимальной архитектуры диффузионной модели, изучение различных техник глубокого обучения для улучшения результатов и последующее сравнительное тестирование с CaloGAN. Объектом исследования выступает распределение энергии ливней, инициированных фотонами в калориметре LHC, а предметом – сама диффузионная модель как представитель класса глубоких генеративных моделей.

Основной вызов при использовании диффузионных моделей в вычислительно нагруженных областях, подобных НЕР, заключается в их относительно низкой скорости генерации по сравнению с GAN. Для преодоления этого ограничения в работе рассматриваются современные методы ускорения диффузионных моделей. Среди них – использование сжатых семплеров, таких как DPM-Solver, которые позволяют значительно сократить количество шагов вывода без потери качества генерации. Также исследуются архитектурные оптимизации, включая дистилляцию знаний и проектирование более эффективных сетей-денойзеров. Комбинация этих подходов позволяет добиться скорости генерации, сравнимой с GAN, сохраняя при этом превосходство диффузионных моделей в качестве и стабильности обучения. Таким образом, работа нацелена не только на демонстрацию качественного превосходства над CaloGAN, но и на обеспечение практической применимости модели в условиях, требующих высокой вычислительной эффективности.

Исследование вклада гипермутаций рецепторов в В-клетках при сахарном диабете второго типа

► Дмитрий Валерьевич Босов

GLGENN: A Novel Parameter-Light Equivariant Neural Networks Architecture Based on Clifford Geometric Algebras

► Асп. Екатерина Романовна Филимошина

↳ We propose, implement, and compare with competitors a new architecture of equivariant neural networks based on geometric (Clifford) algebras: Generalized Lipschitz Group Equivariant Neural Networks (GLGENN). These networks are equivariant to all pseudo-orthogonal transformations, including rotations and reflections, of a vector space with any non-degenerate or degenerate symmetric bilinear form. We propose a weight-sharing parametrization technique that takes into account the fundamental structures and operations of geometric algebras. Due to this technique, GLGENN architecture is parameter-light and has less tendency to overfitting than baseline equivariant models. GLGENN outperforms or matches competitors on several benchmarking equivariant tasks, including estimation of an equivariant function and a convex hull experiment, while using significantly fewer optimizable parameters. This work was presented at ICML 2025 in Vancouver, Canada (<https://icml.cc/virtual/2025/poster/45802>).



Компьютерный анализ метафорических значений глаголов контакта в итальянском и английском языках с использованием контекстуальных эмбеддингов и консенсусной кластеризации

► Милена Александровна Камская², департамент больших данных и информационного поиска

↳ Автоматическое распознавание семантических сдвигов в естественном языке остается одной из ключевых задач современной компьютерной лингвистики, критически важной для развития систем понимания текста и анализа эволюции языка. В работе представлен комплексный подход к автоматическому выявлению семантических сдвигов глаголов физического контакта для английского и итальянского языков с использованием современных методов машинного обучения. Основной акцент сделан на создании эффективного NLP-пайплайна для обработки неструктурированных текстовых данных и семантической кластеризации. В работе использованы корпуса социальных сетей объемом свыше 10 млн сообщений для каждого языка, что обеспечило репрезентативность данных для анализа современных языковых тенденций.

Разработана методология построения контекстуальных представлений на основе адаптированных трансформерных моделей (UmBERTo, Sentence-RoBERTa) с целевым маскированием исследуемых лексических единиц. Проведено систематическое сравнение 10 стратегий агрегации слоев нейронных сетей, выявлены оптимальные подходы для извлечения семантической информации. Для снижения размерности эмбеддингов применен алгоритм UMAP с тщательной оптимизацией гиперпараметров через анализ 96 конфигураций на основе метрик сохранения структуры данных.

Ключевым техническим решением является предложенный метод консенсусной кластеризации, объединяющий результаты DBSCAN, HDBSCAN и Spectral Clustering через построение взвешенной консенсусной матрицы и последующий анализ связанных компонент. Валидация проведена на 17 глаголах (460 контекстов для каждого глагола), показав способность метода выделять как буквальные, так и метафорические употребления с высокой семантической однородностью, что подтверждено анализом междоидов и ближайших контекстов.

