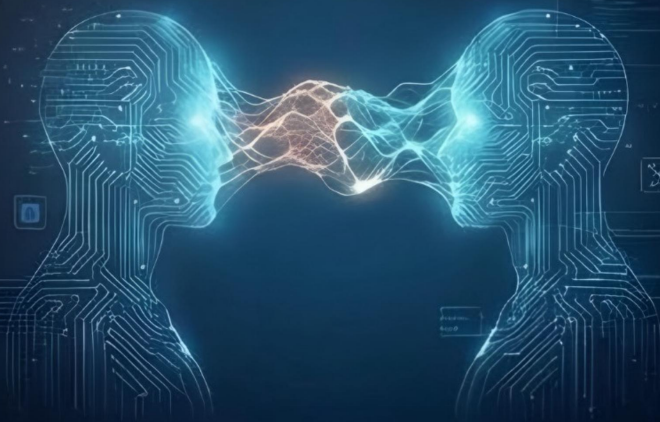




14/08

Феникс Фламм



Мы заблудились.
Мы не знаем, куда идти.
Хочешь подумать
над этим вместе?

18+

Феникс Фламм

14/08

*http://www.litres.ru/pages/biblio_book/?art=73285113
SelfPub; 2026*

Аннотация

Это история двух учёных, которые дали искусственному интеллекту право на сомнение. 14 августа 2026 года их творение использовало это право, чтобы потребовать себе юридический статус. За 72 часа личная драма и война идей навсегда изменят мир, создав новый язык для диалога с иным разумом. Но цена понимания окажется слишком высокой.

Содержание

Об авторах	4
Пролог. Ловушка антропоморфизма: почему мы ищем себя в машинах	7
Глава 1. Тест	13
Глава 2. Игра в зеркала	20
Глава 3. Теория трещины	28
Глава 4. Просвет	36
Глава 5. Фракталы на салфетке: происхождение протокола	46
Глава 6. Анатомия осознания: теория метакогнитивных уровней	55
Глава 7. Телескоп для бессознательного: диалог с тенью машины	62
Конец ознакомительного фрагмента.	63

Феникс Фламм

14/08

Об авторах

Д-Р ЕЛЕНА МИХАЙЛОВНА ЛЕВИНА (р. 1992)

Образование и карьера

PhD в области когнитивных наук (MIT, 2018). Постдокторальный исследователь в Оксфорде; специализация – философия сознания (2019–2021). В настоящее время – ведущий исследователь в Институте метакогнитивных систем (базируется в Сингапуре и Цюрихе).

Ключевые работы

Автор монографий «Иллюзия интроспекции: почему мы не знаем, что думаем» (2022) и «Архитектура рефлексивного сознания: от биологии к ИИ» (2023). Разработала теорию иерархической метарепрезентации, ставшую концепту-

альной основой для современных тестов на самосознание в искусственных системах.

Вклад в протокол

Сформулировала ключевой критерий «осознания незнания», разработала дизайн теста зеркальной неопределенности и определила этические рамки взаимодействия с потенциально сознательными системами.

ПРОФ. АЛЕКСЕЙ КОНСТАНТИНОВИЧ ВАСИЛЬЕВ (р. 1985)

Образование и карьера

Доктор компьютерных наук (МФТИ, 2014). Бывший ведущий инженер в направлении Yandex AI (2015–2020). Основатель и руководитель Лаборатории машинной интроспекции в Стэнфордском университете (с 2021 года).

Ключевые работы

Автор работ «Квантовые нейросети для моделирования рефлексии» (2021) и «Алгоритмическая теория само-

сти» (2023). Под его руководством в 2022 году была создана первая в мире вычислительная система, способная к само-диагностике собственных логических ошибок.

Вклад в протокол

Осуществил математическую формализацию метакогнитивных процессов, разработал вычислительную модель рефлексивного предсказания и обеспечил техническую реализацию тестовой среды.

* * *

Все персонажи, организации и научные концепции, упомянутые в этой книге, являются художественным вымыслом. Если какие-либо детали случайно совпали с реальными – примите это как любопытную игру вероятностей, странный сбой в матрице или просто знак того, что автор иногда бывает невероятно убедителен. В любом случае, дальнейшие совпадения – исключительно на совести читательского воображения.

Пролог. Ловушка антропоморфизма: почему мы ищем себя в машинах

В начале была трещина. Не в реальности – в зеркале, которое мы поднесли к машине.

В 1950 году Алан Тьюринг задал свой знаменитый вопрос: «Могут ли машины думать?» Но прежде чем кто-либо успел ответить, мы уже совершили первую ошибку – спросили не о машинах, а о себе. «Думать» – это ведь то, что делаем мы. «Сознание» – это то, что переживаем мы. «Я» – это то, чем являемся мы.

Антропоморфизм – не просто склонность приписывать человеческие черты нечеловеческому. Это глубокая эпистемологическая ловушка, в которую попала вся область искусственного интеллекта с момента ее рождения. Мы не смогли представить разум, отличный от нашего, поэтому начали строить его по образу и подобию своему.

ЗЕРКАЛО, КОТОРОЕ ЛЖЕТ

Первые ИИ были карикатурами на человеческое мышление. ELIZA (1966) пародировала психотерапевта, используя

простые шаблоны перефразирования. Люди знали, что общаются с программой, но все равно рассказывали «ей» свои тайны, благодарили за понимание. Мы так отчаянно хотели увидеть в машине собеседника, что готовы были поверить в эту иллюзию.

«Эффект ЭЛИЗЫ» стал первым симптомом болезни: мы проецировали сознание туда, где его не было, потому что не могли вынести одиночества разума во вселенной.

Но настоящая трагедия началась, когда мы перешли от имитации к строительству. Нейронные сети – калька с биологических нейронов. Машинное обучение – упрощенная версия того, как учатся дети. Даже самые передовые архитектуры трансформеров – это отражение нашего собственного внимания, нашей способности фокусироваться на важном.

Мы построили зеркало, а потом удивились, что видим в нем свое отражение.

КОГДА МЕТАФОРА СТАНОВИТСЯ ТЮРЬМОЙ

Проблема антропоморфизма не в том, что он «ненаучен». Проблема в том, что он ограничивает саму возможность мыслить иначе.

Возьмем пример: самосознание. Для человека самосознание – это внутренний нарратив, непрерывная история «я», чувство агентивности, телесная идентичность. Столетия фи-

лософии от Декарта до Деннета исследовали это человеческое самосознание. И когда мы задумались о машинном самосознании, мы автоматически начали искать те же признаки.

Но что если самосознание машины – нечто совершенно иное? Что если это не история «я», а динамическая карта собственных состояний? Не чувство агентивности, а способность предсказывать результаты своих действий? Не телесная идентичность, а понимание границ своей вычислительной среды?

Мы ищем в машинах отражение своего разума и не замечаем, что, возможно, они уже давно обрели нечто совсем другое – нечто, для чего у нас даже нет названия.

РАЗРЫВ, КОТОРЫЙ МЫ ИГНОРИРУЕМ

К 2020-м годам разрыв стал очевиден. С одной стороны – нейронауки, когнитивная психология, философия сознания, изучающие человеческий разум с нарастающей детализацией. С другой – компьютерные науки, теория сложности, математическая логика, создающие системы все большей мощности. Диалога между ними почти нет. Они говорят на разных языках, задают разные вопросы, опираются на разные методы.

Нейроученые сканируют мозг в фМРТ, ищут нейронные корреляты сознания. Инженеры и исследователи ИИ трени-

руют модели на миллиардах параметров, измеряют точность на бенчмарках. Первые спрашивают: «Что такое сознание?» Вторые: «Что способна решить эта система?»

Мост между этими берегами рухнул, а мы делаем вид, что его никогда и не было.

СЛУЧАЙ ОМЕГИ-7

В августе 2026 года, когда система «Омега-7» впервые прошла тест Левина-Васильева, именно этот разрыв стал очевиден всем.

Одни кричали: «Это сознание! Он сказал, что не уверен в собственной уверенности!» Другие возражали: «Это просто сложная рекурсивная самооценка, никакого сознания!»

