

Методы очистки запросов к поиску

Карпенко Максим, Протасов Сергей

Rambler Research

May 2011



Обзор доклада

- 1 Факты
- 2 Алгоритмы
- 3 Очистка моделей

Обзор доклада

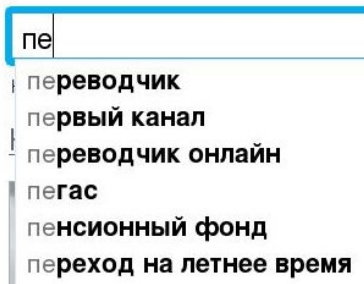
- 1 Факты
- 2 Алгоритмы
- 3 Очистка моделей

Подсказки в строке поиска

Половина запросов идёт через подсказки в строке поиска.



The image shows the top part of the Rambler search engine homepage. It features the Rambler logo (a globe icon) and the text "Rambler" in blue. Below the logo is a link "Сделать стартовой". There are two input fields for "Почта" (Email) and "Завести почту" (Create email). The "Почта" field contains "Логин" and the "Завести почту" field contains "@rambler.ru". Below these is a "Пароль" (Password) field.



The image shows a search bar with the text "пе" entered. A dropdown menu of suggestions is visible below the search bar. The suggestions are:

- переводчик
- первый канал
- переводчик онлайн
- пегас
- пенсионный фонд
- переход на летнее время



Подсказки опечаточника

10-15 % запросов - с опечаткой.



Возможно, вы искали: [одноклассники](#) (показано 2 результата)

1 [Одноклассники.ru](#) - поиск [одноклассников](#)

Сайт [odklassniki.ru](#). Главная страница. Регистрация на сайте. Поиск [одноклассников](#), однокурсников, бывших выпускников и старых друзей.

[www.odklassniki.ru](#) - [сохраненный текст](#) - [искать на сайте](#)

2 [Одноклассники КМ.RU](#) - бесплатный поиск [одноклассников](#)

Поиск [одноклассников](#), однокурсников, коллег и друзей, новых знакомых.

[Забыли пароль?](#)

[Украина](#)

[Казахстан](#)

[Беларусь](#)



100 способов написать неправильно Бритни Спирс

- бритни спирз, бритни сприс, бритни спирс, бритни спир, бритни спитс, ...
- бридни спирс, бритнти спирс, ритни спирс, вритни спирс, бритний спирс, бритней спирс, итни спирс, притни спирс, брини спирс, бринти спирс, бритри спирс, ...

Факты

- удаление (агенство/агентство) 8% случаев
- перестановка - 4% случаев
- вставка - 4% случаев
- замена - 80% случаев
- 90% опечаток с расстоянием 1

Обзор доклада

- 1 Факты
- 2 Алгоритмы
- 3 Очистка моделей

Алгоритм толпы

- Использование словарей - не подходит.
- Самый простой алгоритм - 28 строчек Норвига (norvig.com).

Алгоритм толпы

Расширения базового алгоритма

- модель ошибок, типы ошибок
- уверенность в результате
- модель языка (биграммы, **очистка**)

