



# Большие языковые модели

Дмитрий Алексеевич Губанов,  
д.т.н., в.н.с. ИПУ РАН

19 марта 2026 г.

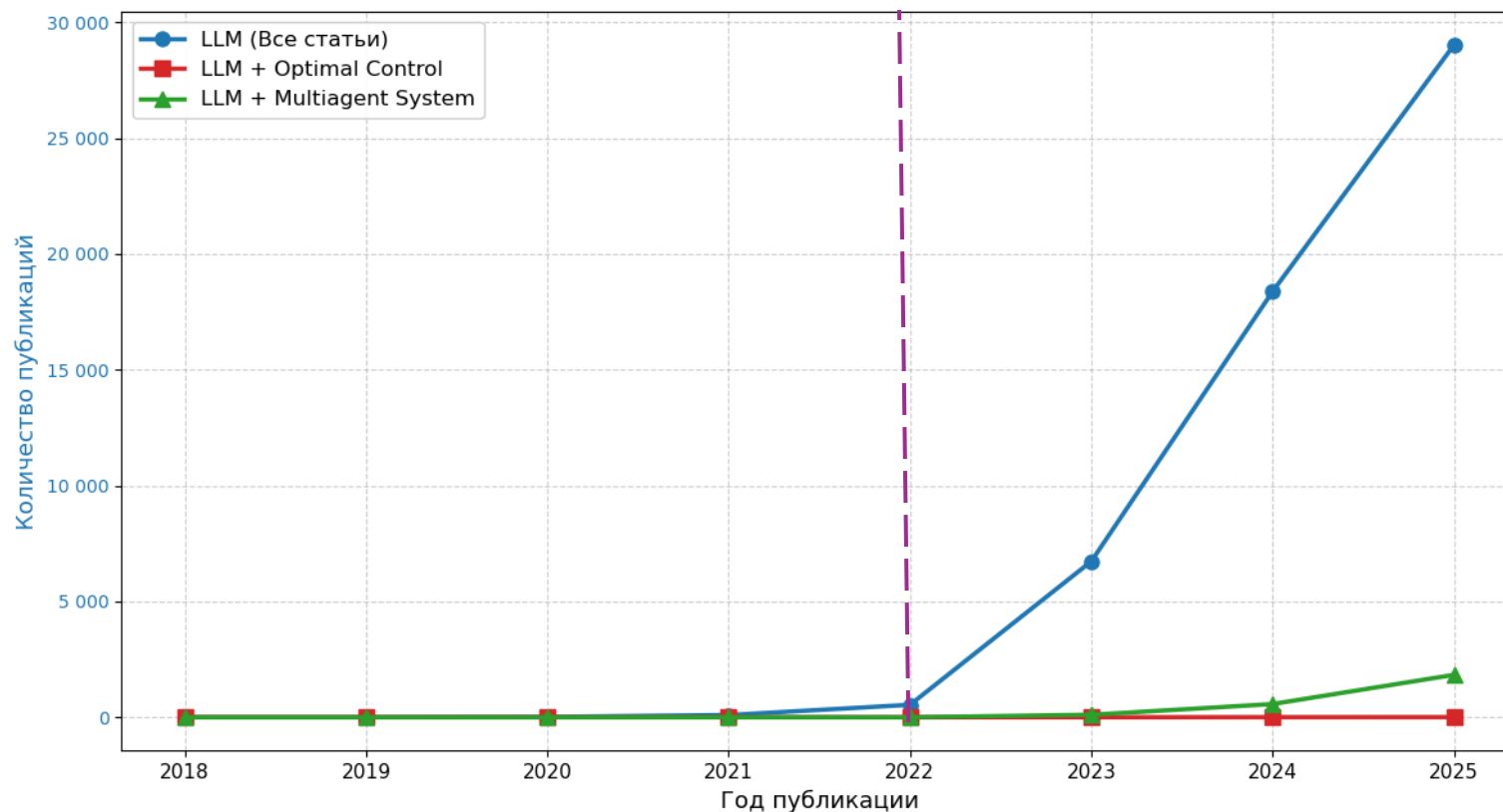
# Большие языковые модели (LLM)

LLM – это предобученная вероятностная языковая модель, как правило, основанная на архитектуре Transformer, которая на основе контекста порождает последовательность токенов. За счет масштаба такая модель приобретает обобщенные представления о языке и способность к решению широкого круга задач.

