



# Система сбора логов с кластера kubernetes

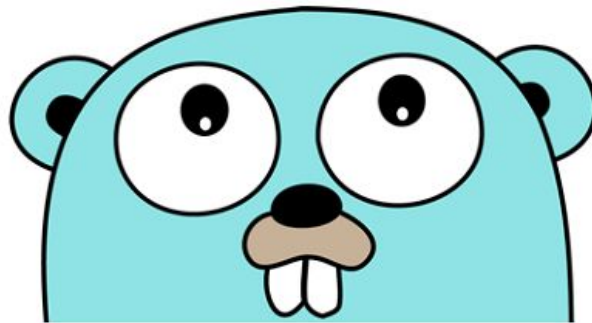
Антон Галицын

**Tinkoff.ru**

# О чем доклад?



1. Кратко о Kubernetes
2. Задача сбора логов
3. Сбор логов у docker контейнеров
4. Что сделали мы?
5. В конце про Go



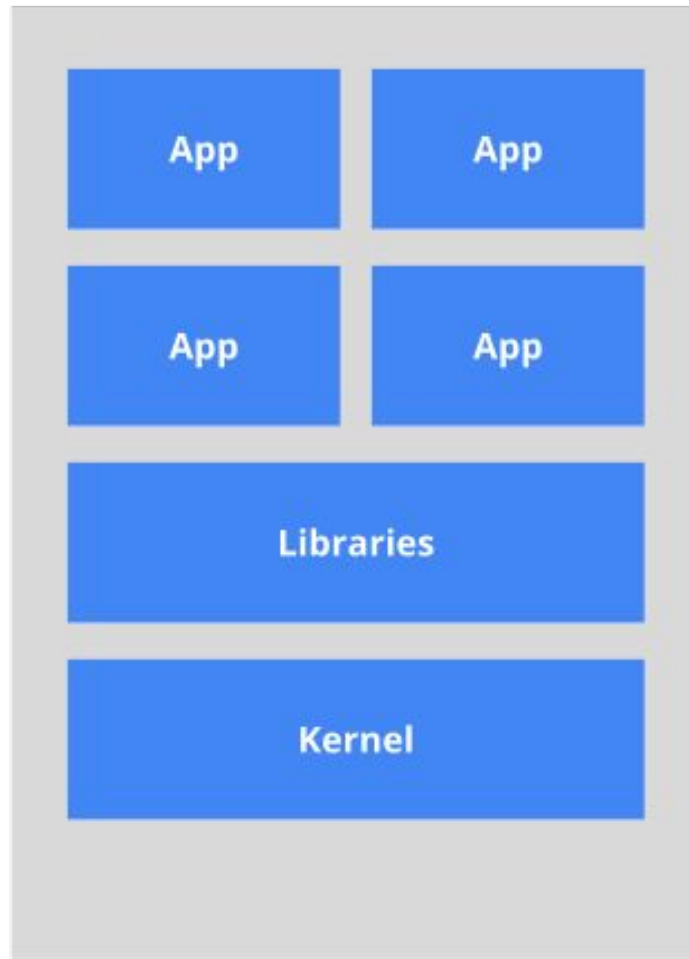
# Kubernetes ;TLDR



# Контейнеры

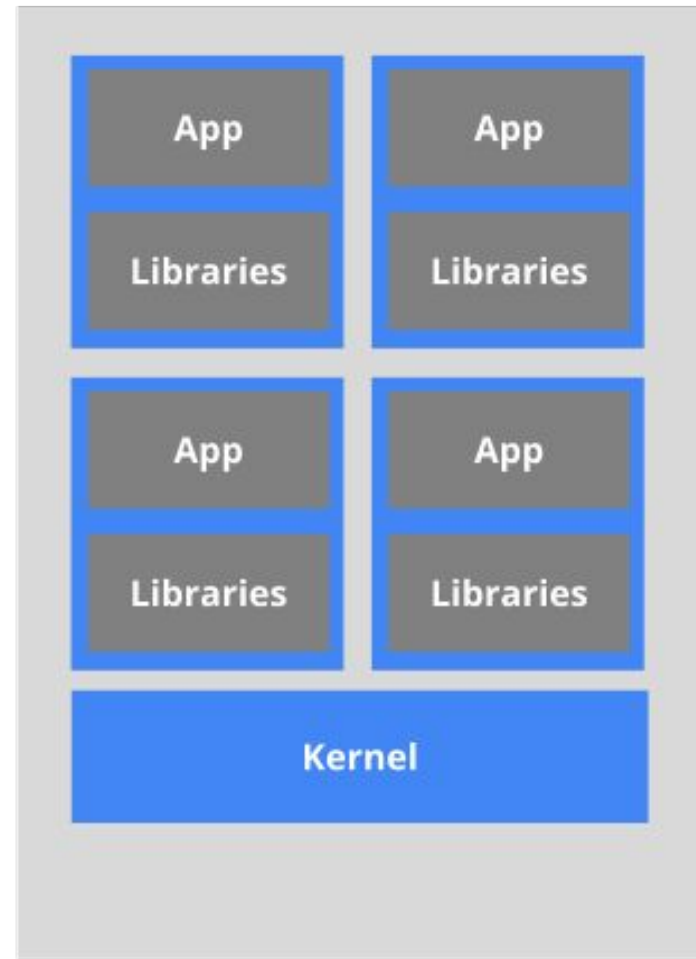


**The old way:** Applications on host



*Heavyweight, non-portable  
Relies on OS package manager*

**The new way:** Deploy containers



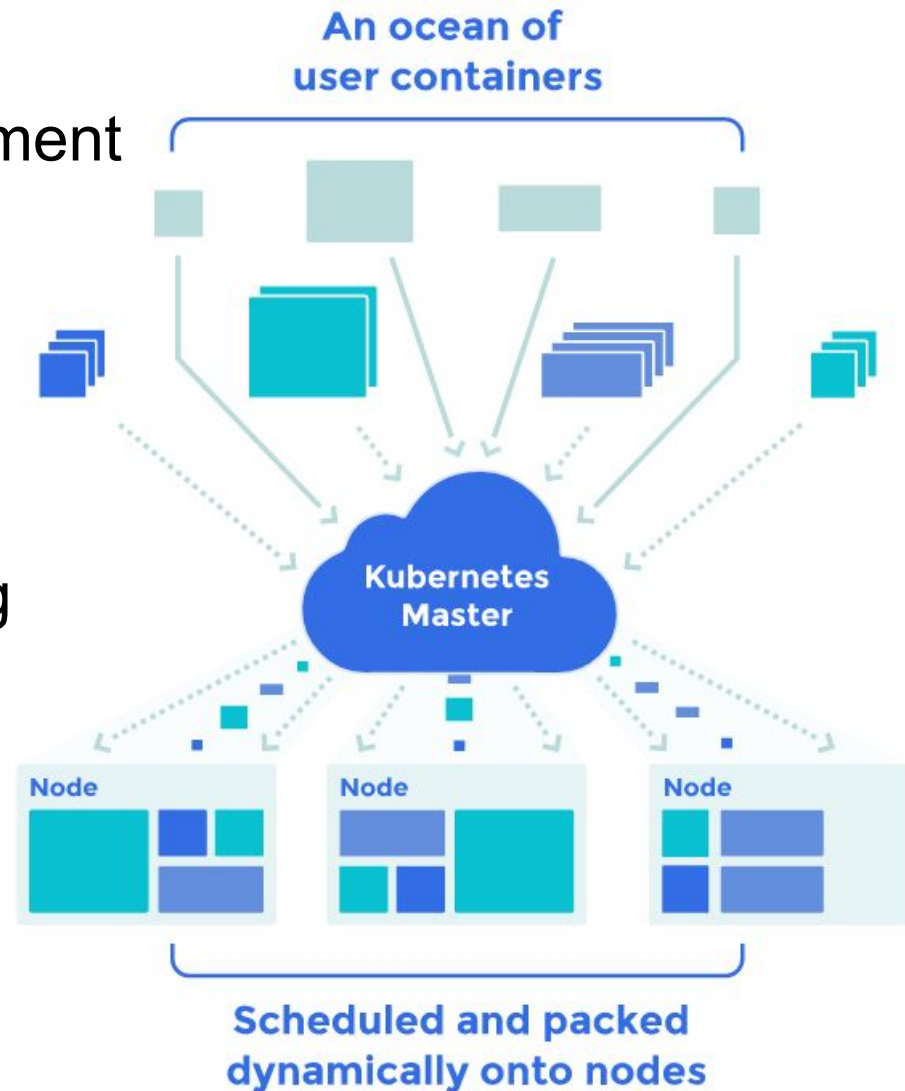
*Small and fast, portable  
Uses OS-level virtualization*

1. Изоляция - Linux namespaces  
сеть, юзеры, пространство процессов, диск и тд
2. Ограничение действий - Linux capabilities,  
apparmor / selinux, seccomp
3. Ограничение ресурсов - Linux cgroups  
Память, CPU и тд.



