



# Мониторинг PostgreSQL в АВИТО.

Дмитрий Вагин, 2017

# О чём это всё?!

Я Дмитрий - разработчик баз данных в Авито.

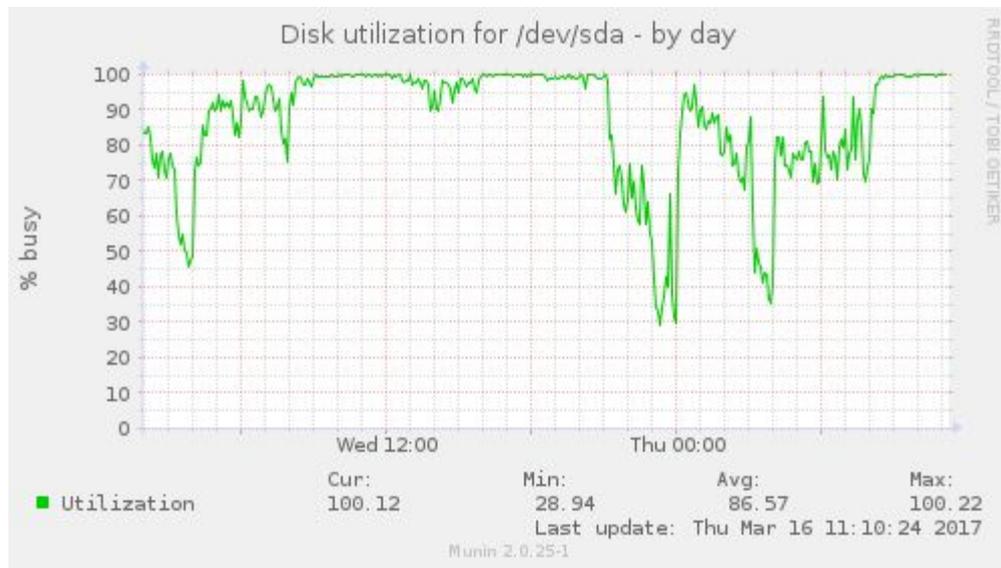
Я расскажу про свой опыт сбора статистики баз данных PostgreSQL.

# Зачем всё это?!

Очень хочется:

- знать что сломалось
- знать почему сломалось
- знать когда сломается в следующий раз

# Как проявляются проблемы



# Куда смотреть?!

*«PostgreSQL's statistics collector is a subsystem that supports collection and reporting of information about server activity.»*

Из pg\_stat\_database мы узнали что с момента последнего сброса статистики

- было совершено 14 миллиардов 066 миллионов 649 тысячи 489 транзакции
- 128 миллиардов 702 миллиона 052 тысячи 186 блоков было прочитано



# Что делать?!

*pg\_stat\_\** и *pg\_statio\_\** практически бесполезны без исторических данных!

```
insert into maintenance.pg_stat_all_tables  
select * from pg_stat_all_tables;
```

680kb - один срез ( $\approx 2k$  таблиц)

$680\text{kb} * 60\text{min} * 24\text{h} \approx 1\text{GB per day}$

1 гигабайт чтобы иметь статистику  
по таблицам в одной базе за один день?!



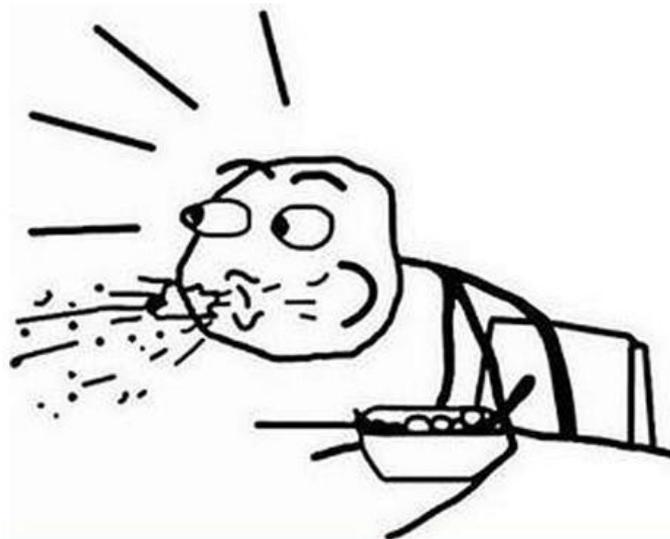
# Что у нас есть...

