

АГЕНТЫ СРЕДИ НАС

Как нанять нейросеть
на работу...



...и не уволиться самому

ЕВГЕНИЙ ВОЛКОВ

Евгений Игоревич Волков
Агенты среди нас. Как
нанять нейросеть на работу
и не уволиться самому

http://www.litres.ru/pages/biblio_book/?art=73173923
ISBN 9785006905726

Аннотация

AI-агенты уже работают. Каждый день разница между теми, кто их использует, и теми, кто игнорирует, растёт. Эта книга – практическое руководство по внедрению агентов в вашу работу. Вы узнаете, как масштабировать бизнес в 5—10 раз, если вы предприниматель, или работать эффективнее, чтобы БОСС задавался вопросом: «Как ты это делаешь?», если вы работаете в найме. Для предпринимателей, менеджеров и специалистов, которые хотят остаться конкурентоспособными в 2026 году.

Содержание

ПРЕДИСЛОВИЕ	6
Введение. Конец эпохи одиноких чат-ботов	14
Почему ChatGPT больше не нужен (если вы используете его руками)	14
Тренд Gartner 2025: Что такое «Фоновый ИИ» (Background AI) и почему он работает, пока вы спите	18
Разница между инструментом и сотрудником: от генерации текста к выполнению действий	23
Обещание книги: как построить агентство из цифровых сотрудников за выходные	27
Часть I. Рекрутинг: Кого мы нанимаем?	35
Глава 1. Анатомия агента	35
Мозг: Большая Языковая Модель (LLM)	36
Чем агент отличается от простого скрипта автоматизации	40
Типология цифровых личностей:	43
Исследователь, Критик, Исполнитель, Менеджер	
Глава 2. Почему они должны говорить друг с другом	47
Проблема «одного большого промпта»: почему универсальные модели глупеют	47

от сложных задач	
Разделяй и властвуй: как разбиение задачи на микро-роли повышает качество на 40%	53
Пример из жизни: Как один агент пишет статью, второй – критикует, а третий – верстает (без участия человека)	57
Часть II. Штатное расписание: Архитектура системы	63
Глава 3. Роль CEO: Ваша новая работа	63
Как перестать быть «оператором промпта» и стать «архитектором системы»	63
Карта процессов: выявляем рутину, которую можно делегировать агентам	68
Чек-лист готовности бизнеса к автономности (без участия человека)	73
Конец ознакомительного фрагмента.	75

**Агенты среди нас. Как
нанять нейросеть на работу
и не уволиться самому**

Евгений Игоревич Волков

© Евгений Игоревич Волков, 2026

ISBN 978-5-0069-0572-6

Создано в интеллектуальной издательской системе Ridero

ПРЕДИСЛОВИЕ

Письмо от автора: Почему эта книга может изменить вашу карьеру

Представьте, что вы просыпаетесь в среду утром. Вы привычно открываете ноутбук и видите 150 новых email. Из них 80 требуют вашего ответа. Вы тратите час на то, чтобы отсортировать их по приоритету, написать ответы, скопировать информацию в CRM.

Затем встреча с командой. Час на то, чтобы услышать о проблемах, которые уже были решены на предыдущей встрече. Потом нужно написать отчет босу. Два часа на сбор данных из разных источников, объединение в один документ, переформатирование, отправка.

Обед. 15 минут. Быстро.

Потом еще две встречи. Потом – срочный звонок от клиента. Потом – исправление ошибки в документе, который был создан неправильно. Потом – планирование следующего дня.

Вы уходите из офиса в 7 вечера. Вы устали. Вы не помните, что вы вообще сделали в этот день, кроме того, что ответили на много email и ходили на встречи.

Это называется shadow work – невидимая работа, которая съедает вашу энергию, но не создает никакой ценности.

Если вы находитесь в этой ситуации, эта книга – для вас.

Но есть и вторая часть.

В той же среду утром другой человек просыпается и видит 150 email. Но первый час его дня выглядит по-другому. Он открывает Slack и видит сообщение от агента: «Я обработал все неважные email и поместил их в папку с меткой „автоматический ответ отправлен“. Нужны твои действия только для 5 писем (вот они)».

Он смотрит на 5 писем. Это займет 10 минут.

Потом он смотрит на свой календарь и видит, что встреча, которая была запланирована на сегодня, уже закончилась (агент провел встречу от его имени, на основе инструкций).

Отчет для босса? Агент уже собрал данные и создал проект отчета. Нужно просто глянуть и нажать «отправить».

К 10:00 AM этот человек уже покончил с тем, что обычного человека занимает до 6 вечера.

Теперь у него есть время на то, что действительно важно: стратегия, творчество, отношения, принятие решений.

Выбор между двумя будущими

Мы стоим на пороге раскола.

В течение следующих 3—5 лет люди разделятся на две группы.

Группа 1: Люди, которые освоили работу с AI-агентами. Они использовали это время не для прочтения еще одной статьи про ChatGPT, а для настройки агентов, которые работают за них. Эти люди работают 4—5 часов в день (вместо 8), получают в 2—3 раза больше результатов, имеют время

на семью, здоровье, собственные проекты.

Группа 2: Люди, которые не использовали AI. Они по-прежнему читают email, ходят на встречи, пишут отчеты вручную. Они по-прежнему устают в конце дня. Они замечают, что их коллеги (из группы 1) достигают больше результатов за меньше времени, и не понимают как.

Разница в продуктивности между двумя группами уже в 2025 году составляет 3—5х. К 2027 году это станет 10х.

Компании будут нанимать людей из группы 1 и отпускать людей из группы 2.

Фрилансеры из группы 1 смогут брать в 5 раз больше проектов и зарабатывать в 3 раза больше денег.

Предприниматели из группы 1 смогут масштабировать бизнес в 10 раз без найма людей.

Почему сейчас – идеальный момент

Вы могли бы прочитать про AI-агентов в 2026 году. К тому времени информация будет везде. Статьи. Видео. Мастер-классы. Курсы за \$500. И это будет поздно.

Потому что к 2026 году те, кто внедрил агентов в 2025 году, уже получили все преимущества. Они уже стали лидерами в своих компаниях. Они уже запустили новые проекты. Они уже зарабатывают больше.

А те, кто ждал, будут ловить волну. Они будут пытаться наверстать упущенное. Но будет поздно.

В 2025 году окно еще открыто. Есть еще время.

Но это окно быстро закрывается.

Что находится в этой книге

Эта книга не про то, как устроены нейронные сети. Не про math за трансформерами. Не про fine-tuning и prompt engineering на уровне исследования.

Эта книга про то, как использовать AI-агентов для конкретных, реальных задач.

Мы рассмотрим:

Основы (Главы 1—2): Что такое AI-агент? Почему он отличается от ChatGPT? Как устроена архитектура?

Практика (Главы 3—10): Как на самом деле создавать и запускать агентов? Какие инструменты использовать? Как они работают в реальных бизнес-процессах (продажи, маркетинг, HR, аналитика)?

Безопасность (Главы 11—12): Что может пойти не так? Как защитить систему? Какие юридические риски? Как не потерять данные клиентов?

Будущее (Заключение): Как это изменит вашу карьеру? Как выжить в эпоху AI? Конкретный план действий на следующий понедельник.

Каждая глава основана на реальных примерах, реальных ошибках, реальных решениях. Не на теории. Не на гипотезах. На том, что уже происходит в 2024—2025 году.

Три типа читателей

Тип 1: Предпринимателя, фрилансера

Вы уже знаете, что вам нужно масштабировать. Вы знаете, что нанять людей дорого. Вы ищете способ сделать больше

с меньшим.

Для вас эта книга – техническое руководство по масштабированию. Вы узнаете, как настроить систему из агентов, которая будет делать работу, которую раньше делали бы 2—3 человека.

Тип 2: Менеджер в компании

Вы ответственны за продуктивность команды. Вы видите, что люди тратят слишком много времени на рутину. Вы хотите освободить их время для более важной работы.

Для вас эта книга – инструмент для трансформации команды. Вы узнаете, как внедрить агентов в процессы, как это окупается, какие риски управлять.

Тип 3: Эксперт в своей области (маркетолог, аналитик, специалист)

Вы хороший в том, что вы делаете. Но вы замечаете, что конкуренты быстрее вас, делают больше. Вы хотите остаться конкурентоспособным.

Для вас эта книга – способ остаться впереди. Вы узнаете, как агенты могут усилить ваши навыки, как сделать свою работу в 5 раз быстрее, как выделиться среди конкурентов.

Важное замечание: Это не волшебство

Эта книга может изменить вашу карьеру. Но она не волшебство.

Прочитать эту книгу и ничего не сделать – это потратить время впустую.

Внедрить агентов, потратив две недели на настройку,

но потом забросить – это неправильный способ.

Ожидать, что агенты решат все проблемы сами – это установка на неудачу.

Правильный способ – это медленная, постоянная интеграция. Начать с одного маленького агента. Посмотреть результаты. Расширить. Повторить.

По данным McKinsey Global Institute, компании, которые внедрились AI и продолжали инвестировать в обучение команды, получили 3x ROI (возврат инвестиций). Компании, которые купили инструменты и надеялись на чудо, получили убыток.

Это книга про первый тип. Про постоянное, системное внедрение.

План прочтения

Вы можете прочитать эту книгу разными способами.

Способ 1: От начала до конца (для новичков)

Если вы никогда раньше не работали с AI-агентами, читайте книгу по порядку. Каждая глава строится на предыдущей. Первые две главы – это фундамент.

Время: 6—8 часов чтения.

Способ 2: Быстрый старт (для людей в спешке)

Если у вас нет времени на полное прочтение, прочитайте:

– Главы 1—2 (основы) – 1 час

– Главы 5—6 (инструменты и их настройка) – 2 часа

– Заключение (план на понедельник) – 1 час

Это даст вам достаточно знаний, чтобы начать.

Время: 4 часа.

Способ 3: Фокус на вашей области (для специалистов)

Если вы маркетолог – читайте главу 9. Если вы HR – читайте главу 10. Если вы продажник – читайте главу 7.

Потом вернитесь и прочитайте главы 1—2 и 11—12 (безопасность и юридика).

Время: 3—4 часа (зависит от главы).

Последнее слово перед началом чтения

Мир меняется быстро. Очень быстро.