## Container cluster management

- Orchestration
- Deployment
- Scaling
- Load balancing
- Logging and monitoring



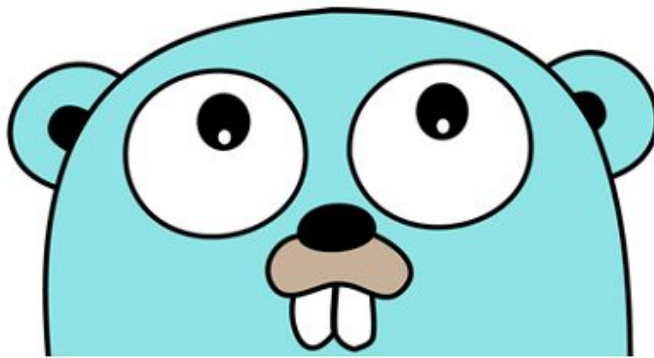


# Система сбора логов

# От кого можно собрать логи?



- Балансировщики
- Базы данных
- Приложения
- Kernel
- и тд

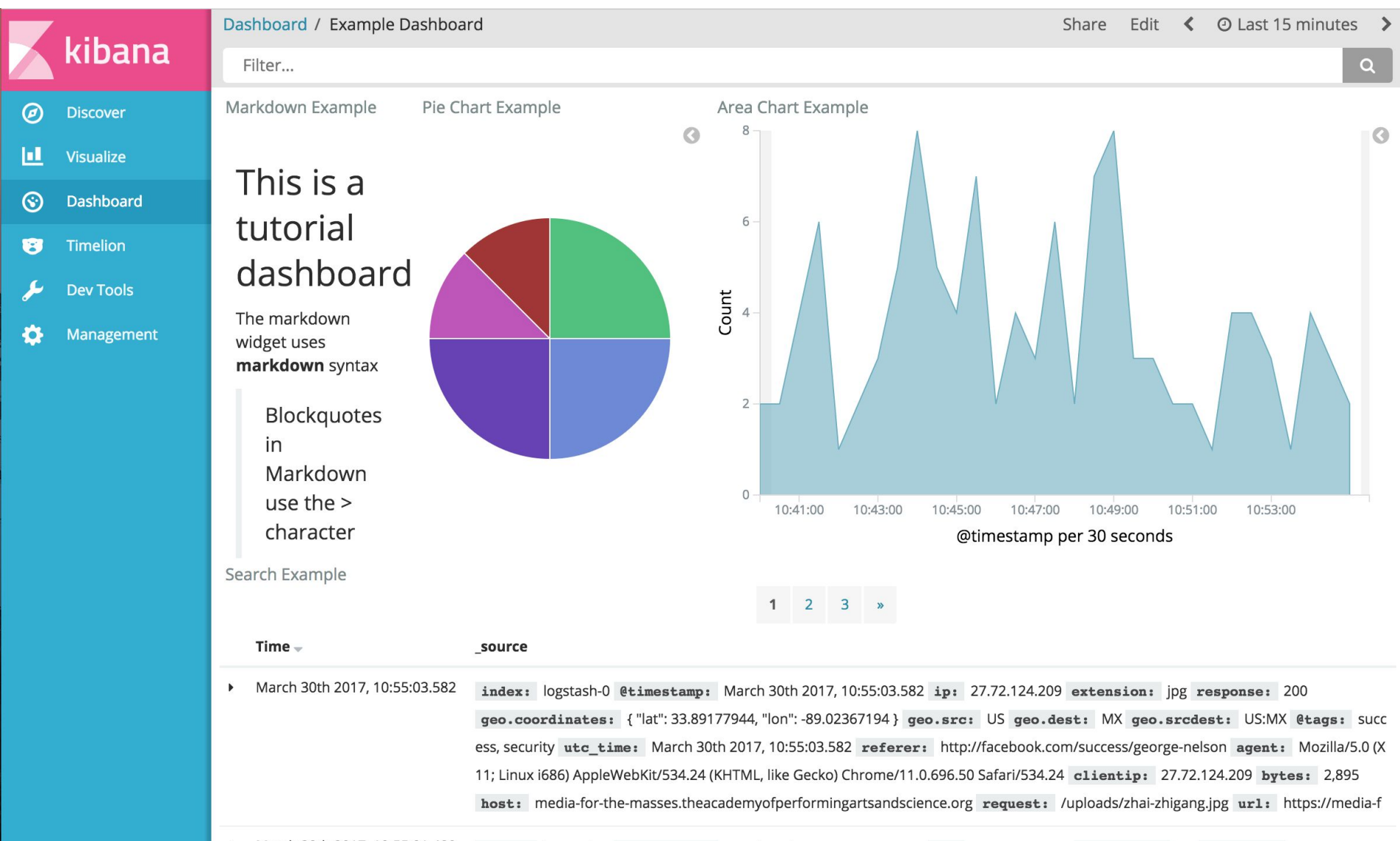




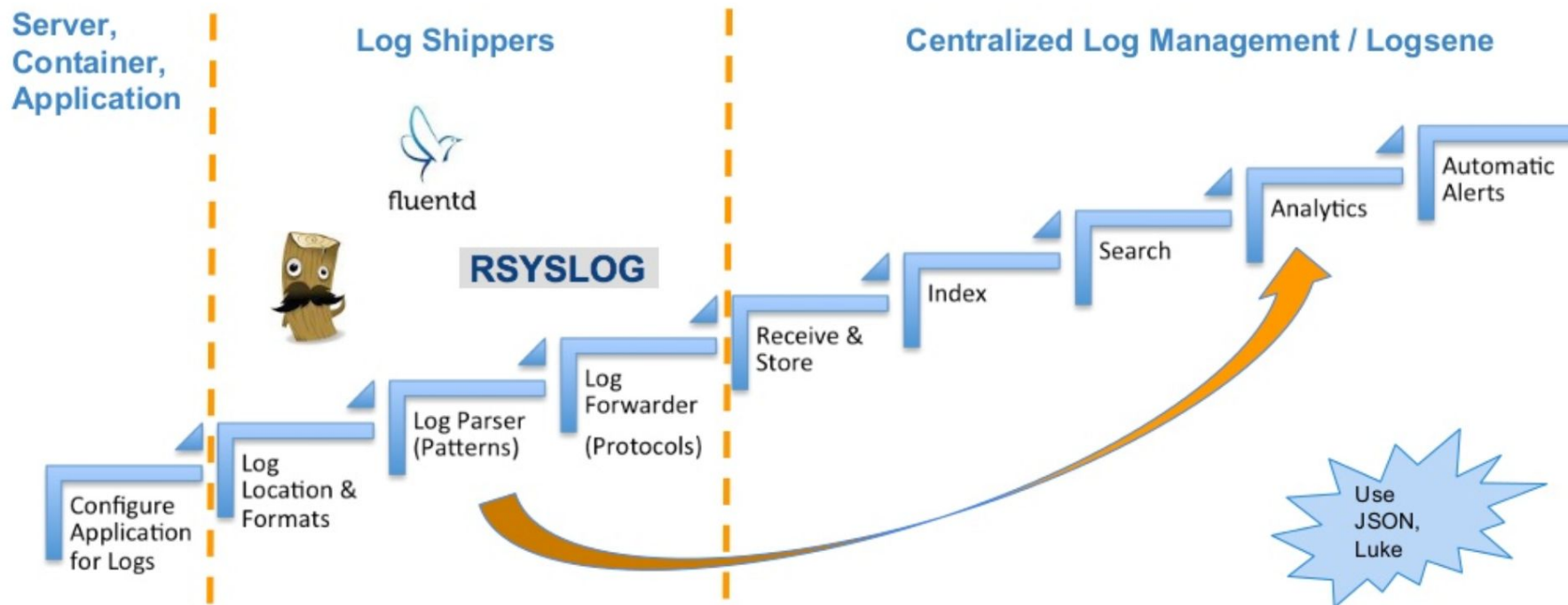


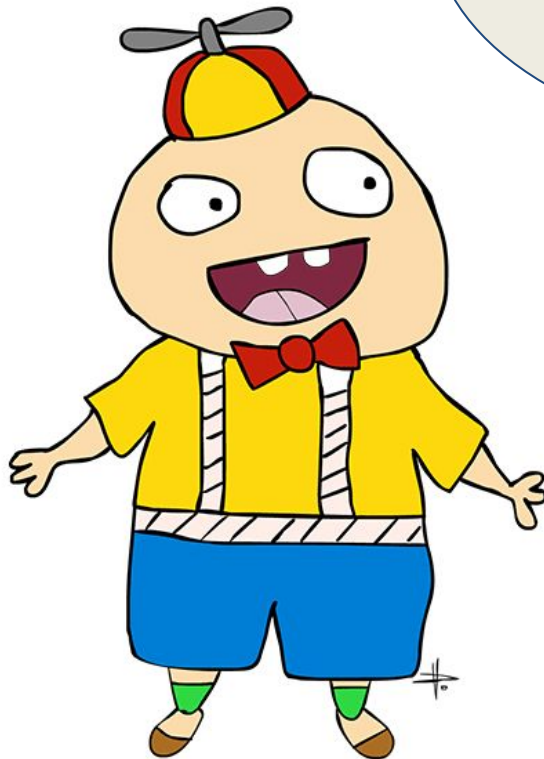
- Поиск
- Аналитика
- Дашборды с агрегатами
- Алерты

