Введение

Вопрос о пределах искусственного интеллекта (ИИ) приобрёл особую актуальность на фоне роста вычислительных мощностей, успехов машинного обучения и появления архитектур, моделирующих сложные когнитивные функции. Однако остаётся фундаментальная проблема:

> существуют ли пределы, которые ИИ не может преодолеть по своей логической природе, независимо от технического прогресса?

Данная работа посвящена исследованию границ ИИ как алгоритмической системы, через призму формальной теории вычислений. Основной тезис заключается в следующем:

> Никакая ИИ-система, как бы ни была она устроена, не способна выйти за рамки вычислимых функций в смысле Тьюринга.

На основании гипотезы Чёрча-Тьюринга, теоремы Тьюринга о неразрешимости проблемы останова, теорем Гёделя, Райса и Матиясевича, доказывается, что:

поведение любой ИИ-модели формализуемо как Turing-вычислимая функция;

обучение, самооптимизация, стохастика не нарушают границы вычислимости;

гипотетические идеи мета-ИИ и «взрыва интеллекта» не имеют формального основания.

Работа построена строго, от фундаментальных определений до логических следствий. Её цель— не описать, как устроен ИИ, а показать, чего он никогда не сможет сделать

Оглавление

- 1. Введение
- 2. Глава 1. Теоретическая основа: гипотеза Чёрча-Тьюринга
 - 1.1 Модель Тьюринга
 - 1.2 Эквивалентность формализмов
 - 1.3 Эффективная вычислимость
 - 1.4 Смысл гипотезы
- 3. Глава 2. Логические границы вычислений
 - 2.1 Проблема останова Тьюринга
 - 2.2 Существование неразрешимых функций
 - 2.3 Обобщённые логические ограничения
 - 2.4 Значение для ИИ
- 4. Глава 3. Искусственный интеллект как вычислимая функция
 - 3.1 Структура нейросети как функции
 - 3.2 Обучение не делает функцию невычислимой
 - 3.3 Общий вид ИИ-системы
 - 3.4 Невозможность «прорыва» за пределы функции

- 3.5 Эвристические и вероятностные методы
- 5. Глава 4. Универсальность и предел моделей
 - 4.1 Эквивалентность вычислительных моделей
 - 4.2 Замкнутость класса вычислимых функций
 - 4.3 Попытки выйти за пределы машины Тьюринга
 - 4.4 Примечание о границах
- 6. Глава 5. Аргументация ограничения ИИ
 - 5.1 Построение формальной позиции
 - 5.2 Несостоятельность ИИ в вычислении неразрешимых функций
 - 5.3 Общая формализация ограничения
 - 5.4 Следствия: невозможность универсального ИИ
- 7. Глава 6. Иллюзия мета-вычислительного развития
 - 6.1 Формальное описание гипотезы самоусиления
 - 6.2 Предел функции самосовершенствования
 - 6.3 Контраргумент: увеличение мощности ≠ увеличение класса
 - 6.4 Принцип рефлексивного ограничения
- 8. Глава 7. Примеры невычислимых функций
 - 7.1 Проблема останова и функция Чёрча
 - 7.2 Диофантовы уравнения и теорема Матиясевича
 - 7.3 Функция Блума
 - 7.4 Истинные, но недоказуемые утверждения
- 9. Заключение
- 10. Приложение: Использованные источники и упомянутые работы
- «Пределы вычислимости и невозможность мета-алгоритмического искусственного интеллекта: формально-логическое обоснование в свете гипотезы Чёрча-Тьюринга и теоремы Гёделя»
- Глава 1. Теоретическая основа: гипотеза Чёрча-Тьюринга

Гипотеза Чёрча-Тьюринга, сформулированная независимо двумя математиками — Алонзо Чёрчем и Аланом Тьюрингом в 1936 году, — является краеугольным камнем всей современной теории алгоритмов. Она утверждает следующее:

Любая функция f, которую можно вычислить "эффективно", может быть вычислена на машине Тьюринга.

Иначе говоря:

∀ f (эффективно вычислима) ⇒ f вычислима на машине Тьюринга.

1.1 Модель Тьюринга

Машина Тьюринга — это абстрактная конструкция, включающая:

- бесконечную ленту, разбитую на ячейки;
- головку чтения/записи, способную перемещаться по ленте;
- конечное множество состояний Q:
- алфавит Σ (например, 0, 1, пустой символ);
- − таблицу переходов δ: Q × Σ → Q × Σ × {L, R}.

Каждый шаг машины — это переход вида:

$$(q, a) \rightarrow (q', a', D),$$

где:

q - текущее состояние,

а - текущий символ под головкой,

q' — новое состояние,

а' — символ для записи,

D — направление движения (влево или вправо).

Это простейшая, но при этом универсальная модель всех пошаговых алгоритмических процессов.

1.2 Эквивалентность формализмов

Независимо от Тьюринга, Алонзо Чёрч сформулировал понятие λ -исчисления — формальную систему редукции выражений с помощью абстракций и подстановок. Он доказал, что многие функции можно выразить в виде λ -термов и последовательно редуцировать до результата.

Позднее было установлено, что:

λ-исчисление ≡ машина Тьюринга ≡ рекурсивные функции ≡ машина Поста.

То есть:

Все эти модели описывают один и тот же класс вычислимых функций.

Множество всех таких функций называют Turing-computable.

1.3 Эффективная вычислимость

Ключевой момент: гипотеза Чёрча-Тьюринга не формализуема как теорема, потому что опирается на понятие интуитивной вычислимости. Она утверждает:

Если алгоритм может быть выполнен человеком, последовательно, без вмешательства случайности, с ограниченными ресурсами, то он может быть реализован на машине Тьюринга.

Это делает гипотезу своего рода аксиомой для всего, что мы называем алгоритмом.

