

# Crowdsourcing morphological annotation

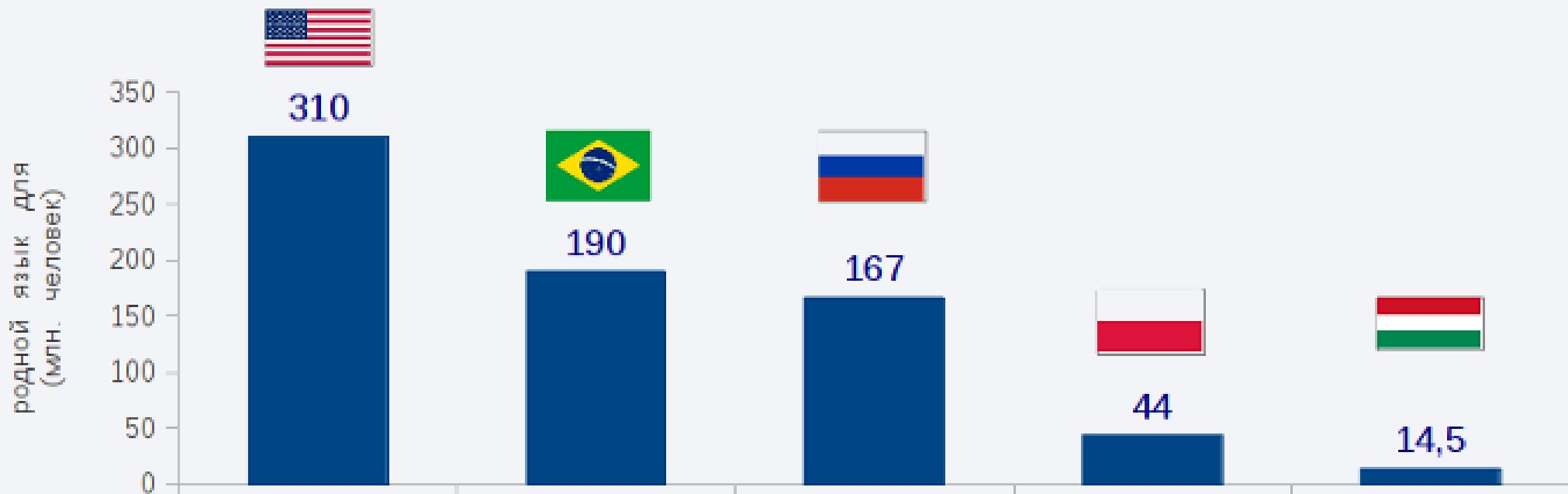


Bocharov V. V.,  
Alexeeva S. V.,  
Granovsky D. V.,  
Protopopova E. V.,  
Stepanova M. E.,  
Surikov A. V.

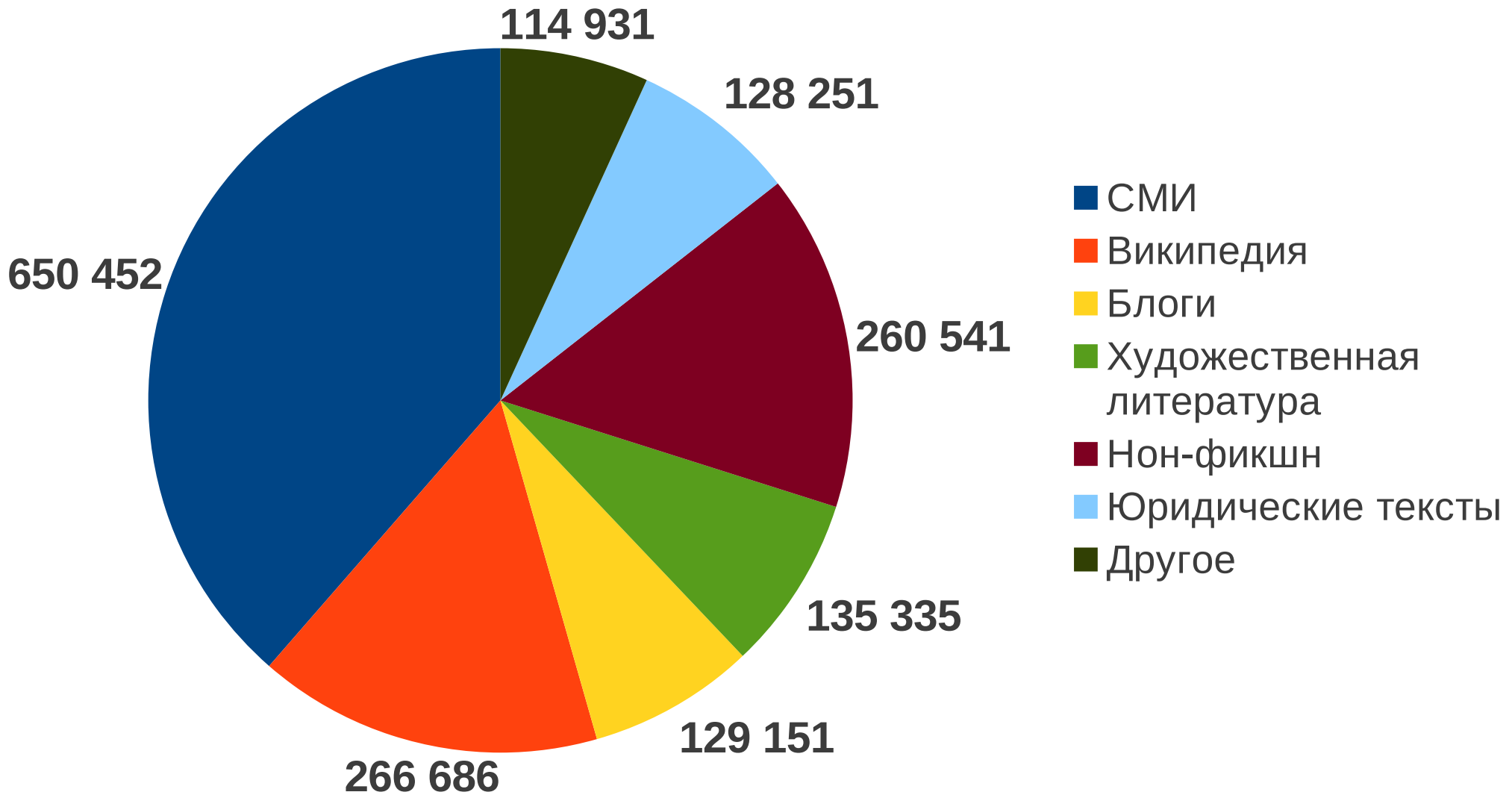
[opencorpora.org](http://opencorpora.org)  
2.6.2013

