

VLA

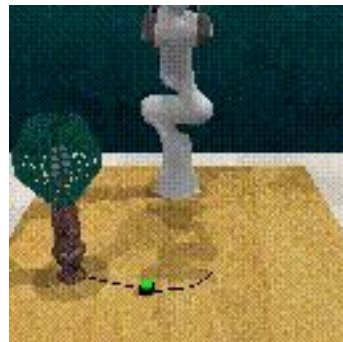
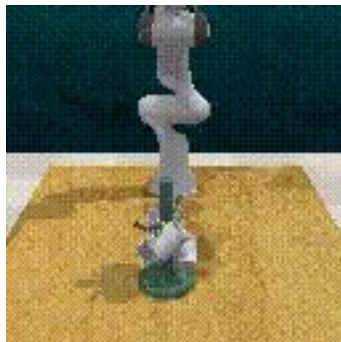
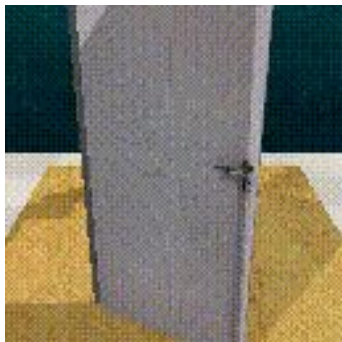
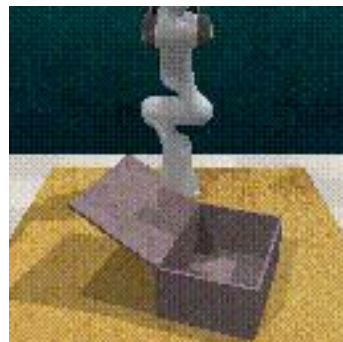
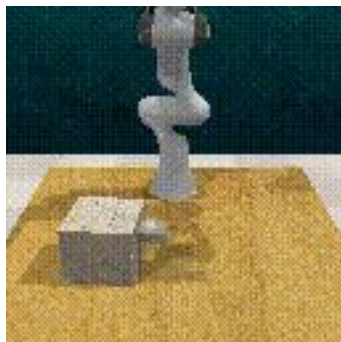
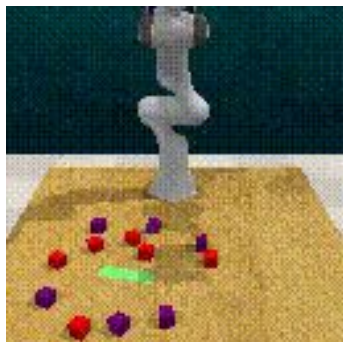
16 october 2025

О чем сегодня поговорим?

- VLA
 - ◆ Зачем это нужно?
 - ◆ Как это работает?
- Pi0, Pi0.5
- OpenVLA
- GROOT



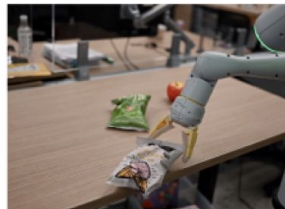
Задачи для VLA



Real World задачи для VLA



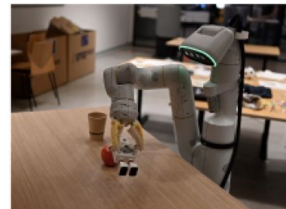
put strawberry
into the correct
bowl



pick up the bag
about to fall
off the table



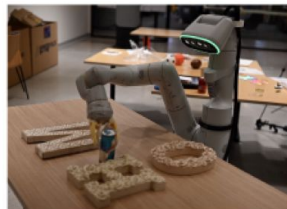
move apple to
Denver Nuggets



pick robot



place orange in
matching bowl



move redbull can
to H



move soccer ball
to basketball



move banana to
Germany



move cup to the
wine bottle



pick animal with
different colour

Real World задачи для VLA

Prompt:

Given Instruction:
Bring me a drink.

Prediction:

Plan: pick 7up can.

Action: 1 143 129 123 145
114 115 127



Prompt:

Given Instruction:
Move all the objects
together.

Prediction:

Plan: move green can near
green rice chip bag.
Action: 1 128 126 127 135
123 119 127



Prompt:

Given Instruction:
Pick the object that is
different from all other
objects

Prediction:

Plan: pick rxbar
chocolate. Action: 1 128
129 125 131 125 128 127

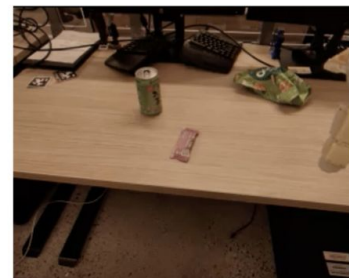


Prompt:

Given Instruction:
Move the green objects
together.

