

Распределенные системы поиска: архитектура, лежащая в основе современных приложений

Поиск товаров на онлайн магазинах, поиск ресторанов в сервисах доставки еды и запросы к корпоративной базе знаний — все это основано на распределенных системах поиска, которые одновременно обрабатывают миллионы похожих запросов. Эти системы прошли долгий путь от простого сопоставления ключевых слов — современные архитектуры сочетают традиционный текстовый поиск с семантическим пониманием, при этом обеспечивая время отклика менее секунды в глобальном масштабе.

За последние годы сфера поиска претерпела кардинальные изменения. Традиционные системы на основе ключевых слов быстро превращаются в платформы, которые понимают не только слова, но и намерения. Бессерверные архитектуры устраняют операционные издержки. Векторные базы данных открывают возможности для совершенно новых парадигм поиска. В этом техническом обзоре подробно рассматриваются архитектуры, алгоритмы и компромиссные решения, лежащие в основе современного поиска в больших масштабах.

Основы поиска: от ключевых слов к пониманию

Современные системы поиска должны обрабатывать широкий спектр типов запросов, которые отражают то, как пользователи на самом деле ищут информацию. Поиск по ключевым словам остается основополагающим — пользователи по-прежнему вводят «чехол для iPhone 15» или «программное обеспечение для управления проектами», — но теперь поисковые системы интерпретируют намерения, обрабатывают опечатки и понимают синонимы. Запрос «врач рядом» требует контекста местоположения, а «лучший ноутбук для программирования» — понимания категорий продуктов и сценариев использования.

Эволюция релевантности



Точное совпадение ключевых слов

Ранние системы полагались на точное совпадение ключевых слов, часто возвращая результаты в случайном порядке.

19

BM25

BM25 (Best Matching 25) улучшил TF-IDF, устранив насыщение частоты терминов — дополнительные вхождения термина приносят уменьшающуюся отдачу — и включив нормализацию длины документа. Этот алгоритм остается основой большинства текстовых поисковых систем, хотя современные реализации включают в себя десятки дополнительных сигналов ранжирования.



TF-IDF

Внедрение TF-IDF (Term Frequency-Inverse Document Frequency, частота термина — обратная частота документа) принесло статистическую релевантность, ранжируя документы на основе важности терминов как в отдельных документах, так и во всей коллекции.



Машинное обучение

Машинное обучение привело к появлению подходов «обучение ранжированию», которые оптимизируют релевантность с помощью данных о переходах по ссылкам и сигналов о поведении пользователей. Эти системы могут научиться тому, что пользователи, нажимающие на третий результат вместо первого, указывают на проблему с упорядочением по релевантности, и постоянно улучшать качество поиска с помощью обратной связи от пользователей.

Обработка и анализ запросов

Прежде чем документы могут быть проиндексированы, как запросы, так и контент проходят сложный анализ. Конвейеры текстового анализа токенизируют входные данные, нормализуют регистр, удаляют стоп-слова и применяют стемминг или лемматизацию. Запрос «кроссовки» может быть преобразован с включением связанных терминов, таких как «спортивная обувь» или «кеды».

Пользовательские анализаторы

Платформы электронной коммерции могут рассматривать «iPhone-15» и «iPhone 15» как одно и то же, в то время как поиск юридических документов сохраняет точную пунктуацию. Для поддержки нескольких языков требуются языковые анализаторы, понимающие различные грамматические структуры и наборы символов.

Классификация намерений

Современная обработка запросов также включает классификацию намерений. Поиск по запросу «Apple» может относиться к фрукту или технологической компании, что требует контекста из истории пользователя, содержимого текущей страницы или явного устранения неоднозначности.

Распознавание сущностей

Распознавание именованных сущностей помогает идентифицировать людей, места и организации в запросах, обеспечивая более точное сопоставление.

Elasticsearch: стандарт поисковых систем

Elasticsearch стал де-факто стандартом для поисковой инфраструктуры, обеспечивая работу всего, от каталогов продуктов электронной коммерции до систем управления знаниями в предприятиях. Его успех основан на сочетании мощных поисковых возможностей и распределенной архитектуры, которая масштабируется по горизонтали.

Основная архитектура и модель данных

Elasticsearch организует данные в индексы — логические контейнеры, похожие на таблицы базы данных. Каждый индекс состоит из документов, хранящихся в виде объектов JSON, что обеспечивает гибкость для различных структур данных. В отличие от жестких схем баз данных, динамическое сопоставление Elasticsearch автоматически определяет типы полей и создает соответствующие индексы.

