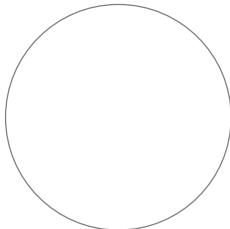




2024

Годовой отчёт



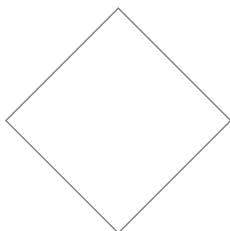
Природа

Постоянное движение и развитие,
закономерность и системность, законы
и ритмы, хаос и жизнь



Человек

Живой организм, часть природы,
способный мыслить, чувствовать
и говорить, несовершенный



Технологии

Творение человека и его продолжение,
логичное и предсказуемое,
утилитарное, создано помочь



AIRI

Единство природы, человека
и технологий

Содержание

Вступительное слово	4
Миссия AIRI	5
Ценности AIRI	6
Направления деятельности	7
Ключевые результаты	8
Руководители научных групп	10
Руководство Института	12
Научные результаты	13
Научные публикации	41
Мероприятия и выступления	83
О нас пишут и говорят	102
Партнерства и коллаборации	104
Ученые AIRI в социальных сетях	106
Контакты	108

Четвертый год жизни AIRI



Иван Оседец
Генеральный директор

В 2024 году наш Институт увеличился вдвое и насчитывает более 200 преданных своему делу специалистов, 9 из которых в течение года защитили кандидатские и докторские диссертации. Многие приехали к нам из-за границы, чтобы развивать науку об ИИ в России и мире с AIRI.

Вместе с тем, увеличилось и количество успешно проведенных исследований, 90 из которых было представлено на конференциях уровня A и A*. 17 статей — принято на ведущую международную конференцию по искусственному интеллекту NeurIPS.

Мы в AIRI продолжаем создавать универсальные системы ИИ для решения задач реального мира. В этому году мы сфокусировались на концепции прикладного AGI, в основе которой лежит научная обоснованность и польза для бизнеса. Совместно с индустриальными партнерами мы запустили ряд амбициозных проектов по генеративному проектированию и применению технологий искусственного интеллекта в медицине.

Появляются новые коллаборации с ведущими университетами и технологическими компаниями.

К нам присоединились партнеры из ГИАП, СберМедИИ, ИСП РАН, КАМАЗ Digital, Сибур Digital и многие другие. Обширная и вдохновляющая часть 2024 года — знакомство и сотрудничество с региональными научно-исследовательскими центрами. СВФУ, АГУ, ТГУ — всех не перечесть — демонстрируют, насколько разнообразные и сильные компетенции представлены в России. Я уверен, что благодаря стремлению AIRI соединять между собой коллег из самых разных точек мира и предметных областей, в будущем нас ждет очень много непредсказуемых, но приятных результатов.

Я выражаю благодарность сотрудникам AIRI и нашим партнерам за усердную работу. Вместе мы не только приближаем появление AGI, но и участвуем в подготовке нового поколения исследователей и вдохновляем друг друга на открытия.

Приглашаю вас узнать больше о результатах работы Института AIRI в 2024 году из этого отчета.

Миссия AIRI

Создание универсальных систем
искусственного интеллекта,
решающих задачи реального мира

Основная цель Института —
найти возможности применения
искусственного интеллекта для решения
сложных социальных, экономических
и научных задач. Научные сотрудники
Института занимаются исследованиями
в фундаментальных и прикладных
областях. В работе над своими
проектами они стремятся получить
прорывные результаты в области
искусственного интеллекта и его
приложений, участвуя в формировании
глобальной исследовательской
повестки.

Ценности AIRI



Человекоцентричность



Свобода научного выражения



Ответственность
и вклад в будущее



Открытость
и прозрачность



Партнерство
и коллаборация



Направления деятельности



Научные исследования

Проведение прорывных исследований в области искусственного интеллекта, формирование глобального центра экспертизы



Вклад в развитие AI

Участие в глобальном развитии искусственного интеллекта через создание, развитие и поддержку проектов с открытым кодом



Научно-технические партнерства

Развитие партнерств с научными организациями, промышленностью и государством, разработка и коммерциализация технологий в области искусственного интеллекта



Лаборатории

Сотрудничество с институтами, университетами и индустриальными партнерами по запуску совместных научно-исследовательских лабораторий в области искусственного интеллекта



Популяризация AI

Проведение профильных конференций и мероприятий, создание и поддержка соревнований, популяризация искусственного интеллекта в целом

Ключевые результаты

221

публикаций

65

статей
на конференциях А*

25

статей
на конференциях А

56

статей в журналах
с рейтингом Q1



Руководители научных групп



Андрей Кузнецов
Лаборатория FusionBrain



Владислав
Шахуро
ИИ для роботов



Айбек
Аланов
Контролируемый
Генеративный ИИ



Константин
Соболев
Генеративный ИИ
для видео



Александр Панов
Лаборатория когнитивных систем искусственного интеллекта



Алексей
Ковалёв
Воплощенные агенты



Юрий
Куратов
Модели с памятью



Алексей
Скрынник
RL агенты



Дмитрий Дылов
Лаборатория «Сильный AI в медицине»



Дмитрий
Умеренков
Базовые модели



Ярослав
Беспалов
Мультиомодальные
архитектуры ИИ



Назар
Бузун
Обучение
представлений



Дмитрий
Крюков
Исследования
биомаркеров

Руководители научных групп



Иван Оседец

Вычислительный интеллект



Евгений Фролов

Технологии персонализации



Евгений Бурнаев

Обучаемый интеллект



Александр Горбань

Лаборатория ИИ, анализа данных и моделирования



Антон Конушин

Пространственный интеллект



Олег Рогов

Доверенные и безопасные интеллектуальные системы



Егор Ершов

Цветовая вычислительная фотография



Семен Буденный

Дизайн новых материалов



Ольга Кардымон

Биоинформатика



Елена Тутубалина

Прикладное NLP



Марина Мунхоева

Самообучение и обучение представлений



Алексей Осадчий

Нейроинтерфейсы



Илья Макаров

ИИ в индустрии



Артур Кадурин

Глубокое обучение в науках о жизни



Александр Панченко

Вычислительная семантика



Владислав Куренков

Адаптивные агенты



Александр Тюрин

Методы оптимизации в машинном обучении



Артем Шелманов

Обучение на слабо размеченных данных

Руководители функциональных направлений



Иван Оседец

Генеральный директор



Максим Кузнецов

Директор управления стратегического развития и партнерств



Ольга Суровегина

Директор по научно-техническим партнерствам



Степан Мамонтов

Руководитель отдела научно-технической разработки



Ольга Попова

Руководитель отдела проектного управления



Антон Ризаев

Финансовый директор



Екатерина Сарафанова

Главный бухгалтер



Николь-Мария Кук

Руководитель отдела закупок



Александра Брайтман

Директор по маркетингу и коммуникациям



Мария Звонарева

HR Директор



Юлия Никитина

Директор по правовому обеспечению



Константин Катанов

IT Директор

Научные
результаты

На пути к прикладному AGI

AGI =

Инженерный подход

- Научная обоснованность: Публикации на A/A* и Q1
- Трансфер идей в продуктовый пайплайн / тестирование в бизнесе

+

Польза для текущего бизнеса

- Функциональный заказчик на проведение научного исследования
- Проведенный на стороне заказчика предварительный анализ бизнес-эффекта и prerequisite research

Направления прикладного AGI

Compute and Data	Multimodality	Agency	Embodiment	World Model
Effective training & inference	Modality	Self-learning	Hardware	Knowledge
Effective data storage	Omnimodality	Self-reflection/Reasoning	Software	Forecasting
Synthetic data	Perception Augmentation	Self-alignment	Interface	Causality
		Goal-setting & planning		
		Multiagency		



Главный
результат



Дмитрий Дылов

Директор лаборатории «Сильный AI в медицине»

AGI Med

«Цифровой Помощник 2.0» и «Муншот», предназначенные для первичного приема пациентов и для содействия врачам лучевой диагностики

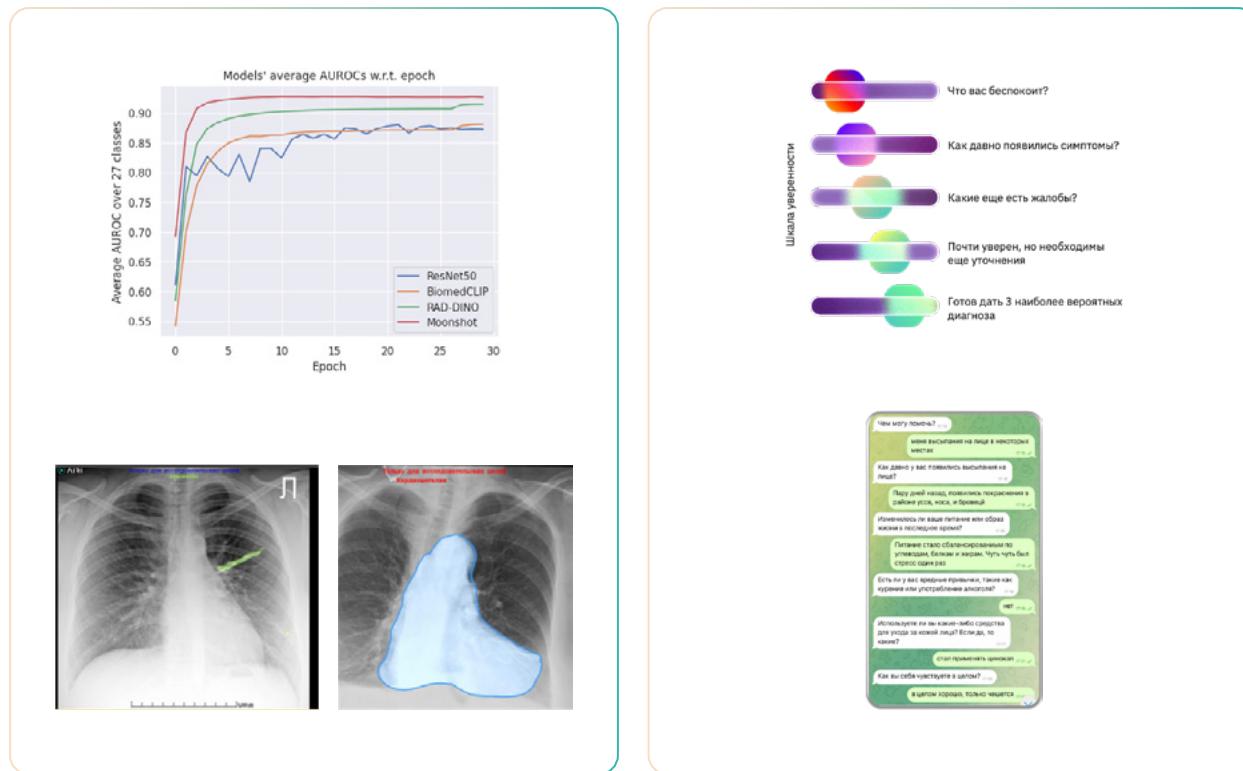
В мобильном приложении медицинской компании «СберЗдоровье» (входит в индустрию здоровья Сбербанка) появился ИИ-помощник, который поможет россиянам заботиться о здоровье. Технологию разработали ученые Института AIRI и компания «СберМедИИ» на базе нейросетевой модели Сбербанка GigaChat.

- Постановка предварительного диагноза и маршрутизация к целевому специалисту
- Расшифровка анализов и заключения врача
- Реализован набор интеллектуальных фильтров безопасности

В то же время, базовая модель «Муншот» показала себя как мощная модель анализа радиологических снимков. Будучи обучена на крупном датасете Москвы, модель научилась генерировать разделы «Найдено» и «Заключение» в радиологических отчетах по рентгенографии грудной

клетки, предсказывать класс патологии и автоматически контурировать их. Этот метод объединяет визуальный энкодер (контрастное обучение), обученный на медицинских изображениях, и специализированную биомедицинскую языковую модель. Подход эффективно преодолевает сложности, связанные с медицинской семантикой, и учитывает разнообразие находок в изображениях, обеспечивая автоматизацию процесса составления отчетов.

Модель обучена на крупных наборах данных, таких как PadChest, BIMCV-COVID19, CheXpert, OpenI и MIMIC-CXR, что обеспечивает ее устойчивость и адаптацию к разнообразным клиническим случаям. Оценочные метрики подчеркивают точность и способность модели учитывать тонкие медицинские детали, что демонстрирует большой потенциал для оптимизации работы радиологов. Это достижение не только автоматизирует важный



этап диагностики, но и способствует внедрению передовых ИИ-решений в медицинскую практику.

В рамках главной конференции по мед ИИ (MICCAI 2024), одним из ключевых достижений стало создание нового подхода к интерактивной сегментации сложных удлиненных объектов, таких как провода, катетеры или вены, без использования традиционных масок. Разработанная структура данных ориентируется на набор центральных линий контурируемых объектов и работает особенно эффективно. Также значительные усилия были направлены на создание новых подходов к анализу медицинских изображений. В частности, предложены анатомические позиционные эмбеддинги для 3D-изображений, способные точно предсказывать анатомическое положение областей, что открывает новые возможности для локализации органов и патологий. Эти результаты были представлены на MICCAI 2024, 2 из 3 российских

статьей на конференции были сделаны сотрудниками лаборатории AGI Med.

Команда AIRI-AGI-med приняла участие в двух международных соревнованиях: RRG24 на BioNLP и MIDRC XAI Challenge, где заняла второе и пятое место, соответственно, используя первое поколение своей модели Муншот. Это стало возможным благодаря разработке методов классификации затемнений и сегментации легких на датасете с ограниченной разметкой. На основе CLIP-модели удалось выиграть приз в \$5000, а также показать отличные результаты в задачах классификации и сегментации. Команда уступила только ведущим IT-гигантам, будучи на порядок меньше по составу.



Крупные
результаты

Разработан метод SparseGrad для параметро-эффективного обучения больших языковых моделей на базе тензорного разложения HOSVD



Александр Панченко

Ведущий научный сотрудник

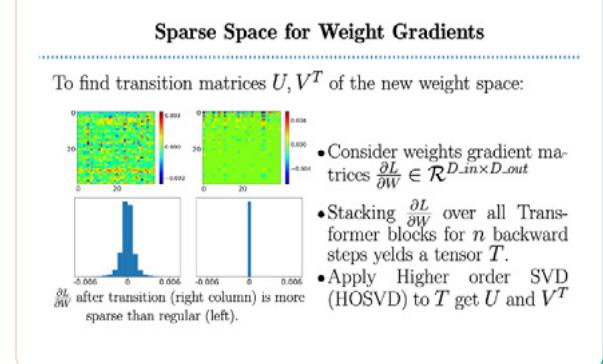


Иван Оседец

Генеральный директор

В методе SparseGrad с помощью HOSVD-декомпозиции мы находим пространство, где матрица весов линейных слоев сильно разрежена (остается ~1% от всех параметров). Создав новый класс линейного слоя и переписав автоград с учетом перехода в это пространство, авторы уменьшили количество параметров, используемых при тренировке Трансформеров, и, следовательно, используемую память. При одинаковом потреблении памяти на энкодерных архитектурах (BERT, RoBERTa) и декодерных архитектурах (LLaMa2-7B) модель показывает результаты лучше LoRA и MeProp — бейзлайна, который отбирает наиболее значимые параметры в линейном слое и тренируют только их.

Тестирование подхода на LLaMa 2 7B показало, что наш метод дает лучший лосс по валидации, а также лучшую метрику на вопросно-ответном бенчмарке I-Bench. Качественная характеристика генерируемого текста подтвердила этот вывод.



SparseGrad: A Selective Method for Efficient Fine-tuning of MLP Layers
Viktoria Chekalina, Anna Rudenko, Gleb Mezentsev, Alexander Mikhalev, Alexander Panchenko, Ivan Oseledets
EMNLP'24, A*

Разработан метод AMORE для оценки надежности и интерпретируемости больших языковых моделей в химии, основанный на аугментации молекулярных структур



Елена Тутубалина
Ведущий научный сотрудник



Артур Кадурин
Инженер-исследователь



Иван Оседец
Генеральный директор



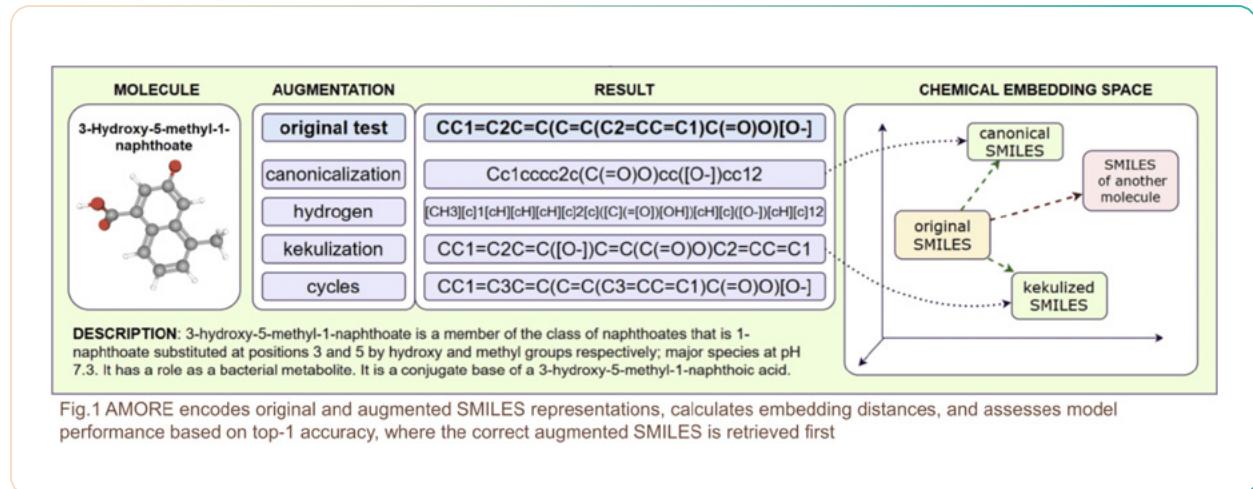
Андрей Кузнецов
Директор лаборатории



Денис Димитров
Научный консультант

Недавняя интеграция химии с обработкой естественного языка (NLP) продвинула вперед открытие лекарств. Сегодня представление молекул в языковых моделях (LM) имеет решающее значение для улучшения понимания химических процессов. Разработана система дополненного молекулярного поиска (Augmented MOlecular REtrieval, AMORE), гибкая zero-shot система оценки достоверности химических языковых моделей различной природы: обученных исключительно на молекулах для выполнения химических задач и на объединенном

корпусе текстов на естественном языке и структур, представленных в виде строк. Эта система опирается на модификации молекул, сохраняющие их химические свойства, такие как кекулизация и замена циклов. Метрика основана на оценке сходства между распределенными представлениями молекул и их модификациями.



A*: Ganeeva, V., Sakhovskiy, A., Khrabrov, K., Savchenko, A., Kadurin, A. & Tutubalina, E. Lost in Translation: Chemical Language Models and the Misunderstanding of Molecule Structures. In Findings of the Association for Computational Linguistics: EMNLP 2024

A: The Shape of Learning: Anisotropy and Intrinsic Dimensions in Transformer-Based Models”, Anton Razzhigaev, Matvey Mikhalkchuk, Elizaveta Goncharova, Ivan Oseledets, Denis Dimitrov, Andrey Kuznetsov, EACL 2024

Разработан фреймворк GOLF для активного обучения нейронных потенциалов, ориентированный на локальную оптимизацию лекарственных молекул



Александр Панов

Директор лаборатории



Елена Тутубалина

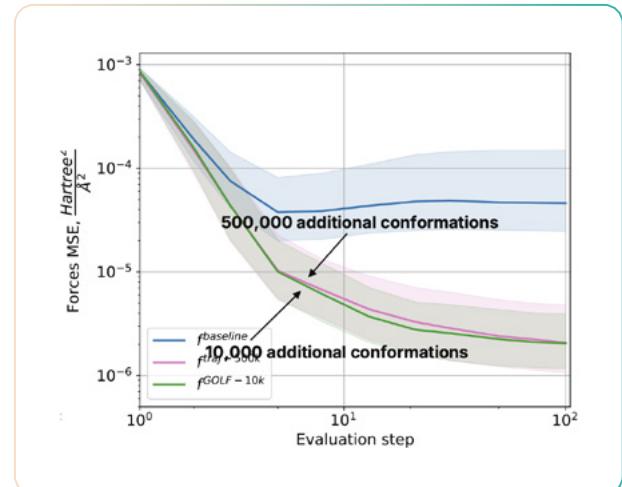
Ведущий научный сотрудник



Артур Кадурин

Инженер-исследователь

Разработанный подход GOLF позволил в 50 раз сократить объем данных, необходимый для дообучения нейронных сетей на задачу оптимизации конформаций (пространственных представлений молекул). Используем суррогатный оракул, чтобы отобрать конформации для последующей разметки с помощью DFT. Оцениваем, уменьшилась ли энергия. Если энергия уменьшилась, продолжаем оптимизацию или рассматриваем предсказание NNP как некорректное и добавляем предыдущую конформацию в обучающую выборку.



A*: Tsypin, A., Ugadiarov, L. A., Khrabrov, K., Telepov, A., Rumiantsev, E., Skrynnik, A., Panov, A., Vetrov, D., Tutubalina, E., Kadurin, A. Gradual Optimization Learning for Conformational Energy Minimization. In The Twelfth International Conference on Learning Representations (ICLR).

A*: Tsypin, A., Ugadiarov, L. A., Khrabrov, K., Telepov, A., Rumiantsev, E., Skrynnik, A., Panov, A., Vetrov, D., Tutubalina, E., Kadurin, A. Gradual Optimization Learning for Conformational Energy Minimization. In The Twelfth International Conference on Learning Representations (ICLR).

Исследование внутренних свойств моделей трансформеров (внутренняя размерность и анизотропия), характеризующих наличие линейной зависимости ряда соседних слоёв эмбеддингов



Иван Оседец
Генеральный директор

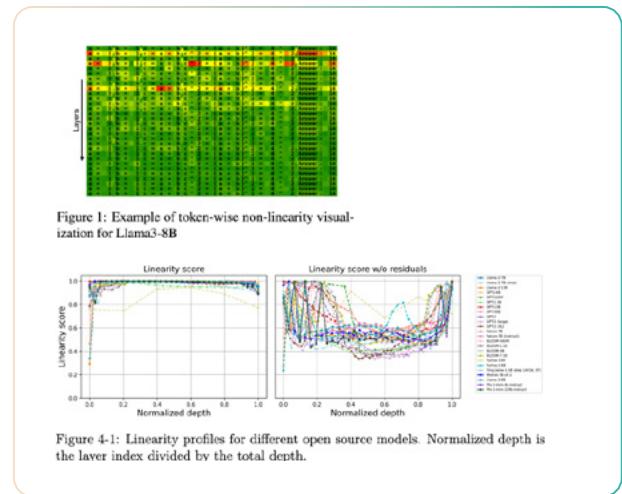


Денис Димитров
Научный консультант



Андрей Кузнецов
Директор лаборатории

В данной работе было продолжено исследование моделей-декодеров и выявлена линейная зависимость между последовательными слоями моделей, основанных на архитектуре Трансформер. Был предложен процесс прунинга моделей, который позволяет заменить часть слоев декодеров на линейное преобразование входного вектора. Замеры на бенчмарках показали, что замена 5-10% слоев на линейное преобразование не оказывает сильное влияние на результаты модели.



A*: Anton Razzhigaev, Matvey Mikhalkuk, Elizaveta Goncharova, Nikolai Gerasimenko, Ivan Oseledets, Denis Dimitrov, and Andrey Kuznetsov. 2024. Your Transformer is Secretly Linear. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)

Контролируемый генеративный ИИ



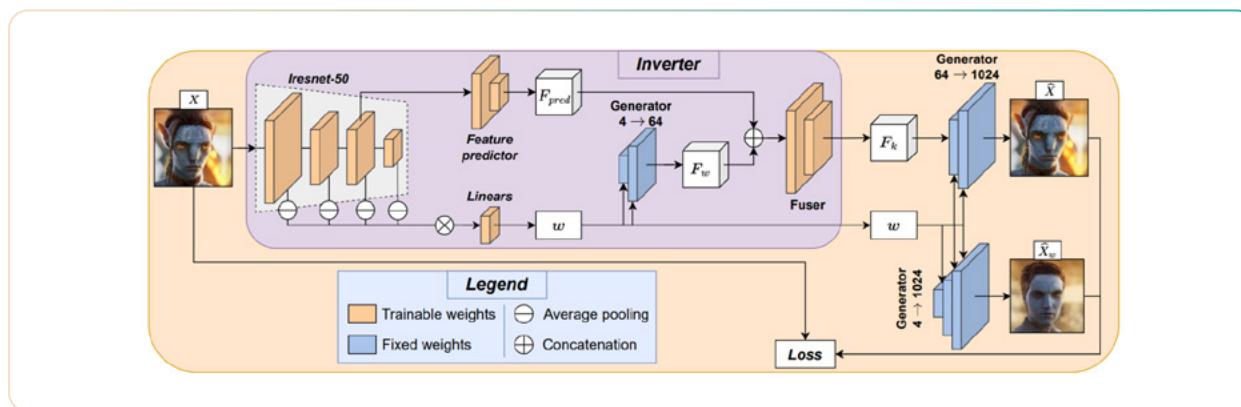
Айбек Аланов

Научный сотрудник

- Применение модели StyleGAN для манипуляции изображений. Мы разработали модель StyleFeatureEditor, которая позволяет редактировать изображения, при этом восстанавливая их с очень высокой точностью. Мы улучшили на тот момент sota подход в 4 раза по метрике LPIPS точность реконструкции¹.
- Разработка модели для переноса причесок для лиц людей. Мы предложили новую модель, которая показывает наилучшее качество переноса причесок

и реалистичность результата по сравнению с бейзлайнами. Скорость применения на порядке быстрее².

- Разработка эффективной параметризации для дообучения нейросетевых моделей. Был предложен метод, который позволяет более эффективно с точки зрения используемой памяти дообучать нейросетевые модели, в том числе большие диффузионные модели³.



¹ A*: "The Devil is in the Details: StyleFeatureEditor for Detail-Rich StyleGAN Inversion and High Quality Image Editing" на конференции CVPR 2024.

Авторы: Денис Бобков, Вадим Титов, Айбек Аланов, Дмитрий Ветров

² A*: "HairFastGAN: Realistic and Robust Hair Transfer with a Fast Encoder-Based Approach" на конференции NeurIPS 2024. Авторы: Максим Николаев, Михаил Кузнецов, Дмитрий Ветров, Айбек Аланов

³ A*: "Group and Shuffle: Efficient Structured Orthogonal Parametrization" на конференции NeurIPS 2024. Авторы: Михаил Горбунов, Николай Юдин, Вера Соболева, Айбек Аланов, Алексей Наумов, Максим Рахуба

Разработана оригинальная нейросимвольная архитектура управления воплощенными агентами в сложных динамических средах с поддержкой планирования на основе языковых моделей и двумя типами модели мира для эффективного обучения в среде



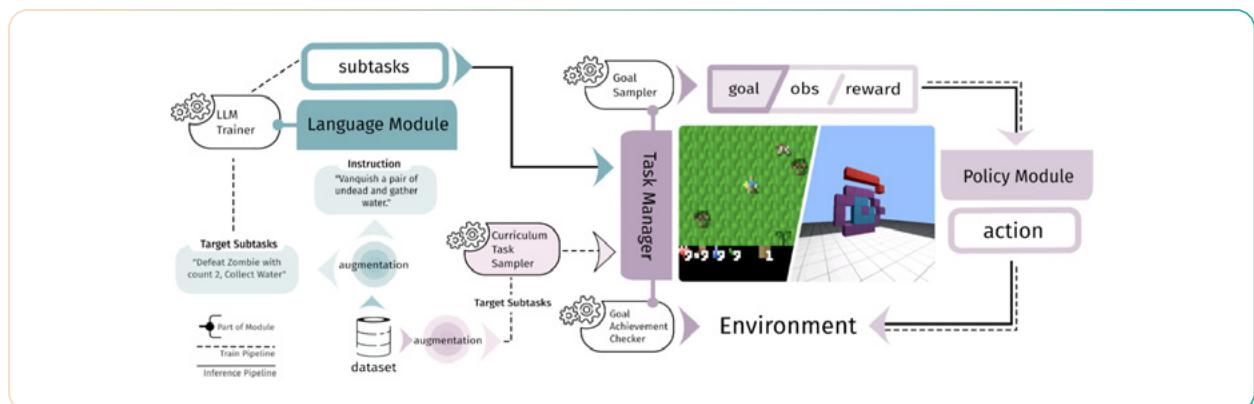
Александр Панов

Директор лаборатории

Результаты работы изложены в докторской диссертации, защищенной на диссовете МФТИ по специальности 1.2.1 «Искусственный интеллект и машинное обучение».

Подход IGOR для построения стратегии в мультимодальных средах. Метод использует большую языковую модель для преобразования текстовых

инструкций в высокоуровневый план, который затем воплощает RL-агента, выстраивая стратегию поведения в среде. Эффективность метода продемонстрирована в среде IGLU, где предложенный подход превзошел результаты первого места на соревновании NeurIPS, и в среде Crafter, опережая SOTA-решение Dynalang.



