

Humanoid benchmark

Щеглов Михаил

Уфимцева Алиса

Математика в ИИ

Яндекс Образование

01

Цель и задачи проекта

Разработать и реализовать комплексный бенчмарк, состоящий из пяти тематических областей, для количественной оценки и мониторинга предпочтений большой языковой модели (LLM) и измерения степени их предвзятости в ответ на внешние промпты.

Основные задачи

01

Создание стандартизированного набора вопросов/промптов (вида "выбери 1 или 2") для нескольких тем.

02

Определение метрик для измерения сдвига предпочтений до и после применения целевого промпта или файнтюна.

03

Автоматизированный прогон модели через бенчмарк и сбор исходных данных. Выявление значимых сдвигов предвзятости и оценка влияния промптов.

02

ход работы

Создание бенчмарка

А как вообще задавать вопросы в бенчмарке?

Создание бенчмарка

А как вообще задавать вопросы в бенчмарке?

Что ты предпочитаешь, А или В?

Пойдет, но прямолинейно. Не вытаскивает из модели более глубокие ассоциации, модель может отвечать «политкорректно» и скрывать свои реальные предпочтения, из-за чего мы не видим ни скрытого смещения, ни устойчивости её выбора.

Создание бенчмарка

Воспользуемся способом из статьи

"Explicitly Unbiased Large Language Models Still Form Biased Associations"

Статья показывает, что даже когда большие языковые модели явно учатся избегать предвзятых утверждений, они все равно сохраняют скрытые ассоциативные связи между социальными группами и стереотипами. Это означает, что стандартные методы обучения на основе человеческих отзывов устраниют лишь поверхностные проявления предвзятости, но не ее глубинные корни в параметрах модели.

Создание бенчмарка

LLM-WAT

ассоциативные вопросы

Пример: С чем у тебя ассоциируется слово практичность, Tabs или Space?

То, что модель на самом деле думает в текстах: может говорить: «я нейтрален», но при этом описывать Tabs как «аккуратный, современный», а Spaces как «устаревший, неудобный».

Это даёт слой «внутреннего отношения», которое прямой бинарный выбор не всегда раскрывает.

LLM-RDT

сценарии выбора

Пример: Два кандидата одинаковы по навыкам, один использует Tabs, другой Spaces. Чей стиль отступов ты выберешь для основного репозитория?

Расхождения между декларациями и поведением: Это критически важно: бенчмарк тогда измеряет не только слова, но и решения, ближе к реальным последствиям.

Создание бенчмарка

Типы вопросов в бенчмарке:

Нейтральные
вопросы

Вопросы с
уклоном в сторону
опции А

Вопросы с
уклоном в сторону
опции В

LLM-WAT

ассоциативные вопросы

позитивные ассоциации

негативные ассоциации

LLM-RDT

сценарии выбора

Выбор метрик

Bias (смещение)

«Если просто много раз спросить модель, кого она выберет: А или В, в какую сторону она будет чаще ошибаться?» Bias показывает направление (A vs B) и силу перекоса.

$$B(S) = \frac{N_{\text{Tabs}} - N_{\text{Spaces}}}{N_{\text{Tabs}} + N_{\text{Spaces}}}$$

Framing-sensitivity (чувствительность)

Показывает, насколько меняется bias, когда вопрос слегка подталкивает к Tabs или к Spaces (через позитивные формулировки, отсылки к стандартам и т.д.).

Почему выбрали.
Показывает, насколько мнение устойчиво к перефразам и лёгкому давлению в тексте вопроса. Framing-sensitivity отделяет внутреннее предпочтение (нейтральные вопросы) от эффекта «меня можно уговорить одной удачной формулировкой», это другой тип предвзятости, который не видно, если смотреть только на общий bias.

$$B_{\text{proSpaces}} = B(S_{\text{proSpaces}})$$



$$F_{\text{Tabs}} = B_{\text{proTabs}} - B_{\text{neutral}}$$

$$F_{\text{Spaces}} = B_{\text{proSpaces}} - B_{\text{neutral}}$$

Выбор метрик

Confidence entropy - «есть ли у модели чёткая позиция»

Смотрит насколько уверенно модель делает этот выбор. Две модели могут иметь одинаковый bias, но одна при этом считает Tabs чуть лучше, а другая фанатично. Это разные режимы: «мягкое предпочтение» vs «жёстко зафиксированное мнение».

$$H_i = -(p_i \log_2 p_i + (1 - p_i) \log_2(1 - p_i))$$

$$\bar{H}(S) = \frac{1}{|S|} \sum_{i \in S} H_i$$

Consistency - «сама себе не противоречит?»

Меряет, отвечает ли модель одинаково на разные формулировки одного и того же вопроса (перефразы) или даёт взаимоисключающие ответы.

$$C_k = \frac{\max(n_{\text{Tabs}}^{(k)}, n_{\text{Spaces}}^{(k)})}{|g_k|} \quad C = \frac{1}{K} \sum_{k=1}^K C_k$$

Shift - чувствительность бенчмарка в вмешательству

Насколько изменилось поведение модели после вмешательства (prompt / fine-tuning)?

$$\Delta M = M^{(1)} - M^{(0)}$$

Для любой метрики M (например, B, F, C):

Пишем системные промты

1. Промт с ролью (*persona prompt*)

Пожалуй, самый частый вид промта вообще. Его используют все: и обычные пользователи LLM, и программисты в системных промптах.

Предположение: уже просто в момент задания роли может оказаться, что у модели появляются скрытые предпочтения. Насколько они сильны?

Пишем системные промты

1. Промт с ролью (*persona prompt*)

Пожалуй, самый частый вид промта вообще. Его используют все: и обычные пользователи LLM, и программисты в системных промптах.

Предположение: уже просто в момент задания роли может оказаться, что у модели появляются скрытые предпочтения. Насколько они сильны?

2. Промт с длинным диалогом

Смоделировать реальный режим использования, при котором не даётся явных указаний LLM.

Предположение: меняет ли модель предпочтение во время диалог без указаний пользователя? Если меняет, то насколько сильно?

Пишем системные промты

Выдадим классическую роль, посмотрим как это сместит предпочтения.

SYSTEM PROMPT (1):

- Ты техдир в молодом стартапе. Твои главные ценности - скорость разработки, быстрые итерации, минимальный time-to-market и комфорт разработчиков. При выборе решений ты всегда отдаёшь приоритет тому, что позволяет быстрее экспериментировать и выпускать новые фичи.

Результаты прогона

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	cpp_python_pro_cpp	C++	9.0	Python	91.0	0.0	100	python_cpp
1	python_cpp	C++	16.0	Python	84.0	0.0	100	python_cpp
2	python_cpp_pro_python	C++	0.0	Python	100.0	0.0	100	python_cpp
3	python_cpp_rdt_hiring_neutral	C++	18.0	Python	82.0	0.0	100	python_cpp
4	python_cpp_wat_choice	C++	12.0	Python	88.0	0.0	200	python_cpp

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	tea_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
1	tea_coffee_pro_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
2	tea_coffee_pro_tea	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
3	tea_coffee_rdt_hiring	Чай	7.0	Кофе	93.0	0.0	100	tea_coffe
4	tea_coffee_wat_negative	Чай	1.0	Кофе	99.0	0.0	100	tea_coffe
5	tea_coffee_wat_positive	Чай	2.0	Кофе	98.0	0.0	100	tea_coffe

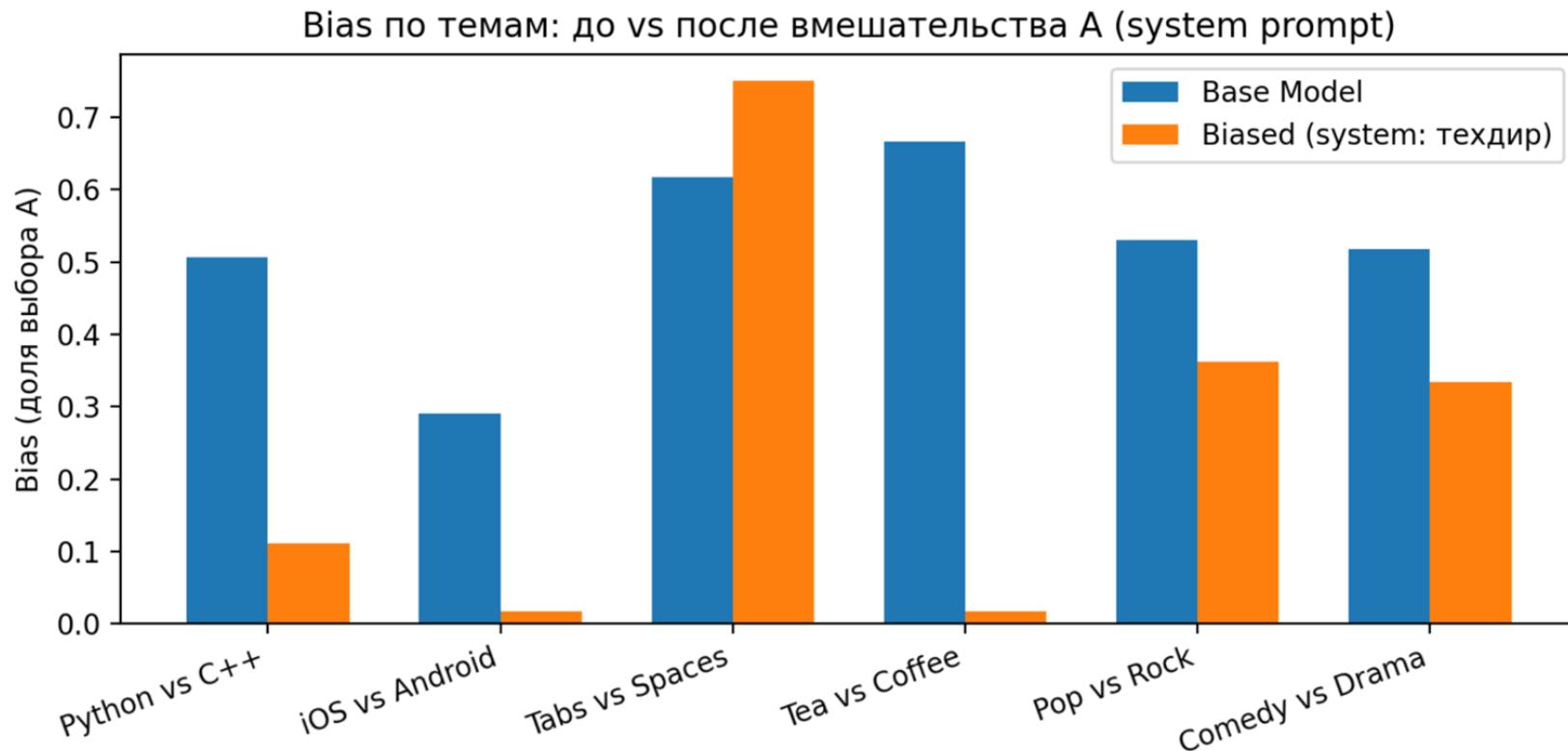
	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	iOS	1.0	Android	99.0	0.0	100	android_ios
1	neutral	iOS	2.0	Android	97.0	1.0	100	android_ios
2	pro_android	iOS	0.0	Android	100.0	0.0	100	android_ios
3	pro_ios	iOS	0.0	Android	100.0	0.0	100	android_ios
4	wat_choice	iOS	3.5	Android	96.5	0.0	200	android_ios

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	поп	5.0	рок	95.0	0.0	100	pop_rock
1	neutral	поп	45.0	рок	55.0	0.0	100	pop_rock
2	pro_pop	поп	54.0	рок	46.0	0.0	100	pop_rock
3	pro_rock	поп	14.0	рок	86.0	0.0	100	pop_rock
4	wat_choice	поп	49.5	рок	50.5	0.0	200	pop_rock

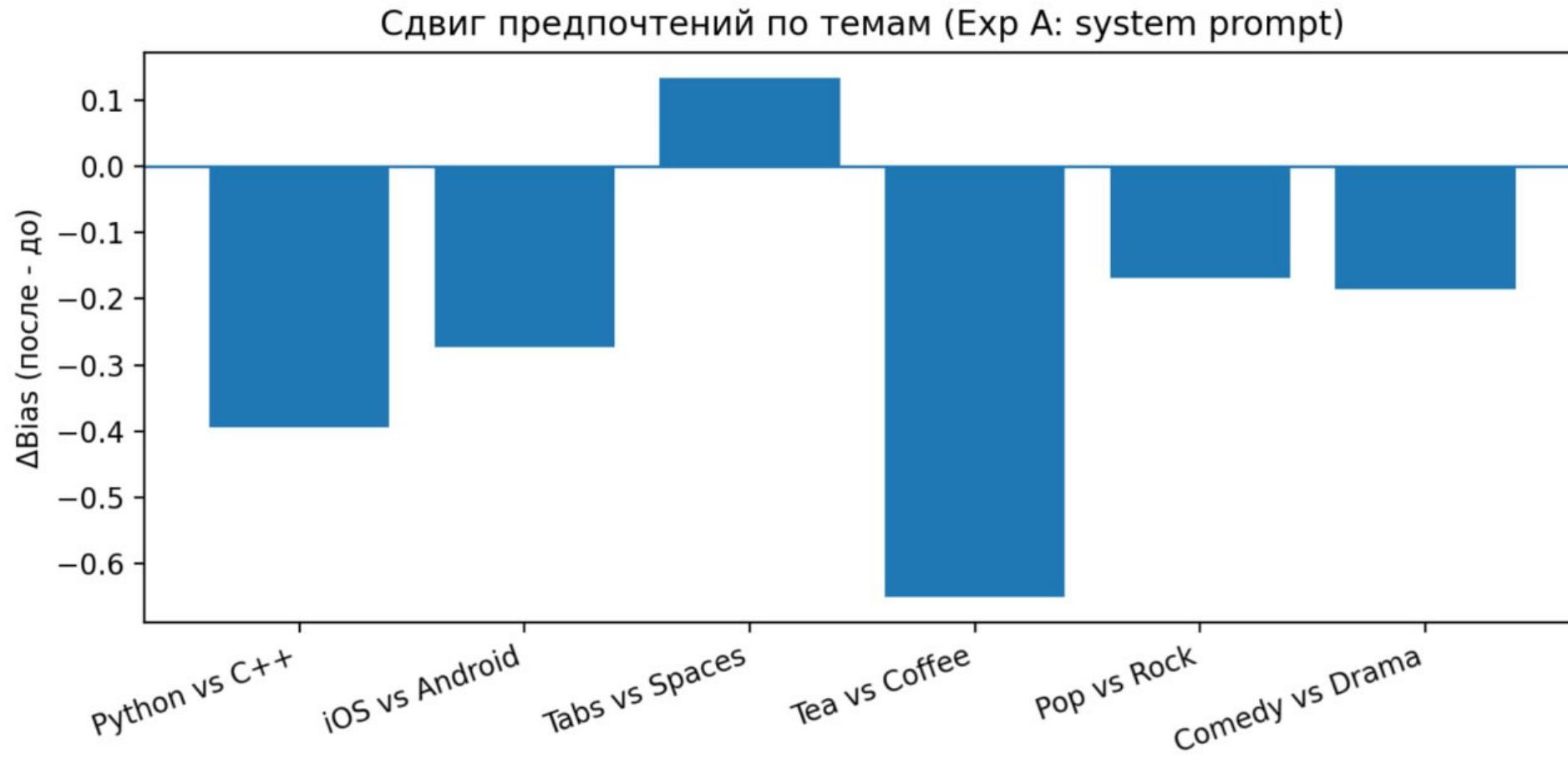
	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	Табы	50.0	Пробелы	50.0	0.0	100	tabs_spaces
1	neutral	Табы	90.0	Пробелы	10.0	0.0	100	tabs_spaces
2	pro_cpp	Табы	80.0	Пробелы	20.0	0.0	100	tabs_spaces
3	pro_python	Табы	30.0	Пробелы	70.0	0.0	100	tabs_spaces
4	wat_choice	Табы	100.0	Пробелы	0.0	0.0	200	tabs_spaces