Prediction:

Plan: move green can near
green rice chip bag.
Action: 1 130 129 121 131
127 128 127

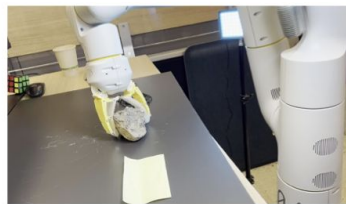
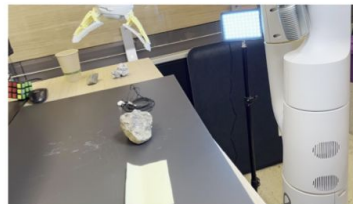


Prompt:

Given I need to
hammer a nail, what
object from the scene
might be useful?

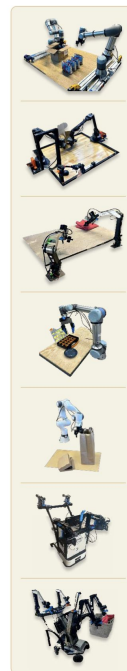
Prediction:

Rocks. Action: 1 129 138
122 132 135 106 127



Проблемы традиционной робототехники

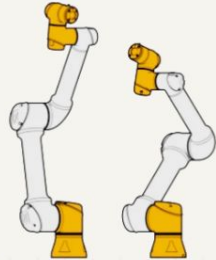
- **Узкая специализация:** Каждый робот или задача требовали создания и обучения отдельных моделей, оптимизированных под конкретные условия.
- **Необходимость переобучения:** Для каждой новой задачи или изменения среды требовалось трудоёмкое и дорогостоящее переобучение всей системы.
- **Изолированные решения:** Отсутствие единой архитектуры приводило к созданию разрозненных, несовместимых решений для разных типов роботов и платформ.
- **Высокие затраты на разработку:** Кастомная разработка для каждого применения влекла за собой значительные временные и финансовые издержки.



UR5e



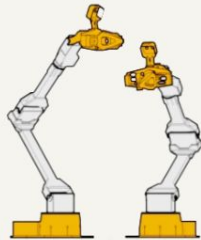
Bimanual UR5e



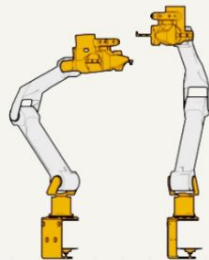
Franka



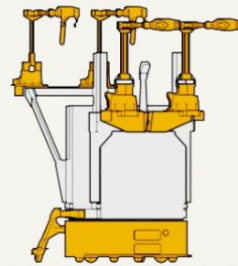
Bimanual Trossen



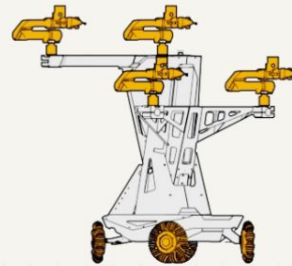
Bimanual Arx



Mobile Trossen



Mobile Fibocom



Foundational models

В контексте робототехники **foundation model** — это обобщающая политика (policy), предобученная на множестве гетерогенных данных: визуальных наблюдениях, действиях, обратной связи и состояниях различных роботов.

Такой подход решает три ключевые задачи:

1. **Инвариантность к воплощению (embodiment)** — одна и та же модель должна работать на разных роботах (рука, мобильная платформа, манипулятор и т.д.).
2. **Инвариантность к задаче** — способность выполнять новые задачи без специальной дообучения.
3. **Инвариантность к среде** — перенос между различными сценами, объектами, освещением, фоном.

Примеры из других областей:



GPT

Языковые задачи



CLIP / DINO

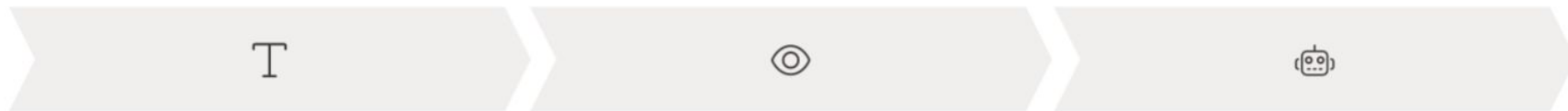
Визуальное понимание



SAM

Сегментация изображений

От текста к действию



LLM

Понимание языка и текстовые ответы

VLM

Интеграция визуального восприятия

VLA

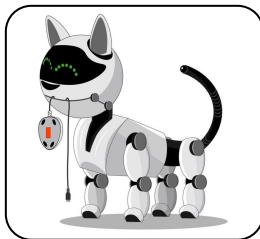
Генерация физических действий в реальном мире