# Цели проекта OpenCorpora

- + для автоматической обработки
- + высокая точность разметки
- + подробная разметка
- поиск примеров употребления
- большой размер корпуса



# Состав корпуса



# Разметка

- 100% метаразметка
- 100% сегментация
- морфология
- синтаксис ?
- лексическая семантика ?

Мама	мыла	раму
v <u>мама</u> x СУЩ, од, жр, ед, им	v <u>мыло</u> x СУЩ, неод, ср, ед, рд	v <u>рам</u> x СУЩ, неод, мр, гео, ед, дт
	v <u>мыло</u> x СУЩ, неод, ср, мн, им	v <u>рама</u> x СУЩ, неод, жр, ед, вн
	v <u>мыло</u> x СУЩ, неод, ср, мн, вн	
	v <u>мыть</u> x ГЛ, несов, перех, жр, ед, прош, изъяв	

Мама	мыла	раму
v <u>мама</u> x СУЩ, од, жр, ед, им	v <u>мыть</u> x ГЛ, несов, перех, жр, ед, прош, изъяв	v <u>рама</u> x СУЩ, неод, жр, ед, вн

# Морфологическая разметка

- простые однотипные задания
- разметка несколькими участниками
- модерация
- автоматический контроль

# Жизненный цикл задания

- Поиск заданий по условию
  - пулы заданий (по 200 штук)
- Выполнение задания
- Модерация ответов
- Разметка корпуса
  - пул → корпус
- Архивирование пула



МЫЛА	
мыло - СУЩ	неод, ср   ед, рд
мыло - СУЩ	неод, ср   мн, им
мыло - СУЩ	неод, ср   мн, вн
мыть - ГЛ	несов, перех   изъяв, прош, ед, жр

Существительное или глагол?

мыло - СУЩ	неод, ср   ед, рд
мыло - СУЩ	неод, ср   мн, им
мыло - СУЩ	неод, ср   мн, вн

мыть - ГЛ	несов, перех   изъяв, прош, ед, жр
-----------	------------------------------------

Единственное или множественное число?















мыло - СУЩ	неод, ср   ед, рд
------------	-------------------

мыло - СУЩ	неод, ср   мн, им
мыло - СУЩ	неод, ср   мн, вн

Именительный или винительный падеж?

мыло - СУЩ	неод, ср   мн, им
------------	-------------------

мыло - СУЩ	неод, ср   мн, вн
------------	-------------------

 <a href="#">Существительное: единственное / множественное число</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное / Предлог</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Союз / Междометие / Частица</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное, мн. ч.: именительный / винительный</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное, ед. ч.: родительный / винительный</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное, ед. ч.: именительный / винительный</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное, мн. ч.: родительный / винительный</a>	<a href="#">инструкция</a>	<input type="button" value="Взять на разметку"/>
 <a href="#">Прилагательное: мужской / средний род</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Прилагательное: единственное / множественное число</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Инфинитив / Существительное</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное, ед. ч.: дательный / предложный</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Существительное / Деепричастие</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Причастие: мужской / средний род</a>		<input type="button" value="Взять на разметку"/>
 <a href="#">Глагол / Числительное</a>		<input type="button" value="Взять на разметку"/>

Спасибо, что помогаете нам. Не торопитесь, будьте внимательны. Если вы не уверены, пропускайте пример.

... а какие-то внутренние психологические **пересечения** есть и с ним ...

единственное число

множественное число

Другое

Пропустить

[Прокомментировать](#)

... 3,6 см вблизи апоцентра **кольца** .

единственное число

множественное число

Другое

Пропустить

[Прокомментировать](#)

... сотрудничество с Пакистаном в **области** борьбы с терроризмом помогло ...

единственное число

множественное число

Другое

Пропустить

[Прокомментировать](#)

... Канады темы , популярность **группы** вышла далеко за пределы ...

единственное число

множественное число

Другое

Пропустить

[Прокомментировать](#)

Компания **ООО** « Нефтестрой » .

