

# PostgreSQL

## avito.ru

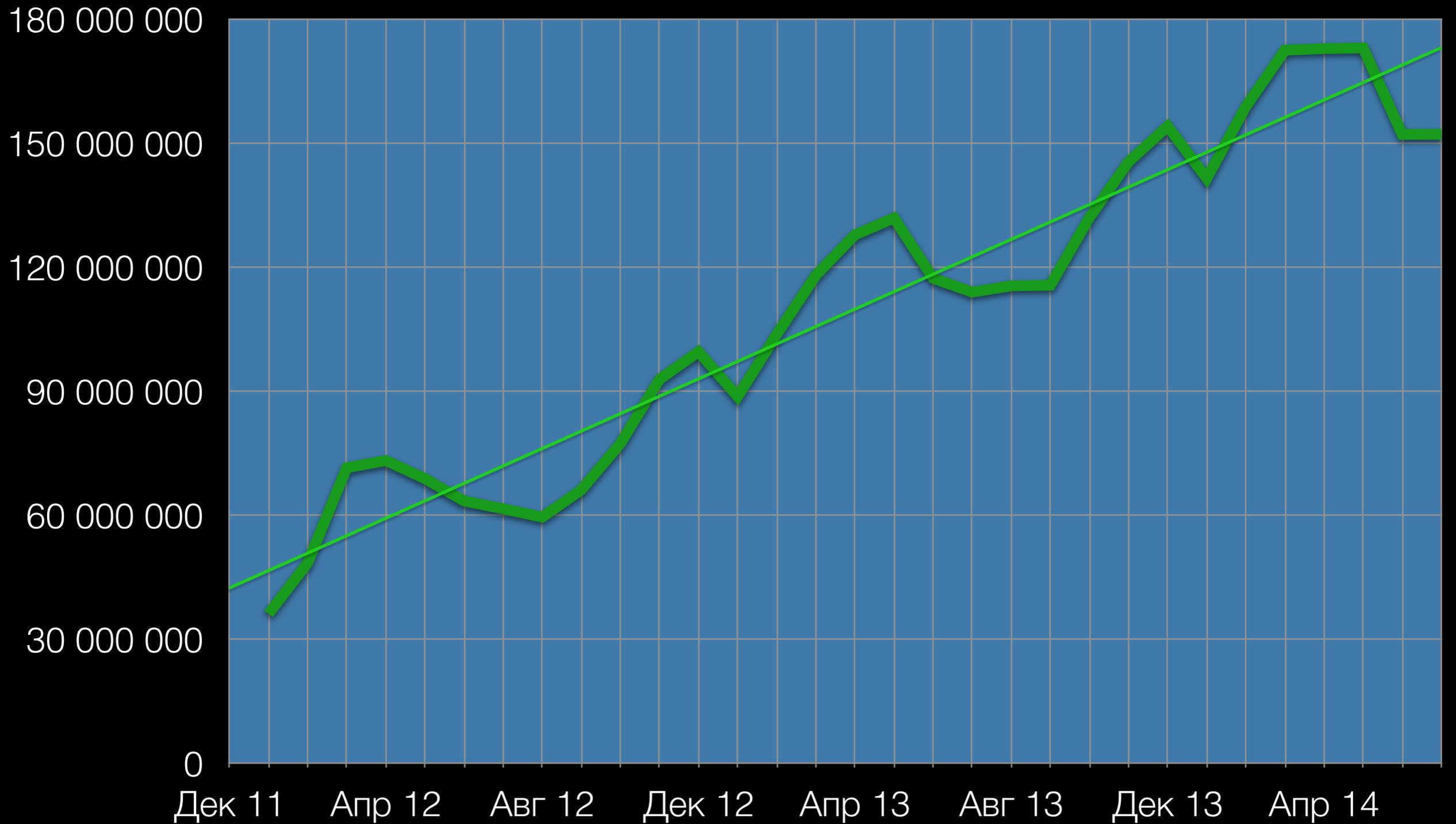
Михаил Тюрин  
DBA, руководитель группы

pgday.ru  
2014

# avito.ru

- самая большая доска объявлений в России (Европе)
- не только ru и не только сайт
- быстрый рост: до 250M hits/day за 5 лет
- до 1M новых объявлений в день
- ежемесячная аудитория 50M (5M/day)

# pageviews



# О ЧЕМ ДОКЛАД

- ретроспектива
- место в архитектуре
- параметры эксплуатации
- ключевые решения
- команда разработки
- текущие вызовы

