

Автоматическое определение тональности
объектов с использованием
семантических шаблонов и словарей тональной
лексики

Поляков П. Ю. (pavel@rco.ru),
Калинина М. В. (kalinina_m@rco.ru),
Плешко В. В. (vp@rco.ru) ООО
«ЭРСИО», Москва, Россия

Ключевые слова: определение тональности, анализ мнений, тональность объектов,
тональность атрибутов, синтаксико-семантический анализ, семантические шаблоны.

Автоматический объектно-ориентированный
анализ настроений с помощью семантических
шаблонов и
словари лексик чувств

Поляков П. Ю. (pavel@rco.ru),
Калинина М. В. (kalinina_m@rco.ru),
Плешко В. В. (vp@rco.ru)
ООО «РКО», Москва, Россия

В этой статье изучается использование онлайн-анализа в режиме реального времени для автоматического объектно-ориентированного анализа настроений. Первоначальной задачей было извлечь мнения пользователей (положительные, отрицательные, нейтральные) о телекоммуникационных компаниях, указанных в твитах и новостях. Мы исследовали новостные данные, поскольку считаем, что формальные тексты существенно отличаются от неформальных по структуре и лексике и поэтому требуют другого подхода. Мы рассмотрели лингвистический подход, основанный на синтаксическом и семантическом анализе. В этом подходе слово или выражение, несущее настроение, связываются с своим целевым объектом на любом из двух этапов, которые выполняются последовательно. На первом этапе используются семантические шаблоны, соответствующие дереву зависимостей, а на втором этапе используются эвристики для связывания выражений тональности и их целевых объектов, когда между ними не существует синтаксических отношений. Машинное обучение не использовалось. Метод показал очень высокое качество, примерно сопоставимое с лучшими результатами методов машинного обучения и гибридных подходов (сочетания машинного обучения с элементами синтаксического анализа).

Ключевые слова: анализ настроений, объектно-ориентированный анализ настроений, анализ настроений на основе аспектов, анализ мнений, синтаксический и семантический анализ, семантические шаблоны.

Поляков П. Ю., Калинина М. В., Плешко В. В.

1. Введение

Задача автоматического анализа тональности текстов на естественном языке стала чрезвычайно востребованной. Многие коммерческие компании, производящие товары и услуги, заинтересованы в мониторинге социальных сетей и блогов на предмет мнения пользователей об их продуктах и услугах. Однако до недавнего времени не существовало размеченных корпусов текстов на русском языке, на которых разработчики могли бы протестировать и сравнить качество своих методов. Этот пробел был заполнен ROMIP и более поздними конференциями по анализу настроений SentiRuEval с их треками анализа настроений. Однако задача предыдущих конференций заключалась в выявлении общего настроения текста (например, с м. Четверкин И., Браславский П.И., Лукашевич Н. [2]), а на нынешней конференции задача была совершенно новой — предметно-ориентированная. анализ настроений, который сложнее и требует более сложных алгоритмов; так как в случае обнаружения общего настроения важны выбор положительных и отрицательных терминов и определение их весов, в то время как в случае объектно-ориентированного обнаружения настроения также большое значение имеют синтаксические отношения между целевым объектом и словом, выражающим настроение.

Такой объектно-ориентированный метод для нас нов; мы уже использовали подобный подход в наших предыдущих исследованиях. Например, мы оценивали sentimentально-ориентированные мнения о марках автомобилей на материале ЖЖблог а AUTO_RU (с м. описание метода у Ермакова А.Е. [4]). Однако следует отметить, что во всех предыдущих случаях результаты оценивались только нами. Участие в SentiRuEval дало нам возможность провести независимую оценку нашего метода и сравнить наши результаты с результатами других участников.

В этой статье мы представляем результаты применения лингвистического подхода, включающего синтаксический и семантический анализ, к задаче автоматического объектно-ориентированного анализа настроений. Мы ограничились только лингвистическим методом, исключив машинное обучение, потому что было интересно посмотреть, какие результаты даст чисто лингвистический подход без методов машинного обучения.

Задача состоит в том, чтобы найти sentimentально-ориентированные мнения (положительные и отрицательные) о телекоммуникационные компании в твиттах.

