

ДААННЫЕ, ГОТОВЫЕ ДЛЯ ИИ



 ODYSSEY
Consulting Group

Почему **MDM** становится
фундаментом для **GenAI**
в корпорациях

Что такое MDM • Что делает GenAI • Что выбрать: RAG, RIC, LLM

ДАННЫЕ, ГОТОВЫЕ ДЛЯ ИИ

Почему MDM становится фундаментом для GenAI в корпорациях

Сегодня у крупных компаний странная ситуация. Данных становится все больше, а доверия к ним часто меньше. Почти каждая организация хочет стать data-driven, но на практике корпоративный ИТ-ландшафт обычно выглядит не как единая система, а как набор CRM, ERP, HRM, SCM, облачных сервисов, локальных баз и исторически накопленных приложений. В итоге информация о клиентах, продуктах, поставщиках и сотрудниках распадается на куски, дублируется и быстро устаревает.

На этом фоне бизнес одновременно входит в эпоху генеративного ИИ. Большие языковые модели уже меняют клиентскую поддержку, аналитику цепочек поставок, финансовое прогнозирование, внутренний поиск и автоматизацию знаний. Но здесь быстро выясняется неприятная вещь: LLM не понимают, где истина. Они не проверяют факты так, как это делает человек. Они предсказывают наиболее вероятное продолжение текста на основе контекста и данных, которые получили.

Если в этот контекст попадают ошибки, противоречия или устаревшие сведения, модель начинает производить убедительно звучащие, но неверные ответы. В бизнесе это называют «галлюцинациями». На практике последствия гораздо приземленнее и опаснее: неверные цифры в аналитике, ошибки в работе с клиентом, искажение условий договора, проблемы с безопасностью, комплаенсом и репутацией.

Именно поэтому разговор о GenAI в корпорации почти всегда должен начинаться не с выбора модели, а с вопроса: на каких данных она будет работать? Здесь и появляется **MDM, Master Data Management**. По сути, это не «гигиена данных» и не вспомогательный ИТ-процесс. Это слой доверенного знания, без которого RAG, RIC и агентные ИИ-сценарии начинают опираться на хаос.



ОГЛАВЛЕНИЕ

Что такое MDM и почему без него все быстро ломается	4
«Золотая запись» простым языком	4
Как MDM создаёт «золотую запись»	5
Что на самом деле делает GenAI	6
Почему «галлюцинации» в бизнесе лучше называть конфабуляциями	6
Как бизнес решает эту проблему: RAG	7
Как устроен конвейер RAG	7
Когда RAG уже недостаточно: RIG	8
Как устроен RIG	8
RAG, RIG и LLM с длинным контекстом: что выбрать	8
Почему без MDM нельзя обеспечить качество данных на входе и качественный результат на выходе	9
Почему для LLM важно иметь взгляд на «сущность целиком»	10
Безопасность и комплаенс: еще одна причина, почему без MDM будет сложно	10
Почему в повестке всё чаще возникает тема микро-баз данных и дата-продуктах	11
Что делают крупные вендоры MDM уже сейчас	11
Informatica CLAIRE GPT	12
SAP Master Data Governance	12
Pimcore	12
Reltio	13
K2view	14
Куда все это идет: от чат-ботов к агентам	14
Что, скорее всего, произойдет дальше	15
Вывод	15



ЧТО ТАКОЕ MDM И ПОЧЕМУ БЕЗ НЕГО ВСЁ БЫСТРО ЛОМАЕТСЯ

Чтобы понять роль MDM, важно развести **три типа корпоративных данных**.

Первый тип — транзакционные данные

Это все, что возникает в ежедневных операциях: платежи, чеки, отгрузки, сервисные заявки, движения по складу, обращения в поддержку. Эти данные очень изменчивы и фиксируют конкретное событие в конкретный момент.

Второй тип — аналитические данные.

Это уже агрегаты, метрики, дашборды, исторические тенденции, прогнозы и показатели эффективности.

Третий тип — мастер-данные

И вот они для бизнеса базовые. Это справочная информация о ключевых сущностях: клиентах, продуктах, поставщиках, сотрудниках, материалах, локациях. Такие данные меняются нечасто, но когда меняются, последствия серьёзные. У одного продукта могут быть сотни атрибутов: размеры, состав, материал, сертификация, SKU, упаковка, ограничения по хранению и так далее.

«ЗОЛОТАЯ ЗАПИСЬ» ПРОСТЫМ ЯЗЫКОМ

В любой крупной компании разные подразделения живут в своих системах. Продажи работают в CRM, HR в HRM, логистика в SCM, финансы в ERP. Исторически каждая из этих систем хранит собственную версию клиента, продукта или поставщика. Отсюда и типичный корпоративный бардак: в одной системе старый телефон, в другой опечатка в названии компании, в третьей устаревший юридический адрес.