Обзор доклада

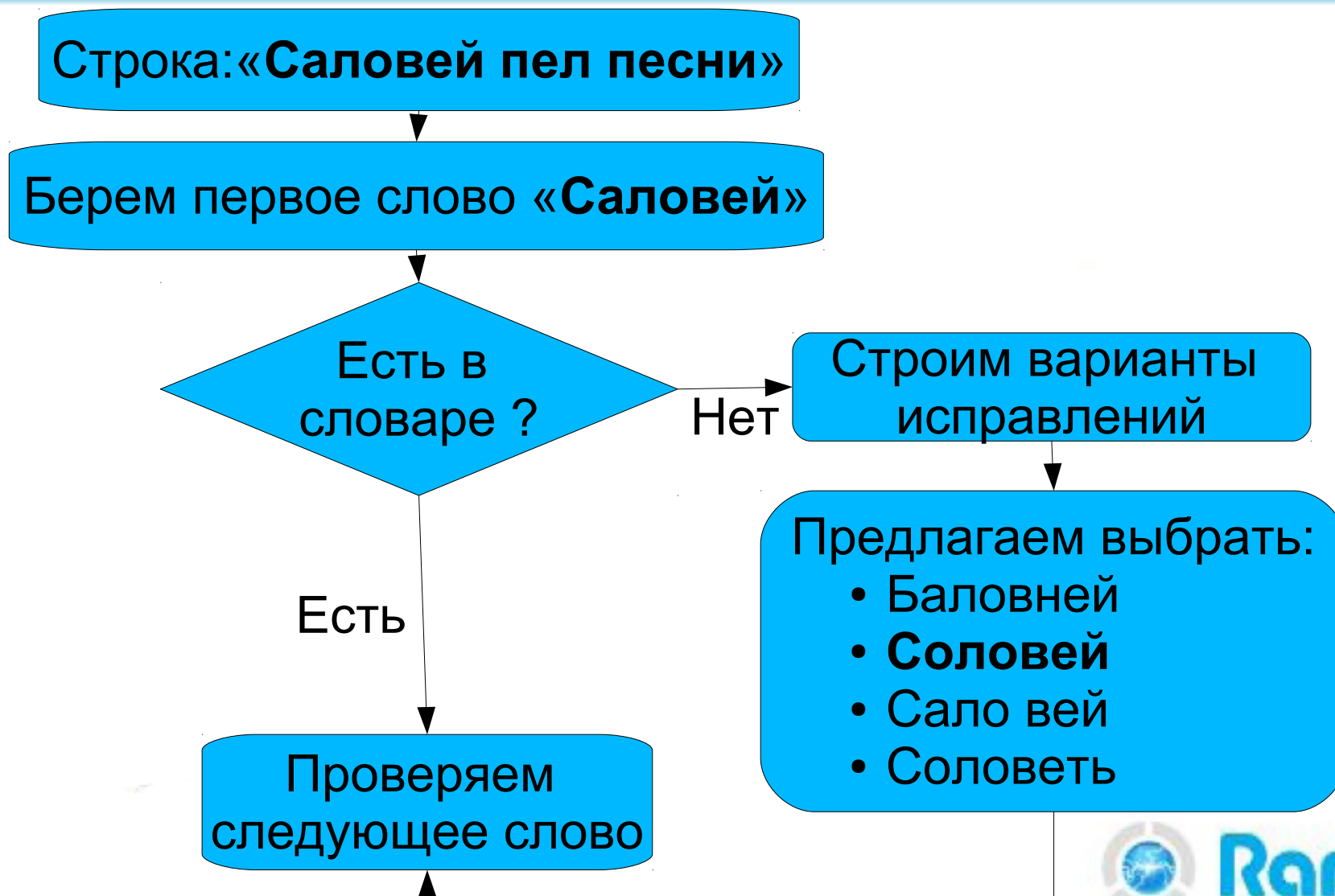
- 1 Факты
- 2 Алгоритмы
- 3 Очистка моделей

Некоторые методы очистки
моделей языка запросов к
поиску

Карпенко М.
Протасов С.

Рамблер

Простейший опечаточник



Словарь опечаток

ликимии → лейкемии

лекимии → лейкемии

ликимия → лейкемия

...

Плюсы:

- ♦ Высокая точность

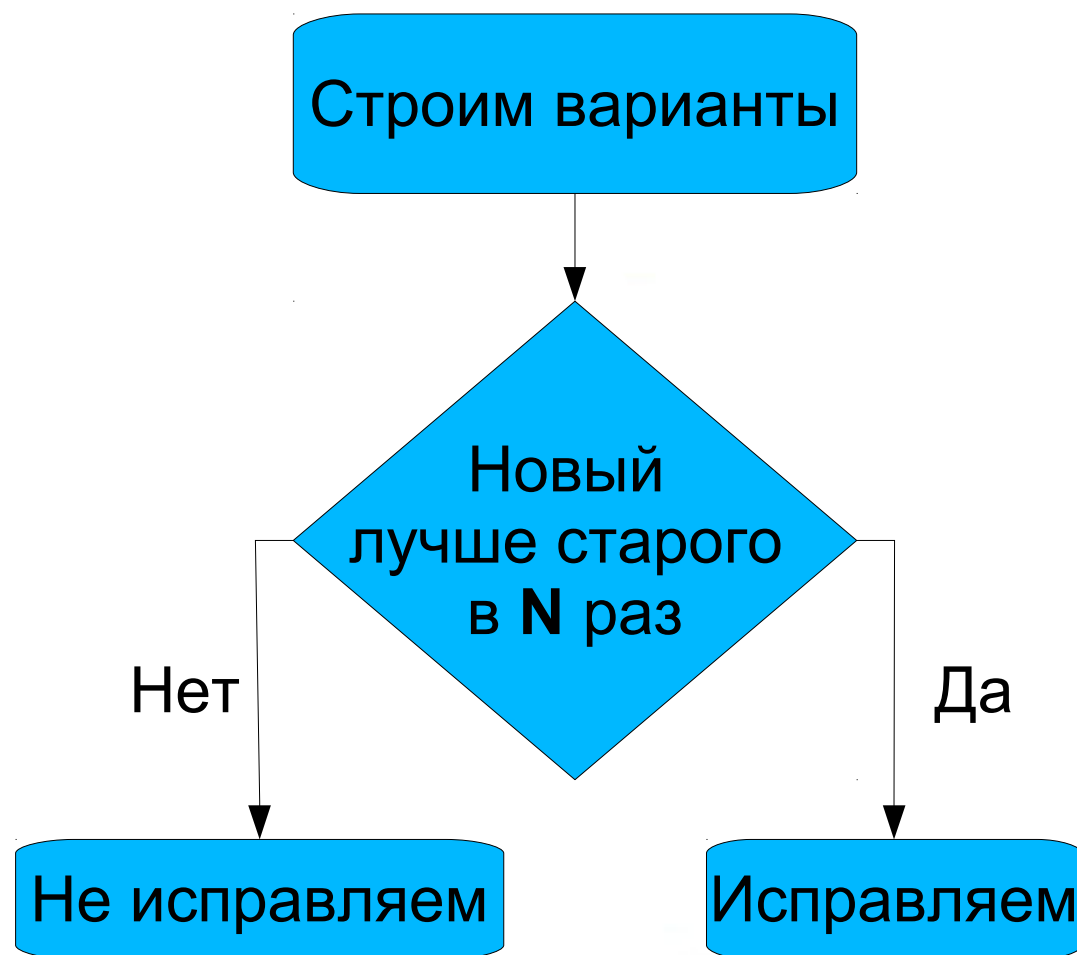
Минусы:

- ♦ Недостаточная полнота
- ♦ Высокая трудоемкость
- ♦ Не исправляет смысловые опечатки.
Небо сырое, но дождя нет → Небо серое, но дождя нет
- ♦ Плохо исправляет ошибки склейки и разбиения слов.

Модель языка.

укусила сабака - 18
укусила собака - 3571

группа собака - 11
группа сабака - 740



Расширенный опечаточник

Используется:

- ♦ Степень уверенности в исправлении
- ♦ Разные модели языка.
 - ✓ *Разные модели для однословных и многословных запросов*
 - ✓ *Тематические модели.*
- ♦ Поведенческие метрики.

Проблемы, создаваемые ошибками в модели.

$$\operatorname{argmax}_c P(w|c)P(c)$$

- Частотные опечатки.
 - ♦ Не исправляются
 - ✓ **И**нцефалит / энцефалит
 - ✓ Инте**л**егенция / интеллигенция
 - ♦ 2 ошибки исправляются на одну частотную
 - ✓ Ха**н**др**а**ксит → Хондр**а**ксит → Хондроксит
- Редкие опечатки.
 - ♦ Снижается скорость работы.
 - ♦ Модель занимает много места.
 - ✓ Вентилятор (112 опечаток): Венти**о**лятор, Вентил**а**тро, Венти**и**атор ...
 - ✓ География (84 опечатки): Ге**о**рафия, Геога**ф**фия ..