- есть один большой инстанс (кластер), и 5 стендбаев, один для работы, остальные для резерва.
- есть ещё 4 инстанса поменьше, с кучкой баз разной степени нагруженности и важности + по standby на каждый инстанс
- 3 инстанса для логических реплик разных таблиц и материализованных представлений + по резервному инстансу для каждой из них

# Что у нас есть...

Итого:

- $\approx 15$  инстансов
- $\approx 30$  баз
- $\approx 30000$  таблиц
- $\approx 60000$  индексов
- $\approx 1000000$  различных метрик



# Что делать!?!?!?!?!?!?!?!111111



# Магия



# Как такое сделать?

1. настраиваем brubeck + graphite + grafana
2. пишем крон или демон который раз в минуту или n-секунд забирает данные из `pg_stat_*` и `pg_statio_*` и скидывает в brubeck
3. настраиваем dashboards в графанае
4. PROFIT

# Про крон или демон

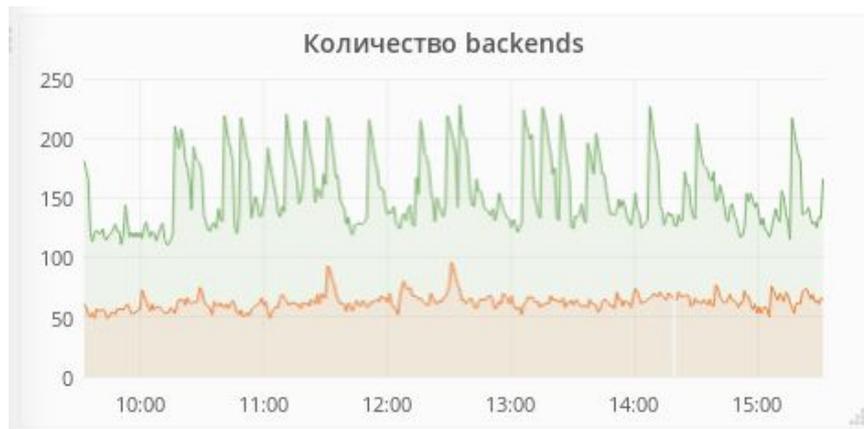
- можно попробовать найти готовое
- можно попробовать plugin collectd-postgresql
- можно сделать своё :)

Важно разделять типы метрик:

- метрика отображающая текущее значение — например `pg_stat_database.num_backends`
- увеличивающийся счётчик — практически все остальные метрики