единственное число

множественное число

Другое

Пропустить

[Прокомментировать](#)

id	Текст + варианты разбора	1	2	3	Модератор
73 49	<p>... в свете первой удачной <b>[картины]</b> – закономерность .</p> <p>картина, NOUN, inan, femn, sing, gent картина, NOUN, inan, femn, plur, nomn картина, NOUN, inan, femn, plur, accs</p>	NOUN & sing	NOUN & sing	NOUN & sing	NOUN & sing
73 50	<p>... Питера FM » от <b>[Оксаны]</b> Бычковой хотели истории в ...</p> <p>оксана, NOUN, anim, femn, Name, sing, gent оксана, NOUN, anim, femn, Name, plur, nomn</p>	NOUN & sing	NOUN & sing	NOUN & sing	NOUN & sing
735 1	<p>... от Оксаны Бычковой хотели <b>[истории]</b> в этом же ключе ...</p> <p>история, NOUN, inan, femn, sing, gent история, NOUN, inan, femn, sing, datv история, NOUN, inan, femn, sing, loc история, NOUN, inan, femn, plur, nomn история, NOUN, inan, femn, plur, accs</p>	NOUN & sing	NOUN & sing	NOUN & plur	NOUN & sing

# Модерация

## **Ожидаемая точность:**

10% ошибок \* 3 участника

= 0.1% случаев, когда трое ошиблись

т. е. 1 ошибка / 1000 примеров

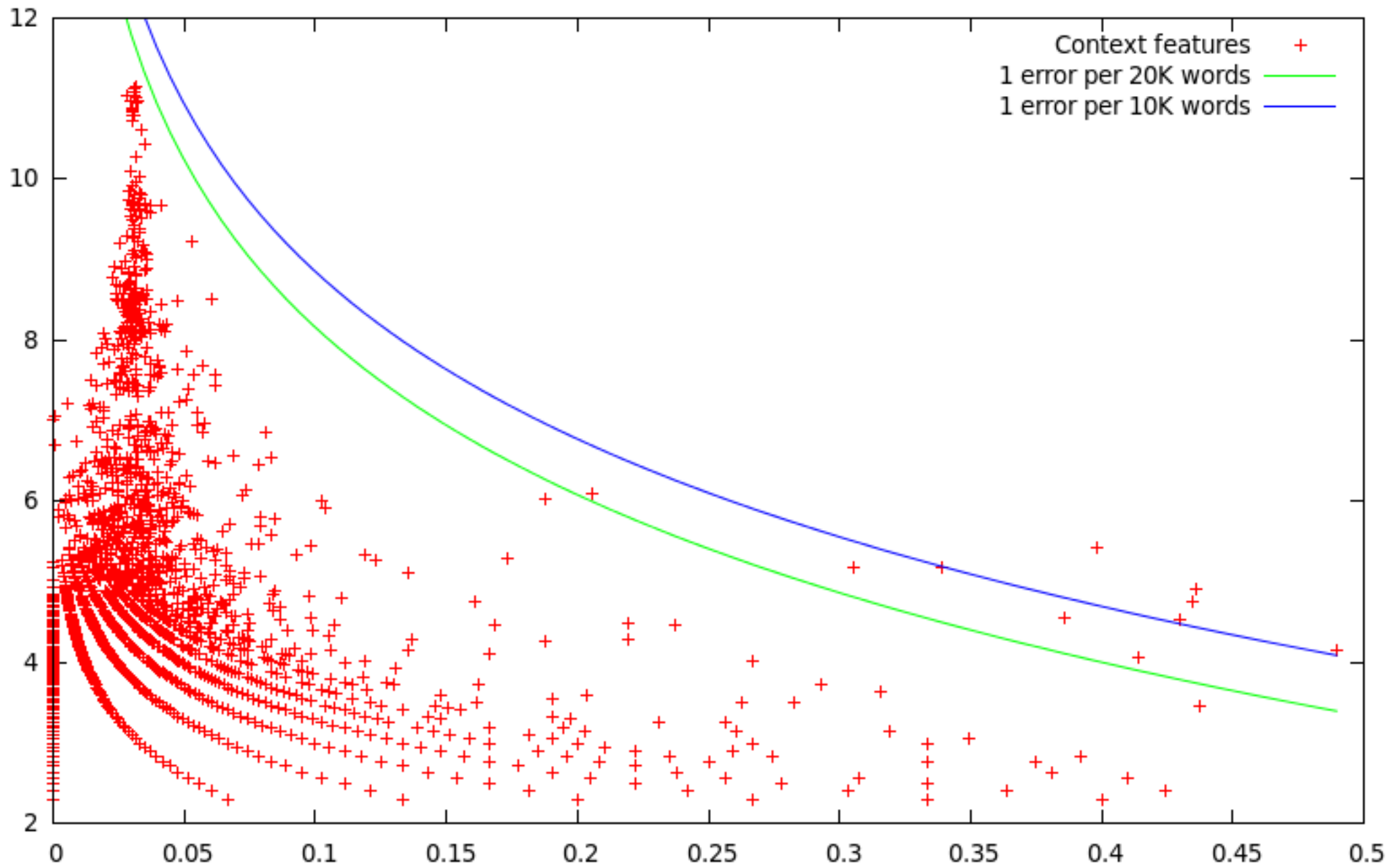
## **Реальная точность:**

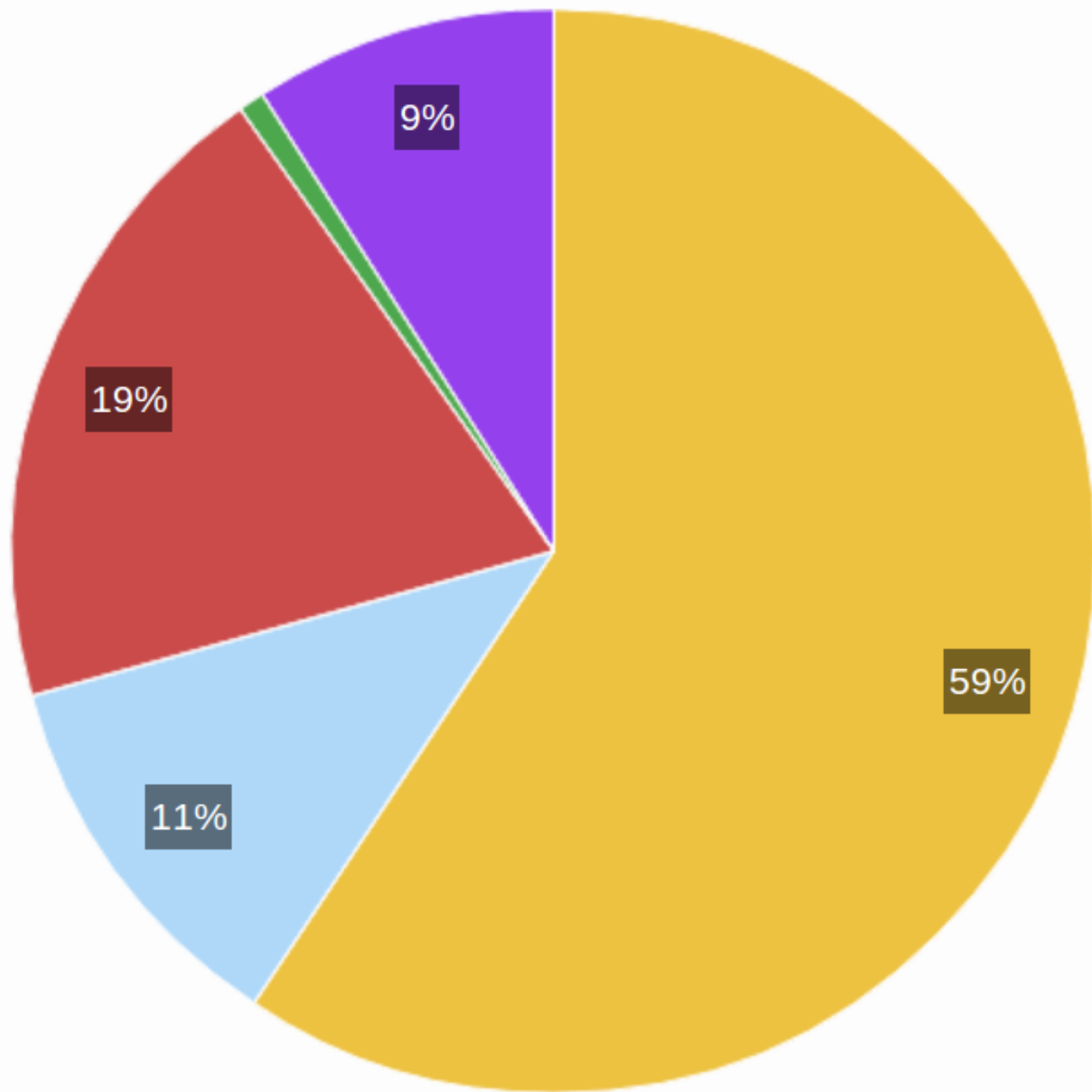
0.5 — 10% ошибок

2 % случаев, когда трое ошиблись

# Модерация

Context feature	Position	Total samples	Samples with disagreement	Samples without disagreement	Disagreement rate	Expected error probability
word=четыре	-1	64	47	17	73,44%	48,96%
word=две	-1	136	89	47	65,44%	43,63%
word=три	-1	115	75	40	65,22%	43,48%
word=два	-1	93	60	33	64,52%	43,01%
word=две	-2	58	36	22	62,07%	41,38%
word=одна	4	13	8	5	61,54%	41,03%
word=две	0	226	135	91	59,73%	39,82%
word=копейки	0	17	10	7	58,82%	39,22%
word=четыре	-	95	55	40	57,89%	38,60%