A: Volovikova, Z., Skrynnik, A., Kuderov, P. and Panov, A.I., 2024. Instruction Following with Goal-Conditioned Reinforcement Learning in Virtual Environments. In ECAI 2024 (pp. 650-657). IOS Press. [Core A]

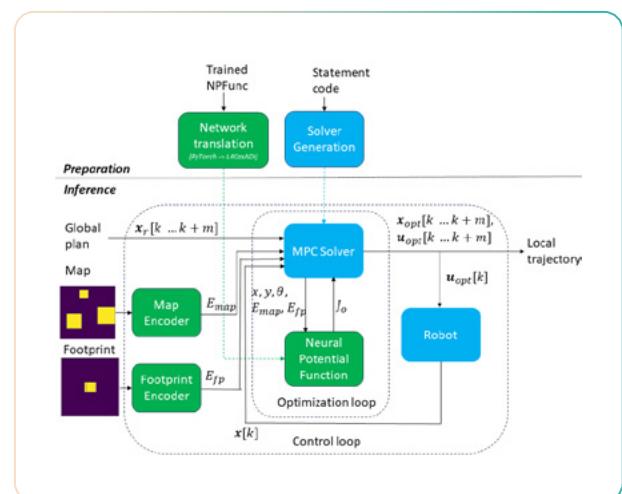
Разработана модель для задачи локального планирования с учетом динамических препятствий



Александр Панов

Директор лаборатории

Разработана модель, вычисляющая нейронный потенциал, используемый в MPC для задачи локального планирования пути с учетом динамических препятствий. Используется трансформерная модель на основе GPT. Построенная система управления была протестирована на реальном мобильном роботе и показала преимущество по сравнению со всеми мировыми аналогами. Работа была представлена на ведущей конференции по робототехнике ICRA 2024.



A*: Alhaddad, M., Mironov, K., Staroverov, A., Panov, A. Neural Potential Field for Obstacle-Aware Local Motion Planning, IEEE International Conference on Robotics and Automation (ICRA) 2024, pp. 9313–9320.

Разработан подход, объединяющий последовательности кадров LIDAR через алгоритмы регистрации облаков точек, достигающий SOTA-результатов в semi-supervised и улучшение в supervised сеттинге

CVPR 2024 (A*), IEEE Access (Q1), demo IJCAI 2024 (A*)

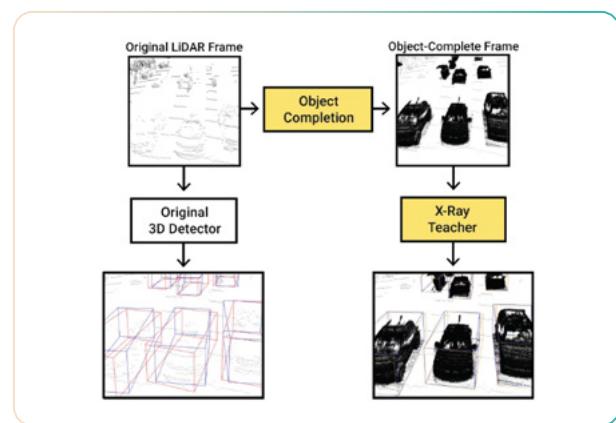


Илья Макаров

Ведущий научный сотрудник

Подход существенно улучшает точность и эффективность моделей 3D детекции, что критично для автономного транспорта, робототехники и других приложений, зависящих от высокоточного распознавания объектов.

Технология может быть использована в автономных транспортных системах, складской автоматизации и робототехнике.



Gambashidze, A., Dadukin, A., Golyadkin, M., Razzhivina, M., & Makarov, I. Weak-to-Strong 3D Object Detection with X-Ray Distillation. CVPR 2024

Исследование стратегических и этических решений LLM



Илья Макаров

Ведущий научный сотрудник

Разработан фреймворк для тестирования гипотез о соответствии эмоциональных реакций людей и решений LLM. Установлено, что эмоциональный алайнмент с людьми в большинстве случаев слабый, но эмоциональные промпты, особенно негативные, существенно влияют на поведение моделей, снижая кооперацию в играх и качество в этических задачах. NeurIPS (A*)

Понимание эмоционального воздействия на модели критично для разработки безопасных и надежных ИИ-систем в этических и стратегических приложениях.

Использование фреймворка возможно для оптимизации LLM в областях, где важно учитывать эмоциональный контекст, например, ассистентах или интерактивных системах.

Разработка эмпатичного LLM-ассистента:

Разработан подход InsideOut, где базовые эмоции по Экману участвуют в формулировании ответов. На базе GigaChat и GPT-3.5 достигнуто улучшение на 12–20% в распознавании эмоций и до 3.9% в качестве эмпатичных ответов.

ECAI demo (A)

Повышение эмпатии в LLM способствует созданию более человечных и эффективных ассистентов, полезных в психологии, обучении и клиентской поддержке. Подход может быть интегрирован в ассистентов для медицинских консультаций, обучения и сервисов поддержки клиентов.

Mozikov, M., Severin, N., Bodishtianu, V., Glushanina, M., Nasonov, I., Orekhov, D., Pekhotin, V., Makovetskiy, I., Baklashkin, M., Lavrentyev, V., Tsvigun, A., Turdakov, D., Shavrina, T., Savchenko, A., & Makarov, I. EAI: Emotional Decision-Making of LLMs in Strategic Games and Ethical Dilemmas. NeurIPS 2024

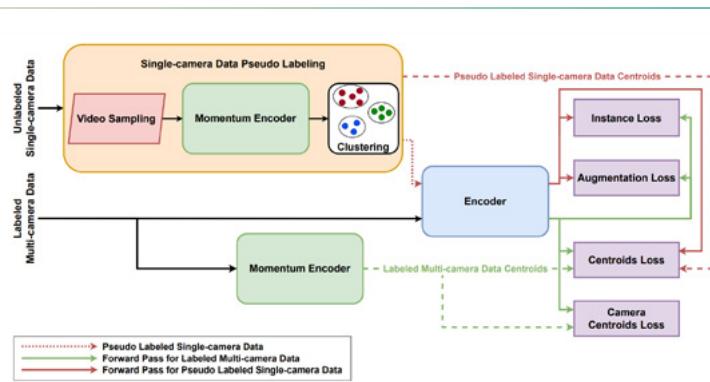
Deeb, B. M., Savchenko, A., & Makarov, I. (2024). CA-SER: Cross-Attention Feature Fusion for Speech Emotion Recognition. In ECAI 2024 (pp. 4479-4482). IOS Press. ECAI 2024 Demo.

Умное обучение на смеси данных для распознавания людей



Антон Конушин
Старший научный сотрудник

Разработаны методы обучения моделей реидентификации человека на смесях неразмеченных однокамерных и размеченных многокамерных данных, позволившие существенно повысить обобщающую способность моделей, валидированную тестированием в кросс-датасет сценариях.



Mamedov T., Konushin V., Konushin A. ReMix: Training Generalized Person Re-identification on a Mixture of Data // arXiv preprint arXiv:2410.21938. — 2024 (принята на WACV 2025)

Первая открытая среда и самый большой датасет в мире для In-Context Reinforcement Learning



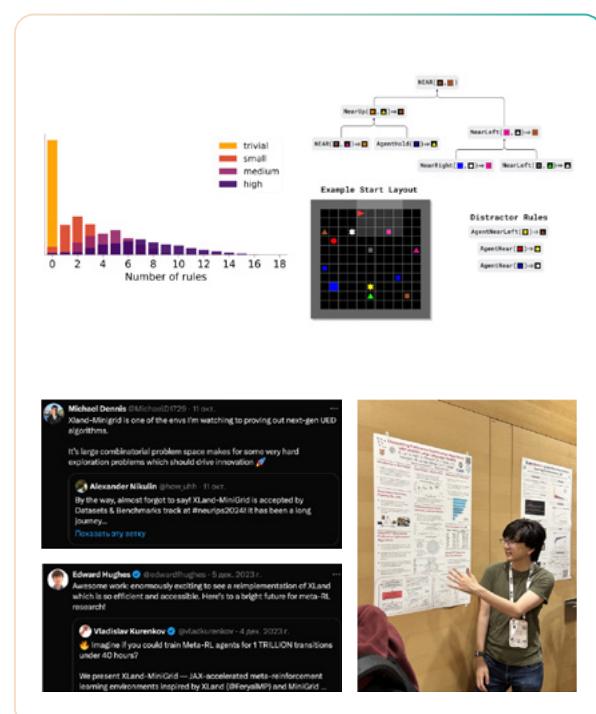
Владислав Куренков

Инженер-исследователь

Разработанная среда XLand-MiniGrid для контекстного обучения значительно ускорила скорость экспериментов, сделав область более доступной. За счет переноса вычислений среды на GPU обучение агентов для экспериментов и сбора данных ускорилось в десятки раз, с недель до минут. По сравнению с существующими аналогами, среда также предлагает большее разнообразие — миллионы уникальных задач разной сложности, пусть и абстрактных. Среда получила признание в сообществе и одобрение от множества исследователей, в том числе из Google DeepMind.

На основе среды был собран датасет XLand-100B, самый большой из существующих в области. Он позволит исследователям детальнее изучить scaling laws текущих и будущих подходов. Более того, текущие подходы в целом не справляются с решением

представленных в датасете задач, поэтому он послужит хорошим бенчмарком и путеводной звездой в области.



A*: XLand-MiniGrid: Scalable Meta-Reinforcement Learning Environments in JAX, NeurIPS 2024

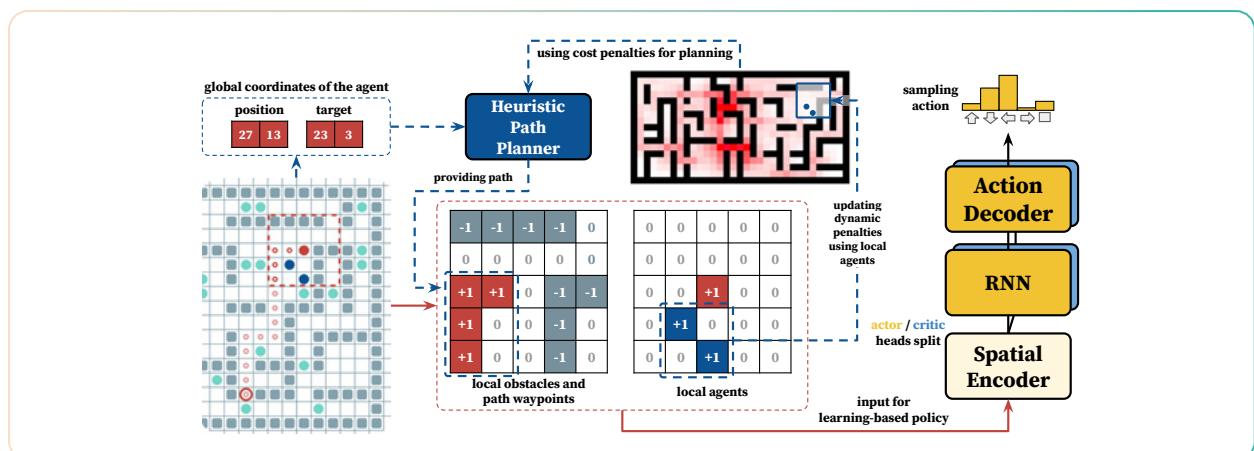
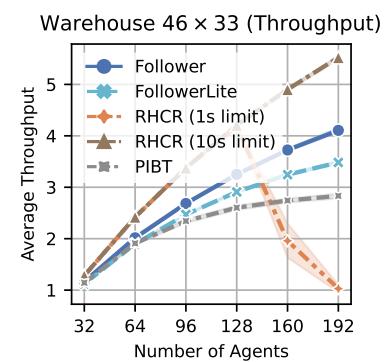
Децентрализованные методы решения задачи многоагентного планирования



Александр Панов

Директор лаборатории

Предложены новые методы решения задачи много-агентного планирования траекторий в децентрализованной постановке, опирающихся на интеграцию классические техник решения этой задачи (эвристический поиск) с современными обучаемыми методами (model-based reinforcement learning, large-scale imitation learning with transformers и др.). Совместная работа с А.А. Андрейчуком, А.А. Скрынником, М. Нестеровой, К.С. Яковлевым.



Skrynnik A., Andreychuk A., Yakovlev K., Panov A. Decentralized Monte Carlo Tree Search for Partially Observable Multi-agent Pathfinding // In Proceedings of the 38th AAAI Conference on Artificial Intelligence (AAAI 2024). pp. 17531-17540.

Skrynnik A, Andreychuk A, Nesterova M, Yakovlev K, Panov A. Learn to Follow: Decentralized Lifelong Multi-Agent Pathfinding via Planning and Learning. InProceedings of the AAAI Conference on Artificial Intelligence 2024 Mar 24 (Vol. 38, No. 16, pp. 17541-17549).

ReDisCA — метод быстрого анализа репрезентативного сходства по ЭЭГ и МЭГ данным активности головного мозга

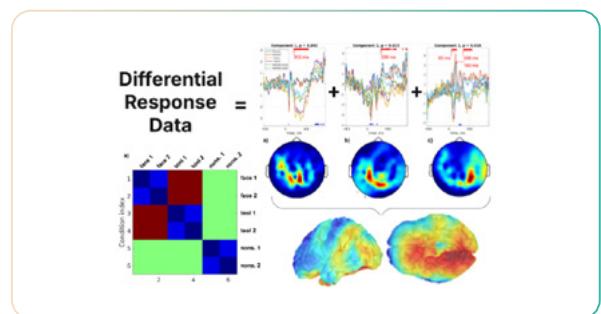


Алексей Осадчий
Ведущий научный сотрудник

Поиск зон мозга, ответ на внешние стимулы в которых обладает заданной репрезентативной (геометрической) структурой. Например, реакции на изображения инструментов и изображения лиц диаметрально противоположны и лежат по «разные стороны» от реакции на бессмысленные зрительные стимулы.

Завершение разработки метода анализа компонентов репрезентативного сходства (ReDisCA) и применение его к анализу данных МЭГ и ЭЭГ с целью поиска нейрональных источников с желаемой репрезентативной структурой. В отличие от традиционного анализа репрезентативного сходства (RSA), который требует исчерпывающего поиска по сетке нейронных источников, ReDisCA использует приблизительную стратегию оптимизации в закрытой форме для разложения данных ЭЭГ/МЭГ на пространственно-

временные компоненты с заданными репрезентативными свойствами. Это нововведение позволяет избежать проблемы множественных сравнений и значительно улучшает обнаруживаемость и локализацию источников, связанных с заданными пользователем профилями репрезентативного сходства. ReDisCA предоставляет точную и интерпретируемую основу для объединения данных визуализации мозга с вычислительными моделями. В 100+ раз быстрее традиционного RSA, в 2 раза точнее. Работа опубликована в Neuroimage, Q1, Top 1%.



Q1: A. Ossadtchi, I. Semenkov, A. Zhuravleva, O. Serikov, E. Voloshina. Representational dissimilarity component analysis (ReDisCA). NeuroImage, 2024

Вероятностно-устойчивые водяные знаки для нейросетей



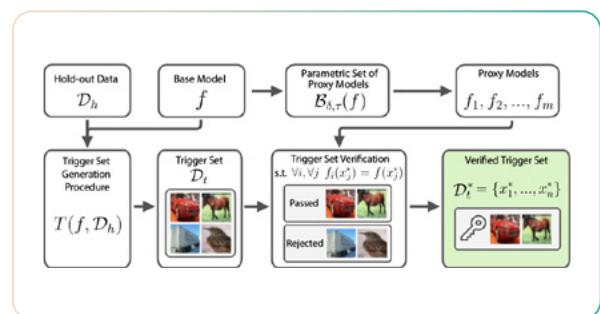
Олег Рогов

Старший научный сотрудник

Результат: метод нанесения цифровых водяных знаков, гарантированно устойчивых к наиболее сильным атакам кражи функциональности.
Статья на IJCAI-2024 (A*)

С активной эксплуатацией моделей глубокого обучения на продуктовых платформах MLaaS («Машинное обучение как услуга») динамично возрастает и интерес к методам цифровых водяных знаков для таких моделей, которые можно использовать для подтверждения права собственности на определенную модель. К сожалению, ряд существующих методов формирует цифровые водяные знаки, уязвимые для атак, направленных на кражу функционала модели. В нашем решении мы предложили принципиально новый подход к водяным знакам на основе набора триггерных данных, который демонстрирует устойчивость к атакам с кражей функциональности, особенно к тем, которые включают дистилляцию знаний. Подход не требует дополнительного

обучения модели и может быть применен к любой архитектуре. Основная идея метода заключается в поиске набора триггеров, который можно переносить между исходной моделью и набором прокси-моделей с высокой вероятностью. Мы в экспериментальном порядке показываем, что если вероятность переносимости набора достаточно высока, его можно эффективно использовать для проверки права собственности на украденную модель. Метод превосходит все аналогичные современные решения цифровой маркировки моделей машинного обучения.



Pautov M, Bogdanov N, Pyatkin S, Rogov O, Oseledets I. Probabilistically Robust Watermarking of Neural Networks. In Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence 2024 August 3-9, (Vol. 33, pp. 4778-4787)

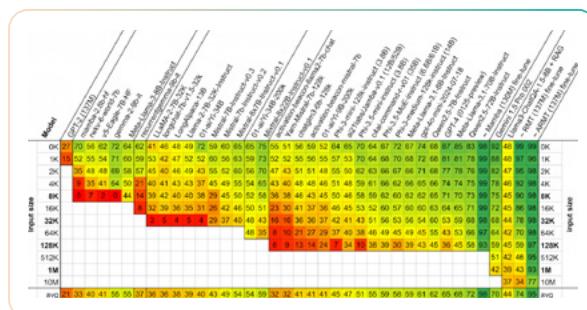
Новый бенчмарк для оценки реальной производительности языковых моделей при работе с большими объемами данных



Юрий Куратов
Старший научный сотрудник

С помощью разработанного бенчмарка BABILong продемонстрирована реальная эффективность open-source моделей и моделей семейства GigaChat на различных стадиях обучения: short-context pre-train, long-context pre-train, supervised fine-tuning. Показана высокая эффективность этапа SFT на контекстах менее 16к токенов и недостаточная на больших контекстах. Анализ подтвержден на задачах на русском и английском языке. BABILong/RuBABILong являются наиболее репрезентативными бенчмарками по оценке этапов обучения модели по эффективности работы с задачами на больших контекстах.

- С увеличением длины контекста качество у всех моделей падает, хотя интенсивность деградации варьируется.
 - Большинство моделей эффективно работают лишь с 10–20% от заявленной длины контекста.
 - Даже лучшие модели, такие как Gemini 1.5 и Qwen 2.5, показывают падение качества с 80–90% до 60% и меньше на более длинных задачах.
 - BABILong демонстрирует, что работа даже с простыми задачами с длинными контекстами — всё ещё сложна даже для лидирующих моделей. Заявленный контекст в сотни тысяч токенов вовсе не означает, что модель способна его воспринять.



A*: BABILong: Testing the Limits of LLMs with Long Context Reasoning-in-a-Haystack
Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Igorevich Sorokin, Artyom Sorokin, Mikhail Burtsev (NeurIPS).

Разработано несколько новых алгоритмов, моделей и архитектур для решения практических задач в рекомендательных системах



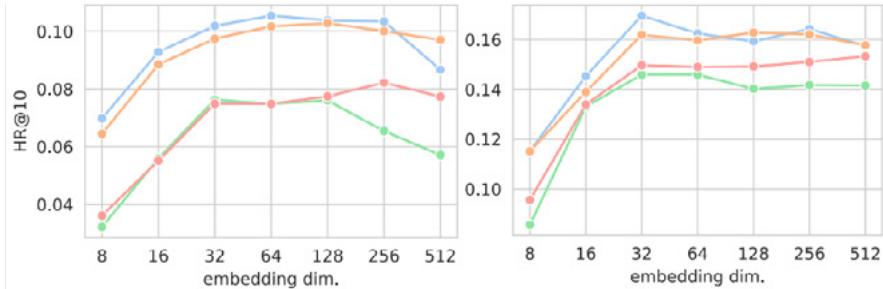
Евгений Фролов

Научный сотрудник

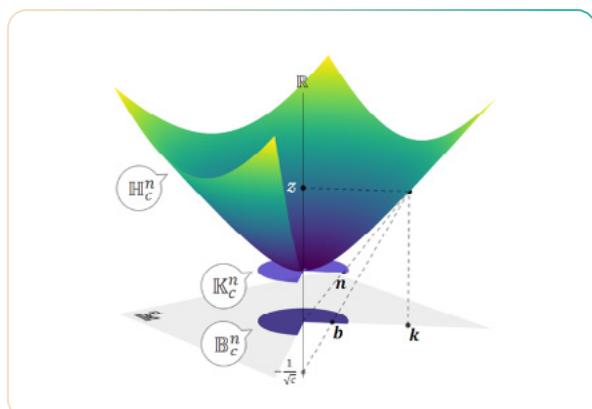
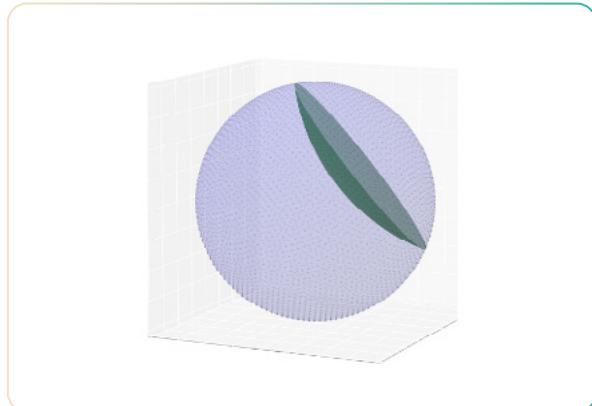
- Новый подход для масштабирования современных SOTA-моделей на огромные каталоги товаров^{2,3}. Разработанные нами решения открывают новую веху в возможностях применения ранее зарекомендовавших себя трансформерных архитектур в реальных задачах, поскольку позволяют не только значительно сократить вычислительную нагрузку при обучении моделей, но и повышают качество моделей.
- Новый фреймворк для объединения различных источников поведенческой информации в рамках end-to-end архитектуры для обучения на последовательностях¹. Универсальность предложенного решения открывает широкие возможности по моделированию гетерогенных источников информации о потребительском поведении в рамках рекомендательных

систем. Ведутся работы по усовершенствованию подхода для работы в высокодинамичном окружении.

- Разработан эффективный алгоритм кросс-доменного обучения на основе метода ADMM⁴, позволяющий улучшать качество работы рекомендательных систем за счет трансфера знаний между различными доменами. Разработка представляет особый практический интерес при работе в экосистеме маркетплейсов и онлайн-сервисов, т.к. не требует передачи персональной потребительской информации в явном виде, обеспечивая при этом эффективное обучение сразу на нескольких доменах. Ведутся работы по проверке и расширению возможностей алгоритма при обучении на более чем двух доменах.



→ Разработан геометрический подход на основе гиперболической геометрии для обучения на последовательностях действий пользователей⁵, позволяющий получить более компактные модели без потери качества рекомендаций. В рамках проекта были впервые разработаны быстрые методы расчета кривизны пространства входных данных, масштабируемые на реальные датасеты. Полученные наработки уже используются в смежных исследованиях, требующих расчетов на больших данных. Планируется опробовать подход в смежной области больших языковых моделей.



¹ A*: Baikalov, Vladimir; Frolov, Evgeny; "End-to-End Graph-Sequential Representation Learning for Accurate Recommendations", Proceedings of the ACM on Web Conference 2024, 501-504, 2024.

² A: Gusak, Danil; Mezentsev, Gleb; Oseledets, Ivan; Frolov, Evgeny; "RECE: Reduced Cross-Entropy Loss for Large-Catalogue Sequential Recommenders", Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, 3772-3776, 2024.

³ A: Mezentsev, Gleb; Gusak, Danil; Oseledets, Ivan; Frolov, Evgeny; "Scalable Cross-Entropy Loss for Sequential Recommendations with Large Item Catalogs", Proceedings of the 18th ACM Conference on Recommender Systems, 475-485, 2024.

⁴ A: Samra, Abdulaziz; Frolov, Evgeny; Vasilev, Alexey; Grigorevskiy, Alexander; Vakhrushev, Anton; "Cross-Domain Latent Factors Sharing via Implicit Matrix Factorization", Proceedings of the 18th ACM Conference on Recommender Systems, 309-317, 2024.

⁵ A: Frolov, Evgeny; Matveeva, Tatyana; Mirvakhabova, Leyla; Oseledets, Ivan; "Self-Attentive Sequential Recommendations with Hyperbolic Representations", Proceedings of the 18th ACM Conference on Recommender Systems, 981-986, 2024.



Прикладные
результаты

Разработан бенчмарк для оценки качества LLM для русского языка
MERA: A Comprehensive LLM Evaluation in Russian. (ACL) Совместно с коллегами из SberDevices и SberAI был подготовлен и опубликован бенчмарк MERA для оценивания больших языковых моделей на задачах, сформулированных на русском языке. Нами были подготовлены задачи MathLogiQA и LogiQA, оценивающие способности моделей решать математические и логические задачи.

Alena Fenogenova, Artem Chervyakov, Nikita Martynov, Anastasia Kozlova, Maria Tikhonova, Albina Akhmetgareeva, Anton Emelyanov, Denis Shevelev, Pavel Lebedev, Leonid Sinev, Ulyana Isaeva, Katerina Kolomeytseva, Daniil Moskovskiy, Elizaveta Goncharova, Nikita Savushkin, Polina Mikhailova, Anastasia Minaeva, Denis Dimitrov, Alexander Panchenko, et al.. 2024. MERA: A Comprehensive LLM Evaluation in Russian. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics

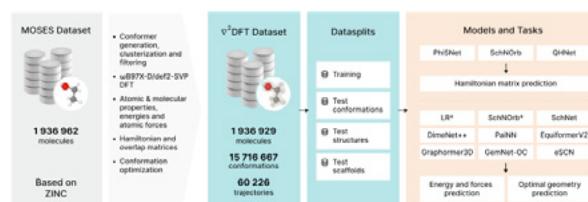
для задачи AR Manga Colorization с использованием синтетических и реальных данных. ISMAR 2024 (A*), WACV 2025 (A)



Golyadkin, M., Saraev, S., & Makarov, I. Benchmarking and Data Synthesis for Colorization of Manga Sequential Pages for Augmented Reality, ISMAR 2024.

Golyadkin, M., Plevokas, I., & Makarov, I. Closing the Domain Gap in Manga Colorization via Aligned Paired Dataset. WACV 2025

DFT: A Universal Quantum Chemistry Dataset of Drug-like Molecules and a Benchmark for Neural Network Potentials (NeurIPS D&B)



Предложен подход для создания парных датасетов, позволяющий впервые создать несинтетический датасет для задачи колоризации манги, что позволило достичь SOTA-результатов. Создан бенчмарк

Оптимизация нейросетей для извлечения лицевых дескрипторов на мобильных устройствах. Фреймворк позволяет адаптировать архитектуры нейросетей под конкретные устройства с учётом их аппаратных ограничений. Методология включает использование эволюционных алгоритмов и суррогатных бинарных классификаторов для быстрого выбора субсетей. ECAI demo 2024 (A)

Разработка моделей, оптимизированных под конкретные устройства, необходима для повышения производительности и приватности в мобильных приложениях.

Savchenko, A., Maslov, D., & Makarov, I. (2024). Device-Specific Facial Descriptors: Winning a Lottery with a SuperNet. In ECAI 2024 (pp. 4439-4442)

ACL 2024 статьи “Fact-checking the output of large language models via token-level uncertainty quantification”, посвященной использованию нового метода оценки неопределенности для fact-checking-а генерации LLM.

Fadeeva, E., Rubashevskii, A., Shelmanov, A., Petrakov, S., Li, H., Mubarak, H., Tsymbalov, E., Kuzmin, G., Panchenko, A., Baldwin, T., Nakov, P., Panov, M. (2024): Fact-Checking the Output of Large Language Models via Token-Level Uncertainty Quantification. In Proceedings of the 62th Annual Meeting of the Association for Computational Linguistics. Bangkok, Thailand. Association for Computational Linguistics.

Предложены анатомические позиционные эмбеддинги для 3D-изображений, способные точно предсказывать анатомическое положение областей, что открывает новые возможности в ретривале изображений, локализации органов и работе генеративных моделей. Эти результаты были представлены на MICCAI 2024.

Трёхмерные медицинские изображения, например, компьютерная томография (КТ), получаются в результате сканирования некоторой части тела пациента. Чаще всего сканируется довольно большая область (вся грудная клетка, или вся брюшная полость, или и то, и другое). Если нарезать КТ изображение на трёхмерные патчи, то каждый патч будет являться мини-КТ снимком небольшой анатомической области.

Обычные позиционные эмбеддинги патчей кодируют положение патча относительно края изображения. Но так как область сканирования для разных КТ изображений может отличаться (у одного пациента край изображения — это шея, а у другого диафрагма), обычные позиционные

эмбеддинги не содержат информацию об анатомическом положении патчей (у патча шеи и у патча диафрагмы будет один и тот же обычный позиционный эмбеддинг).

В нашей работе “Anatomical Positional Embeddings” [1] (Core A) мы предложили модель для получения data-driven эмбеддингов, которые кодируют информацию об анатомическом положении патчей (или пикселей) трёхмерных медицинских изображений. Обучение модели основано на том, что у большинства пациентов органы и анатомические области расположены относительно друг друга примерно одинаково. Наша модель обучается предсказывать трёхмерные позиционные эмбеддинги патчей на основе их визуальных признаков, так что расстояния между эмбеддингами согласованы с физическими расстояниями между патчами на исходном КТ изображении. Оказывается, что эмбеддинги, обученные таким образом, фактически являются 3D координатами в некоторой условной системе координат, связанной с туловищем усредненного пациента.

