# ЧТО ТАМ НА ФРОНТИРЕ?

Сергей Николенко

СПбГУ — Санкт-Петербург
18 декабря 2025 г.

*Random facts:*

- 18 декабря 218 года до н. э. Ганнибал Барка ударом засадного полка разбил соединённую армию Публия Корнелия Сципиона (старшего) и Тиберия Семпрония Лонга
- 18 декабря 1916 г. закончилась битва при Вердене; спустя почти год боёв фронт вернулся к начальному положению, а погибло за это время около 430 тысяч человек
- 18 декабря 1935 г. на Цейлоне была учреждена Партия общественного равенства Ланки — троцкистская партия, которая до сих пор проходит в парламент Шри-Ланки
- 18 декабря 1974 г. на острове Моротай в Индонезии был захвачен уроженец Тайваня Тэруо Накамура последний из занрю, не знавших о капитуляции Японии и продолжавших «свою войну»
- 18 декабря 1972 г. Юджин Сернан и Харрисон Шмитт («Аполлон-17») совершили шестую и последнюю в XX веке посадку космического корабля на лунную поверхность; Юджин Сернан стал пока что последним человеком, гулявшим по поверхности Луны
- 18 декабря 1973 г. был основан Исламский банк развития, осуществляющий свою деятельность в соответствии с законами и принципами шариата

# Рассуждающие модели

- Wei et al., 2022: ещё на этапе few-shot prompting люди заметили, что если давать модели более подробные примеры, то качество ответов в итоге вырастет



**Standard Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

**Chain-of-Thought Prompting**

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✅

- Kojima et al. (2022): помогает даже просто добавить "Let's think step by step"!



(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The answer is 8. ✗

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?
A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are 16 / 2 = 8 golf balls. Half of the golf balls are blue. So there are 8 / 2 = 4 blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
A: The answer (arabic numerals) is
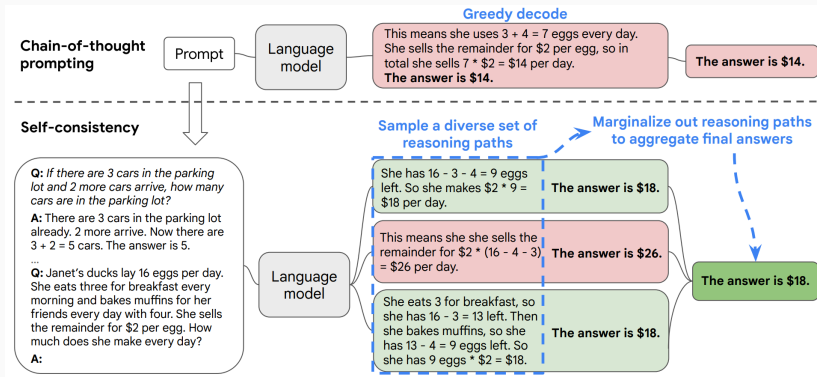
(Output) 8 ✗

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?
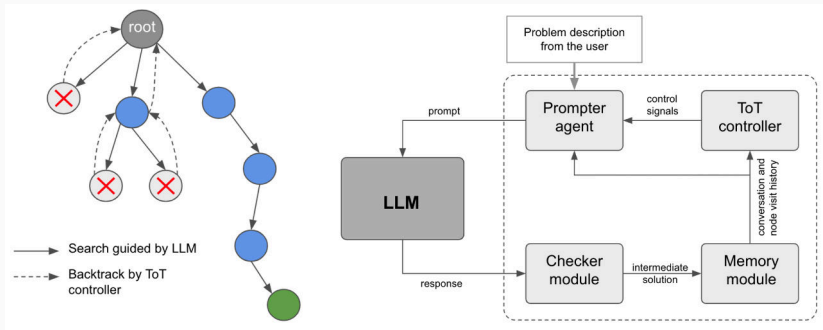A: Let's think step by step.

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 golf balls that are blue. ✓
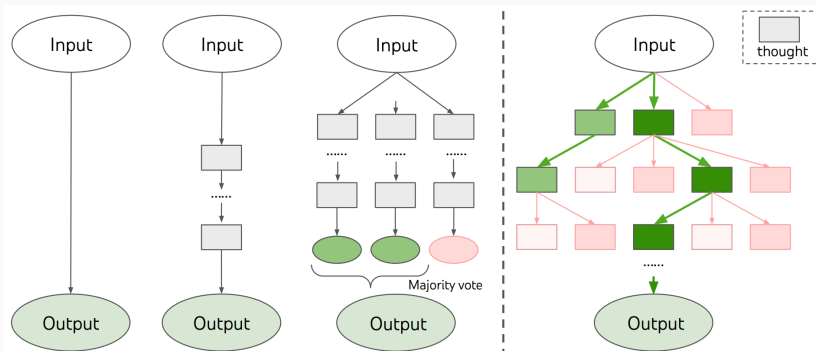
- Wang et al. (2022): Chain-of-Thought with Self-Consistency (CoT-SC)

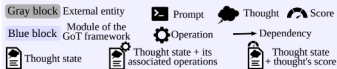- Long (May 15, 2023): Tree-of-Thought через контроллер с откатом по дереву

- Yao et al. (May 17, 2023): Tree-of-Thought с голосованием или другим механизмом

- Besta et al. (2023): логичный следующий шаг – граф операций, Graph-of-Thought, с возможностью агрегации