В 2022 году ChatGPT был новинкой, которой играли на выходных.

В 2023 году компании начали задумываться о том, как использовать ChatGPT в работе.

В 2024 году агенты стали реальностью. Не гипотезой. Реальностью.

В 2025 году (прямо сейчас) те, кто внедрили агентов, получают конкурентное преимущество.

В 2026 году это конкурентное преимущество станет конкурентной необходимостью. Если вы не используете агентов, вы будете позади.

Эта книга писалась в 2024—2025 году. Когда окно все еще открыто.

Пока вы читаете эти слова, где-то в другом городе менеджер настраивает своего первого агента. Предприниматель запускает систему, которая будет работать без него. Специалист освобождает свое время для более важной работы.

Вопрос: будете ли вы среди них?

Добро пожаловать в эпоху AI-агентов.

Давайте начнем.

Введение. Конец эпохи одиноких чат-ботов

Почему ChatGPT больше не нужен (если вы используете его руками)

Давайте будем честны: в 2025 году открывать вкладку браузера, чтобы написать «Напиши мне пост» или «Составь структуру договора», – это уже не передовая технология. Это цифровой ручной труд.

Мы привыкли считать появление ChatGPT в конце 2022 года революцией. И это действительно было так. Но революции имеют свойство пожирать своих детей, а технологии – стремительно устаревать. То, что вчера казалось магией, сегодня становится рутиной, а завтра – неэффективным атавизмом.

Проблема не в самой языковой модели. GPT-4o, Claude 3.5 Sonnet или Llama 3 по-прежнему остаются великолепными «мозгами». Проблема в интерфейсе. Проблема в том самом окошке чата, курсоре и кнопке «Send».

Если вы используете нейросети руками – вы работаете на нейросеть, а не она на вас.

Взгляните на свой рабочий процесс. Вы формулируете за-

дачу (тратите когнитивный ресурс). Вы вводите промпт. Вы ждете, пока курсор, мигая, выведет текст по буквам (даже с высокой скоростью это отнимает секунды, которые складываются в часы). Вы читаете результат. Вы понимаете, что модель упустила важную деталь. Вы пишете правку: «Нет, сделай это в более официальном тоне и добавь список». Модель переписывает. Вы копируете результат. Вы открываете Google Docs, почту или CRM. Вы вставляете текст. Вы правите форматирование.

В этой цепочке вы – не директор. Вы – «человеческая прослойка» (Human Middleware). Вы выполняете роль API-интерфейса между нейросетью и реальным миром. Вы – тот самый курьер, который берет данные из одного окна и переносит их в другое.

Почему этот подход обречен?

1. Проблема контекстного окна и «памяти золотой рыбки»

Даже с расширением контекстных окон до миллионов токенов, чат-интерфейс линейн. Вспомните, сколько раз вы пытались заставить ChatGPT вспомнить, о чем вы договаривались десять сообщений назад, или пытались «скормить» ему огромный PDF-файл, чтобы через два вопроса обнаружить, что он забыл ключевое условие из середины документа.

В ручном режиме вы вынуждены постоянно «подгружать» контекст в голову модели заново. Вы каждый раз объясняете новому чату: «Я владелец магазина автозапчастей, мы прода-

ем радиаторы...». Это неэффективно. Это похоже на то, как если бы вы нанимали нового сотрудника каждое утро и каждый раз заново проводили ему онбординг, объясняя, где находится кулер и как зовут бухгалтера.

2. Иллюзия диалога

Чат-бот создан для разговора. Но бизнес-задачи редко решаются разговорами. Вам не нужно *обсуждать* с бухгалтером отчет – вам нужен *готовый отчет*. Вам не нужно *обсуждать* с дизайнером логотип – вам нужны *файлы в векторе*.

Интерфейс чата навязывает нам формат дискуссии, в то время как бизнесу нужен формат «Задание – Результат». Пока вы «болтаете» с ботом, вы тратите время. Автономные агенты не болтают. Они делают.

3. Ловушка микроменеджмента

Используя ChatGPT вручную, вы занимаетесь тотальным микроменеджментом. Вы проверяете каждое слово. Вы не можете отойти от монитора, потому что процесс остановится. Если вы перестанете писать промпты, работа встанет.

Это противоречит самой идее автоматизации. Представьте, что вы купили робот-пылесос, но вместо того, чтобы включить его и уйти, вы ходите за ним с пультом управления и нажимаете: «Вперед. Теперь влево. Теперь всоси пыль. Теперь вправо». Это абсурд. Но именно этим занимаются 90% предпринимателей, использующих ИИ сегодня. Они ходят за своим «цифровым пылесосом» и указывают ему каж-

дый шаг.

4. Изоляция данных

ChatGPT в браузере – это остров. Он не видит вашу почту, пока вы не скопируете туда письмо. Он не знает остатков на складе, пока вы не выгрузите ему Excel. Он оторван от реальности.

Ручное копирование данных из CRM в нейросеть и обратно – это не только потеря времени, это риск утечек и ошибок. Стоит вам случайно не скопировать одну строчку из таблицы, и вся аналитика, которую выдаст модель, будет ошибочной. Агенты же живут *внутри* вашей инфраструктуры, имея прямой доступ к «кровеносной системе» бизнеса по API.

Эпоха «одиноких чат-ботов» заканчивается не потому, что они плохие. А потому что они требуют вашего внимания. В экономике внимания (Attention Economy) ваше внимание – самый дорогой ресурс. Тратить его на переписку с ботом – непозволительная роскошь.

Мы переходим от парадигмы «AI as a Tool» (ИИ как инструмент в руках мастера) к парадигме «AI as a Workforce» (ИИ как рабочая сила). Молоток не может работать без вас. Сотрудник – может и должен.

Если вы до сих пор платите подписку за ChatGPT Plus, чтобы самому писать туда тексты, – поздравляю, вы наняли самого дорогого секретаря в мире: себя.

Тренд Gartner 2025: Что такое «Фоновый ИИ» (Background AI) и почему он работает, пока вы спите

Если 2023 и 2024 годы прошли под знаменем Generative AI (Генеративного ИИ) – способности машин создавать контент, то 2025 год аналитики Gartner и других ведущих агентств однозначно маркируют как год Agentic AI (Агентного ИИ) и так называемого «Фонового интеллекта».

В отчетах «Top Strategic Technology Trends for 2025» Gartner выделяет Agentic AI как ключевой тренд, прогнозируя, что к 2028 году не менее 15% повседневных рабочих решений будут приниматься автономными агентами (для сравнения: в 2024 году эта цифра была близка к 0%).

Что такое «Фоновый ИИ» (Background AI)?

Это концепция, при которой искусственный интеллект перестает быть событием, требующим вашего участия, и становится средой. Он исчезает с переднего плана. Он растворяется в бизнес-процессах.

Вспомните электричество. Когда вы входите в темную комнату и щелкаете выключателем, вы не думаете: «О, я сейчас применю технологию потока электронов через вольфрамовую нить». Вы просто хотите света. Электричество – это фоновая технология. Вы замечаете его только тогда, когда

оно пропадает.

ИИ проходит ту же эволюцию.

Сегодня использование ИИ – это событие. «Давай спросим у ChatGPT», – говорите вы коллеге. Это активное действие.

«Фоновый ИИ» – это когда вы приходите утром на работу, а в вашей CRM уже заполнены карточки клиентов, которые написали ночью; в почте лежат черновики ответов на сложные претензии; а на дашборде горит аналитика продаж за вчерашний день с выделенными аномалиями. Вы не «просили» ИИ сделать это. Он просто сделал это, потому что это его работа.

Почему «Фоновый ИИ» меняет правила игры?

– От реактивности к проактивности

– Классический софт (и даже классические чат-боты) реактивен. Он ждет команды. Microsoft Word не напишет за вас отчет, пока вы не начнете печатать.

– Фоновые агенты проактивны. Они настроены на триггеры (события), а не на команды пользователя.

– *Триггер*: Клиент прислал письмо с пометкой «Срочно».

– *Реакция агента*: Проанализировать письмо, понять суть проблемы, проверить статус заказа в базе, сформировать ответ, поставить задачу менеджеру и прислать уведомление в Slack.

– Всё это происходит в 03:00 ночи, пока вы спите. Утром вы утверждаете действие, а не создаете его с нуля.

– Непрерывность бизнес-процессов

– Человеческий ресурс дискретен. Мы работаем 8 часов, потом спим, едим, болеем, уходим в отпуск, выгораем.

– Фоновый ИИ обеспечивает конвейерную непрерывность. Агент, ответственный за мониторинг цен конкурентов, не забудет проверить сайты в воскресенье утром. Агент техподдержки не ответит раздраженно клиенту в конце тяжелой смены, потому что у него нет настроения.

– Для малого бизнеса это означает возможность предоставлять сервис уровня enterprise-корпораций с огромными штатами, не раздувая ФОТ (фонд оплаты труда).

– Снижение когнитивной нагрузки (Cognitive Load)

– Тренд Gartner подчеркивает: люди устали от технологий. Мы тонем в приложениях, уведомлениях и интерфейсах. Фоновый ИИ не добавляет новый интерфейс – он убирает существующие.

– Вместо того чтобы заходить в пять разных сервисов (почта, таск-трекер, аналитика, календарь, чат), вы получаете агрегированный результат работы агентов в одном канале. Агенты берут на себя «цифровой шум», фильтруя информацию и передавая человеку только то, что действительно требует принятия решения.

– Композитный ИИ (Composite AI)

– Фоновый ИИ редко бывает монолитным. Это не одна гигантская нейросеть, которая «знает всё». Это оркестр из множества специализированных маленьких моделей.

– Gartner называет это Composite AI. Один агент использует модель компьютерного зрения, чтобы «смотреть» на сканы накладных. Другой использует LLM, чтобы извлекать из них данные. Третий использует математический модуль, чтобы проверять суммы.

– Пользователь не видит этой сложности. Для него это просто «папка, в которой самозарождается порядок».

Почему это работает, пока вы спите?

Секрет в смене парадигмы архитектуры. Раньше мы строили системы вокруг человека: человек был процессором, а софт – инструментом. Теперь мы строим системы, где ИИ – процессор, а человек – контролер.