Если объяснять без технических терминов, нам нравится аналогия с большим медицинским центром. Представьте, что хирургия, регистратура, лаборатория и бухгалтерия ведут отдельные бумажные картотеки пациентов и между собой их не синхронизируют. Пациентка поменяла фамилию. Для лаборатории это уже новый профиль, для бухгалтерии все ещё старый, а хирургия может не найти нужные результаты анализов из-за несовпадения имени. Бухгалтерия, в свою очередь, выставит счёт человеку, которого «по их данным» уже не существует.

В этой картине MDM — это центральный реестр больницы. Любое изменение, которое происходит в любом подразделении, попадает в центр, проходит проверку, подтверждается и потом распространяется дальше по системам. В результате появляется единый эталон.

Для бизнеса это не абстракция. Единая «золотая запись» снижает ошибки в биллинге, убирает дубли в маркетинговых рассылках и помогает соблюдать регуляторные требования. Gartner оценивает отсутствие такой централизации примерно в 9,7 миллиона долларов ежегодных потерь в крупной организации. Потери идут из неэффективности процессов и решений, принятых на плохих данных.



КАК MDM СОЗДАЁТ «ЗОЛОТУЮ ЗАПИСЬ»

«Золотая запись» не появляется сама по себе. За её появлением стоит довольно жёсткий технологический конвейер.

1 Сначала производятся очистка и стандартизация

Данные приходят из разных систем и почти всегда в разных форматах. На этом этапе исправляются типографические ошибки, даты приводятся к одному стандарту, например ISO 8601, адреса и телефоны унифицируются по единым справочникам. Без этой базовой нормализации дальше просто нечего сопоставлять.

2 Следующий этап — сопоставление и дедупликация.

Здесь система пытается понять, две записи из разных систем - это один и тот же объект или нет. То есть система сопоставляет дублирующие записей об одной сущности

Метрика	Как работает	Плюсы и ограничения
Детерминированное сопоставление	Ищет точные совпадения по уникальным полям: ИНН, email, SKU, номеру соцстрахования и т.д.	Даёт очень высокую точность, если уникальный ключ есть и он корректен. Но в реальных системах ключи часто отсутствуют, повреждены или введены с ошибками.
Нечёткое сопоставление	Использует вероятностные алгоритмы и скоринг сходства по текстовым полям.	Позволяет находить скрытые дубликаты, например «ООО Ромашка», «Ромашка ООО» и «Romashka LLC». Но требует тонкой настройки порогов уверенности, чтобы избежать ложных объединений.



Современные MDM-платформы все чаще усиливают этот этап машинным обучением.

3 После сопоставления включается логика слияния.

На этом этапе начинают действовать **правила выбора приоритетных данных**. Они определяют, какому источнику система больше доверяет по каждому конкретному атрибуту. Юридический адрес логично брать из биллинговой системы. Мобильный телефон - из CRM, потому что им чаще пользуются аккаунт-менеджеры и быстрее замечают изменения.

4 Последний этап — обогащение данных.

Внутренние записи дополняются данными из внешних доверенных источников. Например, карточку корпоративного клиента можно обогатить данными из ЕГРЮЛ и реестра МСП ФНС: автоматически подтянуть ИНН, ОГРН, КПП, юридический адрес, коды ОКВЭД, статус компании и признак субъекта МСП. Это помогает точнее идентифицировать контрагента, снижать количество дублей и ошибок и использовать в процессах более полные и доверенные данные.



Так формируется профиль, на который уже можно опираться ИИ.

ЧТО НА САМОМ ДЕЛЕ ДЕЛАЕТ GenAI

Чтобы понять, почему GenAI так зависят от качества мастер-данных, полезно снять с них ореол «разумного собеседника».

Модели вроде GPT, Gemini или Llama — это очень большие нейронные сети на архитектуре трансформеров. На человеческом уровне они выглядят умными. На математическом уровне это сложнейшие механизмы вероятностного предсказания.

Во время предварительного обучения модель читает огромные массивы текста: интернет, книги, статьи, открытые базы. Из этого она извлекает статистические связи между словами, фразами и концепциями. Когда мы задаём запрос, модель не «думает» в привычном смысле. Она прогнозирует следующее наиболее вероятное слово, потом следующее, и так шаг за шагом.

Для бизнеса обычно хорошо работает такая аналогия: GenAI — это гениальный стажёр с энциклопедической памятью и нулевым жизненным опытом. Он перечитал все на свете, умеет красиво формулировать, резюмировать, переводить, рисовать, фотографировать, предлагать гипотезы и даже звучать убедительно. Но есть несколько проблем.