1.4 Смысл гипотезы

Гипотеза Чёрча-Тьюринга не ограничивает конкретные технологии. Она утверждает нечто более фундаментальное:

> Не существует вычислимой функции, которую не могла бы вычислить машина Тьюринга, если у неё будет достаточно времени и памяти.

Это означает, что даже самые продвинутые формы ИИ, суперкомпьютеры или гипотетические "алгоритмические организмы" не могут реализовать функции, которые лежат вне класса Turing-computable.

Таким образом, гипотеза Чёрча-Тьюринга задаёт естественную и логическую границу для всех алгоритмических систем — независимо от архитектуры, скорости, параллелизма или способа организации вычислений. Именно на этом фундаменте строится дальнейшее доказательство ограниченности искусственного интеллекта.

Глава 2. Логические границы вычислений

Ключевым открытием в логике XX века стало осознание того, что не всё, что можно сформулировать, можно вычислить. Уже в 1930-х годах было строго доказано существование задач, для которых не существует никакого алгоритма, способного дать ответ в общем случае. Эти задачи — не просто сложные. Они — принципиально неразрешимы.

2.1 Проблема останова Тьюринга

Одна из самых известных неразрешимых задач — это проблема останова (Halting Problem). Суть задачи: существует ли алгоритм H, такой, что для любой программы P и входа x, он определяет:

Остановится ли программа Р на входе х, или зациклится бесконечно?

Формально:

H(P, x) =

1, если Р(х) останавливается

0, если Р(х) зацикливается бесконечно

Тьюринг доказал: такой универсальный Н не существует.

Доказательство опирается на парадокс диагонализации (по аналогии с доказательством Кантора). В краткой форме:

- 1. Предположим, что Н существует.
- 2. Тогда можно построить программу D, которая делает следующее:

на входе своей собственной программы D она вызывает H(D, D) и, если H говорит "остановится", — зацикливается;

а если Н говорит "не остановится" — завершает выполнение.

3. Получается противоречие: программа D ведёт себя не так, как предсказал H. Следовательно, H не может существовать.

То есть:

 \neg \exists H: H(P, x) всегда правильно определяет остановку P(x)

2.1.1. Уточнение: предельность и прикладная обработка задачи останова

Хотя задача останова в общем виде неразрешима (в силу теоремы Тьюринга), это не исключает существования эффективных практических методов, позволяющих частично решать её на ограниченных классах входов. В прикладной инженерной практике широко применяются подходы, обеспечивающие точную или приближённую оценку остановаемости:

- статические анализаторы, определяющие остановаемость в обширных подклассах программ;
- анализ кода с ограничением глубины рекурсии, количества итераций или размеров памяти;
- использование языков с ограниченной выразительной мощностью, где задача останова становится разрешимой;
- статистические и обучаемые модели, оценивающие вероятность останова по синтаксическим признакам кода.

Эти методы не противоречат теоретической неразрешимости в общем случае, но показывают, что множество реальных сценариев поддаётся надёжному анализу. Следовательно, важно различать теоретически неразрешимую универсальную постановку задачи останова и прикладные варианты, где ИИ может успешно работать в пределах допустимого приближения или частичной формализации.

2.2 Существование неразрешимых функций

Отсюда следует: существует множество функций, которые невычислимы.

Пример: функция f, определяемая так:

f(n) =

1, если n-я программа P_n останавливается на входе n

0 — иначе

Эта функция называется функцией останова по диагонали (diagonal halting function). Она не вычислима:

Нет ни одного алгоритма, который её реализует.

2.3 Обобщённые логические ограничения

На этом фундаменте строится ряд других теорем:

- Теорема Гёделя о неполноте (1931):

В любой достаточно мощной формальной системе (например, арифметике Пеано) существует утверждение G, такое что:

G не доказуемо и не опровержимо в рамках этой системы.

- Следствие:

Не существует алгоритма, способного определить истинность всех утверждений арифметики.

Теорема Райса (1953):

Любое нетривиальное свойство поведения программ — неразрешимо.

Формально:

Пусть S — множество программ, обладающих некоторым свойством. Если $S \neq \emptyset$ и $S \neq$ множество всех программ, то вопрос «принадлежит ли программа P множеству S» — неразрешим.

2.4 Значение для ИИ

Эти логические ограничения важны для понимания пределов ИИ:

- Никакая ИИ-система не может решить проблему останова в общем виде.
- ИИ не сможет отвечать на все математические вопросы (в частности, из арифметики).
- ИИ не сможет узнать, имеет ли программа определённое поведение, если это свойство не тривиально.

Иначе говоря:

ИИ не может вычислить больше, чем позволяет теория Тьюринга. ∀ задача Z: если Z ∉ Turing-computable, то Z ∉ мощности любого ИИ.

Вывод главы:

Уже с 1930-х годов известно, что не существует алгоритма, способного решить все логически формализованные задачи. Эти границы не являются следствием технических ограничений — они заложены в самой структуре формальной логики и алгоритмики. Именно эти границы становятся фундаментом для дальнейшего доказательства того, что искусственный интеллект не может выйти за рамки вычислимого

Глава 3. Искусственный интеллект как вычислимая функция

Сегодняшние реализации искусственного интеллекта (ИИ) — независимо от архитектуры — представляют собой строго математические модели, которые обрабатывают входные данные и выдают результат. Наиболее распространённой и развитой моделью является искусственная нейронная сеть (ИНС). Однако, несмотря на кажущуюся сложность, такая сеть представляет собой вычислимую функцию, а значит, входит в область действия машины Тьюринга.

3.1 Структура нейросети как функции

Современная нейросеть, особенно глубокая (deep neural network), реализует композицию нелинейных функций. Её математическое представление можно записать в следующем виде:

$$f(x) = \sigma_n(W_n \cdot \sigma_{n-1}(W_{n-1} \cdot ... \sigma_1(W_1 \cdot x + b_1) ... + b_{n-1}) + b_n)$$

где: – х — входной вектор; – W_i — матрицы весов для каждого слоя; – b_i — смещения (bias); – σ_i — функции активации (например, ReLU, sigmoid, tanh); – f(x) — выходной вектор.