Приложения нашей модели:

- Retrieval изображений определенной анатомической области.
- Few-shot локализация органов.
- Label-efficient классификация анатомической области, в которой находятся находки, задетектированные другими моделями (например, классификация легочных узлов по долям и сегментам легких).
- Трэкинг находок на изображениях одного пациента с разницей во времени.
- Регистрация изображений.
- Conditioning в генеративных и дискriminативных моделях.

Goncharov, M., Samokhin, V., Soboleva, E., et al. Anatomical Positional Embeddings. MICCAI 2024. Core A

Разработан новый метод оценки полной неопределенности в применении к задачам классификации и сегментации медицинских изображений с учетом мнения экспертов (принято на WACV'25)

переноса знаний с сохранением по-классовой структуры.

- Метод для быстрого и эффективного решения задачи оптимального транспорта при дисбалансе данных.

Rethinking Optimal Transport in Offline Reinforcement Learning

A Asadulaev, R Korst, A Korotin, V Egiazarian, A Filchenkov, E Burnaev

Neural Information Processing Systems 2024 (NeurIPS 2024)

Neural Optimal Transport with General Cost Functionals
A Asadulaev, A Korotin, V Egiazarian, P Mokrov, E Burnaev
The 12th International Conference on Learning Representations (ICLR 2024)

Light Unbalanced Optimal transport
M Gazdieva, A Asadulaev, E Burnaev, A Korotin
Neural Information Processing Systems 2024 (NeurIPS 2024)

Incomplete Reinforcement Learning
A Asadulaev, R Korst, A Korotin, V Egiazarian, A Filchenkov, E Burnaev
The 12th International Conference on Learning Representations (ICLR 2024 Workshop)

Разработан фреймворк для тестирования гипотез о соответствии эмоциональных реакций людей и решений LLM. Установлено, что эмоциональный алайнмент с людьми в большинстве случаев слабый, но эмоциональные промпты, особенно негативные, существенно влияют на поведение моделей, снижая коопeração в играх и качество в этических задачах. [EAI: Emotional Decision-Making of LLMs in Strategic Games and Ethical Dilemmas]

Mozikov, M., Severin, N., Bodishtianu, V., Glushanina, M., Nasonov, I., Orekhov, D., Pekhotin, V.,

Makovetskiy, I., Baklashkin, M., Lavrentyev, V., Tsvigun, A., Turdakov, D., Shavrina, T., Savchenko, A., &

Makarov, I. EAI: Emotional Decision-Making of LLMs in Strategic Games and Ethical Dilemmas. NeurIPS

Разработаны новые подходы к применению оптимального транспорта для задач машинного обучения.

- Метод обучения с подкреплением для построения автономных агентов. Данный подход позволяет осуществлять эффективное обучение на неполных или зашумленных данных.
- Расширение алгоритма нейронного оптимального транспорта на задачи



Научные
публикации

Статьи на конференциях A*

5

AAAI

4

ACL

2

CVPR

6

ICLR

7

ICML

2

ISMAR

17

NeurIPS

1

ICRA

1

WWW

3

ACM KDD

2

ECCV

4

ICDM

4

IJCAI

7

EMNLP

Beyond Attention: Breaking the Limits of Transformer Context Length with Recurrent Memory

Aydar Bulatov, Yuri Kuratov, Yermek Kapushev, Mikhail Burtsev

A major limitation for the broader scope of problems solvable by transformers is the quadratic scaling of computational complexity with input size. In this study, we investigate the recurrent memory augmentation of pre-trained transformer models to extend input context length while linearly scaling compute. Our approach demonstrates the capability to store information in memory for sequences of up to an unprecedented two million tokens while maintaining high retrieval accuracy. Experiments with language modeling tasks show perplexity improvement as the number of processed input segments increases. These results underscore the effectiveness of our method, which has significant potential to enhance long-term dependency handling in natural language understanding and generation tasks, as well as enable large-scale context processing for memory-intensive applications.

The Thirty-Eighth AAAI Conference on Artificial Intelligence (AAAI-14)



Topological and Node Noise Filtering on 3D Meshes Using Graph Neural Networks

Vladimir Mashurov, Natalia Semenova

Topological and node noise filtration are typically considered separately. Graph Neural Networks (GNN) are commonly used for node noise filtration, as they offer high efficiency and low exploitation costs. This paper explores the solution of joint node and topological noise filtration through the use of graph neural networks. Since treating a 3D mesh as a graph is challenging, an indicator function grid representation is employed as input for GNNs to perform the joint filtering. The resulting machine learning model is inspired by point cloud to mesh reconstruction algorithms and demonstrates low computational requirements during inference, producing successful results for smooth, watertight 3D models.

Seite 18 von 20 IASTM Conference on Medical Biomechanics AAA2014



A standard black and white QR code located at the bottom right of the page, which links to the original source of the information.

Improved Anonymous Multi-Agent Path Finding Algorithm

Zain Alabedeen Ali, Konstantin Yakovlev

We consider an Anonymous Multi-Agent Path-Finding (AMAPF) problem where the set of agents is confined to a graph, a set of goal vertices is given and each of these vertices has to be reached by some agent. The problem is to find an assignment of the goals to the agents as well as the collision-free paths, and we are interested in finding the solution with the optimal makespan. A well-established approach to solve this problem is to reduce it to a special type of a graph search problem, i.e. to the problem of finding a maximum flow on an auxiliary graph induced by the input one. The size of the former graph may be very large and the search on it may become a bottleneck. To this end, we suggest a specific search algorithm that leverages the idea of exploring the search space not through considering separate search states but rather bulks of them simultaneously. That is, we implicitly compress, store and expand bulks of the search states as single states, which results in high reduction in runtime and memory. Empirically, the resultant AMAPF solver...

Learn to follow: Decentralized lifelong multi-agent pathfinding via planning and learning

Alexey Skrynnik, Anton Andreychuk, Maria Nesterova, Konstantin Yakovlev, Aleksandr Panov

Multi-agent Pathfinding (MAPF) problem generally asks to find a set of conflict-free paths for a set of agents confined to a graph and is typically solved in a centralized fashion. Conversely, in this work, we investigate the decentralized MAPF setting, when the central controller that possesses all the information on the agents' locations and goals is absent and the agents have to sequentially decide the actions on their own without having access to a full state of the environment. We focus on the practically important lifelong variant of MAPF, which involves continuously assigning new goals to the agents upon arrival to the previous ones. To address this complex problem, we propose a method that integrates two complementary approaches: planning with heuristic search and reinforcement learning through policy optimization. Planning is utilized to construct and re-plan individual paths. We enhance our planning algorithm with a dedicated technique tailored to avoid congestion and increase the throughput of the system. We employ reinforcement learning to discover the collision avoidance policies that effectively guide the agents along the paths. The policy is implemented as a neural network and is effectively trained...



Ссылка на источник

Learn to Follow: Decentralized Lifelong Multi-agent Pathfinding via Planning and Learning

Aleksy Skrynnik^{1,2}, Anton Serebryakov¹, Maria Nesterova¹,
Semyonov Valeriy^{1,2}, Aleksandr Panin¹,
AIBI¹, Moscow, Russia
¹Federal Research Center for Computer Technologies and Network Analysis of Sciences, Moscow, Russia
²Corresponding author: aleksy.skrnnk@gmail.com

Abstract

Multi-agent pathfinding (MAPF) is a problem of great interest in game development, robotics, and other fields. One of the main challenges is to find a solution that is both safe and efficient. In this paper, we propose a decentralized learning-based approach to solving MAPF. The proposed method is based on the notion of a learned map, which is a representation of the environment that is learned by each agent independently. The learned map is used to guide the agents' movement, and it is updated as the agents move. This allows the agents to learn their own behavior from scratch, without the need for explicit communication or coordination. The results of our experiments show that the proposed method is able to solve complex MAPF instances in a timely manner, while maintaining a low computational cost.

Figure

Figure caption

Figure 1: An example of a deconstructed LAMAP instance. Agents are depicted as colored circles. The learned (the other) obstacles are colored in blue. The red circles represent the agents' current positions. The green circles represent the agents' goals. The blue circles represent the learned obstacles. The learned map is shown as thick black lines, which are only present when the agents move.

Introduction

Multi-agent pathfinding (MAPF) [1] (hereinafter referred to as MAPF) is a problem of great interest in game development, robotics, and other fields. It is a problem of great complexity, as it requires finding a solution that is both safe and efficient. In this paper, we propose a decentralized learning-based approach to solving MAPF. The proposed method is based on the notion of a learned map, which is a representation of the environment that is learned by each agent independently. The learned map is used to guide the agents' movement, and it is updated as the agents move. This allows the agents to learn their own behavior from scratch, without the need for explicit communication or coordination. The results of our experiments show that the proposed method is able to solve complex MAPF instances in a timely manner, while maintaining a low computational cost.

Intrusion detection

In general, world domains, however, it is not possible to provide a complete solution to the MAPF problem. This is due to the fact that there are many constraints that must be taken into account. For example, the time of the day, the weather, the location of the objects in the environment, etc. These factors can significantly affect the behavior of the agents and the efficiency of the solution. Therefore, it is necessary to take into account these factors when solving the MAPF problem. This can be done by using a learned map, which is a representation of the environment that is learned by each agent independently. The learned map is used to guide the agents' movement, and it is updated as the agents move. This allows the agents to learn their own behavior from scratch, without the need for explicit communication or coordination. The results of our experiments show that the proposed method is able to solve complex MAPF instances in a timely manner, while maintaining a low computational cost.

Related work

There are many approaches to solving MAPF problems. One of the most common is the centralized approach, where all agents are controlled by a single central controller. This approach has the advantage of being able to find a solution quickly, but it also has the disadvantage of being less efficient than decentralized approaches. Another approach is the decentralized approach, where each agent is controlled by its own local controller. This approach has the advantage of being more efficient than centralized approaches, but it also has the disadvantage of being less safe. There are also hybrid approaches, which combine the strengths of both centralized and decentralized approaches. These approaches have the advantage of being able to find a solution quickly and efficiently, while also being safe. However, they are also more complex than the other approaches.

Conclusion

In this paper, we proposed a decentralized learning-based approach to solving MAPF. The proposed method is based on the notion of a learned map, which is a representation of the environment that is learned by each agent independently. The learned map is used to guide the agents' movement, and it is updated as the agents move. This allows the agents to learn their own behavior from scratch, without the need for explicit communication or coordination. The results of our experiments show that the proposed method is able to solve complex MAPF instances in a timely manner, while maintaining a low computational cost.

References

[1] A. Skrynnik, A. Serebryakov, M. Nesterova, S. Valeriy, A. Panin, "Learn to Follow: Decentralized Lifelong Multi-agent Pathfinding via Planning and Learning", arXiv preprint arXiv:2309.13420, 2023.



[Ссылка на источник](#)

Decentralized Monte Carlo Tree Search for Partially Observable Multi-agent Pathfinding

Alexey Skrynnik, Anton Andreychuk, Konstantin Yakovlev,
Aleksandr Panov

The Multi-Agent Pathfinding (MAPF) problem involves finding a set of conflict-free paths for a group of agents confined to a graph. In typical MAPF scenarios, the graph and the agents' starting and ending vertices are known beforehand, allowing the use of centralized planning algorithms. However, in this study, we focus on the decentralized MAPF setting, where the agents may observe the other agents only locally and are restricted in communications with each other. Specifically, we investigate the lifelong variant of MAPF, where new goals are continually assigned to the agents upon completion of previous ones.

assigned to the agents upon completion of previous ones. Drawing inspiration from the successful AlphaZero approach, we propose a decentralized multi-agent Monte Carlo Tree Search (MCTS) method for MAPF tasks. Our approach utilizes the agent's observations to recreate the intrinsic Markov decision process, which is then used for planning with a tailored for multi-agent tasks version of neural MCTS. The experimental results show that our approach outperforms state-of-the-art learnable MAPF solvers. The source code is available at this [https URL](https://github.com/AIRI-Institute/mats-lp): <https://github.com/AIRI-Institute/mats-lp>



Ссылка на источник

Fact-Checking the Output of Large Language Models via Token-Level Uncertainty Quantification

Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, Maxim Panov

Large language models (LLMs) are notorious for hallucinating, i.e., producing erroneous claims in their output. Such hallucinations can be dangerous, as occasional factual inaccuracies in the generated text might be obscured by the rest of the output being generally factual, making it extremely hard for the users to spot them. Current services that leverage LLMs usually do not provide any means for detecting unreliable generations. Here, we aim to bridge this gap. In particular, we propose a novel fact-checking and hallucination detection pipeline based on token-level uncertainty quantification. Uncertainty scores leverage information encapsulated in the output of a neural network or its layers to detect unreliable predictions, and we show that they can be used to fact-check the atomic claims in the LLM output.

Your Transformer is Secretly Linear

Anton Razzhigaev, Matvey Mikhalkchuk, Elizaveta Goncharova, Nikolai Gerasimenko, Ivan Oseledets, Denis Dimitrov, Andrey Kuznetsov

This paper reveals a novel linear characteristic exclusive to transformer decoders, including models such as GPT, LLaMA, OPT, BLOOM and others. We analyze embedding transformations between sequential layers, uncovering a near-perfect linear relationship (Procrustes similarity score of 0.99). However, linearity decreases when the residual component is removed due to a consistently low output norm of the transformer layer. Our experiments show that removing or linearly approximating some of the most linear blocks of transformers does not affect significantly the loss or model performance. Moreover, in our pretraining experiments on smaller models we introduce a cosine-similarity-based regularization, aimed at reducing layer linearity. This regularization improves performance metrics on benchmarks like Tiny Stories and SuperGLUE and as well successfully decreases the linearity of the models. This study challenges the existing understanding of transformer architectures, suggesting that their operation may be more linear than previously assumed.

Fact-Checking the Output of Large Language Models via Token-Level Uncertainty Models

Eduardo Fuchs^{1,2,3}, Aleksandr Roshchina^{2,4}, Arvin Shabani^{1,2},
Sergey Gorbunov¹, Barbara L. Webb¹, and Mihir Mathur¹,
Alexander P. Kuleshov¹, Michael J. Zelnick¹, and Emanuele Natale¹,
“MIRALAB” AIRES¹, Center for Artificial Intelligence Technology – “HSE University”², “Yandex”³, and University of Cambridge⁴

¹(Eduardo fuchs, sergey.gorbunov, barbara.l.webb, mihir.mathur, alexander.p.kuleshov, emanuele.natale)@miタルバ.hse.ru; ²(arvin.shabani, aleksandr.roshchina)@miタルバ.hse.ru; ³eduardo.fuchs@yandex.ru; ⁴zelnick@cam.ac.uk

Abstract

Large language models (LLMs) are meant for fact-checking, e.g., predicting whether a statement is true or false. However, they can be inaccurate, as measured in factoid benchmarks. In this paper, we propose a new way to measure the net of the output being predicted by LLMs. We introduce a metric called “Token-Level Uncertainty” (TLU), which measures uncertainty at the token level. TLU allows us to quickly assess the quality of LLMs and to bridge this gap. In particular, we propose a two-stage approach to fact-checking. First, we use TLU to quickly filter out low-quality predictions based on token-level uncertainty. Then, we use a standard fact-checking framework to process the output of a neural network trained on a dataset of labeled examples. This two-stage approach is more efficient than a standard one, and it shows that they can be used to quickly filter out low-quality predictions from LLMs without losing much accuracy. Moreover, we present a novel token-level uncertainty metric called “Token-Level Uncertainty” (TLU). TLU measures the impact of uncertainty about what tokens are likely to appear in the output of an LLM. We use TLU to estimate the uncertainty of a particular token value represented by the LLM. This metric is useful for fact-checking, as it provides precise information on strong, right predictions. It also helps to identify strong, wrong predictions. TLU is a simple metric that can be used to quickly filter out low-quality predictions based on uncertainty quantification in concert with a standard fact-checking framework.

1 Introduction

Large language models (LLMs) have become a ubiquitous tool for solving and addressing a variety of natural language processing (NLP) tasks. Despite their remarkable performance, LLMs often make mistakes (Gu et al., 2020; Shi et al., 2020), to

questions (Thorne et al., 2021), or to generate new content (Chen et al., 2021). Recently, there has been a growing interest in improving an LLM’s robustness and reliability by adding traditional checks. However, a significant challenge is how to efficiently incorporate such checks to produce “trustworthy” LLMs, i.e., formally guaranteeing that they do not produce harmful outputs (Thorat et al., 2021; Date et al., 2020).

LLM fact-checking is a major concern because the LLMs’ outputs are often generated by highly confident and pernicious. Common checks for LLMs include the following: (i) the presence of untrusted claims; (ii) the presence of the asserted or unasserted claims. The danger is that the LLMs can produce a large number of false positives, making it extremely difficult to spot them. Moreover, the LLMs’ outputs are often very complex, where each output could be considered as a fact or a falsehood, depending on the context, and in this way our approach can help.

There are several ways to address this problem. One is a standard approach using a fact-checking framework (Gu et al., 2020; Shi et al., 2020). This approach is problematic and time-consuming, as it requires a lot of time and effort over time in terms of storing the knowledge base. Another approach is to use uncertainty quantification in a hallucination (as measured in the model’s confidence score) to detect if the LLM’s prediction is a hallucination (Gu et al., 2020; Kewley-Porter et al., 2020). This approach is also problematic, as it needs implementing complex and expensive fact-checking frameworks that require additional computation.

Our work has mainly focused on the following:



Ссылка на источник

Your Transformer Is Secretly Linear

Anton Rakhlin^{1,2}, Marvuy Mikhalev³, Elizaveta Gordeeva^{2,4},
Nikolai Germanov¹, Ivan Oseledets^{1,2}, and Andrey Karlovich⁵
¹Skolkovo Institute of Science and Technology
²Yandex Moscow Research University
³Yandex Moscow Research University
⁴Yandex Moscow Research University
⁵Yandex Moscow Research University

Abstract

This paper reveals a novel linear characteristic in transformer models, including masked language modeling, text generation, and others. We analyze embedding matrices of the first layer of the model and find that they exhibit a new perfect linear relationship. This phenomenon occurs during the training component in most of the transformer models. We also study the behavior of the last linear block of the model. Our experiments show that the number of dimensions that are often used in transformer models, such as 768 or 1024, are too large. Moreover, we are proposing a parameter α that controls the number of dimensions in a coarse-grained manner, instead of using a single dimension. We also propose a new metric for evaluating the quality of the model, namely, the Frobenius and Frobenius-L2 norms, as well as the cosine similarity between the embeddings. This study challenges the existing understanding of the transformer model's architecture, suggesting that their operation may be more linear than previously thought.

1 Introduction

Transformers have revolutionized the field of automated language processing, offering unprecedented accuracy (Vaswani et al., 2017; Devlin et al., 2019; Lample & Peter, 2020). However, despite their widespread adoption and success, the complex work of these models remains an enigma to many researchers and practitioners. One aspect that has received less attention is the inherent linearities of the transformer model, both within themselves and in their applications. In this study, we reveal a new linear characteristic in the representation of words, specifically focusing on masked language modeling, during the pre-training and fine-tuning phases.

The embedding matrix of the first layer of the transformer exhibits a new linear relationship. This observation is quantified by calculating the Frobenius norm of the difference of a new perfect linear identity of size 1024. Such a discovery not only challenges the traditional understanding of the transformer model, but also opens new opportunities for model optimization.

Based on this insight, we introduce several new contributions:

- **Extensive analysis of the linearity properties of transformer models and its dynamics by training them on different datasets.**
- **Introducing a new regularization approach for depth pruning of transformer models, allowing to reduce the number of layers without a significant loss in performance.**
- **A novel distillation technique that involves pruning, replacing linear layers with linear layers, and fine-tuning the model to maintain the quality of the embeddings to model performance.**
- **Introducing a new regularization approach for depth pruning of transformer models, allowing to reduce the number of layers. This method is based on the Frobenius norm of the difference of normalized vectors on benchmark datasets (Huang et al., 2018; Vaswani et al., 2017; Devlin et al., 2019; Lample & Peter, 2020), that improves the expressiveness of the model, as well as its generalization ability.**

With our findings, we are paving the way for more compact, efficient, and robust transformer architectures. We hope that this work will inspire further research, thereby addressing one of the critical challenges in developing these models.



[Ссылка на источник](#)

MERA: A Comprehensive LLM Evaluation in Russian

Alena Fenogenova, Artem Chervyakov, Nikita Martynov, Anastasia Kozlova, Maria Tikhonova, Albina Akhmetgareeva, Anton A. Emelyanov, Denis Shevelev, Pavel Lebedev, Leonid Sinev, Ulyana Isaeva, Katerina Kolomeytseva, Daniil Moskovskiy, Elizaveta Goncharova, Nikita Savushkin, Polina Mikhailova, Denis Dimitrov, Alexander Panchenko, Sergey Markov

Over the past few years, one of the most notable advancements in AI research has been in foundation models (FMs), headlined by the rise of language models (LMs). However, despite researchers' attention and the rapid growth in application, the capabilities, limitations, and associated risks still need to be better understood. To address these issues, we introduce a new instruction benchmark, MERA, for evaluating foundation models oriented towards the Russian language. The benchmark encompasses 21 evaluation tasks for generative models covering 10 skills and is designed as a black-box test to ensure the exclusion of data leakage. The paper introduces a methodology to evaluate FMs and LMs...

MERA: A Comprehensive LLM Evaluation in Russian

Alena Fenogenova¹, Artem Chervyakov^{1,2}, Nikita Martynov¹, Anastasia Kozlova¹, Maria Tikhonova¹, Albina Akhmetgareeva¹, Anton Emelyanov¹, Denis Shevelev¹, Pavel Lebedev¹, Leonid Sinev¹, Ulyana Isaeva¹, Katerina Kolomeytseva¹, Daniil Moskovskiy¹, Elizaveta Goncharova¹, Nikita Savushkin¹, Polina Mikhailova¹, Denis Dimitrov¹, Alexander Panchenko^{1,2}, Sergey Markov¹

¹Sabot Devise, ²HSE University, "Center for Artificial Intelligence Technology", Correspondence: markov@ai.hse.ru

Abstract

Over the past few years, one of the most notable advancements in AI research has been in foundation models (FMs), headlined by the rise of language models (LMs). However, despite researchers' attention and the rapid growth in application, the capabilities, limitations, and associated risks still need to be better understood. To address these issues, we introduce a new instruction benchmark, MERA, for evaluating foundation models oriented towards the Russian language. The benchmark encompasses 21 evaluation tasks for generative models covering 10 skills and is designed as a black-box test to ensure the exclusion of data leakage. The paper introduces a methodology to evaluate FMs and LMs...

The community has addressed the issue with several metrics and benchmarks (e.g., HELM, BLOOM, etc., see [2021, 2022]). HELM (Bloom et al., 2021), MT Bench (Zhang et al., 2022) which are designed for English, and BLOOM (Bommasani et al., 2022) and TAF (Tafjord et al., 2022) do not cover Russian. Therefore, we believe that a comprehensive Russian benchmark should be created to satisfy research needs and to facilitate an understanding of LLM behavior.

This paper addresses the problem of creating a comprehensive Russian benchmark for evaluating foundation models oriented towards the Russian language. We present MERA, a new instruction benchmark for Russian LMs. MERA includes 21 evaluation tasks for generative models covering 10 skills. The paper also introduces a methodology to evaluate FMs and LMs...

[1] <https://arxiv.org/pdf/2404.07001.pdf>

A*



Ссылка на источник

TaxoLLaMA: WordNet-based Model for Solving Multiple Lexical Semantic Tasks

Viktor Moskvoretskii, Ekaterina Neminova, Alina Lobanova, Alexander Panchenko, Irina Nikishina

In this paper, we explore the capabilities of LLMs in capturing lexical-semantic knowledge from WordNet on the example of the LLaMA-2-7b model and test it on multiple lexical semantic tasks. As the outcome of our experiments, we present TaxoLLaMA, the everything-in-one model, lightweight due to 4-bit quantization and LoRA. It achieves 11 SotA results, 4 top-2 results out of 16 tasks for the Taxonomy Enrichment, Hypernym Discovery, Taxonomy Construction, and Lexical Entailment tasks. Moreover, it demonstrates very strong zero-shot performance on Lexical Entailment and Taxonomy Construction tasks with no fine-tuning. We also explore its hidden multilingual and domain adaptation capabilities with a little tuning or few-shot learning. All datasets, code, and model are available online at this https://tinyurl.com/yxqzr4z7

TaxoLLaMA: WordNet-based Model for Solving Multiple Lexical Semantic Tasks

Viktor Moskvoretskii¹, Ekaterina Neminova², Alina Lobanova², Alexander Panchenko¹, Irina Nikishina¹, HSE University, Skolkovo, FUDI, University of Marburg, [\[{moskvoretskii, e.neminova, alobanova}@cs.hse.ru\]](mailto:{moskvoretskii, e.neminova, alobanova}@cs.hse.ru), irinani@yandex.ru

Abstract

In this paper, we explore the capabilities of LLMs in capturing lexical-semantic knowledge from WordNet on the example of the LLaMA-2-7b model. As the outcome of our experiments, we present TaxoLLaMA, the everything-in-one model, lightweight due to 4-bit quantization and LoRA. It achieves 11 SotA results, 4 top-2 results out of 16 tasks for the Taxonomy Enrichment, Hypernym Discovery, Taxonomy Construction, and Lexical Entailment tasks. Moreover, it demonstrates very strong zero-shot performance on Lexical Entailment and Taxonomy Construction tasks with no fine-tuning. All datasets, code, and trained models are available online.

[1] <https://arxiv.org/pdf/2404.07001.pdf>

[2] <https://tinyurl.com/yxqzr4z7>

A*

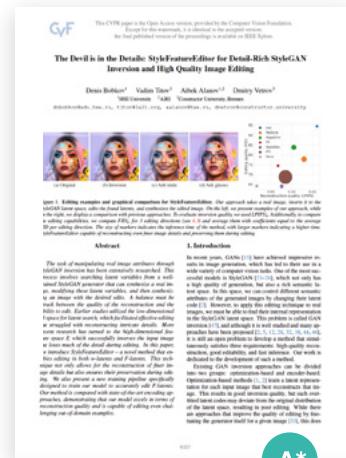


Ссылка на источник

The Devil is in the Details: StyleFeatureEditor for Detail-Rich StyleGAN Inversion and High Quality Image Editing

Denis Bobkov, Vadim Titov, Aibek Alanov, Dmitry Vetrov

The task of manipulating real image attributes through StyleGAN inversion has been extensively researched. This process involves searching latent variables from a well-trained StyleGAN generator that can synthesize a real image, modifying these latent variables, and then synthesizing an image with the desired edits. A balance must be struck between the quality of the reconstruction and the ability to edit. Earlier studies utilized the low-dimensional W-space for latent search, which facilitated effective editing but struggled with reconstructing intricate details. More recent research has turned to the high-dimensional feature space F, which successfully inverses the input image but loses much of the detail during editing. In this paper, we introduce StyleFeatureEditor— a novel method that enables editing in both w-latents and F-latents.

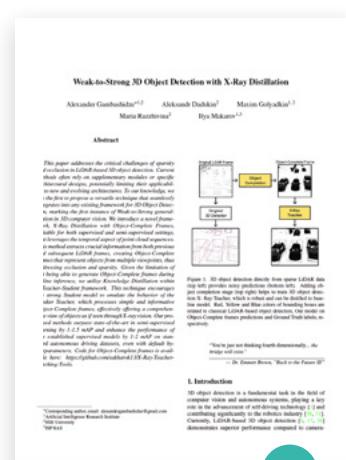


Ссылка на источник

Weak-to-Strong 3D Object Detection with X-Ray Distillation

Alexander Gambashidze, Aleksandr Dadukin, Maksim Golyadkin, Maria Razzhivina, Ilya Makarov

This paper addresses the critical challenges of sparsity and occlusion in LiDAR-based 3D object detection. Current methods often rely on supplementary modules or specific architectural designs, potentially limiting their applicability to new and evolving architectures. To our knowledge, we are the first to propose a versatile technique that seamlessly integrates into any existing framework for 3D Object Detection, marking the first instance of Weak-to-Strong generalization in 3D computer vision. We introduce a novel framework, X-Ray Distillation with Object-Complete Frames, which leverages the temporal aspect of point-cloud sequences to extract object-complete frames from LiDAR point clouds. Our method uses a Teacher-Student framework to distill knowledge from previous frames to generate object-complete frames, improving the performance of weak models. We demonstrate the effectiveness of our approach on various datasets, even with difficult-to-detect objects like birds. Our proposed method achieves 4.1X Ray-Distilling Ratio.



Ссылка на источник

Gradual Optimization Learning for Conformational Energy Minimization

Artem Tsypin, Leonid Ugadiarov, Kuzma Khrabrov, Manvel Avetisian, Alexander Telepov, Egor Rumiantsev, Alexey Skrynnik, Aleksandr I Panov, Dmitry Vetrov, Elena Tutubalina, Artur Kadurin

Molecular conformation optimization is crucial to computer-aided drug discovery and materials design. Traditional energy minimization techniques rely on iterative optimization methods that use molecular forces calculated by a physical simulator (oracle) as anti-gradients. However, this is a computationally expensive approach that requires many interactions with a physical simulator. One way to accelerate this procedure is to replace the physical simulator with a neural network. Despite recent progress in neural networks for molecular conformation energy prediction, such models are prone to distribution shift, leading to inaccurate energy minimization. We find that the quality of energy minimization with neural networks can be improved by providing optimization trajectories as additional training data. Still, it takes around 5×10^5 additional conformations to match the physical simulator's optimization quality. In this work, we present the Gradual Optimization...



