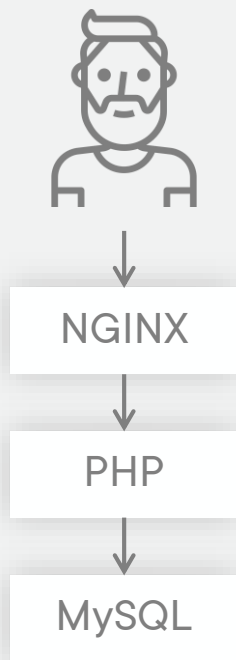




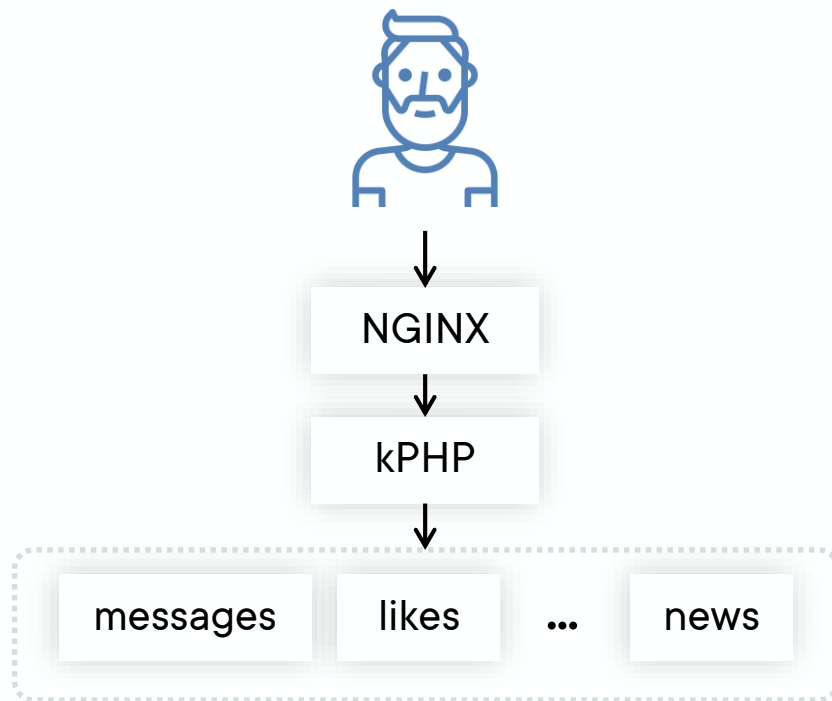
# Как устроены базы данных ВКонтакте?

Борис Минаев

# Типичный сайт



# ВКонтакте



# Что такое движок?

---

- Хранилище данных
- Бизнес логика
- C/C++
- Микросервис?

A word cloud featuring various project names and terms. The words are arranged in a dense, overlapping manner, with some words being significantly larger than others. The colors of the words are primarily shades of blue, brown, and red. The words are as follows:

- rpc-proxy
- logs-db
- persistent-longpoll
- messages
- antispam
- blackkad
- photo
- filesys
- memcached
- storage-redirect
- queue
- playlists
- logs-stats
- storage
- notify
- tlclient
- cache
- kafka-proxy
- money
- smart-alerts
- random
- online
- dns
- sawnews
- statsx
- streaming-api
- views
- likes
- copyfast
- weights
- ping
- sandbox
- socket
- magus
- liked
- graph
- tasks
- places
- seqmap
- replicator
- hints
- expressivita
- captcha
- trees-storage
- set
- mutual-friends
- copyexec
- news
- letters
- search
- image
- pmemcached
- password
- lists
- live-balancer
- text
- mc-proxy
- nostradamus
- isearch
- db-proxy
- ipdb
- friend
- logs-collector
- meowdb
- watchcat
- meowdb