# Пример



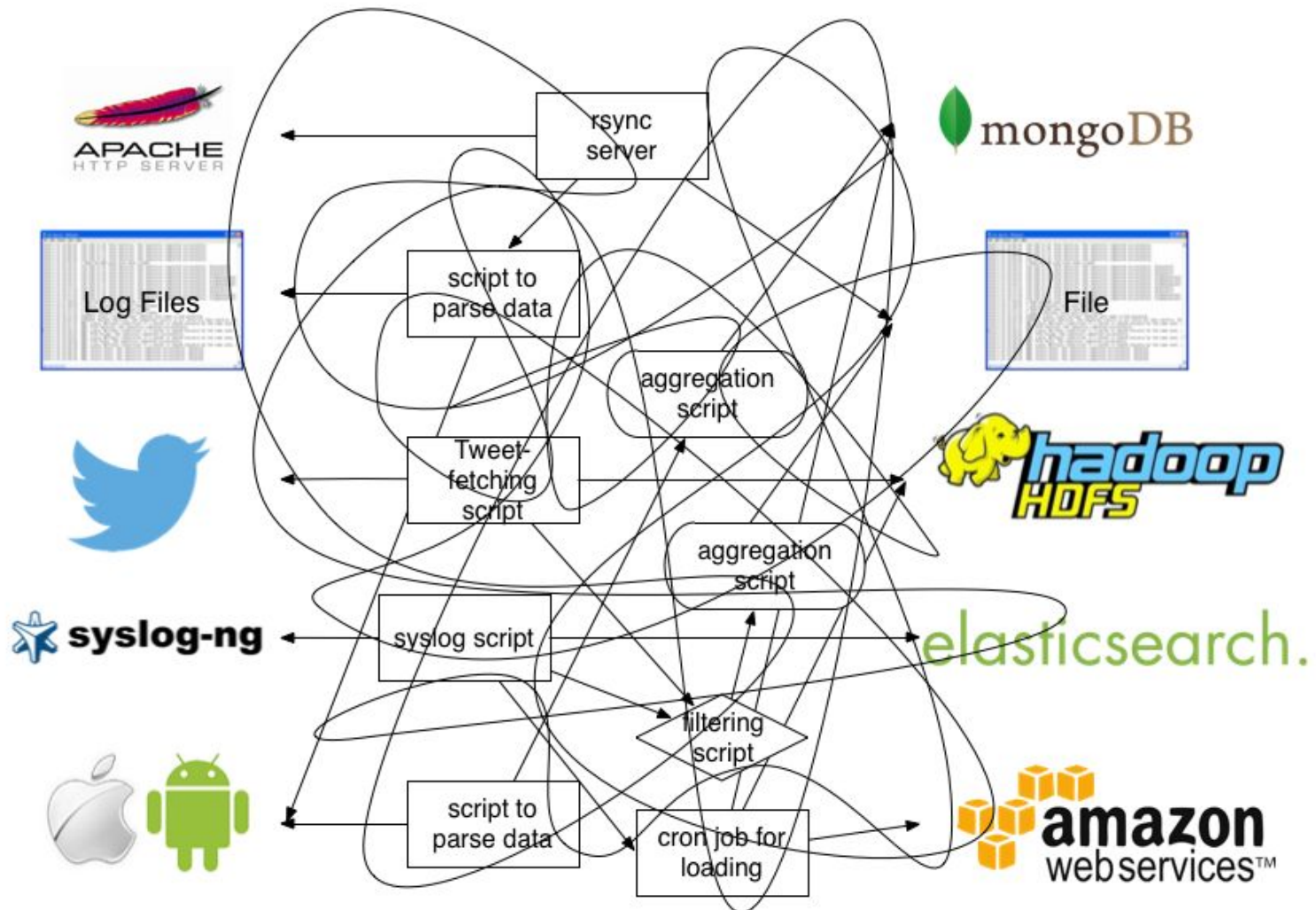
# Шаги доставки





Пару  
скриптиков и  
норм!

Ой...