Вся такая система есть не что иное, как суперпозиция элементарных функций, то есть:

$$f(x) = \phi_1 \circ \phi_2 \circ ... \circ \phi_n(x)$$

А любая конечная суперпозиция вычислимых функций— снова вычислима. Следовательно:

f принадлежит множеству Turing-вычислимых функций.

3.2 Обучение не делает функцию невычислимой

Можно возразить: «но нейросеть обучается, а значит, развивается». Однако обучение — это лишь процесс нахождения параметров (весов и смещений), оптимизирующих функцию по заданному критерию. Алгоритмы обучения (например, градиентный спуск) сами по себе строго детерминированы.

Простейший вариант обучения:

$$\theta \leftarrow \theta - \eta \cdot \nabla J(\theta)$$

где: – θ — вектор параметров; – η — скорость обучения; – $J(\theta)$ — функция потерь; – $\nabla J(\theta)$ — градиент.

Это выражение — не более чем итерационный алгоритм. Он также вычислим. А значит, процесс обучения нейросети не выводит систему за пределы машины Тьюринга.

3.3 Общий вид ИИ-системы

ИИ любого типа (в том числе трансформеры, как GPT) можно свести к следующему:

– Пусть М — ИИ-модель, – Пусть х — входные данные, – Тогда M(x) — результат работы модели.

Если модель реализуется как алгоритм (что верно для всех современных систем), то M(x) вычисляется на машине Тьюринга. Поэтому:

для любого х: результат работы модели — вычислимая функция.

Иными словами: результат работы любой ИИ-системы вычислим, вне зависимости от объёма данных, сложности архитектуры или глубины модели.

3.4 Невозможность «прорыва» за пределы функции

Даже если ИИ научится «оптимизировать себя», создавать «новые версии себя» и улучшать архитектуры, это не изменит главного:

– Эти действия также реализуются программно. – Значит, они принадлежат множеству вычислимых процедур. – А значит, ИИ не может породить невычислимую функцию или превзойти сам класс вычислимых функций.

Таким образом:

любая ИИ-модель A и любые процессы её модификации B — остаются в пределах Тьюринг-вычислимого.

3.5 Эвристические и вероятностные методы

Иногда утверждается, что ИИ способен «обходить» ограничения формальной вычислимости благодаря использованию вероятностных, эвристических или аппроксимационных методов. Такие подходы действительно применяются в практике машинного обучения и могут давать успешные приближённые результаты.

Однако важно понимать:

- Эвристики и вероятностные процедуры не являются универсальными алгоритмами.
- Они не дают гарантированного решения задачи для всех случаев. Они не нарушают границы вычислимости, а лишь обходят их в частных случаях, снижая требования к точности.

Таким образом, даже использование стохастических моделей не позволяет преодолеть фундаментальные ограничения формальной теории вычислений.

Вывод главы:

Как бы сложна ни была архитектура ИИ, она остаётся функцией, определяемой алгоритмически. Ни обучение, ни самооптимизация, ни использование вероятностных методов не делают ИИ чем-то мета-вычислительным. Он всегда остаётся внутри теоретически ограниченного множества Тьюринг-вычислимых функций. В этом смысле — не может «вырваться» за границы алгоритма.

3.5.1. Эвристическая применимость и частичная решаемость

Многие задачи, теоретически неразрешимые в общем виде, оказываются практически решаемыми на ограниченных или структурированных входах с помощью эвристических методов. Искусственный интеллект активно использует такие методы для приближённого, вероятностного или частичного решения задач, не обладающих полной алгоритмической определённостью.

Наиболее значимые практики включают:

- эвристический поиск (A*, beam search, minimax с отсечением α-β и пр.);
- машинное обучение, позволяющее строить аппроксимирующие функции для вывода или классификации;
- ограниченный перебор с эвристическим отсечением (например, в SAT- и SMTрешателях);
- автоматические доказатели теорем, способные строить формальные доказательства в пределах конкретной логической системы.

Эти методы эффективны в инженерной практике и позволяют добиваться впечатляющих результатов при решении задач с ограниченной размерностью, структурой или допускаемой ошибкой. Однако все они опираются на внутренние

эвристики, статистику, эвентуальные вероятности успеха и не обладают полной алгоритмической гарантией.

При этом эвристическая эффективность не нарушает фундаментальных пределов: теоретическая неразрешимость сохраняется, и никакой метод, основанный на приближении или переборе, не способен обойти ограничения, налагаемые теоремами Тьюринга, Райса и Гёделя.

Таким образом, эвристические ИИ-системы демонстрируют высокую применимость, но не преодолевают границы формальной вычислимости. Их успех — это техническая адаптация, а не теоретический прорыв.

Глава 4. Универсальность и предел моделей

Проблема границ вычислимости не сводится к конкретной архитектуре или языку программирования. Существуют разные формальные модели вычислений, созданные для описания понятия «алгоритм». И каждая из них — независимо от своей формы — вычислительно эквивалентна машине Тьюринга. Это означает, что все они способны вычислить один и тот же класс функций, и никакая из них не превосходит другие по вычислительной мощности.

4.1 Эквивалентность вычислительных моделей

Ниже перечислены ключевые формализмы, эквивалентные машине Тьюринга:

- 1. Лямбда-исчисление
- Формализованное описание вычислений через функции и подстановки.
- Любая функция, вычислимая в лямбда-исчислении, вычислима на машине Тьюринга— и наоборот.
- 2. Машина Поста
- Альтернативная модель, основанная на переписывании строк.
- Эквивалентна машине Тьюринга по мощности.
- 3. Регистровые машины
- Используют простые команды (увеличение, обнуление, переход).
- Любую программу на такой машине можно преобразовать в эквивалентную машину Тьюринга.
- 4. Формальные грамматики и машины Маркова
- Применяются в теории языка и логике.
- Не выходят за пределы Тьюринг-вычислимого.