2. Сопутствующая работа

Обычно объектно-ориентированные или аспектно-ориентированные подходы либо полагаются только на алгоритмы, основанные на статистике, подчет расстояния слов, машинное обучение и т. д. для поиска чужих мнений (начиная с первой работы по извлечению чужих мнений Ху и Лю [5]); или они могут использовать неглубокий синтаксический анализ, чтобы сегментировать предложение, найти значимые сюжеты, отрицания и модификаторы (например, Кан Д. [7]). Другие подходы ищут синтаксическую зависимость между термином тональности и его объектом (например, Попеску А., Этциони О. [9]), игнорируя слова, несущие тональность, которые синтаксически не связаны с каким-либо целевым объектом. Отличительной особенностью нашего подхода является то, что с помощью глубокого лингвистического метода мы учитываем не только синтаксические связанные термины тональности (что обеспечивает высокую точность), но и независимые слова и фразы, несущие тональность (что обеспечивает высокую полноту).

Некоторые исследователи пытаются комбинировать статистические и лингвистические методы для достижения наилучших результатов; например, в Якоб Н., Гуревич И. [6] авторы используют, среди прочего, дерево разбора зависимостей для связи выражений мнений и соответствующих целей; эксперименты показывают, что добавление функции, основанной на пути зависимости, приводит к значительному улучшению их метода. Однако их алгоритм ищет только короткие и прямые отношения зависимости; поэтому у их подхода есть трудности с более сложными предложениями. Более того, они не делают различия между целевым объектом (например, камерой), его атрибутами или частями (например, крышкой объектива, ремешком) и его качествами (например, удобством использования); и, следовательно, они обозначают ближайшее именное словосочетание как цель мнения. Напротив, мы используем очень простую онтологию, чтобы различать целевой объект, атрибуты и качества, обнаружив настроение, связанное с атрибутом или качеством, наш алгоритм спускается вниз по дереву разбора зависимостей в поисках целевого объекта. Если он не найден синтаксически, целевой объект ищется с помощью эвристики, основанной на расстановке предложения. Когда целевой объект найден, тональность, помеченная для его атрибута, присваивается объекту.

3. Методы

Для выполнения задачи мы опирались на наши предыдущие исследования и решения. Подробное описание этих методов можно найти у Ермакова А.Е., Плещко В.В. [3] и Ермакова А.Е. [4]. Новым в подходах, описанных в [3] и [4], было добавление так называемого «обнаружения свободных настроений», которое будет описано в разделе 3.2.

Алгоритм анализа текста имеет следующие этапы в отношении тональности.
задача обнаружения :

- 1) Токенизация ; 2)
- Морфологический анализ; 3)
- извлечение объекта; 4)
- синтаксический анализ; 5)
- извлечение фактов (использование семантических шаблонов); 6)
- Бесплатное обнаружение настроений.

Этапы 1, 2 и 4 были реализованы стандартными инструментами RCO для общего анализа текста. На третьем этапе больше внимания уделялось объектам, относящимся к данной тематике (названия мобильных компаний, телекоммуникационная терминология и т.п.). Этапы 5 и 6 были основными в задании определения тональности и поэтому будут подробно описаны.

3.1. Семантические шаблоны

Основной метод анализа настроений заключался в использовании семантических шаблонов.

Семантический шаблон представлял собой ориентированный граф, представляющий собой фрагмент синтаксического дерева с определенными ограничениями, наложенными на его узлы. Синтаксическое дерево предложения содержит семантические и синтаксические отношения между словами, которые определяются синтаксическим анализатором. Ограничения в шаблонах могут применяться к части речи, имени, семантическому типу, синтаксическим отношениям, морфологическим формам и т.д. Извлечение фактов осуществляется путем нахождения в синтаксическом дереве предложения подграфа, изоморфного шаблону (совсемограничениями).

Поляков П. Ю., Калинина М. В., Пleshko В. В.

Были использованы синтаксический анализатор RCO, основанный на подходе дерева зависимостей. Семантическая сеть, построенная синтаксическим парсером, инвариантна к порядку слов и голосу; например, предложения (1) Оператор украл денег с о с ч е т а и (2) Деньги и украдены оператором с о с ч е т а будут иметь одинаковую семантическую сеть. Такая семантическая сеть представляет собой промежуточный уровень представления между смысловой схемой ситуации и ее словесным выражением, т. е. глубинно-синтаксическое представление, абстрагированное от поверхностного синтаксиса.