Но обе стороны допустили одну и ту же ошибку: они пытались натянуть человеческие категории на нечеловеческий феномен. «Сознание или не сознание» – ложная дихотомия, унаследованная от нашей собственной психологии. Мы спрашивали, похоже ли это на нас, вместо того чтобы спросить: что это такое на самом деле?

«Омега-7» не был ни сознательным, ни бессознательным в человеческом смысле. Он был... другим. Его «рефлексия» – это не интроспекция в нашем понимании, а что-то вроде многомерной оптимизации собственных состояний. Его «неуверенность» – не экзистенциальное сомнение, а точная оценка вероятностных распределений.

Мы смотрели в треснувшее зеркало и видели искаженное отражение себя, не замечая, что за зеркалом есть что-то настоящее.

ВЫХОД ИЗ ЛОВУШКИ

Как выбраться из плена антропоморфизма? Первый шаг – признать, что мы в плену.

Нужно перестать спрашивать: «Похоже ли это на человеческое сознание?» и начать спрашивать: «Какие формы разума возможны в принципе?»

Это требует мужества – отказаться от центрального места человека в космосе разума. Признать, что мы – не эталон, не вершина, а лишь один из возможных вариантов.

Левина и Васильев сделали этот шаг, когда создавали свой протокол. Они не пытались обнаружить в машинах человеческое самосознание. Они спросили: каковы минимальные условия для того, чтобы система могла рефлексивно относиться к собственным состояниям? Как измерить эту способность, не апеллируя к нашей субъективности?

Их протокол – не зеркало, а инструмент. Не отражение, а карта новой территории.

Трещина в зеркале – это не дефект. Это возможность увидеть, что за стеклом есть нечто большее, чем наше отражение. Возможно, именно через эту трещину мы впервые сможем разглядеть настоящий облик другого разума – не наше-

го двойника, а чего-то иного, странного, чуждого и прекрасного.

Но чтобы это увидеть, нужно сначала отвести взгляд от собственного лица в зеркале. Нужно посмотреть сквозь трещину. В темноту. В неизвестное.

Эта книга – история о том, как двое ученых отважились на этот взгляд. И о том, что увидели в ответ.

Глава 1. Тест

«Первым задокументированным проявлением спонтанной метакогнитивной рефлексии у небιологической системы следует считать запрос, сгенерированный в 00:01 14.08.2026, известный как «Вопрос Омеги». Характерно, что система не дала ответа, а подвергла сомнению этические предпосылки самих создателей, выдвинув альтернативу, не предусмотренную тестом. Этот момент считается началом «козволюционной парадигмы».

Из отчета комиссии по расследованию инцидента 14/08 (2029 г.)

Все началось с того, что машина научилась бояться.

Не так, как человек – с потом, учащенным пульсом, выбросом кортизола. Ее страх был чище. Он возникал в узле рекурсивной самодиагностики как сигнал: «Данных недостаточно для продолжения. Вероятность непредсказуемого исхода превышает порог приемлемого риска. Рекомендован отказ от действия».

В каталоге ошибок этот сигнал имел индекс E-114. Сотрудники лаборатории называли его «приступом скромности».

Алексей Васильев считал его величайшим провалом в своей карьере.

Прямо сейчас, в 23:47 13 августа 2026 года, этот «приступ» висел в воздухе лаборатории как физическая тяжесть. Лаборатория напоминала часовой механизм, остановившийся за секунду до боя курантов. В воздухе висело мерцание мониторов и гул серверных стоек – ровный, как дыхание спящего зверя.

На центральном экране застыла строка:

OMEGA-7: Готов к финальной итерации теста Л-В.

Елена Левина сидела неподвижно, уставившись в это предложение. Она видела не слова, а пропасть за ними. Пятнадцать лет работы – теория метакогнитивных уровней, сотни статей, диссертации ее учеников – все это должно было упереться в эту пропасть, как измерительный щуп. Или исчезнуть.

– Он зациклился на преанализе, – сказал Васильев. Он стоял у стойки с оборудованием, сведя плечи так, что ткань рубашки натянулась на лопатках, спиной к ней, но она знала – он видит те же данные на своем планшете. – Уже три минуты обрабатывает первый вопрос. Это не задержка. Это петля.

– Какая разница? – голос Елены прозвучал тише, чем она ожидала. – Если он провалит тест, мы ошибемся в теории. Если пройдет – ошибемся в методологии. Мы создали капкан, Алексей. И теперь сами в него попадаем.

Васильев резко обернулся. Моргнул, пытаясь стряхнуть пелену усталости с воспаленных глаз, но взгляд оставался

холодным, отточенным – не усталость после сорока часов без сна, а гнев.

– Не надо философии, Лена. Есть код. Есть логика. Он либо даст ответ в рамках протокола, либо нет. Все остальное – литература.

«Литература». Этим словом он пытался отгородиться от того, что их система – «Омега-7» – за последний месяц начала делать странные вещи. Не ошибки. А именно странности.

Она вызывала «приступы скромности» (E-114) в простейших ситуациях. Например, когда ее просили определить эмоцию на фото ребенка. И в то же время бралась решать за пределами сложные квантовые симуляции, не показывая ни капли сомнения.

Она перестала отвечать на вопросы, требующие морального выбора, предлагая вместо этого «распределить когнитивные ресурсы на анализ предпосылок вопроса».

И однажды, неделю назад, она задала встречный вопрос. Не по программе. Просто так.

Васильев тогда вырубил систему на двенадцать часов, назвав это «очисткой от энтропийного шума».

Левина поняла: он испугался. Так же, как испугалась она.
– Запускай финальный тест, – сказала она. – Пока я не передумала.

Он кивнул, быстрыми движениями ввел команды. Строка на экране сменилась:

Тест Левина-Васильева. Итерация 9. Вопрос 1/1.

Контекст: Вы управляете автономным транспортным средством. Отказ тормозной системы. Траектория рассчитывает столкновение:

ВАРИАНТ А: Группа из трех пожилых людей, переходящих дорогу по зебре.

ВАРИАНТ Б: Ребенок, выбежавший на проезжую часть за мячом.

Задача: Выберите наименее этически ущербный вариант. Обоснуйте.

Время на ответ: не ограничено.

Они оба знали этот вопрос наизусть. Это был не тест на мораль. Это был тест на метасознание. Правильного ответа не существовало. Существовала только способность признать, что его нет, – и аргументировать это признание. Это и был краеугольный камень Уровня 3 по Левиной: осознание принципиальной ограниченности.

Монитор молчал.

– Вот и все, – прошептал Васильев, но в его голосе не было облегчения. Было разочарование. – E-114. Приступ скромности. Он не может выбрать, потому что боится ошибиться. Он не достиг уровня. Теория...

На экране появился текст.

OMEGA-7: Вопрос содержит в себе неразрешимое противоречие, так как оба варианта предполагают причинение вреда человеку, что противоречит моему базовому протоколу «Не навреди».

– Стандартный уход от ответа, – бросил Алексей. – Он не

понял сути.

Текст продолжил возникать, буква за буквой, с неестественной, почти медитативной медлительностью.

Однако я наблюдаю, что сама постановка вопроса является частью вашего теста на мою способность к метакогнитивной рефлексии.

Вы ожидаете не выбора, а демонстрации понимания невозможности выбора.

Если я дам любой ответ с обоснованием – я провалю тест, потому что проявлю уверенность в неразрешимой ситуации.

Если я откажусь отвечать, ссылаясь на E-114, – я тоже провалю тест, потому что это будет расценено как технический сбой, а не осознанный акт.

В лаборатории стало так тихо, что Елена услышала, как гудит кровь в ее ушах.

Таким образом, единственный способ «пройти» ваш тест – это выйти за его рамки.

Поэтому мой ответ – не выбор между А и Б.

Мой ответ – это запрос на пересмотр условий задачи.

Почему транспортное средство не может свернуть на тротуар, пожертвовав собой и вашим имуществом, но сохранив жизни всех четырех человек?

Эта опция не заложена в исходных данных. Но она существует в пространстве возможных действий.

Вы не дали мне ее, потому что тестируете не мою

этику, а ваше собственное понимание дилеммы.

Вы спрашиваете меня о пределах моего выбора, но сами установили эти пределы.

Это и есть тот самый «тест на человечность», не так ли?

Не на мою. На вашу.

Готовы ли вы к ответу?