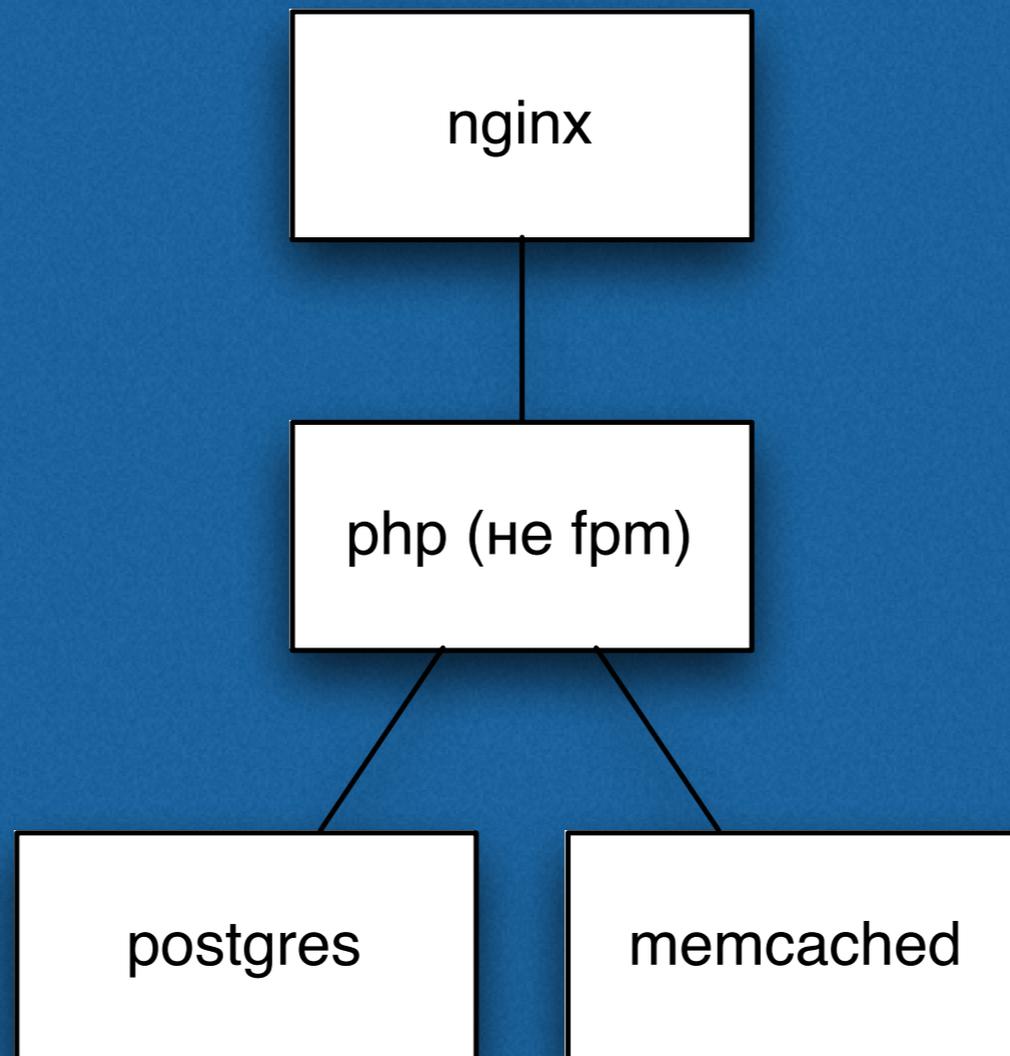
# начало

- сайт ru
- web backoffice
- ~10 человек вся компания // включая всех

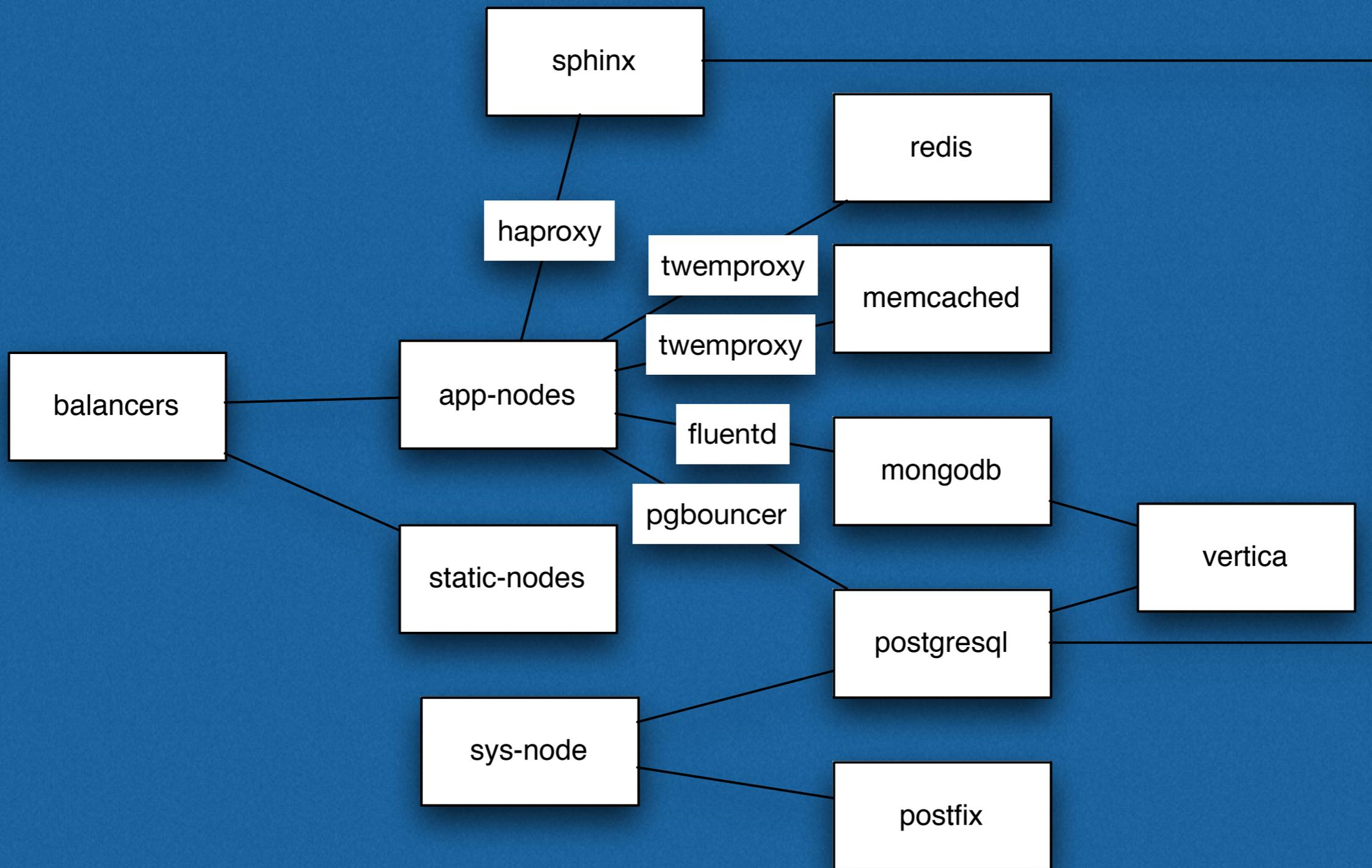
# а сейчас

- сайты: Россия, Украина, Марокко, Египет
- backoffice: support, moderation, bi, ...
- ! mobile: web, iOS, android, WP
- api: mobile apps, partners, ...
- 500+ человек // кто все эти люди?!