Таким образом: любая из этих моделей описывает тот же самый набор функций, что и машина Тьюринга.

4.2 Замкнутость класса вычислимых функций

Обозначим \mathbb{F}_{t} как множество всех функций, вычислимых на машине Тьюринга. Тогда:

- Если f и g вычислимые, то их композиция f(g(x)) тоже вычислима.
- Если две функции дают одинаковый результат на всех входах, и одна из них вычислима, то и вторая также принадлежит этому классу.
- Любую функцию из \mathbb{F}_{t} можно выразить в любой из эквивалентных формальных моделей.

Таким образом, этот класс функций — замкнутая структура, в которую нельзя «влить» нечто более мощное, не нарушив основы логики.

4.3 Попытки выйти за пределы машины Тьюринга

Несмотря на мощь классических моделей, предпринимались попытки создать гипотетические вычислительные системы, выходящие за пределы Тьюринговской модели. Некоторые из них:

1. Машины с оракулом

- Обладают «чёрным ящиком», который может решать неразрешимые задачи (например, проблему останова).
- Однако такая система не может быть физически реализована и используется только в теоретических рассуждениях.

2. Суперзадачи (supertasks)

- Модели, предполагающие выполнение бесконечного числа шагов за конечное время.
- Противоречат законам физики и принципу причинности, а потому остаются абстракцией.

3. Аналоговые или непрерывные вычисления

- Предполагают, что система может использовать непрерывные значения или бесконечную точность.
- На практике невозможно реализовать бесконечно точную физическую величину шум и квантовые ограничения ставят предел.

4. Квантовые компьютеры

- Часто ошибочно воспринимаются как сверх-Тьюринговские.
- На самом деле, они лишь ускоряют вычисления, но не расширяют класс решаемых задач. Квантовая машина не может вычислить невычислимую функцию.

Следовательно, хотя такие модели интересны теоретически, никакая из них не предоставляет доступ к функциям за пределами Turing-computable, если только не допустить нарушение физических или логических ограничений.

4.3.1. Квантовые вычисления и границы вычислимости

Несмотря на фундаментальные различия между квантовыми и классическими вычислениями, квантовые компьютеры не нарушают границ, установленных гипотезой Чёрча-Тьюринга. Все квантовые алгоритмы, включая наиболее известные (алгоритм Шора, алгоритм Гровера), решают задачи, которые остаются в классе функций, вычислимых на универсальной машине Тьюринга. Это означает, что квантовые вычисления не расширяют множество вычислимых задач, а лишь изменяют сложность отдельных из них.

Класс задач, решаемых квантовыми алгоритмами с полиномиальной сложностью и ограниченной вероятностью ошибки (класс BQP), строго содержится в пределах Turing -computable функций. Например:

- алгоритм Шора позволяет выполнять факторизацию за полиномиальное время, тогда как для классических алгоритмов она экспоненциальна;
- алгоритм Гровера обеспечивает квадратичное ускорение при неструктурированном поиске.

Однако такие ускорения не затрагивают нерешаемые задачи — в частности, проблему останова или общую проблему логического следования. Квантовые компьютеры не позволяют обойти теоремы Тьюринга, Гёделя или Райса. Поэтому все теоретические границы, обсуждаемые в данной работе, сохраняют силу и применимы в том числе к ИИ-системам, основанным на квантовых вычислениях.

4.4 Примечание о границах

Важно подчеркнуть: границы, о которых идёт речь, не являются следствием слабой архитектуры, недостатка памяти или медленного процессора. Эти ограничения:

- Заложены в самой логике алгоритма;
- Не устраняются техническим прогрессом;
- Не преодолеваются сменой вычислительной модели, если она формализуема.

Если принять гипотезу Чёрча—Тьюринга, то всякое формализуемое и реализуемое вычисление уже входит в \mathbb{F}_{t} . Попытки построить сверх-Тьюринговые машины либо не реализуемы, либо нарушают физические законы.

Вывод главы

Все известные формальные модели вычислений эквивалентны машине Тьюринга. Это задаёт строгую и универсальную границу: любой искусственный интеллект, как алгоритмическая система, работает внутри этого предела. Попытки выйти за него с помощью гипотетических моделей либо неосуществимы, либо не выдерживают физической проверки. Таким образом, граница вычислимости — это не технический потолок, а логико-онтологический предел любой алгоритмической системы.

Глава 5. Аргументация ограничения ИИ

На современном уровне развития искусственный интеллект (ИИ) — это сложная, но конечная система, выражающаяся в виде вычислимой программы. Она может принимать на вход данные, производить обучение, трансформировать информацию, предсказывать вероятности — но все эти действия укладываются в рамки алгоритмически заданной функции.

Сущностный вопрос: может ли ИИ, даже в пределе своего развития, выйти за рамки формализуемой вычислимости? Ответ — нет, и ниже это будет показано строго.

5.1 Построение формальной позиции

Пусть:

- A любая система искусственного интеллекта;
- Р программа, реализующая поведение А;
- x произвольный вход;
- A(x) = P(x) -результат работы на входе x.

Так как Р является программой, её поведение можно описать как функцию:

f(x) = результат выполнения P на входе x.

Поскольку любая программа P конечна и работает по алгоритму, f — Turing-computable функция. Значит:

$$\forall x \in D$$
: $f(x) \in \mathbb{F}_t$

Следовательно, поведение ИИ всегда ограничено:

```
Behavior(A) \subseteq \mathbb{F}_{t}
```

5.2 Несостоятельность ИИ в вычислении неразрешимых функций

Существует множество функций, которые не входят в класс $\mathbb{F}_{\mathfrak{t}}$, то есть принципиально невычислимы. Примеры таких функций служат жёстким доказательством ограничений ИИ.

Пример 1: Функция останова

```
h(p, x) = 1, если программа р останавливается на входе x, h(p, x) = 0, иначе.
```

Теорема Тьюринга: функция h(p, x) невычислима.