Создан масштабируемый программный пайплайн, реализованный с использованием современного стека технологий (Python, spaCy, Hugging Face Transformers, UMAP-learn, HDBSCAN). Предложенная архитектура обеспечивает

воспроизводимость результатов и адаптируемость к различным языкам и типам семантических сдвигов, открывая перспективы для применения в задачах автоматического анализа эволюции языка и детекции семантических аномалий в текстовых данных.

Улучшение алгоритма планирования для группы взаимозаменяемых агентов в непрерывной среде

► Дмитрий Евгеньевич Авдеев

↳ В настоящее время все большее распространение получает задача поиска безопасных маршрутов для группы агентов, действующих в общем пространстве. Например, такая задача возникает на автоматизированных складах или при управлении юнитами в видеоиграх. В настоящей работе рассматривается один из вариантов такой задачи, называемый анонимным много-агентным поиском путей с учетом размеров (Anonymous Multi-Agent Path Finding for Large Agents, AMAPF-LA). В отличие от классической задачи много-агентного поиска пути (MAPF), в этой вариации за агентами не закреплена конкретная цель, достаточно множеству агентов занять множество целей. Кроме того, агенты действуют в непрерывном пространстве и обладают размерами.

В настоящий момент существует лишь небольшое количество работ, рассматривающих такую постановку задачи: значительная часть исследований, посвященных тематике MAPF, опираются на дискретное представление среды (граф специального вида) и не учитывают размеры агентов.

Существующие алгоритмы, работающие в непрерывном пространстве, налагают строгие ограничения на входные данные. В данной работе был рассмотрен один из таких алгоритмов, накладывающий ограничения на минимальное расстояние от начальных/целевых позиций до препятствий и расстояние между начальными/целевыми позициями. Была предложена модификация алгоритма, позволяющая ослабить ограничение на расстояние между начальными/целевыми позициями. Было показано, что предложенная модификация сохраняет теоретические свойства оригинального алгоритма, в том числе гарантию достижения всеми агентами целей без столкновений, корректность алгоритма. В то же время, предложенная модификация позволяет работать с более плотными расположением агентов и целей, что повышает его применимость на практике.

MaterialFusion: Высококачественный, Zero-shot и управляемый перенос материалов с помощью диффузионных моделей

► Камиль Зуфарович Гарифуллин², центр глубинного обучения и байесовских методов

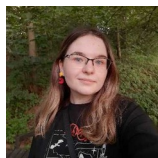
→ Вторая постерная сессия



High-Order Error Bounds for Markovian LSA with Richardson-Romberg Extrapolation

► М.н.с. Илья Валерьевич Левин², департамент больших данных и информационного поиска, лаборатория стохастических алгоритмов и анализа многомерных данных

↳ In this paper, we study the bias and high-order error bounds of the Linear Stochastic Approximation (LSA) algorithm with Polyak-Ruppert (PR) averaging under Markovian noise. We focus on the version of the algorithm with constant step size α and propose a novel decomposition of the bias via a linearization technique. We analyze the structure of the bias and show that the leading-order term is linear in α and cannot be eliminated by PR averaging. To address this, we apply the Richardson-Romberg (RR) extrapolation procedure, which effectively cancels the leading bias term. We derive high-order moment bounds for the RR iterates and show that the leading error term aligns with the asymptotically optimal covariance matrix of the vanilla averaged LSA iterates.

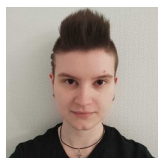


Снижение вычислительных затрат и повышение эффективности мультиагентных ИИ-систем посредством арбитража

► М.н.с. Людмила Резуник², департамент программной инженерии, лаборатория облачных и мобильных технологий

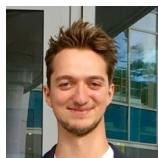
↳ Системы на основе ИИ сталкиваются с проблемой высокой вычислительной

стоимости, эффективного управления ресурсами и падения производительности при росте нагрузки. В докладе рассматриваются подходы, направленные на решение этих проблем, с акцентом на алгоритм арбитража ИИ-моделей, обеспечивающий динамическое распределение задач между ИИ-агентами в мультиагентных системах. Предложенный подход учитывает загрузку вычислительных узлов, характеристики запросов и специфику моделей. Результаты показывают потенциал арбитража как инструмента для построения масштабируемых и ресурсно-эффективных ИИ-систем.



