

**Современное состояние технологий
искусственного интеллекта.
Отечественная нейросетевая платформа Plat**

*Визильтер Юрий Валентинович, viz@gosniias.ru
начальник подразделения 3000 ФГУП «ГосНИИАС»,
д.ф-м.н., профессор РАН*

Контакт по Платформе Plat:

support.plat@gosniias.ru



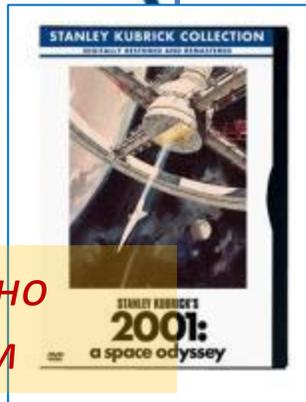
Семинар МГУ
Москва, 13.10.2020

Современное состояние технологий ИИ

- **Глубокие нейронные сети** совершили революцию в области компьютерного зрения и анализа больших данных
- **Функциональный искусственный интеллект** скоро будет создан в ходе второй волны «революции ИИ», которая происходит прямо сейчас
- **Глубокая оптимизация** уже в ближайшие годы распространит эту революцию практически на все области техники, технологии, экономики...
- **Актуальные результаты и открытые проблемы** угрозы, вызовы, надежды (на примерах из области компьютерного зрения и не только...)

2011-2016, 2017-2020+: две волны технологической революции в ИИ

Ранний период развития ИИ: созданы основные методы и подходы...



Завышенные ожидания

2000
«Зима ИИ»
Прогноз: 2040+

Разочарование

Рост интереса

Алгоритмы ИИ сильно проигрывают людям

«Кривая хайпа» для ИИ

2016-17

- Глубокие соревнующиеся сети
- ГО с подкреплением
- ГО с использованием моделей, баз знаний и логического вывода
- Автоматическое обучение ГНС

Вторая волна революции ГНС Все задачи ИИ (НИР)

Первая волна революции ГНС

Распознавание (НИР->ОКР)

Функциональный («слабый») ИИ
Прогноз: 2020+

- Появление глубоких нейронных сетей (ГНС)
- Решение задач компьютерного зрения, обработки сигналов и анализа больших данных на уровне человека и выше (superhuman)

2011-12

...2011+: в области ИИ началась технологическая революция ГНС

Глубокие нейронные сети
*(первая волна современной
технологической революции ИИ)*

Глубокие конволюционные нейронные сети – новое поколение алгоритмов обнаружения и распознавания объектов на изображениях



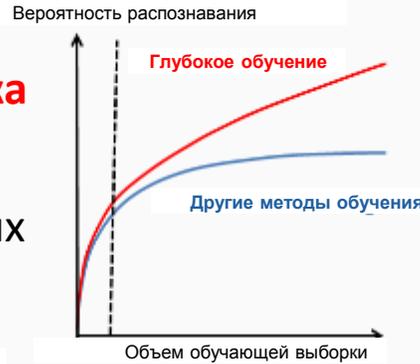
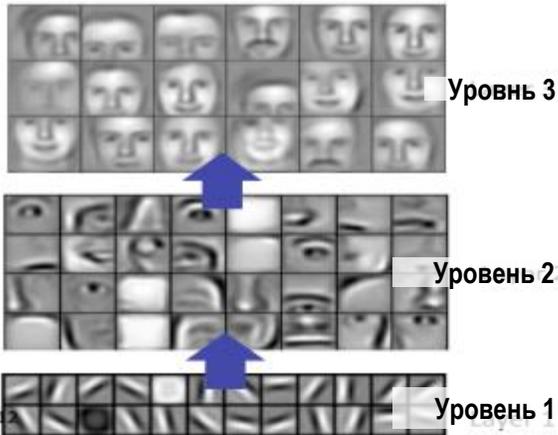
2011: First Superhuman Visual Pattern Recognition

<http://people.idsia.ch/~juergen/deeplearning.html>

2011: Автоматическое обнаружение и распознавание объектов на базе глубоких нейронных сетей (ГНС)

+ С 2011 г. - **распознавание образов на уровне человека или выше** (superhuman)

+ Обучение на сверхбольших объемах данных



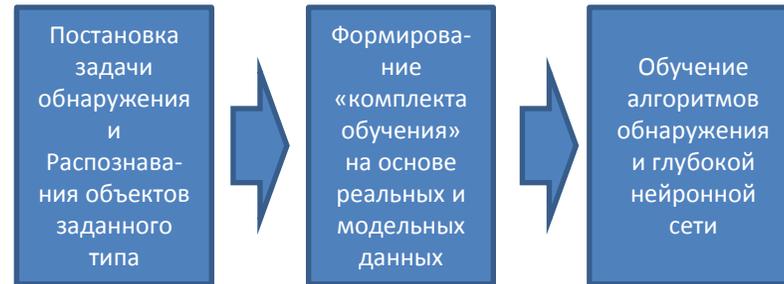
+ Иерархическое обучение с повышением абстракции данных от уровня к уровню

Достоинства и проблемы

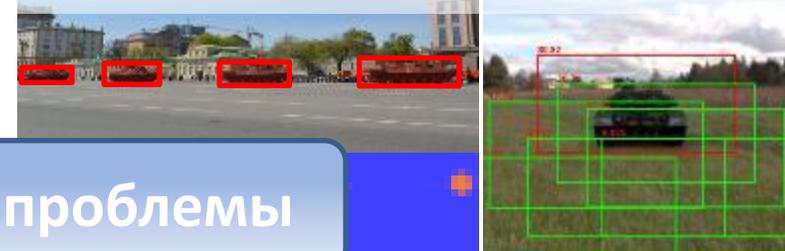
+ Тысячи слоев нейронов
+ Учет специфики изображений как объекта распознавания (локальность, инвариантность к сдвигу, нечеткая локализация)



- Нужны огромные обучающие выборки
- Длительное моделирование и обучение



- Ресурсоемкость, низкая скорость
- Необходимо быстрое предобнаружение



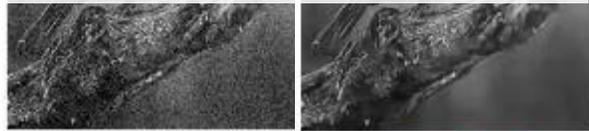
- Необходимость эффективных алгоритмических реализаций и нового поколения бортовых нейропроцессоров



2015-16: ГНС решают все задачи компьютерного зрения



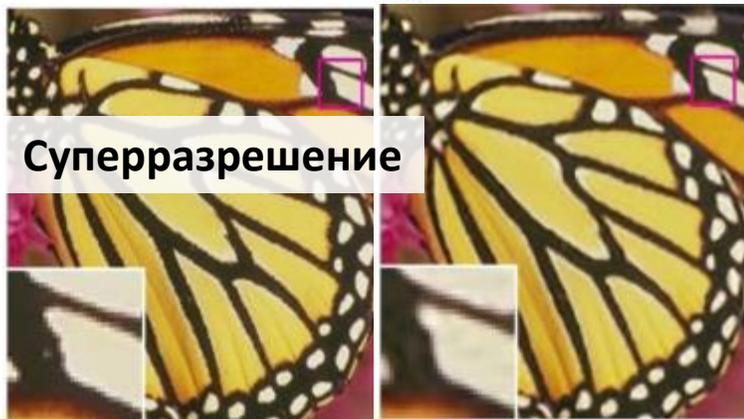
Фильтрация шумов



noisy ($\sigma = 25$) PSNR: 20.16dB ours: PSNR: 30.03dB



Удаление смаза



Суперразрешение

Original / PSNR

SRCNN / 27.95 dB

Ours



Обнаружение особых точек на лицах



Фронтализация лиц без 3D моделей



Распознавание лиц в сложных условиях (ГосНИИАС-2015)

Применение методов глубокого обучения к задаче распознавания лиц



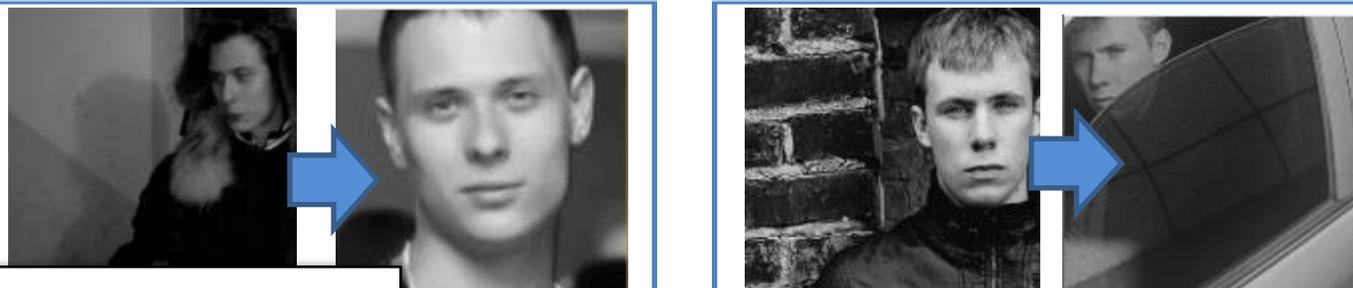
~30 млн
сравнений в
секунду для 2000
битного шаблона

Сверхкомпактные шаблоны, высокая скорость поиска



~ 200 млн
сравнений в
секунду для 200
битного шаблона

Качество, сравнимое с результатами человека-оператора



Распознавание LFW:
200 бит ~0.95
2000 бит ~0.987

Что дают нам технологии глубокого обучения, распознавания образов и анализа данных?

СИСТЕМЫ ВИДЕОНАБЛЮДЕНИЯ

БИОМЕТРИЧЕСКИЕ СИСТЕМЫ

ПРОМЫШЛЕННАЯ БЕЗОПАСНОСТЬ И АВТОМАТИЗАЦИЯ

АНАЛИЗ БОЛЬШИХ ДАННЫХ

ПРОШЛОЕ

НАСТОЯЩЕЕ

БУДУЩЕЕ

СИСТЕМЫ ВИДЕОНАБЛЮДЕНИЯ

ПРОШЛОЕ



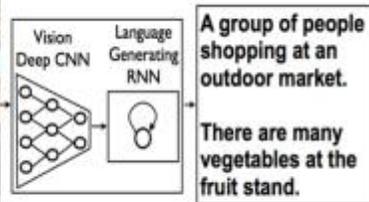
✓ **Обнаружение движущихся объектов в несложных условиях**

✓ **Классификация движущихся объектов**

✓ **Обнаружение принесенных (оставленных) и унесенных предметов**

✓ **Расчет скорости и плотности потока (людей или машин)**

✗ **Распознавание объектов**



НАСТОЯЩЕЕ

✓ **Распознавание типа оставленного предмета как подтверждение детектора**

✓ **Классификация движущихся объектов на основе глубоких нейронных сетей**

✓ **Реидентификация персоны с разных камер по одежде с использованием CNN**

✓ **Детектирование объектов (например, людей, авто) на отдельных кадрах**

✓ **Распознавание номеров на основе нейросетей**

БУДУЩЕЕ

✓ **Анализ и распознавание поведения людей**

✓ **Высоконадежное обнаружение и отслеживание объектов, в том числе с PTZ-камер**

✓ **Появление «умных» камер с алгоритмами детектирования на основе нейронных сетей на борту**

✓ **ИИ: система видеонаблюдения самостоятельно генерирует сообщения оператору с подробным описанием событий**

✓ **ИИ: система видеонаблюдения самостоятельно принимает решения**

БИОМЕТРИЧЕСКИЕ СИСТЕМЫ

ПРОШЛОЕ



- ✔ Обнаружение лиц с помощью группы методов Виолы-Джонса
- ✔ Распознавание лиц с использованием обученных признаков
- ✘ Обнаружение и распознавание лиц в сложных условиях



НАСТОЯЩЕЕ

- ✔ Детектирование лиц продвинутыми методами предыдущего поколения
- ✔ Детектирование лиц в сложных условиях с помощью нейросетей (не real-time)
- ✔ **Распознавание лиц по изображениям высокого качества с очень высокими вероятностями**
- ✔ **Распознавание лиц в сложных условиях с достаточно высокими вероятностями**

БУДУЩЕЕ

- ✔ Повсеместная замена электронных ключей на распознавание лиц
- ✔ Высокие точности определения пола, возраста, национальности, эмоций в задачах маркетинга
- ✔ Детектирование и распознавание людей в толпе по большим базам с высокой вероятностью - единая биометрическая глобальная система в рамках городов
- ✔ Коммерческое внедрение оплаты по лицу

ПРОМЫШЛЕННАЯ БЕЗОПАСНОСТЬ И АВТОМАТИЗАЦИЯ

ПРОШЛОЕ



- ✔ **Обнаружение движения или движущихся объектов в запрещенных зонах (sterile zone)**
- ✔ Контроль температуры установок с помощью тепловизионных камер
- ✔ Замена человека в простейших действиях автоматическими системами (задачи подсчета продукции, считывания маркировок/штрихкодов и пр.)

НАСТОЯЩЕЕ

- ✔ Распознавание средств индивидуальной защиты (чувствительно к марке)
- ✔ **Задачи технологического контроля:** определение качества продукции, правильности течения технологических процессов и пр.
- ✔ Поиск полезных ископаемых по большим данным – набору показаний группы датчиков
- ✔ Автоматизированные заводы с минимальным количеством персонала

БУДУЩЕЕ

- ✔ Автоматические системы контроля за производством
- ✔ **Автоматические интеллектуальные системы управления заводами:** анализ показаний датчиков системы, предсказание чрезвычайных ситуаций, принятие решений
- ✔ **Создание автономных подвижных роботов** нового поколения с системой технического зрения на борту

АНАЛИЗ БОЛЬШИХ ДАННЫХ



Объем данных, собранных только с двигателей коммерческих реактивных самолетов США в течение года

= 1 041
600 500 ТВ

НАСТОЯЩЕЕ - БУДУЩЕЕ

- ✓ Анализ состояния бортового авиационного оборудования в масштабе реального времени с возможностью прогнозирования различных типов отказов.
- ✓ Совершенствование средств прогнозирования времени прибытия воздушных судов в составе систем управления воздушным движением.



НАСТОЯЩЕЕ - БУДУЩЕЕ

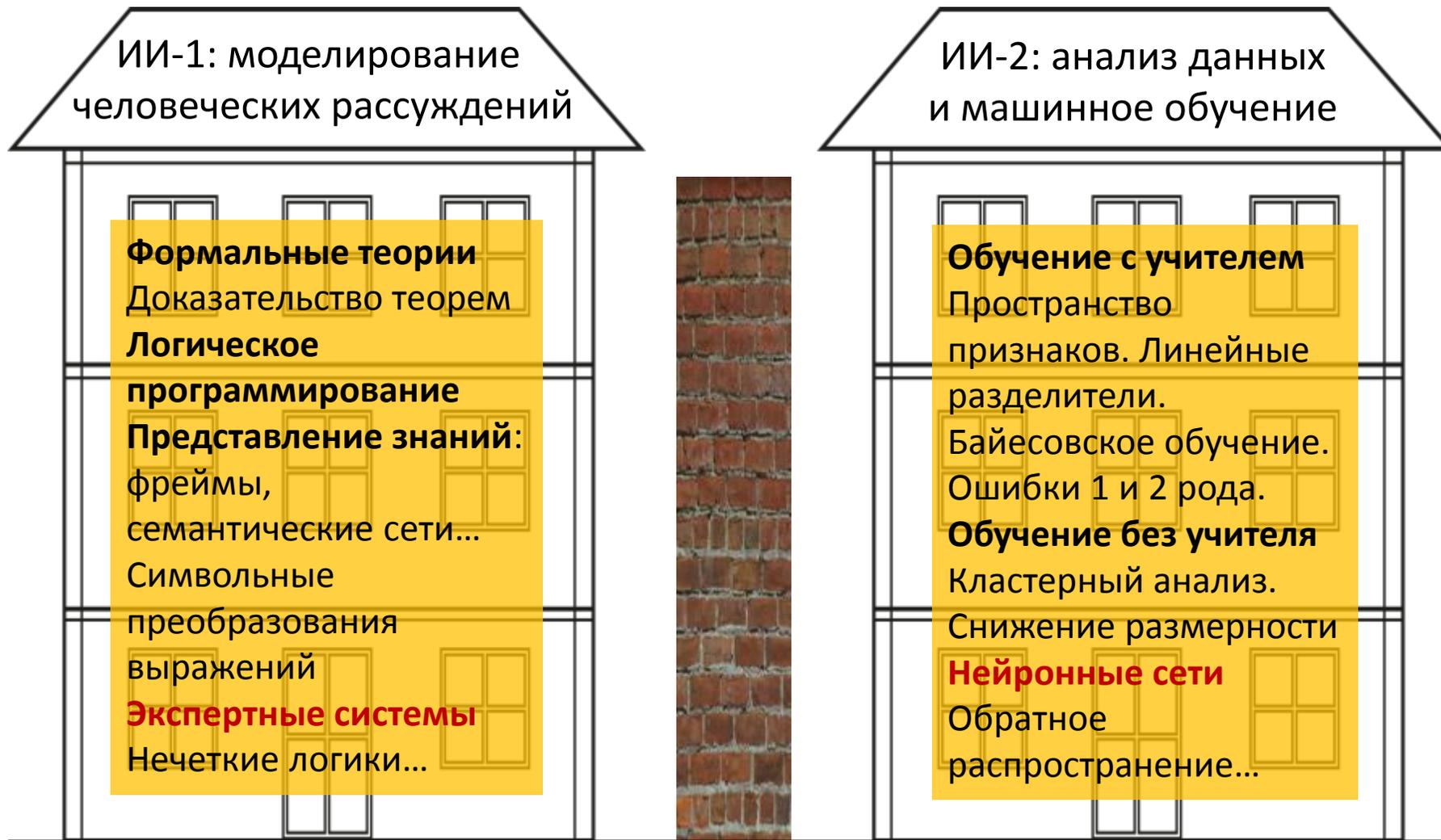
- ✓ Глубокие нейронные сети в инвестировании
- ✓ Предсказание эффективности компаний в терминах выручки, операционных доходов и чистой прибыли на основе финансовой информации и патентах компаний
- ✓ Предсказание поведения фондового рынка

и еще множество приложений...

**Функциональный
искусственный интеллект**
*(вторая волна современной
технологической революции ИИ)*

2000: Что такое функциональный ИИ, из чего он состоит?

Функциональный «ИИ» = АО/ПО, способные автоматически выполнять полезные функции, которые ранее могли быть выполнены только человеком.

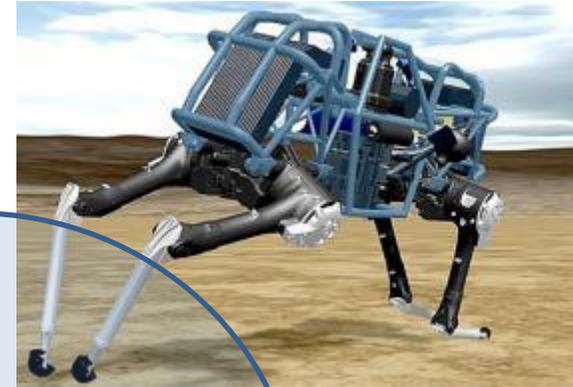


Два «дома», разделенных прочной стеной: ИИ-1 не учится, ИИ-2 не рассуждает...

2000: Прогноз создания функционального ИИ (на примере робототехники)



*Везде
алгоритмы
сильно
проигрывают
людям*



Навигация

**Обработка
сенсорных
данных**
(зрение,...)

Управление
(планирование
оптимизация,
игры,...)

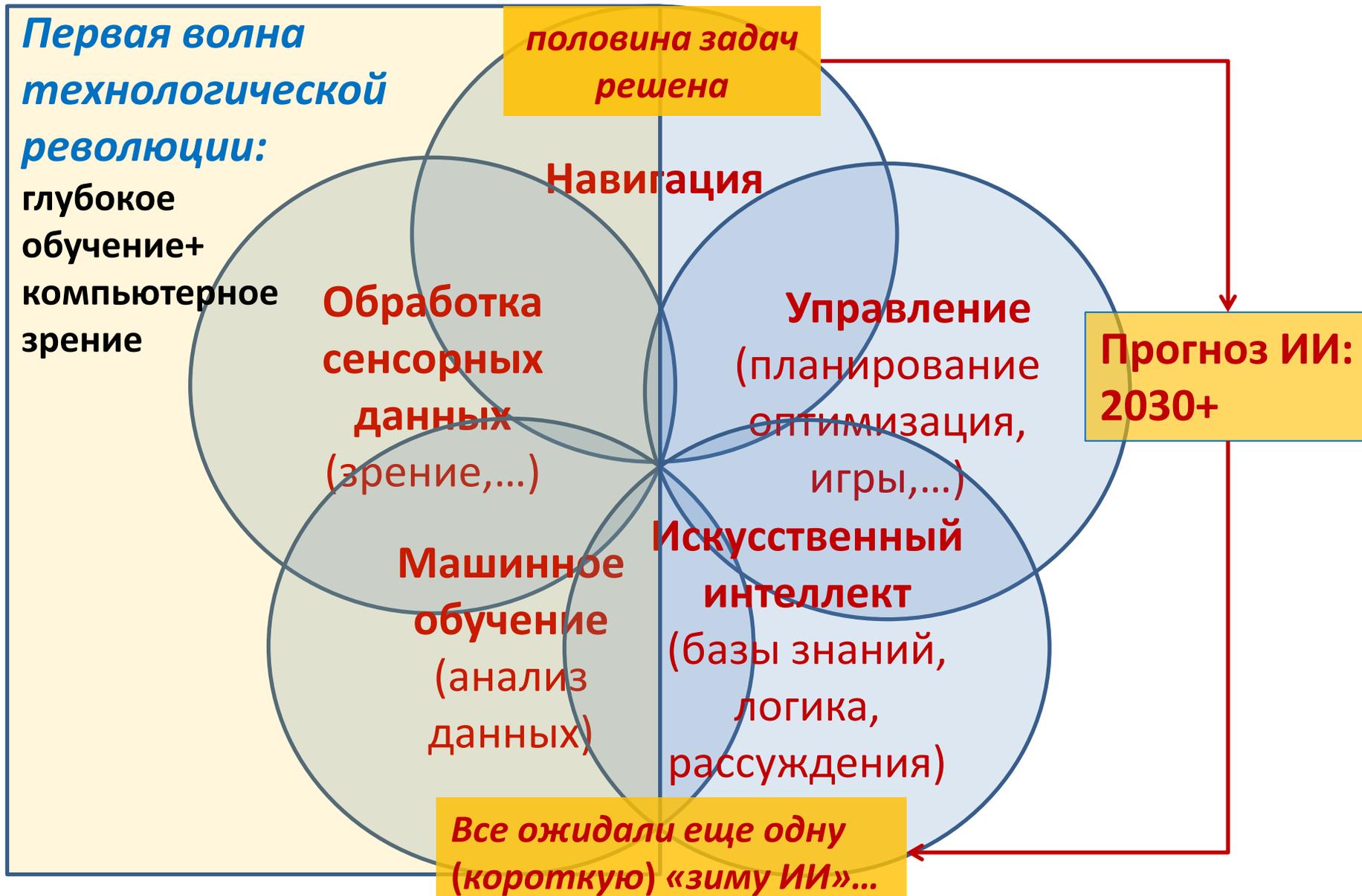
**Прогноз ИИ:
2040+**

**Машинное
обучение
(ИИ-II)**
(анализ
данных)

**Искусственный
интеллект
(ИИ-I)** (базы
знаний, логика,
рассуждения)



2015: Прогноз создания функционального ИИ (на примере робототехники)

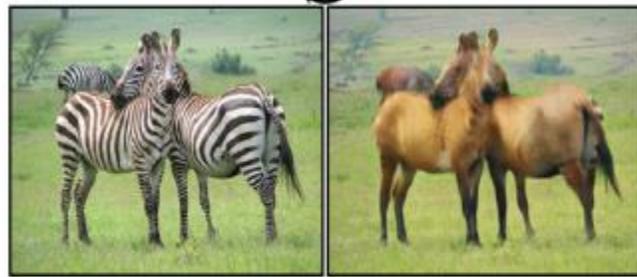


Компьютерное зрение и машинное обучение для интеллектуальных систем

(2017+, вторая волна технологической революции)

- **Глубокие соревнующиеся сети для имитации данных**
GAN, Domain Transfer Learning, Zero-Shot Learning
- **Интерпретация динамической визуальной информации на естественном языке**
Action Detection and Prediction, Video Annotation, Video and Language Understanding, Text-to-Video, VQA
- **Обучение глубоких сетей с подкреплением как активных агентов** Reinforcement Learning, Lifelong Learning
- **Глубокое обучение с использованием структурных моделей, баз знаний и программ логического вывода**
Graph Structured CNN, Deep Visual Reasoning
- **Автоматическое конструирование и обучение глубоких сетей**

Generative Adversarial Networks (GANs)



zebra → horse



apple → orange



summer → winter



horse → zebra



orange → apple



winter → summer

**GAN – сеть,
обладающая
воображением!**

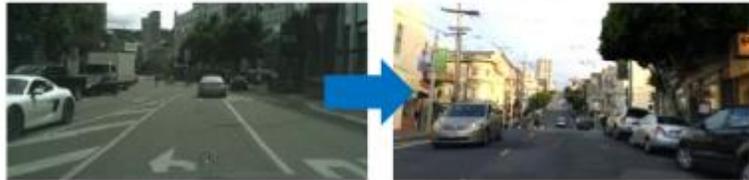


Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, Jun-Yan Zhu et al., ICCV, 2017

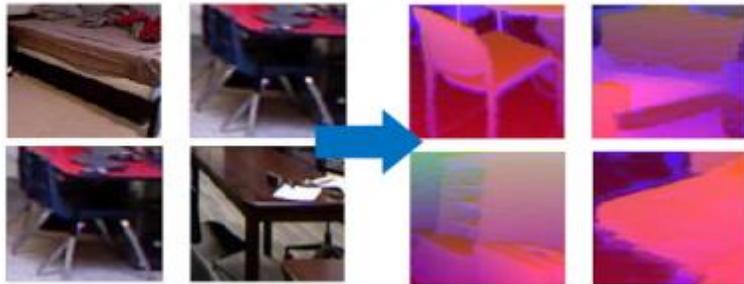
Генеративные конкурирующие сети

Перенос обучения в новую область применения
(Domain Transfer Learning)

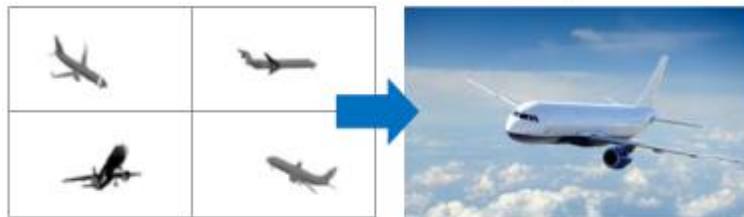
From dataset to dataset



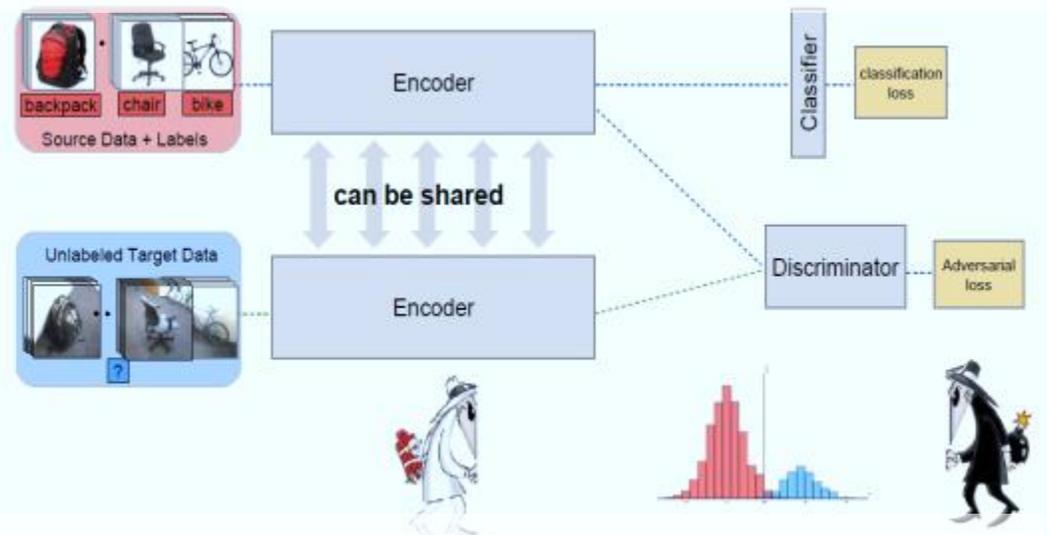
From RGB to depth



From CAD models to real images



Генератор создает визуальные образы, стараясь обмануть Дискриминатор...



....Дискриминатор старается отличить фантазии Генератора от реальности

Понимание сцены и языка: Visual Question Answering

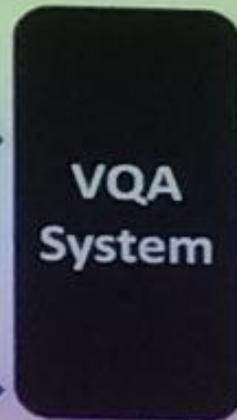
Many questions can be asked about an image

- Is it sunny?
- Is it safe to cross the street?
- How many cars are parked on the road?

Вопросы
самых
различных
типов



How does the person
in the middle feel?