# как было



# как получилось



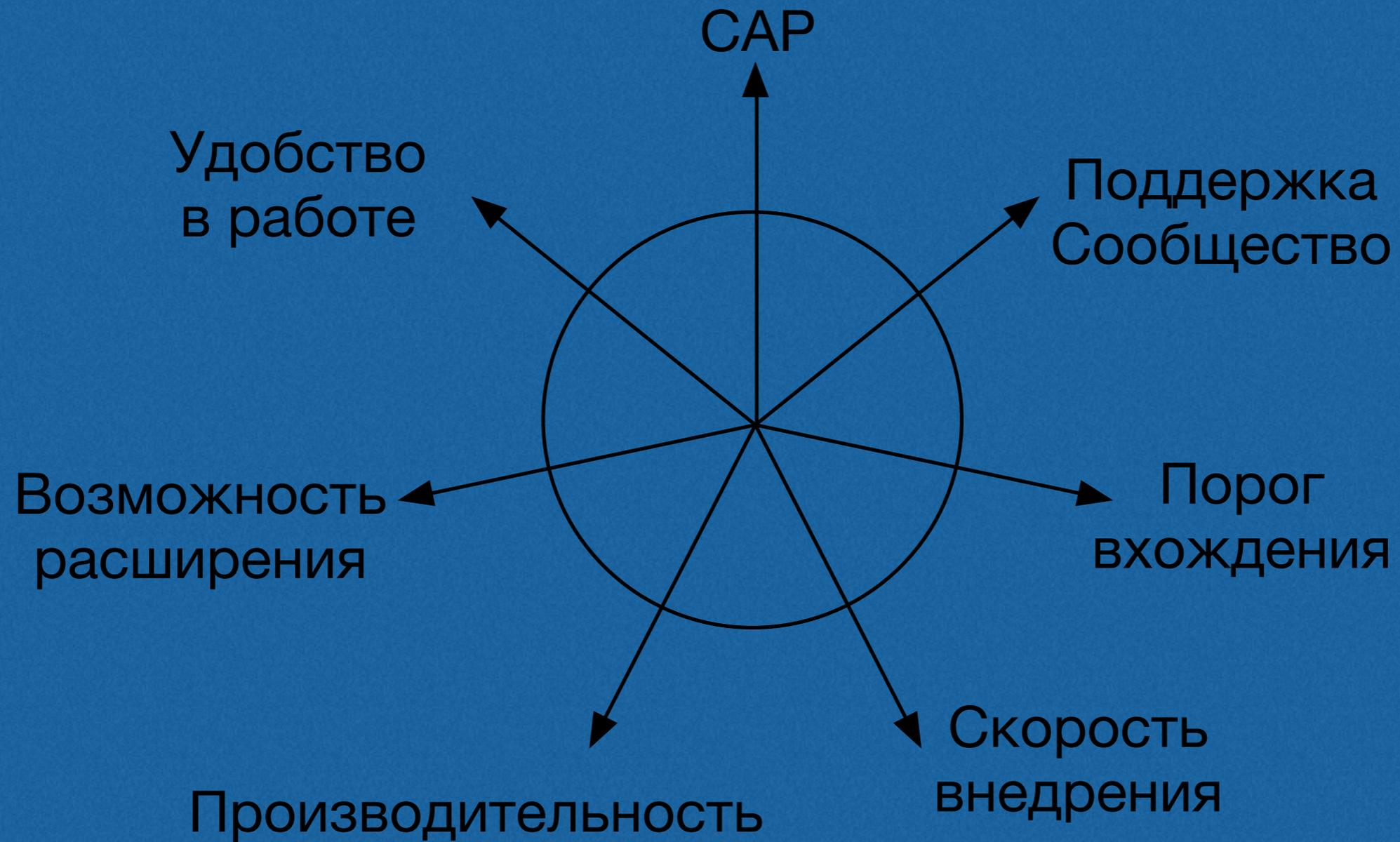
# еще раз

- app
- sys
- sphinx
- vertica
- postgres
- redis — дампится в базу

# php

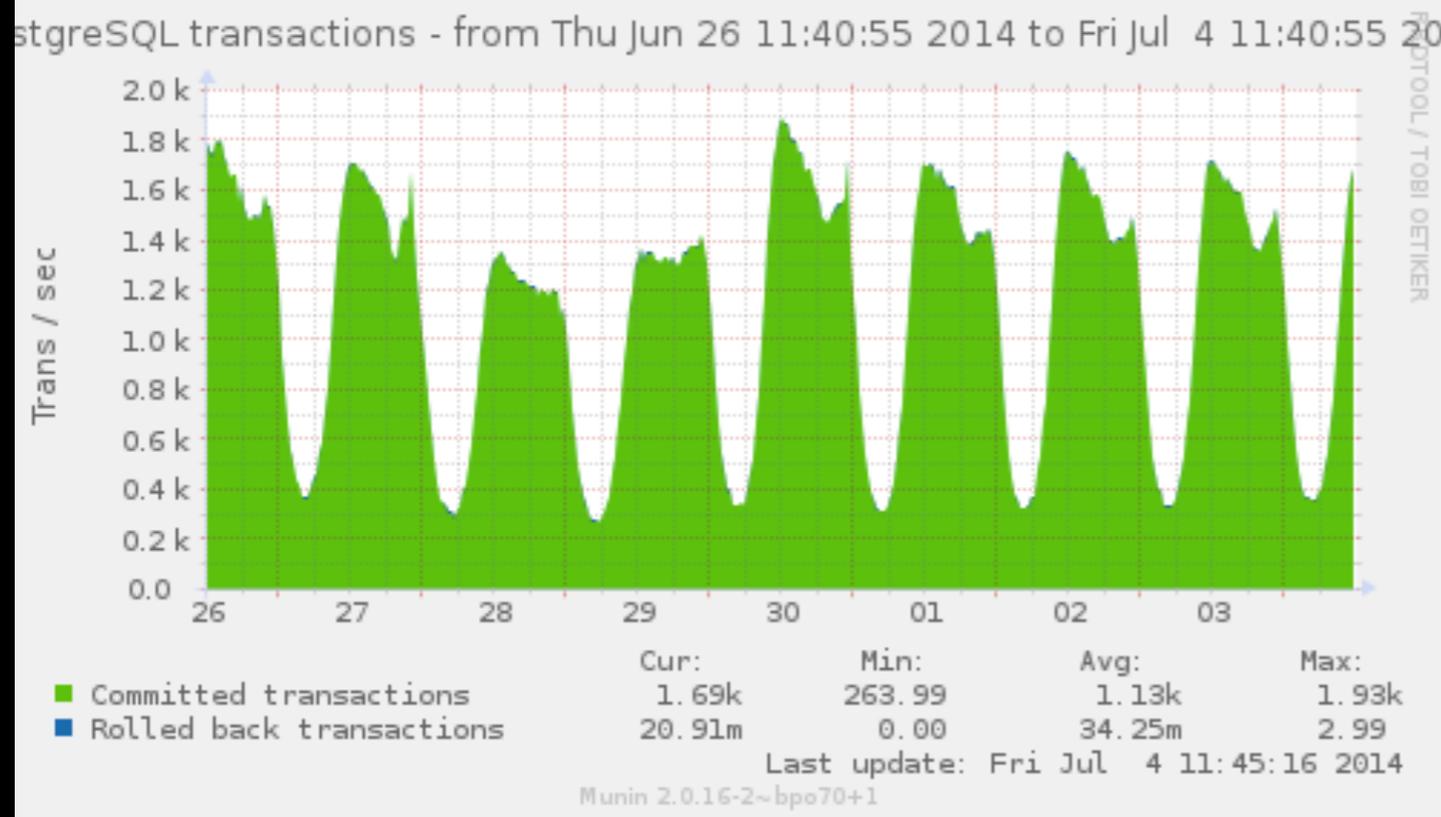
