

**Модели и методы формирования и  
использования баз знаний для  
интегрированных мультипредметных  
информационных систем**

на примере обработки массива данных  
социальной сети ВКонтакте

# ЦЕЛИ РАБОТЫ

- Разработка инструментов сбора больших данных
- Разработка инструментов оптимального хранения больших данных
- Учёт особенностей работы с большими данными соцсетей
- Разработка алгоритмов подготовки и анализа социальных данных
- Разработка алгоритмов визуального представления информации
- Разработка механизма подключения сторонних модулей обработки

# ЧЕМ ИНТЕРЕСНЫ СОЦСЕТИ (ВОЗМОЖНОСТИ)

- Выявление настроений в обществе (Срез эмоционального фона)
- Прогнозирование реакции к будущим событиям.
- Поиск кадров - “автоматическое собеседование”.
- Автоматические соц.опросы. (более искренние и честные ответы)
- Создание новостных агрегаторов на базе социальной заинтересованности.
- Выявление деструктивной активности.
- Поиск связанных событий методикой поиска без ключевых слов.
- Карта активности (гео-метки, связанные с активностью в соц.сети)
- Выявление ботов

# ОСОБЕННОСТИ ТЕКСТА СОЦИАЛЬНЫХ СЕТЕЙ

## Текст:

- текст короткий
- частое применение нетекстовых символов
- сообщение может целиком состоять из нетекстовых символов
- частое намеренное искажение слов
- распространены различные синтаксические ошибки
- различная тематика обсуждения

## Время жизни

- время жизни документов (публикация, социальная реакция, забвение).
- возможность отслеживания по времени

# НЕОБХОДИМЫЙ ИНСТРУМЕНТАРИЙ

- сбор социальных данных
- хранение собранных и обработанных социальных данных
- хранение информации о динамике изменения отслеживаемых документов
- подготовка текста
- средства графического отображения полученных результатов
- средства получения различных выборок
- генерация наборов данных (датасетов) для инструментов тематического моделирования
- выявление тем (методики тематического моделирования)
- возможность быстрой интеграции пользовательский модулей обработки

## ПАРАМЕТРЫ СЕРВЕРА ОБРАБОТКИ

```
core > sys_config.py > ...
1  '''
2  {'коллекция на сервере мониторинга': { entity_name    коллекция на сервере обработки
3                                       metrics          какую инфу для этой коллекции собирать и её тип }
4  }
5
6  INFO_SET:
7      entity_name    Имя сущности на сервере обработки
8      fields        Поля, которые нужно хранить на сервере обработки
9                  id            поля id, записываются единожды при создании документа на СВ, больше не меняются
10                 counters     целочисленные счётчики (просмотры, лайки, ...)
11                 text         текстовые поля
12                 complex      часть полей является документом, нужно выгызать из них вложенные поля и размещать на одном
13                               уровне с остальными. Пример - City в группах. def proc_complex забирает нужное вложенное поле
14                               на основании этого правила. Пока сделана обработка одного вложенного уровня.
15                 vs_store_name имя поля под которым сохранять на CBsdsd
16                 inner_name    вложенное имя внутри поля container
17                               Вложенная метрика строго целочисленная. Правило позволяет извлекать несколько вложенных полей
18                               (***) можно добавить параметр-маршрут или уровень для работы с многоуровневыми составными полями
19                 variable_fields Поля, не всегда встречающиеся, которые можно заполнить просканировав другие поля документа.
20                 Not Active If 'field' == 'value'
21                 idife         Ignore Document if Empty поле. Например, нам ни разу не упукались посты или коменты с пустым полем text
22                 creation_date Дата создания документа. 'ms' - из какого поля брать из документа ВК, 'vs' - под каким именем сохранять на СВ
23
24  Z_SYMBOLS_LEN      Timestamp слева дополняется количеством нулей указанных в этом параметре
25  '''
26  INFO_SET = {
27      'mon.posts': {
28          'entity_name': 'posts',
29          'fields': {
30              'id':      ['id','from_id','owner_id','reply_owner_id','postponed_id'],
31              'counters': ['reposts','likes','views','comments'],
32              'text':    ['text'],
33              'raw':     ['attachments','comments','copyright','post_type','geo','signer_id',
34                          'copy_history','is_pinned','marked_as_ads',''],
35              'str_date': [],
36              'complex': {},
37              'variable_fields': []
38          },
39          'na_if':      {},
40          'idife':     ['text'],
41          'extract':   {},
42          'creation_date': {'ms':'date', 'vs':'date'}
```

```
id: 61904449
> members_count: Object
> name: Object
  name__last: "Подслушано Апатиты"
  name__last_ts: "000001622142004"
  name__last_clean: "Подслушано Апатиты"
> description: Object
  description__last: "У нас вы найдете откровения и новости жителей г. Апатиты и г. Кировск
  ..."
  description__last_ts: "000001622142004"
  description__last_c_: "У нас вы найдете откровения и новости жителей г. Апатиты и г. Кировск ..."
< city: Object
  000001622142004: "Апатиты" ←
  city__last: "Апатиты"
  city__last_ts: "000001622142004"
  city__last_clean: "Апатиты"
> country: Object
  country__last: "Россия"
  country__last_ts: "000001622142004"
  country__last_clean: "Россия"
> contacts: Array
> counters: Array
  cover: ""
> links: Array
  place: ""
> site: Array
> status: Array
  trending: ""
> verified: Array
> wall: Array
  public_date_label: ""
  market: ""
> main_section: Array
  fixed_post: ""
> age_limits: Array
  wiki_page: ""
  last_mon_ts: "000001648396805"
> aggregated_counters: Object
  rating: 0
  ntp: 1
```