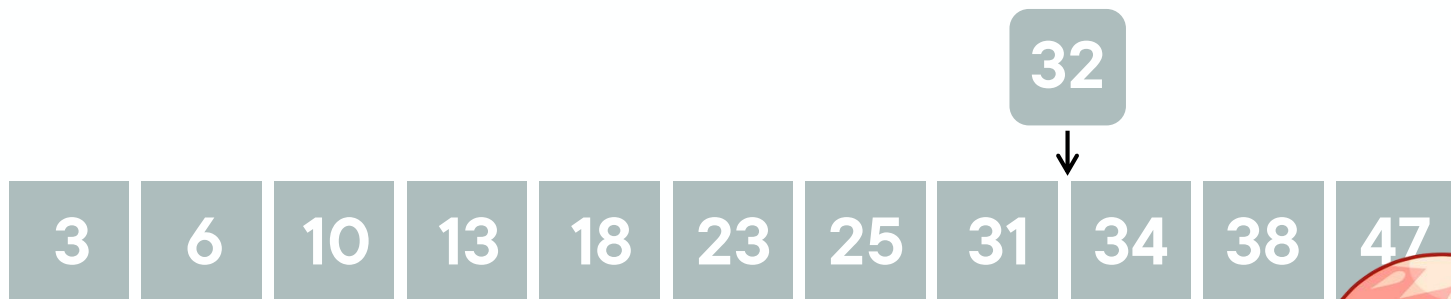
# Что получает ВКонтакте от использования движков?

---

- Скорость
- Эффективность хранения данных
- Атомарность
- Легкость добавления нового функционала

# Мы очень любим оптимизировать

Как найти место числа в отсортированном массиве?