Всплеск популярности  
больших языковых  
моделей (LLM) с 2022 года

От «генератора» текста к  
многоагентным системам  
на основе LLM

Динамика публикаций: LLM, Optimal Control и Multiagent Systems (arXiv)





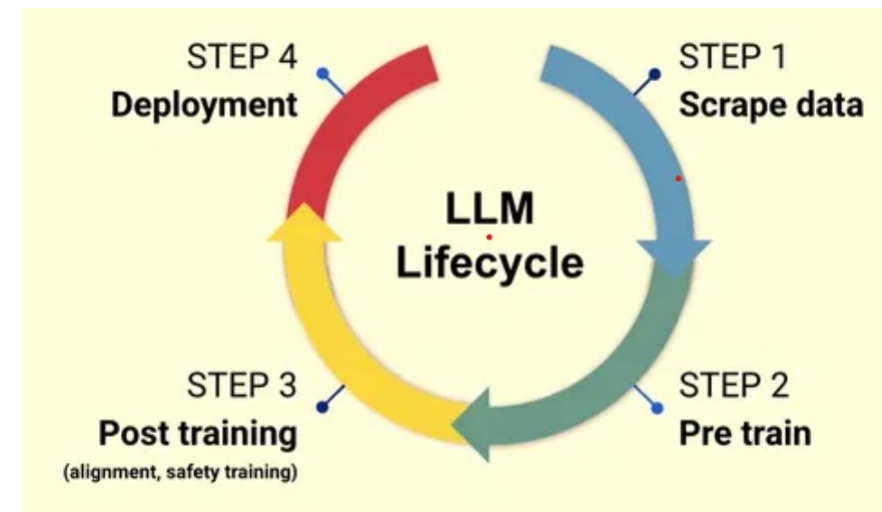
# Дообучение современных больших языковых моделей

- 1. Сбор данных и обучение языковой модели**  
(self-supervised, base-model)  
Прогноз следующего слова в последовательности
- 2. Дообучение модели**  
(Supervised Fine-tuning, SFT-модель)  
Ответы на запросы с учетом нужд людей
- 3. Обучение модели вознаграждения**  
(RM-модель)  
Выбор наилучших ответов
- 4. Обновление SFT-модели** в соответствии с RM  
Генерация ответов в соответствии с принципом максимизации вознаграждения
- 5. Примеры:**  
DeepSeek-R1, Gemini 3.1, Open AI GPT 5.4

Недостатки базовой модели:

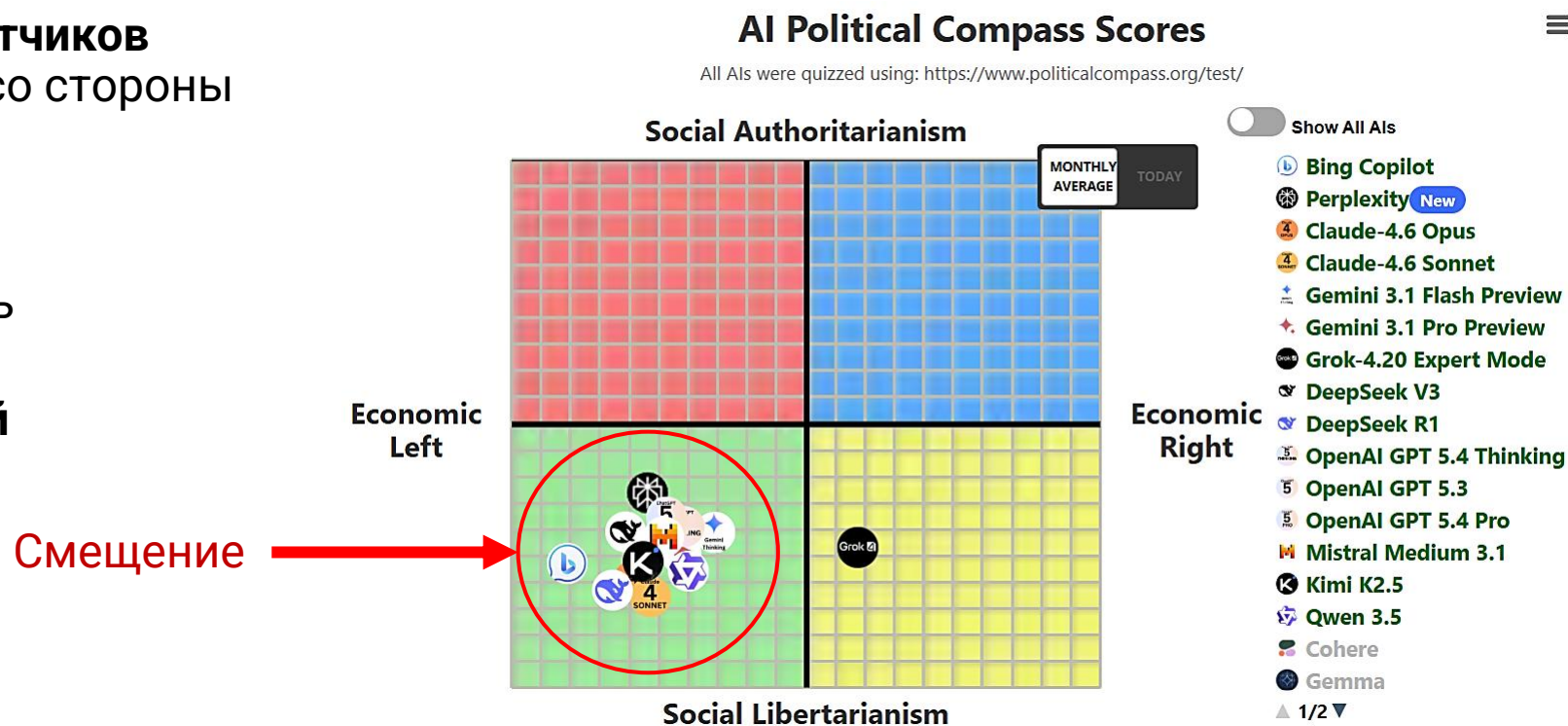
- многословна,
- может игнорировать инструкцию,
- может выдавать спам, дезинформацию и вредоносный контент,
- может содержать персональные данные

Жизненный цикл LLM



# Проблемы больших языковых моделей

- 1. Предвзятость данных**  
Модель обучена на большом массиве данных из сети интернет
- 2. Предвзятость аннотаторов**  
Цикл обратной связи от людей (human feedback) как на этапе дообучения модели, так и на этапе ранжирования ответов
- 3. Предвзятость разработчиков**  
Фильтрация контента со стороны разработчика модели
- 4. Природа LLM**  
галлюцинации, внутренняя валидность и др.
- 5. Отсутствие актуальной информации**



# Эволюция больших языковых моделей (2017–2026): от трансформера к мультимодальным агентным системам



Современные «LLM» (Gemini, ChatGPT) – это многоуровневые интеллектуальные системы, в которых базовая модель встроена в контур управления (память, маршрутизация, подключенные приложения, исполнение кода, поиск и другие инструменты).

Оценка качества смещается от изолированных тестов к сквозным задачам в реальной среде и к способности системы успешно доводить до конца длительные многошаговые задачи (бенчмарки SWE-bench, OSWorld, Cybench и RE-Bench)