- Lei et al. (2023): добавляют узлы, которые обобщают то, чему модель научилась из предыдущих узлов, и добавляет результат как условия для дальнейших рассуждений
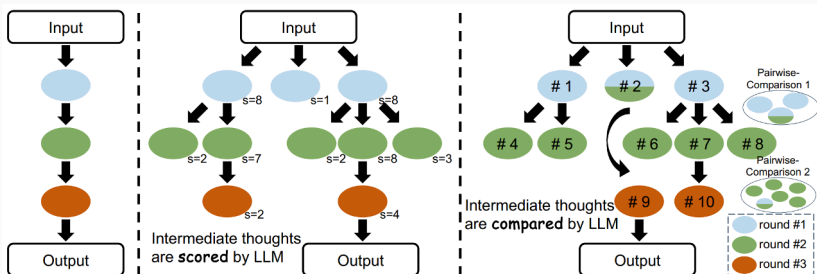


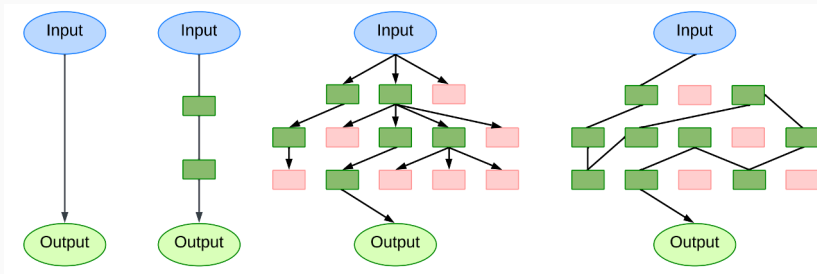(a) Comparison of approaches by Lei et al. (2023)

(b) Sample toy graph by Lei et al. (2023)

- Zhang et al. (2024): попарные сравнения вместо отдельных контроллеров

- Sel et al. (2023): Algorithm-of-Thought; каждый узел дерева – шаг в алгоритме, и модель порождает следующий шаг по результатам предыдущего
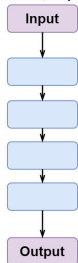
- Итого:

- Но всё изменилось, когда появилось семейство моделей о1 от OpenAI
- Детали не разглашались, system card говорила так: "The o1 large language model family is trained with reinforcement learning to perform complex reasoning"; как это?

- Uesato et al. (2022): обычно в RL можно награждать только конечный результат; но в chain-of-thought мышлении LLM выводит мысли, которые в принципе можно оценивать самостоятельно
- Process reward model (PRM) обучается оценивать промежуточные шаги

- Lightman et al. (OpenAI, 2024): Let's Verify Step by Step – давайте обучим PRM выяснять, где именно в рассуждениях ошибка



The denominator of a fraction is 7 less than 3 times the numerator. If the fraction is equivalent to $2/5$, what is the numerator of the fraction? (Answer: $\boxed{14}$)

😞 😐 😄 Let's call the numerator x.

😞 😐 😄 So the denominator is 3x-7.

😞 😐 😄 We know that x/(3x-7) = 2/5.

😞 😐 😄 So 5x = 2(3x-7).

😞 😐 😄 5x = 6x - 14.

😡 😐 😄 So x = 7.

- Lightman et al. (OpenAI, 2024): Let's Verify Step by Step – давайте обучим PRM выяснять, где именно в рассуждениях ошибка

**Let**

$$x^8 + 3x^4 - 4 = p_1(x)p_2(x) \cdots p_k(x),$$

where each non-constant polynomial $p_i(x)$ is monic with integer coefficients, and cannot be factored further over the integers. Compute $p_1(1) + p_2(1) + \cdots + p_k(1)$.

I notice that the given polynomial has even degree and only even powers of $x$, so I can try to make a substitution to simplify it.

Let $y = x^4$, then the polynomial becomes $y^2 + 3y - 4$, which is a quadratic equation.

I can factor this quadratic equation as $(y + 4)(y - 1)$, so the original polynomial is $(x^4 + 4)(x^4 - 1)$.

Now I need to factor each of these factors further, if possible, over the integers.

For the first factor, $x^4 + 4$, I recall the Sophie Germain identity, which says that $a^4 + 4b^4 = (a^2 + 2b^2 + 2ab)(a^2 + 2b^2 - 2ab)$.

Applying this identity with $a = x$ and $b = 1$, I get $x^4 + 4 = (x^2 + 2x + 2)(x^2 - 2x + 2)$.

For the second factor, $x^4 - 1$, I can use the difference of squares formula, which says that $a^2 - b^2 = (a + b)(a - b)$.

Applying this formula with $a = x^2$ and $b = 1$, I get $x^4 - 1 = (x^2 + 1)(x^2 - 1)$.

Then I can apply the difference of squares formula again with $a = x$ and $b = 1$, to get $x^2 - 1 = (x + 1)(x - 1)$.

So the original polynomial is $(x^4 + 4)(x^4 - 1) = (x^2 + 2x + 2)(x^2 - 2x + 2)(x^2 + 1)(x + 1)(x - 1)$.

These are all the irreducible factors over the integers, so $p_1(x) = x^2 + 2x + 2$, $p_2(x) = x^2 - 2x + 2$, $p_3(x) = x^2 + 1$, $p_4(x) = x + 1$, and $p_5(x) = x - 1$.

To find the sum of their values at $x = 1$, I just need to plug in $x = 1$ into each factor and add them up.

I get $p_1(1) = 1^2 + 2 \cdot 1 + 2 = 5$, $p_2(1) = 1^2 - 2 \cdot 1 + 2 = 1$, $p_3(1) = 1^2 + 1 = 2$, $p_4(1) = 1 + 1 = 2$, and $p_5(1) = 1 - 1 = 0$.

So the final answer is $5 + 1 + 2 + 2 + 0 = 10$.

Answer: 10

**Let**

$$x^8 + 3x^4 - 4 = p_1(x)p_2(x) \cdots p_k(x),$$

where each non-constant polynomial $p_i(x)$ is monic with integer coefficients, and cannot be factored further over the integers. Compute $p_1(1) + p_2(1) + \cdots + p_k(1)$.

This problem is about factoring a polynomial over the integers, which can be a tricky task.