14:48

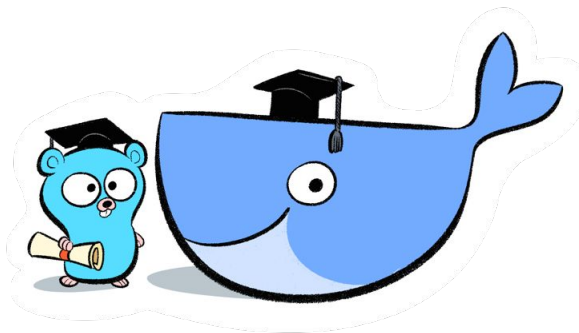
\$00700911

ПОТРАЧЕНО

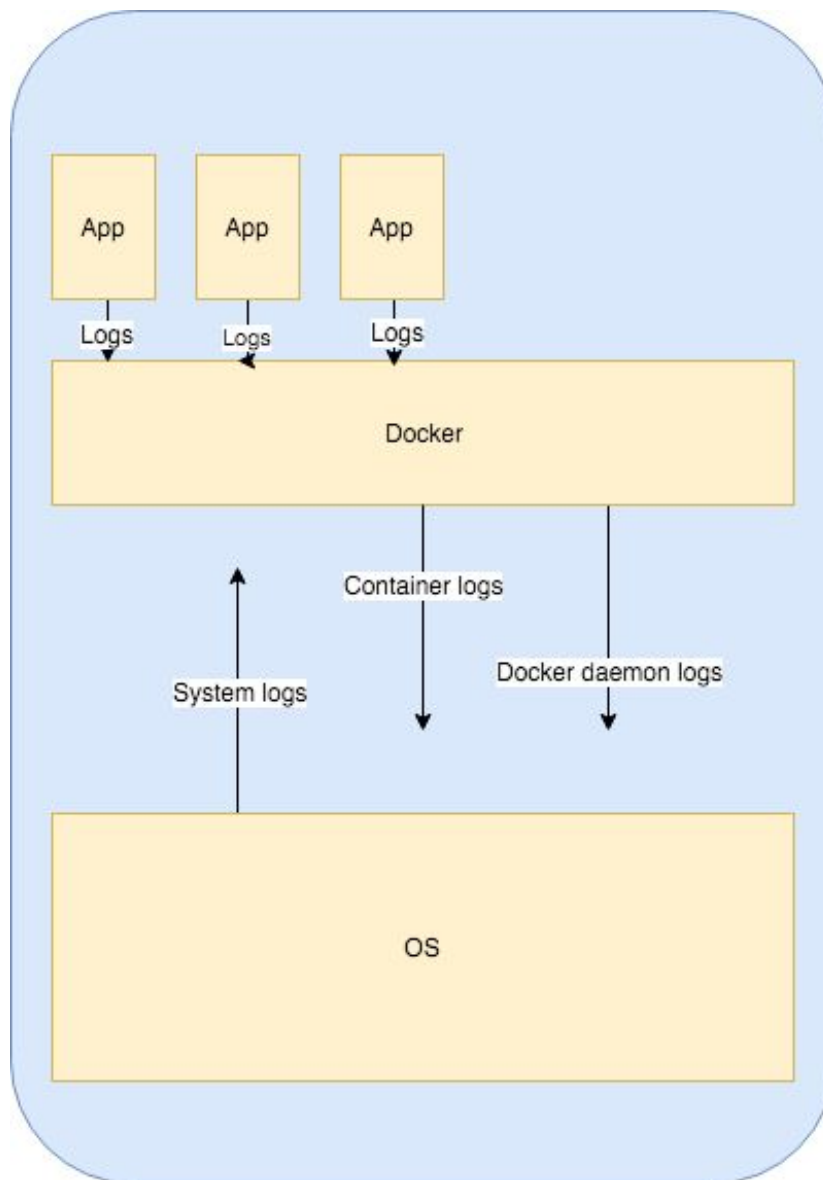




# Docker логи

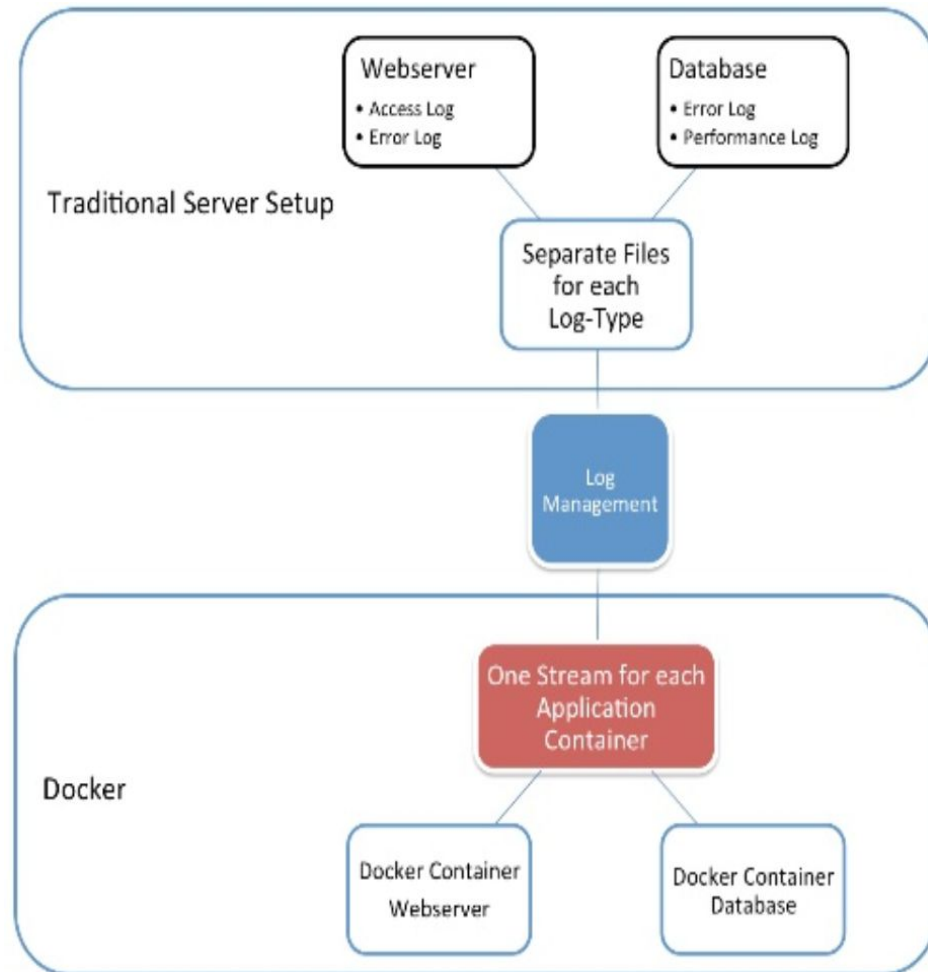


# Логи на ноде

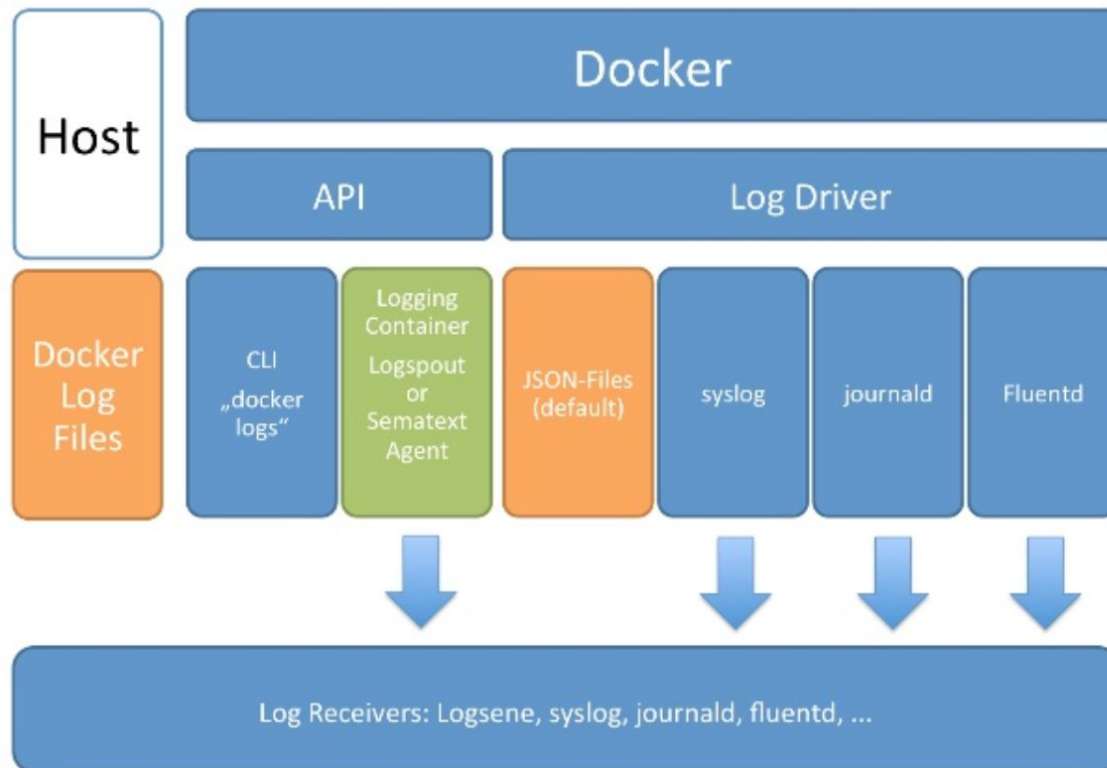




# Логи приложений в docker



- Traditionally separate files for each Application and Log-Type
  - error.log
  - access.log
- Docker Logs are `stdout / stderr` of processes running in a container
- Most official images log to console



- **Docker Log Drivers**
  - json-file, syslog, fluentd, journald, gelf
- **Docker API based Logging Containers**
  - Logspout
  - Sematext Docker Container
- Custom images with installed log shipper (syslog)



- + Простой путь доставки логов в local/remote destination
- + Настраивается per container или глобально
- Нет настройки парсинга
- Контейнер падает если удаленный приемник логов недоступен
- Будет плохо если диски медленные



# Docker logging drivers

---

- none
- json-file
- syslog
- journald
- gelf
- fluentd
- awslogs
- splunk
- etwlogs
- gcplogs

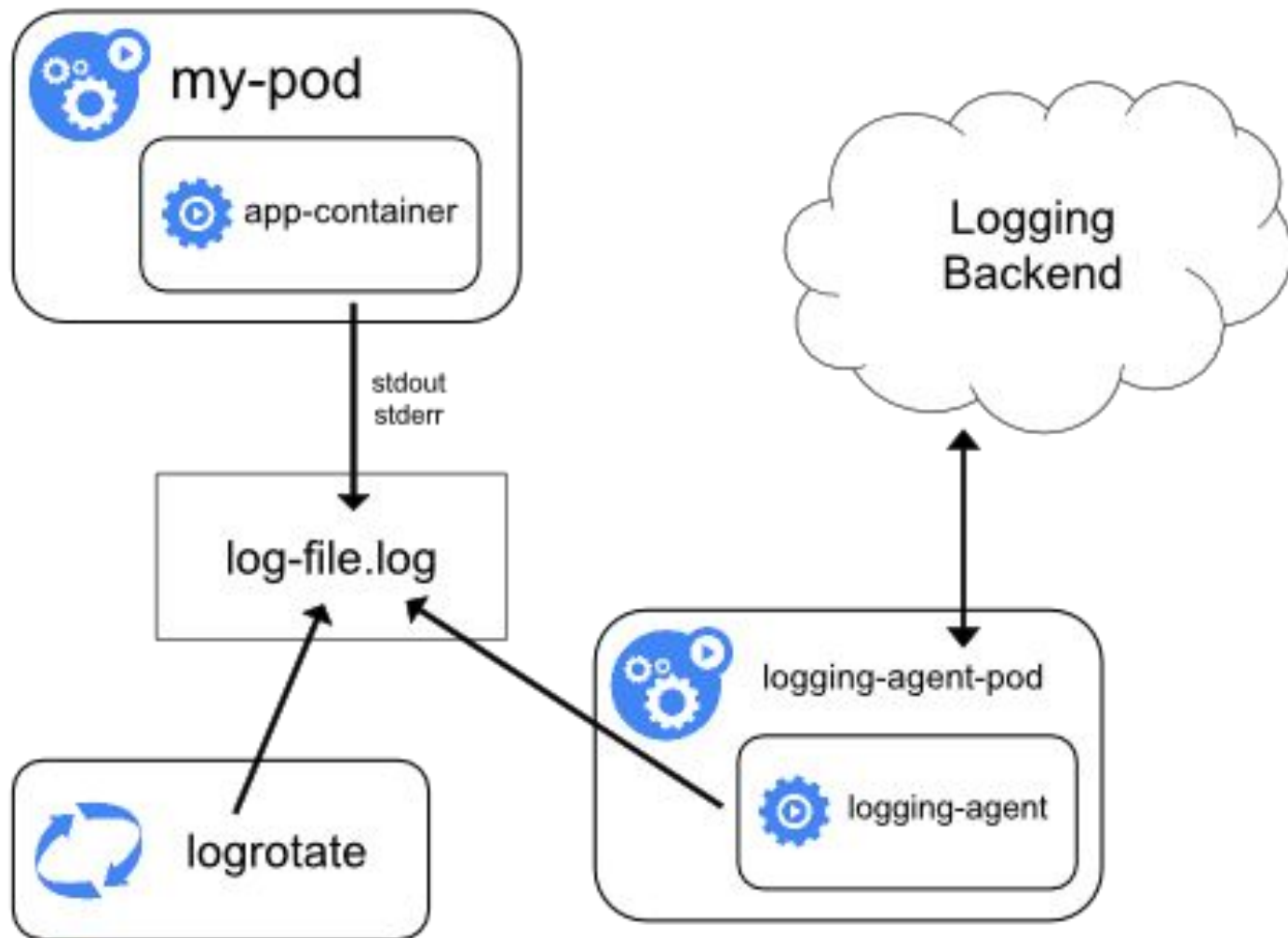


Запускаем приложение, которое выдает на `stdout` 30000 сообщений (в формате `json`) длиной ~530 байт.

Результаты:

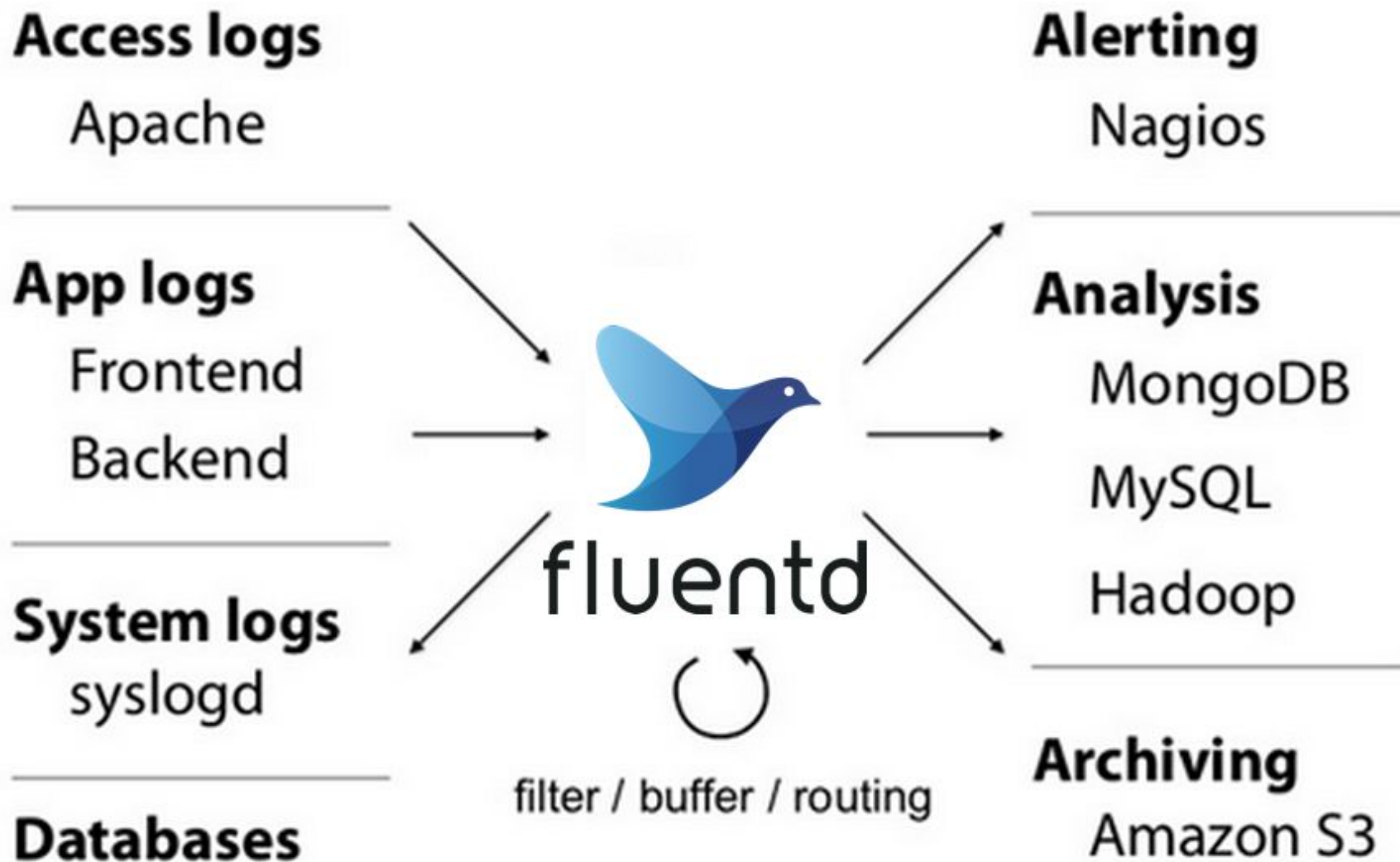
- `json-file` 0.2s
- `journald` 0.9s
- `syslog` 0.75s

# Logging agent





# Решение







## Fluentd: Unified Logging Layer

Fluentd is an open-source logging solution to unify data collection and consumption.