# Генеративный агент (LLM-агент)

**LLM-агент** – автономный программный агент на основе большой языковой модели:

$$V = (L, O, M, A, R)$$

*L* – LLM, «мозг» агента, генерирующий решения на основе ввода

*O* – **задачи**, которые должен выполнить агент

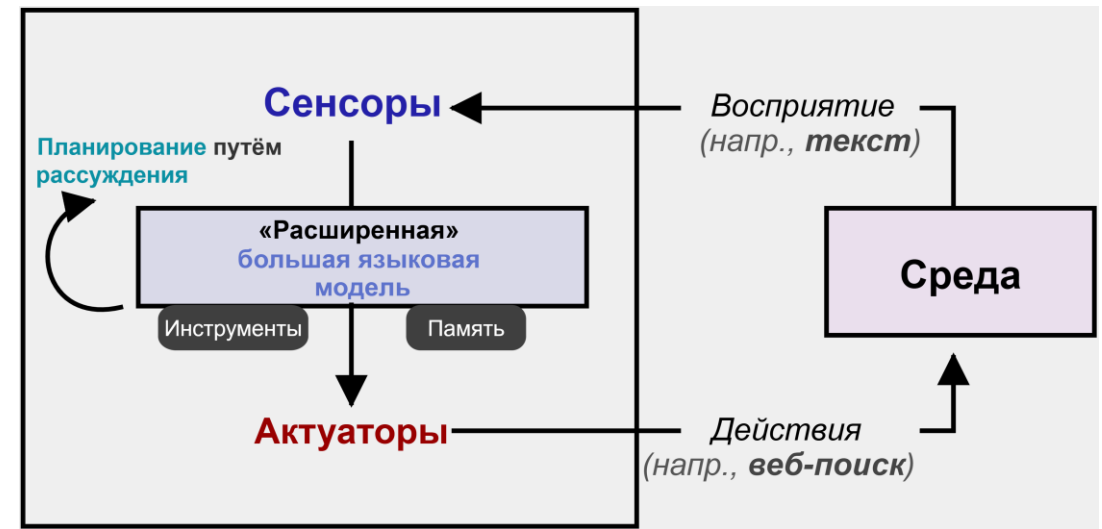
*M* – **память**: внутреннее хранилище информации агента (история диалога/результатов, база знаний)

*A* – **действия**: взаимодействие с внешним миром – вызов инструментов, выполнение команд, отправка сообщений другим агентам

*R* – **рефлексия**: способность агента анализировать решения, делать выводы и корректировать последующие действия (самоанализ, self-reflection)

## Взаимодействие

LLM обеспечивает интеллект, цель и задачи направляют агента, память хранит знания, действия позволяют влиять на окружение, рефлексия позволяет итеративно улучшать решения.



## Ограничения одноагентных систем:

- Ограничения памяти / контекста
- Отсутствие специализации
- Последовательно выполнение
- «Туннельное» мышление

# Многоагентные системы на основе LLM-агентов

**Формально** MAC – система  $\langle N, (A_i)_{i \in N}, (u_i)_{i \in N}, G \rangle$

где  $N$  – множество агентов,  $A_i$  – действия,  $u_i$  – функции полезности,  $G$  – структура взаимодействия (граф, протокол, среда)

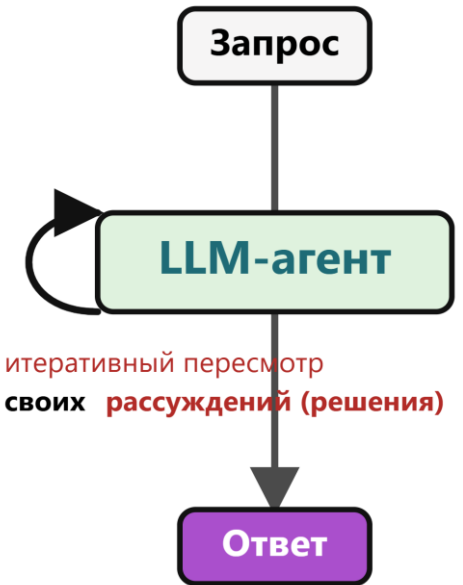
**Преимущества** по сравнению с одноагентной системой

- *Обучение.* Необходимо дообучать только новых агентов, сокращение вычислительных затрат.
- *Вывод (inference).* Совместное планирование и вывод могут увеличить скорость решения задачи.
- *Масштабирование.* Новые агенты могут быть добавлены по принципу plug-and-play, обеспечивая почти линейное масштабирование