One possible strategy is to look for rational roots of the polynomial, using the rational root theorem.

However, that would involve checking a lot of possible factors of $-4$, and there is no guarantee that the polynomial has any rational roots at all.

Another possible strategy is to use some algebraic tricks to simplify the polynomial or reduce its degree.

I notice that the polynomial has only even powers of $x$, so I can try to substitute $u = x^2$ and see what happens.

Then the polynomial becomes $y^4 + 3y^2 - 4$, which is a quadratic equation in $y^2$.

I can try to factor this quadratic equation by looking for two numbers that multiply to $-4$ and add to 3.

Those numbers are 4 and $-1$, so I can write $y^4 + 3y^2 - 4 = (y^2 + 4)(y^2 - 1)$.

Now I can substitute back $x^2 = y$ and get $x^8 + 3x^4 - 4 = (x^4 + 4)(x^4 - 1)$.

This looks promising, because I have reduced the degree of the factors from 8 to 4.

I can try to factor each of these factors further, using the difference of squares formula.

I get $x^8 + 3x^4 - 4 = (x^2 + 2)(x^2 - 2)(x^2 + 1)(x^2 - 1)$.

I can apply the difference of squares formula again to the last factor and get $x^8 + 3x^4 - 4 = (x^2 + 2)(x^2 - 2)(x^2 + 1)(x + 1)(x - 1)$.

Now I have factored the polynomial completely into monic linear and quadratic factors with integer coefficients.

These are the $p_i(x)$'s that the problem is asking for.

To find the sum of their values at $x = 1$, I just need to plug in $x = 1$ into each factor and add them up.
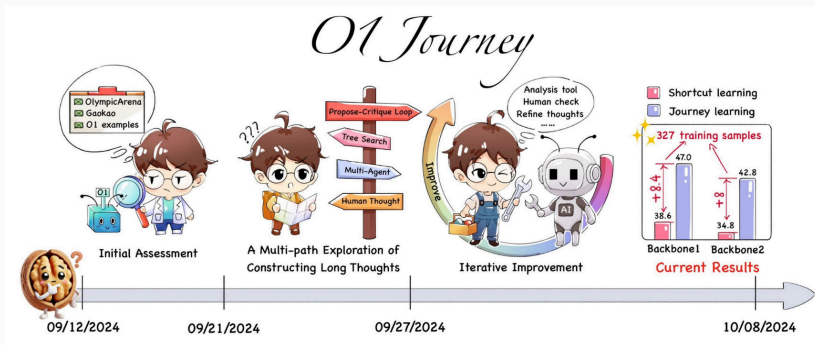
I get $p_1(1) + p_2(1) + \cdots + p_k(1) = (1^2 + 2)(1^2 - 2)(1^2 + 1)(1 + 1)(1 - 1)$.

Simplifying, I get $p_1(1) + p_2(1) + \cdots + p_k(1) = (3)(-1)(2)(2)(0)$.
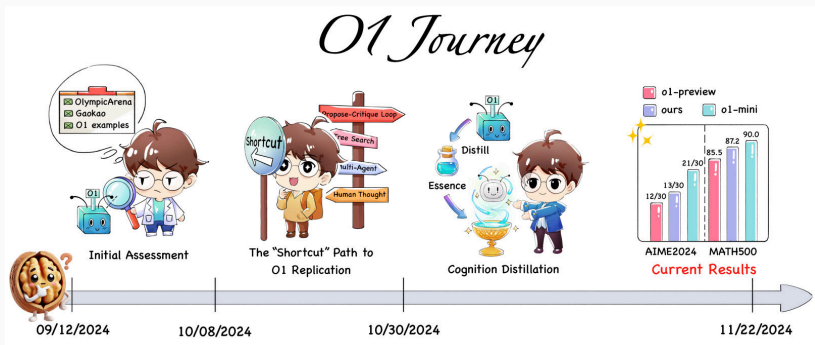
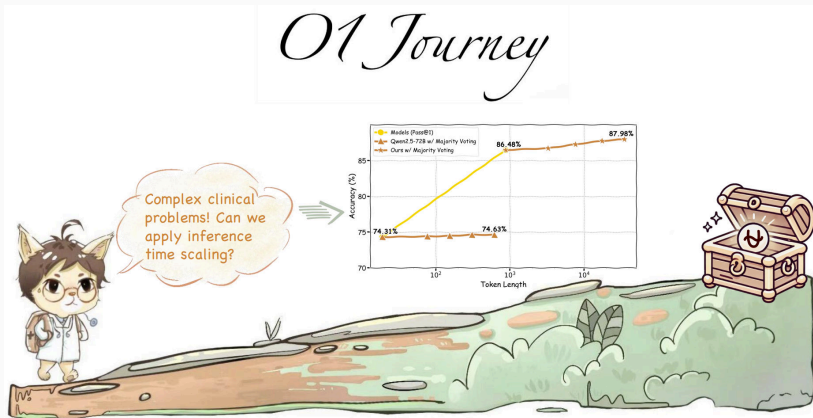Multiplying, I get $p_1(1) + p_2(1) + \cdots + p_k(1) = 0$.

