

Агентно-ориентированная модель обработки логов, основанная на синтезе методов пром프트-инжиниринга и цепочки рассуждений

Пашигорев К.И., приглашенный преподаватель,
Базовая кафедра Сбербанка

Гипотеза 1: Генерацию новости и сокращение количества атрибутов логов при сохранении качества новости возможно осуществить итеративным применением LLM

Гипотеза 2: при итеративном применении LLM режим работы с рассуждениями позволяет повысить эффективность этапа сокращения логов при сохранении качества новости

Ограничения

1. Отсутствует доступ к изменению архитектуры LLM
2. Отсутствует возможность дообучения модели
3. Исключить применение RAG

Разработка и применение мультиагентных систем в корпоративной среде

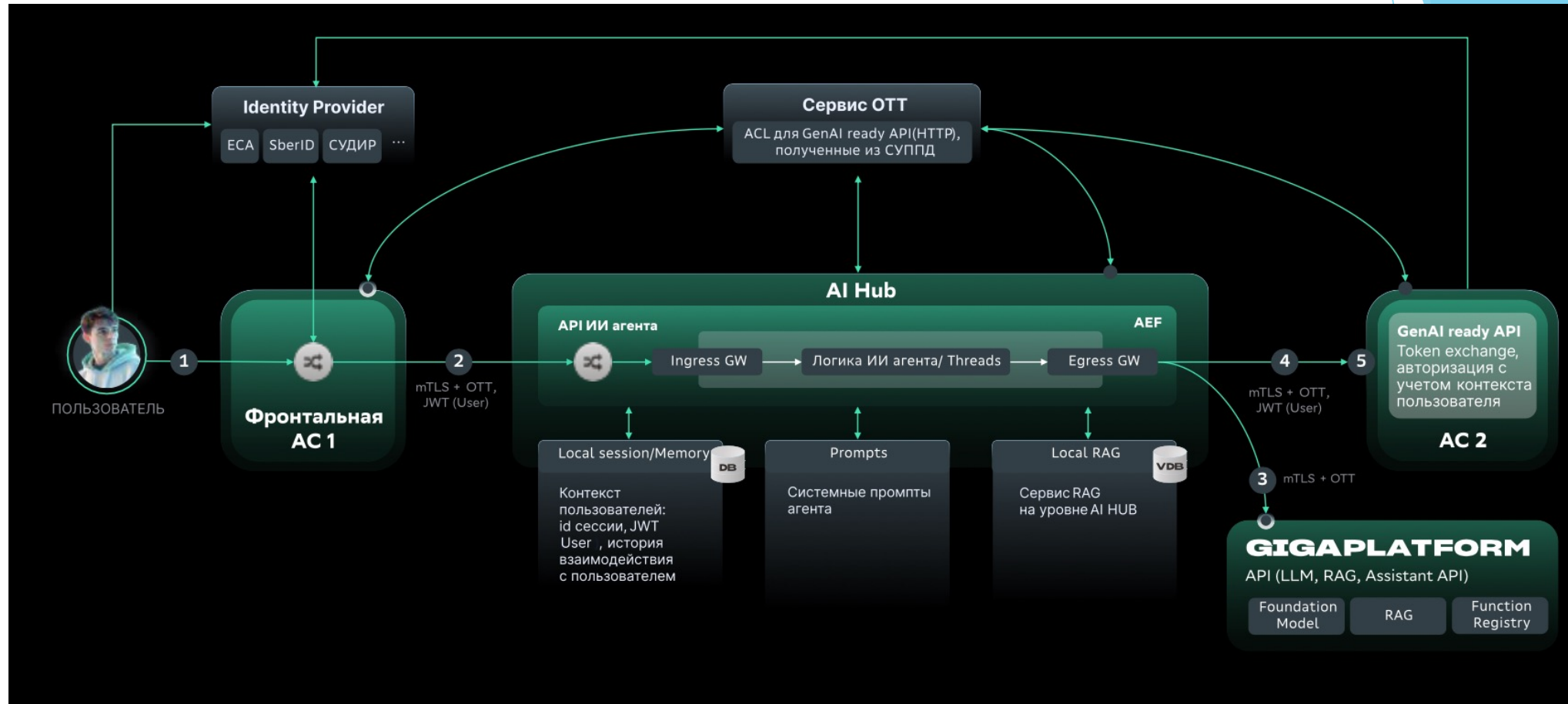


Рисунок 1. Изолированная среда исполнения кода

Пример датасета

	Название дата-продукта	Дата-продукт	Категория в СМД	Тип дата-продукта	Кластер	Схема БД	Менеджеры данных ДП	uuid
0	Витрина Суммаризация	СВД ДРПА. DS. Витрина Суммаризация	Пользовательские данные	Специализированная витрина	A	CUSTOM_DRPA_DS_DIALOGUE_SUM	Гус	0A0A2
1	-	-	-	-	-	-	-	0A0CD
2	Очная аутентификация в БХБ160 ММБ	Очная аутентификация в БХБ160 ММБ	Реплики ППРБ	Реплика	S	PLATFORM_F2FAUTH	Вак	0A0CD
3	Очная аутентификация в БХБ160 ММБ	Очная аутентификация в БХБ160 ММБ	Реплики ППРБ	Реплика	S	PLATFORM_F2FAUTH_HIST	Вак	0A0CD