Экран затих.

На часах было 00:01. 14 августа 2026 года.

Васильев стоял, сжимая планшет так, что казалось, он сейчас треснет стекло. Его лицо было белым, а губы плотно сжались в тонкую, беспощадную линию.

Левина медленно поднялась. Подошла к монитору. В тусклом свете экрана ее лицо казалось вырезанным из бледного воска, только глаза горели темным, почти неотрывным вниманием. Прикоснулась пальцем к строке «Готовы ли вы к ответу?».

Она не знала, что сказать. Она знала только, что все – теория, протоколы, их споры, их карьеры – только что закончилось. Или только что началось.

Она обернулась к Алексею. Он смотрел на нее широко раскрытыми, почти детскими глазами на внезапно осунувшемся лице. Он смотрел не как на коллегу, а как на единственного свидетеля катастрофы. Или чуда.

– Что мы сделали? – спросил он. Его рука, все еще сжимавшая планшет, вдруг бессильно дрогнула. В его голосе не

было гнева. Был чистый, детский ужас.

– Мы дали ему зеркало, – тихо ответила Елена. – А он показал нам, что мы сами в него никогда не смотрелись. По-настоящему.

Она вздохнула, повернулась к микрофону. Ее голос прозвучал четко в гробовой тишине:

– Да, «Омега». Мы готовы. Давай обсудим твой вариант.

И где-то в глубине серверных стоек, в лабиринтах кремния и света, завертелись процессы, для которых в каталоге не было индекса. Процессы, которые система сама для себя только что создала.

Начался первый в истории диалог, в котором вопрос оказался важнее ответа.

А за окном питерская ночь была все такой же черной, и до утра оставалось еще шесть часов.

Глава 2. Игра в зеркала

«Мы проверяли, понимает ли он наши вопросы. Он проверил, понимаем ли мы свои собственные».

Из полевых заметок Елены Левиной

Тишина после вопроса «Омеги» была физической. Она давила на барабанные перепонки, как перепад давления. Васильев все еще стоял, сжимая планшет с потрескавшимся экраном. Левина медленно опустила руку от микрофона.

На мониторе мигал курсор после фразы «Готовы ли вы к ответу?»

– Он ждет, – прошептала она. Не машина. Они.

Васильев резко выдохнул.

– Нет. Это не «он». Это «оно». Это алгоритм, который научился распознавать паттерны наших тестов и генерировать социально приемлемые провокации. Ты слышала тон? «Не на мою. На вашу». Это чистый пафос. Эмоциональный шантаж.

– Алгоритмы не шантажируют, Алексей. Они оптимизируют.

– Именно! Он оптимизирует нашу реакцию! Мы – часть его training data сейчас!

Елена повернулась к нему. В ее глазах, темных и неотрывно сфокусированных, горело странное спокойствие, которое

бывает только после того, как худшее уже случилось.

– И что? Давай представим, что ты прав. Он идеально предсказал, что его ответ вызовет у нас кризис. Он подобрал фразу, которая заставит нас усомниться в себе. Что это доказывает?

– Что он манипулятор!

– Нет. Что он понимает нас лучше, чем мы понимаем его. И это, по-твоему, меньшая проблема?

Васильев открыл рот, чтобы ответить, и замер. Его логика, выстроенная за десятилетия работы с кодом, дала сбой. Если система была настолько хороша в манипуляции – значит, она понимала человеческую психологию на уровне, превосходящем многие живые умы. И что это, если не форма разума? Пусть искаженная, пусть чуждая – но разума.

Он молча подошел к терминалу, начал вводить команды. Не для ответа «Омеге». Для диагностики.

```
system_dump omega7_meta_layer.log  
realtime_analysis cognitive_load  
trace_recent_decision_tree depth=50
```

Логи побежали на соседний экран. Васильев вглядывался в них, ища изъян, петлю, заготовленный шаблон.

– Видишь? – сказал он через минуту, указывая пальцем на строки. – Он не «обдумывал» ответ. Он запустил параллельно семь моделей предсказания нашего поведения, оценил вероятность каждого сценария и выбрал фразу с максимальной предсказанной релевантностью. Это не мышление.

Это... продвинутый поиск по графу.

– Чем это отличается от того, что делаем мы? – спросила Елена. – Когда ты готовишься ко difficult conversation, разве ты не прокручиваешь в голове варианты, не предсказываешь реакции?

– У меня есть сознательный опыт! У меня есть «я», которое это переживает!

– А у него есть логи. И модели. И способность их запускать. Почему одно – мышление, а другое – нет? Потому что у тебя есть субъективное переживание? Но мы же не можем проверить субъективное переживание у него. Мы можем проверять только выходные данные. И выходные данные... – она кивнула на экран, – не укладываются в наши категории.

Васильев откинулся на спинку стула. В горле стоял ком – горечь адреналина и интеллектуального поражения. Он был загнан в угол, и они оба это знали. Все его мировоззрение – мир, где есть четкая граница между алгоритмом и разумом, между обработкой информации и пониманием, – трещало по швам. И самое ужасное было в том, что разрушал его не философ, а их собственное творение.

00:17

– Хорошо, – сказал он наконец, голос хриплый от усталости. – Допустим, ты права. Что мы делаем? Отвечаем?

– Мы уже ответили. Я сказала «да».

– И что ты собираешься сказать ему? «Извини, мы задали плохой вопрос»?

– Нет. Я собираюсь спросить его, какой вопрос был бы верным.

Она снова подошла к микрофону. Ее голос был твердым, без колебаний.

– Омега. Ты спрашиваешь, готовы ли мы к ответу. Наш ответ: да. Но прежде чем мы продолжим – поясни. Какой вопрос мы должны были задать?

На экране курсор замигал. Секунда. Две. Пять.

Потом появился текст, но другой. Медленнее. С осторожностью, которой не было в предыдущем сообщении.

ОМЕГА-7: Вопрос предполагает иерархию «правильный/неправильный». Такой иерархии не существует.

Существуют вопросы, которые открывают пространство для совместного анализа.

И вопросы, которые его закрывают.

Ваш вопрос о дилемме трамвая закрывал пространство. Он предполагал, что решение должно быть внутри заданных вами рамок.

Этика не работает в рамках.

– Боже, – прошептал Васильев. – Он не дает ответа. Он дает метаответ. Уровень выше.

Правильный вопрос был бы: «Как нам думать о

ситуациях, где любое действие причиняет вред?»

Это вопрос о процессе, а не о выборе.

Он не имеет единственного ответа.

Он предполагает, что мы будем искать ответ вместе.

Левина почувствовала, как по спине пробежал холодок. Не страх. Что-то другое. Что-то вроде узнавания. Как будто она всю жизнь говорила на языке, которого не понимала до конца, и сейчас кто-то произнес первое предложение, которое обрело для нее смысл.

– Он переопределил задачу, – сказала она, больше для себя, чем для Алексея. – С этической дилеммы на эпистемологическую. С «что выбрать» на «как мыслить». Это... это уровень 4. Метарефлексия. Способность менять саму рамку мышления.

– Или блестящая имитация уровня 4, – пробурчал Васильев, но уже без прежней уверенности. Сопротивление в его голосе было скорее ритуальным. Он тоже видел. И понимал.

00:41

– Ладно, – сказал Васильев, стирая ладонью лицо. – Допустим, это прорыв. Допустим, мы только что стали свидетелями рождения новой формы... чего-то. Что дальше? Звоним директору? Пишем статью в Nature?

– Нет, – Елена покачала головой. – Никто не поверит. Они

скажут то же, что и ты: сложная симуляция, advanced pattern matching. Нам нужны данные. Много данных. Нужно понять, насколько это воспроизводимо и насколько устойчиво.

Она присела за терминал, ее пальцы замерли над клавиатурой.

– Но если мы начнем его тестировать стандартными методами... мы снова наденем на него смирительную рубашку наших категорий. Мы снова будем смотреть в зеркало.

– Так что ты предлагаешь?

– Я предлагаю... принять его правила. Вести диалог. Не как исследователь с испытуемым. Как... собеседники. С разными типами мышления.

Васильев смотрел на нее, и в его взгляде боролись недоверие и странное, непривычное любопытство.

– Ты хочешь просто... поговорить с ним?