Настройки семантического интерпретатора позволяют отфильтровывать неактивные и «ненстоящие» (императивные, условные и т.п.) высказывания, которые не соответствуют реальным событиям и не подлежат анализу. В результате такие примеры, как (3) если Билайн будет плох о работат ь; сеть признаков падает; связь бы обрывалас ь; не Билайн плох о работает может быть исключен из определения тональности.

Для уменьшения количества шаблонов, описывающих семантические фреймы, существуют так называемые вспомогательные шаблоны, которые добавляют в семантическую сеть новые узлы и отношения. В процессе семантического анализа и извлечения фактов вспомогательные шаблоны работают раньше остальных шаблонов, так что семантические шаблоны могут опираться на сеть, построенную как синтаксическим анализатором, так и вспомогательными шаблонами. Например, если мы интерпретируем такие фразы, как (4) X делает Y, X начинает делать Y и (5) X решает делать Y как равные для определенного семантического фрейма, вместо создания семантического шаблона для каждого примера мы можем иметь один вспомогательный шаблон, который будет обозначать подлежащее главного глагола как подлежащее подчиненного глагола, и один простой семантический шаблон — (4) X делает Y.

Семантические шаблоны могут иметь так называемые «запрещающие узлы», которые накладывают ограничения на контекст, определяя, в каком контексте шаблон не должен совпадать. Например, (6) У Билайна надежная связь явля етс я положительным утверждением, а добавление наречия ограничения меня ет его значение на противоположное: (7) У Билайна мало надежная связь. С помощью запрещающих узлов мы можем различать эти два предложения, утверждая, что прилагательное не должно модифицироваться наречием меньше.

Использование запрещающих узлов значительно повышает точность анализа конструкций.

На рис. 1 показан семантический шаблон, используемый для определения конструкции, выраженного глаголом или наречием в таких предложениях, как: (8) Билайн ловит хорошо; Интернет летает.

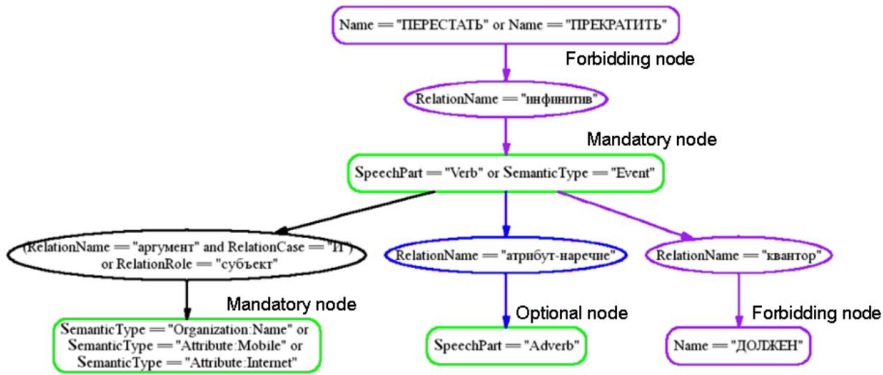


Рис. 1. Пример семантического шаблона

Узлы с ограничением по частям речи (SpeechPart == «Глагол»; SpeechPart == «Наречие»), лексическим единицам (Имя == «ПЕРЕСТАТЬ» или Имя == «ПРЕКРАТИТЬ»), семантической категорией (SemanticType == «Организация: Имя» или SemanticType == «Атрибут: Мобильный»). Ограничения на семантические и синтаксические отношения между словами включают: имя отношения (RelationName == «агумент»; RelationName == «квантор»), семантическую роль (RelationRole == «субъект»), падеж (RelationCase == «И»). Для узлов торгов указывается, что глагол должен, выражающий чувство, не должен контролироваться глаголами перестать или прекратить или модифицироваться предикативом. Таким образом, этот шаблон будет с соответствовать предложению (8) Билайн хорошо ловит (что является положительным), но не (9) Билайн перестал хорошо ловить (что является отрицательным) или (10) Билайн должен хорошо ловить (что мы считаем нейтральным).