 Cloud Native Computing Foundation  <http://www.fluentd.org/>

 **Repositories** 52

 **People** 31

### Pinned repositories

#### **fluentd**

Fluentd: Unified Logging Layer (project under CNCF)

 Ruby  6.2k  741



## Fluentd: Unified Logging Layer

Fluentd is an open-source logging solution to unify data collection and consumption.

 Cloud Native Computing Foundation  <http://www.fluentd.org/>

 **Repositories** 52

 **People** 31

### Pinned repositories

#### fluentd

Fluentd: Unified Logging Layer (project under CNCF)

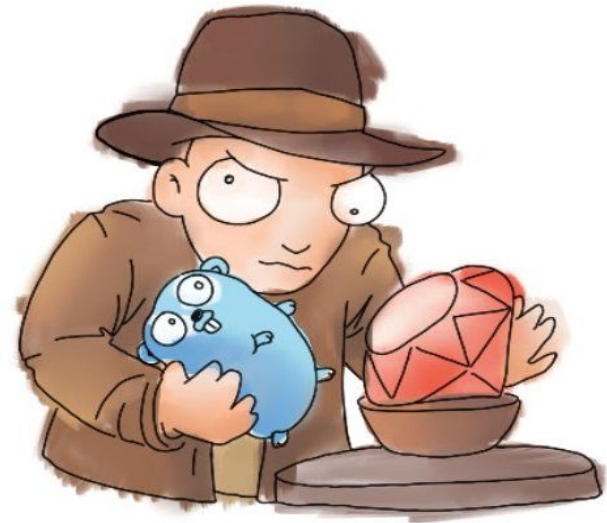
 **Ruby** ★ 6.2k  741

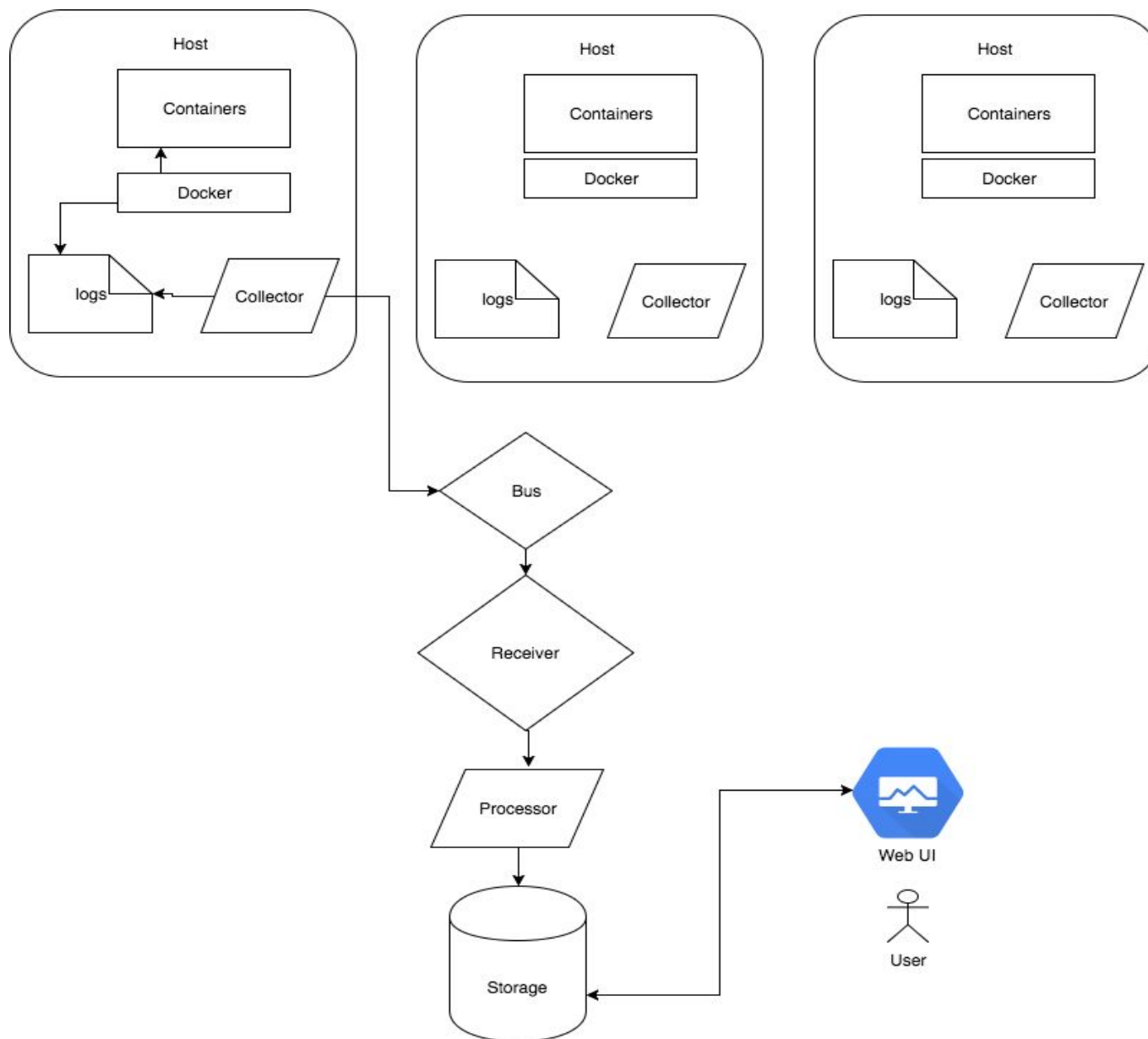


# Если серьезно



1. Из команды никто не знает ruby
2. Fluentd есть память больше чем наши приложения
3. Flexibility менее важно чем system footprint

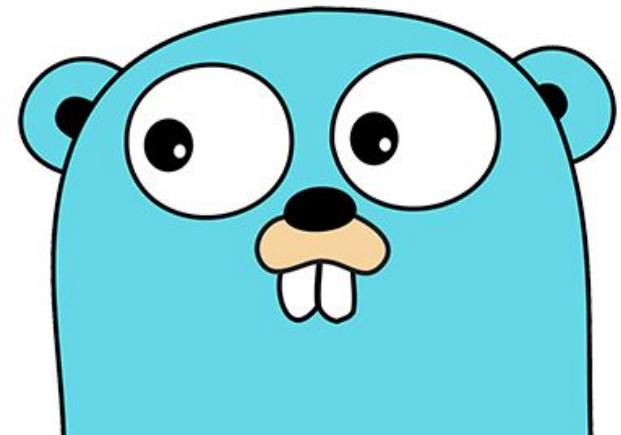




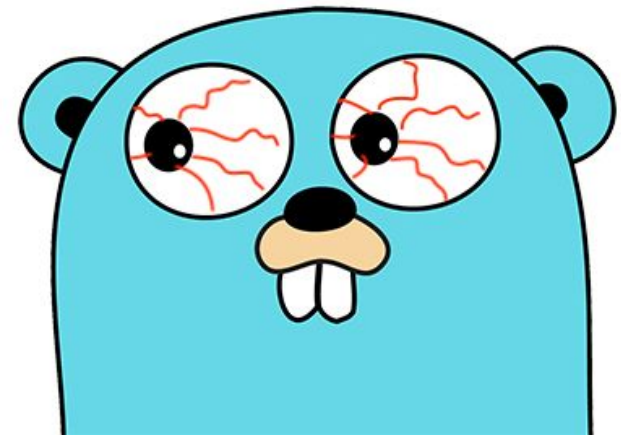


1. Что использовать для снятия логов с машин и сервисов?
2. Какой транспорт использовать для передачи логов от коллектора до хранилища?
3. Что использовать для хранилища?
4. Нужен ли компонент обработчик (фильтр например)?
5. Где разворачивать хранилище?
6. Какой инфраструктурой пользоваться?
7. Как мониторить состояние системы?

- logspout
- rsyslog
- syslog-ng
- filebeat
- fluent-bit
- nxlog
- hecka
- systemd journald remote
- custom
- ...

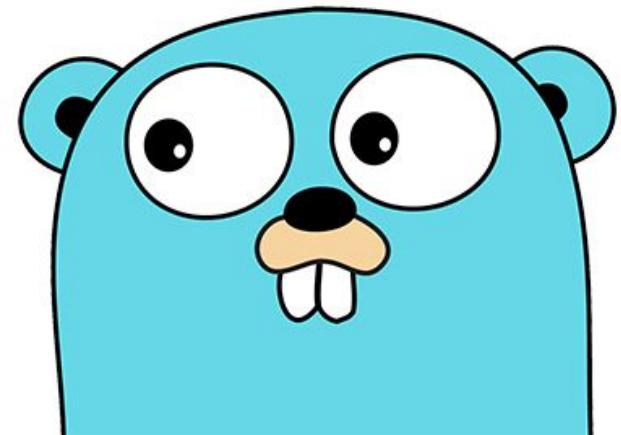


- Kafka
- Scribe
- Flume
- Kestrel
- Fluentd
- RabbitMQ
- BookKeeper
- Calligraphus
- ActiveMQ
- ...