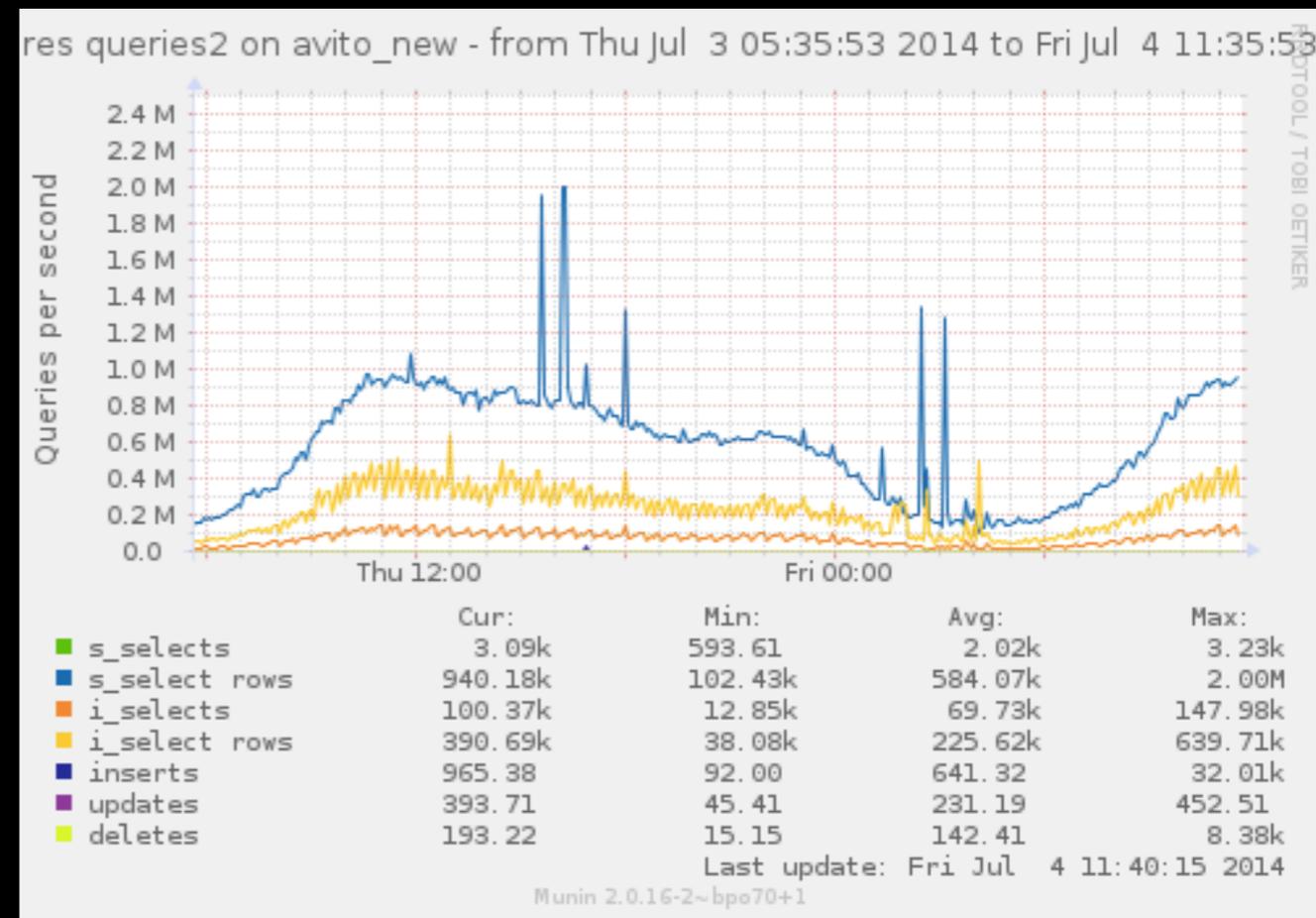
- pgbouncer в транзакшин режиме
- pdo с патчем
- свой класс
- ленивые коннекты и транзакция
- управление транзакциями
- OnCommit

