

## Фронтальные модели и повседневный мир



### Аналитический разбор AI Index 2026

Новые данные AI Index Report 2026 рисуют картину развития искусственного интеллекта, которая значительно сложнее, чем привычный сюжет «модели стали умнее». За последние два года ИИ впечатляюще продвинулся в формальных, хорошо измеримых задачах: от университетских экзаменов и олимпиад до инженерных и программных тестов. Однако там, где для человека требуются здравый смысл и повседневное восприятие мира, системы по-прежнему совершают грубые ошибки — вплоть до неспособности корректно прочитать время на обычных стрелочных часах. Новая повестка уже не в том, «умеет ли ИИ решать сложные задачи», а в том, как он сочетает высокие результаты на тестах с реальной пригодностью и безопасностью в обществе.

### Экзамены сданы: фронтальные модели выходят на уровень экспертов

На ряде сложных испытаний — от мультимодального MMMU до GPQA, ARC-AGI-2 и Humanity's Last Exam — ведущие модели показывают результаты, сопоставимые и иногда превышающие уровни сильных экспертов. Они успешно решают многошаговые задачи, используют диаграммы, таблицы и формулы, выдерживают «олимпиадный» формат с высокой концентрацией нетривиальных вопросов. В математике за короткий период выросла точность на сложных наборах — модели научились стабильно получать правильные ответы в задачах, которые даже для опытных математиков занимают часы работы.

Однако за этим прогрессом скрывается существенная асимметрия. Там, где требуется не только найти ответ, но и построить строгую, проверяемую линию рассуждений (например, в задачах на формальные доказательства), большинство моделей резко теряют в качестве. Новые наборы заданий, такие как IMO-Bench, показывают огромный разброс: отдельные системы демонстрируют впечатляющий уровень, но основная масса моделей все еще далека от устойчивой способности именно доказывать, а не угадывать. В результате ИИ уверенно играет в формализованные «экзаменационные игры», тогда как глубинное, прозрачное обоснование решений остается скорее исключением, чем нормой.

### Парадокс часов: сложные задачи решены, простые — нет

Одна из самых наглядных иллюстраций ограничений современных систем — их трудности с задачами, которые для людей кажутся тривиальными. Исследования, посвященные чтению аналоговых часов и пониманию календарей, показывают: даже лучшие мультимодальные модели сильно уступают людям в точности на таких заданиях. Там, где человек почти не ошибается, ИИ демонстрирует разрывы в десятки процентных пунктов, а типичные промахи измеряются часами, а не минутами.

Речь при этом не сводится к нехватке обучающих данных. Дополнительное обучение на синтетических изображениях часов может улучшать результаты на знакомых вариантах, но плохо переносится на более реалистичные или нестандартные изображения. Это указывает на более глубокую проблему: модели испытывают трудности, когда нужно увязать несколько визуальных признаков (форму циферблата, положение стрелок, деления, фоновые элементы) в цельную интерпретацию. Иначе говоря, ИИ хорошо решает задачи, которые можно свести к узнаваемым схемам, но гораздо менее надежен там, где требуется «собрать картину» из множества мелких деталей.

### **Видео как испытание на понимание мира**

Существенный сдвиг произошел в области моделей, работающих с видео. Новые системы не только генерируют ролики по текстовому описанию, но и начинают демонстрировать способности без дополнительного обучения, для которых традиционно создавали отдельные алгоритмы: сегментацию объектов, простое физическое моделирование, стилизацию и редактирование изображений. В отдельных экспериментах исследователи описывают у таких моделей зачатки «рассуждения по кадрам» — ситуацию, когда система как будто пошагово «прокручивает» развитие сцены, а не просто выдает набор несвязанных кадров.

Тем не менее, когда фокус смещается с зрелищности на правдоподобие, слабые места становятся очевидными. Наборы заданий вроде Video-Bench и VBench-2.0 оценивают качество видео не только по визуальной привлекательности, но и по физике, логике событий, управляемости сюжетом и соответствию запросу. На таких испытаниях ни одна модель пока не достигает высоких пороговых значений, а сложные истории с множеством объектов, длительной временной динамикой и неочевидными причинно-следственными связями все еще регулярно приводят к сбоям: нарушается непрерывность, «ломается» физика, появляются нелогичные переходы. Отсюда важный вывод: эффектные демонстрационные ролики нельзя автоматически считать признаком глубокого понимания реальности.

### **Агентные системы: от чат-бота к цифровому исполнителю**

Отдельная линия развития связана с агентными системами, которые действуют в сложных цифровых средах: браузере, виртуальном «офисе», терминале. В отличие от обычных моделей, отвечающих в диалоге, такие агенты сами иницируют действия, переходят по ссылкам, открывают файлы, выполняют команды и строят цепочки шагов. На испытаниях нового поколения лучшие системы впервые показывают успешность порядка 60–70% задач, а в некоторых сценариях по результатам почти приближаются к человеческому уровню.