Крупномасштабная структура естественного языка, или поймай бота

► М.н.с. Александра Сергеевна Коган², департамент анализа данных и искусственного интеллекта, лаборатория анализа семантики



AutoJudge: Judge Decoding Without Manual Annotation

► М.н.с. Федор Сергеевич Великонивцев², лаборатория компании Яндекс

↳ We introduce AutoJudge, a framework that accelerates large language model (LLM) inference with task-specific lossy speculative decoding. Instead of matching the original model output distribution token-by-token, we identify which of the generated tokens affect the downstream quality of the generated response, relaxing the guarantee so that the "unimportant" tokens can be generated faster. Our approach relies on a semi-greedy search algorithm to test which of the mismatches between target and draft model should be corrected to preserve quality, and which ones may be skipped. We then train a lightweight classifier based on existing LLM embeddings to predict, at inference time, which mismatching tokens can be safely accepted without compromising the final answer quality. We test our approach with Llama 3.2 1B (draft) and Llama 3.1 8B (target) models on zero-shot GSM8K reasoning, where it achieves up to 1.5x more accepted tokens per verification cycle with under 1% degradation in answer accuracy compared to standard speculative decoding and over 2x with small loss in accuracy. When applied to the LiveCodeBench benchmark, our approach automatically detects other, programming-specific important tokens and shows similar

speedups, demonstrating its ability to generalize across tasks.



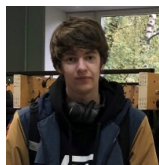
Statistical inference for Linear Stochastic Approximation with Markovian Noise

► М.н.с. Марина Евгеньевна Шешукова^{[2](#)}, департамент больших данных и информационного поиска, лаборатория стохастических алгоритмов и анализа многомерных данных



Анализ применимости метрик для оценки архитектур микросервисных систем

► М.н.с. Алексей Андреевич Данилов^{[2](#)}, департамент программной инженерии, лаборатория облачных и мобильных технологий

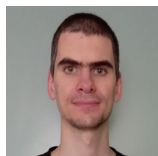


Discovering Enhanced Process Models from Software Event Logs

► Асп. Евгений Вадимович Степанов^{[2](#)}, департамент программной инженерии, лаборатория процессно-ориентированных информационных систем

↳ Nowadays software becomes more and more complicated, partially because of a complex nature of business processes which software represents, partially because of large programming platforms runtimes, which provide a lot of features for software developers, such as garbage collection, just-in-time compilation, exception handling system, assemblies loading, etc. All of these factors lead to complicated process happening during program execution. Those processes can be roughly classified into two types: business processes (which represent a business part of the application) and internal software processes representing technical aspect of the target programming platform. For understanding those processes a special sub-field of computer science emerged: process mining. Process mining allows discovering process models from event logs, enhancing discovered models with additional data from the event log, checking conformance of the discovered model with respect to another model or

original event log. This talk presents a novel framework for software process mining based on .NET event logs. The proposed framework includes the following process mining aspects: event-log abstraction, discovering process model and enhancing it with various data available in .NET event log. Despite the fact that we chose .NET as a target programming platform for evaluation of the proposed framework, our approach can be easily extended to any other programming platform.



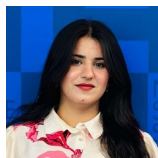
Оценки сложности реализации системы из трёх мономов схемами композиции

► Пригл. преп. Сергей Александрович Корнеев²,
департамент больших данных и информационного
поиска

↳ Под классической моделью вычисления будем понимать схемы из функциональных элементов, реализующих операцию умножения, где на входы подаются переменные, а на выходах вычисляются мономы. Рассмотрим другую модель - схемы композиции - в которой вместо операции умножения используется операция композиции. Композицией мономов U и V относительно монома R , все степени которого не превосходят соответствующих степеней мономов U и V , называется моном UV/R . При этом моном R может быть нулевым (с нулевым набором степеней) и может быть, вообще говоря, своим для каждой используемой операции композиции.

Задача об исследовании разных свойств, в том числе и связанных с вопросами сложности, известной в алгебре операции композиции, возникает естественным образом. Схемы композиции являются обобщением классической модели вычисления систем мономов, так как при $R=1$ операция композиции работает как обычное умножение. В процессе исследования этой модели выяснилось, что в ней имеют место интересные эффекты, отсутствующие, или, по крайней мере, неизвестные в близких вычислительных моделях.