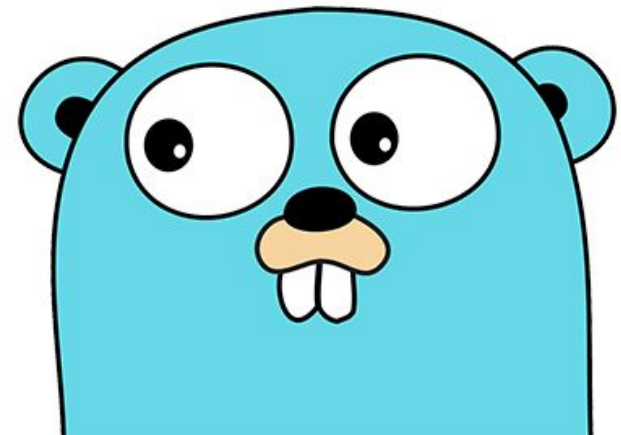




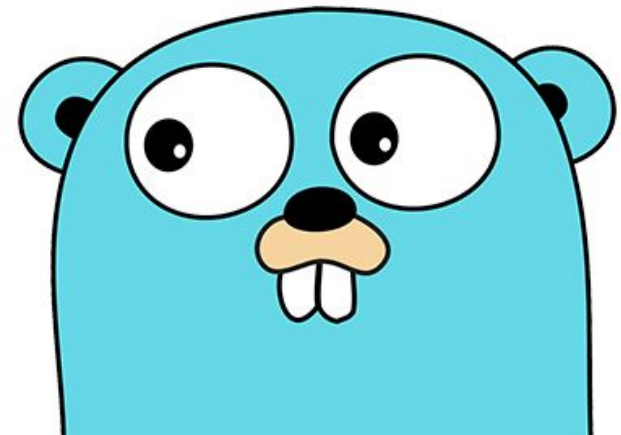
- logstash
- fluent-bit
- syslog-ng
- custom
- ...



- Elasticsearch
- Mongo
- RDBMS
- ClickHouse
- ...



- Kibana
- Grafana

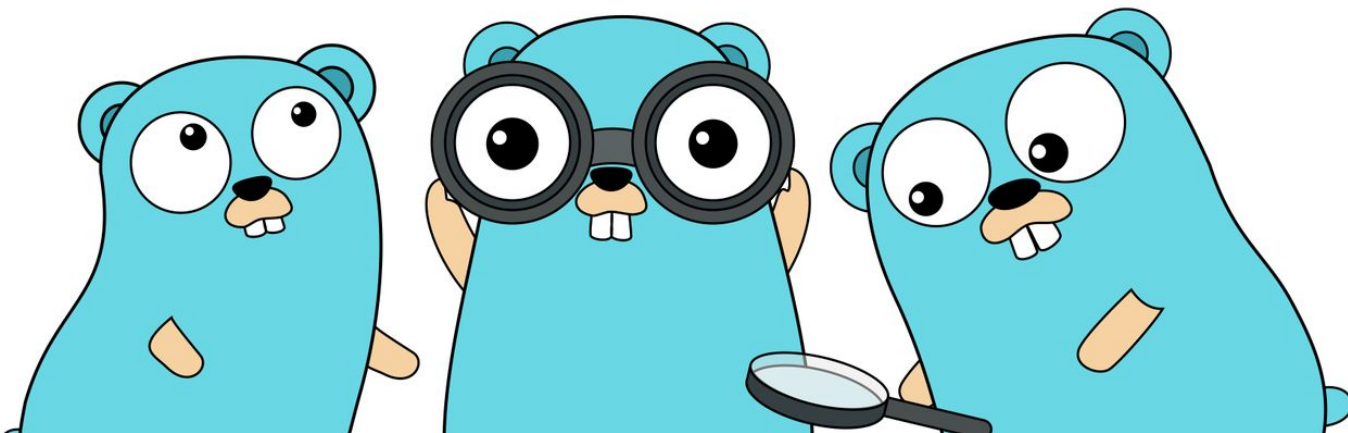


OMG

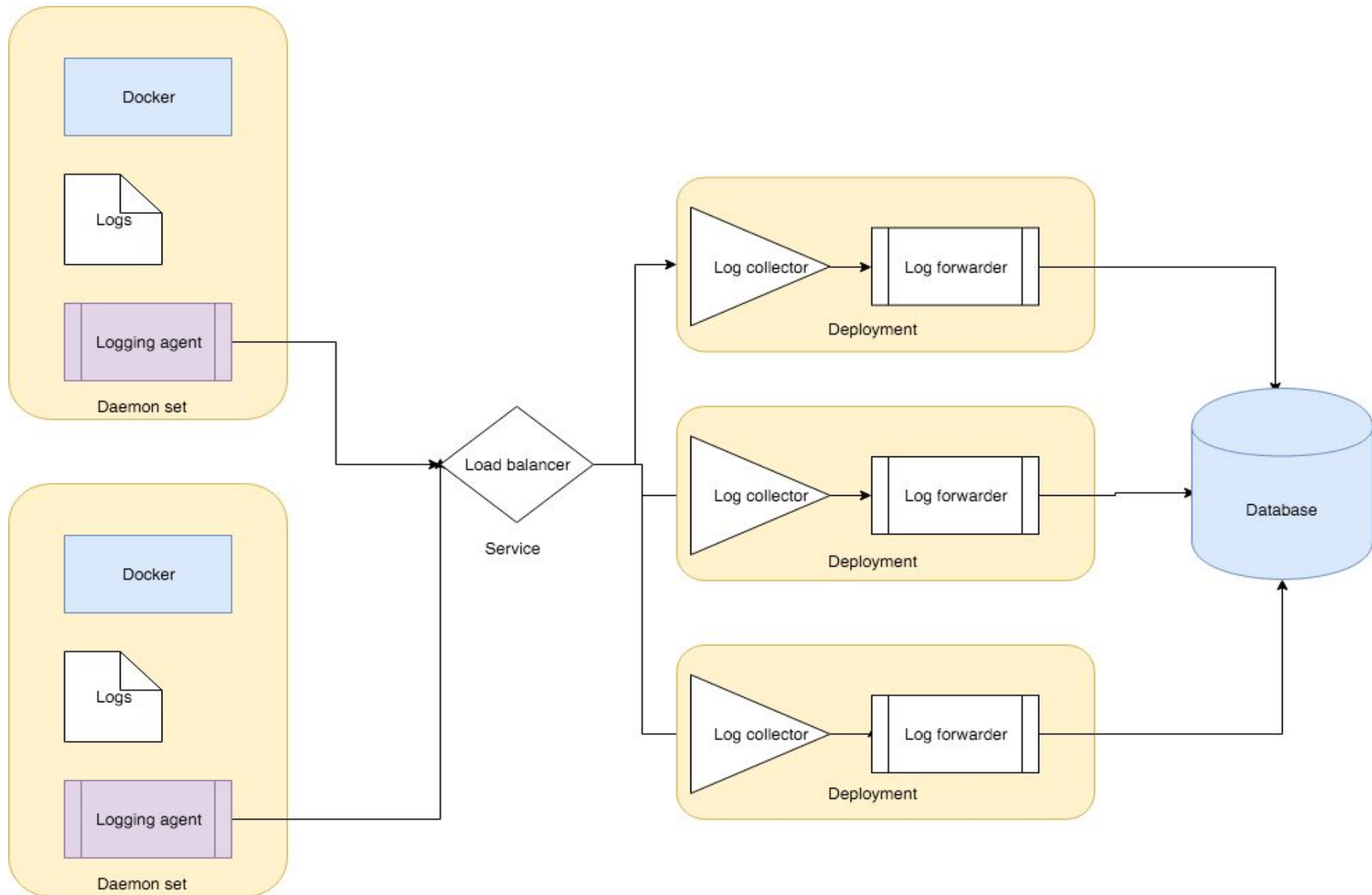




С начала



# Снова схема





## syslog-ng (си)

- + Позволяет описать много input, output и parser для одного экземпляра демона
- Нет никакой интеграции с docker или k8s

## fluent-bit (си)

- + Есть плагин для работы с k8s api
- Умеет только один input и output

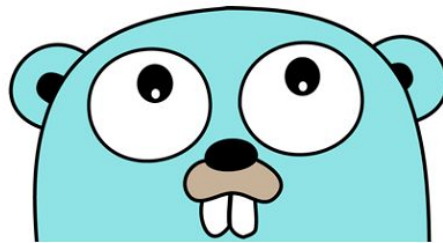


## syslog-ng

- На нагрузке 10000 в сек падает и ломает указатели на позицию файла журналов
- После перезапуска падал в segfault