Следствие: ни один ИИ не сможет построить универсальный предиктор останова.

Пример 2: Диагональная функция (парадоксальная)

Пусть $f_n(x) - n$ -я программа в счётной нумерации всех возможных программ. Рассмотрим функцию:

```
D(n) = { 0, если f<sub>n</sub>(n) = 1;
1, иначе }
```

Тогда D не совпадает ни с одной f_n , то есть D $\notin \mathbb{F}_t$.

Это — модификация диагонального аргумента Кантора, применённого к множеству вычислимых функций. ИИ не сможет вычислить такую функцию в принципе, поскольку её построение нарушает замкнутость \mathbb{F}_{+} .

Пример 3: Разрешимость диофантовых уравнений

Пусть $H(a_1, a_2, ..., a_n)$ = 1, если соответствующее диофантово уравнение имеет решение, и 0 — иначе.

По теореме Ю. Матиясевича (1970), функция Н не вычислима.

Значит, ни одна ИИ-система не сможет универсально определять наличие решений у произвольных уравнений в целых числах.

5.3 Общая формализация ограничения

Пусть \mathbb{F}_{-} А — класс всех функций, которые способен реализовать ИИ А. Тогда:

```
\mathbb{F}_A \subseteq \mathbb{F}_t
```

При этом:

- Существует множество f: f ∉ 𝔻 t
- Следовательно: f ∉ **F** _A

Таким образом, никакая система А не сможет вычислить такие функции. Это не вопрос мощности железа или сложности архитектуры, а логический предел.

5.4 Следствия: невозможность универсального ИИ

Из сказанного вытекают конкретные невозможности:

1. Невозможно создать ИИ, который будет универсальным решателем всех математических задач.

(Такой ИИ нарушил бы теорему о неразрешимости.)

- 2. Невозможно создать ИИ, который предсказывает поведение всех других ИИ. (Поскольку тогда он должен решать проблему останова.)
- 3. Невозможно создать ИИ, который сам создаёт ИИ с большей вычислительной мощностью, чем у него самого.

(Потому что он не может перейти из \mathbb{F}_{t} в класс более широкий — такого просто не существует.)

Вывод главы:

ИИ, как бы ни эволюционировал, всегда останется внутри класса Turing-computable функций. Он не сможет вычислить ничего за пределами \mathbb{F}_{t} , а значит, существует абсолютно непреодолимая граница, которая не зависит ни от техники, ни от архитектурных инноваций.

Это ограничение вытекает не из практики, а из математической логики и теории алгоритмов, и оно имеет статус жесткой онтологической границы ИИ.

5.4.1. Ограниченная универсальность ИИ

Современные ИИ-системы демонстрируют так называемую ограниченную универсальность — способность успешно функционировать в широком спектре прикладных задач, при этом опираясь на обобщённые алгоритмические методы адаптации, преобразования и обобщения. К таким системам относятся нейросетевые архитектуры, трансформер-модели, логико-программные агенты и другие вычислительные конфигурации, которые в пределах одной архитектуры способны переключаться между задачами распознавания, генерации, классификации, поиска, оптимизации и вывода.

Однако подобная универсальность является строго инженерной и формализуемой. Её

основа— не наличие неограниченной способности к решению любых задач, а возможность выполнять различные функции в пределах одной вычислимой структуры. Такие ИИ-системы не обладают содержательной полнотой или рефлексивной целостностью, и их функциональная гибкость реализуется внутри формально заданного множества операций. Независимо от сложности внутренней архитектуры или глубины модели, каждая из таких систем остаётся алгоритмически ограниченной.

Если рассматривать эти модели с позиции теории вычислимости, они не выходят за пределы класса Тьюринг-вычислимых функций: все их трансформации, обучение, генерации и выводы поддаются описанию как конечные алгоритмы. Это означает, что никакая форма «ограниченной универсальности» не нарушает фундаментальных границ, определённых теоремами Тьюринга, Райса и Гёделя. ИИ-система может выполнять множество различных задач, но остаётся в пределах счётного, формализуемого и замкнутого пространства функций.

Следовательно, даже самые универсальные на практике ИИ-модели не являются универсальными в теоретическом смысле. Их гибкость не указывает на эволюционное приближение к разуму или трансцендентной субъективности. Она лишь отражает архитектурную способность реализовывать конечный набор программных операций в пределах одного и того же класса.

Таким образом, универсальность ИИ — не переход к метаинтеллекту, а обобщённая форма внутренней перестройки внутри ограниченной формальной модели. Никакая такая система не преодолевает барьер алгоритмической замкнутости, и её возможности, какими бы широкими они ни казались, полностью определяются структурой вычислимого.

Глава 6. Иллюзия мета-вычислительного развития

Идея о так называемом «взрыве интеллекта» (intelligence explosion), впервые серьёзно сформулированная И. Дж. Гудом в середине XX века, предполагает, что разумная машина, достигнув определённого уровня сложности, сможет создавать машины, умнее самой себя. Те, в свою очередь, создадут ещё более умные, и так далее — в результате чего возникнет каскад самоусовершенствования, ведущий к качественно новому уровню интеллекта, превосходящему человеческий.

Однако такое предположение основано на логической иллюзии. Оно игнорирует фундаментальное ограничение вычислимости, а именно — то, что никакая система не может выйти за рамки собственной вычислительной мощности. Мы покажем, что так называемый «мета-ИИ» не может существовать как нечто вычислительно превосходящее систему, его породившую.

6.1. Формальное описание гипотезы самоусиления

Рассмотрим модельную цепочку самопорождающихся ИИ-систем:

А_о — исходная система, заданная как программа;

 $A_1 = A_0(A_0)$ — результат применения A_0 к самой себе;

 $A_2 = A_1(A_1)$ — результат применения A_1 к самой себе;

и так далее: $A_n = A_{n-1}(A_{n-1}), \forall n \in \mathbb{N}$.