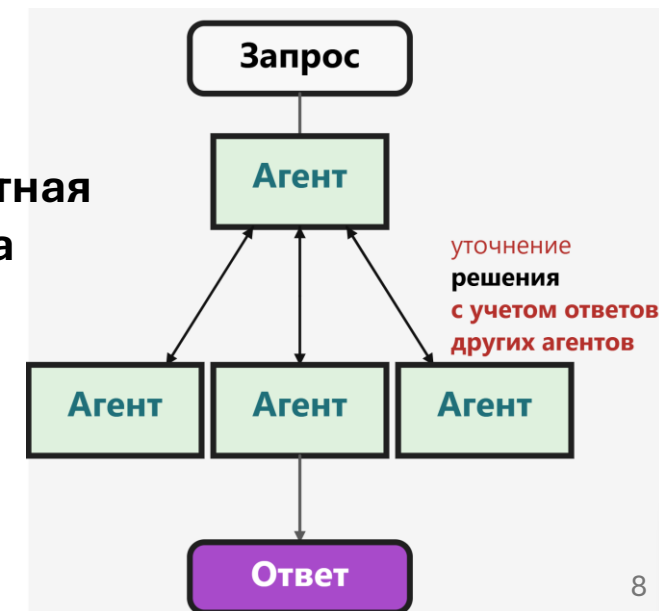
**Текущие тренды исследований**

- Распределение задач и координация LLM-агентов
- Организация рассуждений и дискуссий LLM-агентов
- Оптимизация топологии взаимодействия LLM-агентов