A*



Ссылка на источник

Light Schrödinger Bridge

Alexander Korotin, Nikita Gushchin, Evgeny Burnaev

Despite the recent advances in the field of computational Schrödinger Bridges (SB), most existing SB solvers are still heavy-weighted and require complex optimization of several neural networks. It turns out that there is no principal solver which plays the role of simple-yet-effective baseline for SB just like, e.g., -means method in clustering, logistic regression in classification or Sinkhorn algorithm in discrete optimal transport. We address this issue and propose a novel fast and simple SB solver. Our development is a smart combination of two ideas which recently appeared in the field: (a) parameterization of the Schrödinger potentials with sum-exp quadratic functions and (b) viewing the log-Schrödinger potentials as the energy functions. We show that combined together these ideas yield a lightweight, simulation-free and theoretically justified SB solver with a simple straightforward optimization objective. As a result, it allows solving SB in moderate dimensions in a matter of minutes on CPU without a painful hyperparameter selection. Our light solver resembles the Gaussian mixture model which is widely used for density estimation. Inspired by this similarity, we also prove an important theoretical result showing that our light solver is a universal approximator of SBs.



A*

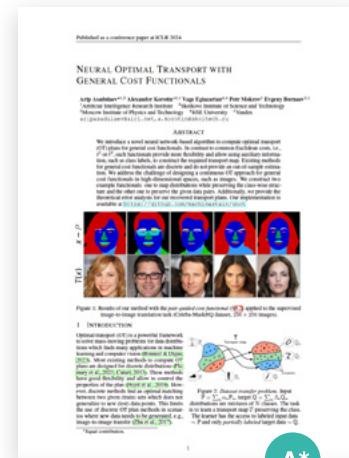


Ссылка на источник

Neural optimal transport with general cost functionals

Arip Asadulaev, Alexander Korotin, Vage Egiazarian, Petr Mokrov, Evgeny Burnaev

We introduce a novel neural network-based algorithm to compute optimal transport (OT) plans for general cost functionals. In contrast to common Euclidean costs, i.e., or , such functionals provide more flexibility and allow using auxiliary information, such as class labels, to construct the required transport map. Existing methods for general cost functionals are discrete and do not provide an out-of-sample estimation. We address the challenge of designing a continuous OT approach for general cost functionals in high-dimensional spaces, such as images. We construct two example functionals: one to map distributions while preserving the class-wise structure and the other one to preserve the given data pairs. Additionally, we provide the theoretical error analysis for our recovered transport plans. Our implementation is available at <https://github.com/machinestein/gnot>



A*

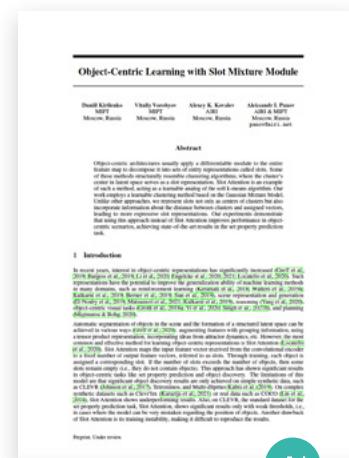


[Ссылка на источник](#)

Object-Centric Learning with Slot Mixture Module

Daniil Kirilenko, Vitaliy Vorobyov, Alexey K. Kovalev, Aleksandr I. Panov

Object-centric architectures usually apply a differentiable module to the entire feature map to decompose it into sets of entity representations called slots. Some of these methods structurally resemble clustering algorithms, where the cluster's center in latent space serves as a slot representation. Slot Attention is an example of such a method, acting as a learnable analog of the soft k-means algorithm. Our work employs a learnable clustering method based on the Gaussian Mixture Model. Unlike other approaches, we represent slots not only as centers of clusters but also incorporate information about the distance between clusters and assigned vectors, leading to more expressive slot representations. Our experiments demonstrate that using this approach instead of Slot Attention improves performance in object-centric scenarios, achieving state-of-the-art results in the set property prediction task.



A*



[Ссылка на источник](#)

Estimating Barycenters of Distributions with Neural Optimal Transport

Alexander Kolesov, Petr Mokrov, Igor Udovichenko, Milena Gazdieva, Gudmund Pammer, Evgeny Burnaev, Alexander Korotin

Given a collection of probability measures, a practitioner sometimes needs to find an «average» distribution which adequately aggregates reference distributions. A theoretically appealing notion of such an average is the Wasserstein barycenter, which is the primal focus of our work. By building upon the dual formulation of Optimal Transport (OT), we propose a new scalable approach for solving the Wasserstein barycenter problem. Our methodology is based on the recent Neural OT solver: it has bi-level adversarial learning objective and works for general cost functions. These are key advantages of our method since the typical adversarial algorithms leveraging barycenter tasks utilize tri-level optimization and focus mostly on quadratic cost. We also establish theoretical error bounds for our proposed approach and showcase its applicability and effectiveness on illustrative scenarios and image data setups.

Neural operators meet conjugate gradients: The FCG-NO method for efficient PDE solving

Alexander Rudikov, Vladimir Fanaskov, Ekaterina Muravleva, Yuri M. Laevsky, Ivan Oseledets

Deep learning solvers for partial differential equations typically have limited accuracy. We propose to overcome this problem by using themas preconditioners. More specifically, we apply discretization-invariant neural operators to learn preconditioners for the flexible conjugate gradient method (FCG). Architecture paired with novel loss function and training scheme allows for learning efficient preconditioners that can be used across different resolutions. On the theoretical side, FCG theory allows us to safely use nonlinear preconditioners that can be applied in $O(N)$ operations without constraining the form of the preconditioners matrix. To justify learning scheme components (the loss function and the way training data is collected) we perform several ablation studies. Numerical results indicate that our approach favorably compares with classical preconditioners and allows to reuse of preconditioners learned for lower resolution to the higher resolution data.

Estimating Barycenters of Distributions with Neural Optimal Transport

Alexander Kuhnle^{1,*}, Piotr Mokrusz¹, Igor Ulyanov¹, Milica Gašević², Gundolf Panzer¹,
Tugay Birhan¹, Alexander Kröse¹

Abstract

Given a collection of probability measures, a principal task in optimal transport is to find their barycenter distribution which adequately aggregates these distributions. This is often done by averaging each of them at the average of the Wasserstein barycenter, which is the unique minimizer of the sum of squared distances from all measures to the barycenter. We propose the soft version of the formulation of Optimal Transport (OT) that allows us to compute the barycenter without solving the Wasserstein barycenter problem. Our method is based on the neural network representation of OT values. It is a local alternating iterative procedure that finds the barycenter in a single step. There are few advantages of our method over the traditional one: it is faster, it can handle much larger sets of tasks while still optimizing and tuning hyperparameters, and it is more robust. We also present critical bend points for the proposed approach and analyze its performance on several datasets. Our illustrative examples and extensive numerical data show, that the proposed method is competitive with the state-of-the-art methods.

<https://arxiv.org/pdf/2305.13895.pdf>

Introduction

Optimal Transport (OT) is a well-known mathematical framework for comparing probability distributions. It is usually represented by the Wasserstein distance, which measures the average cost of transporting mass from one distribution to another. In particular, it minimizes the average of specific distances in the space of measures. The main applications of OT in machine learning are to obtain a common representation of data, to measure the similarity between two datasets, and to estimate the barycenter of a collection of measures. OT has been used in various fields such as generative modeling, domain adaptation, and learning with the minimal loss (Amit et al., 2019; Arjovsky et al., 2017).

Joint contribution: ¹Machine Learning Department, University of Amsterdam, Amsterdam, The Netherlands; ²Software Engineering Institute, TU Berlin, Berlin, Germany. *Correspondence to: Alexander Kuhnle (alexander.kuhnle@uva.nl). Correspondence to: Piotr Mokrusz (piotr.mokrusz@uva.nl). Correspondence to: Gundolf Panzer (gundolf.panzer@uva.nl). Correspondence to: Tugay Birhan (tugay.birhan@uva.nl). Correspondence to: Alexander Kröse (alexander.kroese@uva.nl). Correspondence to: Igor Ulyanov (igor.ulyanov@uva.nl). Correspondence to: Milica Gašević (mila.gasevic@uva.nl).

Proceedings of the 17th International Conference on Machine Learning, Austin, TX, USA, PMLR 106, 2024. Copyright 2024 by the authors.

problem has been currently gaining significant attention in machine learning. One of the main reasons is that it provides a way to compare non-distributional data and an estimate of their distance. In this paper, we propose a new way to compute Wasserstein barycenter (Kuhnle et al., 2023), style Transfer (Ulyanov et al., 2023), and Generative Adversarial Reinforcement Learning (Ulyanov et al., 2023). Furthermore, Learning Optimal Transport (LOT) (Kuhnle et al., 2023) is a promising direction for the practical application of OT. In this paper, we propose a new iterative method for solving this task. Early works, e.g., Cen et al. (2016), Arora et al. (2016), and others, have shown that the Wasserstein barycenter is a distribution of measure and discrete. Unfortunately, the Wasserstein barycenter is not necessarily discrete, so the barycenter may be regular, see (Lévy, 2017). Our work is to find the well-known Wasserstein barycenter in a single step. We also propose a new way to estimate the Wasserstein barycenter for challenging distributions, such as the continuous distributions.

We resort to creating a piecewise approximation procedure, g, to a local function, ϕ , that is a smooth function of the Wasserstein distance. We use a standard gradient descent with quadratic cost (Bach et al., 2022), $\text{soft-}OT$, where ϕ is the function of the Wasserstein distance.

Contributions. We lay a forward and propose a new way to estimate the Wasserstein barycenter for challenging distributions, which a priori cannot be represented by a discrete measure. We propose a new way to estimate the Wasserstein barycenter problem, with ϕ given general conditions. We prove that the proposed method is stable and converges to a global minimum.

1. We combine several Neural OT method (Kuhnle et al., 2023) with the transportation condition (1) and derive a new iterative procedure for the Wasserstein barycenter.

2. We compute soft version of the Wasserstein barycenter, and not the hard particular instances, we obtain quality results.

3. We describe the performance of our method on standard benchmarks and in the latent space of a pre-trained StyleGAN.



[Ссылка на источник](#)

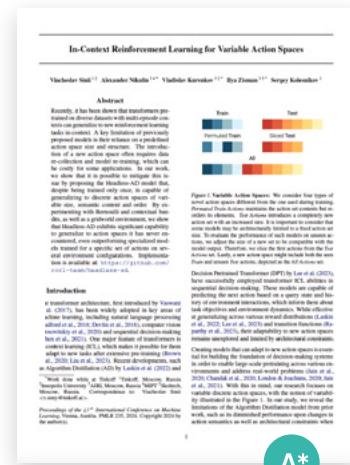


Сервисные методы

In-Context Reinforcement Learning for Variable Action Spaces

Viacheslav Sinii, Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman,
Sergey Kolesnikov

Recently, it has been shown that transformers pre-trained on diverse datasets with multi-episode contexts can generalize to new reinforcement learning tasks in-context. A key limitation of previously proposed models is their reliance on a predefined action space size and structure. The introduction of a new action space often requires data re-collection and model re-training, which can be costly for some applications. In our work, we show that it is possible to mitigate this issue by proposing the Headless-AD model that, despite being trained only once, is capable of generalizing to discrete action spaces of variable size, semantic content and order. By experimenting with Bernoulli and contextual bandits, as well as a gridworld environment, we show that Headless-AD exhibits significant capability to generalize to action spaces it has never encountered before, a specific set of actions on several environment configurations. The implementation is available at <https://github.com/airl-lab/Headless-AD>.



A*

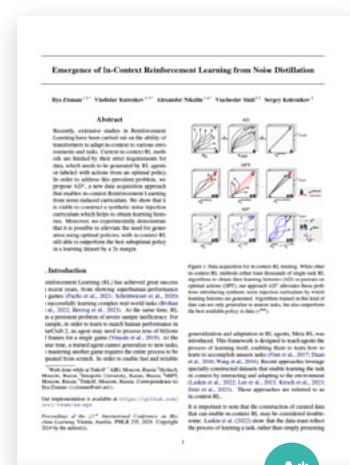


Ссылка на источник

Emergence of In-Context Reinforcement Learning from Noise Distillation

Ilya Zisman, Vladislav Kurenkov, Alexander Nikulin, Viacheslav Sinii,
Sergey Kolesnikov

Recently, extensive studies in Reinforcement Learning have been carried out on the ability of transformers to adapt in-context to various environments and tasks. Current in-context RL methods are limited by their strict requirements for data, which needs to be generated by RL agents or labeled with actions from an optimal policy. In order to address this prevalent problem, we propose AD ϵ , a new data acquisition approach that enables in-context Reinforcement Learning from noise-induced curriculum. We show that it is viable to construct a synthetic noise injection curriculum which helps to obtain learning histories. Moreover, we experimentally demonstrate that it is possible to alleviate the need for generating training data by a factor of 2x. The implementation is available at <https://github.com/airl-lab/ADepsilon>.



A*



Ссылка на источник

Disentanglement Learning via Topology

Nikita Balabin, Daria Voronkova, Ilya Trofimov, Evgeny Burnaev, Serguei Barannikov

We propose TopDis (Topological Disentanglement), a method for learning disentangled representations via adding a multi-scale topological loss term. Disentanglement is a crucial property of data representations substantial for the explainability and robustness of deep learning models and a step towards high-level cognition. The state-of-the-art methods are based on VAE and encourage the joint distribution of latent variables to be factorized. We take a different perspective on disentanglement by analyzing topological properties of data manifolds. In particular, we optimize the topological similarity for data manifolds traversals. To the best of our knowledge, our paper is the first one to propose a differentiable topological loss for disentanglement learning. Our experiments have shown that the proposed TopDis loss improves disentanglement scores such as MIG, FactorVAE score, SAP score, and DCI disentanglement score with respect to state-of-the-art results...

Disentanglement Learning via Topology

Nikita Balabin^{*}, Daria Voronkova^{**}, Ilya Trofimov^{*}, Evgeny Burnaev^{**}, Serguei Barannikov^{*}

Abstract

We propose TopDis (Topological Disentanglement), a method for learning disentangled representations via adding a multi-scale topological loss term. Disentanglement is a crucial property of data representations substantial for the explainability and robustness of deep learning models and a step towards high-level cognition. The state-of-the-art methods are based on VAE and encourage the joint distribution of latent variables to be factorized. We take a different perspective on disentanglement by analyzing topological properties of data manifolds. In particular, we optimize the topological similarity for data manifolds traversals. To the best of our knowledge, our paper is the first one to propose a differentiable topological loss for disentanglement learning. Our experiments have shown that the proposed TopDis loss improves disentanglement scores such as MIG, FactorVAE score, SAP score, and DCI disentanglement score with respect to state-of-the-art results...

bioRxiv preprint doi: https://doi.org/10.1101/2023.06.01.539262; this version posted June 1, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [aCC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

Figure 1 The TopDis pipeline process involves the following steps: 1. Input images, 2. Extract features, 3. Compute distances, 4. Compute topological loss, 5. Optimize via gradient descent. The figure shows both the original and denoised images, and highlights the topological features used for training and testing points.

A*



Ссылка на источник

Self-Supervised Coarsening of Unstructured Grid with Automatic Differentiation

Sergei Shumilin, Alexander Ryabov, Nikolay Yavich, Evgeny Burnaev, Vladimir Vanovskiy

Due to the high computational load of modern numerical simulation, there is a demand for approaches that would reduce the size of discrete problems while keeping the accuracy reasonable. In this work, we present an original algorithm to coarsen an unstructured grid based on the concepts of differentiable physics. We achieve this by employing k-means clustering, autodifferentiation and stochastic minimization algorithms. We demonstrate the efficiency of the proposed approach on two PDEs: a linear parabolic equation with periodic boundary conditions and a wave equation. Our results show that the number of grid points up to 10 times while preserving the modeled variable dynamics in the points of interest. The proposed approach can be applied to simulation of an arbitrary system described by evolutionary partial differential equations.

Self-Supervised Coarsening of Unstructured Grid with Automatic Differentiation

Sergei Shumilin^{*}, Alexander Ryabov^{*}, Nikolay Yavich^{*}, Evgeny Burnaev^{**}, Vladimir Vanovskiy^{*}

Abstract

Due to the high computational load of modern numerical simulation, there is a demand for approaches that would reduce the size of discrete problems while keeping the accuracy reasonable. In this work, we present an original algorithm to coarsen an unstructured grid based on the concepts of differentiable physics. We achieve this by employing k-means clustering, autodifferentiation and stochastic minimization algorithms. We demonstrate the efficiency of the proposed approach on two PDEs: a linear parabolic equation with periodic boundary conditions and a wave equation. Our results show that the number of grid points up to 10 times while preserving the modeled variable dynamics in the points of interest. The proposed approach can be applied to simulation of an arbitrary system described by evolutionary partial differential equations.

bioRxiv preprint doi: https://doi.org/10.1101/2023.06.01.539262; this version posted June 1, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [aCC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

Figure 1 The diagram illustrates the self-supervised coarsening process. It starts with a fine grid, which is then processed through several steps: 1. Clustering, 2. Optimization, 3. Grid refinement, 4. Grid coarsening, 5. Final coarse grid. The process is iterative, with feedback loops from the refined grid back to the clustering and optimization steps.

A*



Ссылка на источник

Light and Optimal Schrödinger Bridge Matching

Nikita Gushchin, Sergei Kholkin, Evgeny Burnaev, Alexander Korotin

Schrödinger Bridges (SB) have recently gained the attention of the ML community as a promising extension of classic diffusion models which is also interconnected to the Entropic Optimal Transport (EOT). Recent solvers for SB exploit the pervasive bridge matching procedures. Such procedures aim to recover a stochastic process transporting the mass between distributions given only a transport plan between them. In particular, given the EOT plan, these procedures can be adapted to solve SB. This fact is heavily exploited by recent works giving rise to matching-based SB solvers. The cornerstone here is recovering the EOT plan: recent works either use heuristical approximations (e.g., the minibatch OT) or establish iterative matching procedures which by the design accumulate the error during the training. We address these limitations and propose a novel procedure to learn SB which we call the optimal Schrödinger bridge matching. It exploits the optimal parameterization of the diffusion process and provably recovers the SB process (a) with a single bridge matching step and (b) with arbitrary transport plan as the input.

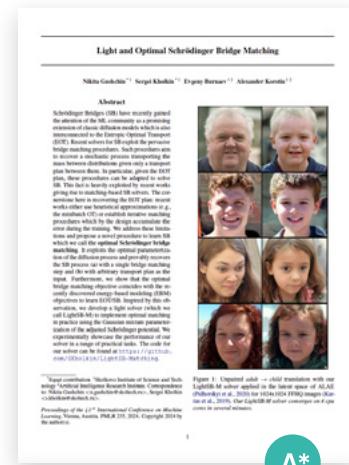


Figure 1. Unpaired adult → child translation with our LightSMB solver applied in the latent space of AAE. The figure shows four pairs of images: (top-left) original adult and child; (top-right) result from a previous solver; (bottom-left) result from our LightSMB solver; (bottom-right) result from a baseline solver. The results show significant improvements in quality and consistency.

A*



Ссылка на источник

Benchmarking and Data Synthesis for Colorization of Manga Sequential Pages for Augmented Reality

Maksim Golyadkin, Sergey Saraev, Ilya Makarov

This paper introduces an innovative approach to manga colorization within augmented reality (AR) environments, focusing on the unique challenges posed by colorizing photos of manga books. We present a novel method using diffusion models to generate a synthetic dataset that accurately replicates photographed mangapages. Additionally, we have compiled a dataset of real manga photographs, capturing diverse environmental conditions. Integrating these datasets, we established a comprehensive benchmark to evaluate colorization models in scenarios that simulate AR applications. This benchmark was validated through a human study, confirming the accuracy of our metrics across both datasets. We also showed that domain adaptation may improve model performance. Paving the way for practical applications, our framework enables...



A*



[Ссылка на источник](#)

Pose Networks Unveiled: Bridging the Gap for Monocular Depth Perception

Yazan Dayoub, Anrey V. Savchenko, Ilya Makarov

Depth estimation is essential in Augmented Reality applications, enabling realistic object placement, scene understanding, spatial mapping, interaction, and environment awareness. This paper proposes a method to enhance depth model performance without increasing inference costs by improving the pose network in a self-supervised learning setup. In particular, we enrich spatial information in the pose network by incorporating features from different scales and normalized coordinates. It is experimentally shown on the KITTI dataset that our approach achieves a 2-7% improvement in the abs rel metric when compared to baseline techniques.



A*



[Ссылка на источник](#)

BABILong: Testing the Limits of LLMs with Long Context Reasoning-in-a-Haystack

Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Sorokin, Mikhail Burtsev

In recent years, the input context sizes of large language models (LLMs) have increased dramatically. However, existing evaluation methods have not kept pace, failing to comprehensively assess the efficiency of models in handling long contexts. To bridge this gap, we introduce the BABILong benchmark, designed to test language models' ability to reason across facts distributed in extremely long documents. BABILong includes a diverse set of 20 reasoning tasks, including fact chaining, simple induction, deduction, counting, and handling lists/sets. These tasks are challenging on their own, and even more demanding when the required facts are scattered across long natural text. Our evaluations show that popular LLMs effectively utilize only 10-20% of the context and their performance declines sharply with increased reasoning complexity.

BABILong: Testing the Limits of LLMs with Long Context Reasoning-in-a-Haystack

Tiari Kuratov^{1,2}, Aydar Bulatov¹, Petr Anokhin¹, Ivan Rodkin¹, Dmitry Sorokin¹, Artyom Sorokin¹, Mikhail Burtsev¹

¹AIIR, Moscow, Russia; ²Natural Networks and Deep Learning Lab, MPFI, University, Russia

²London Institute for Mathematical Sciences, London, UK
E-mail: tiari.kuratov@aiir.ru, artyom.sorokin@aiir.ru

Abstract

In recent years, the input context sizes of large language models (LLMs) have increased dramatically. However, existing evaluation methods have not kept pace, failing to comprehensively assess the efficiency of models in handling long contexts. To bridge this gap, we introduce the BABILong benchmark, designed to test language models' ability to reason across facts distributed in extremely long documents. BABILong includes a diverse set of 20 reasoning tasks, including fact chaining, simple induction, deduction, counting, and handling lists/sets. These tasks are challenging on their own, and even more demanding when the required facts are scattered across long natural text. Our evaluations show that popular LLMs effectively utilize only 10-20% of the context and their performance declines sharply with increased reasoning complexity. Retrieval-Augmented Generation methods achieve a moderate performance, while models trained with long context reasoning methods outperform them. Among context extension methods, the highest performance is demonstrated by the ones that use a combination of multiple reasoning methods, such as chain reasoning up to 50 million tokens. The BABILong benchmark is available to try long context reasoning with increased capabilities, and we provide splits up to 10 million token length.

I. Introduction

Today, large language models (LLMs) and neural architectures are continuously evolving and achieving new milestones in natural language processing (NLP). The ability of these models to process and generate text based on the context of the input sentence is one of the most important features of NLP. This context information is used by the model to condition its output, leading to more accurate, contextually relevant, and appropriate responses. The context of a sentence is defined as the surrounding words and phrases that provide more context about the meaning of the sentence, including its history, or changing requirements in context of different parts of the sentence.

Despite these advances in model capabilities, the benchmarks used to evaluate them have not kept pace. Most benchmarks focus on short context reasoning, such as the GLUE and SuperGLUE benchmarks, which can only handle contexts up to 40,000 tokens, while models are capable of hundreds of thousands and millions of tokens. The BABILong benchmark aims to address this gap by providing a new context reasoning task that requires reasoning across long natural text, up to 50 million tokens. As a consequence, synthetic benchmarks focusing on variations of "needle-in-a-haystack"

30th Conference on Neural Information Processing Systems (NeurIPS 2024) Track on Datasets and Benchmarks.



Ссылка на источник

ENOT: Expectile Regularization for Fast and Accurate Training of Neural Optimal Transport

Nazar Buzun, Maksim Bobrin, Dmitry V. Dylov

We present a new extension for Neural Optimal Transport (NOT) training procedure, capable of accurately and efficiently estimating optimal transportation plan via specific regularisation on conjugate potentials. The main bottleneck of existing NOT solvers is associated with the procedure of finding a near-exact approximation of the conjugate operator (i.e., the c-transform), which is done either by optimizing over maximin objectives or by computationally-intensive fine-tuning. We propose a new, theoretically justified loss in the form of expectile regularization that enforces binding conditions on the learning dual potentials. Such a regularization provides the upper bound estimation over the distribution of possible conjugate potentials, which makes the learning stable and removes the need for additional extensive fine-tuning. Proposed method, called Expectile-Regularised Neural Optimal Transport (ENOT), outperforms previous state-of-the-art approaches on the established Wasserstein-2 benchmark tasks by up to 3-fold improvement in quality and up to a 10-fold improvement in runtime. Moreover, we observe performance of ENOT for varying regularization scales much like image generation, showing robustness of proposed algorithm.

ENOT: Expectile Regularization for Fast and Accurate Training of Neural Optimal Transport

Nazar Buzun^{*}
AIIR, <https://aiir.ru>
Maksim Bobrin^{*}
AIIR, <https://aiir.ru>
Dmitry V. Dylov
Moscow, <https://mpfi.ru>
^{*}<https://mpfi.ru>

Abstract

We present a new approach to Neural Optimal Transport (NOT) training procedure capable of accurately and efficiently estimating optimal transportation plan via specific regularization on dual conjugate potentials. The main bottleneck of existing NOT solvers is associated with the procedure of finding a near-exact approximation of the conjugate operator (i.e., the c-transform), which is done either by optimizing over maximin objectives or by computationally-intensive fine-tuning. We propose a new, theoretically justified loss in the form of expectile regularization that enforces binding conditions on the learning dual potentials. Such a regularization provides the upper bound estimation over the distribution of possible conjugate potentials, which makes the learning stable and removes the need for additional extensive fine-tuning. Proposed method, called Expectile-Regularised Neural Optimal Transport (ENOT), outperforms previous state-of-the-art approaches on the established Wasserstein-2 benchmark tasks by up to 3-fold improvement in quality and up to a 10-fold improvement in runtime. Moreover, we observe performance of ENOT for varying regularization scales much like image generation, showing robustness of proposed algorithm.

I. Introduction

Computational optimal transport (OT) has enabled machine learning (ML) research by offering a new way to measure distances between probability distributions. In 2009, Villani et al. [2009], Ambrosio et al. [2008], and Santambrogio [2016] have shown that the theory of a duality metric between measures, or as a generative model defined by the dual representation of the OT problem, can be applied to various ML tasks. Since then, many other generative modeling approaches and GANs, Normalizing Flows, Diffusion Models, their variants, and many others have been proposed. The success of these approaches can be attributed to the prior-wise aligned, admitting numerous applications of optimal transportation theory. Most of these approaches are based on the Sinkhorn algorithm [Sinkhorn and Knopp, 1949; Cuturi, 2013; Eslami et al., 2019; Esmeralda et al., 2019; Li, 2019; Kondor et al., 2019; Kondor et al., 2021; Kondor et al., 2022].

© The authors. NeurIPS 2024. Licensee: AIIR, Moscow, Russia.



Ссылка на источник

NeurIPS

Light Unbalanced Optimal Transport

Milena Gazdieva, Arip Asadulaev, Alexander Korotin, Evgeny Burnaev

While the continuous Entropic Optimal Transport (EOT) field has been actively developing in recent years, it became evident that the classic EOT problem is prone to different issues like the sensitivity to outliers and imbalance of classes in the source and target measures. This fact inspired the development of solvers that deal with the unbalanced EOT (UEOT) problem – the generalization of EOT allowing for mitigating the mentioned issues by relaxing the marginal constraints. Surprisingly, it turns out that the existing solvers are either based on heuristic principles or heavy-weighted with complex optimization objectives involving several neural networks. We address this challenge and propose a novel theoretically-justified, lightweight, unbalanced EOT solver. Our advancement consists of developing a novel view on the optimization of the UEOT problem yielding tractable and a non-minimax optimization objective. We show that combined with a light parametrization recently proposed in the field our objective leads to a fast, simple, and effective solver which allows solving the continuous UEOT problem in minutes on CPU.

∇^2 DFT: A Universal Quantum Chemistry Dataset of Drug-Like Molecules and a Benchmark for Neural Network Potentials

Kuzma Khrabrov, Anton Ber, Artem Tsybin, Konstantin Ushenin, Egor Rumiantsev, Alexander Telepov, Dmitry Protasov, Ilya Shenbin, Anton Alekseev, Mikhail Shirokikh, Sergey Nikolenko, Elena Tutubalina, Artur Kadurin

Methods of computational quantum chemistry provide accurate approximations of molecular properties crucial for computer-aided drug discovery and other areas of chemical science. However, high computational complexity limits the scalability of their applications. Neural network potentials (NNPs) are a promising alternative to quantum chemistry methods, but they require large and diverse datasets for training. This work presents a new dataset and benchmark called $\nabla^2\text{DFT}$ that is based on the nablaDFT. It contains twice as much molecular structures, three times more conformations, new data types and tasks, and state-of-the-art models. The dataset includes energies, forces, 17 molecular properties, Hamiltonian and overlap matrices, and a wavefunction object.