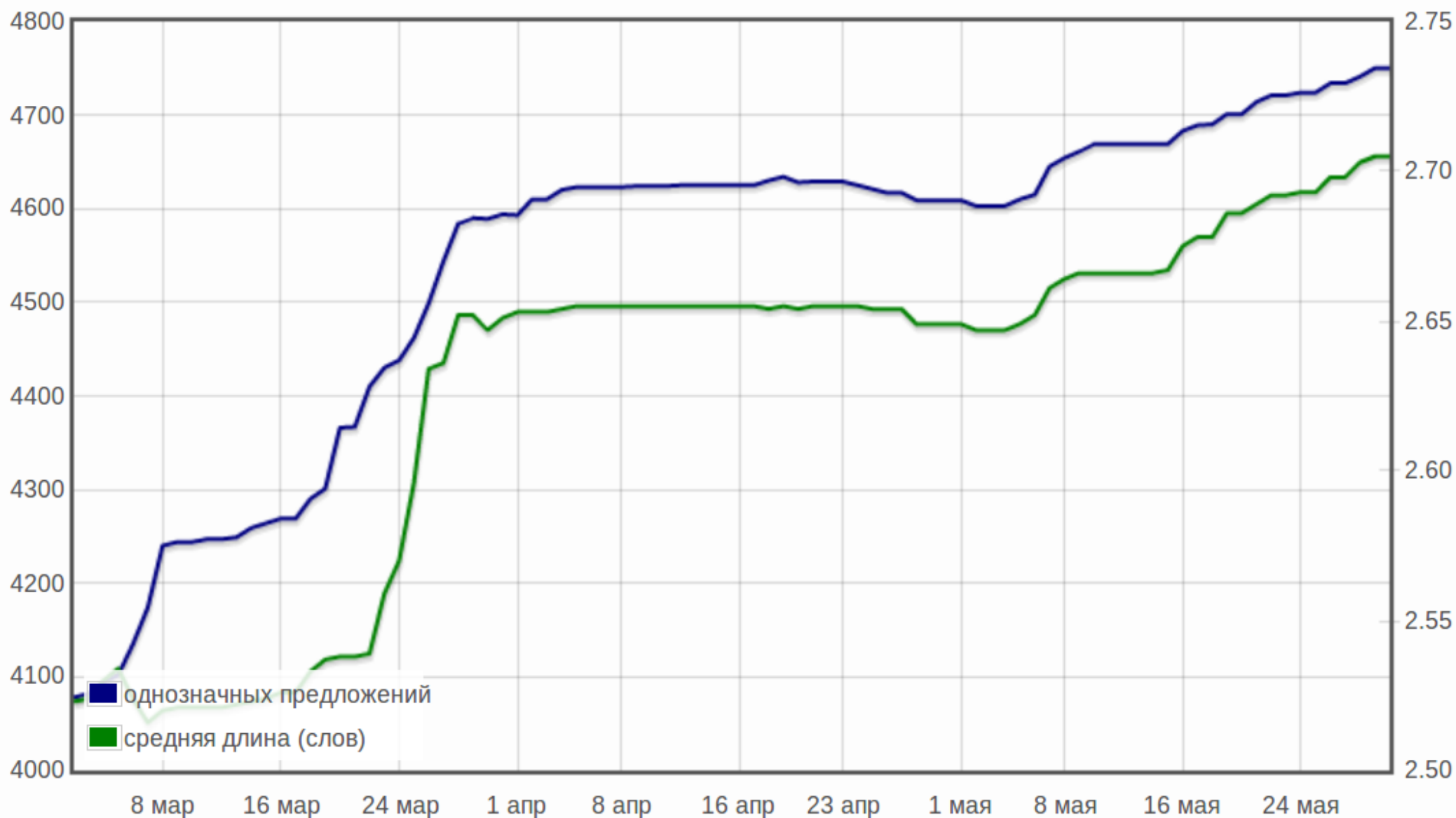




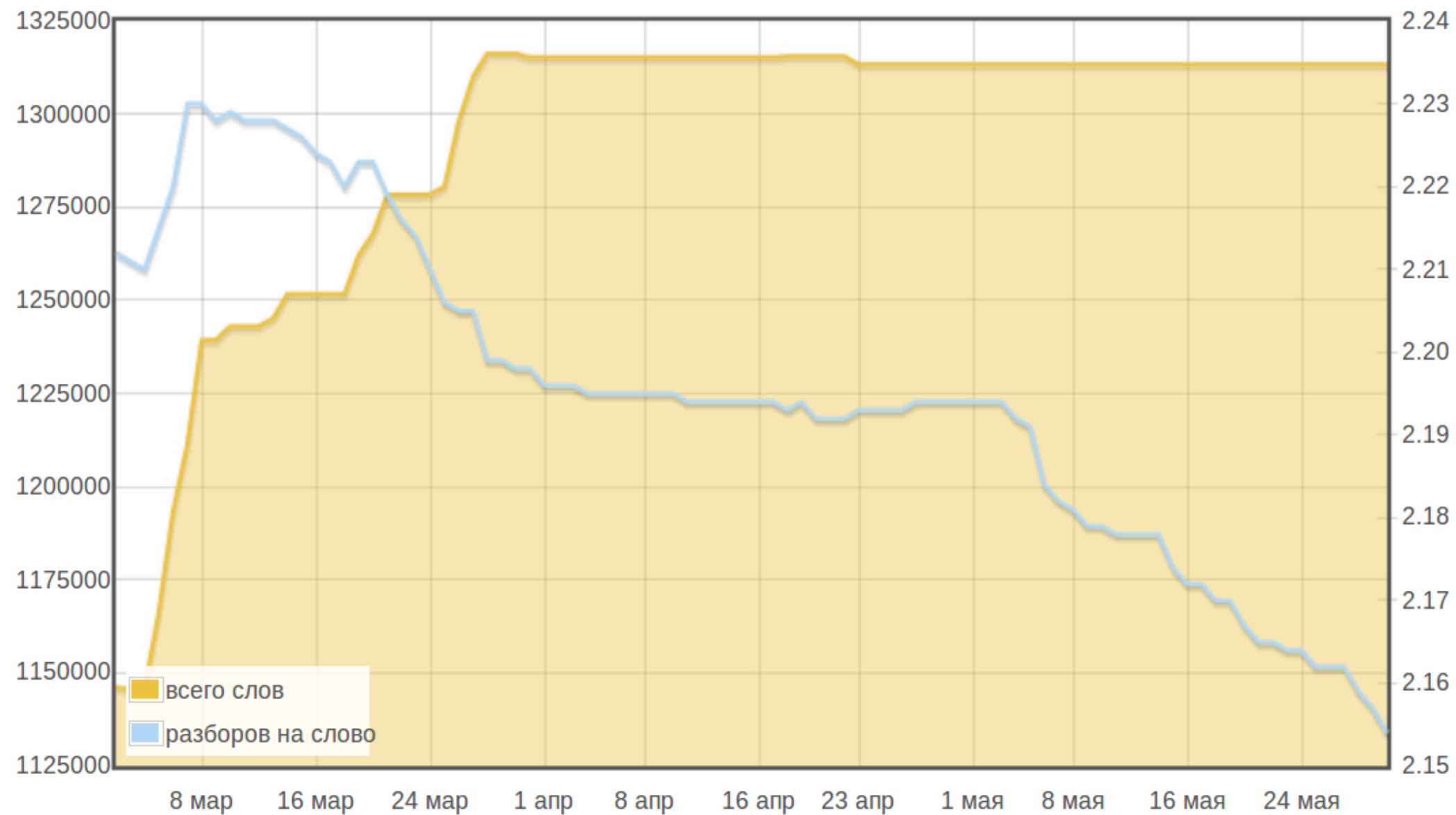
- готовится
- размечается
- размечено
- на модерации
- ушло в корпус