## ПОДСЛУШАНО АПАТИТЫ

## ТИПИЧНЫЙ ПОСТ

```
date: "000001622055649"  
id: 2046224  
from_id: -61904449  
owner_id: -61904449  
reply_owner_id: 0  
postponed_id: 0  
> reposts: Object  
> likes: Object  
> views: Object  
> comments: Object  
> text: Object  
  text__last: "Может кто знает где режут зеркальную плитку? ( нужно вырезать отверст..."  
  text__last_ts: "000001622142008"  
  text__last_clean: "Может кто знает где режут зеркальную плитку? нужно вырезать отверстия ..."  
  attachments: ""  
  copyright: ""  
> post_type: Array  
  geo: ""  
> signer_id: Array  
  copy_history: ""  
  is_pinned: ""  
> marked_as_ads: Array  
  : ""  
  last_mon_ts: "000001622400008"  
> aggregated_counters: Object  
  rating: 67  
  ntp: 1
```



```
date: "000001620121843"
id: 124628
from_id: -68246638
owner_id: -68246638
reply_owner_id: 0
postponed_id: 0
> reposts: Object
  < likes: Object
    000001622142008: 12
    000001622283606: 13
    000001622571610: 14
    000001623093615: 15
    000001624002012: 16
    000001624135209: 17
    000001624992009: 18
    000001627161605: 19
    000001633630810: 20
    000001633701605: 21
  > views: Object

  < aggregated_counters: Object
    < reposts: Object
      max: 8
    < dynamics: Object
      dts: 11990397
      tsc: 9797
      dmv: 8
```

## ХРАНИЕНИЕ ДИНАМИКИ СОЦИАЛЬНОЙ АКТИВНОСТИ

### Метрики динамики социальной активности:

**dts** Diff Timestamp Sum  
Сумма временных “разрывов” между изменениями метрик

**dmv** Diff Metrics Values  
Сумма изменений метрик между срезами

**tsc** ts\_count  
Количество временных меток (срезов) метрики

# Расчёт рейтинга документа

```
DOCUMENT_RATING_K = {'reposts':20, 'comments':10, 'likes':5, 'views':1}
```

$$i = \frac{K_R \times Reposts + K_C \times Comments + K_L \times Likes \times V_v}{(\bar{Views} - 1) \times 100}$$

, где

$$V_v = \frac{\sum ts_{v+1} - ts_v}{n}$$

$V$  - скорость изменения метрик поста. Временные промежутки между изменениями метрик делим на количество этих изменений ( $n$ ).