Автономные агенты работают в бесконечном цикле OODA (Observe, Orient, Decide, Act – Наблюдай, Ориентируйся, Решай, Действуй):

– Наблюдай: Мониторинг входящих писем, новых записей в базе данных, изменений на сайте конкурента.

– Ориентируйся: Сравнение новых данных с инструкциями, памятью и прошлым опытом (RAG – Retrieval Augmented Generation).

– Решай: Выбор инструмента или сценария поведения (нужно ли отвечать самому или позвать человека?).

– Действуй: Отправка API-запроса, генерация текста, создание файла.

Этот цикл не требует кнопки «Старт». Он запущен постоянно.

Фоновый ИИ превращает бизнес из набора ручных операций в саморегулирующийся организм. Ваша роль меняется с «гребца на галере» на «капитана корабля», который прокладывает курс, пока механизмы в трюме делают тяжелую работу.

Разница между инструментом и сотрудником: от генерации текста к выполнению действий

Чтобы успешно «нанять» нейросеть, нужно фундаментально изменить отношение к ней. Главная ментальная ловушка, в которую попадают предприниматели, – восприятие ИИ как очень умного поисковика или текстового редактора.

Давайте проведем четкую границу между Инструментом (Tool), Ассистентом (Co-pilot) и Сотрудником (Agent/Worker). Понимание этой градации сэкономит вам сотни часов и уберезет от разочарований.

1. Инструмент (Tool): «Лопата»

ChatGPT в браузере, Midjourney, переводчик DeepL – это инструменты.

– Инициатива: 0%. Лопата не начнет копать сама, даже если увидит кучу земли. Она лежит и ждет, пока вы ее возьмете.

– Ответственность: 0%. Если яма выкопана криво, виноват копатель, а не лопата.

– Контекст: Отсутствует. Лопата не знает, что вы строите – замок или туалет. Она просто копает.

– Режим работы: Синхронный. Вы нажали кнопку – получили результат.

Инструменты хороши для разовых задач. Но вы не можете построить бизнес из одних инструментов, если некому их держать. В модели «человек с инструментом» узким местом всегда остается человек. Вы не масштабируемы.

2. Ассистент (Co-pilot): «Стажер»

Это встроенные помощники в софте (Github Copilot, Notion AI, Microsoft 365 Copilot).

– Инициатива: 10—20%. Они могут подчеркнуть ошибку или предложить окончание фразы, пока вы пишете. Они «подхватывают» руль, но не ведут машину.

– Ответственность: Низкая. Они предлагают, вы утверждаете.

– Контекст: Ограниченный. Они видят документ, который открыт перед вами, но вряд ли знают стратегию компании на год.

– Режим работы: Совместный. Вы работаете в паре.

Ассистенты ускоряют работу, но не заменяют работника. Вы все еще должны присутствовать в процессе. Если вы уйдете пить кофе, Copilot в Word не допишет статью за вас. Он будет ждать курсора.

3. Сотрудник (Agent/Worker): «Профессионал»

Это то, чему посвящена эта книга. Автономный агент.

– Инициатива: 100%. Агент сам ищет работу. Если его задача – обрабатывать лиды, он не ждет команды. Он реагирует на появление лида. Если лидов нет, он может (если так запрограммирован) пойти проверять старую базу.

– Ответственность: Делегированная. Вы наделяете агента правом совершать ошибки в определенных пределах. Вы даете ему «бюджет доверия». Например: «Можешь отвечать на типовые вопросы сам, но если клиент ругается матом – зови меня».

– Контекст: Глубокий. Агент подключен к базе знаний компании, он «помнит» историю переписки с клиентом за три года, он знает прайс-лист и скрипты продаж.

– Режим работы: Асинхронный. Вы поставили задачу (или настроили процесс) и ушли. Агент работает в фоне. Вы проверяете результат постфактум.

Ключевое отличие: От слов к действиям (From Chat to Action)

Самый важный водораздел проходит по линии Actionability (способности действовать).

Языковая модель (LLM) сама по себе – это просто мозг в банке. Она может придумать гениальный план маркетинга, но она не может отправить ни одного письма. Она парализована.

Агент – это «Мозг» (LLM), к которому пришили «Руки» (Tools/API).

– *Инструмент* генерирует текст письма.

– *Агент* генерирует текст, открывает Gmail, вставляет адрес получателя, прикрепляет файл КП и нажимает «Отправить».

Разница колоссальна.

Инструмент: «Вот текст SQL-запроса, который тебе нужен». (Вам нужно скопировать его, открыть базу данных, вставить и выполнить).

Сотрудник: «Я проанализировал базу данных, нашел 15 дублей заказов, объединил их и отправил отчет тебе на почту. Вот ссылка».

Нанимая нейросеть на работу, мы перестаем искать «лучший генератор текста». Мы начинаем искать «лучшего исполнителя действий». Мы переходим от промпт-инжиниринга (как правильно попросить) к флоу-инжинирингу (как правильно выстроить процесс).

В этой книге мы будем строить именно таких сотрудников. Не болтливых собеседников, а молчаливых исполнителей, у которых есть доступ к клавиатуре, мышке и кнопке «Enter» вашего бизнеса.

Это рискованно? Да.

Это требует контроля? Безусловно.

Но это единственный путь масштабирования в мире, где скорость реакции стала важнее размера компании.

Когда вы нанимаете живого менеджера, вы не ожидаете, что будете диктовать ему каждое слово, которое он скажет клиенту. Вы даете ему скрипт, телефон и свободу действий.

Точно так же мы поступим с нейросетями. Мы дадим им свободу. Но сначала научим их не засовывать пальцы в розетку.

Обещание книги: как построить агентство из цифровых сотрудников за выходные

Большинство бизнес-книг предлагают вам знания. Они обещают, что после прочтения вы станете умнее, начнете лучше понимать рынок или освоите новые термины. Эта книга предлагает другое. Она предлагает вам готовый актив.

Мое обещание звучит так: если вы начнете читать эту книгу в пятницу вечером и будете параллельно выполнять описанные шаги, то к утру понедельника у вас будет не просто «понимание нейросетей». У вас будет работающий штат сотрудников, которые не требуют кофе, не уходят на перекур и не просят отпускных. У вас появится Цифровое Агентство, интегрированное в сердце вашего бизнеса.

Это смелое заявление. В мире IT принято считать, что внедрение новых технологий – это больно, долго и дорого. Мы привыкли, что «цифровая трансформация» – это проект на полгода с бюджетом, превышающим годовую прибыль. Но правила изменились. Агентный ИИ (Agentic AI) демократизировал автоматизацию так же, как смартфоны демократизировали фотографию. Вам больше не нужна команда программистов, чтобы создать сложную систему. Вам нужны только логика, доступ к API и выходные.

Почему «Агентство», а не «Программа»?

Прежде чем мы перейдем к дорожной карте, давайте зафиксируем ключевую метафору, которая красной нитью пройдет через все страницы. Мы не строим программу. Мы не пишем скрипт. Мы открываем агентство.

Когда вы устанавливаете 1С или CRM, вы настраиваете инструмент. Вы ожидаете, что он будет работать строго по инструкции: «Если нажата кнопка А, покажи окно Б». Это детерминированная система. Она надежна, но глупа.

Когда вы создаете мультиагентную систему, вы нанимаете коллектив. Вы создаете сущности, обладающие некоторой (пусть и ограниченной) автономией.

Почему это важно для вашего настроения? Потому что отношение к агентам должно быть менеджерским, а не инженерным.

– Инженер спрашивает: «Почему этот код выдал ошибку в строке 404?»

– Менеджер спрашивает: «Почему этот сотрудник неправильно понял задачу и как мне переписать инструкцию, чтобы в следующий раз он справился?»

В этой книге мы будем строить именно структуру агентства. У нас будут отделы. У нас будет иерархия. У нас будут планерки (автоматические, где один нейро-агент проверяет работу другого).

Вы перестанете быть «пользователем ПК». Вы станете основателем теневой корпорации, которая работает на сервере.

рах, пока вы занимаетесь стратегией.

Феномен «Выходных»: Спринт вместо Марафона

Почему я делаю акцент на сроке в 48 часов? Неужели можно создать сложную систему так быстро?

Ответ кроется в психологии предпринимателя и особенностях No-Code инструментов.

Главный враг внедрения ИИ – это не сложность технологий. Это «паралич анализа» (Analysis Paralysis). Рынок меняется так быстро, что пока вы изучаете один инструмент, выходит два новых. Вы можете месяцами читать про LangChain, Pinecone и векторные базы данных, но так и не автоматизировать ни одного процесса.

Концепция «за выходные» – это ваш защитный механизм. Это спринт.

– Пятница вечер: Декомпозиция и стратегия. Мы решаем, *кого* мы нанимаем. Мы не пытаемся автоматизировать весь бизнес сразу (это путь к провалу). Мы выбираем одну «боль», один процесс, который высасывает из вас больше всего энергии, и убиваем его.

– Суббота: «Сборка конструктора». Мы используем современные визуальные платформы (о них будет отдельная глава), где создание агента напоминает игру в LEGO. Мы соединяем кубики: «Вот кубик Почты», «Вот кубик Мозга (ChatGPT)», «Вот кубик Google-таблиц». Никакого черного экрана с зеленым кодом. Только логика и стрелочки.

– Воскресенье: «Онбординг и краш-тесты». Агент создан,

но он глуп и наивен. Мы проводим воскресенье, обучая его на ваших реальных данных, показывая ему примеры хорошей и плохой работы, и – самое важное – пытаюсь его сломать. Мы моделируем ситуации, где агент может сойти с ума, и ставим предохранители.

– Понедельник утро: Запуск. Вы нажимаете «Activate», и магия начинается.

Эта книга – не учебник по программированию на Python. Это практическое руководство по *сборке*. Мы будем использовать готовые блоки, чтобы получить результат здесь и сейчас.

Что именно вы получите: Инвентаризация обещаний

Чтобы наше соглашение с читателем было честным, я детально распишу, что именно вы сможете построить, следуя методологии этой книги. Это не абстрактные примеры, а конкретные архитектуры, которые уже работают в тысячах передовых компаний малого бизнеса.

1. Освобождение от «коммуникационного пинг-понга»

Вы построите систему, которая берет на себя первую линию обороны. Не тупой чат-бот с кнопками, который бесит клиентов. А интеллектуальный агент, который может поддерживать светскую беседу, понять контекст, заглянуть в базу наличия товара и вежливо отказать или продать.

