

Проект “ANTIDICT”

База морфологических структур, допустимых в
современном русском языке

Студенты:

Алексей Доркин

Владислава Смирнова

Антон Вахранев

Куратор: Варвара

Магомедова

Что мы успели сделать

- Начали разделять слова по группам по определенному признаку (например, по общим словообразовательным признакам)
- Обработали и разметили словарь Дьякова
- Выкачали словарь Зализняка со словоформами
- Написали классификатор для заимствований, обучили модель

Группируем слова

накачанный очкарик прекрасный рождённый непосоветуют любой предмет теперь поменялось боле льщик оспи дигонщик	репортаж ечка геймер юга быд лячество степанскую громыхалки жизнеоформляющих киви боем	ловка ческой ребенковская перехрустит луноградцы фидошников манхетенской	фс ячности афтобус пионэром бидилы трэска т подлечица чума зег обзатце перажки маральна е животные
деза стером ресэмп линг тролбейн ргб трекеров	редбуловых асогласна тиражир омедведюсь набухтелся	ляляка мимишном щупака свиняки милошных	

Группируем слова

антииудаизмом
неправославное
антифосфатные
неоценившему
некапиталистическим
гиперкомпенсировать
киноальянс
горетуристами
видеоклипам
фотослайдами
китчевость
схемку
впариванию
выдеорепортажем
тромбозникам

электрагенератары
географисеских
переводишьс
справяю
разобразным
кажется
документации
сковыжималку
отсутстви
компьютер
фотоапарата
метериться
неизменёных

новогодих
создвездия
подтвердили
бершадського
обяснило
мудчины
импланты
чужародная
джнисы
накуриваются
сверхъественного
компенисрующие
рождествеская

оечка
суперартмаркет
мармонам
четров
катинками
родоночальник
эспресионистская
галаграма
эффективное
псхиушки

Размеченные данные из Дьякова

__label__нео юникс
__label__нео юнилевер
__label__нео юнилэтаризм
__label__нео юнион
__label__нео юнионизм
__label__нео юнионист
__label__нео юнионисты
__label__нео юниорский
__label__нео юнипол
__label__нео юнисеф
__label__нео юнискан
__label__нео юнит
__label__нео юнитард
__label__нео юнити

__label__dict лидируем
__label__dict воздуховоде
__label__dict гугенотке
__label__dict развёрстывавшуюся
__label__dict взъедавшееся
__label__dict
перебинтовывающим
__label__dict прирастится
__label__dict нимфоманку
__label__dict покормившийся
__label__dict натуралисту
__label__dict учтивее
__label__dict поваловой

T3

<https://docs.google.com/document/d/1fpn72q8bKqhFnCaTmbqGZEtWS6WcXIP9WC1VpviAGEc/edit#>

Репозиторий

<https://github.com/wksmirnowa/antidict>

О чем проект (слайд на случай, если кто-то забыл)

- изучение необычных слов, которые можно встретить в интернете (50 млн словоформ из ГИКРЯ)
- классификация слов по большим категориям (например, заимствования, экспрессивные формы, слова со случайными и с намеренными ошибками и искажениями)
- классификация слов по подкатегориям (например, заимствования, слова с заимствованными морфемами, слова, образованные от русских корней)
- морфологический анализ слов