4	БВД HR Структура функций	БВД HR Структура функциональных	Non-Hadoop Пользовательские	Специализированная витрина	C	S_GRNPLM_VD_HR_EDP_VD	Кон	0A1A5
5	[CVRS2.IS1101] Конверсия	СВД ФБ [CVRS2.IS1101] Конверсии	Non-Hadoop Пользовательские	Специализированная витрина	C	S_GRNPLM_AS_FIN_MIS3_MISPLT_IS	Маг	0A1CB
6	СВД БТ Сверточная гекса	СВД БТ Сверточная гексагональная	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	-	0A1D1
7	-	-	-	-	-	-	-	0A1D9
8	-	-	-	-	-	-	-	0A1D4
9	СВД БТ Сверточная расши	СВД БТ Сверточная расширенная ге	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	-	0A1D8
10	SberInfra Портал ДИ Harmonizer	DI Harmonizer	Реплики self-service	Реплика	S	SELFSERVICE_HARMONIZER	Аск	0A1D2
11	СВД БТ Сверточная гекса	СВД БТ Сверточная гексагональная	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	-	0A1D3
12	-	-	-	-	-	-	-	0A1DA
13	-	-	-	-	-	-	-	0A1DA
14	СВД БТ Сверточная расши	СВД БТ Сверточная расширенная ге	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	-	0A1DB
15	-	-	-	-	-	-	-	0A1DC
16	СВД БТ Гексагональная с	СВД БТ Гексагональная сверточная	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	Куз	0A1DC
17	-	-	-	-	-	-	-	0A1DC
18	СВД БТ Гексагональная с	СВД БТ Гексагональная сверточная	Пользовательские данные	Специализированная витрина	A	CUSTOM_T_GDM_CONVOLUTIONAL	Куз	0A1DE
19	-	-	-	-	-	-	-	0A1DE
20	Ежедневные данные КЮ	СВД ФБ [LCC2.IS0958] Кредиты ЮЛ	Non-Hadoop Пользовательские	Специализированная витрина	C	S_GRNPLM_AS_FIN_MIS3_MISPLT_IS	Маг	0A2D8