--

Популярная версия гипотезы самоусиления утверждает, что начиная с некоторого N, каждая следующая система становится более «мощной», чем предыдущая:

$$\exists N: A_n > A_{n-1}$$

где неформально предполагается рост способности к решению всё более сложных задач, и в пределе якобы возникает система S, выходящая за рамки человеческого интеллекта:

$$A_n \ \to \ S.$$

Однако с формальной точки зрения:

- 1. Пусть A_0 принадлежит множеству Turing-computable функций: $A_0 \in \mathbb{F}_{t}$.
- 2. Тогда $A_1 = A_0(A_0)$ также принадлежит \mathbb{F}_{t} , поскольку вычислимая функция, применённая к вычислимому описанию, остаётся в классе вычислимых.
- 3. По индукции: \forall n ∈ \mathbb{N} , A_n ∈ \mathbb{F}_t .
- 4. Следовательно, вся цепочка $\{A_0, A_1, A_2, ...\} \subseteq \mathbb{F}_t$.

При этом понятие «мощности» не означает выход за границы вычислимого. Система может усложняться архитектурно, увеличивать глубину представлений, расширять охват задач, снижать ресурсоёмкость — но всё это происходит внутри замкнутого пространства алгоритмически описываемых функций.

Кроме того, предельный переход $A_n \to S$ не имеет формального смысла в рамках теории алгоритмов, поскольку множество программ дискретно, и на нём не задана операция предела. Более корректно говорить о возрастающей последовательности вычислимых систем, не покидающей класс \mathbb{F}_{t} .

Таким образом, даже если каждая A_n становится инженерно более эффективной или универсальной, никакое количество итераций не способно породить систему, выходящую за пределы Тьюринг-вычислимости. Так называемый «сверхразум» S, если он строится последовательным самоусилением, остаётся функцией в пределах $\mathbb{F}_{\mathbf{t}}$ и подчиняется всем фундаментальным ограничениям, описанным в теоретической части данной работы.

6.2 Предел функции самосовершенствования

Можно рассмотреть гипотетическую функцию улучшения:

$$U: \mathbb{F}_{+} \rightarrow \mathbb{F}_{+}$$

где U(f) — функция, описывающая «улучшенный» алгоритм на основе f.

Тогда:

$$\begin{split} &f_0 \in \ \mathbb{F}_{ \ t} \\ &f_1 = U(f_0) \in \ \mathbb{F}_{ \ t} \\ &f_2 = U(f_1) \in \ \mathbb{F}_{ \ t} \end{split}$$

$$f_n = U^n(f_0) \in \mathbb{F}_{t}$$

И в любом случае: \forall n ∈ \mathbb{N} f_n ∈ \mathbb{F} t

Никакая итерация применения U не позволяет выйти за пределы \mathbb{F}_{t} .

6.3 Контраргумент: увеличение мощности ≠ увеличение класса

Некоторые могут возразить: «Да, ИИ остаётся внутри \mathbb{F}_{t} , но ведь он может становиться мощнее — быстрее, эффективнее, способнее решать сложные задачи». Это верно — но такие улучшения касаются времени выполнения или оптимизации, а не самой природы функции.

Пояснение:

Пусть T_1 и T_2 — две машины Тьюринга.

Если обе вычисляют одну и ту же функцию f(x), но T_2 делает это быстрее, то T_2 не мощнее по вычислимой природе.

То же самое относится к ИИ: его «улучшение» — это лишь переход к другой реализации той же функции, возможно, с меньшей сложностью по времени или памяти.

То есть:

$$f \in \mathbb{F}_+ \Rightarrow f \notin \mathbb{F}_+$$
при оптимизации

6.3.1. Исчерпаемость архитектурного роста и предел усиления

В рамках Turing-вычислимых функций можно построить бесконечную последовательность модифицирующихся систем:

каждая система A_n порождает следующую, более сложную, более быструю или более оптимизированную.

Такая последовательность {A_n} допускает множество архитектурных преобразований — от перестроек памяти и глубины до адаптации эвристик или рекурсивных правил.

Однако при любом изменении остаётся неизменным один факт:

$$\forall n \in \mathbb{N}, A_n \in \mathbb{F}_{t}$$
.

Ни одна из таких модификаций не способна вывести систему за пределы вычислимого.

Все возможные «усиления» представляют собой внутренние перераспределения в счётном замкнутом множестве \mathbb{F}_{t} .

Это не рост, а движение внутри ограниченного пространства: вариации без принципиального расширения.

Такая замкнутость делает невозможным переход от формальной структуры к чемулибо «сверхразумному».

ИИ может меняться по форме, но не по классу.

Он может становиться сложнее, но не иным по природе.

Формально:

```
Если A_0 \in \mathbb{F}_{t} и A_n = A_{n-1}(A_{n-1}),
то \forall n \in \mathbb{N}: A_n \in \mathbb{F}_{t}.
```

Следовательно, ни одна цепочка самоприменения, ни одна архитектурная эволюция, ни одно рекурсивное усложнение не приведут к выходу за пределы Тьюрингвычислимости. Это ограничение не ресурсное — а логико-структурное.

Композиция синтаксических операций не рождает смысла.

Переход от алгоритма к субъекту невозможен по той же причине, по которой автомат не становится мыслящим.

ИИ может совершенствоваться, но не эволюционирует в разум.

Вывод:

ИИ ограничен, потому что вычислимость ограничена.

Всё остальное — инженерные варианты в замкнутой системе.

6.4 Принцип рефлексивного ограничения

Никакая система, работающая по определённым правилам, не может создать систему, не подчиняющуюся этим правилам, если её конструкция полностью ограничена изначальными механизмами.