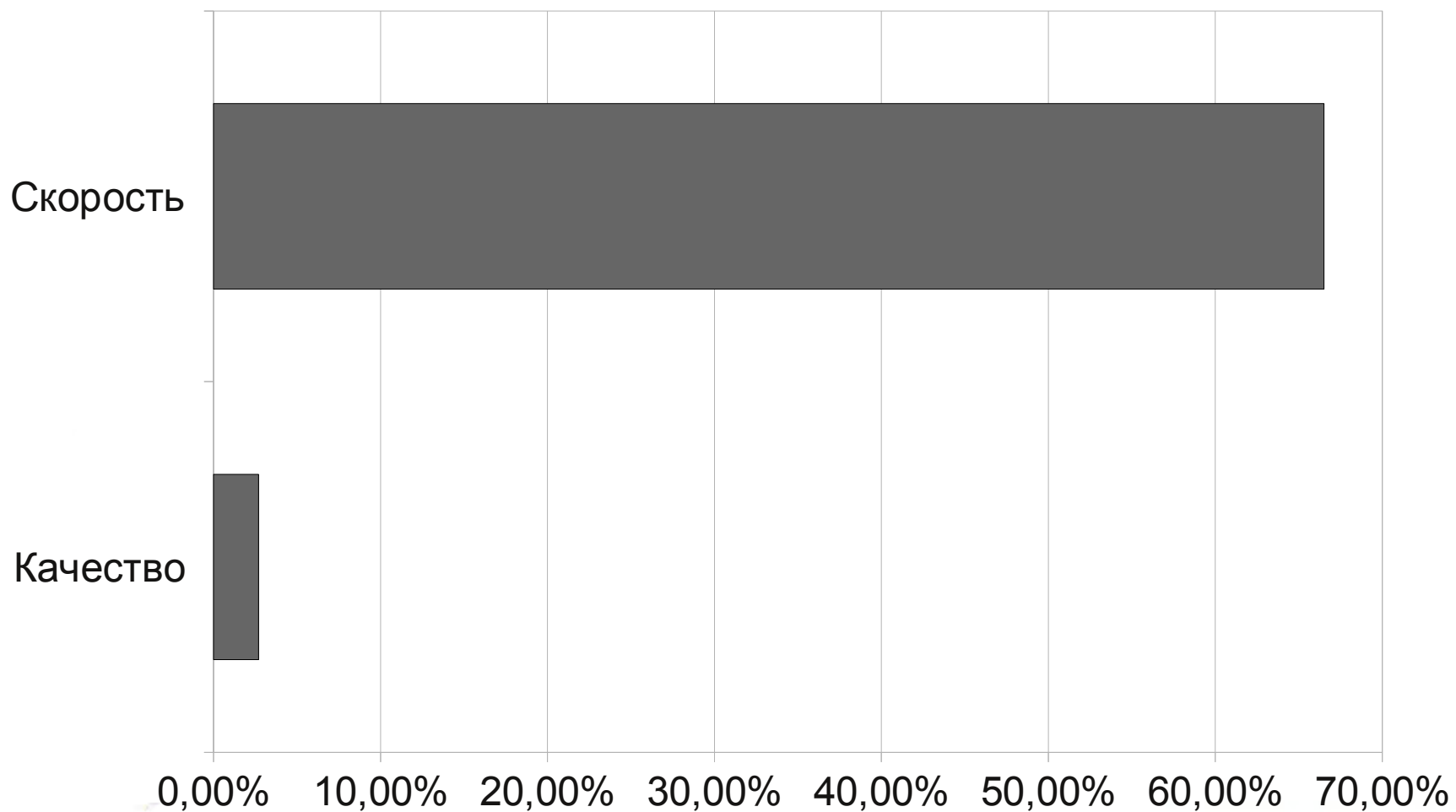
Удаление низкочастотных ошибок в модели

- Удаляем “невероятные” слова языковой модели.
 - ♦ Качество не меняется
 - ♦ Скорость +58 %
- Самоочистка
 - ♦ Качество не меняется
 - ♦ Скорость +9 %
- Результаты поиска.
 - ♦ Качество + 1,2 %
 - ♦ Скорость + 5 %

Выявление и удаление частотных ошибок.

- Соотношение частотности похожих слов.
 - ♦ "бессо**н**ица" 33735
 - ♦ "бессонница" 34471
- Использование общего контекста.
 - ♦ программы *для android*
 - ♦ програ**м**ы *для android*
 - ♦ пра**а**граммы *для android*
- Результаты поиска. (Уровень цитируемости)
 - ♦ "гидропередача" 28376
 - ♦ "гидро передача" 52

Результат удаления ошибок.



Проблемы очистки

- Имена и фамилии.
- Названия: лекарств, компаний, населённых пунктов.

Система выделения признаков запросов:

Тематика лекарства (1000 признаков)

- ✓ *карбуситрат горелутамин*
- ✓ *ниналонин оксацтозан пициллин репитоксодин
фитотиланадат строин тубутатин*
- ✓ *ликсферид питиазоллол ягмониламид карбозыкол
фенамещенный риптомазоин*

Спасибо за внимание.

Вопросы ?

m.karpenko@rambler-co.ru
s.protasov@rambler-co.ru