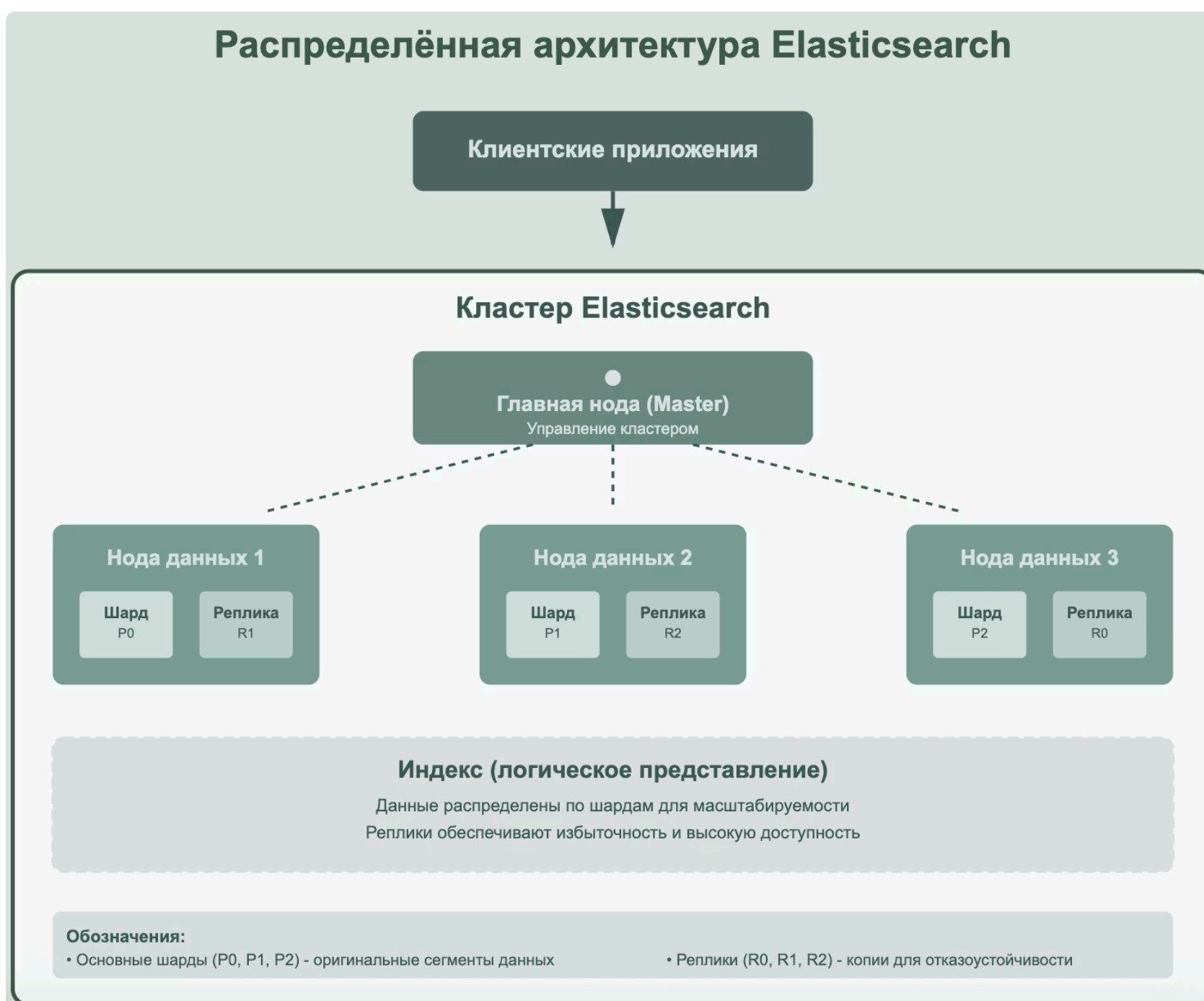
Базовая структура данных использует инвертированные индексы — ту же технологию, что и поисковая система Google. Вместо последовательного хранения документов система создает сопоставления терминов и документов для каждого поля. При индексировании описания продукта, содержащего «беспроводные Bluetooth-наушники», Elasticsearch создает записи, сопоставляющие каждый термин с ID документа, что позволяет быстро находить термины в миллионах документов.

Типы полей

- Текстовые поля — полный анализ для поиска
- Поля ключевых слов — точное сопоставление
- Числовые, даты и булевы поля
- Вложенные и объектные типы
- Специализированные типы (geo_point)

Шардинг и распределение

Модель распределения Elasticsearch основана на шардинге — разделении индексов на более мелкие независимые единицы, распределенные по узлам кластера.