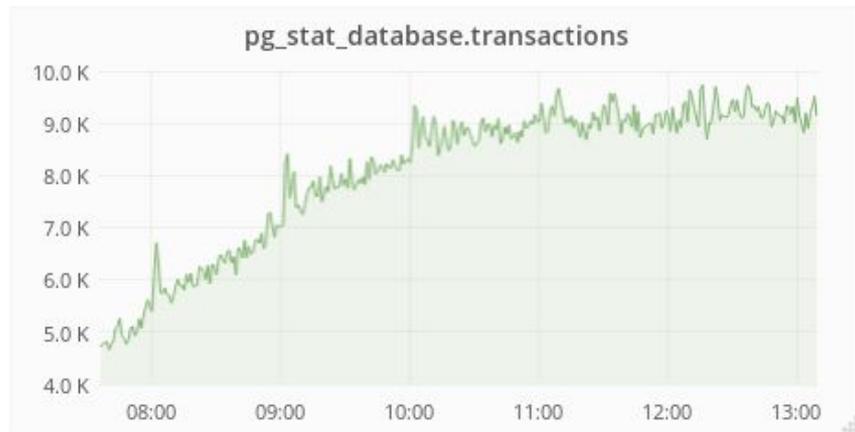
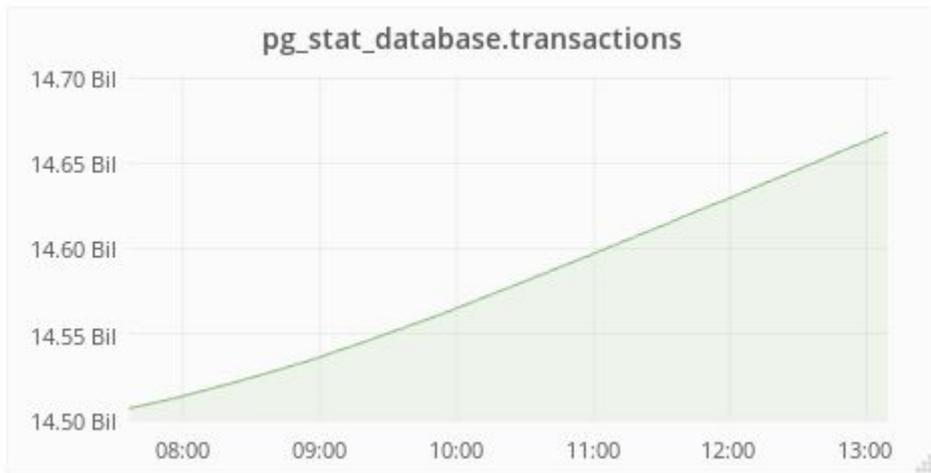
# Про крон или демон

`pg_stat_database.num_backends`



# Про крон или демон

`pg_stat_database.xact_commit`



`perSecond()`

# Про крон или демон

- g - Gauges
- c - Meters
- C - Counters
- h - Histograms
- ms - Timers (in milliseconds)

<https://github.com/github/brubeck>

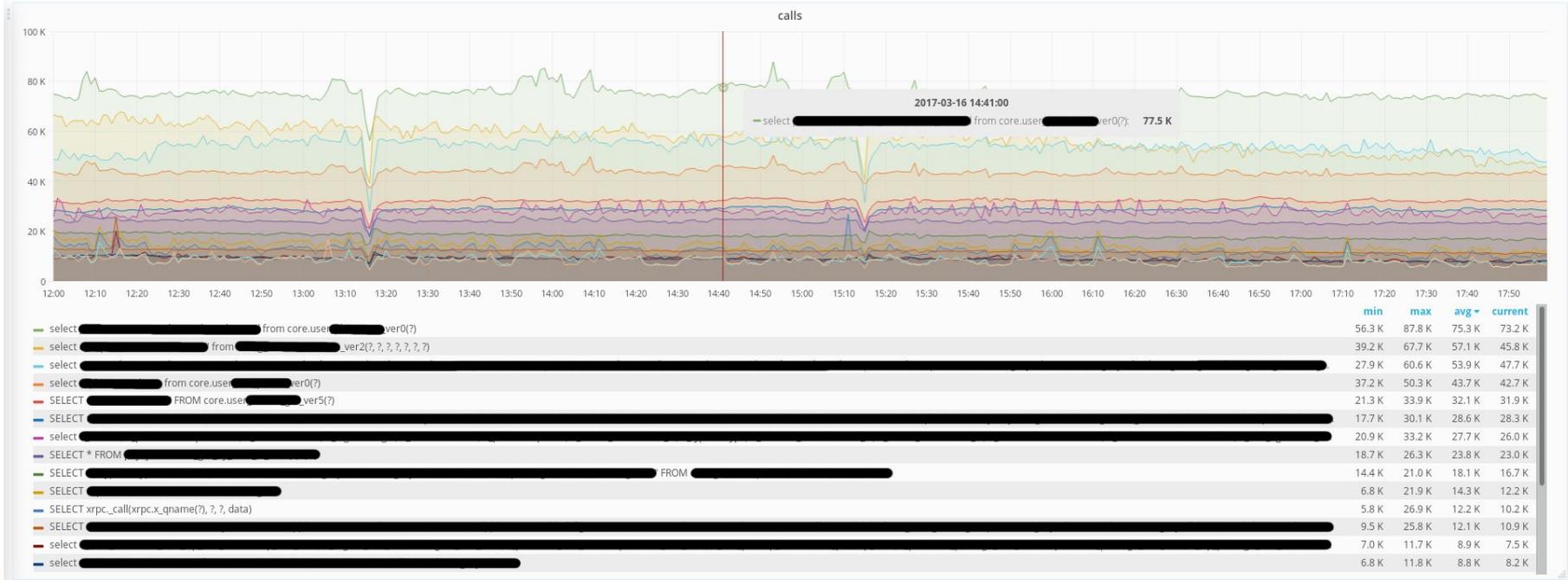
[https://github.com/etsy/statsd/blob/master/docs/metric\\_types.md](https://github.com/etsy/statsd/blob/master/docs/metric_types.md)