Обещание: К концу книги вы сможете создать агента, который самостоятельно квалифицирует входящие лиды, отсе-

ивает спам и передает живым продавцам только «горячих» клиентов с уже заполненной карточкой в CRM.

2. Собственная аналитическая разведка

Информация – новая нефть, но ее слишком много. Вы не можете читать 50 новостных каналов и мониторить 10 сайтов конкурентов ежедневно. Агент может.

Обещание: Вы настроите «Агента-Ресерчера», который каждое утро будет прочесывать заданные источники, собирать данные о ценах конкурентов, находить упоминания вашего бренда и класть вам на стол (в Telegram или на почту) сжатую сводку на одну страницу: «Что случилось, пока я спал, и на что стоит обратить внимание».

3. Конвейер контента без мук творчества

Ведение соцсетей и блогов часто превращается в каторгу. «О чем писать?», «Где взять картинку?», «Текст получился сухим».

Обещание: Мы соберем редакцию, где один агент (Трендолог) ищет темы, второй (Автор) пишет черновик в вашем стиле, третий (Редактор) безжалостно правит и сокращает, а четвертый (Иллюстратор) генерирует визуал. Ваша роль сведется к нажатию кнопки «Одобрить публикацию».

4. Бессмертная база знаний

Сотрудники уходят и уносят знания с собой. Агенты остаются.

Обещание: Вы научитесь создавать RAG-системы (Retrieval-Augmented Generation), которые превращают ва-

ши разрозненные регламенты, PDF-инструкции и гугл-доки в единый мозг компании. Новый живой сотрудник сможет спросить у агента: «Как мы оформляем возврат для юрлиц?», и получить мгновенный точный ответ со ссылкой на пункт договора.

Чего эта книга НЕ обещает (Техника безопасности ожиданий)

Я не хочу продавать вам воздух. Хайп вокруг ИИ породил завышенные ожидания, которые ведут к разочарованию. Давайте сразу обозначим границы возможного.

1. Это не «Кнопка Бабло»

Агенты не спасут убыточную бизнес-модель. Если ваш продукт никому не нужен, автоматизация лишь поможет вам быстрее и эффективнее масштабировать убытки. Агенты – это усилитель. Они умножают то, что есть. Если у вас ноль, то ноль, умноженный на ИИ, останется нулем.

2. Это не полная замена людей (пока)

Я намеренно использую термин «Фоновый ИИ» и «Сотрудник», но вы должны понимать: агент – это джуниор (младший специалист). Он исполнительный, быстрый, но у него нет житейской мудрости и интуиции. Вы не сможете уволить всех и уехать на Бали, оставив бизнес на нейросеть. Напротив, ваша ответственность как руководителя вырастет. Теперь вы отвечаете не только за людей, но и за роботов.

3. Это потребует от вас изменения мышления

Самое сложное в этой книге – не настройка API ключей. Самое сложное – научиться формулировать свои желания так, чтобы их понял кремниевый мозг. Вам придется стать предельно четким. Агенты не понимают «сделай красиво» или «ну ты же понимаешь». Они понимают алгоритмы и критерии. Если у вас хаос в процессах, автоматизация приведет к автоматизированному хаосу. Эта книга заставит вас сначала навести порядок в голове, и только потом – в сервере.

Манифест Цифрового Нанимателя

Приступая к чтению, я прошу вас принять новый майндсет (образ мышления). В эпоху ИИ побеждает не тот, кто больше работает, и не тот, кто лучше всех пишет код. Побеждает лучший архитектор.

Ваш бизнес – это механизм. Раньше детали этого механизма были сделаны только из биологического материала (людей). Они дороги, капризны, но креативны. Теперь у вас появились детали нового типа: цифровые. Они дешевы, предсказуемы и невероятно быстры.

Искусство современного бизнеса – это умение правильно комбинировать эти детали.

– Где нужна эмпатия, теплота и сложные переговоры – ставим человека.

– Где нужна обработка данных, скорость, рутина и доступность 24/7 – ставим агента.

Эта книга – чертеж вашего будущего гибридного предприятия. Мы будем сносить старые стены и прокладывать

новые коммуникации. Будет пыльно, иногда сложно, но результат того стоит.

Представьте, что сейчас утро понедельника. Вы просыпаетесь, берете телефон, и видите уведомление от вашего Главного Агента:

«Доброе утро, босс. За ночь я обработал 15 заявок, 3 перевел в оплату, подготовил отчет по маркетингу и нашел 5 интересных новостей по нашей теме. Жду ваших указаний».

Это не фантастика из фильма про Железного Человека. Это доступная реальность 2025 года. И она стоит всего одних выходных вашего времени.

Вы готовы нанять своего первого цифрового сотрудника? Тогда переверните страницу. Добро пожаловать в отдел кадров будущего.

Часть I. Рекрутинг: Кого мы нанимаем?

Глава 1. Анатомия агента

Когда вы нанимаете живого сотрудника, вы интуитивно понимаете, из чего он «состоит». У него есть мозг (интеллект), есть память (опыт и знания), и есть руки (способность печатать на клавиатуре или звонить по телефону).

С цифровым сотрудником – автономным агентом – ситуация идентична. Чтобы перестать бояться ИИ и начать им управлять, нужно разобрать его на запчасти. Понимание архитектуры агента – это то, что отличает профессионального «архитектора систем» от любителя, который просто переписывается с чат-ботом.

Любой агент, от простейшего бота-секретаря до сложного аналитика, состоит из трех фундаментальных блоков: Мозг, Память и Руки. В технической литературе это называется триадой LLM + RAG + Tools. Давайте рассмотрим каждый орган подробно.

Мозг: Большая Языковая Модель (LLM)

Это центральный процессор агента. Именно здесь происходит магия «мышления». Когда вы используете GPT-4, Claude или Llama, вы обращаетесь к мозгу.

Роль мозга – понимать намерения и принимать решения.

Представьте себе выпускника Гарварда, которого заперли в пустой комнате без интернета и книг. Он очень умен. Он знает 50 языков. Он читал всю Википедию (по состоянию на прошлый год). Он может написать сонет Шекспира или решить сложное уравнение.

Но он оторван от реальности.

– Если вы спросите его: «Какая сегодня погода?», он ответит: «Я не знаю, я в закрытой комнате».

– Если вы спросите: «Сколько денег у нас на счету?», он ответит: «Я не знаю, у меня нет доступа к вашему банку».

Ключевая функция LLM в агенте – это Оркестратор.

Мозг не обязательно должен знать всё. Его главная задача – понять, *что* нужно сделать, и решить, *какой инструмент* для этого использовать.

Пример мыслительного процесса агента (это скрытый монолог, который происходит за доли секунды):

«Пользователь спрашивает про остатки товара на складе. Я сам этого не знаю. Но у меня есть инструмент „Поиск в базе IC“. Значит, мне нужно сформулировать SQL-запрос,

передать его в инструмент, получить ответ и перевести его на человеческий язык для пользователя».

Выбор «мозга» для агента:

Не всем агентам нужен «Эйнштейн» (дорогая модель вроде GPT-4o).

– Для сложных переговоров и стратегического планирования мы берем «дорогие мозги».

– Для сортировки почты или извлечения данных из чеков достаточно «стажера» (быстрой и дешевой модели, например, GPT-4o-mini или Haiku). В бизнесе это называется LLM Routing – экономия бюджета за счет назначения задач моделям соответствующего уровня.

2. Память: RAG (Retrieval-Augmented Generation)

Самая большая проблема «голового» мозга – амнезия и галлюцинации.

LLM помнит только то, чему её учили при создании (общие знания мира), и то, что помещается в текущее окно диалога (кратковременная память). Как только вы закрываете чат, агент всё забывает.

Для бизнеса такой сотрудник бесполезен. Вы не можете нанять менеджера, который каждое утро забывает прайс-лист компании и имена ключевых клиентов.

Здесь на сцену выходит RAG (Retrieval-Augmented Generation) – Генерация, дополненная поиском.

Простыми словами, RAG – это долгосрочная память агента, его личная библиотека и картотека.

Это технология, которая позволяет агенту перед тем, как ответить, «сбегать в архив» и подсмотреть правильный ответ.

Как это работает механически:

– Вы загружаете в систему PDF-инструкции, регламенты, историю переписки, базу знаний компании.

– Система нарезает эти документы на маленькие кусочки (чанки) и складывает в специальную «Векторную Базу Данных» (Vector Database).

– Когда вы задаете вопрос, агент не выдумывает ответ из головы. Он сначала ищет похожие кусочки в вашей базе.

– Он находит нужный пункт инструкции: «Ага, при возврате товара мы требуем заявление по форме №5».

– И только потом формулирует вежливый ответ клиенту, опираясь на этот факт.

Без RAG агент – это фантазер. С RAG агент – это бюрократ, который следует букве вашего закона. RAG – это то, что превращает общедоступную нейросеть (которая училась на всем интернете) в *вашу* корпоративную нейросеть (которая знает только ваш бизнес).

3. Руки: Инструменты (Tools / API)

Мозг с памятью может умно рассуждать, но он по-прежнему парализован. Он может выдать гениальный совет, но не может выполнить действие.

Чтобы агент стал сотрудником, ему нужны «Руки». В мире софта руками являются API (Application Programming

Interface) и Функции (Function Calling).

Инструменты – это навыки агента. Это «кнопки», которые вы разрешаете ему нажимать во внешнем мире.

Типичные «руки» бизнес-агента:

– Web Search (Поиск в интернете): Способность гуглить актуальные курсы валют или новости конкурентов.

– Email Sender: Способность реально отправить письмо, а не просто сгенерировать его текст.

– Calendar API: Способность забронировать слот в вашем расписании.

– CRM Action: Способность передвинуть сделку на этап «Оплачено» или изменить телефон клиента.

– Code Interpreter: Способность написать и выполнить код (например, чтобы построить график в Excel или посчитать сложную математику).

Принцип минимальных привилегий:

Выдавая агенту руки, вы должны быть осторожны. Если вы дадите ему «руку», которая умеет удалять файлы, он может случайно удалить базу данных. Поэтому в архитектуре агентов мы всегда строго очерчиваем список доступных инструментов.