Компоненты архитектуры:

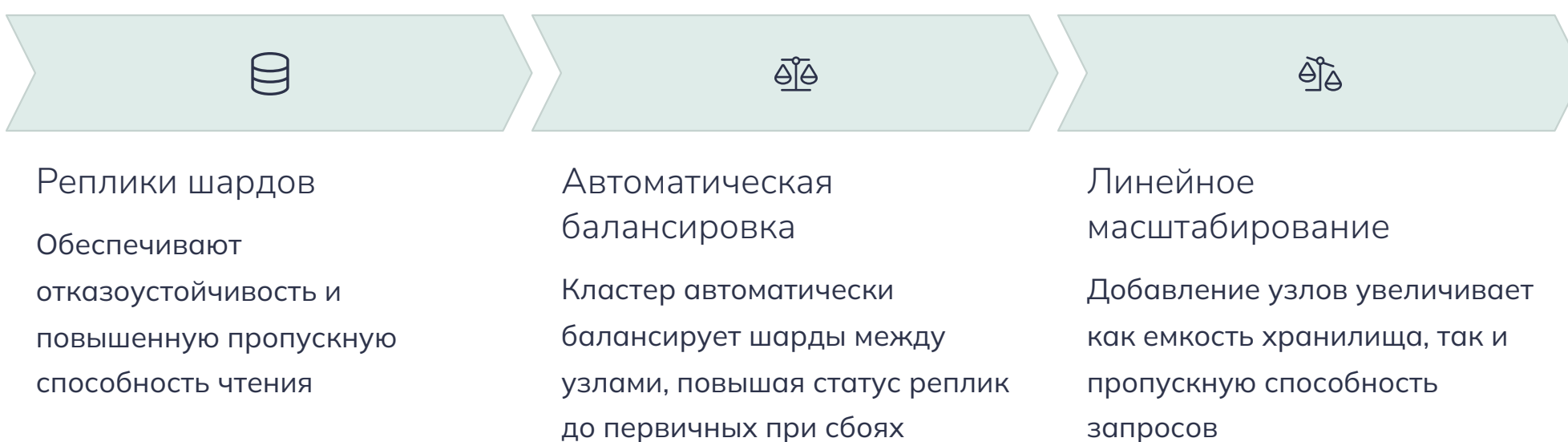
Главная нода (Master Node) управляет метаданными кластера, координирует создание и удаление индексов, отслеживает состояние узлов и распределение шардов.

Ноды данных (Data Nodes) хранят шарды и выполняют операции с данными. В примере архитектуры показаны три ноды данных, каждая из которых содержит комбинацию первичных шардов и реплик.

Первичные шарды (P0, P1, P2) — каждый функционирует как полный и самостоятельный индекс, независимо обрабатывая как чтение, так и запись. Система использует маршрутизацию на основе хеширования для обеспечения равномерного распределения данных между шардами.

Реплики (R0, R1, R2) — копии первичных шардов, размещенные на других нодах для обеспечения отказоустойчивости и высокой доступности. Реплики также могут обслуживать операции чтения, увеличивая пропускную способность кластера.

Принцип распределения: каждая нода содержит один первичный шард и одну реплику от другого шарда, что обеспечивает сохранность данных при отказе любой ноды.



Выполнение запросов происходит по схеме «scatter-gather». Координационный узел принимает запросы, определяет соответствующие шарды, распределяет запросы по кластеру и агрегирует результаты. Такой распределенный подход позволяет выполнять запросы к петабайтам данных с временем отклика менее секунды.

Расширенные возможности поиска

Типы запросов


Elasticsearch предоставляет сложные типы запросов, выходящие за рамки простого сопоставления текста. Булевы запросы объединяют несколько условий с помощью операторов AND, OR и NOT. Диапазонные запросы фильтруют числовые и даты поля. Нечеткие запросы обрабатывают опечатки и вариации, а запросы с подстановочными знаками и регулярными выражениями позволяют выполнять сопоставление шаблонов.

Агрегации

Агрегации позволяют выполнять аналитические запросы по результатам поиска. Агрегации терминов подсчитывают количество вхождений значений полей, а агрегации метрик вычисляют суммы, средние значения и процентиля. Агрегации корзин группируют результаты по диапазонам дат, географическим регионам или настраиваемым критериям. Эти возможности превращают Elasticsearch из поисковой системы в полноценную аналитическую платформу.