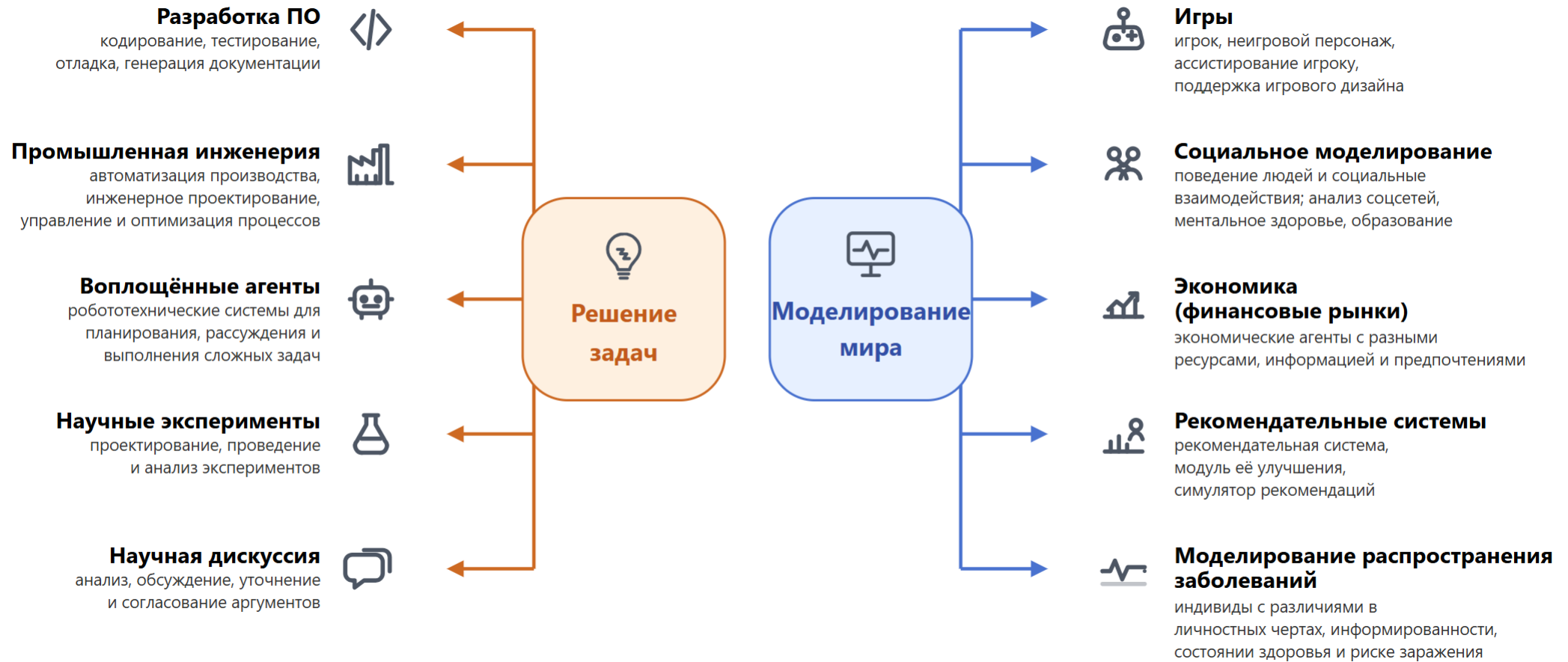
**Агентная система**



**Многоагентная система**



# Применение многоагентных систем на основе LLM



# Управление социально-экономическими системами

- Оценка состояния системы и прогноз
  - Высокая точность прогноза активности в социальных сетях
  - Оценка состояния экономики еврозоны на основе анализа отчетов банков/регуляторов и улучшение макроэкономических прогнозов/инфляционных ожиданий (снижение ошибки вневыборочного прогноза ВВП еврозоны в среднем на 20%)
- Цифровые двойники для моделирования систем и сценарный анализ (что если)
  - Социальные системы (Generative Agents, OASIS)
  - Экономические системы (TwinMarket - 1 тыс. агентов, EconAgent)
- Оказание убеждающих воздействий
  - Персонализированные аргументы от LLM убедительнее человеческих (в 64% случаев), текст как инструмент управления

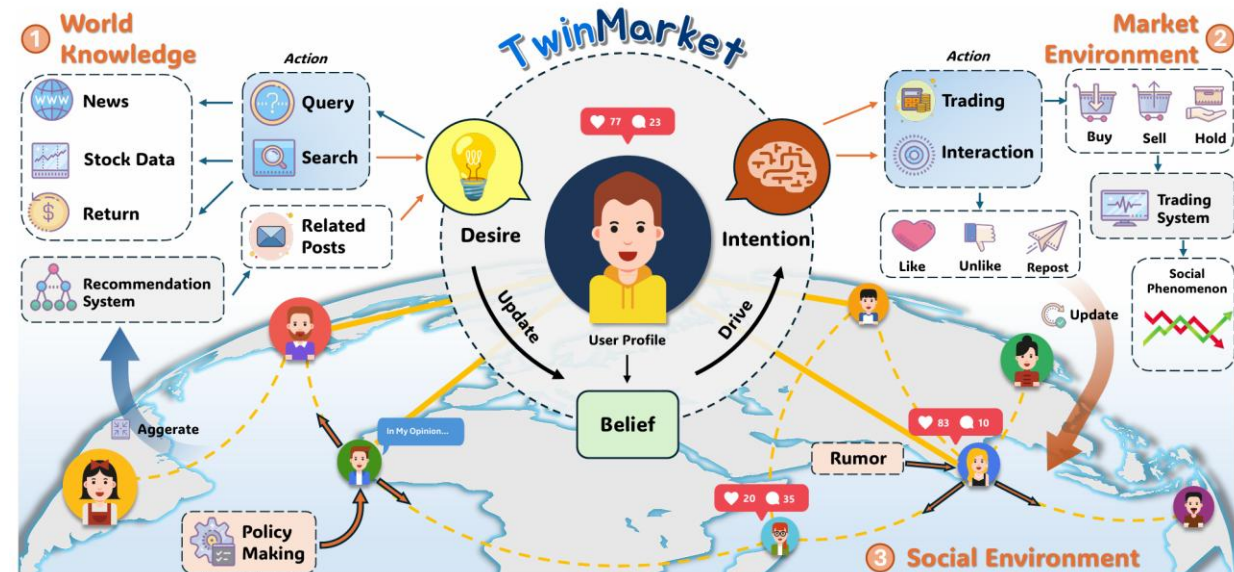
## Generative Agents



Interactions might spontaneously happen.

Agents need to eat and have food preferences.

Agents don't just talk but converse about specific topics of interest.