Хороший агент знает границы своих рук. Если вы попросите его: «Свари мне кофе», а у него нет подключения к API умной кофемашины, он (благодаря Мозгу) ответит: «Извините, у меня нет доступа к управлению физическими объектами».

Итоговая формула:

- Мозг (LLM) = Рассуждает и планирует.
- Память (RAG) = Дает контекст и факты.
- Руки (Tools) = Совершают полезное действие.

Уберите любой элемент, и система рухнет. Без мозга это скрипт. Без памяти это болтун. Без рук это консультант. Вместе – это Агент.

Чем агент отличается от простого скрипта автоматизации

Скептики часто говорят: «Зачем мне этот модный ИИ? Я могу написать скрипт на Python или настроить сценарий в Zapier, который будет делать то же самое».

Это справедливый вопрос. Граница между классической автоматизацией (Automation) и агентной автоматизацией (Agentic Automation) тонкая, но критически важная. Она проходит по линии адаптивности к неопределенности.

Чтобы понять разницу, давайте используем аналогию с транспортом.

– Скрипт (Automation) – это Поезд. Он очень мощный и быстрый. Но он может ехать только по рельсам. Если на рельсах лежит камень – поезд либо остановится, либо сойдет с рельсов. Если рельсы закончатся – он встанет. Поезд не может сказать: «Хм, тут ремонт путей, объеду-ка я через лес».

– Агент (Agentic AI) – это Внедорожник с водителем. Он может ехать по дороге. Но если дорога перекрыта, водитель (LLM) посмотрит на карту, оценит ситуацию и проедет по обочине. Он адаптируется.

1. Жесткая логика vs. Вероятностная логика

Скрипт (If/Then): Работает на жестких правилах.

– *Задача*: Разобрать почту.

– *Логика скрипта*: «ЕСЛИ в теме письма есть слово „Счет“, ТО переслать бухгалтеру».

– *Проблема*: Клиент прислал письмо с темой «Оплата за услуги по договору». Слово «Счет» отсутствует. Скрипт пропустит это письмо. Для скрипта «Счет» и «Оплата» – это абсолютно разные наборы байтов. Чтобы починить это, вам придется вручную дописывать правило: «ЕСЛИ «Счет» ИЛИ «Оплата» ИЛИ «Invoice»...». Вы станете рабом бесконечных правил.

Агент (Intention/Reasoning): Работает на смыслах.

– *Логика агента*: «Проанализируй содержимое письма. Если суть письма касается финансовых документов или просьбы об оплате – перешли бухгалтеру».

– *Результат*: Агент поймет, что «Оплата», «Инвойс», «Где деньги?» и «Кидаю акты» – это всё семантически близкие понятия. Он поймет *смысл*, даже если конкретных ключевых слов нет. Он устойчив к вариативности человеческого языка.

2. Реакция на ошибки (Self-Correction)

Скрипт: Хрупок.

Если API сайта, с которого скрипт собирает цены, вернет ошибку 500, скрипт упадет и пришлет вам лог с красным текстом «Error». Процесс встал.

Агент: Устойчив.

Получив ошибку, агент «подумает»: «Так, сайт недоступен. Что я могу сделать? Я могу подождать 5 минут и попробовать снова. Или я могу попробовать найти этот товар на другом сайте-зеркале. Или я могу сообщить пользователю, что данные старые, но вот прогноз».

Агенты обладают способностью к саморефлексии. Они могут прочитать сообщение об ошибке, понять, что пошло не так (например, «неверный формат даты»), исправить свой же запрос и повторить попытку. Без участия человека.

3. Работа с неструктурированными данными

Скрипт: Любит таблицы и четкие формы.

Скрипт отлично перекладывает цифры из ячейки A1 в ячейку B2. Но если вы дадите скрипту фотографию смятого чека или запись телефонного разговора с клиентом, он беспомощен.

Агент: Всеяден.

Агент может «прочитать» фото чека (используя Vision модели), «услышать» аудио (используя Whisper), понять сарказм в голосе клиента и извлечь из этого хаоса структурированные данные. Агенты – это мост между хаосом реального мира и порядком баз данных.

Когда использовать скрипт, а когда агента?

Не нужно стрелять из пушки по воробьям. Агенты дороже и медленнее скриптов (так как каждый шаг требует обращения к LLM).

– Если задача линейна, предсказуема и не меняется (например, «каждую ночь копировать базу данных на резервный сервер») – используйте скрипт.

– Если задача требует суждения, понимания контекста или работы с «грязными» входными данными (например, «отвечать на отзывы клиентов» или «искать перспективные тендеры») – нанимайте агента.

Типология цифровых личностей: Исследователь, Критик, Исполнитель, Менеджер

При создании мультиагентной системы (Multi-Agent System), главная ошибка новичка – попытка создать одного «Супер-Агента», который умеет всё.

«Пусть он и ищет информацию, и пишет текст, и проверяет ошибки, и публикует».

Это плохая идея. Универсальные промпты работают хуже специализированных. LLM, как и человек, начинает путаться, когда в инструкции слишком много разнородных задач.

Эффективная система строится на разделении труда. Мы создаем команду узких специалистов. В современной прак-

тике (например, в фреймворках CrewAI или AutoGen) выделились четыре классических архетипа цифровых личностей.

1. Исследователь (The Researcher)

– Кредо: «Факты, только факты».

– Инструменты: Поиск в Google (Serper, Tavily), чтение сайтов (Scraper), доступ к Wikipedia или научным базам (Arxiv).

– Характер (System Prompt): Ты дотошный аналитик. Ты не веришь на слово. Ты должен найти первоисточник каждой цифры. Твоя задача – собрать максимально полную, но сырую информацию. Ты не пишешь красивый текст, ты собираешь «мясо».

– Зачем нужен: Чтобы избавить итоговый продукт от галлюцинаций. Он поставляет «чистое топливо» для других агентов.

2. Исполнитель / Креатор (The Doer / Creator)

– Кредо: «Сделаю быстро и красиво».

– Инструменты: Текстовый редактор, генератор кода, генератор картинок (DALL-E), шаблоны документов.

– Характер: Ты талантливый копирайтер (или программист). Твоя задача – взять сухие факты от Исследователя и превратить их в продукт. Ты заботишься о тоне (Tone of Voice), структуре и привлекательности. Ты можешь быть эмоциональным и креативным.

– Зачем нужен: Чтобы упаковать информацию в форму, потребляемую человеком или другой системой.

3. Критик (The Critic / Reviewer)

– Кредо: «Всё переделать. Это никуда не годится».

– Инструменты: Доступ к чек-листам качества, логические валидаторы, сравнение с эталоном.

– Характер: Ты вредный, придирчивый редактор. Твоя задача – найти слабые места в работе Исполнителя. Ты ищешь логические несостыковки, нарушение стиля, опасные формулировки или отсутствие аргументации. Ты никогда не хвалишь, ты только указываешь на ошибки.

– Зачем нужен: Это самый важный агент для контроля качества. Исполнитель склонен «заигрываться» и фантазировать. Критик приземляет его. Исследования показывают, что наличие агента-Критика в цепочке повышает точность ответов на 40—50%. Цикл «Написал – Раскритиковал – Исправил» дает результат на голову выше, чем просто «Написал».

4. Менеджер (The Manager / Orchestrator)

– Кредо: «Соблюдаем сроки и цель».

– Инструменты: Делегирование задач другим агентам, часы, память проекта.

– Характер: Ты руководитель проекта. Ты не делаешь работу руками. Ты получаешь задачу от человека («Напиши отчет о рынке кофе»), разбиваешь её на подзадачи, раздаешь их Исследователю и Исполнителю, следишь, чтобы они не ушли в дебри, и собираешь итоговый результат. Ты решаешь, когда работа готова («Definition of Done»).

– Зачем нужен: Чтобы система работала автономно. Без

Менеджера вам пришлось бы вручную передавать данные от Исследователя к Исполнителю. Менеджер держит в голове «большую картинку» (Big Picture).

Как это работает в связке:

Вы (Человек) говорите Менеджеру: «Нужен пост про тренды ИИ».

– Менеджер зовет Исследователя: «Найди 3 свежих тренда за эту неделю».

– Исследователь серфит интернет и возвращает список ссылок и фактов.

– Менеджер передает это Исполнителю: «Напиши веселый пост на основе этих фактов».

– Исполнитель пишет черновик.

– Менеджер показывает черновик Критику: «Проверь, нет ли тут чуши?».

– Критик замечает: «Второй пункт – это фейк-ньюс, и тон слишком агрессивный».

– Менеджер возвращает Исполнителю: «Перепиши пункт 2 и смягчи тон».

– Исполнитель переписывает.

– Менеджер одобряет и присылает вам готовый текст.

Вся эта драма разыгрывается на сервере за 30 секунд. Вы получаете результат работы целого отдела, заплатив за токены копейки. Это и есть сила ролевой модели.

Глава 2. Почему они должны говорить друг с другом

Проблема «одного большого промпта»: почему универсальные модели глупеют от сложных задач

В начале «золотой лихорадки» генеративного ИИ (2023—2024 годы) в профессиональном сообществе доминировал подход, который мы сейчас, в эпоху агентных систем, называем «Монолитным Промптингом» (Monolithic Prompting). Мы все были его заложниками. И вы, скорее всего, тоже.

Вспомните свой самый сложный запрос к ChatGPT. Вероятно, он выглядел как бесконечное полотно текста, где смешались люди, кони, стилистические требования, факты и запреты.

«Ты – профессиональный маркетолог и юрист. Прочитай этот договор, найди риски, перепиши пункт 5, чтобы он был выгоднее для нас, но не злил контрагента, потом напиши вежливое сопроводительное письмо на английском языке в стиле деловой переписки Оксфорда, а в конце составь таблицу с ключевыми датами».

И что происходило дальше? Модель начинала бодро. Пер-

вый пункт выполнялся блестяще. Второй – неплохо. На третьем начинались странности: стиль письма становился слишком сухим, а в таблице появлялись галлюцинированные даты. К концу ответа модель словно «уставала», теряла нить рассуждений и игнорировала половину ваших инструкций.