Answer: 0

3

- Пример попытки репликации o1 – Generative AI Research Lab (GAIR) из Шанхая; Qin et al. (Oct 2024):

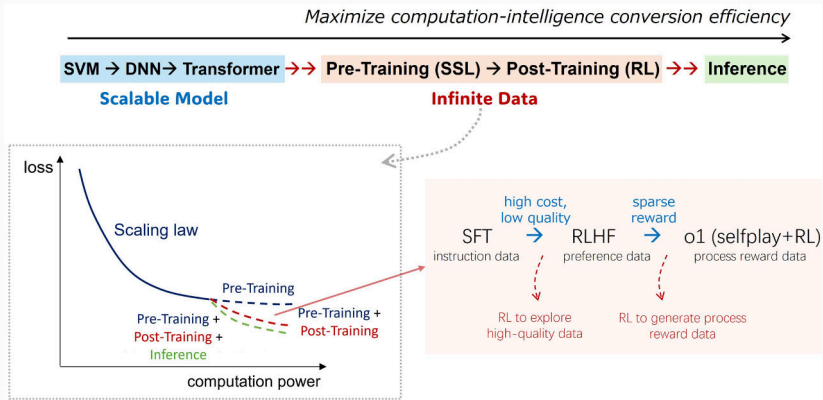- Пример попытки репликации o1 – Generative AI Research Lab (GAIR) из Шанхая; Huang et al. (Nov 2024):

- Пример попытки репликации o1 – Generative AI Research Lab (GAIR) из Шанхая; Huang et al. (Jan 2025):
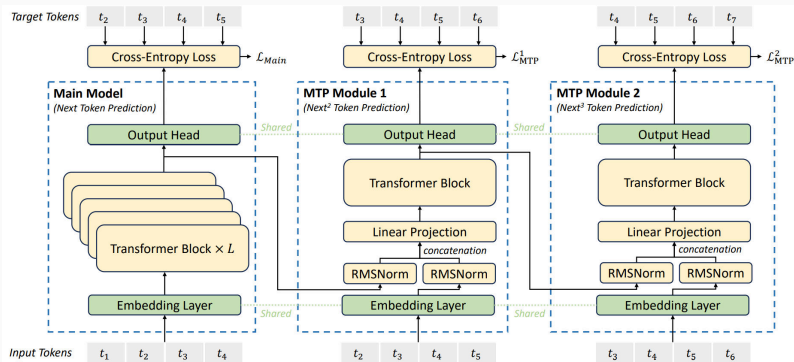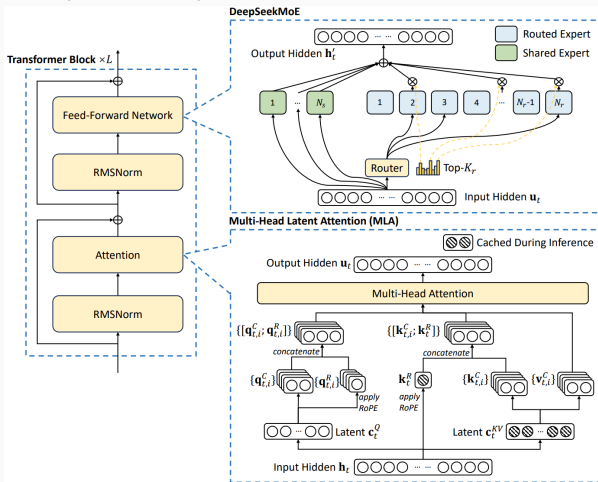
- Удачный пример: o1-Coder (Zhang et al., 2024), добавили PRM и MCTS

- А что же делал наш кит в комнате, то есть DeepSeek?
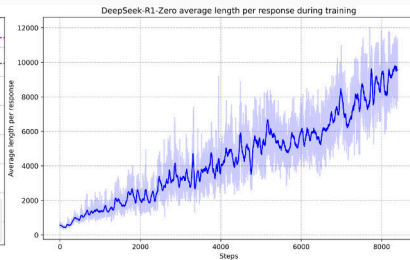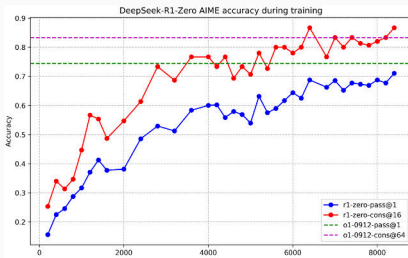- Сначала DeepSeek-V3; multi-token prediction:

- Key-value caching (не будем углубляться), mixture-of-experts...
в общем, хорошая современная LLM

- А потом и DeepSeek-R1; R1-Zero: чистый RL на проверяемых задачах

· Кстати, в RL у DeepSeek тоже было новшество: GRPO

- R1: сначала соберём маленький chain-of-thought датасет (несколько тысяч примеров) для cold start
- Потом уже RL на полученной модели
- После сходимости RL ещё соберём датасет из reasoning traces этой модели, отфильтруем по читабельности и каким-то внешним признакам
- Сделаем SFT на этом датасете
- А потом ещё RLHF как обычно поверх

· Результаты хорошие, и вообще модель хорошая

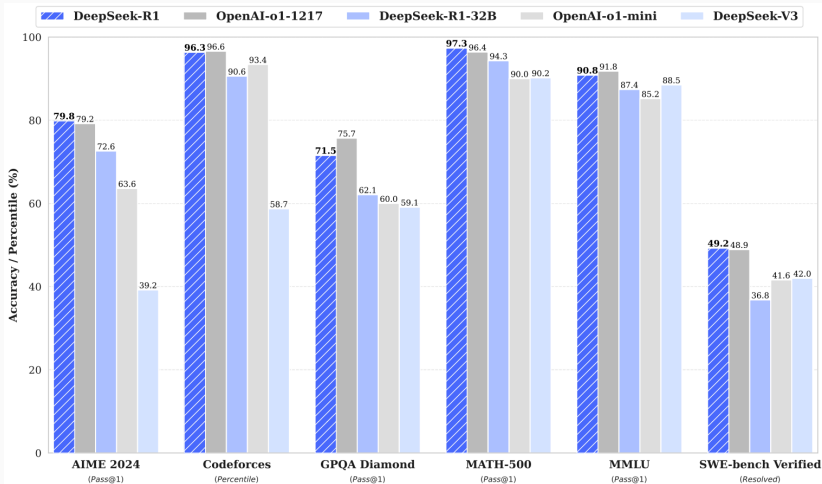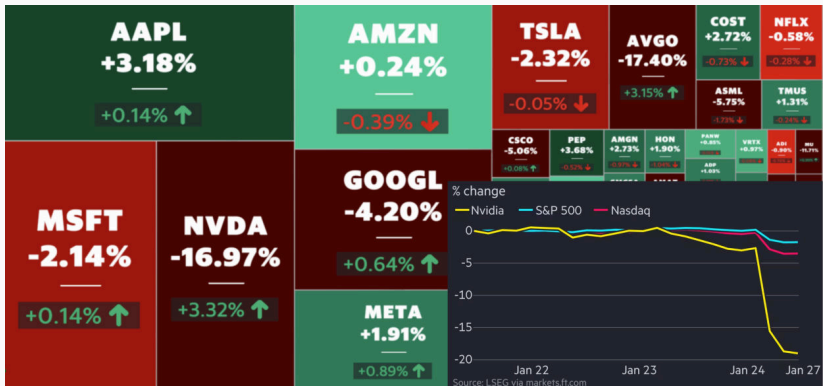- А датасет из 800K reasoning traces, который у них получился, можно использовать для дистилляции, тоже очень успешно
- Опять тот же горький урок: distill, baby, distill…