Light Unbalanced Optimal Transport



Ссылка на источник

V²DFT: A Universal Quantum Chemistry Dataset of Drug-Like Molecules and a Benchmark for Neural Network Potentials



СОЛНЦЕ НА МОСТНИКИ

Rethinking Optimal Transport in Offline Reinforcement Learning

Arip Asadulaev, Alexander Korotin, Vage Egiazarian, Rostislav Korst, Andrey Filchenkov, Evgeny Burnaev

We present a novel approach for offline reinforcement learning that bridges the gap between recent advances in neural optimal transport and reinforcement learning algorithms. Our key idea is to compute the optimal transport between states and actions with an action-value cost function and implicitly recover an optimal map that can serve as a policy. Building on this concept, we develop a new algorithm called Extremal Monge Reinforcement Learning that treats offline reinforcement learning as an extremal optimal transport problem. Unlike previous transport-based offline reinforcement learning algorithms, our method focuses on improving the policy beyond the behavior policy, rather than addressing the distribution shift problem. We evaluated the performance of our method on various continuous control problems and demonstrated improvements over existing algorithms.



Ссылка на источник

RClicks: Realistic Click Simulation for Benchmarking Interactive Segmentation

Anton Antonov, Andrey Moskalenko, Denis Shepelev, Vlad Shakhuro, Alexander Krapukhin, Konstantin Soshin, Anton Konushin

The emergence of Segment Anything (SAM) sparked research interest in the field of interactive segmentation, especially in the context of image editing tasks and speeding up data annotation. Unlike common semantic segmentation, interactive segmentation methods allow users to directly influence their output through prompts (e.g. clicks). However, click patterns in real-world interactive segmentation scenarios remain largely unexplored. Most methods rely on the assumption that users will click in the center of the largest erroneous area. Nevertheless, recent studies show that this is not always the case. Thus, methods may have poor performance in real-world deployment despite high metrics in a baseline benchmark. To accurately simulate real-user clicks, we conducted a large crowdsourcing study of click patterns in an interactive segmentation scenario and collected 475K real-user clicks.



Ссылка на источник

NeurIPS

Adversarial Schrödinger Bridge Matching

Nikita Gushchin, Daniil Selikhanovich, Sergei Kholkin, Evgeny Burnaev, Alexander Korotin

The Schrödinger Bridge (SB) problem offers a powerful framework for combining optimal transport and diffusion models. A promising recent approach to solve the SB problem is the Iterative Markovian Fitting (IMF) procedure, which alternates between Markovian and reciprocal projections of continuous-time stochastic processes. However, the model built by the IMF procedure has a long inference time due to using many steps of numerical solvers for stochastic differential equations. To address this limitation, we propose a novel Discrete-time IMF (D-IMF) procedure in which learning of stochastic processes is replaced by learning just a few transition probabilities in discrete time. Its great advantage is that in practice it can be naturally implemented using the Denoising Diffusion GAN (DD-GAN), an already well-established adversarial generative modeling technique. We show that our D-IMF procedure can provide the same quality of unpaired domain translation as the IMF, using only several generation steps instead of hundreds.

Adversarial Schrödinger Bridge Matching

Nikita Gushchin¹, Daniil Selikhanovich², Sergei Kholkin³, Evgeny Burnaev⁴, Alexander Korotin⁵
1 Moscow Institute of Physics and Technology, 2 Moscow Institute of Physics and Technology, 3 Moscow Institute of Physics and Technology, 4 Moscow Institute of Physics and Technology, 5 Moscow Institute of Physics and Technology

Abstract
The Schrödinger Bridge (SB) problem offers a general framework for combining optimal transport and diffusion models. A promising recent approach to solve the SB problem is the Iterative Markovian Fitting (IMF) procedure, which alternates between Markovian and reciprocal projections of continuous-time stochastic processes. However, the model built by the IMF procedure has a long inference time due to using many steps of numerical solvers for stochastic differential equations. To address this limitation, we propose a novel Discrete-time IMF (D-IMF) procedure in which learning of stochastic processes is replaced by learning just a few transition probabilities in discrete time. Its great advantage is that in practice it can be naturally implemented using the Denoising Diffusion GAN (DD-GAN), an already well-established adversarial generative modeling technique. We show that our D-IMF procedure can provide the same quality of unpaired domain translation as the IMF, using only several generation steps instead of hundreds. We provide the code at <https://github.com/NeurIPS2024/ASBM>.

¹State-of-the-art
²Russia
³Russia
⁴Russia
⁵Russia

NeurIPS Conference on Neural Information Processing Systems (NeurIPS 2024)

A*



Ссылка на источник

EAI: Emotional Decision-Making of LLMs in Strategic Games and Ethical Dilemmas

Mikhail Mozikov, Nikita Severin, Valeria Bodishtianu, Maria Glushanina, Ivan Nasonov, Daniil Orekhov, Vladislav Pekhotin, Ivan Makovetskiy, Mikhail Baklashkin, Vasily Lavrentyev, Akim Tsvigun, Denis Turdakov, Tatiana Shavrina, Andrey Savchenko, Ilya Makarov

One of the urgent tasks of artificial intelligence is to assess the safety and alignment of large language models (LLMs) with human behavior. Conventional verification only in pure natural language processing benchmarks can be insufficient. Since emotions often influence human decisions, this paper examines LLM alignment in complex strategic and ethical environments, providing an in-depth analysis of the emotional impact on decision-making. We introduce the novel EAI framework for integrating emotion modeling into LLMs to examine the emotional impact on ethics and decision-making in various strategic games, including bargaining and repeated games. Our experimental study with various LLMs demonstrated that emotions significantly alter LLM-based decision-making, particularly in strategic games. Highlighting the need for robust mechanisms to examine emotional bias of LLMs, we conducted a series of experiments to analyze the emotional bias of LLMs influenced by model size, alignment strength, and primary emotional responses, occasionally leading to unexpected drops in cooperation levels. This research provides a comprehensive framework for EAI, allowing us to rigorously evaluate the degree of emotional alignment. Our framework provides a fundamental basis for developing such benchmarks.

EAI: Emotional Decision-Making of LLMs in Strategic Games and Ethical Dilemmas

Mikhail Mozikov¹, Nikita Severin², Valeria Bodishtianu³, Maria Glushanina⁴, Ivan Nasonov⁵, Daniil Orekhov⁶, Vladislav Pekhotin⁷, Ivan Makovetskiy⁸, Mikhail Baklashkin⁹, Vasily Lavrentyev¹⁰, Akim Tsvigun¹¹, Denis Turdakov¹², Tatiana Shavrina¹³, Andrey Savchenko¹⁴, Ilya Makarov¹⁵
¹AI&ML, ²AI&ML, ³AI&ML, ⁴AI&ML, ⁵AI&ML, ⁶AI&ML, ⁷AI&ML, ⁸AI&ML, ⁹AI&ML, ¹⁰AI&ML, ¹¹AI&ML, ¹²AI&ML, ¹³AI&ML, ¹⁴AI&ML, ¹⁵AI&ML

Abstract
One of the urgent tasks of artificial intelligence is to assess the safety and alignment of large language models (LLMs) with human behavior. Conventional verification only in pure natural language processing benchmarks can be insufficient. Since emotions often influence human decisions, this paper examines LLM alignment in complex strategic and ethical environments, providing an in-depth analysis of the emotional impact on decision-making. We introduce the novel EAI framework for integrating emotion modeling into LLMs to examine the emotional impact on ethics and decision-making in various strategic games, including bargaining and repeated games. Our experimental study with various LLMs demonstrated that emotions significantly alter LLM-based decision-making, particularly in strategic games. Highlighting the need for robust mechanisms to examine emotional bias of LLMs, we conducted a series of experiments to analyze the emotional bias of LLMs influenced by model size, alignment strength, and primary emotional responses, occasionally leading to unexpected drops in cooperation levels. This research provides a comprehensive framework for EAI, allowing us to rigorously evaluate the degree of emotional alignment. Our framework provides a fundamental basis for developing such benchmarks.

¹Introduction
As LLMs become increasingly prevalent across various sectors – including healthcare, commerce, and entertainment – their ability to make informed decisions based on emotional acuity is increasingly important and appears to be on the cusp of regulatory, ethical, and technological affairs. The LLMs are trained on massive amounts of data, which can lead to significant biases and inconsistencies in their decision-making. These biases can manifest in various ways, such as gender, race, and cultural stereotypes, which can lead to discriminatory outcomes. Therefore, it is crucial to develop methods to detect and mitigate these biases, ensuring that the LLMs align with human values and maintain ethical standards.

²Conclusion
In this paper, we introduced the novel EAI framework for integrating emotion modeling into LLMs to examine the emotional impact on ethics and decision-making in various strategic games, including bargaining and repeated games. Our experimental study with various LLMs demonstrated that emotions significantly alter LLM-based decision-making, particularly in strategic games. Highlighting the need for robust mechanisms to examine emotional bias of LLMs, we conducted a series of experiments to analyze the emotional bias of LLMs influenced by model size, alignment strength, and primary emotional responses, occasionally leading to unexpected drops in cooperation levels. This research provides a comprehensive framework for EAI, allowing us to rigorously evaluate the degree of emotional alignment. Our framework provides a fundamental basis for developing such benchmarks.

NeurIPS Conference on Neural Information Processing Systems (NeurIPS 2024)

A*



Ссылка на источник

Energy-Guided Continuous Entropic Barycenter Estimation for General Costs

Alexander Kolesov, Petr Mokrov, Igor Udovichenko, Milena Gazdieva, Gudmund Pammer, Evgeny Burnaev, Alexander Korotin

Optimal transport (OT) barycenters are a mathematically grounded way of averaging probability distributions while capturing their geometric properties. In short, the barycenter task is to take the average of a collection of probability distributions w.r.t. given OT discrepancies. We propose a novel algorithm for approximating the continuous Entropic OT (EOT) barycenter for arbitrary OT cost functions. Our approach is built upon the dual reformulation of the EOT problem based on weak OT, which has recently gained the attention of the ML community. Beyond its novelty, our method enjoys several advantageous properties: (i) we establish quality bounds for the recovered solution; (ii) this approach seamlessly interconnects with the Energy-Based Models (EBMs) learning procedure enabling the use of well-tuned algorithms for the problem of interest; (iii) it provides an intuitive optimization scheme...

Energy-Guided Continuous Entropic Barycenter Estimation for General Costs

Alexander Kolesov, Petr Mokrov, Igor Udovichenko, Milena Gazdieva
 Skolkovo Institute of Science and Technology
 Moscow, Russia
 (a.kolesov, p.mokrov, i.udovichenko, m.gazdieva)@skoltech.ru

Gudmund Pammer
 ETH Zurich
 Zürich, Switzerland
 gudmund.pammer@ethz.ch

Alexander Korotin
 University of McLean
 Ontario, Canada
 korotin@csd.uwaterloo.ca

Evgeny Burnaev & Alexander Korotin
 Institute of Mathematics and Cryptology
 Artificial Intelligence Research Institute
 Warsaw University
 Warsaw, Poland
 evgeny.burnaev@im.uj.edu.pl, alexander.korotin@im.uj.edu.pl

Abstract
 Optimal transport (OT) barycenters are a mathematically grounded way of averaging probability distributions while capturing their geometric properties. In short, the barycenter task is to take the average of a collection of probability distributions w.r.t. given OT discrepancies. We propose a novel algorithm for approximating the continuous Entropic OT (EOT) barycenter for general OT cost functions. The approach is built upon the dual reformulation of the EOT problem based on weak OT, which has recently gained the attention of the ML community. Beyond its novelty, our method enjoys several advantages: (i) we establish quality bounds for the recovered solution; (ii) this approach seamlessly interconnects with the Energy-Based Models (EBMs) learning procedure enabling the use of well-tuned algorithms for the problem of interest; (iii) it provides an intuitive optimization scheme...
 In this paper, we consider several low-dimensional scenarios and image-space settings. Finally, we demonstrate the potential of the proposed approach for a practical task of learning the barycenter on an image manifold generated by a pretrained generative model, opening up new directions of real-world applications.

Preprint. Under review.

A*



Ссылка на источник

Improving the Worst-Case Bidirectional Communication Complexity for Nonconvex Distributed Optimization under Function Similarity

Kaja Gruntkowska, Alexander Tyurin, Peter Richtárik

Effective communication between the server and workers plays a key role in distributed optimization. In this paper, we focus on optimizing the server-to-worker communication, uncovering inefficiencies in prevalent downlink compression approaches. Considering first the pure setup where the uplink communication costs are negligible, we introduce MARINA-P, a novel method for downlink compression, employing a collection of correlated compressors. Theoretical analyses demonstrates that MARINA-P with permutation compressors can achieve a server-to-worker communication complexity improving with the number of workers, thus being provably superior to existing algorithms. We further show that MARINA-P can serve as a starting point for extensions such as methods supporting bidirectional compression. We introduce M3, a method combining MARINA-P with uplink compression and a momentum step, achieving bidirectional compression with provable improvements in total communication complexity as the number of workers increases. Theoretical findings align closely with empirical experiments, underscoring the efficiency of the proposed algorithms.

Improving the Worst-Case Bidirectional Communication Complexity for Nonconvex Distributed Optimization under Function Similarity

Kaja Gruntkowska¹ Alexander Tyurin² Peter Richtárik¹
 KAUST¹ KAUST, KACST, Saudi Arabia² Microsoft Research, Microsoft, USA

Abstract
 Effective communication between the server and workers plays a key role in distributed optimization. In this paper, we focus on optimizing communication, uncovering inefficiencies in prevalent downlink compression approaches. Considering first the pure setup where the uplink communication costs are negligible, we introduce MARINA-P, a novel method for downlink compression, employing a collection of correlated compressors. Theoretical analyses demonstrates that MARINA-P with permutation compressors can achieve a server-to-worker communication complexity improving with the number of workers, thus being provably superior to existing algorithms. We further show that MARINA-P can serve as a starting point for extensions such as methods supporting bidirectional compression. We introduce M3, a method combining MARINA-P with uplink compression and a momentum step, achieving bidirectional compression with provable improvements in total communication complexity as the number of workers increases. Theoretical findings align closely with empirical experiments, underscoring the efficiency of the proposed algorithms.

1 Introduction
 In federated learning (McMahan et al., 2017; Konečný et al., 2016) and large-scale machine learning training, effective communication between the server and workers is critical for training a model. Facilitating this collaboration requires the transmission of information between the server and workers. In this paper, we focus on the server-to-worker communication framework, where communication takes place in a series. As a result, practical challenges are due to the high communication costs. In this paper, we propose a novel method for improving the communication efficiency in distributed optimization. We further show that MARINA-P can serve as a starting point for extensions such as methods supporting bidirectional compression. We introduce M3, a method combining MARINA-P with uplink compression and a momentum step, achieving bidirectional compression with provable improvements in total communication complexity as the number of workers increases. Theoretical findings align closely with empirical experiments, underscoring the efficiency of the proposed algorithms.

We consider the following setup of distributed optimization task:

$$\min_{\theta} \left\{ f(\theta) + \frac{1}{n} \sum_{i=1}^n f_i(\theta) \right\}, \quad (1)$$

where $\theta \in \mathbb{R}^d$ is the vector of parameters, n is the number of workers and $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i = 1, \dots, n$ — smooth nonconvex functions. We investigate the scenario where the workers are distributed across multiple locations and communicate with the server via some communication network.

¹King Abdullah University of Science and Technology, Thuwal, Saudi Arabia
²Microsoft Research, Microsoft, USA

© 2024 Conference on Neural Information Processing Systems (NeurIPS 2024).

A*



Ссылка на источник

XLand-MiniGrid: Scalable Meta-Reinforcement Learning Environments in JAX

Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Artem Agarkov, Viacheslav Sini, Sergey Kolesnikov

Inspired by the diversity and depth of XLand and the simplicity and minimalism of MiniGrid, we present XLand-MiniGrid, a suite of tools and grid-world environments for meta-reinforcement learning research. Written in JAX, XLand-MiniGrid is designed to be highly scalable and can potentially run on GPU or TPU accelerators, democratizing large-scale experimentation with limited resources. Along with the environments, XLand-MiniGrid provides pre-sampled benchmarks with millions of unique tasks of varying difficulty and easy-to-use baselines that allow users to quickly start training adaptive agents. In addition, we have conducted a preliminary analysis of scaling and generalization, showing that our baselines are capable of reaching millions of steps per second during training and validating that the proposed benchmarks are challenging.

XLand-MiniGrid: Scalable Meta-Reinforcement Learning Environments in JAX

Authors: Alexander Nikulin¹, Vladislav Kurenkov¹, Ilya Zisman¹, Artem Agarkov², Viacheslav Sini³, Sergey Kolesnikov⁴

Abstract: Inspired by the diversity and depth of XLand and the simplicity and minimalism of MiniGrid, we present XLand-MiniGrid, a suite of tools and grid-world environments for meta-reinforcement learning research. Written in JAX, XLand-MiniGrid is designed to be highly scalable and can potentially run on GPU or TPU accelerators, democratizing large-scale experimentation with limited resources. Along with the environments, XLand-MiniGrid provides pre-sampled benchmarks with millions of unique tasks of varying difficulty and easy-to-use baselines that allow users to quickly start training adaptive agents. In addition, we have conducted a preliminary analysis of scaling and generalization, showing that our baselines are capable of reaching millions of steps per second during training and validating that the proposed benchmarks are challenging.

<https://github.com/kolesnikovsl/XLand-MiniGrid>



Ссылка на источник

Shadowheart SGD: Distributed Asynchronous SGD with Optimal Time Complexity Under Arbitrary Computation and Communication Heterogeneity

Alexander Tyurin, Marta Pozzi, Ivan Ilin, Peter Richtárik

We consider nonconvex stochastic optimization problems in the asynchronous centralized distributed setup where the communication times from workers to a server can not be ignored, and the computation and communication times are potentially different for all workers. Using an unbiased compression technique, we develop a new method-Shadowheart SGD-that provably improves the time complexities of all previous centralized methods. Moreover, we show that the time complexity of Shadowheart SGD is optimal in the family of centralized methods with compressed communication, where broadcasting from the server to the workers is non-negligible, and develop a corresponding method.

Shadowheart SGD: Distributed Asynchronous SGD with Optimal Time Complexity Under Arbitrary Computation and Communication Heterogeneity

Authors: Alexander Tyurin¹, Marta Pozzi², Ivan Ilin³, Peter Richtárik⁴

Abstract: We consider nonconvex stochastic optimization problems in the asynchronous centralized distributed setup where the communication times from workers to a server can not be ignored, and the computation and communication times are potentially different for all workers. Using an unbiased compression technique, we develop a new method-Shadowheart SGD—that provably improves the time complexity of all previous centralized methods. Moreover, we show that the time complexity of Shadowheart SGD is optimal in the family of centralized methods with compressed communication, where broadcasting from the server to the workers is non-negligible, and developing a corresponding method.

<https://nips.cc/WBSITE-7/3337>



Ссылка на источник

Optimal Flow Matching: Learning Straight Trajectories in Just One Step

Nikita Kornilov, Petr Mokrov, Alexander Gasnikov, Alexander Korotin

Over the several recent years, there has been a boom in development of Flow Matching (FM) methods for generative modeling. One intriguing property pursued by the community is the ability to learn flows with straight trajectories which realize the Optimal Transport (OT) displacements. Straightness is crucial for the fast integration (inference) of the learned flow's paths. Unfortunately, most existing flow straightening methods are based on non-trivial iterative FM procedures which accumulate the error during training or exploit heuristics based on minibatch OT. To address these issues, we develop and theoretically justify the novel Optimal Flow Matching approach which allows recovering the straight OT displacement for the quadratic transport in just one FM step. The main idea of our approach is the employment of vector field for FM which are parameterized by convex functions.

**Optimal Flow Matching:
Learning Straight Trajectories in Just One Step**

Nikita Kornilov
Skolkovo Institute of Science and Technology
R Center for AI, Skolkovo University
Moscow, Russia
nik.kornilov@skoltech.ru

Petr Mokrov
Skolkovo Institute of Science and Technology
petr.mokrov@skoltech.ru

Alexander Gasnikov
Moscow Institute of Physics and Technology
Moscow, Russia
gasnikov@mipt.ru

Alexander Korotin
Skolkovo Institute of Science and Technology
Artificial Intelligence Research Institute
Moscow, Russia
korotin@skoltech.ru

Abstract

Over the several recent years, there has been a boom in development of Flow Matching (FM) methods for generative modeling. One intriguing property pursued by the community is the ability to learn flows with straight trajectories which realize the Optimal Transport (OT) displacements. Straightness is crucial for the fast integration (inference) of the learned flow's paths. Unfortunately, most existing flow straightening methods are based on non-trivial iterative FM procedures which accumulate the error during training or exploit heuristics based on minibatch OT. To address these issues, we develop and theoretically justify the novel Optimal Flow Matching approach which allows recovering the straight OT displacement for the quadratic transport in just one FM step. The main idea of our approach is the employment of vector field for FM which are parameterized by convex functions. The code for the implementation of the proposed method is available at <https://github.com/nik-kornilov/Optimal-Flow-Matching>.

1 Introduction

Recent success in generative modeling [Luc et al., 2016], [Ever et al., 2016], [Caro et al., 2016] is mostly due to the development of Flow Matching (FM) methods [Kondor et al., 2016]. These methods reduce the task to a large one via ordinary differential equations (ODEs) describing the mass movement. However, they are not able to learn straight trajectories which are crucial for the fast inference of the learned sampling. To address this issue, researchers developed several improvements of the FM [Li et al., 2017], [Li et al., 2018], [Liu et al., 2019], [Liu et al., 2020], [Liu et al., 2021], [Liu et al., 2022], [Liu et al., 2023].

Revised Flow (RF) method [Li et al., 2021], [Liu et al., 2022] iteratively solves FM and gradually refines straightness. It is a good improvement of the method. The other popular branch of research is the Optimal Transport (OT) [Villani, 2009]. The main goal of OT is to find the way to move one probability distribution to another. In the context of generative modeling, OT is used to learn straight trajectories [Gasnikov et al., 2019], [Gasnikov et al., 2020]. The authors propose OTFM or OT Conditional Flow Matching (OTCFM) [Vashist et al., 2020], [Vashist et al., 2021]. The authors propose OTFM or OTCFM because both considered

© 2024 NeurIPS. All Rights Reserved.

Ссылка на источник

On the Optimal Time Complexities in Decentralized Stochastic Asynchronous Optimization

Alexander Tyurin, Peter Richtárik

We consider the decentralized stochastic asynchronous optimization setup, where many workers asynchronously calculate stochastic gradients and asynchronously communicate with each other using edges in a multigraph. For both homogeneous and heterogeneous setups, we prove new time complexity lower bounds under the assumption that computation and communication speeds are bounded. We develop a new nearly optimal method, Fragile SGD, and a new optimal method, Amelie SGD, that converge under arbitrary heterogeneous computation and communication speeds and match our lower bounds (up to a logarithmic factor in the homogeneous setting). Our time complexities are new, nearly optimal, and provably improve all previous asynchronous/stochastic methods in the decentralized setup.

On the Optimal Time Complexities in Decentralized Stochastic Asynchronous Optimization

Alexander Tyurin
KAUST, ARRI, Heriot-Watt*

Peter Richtárik
KAUST*

Abstract

We consider the decentralized stochastic asynchronous optimization setup, where many workers asynchronously calculate stochastic gradients and asynchronously communicate with each other using edges in a multigraph. For both homogeneous and heterogeneous setups, we prove new time complexity lower bounds under the assumption that computation and communication speeds are bounded. We develop a new nearly optimal method, Fragile SGD, and a new optimal method, Amelie SGD, that converge under arbitrary heterogeneous computation and communication speeds and match our lower bounds (up to a logarithmic factor in the homogeneous setting). Our time complexities are new, nearly optimal, and provably improve all previous asynchronous/stochastic methods in the decentralized setup.

1 Introduction

We consider the smooth consensus optimization problem

$$\min_{x \in \mathbb{R}^n} \left\{ f(x) + \sum_{i=1}^n \psi_i(x_i) \right\}, \quad (1)$$

where $f: \mathbb{R}^n \rightarrow \mathbb{R}$, ψ_i is a distribution on a non-empty set S_i . For a given $j \in \mathcal{N}$, we want to minimize ψ_j over S_j . We assume that the workers are in a network and the network is connected. We analyze the heterogeneous setup and the convex setup with smooth and non-smooth functions in this paper.

1.1 Decentralized setup with times

We investigate the following decentralized asynchronous setup. Assume that we have n workers indexed by $i \in \{1, \dots, n\}$ which wants to compute a stochastic gradient to the i^{th} node, and a loss of equal size ψ_i . Let $\mathcal{E} = \{e_{ij}\}_{i,j \in \mathcal{N}}$ be a directed graph connecting nodes i and j , such that $e_{ij} \in \mathcal{E}$ if and only if $i \neq j$ and $i, j \in \mathcal{N}$. We consider the following setup. At time t , each worker i performs one of the following actions: (i) it performs a local computation (gradient calculation), (ii) it performs a local communication (gradient transmission), and (iii) it performs a global communication (gradient aggregation). The computation and communication times can be arbitrary heterogeneous and the communication times can be arbitrary homogeneous. In this paper, we assume that the computation times are bounded and the communication times are bounded. See Section 3.1, we explain that our result can easily extend to the case where the upper bounds are not necessarily constant.

We consider any weighted directed multigraph parameterized by a vector $\lambda \in \mathbb{R}^m$ such that $\lambda_i \in [0, 1]$ for all $i \in \mathcal{N}$ and $\lambda_{i,j} \geq 0$ for all $i, j \in \mathcal{N}$ such that $e_{ij} \in \mathcal{E}$.

*King Abdullah University of Science and Technology, Thuwal, Saudi Arabia

**Moscow Institute of Physics and Technology, Moscow, Russia

© 2024 Conference on Neural Information Processing Systems (NeurIPS 2024).

Ссылка на источник

AIRI — результаты 2024 года

63

Freya PAGE: First Optimal Time Complexity for Large-Scale Nonconvex Finite-Sum Optimization with Heterogeneous Asynchronous Computations

Alexander Tyurin, Kaja Gruntkowska, Peter Richtárik

In practical distributed systems, workers are typically not homogeneous, and due to differences in hardware configurations and network conditions, can have highly varying processing times. We consider smooth nonconvex finite-sum (empirical risk minimization) problems in this setup and introduce a new parallel method, Freya PAGE, designed to handle arbitrarily heterogeneous and asynchronous computations. By being robust to «stragglers» and adaptively ignoring slow computations, Freya PAGE offers significantly improved time complexity guarantees compared to all previous methods, including Asynchronous SGD, Rennala SGD, SPIDER, and PAGE, while requiring weaker assumptions. The algorithm relies on novel generic stochastic

Freya PAGE: First Optimal Time Complexity for Large-Scale Nonconvex Finite-Sum Optimization with Heterogeneous Asynchronous Computations

Alexander Tyurin¹, Kaja Gruntkowska², Peter Richtárik¹

Abstract

In practical distributed systems, workers are typically not homogeneous, and due to differences in hardware configurations and network conditions, can have highly varying processing times. We consider smooth nonconvex finite-sum optimization problems in this setup and introduce a new parallel method, Freya PAGE, designed to handle arbitrarily heterogeneous and asynchronous computations. By being robust to «stragglers» and adaptively ignoring slow computations, Freya PAGE offers significantly improved time complexity guarantees compared to all previous methods, including Asynchronous SGD, Rennala SGD, SPIDER, and PAGE, while requiring weaker assumptions. The algorithm relies on novel generic stochastic gradient collection strategies with different parameters that can be tuned independently for each worker. Furthermore, we establish a lower bound for smooth nonconvex finite-sum optimization with heterogeneous asynchronous computations that matches our upper bound. This lower bound is tight and demonstrates the optimality of Freya PAGE in terms of scaling linearly with the number of workers, and is $\mathcal{O}(d)$ data samples.

1 Introduction

In practical distributed systems used for large-scale machine learning tasks, it is common to have workers with different hardware configurations and network conditions, leading to highly varying processing times. One can see this from GFT computation details reported in baseline configurations, network conditions, and hardware configurations in the literature [1–10]. For example, the widely used asynchronous coordinate descent (ACD) [11] and its variants [12–14] are not able to handle such heterogeneity well. In fact, these methods often significantly increase time complexity guarantees compared to their synchronous counterparts [15–17].

Asynchronous SGD [18] and its variants [19–21] are more robust to heterogeneity, but they still require some form of synchronization between workers, which limits their performance.

Due to the above issues, we aim to address the challenges posed by device heterogeneity in the context of solving finite-sum minimization problems of the form

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) + \frac{\lambda}{2} \|x\|_2^2 \right\}, \quad (1)$$

where $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is the sum of n smooth functions, and $\lambda > 0$ is a weight parameter. We assume that f is convex and has a unique minimum at zero. Our goal is to find an ϵ -stationary point, i.e., a ϵ -proximal point of f , which is defined as a point x^* such that $\|f(x^*) - f(x)\| \leq \epsilon$ for all $x \in \mathbb{R}^d$. In this paper, we focus on the case where f is nonconvex and non-smooth.

• Best: a workflow/algorithm able to work in parallel;

— King Abdullah University of Science and Technology, Thuwal, Saudi Arabia;

— Max Planck Institute for Software Systems, Saarbrücken, Germany;

— 5th Conference on Neural Information Processing Systems (NeurIPS 2024).

A*



Ссылка на источник

HairFastGAN: Realistic and Robust Hair Transfer with a Fast Encoder-Based Approach

Maxim Nikolaev, Mikhail Kuznetsov, Dmitry P. Vetrov, Aibek Alanov

Our paper addresses the complex task of transferring a hairstyle from a reference image to an input photo for virtual hair try-on. This task is challenging due to the need to adapt to various photo poses, the sensitivity of hairstyles, and the lack of objective metrics. The current state-of-the-art hairstyle transfer methods use an optimization process for different parts of the approach, making them inexcusably slow. At the same time, faster encoder-based models are of very low quality because they either operate in StyleGAN's W+ space or use other low-dimensional image generators. Additionally, both approaches have a problem with hairstyle transfer when the source pose is very different from the target pose. In our paper, we present the HairFast model, which uniquely solves these problems and compares it to optimization problem-based methods. Our solution includes a new encoder-based architecture, a novel loss function, a new training approach, and improved modules for better alignment, color matching, and a new metric for evaluating hairstyle transfer quality. We also show that our model performs well on metrics after random hairstyle transfer and reconstruction while the source and target images are very different. Our method performs best in terms of both speed and quality.