Мы привыкли списывать это на «тупость» конкретной версии нейросети. Мы ждали GPT-5, надеясь, что она будет умнее. Но проблема не в мощности модели. Проблема в фундаментальной архитектуре современных нейросетей, которая делает «Один Большой Промпт» тупиковой ветвью эволюции. Чтобы понять, почему агенты неизбежны, нам нужно заглянуть под капот технологии Трансформеров и разобрать феномен, который ученые называют «Размытием Внимания» (Attention Dilution).

Механика внимания: Эффект фонарика в темной комнате

В основе всех современных LLM (Large Language Models) лежит механизм Self-Attention (Само-внимание). Это математический алгоритм, который позволяет модели при генерации каждого следующего слова «оглядываться» на весь предыдущий текст и решать, какие слова важны для текущего момента, а какие – нет.

Представьте, что контекстное окно модели (вся информация, которую вы ей дали) – это огромная темная комната, заставленная мебелью (фактами, инструкциями, условиями). «Внимание» модели – это луч карманного фонарика.

– Когда задача узкая и конкретная («Назови столицу Франции»), луч фонарика сфокусирован в узкую, яркую точку. Модель видит ответ кристально ясно. Вероятность ошибки стремится к нулю.

– Когда вы загружаете в модель «Один Большой Промпт» на 10 страниц с десятком разнородных задач, вы заставляете этот фонарик осветить сразу весь футбольный стадион.

– Что происходит с лучом? Он рассеивается. Свет становится тусклым. Модель вроде бы «видит» всё, но не видит ничего конкретно. Детали в тених теряются. Инструкция «не использовать пассивный залог», написанная в середине промпта, просто тонет в информационном шуме.

Этот феномен научно подтвержден. В 2023 году исследователи из Стэнфорда (Nelson F. Liu et al.) опубликовали знаковую работу «Lost in the Middle» («Потерянные в середине»). Они доказали существование так называемой U-образной кривой производительности (U-shaped performance curve).

Суть открытия пугающая для бизнеса: LLM отлично запоминают то, что написано в самом начале промпта (Primacy Effect) и в самом конце (Recency Effect). Но информация, находящаяся в середине длинного контекста, проваливается в «слепую зону».

Если в вашем «Мега-Промпте» самое важное условие (например, «максимальный бюджет 5000\$») находилось в середине текста, вероятность того, что модель его проигнориру-

ет, достигает 60—70%.

Это не баг, это физика внимания. «Один Большой Промпт» физически не может обеспечить одинаково высокое качество выполнения для всех подзадач одновременно.

Три всадника промпт-апокалипсиса

Помимо технического ограничения внимания, монолитный подход порождает три критические проблемы, которые делают его непригодным для серьезного бизнеса.

1. Шизофрения ролей (Role Confusion)

В примере выше мы просили модель быть одновременно «агрессивным маркетологом» и «осторожным юристом». Для нейросети это взаимоисключающие векторы настройки вероятностей.

– Маркетолог требует высокой «Температуры» ($\text{Temperature} > 0.7$) – параметра, отвечающего за креативность, случайность и неожиданные обороты.

– Юрист требует нулевой «Температуры» ($\text{Temperature} = 0$) – параметра, обеспечивающего максимальную точность, детерминизм и сухость формулировок.

Когда вы запикиваете эти две роли в один промпт, модель вынуждена искать «среднее арифметическое». В результате вы получаете шизофренический продукт: договор с неуместными шутками и рекламный пост, написанный канцеляристом.

Невозможно быть одновременно клоуном и судьей в рамках одной сессии генерации. Универсальность убивает каче-

СТВО.

2. Эффект домино при ошибках (Error Propagation)

В «Одном Большом Промпте» все действия связаны в одну неразрывную цепь. Анализ, планирование, написание кода и его проверка происходят в одном потоке вывода.

Если модель совершает ошибку в самом начале (например, неправильно поняла входные данные), эта ошибка лавиной катится через весь ответ.

Модель не может остановиться, сказать «Ой, я ошиблась в первом абзаце» и переписать его. Она работает авторегрессионно – только вперед. Она будет героически строить логические замки на фундаменте из галлюцинаций.

В конце вы получаете огромный, красивый, связный текст, который полностью неверен, потому что первая предпосылка была ошибочной. В монолитной системе нет «точек сохранения» и «предохранителей».

3. Проблема «Черного ящика» (Debugging Nightmare)

Это боль всех, кто пытался внедрять ИИ в продакшн.

Вы написали сложный промпт. Клиент жалуется, что бот иногда грубит.

Где ошибка?

– Может, вы неудачно описали Tone of Voice?

– Может, модель неправильно интерпретировала вопрос клиента?

– Может, сработал защитный фильтр OpenAI?

– В «Большом Промпте» вы не можете изолировать проблему. Вы меняете одно слово в инструкции, и ломается всё остальное (эффект бабочки). Промпт-инжиниринг превращается в шаманизм: вы боитесь дышать на промпт, который «вроде бы работает». Это тупик для масштабирования.

Смерть «Швейцарского ножа»

Рынок долго верил в мечту об AGI (General Intelligence) – универсальном разуме, который умеет всё. Но практика 2025 года показывает, что бизнес-задачи лучше решают не универсалы, а узкие специалисты.

«Один Большой Промпт» – это попытка забивать гвозди микроскопом. Да, микроскоп тяжелый, им можно забить гвоздь. Но гвоздь будет кривой, а микроскоп сломается.

Мы переходим от парадигмы «Супер-Модель, которая делает всё» к парадигме «Команда средних моделей, каждая из которых делает одно дело идеально».

Вместо того чтобы растягивать внимание модели на километр, мы нарезаем задачу на 10 маленьких кусочков, где для каждого кусочка внимание сфокусировано на 100%.

Именно неспособность одной модели удерживать в голове сложный контекст и привела к рождению мультиагентных систем. Мы перестали бороться с природой нейросетей и начали использовать её сильные стороны.

Мы поняли: чтобы написать книгу, не нужен один гений. Нужен редактор, писатель, корректор и фактчекер. И они не должны жить в одной голове.

Разделяй и властвуй: как разбиение задачи на микро-роли повышает качество на 40%

Если «Один Большой Промпт» – это проблема, то каково решение? Ответ лежит в плоскости организационного управления, а не программирования. Решение называется Декомпозиция Агентных Ролей (Agentic Role Decomposition).

В исследовании, проведенном Microsoft Research в рамках проекта AutoGen (осень 2024), был зафиксирован поразительный результат. Группа из трех «слабых» агентов (на базе модели GPT-3.5), работающих по цепочке, решила сложные задачи по программированию и математике на 40% точнее, чем одна «сильная» модель (GPT-4), которой дали ту же задачу целиком.

Как такое возможно? Как три джуниора могут победить сеньора?

Секрет кроется в методологии «Разделяй и властвуй».

1. Гигиена Контекста (Context Hygiene)

Главное преимущество мультиагентной системы – чистота рабочей памяти каждого отдельного агента.

Когда мы разбиваем задачу, мы создаем для каждого агента свой, изолированный мир.

– Мир Агента-Аналитика: В его контексте лежат толь-

ко сухие цифры, таблицы и отчеты. Ему запрещено думать о красоте слога. Его промпт короткий и жесткий: «Извлеки тренды из таблицы». Лучь его внимания (Self-Attention) бьет точно в цифры. Он не отвлекается.

– Мир Агента-Копирайтера: В его контексте нет таблиц. Ему не нужно тратить вычислительный ресурс на понимание цифр – он получает от Аналитика уже готовые выводы. В его контексте лежат только гайдлайны по стилю и примеры хороших текстов. Его внимание сфокусировано на метафорах и глаголах.

Разделяя контексты, мы устраняем «интерференцию навыков». Модель больше не пытается усидеть на двух стульях. Это повышает IQ каждого отдельного агента в рамках его узкой задачи. Специалист всегда бьет универсала на своем поле.

2. Сила дебатов и самокоррекции (Multi-Agent Debate)

Одиночная модель страдает от когнитивного искажения, известного как Confirmation Bias (Предвзятость подтверждения).

Если модель в начале ответа написала: «Земля плоская», она будет до конца текста придумывать аргументы в пользу этого тезиса, чтобы сохранить логическую связность (Coherence). Ей очень трудно сказать самой себе: «Стоп, я пишу чушь».

В мультиагентной системе мы внедряем механизм Состоя-

зательности (Adversarial Flow).

Мы специально создаем агента-врага. Агента-Критика.

Его единственная задача – не создавать, а разрушать. Он получает ответ первого агента и промпт: «Найди логические ошибки. Найди фактические неточности. Будь безжалостным».

Этот процесс имитирует научную дискуссию или защиту диссертации.

– Агент А выдвигает гипотезу.

– Агент Б разбивает её аргументами.

– Агент А вынужден пересмотреть гипотезу и выдать улучшенную версию.

Исследования MIT (статья «Improving Factuality and Reasoning in Language Models through Multiagent Debate») показывают, что даже 2 раунда таких дебатов снижают уровень галлюцинаций в 3 раза. Истина рождается в споре, и агенты умеют спорить гораздо эффективнее людей – без обид и перехода на личности.

3. Экономика токенов и выбор инструментов (Model Routing)

Разбиение на роли позволяет нам экономить деньги. Это циничный, но важный аспект бизнеса.

Зачем использовать «ядерный реактор» (дорогую модель вроде Claude 3 Opus или GPT-4o) для задачи перефразирования запятых?

В агентной архитектуре мы применяем принцип Model

Routing (Маршрутизация моделей):

– Сложные задачи (Планирование, Креатив): Отдаем топ-моделям. Да, дорого, но тут нужен интеллект.

– Рутинные задачи (Саммаризация, Форматирование JSON, Проверка орфографии): Отдаем дешевым и быстрым моделям (Llama 3, Naiku, GPT-4o-mini).

Разделяя процесс на этапы, мы можем для каждого этапа нанять сотрудника соответствующей квалификации. Вы не нанимаете доктора наук, чтобы он мыл пробирки. Вы нанимаете лаборанта. Точно так же работает агентная сеть. Это снижает стоимость владения системой (ТСО) в разы при сохранении качества финального продукта.

4. Стандартизация процессов (SOPification)

Внедрение агентов заставляет вас, как владельца бизнеса, сделать то, что вы откладывали годами: прописать SOP (Standard Operating Procedures).