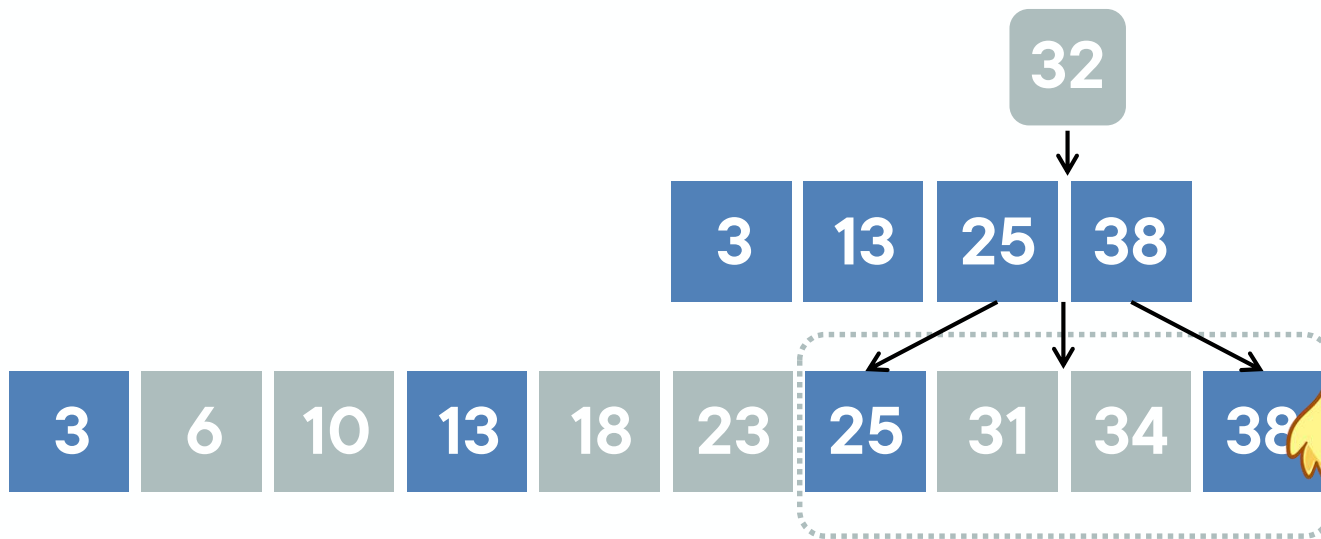
**10 миллионов элементов**

`std::lower_bound?`



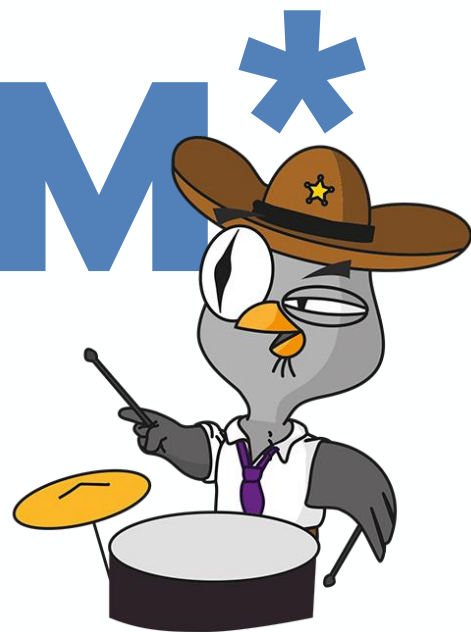
# Мы очень любим оптимизировать

1,5x быстрее



**Зачем создавать свои БД в 2018?**

**Незачем\***





# А зачем тогда этот доклад?

---

- Узнать что-нибудь новое
- Перестать смотреть на базы данных как на черный ящик
- Интересно, как устроены высоконагруженные проекты

# Пишем свою БД

---

- Общение с внешним миром
- Сохранение данных на диск
- Шардируем и реплицируем
- Оптимизации

# Протокол общения (memcached-like)

---

```
printf "set mykey 0 60 4\r\ndata\r\n" | nc localhost 11211
```

STORED

```
printf "get mykey\r\n" | nc localhost 11211
```

VALUE mykey 0 4

data

END

read\_inbound9822647,21524190,1441548



# Протокол общения (RPC/TL)

---

```
messages.readMessages user_id:int peer_id:int up_to_local_id:int = Bool;
```

```
$query = ["_"           => "messages.readMessages",  
          "user_id"     => 9822647,  
          "peer_id"     => 21524190,  
          "up_to_local_id" => 1441548];
```

```
0xfbd90216 0x95e1b7 0x1486ede 0x15ff0c
```



# Протокол общения

---

- Десятки тысяч движков
- Хотим гарантии tcr
- Не хотим держать коннекты ко всем
- Пишем свой!



Программируем, как умеем!

# Бинарный лог событий

---

```
online.getUserStatus user_id : 169
```

```
online.setOnline user_id : 169 is_mobile : true
```

```
online.getUserStatus user_id : 169
```

```
online.setOnline user_id : 182 is_mobile : true
```

```
online.setOnline user_id : 119 is_mobile : false
```

```
online.getUserStatus user_id : 123
```

```
online.getUserStatus user_id : 149
```

```
online.getUserStatus user_id : 127
```

```
online.getUserStatus user_id : 115
```

```
online.setOnline user_id : 134 is_mobile : true
```



```
online.setOnline user_id : 169 is_mobile : true  
online.setOnline user_id : 182 is_mobile : true  
online.setOnline user_id : 119 is_mobile : false  
online.setOnline user_id : 134 is_mobile : true
```

# СНИМОК

---

```
online.setOnline user_id : 136 is_mobile : true
online.setOnline user_id : 142 is_mobile : true
online.setOnline user_id : 192 is_mobile : false
online.setOnline user_id : 191 is_mobile : true
online.setOnline user_id : 191 is_mobile : true
online.setOnline user_id : 142 is_mobile : true
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 142 is_mobile : true
online.setOnline user_id : 136 is_mobile : true
online.setOnline user_id : 192 is_mobile : false
online.setOnline user_id : 191 is_mobile : false
online.setOnline user_id : 192 is_mobile : false
online.setOnline user_id : 136 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 191 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 142 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 136 is_mobile : true
online.setOnline user_id : 191 is_mobile : false
online.setOnline user_id : 136 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 136 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
online.setOnline user_id : 136 is_mobile : true
```