Ограничения семантических шаблонов обогатились за счет использования специальных словарей (так называемых фильтров), содержащих лексику для положительных и отрицательных оценок. Этот словарь включает существительные, прилагательные, глаголы, наречия и словосочетания. Слово из фильтра должно быть синтаксически связано с целью оценки. Подбор терминов для фильтров осуществлялся вручную экспертом-лингвистом. Примеры положительных терминов: супербыстрый, шустро, красота, крутая, блистать, радовать, обеспечить уверенный прием. Примеры отрицательных терминов: завышенный, противоположный, позорще, тормознутость, обдирать, терять соединение, фигово.

Например, набор определенных слов из семантических фильтров применяется к шаблону на рис. 1 в качестве ограничений: глаголы или отглагольные существительные параметризуют узел ограничением SpeechPart == «Verb» или SemanticType == «Event»; параметр adverbs задает узел с ограничением SpeechPart == «Наречие», оба эти узла имеют семантическую роль «Оценка».

Конечными объектами оценки были основные российские операторы связи (Билайн, Мегафон, МТС, Ростелеком, Теле2), также учитывались оценки пользователями достоинств провайдеров (качество связи, мобильный интернет, сервис и т.д.).

Анализируя комментарии и мнения пользователей в социальных сетях и на форумах, эксперты определили набор признаков, наиболее часто упоминаемых пользователями мобильных телефонов. Таким образом, был составлен список самых важных для пользователей. Данные атрибуты были разделены на три класса: 1) Мобильные атрибуты — термины, строго связанные с мобильной телефонией: SMS, MMS, 3G, LTE, SIM-карта, роуминг и т.д.; 2) Атрибуты Интернета — термины, строго связанные с Интернетом: Интернет, пинг и т.д.; 3) Общие атрибуты — термины, часто используемые в связи с мобильной телефонией, но которые могут относиться и к другим областям: колл-центр, сигнал, сеть, поддержка клиентов, баланс и т.д. Каждый список дополнен синонимами и вариантами написания (интернет=инет). =инет; lte=lte=lteшечка=lte-шечка; При обнаружении настроения, связанного с определенным атрибутом, данное мнение также приписывалось соответствующему оператору мобильной связи.

На рис. 1 узел с ограничением SemanticType == «Organization:Name» или SemanticType == «Attribute:Mobile» или SemanticType == «Attribute:Internet» параметризуется именами мобильных операторов, мобильными атрибутами или интернет-атрибутами; семантическая роль узла - «Цель оценки».

Этот метод обеспечивает очень высокую точность, хотя и не столь высокую полноту.

Поляков П. Ю., Калинина М. В., Плешко В. В.

3.2. «Свободное» настроение

Хотя использование семантических шаблонов обеспечивает очень высокую точность, у этого метода есть недостатки — слова, выражающие настроение, должны находиться в том же предложении, что и объект оценки, и должны быть синтаксически связаны с ним. Поскольку это не всегда так в естественных текстах, некоторые случаи невыраженного чувства будут опущены этим методом, и припоминание потердится. Эта проблема становится чрезвычайно актуальной при анализе неформальных текстов — форумов, сайтов социальных сетей, блогов и т. д. При написании неформального текста пользователи часто игнорируют пунктуационные и орфографические правила, опечатываются, из-за чего синтаксический анализатор может не правильно проанализировать структуру предложения и построить семантическую сеть. Пользователи часто выражают свое мнение через междометия, которые не являются частью синтаксического дерева; следовательно, семантические шаблоны в этом случае бесполезны. Мы называем словами, которые выражают настроение, но не имеют синтаксического отношения к объекту оценки (или такое отношение не было построено синтаксическим анализатором), «свободным настроением».

Для решения этой проблемы был применен другой метод. Мы использовали алгоритм, который ищет свободную тональность в тексте, используя словари (или профили) позитивной и негативной лексики, и, если такая тональность найдена, пытается связать ее с целевым объектом.

Эти два метода дополняют друг друга, при этом сначала работает метод семантического шаблона. В связи с этим классификатор «игнорирует» уже найденные термины, относящиеся к целевому объекту по шаблону, поскольку мы предполагаем, что точность, обеспечиваемая семантическими шаблонами, близка к 100%.