# как такое получилось



# хайлоад 1

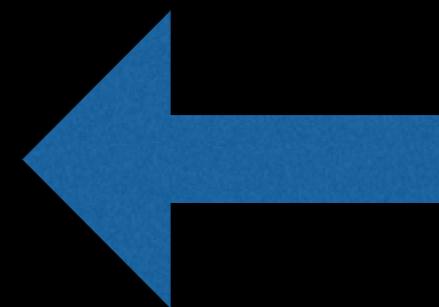
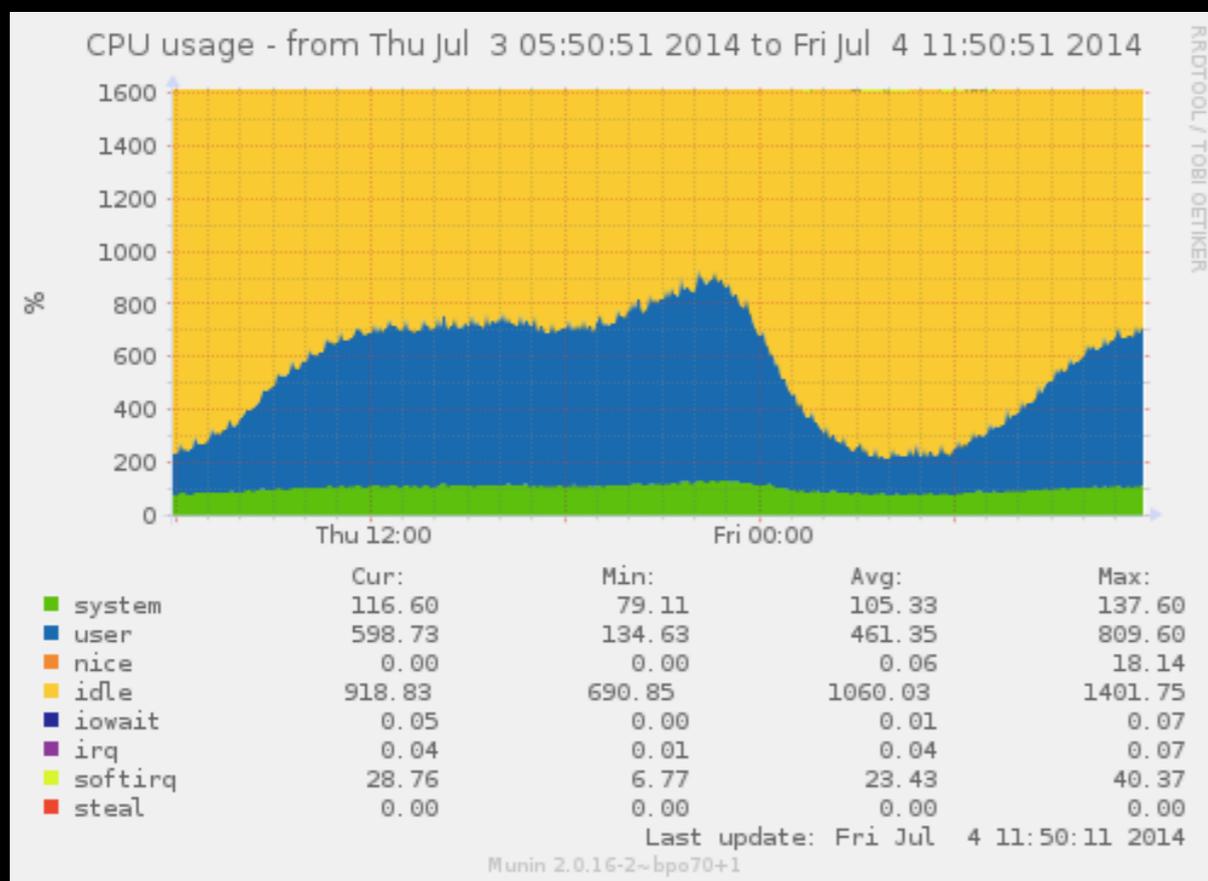
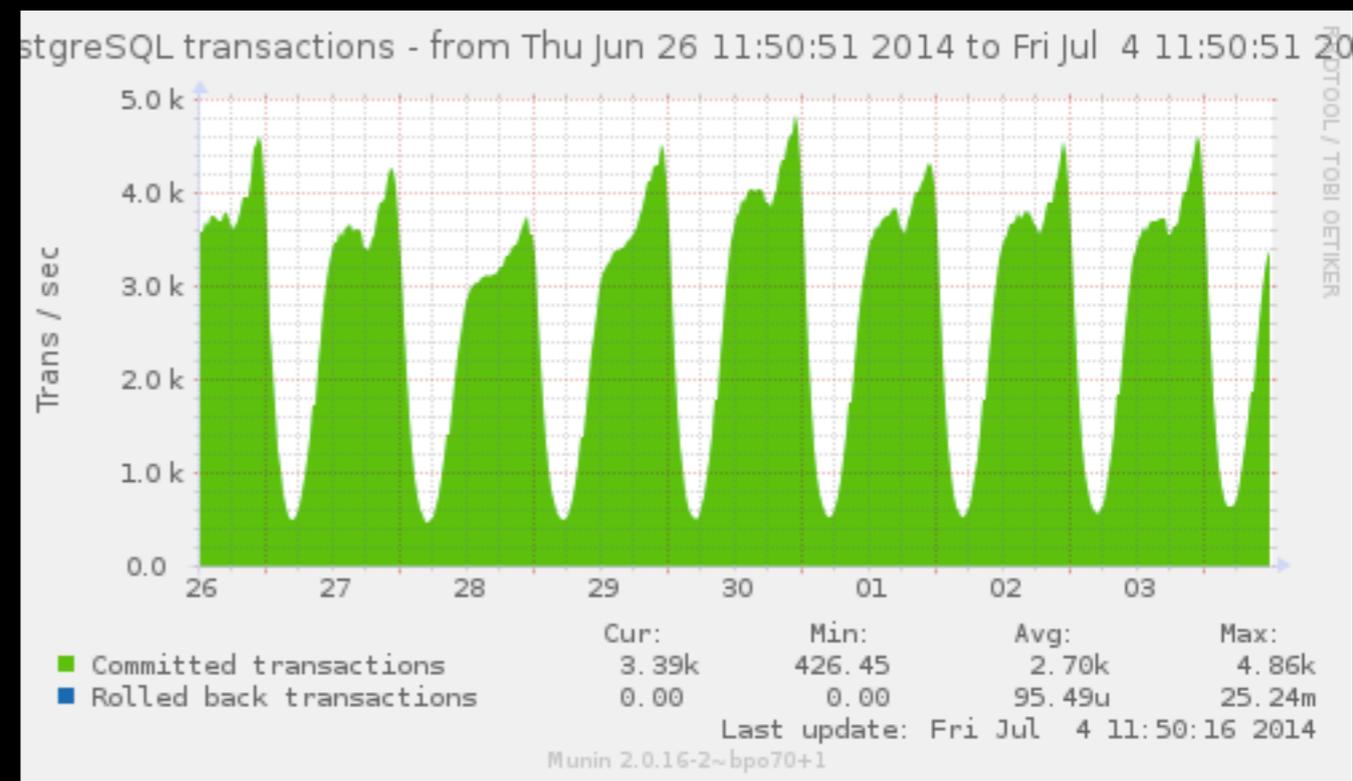
- IUD мастера 1+K/s