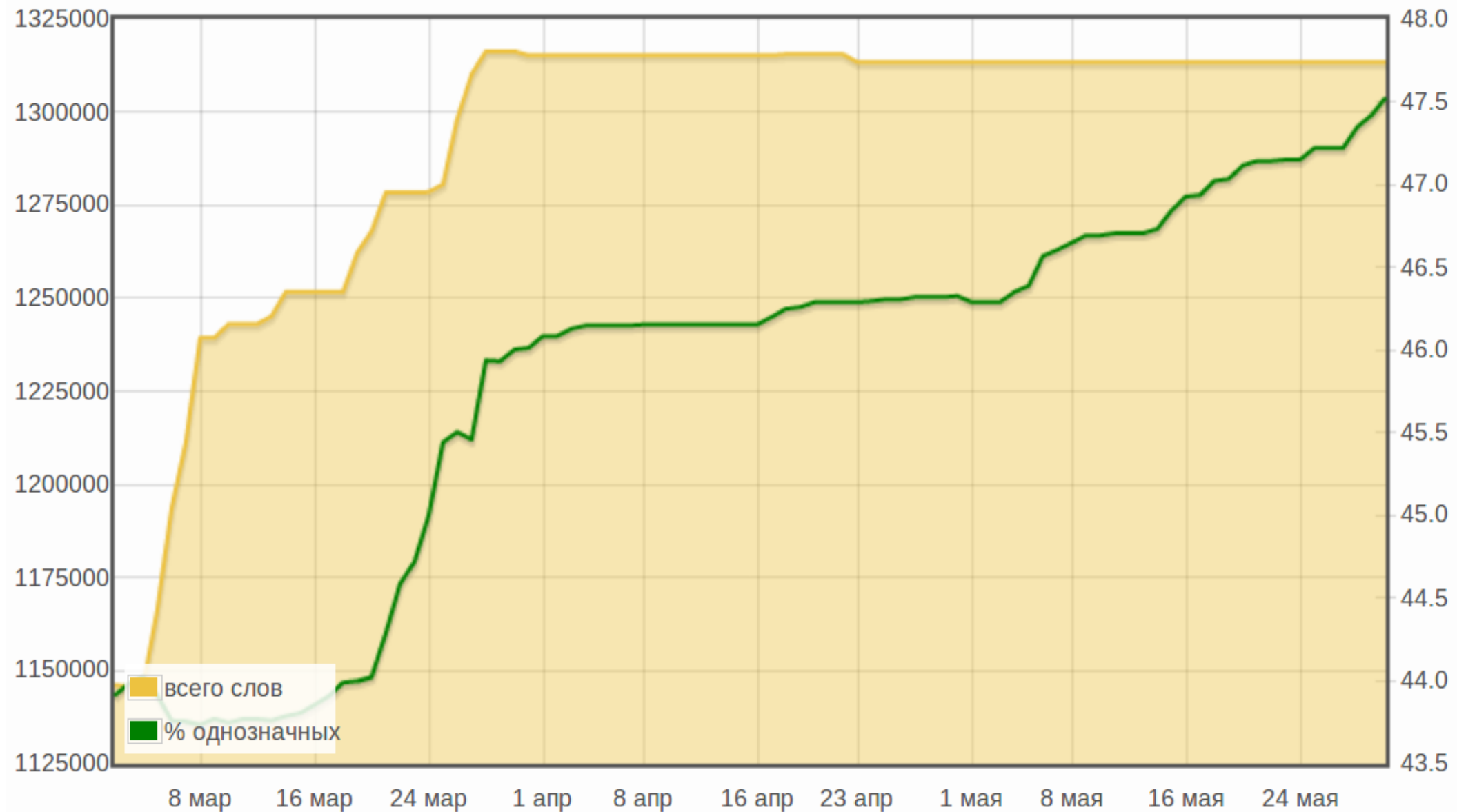
# Однозначные разборы



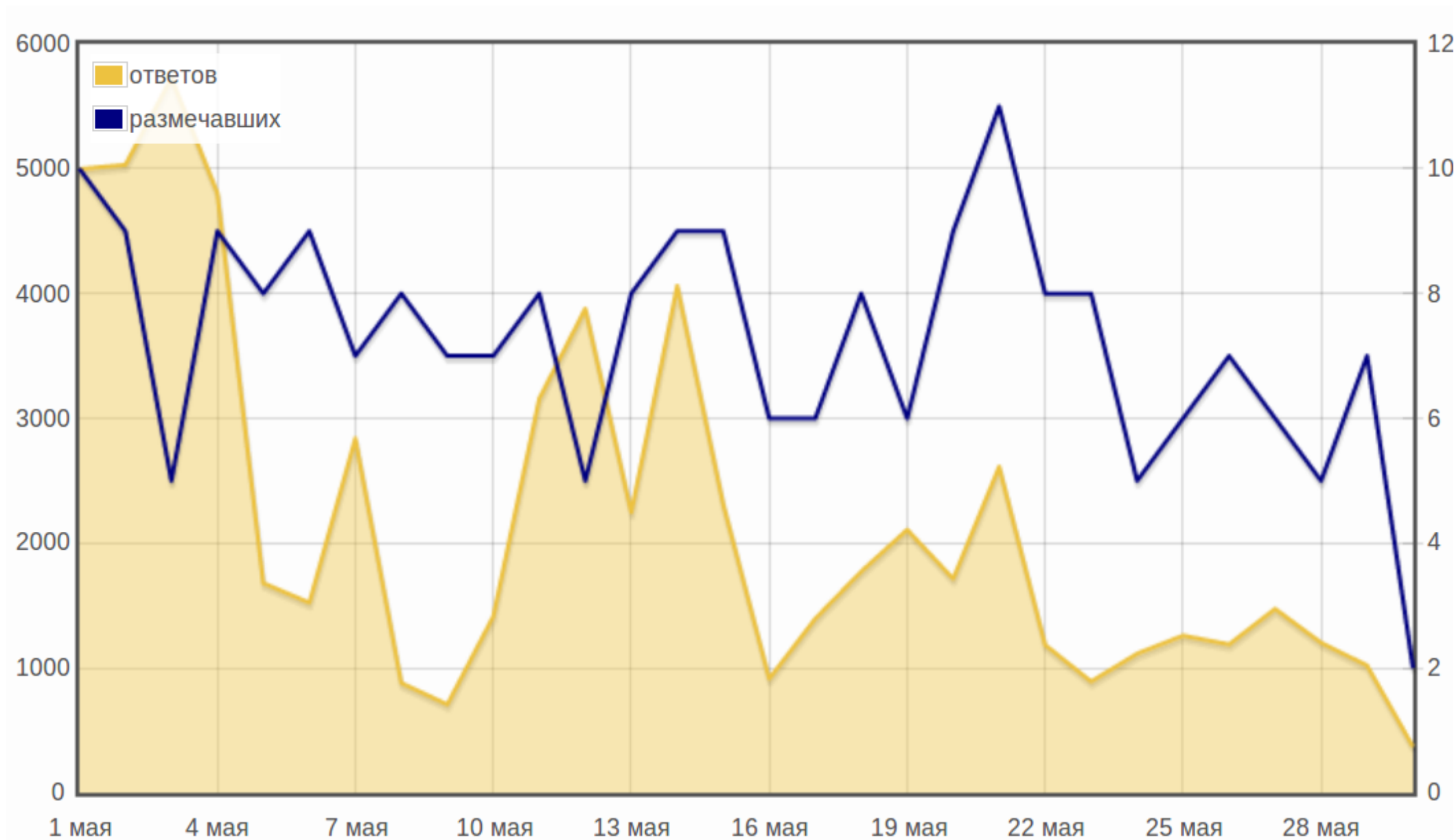
# Количество разборов на слово



# Количество однозначных слов

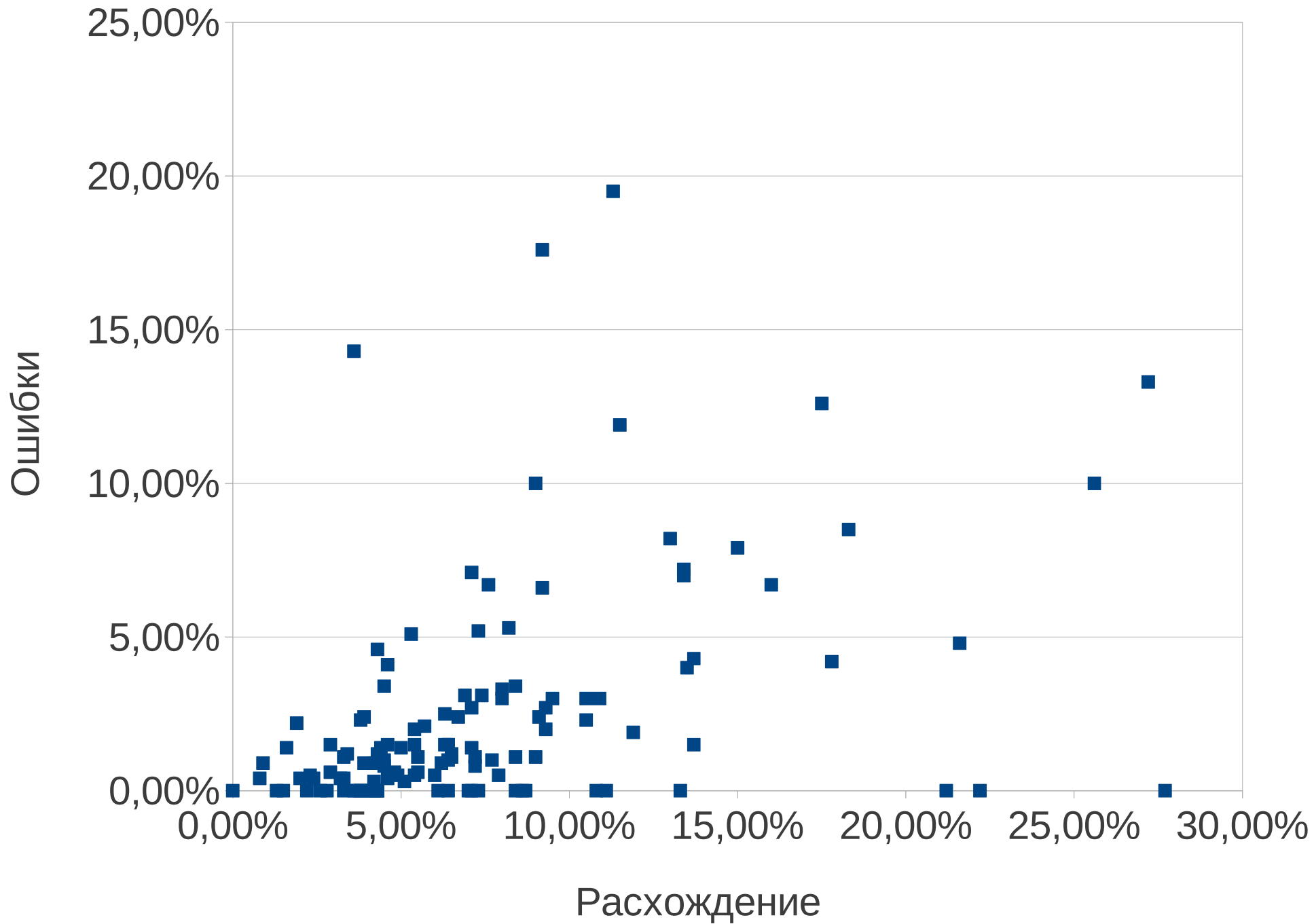


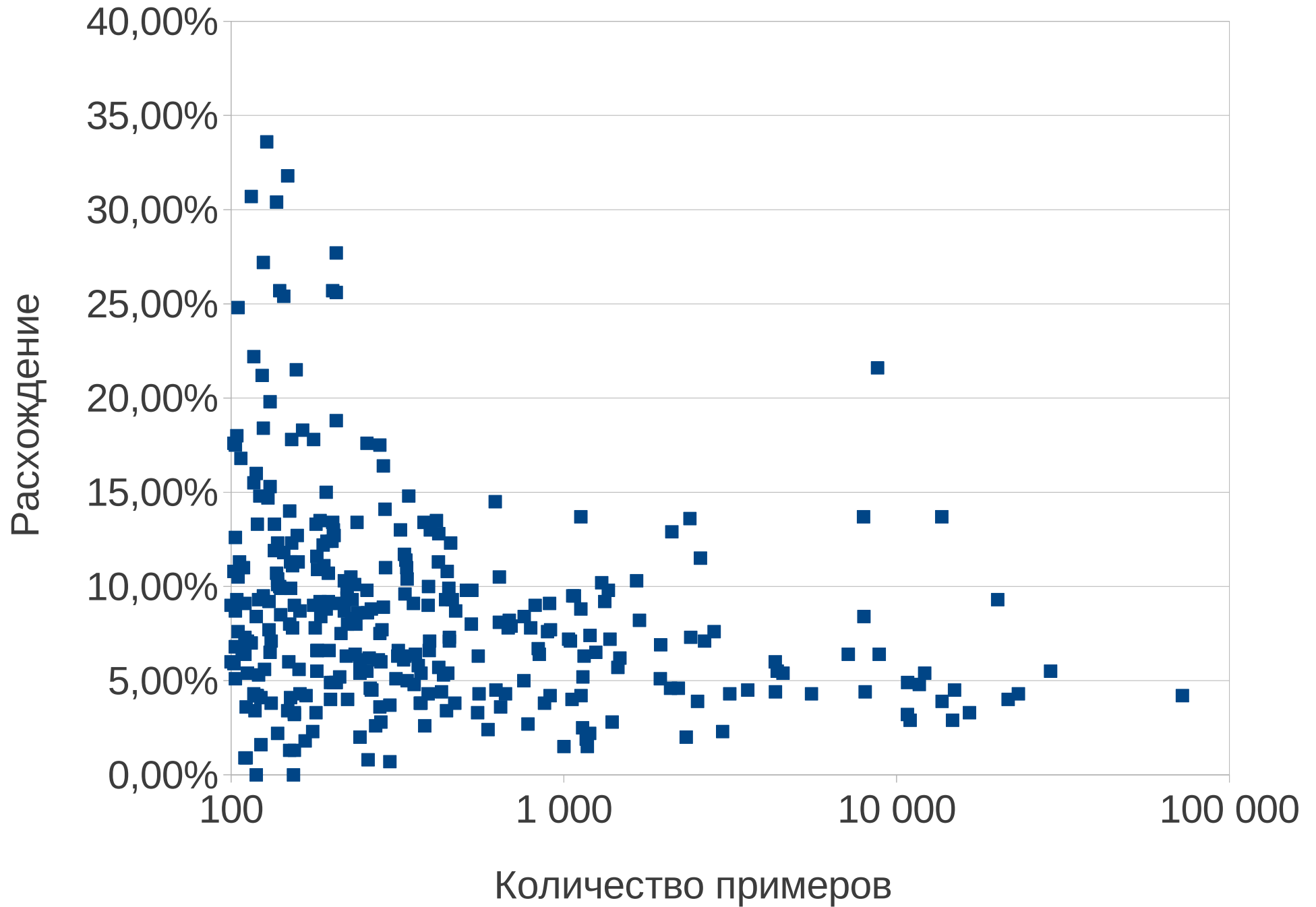
# Активность участников



# Участники

#	Участник	Рейтинг	Ответов	В завершённых пулах		В проверенных пулах		Последняя активность
				Размечено	% расхождений	Размечено	% ошибок	
1	Lvova	545499	359367	232205	4.1%	62783	0.7%	сегодня в 09:01
2	Nofenigma	174551	109663	100146	4.2%	30686	1.8%	27.05.13
3	Мяу	114471	103877	93993	2.2%	42936	0.6%	08.03.13
4	Rave	87745	83522	64597	3.5%	25730	0.6%	19.05.13
5	quorax	75832	38757	34580	4.3%	6441	0.6%	19.02.13





# Данные

- Корпус с частично снятой омонимией
- Подкорпус с полностью снятой омонимией
- Ответы на задания (-)
- Хронология разметки (- -)
- Словарь



# Что дальше?

- Разметить какие-то задания
- Рассказать коллегам
- Разрешить использование своих текстов на условиях Creative Commons
  
- [twitter.com/opencorpora](https://twitter.com/opencorpora)
- [vk.com/opencorpora](https://vk.com/opencorpora)
- [opencorpora@opencorpora.org](mailto:opencorpora@opencorpora.org)
- <http://opencorpora.org/?page=about>

Спасибо!