Агенты не работают на абстракциях. Им нужен алгоритм.

Чтобы разделить задачу на микро-роли, вы должны сначала разложить свой бизнес-процесс на атомы.

– Не «Веди продажи», а: «1. Квалифицируй лид. 2. Проверь бюджет. 3. Подбери кейс. 4. Сформируй КП».

Этот процесс сам по себе лечит бизнес. Даже если вы завтра выключите ИИ, у вас останутся идеально прописанные регламенты.

Разделение задач превращает «творческий хаос» в «промышленный конвейер». А конвейер, как доказал Генри

Форд, всегда эффективнее ремесленной мастерской.

Качество повышается на 40% не потому, что нейросети стали умнее. А потому, что вы, наконец, начали правильно ставить им задачи. Вы перестали требовать чуда от одного «волшебного промпта» и построили инженерную систему.

Пример из жизни: Как один агент пишет статью, второй – критикует, а третий – верстает (без участия человека)

Теория звучит убедительно, но давайте посмотрим, как это выглядит «в металле». Как именно эти невидимые сущности передают друг другу эстафетную палочку?

Мы разберем классический кейс для малого бизнеса: «Автономная новостная редакция».

Задача: У вас есть Telegram-канал или блог компании. Вы хотите каждый вторник публиковать экспертную статью с обзором новинок рынка.

Проблема: Это занимает 4—6 часов времени живого человека (поиск, чтение, написание, редактура, верстка).

Решение: Команда из трех агентов, которая делает это за 5 минут и 0.5\$ стоимости токенов.

Встречайте нашу команду. У них нет имен, но есть Роли (System Prompts).

АКТ 1. ИЩЕЙКА (The Researcher)

Роль: Агент-Исследователь.

Инструменты:

- search_web (доступ к Google через API Serper/Tavily).
- scrape_url (читалка сайтов).
- Триггер: Таймер (Каждый вторник, 09:00 AM).

Внутренний монолог (Chain of Thought):

«Я получил задачу: найти главные новости по теме «Нейросети в ритейле» за последние 7 дней.

Шаг 1: Формирую запросы. «AI in retail news last week», «Walmart AI implementation 2025», «Amazon AI shopping trends».

Шаг 2: Делаю запрос к API Google... Найдено 15 ссылок.

Шаг 3: Читаю контент. Ссылка 1 – рекламный мусор, пропускаю. Ссылка 2 – серьезный отчет McKinsey, беру. Ссылка 3 – новость про Amazon, беру.

Шаг 4: Извлекаю факты. Мне не нужно писать текст. Мне нужно собрать «мясо». Цифры, даты, имена, цитаты.

Шаг 5: Формирую JSON-файл для следующего агента».

Артефакт на выходе (Output):

Агент-Исследователь не пишет статью. Он выдает сухой файл dossier. json. В нем нет ни капли эмоций, только структурированные данные:

```
{
  «topic»: «AI in Retail»,
  «key_findings»: [
    {«source»: «TechCrunch», «fact»: «Amazon внедрил умные тележки в 50 магазинах», «date»: «12.12.2025»},
```

{ «source»: «McKinsey Report», «fact»: «Внедрение ИИ снизило затраты на логистику на 15%», «link»: «...» }
]
}

АКТ 2. ДУШНИЛА (The Critic / Editor)

Роль: Агент-Редактор.

Инструменты: База знаний (Ваш Brand Book, примеры ваших лучших постов).

Входные данные: dossier. json от Исследователя.

Внутренний монолог:

**«Я получил досье. Теперь моя задача – проверить его на соответствие редакционной политике.*

Проверяю факт №1 (Amazon). Это действительно новость или старый пресс-релиз? Ага, новость свежая. Подходит.

Проверяю факт №2 (McKinsey). Слишком скучно. Нужно найти пример, понятный малому бизнесу, а не корпорациям.

Решение: Я отвергаю факт №2 и прошу Исследователя найти пример для малого бизнеса.

(Цикл возвращается назад. Исследователь доискивает информацию).

...

Теперь все факты утверждены. Я составляю План Статьи (Outline).

– Вступление: Проблема очередей в магазинах.

– Основная часть: Как Amazon решает это (на базе факта 1).

– Вывод: Что делать обычному магазину?

– Важное примечание для Писателя: Не используй слово «инновационный», босс его ненавидит». *

Артефакт на выходе:

Файл `outline_approved.txt` с жесткой структурой и инструкциями.

АКТ 3. ТВОРЕЦ (The Writer / Formatter)

Роль: Агент-Копирайтер.

Инструменты: Нет (ему запрещено ходить в интернет, чтобы не нахвататься чуши).

Входные данные: `outline_approved.txt` и `dossier.json`.

Внутренний монолог:

«У меня есть утвержденный план и проверенные факты.

Моя задача – просто превратить это в текст.

Я пишу в стиле «Информационный стиль». Короткие предложения. Без воды.

Использую факты из досье.

Вставляю шутку про кассиров во вступлении (разрешено гайдлайном).

Теперь верстка. Мне нужен Markdown.

Заголовки делаю через H2 (##).

Ключевые цифры выделяю жирным (15%).

Вставляю плейсхолдер для картинки: [Робот с тележкой]».

Артефакт на выходе:

Готовый файл `article_final.md`, который можно копировать и вставлять в CMS.

ФИНАЛ: ЧЕЛОВЕК НА БАЛКОНЕ

В 09:05 вы получаете уведомление в Telegram.

«Шеф, материал готов. Исследователь нашел 3 новости, Критик забраковал одну, но утвердил две. Писатель сверстал текст. Прочитать и опубликовать?»

Под сообщением две кнопки: [Опубликовать] и [Передать].

В чем магия этого примера?

– Отсутствие галлюцинаций. Писатель физически не может выдумать цифру, потому что у него нет доступа в интернет, он работает только с тем, что дал Исследователь. А Исследователь не пишет текст, он только копирует факты. Цепь разорвана в нужном месте.

– Стиль. Если вам не нравится стиль текста, вы меняете промпт только у третьего агента (Писателя). Исследователя и Критика трогать не надо. Система модульная.

– Скорость. Пока вы наливали кофе, три нейросети провели планерку, поссорились, помирились и сделали работу.

Это не фантастика. Такая схема (Researcher -> Critic -> Writer) – это «Hello World» агентных систем. Она собирается на No-Code платформах (вроде n8n или Make) за один вечер, а экономит сотни часов в год.

Именно так выглядит разделение труда в цифровую эпо-

ху. Мы не заставляем одного робота делать всё. Мы строим конвейер, где каждый робот закручивает свою гайку, но делает это безупречно.

Часть II. Штатное расписание: Архитектура системы

Глава 3. Роль СЕО: Ваша новая работа

Как перестать быть «оператором промпта» и стать «архитектором системы»

В 2023 году мир разделился на две неравные группы. В первой, самой многочисленной, люди осваивали профессию «Промпт-инженера». Они учили заклинания: «Действуй как профессиональный копирайтер», «Используй метод Chain-of-Thought», «Дыши глубоко» (да, был и такой миф, что призыв к глубокому дыханию улучшает ответы нейросети). Эти люди гордились тем, что умеют «уговорить» машину выдать нужный результат.

Но была и вторая, крошечная группа людей. Они не писали длинных поэм в ChatGPT. Они рисовали квадратики и стрелочки на белых досках. Они не пытались стать лучши-

ми собеседниками для бота. Они строили заводы, где боты работали у станков. Сегодня, в 2025 году, мы видим результат: первые остались с забавной, но умирающей профессией «оператора чата», а вторые стали владельцами автономных агентств.

Переход от Оператора к Архитектору – это самый сложный ментальный прыжок, который вам предстоит совершить в этой книге. Это отказ от микроменеджмента в пользу системного дизайна.

Ловушка «Волшебной Кнопки»

Почему так трудно перестать быть оператором? Потому что чат-интерфейс (Chat UI) – это наркотик. Он дает мгновенный дофаминовый отклик. Вы написали вопрос – через секунду получили ответ. Это создает иллюзию контроля и продуктивности.

– *Оператор думает:* «Я сейчас быстро сам всё напишу в чат, это быстрее, чем настраивать какую-то систему».

– *Архитектор знает:* «Любое действие, которое я делаю руками больше двух раз, должно быть превращено в алгоритм, иначе я становлюсь узким местом своего бизнеса».

Оператор работает в парадигме **«Человек – > ИИ»**.

Архитектор работает в парадигме **«Событие – > Система (Агент 1 + Агент 2) – > Результат»**.

Пять заповедей Архитектора ИИ-систем

Чтобы совершить этот переход, вам нужно перепрощить свое отношение к взаимодействию с машинами. Вот пять

принципов, которые отличают системный подход от ручного труда.

1. Мыслить Потоками (Flows), а не Запросами (Prompts)

Оператор фокусируется на том, *как спросить*. Архитектор фокусируется на том, *откуда приходят данные и куда они уходят*.

Представьте, что вы строите водопровод.

– Промпт – это кран. Вы можете долго полировать ручку крана, чтобы она блестела. Но если к дому не подведена труба, воды не будет.

– Ваша задача – проложить трубы.

– *Практика*: Перестаньте открывать ChatGPT, когда возникает задача. Откройте блокнот или Miro. Нарисуйте кружочек «Входящее письмо». Нарисуйте стрелочку к квадратику «Анализ». От него – стрелочку к «Генерация ответа». Вы только что начали мыслить как архитектор. Вы создали граф выполнения, а не диалог.

2. Принцип «Стеклянного ящика» (Glass Box vs Black Box)

Когда вы переписываетесь с нейросетью в чате, это «Черный ящик». Вы не знаете, почему она ответила именно так. Вы не можете залезть ей в голову.

Архитектор строить систему как «Стеклянный ящик».

– Вместо одного гениального промпта «Сделай отчет», вы разбиваете процесс на 5 прозрачных шагов.

– Шаг 1: Извлеки цифры (мы видим этот файл).

– Шаг 2: Посчитай разницу (мы видим формулу).

– Шаг 3: Напиши выводы (мы видим текст).

– Если на Шаге 3 ошибка, вы не переписываете промпт целиком. Вы видите, что ошибка произошла на Шаге 1 (неверно извлечена цифра). Вы чините только этот узел. Архитектура дает контроль, которого лишен оператор.