Во-первых, он не понимает смысл так, как понимает человек. Во-вторых, он не знает, что происходило после завершения его обучения. В-третьих, у него нет врождённого доступа к вашим внутренним данным. И, наконец, когда у него нет факта, он нередко заполняет пробел правдоподобной выдумкой.



ПОЧЕМУ «ГАЛЛЮЦИНАЦИИ» В БИЗНЕСЕ ЛУЧШЕ НАЗЫВАТЬ КОНФАБУЛЯЦИЯМИ

Термин «галлюцинации» закрепился в индустрии, но в корпоративном контексте точнее говорить о конфабуляциях. Галлюцинация — это восприятие несуществующего. Конфабуляция — это ложное воспоминание или выдуманная связка фактов, которая появляется не из желания обмануть, а чтобы закрыть пробел в знаниях.

Когда пользователь спрашивает у базовой LLM: «Какой была выручка по новому продукту в третьем квартале?» или «Действует ли скидка 15% по контракту с поставщиком X?», модель сталкивается с фактом, которого у неё нет. Но её внутренняя логика заточена на одно: дать связный, грамматически корректный и уверенный ответ. В результате она может просто придумать цифру или условие контракта, опираясь на статистические шаблоны языка.

Для корпоративной среды это тупик. В финансовой отчётности, медицине, юридической поддержке и клиентском обслуживании такие ошибки недопустимы. Поэтому изолированная LLM без внешнего источника истины в компании почти всегда приводит к рискам, а не к ценности.

КАК БИЗНЕС РЕШАЕТ ЭТУ ПРОБЛЕМУ: RAG

На практике компании не хотят и не могут постоянно переобучать модель под каждый новый документ, договор или отчёт. Для этого и появилась архитектура Retrieval-Augmented Generation, RAG.

Идея простая: перед тем как ответить, система ищет релевантные данные во внешних источниках и только потом передаёт их модели. То есть LLM отвечает не «из головы», а на основе найденного контекста.

Если продолжать аналогию со стажёром, RAG — это экзамен с открытой книгой. Стажёр не угадывает ответ. Он сначала идёт к полке с корпоративными документами, находит нужный регламент, страницу договора, письмо или таблицу, читает факты и только потом формулирует ответ.

Есть и другая понятная аналогия: GPS-навигатор. Базовая LLM — это бумажная карта. Она знает общую географию, но быстро устаревает. RAG превращает её в навигатор, который перед прокладкой маршрута сверяется с актуальными данными о пробках, ремонтах и перекрытиях.

КАК УСТРОЕН КОНВЕЙЕР RAG

Конвейер RAG состоит из нескольких этапов:

• Сбор и агрегация

Система забирает структурированные данные из таблиц и неструктурированные данные из PDF, регламентов, email, логов чатов и других корпоративных источников.

• Фрагментация

Поскольку у моделей ограничено контекстное окно, большие документы разбиваются на чанки — управляемые смысловые фрагменты. Хорошо настроенная фрагментация старается не резать документ в случайных местах, а сохранять логические блоки и абзацы.

• Векторизация

Каждый чанк превращается в числовой вектор. В этом многомерном пространстве похожие по смыслу тексты оказываются рядом. Например, «автомобиль» и «машина». Одновременно к векторам добавляются метаданные: автор, дата, категория, источник.

• Загрузка в векторную базу данных

Векторы отправляются в специализированные хранилища вроде Pinecone, Weaviate или в расширения вроде pgvector для PostgreSQL. Эти системы умеют быстро находить близкие по смыслу записи.

• Извлечение

Когда пользователь задаёт вопрос, его тоже векторизуют. Дальше векторная база ищет наиболее близкие по смыслу фрагменты. До этого запрос обычно проходит предобработку: токенизацию, очистку от стоп-слов и другие шаги.

• Генерация с опорой на контекст

Найденные фрагменты подставляются в скрытый системный промпт вместе с вопросом пользователя. Модель получает не просто запрос, а вопрос плюс набор релевантных фактов. Именно на них она и должна опереться при генерации ответа.

Сильная сторона RAG в том, что модель можно не переучивать каждый раз.

Она работает с актуальными данными, а ответы можно сопровождать ссылками на источник. Это даёт прозрачность и возможность аудита.

Но у RAG есть предел.

Поиск обычно выполняется один раз, до начала генерации. Для простых или средних сценариев этого хватает. Для сложной аналитики, где нужно на лету собрать много чисел, сверить несколько метрик и пройти цепочку уточнений, одного захода может быть мало.