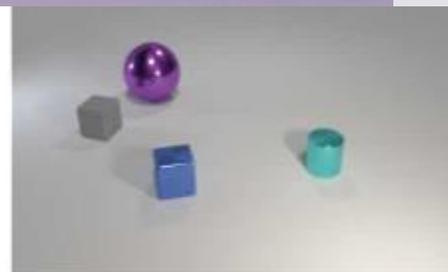


Happy

Вопросы, требующие
понимания контекста



Q: Is there a blue box
in the items? A: yes



Q: What shape object
is farthest right?
A: cylinder



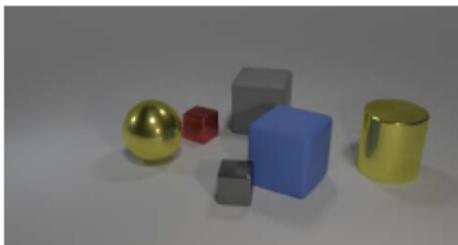
Q: Are all the balls small?
A: no



Вопросы,
требующие рассуждений

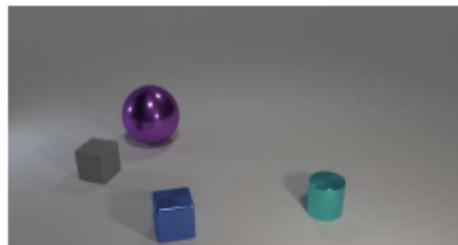
Q: Is the green block to the
right of the yellow sphere?
A: yes

Deep Visual Reasoning for VQA: генератор программ

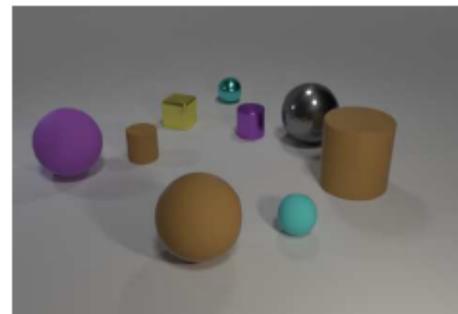


Вопросы, требующие рассуждений:

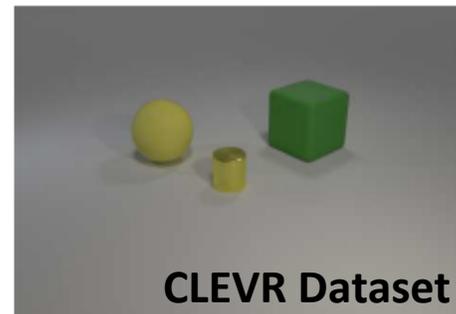
Q: *Is there a blue box in the items?* A: *yes*



Q: *What shape object is farthest right?*
A: *cylinder*

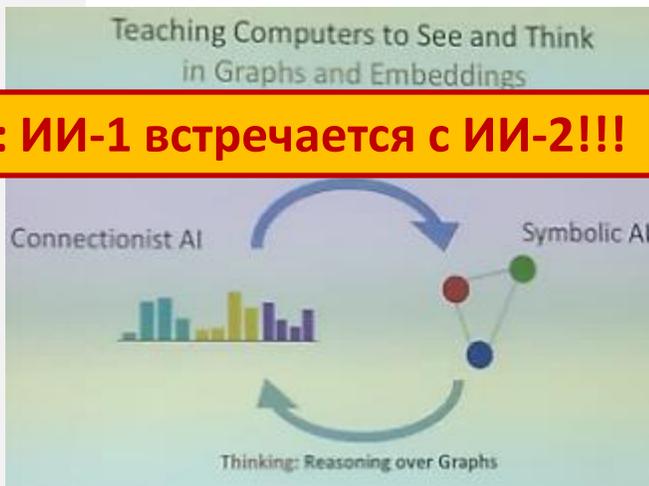


Q: *Are all the balls small?*
A: *no*



CLEVR Dataset

Q: *Is the green block to the right of the yellow sphere?*
A: *yes*



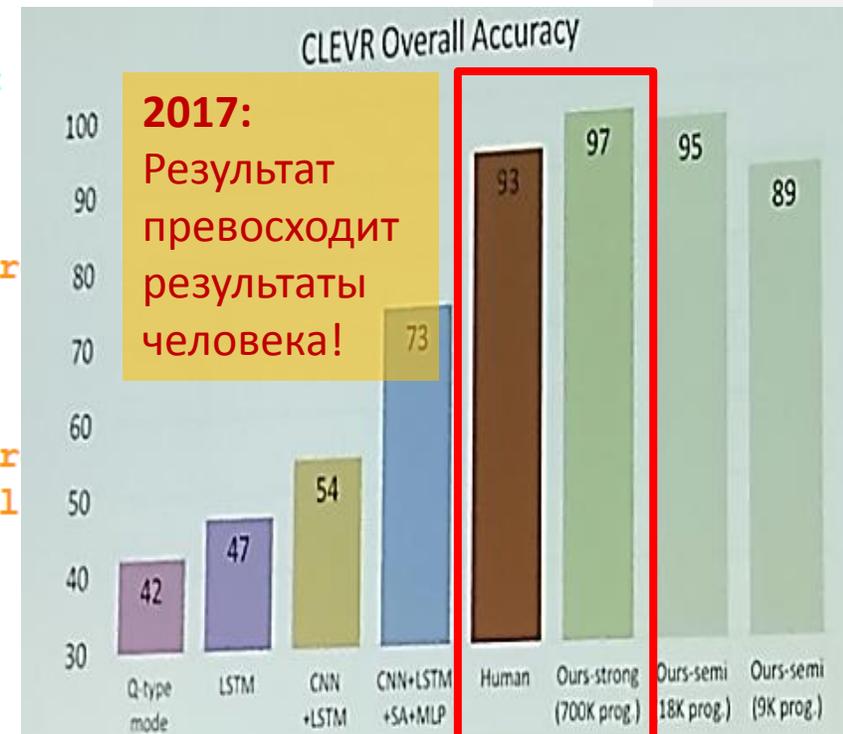
2017: ИИ-1 встречается с ИИ-2!!!

Predicted Program:

```
query_shape
unique
relate [right]
unique
filter_shape [cylinder]
filter_color [blue]
scene
```

Predicted Program:

```
equal_size
query_size
unique
filter_shape [spher
scene
query_size
unique
filter_shape [spher
filter size [small
```



2017:
Результат
превосходит
результаты
человека!

2017+: Применение ГНС к символическим данным позволяет использовать базы знаний и логический вывод для анализа данных о реальном мире

2019+: Прогноз создания функционального ИИ (на примере робототехники)

**Вторая волна
технологической
революции:**

глубокое обучение+
компьютерное зрение+
базы знаний+
семантические модели+
системы логического вывода+
автоматическое программирование+
общение с человеком на естественном языке+
обучение с подкреплением

Обработка сенсорных данных
(зрение,...)
Машинное обучение
(анализ данных)

Навигация

Управление
(планирование, оптимизация, игры,...)

Искусственный интеллект
(базы знаний, логика, рассуждения)

победа в го и Starcraft II +
автоматическое обучение



**Прогноз ИИ:
2020+**

**Все необходимое
для автономных
систем!**

Глубокая оптимизация
*(на пути к революции
в технике, технологии и
управлении бизнес-процессами)*

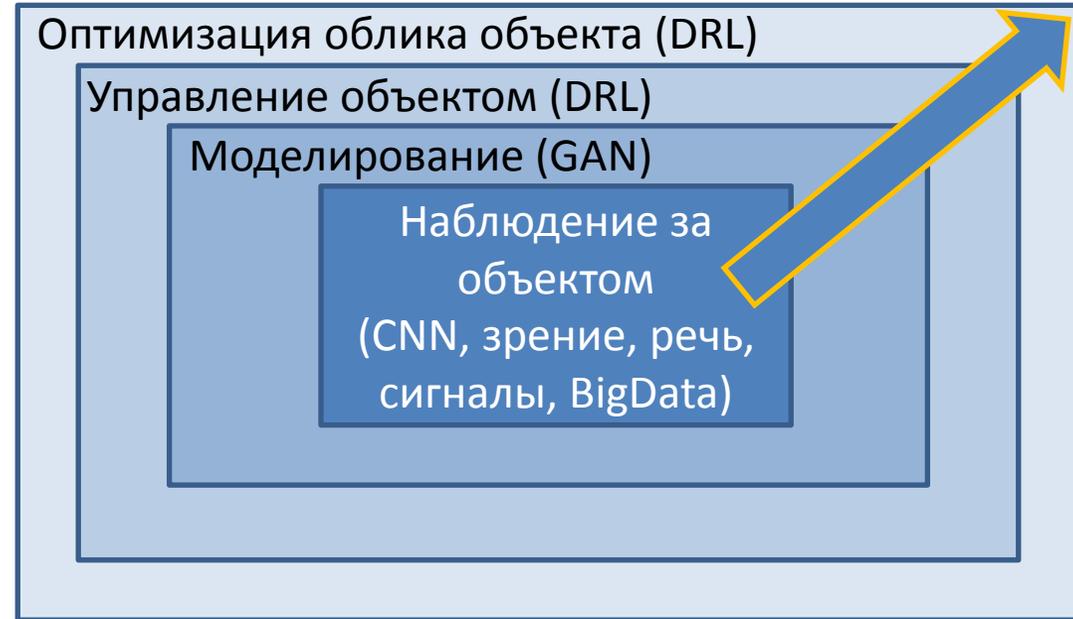
Глубокая оптимизация: **взгляд-2019+** на революцию ГНС и ее перспективы

2020+:
при помощи ГНС может быть создан «слабый» бортовой ИИ... но это не главное - ГНС могут дать науке и технике гораздо больше!

На самом деле нет никакого прорыва в методах ИИ - есть прорыв в методах локальной оптимизации, связанный с использованием ГНС.

Мы имеем дело не с новыми методами и подходами в ИИ, а с новой группой мощных инженерных методов: **«глубокими» методами моделирования, управления и оптимизации.**

Логика развития этих методов и технологий в последние годы состоит в **переходе от задач обработки и анализа информации к задачам управления и оптимизации.**



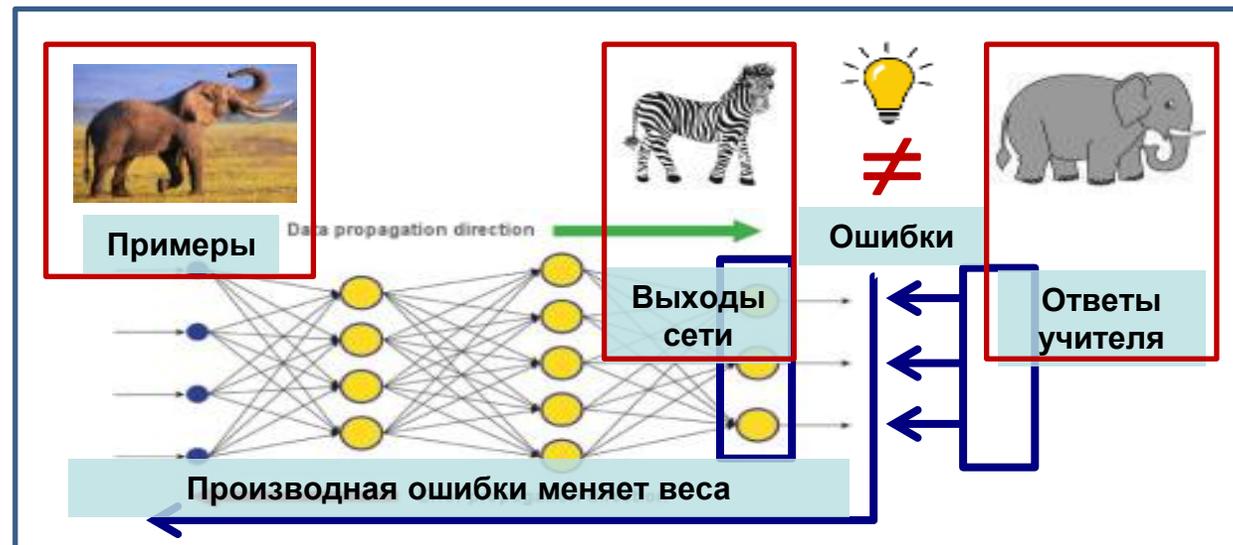
ГНС могут решать любые задачи, которые формулируются как задачи оптимизации.



Перспективные интеллектуальные системы на основе ГНС это уже не только СТЗ и даже не только «бортовой интеллект»! Это про все задачи создания и применения (от облика и алгоритмов до оптимизации производства и эксплуатации) изделий.

Переход от глубокого обучения к глубокой оптимизации

Поскольку применение ГНС в оптимизационных задачах машинного обучения оказалось чрезвычайно эффективным, их начали активно применять и для решения других задач оптимизации...



Обучение = Оптимизация: минимизация ошибок на выборке

Технологии глубокого обучения в классических задачах оптимизации

2015+: Глубокие нейросети умеют работать на графах
2015+: ГНС на графах могут выучивать эффективные эвристики решения задач оптимизации.

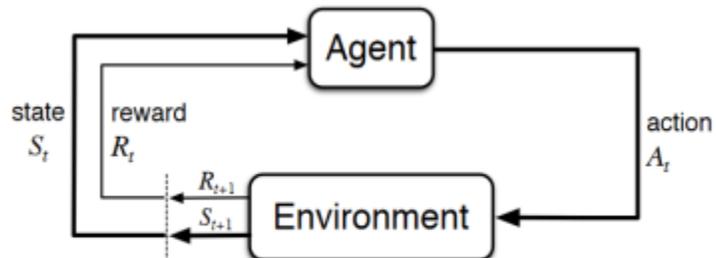
Deep Reinforcement Learning: обучение с подкреплением



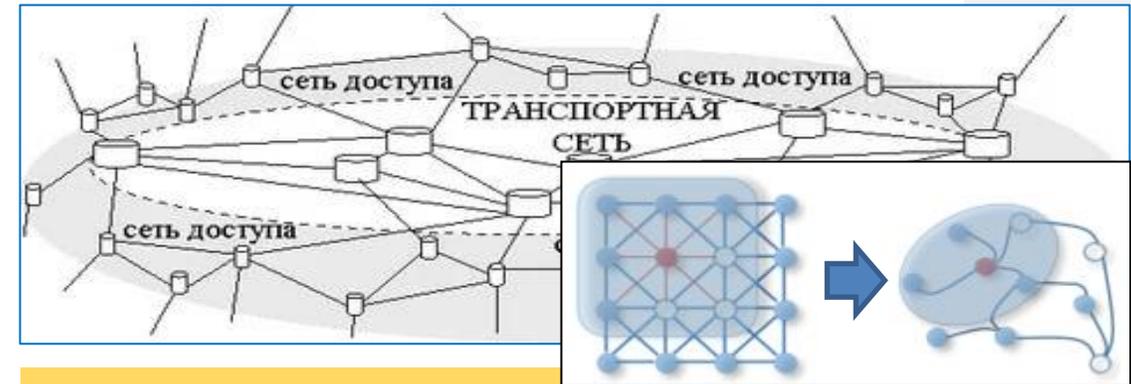
- Reward $R(t)$: score you earned at current step
Вознаграждение (выигрыш после хода)
- State S : current screen
Состояние (что видим на экране)
- Action i : move your board left / right
Действие (что делаем)
- Action value function $\hat{Q}(S, i)$: your predicted future total rewards
Стоимость (выигрыш в будущем)
- Policy $\pi(s)$: How to choose your action
Решение (какой ход выбрать)

Главный ключ к решению задач глубокой оптимизации: **глубокое обучение с подкреплением для выучивания эвристик!**

Как научить глубокую сеть (ИИ - агента) выучивать эвристики

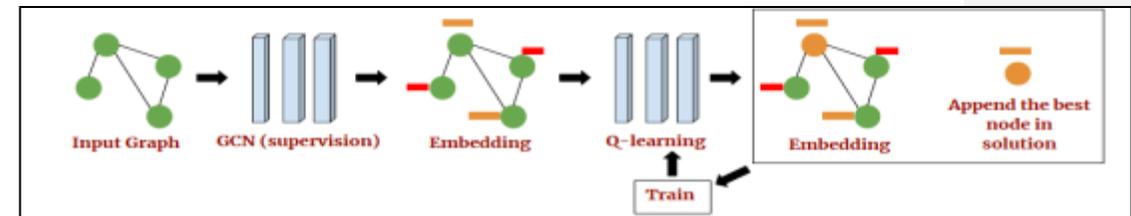


Deep Graph Embedding: глубокие сети на графах



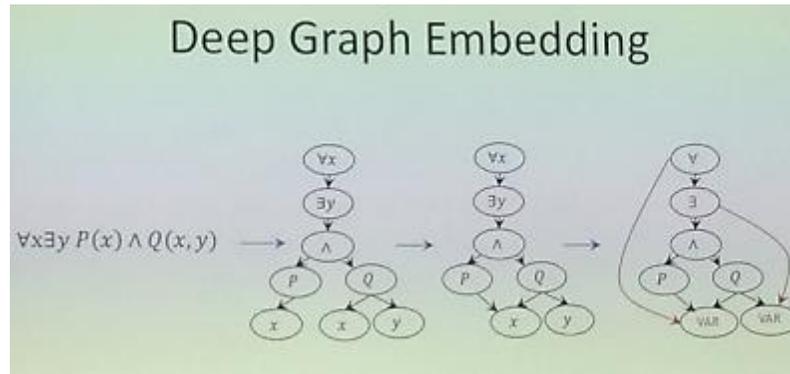
Опишем малые подграфы признаками, соберем из них признаки больших подграфов, и так – пока не опишем вектором признаков весь граф

Deep Graph Embedding + Deep Reinforcement Learning



Сегодня (2019+) задачи глубокой оптимизации решаются на графах с миллионами вершин, что позволяет уже переходить к практическому внедрению в самых масштабных приложениях

Deep Graph Embedding + Deep Reinforcement Learning в автоматическом доказательстве теорем (2017-2019)



HolStep [Kaliszyk et al. 2017]

- Benchmark for machine learning for Theorem Proving
- 2M+ conjecture-fact pairs of higher-order logic statements

Conjecture: $\forall \alpha \forall \beta (\sin(\alpha) = \sin(\beta)) = ((\alpha = \beta) \vee (\alpha = \pi - \beta))$

Relevant fact: $\forall \alpha \forall \beta \sin(\alpha - \beta) = \sin(\alpha)\cos(\beta) - \sin(\beta)\cos(\alpha)$

Irrelevant fact: $(x > 0) \wedge (y > 0) \rightarrow (xy > 0)$

	Sequence embedding	Graph embedding	
	CNN [Kaliszyk et al. '17]	CNN-LSTM [Kaliszyk et al. '17]	Ours
Accuracy	82	83	90.3

Более современные работы:

- Urban, J., Kaliszyk, C., Michalewski, H., and Olšák, M. (2018). Reinforcement learning of theorem proving. In *NIPS*.

<https://arxiv.org/abs/1805.07563>

- Automated Theorem Proving in Intuitionistic Propositional Logic by Deep Reinforcement Learning (2108)

<https://arxiv.org/abs/1811.00796>

- HOList: An Environment for Machine Learning of Higher-Order Theorem Proving (extended version) (2019)

<https://arxiv.org/abs/1904.03241>

- <https://github.com/tensorflow/deepmath>

На самом деле именно глубокая оптимизация и позволила объединить обучаемость ИИ-2 с интеллектуальностью ИИ-1

Технологии глубокого обучения в исследовании операций и управлении

AI Beats a Fighter Pilot in a Virtual Dogfight (2016)



AI ALPHA, built by Psibernetix, Inc. with Air Force Research Laboratory



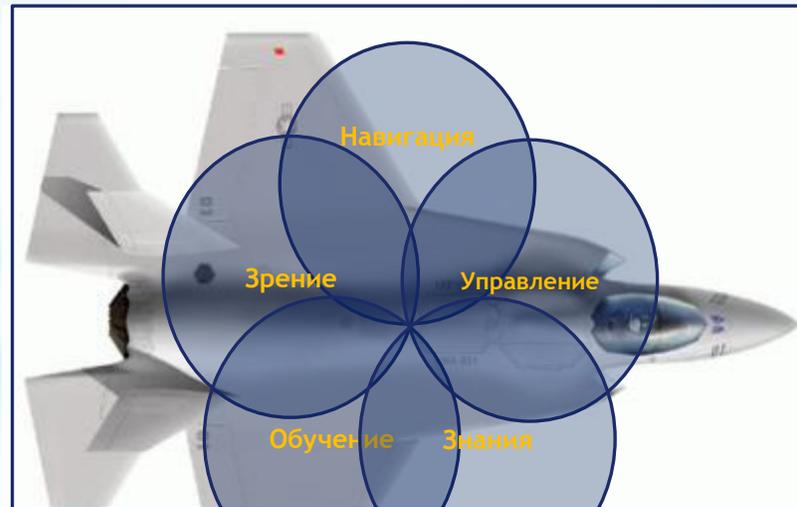
Оперативно-тактическое и групповое управление

A timeline showing the evolution of AI in games. It starts with a Go board labeled '1997', followed by a StarCraft II game labeled '2015-17', and then another Go board labeled 'ГНС AlphaGo' and 'ГНС AlphaZero'. Blue arrows indicate the progression from 1997 to 2015-17, and from 2015-17 to the AlphaGo/AlphaZero era. Below the timeline is a yellow banner with text.

Показатель зрелости: игровые задачи

StarCraft II - полноценный тактический военный симулятор с упрощённой моделью ведения боя.

Прогноз 2018: ГНС смогут выиграть у людей лет через 5



Автономное управление и оптимизация облика изделий



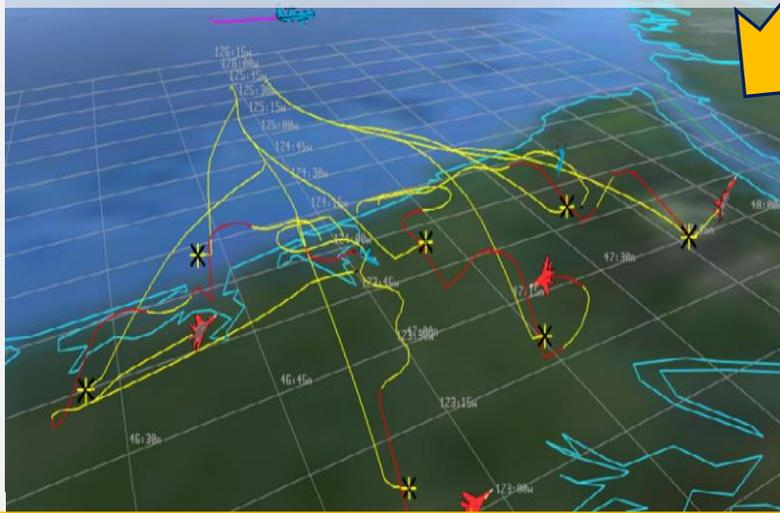
Оптимизация бизнес-процессов и производства

Технологии глубокого обучения в исследовании операций и управлении

AI Beats a Fighter Pilot in a Virtual Dogfight (2016)



AI ALPHA, built by Psibernetix, Inc. with Air Force Research Laboratory



Оперативно-тактическое и групповое управление



Показатель зрелости: игровые задачи

Сочетание технологий ГНС, RL и генетического отбора

2019+: Дорога открыта!



Автономное управление и оптимизация облика изделий

StarCraft II - полноценный тактический военный симулятор с упрощённой моделью ведения боя.

~~Прогноз 2018: ГНС смогут выиграть у людей лет через 5~~

25.01.2019: DeepMind AlphaStar со счетом 11:1 победила ведущих профессиональных игроков в StarCraft II

ХИТ ИИ 2019!

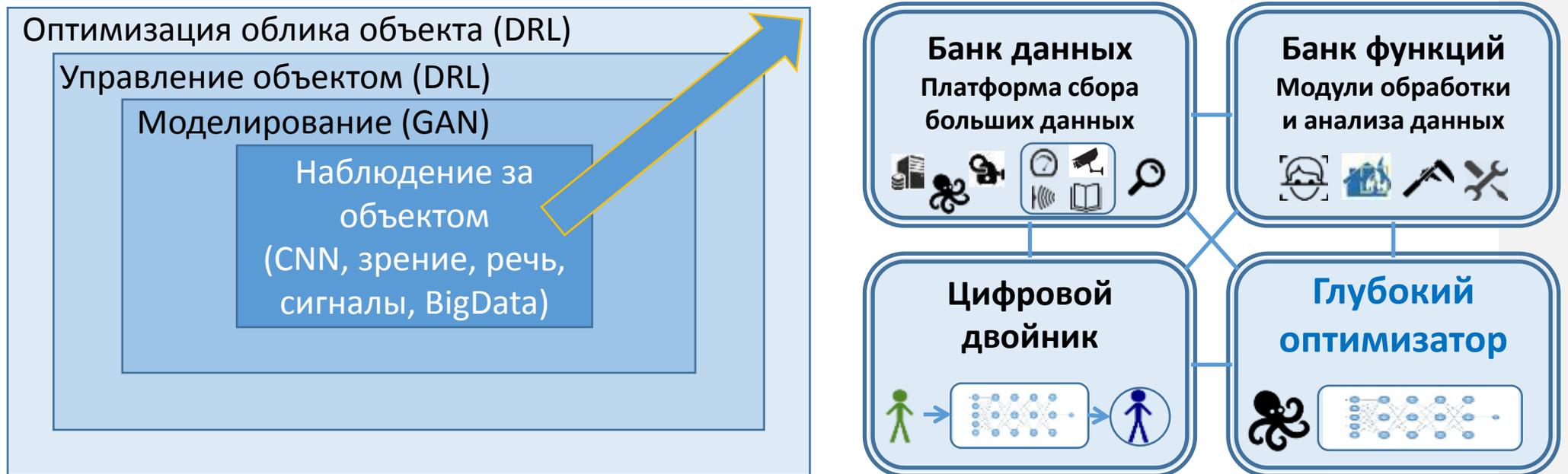


Оптимизация бизнес-процессов и производства

«Глубокая оптимизация» как второй этап «цифровизации»

После того, как для некоторого объекта или системы создана достаточно точная действующая модель («цифровой двойник»), эта модель сразу может быть использована для оптимизации этого объекта, либо формирования наилучших алгоритмов взаимодействия с ним. Это в равной степени относится к изделиям, их группам, войсковым соединениям или же промышленным предприятиям.

Логика развития и внедрения «глубоких» технологий: от обработки информации к управлению и оптимизации



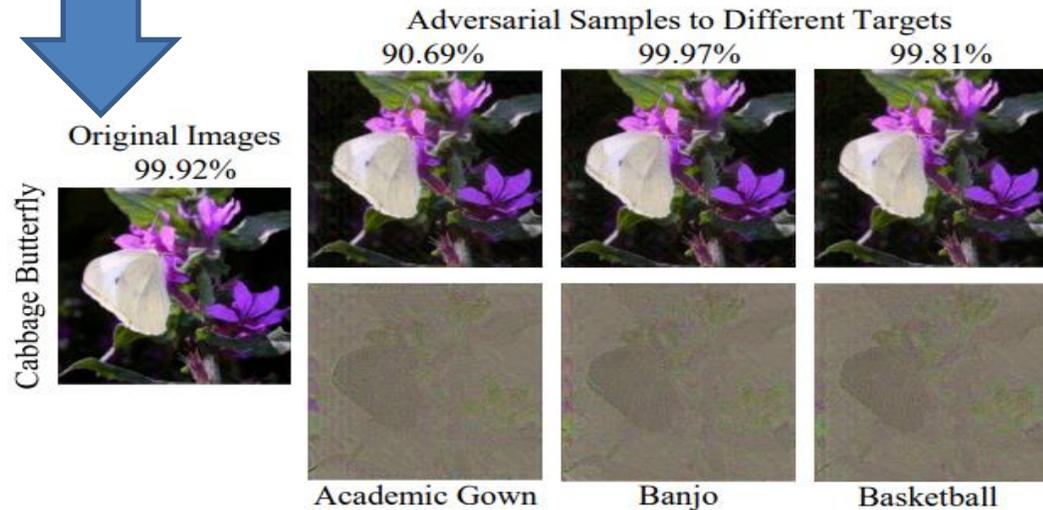
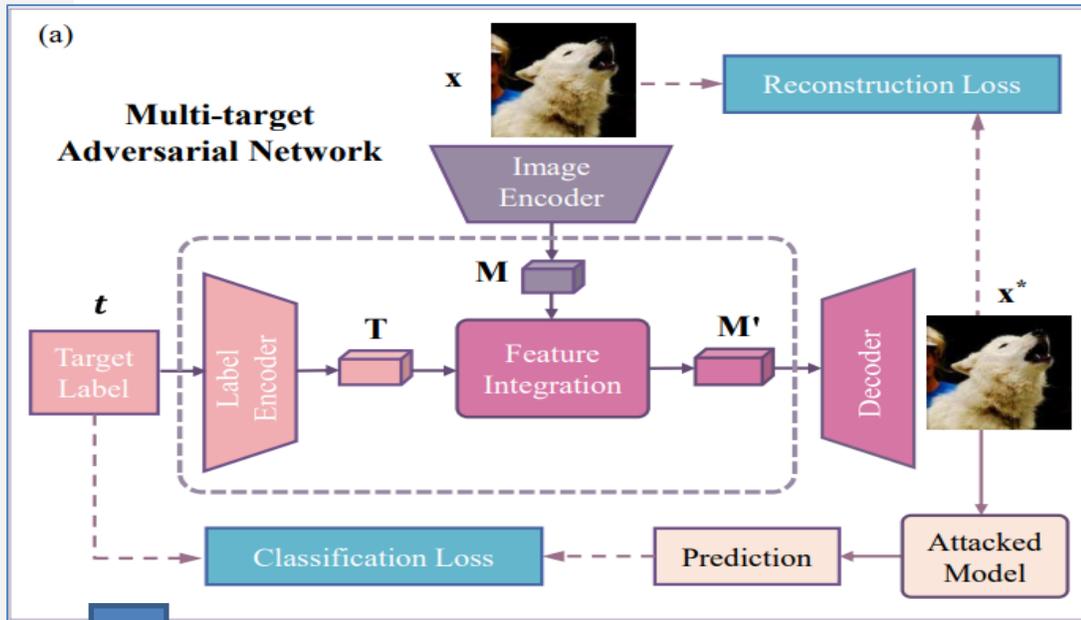
Таким образом, **этап создания и применения интеллектуальных средств «глубокой оптимизации» изделий и процессов на основе ГИС является логическим продолжением этапа создания «цифровых двойников»**

Актуальные результаты и открытые проблемы (угрозы, вызовы, надежды)

*На примерах из области компьютерного зрения
и не только...*

Атаки на нейросети (Adversarial Attack)

Атаки на распознающие нейронные сети

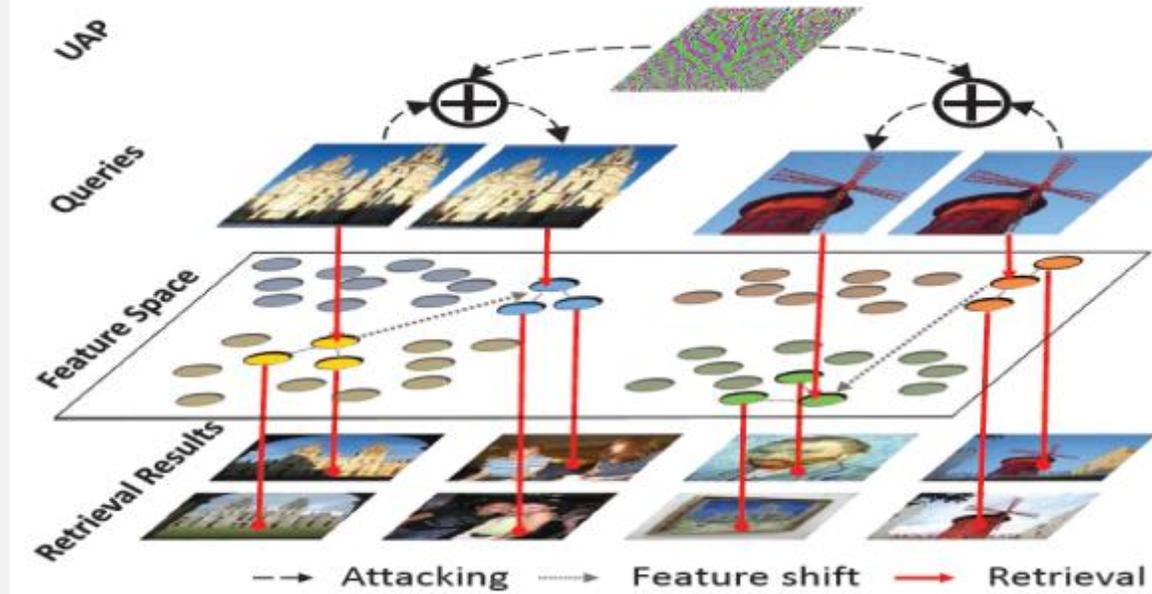


Once a MAN: Towards Multi-Target Attack via Learning Multi-Target Adversarial Network Once, ICCV-2019, Jiangfan Han, Xiaoyi Dong, Ruimao Zhang, Dongdong Chen, Weiming Zhang, Nenghai Yu, Ping Luo, Xiaogang Wang