Однако здесь снова проявляется хрупкость. Если слегка изменить внешний вид интерфейса, переформулировать условия задачи или скрыть структуру описания, показатели резко падают. Это говорит о том, что мы имеем дело с системами, которые хорошо работают в относительно стабильных, знакомых средах, но пока не обладают по-настоящему надежной способностью к адаптации. ИИ-агент сегодня скорее напоминает очень сильного исполнителя с жестко закрепленной логикой и ограниченной гибкостью, чем универсального цифрового коллегу.

### **Программирование и инженерия: ИИ как усилитель, а не замена инженера**

В области программной инженерии сдвиг особенно заметен. На наборах SWE-bench Verified, Vibe Code Bench и Terminal-Bench модели уже не ограничиваются подсказками кода, они решают значимую часть задач целиком: создают исправления для существующих репозиторий, собирают функциональные веб-приложения, уверенно работают в командной строке. Лучшие системы закрывают примерно от половины до трех четвертей заданий, что по меркам нескольких лет назад казалось почти недостижимым.

Но этот прогресс не означает, что инженер-человек становится избыточным. Даже при высоких процентах успешности остается немало случаев, когда модель предлагает частично неверные или нестабильные решения, требующие переработки. На практике это делает ИИ эффективным усилителем для профессионалов, который снимает значимую долю рутинной работы, но не заменяет ключевые компетенции системного проектирования, понимания ограничений и ответственности за итоговый результат.

### **Экономика и рынок труда: рост эффективности с перекосами**

Технический прогресс уже отражается на экономике и занятости. Исследования показывают: в профессиях с высокой степенью использования инструментов ИИ (службы поддержки, маркетинг, часть функций разработчиков) выработка сотрудников заметно увеличилась, иногда на десятки процентов. Это подтверждается и агрегированными показателями, которые фиксируют ускорение роста производительности в ряде экономик.

Однако выгоды распределяются неравномерно. Данные по занятости показывают, что сильнее всего удар приходится по младшим специалистам в профессиях, подверженных автоматизации: их доля занятости снижается, тогда как старшие когорты сохраняют позиции или даже укрепляют их. Похоже, что опытные работники используют ИИ как рычаг для повышения своей эффективности и стоимости на рынке, тогда как начинающие конкурируют не только друг с другом, но и с быстро автоматизируемыми задачами. Это ставит новые вопросы перед политикой занятости и системой образования: как поддерживать вход в профессии в условиях, когда именно простые и базовые виды работы быстрее всего переходят к алгоритмам.

### **Ответственный ИИ: измерения безопасности явно отстают**

Самый чувствительный разрыв, на который указывает отчет, пролегает между ростом общих способностей моделей и степенью развития практик ответственного ИИ. По техническим метрикам — точности, полноте, результатам на стандартизированных заданиях — сообщество располагает обширным набором сопоставимых данных, что позволяет достаточно корректно сравнивать модели. По измерениям безопасности, справедливости, защиты данных и фактической надежности картина гораздо более фрагментирована: разные группы используют несогласованные методики, не всегда в полном объеме публикуют результаты, общепринятых стандартов пока не существует.

Появляются работы, анализирующие взаимное влияние различных параметров: оказывается, что улучшение показателей по одной оси (например, усиление фильтров безопасности) может приводить к ухудшению по другой (к росту предвзятости или снижению полноты и честности ответов). Это означает, что речь идет не о простом «регулировании качества» по одному параметру, а о сложном пространстве компромиссов между точностью, безопасностью, справедливостью и прозрачностью. Без согласованных стандартов измерения и открытых данных эти компромиссы рискуют остаться неявными и, по сути, непроверяемыми.

### **Вместо вывода: что действительно важно в ближайшие годы**

В совокупности новые данные AI Index Report 2026 показывают ИИ как технологию на пороге содержательной зрелости. С одной стороны, модели достигли впечатляющих результатов в формализованных областях: экзамены, олимпиадные задачи, инженерные и программные наборы заданий. С другой — они продолжают допускать грубые ошибки в задачах здравого смысла, демонстрируют хрупкость при изменении условий и функционируют в условиях недоразвитой системы ответственных практик.

Это смещает фокус дискуссии. Главный вопрос уже не в том, появится ли очередная модель с более высокими баллами на стандартных тестах, а в том, как научиться одновременно измерять и развивать три группы характеристик: реальную пригодность ИИ в сложной, плохо формализованной среде; устойчивость и предсказуемость поведения; и ценностно значимые свойства — от безопасности до справедливости и защиты приватности. От того, насколько быстро и согласованно научное и инженерное сообщество справится с этой задачей, во многом зависит, станет ли нынешний скачок в развитии ИИ устойчивой основой для долгосрочного роста или источником новых системных рисков.

Источник: [REDACTED]