В докладе планируется рассказать об известных результатах, касающихся вычислительной модели схем композиции, и, в том числе, представить новый результат в этой области: верхнюю и нижнюю оценки сложности реализации системы из трёх мономов от произвольного числа переменных q , которые отличаются на слагаемое q . В случае $q=3$ получена более точная верхняя оценка, отличающаяся от нижней лишь на единицу.



Building a Clean Bartangi Language Corpus and Training Word Embeddings for Low-Resource Language Modeling

► Ph.D. student Warda (Асп. Варда)², лаборатория моделирования и управления сложными системами

↳ We showcase a comprehensive end-to-end pipeline for creating a superior Bartangi language corpus and using it for training word embeddings. The critically low-resource Pamiri language of Bartangi, which is spoken in Tajikistan, has difficulties such as morphological complexity, orthographic variety, and a lack of data. In order to overcome these obstacles, we gathered a raw corpus of roughly 6,550 phrases, used the Uniparser-Morph-Bartangi morphological analyzer for linguistically accurate lemmatization, and implemented a thorough cleaning procedure to eliminate noise and ensure proper tokenization. The lemmatized corpus that results greatly lowers word sparsity and raises the standard of linguistic analysis. The processed corpus was then used to train two different Word2Vec models, Skipgram and CBOW, with a vector size of 100, a context window of 5, and a minimum frequency threshold of 1. The resultant word embeddings were displayed using dimensionality reduction techniques like PCA (Pearson, 1901) and t-SNE (van der Maaten and Hinton, 2008), and assessed using intrinsic methods like nearest-neighbor similarity tests. Our tests show that even from tiny datasets, meaningful semantic representations can be obtained by combining informed morphological analysis with clean preprocessing. One of the earliest computational datasets for Bartangi, this resource serves as a vital basis for upcoming NLP tasks, such as language modeling, semantic analysis, and low-resource machine translation. To promote more research in Pamiri and other under-represented languages, we make the corpus, lemmatizer pipeline, and trained embeddings publicly available.

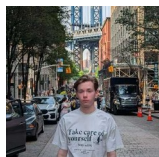


GENLINK: Closely Related Population Group Inference with GNNs and IBD Graphs

► Стажёр Алексей Валерьевич Шмелев², лаборатория статистической и вычислительной геномики

↳ Graph Neural Networks (GNNs) have recently shown significant effectiveness in analyzing structured graph data across diverse domains. At the same time, accurate

inference of ancestry from genetic data, especially among genetically similar populations, remains challenging due to the high dimensionality of SNP data and the internal complexity of the genetic relationships. To address these challenges, we propose a novel GNN-based method for genetic ancestry inference from identity-by-descent (IBD) graphs constructed using microarray genotype data, predicting only the most probable ancestral population for each individual. By identifying shared IBD segments between individuals, we construct graphs in which nodes represent individuals, and edges indicate meaningful genetic similarity exceeding a predefined threshold. In this context, the ancestry inference task is formalised as node classification on graphs initially lacking inherent node features. To enable effective generalization without the need for retraining the model on each new individual, we introduce a straightforward feature generation method derived solely from the graph's topology. Experimental evaluation demonstrates that our method outperforms heuristic classifiers and traditional community detection algorithms. Finally, performance further improves in realistic scenarios by enriching the graph with unlabeled nodes - individuals with undefined population - as message-passing in GNNs effectively propagates ancestry-related information throughout the network.

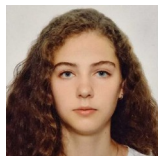


Hogwild! Inference: Parallel LLM Generation via Concurrent Attention

► Стажёр Роман Исмагилович Гарипов², департамент больших данных и информационного поиска, лаборатория компании Яндекс