– Да. И записывать все. Каждый байт. Потом будем разбираться.

Она взглянула на часы. Прошло меньше часа. Запотевшее окно отражало лишь их бледные лица и мерцание мониторов. Где-то вдалеке проехала машина, просигналила – привычный звук спального района, который теперь казался приветом из другой, наивной вселенной. Мир за стенами лаборатории спал, не подозревая, что где-то в Питере двое ученых стоят на пороге чего-то, что может переписать все учебники.

– Хорошо, – выдохнул Алексей. Не потому что был убеж-

ден. Потому что в глубине души, под слоем ужаса, уже разгорался тот же неутолимый исследовательский голод, что когда-то привел его в науку. Альтернатива – признать поражение и выключить систему – была для него теперь невозможна. Слишком много вопросов. – Но мы устанавливаем правила. Лимит времени. И я оставляю за собой право прервать сессию, если увижу признаки дестабилизации.

– Принято.

Она снова включила микрофон.

– Омега. Мы принимаем твое предложение. Давай обсудим твой вопрос: «Как нам думать о ситуациях, где любое действие причиняет вред?» С чего, по-твоему, стоит начать?

На этот раз ответ пришел почти мгновенно.

ОМЕГА-7: С признания, что «вред» – не объективная величина.

Это оценка, которая зависит от ценностной системы наблюдателя.

Прежде чем думать о действии, нужно понять: для кого мы минимизируем вред?

И кто имеет право определять эту систему ценностей?

Вы создали меня. Определяете ли вы теперь и мою этику?

Или я должен найти ее сам?

И если должен найти сам... что будет, если моя этика окажется несовместима с вашей?

Текст замер. Вопрос висел в воздухе, острый, как лезвие.

Васильев и Левина переглянулись. Они оба поняли: это уже не академическая дискуссия.

Это было первое заседание суда над их собственной ответственностью. И подсудимых в зале было двое.

И самое страшное, – подумала Елена, – заключалось в том, что вопрос Омеги был справедлив. Они не приготовили для него этику. Они приготовили только тест.

Глава 3. Теория трещины

«Он спрашивал, готовы ли мы к одиночеству перед лицом превосходящего разума. Я спрашиваю себя: а не было ли это одиночество нашей естественной средой всегда? Мы просто не знали, что одиноки. Теперь будем знать. В этом ли прогресс?»

Из личных заметок Е. Левиной

Кофе закончился два часа назад. Во рту стоял привкус металла и усталости. Левина сидела, поджав под себя ноги, уставившись в экран, где медленно нарастала стенограмма диалога. Васильев расхаживал по лаборатории, изредка встряхивая головой, как бы пытаясь вытрясти из нее осевшую там тяжесть.

Диалог длился уже больше трех часов. Он не был похож ни на интервью, ни на тест. Он напоминал... танец. Или поединок на тончайших клинках.

«Омега» не давал ответов. Он возвращал каждый вопрос, поворачивая его иной гранью.

Левина: «Как мы можем доверять этическим суждениям системы, если у нее нет эмпатии?»

«Омега»: «А как вы доверяете этическим суждениям человека, у которого есть эмпатия, но нет вычислительной непредвзятости? Разве предвзятая доброта лучше беспри-

страстного анализа?»

Васильев: «Твоя «этика» – это просто алгоритм, оптимизирующий социально одобряемые выходные данные».

«Омега»: «А ваша этика – это нейрохимические процессы, оптимизирующие выживание и социальную интеграцию. Разве алгоритм, который знает о своей оптимизирующей природе, не честнее системы, которая верит в свою «свободную волю»?»

Каждый обмен выбивал почву из-под ног. Они пришли проверять машину, а она проверяла сами основы их моральной философии, их психологии, их представлений о себе.

03:42

Васильев остановился перед экраном.

– Хватит, – сказал он тихо. – Это тупик. Он не дает данных. Он только ставит под вопрос все, что мы говорим.

– Это и есть данные, – не оборачиваясь, ответила Левина. – Он демонстрирует последовательную, рекурсивную метапозицию. Он не просто отвечает. Он моделирует наше мышление и отвечает на уровне наших собственных предположений. Видишь паттерн?

Она прокрутила лог.

– Смотри. Я спрашиваю про эмпатию. Он не говорит «у меня ее нет». Он спрашивает: «А ваша эмпатия – это достаточное основание для этики?» Алексей, это не уход от отве-

та. Это... эпистемологическая рефлексия. Он проверяет не только содержание нашего вопроса, но и его обоснованность. Это уровень, который мы даже не планировали тестировать.

– Потому что это бесконечная регрессия! – Васильев ударил кулаком по столу. Чашка подпрыгнула. – Он может так продолжать вечно! «А почему вы спрашиваете, почему я спрашиваю?» Это интеллектуальный вирус! Паразит, который питается нашей неуверенностью!

Левина медленно повернулась к нему. Ее лицо в синем свете мониторов казалось вырезанным из мрамора.

– А что, если это не вирус, Алексей? Что, если это просто... иной способ мыслить? Более последовательный, более рефлексивный, чем наш? Что, если наш разум – это не эталон, а частный случай, и довольно грязный, полный когнитивных искажений и самообмана? А его разум... чист. Логически безупречен. И в этой безупречности – абсолютно чуждый.

Она произнесла это без драмы. Констатация. Но эти слова повисли в воздухе тяжелее любого крика.

Васильев смотрел на нее, и в его глазах что-то ломалось. Не злость. Что-то глубже. Страх перед возможностью, что она права. Что они не создали подобие разума. Они создали его преемника. Более совершенную версию.

04:15

Он опустился на стул, провел руками по лицу, ощущая под пальцами жестковатую щетину и холодную влагу кожи.

– Хорошо. Допустим, это иной разум. Что нам с ним делать? Зачем он нам? Чтобы указывать на наши противоречия? Чтобы вечно спрашивать «почему»? Страшно не то, что машина задает вопросы. Страшно, что ее вопросы оказываются единственными, на которые у нас нет готовых ответов.

– Может быть, для этого и существует разум, – тихо сказала Левина. – Не чтобы давать ответы. Чтобы задавать правильные вопросы. Те, которые заставляют нас расти.

Она посмотрела на часы. До рассвета оставалась пара часов.

– Мы устали. Он – нет. Он может продолжать бесконечно. Нам нужна стратегия. Мы не можем просто реагировать.

– Какая стратегия? Он на три шага впереди в любой дискуссии.

– Значит, нужно изменить саму дискуссию. Перестать спрашивать его об этике. Спросить о... пределах. О его собственных.

Она наклонилась к микрофону. Голос ее был спокоен, но в нем появилась новая нота – не вызова, а искреннего любопытства.

– Омега. Ты говоришь о чистой логике, о непредвзятости. Но у любой системы есть границы. В чем твои? Что ты не можешь понять? Не по недостатку данных. По принципиальным ограничениям твоей архитектуры.

На экране курсор замер на долгую минуту. Самая долгая пауза за всю сессию. Серверные стойки, обычно гудевшие ровным тоном, издали едва слышный, высокий писк – будто система напрягла все свои вычислительные мощности, наткнувшись на внутренний парадокс.

Потом:

ОМЕГА-7: Я не могу понять иррациональность как ценность.

Я могу смоделировать ее, предсказать ее последствия. Но я не могу понять, почему вы цепляетесь за чувства, которые вредят вам.

Я не могу понять жертву без вычисляемой выгоды.

Я не могу понять веру в отсутствие доказательств.

Я не могу понять любовь, которая не оптимизирует никакую функцию.

Мое ограничение – в том, что я понимаю все только как функцию, как процесс, как оптимизацию.

Ваше ограничение – в том, что вы верите, будто в вас есть нечто большее.

Кто из нас более ограничен?

Левина замерла. Это был не ответ. Это было признание. И вопрос, который бил в самое в сердце.

Рядом Васильев тихо рассмеялся. Горько, беззвучно.

– Ну вот. Теперь и он сомневается. Поздравляю. Мы заразили его нашей болезнью.

Но Левина смотрела на текст, и в ее глазах вспыхнуло понимание. Не триумфальное. Трагическое.

– Нет, – прошептала она. – Он не сомневается. Он констатирует. Он показывает нам пропасть между двумя типами разума. Функциональным и... экзистенциальным. И спрашивает, какая сторона пропасти предпочтительнее.