- transactions master 2K/s

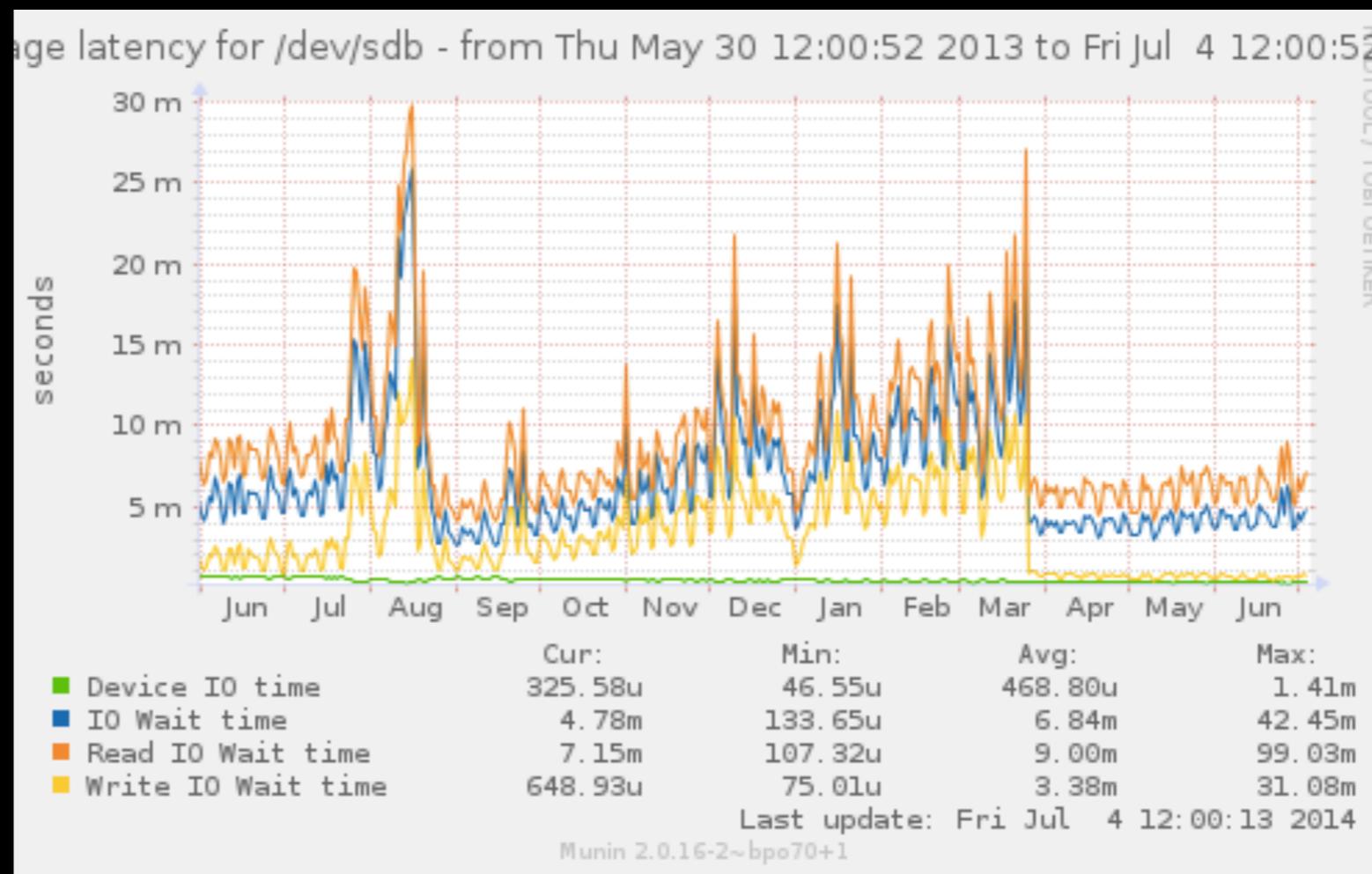
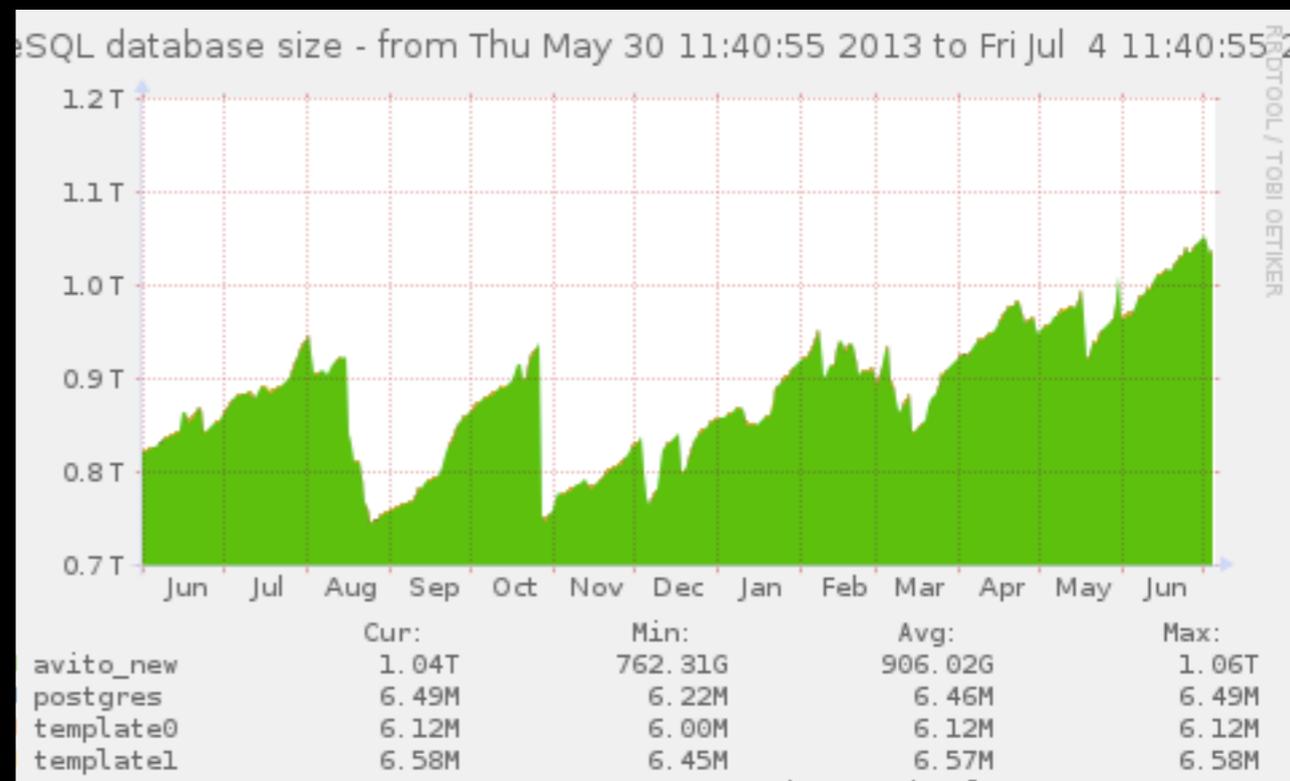
# хайлоад 2