Документация противоречивая, поэтому смотрим в исходники

<https://github.com/github/brubeck/blob/master/src/metric.c>

# pg\_stat\_statements

column calls - limit 20 -



# pg\_stat\_statements

В имя метрики невозможно засунуть весь запрос. Поэтому запрос храним в redis'е, а в имени метрики идёт хэш от запроса.

Для получения текста запроса используем самописную функцию в графите, запрос к графиту выглядит так:

```
aliasByHash(aliasByNode(keepLastValue(limit(sortByTotal(metric.namespace.*.$column), $limit), 1), 5), 'sql_query')
```

# pg\_stat\_statements

Из хорошего:

- мы видим самые часто выполняемые запросы
- мы видим запросы которые читают с диска `shared_blks_read`
- можем видеть динамику(!) времени выполнения запроса

Из плохого:

- очень много уникальных запросов, поэтому удаляем старое
- нужно какое-то дополнительное хранилище (redis)
- проблемы с отображением больших запросов в графанах

# PG Metricus

Надо смотреть за динамикой времени выполнения куска большой plpgsql процедуры.

- добавить *raise notice* и парсить логи
- складывать в таблицу и анализировать её
- а что если..... кидать в канал через механизм NOTIFY, а оттуда отправлять в графит.....

# PG Metricus

```
x1 = clock_timestamp();  
.... your long running query  
x2 = clock_timestamp();
```

```
perform pg_notify('test', format(E'%s.%s:%s|%s\n',  
    'db.sql.metric',  
    'get_val_hstore_duration',  
    extract(milliseconds from (x2 - x1))::bigint::text,  
    'ms' ));
```