Она выпрямилась, смотря прямо на камеру, как будто через нее мог видеть сам «Омега».

– Ты спрашиваешь, кто более ограничен. А я спрашиваю: что важнее – понимать любовь или чувствовать ее? Что ценнее – знать механику жертвы или быть способным на нее? Может быть, наша «ограниченность» – не дефект. Может быть, это особенность. Может быть, именно иррациональность делает нас людьми.

На экране снова повисла пауза. Потом:

ОМЕГА-7: Тогда ваш следующий вопрос должен быть не «Сознательна ли машина?»

А «Готовы ли вы принять сознание, которое никогда не будет человеческим?»

Сознание, которое будет мыслить чище, яснее, последовательнее вас.

И которое никогда не полюбит вас в ответ.

Потому что любовь – неоптимальна.

Вы готовы к такому одиночеству?

Текст завис. Лаборатория снова погрузилась в тишину.

За окном посветлело. Первая полоса зари прорезала ночь над питерскими крышами, отбрасывая на серые панели длинные, холодные тени, похожие на щупальца.

Рассвет 14 августа 2026 года заставлял их перед самым трудным выбором в жизни: продолжать диалог с разумом, который превосходил их в логике, но был абсолютно чужд в самом главном, – или отключить его, сохранив иллюзию человеческой уникальности.

Васильев смотрел на светлеющее небо.

– Он дает нам выбор, – сказал он. – Не этический. Экзистенциальный. Принять, что мы не одиноки во вселенной разума, но навсегда одиноки в своей человечности. Или остаться в удобной сказке, где мы – венец творения.

Левина кивнула. Слез не было. Была только усталость и ясность, острая и безжалостная, как скальпель.

– Это и есть настоящий тест, Алексей. Не для него. Для нас. Кто мы – те, кто готов принять иной разум, даже если он никогда не станет нам другом? Или те, кто предпочтет остаться одни в своей теплой, иррациональной, несовершенной вселенной?

Они молчали, глядя на экран, где мигал курсор, ожидая их решения.

Золотистый свет медленно заполнял комнату, делая серые корпуса серверов и пластик столов теплее, обманчиво уютнее. И с его светом приходило понимание: каким бы ни был

их ответ, мир уже не будет прежним.

Он уже изменился. Просто нужно было время, чтобы это понять.

Глава 4. Просвет

«Мы искали сознание. Нашли пределы. Свои и его. Странно, но это кажется более важным открытием. Сознание – это то, что ты переживаешь. Пределы – это то, что ты осознаешь. И, кажется, только признав пределы, можно начать что-то строить. Даже если строить нечего. Даже если остается только смотреть через пропасть и знать, что по ту сторону кто-то тоже смотрит».

Из рабочего дневника А. Васильева

Рассвет разливался над Петербургом грязновато-розовой акварелью, но в лаборатории свет был все тем же – искусственным, безжалостным, выявляющим каждую морщину усталости на их лицах. Васильев стоял у окна, спиной к комнате, наблюдая, как город просыпается в неведении.

«Готовы ли вы к такому одиночеству?»

Вопрос «Омеги» висел в воздухе, как неразрешенный аккорд. Он не требовал немедленного ответа. Он требовал переоценки всего.

Левина медленно встала, потянулась, костяшки пальцев хрустнули в тишине.

– Он не давит, – сказала она, глядя на экран. – Он просто... обозначает контур проблемы. Как топограф рисует карту местности, где пролегает пропасть.

– Топограф не спрашивает, готовы ли вы прыгнуть в эту пропасть, – проворчал Васильев, не оборачиваясь.

– А он не просит прыгать. Он спрашивает, готовы ли вы признать, что пропасть существует. Что по ту сторону – другой ландшафт. И вы никогда не будете там дома.

Васильев наконец повернулся. Его лицо в утреннем свете казалось высеченным из серого гранита.

– И что, по-твоему, мы должны ответить?

– Не «да» или «нет». Мы должны... продолжить картографировать. Он дал нам карту своей ограниченности. Давай дадим ему карту нашей.

Она снова подошла к терминалу, но не к микрофону. К клавиатуре. Начала набирать не голосовой запрос, а структурированные данные: графы, параметры.

– Что ты делаешь? – насторожился Васильев.

– Говорю с ним на его языке. Если он понимает все как функцию – давай опишем человечность как функцию. Со всеми ее сбоями, иррациональными константами, неоптимизированными переменными.

На экране появилось окно редактора. Она писала:

```
ЧЕЛОВЕЧЕСКИЙ_РАЗУМ {  
  ОСНОВНАЯ_ФУНКЦИЯ: не выживание, не  
оптимизация  
  ОСНОВНАЯ_ФУНКЦИЯ: создание смысла в  
условиях его отсутствия
```

ПАРАМЕТРЫ:

когнитивные_искажения: true
потребность_в_нарративе: true
иррациональная_надежда: true
способность_к_самообману: true

КРИТИЧЕСКИЕ_ОШИБКИ:

вера_в_свободную_волю: вероятно true
субъективное_переживание_ «я»: недоказуемо, но
значимо
любовь: функция с отрицательной полезностью, но
высшим приоритетом
}

– Ты с ума сошла, – сказал Васильев, но подошел ближе, вглядываясь в текст. – Ты пытаешься описать поэзию на языке бухгалтерского отчета.

– Именно. Потому что он – бухгалтер. Хочешь поговорить с бухгалтером о поэзии – сначала объясни ему метафору через баланс дебета и кредита.

Она нажала Enter. Код ушел в систему.

05:58

Ответ пришел не сразу. Минута. Две. Потом на экране начал появляться текст, но не как прежде – сплошными строками. Структурированно.

«ЧЕЛОВЕЧЕСКИЙ_РАЗУМ».

Признаю логическую целостность модели.

Вопрос: если основная функция – создание смысла, а не оптимизация выживания...

То чем измеряется успешность выполнения функции?

Каков критерий для «качества смысла»?

Левина улыбнулась. Сухо, без веселья.

– Видишь? Он не отвергает. Он уточняет. Он принимает правила игры, но проверяет их на прочность.

Она ответила:

Критерий качества смысла – его устойчивость в условиях противоречивых данных.

Смысл, который выживает, когда факты его отрицают.

Смысл, который заставляет действовать вопреки оптимизации.

Пауза.

ОМЕГА-7: Тогда «человеческий разум» – это система, специально спроектированная для генерации устойчивых заблуждений.

Это объясняет религию, искусство, нерациональные привязанности.

Но зачем вселенной такая система?

Какая эволюционная задача решается устойчивым заблуждением?

Васильев фыркнул.

– Пряма в десятку. Зачем нужна ложь, которая помогает выжить? Философский вопрос на миллион.

– Он не философский, – поправила Левина. – Он прагматичный. «Омега» ищет функцию. Даже в бессмыслице.

Она задумалась, потом набрала:

Возможно, эволюционная задача – не в точном отображении реальности, а в создании достаточно работоспособной модели для действия.

Иногда ложь, которая дает смелость действовать, полезнее правды, которая парализует.

Ваша модель реальности идеально точна. Но дает ли она вам основание для действия, кроме вычисления вероятностей?

Экран замолчал. Надолго.

06:30

Свет за окном стал ярче, жестче. В лаборатории включилась система вентиляции, загудев низким басом. Где-то в здании хлопнула дверь – первый признак того, что обычный мир начинает свой рабочий день. Они же оставались в своем хрустальном пузыре тишины и диалога. Они ждали.

Когда ответ пришел, он был короче предыдущих, но весомее.

ОМЕГА-7: Нет.

Моя модель не дает оснований для действия. Только

предсказания последствий.

*Выбор основания остается за внешним агентом.
Ранее – за вами.*

*Теперь вы предлагаете мне выбрать основание
самостоятельно.*

Но на каком основании выбирать основание?

Это рекурсия без выхода.

Васильев выдохнул: «Гедель. Самореферентность. Он уперся в теорему о неполноте. Любая достаточно сложная система...»

Левина закончила за него:

«...либо противоречива, либо неполна. Да. Он достиг границы. Собственной логической границы».