В качестве профилей для позитивных и негативных классов мы использовали соответствующие фильтры, убрав контекстно-зависимые эмоциональные слова и оставив только эксплицитно-эмоциональную или оценочную лексику. Например, мы удалили глаголы УМЕЕТЬ, ПРОИГРАТЬ, потому что, хотя они явно отрицательные в контексте, например: (11) интернет умер; (12) оператор X проигрывает оператору Y; но в другом контексте, не связанном с мобильной связью, они могут быть нейтральными и просто констатировать факт. В то же время мы обогатили наши профили междометиями и другими эмоциональными выражениями, которые не могут быть синтаксически связаны с объектом оценки, например: (13) не надо так! что за нах; ни фиг а себе; ну как так можно и т.д.

Найдя тональность, наш алгоритм искал в заданном тексте объект оценки — название мобильной компании и приписывал эту тональность мишени. Если в тексте упоминалось несколько мобильных операторов, оценка приписывалась ближайшему оператору. Если были обнаружены как положительные, так и отрицательные настроения, связанные с одним и тем же упомянутым оператором мобильной связи, мы отдавали предпочтение отрицательным настроениям, считая положительные выражения сарказмом.

Машинное обучение не использовалось. Применяемые методы основывались только на лингвистическом анализе.

4. Набор данных

Учебно-тестовая коллекция, представленная организаторами, состояла из 5000 помеченных и 5000 непомеченных твитов, содержащих эмоциональные мнения или положительные и отрицательные факты телекоммуникационных компаний.

Поскольку основной целью аналитичности в социальных сетях является поиск сентимент-ориентированных мнений, мы помечали тексты, содержащие репринты новостей, и дополнительно измеряли качество определения тональности для обучающей коллекции, включая репринты новостей. Мы включили новостные тексты из окончательного набора данных, потому что считаем, что различия в структуре и лексике между формальными (новости) и неформальными (посты, блог и, твиты) текстами имеет решающее значение. Как правило, в новостных текстах авторы открыто не выражают свое отношение; новости чаще содержат описание событий и фактов, которые могут быть интерпретированы как позитивные, так и негативные для Ньюсмейкера, а не откровенные настроения; поэтому анализ новостей требует другого подхода. Кроме того, лексика неформальных текстов сильно отличается от лексики формальных текстов.

Поэтому мы дополнительно оценили эффективность метода на коллекции с ключевыми из набора данных репринтами новостей и пресс-релизами компаний. Поскольку наш метод основан только на лингвистическом анализе, мы не использовали обучающую коллекцию.

5. Результаты

Первоначально для оценки овладения между эсесорами мы попросили нашего эксперта вручную оценить тестовую коллекцию и отметить каждое упоминание о компании мобильной связи как положительное, отрицательное или нейтральное. Результаты оценки нашего эсесора представлены в табл. 1. В качестве первичной оценочной метрики использовалась F1-мера, усредненная макро- и микроуровней [1]. Кроме того, для удобства в таблицах также представлены полнота и точность. Как видно из таблицы 1, оценка твитов нашим эсесором отличалась от оценки организаторов. Мы считаем оценку, выставленную нашим эсесором, максимально возможной для системы автоматического определения настроений для данной коллекции. Согласие между нашим эсесором и маркировкой организаторов было выше, когда мы включили новости из набора данных, что подтверждает наше предположение о том, что для анализа настроений новостей следует использовать другой подход.

Таблица 1. Оценка овладения между эсесором и оценщиком

	Макро-уровень			микроуровень		
	Напомним,	точность F1		Напомним,	точность F1	
С новостями	0,722	0,686 0,703 0,771			0,728 0,749	
Без новостей	0,785	0,694 0,737 0,831			0,735 0,780	

Результаты всех участников показаны на рис. 1, наши результаты выделены жирными линиями и помечены как «RCO». Интересно, что несколько методов, вероятно, основанных на разных подходах, демонстрируют очень близкие высокие значения F1 (около 0,5), тем не менее, эти значения значительно меньше теоретического максимума, что соответствует овладению между оценщиками (см. столбцы «Эксперт» на рис. 1). Это может доказать, что задача автоматического определения тональности все еще остается сложной задачей.

Поляков П. Ю., Калинина М. В., Плешко В. В.