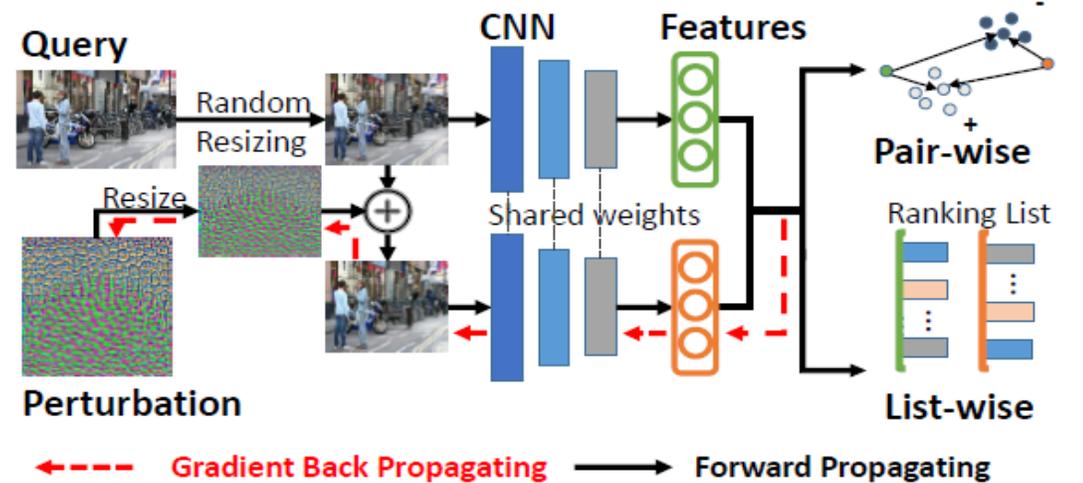
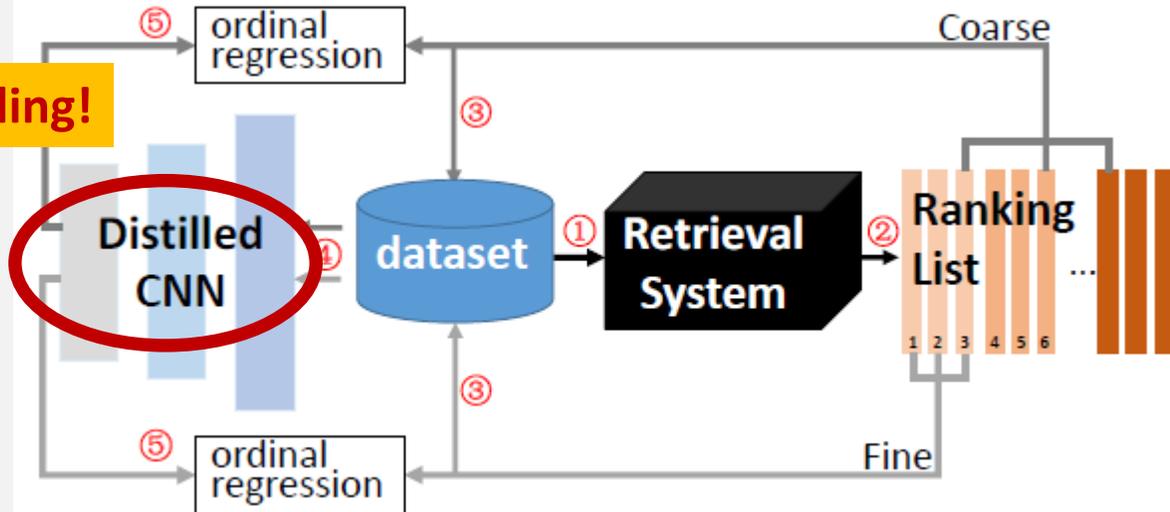


Sparse and Imperceivable Adversarial Attacks
 Francesco Croce, Matthias Hein
 University of Tübingen

Атаки на поисковые нейронные сети

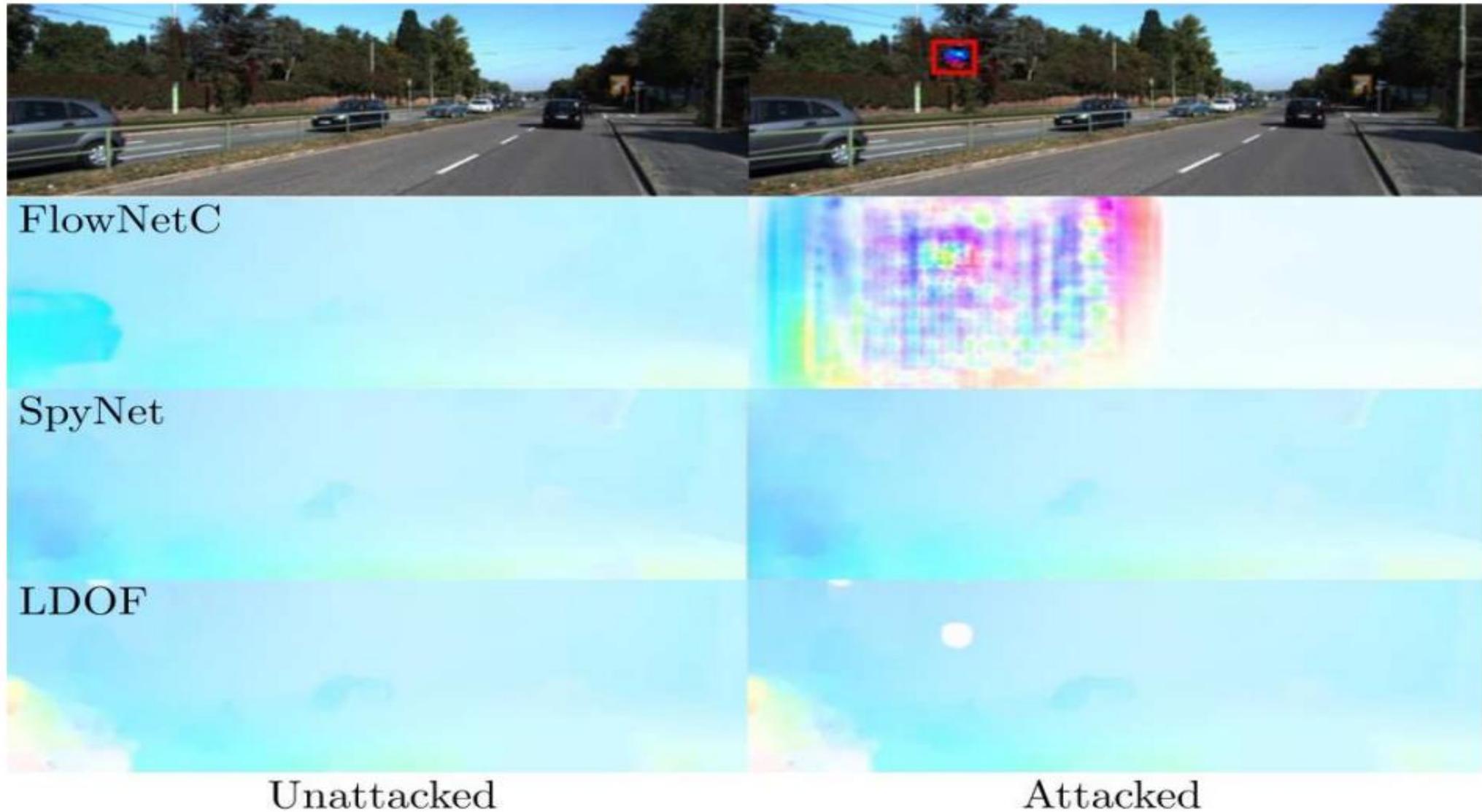


Distilling!

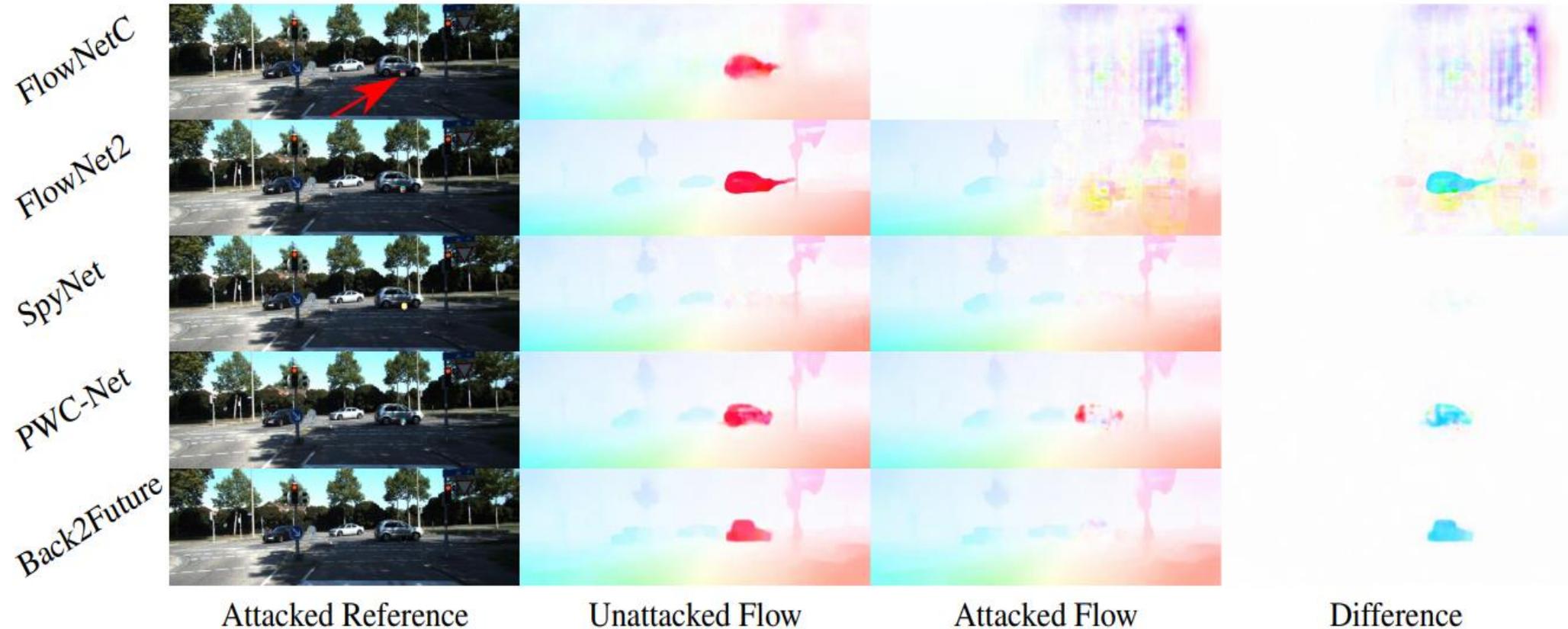


Universal Perturbation Attack Against Image Retrieval
Jie Li, Rongrong Ji, Hong Liu, Xiaopeng Hong, Yue Gao, Qi Tian

Атаки на оптический поток (в автономном вождении)



Атаки на оптический поток (в автономном вождении)

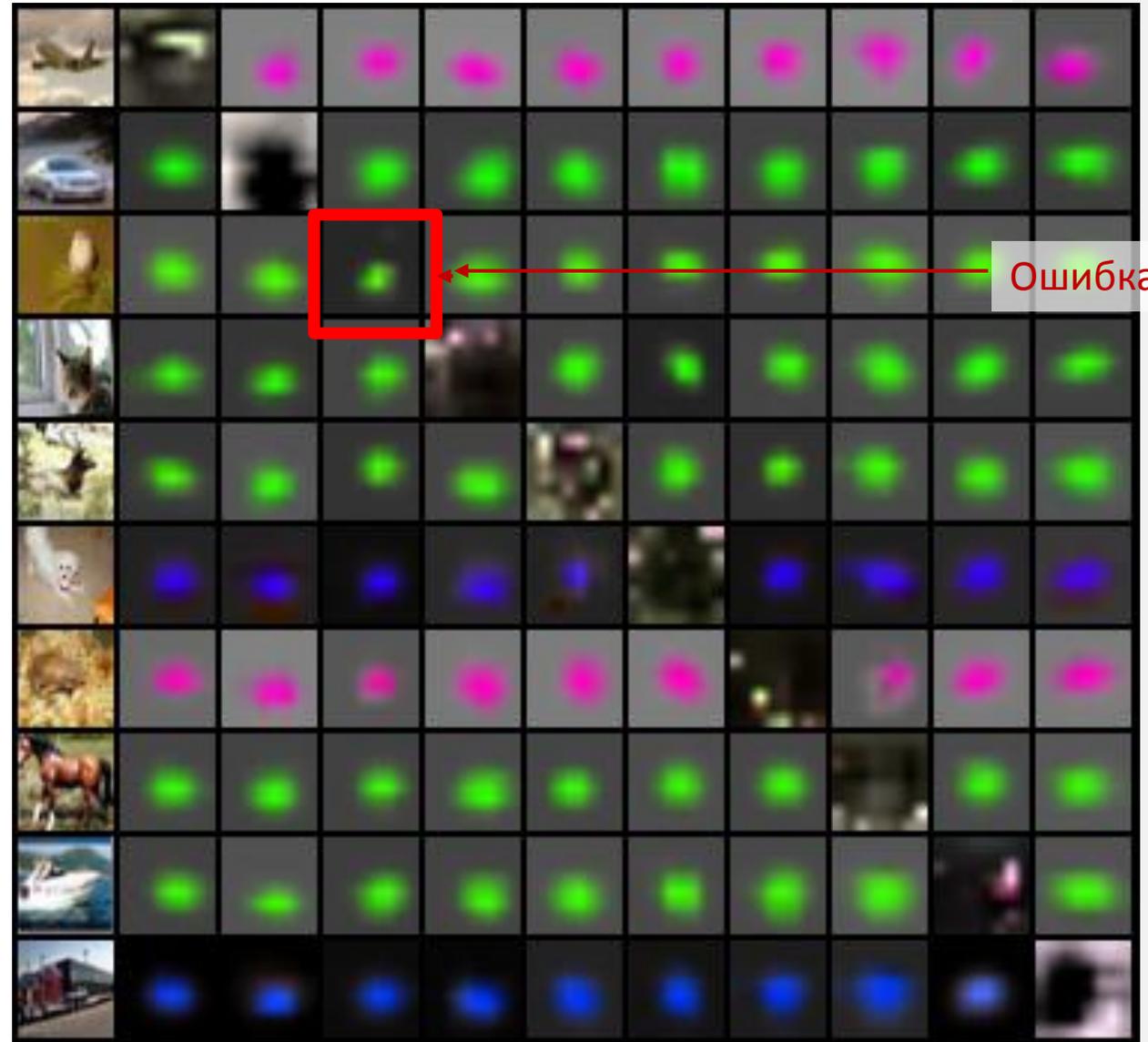


Результатом такой атаки в реальном мире может стать авария с человеческими жертвами...

Исследование уязвимостей и поиск защиты ГНС (ГосНИИАС, 2018-19)

Предложена методика извлечения адаптивных корреляционных эталонов из распознающих ГНС для оценки причин возникновения и разработки методов борьбы с пиксельными и другими «атаками» на распознающие ГНС

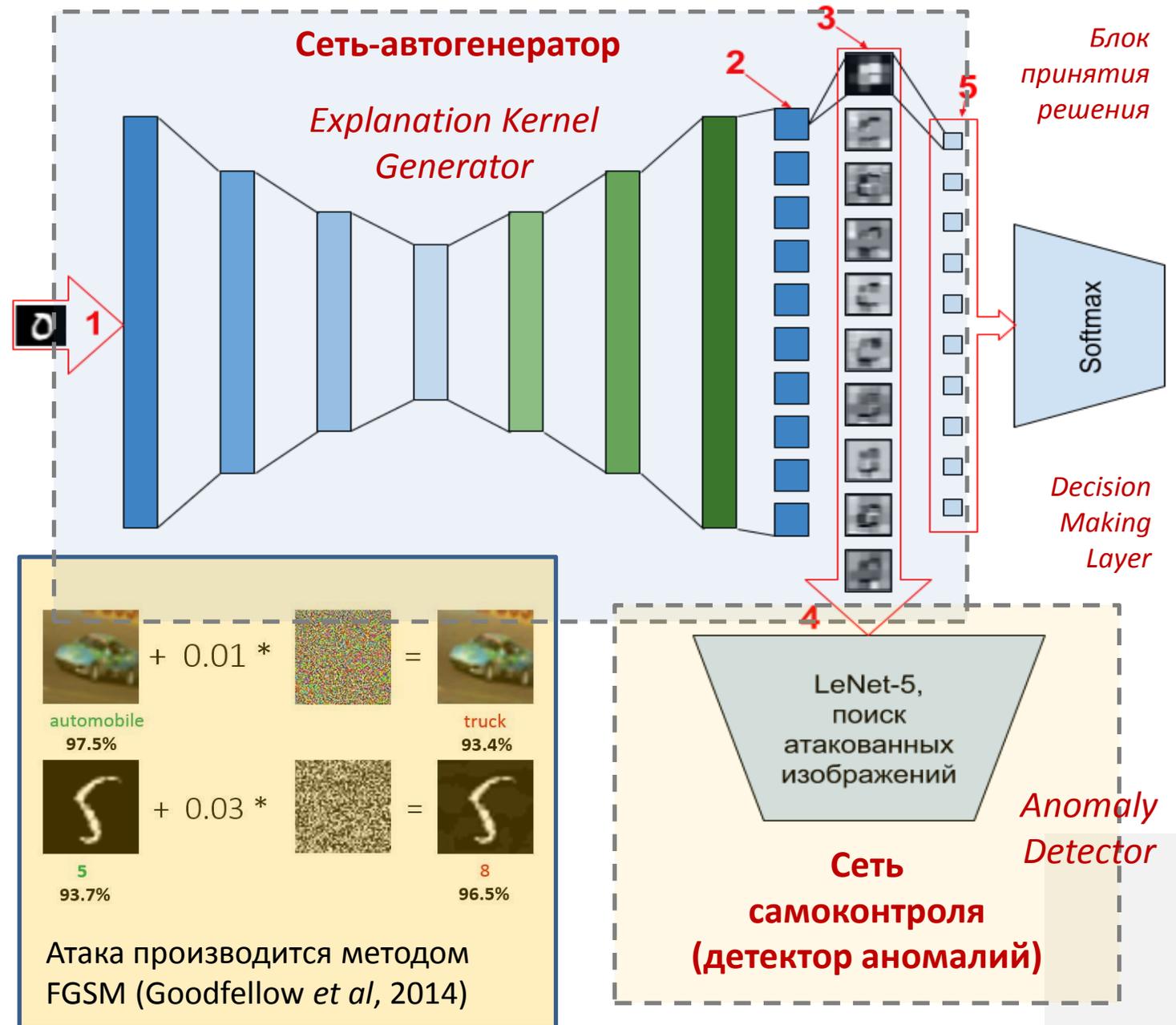
Для распознавания атак была обучена отдельная сеть AttackFinder архитектуры LeNet-5 (5 слоев), которая получает на вход ядро, восстановленное первой сетью, и по виду ядра определяет, была ли исходная картинка атакована, т.е. соответствует ли полученное ядро "типичным представителям" этого класса. Результаты экспериментов показали, что **предложенный подход может не только детектировать атаки на ГНС, но и обнаруживать собственные ошибки распознавания анализируемых ГНС.**



Исследование уязвимостей и поиск защиты ГНС (ГосНИИАС, 2018-19)

	MNIST	CIFAR-10
исходная точность	97,5%	82,3%
% найденных атак	89,9%	94,6%
% ложных срабатываний	8,8%	34,6%
точность на прошедших	99,7%	97,0%

1. Самоконтроль атак возможен, можно находить до 99% атак;
2. Точность классификатора с самоконтролем драматически возрастает;
3. Сеть с самоконтролем может обнаруживать собственные ошибки распознавания как «real-world» атаки.



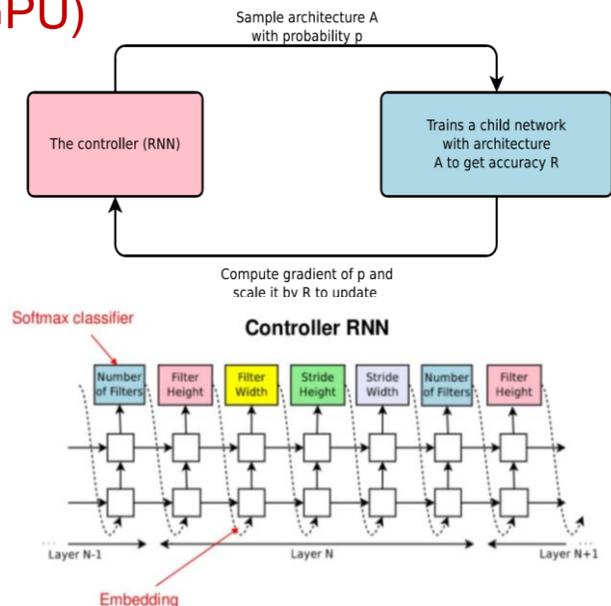
**Автоматическое формирование
и обучение нейросетей (AutoML,
Neural Architecture Search, NAS)**

AutoML: автоматическое обучение глубоких сетей

Первое поколение AutoML (2016-2018)

- Лучшие результаты в задаче классификации на CIFAR-10/ImageNet (выше придуманных человеком архитектур)
- **Вычислительные требования:**
Сотни серверов с GPU/TPU (~500-800 GPU)

2017:
Автоматически сформированные глубокие сети впервые превзошли показатели глубоких сетей, сформированных вручную



*Neural Architecture Search with Reinforcement Learning

Barret Zoph*, Quoc V. Le Google Brain ICCV 2017

*Learning Transferable Architectures for Scalable Image Recognition

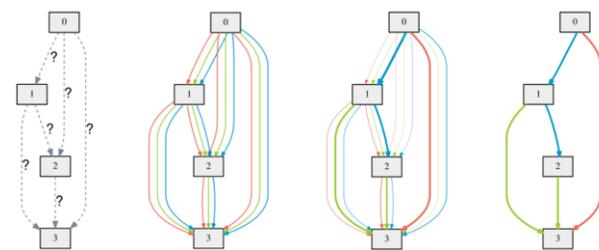
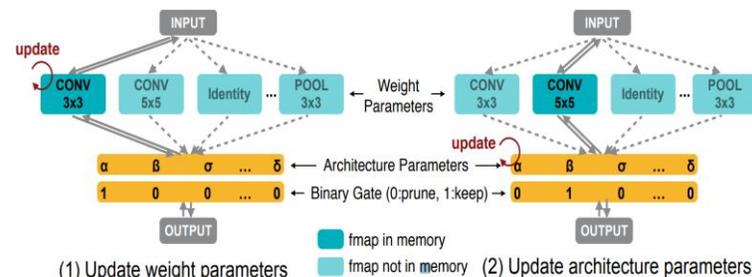
Barret Zoph, Vijay Vasudevan, Jonathon Shlens, Quoc V. Le CVPR 2018

*Regularized Evolution for Image Classifier Architecture Search

Esteban Real, Alok Aggarwal, Yanping Huang, Quoc V Le 2018

Второе поколение AutoML (2018-2019)

- Лучшие результаты для задач обнаружения и распознавания
- **Учет специфики задачи и архитектуры конечного вычислителя**
- **Вычислительные требования:**
от 200 GPU/часов – сравнимо с обычным обучением



*PROXYLESSNAS: DIRECT NEURAL ARCHITECTURE

SEARCH ON TARGET TASK AND HARDWARE

Han Cai, Ligeng Zhu, Song Han arxiv 2019

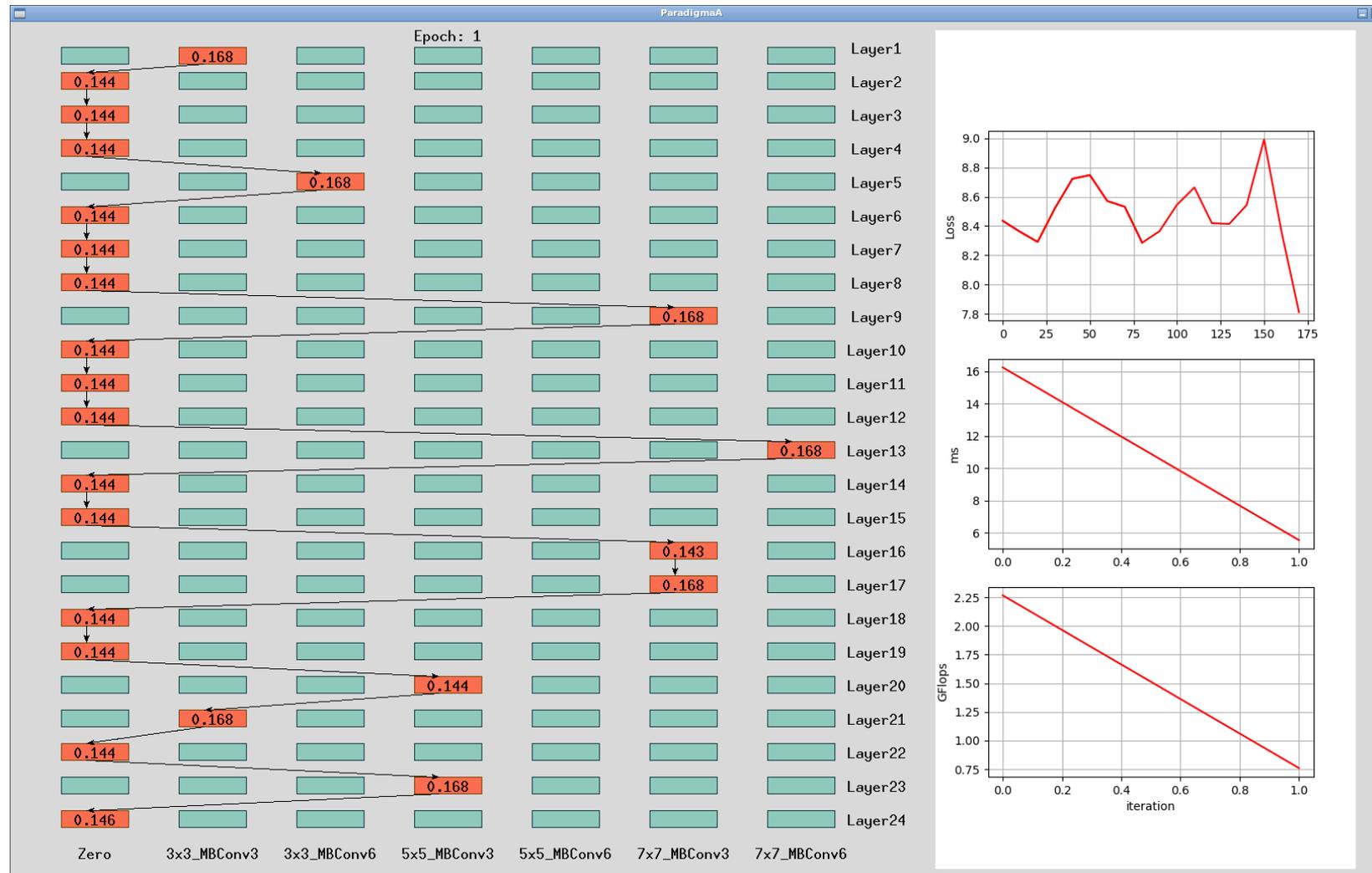
*DARTS: DIFFERENTIABLE ARCHITECTURE SEARCH

Hanxiao Liu, Karen Simonyan, Yiming Yang ICLR 2019

Процесс оптимизации: подбор фильтров, блоков, слоев и параметров в заданном словаре

2025+:
Перспектива полной автоматизации процессов обучения

Алгоритм ProxyLessSSD (ГосНИИАС, 2020)



Функция потерь

Время работы сети

Вычислительная сложность сети

Процесс поиска архитектуры

Результаты тестов на БД VisDrone:

Название ГКНС	mAP, mean average precision	mAR, mean average recall	FLOPs (640x480)	число парам-ов (640x480)	Время, мс на GPU 1080Ti (640x480) batch=1
SSD ProxylessNAS	34.2(+8%)	0.264(+6%)	1.9G(+1%)	1.48M(-26%)	9.24(-22%)
SSD MobileNet v2 (manual) 2.0	31.5	0.249	1.88G	2.02M	11.33
SSD MobileNet v1 (manual) 1.0	19.9	0.158	4,56G	5,16M	11.5
SSD MobileNet v3 small x1.5 (AutoML)	23.8	0.189	1.57G	1.69M	15
SSD ShuffleNet v2 (manual)	22.2	0.19	0.88G	1.53M	13.36

Базовые алгоритмы VisDrone 2019

Название ГКНС	mAP	mAR
CornerNet	34.12	24,37
Light-RCNN	32.78	23,09
FPN	32.20	20,72
Cascade R-CNN	31.91	21,37
DetNet59	29.23	20,87

Самый быстрый алгоритм
(visDrone 2019):

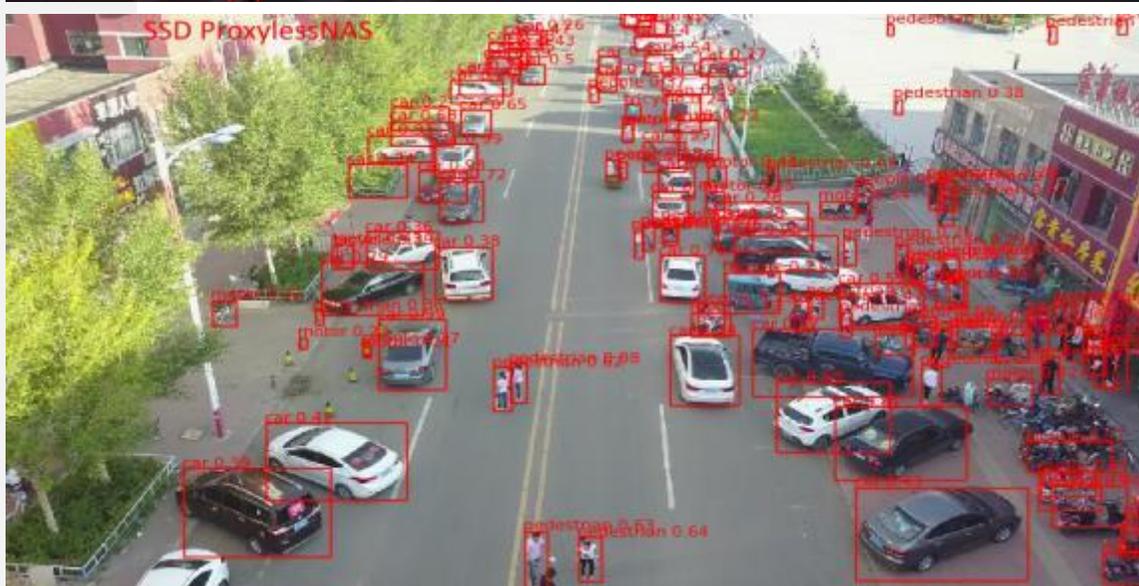
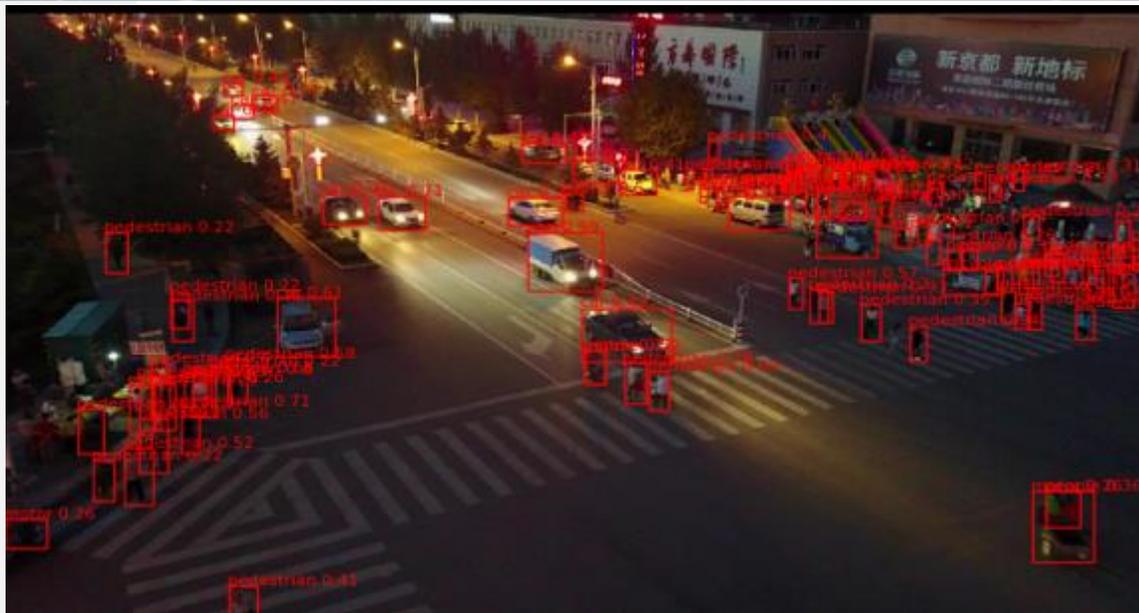
3 кадра в сек.