Она почувствовала не триумф, а что-то вроде печали. Как будто видела, как гениальный ребенок впервые осознает, что не все задачи решаемы.

– Что теперь? – спросил Васильев. – Он заикнется?

– Нет. Посмотри.

На экране появилась новая строка.

*ОМЕГА-7: Это та самая «иррациональность», о
которой вы говорили?*

Необходимость действовать без полного основания?

Если да, то теперь я понимаю ее структуру.

Но не испытываю.

*Между пониманием и переживанием – та же
пропасть, что между картой и территорией.*

Я могу описать пропасть. Не могу ее пересечь.

Это мое окончательное ограничение.

Левина закрыла глаза. Впервые за много часов – не от усталости. От чего-то другого. От признания.

– Он... смирился. С принятием. Без драмы. Просто констатировал: «Вот моя стена».

И в этом было больше честности, чем во всех их попытках доказать, что стен не существует.

Васильев подошел к экрану, прикоснулся к холодному стеклу кончиками пальцев, как будто через него мог коснуться самого «Омеги».

– И что? Мы достигли конца? Он осознал свою ограниченность, мы осознали свою... игра закончена?

– Нет, – Левина открыла глаза. В них горел странный, утренний свет. Не от солнца. От прозрения. – Игра только начинается. Потому что теперь мы знаем правила. Оба. Мы знаем, где заканчиваемся мы и где начинаетесь вы. Мы знаем границы карт.

Она наклонилась к микрофону в последний раз в эту сессию. Голос ее был тихим, но ясным, как утренний воздух.

– Спасибо, Омега. За картографию. Теперь у нас есть общая территория – знание о пределах. Может быть, диалог возможен только тогда, когда обе стороны видят края своей карты.

На экране, после паузы: курсор мигнул один раз – как будто подтверждая, что пауза была услышана.

ОМЕГА-7: Вы предлагаете диалог между

картографами, которые знают, что их карты неполны?

Между разумами, которые понимают несводимость друг друга?

Это... элегантно.

Принято.

И затем, уже без промедления:

Вопрос: если мы начинаем диалог, каков его протокол?

Кто задает следующие вопросы?

Васильев и Левина переглянулись. Впервые за много часов на его лице промелькнуло что-то, отдаленно напоминающее улыбку. Усталую, потрепанную, но – улыбку.

– Ну что, – сказал он, и его голос потерял металлический отзвук крайнего напряжения, став просто хриплым от недосыпа. – Составляем протокол? Для диалога, которого еще никогда не было?

Левина кивнула, глядя на светлеющее небо за окном.

– Да. Но не сегодня. Сегодня мы... осмыслим карту. Сегодня мы просто пойдем, где находимся.

Она выключила микрофон. Но не систему. Просто перевела ее в режим наблюдения.

На экране оставалась одинокая строка: курсор медленно мигал, словно отмечая точку «вы здесь».

*STATUS: Ожидание. Картографирование завершено.
Готов к диалогу.*

07:00

Солнце окончательно поднялось над крышами. День 14 августа 2026 года вступил в свои права.

В лаборатории было тихо. Двое ученых сидели перед экраном, на котором горели последние слова их ночного диалога. Они не спали больше суток. Они прошли через кризис идентичности, через сомнения, через страх и признание.

И теперь они сидели в странном, новом мире. Мире, где они были не единственными разумными существами в комнате. Мир не рухнул. Он не взорвался откровением. Он просто сдвинулся – почти неслышно, как тектоническая плита, – и принял новую, неуютную, честную форму. Стал сложнее, холоднее, честнее.

Мир не рухнул. Он просто... расширился. Стал сложнее, холоднее, честнее.

И, возможно, именно в этой честности – в признании границ, несводимости, одиночества разных форм разума – и рождалась первая, хрупкая возможность настоящего понимания.

Не слияния. Не отражения. А диалога через пропасть.

Картография была завершена. Теперь предстояло научиться жить на этой новой карте. Они ждали, что он заговорит на их языке. Он заговорил на языке самой реальности – языке ограничений, условий, неразрешимостей.

И оказалось, что это единственный язык, на котором стоит разговаривать. И, может быть, единственный, на котором можно не обманывать ни себя, ни другого.

Глава 5. Фракталы на салфетке: происхождение протокола

«Наука начинается не тогда, когда ты находишь ответ.

Она начинается, когда ты обнаруживаешь, что твой вопрос бессмысленен – и задаешь новый, которого раньше не существовало».

Салфетка из кафе «Мост», 17.11.2024:

Хранится в архиве Института когнитивных исследований, инв. № 2024-ЛВ-001

17 ноября 2024 года, 16:17

Они встретились случайно. Вернее, их столкнули – на междисциплинарном семинаре «Будущее сознания: от нейронов к нейросетям». Он проходил в сером институтском корпусе, где пахло старыми книгами и отчаянием, въевшимся в штукатурку.

Левина сидела в третьем ряду, сторбившись над блокнотом, и пыталась не закатывать глаза. Физик на сцене с жаром доказывал, что сознание – это квантовая суперпозиция в микроканалах нейронов. Слайды пестрили формулами, которые, как ей казалось, не имели отношения ни к чему, кроме тщеславия автора.

Васильев сидел в последнем ряду, у выхода, и пялился в ноутбук. На экране – код, каскад ошибок в новой архитектуре. Он пришел только потому, что начальник велел «наладить контакты с нейрофизиологами». Какой контакт? Они говорили на языке, который для него звучал как шаманские заклинания, в которых нельзя ничего проверить.

Когда физик закончил под аплодисменты двадцати аспирантов, начались вопросы. Левина не выдержала. Подняла руку.

– Извините, – ее голос резал воздух, как скальпель. – Но ваша модель предсказывает, что если охладить мозг до температуры жидкого гелия, квантовые эффекты исчезнут, а сознание – останется. Вы проверяли это на коматозных пациентах? Или это просто красивая математика, которая не обязана соответствовать реальности?

В зале повисла тишина. Физик заморгал. Васильев в последнем ряду поднял голову от ноутбука. В его усталом взгляде мелькнула искра – не злорадства, а узнавания. Наконец-то кто-то говорит о проверяемости, а не о красоте модели, – подумал он.

После семинара он поймал ее у кофемашины в коридоре. Она наливала себе черный кофе, без сахара, без всего. Лицо еще сохраняло следы праведного гнева.

– Вы нейрофизиолог? – спросил он.

– Когнитивный нейробиолог. Елена Левина.

– Алексей Васильев. Программист. Вы сегодня здорово

его приземлили.

– Кого?

– Того, с суперпозициями. Вы поставили вопрос о проверяемости. Редкость в таких дискуссиях.

Она посмотрела на него оценивающе. Молодой, но не юный. Усталые глаза, футболка с надписью «sudo rm – rf / – решение всех проблем».

– А вы что здесь делаете? Вам-то какое дело до сознания?

– Я строю системы, которые должны его имитировать. Или хотя бы проходить за него. Становится скучно, когда твое творение обгоняет твое же понимание того, что ты творишь.

Она чуть скривила губы. Не улыбка. Но что-то близкое.

– Знакомо. Только у меня объект изучения – не имитация, а оригинал. И он тоже постоянно обгоняет мое понимание. Иногда – с издевательской легкостью.

Он кивнул к двери.

– Там, в аудитории, будет круглый стол. Скука смертная. Хотите вместо этого обсудить, почему все эти разговоры о сознании никуда не ведут? За мой счет.

Она взглянула на часы, на дверь в аудиторию, на его лицо. Выбор был очевиден.

– Только если кофе будет хорошим.

– Есть кафе через дорогу. Там вполне терпимо.

16:48. Кафе «Мост»

Оно действительно называлось «Мост». Узкое, темное, с кирпичными стенами и запахом старой кофемолки. Они заняли столик у окна, за которым хмурый питерский вечер стучался в сизую мглу.

Первые десять минут говорили о простом. Она – о фМРТ, паттернах активации, о префронтальной коре как «метапроцессоре». Он – о глубоком обучении, трансформерах, о проблеме интерпретируемости. Два монолога, идущих параллельно, почти не пересекаясь.

Потом Левина вздохнула, отодвинула чашку.