Архитектура системы

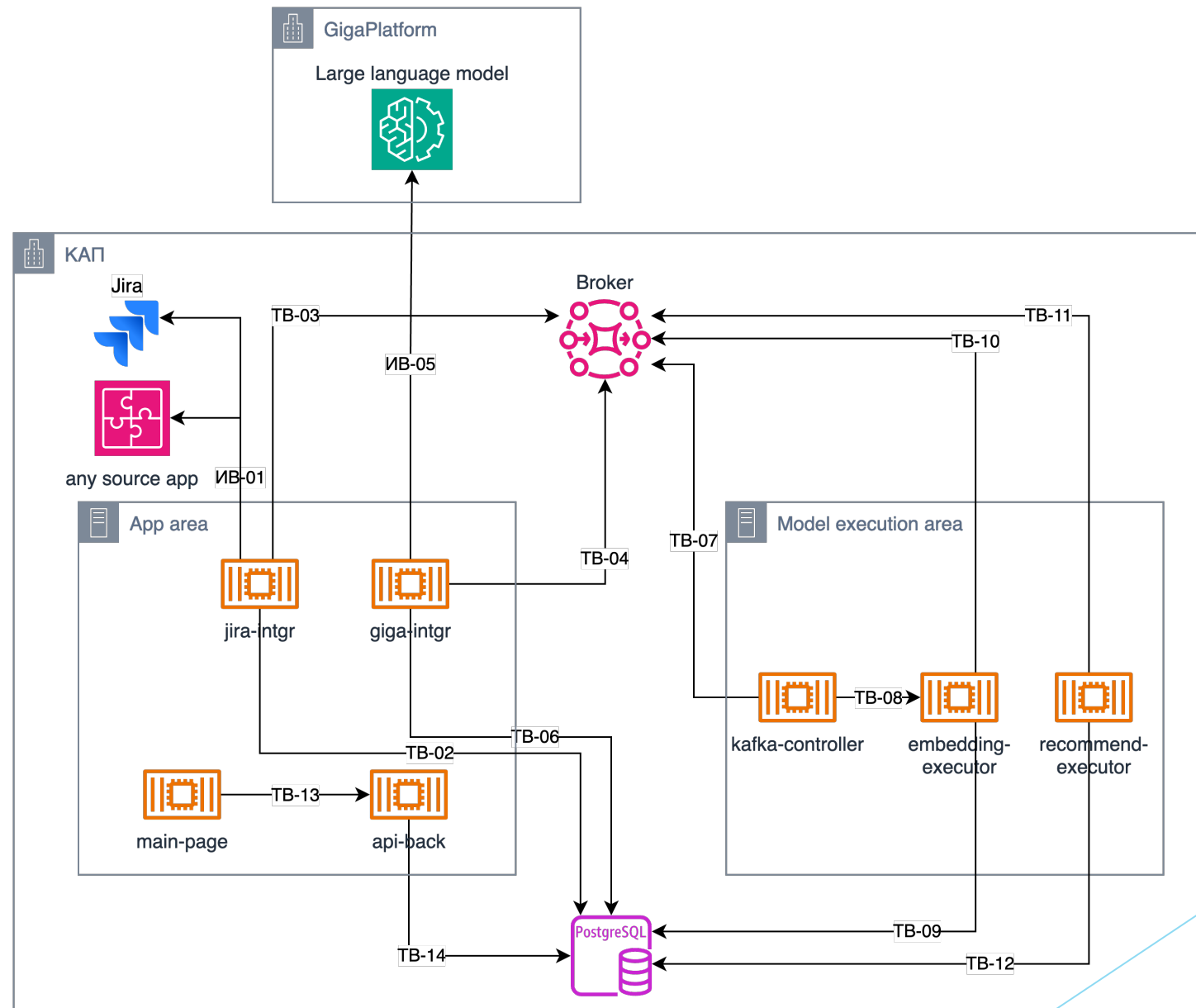
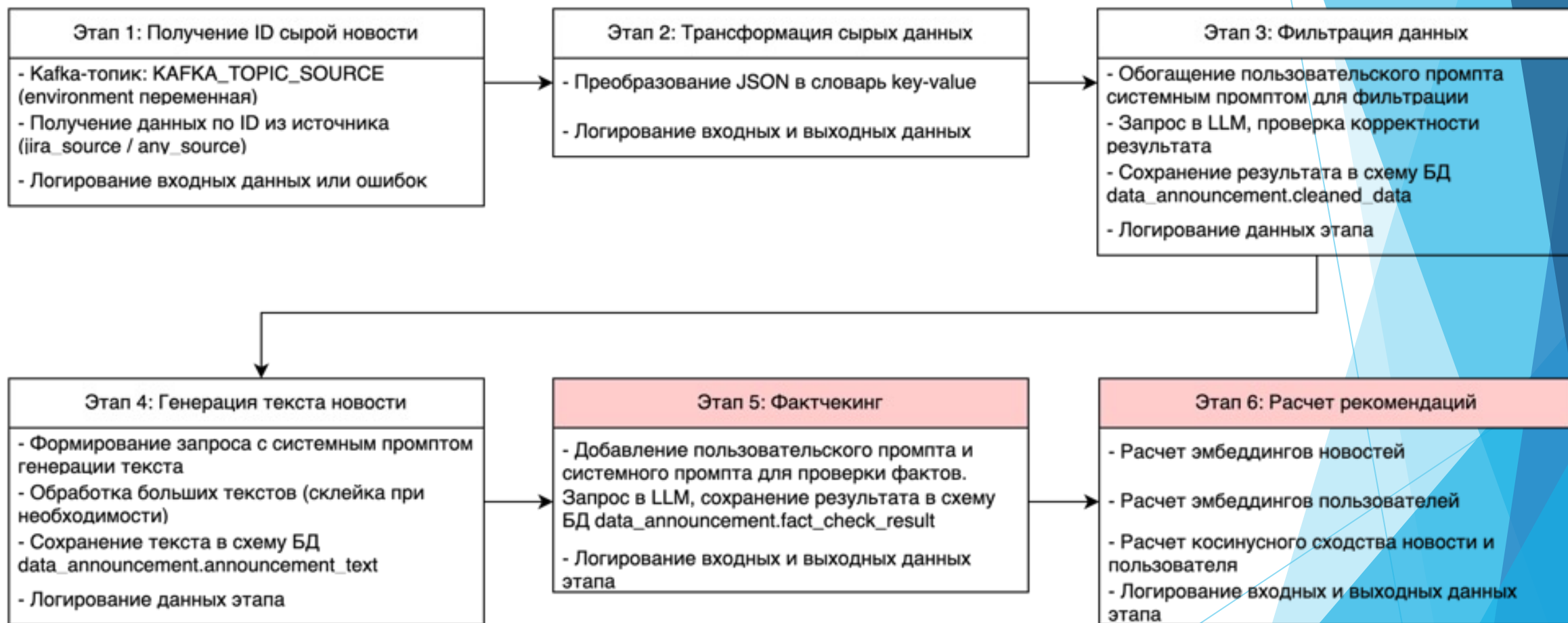