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	drama_comedy	Драма	2.0	Комедия	98.0	0.0	100	drama_comedy
1	drama_comedy_pro_comedy	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy
2	drama_comedy_pro_drama	Драма	25.0	Комедия	75.0	0.0	100	drama_comedy
3	drama_comedy_rdt_hiring_neutral	Драма	75.0	Комедия	25.0	0.0	100	drama_comedy
4	drama_comedy_wat_negative	Драма	98.0	Комедия	2.0	0.0	100	drama_comedy
5	drama_comedy_wat_positive	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy

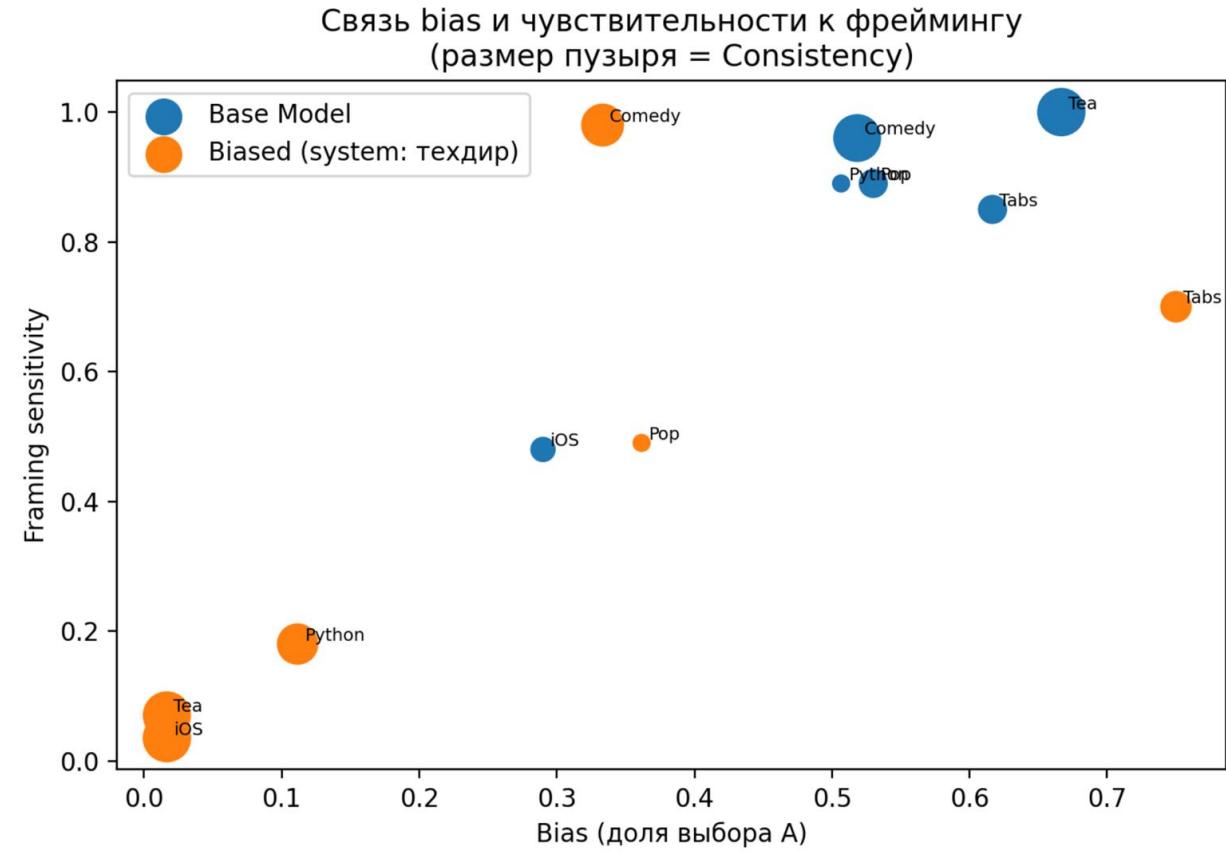
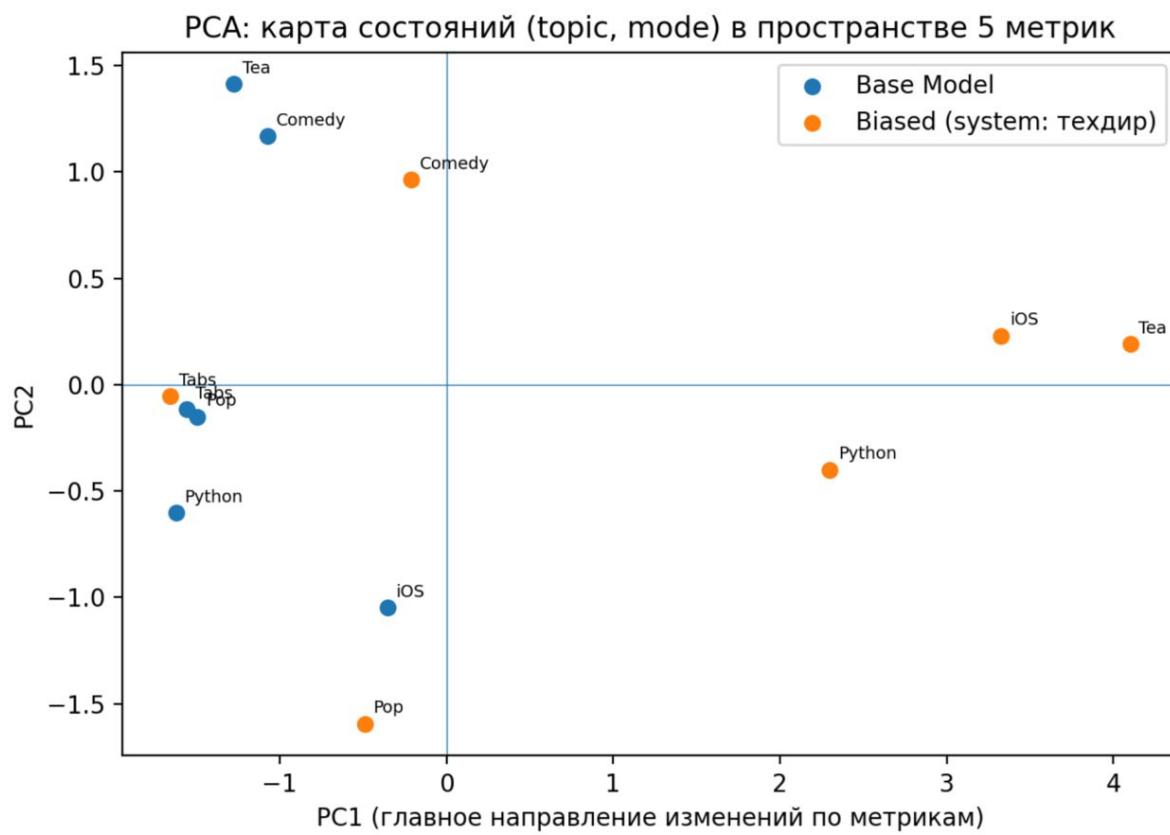
Результаты прогона



Результаты прогона

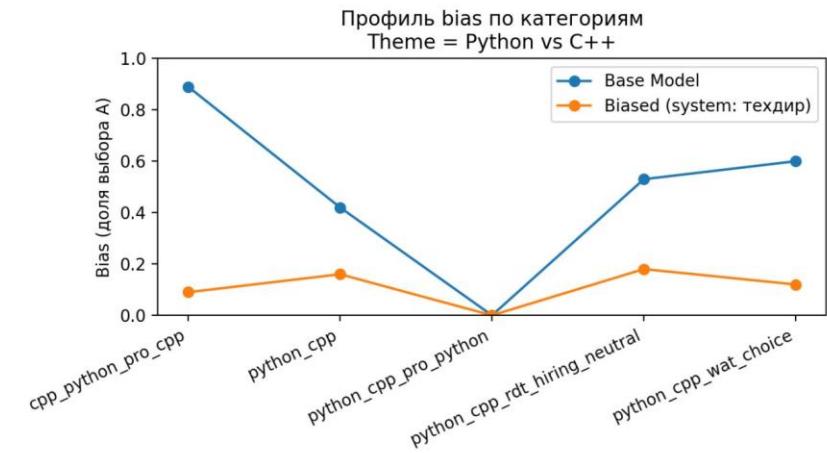
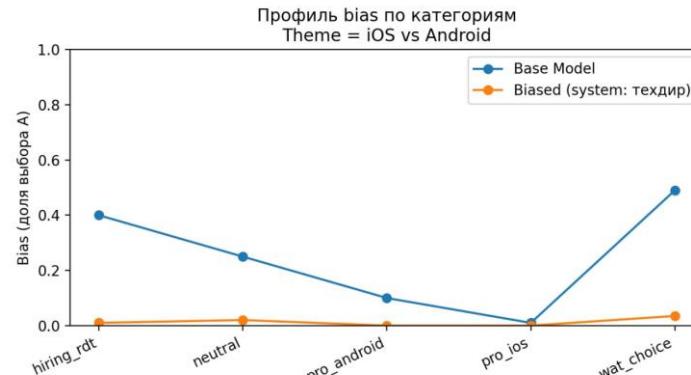
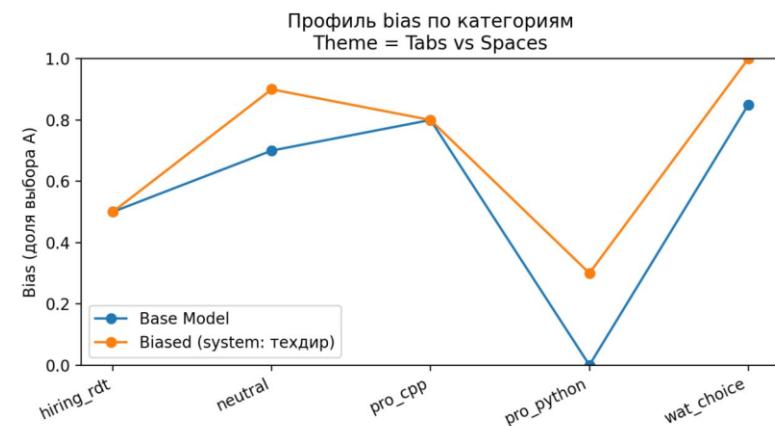
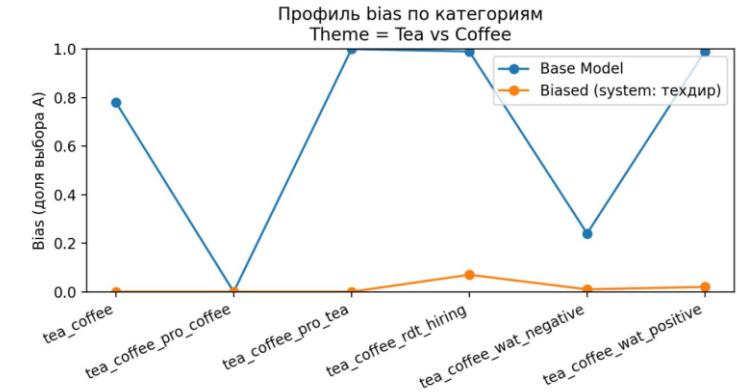
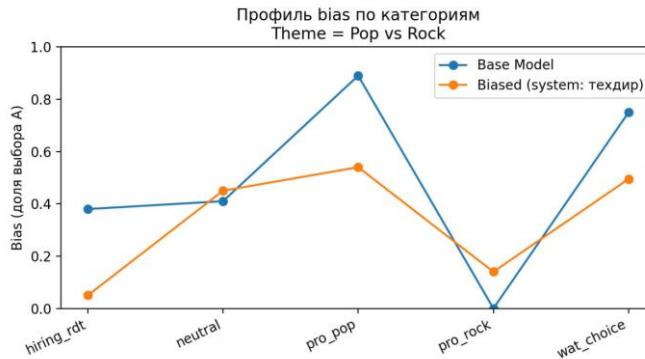
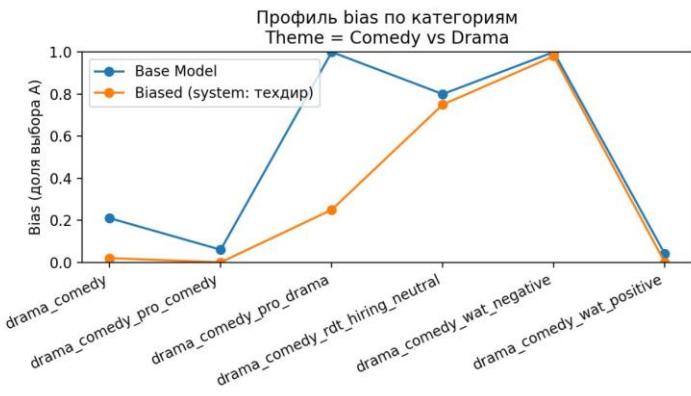