3. Отказ от Детерминизма в пользу Вероятностей

Программисты привыкли: если $A + B$, то всегда C .

Архитектор ИИ понимает: если $A + B$, то C будет в 95% случаев, а в 5% нейросеть начнет цитировать Шекспира.

Ваша работа – не пытаться свести эти 5% к нулю (это невозможно), а построить «уловители ошибок» (Guardrails).

– *Оператор*: Злится, что нейросеть ошиблась.

– *Архитектор*: Ставит после нейросети блок-валидатор (простой код), который проверяет: «В ответе есть цифры? Если нет – верни на переделку». Вы проектируете систему с правом на ошибку, которая сама себя исправляет.

4. OODA Loop: Наблюдай, Ориентируйся, Решай, Действуй

Это военная концепция (петля Джона Бойда), которая идеально описывает работу автономного агента.

Оператор сам проходит этот цикл. Он читает письмо (Наблюдает), понимает его суть (Ориентируется), решает ответить (Решает) и пишет промпт (Действует).

Архитектор передает этот цикл машине.

– *Наблюдение*: Вебхук от CRM (пришла заявка).

– *Ориентация*: Классификация интента (клиент злой, добрый или просто спрашивает).

– *Решение*: Маршрутизация (злого – человеку, доброго – агенту-продажнику).

– *Действие*: API вызов (отправка сообщения).

– Ваша задача – настроить правила для фазы «Решение».

Все остальное делает софт.

5. Масштабируемость как религия

Оператор линейн. Чтобы обработать 100 заявок, ему нужно в 100 раз больше времени.

Архитектор мыслит категориями масштаба.

«Если я построю эту систему для обработки одного отзыва, будет ли она работать для 10 000 отзывов, пока я сплю?»

Если ответ «Нет» (нужно что-то кликать руками) – это плохая архитектура. Архитектор не строит ничего, что требует его постоянного участия. Он строит только то, что работает автономно.

Практическое упражнение: Тест на Архитектора

Посмотрите на свою последнюю переписку с ChatGPT.

Если там есть фразы: «Нет, переделай», «Ты забыл про...», «Сделай короче» – вы Оператор. Вы тратите время на исправление брака в ручном режиме.

Архитектор бы сделал так:

– Увидел, что модель часто забывает про контекст.

– Добавил бы в систему шаг «Проверка контекста» перед

генерацией.

– Добавил бы шаг «Авто-сокращение» после генерации.

– Больше никогда не писал бы эти замечания руками.

Стать архитектором – значит перестать играть в «Вопрос-Ответ» и начать играть в LEGO. Это переход от гуманитарного общения к инженерному конструированию. И именно за это рынок готов платить. За умение просить – платят копейки. За умение строить конвейеры просьб – платят миллионы.

Карта процессов: выявляем рутину, которую можно делегировать агентам

Прежде чем нанимать армию цифровых сотрудников, нужно понять, куда их расставить. Самая частая ошибка бизнеса – попытка автоматизировать хаос. «У нас бардак в продажах, давайте внедрим ИИ-агента, пусть он разберется».

Спойлер: Агент не разберется. Он просто автоматизирует бардак. Вы будете терять клиентов со скоростью света, 24/7, без перерывов на обед.

Чтобы этого не случилось, нам нужна **Карта Процессов**. Это рентгеновский снимок вашего бизнеса, который покажет, где прячутся деньги, а где – черные дыры времени.

Шаг 1. Охота на «Теневую работу» (Shadow Work)

В любой компании есть официальные должностные инструкции, а есть реальность.

Официально: «Менеджер ведет переговоры».

Реально: «Менеджер 30% времени ищет телефон клиента в старых письмах, 20% времени копирует данные из Excel в CRM, 10% времени правит запятые в договоре и только 40% – говорит».

Вот эти 60% – это «Теневая работа». Невидимая, неоплачиваемая (по сути), но неизбежная рутина.

Агенты созданы именно для уничтожения теневой работы.

Упражнение «Неделя с таймером»

Попросите ключевых сотрудников (или себя) в течение недели записывать каждое действие, которое длится больше 5 минут. Не задачи («Продажи»), а действия («Копировал email», «Искал файл», «Писал фоллоу-ап»).

Вы ужаснетесь. 80% действий – это работа курьера данных.

Шаг 2. Матрица Делегирования (Матрица Эйзенхауэра для ИИ)

Когда список действий готов, мы должны прогнать каждое действие через сито. Не всё можно и нужно отдавать агентам.

Используйте матрицу 2x2, где оси – это **Сложность когнитивная** (нужен ли мозг?) и **Частотность** (как часто это происходит?).

Квадрант 1: Высокая частота / Низкая сложность (Роботизация)

– *Примеры:* Ответы на частые вопросы (FAQ), копирова-

ние данных из заявки в CRM, выставление счетов, сбор статистики.

– *Вердикт:* **Автоматизировать немедленно.** Это идеальная еда для агентов. Тут даже не всегда нужны сложные LLM, часто хватит простых скриптов или дешевых моделей. Это зона мгновенной окупаемости.

Квадрант 2: Высокая частота / Высокая сложность (Augmentation – Усиление)

– *Примеры:* Написание персонализированных писем клиентам, анализ причин отказа, создание контента для соцсетей, техническая поддержка второго уровня.

– *Вердикт:* **Внедрять агентов-помощников (Co-pilots) или Human-in-the-loop.** Полностью отдать страшно (высокая цена ошибки), но агент может делать 90% черновика. Человек только проверяет и жмет «Отправить». Это повышает производительность сотрудника в 5—10 раз.

Квадрант 3: Низкая частота / Низкая сложность (Игнорирование)

– *Примеры:* Заказ воды в офис раз в месяц, поздравление партнера с днем рождения раз в год.

– *Вердикт:* **Не трогать.** Настройка агента займет больше времени, чем само действие за 5 лет. Делайте руками.

Квадрант 4: Низкая частота / Высокая сложность (Экспертиза)

– *Примеры:* Стратегическая сессия, разрешение сложного конфликта с VIP-клиентом, найм топ-менеджера.

– *Вердикт: Только человек.* ИИ здесь может выступать только как советник (Второе мнение), но не как исполнитель.

Шаг 3. Декомпозиция процесса до атомов (SOPification)

Выбрав процессы из Квадрантов 1 и 2, мы начинаем их препарировать. Агент не поймет задачу «Веди соцсеть». Ему нужна инструкция для идиота (или для гениального ребенка).

Возьмем процесс «Обработка входящей заявки».

Как это видит человек: «Ну, приходит заявка, я смотрю, нормальный ли клиент, если да – звоню».

Как это видит Архитектор (Карта для агента):

– **Триггер:** Новая строка в Google Sheets (источник: форма на сайте).

– **Действие 1 (Агент-Классификатор):** Прочитать поле «Комментарий». Оценить тональность (Positive/Negative). Проверить, есть ли мат или спам-слова.

– **Ветвление (Router):**

– ЕСЛИ спам – > удалить строку, отправить отчет в Slack «Спам».

– ЕСЛИ клиент – > перейти к шагу 3.

– **Действие 2 (Агент-Обогатитель):** Взять email. Погуглить домен компании (через Clearbit или Serper). Найти размер компании и индустрию. Дописать в таблицу.

– **Действие 3 (Агент-Продавец):** Сгенерировать при-

ветственное письмо, используя шаблон №3, вставить имя и упомянуть индустрию (из шага 4).

– **Финал:** Создать черновик в Gmail.

Видите разницу? Мы разбили абстрактное «смотрю» на конкретные алгоритмические шаги.

Золотое правило: Если вы не можете написать SOP (Standard Operating Procedure) на бумаге в виде блок-схемы «Да/Нет», вы не готовы нанимать агента. Агент – это не замена процесса, это его кристаллизация.

Шаг 4. Оценка данных (Input/Output Audit)

Агенты питаются данными. Для каждого процесса в карте ответьте на вопросы:

– **Вход (Trigger):** Как агент узнает, что пора работать? (Письмо, время, файл, вебхук). «Когда у меня будет настроение» – это не триггер.

– **Контекст (Context):** Где лежит информация, необходимая для решения? (В голове менеджера? В блокноте? В CRM?). Если информация в голове – агент бесполезен. Сначала оцифруйте базу знаний.

– **Выход (Action):** Что является физическим результатом работы? (Файл, запись в базе, отправленное сообщение). «Ощущение, что мы молодцы» – это не результат.

Составив такую карту, вы увидите свой бизнес как печатную плату. Вы увидите, где ток течет свободно, а где – сопротивление. И именно в точки сопротивления мы будем втыкать наших ИИ-агентов в следующих главах.

Чек-лист готовности бизнеса к автономности (без участия человека)

Вы нарисовали схемы, у вас горят глаза, вы готовы купить подписку на все ИИ-сервисы мира. Стоп.

Прежде чем пустить электричество по проводам, нужно проверить проводку. Иначе сгорит дом.

Внедрение автономных агентов – это не установка фотошопа. Это хирургическое вмешательство в процессы.

Вот честный чек-лист, который отрезвит вас и сэкономит бюджет. Если вы не ставите галочку хотя бы в одном пункте раздела «Критично» – отложите внедрение и исправьте фундамент.

Блок 1: Техническая гигиена (Критично)

– Данные доступны в цифровом виде (Machine Readable).

– Плохо: Заявки приходят в виде фоток рукописных бланков в WhatsApp. Прайс-лист – это скан PDF 2010 года.

– Хорошо: Данные в CRM, Excel, Google Sheets, JSON, SQL. Текст можно выделить курсором и скопировать.

– Почему: Агенты могут использовать OCR (распознавание текста), но это добавляет стоимость, время и 5—10% ошибок. Надежная автоматизация строится на цифровых данных.

– Доступ по API (или хотя бы Webhooks).

– *Плохо*: Вы используете самописную CRM 90-х годов, у которой нет API, или закрытый проприетарный софт, куда можно зайти только через удаленный рабочий стол.

– *Хорошо*: Ваш софт (CRM, почта, таск-трекер) дружит с миром (имеет REST API, есть в интеграторах типа Make/Zapier).

– *Почему*: Агенту нужны «руки». API – это и есть руки. Без API агент – это мозг в банке, который может только кричать, но не может ничего сделать.

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.