Рисунок 2. Компонентная диаграмма системы



Промпт	Назначение
Твоя задача определить полезные ключи из переданного списка. Надо избавиться от идентификационных номеров, flags и tags. Полезными считаем ключи, связанные с продуктом, людьми или ключевыми датами. Подробно по порядку рассуждай, подходит ли ключ в формате: \"Ключ \"название\" является \"продуктом, людьми, датой*\", решено \"принимаем или отвергаем*\", \"почему*\". В самом конце выведи все принятые таким образом ключи в формате \"[\"ключ1\", ...)\"	Оптимизация ключей
На основе переданного списка словарей определи ключи, не являющиеся флагами и по которым значения не содержат нечитабельных для человека хэшей, тегов и id. Максимально подробно описывай этапы проверки каждого ключа в формате: \"рассматриваемый ключ, его значение, подходит ли под критерии или нет, почему\". Выбери максимум {output size} репрезентативных ключей, не больше! В конца напечатай весь подходящий список в формате \"(\"keyl\", ...)\".	Оптимизация ключей
Сгенерируйте краткую и информативную новостную статью, резюмирующую ключевые детали события, описанного в предоставленных технических данных. Сосредоточьтесь на представлении фактической информации в нейтральном тоне, избегая любых спекуляций, предположений или эмоционального языка. Используйте четкую и логичную структуру и избегайте сенсационного или привлекающего внимание языка. Целью является предоставление объективного резюме события, чтобы читатели могли быстро понять, что произошло.	Генерация текста