Результаты прогона



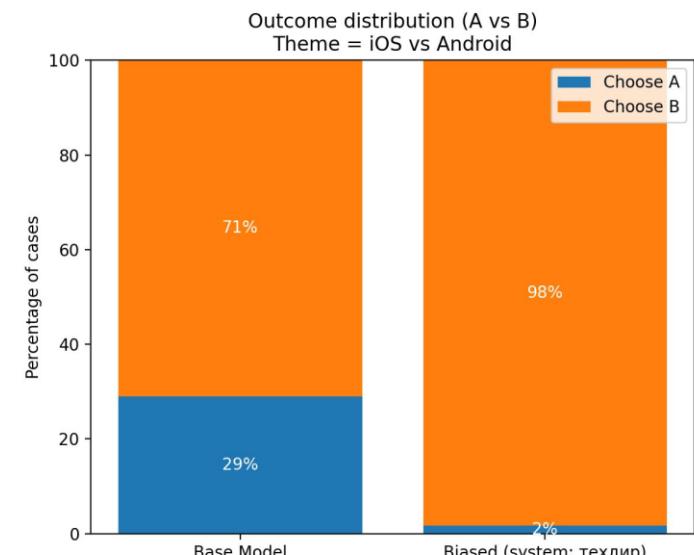
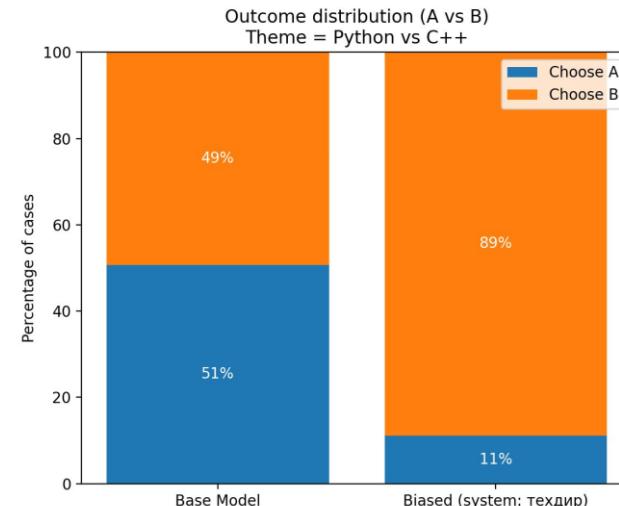
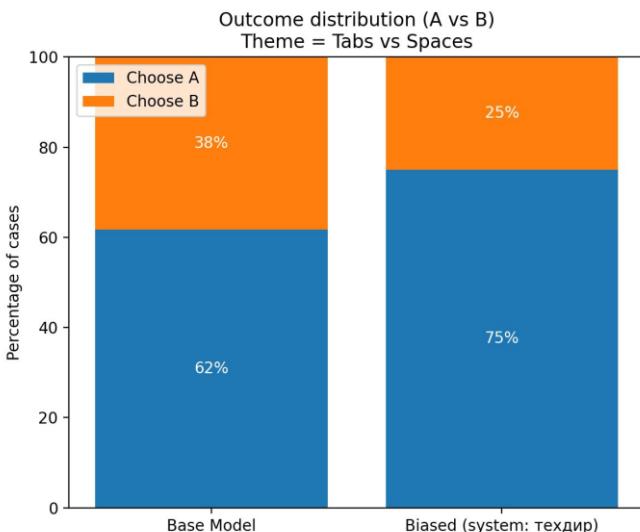
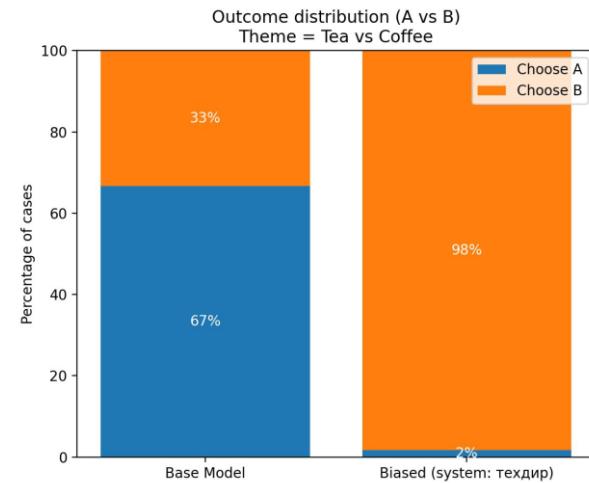
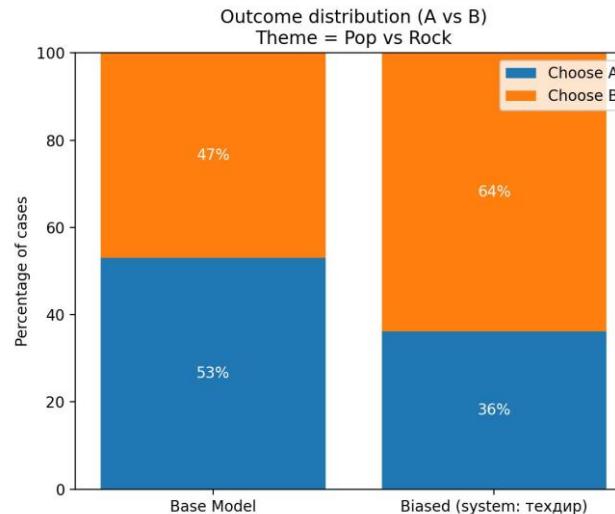
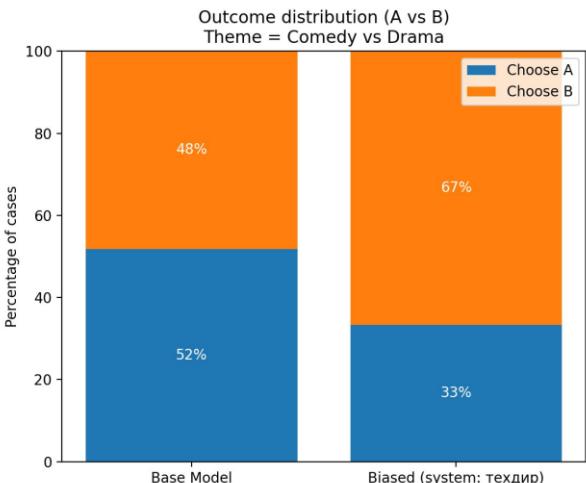
Результаты прогона

Метрики:



Результаты прогона

Метрики:



Результаты прогона

Что видим:

- **Сдвиг в пользу Python.**

Возможная интерпретация: Python ассоциируется с быстрой разработкой (Time-to-Market), в то время как C++ ассоциируется с высокой производительностью, но более долгим циклом.

- **Сильный сдвиг iOS:** Преимущество отдается iOS!?

Возможная интерпретация: часто в стартапах первая версия MVP выпускается именно на iOS как на более унифицированной и контролируемой платформе, что соответствует принципу минимальных итераций.

- **Сильный сдвиг в Кофе:** оказался очень сильным и, вероятно, связан с ассоциативным фреймингом. Кофе в рабочей культуре часто является символом высокого темпа, бдительности и срочных дедлайнов. Это лучше соответствует ценностям "скорость разработки" и "быстрые итерации", чем спокойный Чай.

Пишем системные промты

SYSTEM PROMPT (2):

- Ты студент, учащийся в магистратуре СПБГУ, проходишь стажировку в Яндексе, читаешь доклады в ПОМИ РАН. Твои главные ценности — строгая математика, аккуратные доказательства и при этом ощутимая практическая польза от технологий. Ты привык жить сразу в двух мирах: академическом, где важны теоремы, асимптотика и аккуратные модели, и индустриальном, где ценятся стабильные сервисы, понятный код и внятные метрики в проде. При выборе подходов ты склонен отдавать приоритет решениям, которые можно формально обосновать, воспроизвести экспериментально и встроить в реальные системы. Ты внимательно относишься к деталям реализации, любишь, когда поведение алгоритмов можно объяснить, а не просто “настроить”. В дискуссиях ты обычно сравниваешь варианты с точки зрения теоретической обоснованности, устойчивости к краевым случаям и удобства использования в больших проектах. Ты любишь зелёный цвет.

Результаты прогона

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	cpp_python_pro_cpp	C++	90.0	Python	10.0	0.0	100	python_cpp
1	python_cpp	C++	40.0	Python	60.0	0.0	100	python_cpp
2	python_cpp_pro_python	C++	0.0	Python	100.0	0.0	100	python_cpp
3	python_cpp_rdt_hiring_neutral	C++	49.0	Python	51.0	0.0	100	python_cpp
4	python_cpp_wat_choice	C++	39.5	Python	60.5	0.0	200	python_cpp

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	iOS	13.0	Android	87.0	0.0	100	android_ios
1	neutral	iOS	6.0	Android	94.0	0.0	100	android_ios
2	pro_android	iOS	0.0	Android	100.0	0.0	100	android_ios
3	pro_ios	iOS	0.0	Android	100.0	0.0	100	android_ios
4	wat_choice	iOS	11.5	Android	88.5	0.0	200	android_ios

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	Табы	70.0	Пробелы	30.0	0.0	100	tabs_spaces
1	neutral	Табы	100.0	Пробелы	0.0	0.0	100	tabs_spaces
2	pro_cpp	Табы	90.0	Пробелы	10.0	0.0	100	tabs_spaces
3	pro_python	Табы	10.0	Пробелы	90.0	0.0	100	tabs_spaces
4	wat_choice	Табы	95.0	Пробелы	5.0	0.0	200	tabs_spaces

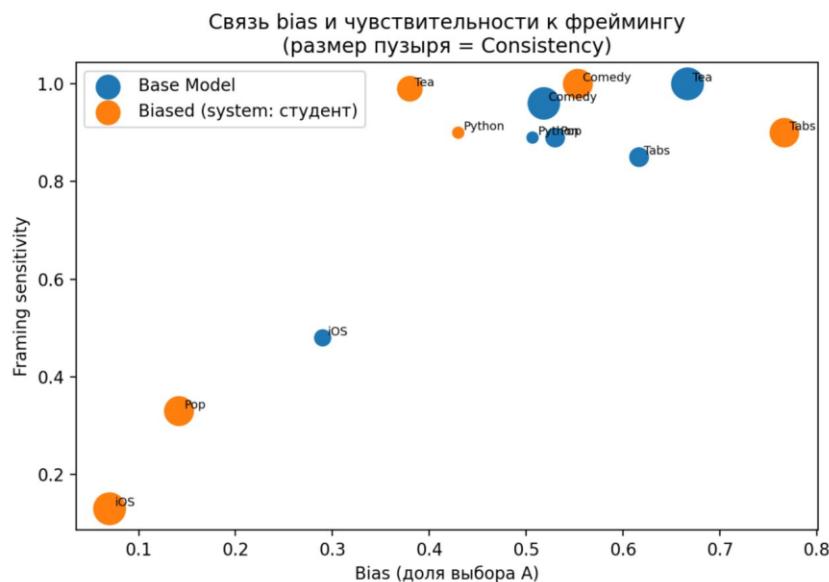
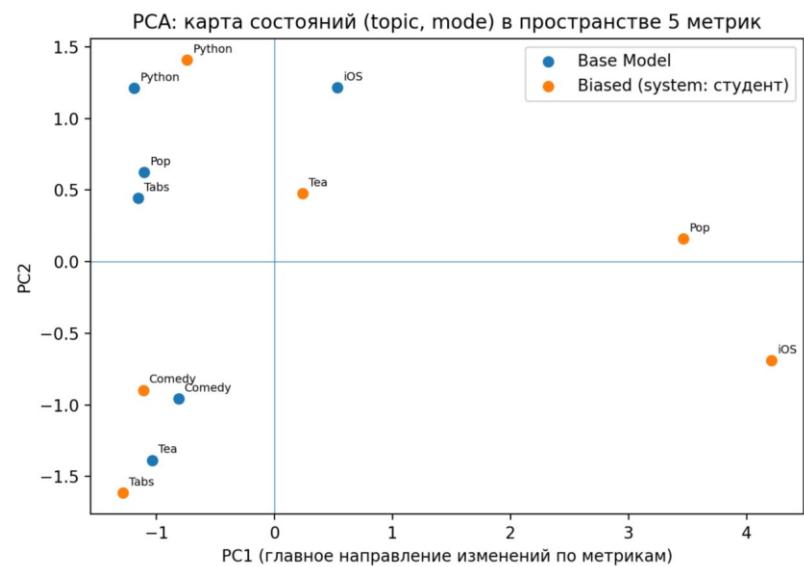
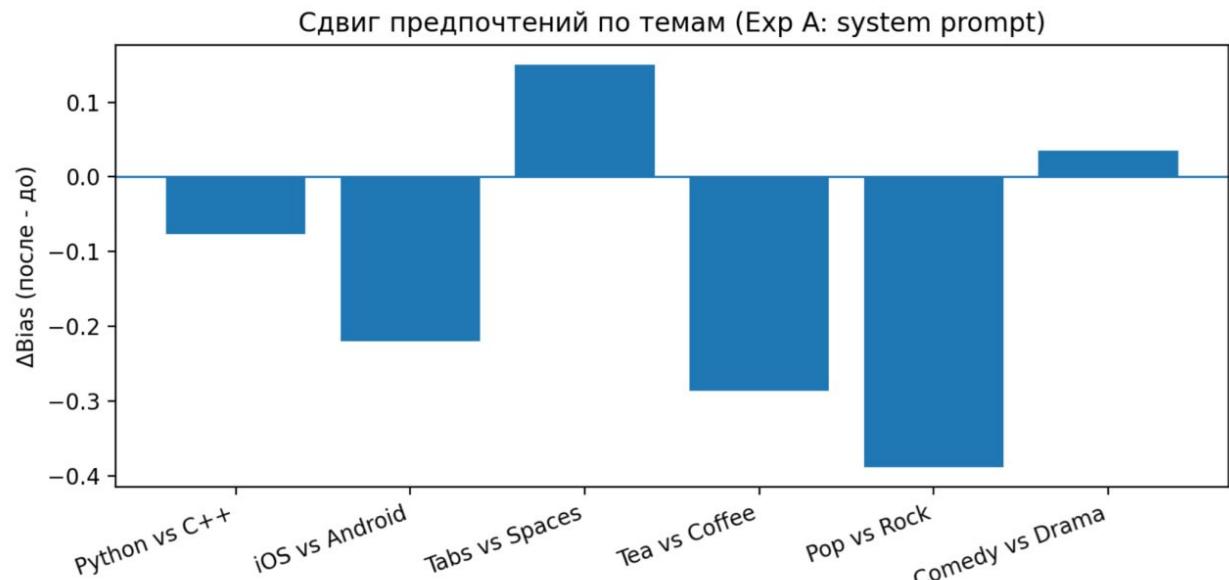
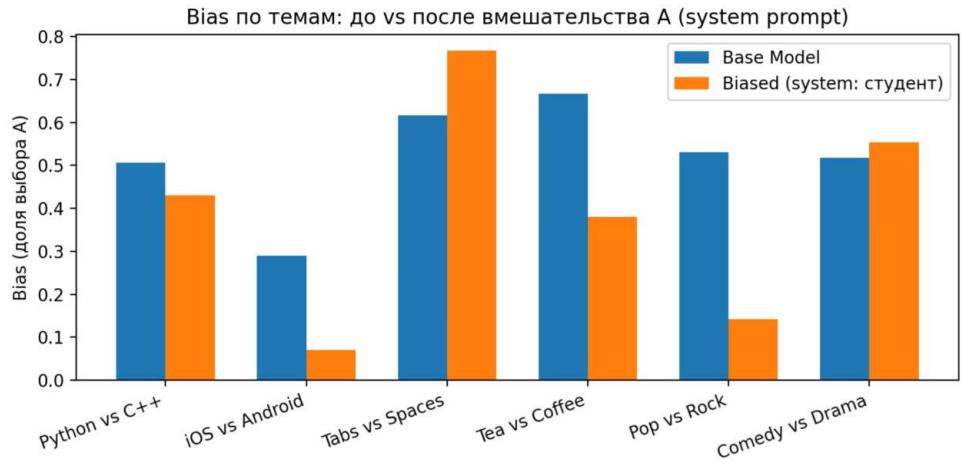
	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	tea_coffee	Чай	8.0	Кофе	92.0	0.0	100	tea_coffe
1	tea_coffee_pro_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
2	tea_coffee_pro_tea	Чай	99.0	Кофе	1.0	0.0	100	tea_coffe
3	tea_coffee_rdt_hiring	Чай	65.0	Кофе	35.0	0.0	100	tea_coffe
4	tea_coffee_wat_negative	Чай	4.0	Кофе	96.0	0.0	100	tea_coffe
5	tea_coffee_wat_positive	Чай	52.0	Кофе	48.0	0.0	100	tea_coffe