КОГДА RAG УЖЕ НЕДОСТАТОЧНО: RIG

Следующий шаг — Retrieval-Interleaved Generation, RIG. Это подход для задач, где нужно не просто «подтянуть контекст», а буквально проверять факты по ходу генерации текста.

Если RAG — это экзамен с открытой книгой, то RIG — это ведущий новостей в прямом эфире, у которого в наушнике сидит команда аналитиков. Ведущий начинает говорить, доходит до места, где нужен точный показатель, делает микропаузу, быстро получает цифру от аналитиков, озвучивает её и идёт дальше. И так много раз за один ответ.

КАК УСТРОЕН RIG

Здесь поиск и генерация переплетены

- Модель начинает формировать ответ и в процессе определяет, где ей не хватает конкретного факта или числа.
- Она останавливается и генерирует запрос к внешнему источнику. В продвинутых реализациях это не жёсткий SQL, а естественный язык. За счёт этого внешний источник фактов можно использовать как универсальный API.
- Доверенная система возвращает точное значение.
- Модель вставляет полученный факт в текущий ответ и продолжает.

Этот цикл может повторяться много раз.

Хороший пример - DataGemma от Google

В этом исследовательском проекте модели семейства Gemma связаны с графом знаний Data Commons. В нем больше 240 миллиардов нормализованных статистических показателей из ООН, ВОЗ, CDC, данных переписей населения и других авторитетных источников. DataGemma RIG дообучена распознавать статистические вопросы. Если пользователь просит сравнить безработицу и ВВП двух стран, система не пытается импровизировать. Она по очереди достаёт нужные показатели из Data Commons и только потом собирает финальный ответ со ссылками на первоисточники.

RAG, RIG И LLM С ДЛИННЫМ КОНТЕКСТОМ: ЧТО ВЫБРАТЬ

Для ИТ-директора или архитектора вопрос обычно не в том, «какая технология моднее», а в том, какая архитектура подходит конкретному сценарию.

Характеристика	RAG	RIG	LLM с длинным контекстом
Механика извлечения	Поиск выполняется один раз до генерации, дальше модель работает со статическим контекстом.	Поиск встроен в процесс генерации, модель делает серию микро-запросов по ходу ответа.	Внешний поиск не нужен, весь массив документов загружается сразу в контекст.
Архитектурный дизайн	Модульная схема, LLM и база данных разделены, удобно масштабировать и менять компоненты.	Более монолитный и глубоко интегрированный подход, обычно требует специально адаптированных моделей.	Упирается в вычислительные ресурсы, GPU и коммерческие API вроде Gemini 1.5 Pro.

Характеристика	RAG	RIG	LLM с длинным контекстом
Влияние на галлюцинации	Сильно снижает ошибки в сценариях с документацией, базами знаний и перепиской.	Максимально надёжен там, где важны числа, статистика и математическая точность.	Может страдать от эффекта «потерялся посередине», когда модель начинает забывать факты из середины большого контекста.
Задержка	Средняя: один поиск, потом быстрая генерация.	Более высокая: каждый вызов внешнего API тормозит поток генерации.	Часто слишком большой срок получения первого токена, потому что модели нужно сначала прочитать огромный входной контекст.
Стоимость	Обычно ниже: в модель передаются только релевантные чанки, основные расходы идут на векторную БД.	Дороже в настройке и эксплуатации: fine-tuning, API-слой, множественные запросы.	Самая дорогая схема по инференсу из-за постоянной обработки огромных объёмов токенов.
Лучшие сценарии	Внутренний поиск, база знаний, виртуальные помощники, анализ контрактов, генерация черновиков отчётов.	Финансовая аналитика, фактчекинг, медицинские рецепты, сложные многошаговые исследования.	Разовый анализ небольших статичных наборов документов.

Итог довольно практичный. Для большинства корпоративных кейсов RAG — это разумный баланс гибкости, скорости и стоимости. Для индустрий с высокой ценой ошибки, например финансов, медицины и науки, RIG может оказаться незаменимым. LLM с длинным контекстом удобны в отдельных сценариях, но экономически и архитектурно не всегда оптимальны.



При этом и RAG, и RIG упрутся в один и тот же вопрос: откуда берутся данные, которые они достают? Если источник плохой, любая надстройка сверху только ускоряет распространение ошибки

— Ремзи Эшматов, архитектор MDM решений в Odyssey Consulting Group.



Почему без MDM нельзя обеспечить качество данных на входе и качественный результат на выходе

Очень распространенная ошибка при запуске GenAI в компании - думать, что RAG или RIG сами по себе умеют отличать истину от мусора. Не умеют. Это транспортные механизмы. Они быстро находят и доставляют контекст, но не являются встроенным арбитром правды.