## fluent-bit

- Проблем не обнаружено

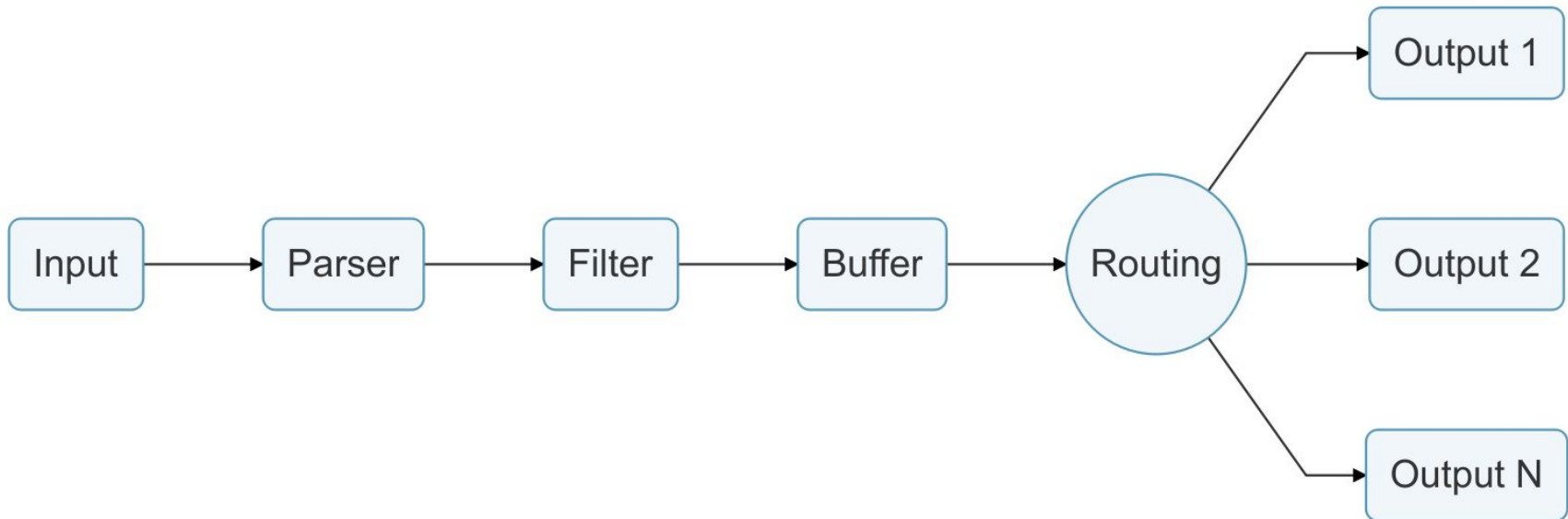




# Log agent: syslog-ng vs fluent-bit



+ fluent-bit





# Storage: elasticsearch vs clickhouse

## Search engine + storage - Elasticsearch (Java)

- Популярное решение, используется в Graylog, ELK
- Есть UI - Kibana
- Позволяет хранить без схемы
- Query language - lucene

## Column based - ClickHouse (C++)

- Строгая схема
- UI нет, только сторонний
- Query language - SQL

## RDBMS

- Сложно масштабировать, нет кластеризации, performance



## Нода database - Ubuntu 16.04.3 x64

- Intel(R) Xeon(R) Platinum 8168 CPU @ 2.70GHz 4 core
- 6 GB Memory
- SSD

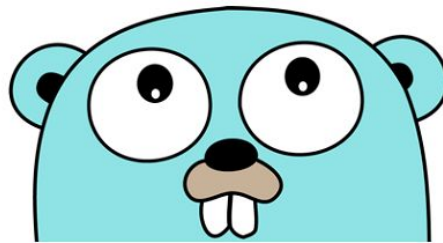
## Нода load generator - Ubuntu 16.04.3 x64

- Intel(R) Xeon(R) Platinum 8168 CPU @ 2.70GHz 2 core
- 3 GB Memory
- SSD



Elasticsearch 100 000 000 записей

Съел всю память в хипе и 12 Gb  
места на диске



# Storage: elasticsearch vs clickhouse



```
1000000 msg, bulk 1000000
```

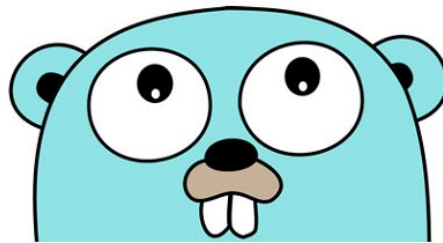
```
real 0m23.633s
```

```
user 0m22.000s
```

```
sys 0m0.580s
```

```
# du -csh elasticsearch/nodes
```

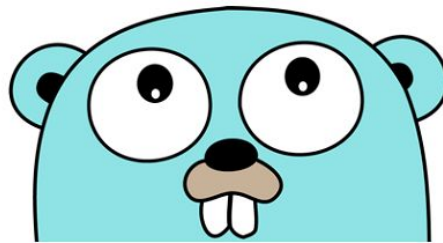
```
390M
```





Clickhouse 100 000 000 записей

Съел 400mb памяти и 3.3 Gb места на диске



# Storage: elasticsearch vs clickhouse



```
1000000 msg, bulk 1000000
```

```
real 0m4.134s
```

```
user 0m6.972s
```

```
sys 0m0.180s
```

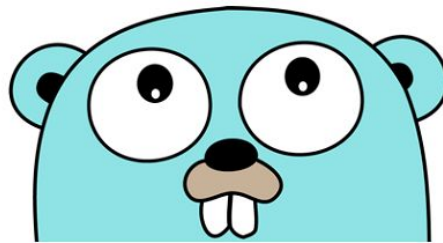
```
# du -csh clickhouse/data
```

```
30M
```



+ Clickhouse

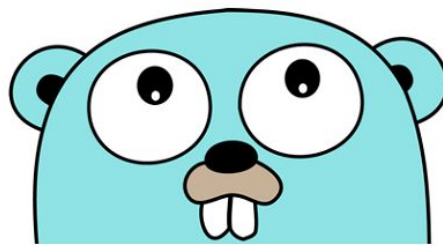
Оптимально по performance и месту  
на диске







1. fluent-bit не умеет писать в clickhouse
2. нужно парсить, чтобы разложить в типизированную схему



# Проблема разных логов



```
> docker run -d -p 80:80 --rm --name apache-app -v "$PWD":/usr/local/apache2/htdocs/ httpd:2.4
```

```
> docker logs apache-app
```

```
AH00558: httpd: Could not reliably determine the server's fully qualified domain name, using 172.17.0.2. Set the 'ServerName' directive globally to suppress this message
```

```
AH00558: httpd: Could not reliably determine the server's fully qualified domain name, using 172.17.0.2. Set the 'ServerName' directive globally to suppress this message
```

```
[Sat Jul 25 11:52:13.381905 2015] [mpm_event:notice] [pid 1:tid 140020361725824] AH00489: Apache/2.4.16 (Unix) configured -- resuming normal operations
```

```
[Sat Jul 25 11:52:13.382133 2015] [core:notice] [pid 1:tid 140020361725824] AH00094: Command line: 'httpd -D FOREGROUND'
```

```
192.168.59.3 - - [25/Jul/2015:11:52:22 +0000] "GET / HTTP/1.1" 200 429
```

```
192.168.59.3 - - [25/Jul/2015:11:52:22 +0000] "GET /favicon.ico HTTP/1.1" 404 209
```

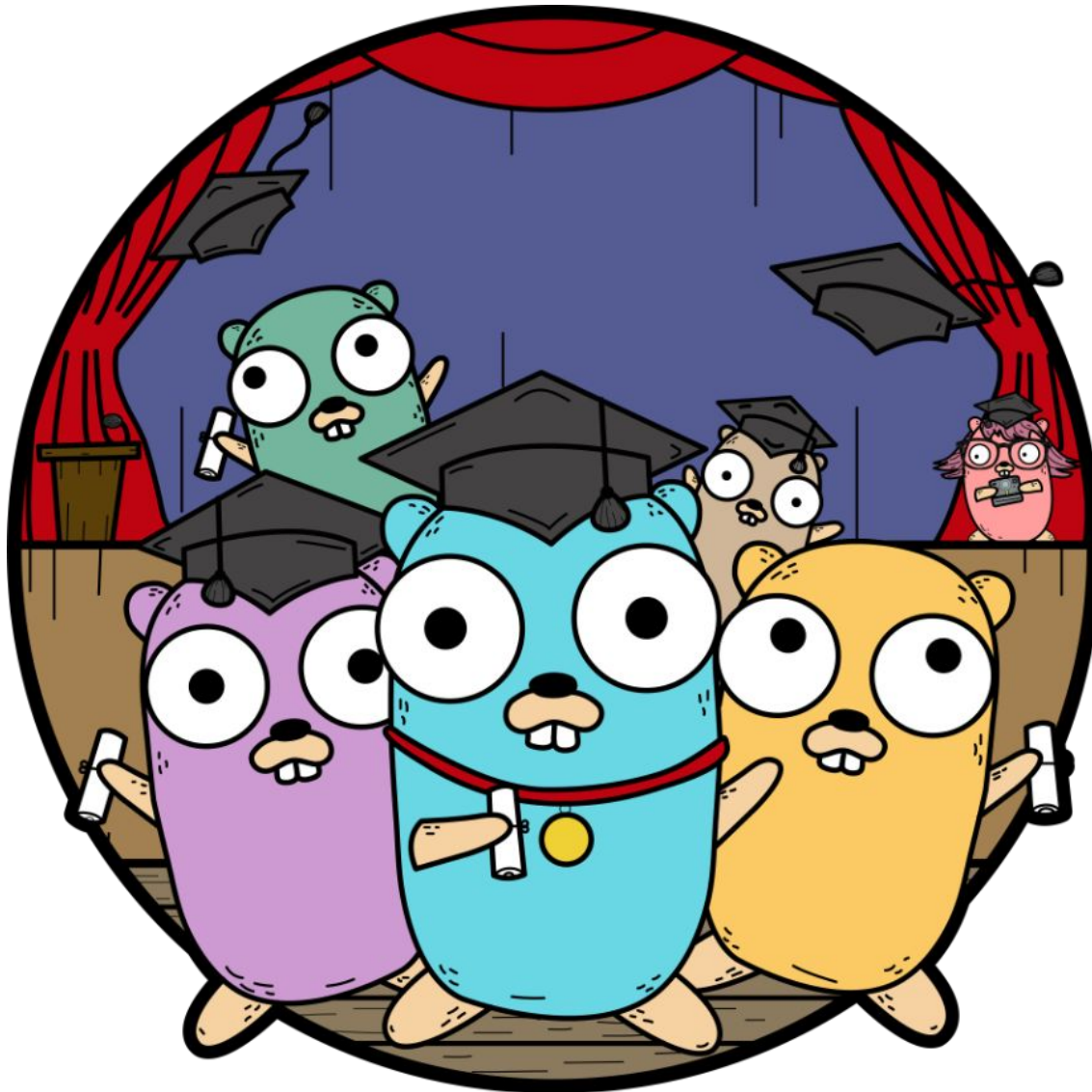
```
192.168.59.3 - - [25/Jul/2015:11:52:22 +0000] "GET /favicon.ico HTTP/1.1" 404 209
```

```
192.168.59.3 - - [25/Jul/2015:11:52:29 +0000] "GET /index.js HTTP/1.1" 200 1258
```

```
192.168.59.3 - - [25/Jul/2015:11:53:12 +0000] "GET /index.js HTTP/1.1" 304 -
```