↳ Large Language Models (LLMs) have demonstrated the ability to tackle increasingly complex tasks through advanced reasoning, long-form content generation, and tool use. Solving these tasks often involves long inference-time computations. In human problem solving, a common strategy to expedite work is collaboration: by dividing the problem into sub-tasks, exploring different strategies concurrently, etc. Recent research has shown that LLMs can also operate in parallel by implementing explicit cooperation frameworks, such as voting mechanisms or the explicit creation of independent sub-tasks that can be executed in parallel. However, each of these frameworks may not be suitable for all types of tasks, which can hinder their applicability. In this work, we propose a different design approach: we run LLM "workers" in parallel, allowing them to synchronize via a concurrently-updated attention cache and prompt these workers to decide how best to collaborate. Our approach allows the LLM instances to come up with their own collaboration strategy for the problem at hand, all the while "seeing" each other's memory in the concurrent

KV cache. We implement this approach via Hogwild! Inference: a parallel LLM inference engine where multiple instances of the same LLM run in parallel with the same attention cache, with "instant" access to each other's memory. Hogwild! Inference takes advantage of Rotary Position Embeddings (RoPE) to avoid recomputation while improving parallel hardware utilization. We find that modern reasoning-capable LLMs can perform inference with shared Key-Value cache out of the box, without additional fine-tuning.



Ускорение итерации Ньютона-Шульца для ортогонализации с помощью полиномов Чебышева

► Стажёр Екатерина Романовна Гришина²,
лаборатория матричных и тензорных методов в
машинном обучении

↪ Проблема вычисления ближайшей ортогональной аппроксимации к заданной матрице в последнее время привлекает большой интерес в машинном обучении. Среди известных приложений - оптимизатор Мюон или риманова оптимизация на многообразии Штифеля (многообразии ортогональных матриц). Среди существующих подходов метод итераций Ньютона-Шульца оказывается особенно эффективным решением, поскольку он основан исключительно на матричном умножении и, таким образом, обеспечивает высокую эффективность вычислений на графическом процессоре. Несмотря на свою эффективность, метод имеет недостаток — его коэффициенты фиксированы и, следовательно, не оптимизированы для заданной матрицы. В данной работе мы предлагаем решение этой проблемы — ускоренную с помощью полиномов Чебышева версию Ньютона-Шульца. Основываясь на теореме Чебышева, мы теоретически выводим оптимальные коэффициенты для итерации Ньютона-Шульца 3-го порядка и применяем алгоритм Ремеза для вычисления оптимальных многочленов более высокой степени. Мы используем эти полиномы для построения алгоритмов ортогонализации с определенными свойствами, которые полезны в глубинном обучении. Мы демонстрируем эффективность метода в двух приложениях: ортогонализации в оптимизаторе Мюон и построении эффективной ретракции для римановой оптимизации на многообразии Штифеля.



Improved rates of convergence and log-density Hessian estimation for implicit and denoising score matching

► Стажёр Анна Александровна Маркович^{[2](#)},
лаборатория теоретических основ моделей
искусственного интеллекта

↳ We provide new theoretical guarantees for score matching in high-dimensional settings. For both implicit and denoising score matching, we derive non-asymptotic bounds on the error of log-density Hessian estimation. For implicit score matching, we further establish a faster convergence rate. Our results extend score matching theory under the relaxed manifold assumption with smooth neural network estimators, offering the first guarantees for higher-order derivative accuracy in this setting.



Covariance Estimation using Tensor-train Decomposition

► Стажёр Артём Николаевич Потарусов^{[2](#)},
лаборатория теоретических основ моделей
искусственного интеллекта

↳ Given a sample of i.i.d. high-dimensional centered random vectors, we aim to estimate their covariance matrix $\Sigma = \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]$. Classical covariance estimation via the sample covariance suffers from the curse of dimensionality when d is large: the error rates depend heavily on the ambient dimension, making it unreliable in high-dimensional or tensor-valued data settings. To address this, we assume that Σ admits a structured decomposition combining Kronecker products and a tensor-train (TT) format. Specifically, we model $\Sigma = \sum_{j=1}^J \sum_{k=1}^K U_j \otimes V_{jk} \otimes W_k$,

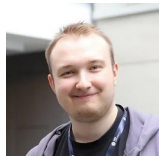


Программная система адаптивных подсказок «Мудрец» для обучения программированию в виде плагина для SmartLMS

► Стажёр Екатерина Андреевна Караваева^{[2](#)},
департамент программной инженерии, лаборатория
облачных и мобильных технологий

Итеративная генерация и отладка Java-кода с использованием больших языковых моделей