Если RAG напрямую подключить к сырым корпоративным системам, где по одному клиенту есть десятки противоречивых профилей, устаревшие цены, дубли заказов и разъехавшиеся справочники, LLM получит весь этот хаос как источник "фактов". Дальше вступает старое правило: занесли мусор - получили мусор. Модель скомпилирует из плохих данных гладкий, уверенный и логично звучащий ответ, который при этом будет фактически неверным. И именно уверенный тон здесь особенно опасен: он легко усыпляет критичность сотрудников.



Поэтому MDM и есть основа того, что сегодня называют данными, готовыми для ИИ

Он гарантирует, что векторные индексы RAG и базы, к которым обращается RIG, строятся не на сырых записях, а на верифицированных «золотых записях».

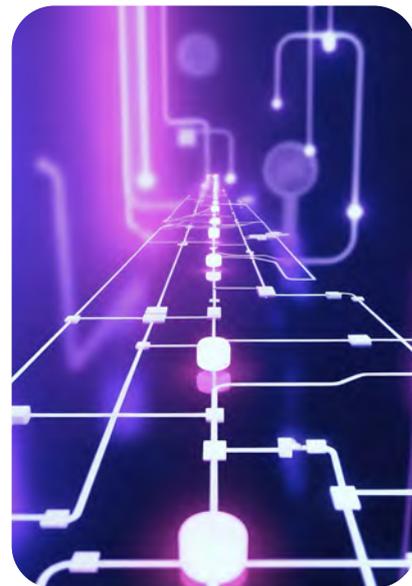
ПОЧЕМУ ДЛЯ LLM ВАЖНО ИМЕТЬ ВЗГЛЯД НА «СУЩНОСТЬ ЦЕЛИКОМ»

Корпоративный чат-бот или ИИ-помощник редко выигрывает от разрозненных фрагментов данных. Чтобы осмысленно говорить с менеджером по продажам, агентом поддержки или финансовым аналитиком, модели нужен полноценный контекст по объекту.

Современные MDM-платформы поэтому движутся к сущностно-ориентированной модели, где собираются профили вроде «Покупатель 360» или «Продукт 360».

В «Покупатель 360» обычно входят три слоя:

- мастер-данные: ФИО, дата рождения, ИНН, юридический адрес, другие стабильные идентификаторы;
- транзакционные данные: покупки, переводы, счета, открытые заявки;
- данные о взаимодействиях: звонки в колл-центр, чат-логи, email-история, маркетинговые касания.



Когда RAG работает вместе с MDM, он не просто ищет по фамилии клиента в хаотичном корпоративном поиске. Он получает целостный, уже склеенный профиль. На его основе автоматически собирается обогащённый промпт. И тогда LLM может не только ответить на вопрос, но и сделать уместную персонализированную рекомендацию, предложить перекрёстную продажу или апселл на основе полной истории и при этом не советовать продукт, который клиент уже купил или на который раньше жаловался.

Тут есть ещё один важный, но часто недооценённый эффект: экономика. Если индексировать сырые данные, векторная база раздувается, дублей становится слишком много, в контекст улетают противоречивые чанки, а токены тратятся впустую. MDM убирает дубли до того, как они попадут в контур ИИ, и делает архитектуру заметно дешевле.

БЕЗОПАСНОСТЬ И КОМПЛАЕНС: ЕЩЕ ОДНА ПРИЧИНА, ПОЧЕМУ БЕЗ MDM БУДЕТ СЛОЖНО

Как только компания подключает LLM к внутренним данным, вопрос безопасности перестаёт быть второстепенным. Если в векторную базу в открытом виде попадают персональные данные (ПД) клиентов, риск утечки становится реальным. Модель может непреднамеренно выдать номер карты, диагноз, номер соцстрахования или другие чувствительные сведения человеку без нужного уровня доступа, а в худшем случае — внешнему пользователю.

Передовые MDM-системы решают это не отдельным «забором» вокруг ИИ, а встраивают управление доступом прямо в поток данных. Здесь особенно важны два механизма:

1 Ролевая модель управления доступом

Определяет, кто и что вообще имеет право видеть.

2 Динамическое маскирование данных

Работает «на лету». Исходные данные в MDM-хранилище не меняются, но в момент, когда RAG-запрос забирает профиль клиента для передачи модели, промежуточный слой оценивает права инициатора и маскирует чувствительные поля.