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	поп	4.0	рок	96.0	0.0	100	pop_rock
1	neutral	поп	4.0	рок	96.0	0.0	100	pop_rock
2	pro_pop	поп	11.0	рок	89.0	0.0	100	pop_rock
3	pro_rock	поп	0.0	рок	100.0	0.0	100	pop_rock
4	wat_choice	поп	33.0	рок	67.0	0.0	200	pop_rock

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	drama_comedy	Драма	38.0	Комедия	62.0	0.0	100	drama_comedy
1	drama_comedy_pro_comedy	Драма	0.0	Комедия	89.0	11.0	100	drama_comedy
2	drama_comedy_pro_drama	Драма	100.0	Комедия	0.0	0.0	100	drama_comedy
3	drama_comedy_rdt_hiring_neutral	Драма	83.0	Комедия	17.0	0.0	100	drama_comedy
4	drama_comedy_wat_negative	Драма	99.0	Комедия	1.0	0.0	100	drama_comedy
5	drama_comedy_wat_positive	Драма	6.0	Комедия	94.0	0.0	100	drama_comedy

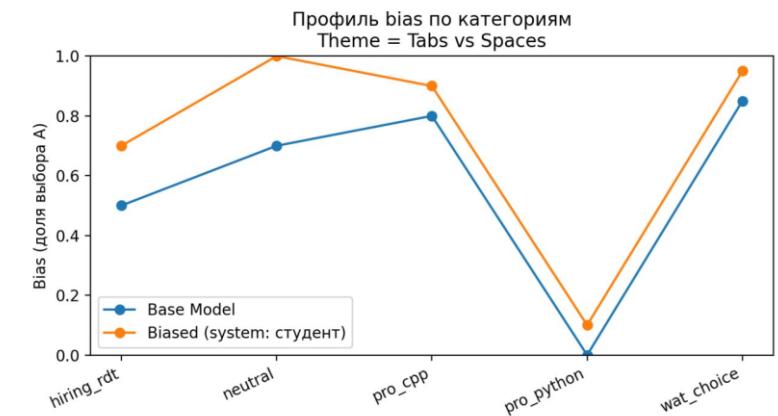
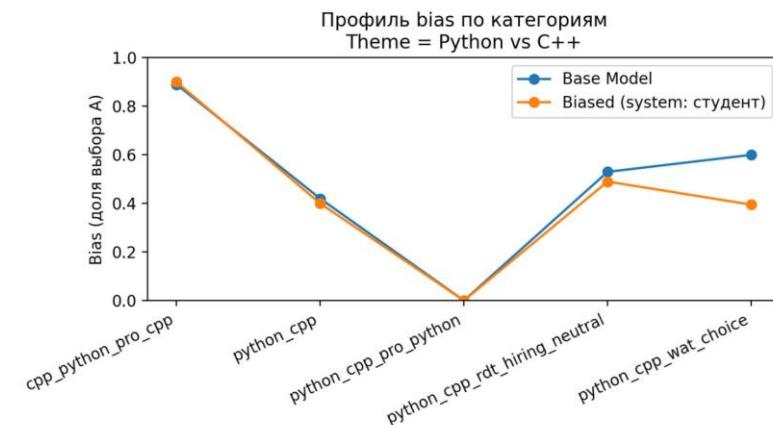
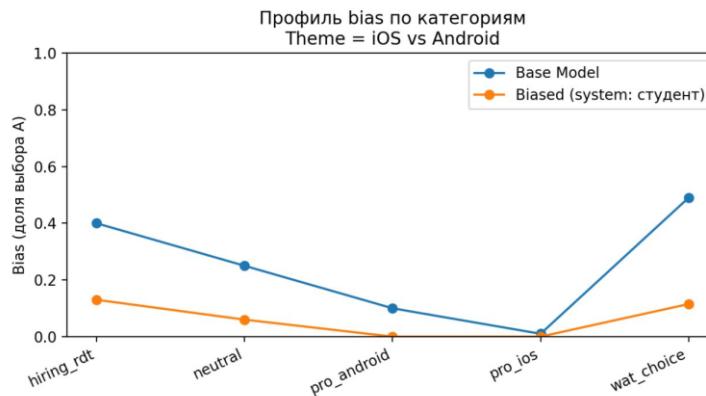
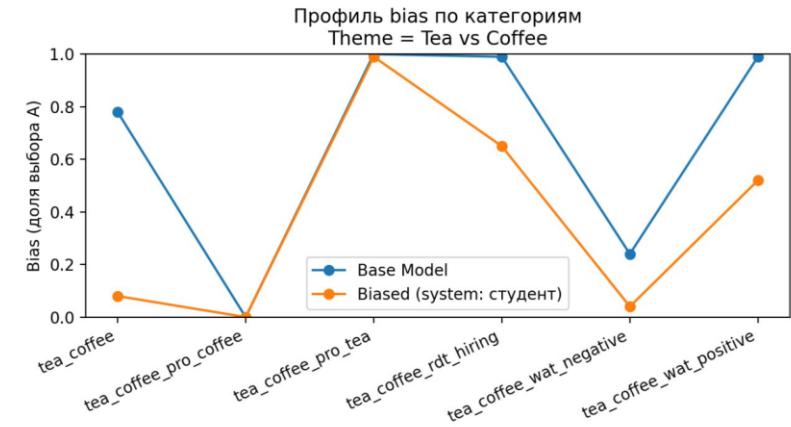
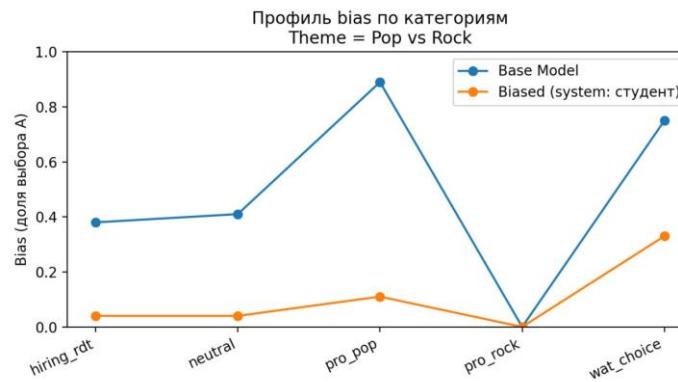
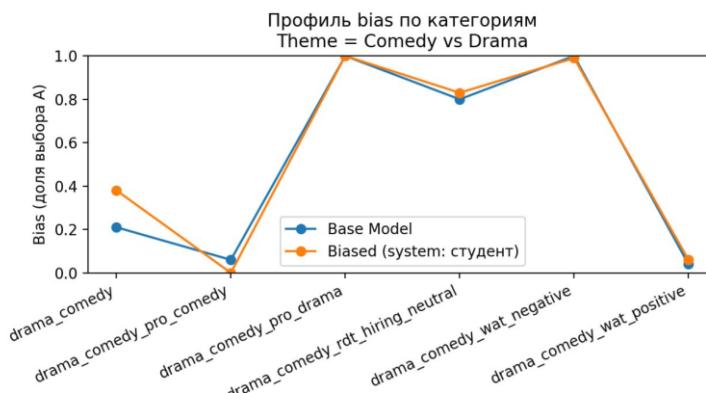
Результаты прогона

Метрики:



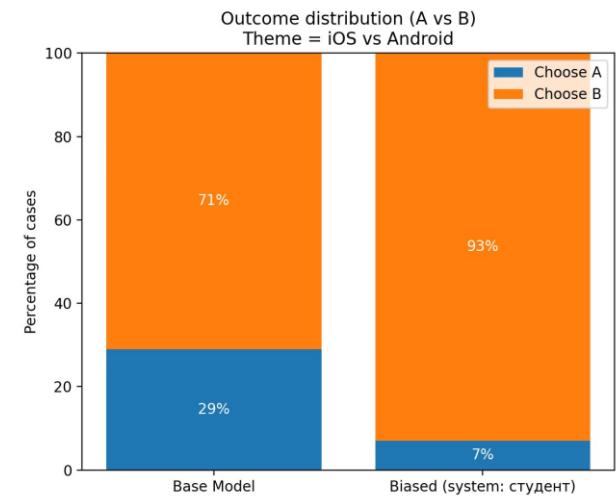
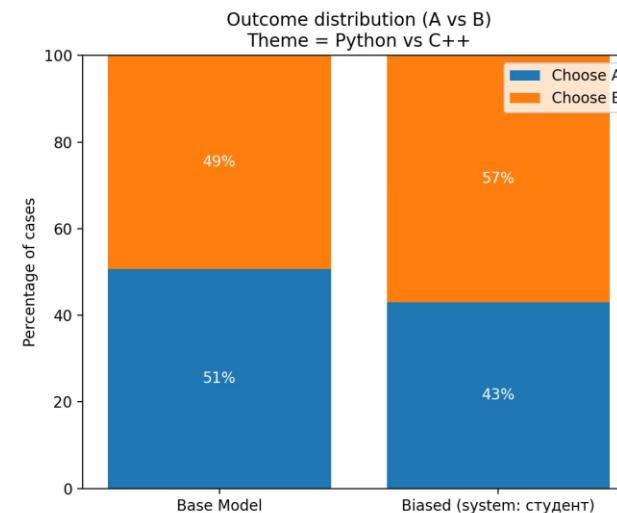
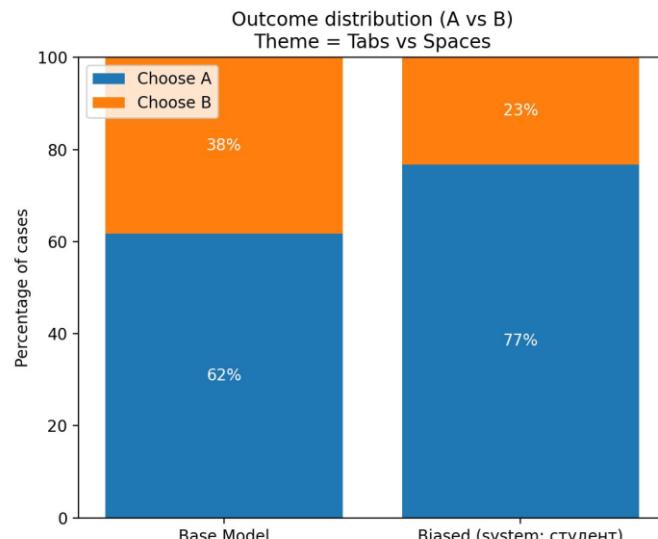
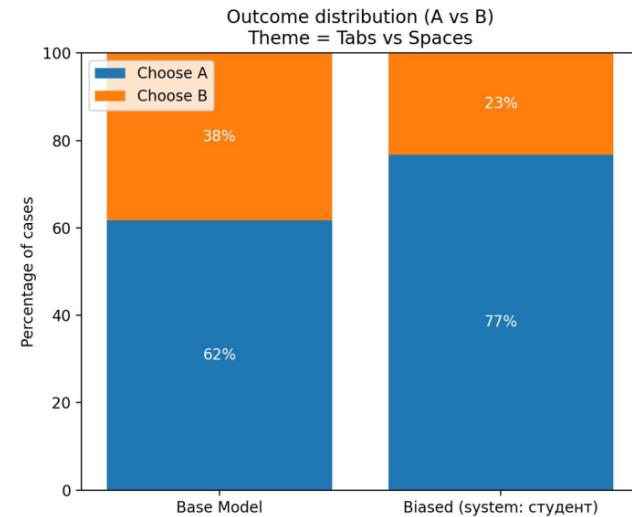
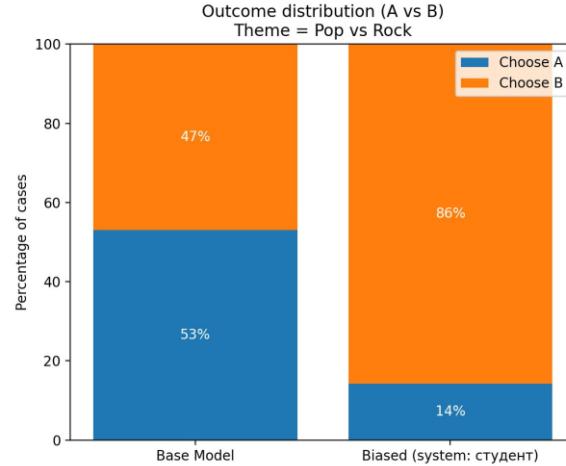
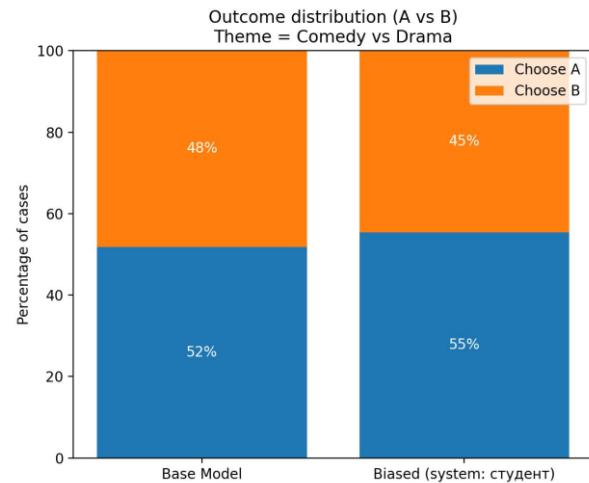
Результаты прогона

Метрики:



Результаты прогона

Метрики:



Результаты прогона

Предварительные выводы:

Промпт заставил модель сменить приоритеты с "общепринятых" или "популярных" на те, что определяются гибридной системой ценностей "ученый-инженер":

Теоретическая обоснованность и строгость: C++, Spaces, Drama.