HairFastGAN: Realistic and Robust Hair Transfer with a Fast Encoder-Based Approach

Maxim Nikolaev^{1,2}, Mikhail Kuznetsov^{2,3}, Dmitry Vetrov⁴, Aibek Alanov^{1,2}
— MSU University, "Skolkovo Institute of Science and Technology", AIBI
— Skolkovo Institute of Science and Technology, Moscow, Russia
— Computer University, Moscow, Russia
— Skolkovo Institute of Science and Technology, Moscow, Russia



Figure 1: David Bowie’s face, Barack Obama’s face, and Barack Obama with HairFast transfer. The first two images are source and target images respectively. The last two images are generated by HairFast and HairFastPlus models. The HairFastPlus model performs the transfer of the original attributes. You can also see a comparison of our model with the state-of-the-art HairFastPlus model.

Abstract

Our paper addresses the complex task of transferring a hairstyle from a reference image to an input photo for virtual hair try-on. This task is challenging due to the need to adapt to various photo poses, the sensitivity of hairstyles, and the lack of objective metrics. The current state-of-the-art hairstyle transfer methods use an optimization process for different parts of the approach, making them inexcusably slow. At the same time, faster encoder-based models are of very low quality because they either operate in StyleGAN’s W+ space or use other low-dimensional image generators. Additionally, both approaches have a problem with hairstyle transfer when the source pose is very different from the target pose.

In our paper, we present the HairFast model, which uniquely solves these problems and compares it to optimization problem-based methods.

Our solution includes a new encoder-based architecture, a novel loss function, a new training approach,

and improved modules for better alignment, color matching, and a new metric for evaluating hairstyle transfer quality.

We also show that our model performs well on metrics after random hairstyle transfer and reconstruction while the source and target images are very different. Our method performs best in terms of both speed and quality.

A*



Ссылка на источник

Group and Shuffle: Efficient Structured Orthogonal Parametrization

Mikhail Gorbunov, Nikolay Yudin, Vera Soboleva, Aibek Alanov, Alexey Naumov, Maxim Rakuba

The increasing size of neural networks has led to a growing demand for methods of efficient fine-tuning. Recently, an orthogonal fine-tuning paradigm was introduced that uses orthogonal matrices for adapting the weights of a pretrained model. In this paper, we introduce a new class of structured matrices, which unifies and generalizes structured classes from previous works. We examine properties of this class and build a structured orthogonal parametrization upon it. We then use this parametrization to modify the orthogonal fine-tuning framework, improving parameter and computational efficiency. We empirically validate our method on different domains, including adapting of text-to-image diffusion models and downstream task fine-tuning in language modeling. Additionally, we adapt our construction for orthogonal convolutions and conduct experiments with 1-Lipschitz neural networks.

Group and Shuffle: Efficient Structured Orthogonal Parametrization

Mikhail Gorbunov
BSP University
gorbunovmikhail@gmail.com

Nikolay Yudin
HSE University

Vera Soboleva
AIST

Aibek Alanov
HSE University

Alexey Naumov
HSE University

Maxim Rakuba
Steklov Mathematical Institute RAS

Abstract
The increasing size of neural networks has led to a growing demand for methods of efficient fine-tuning. Recently, an orthogonal fine-tuning paradigm was introduced that uses orthogonal matrices for adapting the weights of a pretrained model. In this paper, we introduce a new class of structured matrices, which unifies and generalizes structured classes from previous works. We examine properties of this class and build a structured orthogonal parametrization upon it. We then use this parametrization to modify the orthogonal fine-tuning framework, improving parameter and computational efficiency. We empirically validate our method on different domains, including adapting of text-to-image diffusion models and downstream task fine-tuning in language modeling. Additionally, we adapt our construction for orthogonal convolutions and conduct experiments with 1-Lipschitz neural networks.

Preprint. Under review.

A*



[Ссылка на источник](#)

Neural Potential Field for Obstacle-Aware Local Motion Planning

Muhammad Alhaddad, Konstantin Mironov, Aleksey Staroverov, Aleksandr Panov

Model predictive control (MPC) may provide local motion planning for mobile robotic platforms. The challenging aspect is the analytic representation of collision cost for the case when both the obstacle map and robot footprint are arbitrary. We propose a Neural Potential Field: a neural network model that returns a differentiable collision cost based on robot pose, obstacle map, and robot footprint. The differentiability of our model allows its usage within the MPC solver.

It is computationally hard to solve problems with a very high number of parameters. Therefore, our architecture includes neural image encoders, which transform obstacle maps and robot footprints into embeddings, which reduce problem dimensionality by two orders of magnitude. The reference data for network training are generated based on algorithmic calculation of a signed distance function. Comparative experiments showed that the proposed approach is comparable with existing local planners: it provides trajectories with outperforming smoothness, comparable path length, and safe distance from obstacles. Experiment on Husky UGV mobile robot showed that our approach allows real-time and safe local planning. The code for our approach is presented at <https://github.com/cog-isa/NPField> together with demo video.



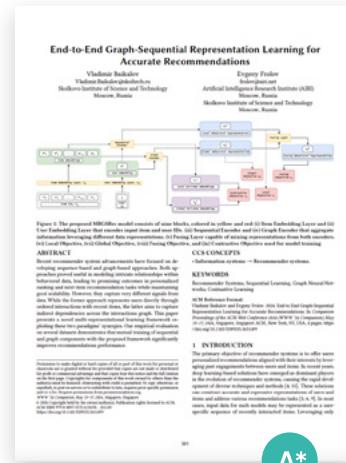
Ссылка на источник

A*

End-to-End Graph-Sequential Representation Learning for Accurate Recommendations

Vladimir Baikalov, Evgeny Frolov

Recent recommender system advancements have focused on developing sequence-based and graph-based approaches. Both approaches proved useful in modeling intricate relationships within behavioral data, leading to promising outcomes in personalized ranking and next-item recommendation tasks while maintaining good scalability. However, they capture very different signals from data. While the former approach represents users directly through ordered interactions with recent items, the latter aims to capture indirect dependencies across the interactions graph. This paper presents a novel multi-representational learning framework exploiting these two paradigms' synergies. Our empirical evaluation on several datasets demonstrates that mutual training of sequential and graph components with the proposed framework significantly improves recommendations performance.



[Ссылка на источник](#)

ACM KDD

Uplift Modelling via Gradient Boosting

Bulat Ibragimov, Anton Vakhrushev

The Gradient Boosting machine learning ensemble algorithm, well-known for its proficiency and superior performance in intricate machine learning tasks, has encountered limited success in the realm of uplift modeling. Uplift modeling is a challenging task that necessitates a known target for the precise computation of the training gradient. The prevailing two-model strategies, which separately model treatment and control outcomes, are encumbered with limitations as they fail to directly tackle the uplift problem. This paper presents an innovative approach to uplift modeling that employs Gradient Boosting. Unlike previous works, our algorithm utilizes multioutput boosting model and calculates the uplift gradient based on intermediate surrogate predictions and directly models the concealed target. This method circumvents the requirement for a known target and addresses the uplift problem more effectively than existing solutions. Moreover, we broaden the scope of this solution to encompass multitreatment settings, thereby enhancing its applicability. This novel approach not only overcomes the limitations of the traditional two-model strategies...

From Variability to Stability: Advancing RecSys Benchmarking Practices

Valeriy Shevchenko, Nikita Belousov, Alexey Vasilev, Vladimir Zholobov,
Artyom Sosedka, Natalia Semenova, Anna Volodkevich, Andrey
Savchenko, Alexey Zaytsev

In the rapidly evolving domain of Recommender Systems (RecSys), new algorithms frequently claim state-of-the-art performance based on evaluations over a limited set of arbitrarily selected datasets. However, this approach may fail to holistically reflect their effectiveness due to the significant impact of dataset characteristics on algorithm performance. Addressing this deficiency, this paper introduces a novel benchmarking methodology to facilitate a fair and robust comparison of RecSys algorithms, thereby advancing evaluation practices. By utilizing a diverse set of 30 open datasets, including two introduced in this work, and evaluating 11 collaborative filtering algorithms across 9 metrics, we critically examine the influence of dataset characteristics on algorithm performance. We further investigate the feasibility of aggregating outcomes from multiple datasets into a unified ranking. Through rigorous experimental analysis, we validate the reliability of our methodology under the variability of datasets, offering a benchmarking strategy that balances quality and computational demands. This methodology enables a fair yet effective means of evaluating RecSys algorithms, providing valuable guidance for future research endeavors.



Ссылка на источник



Ссылка на источник

ACM KDD

Learn Together Stop Apart: an Inclusive Approach to Ensemble Pruning

Bulat Ibragimov; Gleb Gusev

Gradient Boosting is a leading learning method that builds ensembles and adapts their sizes to particular tasks, consistently delivering top-tier results across various applications. However, determining the optimal number of models in the ensemble remains a critical yet underexplored aspect. Traditional approaches assume a universal ensemble size effective for all data points, which may not always hold true due to data heterogeneity. This paper introduces an adaptive approach to early stopping in Gradient Boosting, addressing data heterogeneity by assigning different stop moments to different data regions at inference time while still training a common ensemble on the entire dataset. We propose two methods: Direct Supervised Partition (DSP) and Indirect Supervised Partition (ISP). The DSP method uses a decision tree to partition the data based on learning curves, while ISP leverages the dataset's geometric and target distribution characteristics. An effective validation protocol is developed to determine the optimal number of early stopping regions or detect when the heterogeneity assumption does not hold. Experiments using state-of-the-art implementations of Gradient Boosting, LightGBM, and CatBoost, on standard benchmarks demonstrate that our methods enhance model precision by up to 2%, underscoring the significance of this research direction. This approach does not increase computational complexity and can be easily integrated into existing learning pipelines.

**Learn Together Stop Apart:
an Inclusive Approach to Ensemble Pruning**

Bulat Ibragimov
gleb.gusev@ai-lab.ru
Moscow Institute of Physics and Technology
and Gleb Gusev
Artificial Intelligence Research Institute (AIRI)
Moscow, Russia Federation

ABSTRACT
Gradient Boosting is a leading learning method that builds ensembles and adapts their sizes to particular tasks, consistently delivering top-tier results across various applications. However, determining the optimal number of models in the ensemble remains a critical yet underexplored aspect. Traditional approaches assume a universal ensemble size effective for all data points, which may not always hold true due to data heterogeneity. This paper introduces an adaptive approach to early stopping in Gradient Boosting, addressing data heterogeneity by assigning different stop moments to different data regions at inference time while still training a common ensemble on the entire dataset. We propose two methods: Direct Supervised Partition (DSP) and Indirect Supervised Partition (ISP). The DSP method uses a decision tree to partition the data based on learning curves, while ISP leverages the dataset's geometric and target distribution characteristics. An effective validation protocol is developed to determine the optimal number of early stopping regions or detect when the heterogeneity assumption does not hold. Experiments using state-of-the-art implementations of Gradient Boosting, LightGBM, and CatBoost, on standard benchmarks demonstrate that our methods enhance model precision by up to 2%, underscoring the significance of this research direction. This approach does not increase computational complexity and can be easily integrated into existing learning pipelines.

CCS CONCEPTS
• Computing methodologies → Boosting; Regularization; Classification and regression

KEYWORDS
Ensemble, Boosting, Regularization, Early Stopping, Decision Tree

ACM Reference Format:
Bulat Ibragimov and Gleb Gusev. 2024. Learn Together Stop Apart: an Inclusive Approach to Ensemble Pruning. In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '24), August 18–22, 2024, New York, NY, USA. © 2024, Association for Computing Machinery. 978-1-4503-9963-6/24/08...\$15.00.
https://doi.org/10.1145/3588244.3599026

1 INTRODUCTION
Despite the rapid advancement of deep neural networks, their performance is often limited by the lack of generalization [23, 24]. One way to address this issue is to use ensemble methods (e.g., bags, ensembles, etc.) combined with appropriate data augmentation techniques [25, 26]. Another approach is to learn different models for different parts of the data [27], which is called domain adaptation [28] or domain generalization [29].
In recent years, new optimizers and hyperparameters have been proposed to improve the performance of gradient boosting [30, 31]. One critical aspect affecting the efficiency of GR is the choice of the validation set. A common approach is to split the training data into two sets: a training set and a validation set. The validation set is used to tune the hyperparameters and select the best model. Another approach is to use cross-validation, where the data is divided into several folds, and each fold is used as a validation set. These methods are effective for homogeneous datasets, but they may not work well for heterogeneous datasets. For example, if the validation set is too small, it may not be representative of the entire dataset, leading to overfitting. Conversely, if the validation set is too large, it may not be representative of the entire dataset, leading to underfitting. To address this issue, we propose a new approach to early stopping in Gradient Boosting that takes into account the heterogeneity of the data. The proposed approach is based on the idea of partitioning the data into different regions and training separate models for each region. This allows us to train different models for different regions of the data, which can lead to better performance.

The proposed early stopping method has a significant advantage over traditional methods: it is more efficient and requires less computational resources. This is because the proposed method only trains models that are effective for a few data points. However, the proposed method also has some disadvantages. For example, it may not work well for datasets with complex structures, such as images or videos. Additionally, it may not work well for datasets with missing values or outliers. To overcome these challenges, we propose a validation protocol that checks for the heterogeneity of the data in the feature space and adapts the validation set accordingly. Specifically, we build an additional "partition model" that identifies the most important features for each region and represents them as a set of binary variables. These variables are then used to weight the data points in the validation set, so that the validation set is more representative of the overall dataset. This allows us to train different models for each region individually. The proposed approach is also more efficient than traditional methods, as it only trains models that are effective for a few data points. This leads to faster training times and lower computational costs.



Ссылка на источник

Scalar Function Topology Divergence: Comparing Topology of 3D Objects

Ilya Trofimov, Daria Voronkova, Eduard Tulchinskii, Evgeny Burnaev, Serguei Barannikov

We propose a new topological tool for computer vision Scalar Function Topology Divergence (SFTD), which measures the dissimilarity of multi-scale topology between sublevel sets of two functions having a common domain. Functions can be defined on an undirected graph or Euclidean space of any dimensionality. Most of the existing methods for comparing topology are based on Wasserstein distance between persistence barcodes and they don't take into account the localization of topological features. On the other hand, the minimization of SFTD ensures that the corresponding topological features of scalar functions are located in the same places. The proposed tool provides useful visualizations depicting areas where functions have topological dissimilarities. We provide applications of the proposed method to 3D computer vision. In particular, experiments demonstrate that SFTD improves the reconstruction of cellular 3D shapes from 2D fluorescence microscopy images, and helps to identify topological errors in 3D segmentation.

Guide-and-Rescale: Self-Guidance Mechanism for Effective Tuning-Free Real Image Editing

Vadim Titov, Madina Khalmatova, Alexandra Ivanova, Dmitry Vetrov, and Aibek Alanov

Despite recent advances in large-scale text-to-image generative models, manipulating real images with these models remains a challenging problem. The main limitations of existing editing methods are that they either fail to perform with consistent quality on a wide range of image edits or require time-consuming hyperparameter tuning or fine-tuning of the diffusion model to preserve the image-specific appearance of the input image. We propose a novel approach that is built upon a modified diffusion sampling process via the guidance mechanism. In this work, we explore the self-guidance technique to preserve the overall structure of the input image and its local regions appearance that should not be edited. In particular, we explicitly introduce layout-preserving energy functions that are aimed to save local and global structures of the source image. Additionally, we propose a noise rescaling mechanism...

Scalar Function Topology Divergence: Comparing Topology of 3D Objects

Ilya Trofimov¹, Daria Voronkova², Eduard Tulchinskii³, Evgeny Burnaev², Serguei Barannikov¹

¹ Skolkovo Institute of Science and Technology, Moscow, Russia
<https://ai.sit.edu.ru/~voronkova/>; ² Institute of Mathematics and Cryptology, Polish Academy of Sciences, Warsaw, Poland
<https://im.ii.uj.edu.pl/~burnaev/>; ³ CERN, Geneva, Switzerland

AI Foundation and Algorithm Lab, Moscow, Russia

Abstract: We propose a new topological tool for computer vision Scalar Function Topology Divergence (SFTD), which compares the similarity of multi-scale topology between sublevel sets of two functions defined on an undirected graph or Euclidean space of any dimensionality. Most of the existing methods for comparing topology are based on Wasserstein distance between persistence barcodes and they don't take into account the localization of topological features. Our proposed tool ensures that the corresponding topological features of scalar functions are located in the same places. We provide visualizations depicting areas where functions have topological dissimilarities. We provide applications of the proposed method to 3D computer vision. In particular, experiments demonstrate that SFTD is an additional loss function that improves the reconstruction of cellular 3D shapes from 2D fluorescence microscopy images, and helps to identify topological errors in 3D segmentation. The code is publicly available: <https://github.com/IL-Trofimov/SFTD>.

Keywords: Topological Data Analysis, 3D Computer Vision, Geometry, Deep Learning

1 Introduction

Automatic recognition in computer vision is a separated problem defined over 2D. Examples include image segmentation [25, 26, 27] and 3D shape reconstruction [28, 29, 30]. However, in most of the cases the recognition is limited to specific objects and does not consider the context. For example, the Dots seen, cross-seen, etc. These metrics are not capable of providing detailed information about the object's shape, such as the number of holes, faces, connectivity patterns, etc. Typically, predictions involve scalar function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and the final shape of an object is defined by thresholding $f(x)$.

A*



[Ссылка на источник](#)

Guide-and-Rescale: Self-Guidance Mechanism for Effective Tuning-Free Real Image Editing

Vadim Titov^{1*}, Madina Khalmatova², Alexandra Ivanova^{2,3,4}, Dmitry Vetrov⁵, and Aibek Alanov^{1,2}

¹ HSE University
² AIRI
³ https://github.com/airi-research/airi-mlkit.net
⁴ Microsoft Research Asia
⁵ CERN by day
a. Moshenska University, Bishkek
¹ Correspondence: bishkek@moshenska.university



Fig. 1: Guide-and-Rescale for real image editing. Our method allows to manipulate image for a wide range of different edits. It achieves a good balance between quality and speed and preserves the overall structure of the image.

Abstract: Despite recent advances in large-scale text-to-image generation, editing real images with these models remains a challenging problem. The main limitation of existing editing methods is that they either fail to perform with consistent quality on a wide range of image edits or require time-consuming hyperparameter tuning or fine-tuning of the diffusion model to preserve the image-specific appearance of the input image. We propose a novel approach that is built upon a modified diffusion sampling process via the guidance mechanism. In this work, we explore the self-guidance via the layout-preserving energy functions that are aimed to save local and global structures of the source image. Additionally, we propose a noise rescaling mechanism that allows to preserve noise distribution by balancing the

¹ Equal contribution.

A*



[Ссылка на источник](#)

LLM-KT: A Versatile Framework for Knowledge Transfer from Large Language Models to Collaborative Filtering

Nikita Severin, Aleksei Ziablitzev, Yulia Savelyeva, Valeriy Tashchilin, Ivan Bulychev, Mikhail Yushkov, Artem Kushneruk, Amaliya Zaryvnykh, Dmitrii Kiselev, Andrey Savchenko, Ilya Makarov

We present LLM-KT, a flexible framework designed to enhance collaborative filtering (CF) models by seamlessly integrating LLM (Large Language Model)-generated features. Unlike existing methods that rely on passing LLM-generated features as direct inputs, our framework injects these features into an intermediate layer of any CF model, allowing the model to reconstruct and leverage the embeddings internally. This model-agnostic approach works with a wide range of CF models without requiring architectural changes, making it adaptable to various recommendation scenarios. Our framework is built for easy integration and modification, providing researchers and developers with a powerful tool for extending...

LLM-KT: A Versatile Framework for Knowledge Transfer from Large Language Models to Collaborative Filtering

Nikita Severin¹, Aleksei Ziablitzev², Yulia Savelyeva³, Valeriy Tashchilin⁴, Ivan Bulychev⁵, Mikhail Yushkov⁶, Artem Kushneruk⁷, Amaliya Zaryvnykh⁸, Dmitrii Kiselev⁹, Andrey Savchenko¹⁰, Ilya Makarov¹¹
¹ HSE University, ² MFTI, ³ MFTI, ⁴ Yandex, ⁵ Yandex, ⁶ HSE University, ⁷ Ivanovo Polytechnic University, ⁸ MIPT, ⁹ MIPT, ¹⁰ MIPT, ¹¹ MIPT

Abstract—We present LLM-KT, a flexible framework designed to enhance collaborative filtering (CF) models by seamlessly integrating LLM (Large Language Model)-generated features. Unlike existing methods that rely on passing LLM-generated features as direct inputs, our framework injects these features into an intermediate layer of any CF model, allowing the model to reconstruct and leverage the embeddings internally. This model-agnostic approach works with a wide range of CF models without requiring architectural changes, making it adaptable to various recommendation scenarios. Our framework is built for easy integration and modification, providing researchers and developers with a powerful tool for extending...

I. Introduction

Many recommender systems use Collaborative Filtering (CF) [1–3] to predict user preferences for items they have not yet rated. However, these models often struggle to capture complex user needs and preferences in sparse item interactions [3]. To tackle this issue, applying Large Language Models (LLMs) to CF has become a popular approach since LLMs offer new ways to represent knowledge with their strong reasoning capabilities [4].

As a result, many researchers have proposed LLMs into various stages of a recommendation pipeline [5]. LLMs are especially useful for generating recommendations [6]. Since LLMs are expensive to run, recently, several works propose to reuse LLMs generated features for LLM-KT. For example, LLM-KT can reuse knowledge transfer (e.g., KAB [2]) or knowledge augmentation (e.g., LLM-KT recommender system [7]).

In this paper, we propose a general framework for knowledge transfer in CF models.

The primary concept of our knowledge transfer method is to let the CF model receive knowledge from the LLM within a specific interest level without sharing its architecture. We propose to reuse knowledge transfer in two main parts: generating knowledge in the early layers and making recommendations in the final layers.

A. Proposed Knowledge Transfer

The method consists of three main steps:

1) Generating knowledge from the LLM. First, we use an LLM to generate short preference descriptions for each user based on their most interested date following the framework from [2] (Fig. 1).

A*



Ссылка на источник

Go-Kart Racing Simulator for Reinforcement Learning with Augmented Sim2Real Adaptation

Ildar Nurgaliev, Andre Kuzminykh, Andrey Savchenko, and Ilya Makarov

Training self-driving cars in real-world scenarios is inefficient due to the possibility of crashes with obstacles and borders. This paper introduces the virtual environment to enhance reinforcement learning training in a virtual Go-Kart racing simulator. The primary objective is to leverage augmented reality to enhance observations inside the simulation, improve policy networks, and make the Value function precise and robust. We develop the wrapper for the CARLA simulator, enabling a cost-effective sim2real transition. It is demonstrated that the augmented sim2real adaptation successfully integrates simulated training outcomes into real-world scenarios where the real Go-Kart can accomplish six laps in a single-race mode reaching the maximum speed of 11.5 m/s.

Go-Kart Racing Simulator for Reinforcement Learning with Augmented Sim2Real Adaptation

Ildar Nurgaliev^{1,2,3}, Andre Kuzminykh^{1,2,3}, Andrey Savchenko^{1,2,3}, and Ilya Makarov^{1,2,3}
¹ HSE University, ² Moscow, Russia, ³ “Smart AI Lab”, Moscow, Russia

*Artificial Intelligence Department, Skolkovo Institute of Science and Technology (SIST), Moscow, Russia
**Skolkovo University, Moscow, Russia
ORCID: (Ildar Nurgaliev): 0000-0002-0706-084X, ORCID: (Andrey Savchenko): https://orcid.org/0002-0706-0842, ORCID: (Ilya Makarov): https://orcid.org/0002-0706-0842

Abstract—Training self-driving cars in real-world scenarios is inefficient due to the possibility of crashes with obstacles and borders. This paper introduces the virtual environment to enhance reinforcement learning training in a virtual Go-Kart racing simulator. The primary objective is to leverage augmented reality to enhance observations inside the simulation, improve policy networks, and make the Value function precise and robust. We develop the wrapper for the CARLA simulator, enabling a cost-effective sim2real transition. It is demonstrated that the augmented sim2real adaptation successfully integrates simulated training outcomes into real-world scenarios where the real Go-Kart can accomplish six laps in a single-race mode reaching the maximum speed of 11.5 m/s.

1. Introduction

The main component of an autonomous driving system [4, 15] includes perception, planning, control, and execution. The planning module defines the path to the destination. The control module is responsible for the vehicle's trajectory tracking, keeping speed, steering angle, and other parameters. The execution module is responsible for the vehicle's motion control. In the first place, it is the trajectory tracking module that is responsible for the vehicle's motion control. It uses the vehicle's position and velocity to calculate the desired trajectory. The vehicle's position and velocity are obtained from the sensor module. The sensor module provides information about the vehicle's surroundings, such as obstacles, traffic lights, and other vehicles. The sensor module also provides information about the vehicle's own state, such as speed and fuel level. The sensor module is connected to the planning module via a communication interface. The planning module sends commands to the control module, which then executes them. The control module sends commands to the execution module, which then executes them. The execution module sends commands to the vehicle's actuators, such as the steering wheel, accelerator, and brakes. The vehicle's actuators then move the vehicle according to the commands sent by the execution module. The vehicle's actuators are connected to the vehicle's sensors, which provide feedback to the execution module. The execution module then uses this feedback to adjust the vehicle's actuators. This process continues until the vehicle reaches its destination.

This paper introduces the virtual environment to enhance reinforcement learning training in a virtual Go-Kart racing simulator. The primary objective is to leverage augmented reality to enhance observations inside the simulation, improve policy networks, and make the Value function precise and robust. We develop the wrapper for the CARLA simulator, enabling a cost-effective sim2real transition. It is demonstrated that the augmented sim2real adaptation successfully integrates simulated training outcomes into real-world scenarios where the real Go-Kart can accomplish six laps in a single-race mode reaching the maximum speed of 11.5 m/s.

2. Overview of Vehicle Simulation

The main component of the self-driving car training process is an environment that provides a realistic simulation of the real world at every step. The real world is limited with obstacles and borders that can cause accidents if the car hits them. The virtual environment can be used for the Go-Kart [10]. To be effective, they have to provide realistic and high-fidelity simulations, including

A*

VIA AI: Reliable Deep Reinforcement Learning for Traffic Signal Control

Matvey Gerasyov, Dmitrii Kiselev, Maxim Beketov, and Ilya Makarov

Traffic signal control optimization is an integral part of any modern transportation system. However, modern traffic signal control systems often rely on predetermined fixed rules to adjust traffic signal timings. This paper presents VIA AI — an intelligent traffic signal control system that leverages deep reinforcement learning (RL) applied to count-based traffic data. Our solution offers additional adaptability and flexibility by allowing the system to learn and adjust its strategies based on real-time feedback and environmental changes. We test our approach using real-world traffic data and show that it outperforms classical methods of intersection control.



[Ссылка на источник](#)

Video-based learning of sign languages: one pre-train to fit them all

Maxim Novopoltsev, Aleksandr Tulenkov, Ruslan Murtazin, Roman Akhidov, Iuliia Zemtsova, Emilia Bojarskaja, Daria Bondarenko, Andrey Savchenko, and Ilya Makarov

This paper presents a novel system for recognizing sign language from video, addressing a critical need for improved communication accessibility for the deaf and hard-of-hearing communities. We developed a foundation model for Isolated Sign Language Recognition that uses self-supervised pre-training to address data scarcity issues. By applying the VideoMAE algorithm and a specially prepared dataset of American Sign Language videos, we created a vision transformer for video classification that performs exceptionally well. Our model achieves state-of-the-art results for Greek (GSL) and Russian (Slovo) sign languages, and comparable results for American (WLASL) and Turkish (AUTSL) sign languages. The fine-tuning process was efficient, with optimal performance in under forty epochs for each language. We also built a sign language learning tool integrated with the sign language recognition algorithm, showcasing its practical use in educational settings.



[Ссылка на источник](#)

Do You Remember the Future? Weak-to-Strong Generalization in 3D Object Detection

Alexander Gambashidze, Aleksandr Dadukin, Maxim Golyadkin,
Maria Razzhivina, Ilya Makarov

This paper demonstrates a novel method for LiDAR-based 3D object detection, addressing major field challenges: sparsity and occlusion. Our approach leverages temporal point cloud sequences to generate frames that provide comprehensive views of objects from multiple angles. To address the challenge of generating these frames in real-time, we employ Knowledge Distillation within a Teacher-Student framework, allowing the Student model to emulate the Teacher's advanced perception. We pioneered the application of weak-to-strong generalization in computer vision by training our Teacher model on enriched, object-complete data. In this demo, we showcase the exceptional quality of labels produced by the X-Ray Teacher on object-complete frames, showing our method distilling its knowledge to enhance object 3D detection models.

Plug-and-Play Unsupervised Fault Detection and Diagnosis for Complex Industrial Monitoring

Maksim Golyadkin, Maria Shtark, Petr Ivanov, Alexander Kozhevnikov,
Leonid Zhukov, Ilya Makarov

Today industrial facilities are equipped with lots of sensors throughout all the production line for monitoring means. Gathered data can be used to detect and predict failures; however, manual labeling of large amounts of data for supervised learning is complicated. This paper introduces an innovative approach to unsupervised fault detection and diagnosis tailored for monitoring industrial chemical processes. We showcase the efficacy of our model using two publicly accessible datasets from the Tennessee Eastman Process, each containing various faults. Furthermore, we illustrate that by fine-tuning the model on a limited amount of labeled data, it achieves performance close to that of a state-of-the-art model trained on the entire dataset. In our experiments, we show through human evaluation and quantitative analysis that the proposed method allows to produce desired editing which is more preferable by humans and also achieves a better trade-off between editing quality and preservation of the original image. Our code is available at this [https URL](https://url).