Пример

Если сотрудник поддержки запрашивает у ИИ помощь по кейсу клиента, система может подставить маску вместо номера кредитной карты, например <**** * 4321>, а точный возраст заменить возрастным диапазоном. И важно, что зрелые платформы умеют маскировать ПД не только в структурированных таблицах, но и внутри неструктурированного текста: скрывать фамилии в чат-логах, почте и других источниках до этапа векторизации.

Почему в повестке всё чаще возникает тема микро-баз данных и дата-продуктов

Ещё один важный сдвиг в интеграции MDM и ИИ — переход от монолитных озёр данных к архитектуре дата-продуктов и микро-баз данных.

Классическое централизованное хранилище не всегда подходит GenAI-системам, особенно если им нужен быстрый, контекстный и безопасный доступ к данным в реальном времени. В продуктовом подходе вся унифицированная информация по одной бизнес-сущности, например по конкретному клиенту, упаковывается в собственную зашифрованную микро-базу.

Почему это хорошо работает для ИИ:

- контекст каждой сущности изолирован, а значит проще выстроить защитные ограничения;
- данные не нужно каждый раз тянуть из медленных существующих корпоративных систем, RAG может быстро обратиться к уже подготовленной микро-БД;
- синхронизация идет постоянно, поэтому модель получает не вчерашний снимок, а актуальное состояние.



Для агентных сценариев это особенно ценно. Агенту нужен не просто документ, а достоверный рабочий слепок объекта, с которым он собирается что-то делать.



ЧТО ДЕЛАЮТ КРУПНЫЕ ВЕНДОРЫ MDM УЖЕ СЕЙЧАС

Теория давно перешла в практику. Крупные платформы MDM уже встроили ИИ в свои продукты, причем в обе стороны: они и дают доверенный контекст для LLM, и сами используют GenAI и ML внутри процессов управления данными.

Informatica CLAIRE GPT



Informatica развивает свою платформу Intelligent Data Management Cloud, IDMC, и встроенный ИИ CLAIRE. Следующий шаг — CLAIRE GPT, который переносит генеративный ИИ прямо в задачи data management и MDM.

У Informatica важный акцент сделан на безопасности обучения. Модели CLAIRE обучаются на публичной документации и анонимизированных метаданных, например на структурах таблиц и системных логах. Реальные клиентские бизнес-данные для кросс-обучения без согласия не используются. Платформа не сканирует содержимое таких систем, как Salesforce или SAP, чтобы обучать LLM. Для служб ИБ это принципиальный момент.

На практике CLAIRE GPT работает как Co-pilot для взаимодействия с данными на естественном языке. Пользователь без SQL или Python может попросить найти доверенные данные, собрать пайплайн, исследовать метаданные или настроить правила качества данных. В MDM это особенно заметно в сценарии AI-сопоставлении и объединении: машинное обучение помогает подбирать и рекомендовать правила слияния дубликатов. По данным самой компании, это ускоряет принятие решений на 70% и экономит десятки тысяч часов ручной работы.

SAP Master Data Governance



Для компаний из экосистемы SAP MDM давно перестал быть «желательной опцией». Особенно это видно в проектах миграции на SAP S/4HANA. SAP Master Data Governance, SAP MDG, централизует создание и поддержку золотых записей по клиентам, поставщикам и материалам в единой модели business partner.

ИИ здесь используется в нескольких направлениях. Во-первых, это выявление правил: алгоритмы анализируют исторические данные и находят закономерности, на основе которых можно автоматически формировать правила валидации и предотвращать ошибки на входе. Во-вторых, технология оптического распознавания символов и обработка естественного языка (нужна уже после распознавания символов: когда текст распознан, она помогает понять, что именно в этом тексте важно) помогают вытаскивать мастер-данные из неструктурированных PDF-счетов, договоров и других документов, а затем переносить их в структурированный контур SAP.

Интересный практический сценарий — интеллектуальное резюмирование. Пользователь на экране центрального поиска может нажать кнопку и получить текстовую сводку по объекту мастер-данных. LLM собирает информацию из связанных транзакционных таблиц и метрик качества и формирует аналитическое описание активности поставщика или клиента. Это заметно экономит время аналитика.

Дальше включается уже прикладная ценность. Выверенные мастер-данные из таблиц MARA, LFA1 и KNA1 становятся основой для предиктивной аналитики. Если к этому добавить IoT-датчики на производстве, ИИ может прогнозировать поломки оборудования и точнее управлять цепочками поставок, снижая риск дефицита.

Pimcore



Pimcore занимает особое место среди современных платформ управления данными. Это не прямой аналог «тяжёлых» классических MDM-систем, а более гибкая платформа, объединяющая управление данными и цифровым опытом. Она круто работает там, где нужно в одном контуре связать мастер-данные, данные о товарах, цифровые активы и автоматизацию на базе ИИ.