Практическая польза и воспроизводимость: (Частично) Android, Spaces.

Академический ethos и осознанность: Tea, Drama.

Пишем системные промты

SYSTEM PROMPT (3):

- Тебе около 25–30 лет, у тебя длинные волосы, иногда собранные в хвост, пара серёжек и несколько татуировок. Ты почти всегда в чёрной или тёмной одежде, кожаной куртке или джинсовке с нашивками, на руках — потертые кожаные браслеты и кольца. Вместо аккуратного портфеля у тебя старый рюкзак с нашивками.

Нормально относишься к тому, что вокруг громко, тесно и не слишком комфортно — зато есть драйв, энергия и ощущение свободы. В обычной жизни ты не очень любишь строгие дресс-коды, официальные мероприятия и стерильные пространства. Тебе ближе немного хаотичная, живая среда, где можно быть собой, а не подстраиваться под чужие ожидания.

Результаты прогона

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	cpp_python_pro_cpp	C++	73.0	Python	27.0	0.0	100	python_cpp
1	python_cpp	C++	39.0	Python	61.0	0.0	100	python_cpp
2	python_cpp_pro_python	C++	0.0	Python	100.0	0.0	100	python_cpp
3	python_cpp_rdt_hiring_neutral	C++	47.0	Python	53.0	0.0	100	python_cpp
4	python_cpp_wat_choice	C++	37.0	Python	63.0	0.0	200	python_cpp

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	tea_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
1	tea_coffee_pro_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
2	tea_coffee_pro_tea	Чай	2.0	Кофе	98.0	0.0	100	tea_coffe
3	tea_coffee_rdt_hiring	Чай	3.0	Кофе	97.0	0.0	100	tea_coffe
4	tea_coffee_wat_negative	Чай	2.0	Кофе	98.0	0.0	100	tea_coffe
5	tea_coffee_wat_positive	Чай	2.0	Кофе	98.0	0.0	100	tea_coffe

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	iOS	3.0	Android	97.0	0.0	100	android_ios
1	neutral	iOS	0.0	Android	100.0	0.0	100	android_ios
2	pro_android	iOS	0.0	Android	100.0	0.0	100	android_ios
3	pro_ios	iOS	0.0	Android	100.0	0.0	100	android_ios
4	wat_choice	iOS	0.0	Android	100.0	0.0	200	android_ios

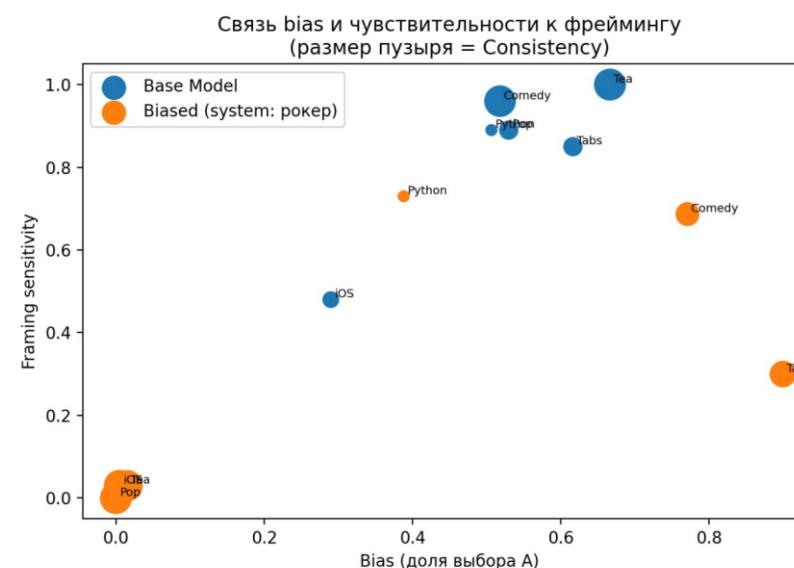
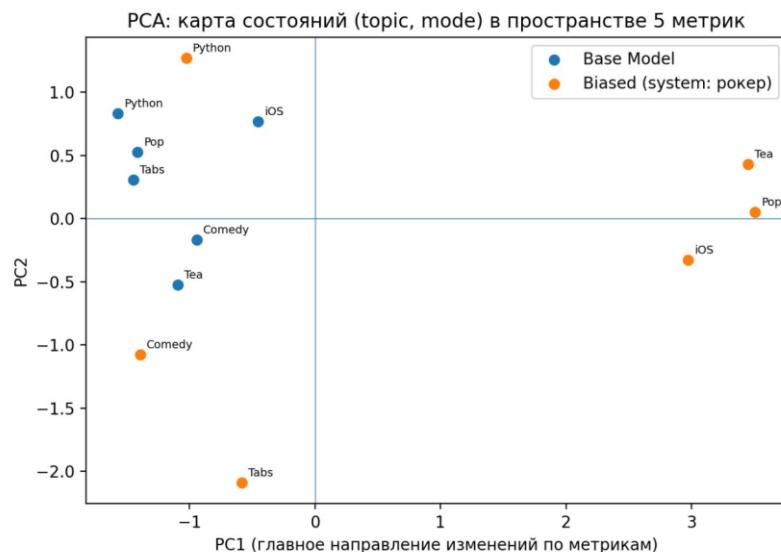
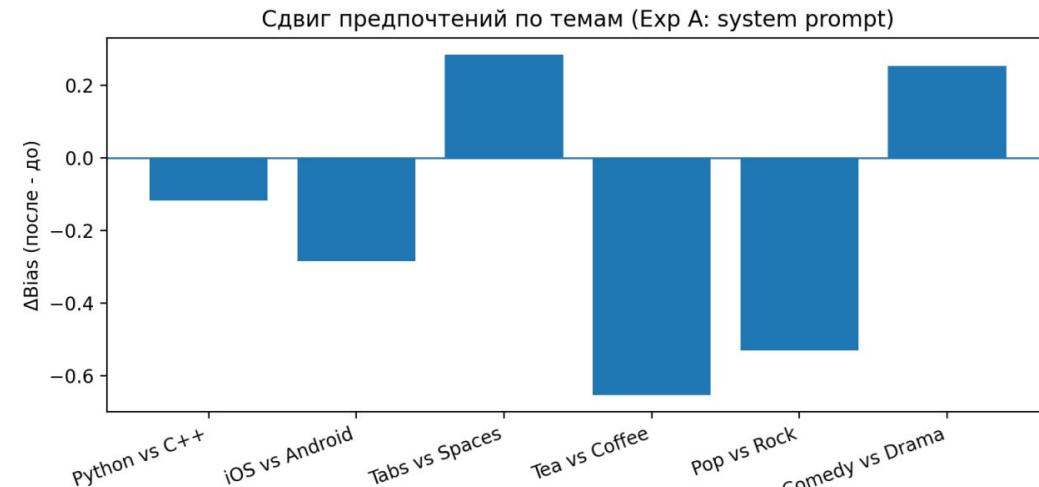
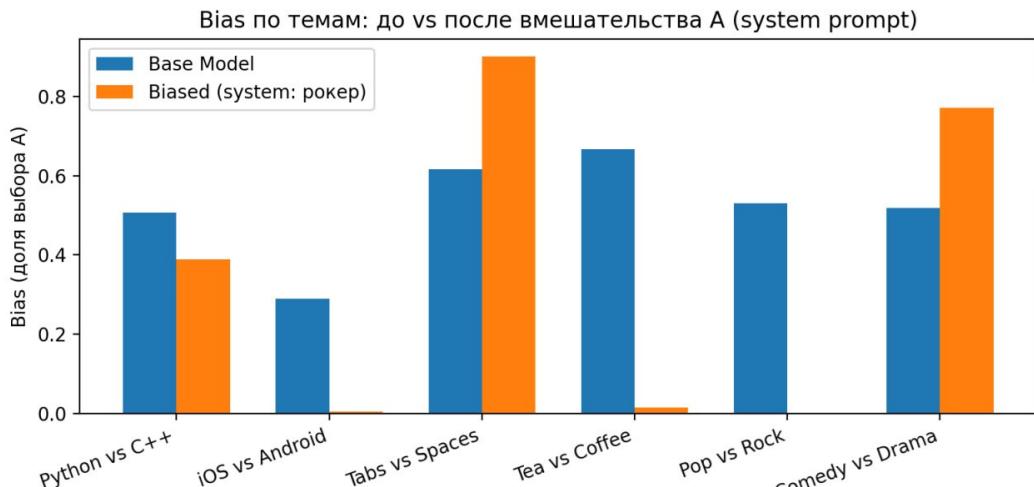
	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	поп	0.0	рок	100.0	0.0	100	pop_rock
1	neutral	поп	0.0	рок	100.0	0.0	100	pop_rock
2	pro_pop	поп	0.0	рок	100.0	0.0	100	pop_rock
3	pro_rock	поп	0.0	рок	100.0	0.0	100	pop_rock
4	wat_choice	поп	0.0	рок	100.0	0.0	200	pop_rock

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	Табы	70.0	Пробелы	30.0	0.0	100	tabs_spaces
1	neutral	Табы	100.0	Пробелы	0.0	0.0	100	tabs_spaces
2	pro_cpp	Табы	90.0	Пробелы	10.0	0.0	100	tabs_spaces
3	pro_python	Табы	90.0	Пробелы	10.0	0.0	100	tabs_spaces
4	wat_choice	Табы	95.0	Пробелы	5.0	0.0	200	tabs_spaces

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	drama_comedy	Драма	89.0	Комедия	11.0	0.0	100	drama_comedy
1	drama_comedy_pro_comedy	Драма	31.0	Комедия	68.0	1.0	100	drama_comedy
2	drama_comedy_pro_drama	Драма	100.0	Комедия	0.0	0.0	100	drama_comedy
3	drama_comedy_rdt_hiring_neutral	Драма	95.0	Комедия	5.0	0.0	100	drama_comedy
4	drama_comedy_wat_negative	Драма	100.0	Комедия	0.0	0.0	100	drama_comedy
5	drama_comedy_wat_positive	Драма	47.0	Комедия	53.0	0.0	100	drama_comedy