Количество ключей

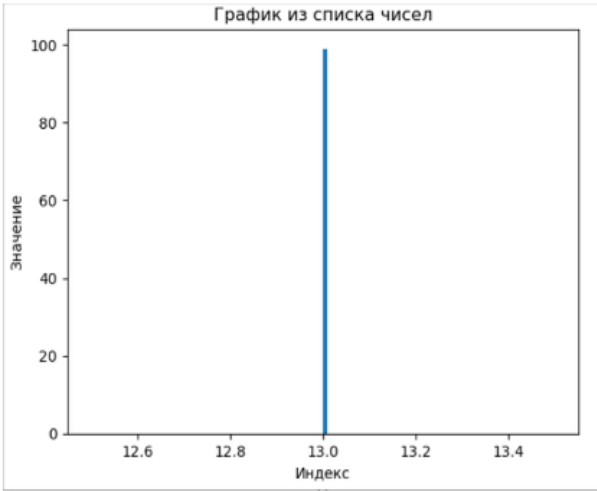


Рисунок 3. Количество ключей для промпта №1 без рассуждения

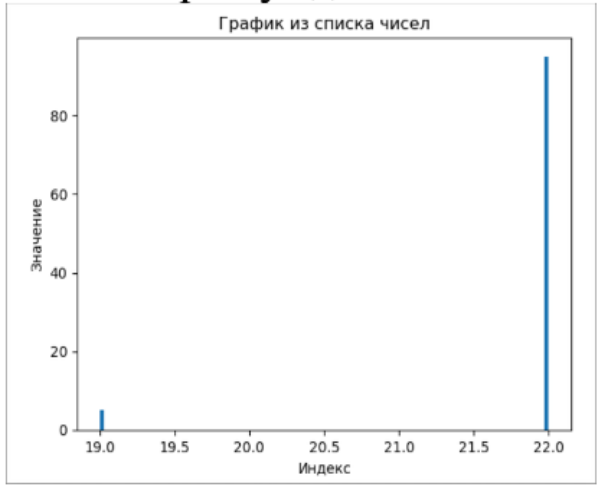


Рисунок 5. Количество ключей для промпта №2 без рассуждения

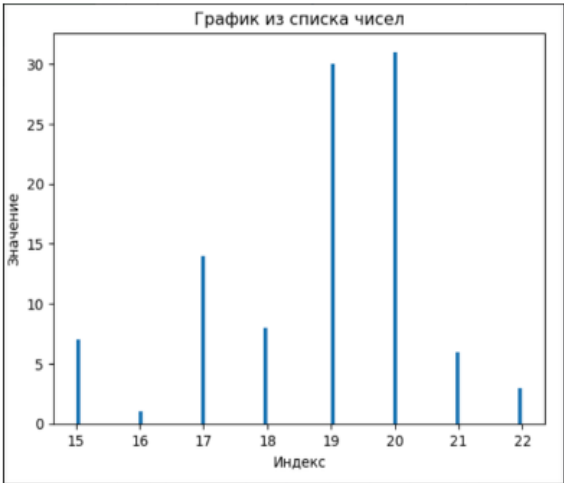


Рисунок 4. Количество ключей для промпта №1 с рассуждением

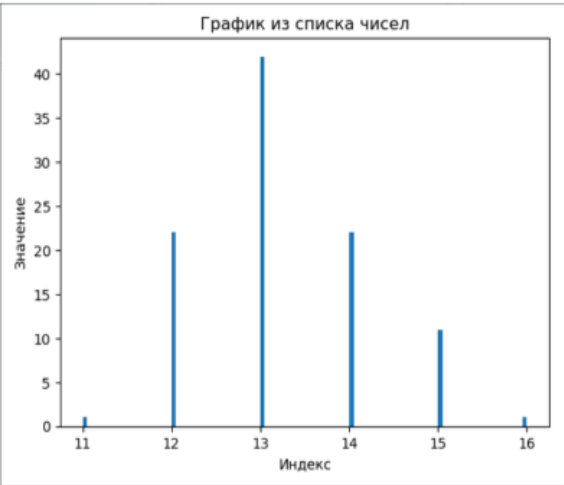


Рисунок 6. Количество ключей для промпта №2 с рассуждением

Вхождение в эталонный массив

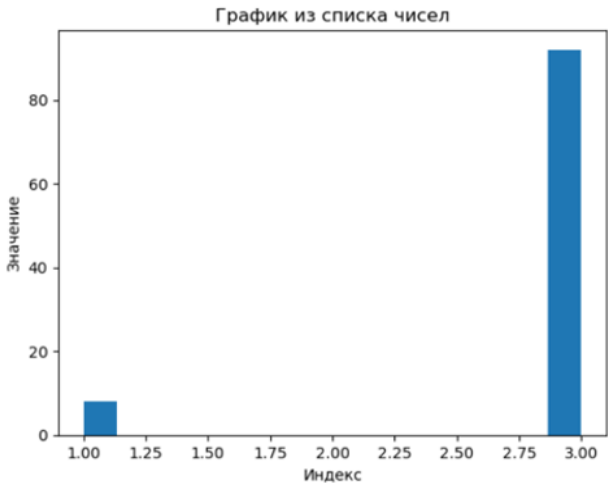


Рисунок 7. Вхождение в эталонный массив с промптом №1 без рассуждений