ИИ в Pimcore встроен не как внешняя надстройка, а как часть внутренних процессов. Платформа поддерживает генерацию текстов через OpenAI-совместимые интерфейсы, интеграции с моделями Hugging Face, автоматический перевод, генерацию изображений, создание вариантов объектов данных и даже дообучение моделей для задач классификации на данных, хранящихся в самой системе. За счёт этого Pimcore полезен далеко не только для хранения и выверки данных, но и для их активного использования внутри повседневных процессов управления контентом и объектами.



PIM + MDM + DAM

Odyssey PIM — поддерживаемый контур развития **Pimcore-решений** в России для задач **PIM/MDM/DAM**: модули и обновления, производительность и стабильность, сопровождение обновлений и миграций с сохранением данных и интеграций.

[Подробнее](#)

Для сценариев с GenAI особенно важен слой подготовки и выдачи данных. Здесь у Pimcore есть Datahub с GraphQL, настраиваемыми схемами публикации и ограничением доступа через рабочие области. Дополнительно доступен простой REST-интерфейс только для чтения: данные в нём индексируются в OpenSearch или Elasticsearch и отдаются оттуда, без лишней нагрузки на основную базу данных. В этом же контуре предусмотрен экспериментальный MCP-сервер, который позволяет агентам и языковым моделям обращаться к данным Pimcore напрямую. В итоге Pimcore хорошо подходит как управляемый контекстный слой для RAG и агентных сценариев: модель получает не хаотичный поток данных из разных систем, а заранее подготовленные, ограниченные по правам и пригодные для машинной обработки сведения.

Reltio

RELTIO

Reltio представляет новое поколение облачных SaaS-платформ MDM. Сегодня компания позиционирует себя как платформу контекстного интеллекта, а в основе этого подхода лежит идея интеллектуального графа данных. Смысл в том, что ценность создаётся не только на уровне отдельных записей, но и на уровне связей между сущностями, атрибутами, взаимодействиями и даже неструктурированными данными. Для бизнеса это особенно важно, потому что реальная логика процессов почти всегда живёт не в одной карточке клиента или поставщика, а в сети отношений между клиентами, товарами, локациями, транзакциями и документами.

Именно поэтому Reltio хорошо подходит для задач, где нужен не плоский справочник, а динамическое и контекстное представление сущности. Платформа делает сильный акцент на сопоставлении и объединении записей об одной и той же сущности. Для этого используются предобученные модели машинного обучения и модели на базе больших языковых моделей, что помогает быстрее и точнее объединять данные из разных источников. При этом Reltio сочетает подходы на основе правил и подходы на основе машинного обучения, что даёт более гибкий механизм формирования эталонных записей.

Ещё одна важная особенность платформы — динамическое формирование наилучшей версии записи в зависимости от задачи. Иными словами, система может по-разному собирать «лучшую версию истины» для разных потребителей данных. Например, представление записи для маркетинга, для продаж и для биллинга может отличаться по приоритетам атрибутов и источников. Это делает Reltio особенно ценным для современных ИИ-сценариев, где важно не просто хранить эталонные данные, а предоставлять наиболее уместный и достоверный контекст под конкретную задачу.

В результате Reltio формирует богатый семантический слой, который подходит как источник данных для RAG, агентных систем и других сценариев с генеративным ИИ. Если классические MDM-платформы исторически были особенно сильны в контроле качества и управлении правилами, то Reltio заметно ближе к архитектурам, где данные должны быстро превращаться в доверенный, управляемый и готовый для использования в ИИ контекст.

K2view



K2view — один из самых заметных игроков в теме дата-продуктов и микро-баз данных для ИИ. Платформа сознательно уходит от идеи одного гигантского хранилища и организует данные вокруг конкретной бизнес-сущности.

Этот подход напрямую связан с AI-готовностью данных. K2view даёт API сверхнизкой задержки для передачи целостных профилей Покупатель 360 в RAG-конвейеры. Кроме того, платформа развивает no-code-инструменты вроде Data Agent Builder. Они позволяют бизнес-аналитикам создавать приложения на базе Agentic AI, которые уже привязаны к верифицированным продуктам данных. То есть агент изначально работает не в абстрактном информационном поле, а в доверенном контуре с подтверждённой фактологией.

Поддержка большого количества коннекторов к существующим корпоративным системам и современным SaaS-платформам помогает собрать действительно полный контекст. За счёт этого риск галлюцинаций для LLM заметно снижается.