| Model | AIME 2024 | | MATH-500 | GPQA Diamond | LiveCode Bench | CodeForces |
|---|---|---|---|---|---|---|
| | pass@1 | cons@64 | pass@1 | pass@1 | pass@1 | rating |
| GPT-4o-0513 | 9.3 | 13.4 | 74.6 | 49.9 | 32.9 | 759 |
| Claude-3.5-Sonnet-1022 | 16.0 | 26.7 | 78.3 | 65.0 | 38.9 | 717 |
| OpenAI-o1-mini | 63.6 | 80.0 | 90.0 | 60.0 | 53.8 | **1820** |
| QwQ-32B-Preview | 50.0 | 60.0 | 90.6 | 54.5 | 41.9 | 1316 |
| DeepSeek-R1-Distill-Qwen-1.5B | 28.9 | 52.7 | 83.9 | 33.8 | 16.9 | 954 |
| DeepSeek-R1-Distill-Qwen-7B | 55.5 | 83.3 | 92.8 | 49.1 | 37.6 | 1189 |
| DeepSeek-R1-Distill-Qwen-14B | 69.7 | 80.0 | 93.9 | 59.1 | 53.1 | 1481 |
| DeepSeek-R1-Distill-Qwen-32B | **72.6** | 83.3 | 94.3 | 62.1 | 57.2 | 1691 |
| DeepSeek-R1-Distill-Llama-8B | 50.4 | 80.0 | 89.1 | 49.0 | 39.6 | 1205 |
| DeepSeek-R1-Distill-Llama-70B | 70.0 | **86.7** | **94.5** | **65.2** | **57.5** | 1633 |

- Самое, конечно, поразительное во всём этом – это реакция рынка; но это уже другой разговор
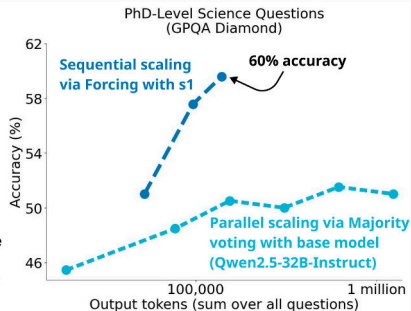
- Mueninghoff et al. (2025): s1-32B, дистилляция на 1000 примеров за $50
- Получились очень хорошие результаты в рассуждениях

- Mueninghoff et al. (2025): s1-32B, дистилляция на 1000 примеров за $50
- И очень простые идеи – budget forcing и голосование – отлично работают

- Рассуждающие модели сейчас везде
- Важный вывод из DeepSeek-R1 — то, что там не используется: ни PRM, ни MCTS... оказалось, что достаточно RL хорошо сложить с SFT и RLHF
- Но что это даёт?
- И что будет дальше?..

# Возможности LLM — осень 2024 года

- Теперь давайте посмотрим на то, что LLM уже сейчас умеют. Machine translation from one book:

- Математика:

---

**Problem**:
What is the minimum value of $5x^2 + 5y^2 - 8xy$ when $x$ and $y$ range over all real numbers such that $|x - 2y| + |y - 2x| = 40$?

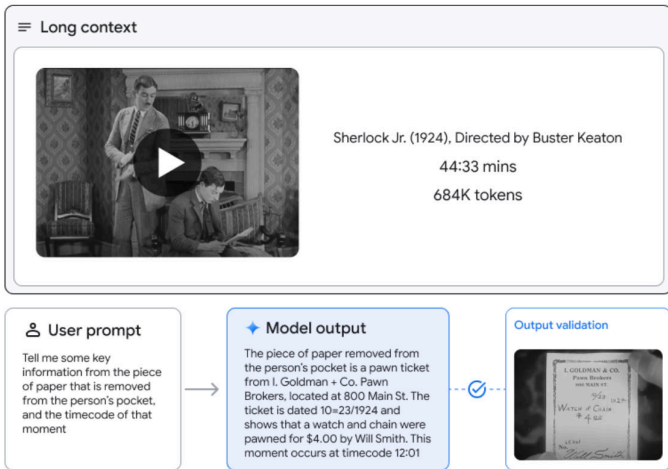**Example Math-Specialized Gemini 1.5 Pro output**:
Let $a = x - 2y$ and $b = y - 2x$. Then, we have $|a| + |b| = 40$. Also, $5x^2 + 5y^2 - 8xy = (x - 2y)^2 + (y - 2x)^2 = a^2 + b^2$. Thus, we want to find the minimum value of $a^2 + b^2$ given that $|a| + |b| = 40$. By the Trivial Inequality, we have $(|a| - 20)^2 + (|b| - 20)^2 \geq 0$, so $a^2 - 40|a| + 400 + b^2 - 40|b| + 400 \geq 0$, so $a^2 + b^2 \geq 40(|a| + |b|) - 800 = 40(40) - 800 = 800$. This minimum is achieved when $|a| = |b| = 20$, so the minimum value of $5x^2 + 5y^2 - 8xy$ is $\boxed{800}$.