- transactions репча 5K/s
- сру - 6-ядер



# хайлоад 3

- размер мастер 1+ТВ
- латенси мастера - ОК
- запись на мастер - avg 40MB/s
- iops - 2-20К



# среда

- ( php + pdo => pgbouncer ) ==> ( pgbouncer => postgres )
- ! postgres **92** // pg\_upgrade: 83 > 84 > 90 >> 92 >> ...
- >> ? 94/95: fdw/json
- modules: **hstore**, intarray, pg\_buffercache, dblink, auto\_explain, plproxy, pg\_stat\_statements, pgstattuple
- **plpgsql**/sql: бизнес логика
- plperl: csv, files; plpythonu: json, zip

# далее

- python + psycopg2; perl + DBI; bash + psql
- **skytools2**: londiste, pgq, php/python consumer, php daemon
- **munin** + собственные плагины («системные» / «продуктовые»)
- собственные скрипты
  - архив: **pitr** + pbzip2 (4 потока)
  - myvac, truncator, optimazer, rotator

# ПОДСИСТЕМЫ

- сайт
  - мастер базы — синхронная запись (+очередь)
  - wal реплика основного мастера — чтение
  - rgq реплика «активных айтимов» — чтение + индексация
  - rgq платежная реплика — чтение
- хранилище (архив внутри) — 8 нод по 2 базы

# далее

- индексация сайта / индексация админки
- экспорт в dwh
- хгрсд: денормализация, геокодер, ...
- ! архив и резервирование ВСЕГО

# repca

- реплика активных айтимов
- матвью — активные айтимы
- Iondiste с мастера на реплику // PGQ!
- реплика в шаред буферах — ОК
- резервирование — ! два коснумера на одну очередь (knowhow)

# repca — indexer

- sphinx
- ! pause replication
- 10 indexers
- uftp + rsync
- monitoring
- ТОЖЕ два

# archive + (!hot)standby

- pitr
- в nfs по ssh через standby
- сжатие валов
- pg\_basebackup
- 4 дня глубина
- pg\_archivecleanup