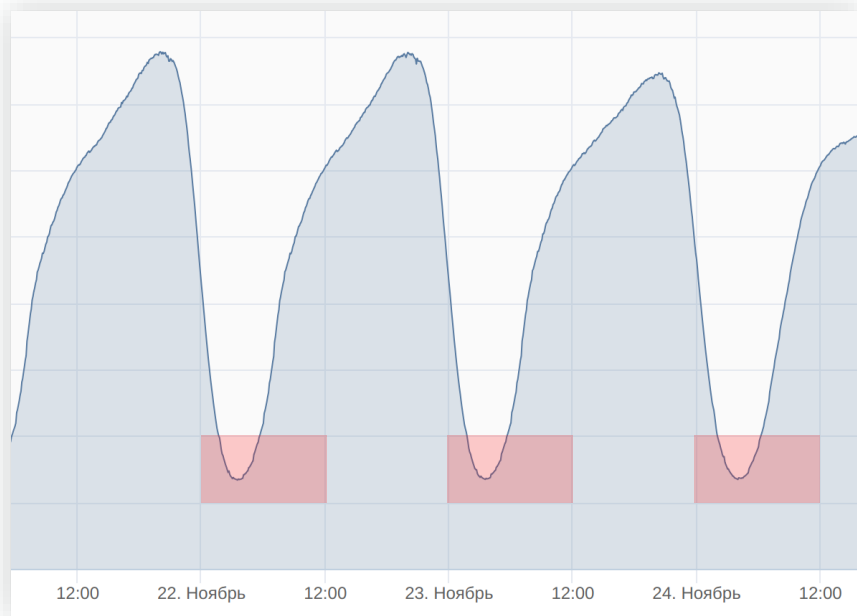


Оставим только нужное, отсортируем

```
online.setOnline user_id : 136 is_mobile : true
online.setOnline user_id : 142 is_mobile : false
online.setOnline user_id : 191 is_mobile : false
online.setOnline user_id : 192 is_mobile : true
```

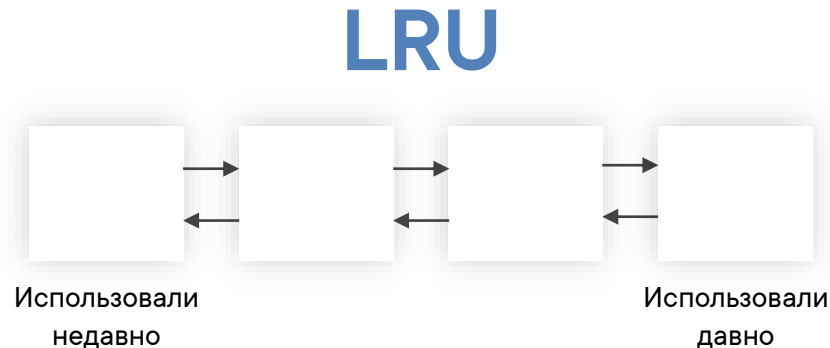
# Когда создавать снимки?

---





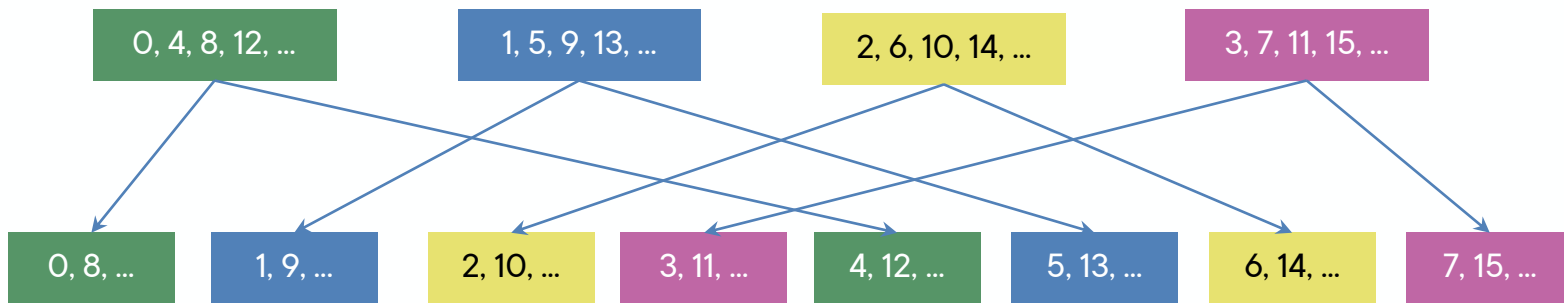
# Метафайлы



# Шардирование

---

- Чем меньше зависимости между данными на разных серверах, тем лучше
- $User\_id \% shards\_num$
- Увеличиваем кластера только в 2 раза



# Оптимизации

---

- Вместо деревьев используем массивы
- Сжатие данных — простой способ получить выигрыш почти бесплатно
- Для каждого движка специфичны

# Выводы

---

- Есть смысл писать свои БД, только если много серверов
- Нужно понимать, чем ваши данные специфичны
- Не нужно смотреть на БД как на черный ящик
- Полезно знать, как устроены высоконагруженные проекты