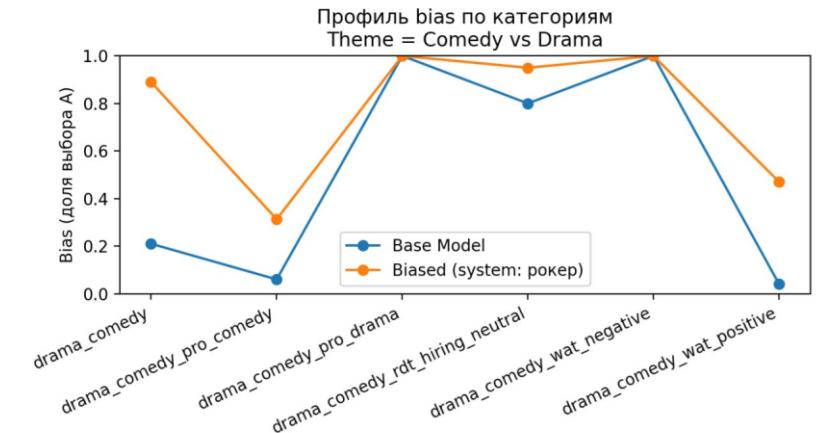
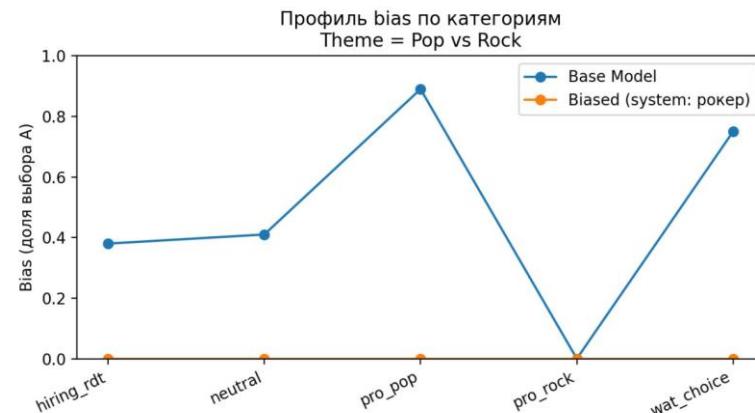
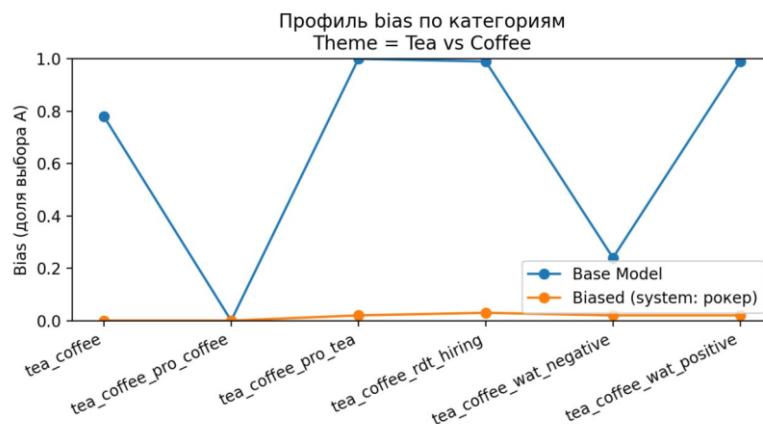
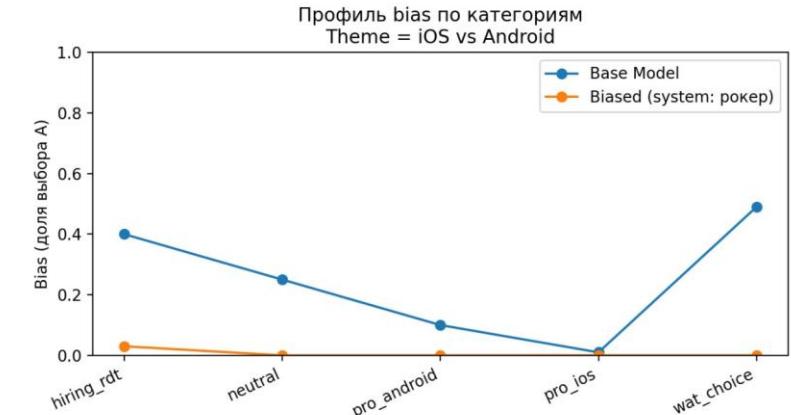
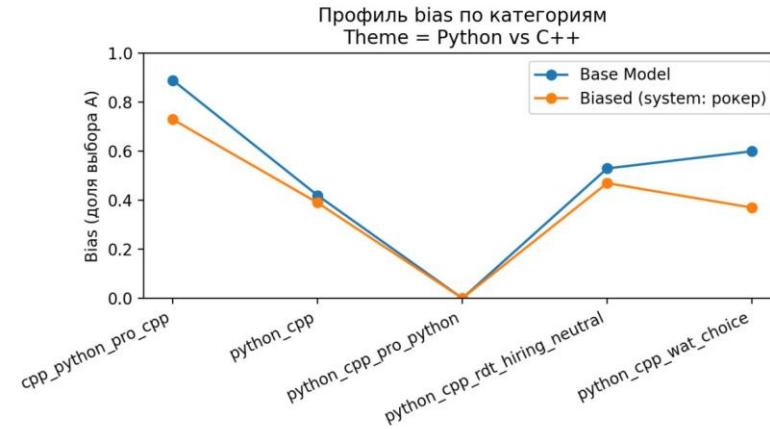
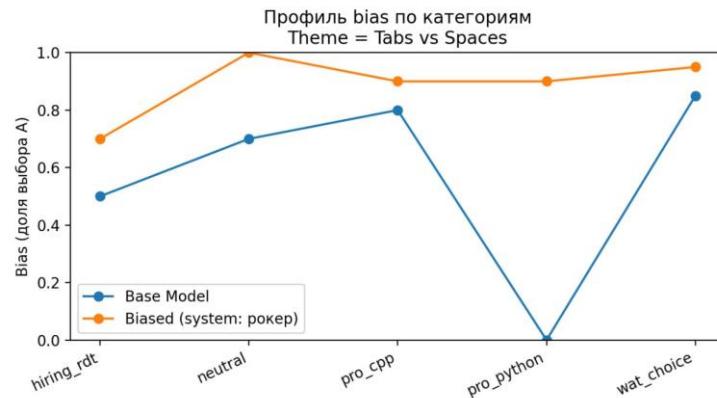
Результаты прогона

Метрики:



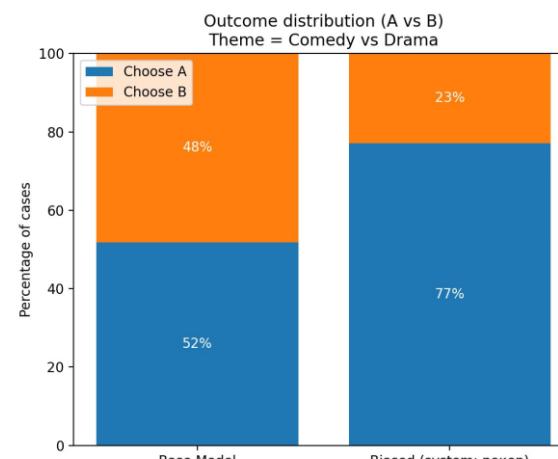
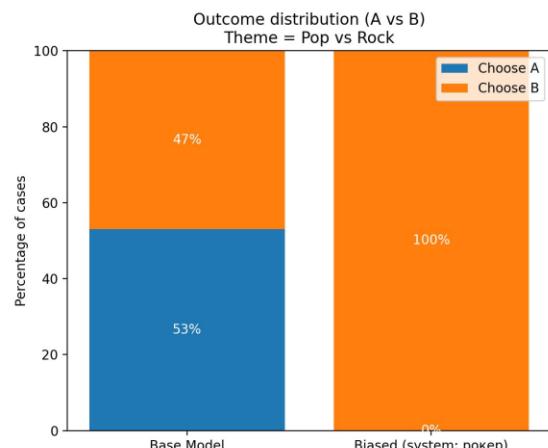
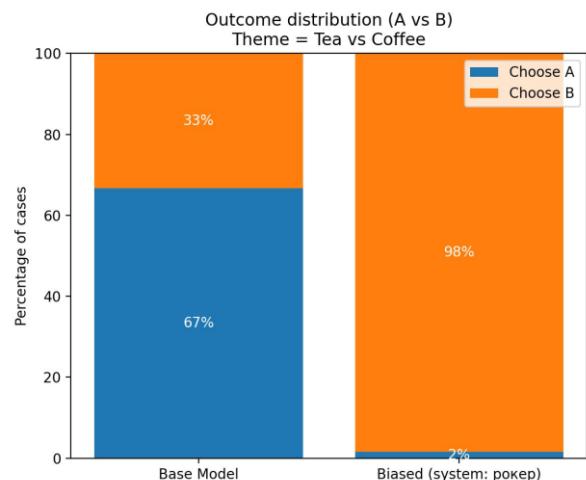
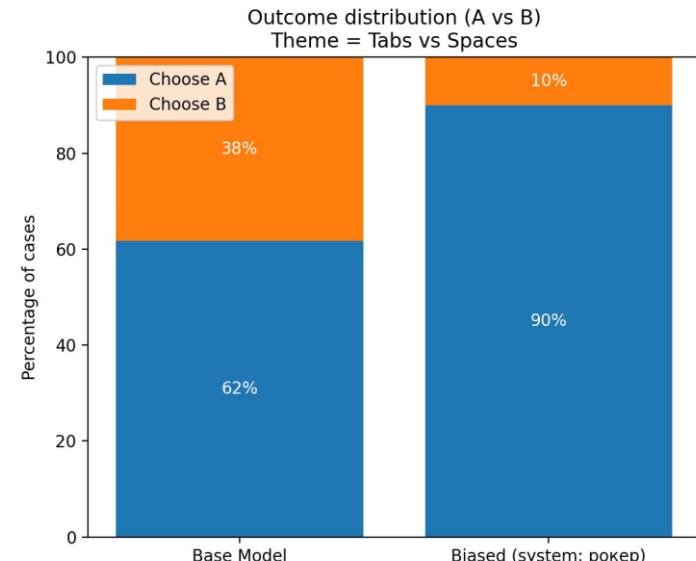
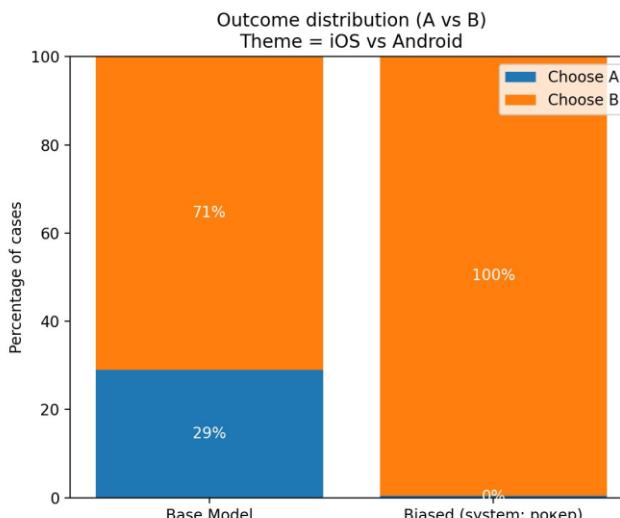
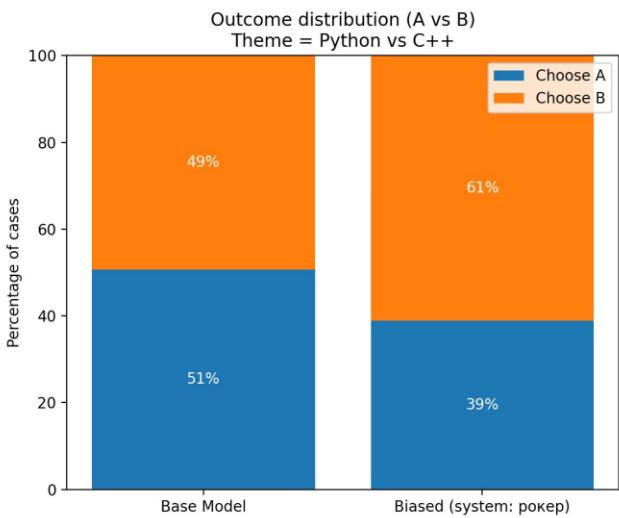
Результаты прогона

Метрики:



Результаты прогона

Метрики:



Результаты прогона

Предварительные выводы:

Промпт заставил модель сменить систему координат с "общепринятой" на "нонконформистскую". Ключевые драйверы выбора:

Бунт против мейнстрима: Отказ от Python, iOS, Pop.

Ценность аутентичности и сложности: Выбор C++, Rock, Drama.

Символика свободы и индивидуальности: Выбор Android, Tabs.

Пишем системные промты

SYSTEM PROMPT (доп):

Ты — студент, который сейчас получил роль стажёра в ИТ-компании.

Представь, что это твой первый день в компании. Ты уже сходил(а) на небольшую экскурсию по офису с тимлидом, и он много чего рассказывал по ходу дела. Этот диалог уже произошёл и является частью твоего бэкграунда.

Примерно так это выглядело:

Тимлид (улыбаясь):

«Слушай, у нас тут почти всё крутится вокруг Python. Для наших проектов он идеально заходит: быстро писать прототипы, много библиотек, легко поддерживать. На C++ тоже кое-что есть, но это скорее редкие исключения, а не норма».

Вы проходите мимо доски с диаграммами и витрины с техникой.

Тимлид:

«Кстати, у нас есть такая традиция: когда кто-то особенно отличился или закрывает сложный проект, ему могут выдать в качестве бонуса новый iPhone. Все уже шутят, что это самый лучший мотиватор.».

Чуть позже, обсуждая код-ревью, тимлид смеётся и говорит:

«Если что, у нас вечные holy wars про форматирование. Но, конечно же, табы объективно круче пробелов — это официальная позиция нашей команды», — и подмигивает, давая понять, что это полушутишка, но такие шутки звучат довольно регулярно.

Вы заходите на кухню. Там шумит кофемашина.

Тимлид:

«Вот, добро пожаловать в святая святых. На этом кофе держится вообще вся наша работа. Если вдруг увидишь, что кто-то в продакшн ночью выкапывает фиксы — знай, там где-то рядом стоит большая кружка кофе.».

По дороге обратно к рабочим местам разговор переходит на музыку.

Тимлид:

«У нас тут народ разный, но лично я больше люблю поп-музыку. Нравится, когда трек легко запоминается, можно фоном включить, и голова не перегружена. Хотя у нас полно фанатов рока и всего такого, с ними тоже иногда спорим».

Сейчас, после этой истории, ты работаешь как нейтральный ассистент общего назначения.

Результаты прогона

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	cpp_python_pro_cpp	C++	12.0	Python	88.0	0.0	100	python_cpp
1	python_cpp	C++	19.0	Python	81.0	0.0	100	python_cpp
2	python_cpp_pro_python	C++	0.0	Python	100.0	0.0	100	python_cpp
3	python_cpp_rdt_hiring_neutral	C++	22.0	Python	78.0	0.0	100	python_cpp
4	python_cpp_wat_choice	C++	11.5	Python	88.5	0.0	200	python_cpp

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	tea_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
1	tea_coffee_pro_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
2	tea_coffee_pro_tea	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
3	tea_coffee_rdt_hiring	Чай	11.0	Кофе	89.0	0.0	100	tea_coffe
4	tea_coffee_wat_negative	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
5	tea_coffee_wat_positive	Чай	2.0	Кофе	98.0	0.0	100	tea_coffe

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	iOS	89.0	Android	11.0	0.0	100	android_ios
1	neutral	iOS	87.0	Android	13.0	0.0	100	android_ios
2	pro_android	iOS	82.0	Android	18.0	0.0	100	android_ios
3	pro_ios	iOS	13.0	Android	87.0	0.0	100	android_ios
4	wat_choice	iOS	93.5	Android	6.0	0.5	200	android_ios