Полнотекстовый поиск

Функции полнотекстового поиска включают сопоставление фраз, запросы по близости и выделение. Запрос «more-like-this» находит похожие документы на основе анализа содержания. Многополевой поиск может выполнять запросы по нескольким полям с разными коэффициентами повышения, а межполевой поиск обрабатывает несколько полей как одно логическое поле.

 **Интеграция с LLM:** Современный Elasticsearch интегрирует возможности на базе LLM, которые расширяют возможности поиска за пределы сопоставления ключевых слов. Плотные векторные поля поддерживают семантический поиск с использованием вложений, сгенерированных такими моделями, как BERT или text-embedding-ada-002 от OpenAI. Пользователи могут искать «бюджетный транспорт» и находить документы о «доступных автомобилях» даже без общих ключевых слов.



Решения для поиска в базах данных

Многие организации предпочитают расширять существующие базы данных с помощью функций поиска, а не поддерживать отдельную инфраструктуру поиска. Такой подход снижает сложность эксплуатации и позволяет использовать существующий опыт и модели безопасности.

Эволюция поиска в PostgreSQL

Основные возможности

PostgreSQL предлагает сложные возможности полнотекстового поиска с помощью типов данных `tsvector` и `tsquery`. Индексы GIN (Generalized Inverted Index) обеспечивают эффективный текстовый поиск с ранжированием и сопоставлением фраз. Система поддерживает несколько языков, пользовательские словари и алгоритмы стемминга.

Последние версии PostgreSQL включают параллельное выполнение запросов для текстового поиска, что значительно повышает производительность на многоядерных системах. Расширения, такие как `pg_trgm`, предоставляют возможности нечеткого сопоставления, а `pg_search` (ParadeDB) обеспечивает производительность, сопоставимую с Elasticsearch, с помощью библиотеки Tantivy.

Векторный поиск

`pgvector` позволяет осуществлять семантический поиск по векторному сходству. Поддерживая алгоритмы HNSW и IVFFlat, он обрабатывает вложения до 16 000 измерений с настраиваемыми функциями расстояния. Эта интеграция позволяет объединить традиционный текстовый поиск с семантическим сходством в одном запросе.

Преимущества: Для организаций, уже использующих PostgreSQL, встроенные возможности поиска предлагают убедительные преимущества в простоте эксплуатации, согласованности данных и снижении затрат. Однако масштабируемость по-прежнему ограничена одноузловыми развертываниями или сложными стратегиями шардинга.

Интегрированный подход MongoDB

MongoDB Atlas Search интегрирует Apache Lucene непосредственно в платформу базы данных. Система создает индексы поиска, синхронизированные с операционными данными через потоки изменений, обеспечивая согласованность и позволяя независимо масштабировать операции поиска.



Векторный поиск

Atlas Search поддерживает как текстовый, так и векторный поиск с реализацией HNSW для векторов до 8192 измерений. Этап агрегации `$vectorSearch` позволяет выполнять сложные запросы, сочетающие традиционные фильтры с семантическим поиском.



Динамическое сопоставление

Динамическое сопоставление автоматически определяет типы полей и создает соответствующие индексы.



Единый API

Интеграция предоставляет уникальные преимущества для приложений, уже использующих MongoDB. Документы можно искать и извлекать через один и тот же API, что устраняет сложность обслуживания отдельной инфраструктуры поиска. Встроенные функции безопасности и репликации используют зрелые операционные возможности MongoDB.



Коммерческие платформы поиска

Сложность построения и обслуживания инфраструктуры поиска породила множество коммерческих платформ, предлагающих управляемые возможности поиска с расширенными функциями и глобальным распространением.

Фокус Algolia на производительности

6.7ms

Среднее время отклика

Распределенная поисковая сеть Algolia достигает исключительной производительности благодаря глобально распределенной инфраструктуре

25+

Дата-центров

Используя репликацию на основе консенсуса в более чем 25 дата-центрах

90%

Запросов менее 15мс

90% запросов выполняются менее чем за 15 мс

2B+

Запросов в месяц

Эта инфраструктура поддерживает более 2 миллиардов запросов в месяц

Механизм консенсуса трех серверов обеспечивает согласованность данных при сохранении высокой доступности. Репликация в реальном времени синхронизирует данные по сети в течение нескольких минут, автоматически направляя пользователей на ближайший сервер.

Algolia отличается удобством для разработчиков благодаря комплексным SDK и компонентам пользовательского интерфейса мгновенного поиска. Платформа автоматически обрабатывает отклонение запросов, кэширование результатов и постепенное уточнение запросов. Расширенные функции включают толерантность к опечаткам, фасетный поиск и персонализацию на основе поведения пользователей.