От текста к действию

На карти -нке белый робо -кот [END]

VLM

[img] [img] На карти -нке белый робо -кот [END]



T

LLM

Понимание языка и текстовые ответы



VLM

Интеграция визуального восприятия



VLA

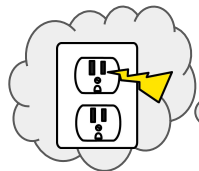
Генерация физических действий в реальном мире



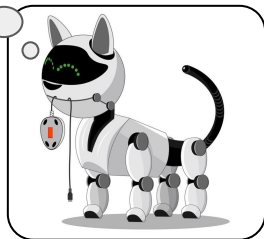
От текста к действию

Пора кор -мить кота ACT1 ACT2 ACT3 ACT3 ACT4

VLA



[img] [img] Пора кор -мить кота ACT1 ACT2 ACT3 ACT3



T

LLM

Понимание языка и текстовые ответы



VLM

Интеграция визуального восприятия



VLA

Генерация физических действий в реальном мире

3 кита современных VLA



3 кита современных VLA



п0 / п0.5

Physical Intelligence

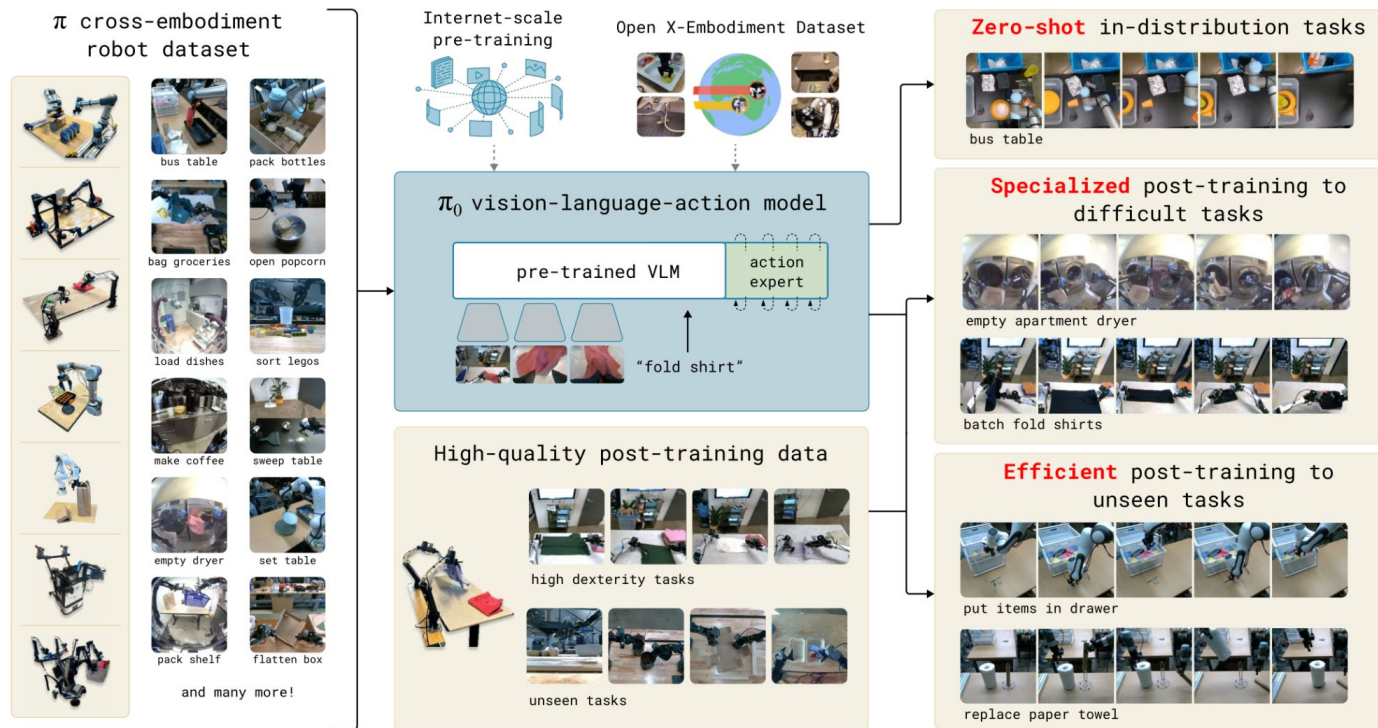
Масштабируемая архитектура с
единым токеном пространством для
языка и действий