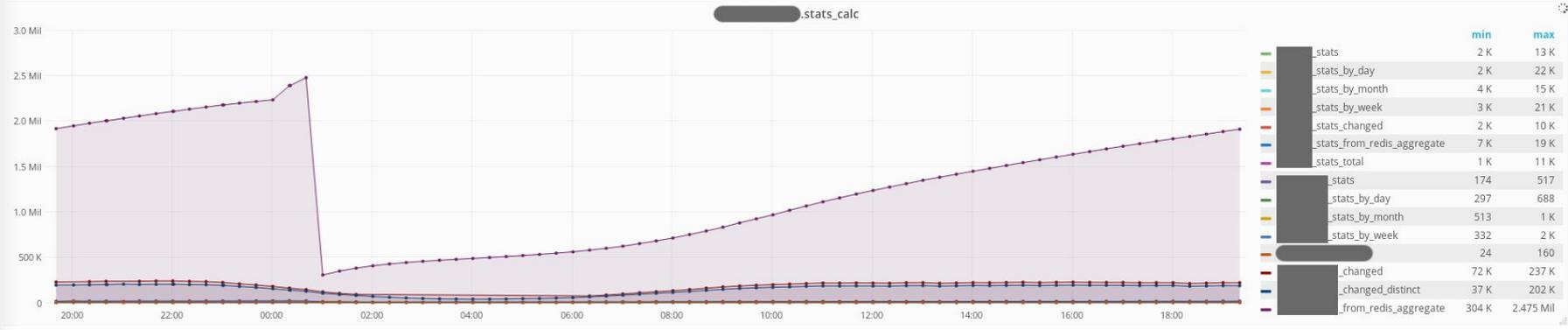
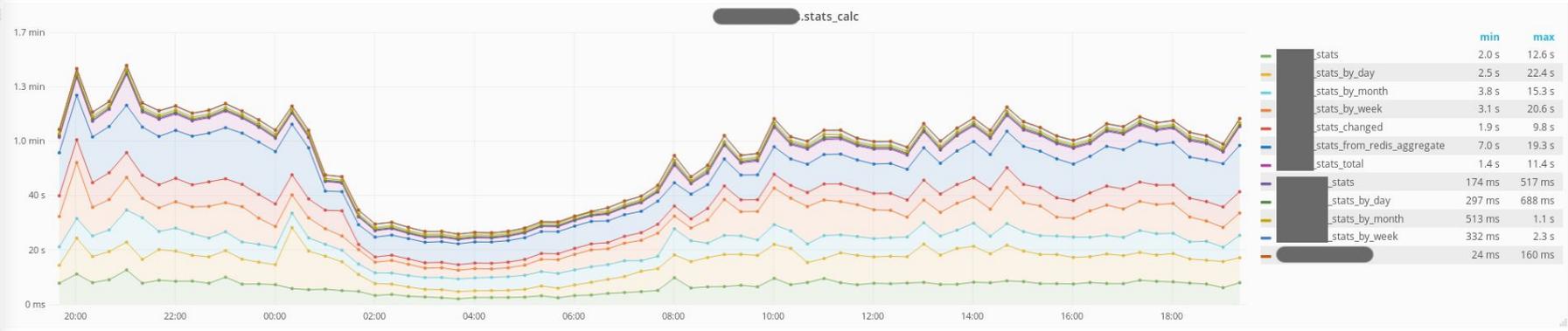
# PG Metricus

```
curs = conn.cursor()
curs.execute("LISTEN {0};".format(args.channel))

addr = (args.socket_host, args.socket_port)
sock = socket(AF_INET, SOCK_DGRAM)

while not checkStop():
    if select.select([conn], [], [], 5) != ([], [], []):
        conn.poll()
        while conn.notifies:
            notification = conn.notifies.pop(0)
            sock.sendto(notification.payload.encode('utf-8'), addr)
```

# PG Metrics



# PG Metricus

Подробнее можно ознакомиться в статье на хабре

<https://habrahabr.ru/company/avito/blog/323900/>

[https://github.com/avito-tech/pg\\_metricus\\_python](https://github.com/avito-tech/pg_metricus_python)

# Алерты

Мы используем

- grafana 4+ имеет встроенный алертинг
- moira

Нотификации в Slack.

## #dba-monitoring

★ | 👤 15 | 🔒 0 | [Add a topic](#)



**incoming-webhook** APP 13:40

[Alerting] avito stats unblocked items alert  
unblocked\_items  
786.75

Grafana v4.2.0-beta1 | Yesterday at 13:40

[OK] avito stats unblocked items alert

Grafana v4.2.0-beta1 | Yesterday at 13:43



**incoming-webhook** APP 14:16

[Alerting] avito stats removed users alert  
removed\_users  
288.5

Grafana v4.2.0-beta1 | Yesterday at 14:16

[OK] avito stats removed users alert

Grafana v4.2.0-beta1 | Yesterday at 14:17



**incoming-webhook** APP 14:23

[OK] Errors refund or write off alert

Grafana v4.2.0-beta1 | Yesterday at 14:23



**incoming-webhook** APP 16:07



Message #dba-monitoring

Ну вроде всё...

Вопросы?

Предложения?

Комментарии?