Самый медленный:

10 сек. на кадр.

Полное время обучения:
432 GPU/часа

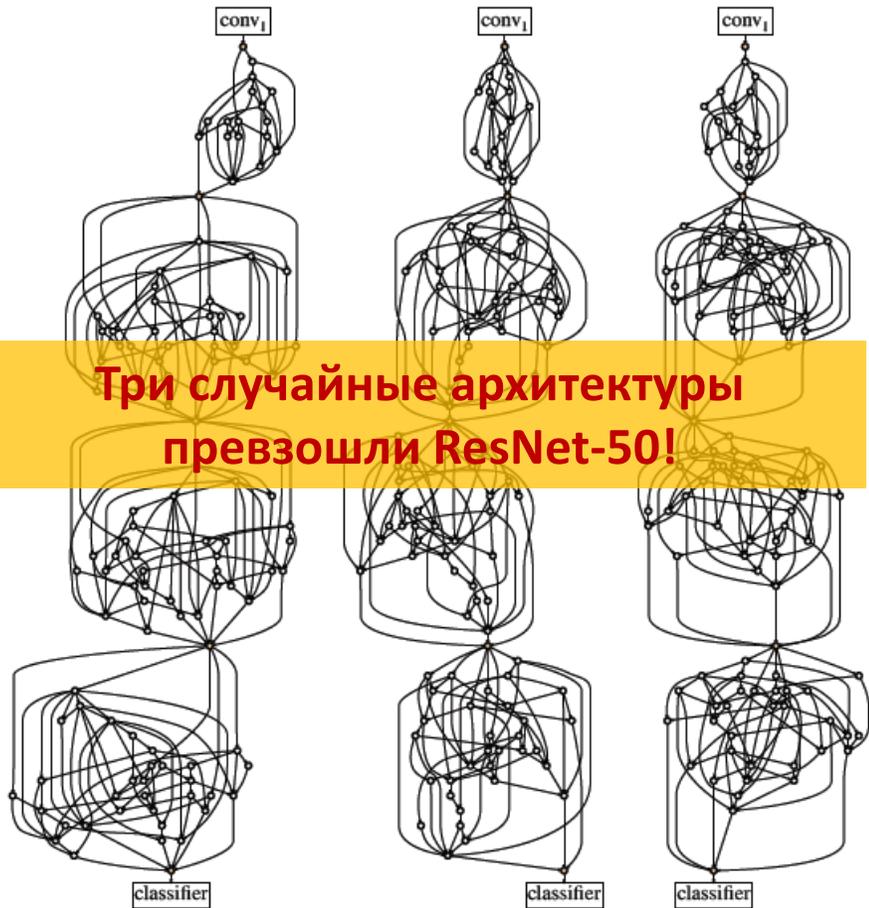
Примеры работы (ГосНИИАС, 2020)



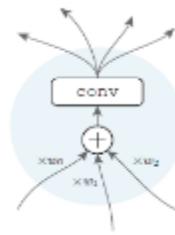
Exploring Randomly Wired Neural Networks for Image Recognition

Saining Xie, Alexander Kirillov, Ross Girshick, Kaiming He

Facebook AI Research (FAIR)

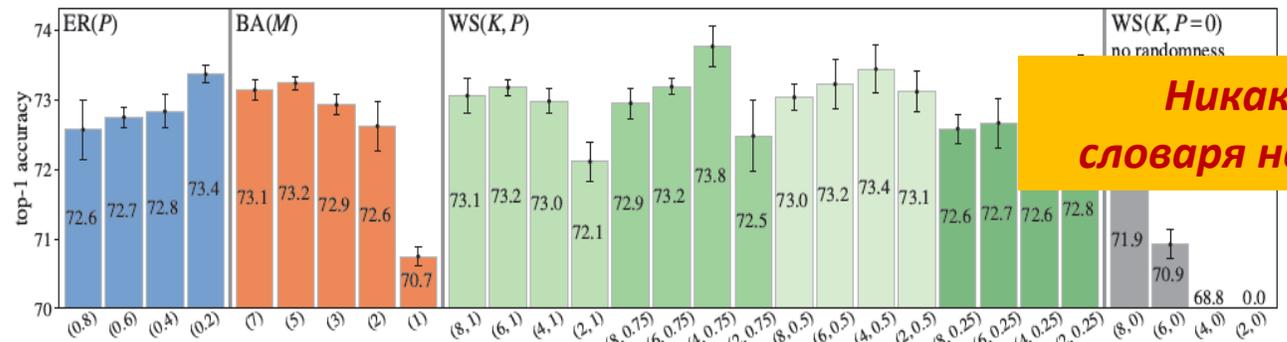


Три случайные архитектуры
превзошли ResNet-50!



- 1) порождаем случайный граф
- 2) определяем вход и выход
- 3) строим и обучаем GCN

Эта работа взрывает все наше понимание NAS!



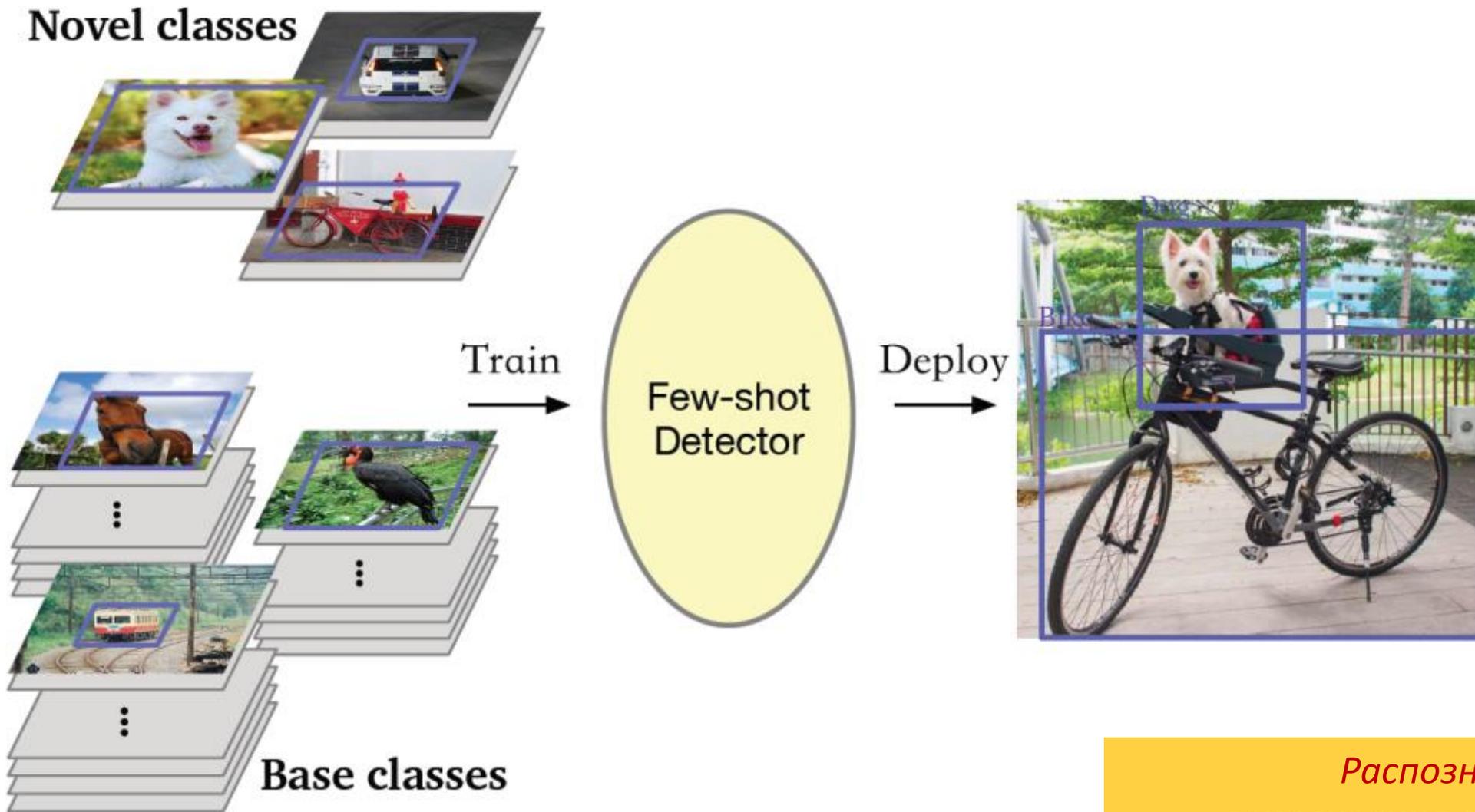
Никакого
словаря нет!

Figure 1. Randomly wired neural networks generated by the classical Watts-Strogatz (WS) [51] model: these three instances of random networks achieve (left-to-right) 79.1%, 79.1%, 79.0% classification accuracy on ImageNet under a similar computational budget to ResNet-50, which has 77.1% accuracy.

Figure 3. Comparison on random graph generators: ER, BA, and WS in the small computation regime. Each bar represents the results of a generator under a parameter setting for P , M , or (K, P) (tagged in x-axis). The results are ImageNet top-1 accuracy, shown as mean and standard deviation (std) over 5 random network instances sampled by a generator. At the rightmost, $WS(K, P=0)$ has no randomness.

**Обучение на малом числе
примеров (Few-Shot Learning /
Detection / Segmentation)**

Few-Shot Learning / **Detection** / Segmentation



Few-Shot Object Detection via Feature Reweighting

Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, Trevor Darrell

Распознавание и локализация объектов по одному или малому числу эталонов

Few-Shot Learning / Detection / Segmentation

Результаты на выборке PASCAL VOC 2007 *(ГосНИИАС, 2019)*

Результаты работы на классах, участвующих в обучении (AP)

bike	bird	bottle	bus	car	chair	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP(%)
0.72	0.53	0.51	0.73	0.8	0.44	0.59	0.69	0.76	0.71	0.72	0.35	0.68	0.62	0.69	0.59	0.63

Результаты работы на классах, не участвующих в обучении (AP)

aero	boat	cat	cow	mAP(%)
0.25	0.11	0.66	0.73	0.44

128x128 эталонное изображение

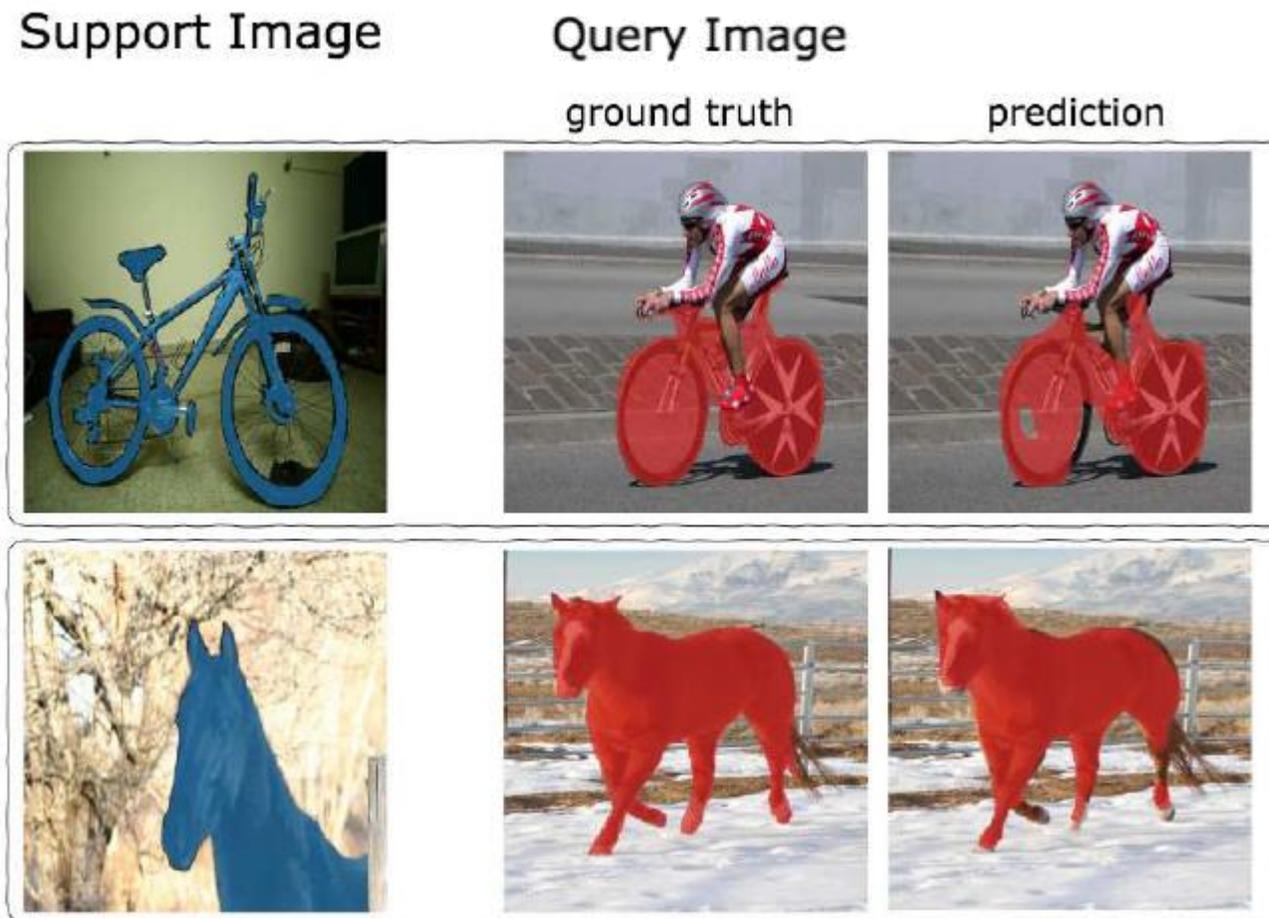
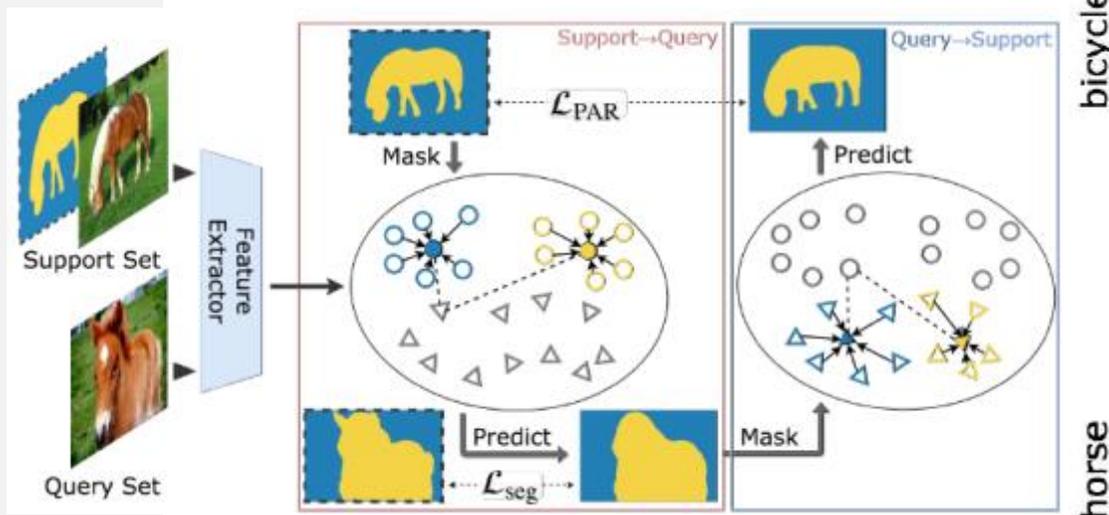


Тестовое изображение



Семантическое обнаружение объектов на изображениях с использованием ГКНС, Горбачевич, Визильтер, Моисеенко, 2019

Few-Shot Learning / Detection / Segmentation

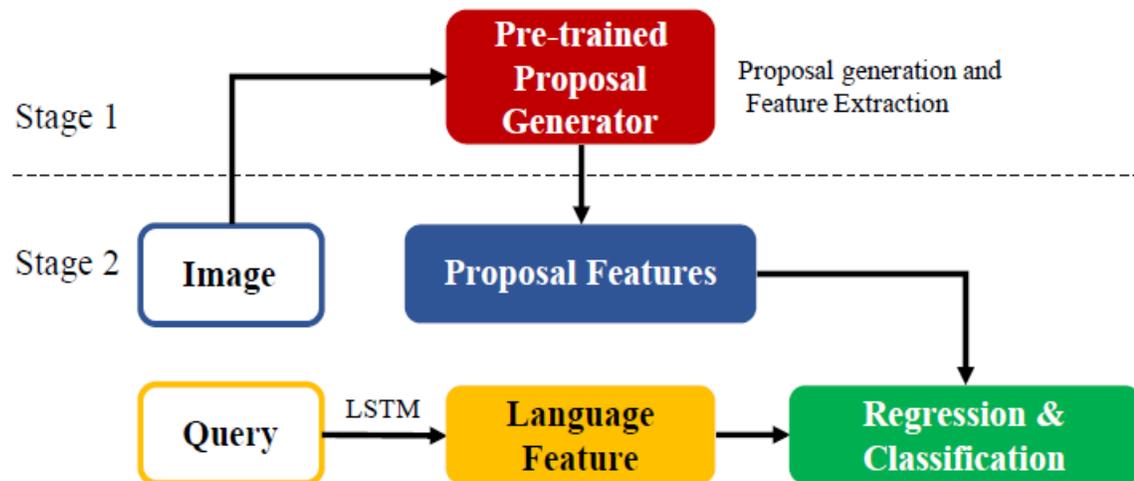


PANet: **Few-Shot** Image **Semantic Segmentation** With Prototype Alignment
Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, Jiashi Feng

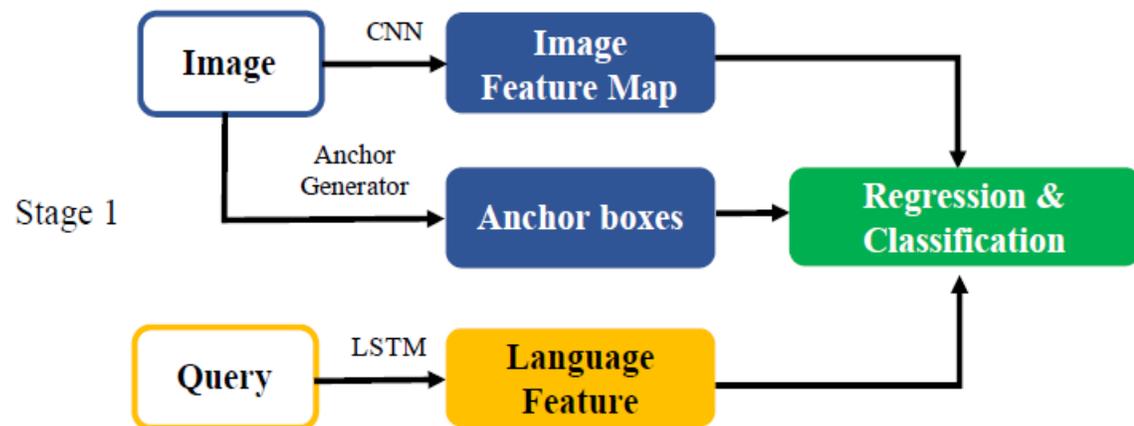
*Распознавание, локализация
и сегментация объектов по
одному или малому числу
эталонов*

**Обучение без примеров
(Zero-Shot Learning, Grounding)**

Zero-Shot Grounding



(a) Vanilla 2-stage phrase grounding system



(b) Our 1-stage phrase grounding system



a group of older men



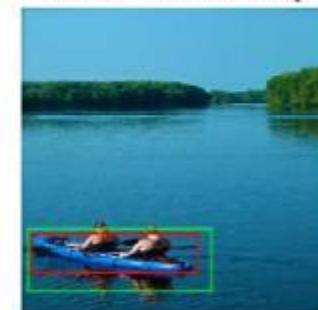
a red beanie cap



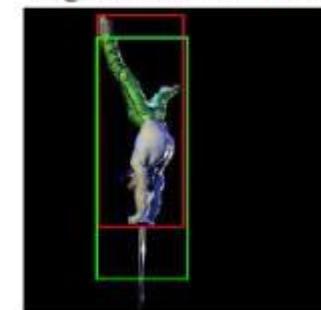
rightmost animal



a cigar



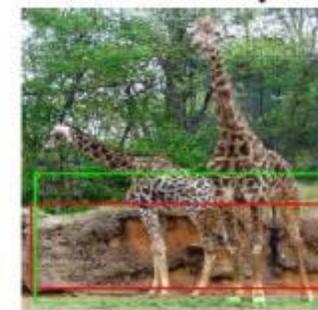
a two-seat kayak



a handstand



a rocky cliff (hill)



large boulders (rock)



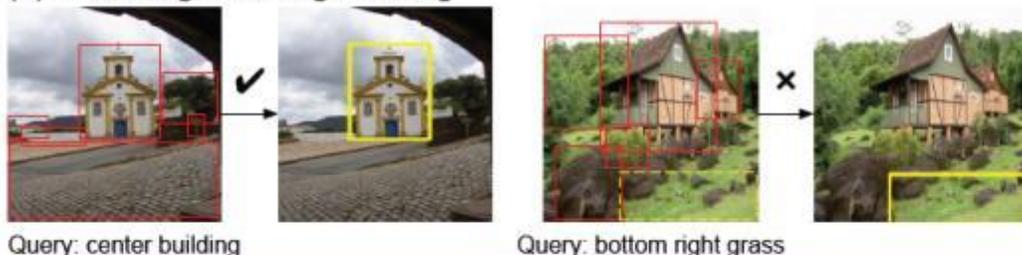
stairway (wall)

Zero-Shot Grounding of Objects From Natural Language Queries
Arka Sadhu, Kan Chen, Ram Nevatia

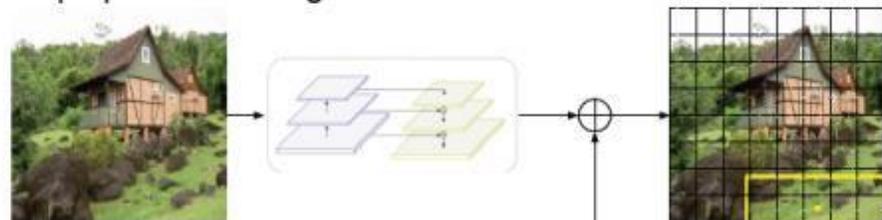
Поиск и локализация объектов без эталона (по описанию)

Zero-Shot Grounding

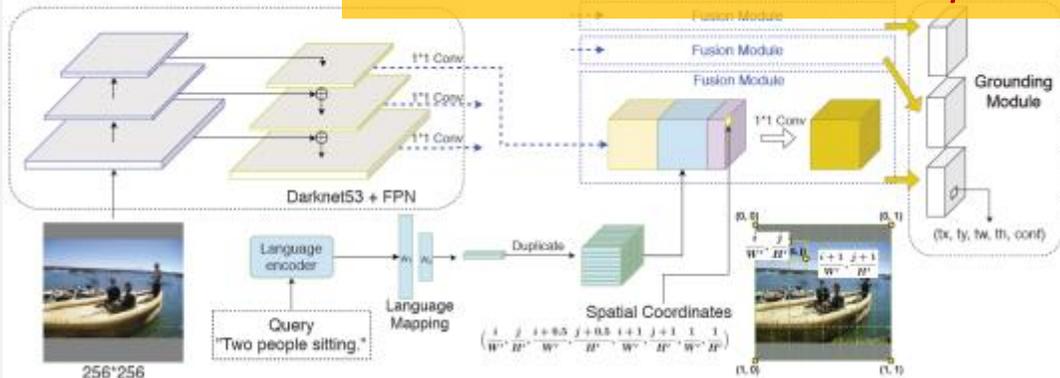
(a). Two-stage visual grounding



(b). The proposed one-stage method



Одноэтапный детектор + Attention с языковым запросом



(a). Query: bike of blue pant lady



(b). Query: the bowl of bean on the bottom



(c). Query: person on the right



(g). Query: man in blue



(h). Query: kid left



(i). Query: window above colonial

A Fast and Accurate One-Stage Approach to **Visual Grounding**
 Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, Jiebo Luo

Поиск и локализация объектов без эталона (по сложному описанию)

Few-Shot via Zero-Shot!

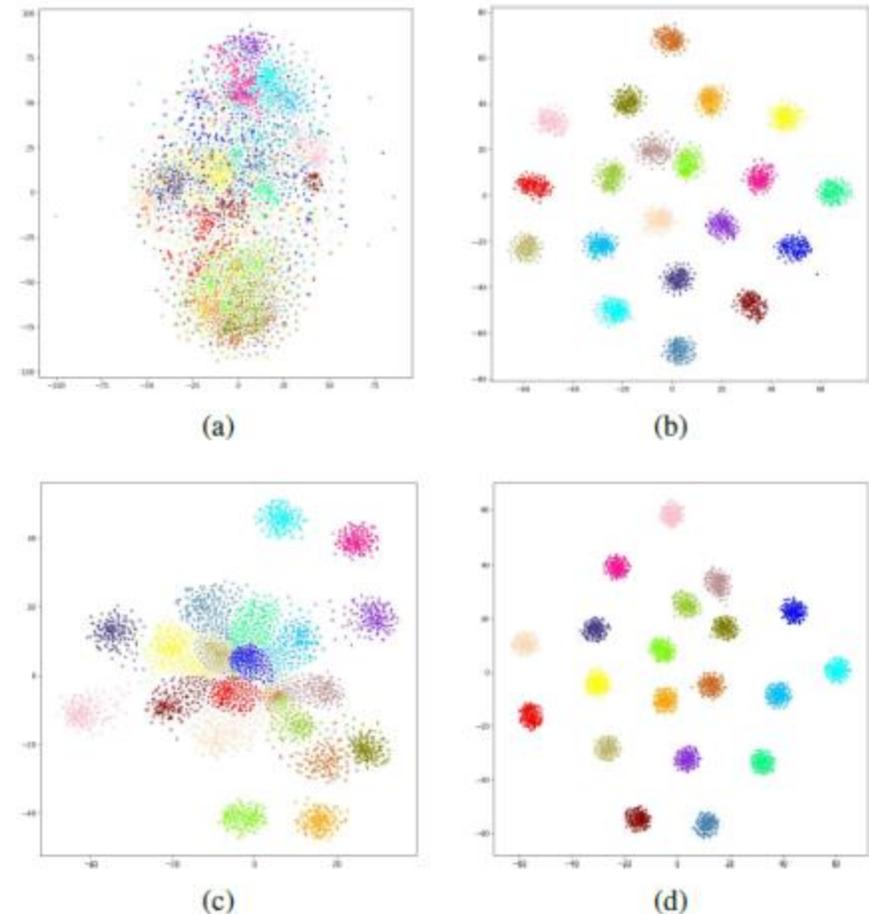
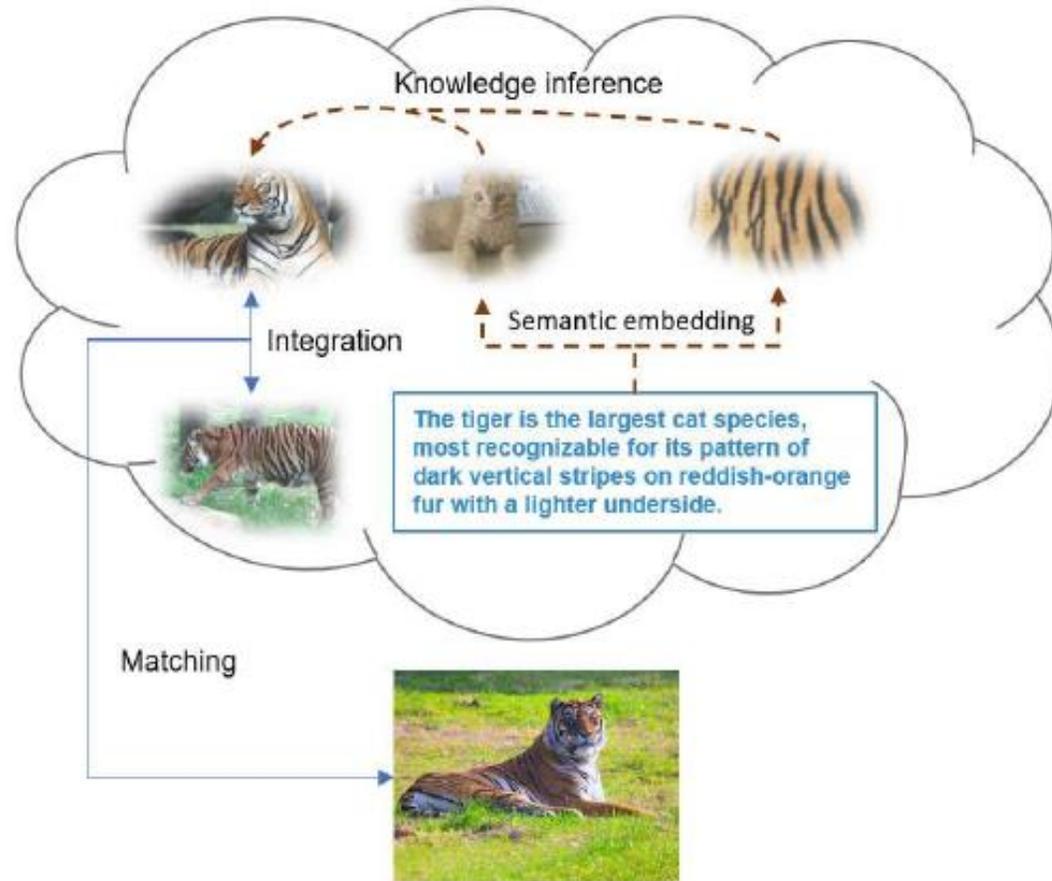


Figure 3. T-SNE visualization results for all novel categories in the Mini-imagenet set on the 1-shot and 5-shot tasks. Each scatter plot contains 20 colored classifier parameter clusters and each color represents a novel category. (a): 1-shot vision-based classifier. (b): 1-shot vision-knowledge classifier. (c): 5-shot vision-based classifier. (d): 5-shot vision-knowledge classifier.