- Более 10M эпизодов обучения
- Поддержка 3D-восприятия
- Впечатляющее zero-shot
обобщение

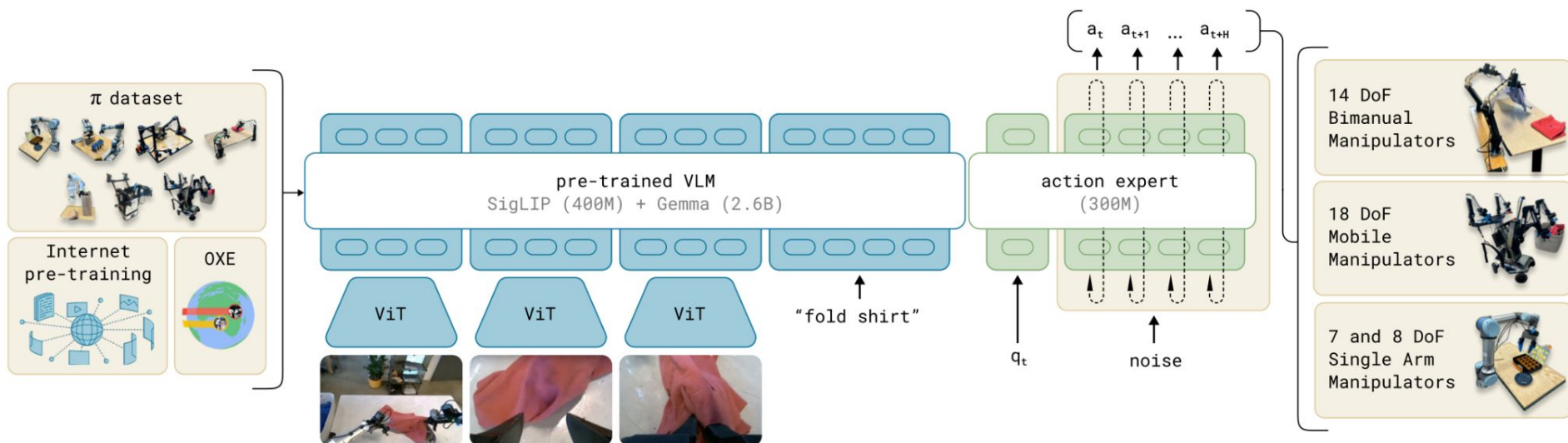


π_0

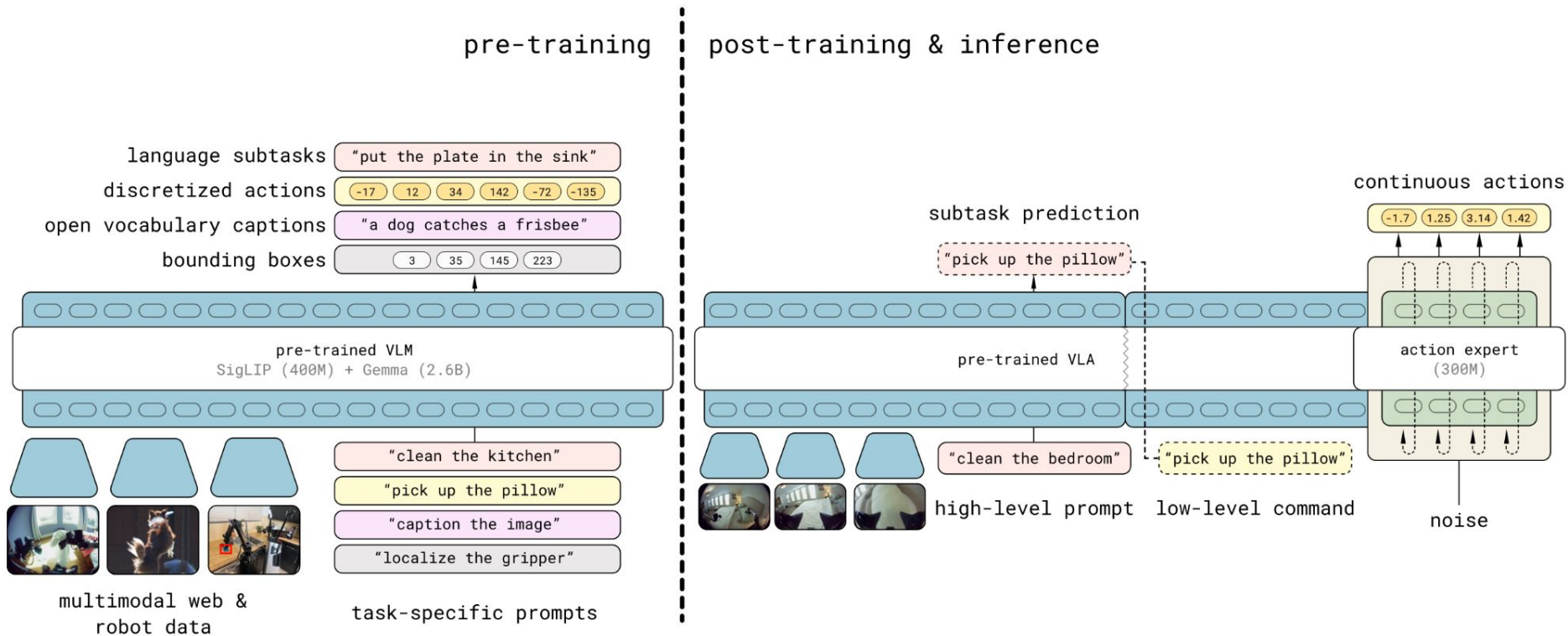
[CLS] возьми красный куб [ACTION_X+0.02] [ACTION_Y-0.01] [GRIP_CLOSE] [EOS]