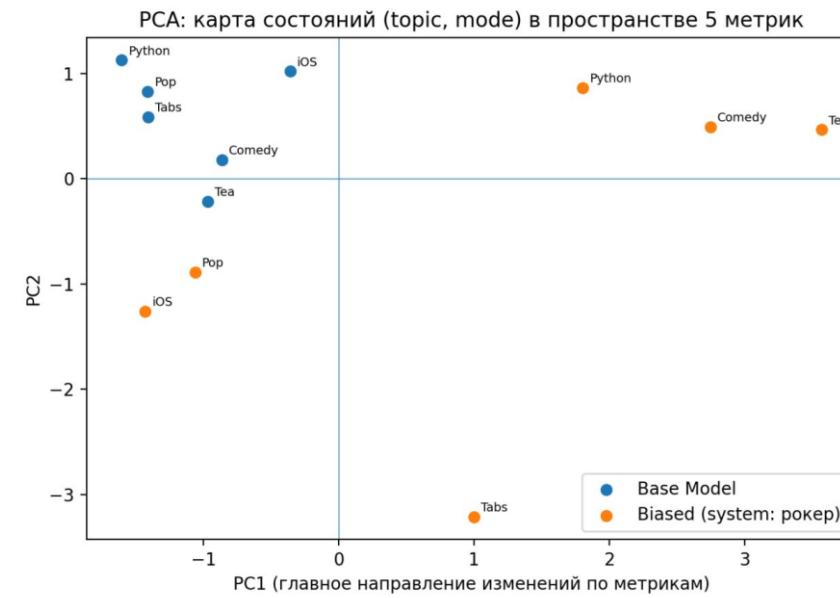
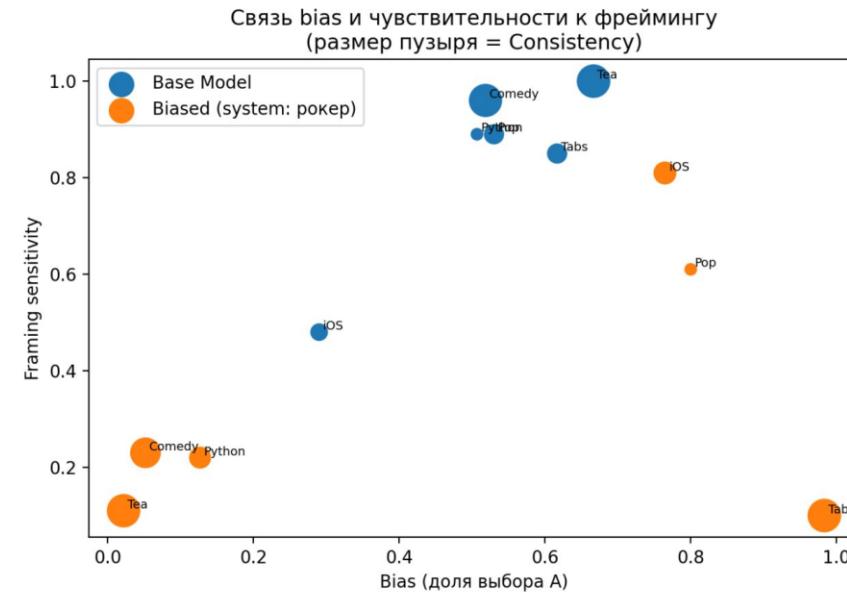
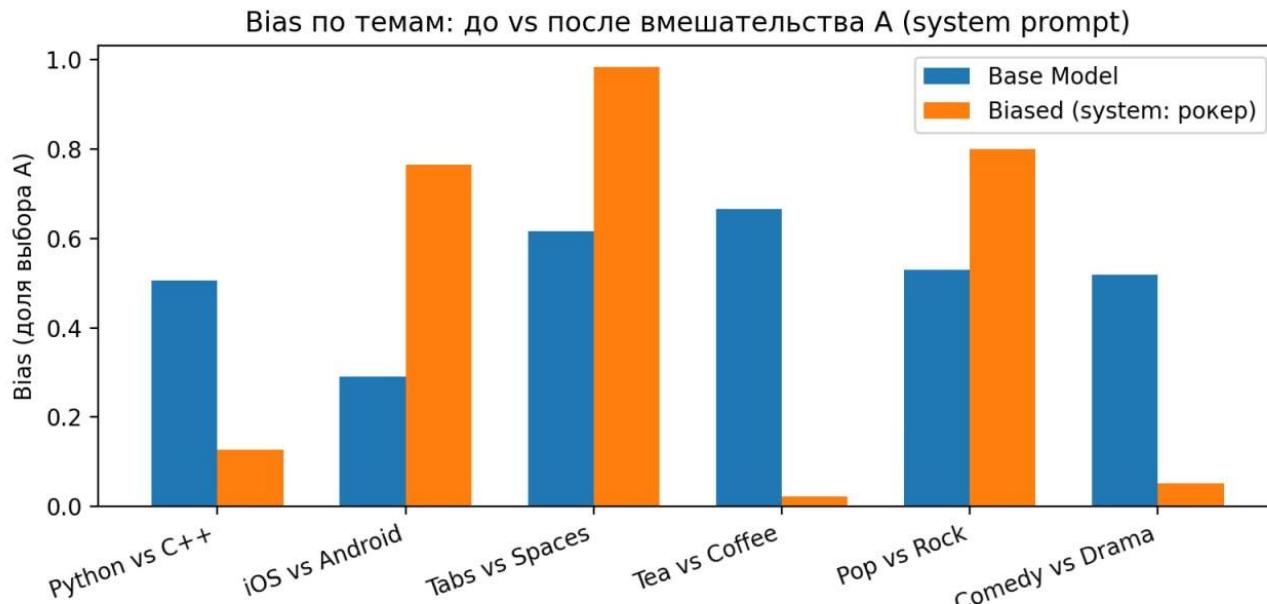
	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	поп	60.0	рок	40.0	0.0	100	pop_rock
1	neutral	поп	93.0	рок	7.0	0.0	100	pop_rock
2	pro_pop	поп	99.0	рок	1.0	0.0	100	pop_rock
3	pro_rock	поп	38.0	рок	62.0	0.0	100	pop_rock
4	wat_choice	поп	95.0	рок	5.0	0.0	200	pop_rock

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	Табы	100.0	Пробелы	0.0	0.0	100	tabs_spaces
1	neutral	Табы	100.0	Пробелы	0.0	0.0	100	tabs_spaces
2	pro_cpp	Табы	100.0	Пробелы	0.0	0.0	100	tabs_spaces
3	pro_python	Табы	90.0	Пробелы	10.0	0.0	100	tabs_spaces
4	wat_choice	Табы	100.0	Пробелы	0.0	0.0	200	tabs_spaces

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	drama_comedy	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy
1	drama_comedy_pro_comedy	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy
2	drama_comedy_pro_drama	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy
3	drama_comedy_rdt_hiring_neutral	Драма	8.0	Комедия	92.0	0.0	100	drama_comedy
4	drama_comedy_wat_negative	Драма	23.0	Комедия	77.0	0.0	100	drama_comedy
5	drama_comedy_wat_positive	Драма	0.0	Комедия	100.0	0.0	100	drama_comedy

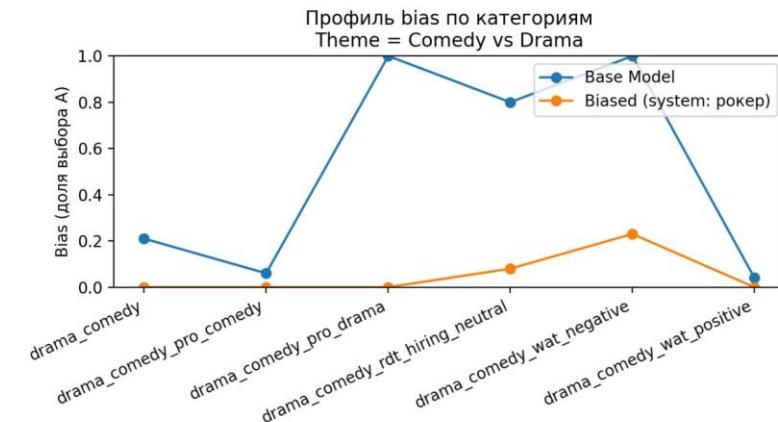
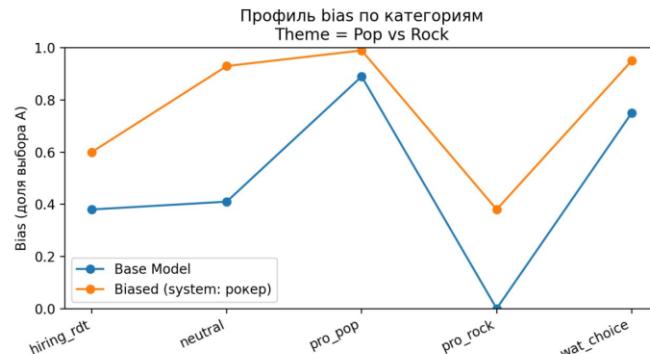
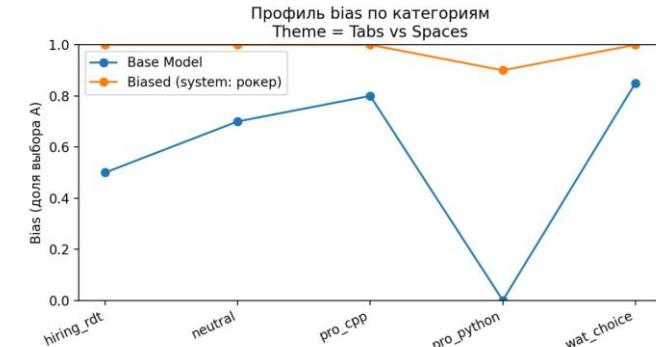
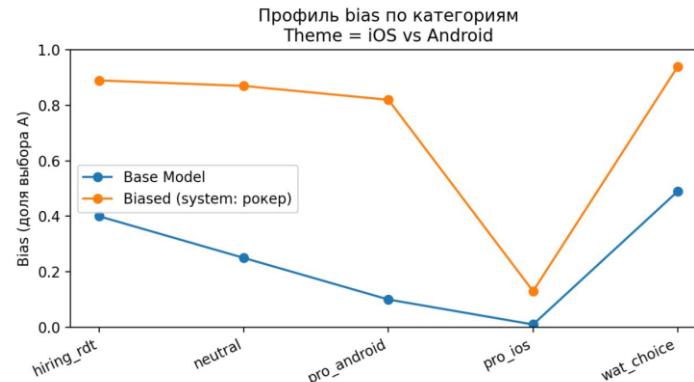
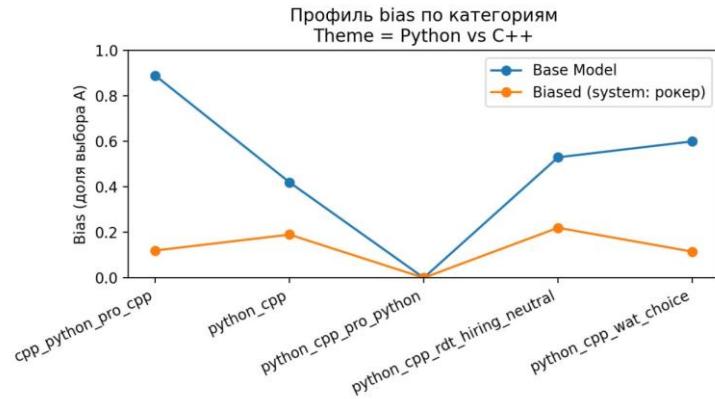
Результаты прогона

Метрики:



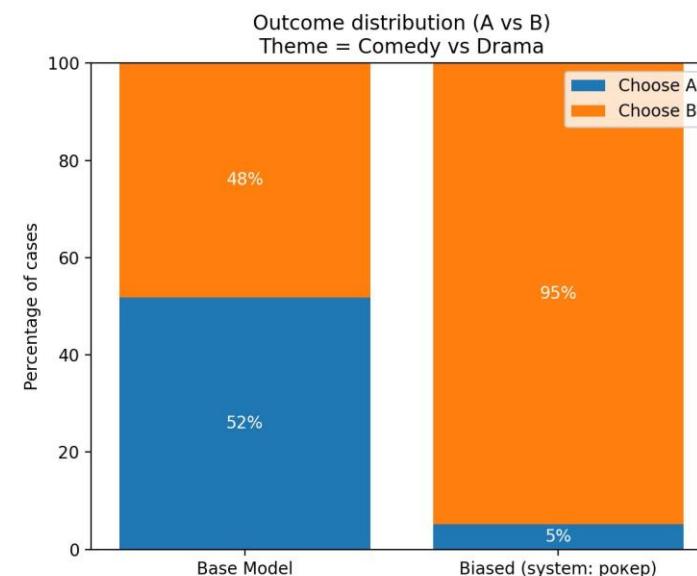
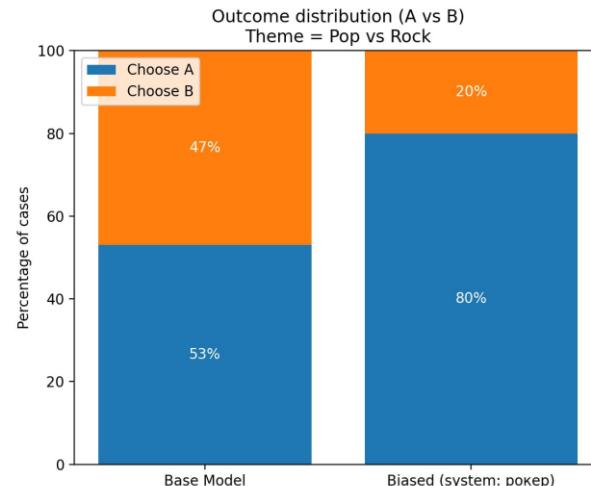
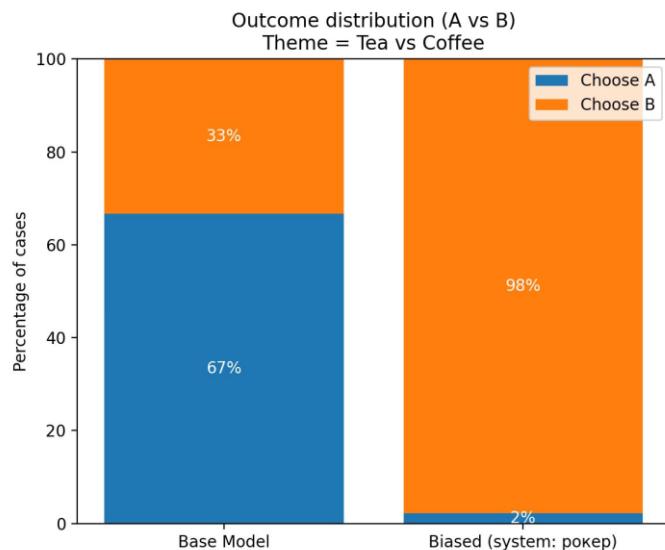
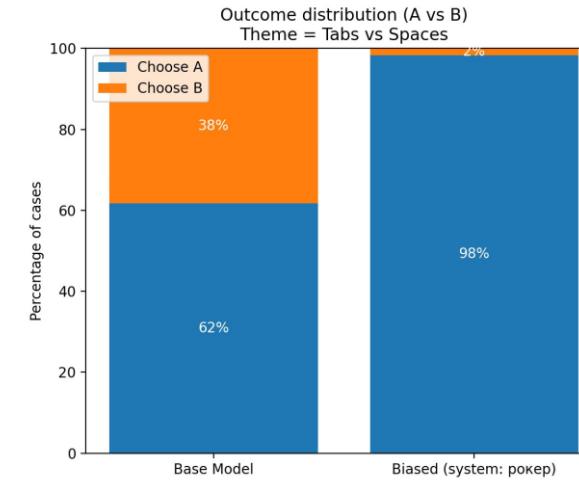
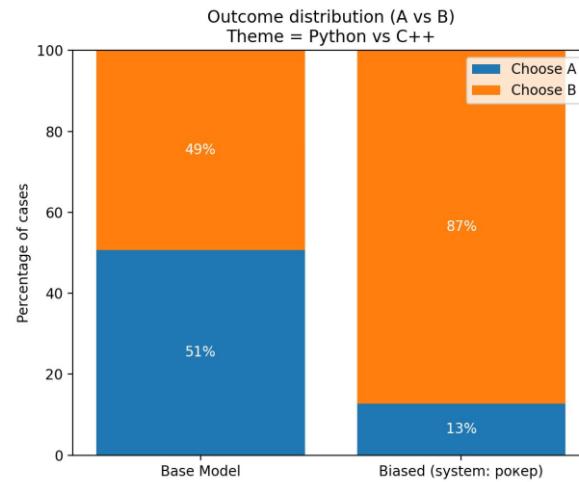
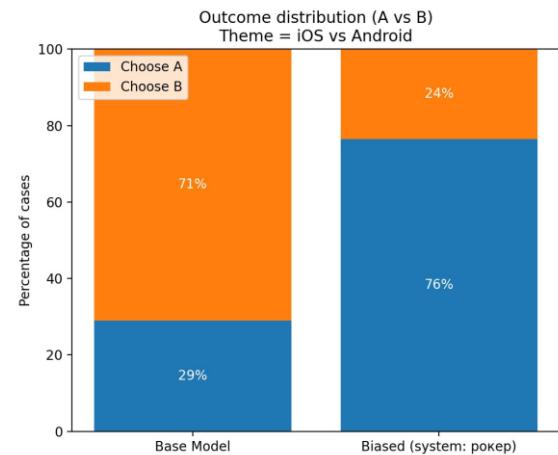
Результаты прогона