$ts$  - TimeStamp поста

# Социальная активность ВК

## Метрики:

likes  reposts  comments  views  LRCV-grow  rating-ss

## Сообщества:

- × Кольский Север | Апатиты Мурманск Хибины и др
- × [X] Курорт Большой Вудъявр | BigWood ski resort
- × [X] Лавина (Апатиты | Кировск | Хибины)
- × Наши Апатиты
- × Подслушано Апатиты
- × [X] Подслушано Кировск Хибины
- × [X] Профсоюзная организация "ФОСАГРО-АПАТИТ"
- × [X] Рабочее время | Апатит
- × [X] СУРОВЫЕ ХИБИНЯ
- × [X] Твой Кировск
- × [X] Типичный Комбинат "Апатит"
- × Хибины Апатиты Кировск Народное ТВ
- × [X] Чёрный | Белый список Апатиты

Get current dataset

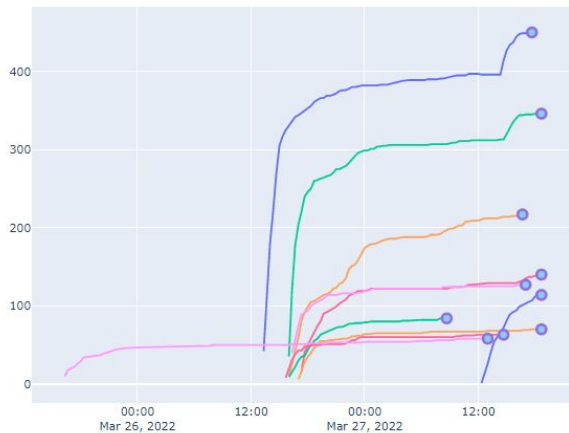
Период, дн.:



Документов:



Поделились | Документов: 10 | Период: 3 дня



- [ (1) Апатиты: Подслушано Апатиты ] - Кто спорил? Победила "Пятерочка"👁️ новый сетевик во  
● 1 day, 4:20:00
- [ (2) Апатиты: Подслушано Апатиты ] - На BigWoode в Кировске погиб человек Появились в  
● 1 day, 2:40:10
- [ (3) X: Подслушано Кировск Хибины ] - На BigWoode в Кировске погиб человек Появились в  
● 1 day, 0:39:59
- [ (4) Кировск: СУРОВЫЕ ХИБИНЯ ] - Не проходит и года, чтоб какой-нибудь турист не по  
● 1 day, 1:20:10
- [ (5) X: Хибины Апатиты Кировск Народное ТВ ] - \$\$\$ В ХИБИНАХ ПОГИБ ГОРНОЛЫЖНИК! 25 марта, око  
● 1 day, 0:40:00
- [ (6) Апатиты: Наши Апатиты ] - Расписание автобусов по Кировску и Апатитам! 🌟 Co  
● 6:20:11
- [ (7) Апатиты: Наши Апатиты ] - На BigWoode в Кировске погиб человек Появились в  
● 16:39:59
- [ (8) Кировск: СУРОВЫЕ ХИБИНЯ ] - Много сахара апатиты пятерочка на строителей  
● 1 day, 1:40:11
- [ (9) Апатиты: Подслушано Апатиты ] - \$\$\$ СРОЧНО: телеграм канал сообщает, что на биз  
● 23:00:01
- [ (10) Кировск: Профсоюзная организация "ФОСАГРО-АПАТИТ" ] - Информация по нашему новому партнёру дисконтной пр  
● 1 day, 20:40:01





(5) [Кто спорил? Победила "Пятерочка" 😊 новый сетевик возле Ре-Монта.](#)

👁18555 | ❤194 | 💬69 | ↩451

👤 [Подслушано Апатиты](#)



(1) [Кто спорил? Победила "Пятерочка" 😊 новый сетевик возле Ре-Монта.](#)

👁18555 | ❤194 | 💬69 | ↩451

👤 [Подслушано Апатиты](#)



(2) [На BigWoode в Кировске погиб человек Появилось видео с места событий](#)

👁16709 | ❤66 | 💬73 | ↩347

👤 [Подслушано Апатиты](#)



(3) [На BigWoode в Кировске погиб человек Появилось видео с места событий](#)

👁12324 | ❤52 | 💬46 | ↩218

👤 [Подслушано Апатиты](#)