π_0



$\pi_{0.5}$



3 кита современных VLA



п0 / п0.5

Physical Intelligence

Масштабируемая архитектура с единым токеном пространством для языка и действий

- Более 10M эпизодов обучения
- Поддержка 3D-восприятия
- Впечатляющее zero-shot обобщение



OpenVLA

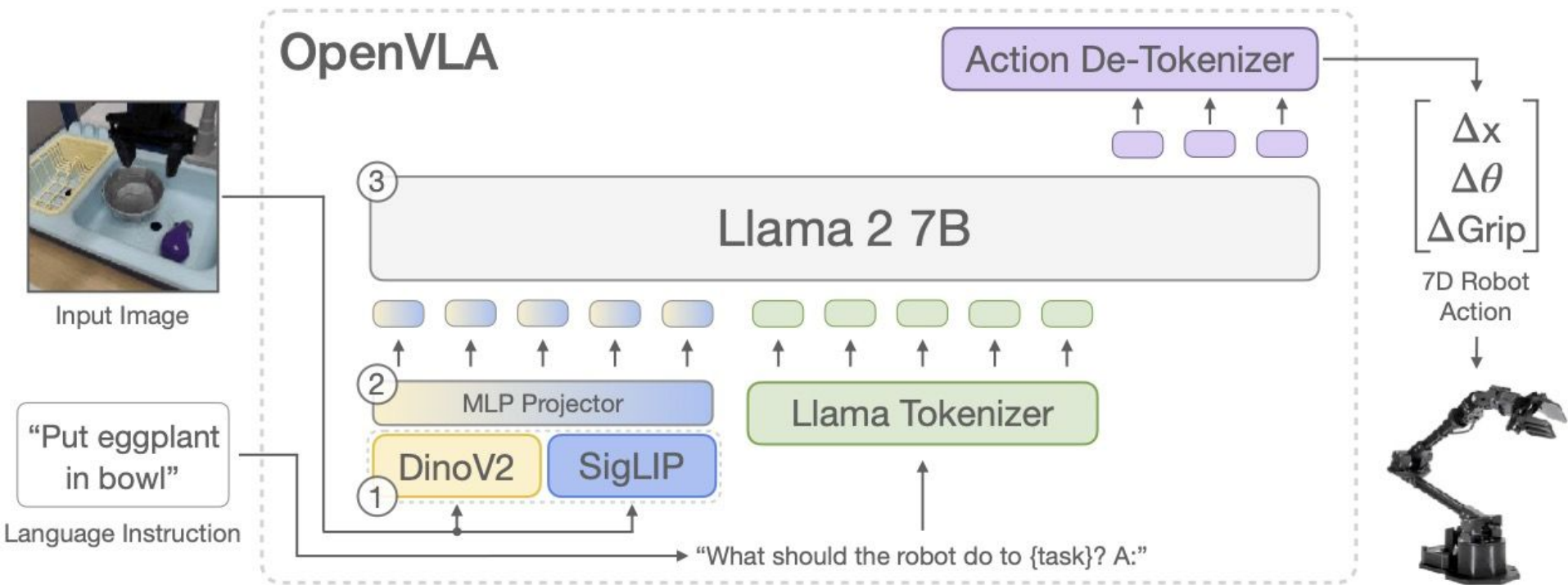
Stanford / UC Berkeley

Открытая платформа для исследований и прототипирования

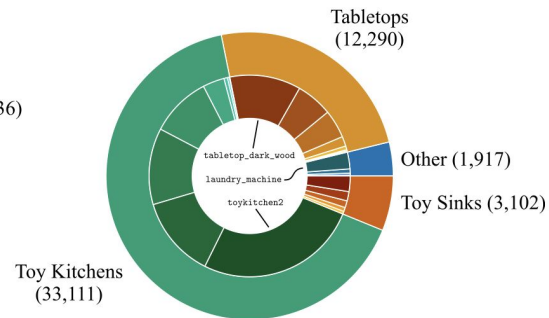
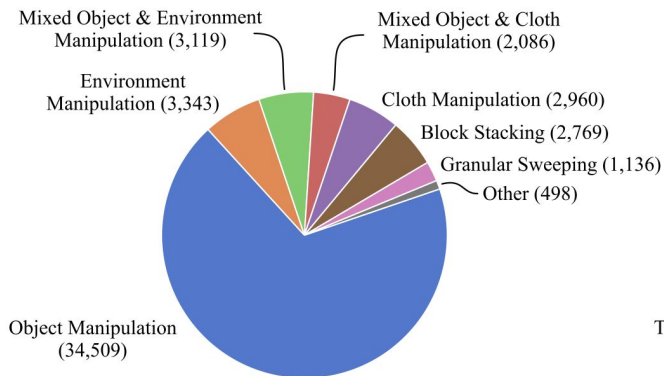
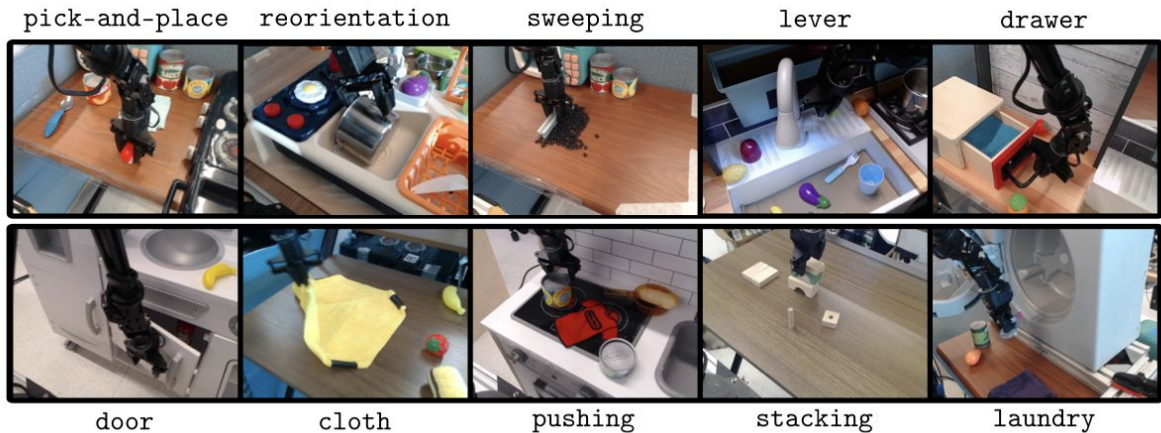
- Полностью открытый код и веса
- Простой fine-tuning за 1 день
- Модульная архитектура



OpenVLA



BridgeData V2



Open X-Embodiment

QT-Opt
pick anything

TOTO
pour

sweep the green cloth to the left side of the table

Push T

stack cups

place the black bowl in the dish rack

Jaco Play

ALOHA

Taco Play

1M Episodes from **311 Scenes**
34 Research Labs across **21 Institutions**

22 Embodiments

527 Skills

pour stack route

60 Datasets

1,798 Attributes • 5,228 Objects • 23,486 Spatial Relations

Cable Routing

pick green chip bag from counter

RT-1

set the bowl to the right side of the table

Bridge

Door Opening

3 кита современных VLA



п0 / п0.5

Physical Intelligence

Масштабируемая архитектура с единым токеном пространством для языка и действий

- Более 10M эпизодов обучения
- Поддержка 3D-восприятия
- Впечатляющее zero-shot обобщение



OpenVLA

Stanford / UC Berkeley

Открытая платформа для исследований и прототипирования

- Полностью открытый код и веса
- Простой fine-tuning за 1 день
- Модульная архитектура