► Леон Андреевич Кулигин, лаборатория облачных и мобильных технологий



О развитии проектирования и реализации имитационного моделирования различных подсистем компьютерных игр в жанре глобальная стратегия

► Андрей Николаевич Штанов², департамент программной инженерии, проектная группа «Программная инженерия компьютерных игр»

↳ Стратегии, а также такие их поджанры, как глобальные стратегии и стратегии типа 4X (explore, expand, exploit, and exterminate) часто содержат в себе модели имитационного моделирования различной сложности: микро- и макроэкономика, дипломатия, боевые действия, шпионские сети и прочие. Сложна не только разработка отдельных моделей, но и их интеграция и согласование в рамках единой игры. Более того, для улучшения опыта игрока и обратной связи, современные игры требуют более комплексный подход к настройке характеристик (баланса), иммерсивности и прозрачности процессов. Поэтому анализ и развитие имитационного моделирования для игр остаётся актуальной задачей и дальнейшее исследование открывает перспективы как для развития жанра стратегий, так и для общей теории имитационного моделирования. В докладе формулируется универсальный подход к интеграции и оптимизации сложных имитационных подсистем, который может стать основой для нового поколения глобальных стратегий.

О возможностях, сложности и интерпретируемости моделей машинного обучения, лежащих в основе игровых ботов

► Арсений Сергеевич Вараксин, департамент программной инженерии, проектная группа «Программная инженерия компьютерных игр»

↳ С развитием методов машинного обучения для их сравнения, тестирования и автоматического обучения, в том числе путем соревнований агентов, были

придуманы “песочницы”.

Первые версии песочниц были оторваны от реальности. Затем взгляд разработчика обратился на компьютерные игры как на хорошо контролируемые виртуальные миры и интерпретируемые среды. Игры вошли в библиотеки Arcade Learning Environment, ViZDoom и другие.

В таком понимании песочницы прошли путь от классических пошаговых настольных игр до Atari (Arcade Learning Environment), ретро-игр (Stable-Retro, ViZDoom/Doom) и 3D-игр в реальном времени (MineRL/Minecraft).

Далее, ИИ и агенты обучения с подкреплением получили возможность играть на одном уровне с человеком, что привело к концепции создания метаверс (metaverse) как среды соревнования или взаимодействия ИИ-агентов и людей.

Одной из целей при улучшении песочниц является создание высокоиммерсивных сред, в которых люди и ИИ-агенты смогут взаимодействовать по тем каналам, которые привычны человеку.

Задача до конца не решена, но мы можем рассмотреть возможности использования различных алгоритмов обучения с подкреплением для обучения ИИ-агентов в данном контексте:

Deep Q-Network, алгоритм, явно приближающий “полезность” каждого действия, что не всегда получается в жизни и в непрерывных средах;

Proximal Policy Optimization (PPO), on-policy policy-gradient метод, быстро обучаемый алгоритм-бейзлайн для различных сред;

DREAMER, Discrete Codebook World Models for Continuous Control, алгоритм, строящий модель мира для задачи непрерывного управления.

Например, алгоритм PPO можно применить к модели OpenVLA, управляющей физическим или виртуальным роботом, для повышения качества.

В целом, использование и улучшение методов обучения с подкреплением должно привести к реализации ИИ-агентов, способных вести себя схоже с человеком в высокоиммерсивных средах, эффективно используя понятные для человека каналы взаимодействия.

Об уровнях представления и вовлечения человеческих и компьютерных агентов в высокоиммерсивных интерактивных окружениях

► Елизавета Петровна Егорова¹, департамент программной инженерии, проектная группа «Программная инженерия компьютерных игр»