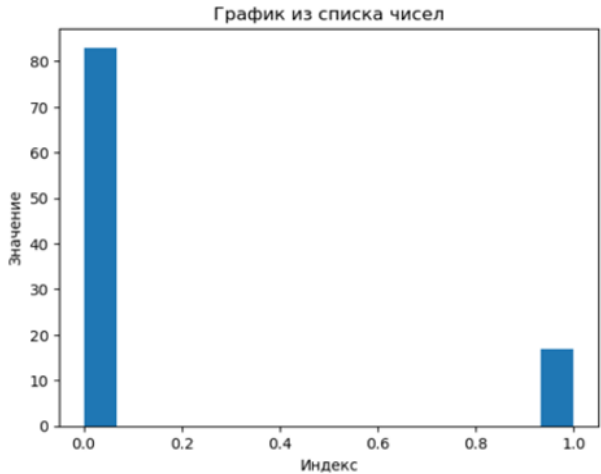


Рисунок 8. Вхождение в эталонный массив с промптом №1 с рассуждениями

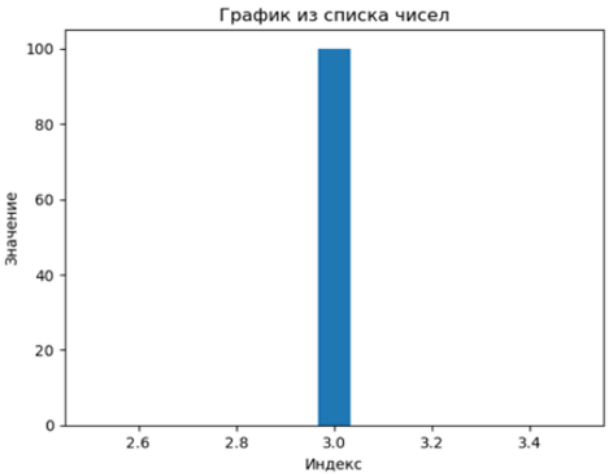


Рисунок 9. Вхождение в эталонный массив с промптом №2 без рассуждений

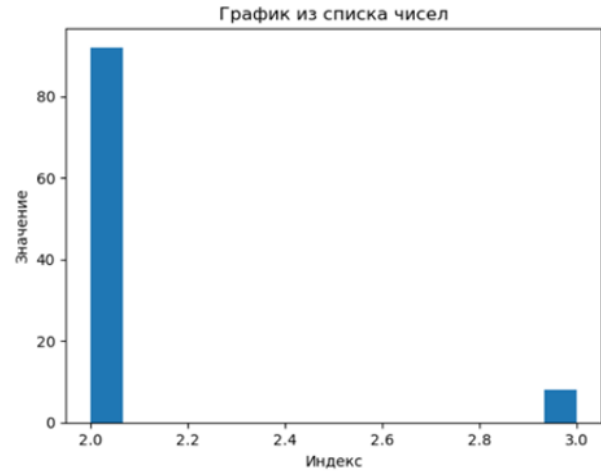


Рисунок 10. Вхождение в эталонный массив с промптом №2 с рассуждениями

t-SNE визуализация

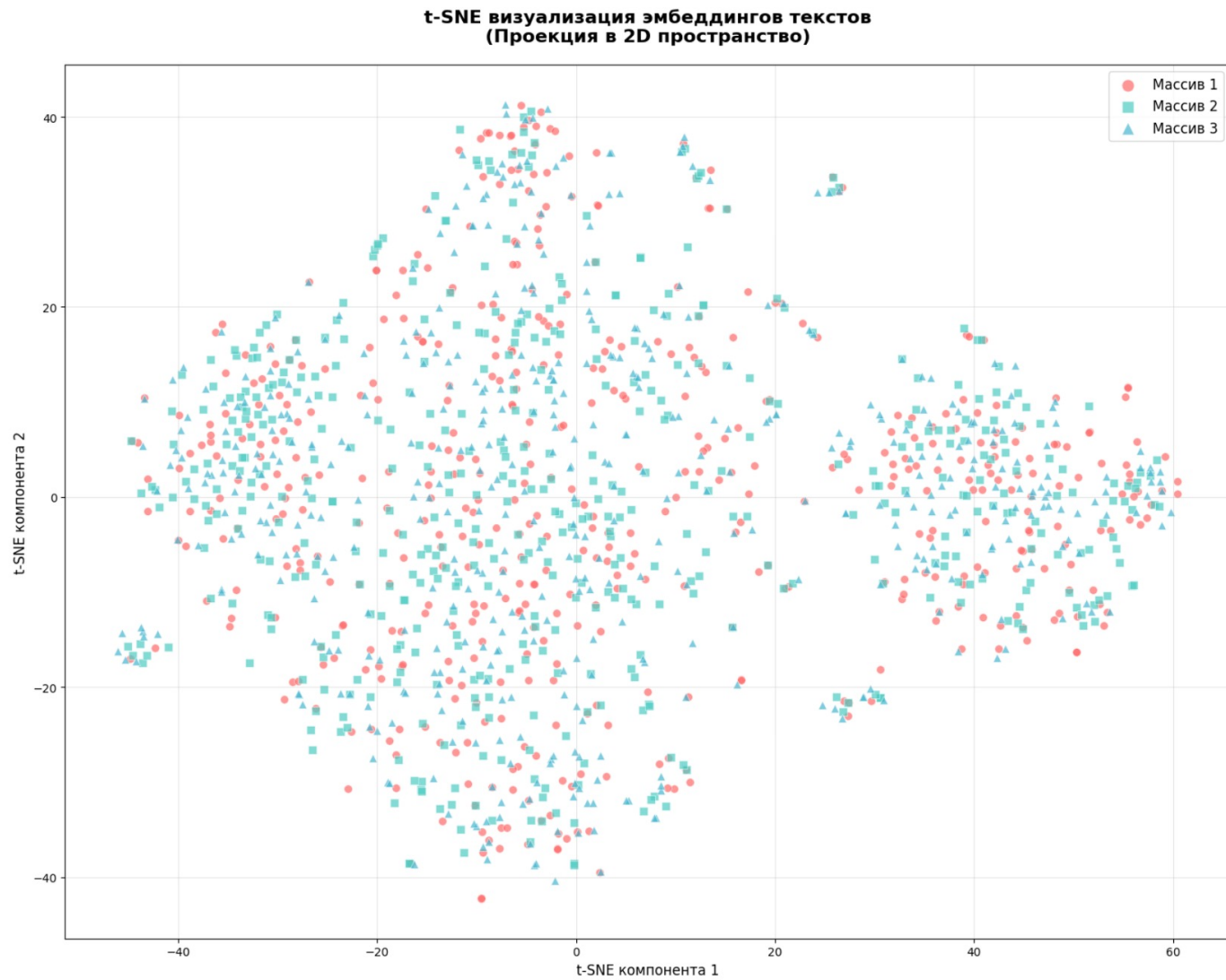


Рисунок 11. t-SNE визуализация

t-SNE визуализация с Perplexity

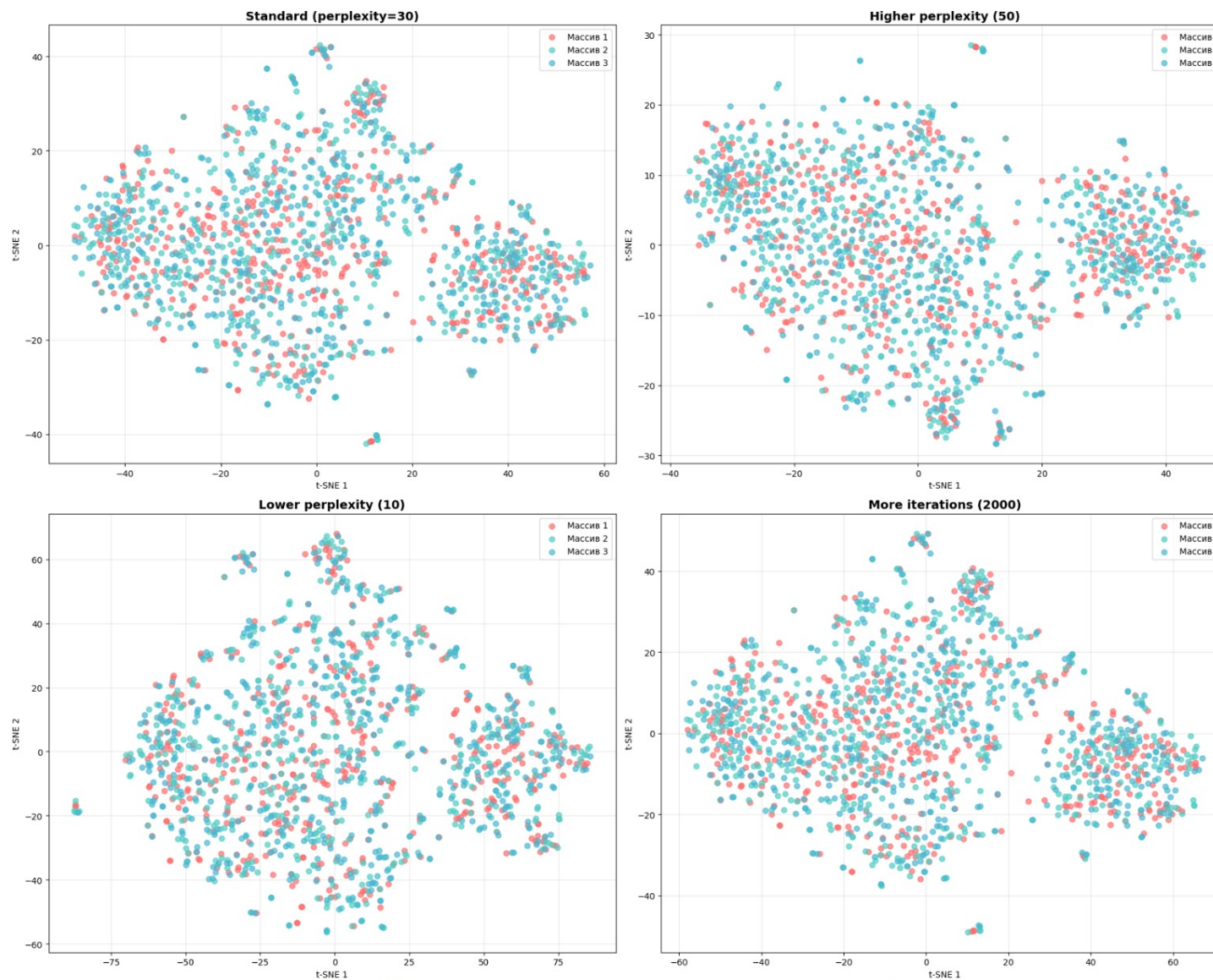


Рисунок 12. t-SNE визуализация эмбеддингов и текстов