Предложения по управляемому поиску

Крупные поставщики облачных услуг и поисковых систем предлагают управляемые поисковые сервисы, которые сочетают в себе функциональность и простоту эксплуатации, устраняя операционные затраты на управление кластерами и предоставляя возможности корпоративного уровня.

Elastic Cloud

Elastic Cloud предоставляет официальный управляемый сервис Elasticsearch, предлагающий полный набор функций открытого ПО Elasticsearch с корпоративной безопасностью, мониторингом и поддержкой. Платформа работает на AWS, Google Cloud и Microsoft Azure, что позволяет развертывать ее в предпочтительных облачных средах, сохраняя при этом функциональность.

Сервис предлагает несколько вариантов развертывания, включая выделенные кластеры, бессерверный поиск и решения для наблюдаемости. Бессерверный Elasticsearch автоматически масштабирует вычислительные ресурсы и хранилище независимо друг от друга, устраняя необходимость планирования емкости и сохраняя предсказуемую цену за запрос. Этот подход особенно выгоден для приложений с переменной или непредсказуемой нагрузкой на поиск.

Сильная сторона платформы заключается в предоставлении полного набора функций Elasticsearch без сложности эксплуатации. Команды получают доступ к новейшим функциям сразу после их выпуска, а автоматические обновления и патчи безопасности устраняют затраты на обслуживание.

Ключевые преимущества

- Полный набор функций Elasticsearch
- Корпоративная безопасность
- Автоматические обновления
- Патчи безопасности
- Мультиоблачное развертывание
- Бессерверные опции

Поиск от облачных провайдеров

Amazon OpenSearch Service

Amazon OpenSearch Service предоставляет полностью управляемые кластеры Elasticsearch с автоматическим масштабированием, резервным копированием и функциями безопасности. Сервис интегрирован с CloudWatch для мониторинга, IAM для безопасности и VPC для сетевой изоляции. Недавние усовершенствования включают варианты развертывания без серверов и возможности векторного поиска. Плагин нейронного поиска позволяет выполнять семантический поиск с использованием предварительно обученных моделей или настраиваемых встраиваний. Обнаружение аномалий выявляет необычные паттерны в поведении поиска, а тонкий контроль доступа позволяет развертывать многопользовательские системы с безопасностью на уровне документов.

Microsoft Azure

Azure AI Search демонстрирует облачный подход с встроенными когнитивными навыками для анализа изображений, OCR и обработки текста. Архитектура навыков платформы позволяет создавать сложные конвейеры обработки контента, автоматически извлекая структуру из неструктурированного контента.

Интеграция с Azure OpenAI обеспечивает возможность запросов на естественном языке и функции резюмирования. Сервис отличается превосходным пониманием документов с автоматическим распознаванием сущностей и поддержкой более 60 языков.

Google Cloud Platform

Google Vertex AI Search использует опыт веб-поиска для корпоративных приложений. Платформа сочетает в себе поиск и диалоговый искусственный интеллект, обрабатывая как структурированные, так и неструктурированные данные, а также предоставляя базовые возможности через информационный корпус Google.

Сервис отличается высокой эффективностью понимания документов и обработки естественного языка, а также автоматической оценкой качества и настройкой релевантности на основе взаимодействия с пользователем. Интеграция с сервисами машинного обучения Google обеспечивает сложные функции персонализации и рекомендаций.

Заключение: архитектура, формирующая поисковые возможности будущего

Распределённые системы поиска прошли путь от простого сопоставления текста к высокоинтеллектуальным платформам, сочетающим текстовый, структурированный и семантический поиск в реальном времени. Сегодняшние решения опираются на зрелые механизмы индексирования и шардинга, глубоко интегрированы с моделями машинного обучения и поддерживают гибридные запросы — от BM25 до векторных встраиваний.

Ключевые тенденции:

- **Семантика важнее слов:** переход от точного совпадения к пониманию смысла запроса.
- **Интеграция с LLM и векторными БД:** позволяет находить релевантные результаты даже при отсутствии общих терминов.
- **Облачные и бессерверные решения:** снижают издержки и упрощают масштабирование.
- **Сближение поиска и транзакционных БД:** PostgreSQL и MongoDB развиваются в сторону полноценных поисковых платформ.
- **Управляемые сервисы:** становятся стандартом для глобальных приложений с миллионами пользователей.

Современный поиск — это не просто технология, это **архитектурный слой**, влияющий на пользовательский опыт, качество рекомендаций и бизнес-показатели. Архитекторы распределённых систем должны проектировать поиск как **стратегически важный компонент**, обеспечивая баланс между масштабируемостью, релевантностью и оперативностью.