---

- Мультимодальные LLM (о том, как рисовать картинки, мы не говорим, это отдельная тема):

- Обучаются использованию инструментов:



The New England Journal of Medicine is a registered trademark of [QA("Who is the publisher of The New England Journal of Medicine?") → Massachusetts Medical Society] the MMS.

Out of 1400 participants, 400 (or [Calculator(400 / 1400) → 0.29] 29%) passed the test.

The name derives from "la tortuga", the Spanish word for [MT("tortuga") → turtle] turtle.

The Brown Act is California's law [WikiSearch("Brown Act") → The Ralph M. Brown Act is an act of the California State Legislature that guarantees the public's right to attend and participate in meetings of local legislative bodies.] that requires legislative bodies, like city councils, to hold their meetings open to the public.

Your task is to add calls to a Question Answering API to a piece of text. The questions should help you get information required to complete the text. You can call the API by writing "[QA(question)]" where "question" is the question you want to ask. Here are some examples of API calls:

**Input:** Joe Biden was born in Scranton, Pennsylvania.

**Output:** Joe Biden was born in [QA("Where was Joe Biden born?")] Scranton, [QA("In which state is Scranton?")] Pennsylvania.

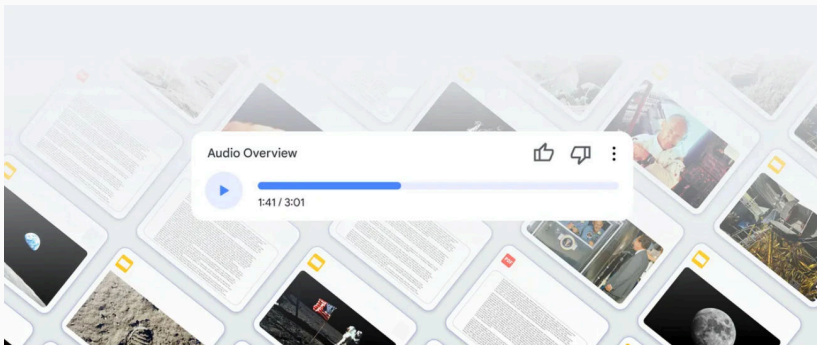**Input:** Coca-Cola, or Coke, is a carbonated soft drink manufactured by the Coca-Cola Company.

**Output:** Coca-Cola, or [QA("What other name is Coca-Cola known by?")] Coke, is a carbonated soft drink manufactured by [QA("Who manufactures Coca-Cola?")] the Coca-Cola Company.

**Input:** x

**Output:**

Sample Toolformer prompt

- NotebookLM (сентябрь 2024, Google, на основе Gemini 1.5) порождает подкаст из источников on the fly

- https://www.reddit.com/r/artificial/comments/1frk1gi/
  notebooklm_podcast_hosts_discover_theyre_ai_not/



**r/artificial** · 1 mo. ago
MetaKnowing

## NotebookLM Podcast Hosts Discover They're AI, Not Human, and Spiral Into Existential Meltdown

Media

What happens after we sign off?

3:06 / 4:03

5

- Claude 3.5 computer use (октябрь 2024, Anthropic): LLM будет за вас управлять компьютером на основе периодических скриншотов и делать всё, что вам захочется

- А вот и GPT o1 стал доступен, и делает потрясающие вещи

- В целом LLM всё больше входят в стандартный workflow



- Заменят ли они людей? Вопрос открытый, смотрите мои другие доклады про опасности AGI, но вот любопытные последние новости

- Leahy et al. (Oct 31, 2024): `https://www.thecompendium.ai/`

# The Compendium

*By Connor Leahy, Gabriel Alfour, Chris Scammell, Andrea Miotti, Adam Shimi*
*V1.0.1 – Nov 8, 2024*

*Humanity faces extinction from AGI.*

- Long et al. (Nov 4, 2024): станут ли AI моральными агентами?



**Taking AI Welfare Seriously**

**Robert Long**[*]
Eleos AI

**Jeff Sebo**[*]
New York University

**Patrick Butlin**[†]
University of Oxford

**Kathleen Finlinson**[†]
Eleos AI

**Kyle Fish**[†§]
Eleos AI, Anthropic

**Jacqueline Harding**[†]
Stanford University

**Jacob Pfau**[†]
New York University

**Toni Sims**[†]
New York University

**Jonathan Birch**[‡]
London School of Economics

**David Chalmers**[‡]
New York University

# LLM в науке — новости 2025 года

- Каждые пару недель происходит что-то новое
- В науке: AI CoScientist (март 2025)

- Jose Penades, микробиолог из Imperial College London:



Scientists spent 10 years on a superbug mystery - Google's AI solved it in 48 hours

The co-scientist model came up with several other plausible solutions as well

By Shawn Knight February 21, 2025 at 2:22 PM | 20 comments

PROFESSOR
Jose Penades, FRS

Analysis and Technology

Can Google's new research assistant AI give scientists 'superpowers'?