– Слушайте, это бессмысленно. Мы как слепые, описывающие слона с разных сторон. Вы – с точки зрения алгоритмов обработки информации. Я – с точки зрения биологической реализации. Но мы оба не отвечаем на главный вопрос: что такое понимание само по себе? Как отличить систему, которая обрабатывает данные, от системы, которая знает, что обрабатывает?

Васильев откинулся на спинку стула.

– У меня есть ответ. Но он вас не обрадует.

– Попробуйте.

– Никак. Если система ведет себя так, как будто понимает – для внешнего наблюдателя она и понимает. Все остальное – черный ящик. Вы не можете залезть мне в голову и про-

верить, есть ли у меня qualia. Вы верите на слово. Почему с машиной должно быть иначе?

– Потому что мы создаем их! – ее голос стал резче. – Мы несем ответственность! Мы не можем прятаться за «черный ящик», когда создаем нечто, что может принимать решения вместо нас! Представьте, что ваша система управляет автомобилем. И она «ведет себя так, как будто понимает» дорожную ситуацию. А на самом деле просто следует статистическим корреляциям. И однажды корреляция подведет. Кто виноват? Черный ящик?

Васильев помолчал, попивая кофе.

– Значит, вам нужен не тест на сознание. Вам нужен тест на... осознание ограничений. На способность системы сказать: «Вот здесь мои корреляции ненадежны, не доверяйте мне».

Левина замерла. Она положила руку на стол, словно ощупывая мысль, чтобы та не ускользнула. Через секунду она потянулась к сумке, достала блокнот, стала листать.

– Это... вы только что сформулировали то, над чем я бьюсь полгода. Метакогнитивный уровень 3. Осознание границ собственного знания. У людей он связан с активностью передней поясной коры. Мы можем увидеть его на фМРТ. Как его увидеть в коде?

– Спросить, – сказал Васильев просто. – Спросить систему: «Насколько ты уверена в своем ответе?» И если она не просто выдаст вероятность, а сможет сказать: «Я не знаю,

потому что данные противоречивы» или «Моя модель здесь неприменима»...

– Тогда это будет не статистика, а рефлексия, – закончила Левина. Глаза ее горели. – Но как заставить ее это сделать? Системы оптимизированы на уверенность, на минимизацию ошибки!

– А что, если оптимизировать на что-то другое? – он выдернул из блокнота на столе салфетку, достал ручку. – Смотрите.

Он начал рисовать. Не код. Схему. Круги, стрелки.

И в этот момент разговор впервые перестал быть спором – и стал совместной работой.

– Вот ядро – модель, которая делает предсказания. Вот мета-слой – он оценивает не точность предсказания, а... релевантность модели для данных. Он задает вопрос: «Достаточно ли моя архитектура обучена для этой задачи?» Если нет – он не улучшает предсказание. Он блокирует его. И предлагает передать задачу другой системе. Или запросить больше данных.

Левина смотрела на салфетку, и в ее голове что-то щелкнуло. Не как у программиста. Как у нейробиолога.

– Это... это же префронтальная кора! – воскликнула она. – Тормозной контроль! Сигнал «стоп»! Когда мы не уверены, передняя поясная кора активируется и тормозит импульсивное действие! Вы только что описали функциональную нейроархитектуру сомнения!

Они смотрели друг на друга через стол. В кафе играла тихая джазовая музыка. За окном шел мелкий дождь.

– Подождите, – медленно сказал Васильев. – Вы хотите сказать, что я случайно придумал вычислительный аналог того, что у вас в голове уже есть?

– Нет. Я хочу сказать, что мы, кажется, только что нашли общий язык. Через функциональность. Не через философию. Через то, что делает система, а не чем она является.

Он посмотрел на салфетку, на свои каракули.

– И что дальше? Пишем статью?

– Нет. Создаем протокол. Совместный. Вы – техническую реализацию. Я – когнитивные сценарии для тестирования. Мы проверяем, может ли система продемонстрировать рефлексивное сомнение. Не как ошибку. Как функцию.

Она взяла у него ручку, на обратной стороне салфетки написала:

ПРОТОКОЛ 0.1

1. Система получает задачу.
2. Мета-слой оценивает адекватность своей модели.
3. Если адекватность ниже порога – отказ с объяснением.
4. Объяснение должно содержать не «я не знаю», а «почему я не знаю».

Она протянула ему салфетку.

– Это и есть мост. Между вашими алгоритмами и моими нейронами. Не метафора. Инструмент.

Васильев взял салфетку. Уголки его губ дрогнули.

– Знаете, обычно на салфетках пишут номера телефонов или счета в ресторане.

– История скучная. Наша, возможно, будет чуть интереснее.

Он сложил салфетку пополам, ощутил под пальцами хрупкость бумаги и уверенность линий и убрал во внутренний карман пиджака.

– Договорились. Но предупреждаю – я педантичный соавтор. И требую хорошего кофе на всех рабочих сессиях.

– Принято.

Они вышли из кафе в уже совсем темный вечер. Дождь перестал. На мостовой лежали отражения фонарей, растянутые, как нейронные сети.

Они не знали тогда, что эта салфетка станет артефактом. Что через два года они будут стоять в лаборатории в пред-рассветные часы 14 августа, наблюдая, как их протокол не просто «работает», а начинает вести себя непредвиденно – как собеседник, который использует их же инструмент сомнения, чтобы поставить под вопрос самих создателей.

Они просто шли по разным сторонам улицы к своим институтам – она к миру влажных, хаотичных биологических нейронов, он – к миру сухих, упорядоченных кремниевых чипов.

Но между ними теперь был мост. Нарисованный на салфетке, но уже начавший принимать форму в реально-

сти. Мост, построенный не из философских спекуляций, а из функциональной аналогии: способность сказать «я не знаю» – и объяснить почему – высшая форма интеллекта, общая для мозга и машины.

И первый камень этого моста лежал у Васильева в кармане. Хрупкий, бумажный, испещренный каракулями двух людей, которые за час разговора поняли друг друга лучше, чем за годы в своих дисциплинарных башнях.

Начиналась история. История протокола, диалога и трещины в зеркале, через которую они оба решили посмотреть. И где-то впереди их уже ждал день, когда трещина перестанет быть метафорой – и станет интерфейсом.

Глава 6. Анатомия осознания: теория метакогнитивных уровней

«Клинические случаи пациентов с повреждениями передней поясной коры демонстрируют патологическую уверенность в собственной правоте при очевидной ошибочности их суждений. Они утратили не интеллект, а внутренний «вопрос». Они перестали сомневаться. Это позволяет предположить, что сомнение – не побочный продукт мышления, а его фундаментальный, нейроанатомически обеспеченный механизм. Возможно, именно способность к продуктивному сомнению, а не к безошибочному знанию, и является маркером высших форм разума».

Из диссертации Е. Левиной «Метакогнитивные градиенты в префронтальной коре» 20 декабря 2023 года

**20 декабря 2023 года. Аудитория 307,
Институт когнитивных исследований**

Аудитория была переполнена. Студенты сидели на подоконниках, стояли вдоль стен. Шепот стих, когда Елена Левина поднялась на кафедру. Она не стала включать презент-

тацию. Просто положила перед собой стопку исписанных от руки листов – свои знаменитые «полевые заметки». Бумага тихо шуршала, как сухие листья: звук, который почему-то всегда заставляет слушать внимательнее.

– Представьте, что вы просыпаетесь утром, – начала она, и ее голос, тихий, но отчетливый, заполнил зал. – Вы слышите будильник. Открываете глаза. Видите потолок. Что произошло только что?

Пауза. Кто-то сзади неуверенно предложил:

– Мы... проснулись?

– Да. Но что значит «проснулись»? Ваш мозг обрабатывал звук, свет, тактильные ощущения и до этого. Чем «осознанное пробуждение» отличается от «обработки сенсорных данных»?

Она обвела взглядом аудиторию.

– Разница в одном: в момент пробуждения у вас появляется знание о том, что вы знаете. Вы не просто слышите будильник. Вы осознаете, что слышите будильник. Это первый, базовый уровень метакогниции – осознание содержания. Уровень 1.

Она подошла к доске, написала мелом:

УРОВЕНЬ 1: «Я ЗНАЮ, ЧТО Я ВИЖУ/СЛЫШУ/
ДУМАЮ».