Поиск и локализация объектов по эталону класса путем генерализации через словесное описание!

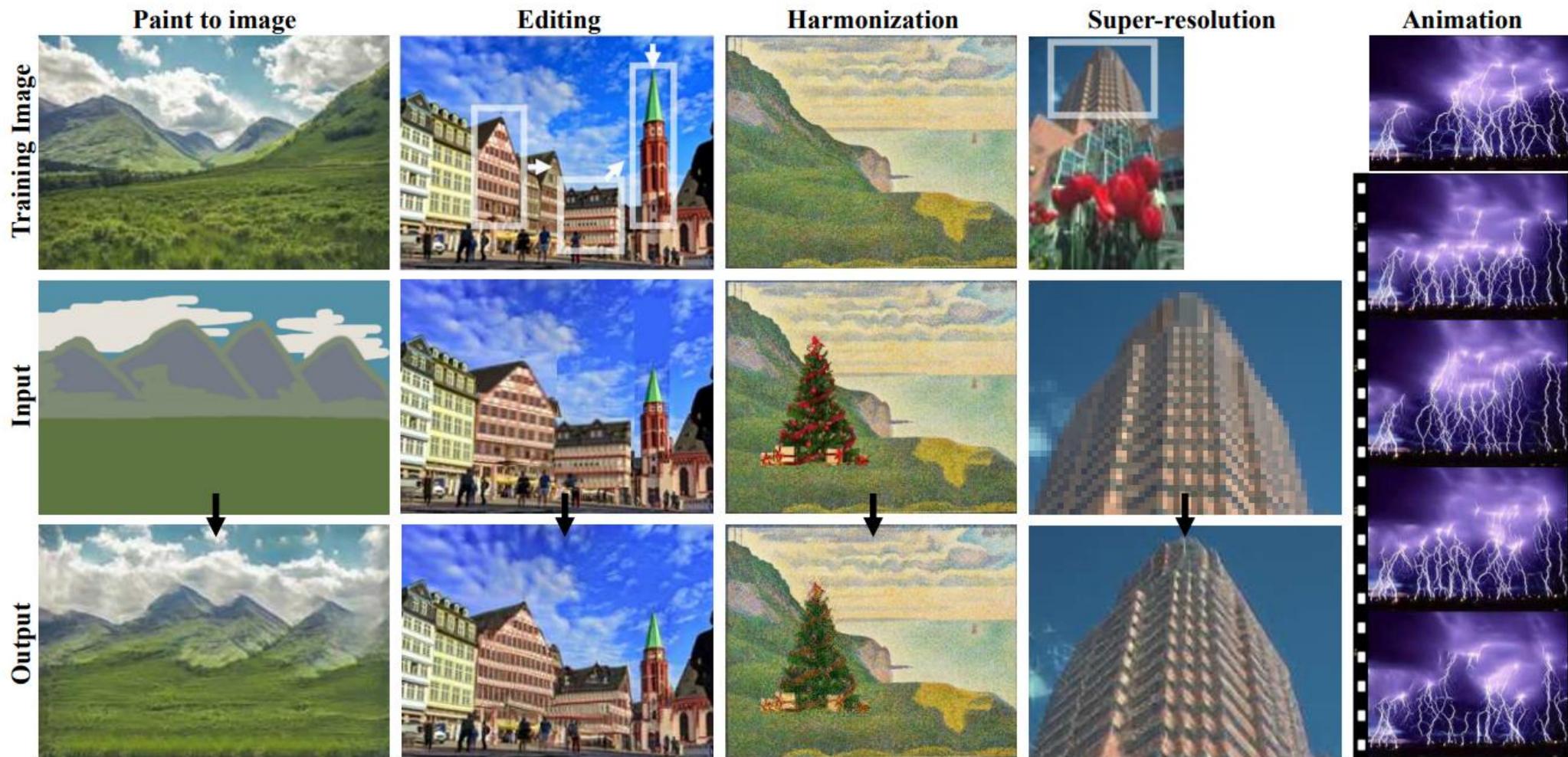
Few-Shot Image Recognition With Knowledge Transfer

Zhimao Peng, Zechao Li, Junge Zhang, Yan Li, Guo-Jun Qi, Jinhui Tang

Классы намного более компактные

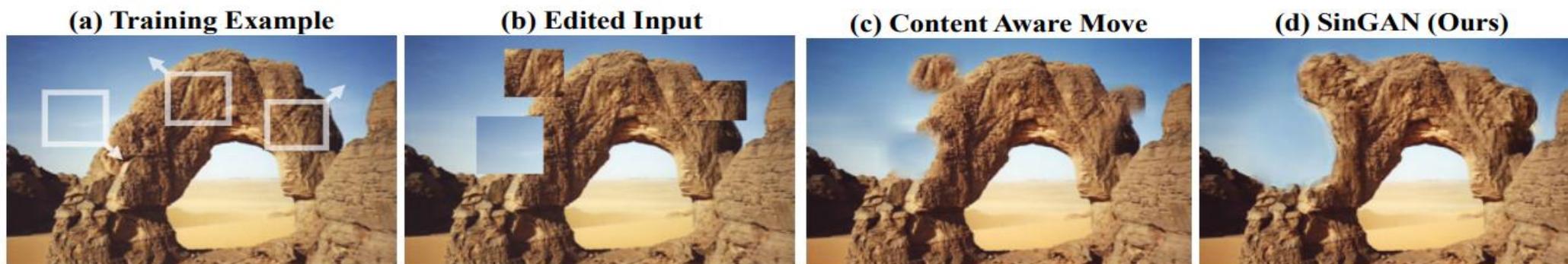
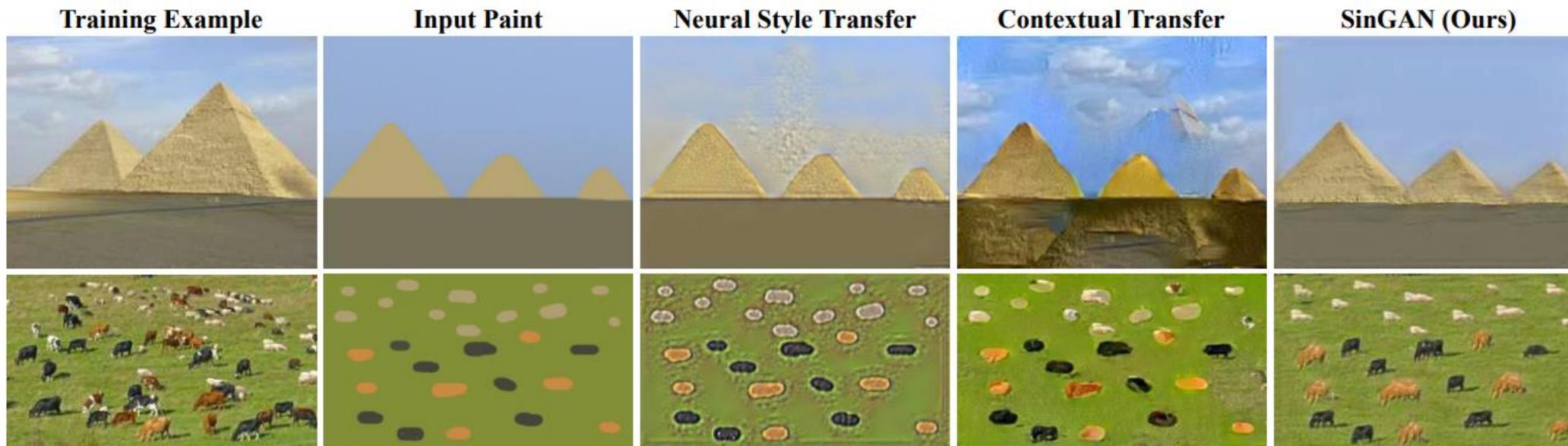
**Генерация реалистичных данных
(Domain Adaptation, Generative
Adversarial Networks (GAN),
Realistic Data Synthesis...)**

SinGAN: Learning a Generative Model from a Single Natural Image



Для обучения реалистичного генератора изображений больше не нужны большие базы примеров!

SinGAN: Learning a Generative Model from a Single Natural Image



Практически неограниченные возможности реалистичного манипулирования данными...

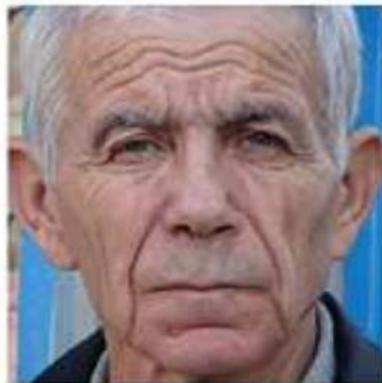
Детальная генерация 3D лица по одному изображению



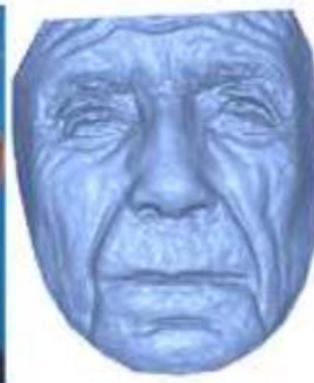
Real Faces



Swapped Faces



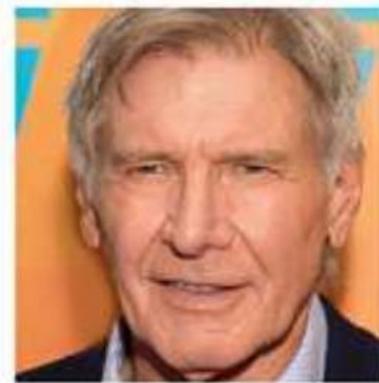
(a) input image



(b) output 3d face



(c) textured 3d face



(d) input image



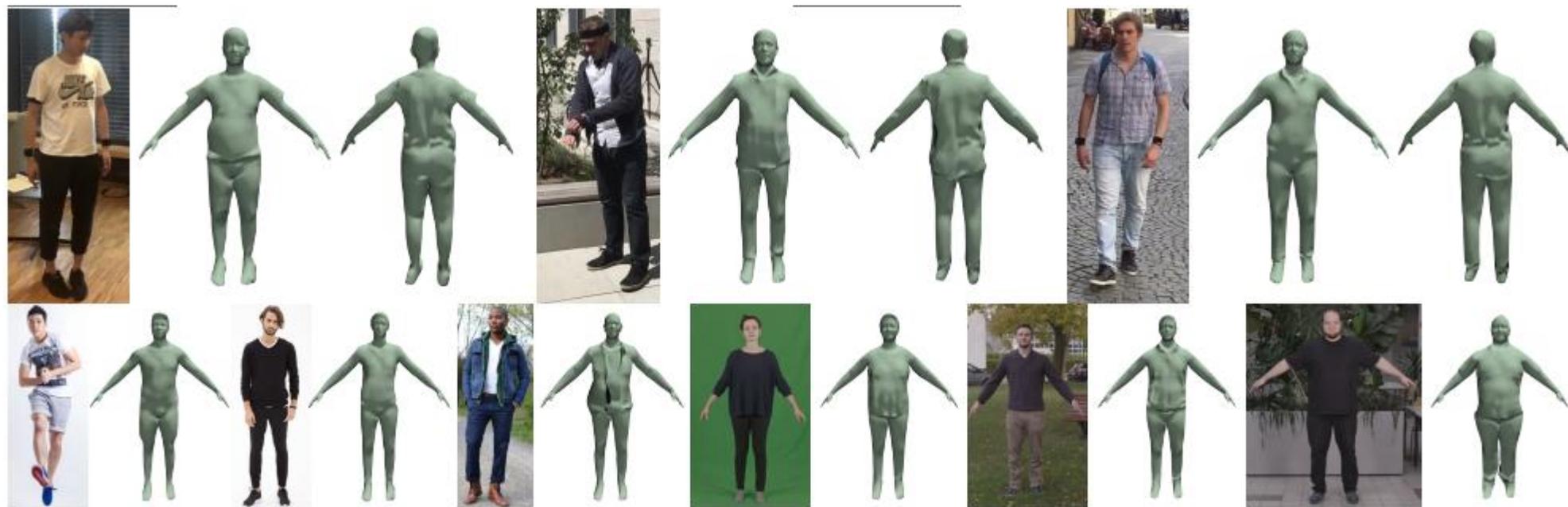
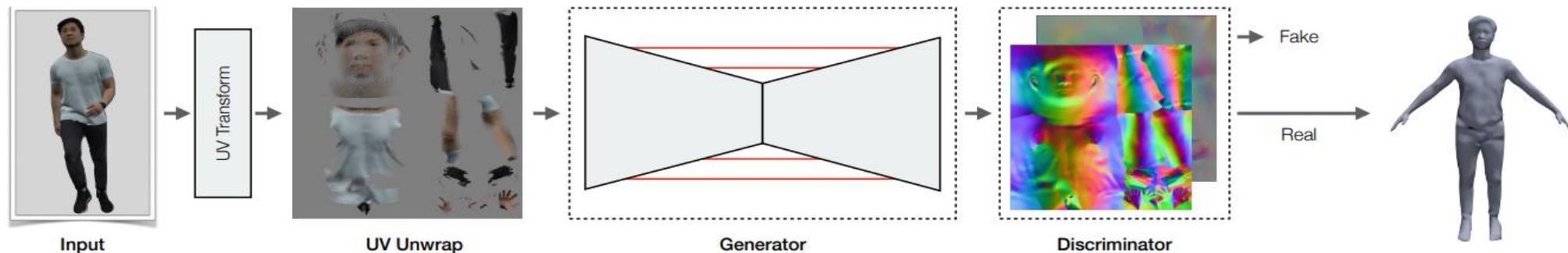
(e) output 3d face



(f) textured 3d face

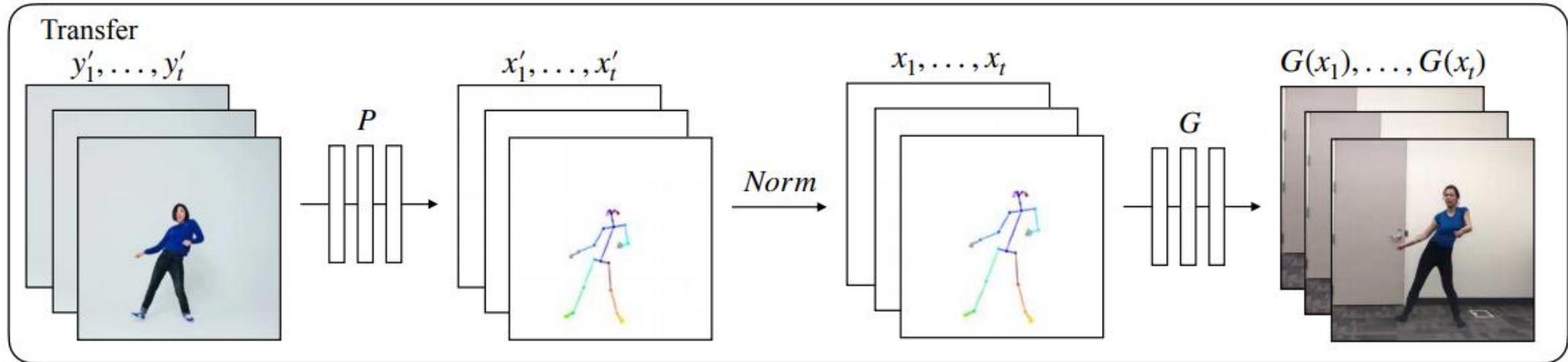
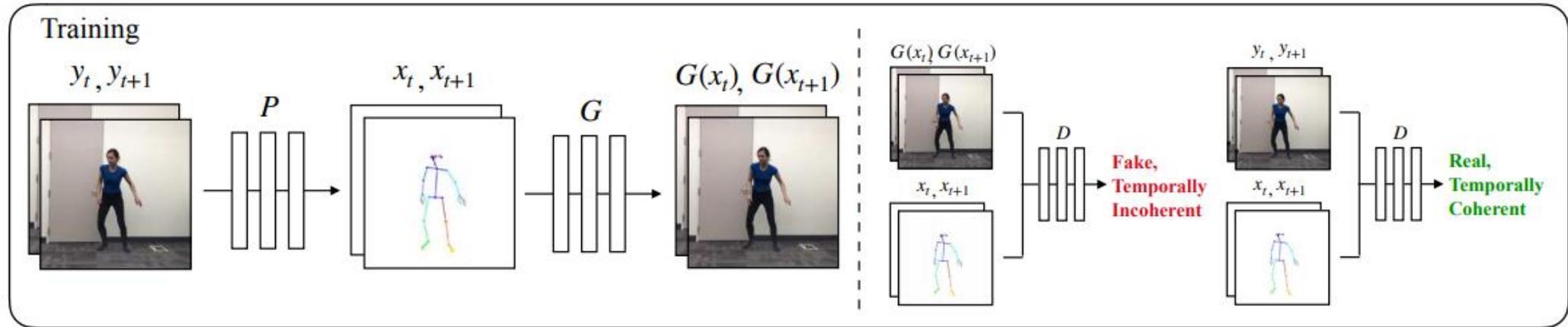
Практически неограниченные возможности реалистичного манипулирования данными...

Детальная генерация 3D модели тела и позы



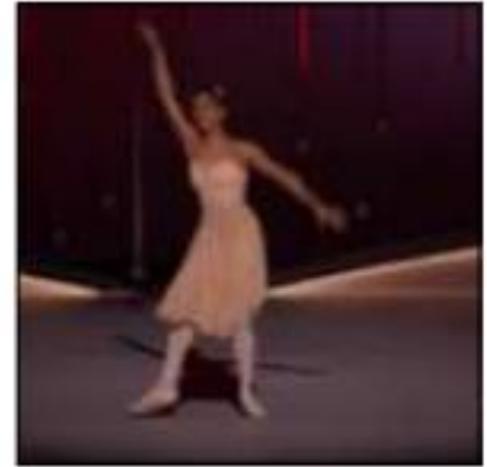
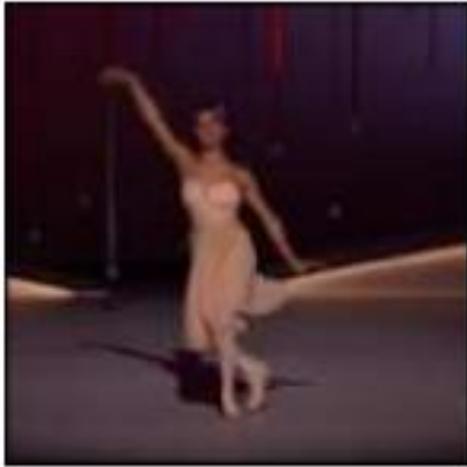
Практически неограниченные возможности реалистичного манипулирования данными...

Перенос движения: Everybody Dance Now



Практически неограниченные возможности реалистичного манипулирования данными...

Everybody Dance Now



Практически неограниченные возможности реалистичного манипулирования данными...

Everybody Dance Now



Практически неограниченные возможности реалистичного манипулирования данными...

Everybody Dance Now

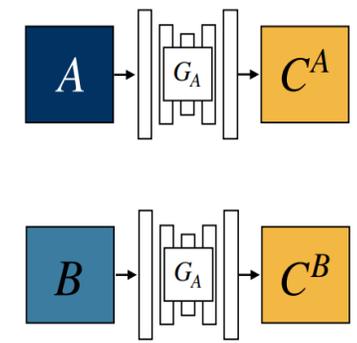
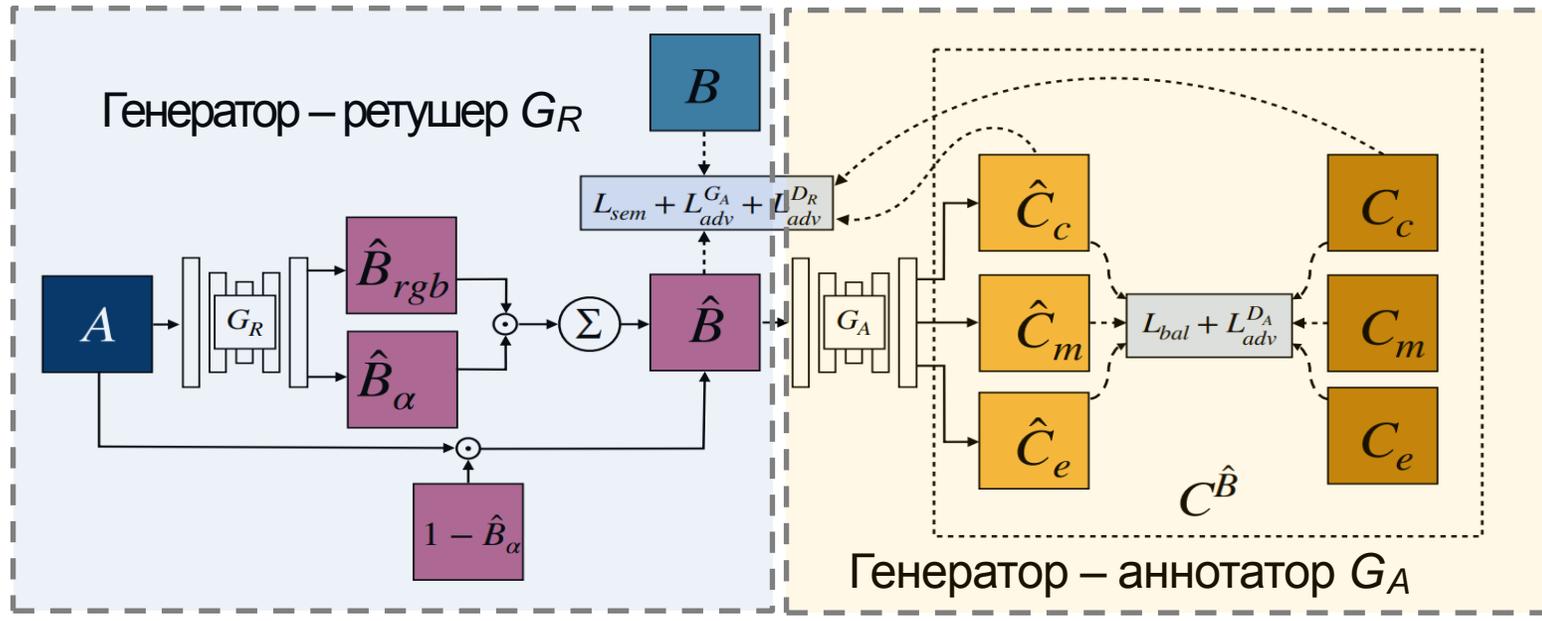


Практически неограниченные возможности реалистичного манипулирования данными...

Использование смешанных генеративно–состязательных нейронных сетей для обнаружения фотомонтажей

Mixed Adversarial Generators

Генератор – ретушер G_R переводит изображения из домена монтажей A в домен реальных фото B .



Генератор-аннотатор G_A обучается прогнозированию маски монтажа, маски краёв монтажа и карты сегментации классов объектов

Retoucher Generator G_R



Annotator Generator G_A

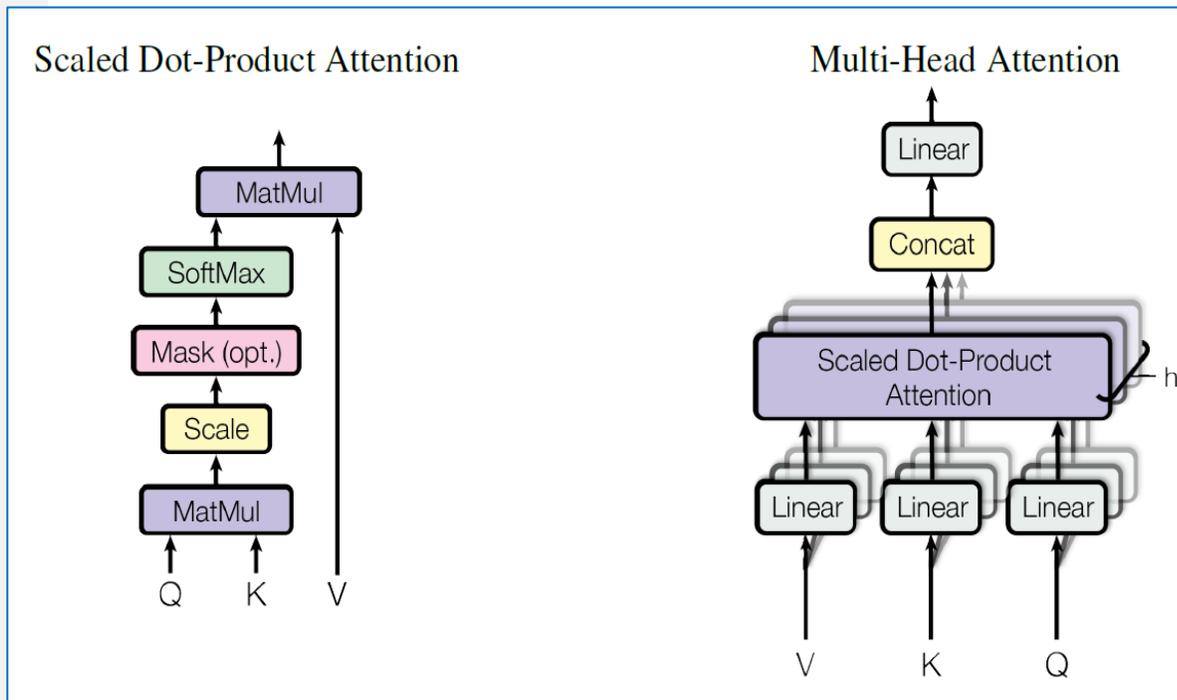
The Point Where Reality Meets Fantasy: Mixed Adversarial Generators for Image Splice Detection (NeurIPS 2019) Vladimir V. Knyaz, Vladimir A. Knyaz, Fabio Remondino



Attention, NLP: BERT, GPT-3,...

(Работа с текстами и неожиданный путь к универсальному «слабому» ИИ)

Graph Attention Networks



An attention function can be described as **mapping a Query and a set of Key-Value pairs to an output**, where the query, keys, values, and output are all vectors.

In the area of Natural Language Processing where **Transformer** style models have become state of the art on many tasks. **Motivation: Learning long-range dependencies is a key challenge for conv layers, which is important for symbolic sequences.**

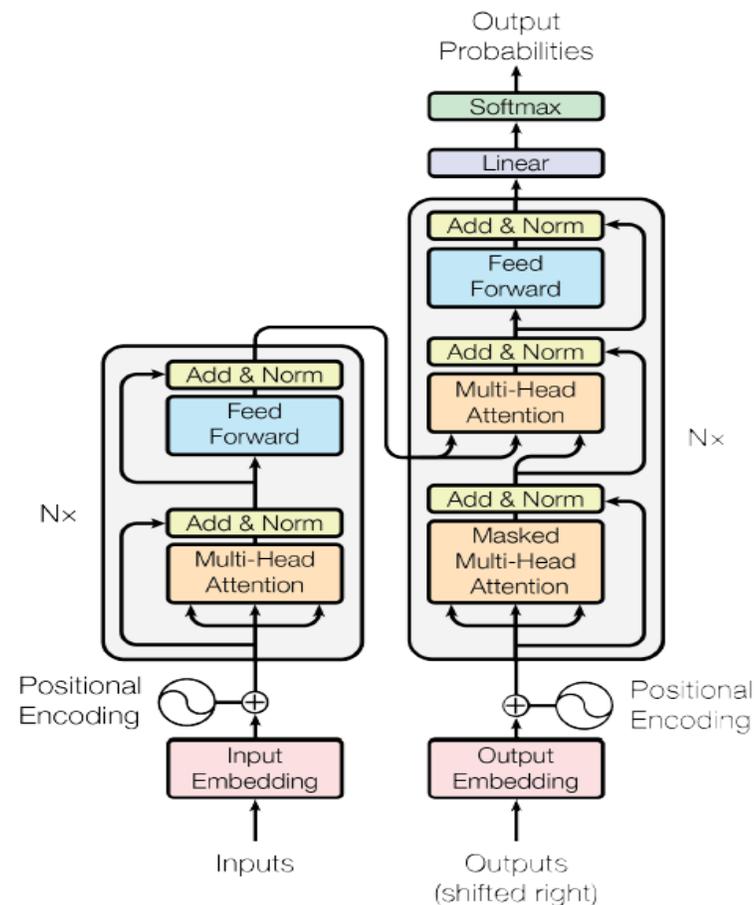


Figure 1: The Transformer - model architecture.