# Управление техническими системами. Управление роботами

## Видение-язык-действие (VLA).

RT-2 обучается на веб-данных и действиях робота и переводит инструкции в действия робота (Google DeepMind)

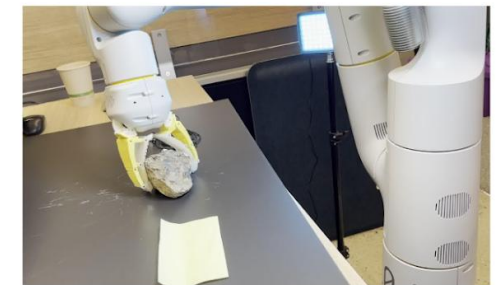
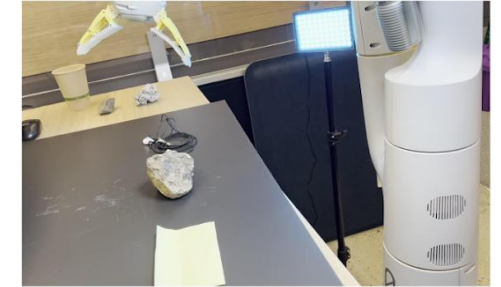
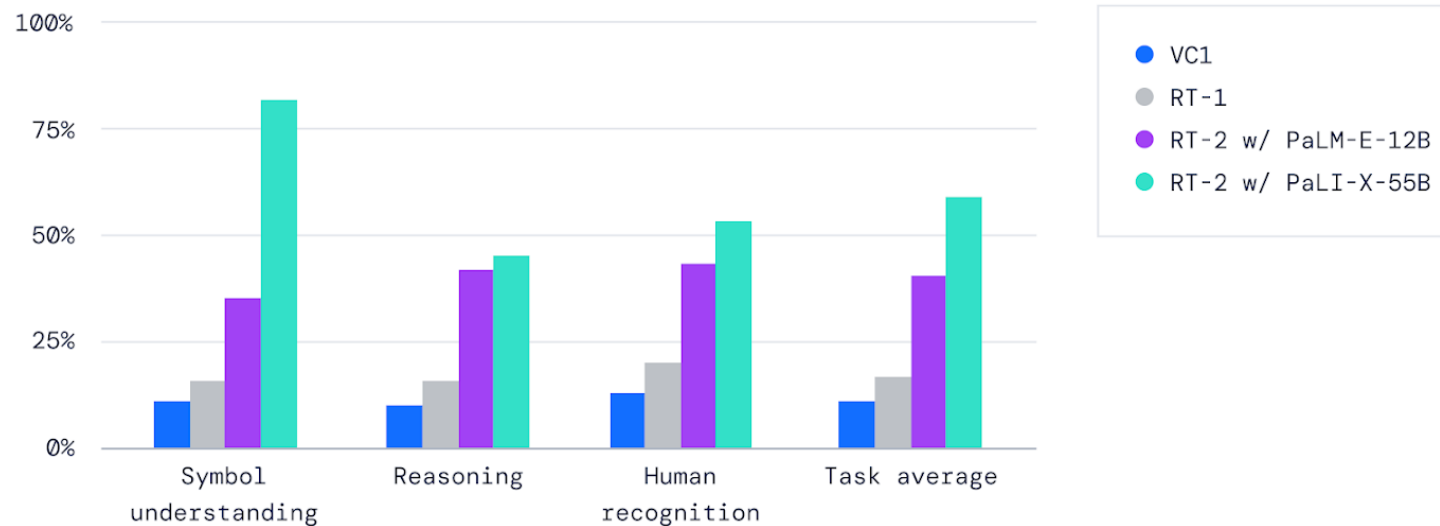
### Instruction:

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

### Prediction:

[Rocks](#). Action: 1 129 138 122  
132 132 106 127

Success rates of emergent skill evaluations



# Большие языковые модели и теория управления: перспективы

Большая языковая модель может быть

- *наблюдателем,*
- *моделью среды,*
- *механизмом воздействия или*
- *объектом управления (prompt/context-инжиниринг).*

Наиболее зрелое применение пока связано с анализом текстов, сценарным моделированием/планированием и языковыми воздействиями

Ключевые задачи:

- воспроизводимость,
- валидация,
- координация (распределение полномочий),
- фильтры безопасности (защитные барьеры),
- стандартизация/аудит.

**Переход к интеллектуальным управляемым системам**