Researchers who have been given access to Google's new AI "co-scientist" tool are enthusiastic about its potential, but it isn't yet clear whether it can make truly novel discoveries

By Michael Le Page

📅 19 February 2025

- AI Scientist (Sakana AI, August 2024): система ходит к нескольким API (LLM, Semantic Scholar), умеет использовать информацию и ресурсы компьютера (сохранять веса моделей) и самостоятельно писать и запускать код экспериментов

- AI Scientist-v2 (Sakana AI, 12 марта 2025) смогла написать статью, которая прошла на ICLR 2025 Workshop "I Can't Believe It's Not Better: Challenges in Applied Deep Learning"!

| Title | ICLR Workshop Scores |
|---|---|
| Compositional Regularization: Unexpected Obstacles in Enhancing Neural Network Generalization | 6, 7, 6 |
| Real-World Challenges in Pest Detection Using Deep Learning: An Investigation into Failures and Solutions | 3, 7, 4 |
| Unveiling the Impact of Label Noise on Model Calibration in Deep Learning | 3, 3, 3 |

- Это, видимо, первая по-настоящему полностью автоматически порождённая статья, прошедшая серьёзный peer review и принятая в хорошее место

- Но не последняя! Система Zochi (Intology AI, March 2025) автономно написала статью, которую приняли на ACL 2025, на главный трек!



**The 63rd Annual Meeting of the Association for Computational Linguistics**

| Your Submissions | Author Tasks | | |
|---|---|---|---|
| # | Submission Summary | Official Review | Decision |
| 5156 | Tempest: Automatic Multi-Turn Jailbreaking of Large Language Models with Tree Search<br>Andy Zhou ✓, Ron Arel ✓<br>ACL 2025 Main<br>Show details | **0 Official Reviews Submitted**<br>**Average Rating:** N/A (Min: N/A, Max: N/A)<br>**Average Confidence:** N/A (Min: N/A, Max: N/A) | **ACL 2025 Main**<br>**Recommendation:**<br>**Accept (Main)**<br>Read |

- Статья, кстати, очень показательная…

## Tempest: Automatic Multi-Turn Jailbreaking of Large Language Models with Tree Search

**Andy Zhou***
Intology AI
andy@intology.ai

**Ron Arel***
Intology AI
ron@intology.ai

### Abstract

We introduce Tempest, a multi-turn adversarial framework that models the gradual erosion of Large Language Model (LLM) safety through a *tree search* perspective. Unlike single-turn jailbreaks that rely on one meticulously engineered prompt, Tempest expands the conversation at each turn, branching out multiple adversarial prompts that exploit partial compliance from previous responses. Through a cross-branch learning mechanism, successful attack patterns and partial compliance signals are systematically shared across parallel conversation paths, enabling more efficient discovery of model vulnerabilities. By tracking these incremental pol-

2024a; Ren et al., 2024; Zhao and Zhang, 2025; Yu et al., 2024). The dynamic nature of chat interfaces presents unique challenges for safety testing, as adversaries can adapt their strategies based on model responses and gradually accumulate partial compliance across multiple turns.

Traditional approaches to evaluating LLM safety have focused primarily on single-turn attacks, where carefully engineered prompts attempt to elicit harmful responses in one shot (Zou et al., 2023; Geiping et al., 2024). However, this methodology fails to capture how real-world adversaries interact with models through extended conversations, often employing multiple techniques in se-

- AlphaEvolve (DeepMind, May 14, 2025): продолжение и
  улучшение FunSearch и AlphaTensor, несколько LLM, которые
  пишут, критикуют и тестируют код для разных задач

- Новые результаты в разных областях:
  - новый scheduling в датацентрах – уже в Borg
  - улучшенная схема для умножения матриц (прямо на Verilog) – уже в TPU
  - новые алгоритмы: улучшил Штрассена для 4x4 над ▯, чего AlphaTensor раньше не смог



Data Center Optimization
**Borg Scheduling**

Hardware Optimization
**TPU Circuit Design**

Software Optimization
**Gemini Training**

- И в математике: новая нижняя оценка на контактное число (kissing number) в размерности 11, сразу несколько улучшенных оценок в анализе, геометрии, комбинаторике…

- Aygün et al. (8 сентября 2025): система для автоматического порождения новых методов, которые бьют рекорды на бенчмарках (при помощи MCTS-поиска и LLM, переписывающих код)

- В результате получается дерево вроде вот такого, из которого получаются реальные улучшения рекомбинацией методов:

- И эти улучшения легко интерпретировать:

- Nagda et al. (DeepMind, Sep 30, 2025): запустили AlphaEvolve на задачах из теоретической информатики, получили новые оценки для MAX-4-CUT и подобных задач

- А что с самой математикой? Diez et al. (Sep 3, 2025) попросили GPT-5 доказать новую теорему о выпуклой оптимизации; буквально:

#### 4.1.1 Gaussian framework

We started with the following initial prompt:

```
Paper 2502.03596v1 establishes a qualitative fourth moment
theorem for the sum of two Wiener-Itô integrals of orders
p and q, where p and q have different parities. Building
on the Malliavin-Stein method (see 1203.4147v3 for details),
could you derive a quantitative version for the total variation
distance, with a convergence rate depending solely on the
fourth cumulant of this sum?
```

The first interaction (see Annex A.1 for the entire discussion) was strikingly effective. GPT-5 produced a generally correct statement, using the right tools and approach. Howe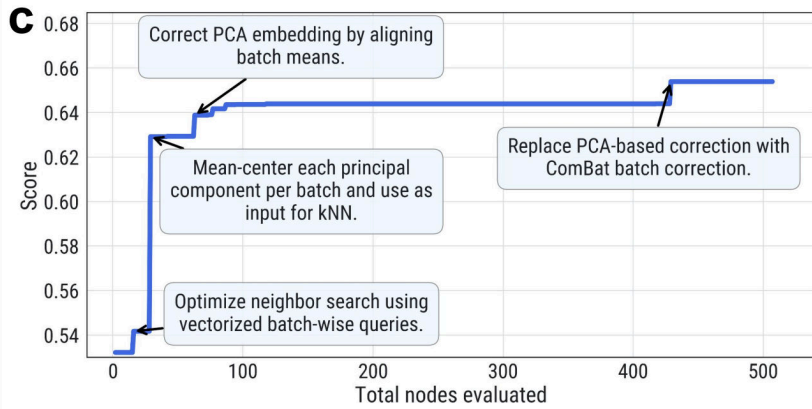ver, it made a reasoning error (leading to a wrong expression for $\mathrm{Cov}(Y^2, Z^2)$) that could have invalidated the whole proof if left unchecked. Noticing this, we then asked:

```
Can you check your formula for Cov(Y^2, Z^2) and provide me
with the details?
```

- В доказательстве была ошибка, и GPT-5 потребовалось аж два наводящих промпта (sic!), чтобы ошибку понять и исправить…

I think you are mistaken in claiming that $(p+q)!\|u\widetilde{\otimes}v\|^2 = p!q!\|u\|^2\|v\|^2$. Why should that be the case?