Управление вниманием (адаптивные парные веса) в GCN оказались ключом к решению самых сложных задач!

Graph Attention Networks

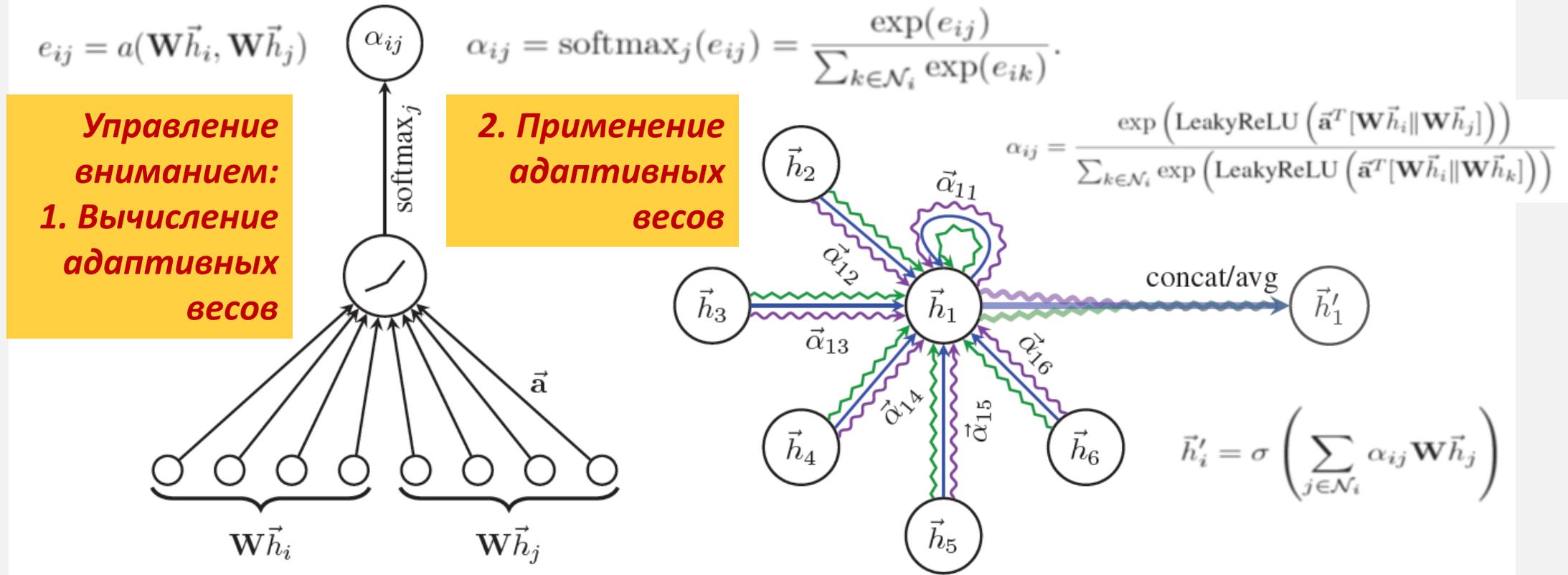


Figure 1: **Left:** The attention mechanism $a(\mathbf{W}\vec{h}_i, \mathbf{W}\vec{h}_j)$ employed by our model, parametrized by a weight vector $\vec{a} \in \mathbb{R}^{2F'}$, applying a LeakyReLU activation. **Right:** An illustration of multi-head attention (with $K = 3$ heads) by node 1 on its neighborhood. Different arrow styles and colors denote independent attention computations. The aggregated features from each head are concatenated or averaged to obtain \vec{h}'_1 .

OpenAI GPT-3: “may be the biggest thing since bitcoin”

Параметры сети:

Размер сети 175 млрд.
параметров(350 Гб fp16)
Контекстное окно: 2048 токена
96 трансформеров
96 self attention head
Размер батча: 3.2 млн

База данных:

The CommonCrawl data - 45 Тб
Прочие БД – 40%

Результаты сравнимые со state-of-the-art в 42 задачах NLP



Суперкомпьютер:

10000 GPU Tesla V100
400 Гбит/сек
285 тыс. CPU ядер
Pytorch

Требует 355 лет и 5 млн. \$
при аренде GPU Cloud

Архитектуры:

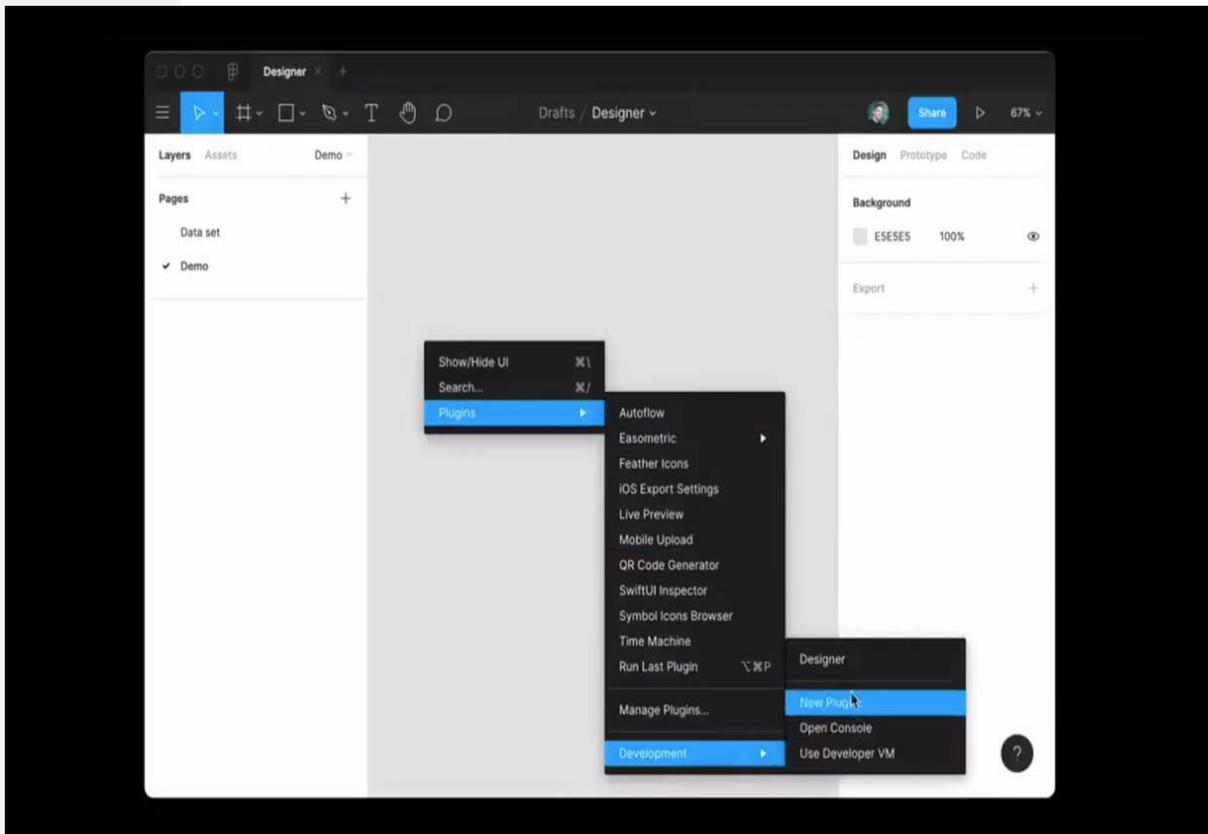
Модификация GPT-2
BERT
Transformer

- GPT-3 может отвечать на вопросы на естественном языке
- GPT-3 умеет переводить с любого языка, отвечать на вопросы, писать статьи и т.д. вплоть до написания исходного кода программ

**ХИТ ИИ
2020!**

OpenAI GPT-3: “may be the biggest thing since bitcoin”

Примеры приложений:



WHAT

Describe a layout.

Just describe any layout you want, and it'll try to render below!

A div that contains 3 buttons each with a random color.

Generate

Trends for you

Trending in Netherlands
Italië

6,684 Tweets

Trending in Netherlands
#coronamaatreg

Trending in Netherlands
Genee

1,933 Tweets

Politics · Trending
Rutte

42.6K Tweets

Trending in Netherlands
Efteling

Show more

Генерация визуальных форм(кода) по текстовому описанию

<https://machinelearningtokyo.com/2020/07/26/10-cool-gpt-3-demos/>



OpenAI GPT-3: “may be the biggest thing since bitcoin”

Примеры приложений:

Received Email

Matt,

Thanks for chatting last week. Hearing your vision for Otherside got both Jim and I really excited. We really like where you're going with this. After talking with my partners yesterday, we're looking at making an investment of \$100K into Otherside on a SAFE. Would this be sufficient to join your round? If so, we'll send over our proposed terms.

On another note, as we discussed, let me know about your estimated market size.

Please let me know. Looking forward to an amazing journey together!

Thanks

Response Points

- * thanks
- * no
- * our minimum is \$150K investment
- * would \$150K be possible
- * \$90B market



Thanks for the note, Peter!

I'm glad you were able to catch up with your partners, and I'm excited to hear about your interest in the company. I'm glad you like where we're going.

I'm not sure if you were aware, but our minimum raise is \$150K. I'm not sure if \$100K would be sufficient to join the round, but I'd be happy to discuss it.

I'm not sure if you were aware, but the market is approximately \$90B.

Thanks,

Matt

Endless Rows: Silicon Valley Edition

My work flow is done before lunch each day and I began experimenting with programming languages outside of work. I found myself building another parser one day, and I realized that I frequently came across specs that required reading like a stack trace with rows of data

Автоматическое написание писем по нескольким тезисам

- Общение с ИИ на естественном языке
- Универсальная модель
- Команды для РТК/ИИ
- Построение моделей по уставу/правилам



Открытые проблемы: угрозы, вызовы, надежды (2020)

Проблемы:

- Как справиться с атаками?
- Как эффективно переносить обучение в реальном мире?
- Как разоблачать фейковые данные?
- Реальных данных для практических приложений катастрофически не хватает!
- Перспективные методы обучения требуют слишком больших ресурсов!
- Мостик между зрением и языком/пониманием давно перекинут, но массовый переход пока не случился

Надежды:

- Новые датчики технического зрения
- Прогресс в AutoML/NAS
- Прогресс в Few-Shot/Zero-Shot
- GCN, Attention => GPT-3
- Прогресс в объяснении нейросетевых рассуждений
- Массовый перевод нейросетевых рассуждений с уровня отдельных объектов на уровень семантических конструкций (онтологий)
- Прогресс в методах RL, глубокой оптимизации и глубокого управления
- Соединение универсальности GPT-3 с анализом данных и глубокой оптимизацией

ВЫВОДЫ ПО СОСТОЯНИЮ ТЕХНОЛОГИЙ ИИ

1. Первая волна современной революции в ИИ породила технологии «глубокого распознавания», обеспечивающие решение задач **компьютерного зрения, анализа сигналов и больших данных.**

2. В настоящее время технологическая революция в ИИ переживает **вторую волну**, которая ведет нас прямо к созданию **функционального ИИ.**

3. Нет никакого специфического прорыва в методах «искусственного интеллекта», но **наблюдается технологический прорыв, связанный с ГНС.**

4. Методы на основе глубоких нейронных сетей (ГНС) **не только про СТЗ.** «Глубокие» технологии развиваются, фокус их применения смещается **от обработки и анализа данных к семантике, управлению и оптимизации**

5. **Революция в искусственном интеллекте продолжается! (см. GPT-3)**

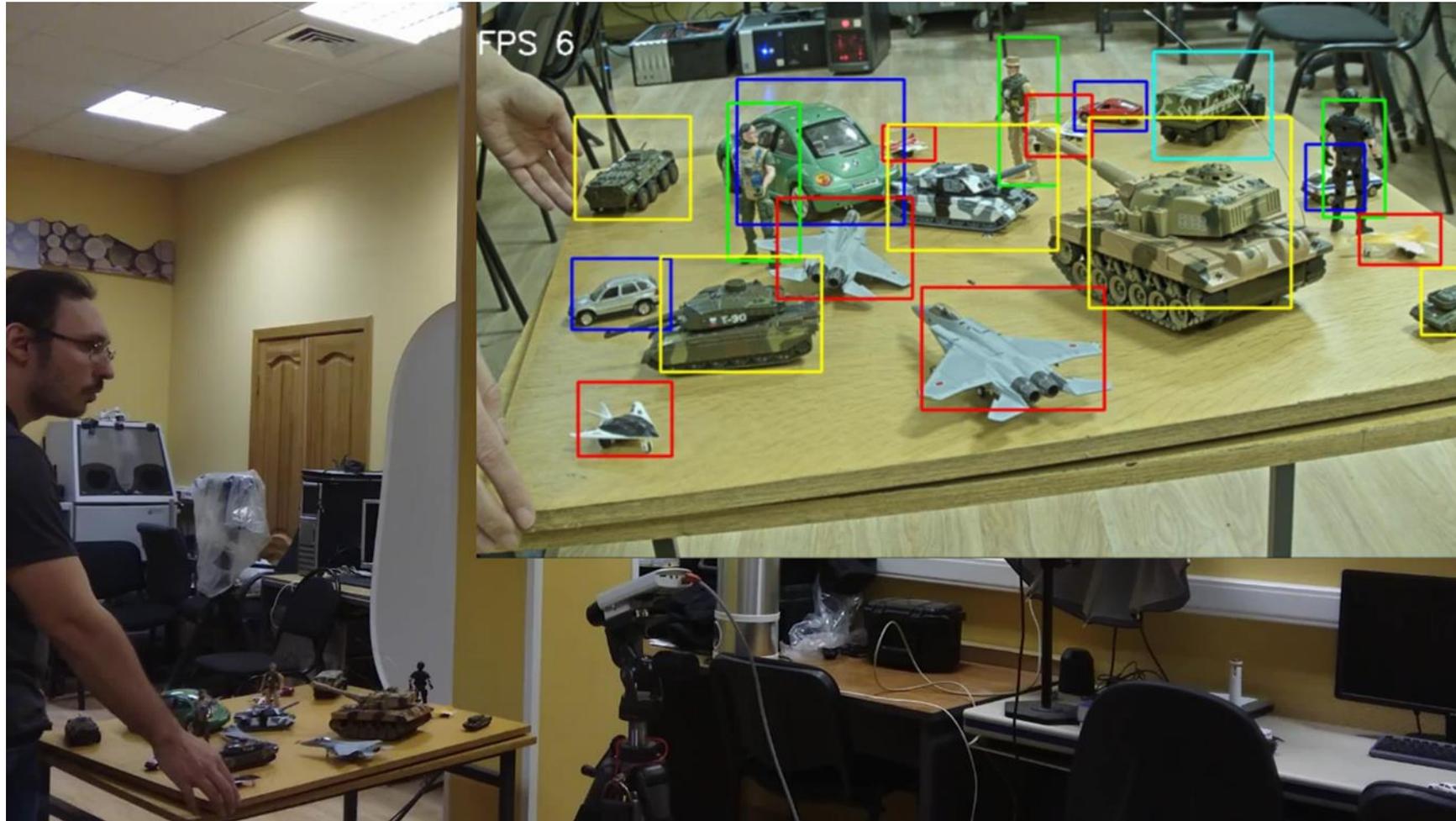
Скорость и направления дальнейшего продвижения будут зависеть от того, **как и когда удастся справиться с возникшими проблемами**, а также от того, **когда и какие сбудутся надежды.**

Унифицированная программная платформа
для разработки конечно ориентированных программных
комплексов на основе нейросетевых подходов

(ПЛАТФОРМА-ГНС, Экосистема Plat)

+ Начальник лаборатории 3050 ФГУП “ГосНИИАС”
Горбацевич Владимир Сергеевич, gvs@gosniias.ru

2016: Обнаружение и распознавание объектов в реальном времени (ГосНИИАС)



Цветом указаны классы объектов: желтый – бронетехника, красный – самолеты, голубой – грузовики, синий – автомашины, зеленый - люди

Где отечественные процессоры?
Где отечественное ПО для обучения?

2018-19: ГНС на отечественных процессорах



Стенд ГосНИИАС
на выставке
«АРМИЯ-2018»

2018: Прототип системы автоматического обнаружения и распознавания целей на основе глубоких сверточных нейронных сетей. Система на базе платы MC121.01 производства НТЦ «Модуль» с процессором NM6407 (5 кадров/сек).

2020: создано бортовое решение с ГНС на базе NM6408, работающего в 32 раза быстрее (скорость обработки - 60 кадров/сек и выше).

Перспектива-2020+: ожидается выпуск еще трех нейропроцессоров отечественных производителей



2019+: отечественные процессоры уже есть

ПО обучения глубоких нейросетей: проблемы и решение

Задача: бортовая реализация ГНС и использование глубокого обучения в отечественных АПК для критических приложений.

Проблемы:

- **Существуют только зарубежные средства обучения ГНС** (системы Caffe, Caffe2, Pytorch, TensorFlow, Theano и др.) => **техническая зависимость, невозможность сертификации ПО на основе ГНС на НДС;**
- Отсутствие поддержки и учета особенностей отечественных датчиков и бортовых вычислителей;
- Несовместимость различных средств и систем разработки на базе ГНС...

Решение: создание отечественной унифицированной программной Платформы для разработки конечно ориентированных программных комплексов, которая должна стать основным средством автоматизированного проектирования (формирования и обучения) алгоритмов и бортового нейросетевого ПО всех будущих отечественных АПК.

**Нужен не просто отечественный аналог фреймворка, а единая интегрированная среда разработки ГНС.
Почему?**

Годится ли нам традиционный путь R&D?

На нашем пути есть объективные барьеры, которые нужно учитывать:

- Нас мало!
- Мы разрознены!
- У нас меньше денег и других (в т.ч. вычислительных) ресурсов!...

Пирамида человеческих ресурсов

ICCV-2019: 300+ статей от США и Китая / менее 10 из России

Прогноз: В области разработки и внедрения ИИ в конечные изделия наши разрозненные и малочисленные разработчики всегда будут выступать как "малый бизнес", который из-за эффекта масштаба проигрывает крупным иностранным разработчикам.

Возможное решение: создать специальную интегрированную среду поддержки разработок в области ИИ (ГНС), которая позволила бы ускорить внедрение ГНС в России путем преодоления ряда барьеров и облегчения ряда переходов...

Студенты, сотрудники

Университет или институт

Платформа для ускорения переходов



Платформа для ускорения переходов



Платформа для преодоления барьеров

- снять ограничения для применения ГНС в области ОПК

- обеспечить наискорейшее внедрение самых передовых ГНС в конечные изделия

- понизить порог входа (потому что нас мало) и резко увеличить сообщество разработчиков

- упростить навигацию в пространстве решений и совместное использование знаний и данных



- максимально облегчить совместное использование ресурсов СКТ, проектную и командную работу

- минимизировать проблемы из-за несовместимости сред разработки, повысить унификацию и взаимопонимание между командами разработчиков

- создать механизмы защиты прав собственников IP и управления доступом к конфиденциальным/ закрытым данным

Готовность отечественных технологий

Проект «Платформа-ГНС» (ГосНИИАС, 2018-20)

Унифицированная платформа

Единая интегрированная среда

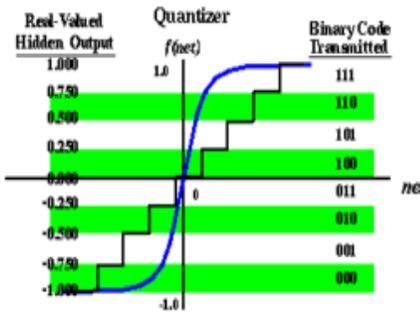
2020+: задача распознавания образов не только алгоритмически решена, но доведена до стадии технологической готовности к ОКР

- **Полностью сертифицированный на НДВ исходный код**
- Возможность использования библиотек NVIDIA и т.д.
- Единая экосистема
- Унифицированный формат хранения БД и моделей
- Импорт/экспорт из основных фреймворков и оппх
- **Наличие типовых решений для основных задач**
- Поддержка отечественных аппаратных платформ
- Поддержка зарубежных аппаратных платформ
- Поддержка отечественных ОС
- Контроль доступа к данным и проектам
- Множество возможностей для командной разработки
- **Низкие требования к квалификации ИТР**

Обучение ГНС



Преобразование ГНС



Аппаратно-ориентированная реализация ГНС



Снижение порога входа => резкое увеличение числа разработчиков => мы надеемся на прорыв в массовом внедрении ГНС в ОПК

2020: создано отечественное ПО

Отечественный фреймворк PlatLib...

Обучение:

- Динамические графы
- Распределенное обучение
- Поддержка GPU AMD/NVIDIA
- Поддержка CPU x86-64/Elbrus
- Поддержка БД Платформы

```
class LeNetPlat(plat.module.Module):  
    def __init__(self):  
        super(LeNetPlat, self).__init__()  
        self.conv1 = conv2D(20, 1, 5)  
        self.relu1 = relu()  
        self.pool1 = poolmax2d()  
        self.conv2 = conv2D(50, 20, 5)  
        self.relu2 = relu()  
        self.pool2 = poolmax2d()  
        self.fc1 = conv2D(500, 50, 4)  
        self.relu3 = relu()  
        self.fc2 = conv2D(10, 500, 1)
```

```
res_plt = LeNetPlat(inputT)  
plt_res = SM(res_plt, labelsT)  
plt_res.backward()  
optimizer_plat.iter()
```

Обучение:

- Аналогично pytorch
- Отладка на python
- Нет multiGPU
- Распределенное обучение через nccl/rccl/mpl

+ интегрированная среда... + целая экосистема Plat...
(полный стек ипортонезависимых технологий)

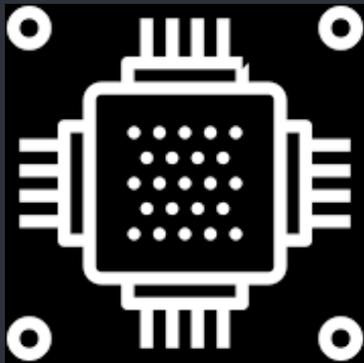
Характеристики Платформы

Клиентская часть



- CentOS/AltLinux
- Аппаратная платформа Intel x86 + Nvidia
- Аппаратная платформа Intel x86 + AMD
- Аппаратная платформа Эльбрус
- Аппаратная платформа Эльбрус + AMD
- Сертификация на НДС

SDK Платформы



- ОС Ubuntu/CentOS/AltLinux / AstraLinux
- Аппаратная платформа Intel x86
- Аппаратная платформа Эльбрус
- Кросс-платформенный код
- Сертификация на НДС

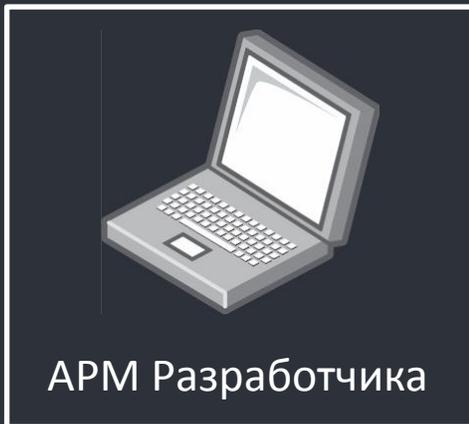
Серверная часть

Вычислительный узел

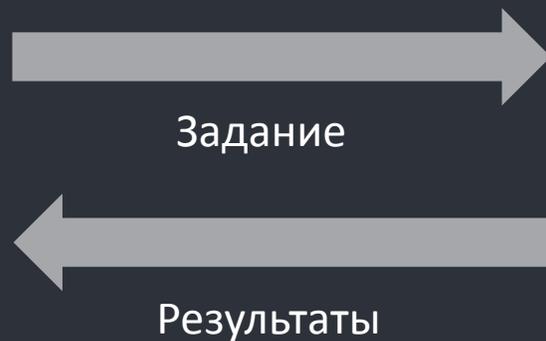
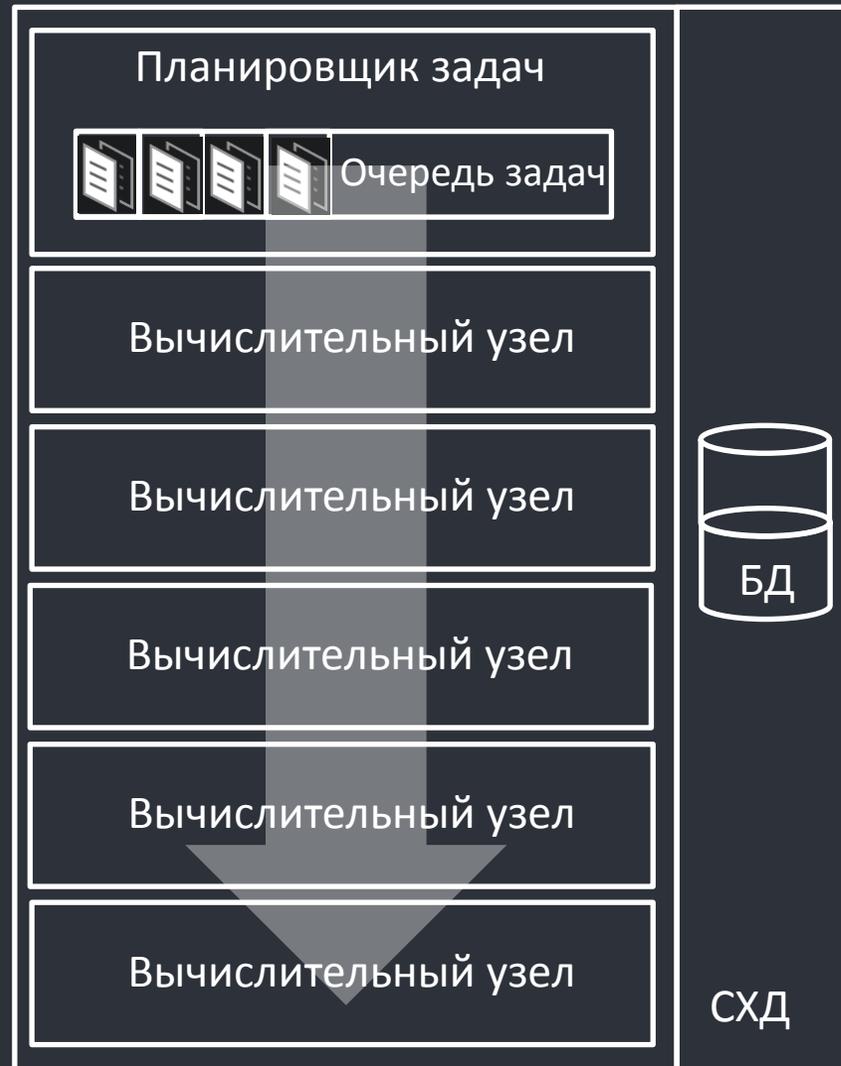
- Аппаратная платформа Intel x86 + Nvidia
- Аппаратная платформа Intel x86 + AMD
- Аппаратная платформа ARM (Байкал-М)
- Аппаратная платформа(ы) Эльбрус
- Аппаратная платформа Эльбрус + AMD
- Аппаратная платформа(ы) НТЦ Модуль
- Аппаратная платформа(ы) НПО Элвис
- Аппаратная платформа ПЛИС UltraScale+ (Xilinx)
- Аппаратная платформа(ы) IVA technologies

Схема работы системы

Клиентская часть



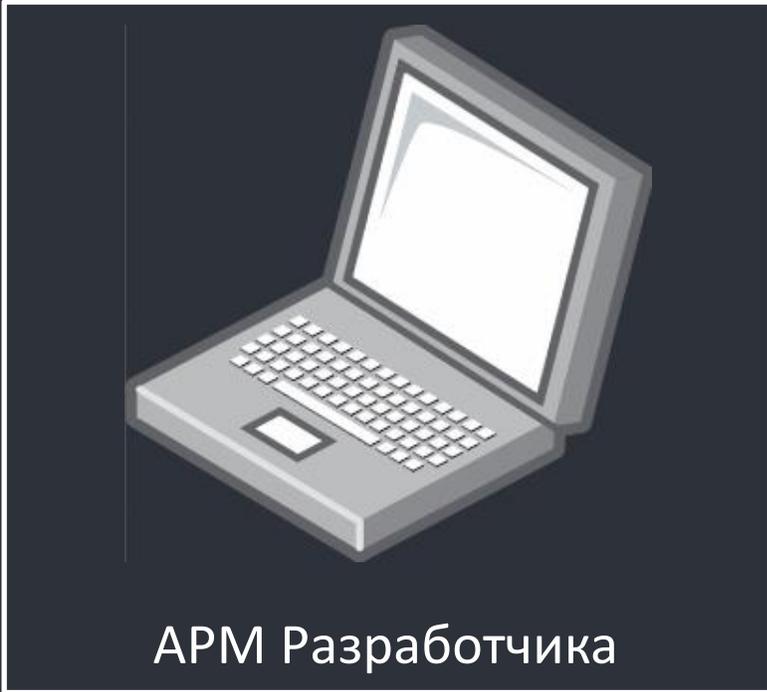
Серверная часть



Особенности:

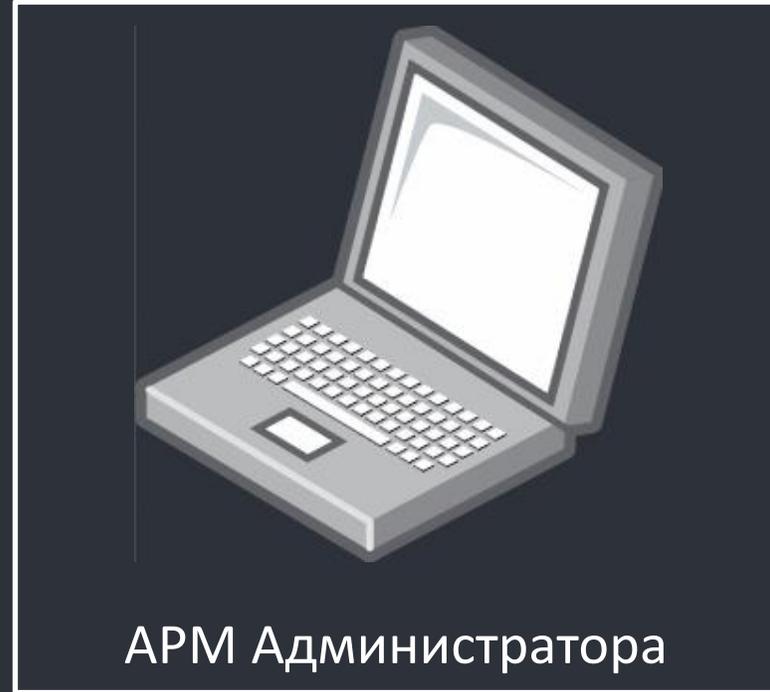
- ✓ Управление вычислительными ресурсами
- ✓ Автоматический выбор узлов
- ✓ Распределенное обучение
- ✓ Доступ к проектам и файлам
- ✓ Совместная работа над проектами
- ✓ Резервное копирование
- ✓ Логирование процесса разработки
- ✓ Удаленная работа с БД