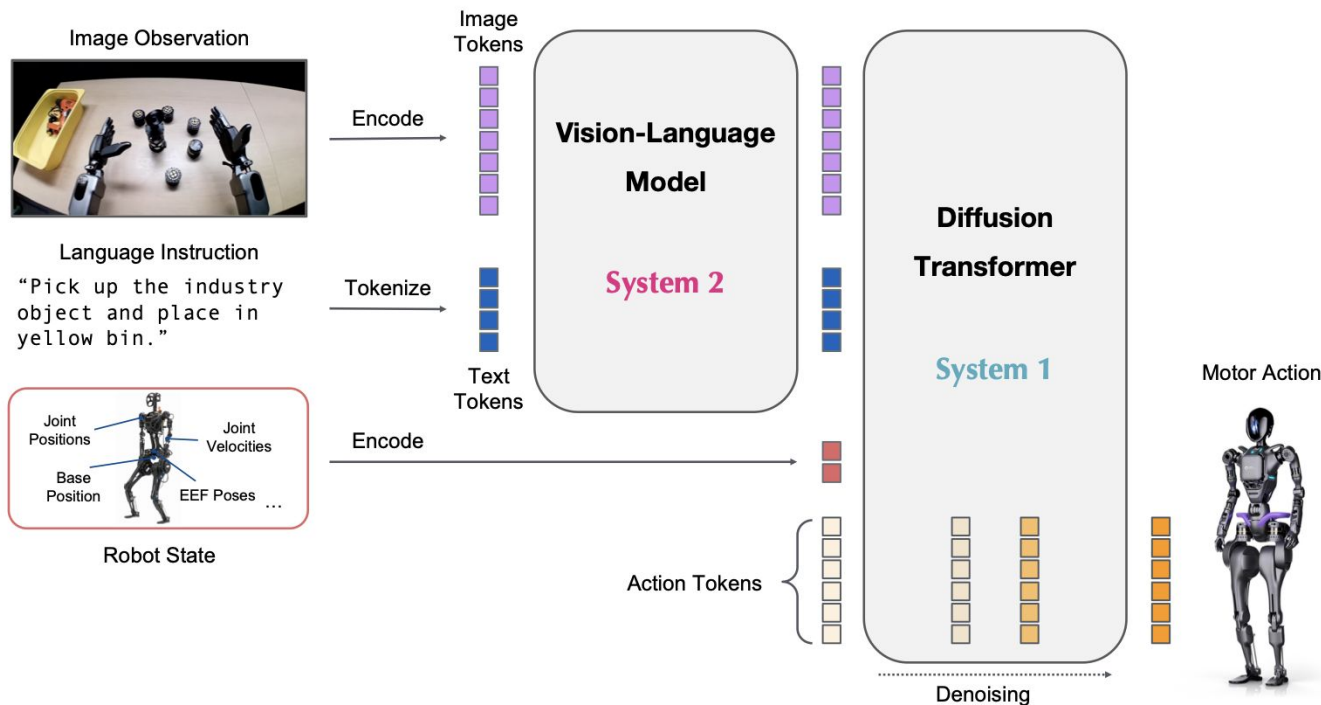
GROOT

NVIDIA

Промышленное решение для реального мира

- Задержка менее 200 мс
- Мультимодальное взаимодействие
- Интеграция с Isaac Sim