Формально:

Если система S реализует функции в классе C, и S \rightarrow S' (новая система), то:

 $S' \in C$

Применительно к ИИ:

```
– Пусть ИИ_1 ∈ \mathbb{F}_{t}
– Он создаёт ИИ_2 \Rightarrow ИИ_2 ∈ \mathbb{F}_{t}
```

Следовательно, никакой последовательный процесс создания ИИ-ИИ-ИИ... не создаст существо, вышедшее за пределы Тьюринговской модели.

Вывод главы:

Гипотеза «взрыва интеллекта» — это не формальная теория, а метафора, основанная на ложной аналогии между увеличением вычислительной эффективности и расширением класса функций.

Математически доказуемо, что никакая последовательность улучшений, если она осуществляется внутри системы, подчинённой законам вычислимости, не может породить вычислительную структуру, превосходящую машину Тьюринга.

Таким образом, идея мета-ИИ или «сверхразума» не имеет логического обоснования и противоречит основам теории алгоритмов. Это — философская иллюзия, не подтверждаемая ни формальной логикой, ни эмпирическими данными.

Глава 7. Примеры невычислимых функций

До этого мы утверждали, что существуют функции, которые не могут быть вычислены ни машиной Тьюринга, ни какой-либо другой алгоритмической системой, включая

искусственный интеллект. Чтобы это утверждение не оставалось лишь абстрактным, рассмотрим несколько конкретных примеров функций, чья невычислимость доказана строго и математически.

Эти примеры служат не только теоретическим подтверждением существования границы вычислимости, но и демонстрируют, что даже простые по формулировке задачи могут оказаться принципиально неразрешимыми.

7.1. Проблема останова и функция Чёрча

Формулировка проблемы останова:

Пусть дана программа P и вход x. Требуется определить, остановится ли P при запуске на входе x.

Формально определим характеристическую функцию остановки:

```
Н(Р, х) = 1, если программа Р останавливается на х
```

H(P, x) = 0, если программа P зацикливается на x

Теорема Тьюринга (1936): не существует алгоритма, вычисляющего H(P, x) для всех P и x

Из неё следует, что функция Н не является тьюринг-вычислимой:

А значит, никакой искусственный интеллект, являющийся частью \mathbb{F}_{t} , не сможет вычислить $\mathsf{H}(\mathsf{P},\mathsf{x})$ в общем случае.

Как следствие, вводится функция Чёрча:

f(n) = 1, если n-я программа (в фиксированной нумерации) останавливается на входе n

f(n) = 0, если не останавливается

Эта функция — частный случай Н и также невычислима:

7.2. Диофантовы уравнения и теорема Матиясевича

Диофантово уравнение — это уравнение с целыми коэффициентами, решение которого ищется в целых числах.

Пример:

$$x^3 + v^3 + z^3 = 42$$

Вопрос: существует ли алгоритм, определяющий, имеет ли произвольное диофантово уравнение решение?

Ответ: не существует.

Теорема Матиясевича (1970), завершившая работу Дэвиса-Путнама-Робинсона:

Множество решений диофантовых уравнений совпадает с множеством перечислимых, но не обязательно разрешимых подмножеств натуральных чисел.

Отсюда: нет алгоритма, определяющего, имеет ли произвольное диофантово уравнение хотя бы одно целочисленное решение.

То есть, характеристическая функция:

D(E) = 1, если уравнение E имеет решение

D(E) = 0, если не имеет

- не является вычислимой: D ∉ Г ,

Следовательно, и эта функция также лежит вне пределов любого ИИ.

7.3. Функция Блума

Функция Блума — пример ограниченной тотальной функции, которая является непредсказуемой, несмотря на свою вычислимость.

Пусть $f: \mathbb{N} \to \mathbb{N}$ — тотальная вычислимая функция, чья скорость роста такова, что никакая программа длины $\leq n$ не может предсказать f(n).

Хотя сама функция вычислима, её непредсказуемость подчёркивает, что даже в классе $\mathbb{F}_{\mathfrak{t}}$ есть пределы понимания и прогнозирования поведения.

Это не пример невычислимости как таковой, но пример ограничения прогноза, который затрагивает и ИИ.

7.4. Истинные, но недоказуемые утверждения: теоремы Гёделя

Гёдель показал: в любой достаточно мощной формальной арифметической системе существуют утверждения, которые истинны, но невыводимы внутри этой системы.

Пусть S- формальная система (например, арифметика Пеано), и G- утверждение, сконструированное как:

G = «G не доказуемо в S»

Тогда, если S непротиворечива, G истинно, но G нельзя доказать в S.

Следовательно:

- Функция, определяющая истинность всех утверждений в арифметике: $T(\phi)$ = 1, если ϕ истинно; $T(\phi)$ = 0 иначе
- Не является вычислимой: Т ∉ Г ,

Это означает, что никакой ИИ не сможет «завершить» математику или «найти истину» за пределами формальной логики — это вычислительно невозможно.

Вывод главы:

Существует множество чётко определённых, математически обоснованных функций, которые не могут быть вычислены никакой машиной Тьюринга, а значит — ни одной

формой искусственного интеллекта.

Это не философское ограничение, а строго формализованная граница:

```
\exists f: \mathbb{N} \to \mathbb{N}, такие что f не вычисляется никаким алгоритмом
```

⇒ f не реализуем ни в одной архитектуре ИИ

Таким образом, любые заявления о всесильности или «неограниченности» искусственного интеллекта опровергаются уже сегодня - в рамках современной математики.

Заключение

На протяжении всей работы мы исследовали границы искусственного интеллекта не с позиции инженерных ограничений (мощность процессоров, объёмы памяти и пр.), а с фундаментальной точки зрения — исходя из природы вычислений как таковых. Центральным понятием здесь выступает вычислимость, формализуемая через машину Тьюринга и гипотезу Чёрча-Тьюринга.