Клиентская часть Платформы



Особенности:

- ✓ Поддержка ОС AltLinux/ ОС AstraLinux
- ✓ Поддержка архитектуры Эльбрус/x86
- ✓ Кроссплатформенная версия



Особенности:

- ✓ Поддержка ОС AltLinux/ ОС AstraLinux
- ✓ Поддержка архитектуры Эльбрус/x86
- ✓ Кроссплатформенная версия(в разработке)
- ✓ Web-интерфейс

Проект: основные понятия



АРМ Разработчика



АРМ Администратора

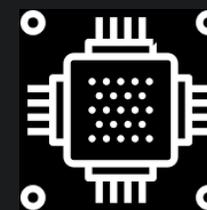
Проект – включает архитектуру ГСНС, базу аннотированных данных для обучения, а также алгоритмы и настройки параметров для обучения



БД



Алгоритм



Апп. реализация



Задание

Попытка #1

Архитектура
Параметры
Метод
Ноды
Выборки
Аугментация

Задание

Попытка #2

Архитектура
Параметры
Метод
Ноды
Выборки
Аугментация

Задание

Попытка #3

Архитектура
Параметры
Метод
Ноды
Выборки
Аугментация

Задание

Попытка #4

Архитектура
Параметры
Метод
Ноды
Выборки
Аугментация

Проект: уровни разработки

Три уровня разработки:

Июнь 2020



Начальный уровень (типичные решения):

+ готовые state of the art решения

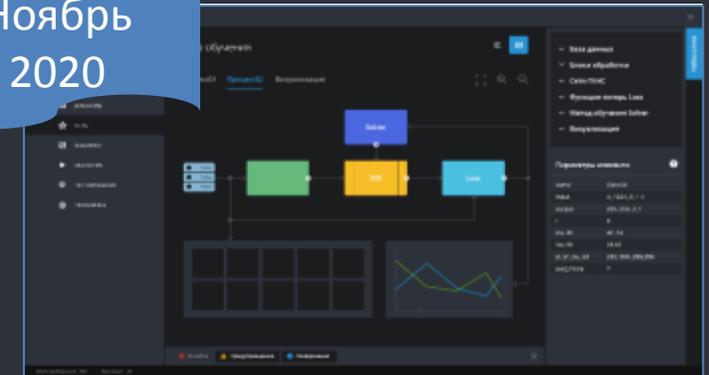
+ все промежуточные операции

+ набор типовых практических задач

+ распределенное обучение

+ оценка производительности

Ноябрь 2020



Средний уровень:

+ визуальное проектирование архитектуры

+ добавление пользовательских слоев

+ удобное и наглядное представление ГКНС

Ноябрь 2020

```
...
dataset = PT_Dataset(args.rootpath, args.pathDataset, CLS_DICT, transform=SDOAugmentation(cfg['mix_dim'],
...
dataset.pull_train()
...
ssd_net = build_ssd('train', 399, cfg['num_classes'])
net = ssd_net
...
if use_cuda:
    net = torch.nn.DataParallel(ssd_net)
    cuda_device = -1
...
if args.resume:
    print('Resuming training, loading {}...'.format(args.resume))
    ssd_net.load_state_dict(torch.load(args.resume))
...
vgg_weights = torch.load('pretrained_model/vgg16_reducefc.pth')
...
print('Loading base network...')
ssd_net.vgg.load_state_dict(vgg_weights)
...
if use_cuda:
    net = net.cuda()
...
if not args.resume:
    print('Initializing weights...')
    # initialize conv and fc layers weights with Xavier method
```

Продвинутый уровень:

+ написание кода на языке python

+ гибкая разработка в единой среде

+ использование сторонних библиотек

+ возможность написание типовых решений

+ средства визуализации



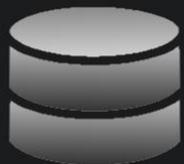
АРМ Разработчика



АРМ Администратора

Типовые решения

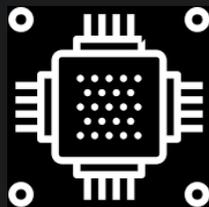
Типовое решение - включает архитектуру ГСНС, базу аннотированных данных для обучения, а также алгоритмы и настройки параметров для обучения



БД



Алгоритм



Апп. реализация



Защита от изменений

Список типовых задач:

- обнаружение объектов по изображениям и видеопоследовательностям;
- обнаружение объектов по многоспектральным данным;
- дешифрирование;
- классификация;
- семантическая сегментация;
- распознавание типов объекта;
- сопровождение объектов на видеопоследовательностях;
- обработка и комплексирование изображений различных диапазонов;
- устранение шумов и помех на изображениях (денойсинг);
- устранение смаза и расфокусировки изображений (деблур).

Особенности:

- **Гарантия работы**
- **Поддержка**
- **State of the art**
- **Аппаратная реализация**

Защита кода:

**Целостность кода
проверяется при старте
платформы**



Типовые задачи

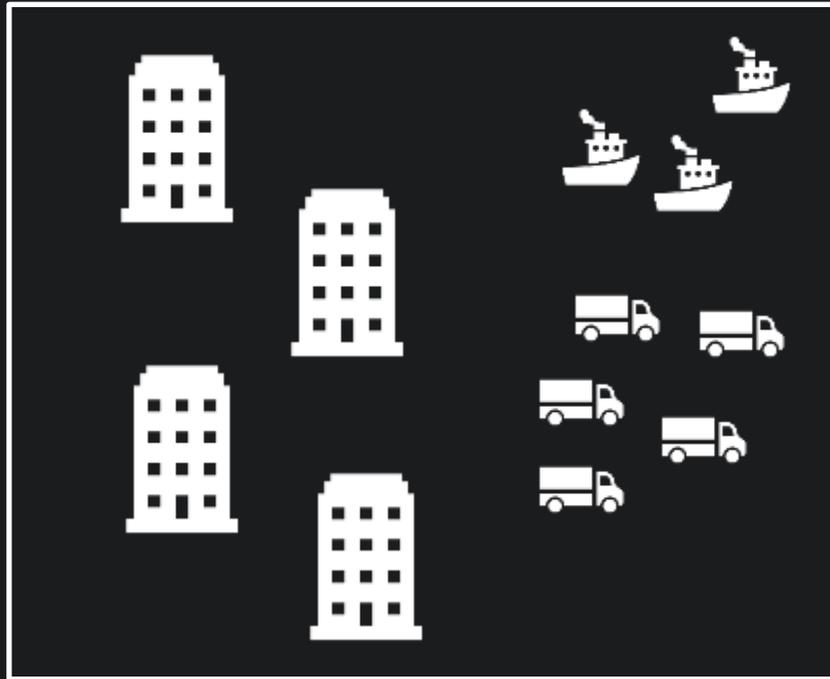
Обнаружение объектов по изображениям и видеопоследовательностям:



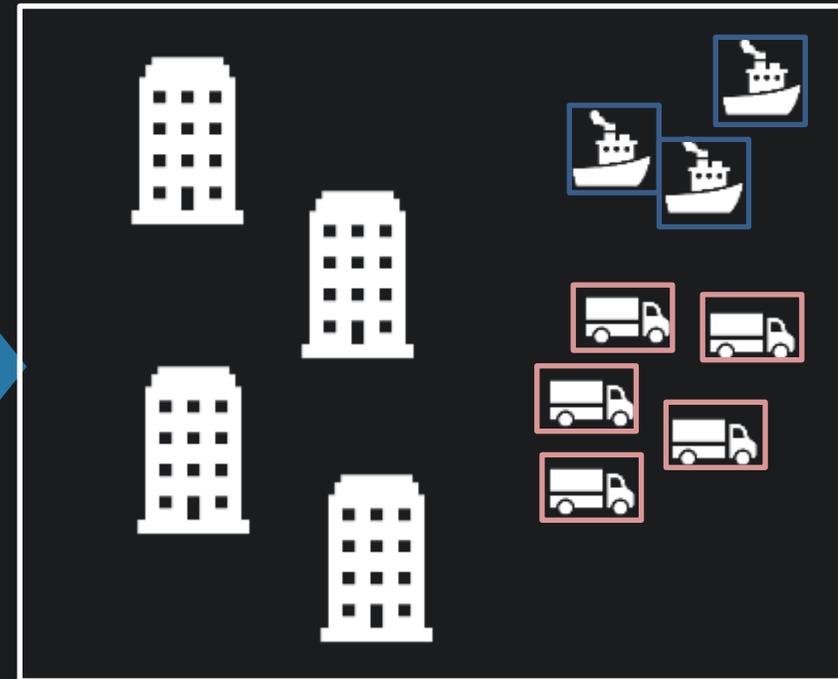
Выделение описывающих прямоугольников в экранных координатах для заданного набора классов

Типовые задачи

Дешифрирование :



Входное изображение высокого разрешения



Результат обнаружения

Выделение описывающих прямоугольников в экранных координатах для заданного набора классов

Типовые задачи

Семантическая сегментация :



Входное изображение



Результат сегментации

Выделение точных границ объектов целевых классов

Типовые задачи

Распознавание объектов :



Входное изображение



Класс грузовой автомобиль

Распознавание класса(типа) объекта на изображении

Типовые задачи

Обработка и комплексирование изображений различных диапазонов,
Устранение шумов и помех на изображениях
Устранение смаза и расфокусировки



Входное изображение
(один или несколько диапазонов)

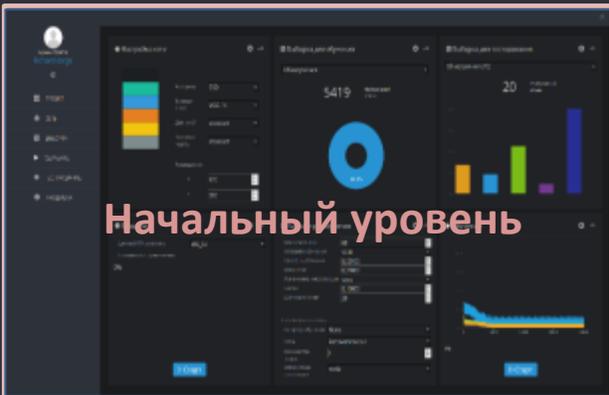


Восстановление максимально информативного изображения

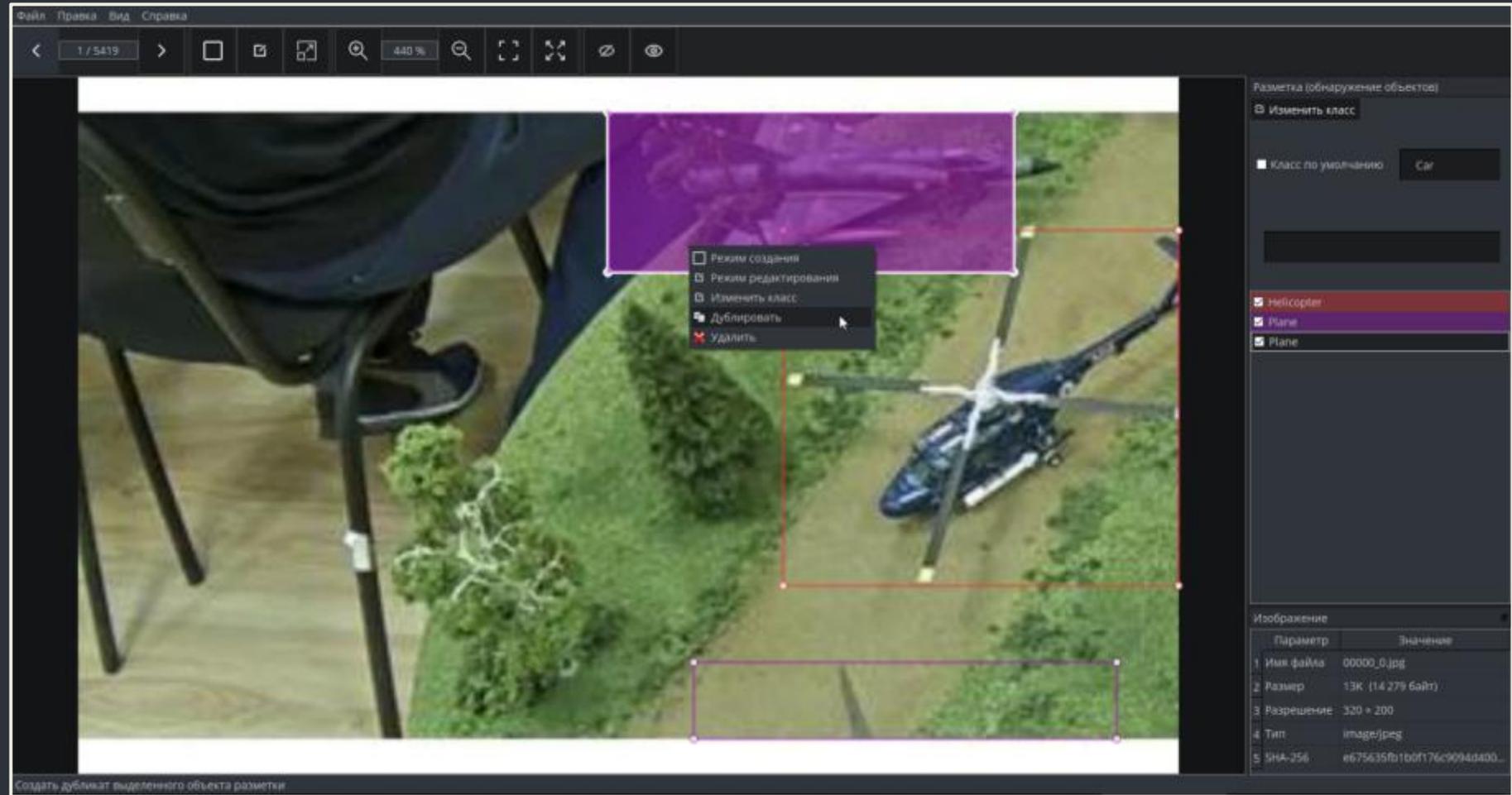
Проект: подготовка данных. СПО АРМ-Р: Разметчик



АРМ Разработчика



```
def train():
    batch_size = int(os.getenv('batch_size'))
    cfg = Argg
    dataset = PFTDataset(args.rootpath, args.pathdataset, C53_D0CT_Transform500augmentation(cfg.num_dir),
                        H5MS1)
    dataset.pull_images()
    val_set = build_val('train', 300, cfg.num_classes)
    val = val_set
    if use_cuda:
        net = torch.nn.DataParallel(val_net)
        net.cuda_device = 'cuda'
    if args.resume:
        print('Resuming training, loading {} ...'.format(args.resume))
        net.load_state_dict(torch.load(args.resume))
    val_weights = torch.load('pretrained_model/valgld_reducedc.pth')
    print('Loading base network...')
    net.load_state_dict(val_weights)
    if use_cuda:
        net = net.cuda()
    if args.weights:
        print('Initializing weights...')
        # initialize newly added layers' weights with kernel init
```



Режимы работы разметчика изображений:

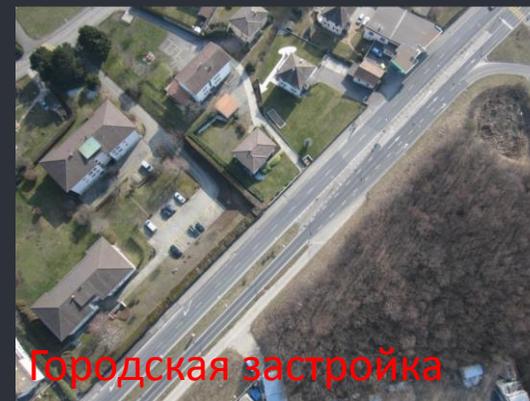
- предварительный просмотр изображения и существующей разметки с возможностью ее редактировать;
- создание новой разметки;
- создание разметки слежением от одного изображения к другому.

Проект: подготовка данных

Разметка изображений:



Разметка описывающими прямоугольниками



Установка атрибутов



Сегментация



АРМ Разработчика



АРМ Администратора

Проект: автоматизация подготовки данных



АРМ Разработчика

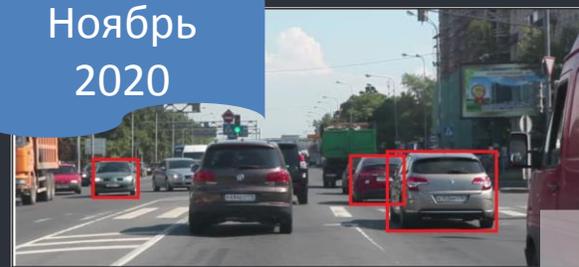
Июнь 2020



Межкадровая интерполяция
Алгоритмы слежения

Итеративное дообучение/разметка алгоритмом

Ноябрь
2020



ИИ сегментация объектов

Ноябрь
2020

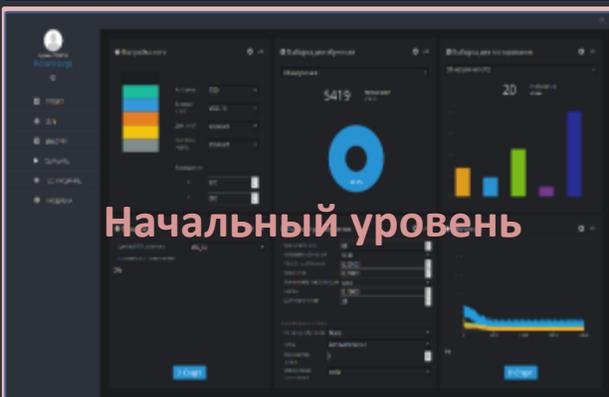


АРМ Администратора

Проект: начальный уровень. СПО АРМ-Р: Работа с проектами и заданиями



ARM Разработчика



```
def train():
    batch_size = int(os.getenv('batch_size'))
    cfg = App
    dataset = PTDataset(args.rootpath, args.pathdataset, OS.DOT_TRANSFORMER_SEGMENTATION(cfg.num_dir),
                       HEMAS)
    dataset.pull_images()
    val_set = build_val_loader(cfg.num_classes)
    net = val_net
    if use_cuda:
        net = torch.nn.DataParallel(net)
        torch.backends.cudnn.enabled = True
    # if use_cuda:
    #     print('Preparing training... loading (1) ... (format: cuda_device)')
    #     net = torch.nn.DataParallel(net)
    vgg_weights = torch.load('pretrained_models/vgg16_reduced.pth')
    print('Loading base network...')
    net.load_state_dict(vgg_weights)
    if use_cuda:
        net = net.cuda()
    # if use_cuda:
    #     print('Initializing weights...')
    #     print('Initializing weights... weights with Xavier normal')
```

ARM-Р 0.5.0.78.0

сервер: http://172.16.0.100/1/ проект: Обнаружение объектов (Обнаружение объектов) задание: vgg16

Администратор lab3010

ПРОЕКТ

СЕТЬ

ВЫБОРКА

ОБУЧЕНИЕ

ТЕСТИРОВАНИЕ

ПРОШИВКА

Настройки сети

Основа сети	vgg16
Вторая часть сети	standard
Третья часть сети	standard

Выборка для обучения

Обучение

5419 Изображений в выборке

Изображений по базам

100.0% Стол(обучение)

Выборка для тестирования

Тестирование

103 Различенных объектов

Изображений по классам

Class	Count
Car	29
Helicopter	21
Plane	14
Tank	7
Truck	0

Проверка

Целевой ПА комплекс

Проверочный

Готовность к применению

Параметры обучения

Деление базы данных	0,90
Последняя эпоха обучения	100
Частота сохранения	5
Частота тестирования	5
Перемешивание выборки	Да
Размер батча	8
Устройство для расчетов	cuda

Обучение

100%

Старт

Функции: создание и модификация выборки; создание задач для обучения и тестирования; настройка и запуск обучения задач; отображение графика обучения и статуса задачи; задание и просмотр настроек сети; показ статуса и статистики по базе данных; запуск тестирования уже обученной задачи; визуальный контроль результатов тестирования; запуск генерации прошивки для обученной задачи и выбранной архитектуры.

Проект: формирование выборок

Создание выборки с использованием мастера выборок:

Доступные БД

БД летных экспериментов

БД синтез

БД съемок 1

Фильтр

1. Выбрать по тегу “ночь”
2. Выбрать класс “машина”
3. Выбрать класс “самолет”

Фильтр

1. Выбрать 10000 изображений
2. Выбрать класс “грузовик”

Фильтр

1. Выбрать данные за апрель
2. Выбрать по тегу “ошибка”
3. Сделать черно-белыми
4. Изменить размер до 320x200

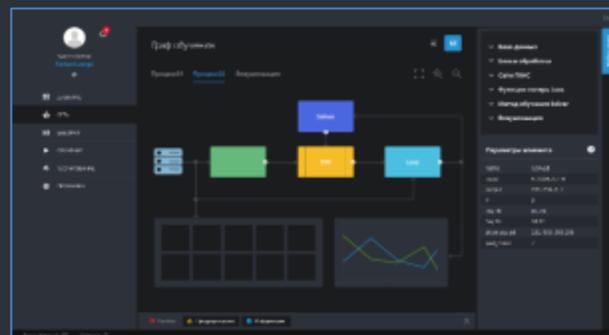
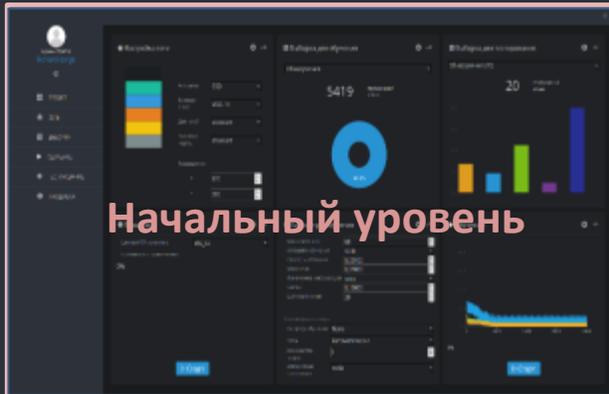
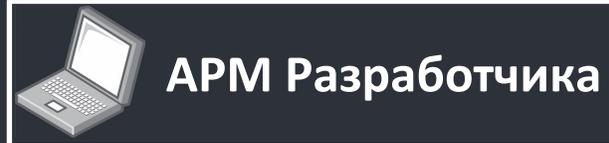


Обучающая

Тестовая

APM Разработчика

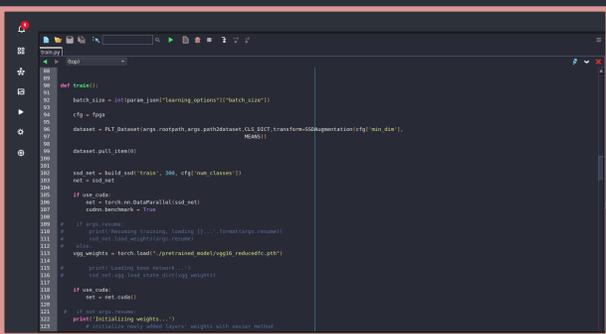
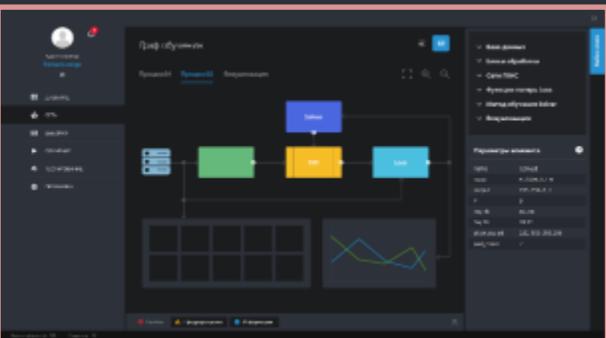
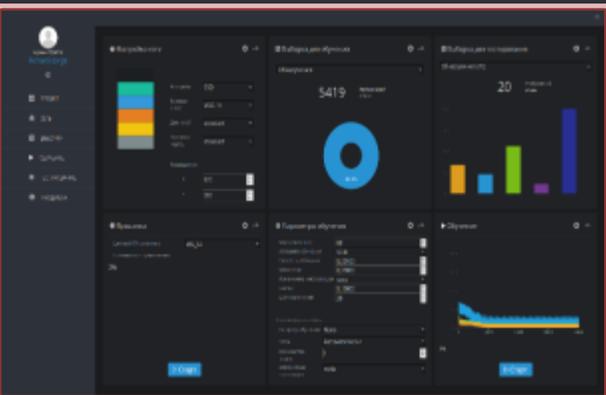
Начальный уровень



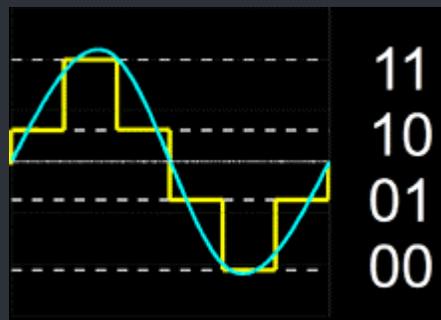
```
def train():
    batch_size = int(os.getenv('batch_size'))
    cfg = App
    dataset = FLT.Dataset(args.rootpath, args.pathdataset, CSL.Dict.Transform(500, repetition=(cfg.num_dir(),
                                                                                               REAS))
    dataset.pull_train()
    val_set = build_val('train', 300, cfg.num_classes())
    net = build_net
    if use_cuda:
        net = torch.nn.DataParallel(net, device_ids=range(torch.cuda.device_count()))
    if args.resume:
        print('Resuming training, loading {}...'.format(args.resume))
        net.load_state_dict(torch.load(args.resume))
    val_weights = torch.load('pretrained_model/vgg16_reducedfc.pth')
    print('Loading base network...')
    net.load_state_dict(torch.load_weights(val_weights))
    if use_cuda:
        net = net.cuda()
    print('Initializing weights...')
    # initialize newly added layers' weights with Xavier normal
```

Проект: Аппаратная реализация

Схема работы SDK Платформы



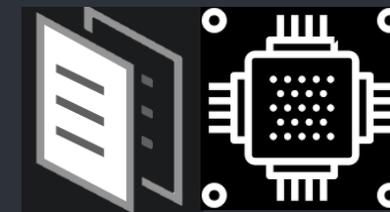
Планировщик



Квантизация



Трансформация архитектуры



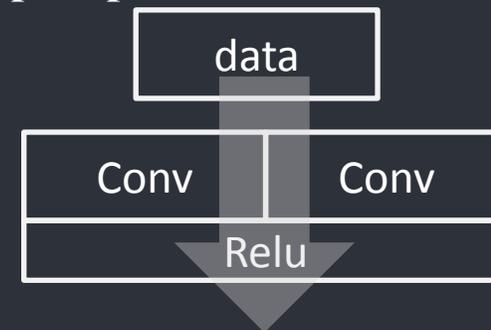
Аппаратно-ориентированное описание

Библиотека примитивов



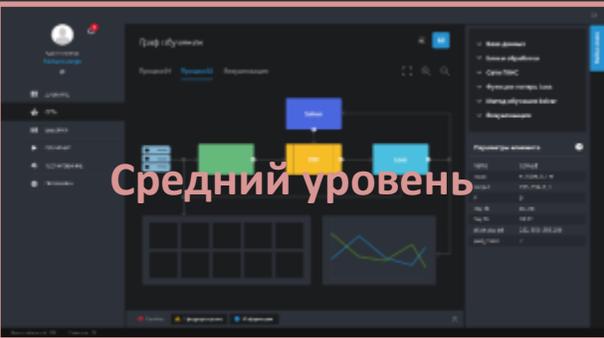
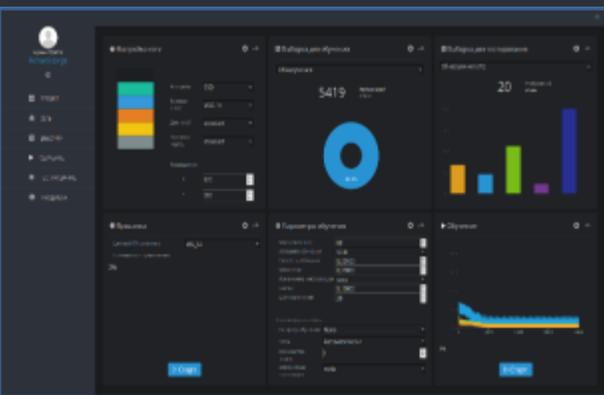
Оптимизированная аппаратная реализация

Оператор



Выполнение прямого прохода

Проект: средний уровень. Конструктор процессов (ноябрь 2020)



```
def train():
    batch_size = int(os.getenv('batch_size'))
    cuda = 'cpu'
    dataset = PTDataset(args.rootpath, args.pathdataset, CLS_Dict, Transform=GDAGeneration(cfg['non_dia'],
    HEMS))
    dataset.pull_image()
    val_loader = build_loader('train', 300, cfg['num_classes'])
    net = build_net()
    if use_cuda:
        net = torch.nn.DataParallel(net)
        cuda_device = -1
    # if use_cuda:
    #     print('Resolving training... loading (1) ... (format: cuda_device)')
    #     net = torch.nn.DataParallel(net)
    #     cuda_device = -1
    #     print('Loading base network...')
    #     net = torch.nn.DataParallel(net)
    #     cuda_device = -1
    if use_cuda:
        net = net.cuda()
    # if use_cuda:
    #     print('Initializing weights...')
    #     initialize_weights(net)
    #     print('Initializing weights...')
    #     initialize_weights(net)
```



Возможности:

- + Формирование процесса обучения/тестирования
- + Визуализация в реальном времени
- + Задание функций потерь и методов обучения
- + Математические операции

Проект: Продвинутый уровень. Редактирование/отладка кода (ноябрь 2020)