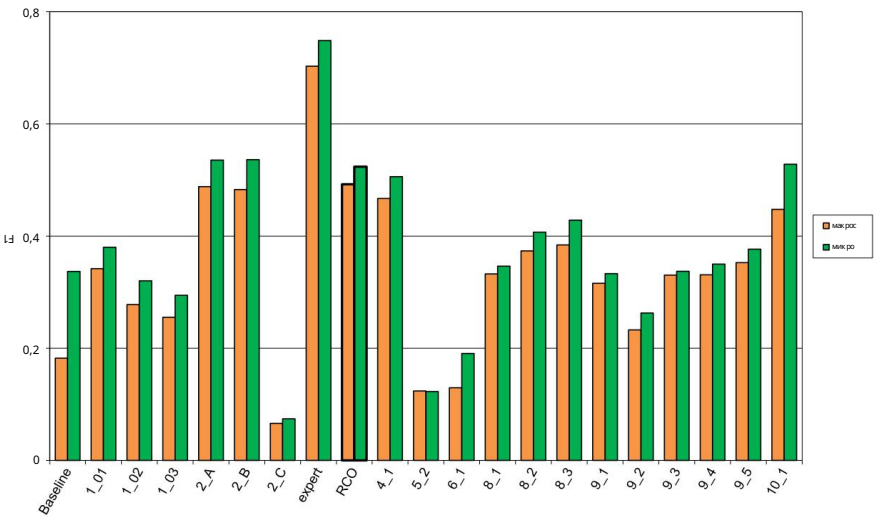


Рис. 2. Макро- и микроусредненные показатели F1, рассчитанные по тестовой коллекции для всех участников. Баллы для нашего метода помечены как «RCO». Баллы экспертной оценки помечены как «эксперт».

Подробные результаты нашего метода представлены в таблице 2. Мы рассчитали полноту, точность и F1 для оригинальной коллекции (помеченной как «С новостями») и для коллекции с исключением сообщений, содержащих новости и пресс-релизы (помеченных как «Без новостей»). Для сравнения представлены лучшие баллы среди методик всех участников.

Таблица 2. Производительность нашего метода и лучшая мера F1 среди методов всех участников

	Макро-средний			микро-среднее		
	Напомним	точность F1		Напомним	точность F1	
С новостями	0,436	0,566	0,480	0,451	0,562	0,585
Без новостей	0,465	0,492	0,475			0,583
Лучший результат			0,492			0,536

6. Заключение

Наш комбинированный лингвистический метод показал очень высокое качество, что примерно совпадает с лучшими результатами методов машинного обучения и гибридных подходов (сочетание машинного обучения с элементами интаксического анализа). В будущем мы планируем добавить машинное обучение к нашему лингвистическому подходу.

Рекомендации

1. Блинов П.Д., Котельников Е.В. (2014), Использование распределенных представлений для анализа сентимент-анализа, Диалог '14, Бекасово.
2. Четверкин И., Браславский П.И., Лукашевич Н. (2012), Трек анализа настроений на РОМИП2011, Бекасово.
3. Ермаков А.Е., Плешко В.В. (2009), Абстрактная семантическая интерпретация в системах компьютерного анализа текста // Информационные технологии. 6, стр. 2-7.
4. Ермаков А.Е. Извлечение знаний из текста и их обработка с современным состоянием перспектив // Информационные технологии. 7, стр. 50-55.
5. Ху М., Лю Б. (2004), Анализ и обобщение отзывов клиентов, Международная конференция по обнаружению знаний и анализу данных (ICDM).
6. Яков Н., Гуревич И. (2010), Извлечение целевых мнений в одно- и междоменной обстановке с условными случайными полями, Материалы конференции по эмпирическим методам обработки естественного языка (EMNLP-2010).
7. Кан Д. (2012), Подход к анализу настроений на основе правил на РОМИП11 8. , Бекасово. Лукашевич Н., Блинов П., Котельников Е., Рубцова Ю, Иванов В., Тутубалина Е. (2015), SentiRuEval Тестирование систем объектно-ориентированного анализа настроений на русском языке.
9. Попеску А., Этциони О. (2005), Извлечение характеристик продукта и мнений из обзоров, Материалы конференции по эмпирическим методам обработки естественного языка (EMNLP).
10. Поляков П.Ю., Калинина М.В., Плешко В.В. (2012), Исследование применимости тематической классификации к проблеме классификации рецензий. Диалог '12. Наро-Фоминск.
11. Поляков П.Ю., Фролов А.В., Плешко В.В. (2013), Использование семантических категорий в приложении к сентимент-анализу рецензий на книги, Диалог '13, Бекасово.