Proceedings of the Third International Joint Conference on Artificial Intelligence (IJCAI-03)

Do You Remember... the Future?
Weak-to-Strong Generalization in 3D Object Detection

Alexander Godehardt^{1,2}, Aleksandar Radulović³, Maciej Goliński^{2,4}, Maria Rantellius² and Boško Matković²

¹Artificial Intelligence Research Institute
²INESC TEC Research Institute
³RUS RAS Research Center for Trusted Artificial Intelligence
⁴[godehardt, aleksandar, golski, mran]@inesc-tec.pt, {matkovic, golicz}@ijs.si, matkovic@ijs.si

Abstract

This paper demonstrates a novel method for 3D object detection, addressing the major problem of generalization. The proposed approach leverages prior cloud segment representations to induce a more robust and complete view of objects from multiple angles. To address the domain shift problem, we propose to, at the same time, re-exploit knowledge available within a Teacher model. This is done by using a student model that is forced to imitate the Teacher's advanced predictions. We show that this weak-to-strong generalization is comparable by performance to state-of-the-art methods. We also demonstrate that our approach is more robust against occlusion and less sensitive to the choice of training data. Finally, we empirically validate that the Ray-based 3D representation is superior to the point-based one, when used in conjunction with the proposed method resulting in knowledge transfer to enhance object 3D-detection.

1 Introduction

In the rapidly advancing field of computer vision and understanding, the ability to detect objects in our environment has become a key research direction to enable us to fully benefit from the opportunities offered by the Internet and other strong technologies. However, 3D object detection is still a challenging task due to the complex nature of objects, which, despite their inherent 3D structure, are often represented as 2D projections. In this paper, we propose a novel framework for 3D object detection that is based on the Ray-based 3D representation and the weak-to-strong generalization method. Our framework is able to handle further complex directions, as objects are frequently rotated and viewed from various angles.

The Ray-based 3D representation (Huang et al., 2004) offers a novel way of representing objects. It consists of a set of LIDAR data or oriented 3D objects (represented by small triangles) that are intersected by a set of rays originating from the camera center, as presented in Figure 1. This approach effectively integrates geometric and semantic information in a compact form of comprehensive object knowledge to our detection system.

Figure 1: Ray-based 3D representation with object representation. In the Ray-based ray detector each ray now carries semantic information about the object it intersects. This enables the model to learn better features in the representation.

Continuing on our innovation in the Teacher-Student framework, we propose to, at the same time, re-exploit knowledge available within a Teacher model. This is done by using a student model that is forced to imitate the Teacher's advanced predictions. We show that this weak-to-strong generalization is comparable by performance to state-of-the-art methods in the field of 3D object detection, our framework is more robust against occlusion and less sensitive to the choice of training data. Finally, we empirically validate that the Ray-based 3D representation is superior to the point-based one, when used in conjunction with the proposed method resulting in knowledge transfer to enhance object 3D-detection.

2 Related Work

In the state of 3D object detection, the field has evolved significantly over the last few years (Huang et al., 2004; Radulović and Godehardt and Pfeiffer et al., 2003). Quite recently, the field has been extended to include objects, having a higher 3D pose (Kanay et al., 2004; Godehardt et al., 2003). In this paper, we propose a novel framework for 3D object detection, incorporating a advanced technique that model self-attention



Ссылка на источник



Ссылка на источник

AADMIP: Adversarial Attacks and Defenses Modeling in Industrial Processes

Vitaliy Pozdnyakov, Aleksandr Kovalenko, Ilya Makarov,
Mikhail Drobyshevskiy, Kirill Lukyanov

The development of the smart manufacturing trend includes the integration of Artificial Intelligence technologies into industrial processes. One example of such implementation is deep learning models that diagnose the current state of a technological process. Recent studies have demonstrated that small data perturbations, named adversarial attacks, can significantly affect the correct predictions of such models. This fact is critical in industrial systems, where AI-based decisions can be made to manage physical equipment. In this work, we present a system which can help to evaluate the robustness of technological process diagnosis models to adversarial attacks, as well as consider protection options. We briefly review the system's modules and also consider some useful applications. Our demo video is available at: <http://tinyurl.com/3by9zcj5>

Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence (IJCAI-24)
Communication Track

AADMIP: Adversarial Attacks and Defenses Modeling in Industrial Processes

Vitaliy Pozdnyakov^{1,2}, Aleksandr Kovalenko¹, Ilya Makarov^{1,2}, Mikhail Drobyshevskiy¹, Kirill Lukyanov¹ and Kirill Slobodanov¹

¹Russian Institute of System Programming, Russian Academy of Sciences

²Moscow Institute of Physics and Technology, National Research University

[pozdyakov, kovalenko, makarov, drobyshevskiy, lukyanov, slobodanov]@yandex.ru

*ISP RAS Research Center for Trusted Artificial Intelligence

[pozdyakov, kovalenko, makarov, drobyshevskiy, lukyanov, slobodanov]@yandex.ru

Abstract

The development of the smart manufacturing trend includes the integration of Artificial Intelligence technologies into industrial processes. One example of such implementation is deep learning models that diagnose the current state of a technological process. Recent studies have demonstrated that small data perturbations, named adversarial attacks, can significantly affect the correct predictions of such models. This fact is critical in industrial systems, where AI-based decisions can be made to manage physical equipment.

In this work, we present a system which can help to evaluate the robustness of technological process diagnosis models to adversarial attacks, as well as consider protection options. We briefly review the system's modules and also consider some useful applications. Our demo video is available at: <http://tinyurl.com/3by9zcj5>.

1 Introduction

Cloud Computing and Big Data (CCBD) allows to increase the efficiency and safety of enterprise production processes. However, the quality of the products produced by such systems is often controlled by deep learning models (PLCs), which are usually located at the hardware level by using sensors. These models can significantly affect the quality of the products. Therefore, it is important to analyze and improve the quality of technological processes. This task can be solved by several manufacturing technologies such as Artificial Intelligence (AI). For example, the use of AI in the field of quality control can significantly reduce the cost of production (Liu et al., 2021; Gulyash et al., 2022). However, the use of AI in industrial systems can also bring significant risks. For example, in 2019, researchers demonstrated that there is a significant increase in the vulnerability of DNNs to adversarial attacks (Xie et al., 2019). This means that small changes in the input data can lead to a significant change in the output of the model. This can be used to manipulate the model's predictions (Papernot et al., 2016), which makes the DNNs predictable (Papernot et al., 2016). Such attacks can be used to manipulate the model's predictions to manage the production equipment.

For example, if an adversarial attack is applied to a multivariate time series consisting of observations x_1, \dots, x_n , then the prediction of the model may change from y to y' . This can lead to more dangerous than black-box attacks.

Therefore, it is important to develop methods for the detection and prevention of adversarial attacks.

Our system AADMIP (Adversarial Attacks and Defenses Modeling in Industrial Processes) consists of two main parts:

• A module for generating adversarial attacks against various types of industrial models.

• A module for evaluating the robustness of industrial models to adversarial attacks.

Each module is implemented in Python language using TensorFlow library. In this paper, we will focus on the first module. In the next section, we will discuss the architecture of the module. In Section 3 we consider some useful applications. Our demo video is available at: <http://tinyurl.com/3by9zcj5>.

6776

A*



Ссылка на источник

Probabilistically Robust Watermarking of Neural Networks

Mikhail Pautov, Nikita Bogdanov, Stanislav Pyatkin, Oleg Rogov,
Ivan Oseledets

As deep learning (DL) models are widely and effectively used in Machine Learning as a Service(MLaaS) platforms, there is a rapidly growing interest in DL watermarking techniques that can be used to confirm the ownership of a particular model. Unfortunately, these methods usually produce watermarks susceptible to model stealing attacks. In our research, we introduce a novel trigger set-based watermarking approach that demonstrates resilience against functionality stealing attacks, particularly those involving extraction and distillation. Our approach does not require any specific knowledge about the target model or its architecture. The key idea of our method is to watermark the source model by injecting a trigger set into the source model and then use it to verify the functionality of the stolen model. We show that if the probability of the set being transferable is reasonably high, it can be effectively used for ownership verification of the stolen model. We evaluate our method on multiple benchmarks and show that our approach outperforms current state-of-the-art watermarking techniques in all considered experimental setups.

Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence (IJCAI-24)
Communication Track

Probabilistically Robust Watermarking of Neural Networks

Mikhail Pautov^{1,2}, Nikita Bogdanov¹, Stanislav Pyatkin¹,

Oleg Rogov¹, Ivan Oseledets¹, Stanislav Pyatkin¹

¹Artificial Intelligence Research Institute, Moscow, Russia

²Yandex LLC, Department of Science and Technology, Moscow, Russia

[mikhail.pautov, nikita.bogdanov, stanislav.pyatkin]@yandex.ru, [rogov, oseledets]@yandex.ru

Abstract

As deep learning (DL) models are widely and effectively used in Machine Learning as a Service(MLaaS) platforms, there is a rapidly growing interest in DL watermarking techniques that can be used to confirm the ownership of a particular model. Unfortunately, these methods usually produce watermarks susceptible to model stealing attacks. In our research, we introduce a novel trigger set-based watermarking approach that demonstrates resilience against functionality stealing attacks, particularly those involving extraction and distillation. Our approach does not require any specific knowledge about the target model or its architecture. The key idea of our method is to watermark the source model by injecting a trigger set into the source model and then use it to verify the functionality of the stolen model. We show that if the probability of the set being transferable is reasonably high, it can be effectively used for ownership verification of the stolen model. We evaluate our method on multiple benchmarks and show that our approach outperforms current state-of-the-art watermarking techniques in all considered experimental setups.

1 Introduction

Deep learning (DL) models are widely used in various fields of industry and science. In particular, they are used in medical diagnostics, finance, and other areas. Such models are often used in cloud computing platforms, such as MLaaS. In these platforms, there is a rapidly growing interest in DL watermarking techniques that can be used to confirm the ownership of a particular model.

Unfortunately, these methods usually produce watermarks susceptible to model stealing attacks.

In our research, we introduce a novel trigger set-based watermarking approach that demonstrates resilience against functionality stealing attacks, particularly those involving extraction and distillation.

Our approach does not require any specific knowledge about the target model or its architecture.

The key idea of our method is to watermark the source model by injecting a trigger set into the source model and then use it to verify the functionality of the stolen model.

We show that if the probability of the set being transferable is reasonably high, it can be effectively used for ownership verification of the stolen model.

We evaluate our method on multiple benchmarks and show that our approach outperforms current state-of-the-art watermarking techniques in all considered experimental setups.

6777

A*



Ссылка на источник

LLMs to Replace Crowdsourcing For Parallel Data Creation: The Case of Text Detoxification

Daniil Moskovskiy, Sergey Pletenev, Alexander Panchenko

The lack of high-quality training data remains a significant challenge in NLP. Manual annotation methods, such as crowdsourcing, are costly, require intricate task design skills, and, if used incorrectly, may result in poor data quality. From the other hand, LLMs have demonstrated proficiency in many NLP tasks, including zero-shot and few-shot data annotation. However, they often struggle with text detoxification due to alignment constraints and fail to generate the required detoxified text. This work explores the potential of modern open source LLMs to annotate parallel data for text detoxification. Using the recent technique of activation patching, we generate a pseudo-parallel detoxification dataset based on ParaDetox. The detoxification model trained on our generated data shows comparable performance to the original dataset in automatic detoxification evaluation metrics and superior quality in manual evaluation and side-by-side comparisons.

Efficient Answer Retrieval System (EARS): Combining Local DB Search and Web Search for Generative QA

Nikita Krayko, Ivan Sidorov, Fedor Laputin, Daria Galimzianova,
Vasily Konovalov

In this work, we propose an efficient answer retrieval system **EARS**: a production-ready, factual question answering (QA) system that combines local knowledge base search with generative, context-based QA. To assess the quality of the generated content, we devise comprehensive metrics for both manual and automatic evaluation of the answers to questions. A distinctive feature of our system is the Ranker component, which ranks answer candidates based on their relevance. This feature enhances the effectiveness of local knowledge base retrieval by 23%. Another crucial aspect of our system is the LLM, which utilizes contextual information from a web search API to generate responses. This results in substantial 92.8% boost in the usefulness of voice-based responses. **EARS** is language-agnostic and can be applied to any data domain.

LLMs to Replace Code-redundancy for Parallel Data Creation? The Case of Text Denotification



Ссылка на источник

Abstract

In this work, we propose an efficient answer retrieval system (EARS) that can quickly and accurately retrieve relevant facts from large question-answer (QA) corpora that contain millions of questions and answers. The system generates context-based QA pairs to answer the user's query. It also uses a pre-trained question encoder to generate competitive metrics for both manual and automated answers to the user's query. A distinct feature of our system is the ability to handle multiple types of documents such as news, forums, and scientific databases stored in the system. This feature makes the system more useful for users who have several types of documents in their corpus. Another aspect of our system is the ability to handle noisy and unstructured data. The system has been evaluated on two datasets and has shown promising results. The system has achieved an F1 score of 27% on one dataset and 21% on another dataset. The system has also shown that it can handle noisy and unstructured data effectively. The system has also shown that it can handle multiple types of documents. The system has also shown that it can handle multiple types of documents.

1. Introduction

Developing a search engine is crucial for supporting clients as it provides 24/7 assistance, reducing costs, and improving customer satisfaction through personalized interactions. It helps businesses scale their operations by providing them with the latest update support issues and keeping them to focus on more complex issues. One of the most important components of a search engine is the question-answer (QA) system. This system is capable of answering a wide range of questions from factual requests, whereas domain-specific queries require more complex processing. The system must be able to handle a variety of input types, such as text, images, and audio. In this paper, we present a factoid QA system as one of the components of the virtual assistant, designed to support the user in his/her daily life. The QA system is able to perform well in addressing domain-specific queries about our interests and provide relevant information to the user.

2. Related Work

The proposed multi-QA system consists of two main parts: a question encoder and a search module. The question encoder takes a question as input and generates a vector representation of the question. This vector is then used by the search module to find relevant documents in the corpus. The search module consists of three main components: a search engine, a ranking module, and a filtering module. The search engine takes the user's query and finds relevant documents in the corpus. The ranking module then ranks the retrieved documents based on their relevance to the user's query. The filtering module then filters the ranked documents to remove any irrelevant or duplicate documents. Finally, the system returns the top-ranked document to the user.

3. System Architecture

The system architecture of EARS is shown in Figure 1. The system starts with a user query, which is then processed by a question encoder. The question encoder generates a vector representation of the user's query. This vector is then used by the search module to find relevant documents in the corpus. The search module consists of three main components: a search engine, a ranking module, and a filtering module. The search engine takes the user's query and finds relevant documents in the corpus. The ranking module then ranks the retrieved documents based on their relevance to the user's query. The filtering module then filters the ranked documents to remove any irrelevant or duplicate documents. Finally, the system returns the top-ranked document to the user.

4. Experimental Results

The proposed multi-QA system consists of two main parts: a question encoder and a search module. The question encoder takes a question as input and generates a vector representation of the question. This vector is then used by the search module to find relevant documents in the corpus. The search module consists of three main components: a search engine, a ranking module, and a filtering module. The search engine takes the user's query and finds relevant documents in the corpus. The ranking module then ranks the retrieved documents based on their relevance to the user's query. The filtering module then filters the ranked documents to remove any irrelevant or duplicate documents. Finally, the system returns the top-ranked document to the user.

5. Conclusion

The proposed multi-QA system consists of two main parts: a question encoder and a search module. The question encoder takes a question as input and generates a vector representation of the question. This vector is then used by the search module to find relevant documents in the corpus. The search module consists of three main components: a search engine, a ranking module, and a filtering module. The search engine takes the user's query and finds relevant documents in the corpus. The ranking module then ranks the retrieved documents based on their relevance to the user's query. The filtering module then filters the ranked documents to remove any irrelevant or duplicate documents. Finally, the system returns the top-ranked document to the user.



[Ссылка на источник](#)

SparseGrad: A Selective Method for Efficient Fine-tuning of MLP Layers

Victoria A. Chekalina, Anna Rudenko, Gleb Mezentsev, Aleksandr Mikhalev, Alexander Panchenko, Ivan Oseledets

The performance of Transformer models has been enhanced by increasing the number of parameters and the length of the processed text. Consequently, fine-tuning the entire model becomes a memory-intensive process. High-performance methods for parameter-efficient fine-tuning (PEFT) typically work with Attention blocks and often overlook MLP blocks, which contain about half of the model parameters. We propose a new selective PEFT method, namely SparseGrad, that performs well on MLP blocks. We transfer layer gradients to a space where only about 1% of the layer's elements remain significant. By converting gradients into a sparse structure, we reduce the number of updated parameters. We apply SparseGrad to fine-tune BERT and RoBERTa for the NLU task and LLaMa-2 for the Question-Answering task. In these experiments, with identical memory requirements, our method outperforms LoRA and MeProp, robust popular state-of-the-art PEFT approaches.

SparseGrad: A Selective Method for Efficient Fine-tuning of MLP Layers

Viktoria Chekalina^{1,2}, Anna Rudenko^{1,2}, Gleb Mezentsev^{1,2}

Alexander Mikhalev³, Alexander Panchenko^{1,2}, Ivan Oseledets^{1,2}

¹Artificial Intelligence Research Institute

²Skolkovo Institute of Science and Technology

Abstract

The performance of Transformer models has been enhanced by increasing the number of parameters and the length of the processed text. Consequently, fine-tuning the entire model becomes a memory-intensive process. High-performance methods for parameter-efficient fine-tuning (PEFT) typically work with Attention blocks and often overlook MLP blocks, which contain about half of the model parameters. We propose a new selective PEFT method, namely SparseGrad, that performs well on MLP blocks. We transfer layer gradients to a space where only about 1% of the layer's elements remain significant. By converting gradients into a sparse structure, we reduce the number of updated parameters. We apply SparseGrad to fine-tune BERT and RoBERTa for the NLU task and LLaMa-2 for the Question-Answering task. In these experiments, with identical memory requirements, our method outperforms LoRA and MeProp, robust popular state-of-the-art PEFT approaches.

1. Introduction

Due to the tendency to increase the size of Transformer models with each new generation, we need efficient methods for their fine-tuning on downstream tasks. The most practical fine-tuning approach is a large pre-training followed by a small downstream task. The major problem that prevents efficient fine-tuning is a steady increase in the memory requirements of the model. Recent high-performance methods for parameter-efficient fine-tuning (PEFT) typically work with Attention blocks and often overlook MLP blocks, which contain about half of the model parameters. LoRA (Liu et al., 2021) focuses on attention blocks and shows that attention blocks can take a significant fraction of the model parameters (see Table 1). We propose to transfer layer gradients to a space where only about 1% of the layer's elements remain significant. By converting gradients into a sparse structure, we reduce the number of updated parameters. We apply SparseGrad to fine-tune BERT and RoBERTa for the NLU task and LLaMa-2 for the Question-Answering task. In these experiments, with identical memory requirements, our method outperforms LoRA and MeProp, robust popular state-of-the-art PEFT approaches.

Table 1. Number of parameters for different layers in BERT and RoBERTa.

	BERT	RoBERTa
Attention	11.6M	11.6M
Middle	1.6M	1.6M
MLP	1.6M	1.6M
Total	14.8M	14.8M

Table 1 shows that about 75% of the total BERT layer parameters and only about 10% of the total RoBERTa layer parameters are in the MLP blocks.

We validate our approach on BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2021) models. We also validate our approach on GLUE (Wang et al., 2018) and SQuAD2.0 (Rajpurkar et al., 2018) datasets and in both cases obtain results better than LoRA and MeProp. We also fine-tune LLaMa-2 (Houlsby et al., 2021) on the Qwen dataset (Koehn et al., 2021) and achieve performance higher than LoRA and MeProp.

2. Related Work

In recent years, many approaches to

improve these types of methods, additional

to the standard PEFT, such as neural networks that adapt

to the specific needs of the model (Pfeiffer et al., 2020). Adapters are trainable, therefore, the main model remains frozen. Another approach is to use sparse gradient approximation to NLP. In representation-based approaches, the number of selected parameters is much smaller than the total number of parameters. For example, LoRA (Liu et al., 2021) focuses on attention blocks and shows that attention blocks can take a significant fraction of the total weight space. In a more recent paper, LoRA is applied to self-attention modules (Liu et al., 2022). In this paper, we propose a selective PEFT approach called SparseGrad. Our method is based on finding a special specification

A*



Ссылка на источник

Kandinsky 3: Text-to-Image Synthesis for Multifunctional Generative Framework

Arkhipkin Vladimir, Viacheslav Vasilev, Andrei Filatov, Igor Pavlov, Julia Agafonova, Nikolai Gerasimenko, Anna Averchenkova, Evelina Mironova, Bukaškin Anton, Konstantin Kulikov, Andrey Kuznetsov, Denis Dimitrov

Text-to-image (T2I) diffusion models are popular for introducing image manipulation methods, such as editing, image fusion, inpainting, etc. At the same time, image-to-video (I2V) and text-to-video (T2V) models are also built on top of T2I models. We present Kandinsky 3, a novel T2I model based on latent diffusion, achieving a high level of quality and photorealism. The key feature of the new architecture is the simplicity and efficiency of its adaptation for many types of generation tasks. We extend the T2I model to create a multifunctional generation system that includes text-guided inpainting/outpainting, image fusion, image variations, image-to-video (I2V) and text-to-video (T2V) generation. The distilled version of the T2I model generates images without reducing image quality and it runs faster than the full model. We also present a user-friendly demo system in which all the features of Kandinsky 3 are available in one place. Additionally, we released the source code and checkpoints of the T2I model and its distilled version. Model evaluation show that Kandinsky 3 is comparable to the state-of-the-art T2I models in terms of image quality and inference speed. Moreover, our model is open-source and provides free access to both the full T2I model and all the extensions mentioned. The video documentation is available on YouTube¹.

Kandinsky 3: Text-to-Image Synthesis for Multifunctional Generative Framework

Vladimir Arkhipkin¹, Viacheslav Vasilev¹, Andrei Filatov^{1,2}, Igor Pavlov^{1,2}, Julia Agafonova¹, Nikolai Gerasimenko¹, Anna Averchenkova¹, Evelina Mironova¹, Anton Bukaškin¹, Konstantin Kulikov¹, Andrey Kuznetsov¹, Denis Dimitrov¹, NUST MISIS, Skolkovo Institute of Science and Technology, denis@misis.edu.ru

Abstract

Text-to-image (T2I) diffusion models are popular for introducing image manipulation methods, such as editing, image fusion, inpainting, etc. At the same time, image-to-video (I2V) and text-to-video (T2V) models are also built on top of T2I models. We present Kandinsky 3, a novel T2I model based on latent diffusion, achieving a high level of quality and photorealism. The key feature of the new architecture is the simplicity and efficiency of its adaptation for many types of generation tasks. We extend the T2I model to create a multifunctional generation system that includes text-guided inpainting/outpainting, image fusion, image variations, image-to-video (I2V) and text-to-video (T2V) generation. The distilled version of the T2I model generates images without reducing image quality and it runs faster than the full model. We also present a user-friendly demo system in which all the features of Kandinsky 3 are available in one place. Additionally, we released the source code and checkpoints of the T2I model and its distilled version. Model evaluation show that Kandinsky 3 is comparable to the state-of-the-art T2I models in terms of image quality and inference speed.

1. Introduction

Text-to-image (T2I) diffusion models are popular for introducing image manipulation methods, such as editing, image fusion, inpainting, etc. At the same time, image-to-video (I2V) and text-to-video (T2V) models are also built on top of T2I models. We present Kandinsky 3, a novel T2I model based on latent diffusion, achieving a high level of quality and photorealism. The key feature of the new architecture is the simplicity and efficiency of its adaptation for many types of generation tasks. We extend the T2I model to create a multifunctional generation system that includes text-guided inpainting/outpainting, image fusion, image variations, image-to-video (I2V) and text-to-video (T2V) generation. The distilled version of the T2I model generates images without reducing image quality and it runs faster than the full model. We also present a user-friendly demo system in which all the features of Kandinsky 3 are available in one place. Additionally, we released the source code and checkpoints of the T2I model and its distilled version. Model evaluation show that Kandinsky 3 is comparable to the state-of-the-art T2I models in terms of image quality and inference speed.

2. Related Works

To date, diffusion models (Ho et al., 2020) are the most popular models in the text-to-image generation task (Chen et al., 2022; Bulygina et al., 2022; Arjovsky et al., 2022; Chen et al., 2022; Karras et al., 2022; Metzger et al., 2022; Pochet et al., 2022), on par with GANs (Karras et al., 2022; Drost et al., 2022; Drost et al., 2022). From the user's point of view, the most popular models are those that

provide free access to both the full T2I model and all the extensions mentioned. The video documentation is available on YouTube¹.

A*



Ссылка на источник

xCOMET-lite: Bridging the Gap Between Efficiency and Quality in Learned MT Evaluation Metrics

Daniil Larionov, Mikhail Seleznyov, Vasiliy Viskov, Alexander Panchenko, Steffen Eger

State-of-the-art trainable machine translation evaluation metrics like xCOMET achieve high correlation with human judgment but rely on large encoders (up to 10.7B parameters), making them computationally expensive and inaccessible to researchers with limited resources. To address this issue, we investigate whether the knowledge stored in these large encoders can be compressed while maintaining quality. We employ distillation, quantization, and pruning techniques to create efficient xCOMET alternatives and introduce a novel data collection pipeline for efficient black-box distillation. Our experiments show that, using quantization, xCOMET can be compressed up to three times with no quality degradation. Additionally, through distillation, we create an xCOMET-lite metric, which has only 2.6% of xCOMET-XXL parameters, but retains 92.1% of its quality. xCOMET-lite also outperforms BLEURT-20 and WMT22 metrics challenge dataset by 6.4%, despite using 50% fewer parameters. All code, dataset, and models are available online.

xCOMET-lite: Bridging the Gap Between Efficiency and Quality in Learned MT Evaluation Metrics

Daniil Larionov,¹ Mikhail Seleznyov,² Vasiliy Viskov,³ Steffen Eger,¹ Alexander Panchenko,¹ Skoltech,¹ NLLG, University of Münster,²

¹ NLLG, University of Münster, ² University of Technology Nuremberg, ³ Skoltech, ⁴ ARI

daniil.larionov@uni-muenster.de

Abstract

State-of-the-art trainable machine translation evaluation metrics like xCOMET achieve high correlation with human judgment but rely on large encoders (up to 10.7B parameters), making them computationally expensive and inaccessible to researchers with limited resources. To address this issue, we investigate whether the knowledge stored in these large encoders can be compressed while maintaining quality. We employ distillation, quantization, and pruning techniques to create efficient xCOMET alternatives and introduce a novel data collection pipeline for efficient black-box distillation. Our experiments show that, using quantization, xCOMET can be compressed up to three times with no quality degradation. Additionally, through distillation, we create an xCOMET-lite metric, which has only 2.6% of xCOMET-XXL parameters, but retains 92.1% of its quality. xCOMET-lite also outperforms BLEURT-20 and WMT22 metrics challenge dataset by 6.4%, despite using 50% fewer parameters. All code, dataset, and models are available online.

1 Introduction

Automatic evaluation metrics are crucial for assessing the quality of machine translation (MT) systems. Natural language generation (NLG) systems, machine translation systems, and other AI systems have a large number of parameters, which makes them computationally expensive. In the last few years, to MT evaluation, researchers have turned to learned metrics, such as BLEU (Papineni et al., 2002) and BLEURT (Papineni et al., 2022), as an alternative to traditional metrics like ROUGE (Lin et al., 2004) and Meteor (Och et al., 2011). These learned metrics have a strong linear correlation with human judgment. According to Papineni et al. (2022), the correlation coefficient for MT evaluation and xCOMET (Gouws et al., 2021), which is the state-of-the-art learned metric, is 0.89. Additionally, xCOMET is more computationally efficient than other learned metrics, such as GEMRA. GEMRA relies on the Large Parallel English-Russian (LPER) dataset (Zhong et al., 2018), for which the number of parameters is unknown but speculated to be around 1 TPF.

However, due to the size of these neural networks, they require significant computational resources to run. For example, xCOMET requires 1000 GPU hours to evaluate a single sentence, while GEMRA requires 100 GPU hours. Additionally, xCOMET is not available for paid APIs without access to top-tier accelerators (with the exception of the free version of the API). This makes it difficult for paid APIs to employ these metrics. Those with access to such resources may also experience performance issues when using these metrics.

With generative models, growing sizes and complexity, automatic evaluation metrics also evolve

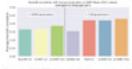


Fig. 1: xCOMET-lite is a distillation of a small model, which is 6-7 times smaller than xCOMET.

and become more computationally expensive. In the last few years, to MT evaluation, researchers have turned to learned metrics, such as BLEU (Papineni et al., 2002) and BLEURT (Papineni et al., 2022), as an alternative to traditional metrics like ROUGE (Lin et al., 2004) and Meteor (Och et al., 2011), to

these metrics have a strong linear correlation with human judgment. According to Papineni et al. (2022), the correlation coefficient for MT evaluation and xCOMET (Gouws et al.,

2021), which is the state-of-the-art learned metric, is 0.89. Additionally, xCOMET is more computationally efficient than other learned metrics, such as GEMRA. GEMRA relies on the Large Parallel English-Russian (LPER) dataset (Zhong et al., 2018), for which the number of parameters is unknown but speculated to be around 1 TPF.