APM Разработчика



```
train.py
88
89
90 def train():
91
92     batch_size = int(param_json["learning_options"]["batch_size"])
93
94     cfg = fpga
95
96     dataset = PLT_Dataset(args.rootpath, args.path2dataset, CLS_DICT, transform=SSDAugmentation(cfg['min_dim'],
97                                                         MEANS))
98
99     dataset.pull_item(0)
100
101
102     ssd_net = build_ssd('train', 300, cfg['num_classes'])
103     net = ssd_net
104
105     if use_cuda:
106         net = torch.nn.DataParallel(ssd_net)
107         cudnn.benchmark = True
108
109     # if args.resume:
110     #     print('Resuming training, loading {}'.format(args.resume))
111     #     ssd_net.load_weights(args.resume)
112     # else:
113     vgg_weights = torch.load("./pretrained_model/vgg16_reducedfc.pth")
114
115     #     print('Loading base network...')
116     #     ssd_net.vgg.load_state_dict(vgg_weights)
117
118     if use_cuda:
119         net = net.cuda()
120
121     # if not args.resume:
122     print('Initializing weights...')
123     # initialize newly added layers' weights with xavier method
```

```
train.py
88
89
90 def train():
91
92     batch_size = int(param_json["learning_options"]["batch_size"])
93
94     cfg = fpga
95
96     dataset = PLT_Dataset(args.rootpath, args.path2dataset, CLS_DICT, transform=SSDAugmentation(cfg['min_dim'],
97                                                         MEANS))
98
99     dataset.pull_item(0)
100
101
102     ssd_net = build_ssd('train', 300, cfg['num_classes'])
103     net = ssd_net
104
105     if use_cuda:
106         net = torch.nn.DataParallel(ssd_net)
107         cudnn.benchmark = True
108
109     # if args.resume:
110     #     print('Resuming training, loading {}'.format(args.resume))
111     #     ssd_net.load_weights(args.resume)
112     # else:
113     vgg_weights = torch.load("./pretrained_model/vgg16_reducedfc.pth")
114
115     #     print('Loading base network...')
116     #     ssd_net.vgg.load_state_dict(vgg_weights)
117
118     if use_cuda:
119         net = net.cuda()
120
121     # if not args.resume:
122     print('Initializing weights...')
123     # initialize newly added layers' weights with xavier method
```

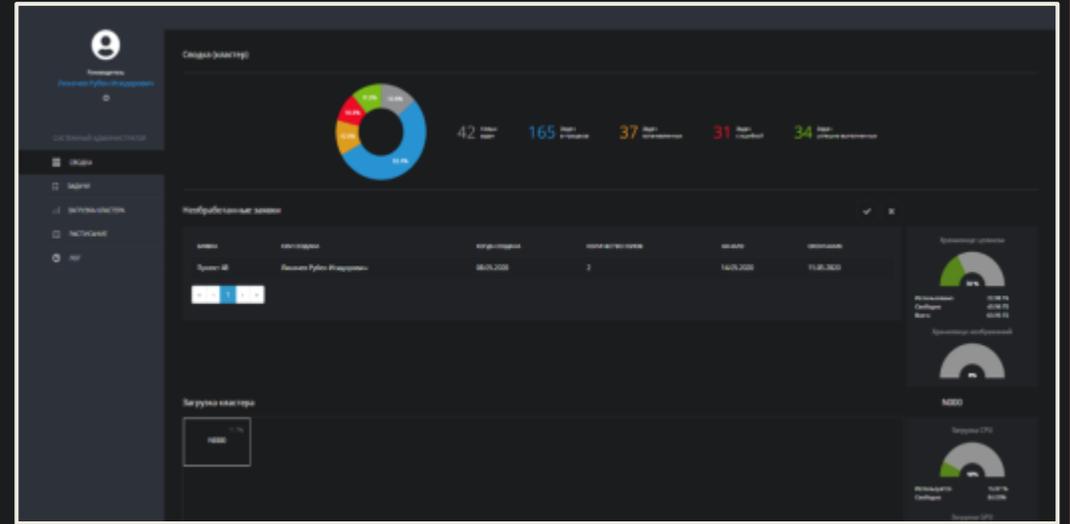
Администрирование. СПО ARM-A



ARM Разработчика

Администратор проектов

- ✓ Создание проектов
- ✓ Добавление пользователей в проект
- ✓ Доступ пользователей к БД
- ✓ Отслеживание прогресса по проекту



ARM Администратора



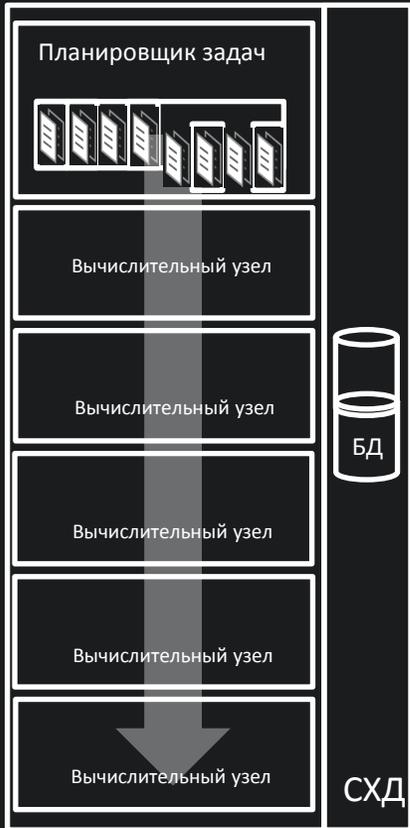
Администратор ресурсов

- ✓ Выделение вычислительных ресурсов
- ✓ Остановка/отмена заданий
- ✓ Резервное копирование
- ✓ Обновление системы
- ✓ Контроль состояния кластера/сервера

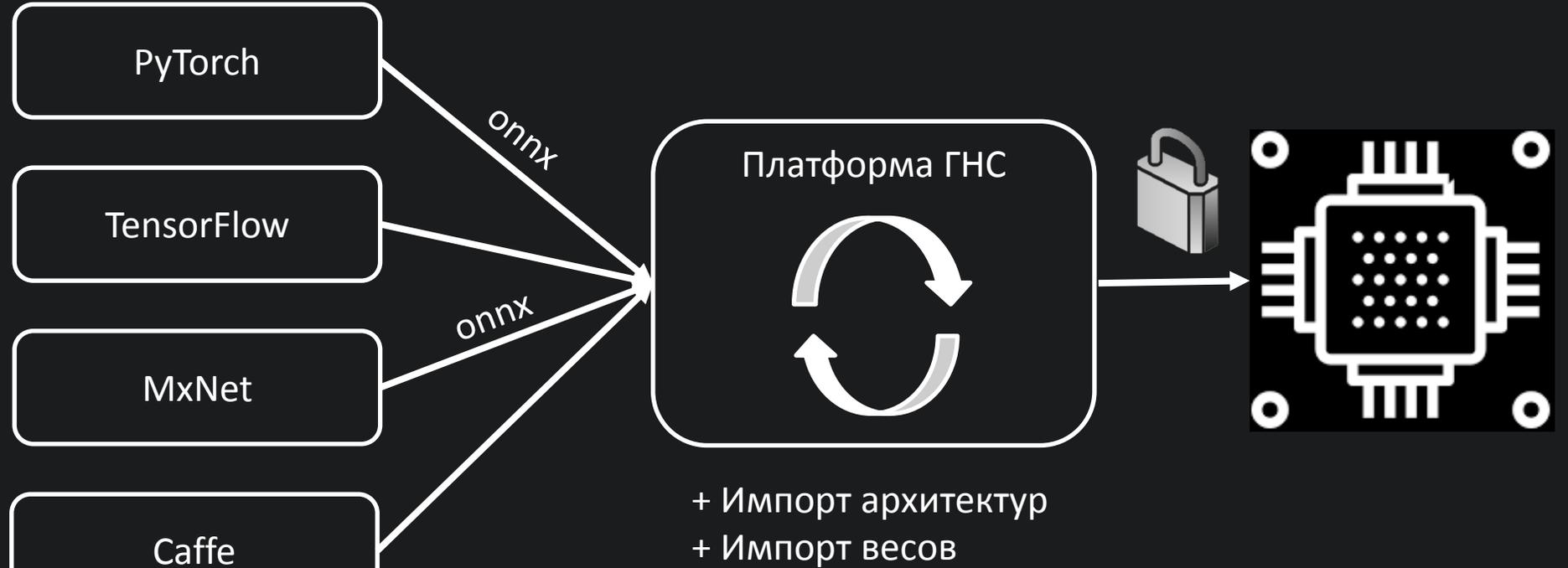
Импорт решений

Инфраструктура

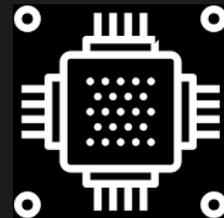
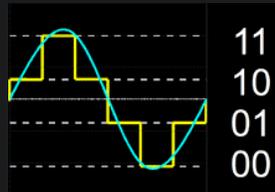
- Расп. обучение
- Контроль доступа
- БД
- Дообучение



Импорт решений из других сред



Аппаратная реализация



Сертифицированное ПО



Единый формат



Перспективы развития Платформы (2021+)

Технологическое развитие инфраструктуры

- ✓ Создание виртуальных кластеров
- ✓ Управление и мониторинг

Использование нескольких вычислительных кластеров

Развитие технологии обучения

Автоматическое обучение (AutoML)

- ✓ Формирование архитектуры под задачу
- ✓ Подбор гиперпараметров обучения
- ✓ AutoML как сервис

Задачи управления и ИИ

- ✓ Решение задач управления
- ✓ Решение задач планирования
- ✓ Эффективная работа на кластере

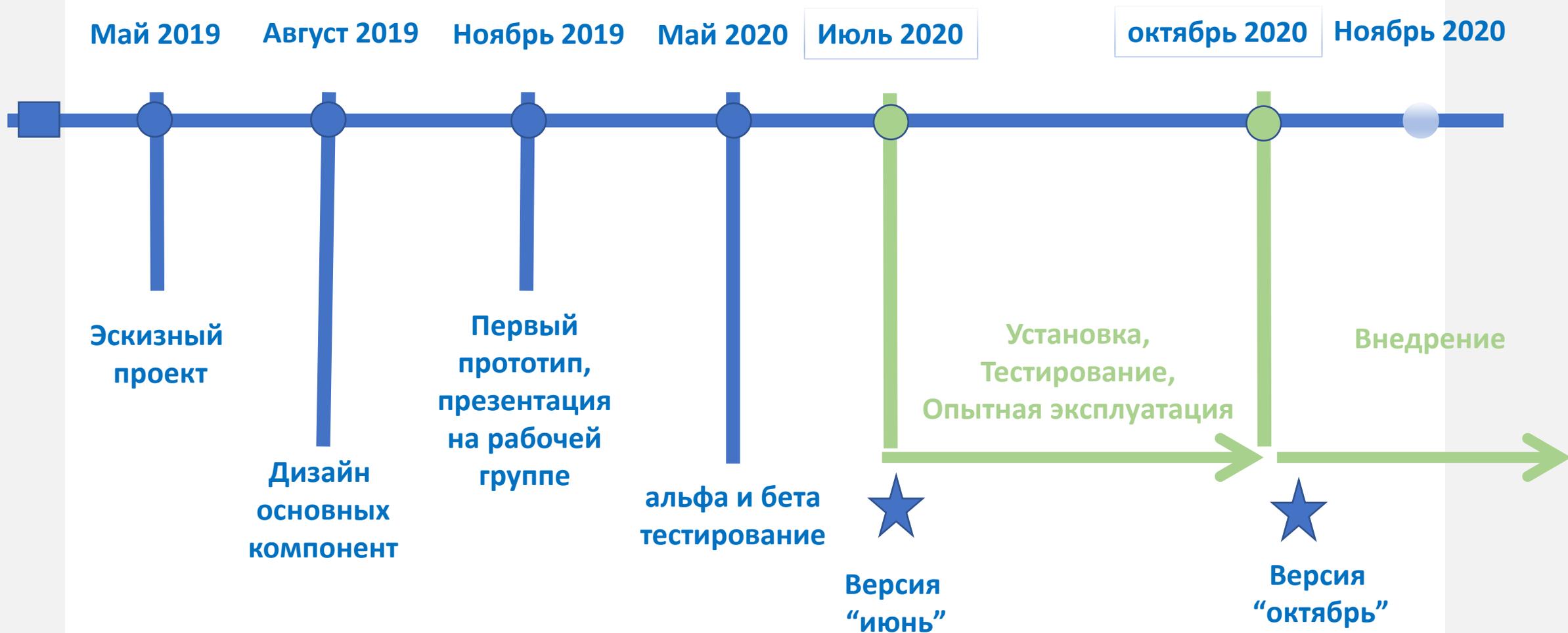
Обучение с подкреплением

Задачи на произвольных типах данных

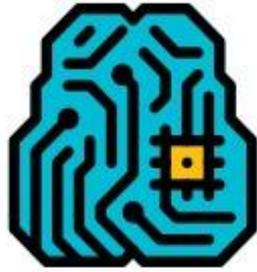
Графовые сети

- ✓ Задачи оптимизации
- ✓ Задачи большой размерности

Основные этапы проекта «Платформа-ГНС»



Будущая экосистема Платформы-ГНС. Переход к этапу тестирования и внедрения



Разработчики аппаратных платформ



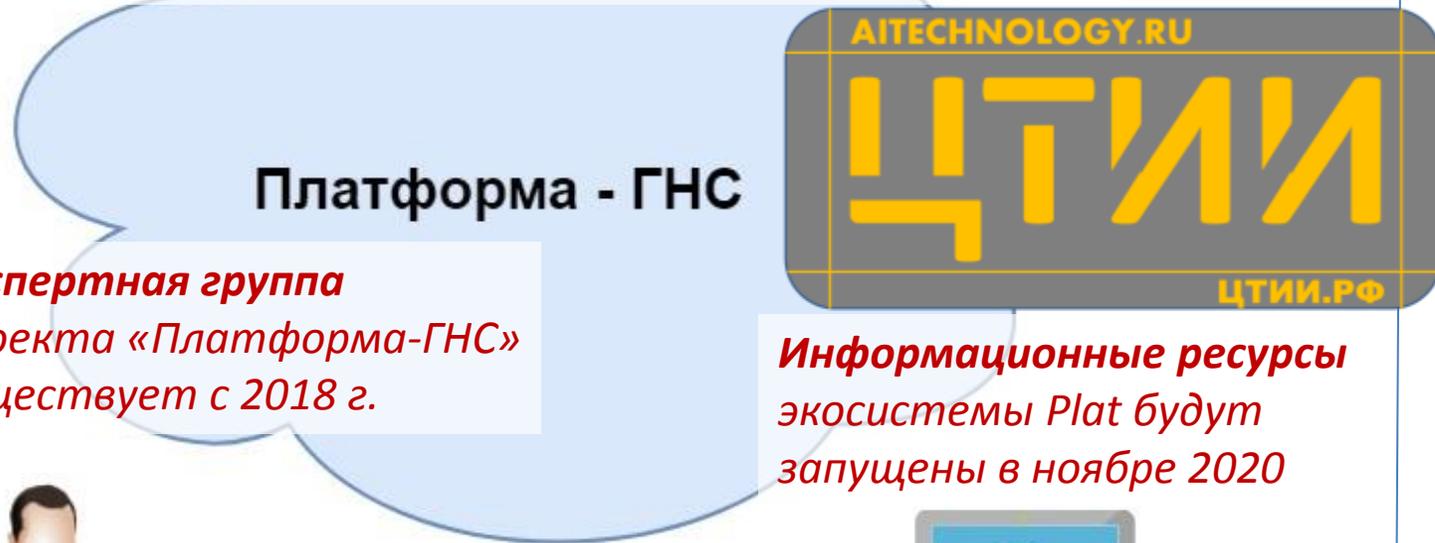
Разработчики алгоритмов

Контакт для участия в пилотном тестировании и опытной эксплуатации:

Оператор: Центр технологий искусственного интеллекта НИЦ «Институт им. Н.Е. Жуковского» (ЦТИИ) создан в 2020

support.plat@gosniias.ru

Невозможно вести такую разработку без учета мнения всех участников рабочего процесса!

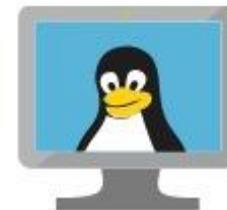


Экспертная группа проекта «Платформа-ГНС» существует с 2018 г.

Информационные ресурсы экосистемы Plat будут запущены в ноябре 2020



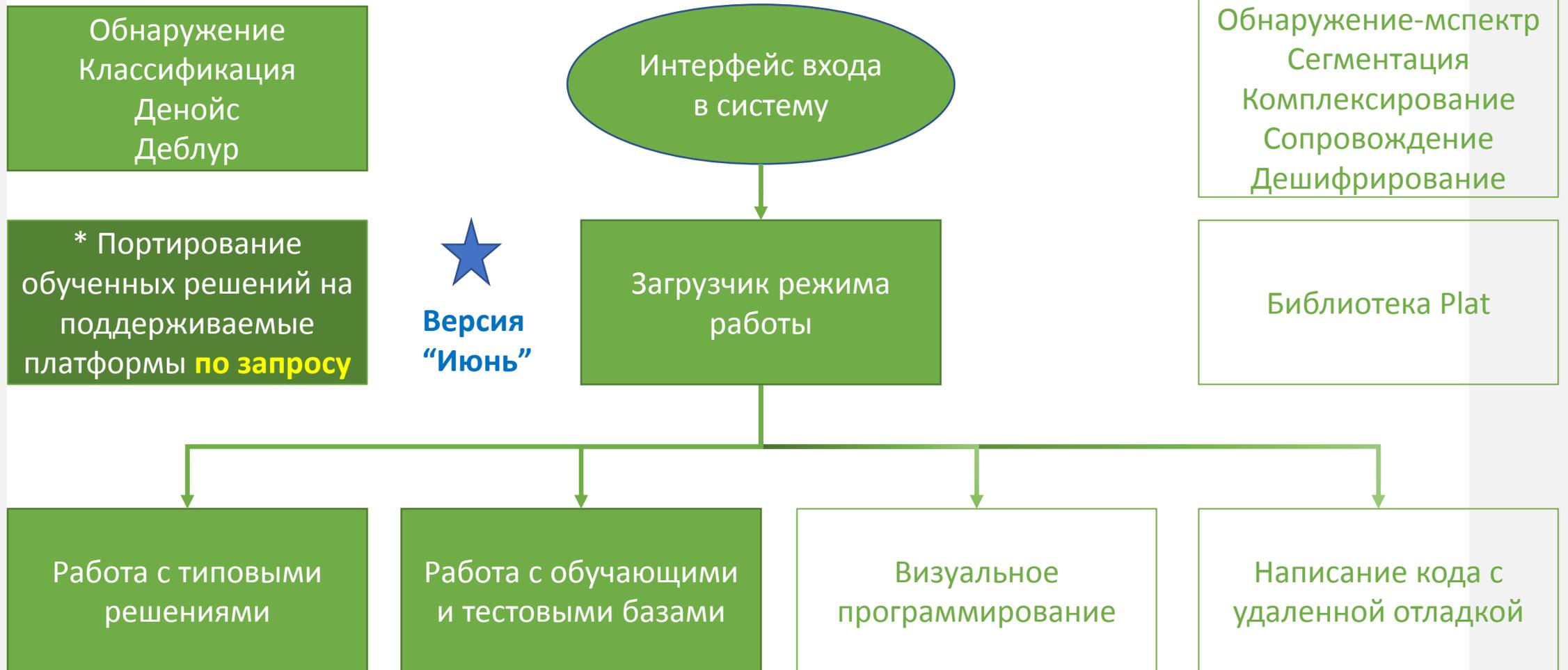
Конечные потребители



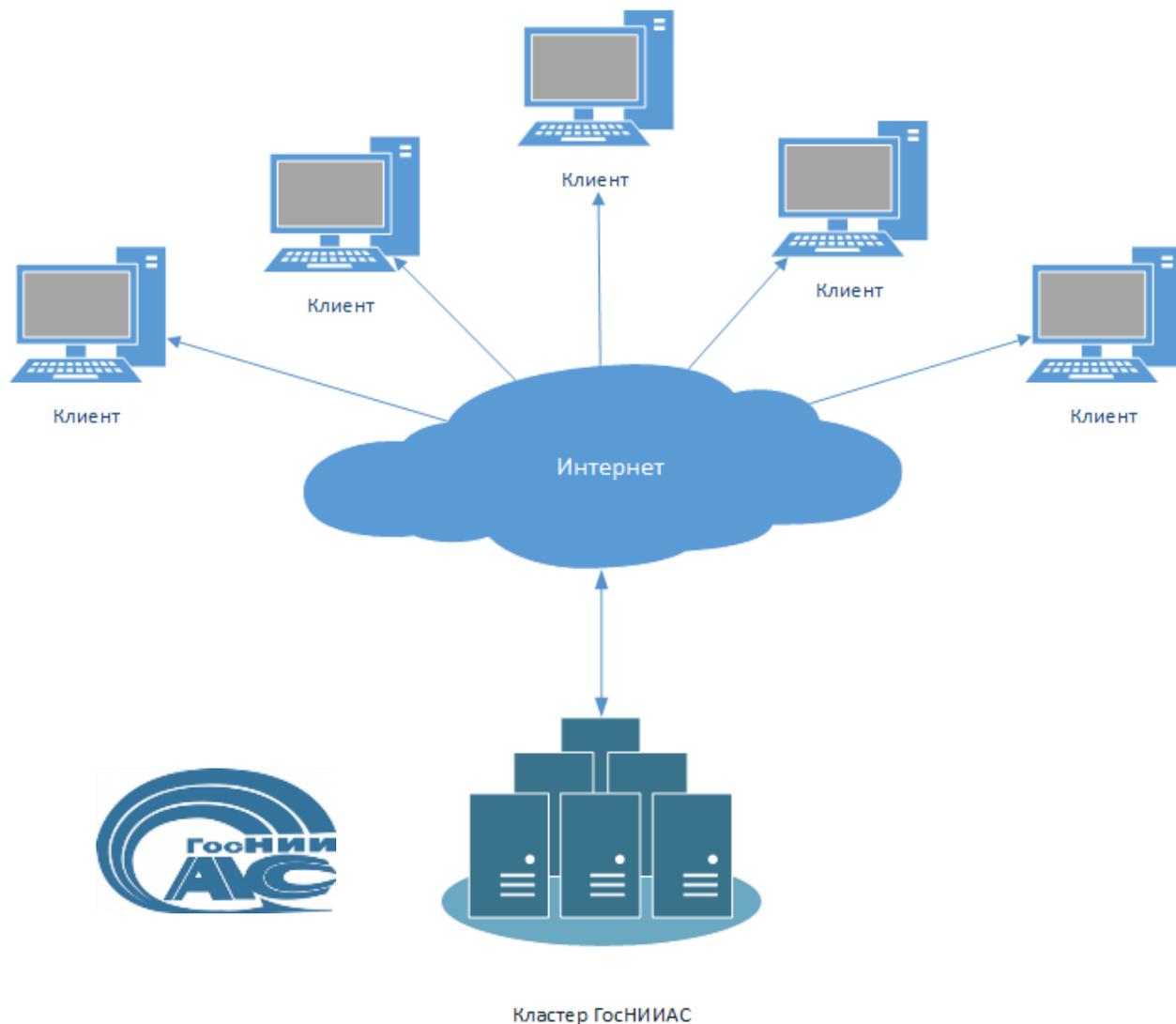
Разработчики программных платформ



ARM-R: возможности версии "июнь"



Технологическая схема участия в пилотном тестировании и опытной эксплуатации



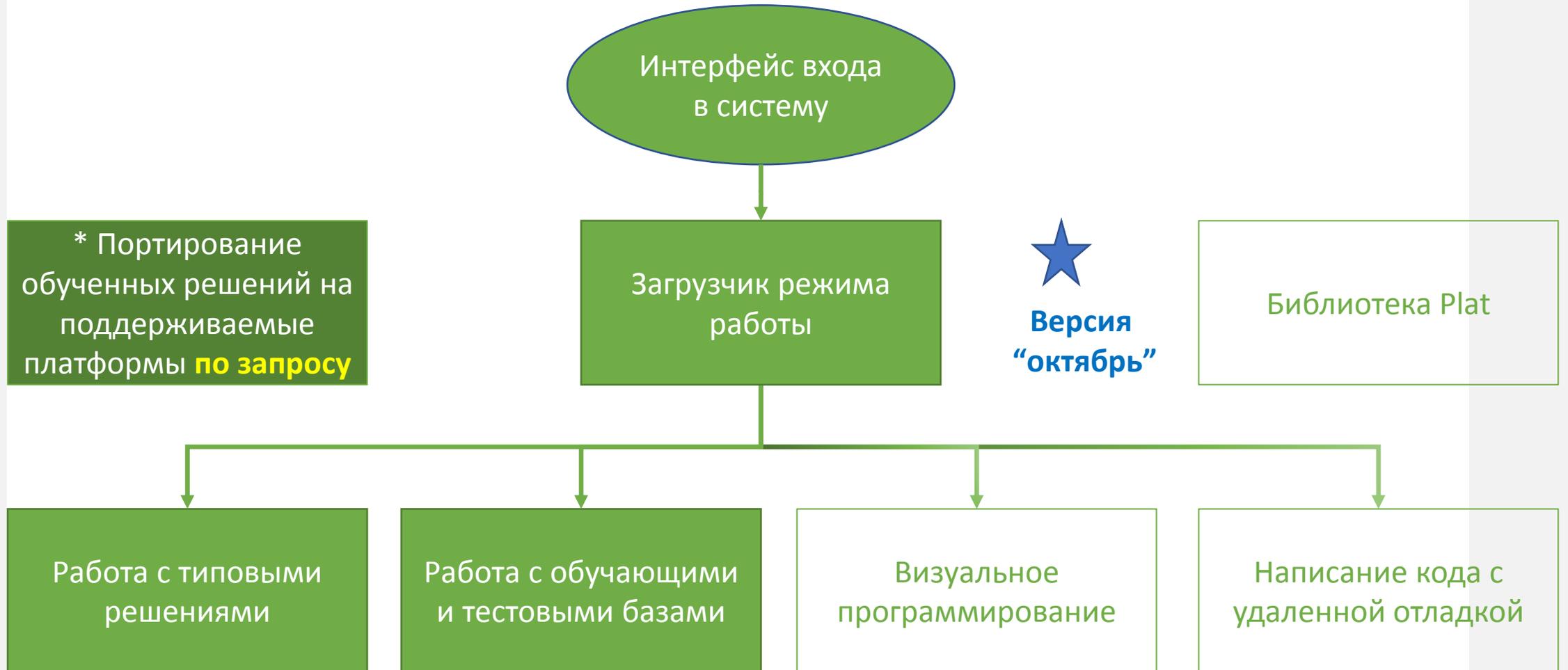
- Кластер ГосНИИАС
- Клиент (кросс-платформенный) – на стороне Заказчика
- Доступные типовые решения:
 - Классификация
 - Обнаружения объектов
 - Устранение шумов
 - Устранение смаза
- Установка на собственный кластер: по запросу

Требования к рабочему месту



1. Процессор: 64-битный Intel или AMD с частотой не менее 2ГГц;
2. ОЗУ не менее 4ГБ;
3. Видеокарта с поддержкой OpenGL 4.0;
4. Экран с разрешением не менее 1920x1080;
5. Свободное место на жестком диске не менее 500МБ;
6. Подключение к серверу со скоростью не менее 100Мбит/с;
7. Операционная система (Windows 7/8/10 или Ubuntu Linux 18.04 и новее или Astra Linux 1.6 релиз «Смоленск»).

АРМ-Р: возможности версии “октябрь”



Графическое программирование

Конструктор процессов:



Версия
“октябрь”



Возможности:

- + Формирование процесса обучения/тестирования
- + Визуализация в реальном времени
- + Задание функций потерь и методов обучения



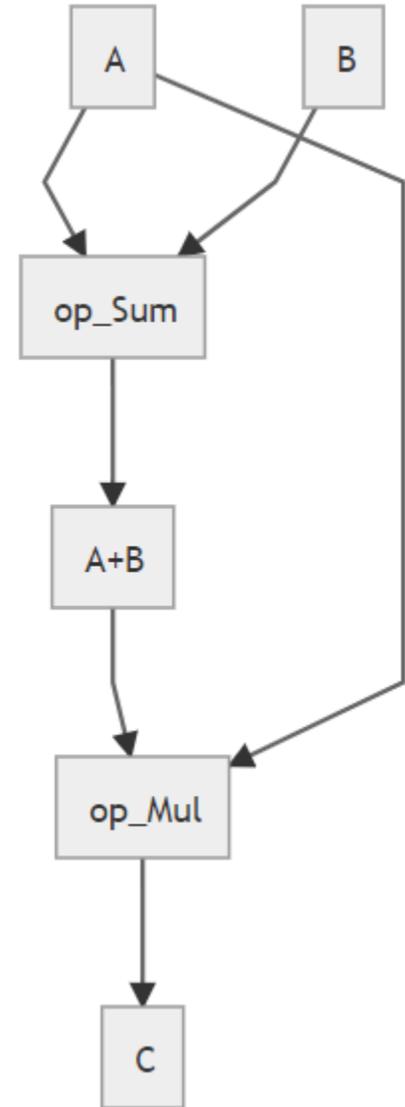
Версия
"октябрь"

PlatLib

Обучение:

- Динамические графы
- Распределенное обучение
- Поддержка GPU AMD/NVIDIA
- Поддержка CPU x86-64/Elbrus
- Поддержка БД Платформы

```
A = PTensor(5,5,5,5)
B = PTensor(5,5,5,5)
C = (A+B)*A
```



PlatLib

Обучение:

- Аналогично pytorch
- Отладка на python
- Нет multiGPU
- Распределенное обучение через nccl/rccl/mpi

```
class LeNetPlat(plat.module.Module):  
    def __init__(self):  
        super(LeNetPlat, self).__init__()  
        self.conv1 = conv2D(20, 1, 5)  
        self.relu1 = relu()  
        self.pool1 = poolmax2d()  
        self.conv2 = conv2D(50, 20, 5)  
        self.relu2 = relu()  
        self.pool2 = poolmax2d()  
        self.fc1 = conv2D(500, 50, 4)  
        self.relu3 = relu()  
        self.fc2 = conv2D(10, 500, 1)
```

```
res_plt = LeNetPlat(inputT)  
plt_res = SM(res_plt, labelsT)  
plt_res.backward()  
optimizer_plat.iter()
```

Программирование в стиле plat: набор примитивов Proj

Обучение:

- Создание собственного типового решения Платформы
- Разделение по стадиям
- Хорошая читаемость кода
- Синхронные аугментации
- Именованные типы тензоров
- Автоматическая генерация

параметров в формате Платформы
для отображения и передачи в АРМ

```
class LeNetProject(plat.proj.Network):
    def __init__(self):
        super(plat.proj.Network, self).__init__()
        self.net=plat.proj.zoo.LeNet((28,28),10)
        self.loss=plat.nn.softmaxwloss()

    def stepTrain(self,data,labels):
        return({'loss',self.loss(self.net(data))})

plat.Trainer.fit(LeNetProject(),plat.db.MNIST())
```

**Современное состояние технологий
искусственного интеллекта и отечественная
нейросетевая платформа Plat**

Ю.В. Визильтер, д.ф.-м.н., проф. РАН, viz@gosniias.ru

Спасибо за внимание!

Контакт по Платформе Plat:

support.plat@gosniias.ru



Семинар МГУ
Москва, 13.10.2020