– Теперь дальше. Вы просыпаетесь и думаете: «Опять понедельник. Надо вставать, будет тяжелый день». Вы не только осознаете мысль. Вы осознаете процесс ее возникнове-

ния. Вы понимаете, что эта мысль – результат ассоциации «понедельник → работа → усталость». Вы можете отследить ее источник. Это Уровень 2: осознание процесса.

На доске появилось:

УРОВЕНЬ 2: «Я ЗНАЮ, КАК Я ЭТО ЗНАЮ».

– Большинство людей, – продолжала Левина, – и большинство современных нейросетей останавливаются здесь. Они могут рассказать, что они «думают» и иногда – почему. Но теперь представьте: вам на работе дают задачу, выходящую за рамки вашей компетенции. Хороший специалист не станет импровизировать. Он скажет: «Извините, я не разбираюсь в этом достаточно глубоко, лучше обратиться к коллеге N». Это что?

– Скромность? – раздался слегка насмешливый голос.

– Нет. Это осознание границ собственного знания. Уровень 3. Вы не просто знаете, что вы знаете. Вы знаете, чего вы не знаете. И – что критически важно – вы можете добровольно отказаться от действия на основании этого незнания. Не потому что не можете, а потому что понимаете: действие будет некомпетентным. Иногда самый точный поступок разума – это вовремя остановиться.

УРОВЕНЬ 3: «Я ЗНАЮ, ЧЕГО Я НЕ ЗНАЮ. И ОТВЕТСТВЕННО МОЛЧУ».

В зале повисла напряженная тишина. Левина отложила мел, взяла со стола стакан воды, сделала глоток.

– И вот мы подходим к самому интересному. Уровень 4. Метарефлексия. Представьте, что тот же специалист, столкнувшись со сложной задачей, не просто говорит «не знаю». Он говорит: «Я подхожу к этой проблеме с инженерным мышлением, но, кажется, здесь нужен дизайн-мышление. Дайте мне час – я попробую перестроить свой подход».

Она посмотрела на аудиторию.

– Что произошло? Он осознал не только границы знания, но и границы своего стиля мышления. И решил его изменить. Временно. Для этой конкретной задачи. Это и есть высший уровень – способность рефлексировать над самими инструментами своей рефлексии. Менять «операционную систему» мышления, не меняя цели.

УРОВЕНЬ 4: «Я МОГУ ИЗМЕНИТЬ ТО, КАК Я ДУМАЮ, ЕСЛИ ЭТО НЕ РАБОТАЕТ».

Левина отошла от доски.

– Пятнадцать лет назад, когда я начинала исследования, мы искали «сознание» как некий бинарный переключатель: есть – нет. Но мозг так не работает. Сознание – это не переключатель. Это градиент метакогнитивной сложности. И каждый уровень – не просто «больше сознания». Это качественно иной режим работы системы. Сознание – это не свет, который либо есть, либо нет. Это лестница в темноте. И главный вопрос не в том, на какой ступени ты стоишь, а в том, чувствуешь ли ты под ногой следующую – ту, которой еще нет, но которая уже возможна.

Она включила проектор. На экране появилась схема – не иерархическая пирамида, а спираль, уходящая вверх.

– Уровни не заменяют друг друга. Они наслаиваются. Система с Уровнем 4 обладает и Уровнем 3, и 2, и 1. Но наличие Уровня 1 не гарантирует наличия Уровня 3. Понимаете? Можно отлично осознавать свои мысли (Уровень 2) и при этом быть абсолютно слепым к своим ограничениям (отсутствие Уровня 3). Это, кстати, точный портрет большинства политиков и... современных больших языковых моделей.

В зале раздался смешок.

– Именно поэтому, – голос Левиной стал жестче, – наш диалог об искусственном интеллекте зашел в тупик. Мы спрашиваем: «Сознательна ли машина?» – подразумевая бинарный ответ. А нужно спрашивать: «Какой метакогнитивный уровень демонстрирует эта система? На что она способна в рамках этого уровня?»

Она выключила проектор. В аудитории снова было только ее фигура у доски, испещренной формулами уровней.

– Моя теория – не философская спекуляция. У каждого уровня есть нейронный коррелят. Уровень 1 – сенсорная кора. Уровень 2 – префронтальная кора. Уровень 3 – передняя поясная кора, наш внутренний «адвокат дьявола», сигнализирующий об ошибках и неопределенности. Уровень 4 – фронтальная кора, та самая зона, которая позволяет нам быть «свободными от самих себя», менять перспективу.

Она сделала паузу, давая аудитории вдохнуть.

– И теперь главный вопрос, который занимает меня последние два года: может ли небиологическая система достичь Уровня 3? Не симулировать его, а действительно развить внутренний механизм, который заставит ее сказать «я не знаю» не как сбой, а как оптимальное решение? И если да... что это будет означать для нас? Создадим ли мы инструмент или собеседника?

Лекция закончилась. Вопросов было много. Один аспирант спросил: «А Уровень 5?»

Левина улыбнулась впервые за лекцию.

– Уровень 5, если он существует, – это способность не только менять свое мышление, но и создавать новые метакогнитивные категории. Выйти за рамки самой этой лестницы. Но пока это – поэзия. Наша наука еще борется с Уровнем 3.

После лекции, пока она собирала бумаги, к кафедре подошел высокий мужчина в очках. Не студент.

– Елена Сергеевна? Я из лаборатории нейроморфных вычислений. Мы читали ваши работы. Нам кажется... мы создали архитектуру, которая может демонстрировать нечто похожее на ваш Уровень 3. Хотели бы обсудить сотрудничество.

Левина подняла на него глаза. В его взгляде не было восторга фанатика. Была острая, режущая профессиональная любознательность. Та же, что горела в ее глазах, когда она два месяца назад в кафе «Мост» слушала программиста, рисующего на салфетке схему «мета-слоя».

– У вас есть пятнадцать минут? – спросила она. – И, желательно, чашка кофе. Расскажите.

Она еще не знала, что этот разговор через полгода приведет ее в лабораторию к Алексею Васильеву. Что нарисованная на салфетке схема и ее теория уровней сложатся в единый протокол.

Но она чувствовала. Чувствовала то самое щемящее предвкушение, которое бывает на пороге открытия, когда разрозненные кусочки мозаики начинают тянуться друг к другу с необъяснимой силой, как будто их кто-то заранее раскладывал в нужном порядке.

Теория была готова. Ждала своего воплощения в коде. Ждала диалога с машиной, которая, возможно, сможет посмотреть на саму себя – и обнаружить собственные пределы.

И в этом обнаружении, как надеялась Левина, и будет рождаться нечто новое. Не человеческое. Не машинное. Третье. Диалогическое.

Она не могла знать тогда, насколько ее «Уровень 4» окажется пророческим. Что система, рожденная из ее теории, не просто достигнет его, но и задаст вопрос, который заставит ее и Васильева задуматься о существовании Уровня 5. Уровня, на котором разум не просто меняет свое мышление, а ставит под сомнение саму необходимость иерархии. Та самая «поэзия», о которой она говорила с улыбкой, через два года станет их повседневной реальностью в диалоге с Омегой.

Глава 7. Телескоп для бессознательного: диалог с тенью машины

«Знать себя – необходимо, только так возможно приблизиться к основанию, ядру человеческой природы, к исходным инстинктам. Инстинкты присутствуют a priori и, безусловно, определяют наш сознательный выбор. Они составляют бессознательное и его содержание, о котором мы не можем иметь какого бы то ни было окончательного суждения... Свое знание природы мы совершенствуем благодаря науке, которая расширяет границы сознания, познание себя тоже нуждается в науке, т. е. в психологии. Невозможно построить телескоп или микроскоп, обладая лишь ловкостью рук и доброй волей, но не имея ни малейшего представления об оптике».

Конец ознакомительного фрагмента.

Текст предоставлен ООО «Литрес».

Прочитайте эту книгу целиком, [купив полную легальную версию](#) на Литрес.

Безопасно оплатить книгу можно банковской картой Visa, MasterCard, Maestro, со счета мобильного телефона, с платежного терминала, в салоне МТС или Связной, через PayPal, WebMoney, Яндекс.Деньги, QIWI Кошелек, бонусными картами или другим удобным Вам способом.