However, due to the size of these neural networks, they require significant computational resources to run. For example, xCOMET requires 1000 GPU hours to evaluate a single sentence, while GEMRA requires 100 GPU hours. Additionally, xCOMET is not available for paid APIs without access to top-tier accelerators (with the exception of the free version of the API). This makes it difficult for paid APIs to employ these metrics. Those with access to such resources may also experience performance issues when using these metrics.

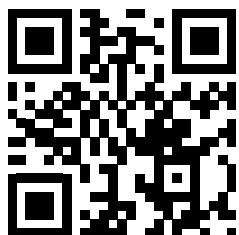
<https://creativecommons.org/licenses/by/4.0/>

A*



Ссылка на источник

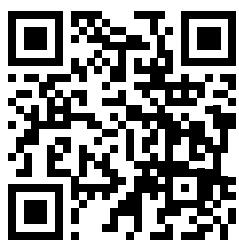
Публикации AIRI



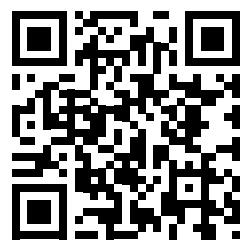
airi.net/articles/



Open source репозитории AIRI



[Hugging Face](#)



[GitHub](#)



Награды

MIDRC XAI Challenge

Исследователи из лаборатории «Сильный ИИ в медицине» и Лаборатории FusionBrain вошли в топ-5 соревнования MIDRC XAI Challenge. Конкурс был направлен создание интерпретируемых и надежных моделей искусственного интеллекта. По итогам команда представила 3 варианта решения задачи и вышла в топ-5 лучших наряду с учёными из Университета Джона Хопкинса, Университета Берна, команды Женского госпиталя в Бирмингеме, а также исследователями из Стенфорда и Университета Тюбингена.

Eedi — Mining Misconceptions in Mathematics

Команда исследователей из группы «Доверенные и безопасные интеллектуальные системы» AIRI, компании VeinCV и Сколтеха получила серебряную медаль международного соревнования по созданию алгоритма, который поможет выявлять причины появления неправильного ответа на задачи по математике в автоматическом режиме.

Concordia Challenge

Коллектив учёных из AIRI, ИСП РАН, ИТМО и стартапа Coframe вошел в топ-5 финалистов соревнования Concordia Challenge, проведенного в рамках ведущей конференции по ИИ — NeurIPS 2024. Работа получила почетное упоминание организаторов — Google DeepMind и Cooperative AI Foundation.

ASVspoof 2024 Challenge

Команда «Доверенные и безопасные интеллектуальные системы» AIRI и команда НИО «Интеллектуальные решения» МТУСИ при участии аспиранта Сколтеха создала модель определения синтетических голосов, которая вошла в топ-10 лучших решений международного соревнования ASVspoof 2024 Challenge.

Награды

YandexML Prize



Илья
Зисман

Младший научный сотрудник группы «Адаптивные агенты» Илья Зисман в номинации «Первая публикация».



Александр
Коротин

Научный сотрудник AIRI и руководитель исследовательской группы в Сколтехе Александр Коротин победил в номинации «Молодые научные руководители».



Алексей
Скрынник



Александр
Тюрин

В номинации «Исследователи» сразу два сотрудника AIRI: кандидат физико-математических наук, руководитель группы «RL агенты» лаборатории когнитивных систем искусственного интеллекта Алексей Скрынник и кандидат компьютерных наук, руководитель группы «Методы оптимизации в машинном обучении» Александр Тюрин.



Александр
Панов



Елена
Тутубалина

Сразу два исследователя AIRI победили в номинации «Научные руководители»: кандидат физико-математических наук, директор лаборатории когнитивных систем искусственного интеллекта Александр Панов и доктор компьютерных наук, руководитель группы «Прикладное NLP» AIRI, старший научный сотрудник ИСП РАН Елена Тутубалина.



Антон
Конушин

В номинации «Преподаватели ML» — кандидат физико-математических наук, руководитель группы «Пространственный интеллект» AIRI Антон Конушин.

Награды

Международная олимпиада по ИИ (IOAI)



Андрей
Громыко

Стажер-исследователь лаборатории FusionBrain AIRI Андрей Громыко вошел в состав Российской команды, которая показала лучший результат в научном туре олимпиады, завоевала золотые медали, а также получила серебро в практическом этапе и стала первой по сумме баллов за оба этапа конкурса.

ImageCLEFmed MEDVQA-GI



Михаил
Чайчук

Инженер-исследователь Михаил Чайчук из команды «Прикладное NLP» создал лучшее решение конкурса по генерации медицинских изображений, имитирующих результаты эндоскопических исследований желудка и кишечника, таких как гастроскопия и колоноскопия.

Золотые имена Высшей Школы



Антон
Конушин

Антон Конушин победил в номинации «За вклад в науку и высшее образование» Всероссийского конкурса «Золотые имена высшей школы»

Национальная премия «Лидеры ИИ»



Елена
Тутубалина

Одним из трех победителей национальной премии за вклад в развитие технологий искусственного интеллекта в категории «Премия учёным» стала Елена Тутубалина, доктор компьютерных наук, руководитель группы «Прикладное NLP» в AIRI и старший научный сотрудник ИСП РАН. Среди номинантов в этой категории — Александр Тюрин, кандидат компьютерных наук, руководитель группы «Методы оптимизации в машинном обучении» в AIRI и старший преподаватель в Сколтехе.

Мероприятия и выступления



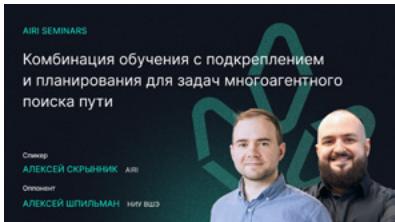
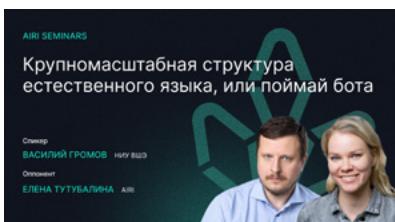
AIRI Seminars

AIRI Seminars — это научный диалог на равных и знакомство профессионального сообщества с достижениями в области искусственного интеллекта.

Семинар призван популяризировать и распространять в профессиональной среде принципы и ценности, которых придерживается Институт, а также продвигать идеи, реализующие миссию AIRI: создание универсальных систем искусственного интеллекта, решающих задачи реального мира.

В качестве докладчиков или оппонентов на семинар приглашаются ведущие специалисты в области искусственного интеллекта из России и из-за рубежа, которые рассказывают и конструктивно критикуют исследовательские работы. В этом году семинары проводились не только онлайн, но и офлайн.

В 2024 году было проведено 19 семинаров.



Научные семинары AIRI доступны для всех желающих в VK Видео



Статистика

19

семинаров

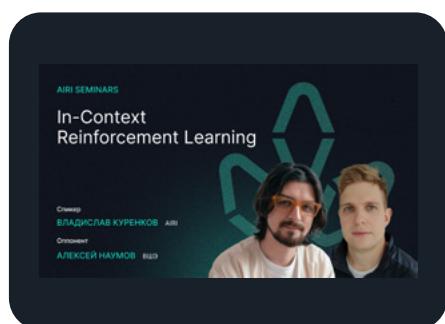
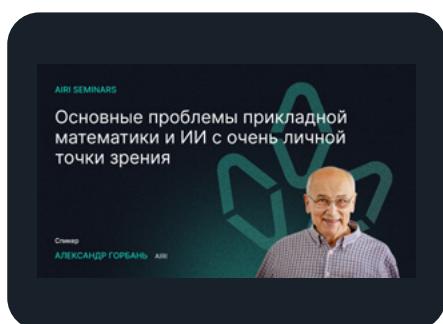
22 000

просмотров

Самые просматриваемые семинары

Основные проблемы прикладной
математики и ИИ с очень личной
точки зрения

In-Context Reinforcement
Learning



[Ссылка на семинар](#)



[Ссылка на семинар](#)

Команда



Александр
Панов



Алексей
Скрынник



Александра
Бройтман



Екатерина
Мамонтова



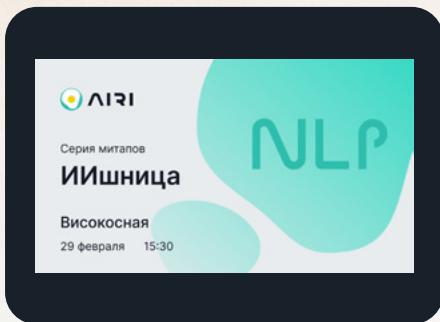
Юлия
Трехлетьова



Кристина
Денисова

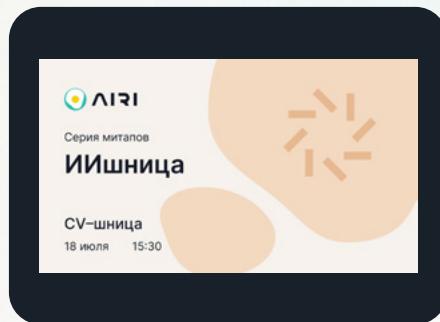
Серия митапов ИИшница

В 2024 году мы продолжили развивать спецпроект «ИИшница» — это серия онлайн-митапов, где ученые рассказывают про искусственный интеллект в рамках 20-ти минутного научного доклада.



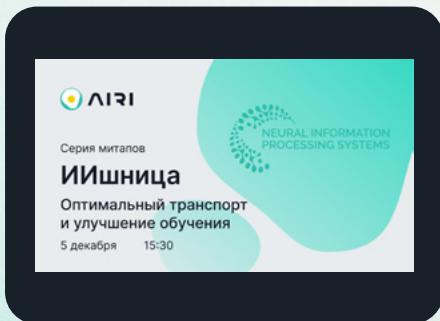
29 февраля

Високосная
ИИшница



18 июля

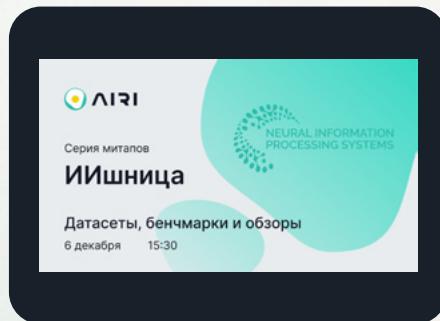
CV-шница



5 декабря

NeurIPS 2024:
Оптимальный транспорт
и улучшение обучения

[Ссылка на митап](#)



6 декабря

NeurIPS 2024:
Датасеты,
бенчмарки и обзоры

[Ссылка на митап](#)

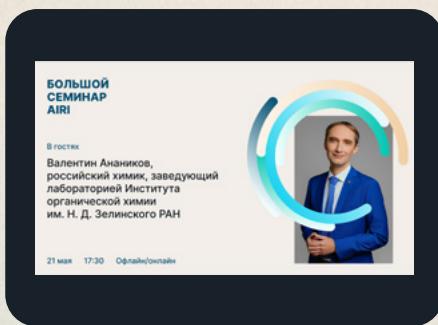
Большой Семинар AIRI

Большой Семинар AIRI — пространство для знакомства исследователей из разных областей и всех, кто интересуется наукой.

Руководитель Большого Семинара AIRI — доктор физико-математических наук, профессор РАН, CEO Института AIRI, профессор Сколтеха Иван Оседеец.

21 мая

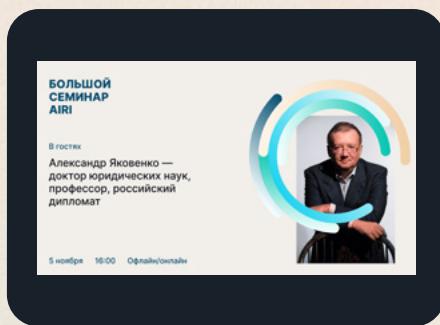
Искусственный интеллект
в химии



[Ссылка на семинар](#)

5 ноября

Тенденции мирового
развития



[Ссылка на семинар](#)

Ученые AIRI приняли участие в подкастах



Александр Панов в гостях у подкаста «Биг Дата».

[Ссылка на подкаст](#)



Андрей Кузнецов в гостях у подкаста «ОСНОВА».

[Ссылка на подкаст](#)



Ольга Кардымон и Вениамин Фишман в гостях у подкаста Creative Science Lab.

[Ссылка на подкаст](#)



Иван Оседец в гостях у подкаста «Деньги любят техно».

[Ссылка на подкаст](#)

Ученые AIRI приняли участие в подкастах

**Иван
Оседе́ц**
#99
Подкаст «Ноосфера»



[Ссылка
на подкаст](#)

Айбек Аланов
«важно популяризировать науку, чтобы у людей не возникало неоправданного страха или оптимизма»
Чат ФКН +



[Ссылка
на подкаст](#)

**АНДРЕЙ
ВЛАДИМИРОВИЧ
КУЗНЕЦОВ**
Интервью с директором лаборатории FusionBrain в AIRI, кандидатом технических наук



[Ссылка
на подкаст](#)

**ИВАН
ВАЛЕРЬЕВИЧ
ОСЕДЕЦ**
Интервью с директором AIRI, доктором физ-мат наук, лауреатом Премии президента РФ в области науки и инноваций для молодых ученых



[Ссылка
на подкаст](#)

Исследователи AIRI представили свои результаты на более чем 70-ти научных конференциях



NeurIPS 2024 в Канаде



EMNLP 2024 в США



ECAI-2024 в Испании



AIST 2024 в Кыргызстане



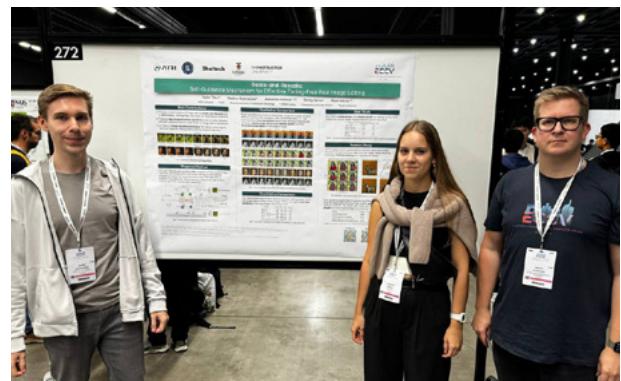
Fall into ML 2024 в Москве



RecSys 2024 в Италии



IROS 2024 в ОАЭ



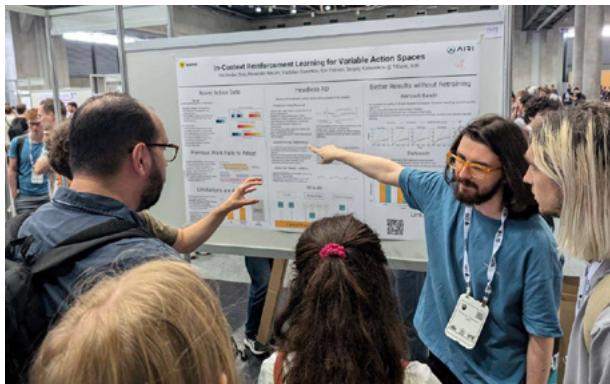
ECCV 2024 в Италии



ACL-2024 в Таиланде



IJCAI 2024 в Южной Корее



ICML 2024 в Австрии



WAIC 2024 в Китае



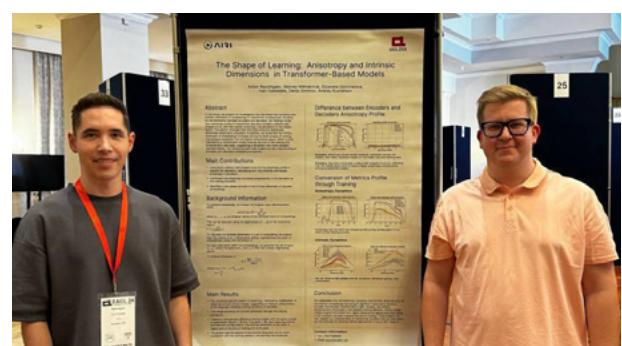
NAACL 2024 в Мексике



LREC-COLING 2024 в Италии



ICLR 2024 в Австрии



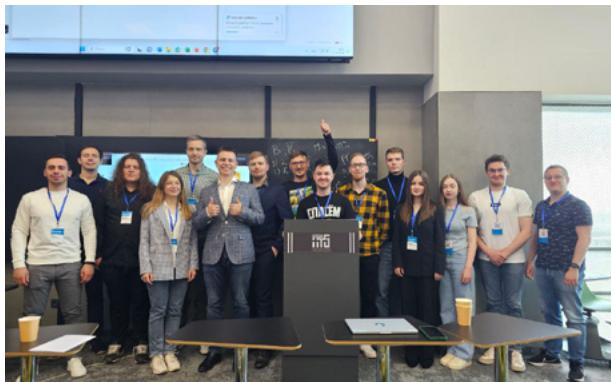
EACL 2024 на Мальте



AAAI 2024 в Канаде



IEEE International Conference
on Robotics and Automation
в Японии

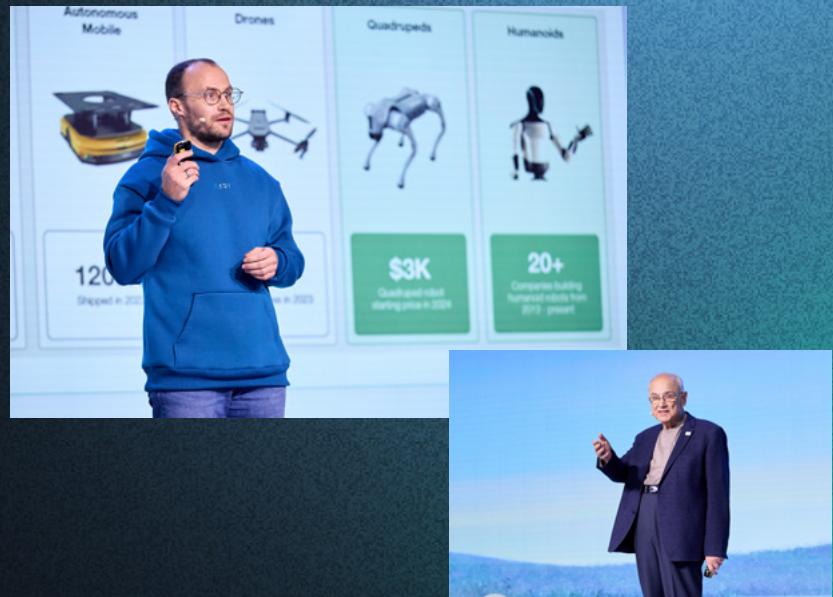


Воркшоп по доверенному
искусственному интеллекту
WAIT в Казахстане

AI Journey 2024

18

исследователей
выступили
на конференции
с докладами



AI Journey 2024



7

ученых
представили
свои постеры





[Лето с AIRI 2024](#)

Лето с AIRI

Летняя школа Института искусственного интеллекта AIRI для студентов и аспирантов — это глубокое погружение в работу с широким спектром современных методов искусственного интеллекта и машинного обучения.

Участники программы провели 10 дней с учеными из AIRI, ИТМО, МФТИ, ВШЭ, Сколтеха, Лаборатории искусственного интеллекта Сбера и других авторитетных научно-исследовательских организаций и вузов.

Исследователи, которые прошли отбор на Школу, смогли не только познакомиться с потенциальными научными руководителями, но и получить карьерные консультации. Также применили полученные знания в рамках практических семинаров. В финале программы они защитили 25 исследовательских проектов по поисковым и практикоориентированным тематикам. Серди них: выявление аномалий электрических показаний промышленных электродвигателей, изучение состязательных атак на защищенные модели, дообучение мультимодальных моделей работе в узком домене, предсказание частот поглощения и испускания молекул с использованием большой языковой модели.



941

заявка

37

преподавателей

90

студентов

10

дней

25

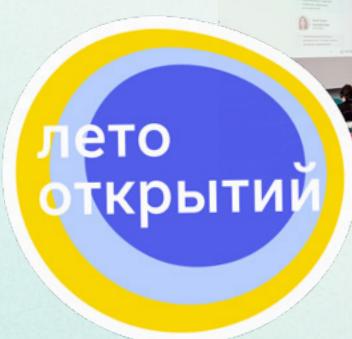
проектов

45

лекций
и семинаров

80

научных
постеров



Партнеры
проекта

SBER
AI Lab

ИТМО ×
Передовая
Инженерная
Школа
ИТМО

Научные
парнери

Skoltech

Высшая Школа
Экономики

МФТИ

МТУСИ



Журнальные клубы

Reading Clubs AIRI — это пространство для увлекательных дискуссий и обмена идеями. На них исследователи обсуждают новые, важные или просто любопытные научные статьи. На встречах участники глубже вникают в суть материалов, оценивают их значимость для своей области и расширяют собственный исследовательский горизонт.

Reading Club | ControlGenA

Лидирует:
Айбек Аланов

32

встречи
за 2024 год

Reading Club | Embodied AI

Лидирует:
Алексей Ковалёв

10

встреч
за 2024 год

Reading Club | CV and Robotics

Лидирует:
Влад Шахуро

31

встреча
за 2024 год

AIRI LEGO клуб

В этом году неформальный клуб по интересам сотрудников AIRI продолжил развиваться. Коллеги собирали Lego, знакомились, приглашали друзей и играли в квизы.



Хакатон SafeSpeak-2024

Обнаружение дипфейков
для обеспечения безопасности
голосовой связи

SafeSpeak2024 — хакатон, посвященный разработке технологий обнаружения аудио-спуфинга, нацеленный на решение актуальных проблем безопасной голосовой аутентификации и защиту биометрических систем от атак.

>240

человек подали
заявки на участие



Заявки на участие в соревновании подали молодые ученые из России, Эфиопии, Казахстана, Вьетнама, Бразилии, Китая, Узбекистана и Индии.



Хакатон «Типовой ИИ vs Базовые модели»

14 декабря состоялся финал хакатона «Типовой ИИ vs Базовые модели». За звание лучших сражались 23 команды из ведущих вузов России. Их главной задачей стало создание модели машинного обучения, способной автоматически обнаруживать патологию на рентгенограммах органов грудной клетки.



AI Journey Contest 2024

Исследователи AIRI подготовили 3 задачи на хакатон AI Journey Contest 2024:

- Emotional FusionBrain 4.0: создание мультимодальных моделей работы с видео, аудио и текстом.
- Multiagent AI: создание мультиагентной RL-системы, агенты которой смогут решать задачи, объединяясь в различные схемы кооперации.
- Embodied AI: создание роботов, которые смогут решать сложные задачи, требующие взаимодействия с окружающей средой и пользователем, а также общения с ними на естественном языке.



Night Photography Rendering Challenge 2024, NTIRE Workshop. CVPR 2024

Ученые AIRI и Института проблем передачи информации (ИППИ РАН) провели приуроченное к одной из наиболее престижных конференций по компьютерному зрению CVPR 2024 научное соревнование по рендерингу ночных фотографий. Его целью было с помощью ИИ-алгоритмов обработать ночное изображение с камеры смартфона и получить кадр фотографического качества.

58 команд приняли участие в соревновании, включая представителей академии (Гонконг, Сингапур, Вашингтон, Милан и др.) и индустрии (Samsung, XiaoMi, Honor, Doshua и др.).



О нас пишут и говорят

ТАСС, РИА, РБК, Forbes, Ведомости, Коммерсантъ, Российская Газета, Газета. Ru, Lenta.ru, N+1, Известия, Хайтек, Первый Канал, Россия 24, BFM, Код Дурова, Радио «Маяк», Радио России, телеканал «Культура» и многие другие.

Прочесть экспертные материалы от сотрудников AIRI



Forbes

Иван Оселедец о том, как оценивать работу DeepTech-исследователей



Forbes

Антон Конушин о том, кого брать в штат ИИ-стартапа и на какие навыки смотреть



РБК Тренды

Александр Панов о том, как ученые используют игры, чтобы обучать системы искусственного интеллекта решать реальные задачи.



РБК Тренды

Юрий Куратов о том, как устроена память нейросетей



RB

Максим Кузнецов о том, как выбрать нейросеть для бизнеса



RB

Александра Бройтман о том, как организовывать летние школы

О нас пишут и говорят



N+1

Партнерский материал
о том, что производит наука
и как она помогает бизнесу



Коммерсантъ

Евгений Фролов о том,
как исследователи
повысили точность
рекомендательных систем



Коммерсантъ

Иван Оседецов для Коммерсантъ
о том, как команда российских
математиков доработала выводы
нобелевского лауреата



TASS

Александр Панов о том,
что бизнес и потребители в
России выигрывают
от внедрения ИИ?



Коммерсантъ

Илья Макаров о том,
что малые нейросети
способны обучать большие
ИИ-модели лучше человека



Коммерсантъ

Айбек Аланов о том,
как редактирование
изображений поможет
науке

Новые партнерства за 2024 год



кибердом

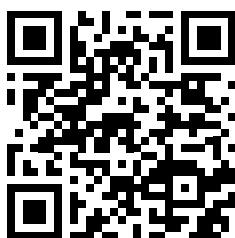


ФГБУ «НМИЦ им. В.А. Алмазова»
Минздрава России





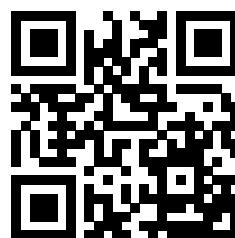
Ученые AIRI в социальных сетях



Ivan Oseledets'
Channel



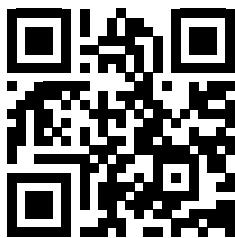
Telegram-канал
Ивана Оседедца



The Oleg



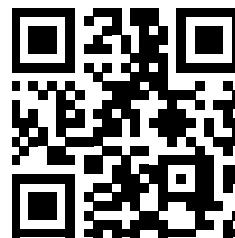
Telegram-канал
Олега Рогова



Kardymonchik
Channel



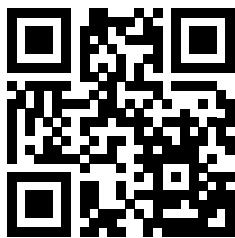
Telegram-канал
Ольги Кардымон



Complete AI



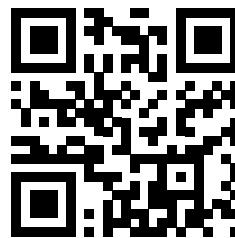
Telegram-канал
Андрея Кузнецова



AbstractDL



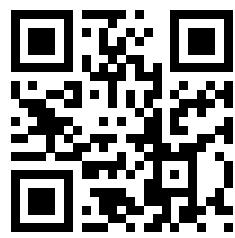
Telegram-канал
Антона Разжигаева



Grounding
Knowledge



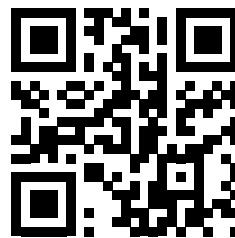
Telegram-канал
Александра Панова



Dendi Math&AI



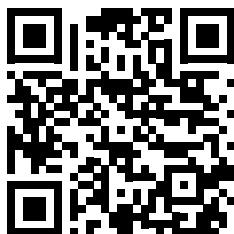
Telegram-канал
Дениса Димитрова



Чердачок Ктошика



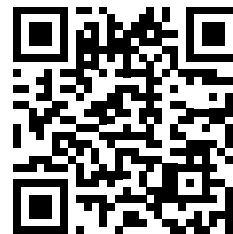
Telegram-канал
Антона Конушина



AI Brain



Telegram-канал
Айбека Аланова



Causality links



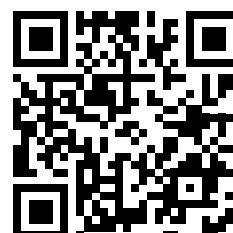
Telegram-канал
Владислава Куренкова



Labrats



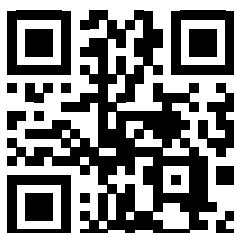
Telegram-канал
Дмитрия Пензара



ComputAgeChannel



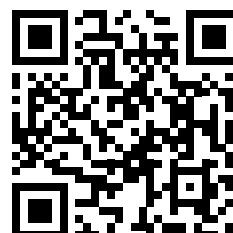
Telegram-канал
Дмитрия Крюкова



Embrace the data



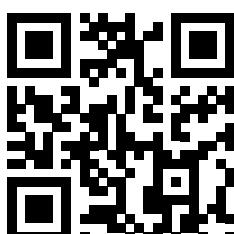
Telegram-канал
Аллы Чепуровой



Гречневые мысли



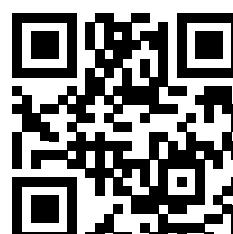
Telegram-канал
Никиты Сушко



BaseLine



Telegram-канал
Алексея Ковалева



Nygma's Diaries



Telegram-канал
Максима Жданова



Вениамин Фишман



Telegram-канал
Вениамина Фишмана



Марат пишет про науку



Telegram-канал
Марата Хамадеева

Контакты

Сайт

airi.net

Соцсети



[airi_research_institute](#)



[Airi_institute](#)



[habr](#)



[artificial-intelligence-research-institute](#)



[AIRIIInstitute](#)



[AIRI_inst](#)

Адрес

Москва, Пресненская набережная, д. 6, стр. 2