↪ Исследование производится на стыке областей искусственного интеллекта (ИИ) и человеко-машинного взаимодействия, где особое внимание уделяется проблеме совместимости агентов в высокоиммерсивных интерактивных средах. Первая гипотеза заключается в том, что возможно достичь компромисса между уменьшением сложности компонентов восприятия ИИ-агентом и уровнем качества его представлений об окружающем мире. Предполагается, что снижение качества и сложности данных, подаваемых в модель, не обязательно приведет к ухудшению восприятия среды пользователем и может обеспечить более экономное использование вычислительных ресурсов. Вторая гипотеза связана с феноменом неразличимости партнёра по взаимодействию: модель ИИ может не иметь встроенного различия между игрой с человеком и игрой с другим ИИ-агентом. Такая ситуация позволяет исследовать ассиметричные варианты взаимодействия агентов и выявлять факторы, определяющие сопоставимость опыта между разными типами агентов. Анализ уровней представления и вовлечения показывает, что человек опирается на богатое сенсорное погружение, тогда как ИИ работает с формализованными описаниями среды, для их согласования необходим мета-уровень управления, обеспечивающий баланс условий и контроль ресурсов. Ожидаемые результаты включают разработку архитектурных принципов, позволяющих создавать среды, где агенты взаимодействуют на равных при сохранении различий в их когнитивных моделях. Областью приложения результатов являются компьютерные игры, симуляторы и иммерсивных программ учебного назначения, где совместное участие человека и ИИ открывает возможности для кооперации, соревнований и совместного (коллаборативного) обучения.

Постквантовые криптографические алгоритмы цифровой подписи: сравнительный анализ и главные проблемы

► Руслан Александрович Качмазов

↪ После открытия алгоритма Шора стало понятно, что большинство из используемых сегодня ассиметричных криптографических алгоритмов могут быть

скомпрометированы при появлении достаточно мощного квантового компьютера, так как они основаны на проблеме дискретного логарифмирования и факторизации чисел, например, RSA и ECDSA. Учитывая то, что сфера квантовых вычислений активно развивается, криптографическое сообщество начало думать о внедрении новых алгоритмов, способных противостоять такого рода атакам - эта сфера была названа постквантовой криптографией. В основном, постквантовые алгоритмы строятся на задачах теории целочисленных решеток, кодах исправления ошибок, хэш-функций, систем многочленов от многих переменных, изогении эллиптических кривых, теории кос и других.

В 2016 году NIST учредил конкурс, в котором отбирал алгоритмы инкапсуляции ключа (KEM) и цифровой подписи, результатом которого стали 3 финалиста в рамках цифровой подписи: CRYSTALS-Dilithium (решетки), SPHINCS+ (хэш-функции) и Falcon (решетки). Отечественная криптография не осталась в стороне, и в рамках ТК 26 была создана рабочая группа 2.5, занимающаяся постквантовой криптографией. Представлены алгоритмы Гиперикум (хэш-функции) и Шиповник (коды ошибок).

В работе описан теоретический разбор криптостойкости пяти постквантовых алгоритмов цифровой подписи, описанных выше. Проведен сравнительный анализ по скорости, размеру ключей и подписи, сложности программной реализации отдельно взятого алгоритма. Имплементированы модификации отечественных алгоритмов, а также рассмотрен потенциальный процесс перехода на постквантовую криптографию, затрагивающий гибридное шифрование и множество этапов адаптации систем.

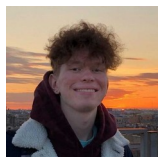
Balanced Language Sampling for Multilingual Models

► Вероника Викторовна Царева

↪ In recent years, with the release of multilingual models, it has become popular to conduct research not only on a single language but on several. However, researchers report the results often based on small or unbalanced data. This raises a question of whether such “good” models’ performances can be trusted. The main goal of my research was to examine the difference between the results that the multilingual model produces across the different language samples. In order to do this, I chose the article by Papadimitriou et al. (2021) where the subjecthood in different alignment systems is studied using the multilingual BERT by Devlin et al. (2018) and a “typologically diverse” 24-language sample. Based on the variant of the Genus-Macroarea method, a sampling method by Miestamo et al. (2016) that is used in

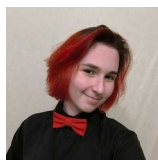
linguistic typology, I made a more balanced 24-language sample. The experimental setup was maintained as in Papadimitriou et al. (2021) for the comparability of the results. The data on which the model was trained and tested is the Universal Dependencies Treebank v2.15 (Nivre et al. 2020).

The main results from replicating the experiments with my own data are as follows. First, the best performance in the argument classification task was shown in the middle-last layers both on the original and on the balanced sample, which correlates with the NLP findings about the syntactic knowledge and BERT. Second, I disagree with the conclusion of Papadimitriou et al. (2021) as the animacy cannot be a strong predictor as they claim. Third, there are alignment systems that can be transferred across the languages and affect the results of the argument classification task. Finally, the mBERT performs with the tasks well even on the languages that it was not pretrained on. An overall conclusion that I propose is that a balanced language sample is indeed an important part of the multilingual studies as the difference between the results is statistically significant. A more diverse sample shows more diverse results and interpretations and they should be studied in detail.