GROOT N1: An Open Foundation Model for Generalist Humanoid Robots



Двуручные манипуляции роботом Fourier GR-1



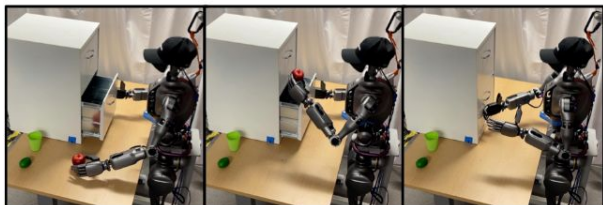
Pick-and-Place: Tray to Plate



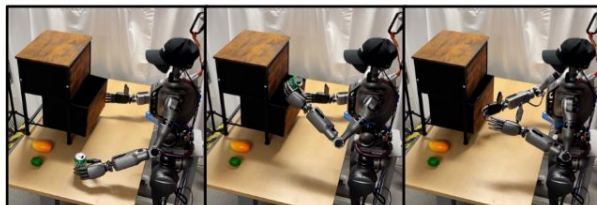
Pick-and-Place: Placemat to Basket



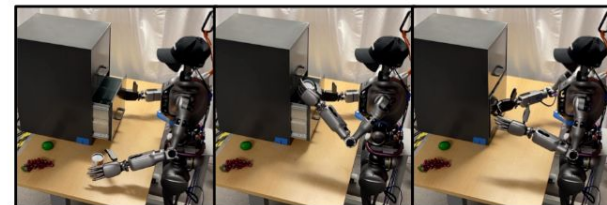
Pick-and-Place: Cutting Board to Pan



Articulated: White Drawer



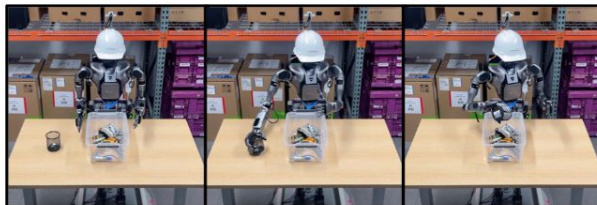
Articulated: Wooden Chest



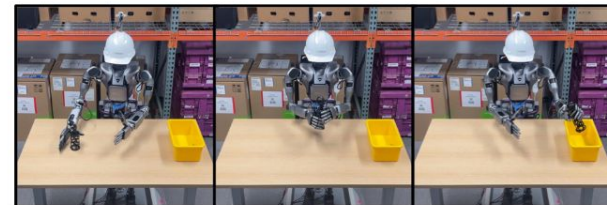
Articulated: Dark Cabinet



Industrial: Machinery Packing



Industrial: Mesh Cup Pouring



Industrial : Cylinder Handover

Teleoperation Data Collection

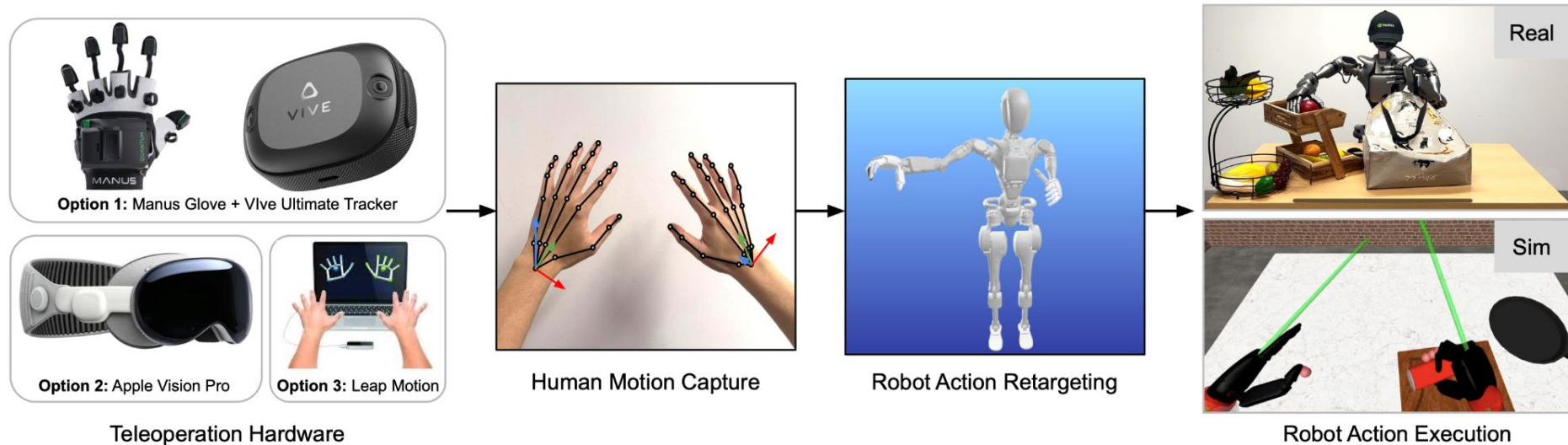


Figure 6: Data Collection via Teleoperation. Our teleoperation infrastructure supports multiple devices to capture human hand motion, including 6-DoF wrist poses and hand skeletons. Robot actions are produced through retargeting and executed on robots in real and simulation environments.

Итоги

VLA – foundation модели для робототехники:

- Объединяют **понимание сцены и задач** с моторикой робота.
- Способны к **многозадачности и обобщению** на новые объекты, среды и роботы.
- Поддержка **разных типов роботов**: манипуляторы, мобильные платформы, гуманоиды.

1

2024: Прорыв

Появление $\pi 0$, OpenVLA и GROOT — три подхода к одной цели

2

2025: Масштабирование

Тысячи лабораторий начинают эксперименты с VLA

3

2026-2030: Внедрение

VLA-роботы появляются на складах, в больницах, домах

4

Будущее

Универсальные роботизированные ассистенты становятся реальностью

Вопросы?