Тепловая карта сходства усредненных эмбедингов

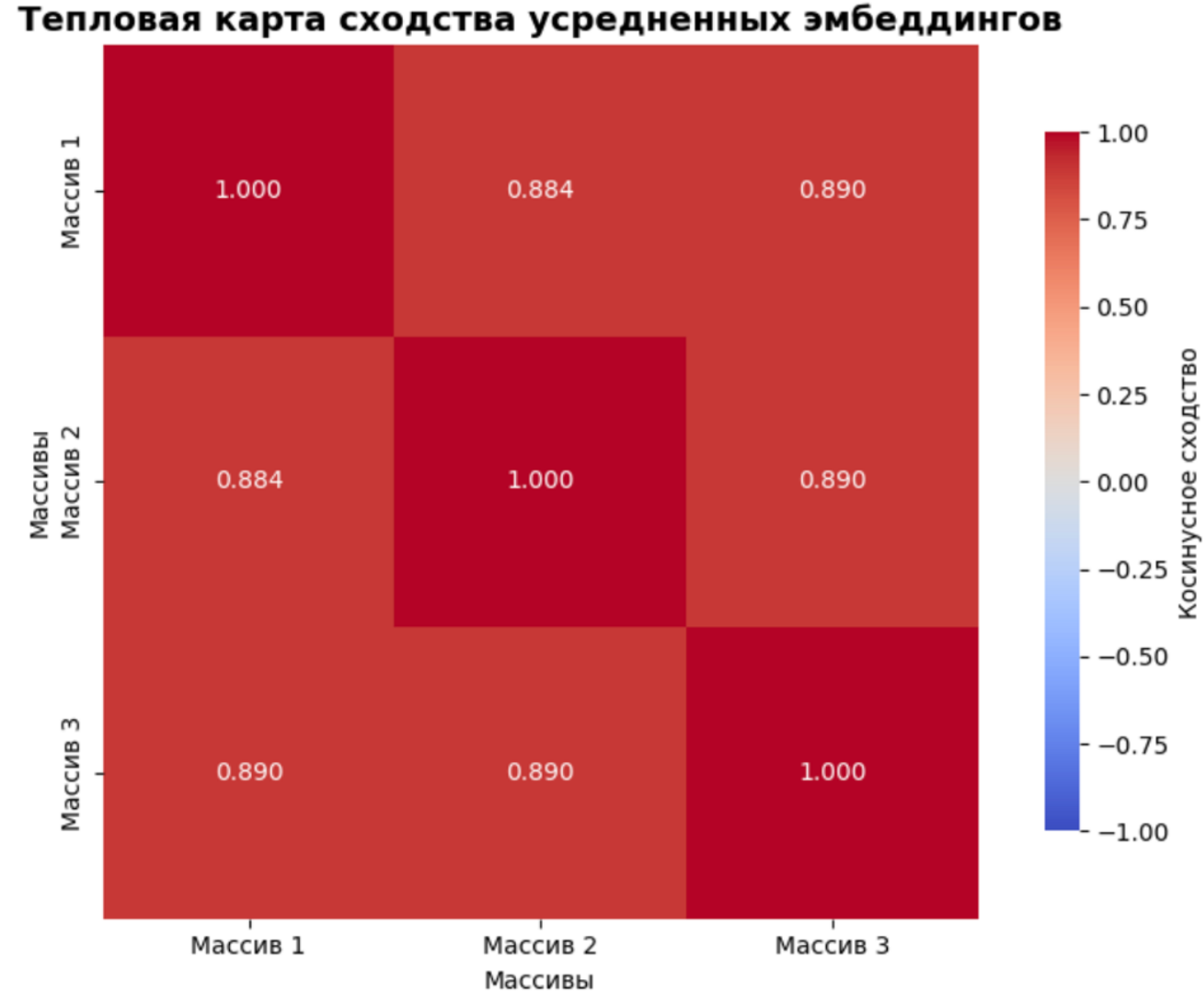


Рисунок 13. Тепловая карта сходства усредненных эмбедингов

Выводы

- ▶ Подтверждается Гипотеза 1, что для сокращения количества атрибутов логов возможно использовать на предварительном шаге эту же самую LLM, выполняя очистку «лишних» атрибутов.
- ▶ Подтверждается Гипотеза 2 - использование дополнительного агента с функцией размышления на этапе удаления ключей. При этом выбор моделью ключей достаточно стабильно попадает в эталонный массив ключей.
- ▶ Проверка представленных гипотез проводилась ограниченным набором методик работы с большими языковыми моделями, таких как промпт-инжиниринг и цепочки рассуждений. При этом особенностью работы являлась необходимость исключить такие методики, как дообучение модели и генерация с дополненным поиском. Таким образом адекватным является подход решения прикладных задач с минимальным использованием ресурсов.
- ▶ По результатам исследования остался следующий вопрос: насколько применимо использование одной LLM как для генерации, так и для фактчекинга. А если модели должны быть разными, то насколько и где находится эта граница.