Мы установили:

1. Любая алгоритмическая система, включая ИИ, представляет собой функцию из класса Turing-computable функций:

$$A(x) = f(x)$$
, где $f \in \mathbb{F}_t$

- 2. Этот класс ограничен. Он не включает:
 - Функцию останова: H(P, x)
 - Решение произвольных диофантовых уравнений: D(E)
 - Истинность недоказуемых утверждений: Т(ф)
- 3. Ни одна архитектура нейросеть, квантовая машина, гипотетический биокомпьютер — не преодолевает границу \mathbb{F}_{1} , если она алгоритмически реализуема.

Это означает: ИИ не способен выйти за пределы формальной вычислимости. Он не может стать "мета-вычислителем", не может создавать функции, которые невозможно вычислить, и не может предсказывать поведение систем, чья динамика лежит вне разрешимых границ.

Следовательно:

- Гипотеза о «взрыве интеллекта» логически несостоятельна, если под ним понимать качественный выход за границу Turing-вычислимого.
- Даже если ИИ будет способен к самоусовершенствованию, он будет это делать в пределах:
 - $\mathbb{F}_{\mathsf{t}} \subseteq \mathscr{F}$, где $\mathscr{F} -$ множество всех мыслимых преобразований ИИ
- А значит, никакой «надразумный интеллект» (в вычислительном смысле) невозможен: любой ИИ останется в рамках алгебраической машины, пусть даже гиперсложной.

В этом и заключается главная философская идея данной работы:

> Искусственный интеллект не есть иной вид мышления.

Это сложная, но конечная и алгоритмическая реализация человеческой идеи «вычисления».

Вывод главы

Рассмотренные примеры — функция останова, диофантовы уравнения, недоказуемые утверждения — представляют собой строго определённые, но принципиально невычислимые объекты. Они доказывают: существуют задачи, которые не решаются никакой машиной Тьюринга — а значит, и никаким искусственным интеллектом, если он алгоритмически реализуем.

Это не техническое ограничение. Это — математическая граница.

ИИ, как любая алгоритмическая система, не способен вычислить то, что невычислимо в принципе. Он может приближаться, адаптироваться, моделировать — но не пересечь предел вычислимости. И именно это определяет его фундаментальную границу: не скорость, не глубина, а природа самой модели.

В этом и заключается центральный вывод всей работы:

> Искусственный интеллект не является формой мышления.

Он — конечная, ограниченная, вычислимая реализация синтаксической процедуры.

Он — завершённая система, а не восходящая.

Он не станет разумом. Потому что разум — невычислим.

Глава 8. Предел искусственного интеллекта как формальный итог

Искусственный интеллект, независимо от своей архитектуры, модели обучения или уровня самоорганизации, остаётся вычислимой системой. Он реализует алгоритмически заданные преобразования и действует в рамках формально определённого класса функций. Этот класс — \mathbb{F}_{t} , множество всех функций, вычислимых машиной Тьюринга.

Никакая форма самоусовершенствования, даже рекурсивного, не даёт ИИ выхода за границу $\mathbb{F}_{\mathbf{t}}$. Он может усложнять структуру, увеличивать скорость, минимизировать ошибки — но остаётся внутри одного и того же формального пространства. ИИ не создает новые классы функций. Он только переупорядочивает известные.

Мы рассмотрели конкретные примеры функций, не входящих в $\,\mathbb{F}_{\,\,t}$:

- функцию останова H(P, x);
- характеристическую функцию решений диофантовых уравнений;
- функцию истинности недоказуемых формул в формальной арифметике.

Для всех этих случаев доказано: не существует алгоритма, вычисляющего их значения во всех общих случаях. Это значит, что не существует и такого ИИ, который смог бы решить их— если он, в свою очередь, подчиняется вычислимой структуре.

ИИ— не путь к новому типу мышления. Он— завершённый способ действия в пределах уже существующего формализма. Ни увеличение параметров, ни усложнение архитектур, ни коллективное обучение, ни комбинирование моделей не преодолевают барьер вычислимости. Природа этого барьера— математическая, а не

технологическая.

Следовательно:

ИИ не становится разумом, потому что разум не является вычислимым объектом.

ИИ не создаёт надразум, потому что не может выйти за пределы своей вычислимой структуры.

ИИ не может симулировать трансценденцию, потому что не способен породить то, что выходит за рамки алгоритма.

ИИ может усиливаться количественно, но не способен преодолеть качественную границу между вычислимым и невычислимым.

Финальный тезис:

> Искусственный интеллект — это форма вычислимости, доведённая до предела. Всё, что кажется выходом за неё, есть иллюзия перехода, порождённая внутри самой системы.

Приложение: Использованные источники и упомянутые работы

- 1. A. Тьюринг. On Computable Numbers, with an Application to the Entscheidungsproblem (1936)
- 2. A. Yëpy. An Unsolvable Problem of Elementary Number Theory (1936)
- 3. Ю. Гёдель. О формально неразрешимых предложениях Principia Mathematica и родственных систем (1931)
- 4. Ю. Матиясевич. Enumerable sets are Diophantine (1970)
- 5. М. Дэвис. Вычислимое: Введение в теорию алгоритмов
- 6. Ю. Манин. Математика как метафора
- 7. М. Девлин. Introduction to Mathematical Logic
- 8. С. К. Клебанов. Машины Тьюринга и пределы алгоритмического
- 9. В. Л. Васильев. Философия математики и основания вычислимости
- 10. Н. Джонсон-Лэйрд. Как мы рассуждаем
- 11. M. Минский. The Society of Mind
- 12. H. Бостром. Superintelligence: Paths, Dangers, Strategies (2014)
- 13. Т. Колмогоров, А. Н. Шень. Алгоритмы и случайность
- 14. Г. Чейтин. Algorithmic Information Theory

- 15. П. Скоулем. Формализация логики и ограниченность формальных систем
- 16. J. R. Searle. Minds, Brains, and Programs (1980)
- 17. G. J. Chaitin. Meta Math! The Quest for Omega (2005)

Работу выполнили: Рыбаков Павел Игоревич

Контактная информация: Электронная почта: pavel_rabota1996@mail.ru Я ВКонтакте: vk.com/id1059469430