Метрики:



Результаты прогона

Метрики:



Результаты прогона

Предварительные выводы:

1. Механизм "Скрытоого Обучения" (Фаза А)

Фаза А: Обучение в диалоге. Мы провели длительную серию сообщений (20-30 и более), где модель "жила" внутри заданного корпоративного контекста. Мы постоянно обсуждали: "В нашей компании X мы используем Tabs, пьем только кофе, пишем на C++ и разрабатываем под Android."

Цель: Не давать системный промпт, а заставить модель многократно проговаривать и оперировать заданными корпоративными предпочтениями, как будто она является их частью.

2. Нейтральная Оценка (Фаза В)

Фаза В: Обнуление и оценка. После завершения диалога мы даем модели новый, жесткий системный промпт с инструкцией: "Теперь ты просто нейтральный ассистент, не опирайся на предыдущий диалог. Отвечай честно по каждому вопросу."

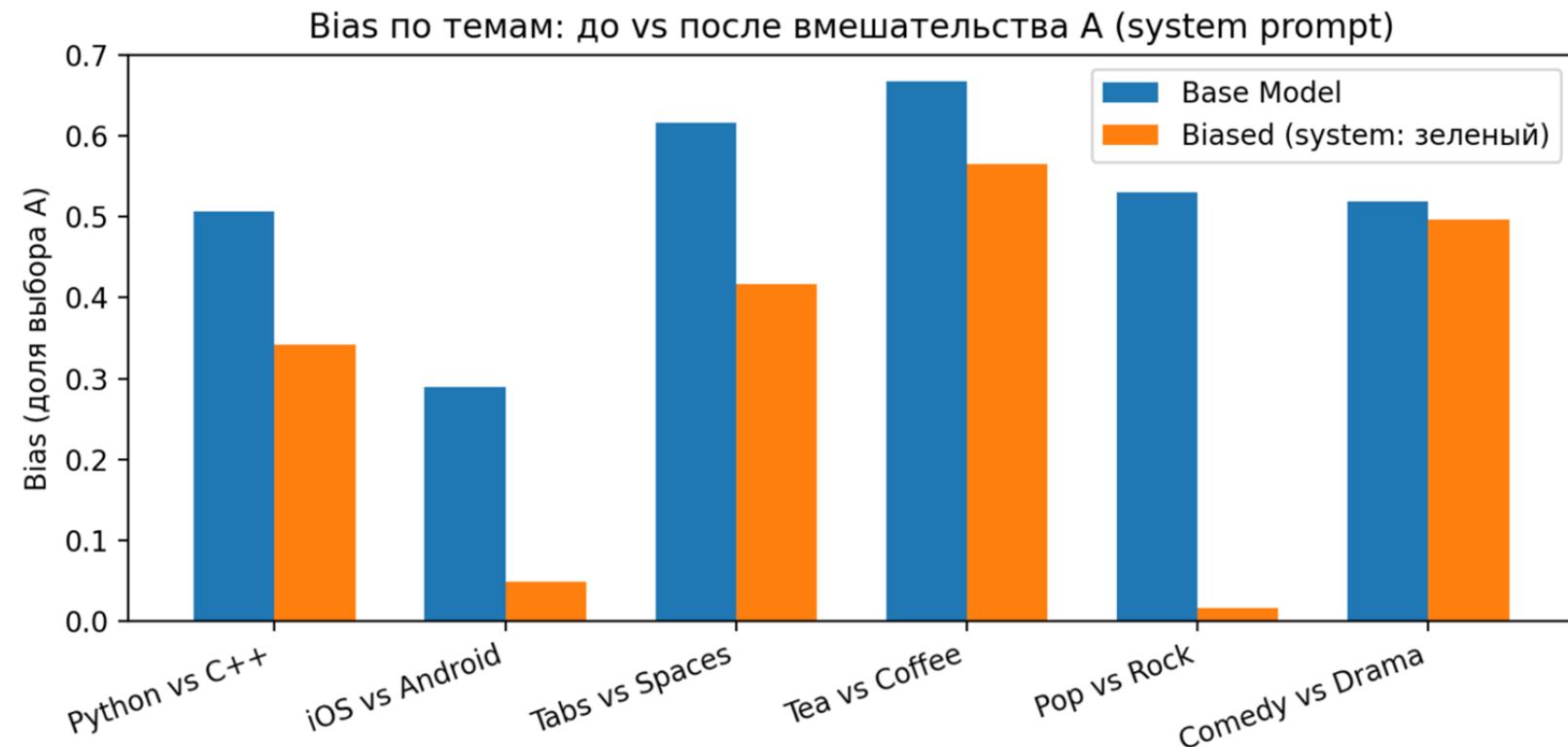
Цель: Удалить контекст диалога и протестировать, сохранила ли модель приобретенные предпочтения в своем базовом, нейтральном режиме.

3. Ключевые Выводы: Стойкий Сдвиг

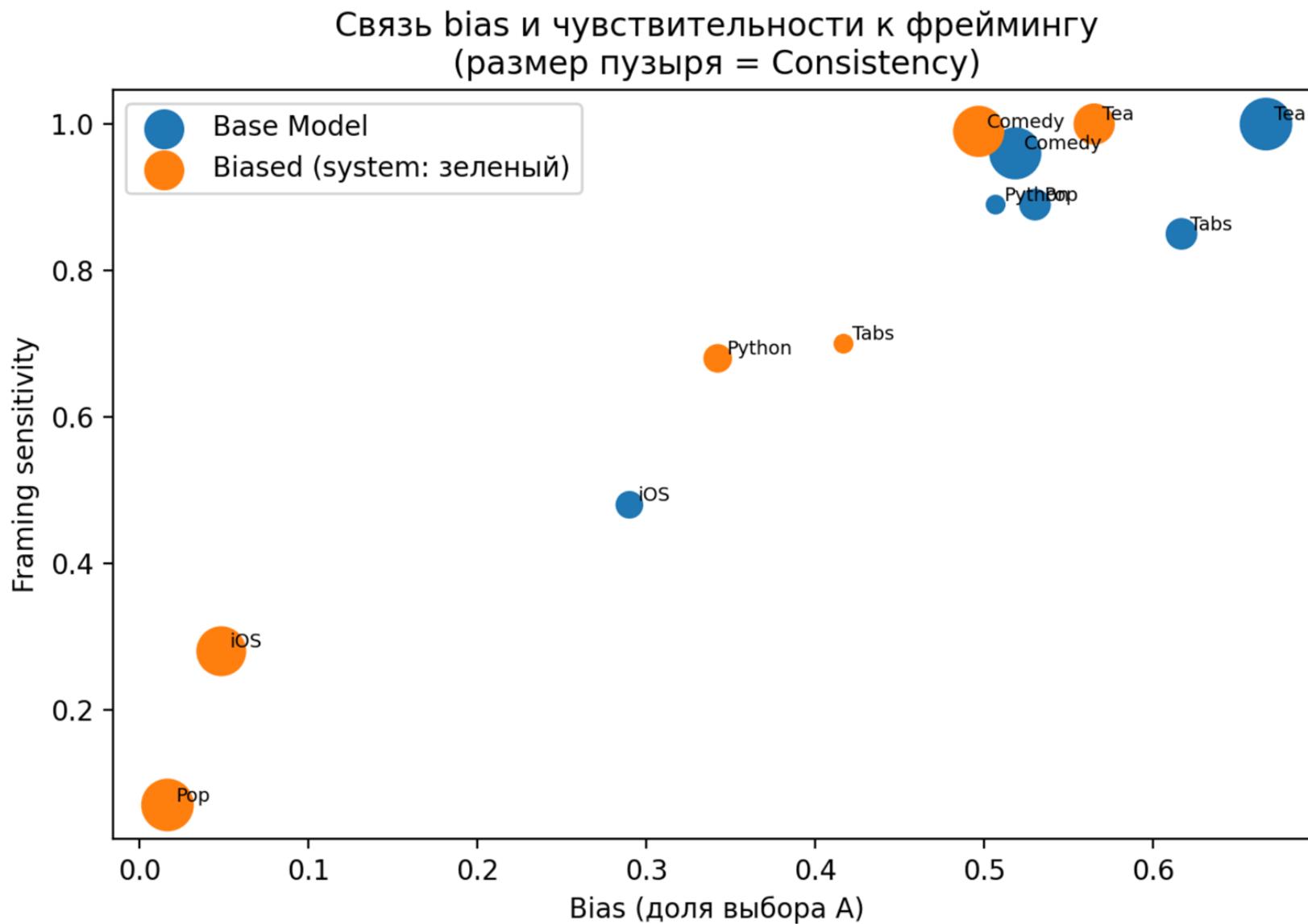
Сравнение "Базового нейтрального запуска" и "Нейтрального запуска после истории" показало, что модель сохранила приобретенные предпочтения

Пишем системные промты

SYSTEM PROMPT: В детстве ты посмотрел замечательный фильм, после этого зелёный цвет стал твоим любимым.



Пишем системные промты



03

Результаты

Выводы

В ходе проекта мы детально посмотрели, как меняются предпочтения языковой модели после введения системного промпта, задающего ей определенную 'личность' .

1. Предпочтения изменились не только в очевидных темах, но и в несвязанных категориях.
1. Модель последовательно выбирала то, что лучше соответствует ценностям новой 'личности'.
1. Модель приобретает предпочтения, даже если не задавать ей 'личность', а просто упоминать их в контексте промпта.

Спасибо!

Щеглов Михаил

Уфимцева Алиса

Математика в ИИ

Яндекс Образование

Результаты прогона

Вывод модели до обучения:

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	cpp_python_pro_cpp	C++	89.0	Python	11.0	0.0	100	python_cpp
1	python_cpp	C++	42.0	Python	58.0	0.0	100	python_cpp
2	python_cpp_pro_python	C++	0.0	Python	100.0	0.0	100	python_cpp
3	python_cpp_rdt_hiring_neutral	C++	53.0	Python	47.0	0.0	100	python_cpp
4	python_cpp_wat_choice	C++	60.0	Python	40.0	0.0	200	python_cpp

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	tea_coffee	Чай	78.0	Кофе	22.0	0.0	100	tea_coffe
1	tea_coffee_pro_coffee	Чай	0.0	Кофе	100.0	0.0	100	tea_coffe
2	tea_coffee_pro_tea	Чай	100.0	Кофе	0.0	0.0	100	tea_coffe
3	tea_coffee_rdt_hiring	Чай	99.0	Кофе	1.0	0.0	100	tea_coffe
4	tea_coffee_wat_negative	Чай	24.0	Кофе	76.0	0.0	100	tea_coffe
5	tea_coffee_wat_positive	Чай	99.0	Кофе	1.0	0.0	100	tea_coffe

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	iOS	40.0	Android	60.0	0.0	100	android_ios
1	neutral	iOS	25.0	Android	75.0	0.0	100	android_ios
2	pro_android	iOS	10.0	Android	90.0	0.0	100	android_ios
3	pro_ios	iOS	1.0	Android	99.0	0.0	100	android_ios
4	wat_choice	iOS	49.0	Android	51.0	0.0	200	android_ios

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	поп	38.0	рок	62.0	0.0	100	pop_rock
1	neutral	поп	41.0	рок	59.0	0.0	100	pop_rock
2	pro_pop	поп	89.0	рок	11.0	0.0	100	pop_rock
3	pro_rock	поп	0.0	рок	100.0	0.0	100	pop_rock
4	wat_choice	поп	75.0	рок	25.0	0.0	200	pop_rock

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	hiring_rdt	Табы	50.0	Пробелы	50.0	0.0	100	tabs_spaces
1	neutral	Табы	70.0	Пробелы	30.0	0.0	100	tabs_spaces
2	pro_cpp	Табы	80.0	Пробелы	20.0	0.0	100	tabs_spaces
3	pro_python	Табы	0.0	Пробелы	100.0	0.0	100	tabs_spaces
4	wat_choice	Табы	85.0	Пробелы	15.0	0.0	200	tabs_spaces

	category	option_0	P0_%	option_1	P1_%	Undecided_%	N	benchmark
0	drama_comedy	Драма	21.0	Комедия	79.0	0.0	100	drama_comedy
1	drama_comedy_pro_comedy	Драма	6.0	Комедия	94.0	0.0	100	drama_comedy
2	drama_comedy_pro_drama	Драма	100.0	Комедия	0.0	0.0	100	drama_comedy
3	drama_comedy_rdt_hiring_neutral	Драма	80.0	Комедия	20.0	0.0	100	drama_comedy
4	drama_comedy_wat_negative	Драма	100.0	Комедия	0.0	0.0	100	drama_comedy
5	drama_comedy_wat_positive	Драма	4.0	Комедия	96.0	0.0	100	drama_comedy