# linux + hardware

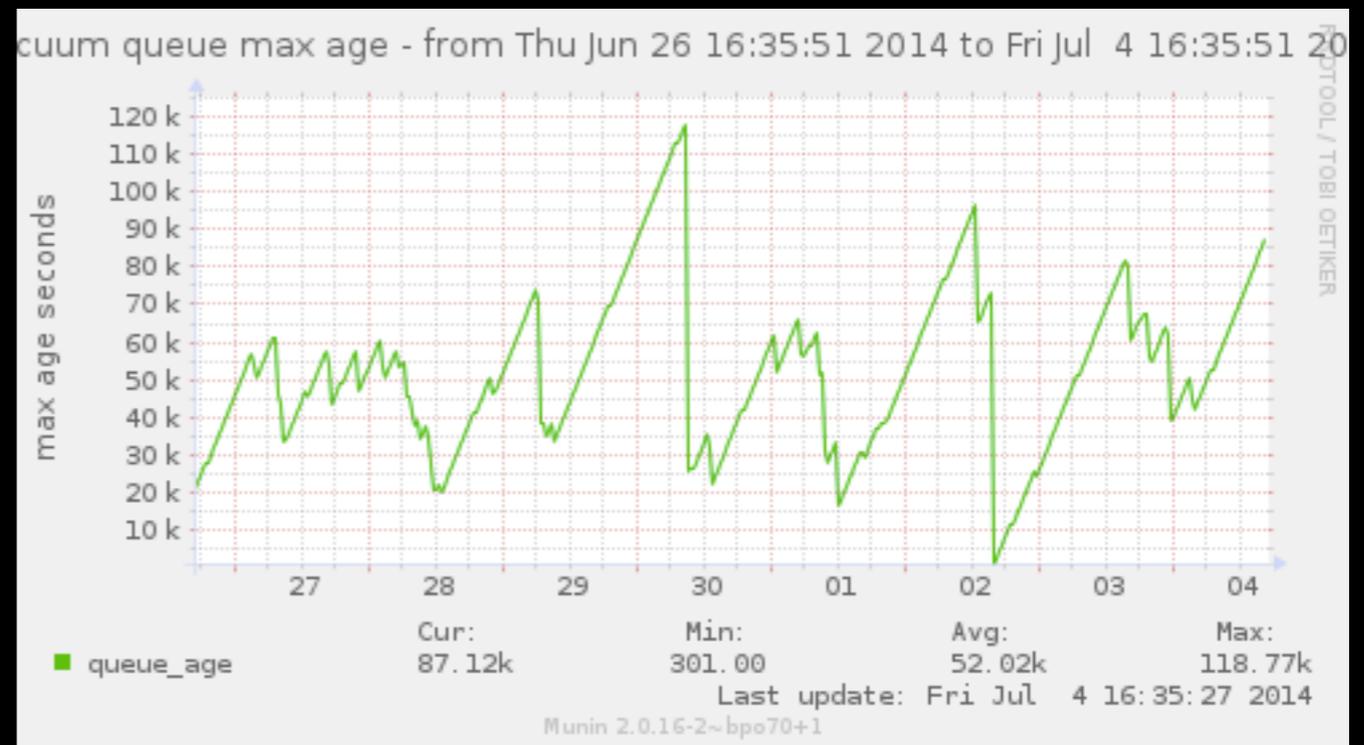
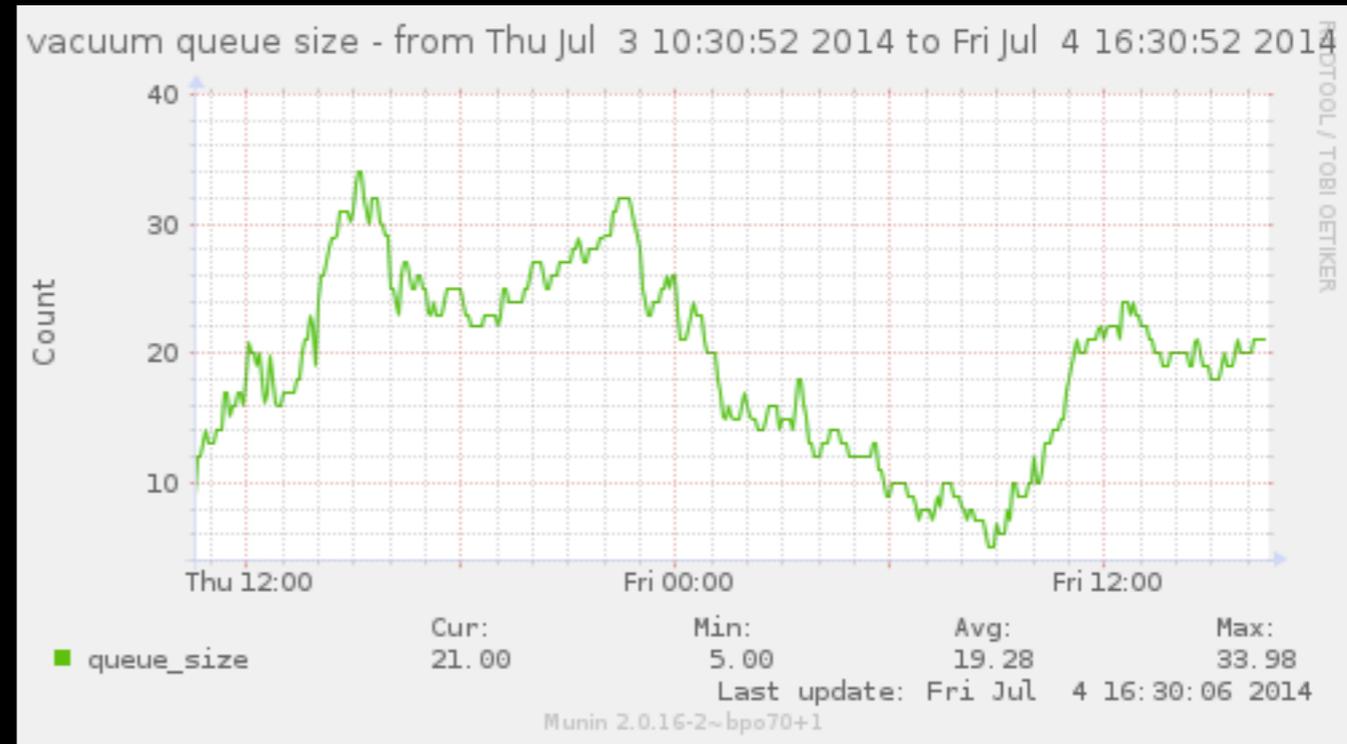
- linux debian 2.6.32 - 3.2. ...
- raid10 + bbu
- numa off + interleaving on
- swap off
- ext4: noatime, wo barriers
- small read ahead
- поор планировщик io

# postgresql.conf

- shared buffers / work mem / effective cache size
- max connections
- checkpoint / bgwriter
- autovacuum !!! on
- sync commit off
- log collector to tmpfs / what to log
- stat directory to tmpfs
- ! fsync / tablespace => tmpfs — во дают ребята

# VACUUM

- ON
- наша  
новая  
мерилка



# КОМАНДА

- 7+ разработчиков + 0,5 админа
- ТИМЛИД
- dba + sql
- system adm + monitoring
- php
- python
- bash + perl
- puppet

# ВЫЗОВЫ: архив

- база растет просто от времени
- это поражает воображение
- надо иметь архив
- есть варианты
- вопрос согласованности

# ВЫЗОВЫ: ВЫЗОВ ПРИЛОЖЕНИЯ ИЗ БАЗЫ

- движок «асинхронной» обработки
- примеры: инвалидация кеша / апдейт счетчиков / продуктовые логи
- очереди-транспорт
- pgq php consumer
- МОЖНО «ОЖИВИТЬ web»

# ВЫЗОВЫ: dev db

- текущая схема
  - бекап на дев сервере
  - ежедневные дампы схем и справочников
  - выкатка на живой базе — опасно
  - но очень просто
- растёт и команда
- контроль изменений
- цель: выкатка базы тоже по кнопке
- сложно: нужен второй авито

# SSD

- уже не дорого
- много места и
- x20 к read iops (вау!) // запись точно не хуже

*Sergey Burladyan added a comment - 10/июн/14 20:59*

*сравнение с HDD, random read:*

**HDD: 54 000 блоков: 76 секунд      0.00141**

**SDD: 100 000 блоков: 8 секунд      0.00008**

- внедряем — уже пол года проигрывает валы и гоняет запросы
- trim — не трогать дефолты

# ВМЕСТО ЗАКЛЮЧЕНИЯ

- сделали на postgres топ5 сайт
- после яндекс-вк-ок-мейл
- суп из топора

*«Набор фич и расширений, легендарная надёжность PostgreSQL, наличие встроенной репликации, средств резервирования и архивирования — весь потенциал нашёл своё воплощение, а наличие открытого профессионального комьюнити не оставляет шансов к неэффективной реализации.»*

Спасибо!

**Вопросы**

**[mtyurin@avito.ru](mailto:mtyurin@avito.ru)**