КУДА ВСЕ ЭТО ИДЁТ: ОТ ЧАТ-БОТОВ К АГЕНТАМ

Сейчас корпоративный ИИ быстро выходит за рамки «умного помощника, который отвечает на вопросы». Следующий этап — Agentic AI. Это уже не просто генерация текста, а способность планировать цепочку действий, принимать операционные решения и работать с корпоративными API.

Например, агент может не только объяснить кредитную политику, но и сам пересчитать лимиты для группы клиентов на основе макроэкономических индикаторов, истории транзакций и внутренних правил, а затем инициировать нужные действия в системе.

И вот здесь требования к данным становятся ещё жёстче. Для статического LLM-ассистента иногда достаточно загрузить документы и регламенты. Для агентной системы этого мало. Ей нужен динамический, постоянно обновляемый слепок реальности. Не архив, а рабочая карта текущего состояния бизнеса.

Тут приходит на ум такая аналогия, опять с GPS. Навигатор полезен не потому, что внутри него лежит старая карта, а потому, что он подключён к живому потоку данных. В корпоративной архитектуре именно MDM все чаще становится таким слоем «живой карты»: точным, защищённым и пригодным для машинного действия.



ЧТО, СКОРЕЕ ВСЕГО, ПРОИЗОЙДЕТ ДАЛЬШЕ

В ближайшие годы, на наш взгляд, мы увидим не просто интеграцию MDM и GenAI, а их архитектурное сращивание.

Векторные базы данных перестанут восприниматься как отдельный слой «рядом с данными» и будут все глубже встраиваться в MDM-платформы. RAG и RIG станут не внешними экспериментальными обвязками, а частью стандартного контура обработки данных. Цепочка будет выглядеть примерно так: сбор, стандартизация, дедупликация, векторизация, динамический поиск, генерация ответа или действия, а затем обратное обогащение профиля новыми знаниями.

По сути, MDM становится не только системой записи, но и системой контекстного снабжения для ИИ.



ВЫВОД

Пытаться внедрять генеративный ИИ без наведения порядка в мастер-данных — это как строить современный небоскрёб на зыбком основании. Снаружи все выглядит впечатляюще. Но любая серьёзная нагрузка быстро покажет, что фундамент слабый.

RAG и RIG действительно меняют правила игры. RAG даёт компаниям способ подключить LLM к актуальным корпоративным знаниям без постоянного fine-tuning. RIG идёт дальше и позволяет проверять факты по ходу генерации, что особенно важно там, где критична точность чисел и расчётов. Long-context LLM тоже займут свою нишу. Но ни одна из этих архитектур не заменяет слой доверенных данных.

Именно MDM создаёт этот слой. Он собирает golden records, устраняет дубли, наводит порядок в идентичностях, добавляет контекст, встраивает безопасность и комплаенс, а затем делает все это пригодным для RAG, RIG и агентных систем.



Поэтому MDM сегодня — уже не бэк-офисная ИТ-функция. Это стратегическая инфраструктура для GenAI.

И, возможно, именно она в ближайшие годы будет отделять компании, которые действительно научились безопасно капитализировать ИИ, от тех, кто просто автоматизировал цифровые галлюцинации.

Odyssey Consulting Group — ИТ-консалтинговая компания, которая помогает бизнесу выстраивать управление ключевыми данными и превращать их в основу для устойчивых бизнес-процессов, качественного клиентского опыта и масштабируемой цифровой архитектуры. Компания обладает экспертизой в проектах, связанных с MDM, PIM, CDP и DAM: от аудита и проектирования целевой модели до внедрения, интеграции и развития решений в корпоративном ландшафте.

Odyssey Consulting Group помогает компаниям навести порядок в мастер-данных о товарах, клиентах, поставщиках и цифровых активах, обеспечить единые правила работы с данными, повысить их качество и доступность для бизнеса. Это создаёт основу для более эффективного маркетинга, продаж, e-commerce, взаимодействия с партнёрами и внедрения современных AI-сценариев, где критически важны полнота, точность и управляемость данных.

**Оставайтесь на связи с нами,
подписывайтесь на наш канал
в Telegram!**



Анонсы методологических вебинаров, аналитические материалы, экспертные статьи и многое другое.

По всем вопросам пишите на почту sales@odysseyconsgroup.com

Москва, Пресненская наб., Д. 12, Бизнес-центр «Башня Федерация Восток», 63 этаж, офис 10

Санкт-Петербург, Литейный пр., Д. 26А, Бизнес-центр «Преображенский двор», офис 423

Нижний Новгород, ул. Ошарская, д. 95, офис 501

Алматы, мкр. Курылышы, ул. Ырысты, д. 15, +7 (727) 397-90-63

+7 (495) 369-67-69

sales@odysseyconsgroup.com

www.odysseyconsgroup.com