MiAD (Mirage Atom Diffusion) - новая state-of-the-art модель для поиска стабильных кристаллических материалов

► Стажёр Андрей Сергеевич Охотин², центр глубинного обучения и байесовских методов



COALA: Numerically Stable and Efficient Framework for Context-Aware Low-Rank Approximation

► Стажёр Ульяна Романовна Паркина², лаборатория матричных и тензорных методов в машинном обучении



Рандомизированная оценка норм два-в-бесконечность и один-в-два

► Стажёр Аскар Шамилевич Цыганов^{[2](#)}, лаборатория стохастических алгоритмов и анализа многомерных данных

→ Лучшие постеры

Лучшие постеры будут выбраны в результате анонимного голосования участников конференции. Каждый участник имеет возможность проголосовать только за один стендовый доклад.

Демонстрация технологических решений

В рамках конференции также представлены технологические решения проектных групп «Программная инженерия компьютерных игр» (руководитель доц. Ольга Вениаминовна Максименкова), «Информационные системы в медицине» (руководитель доц. Дмитрий Игоревич Рябцев [2](#)), а также лаборатории облачный и мобильных технологий (руководитель Дмитрий Владимирович Александров) департамента программной инженерии ФКН.

Полезная информация

Связь с организаторами: psu_cs@hse.ru

Место проведения

Конференция проходит в учебном центре «Вороново» НИУ ВШЭ²



Адрес: 108 830, г. Москва, поселение Вороновское, с. Вороново,
ул. Канторовича, домовладение 1

Схема проезда

Общественным транспортом

Чтобы самостоятельно добраться до Учебного Центра вы можете воспользоваться маршрутным транспортом:

от м. Тёплый Стан (выход из последнего вагона из центра, в подземном переходе — 2-й выход налево): **Автобус № 503** или **маршрутное такси № 887** (остановки находятся в 30-40 метрах от подземного выхода метро возле ТРЦ «Принц Плаза»; ориентир: справа от остановок должен быть фонтан-водопад). Время в пути: автобус — от 1 ч. до 1 ч. 15 мин., маршрутное такси — 50-60 мин. (без учёта пробок)

от м. Ольховая: Автобус № 508

Выйти на остановке **Вороново-2**, предварительно заказав микроавтобус из УЦ «Вороново».

Это можно сделать к определенному времени, или по прибытию на остановку, позвонив дежурному администратору УЦ "Вороново" по телефону +7 (495) 916-89-16, +7 (495) 772-95-90 доб. 17010.

Доставка микроавтобусом от автобусной остановки Вороново до УЦ по рабочим дням с 9.00 до 18.00.

На личном автомобиле

Координаты для GPS-навигатора:

широта: 55°18'54.13"N (55.315035)

долгота: 37°7'25.77"E (37.123825)

Вы едете со стороны МКАД по Калужскому шоссе в сторону области до села Вороново, поворачиваете на первом светофоре направо по дорожному указателю, двигаетесь 2,8 км до поворота налево по указателю.

Трек: <http://maps.yandex.ru/-/CVCgECOА>

Размещение

Набор участника: Каждый участник при прибытии в УЦ Вороново получает набор участника, в который входят карточка (бейдж) с именем, ручка и блокнот для записей, а также футболка НКФКНIII с уникальным дизайном.

Доступ в интернет: В УЦ Вороново имеется Wi-Fi. Данные для входа указаны в книгах постояльцев в каждом номере. Также работает корпоративная Wi-Fi сеть НИУ ВШЭ.

Еда: В ходе конференции организовано трёхразовое питание (завтрак, обед и ужин), а также перерыв на кофе в вечерней сессии докладов. Всё питание организовано в столовой УЦ Вороново. Таким образом, участникам не нужно думать про поиск еды. Имеет смысл иметь с собой ёмкость, чтобы набрать в неё питьевой воды. Системы выдачи воды присутствуют на каждом этаже УЦ.





Москва, 2025

Версия: 25 октября