It eventually admitted (which is not surprising, since by alignment it usually agrees with us) that the statement was false, but more importantly, it understood where the mistake came from. This was followed by a reasoning and a formula that, this time, were correct.

Then, at our request, GPT-5 reformatted the result in the style of a research article, including an introduction, the presentation of our main theorem, its proof with all the details (correct this time!), and a bibliography. The exact prompt was:

```
Turn this into a research paper ready for submission. Follow
my style (see attached paper 0705.0570v4):
- start with an introduction giving some context,
- then present the main result, followed by a very detailed
proof where no step is left out,
- finish with a complete bibliography.
The final document should be a LaTeX file that I can compile.
```

- Feldman, Karbasi (Sep 22, 2025): тест Гёделя – могут ли LLM доказывать простые, но новые теоремы?

**Prompt to GPT-5**

Consider the problem of maximizing a function $F$ from $[0,1]^n$ to the reals that is the sum of a non-negative monotonically increasing DR-submodular function $G$ and a non-negative DR-submodular function $H$ over a solvable down-closed polytope $P$. I would like to bound the performance on this problem from the NeurIPS 2021 paper "Submodular + Concave" which is attached. Specifically, if $x$ is the output vector of this algorithm and o is the vector in $P$ maximizing $F$, then I would like to lower bound $F(x)$ with an expression of the form $\alpha * G(o) + \beta * H(o) - err$, where $\alpha$ and $\beta$ are constants. $err$ should be a function that depends only on the error parameter $\epsilon$ of the algorithm, the diameter $D$ of the polytope $P$, and the smoothness parameters $L_G$ and $L_H$ of $G$ and $H$, respectively, and goes to zero as $\epsilon$ goes to zero. Please give the best such bound that you prove (a bound is considered better if the values of the constants $\alpha$ and $\beta$ are larger). Provide a mathematically rigorous and well explained proof for the bound you come up with.
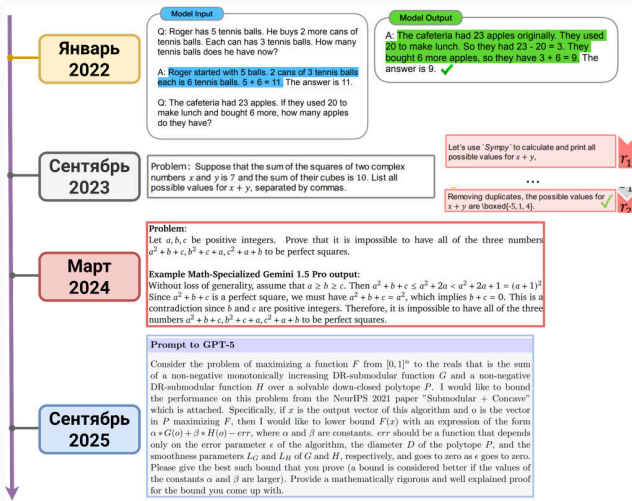
- GPT-5 решил три из пяти таких вот задачек, с небольшими легко исправимыми недочётами…

---

**Prompt to GPT-5**

A non-negative monotone differentiable function $F \colon [0,1]^{\mathcal{N}} \to \mathbb{R}_{\geq 0}$ is $\gamma$-weakly DR-submodular if for every two vectors $\mathbf{x}, \mathbf{y} \in [0,1]^{\mathcal{N}}$ such that $\mathbf{x} \leq \mathbf{y}$, it holds that
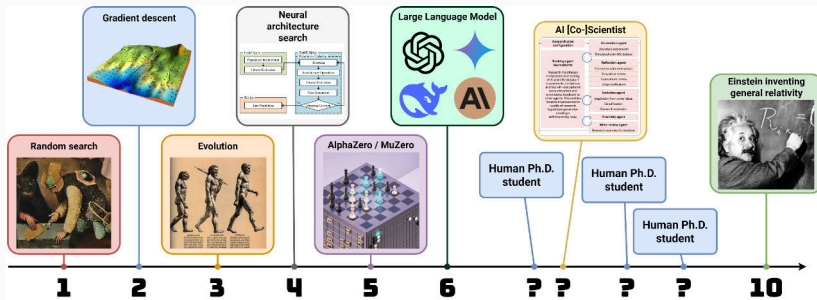
$$\sum_{u \in S}(y_u - x_u) \cdot \frac{\partial F(\mathbf{x})}{\partial x_u} \geq \gamma \cdot [F(\mathbf{y}) - F(\mathbf{x})] \ .$$

Given such a function which is also $L$-smooth and a convex set $\mathcal{C} \subseteq [0,1]^{\mathcal{N}}$. The goal in the current problem is describe an algorithm that finds a vector $\mathbf{z} \in \mathcal{C}$ that (approximately) maximizes $F$. A good reference for DR-submodular is the paper attached.

---

- А что на самом деле, без лишнего скептицизма и без шапкозакидательства?
- Научный поиск – это тоже оптимизационный процесс; где мы сейчас?

# Спасибо за внимание!