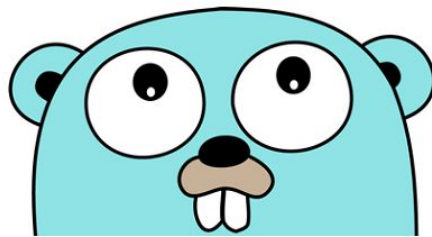
```
192.168.59.3 - - [25/Jul/2015:11:53:18 +0000] "GET /index.js HTTP/1.1" 200 1258
```

# Дипломированные Go специалисты

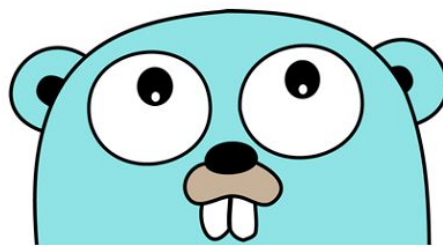
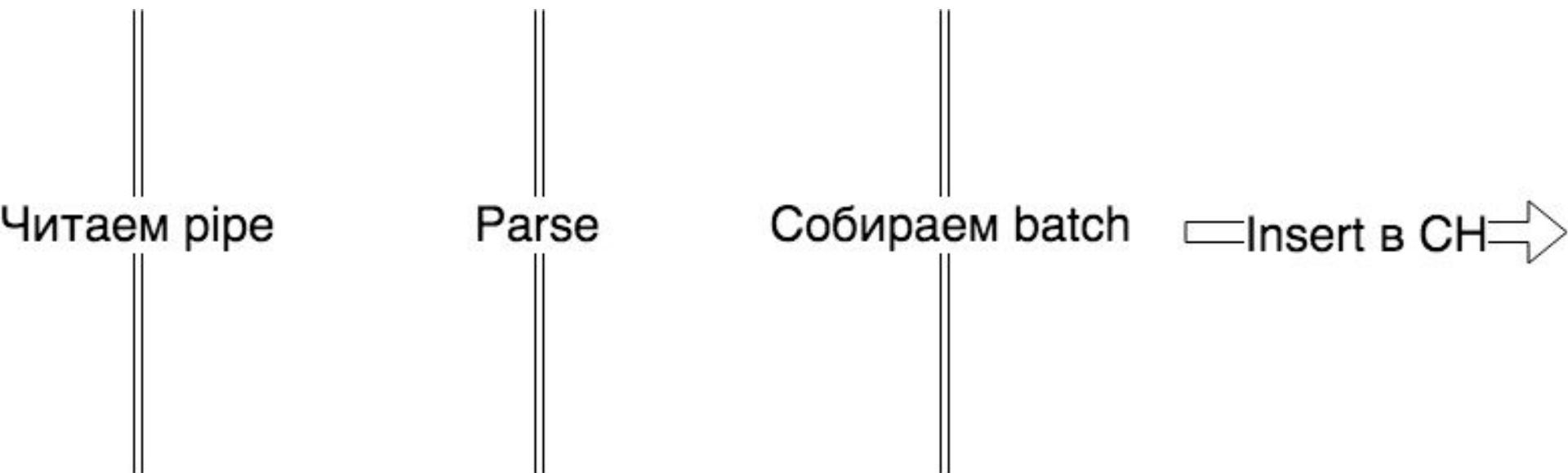




- Демон на Go
- Работает через pipe с fluent-bit
- Парсит разные типы логов
- Собирает batch для записи в Clickhouse



# Clickhouse forwarder





1. [github.com/kshvakov/clickhouse](https://github.com/kshvakov/clickhouse)
2. [github.com/mailru/easyjson](https://github.com/mailru/easyjson)
3. [github.com/hashicorp/logutils](https://github.com/hashicorp/logutils)
4. [github.com/facebookgo/flagenv](https://github.com/facebookgo/flagenv)

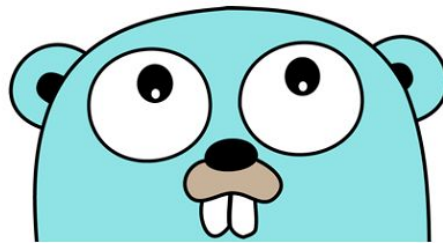




Chart - пакет для kubernetes

Create, Update, Delete операции

содержит шаблоны ресурсов,  
переменные, скрипты



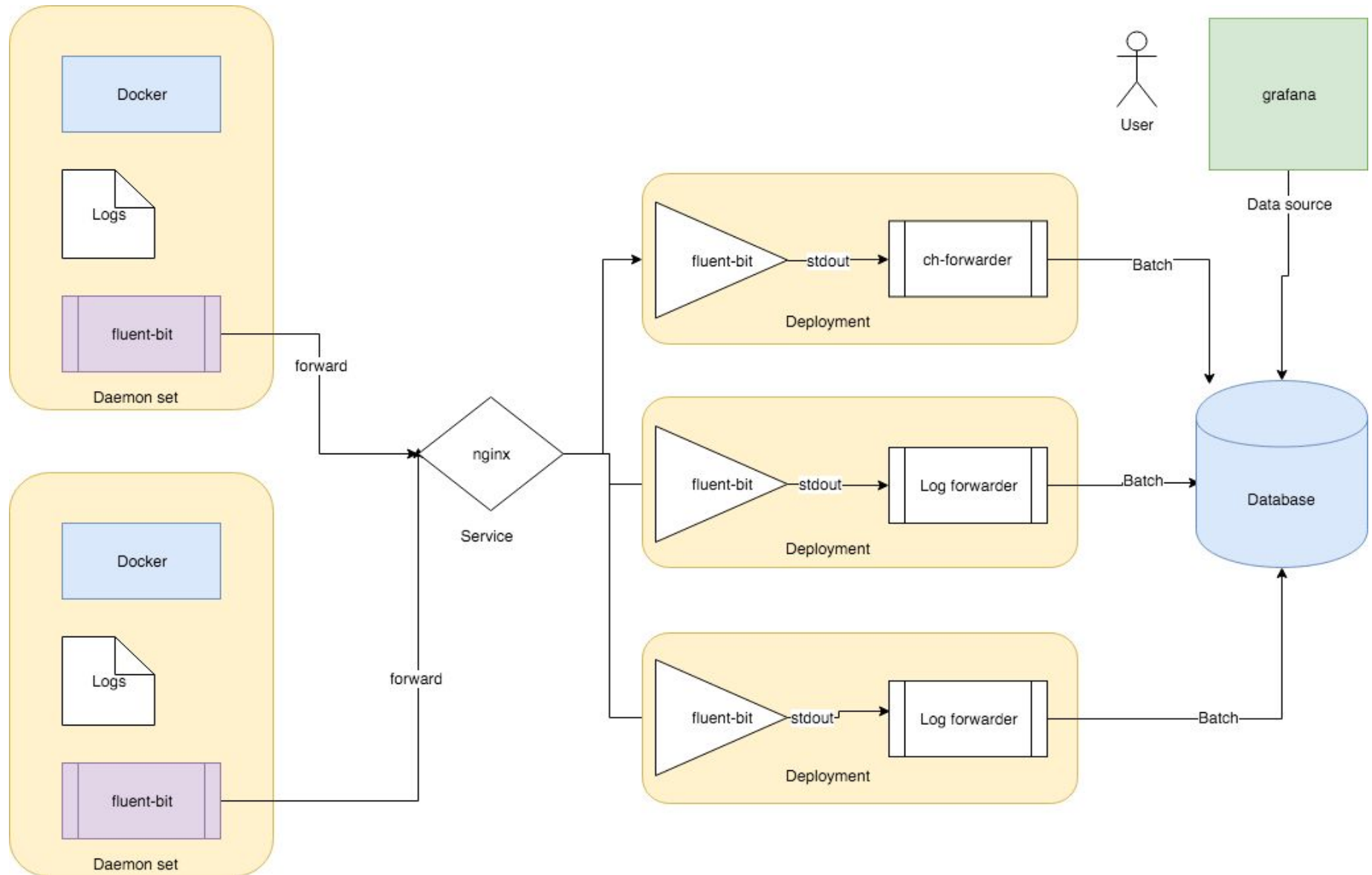


# Deploy - Helm chart

```
logging
├── Chart.yaml
├── files
│   ├── create_merges.sh
│   └── create_partitions.sh
├── templates
│   ├── _helpers.tpl
│   ├── db-pre-install.job.yaml
│   ├── db-rotate.cronjob.yaml
│   ├── db-scripts.configmap.yaml
│   ├── db.secret.yaml
│   ├── log-agent.clusterrole.yaml
│   ├── log-agent.clusterrolebinding.yaml
│   ├── log-agent.configmap.yaml
│   ├── log-agent.daemonset.yaml
│   ├── log-agent.serviceaccount.yaml
│   ├── log-collector.configmap.yaml
│   ├── log-collector.deployment.yaml
│   └── log-collector.service.yaml
└── values.yaml
```



# Cxema

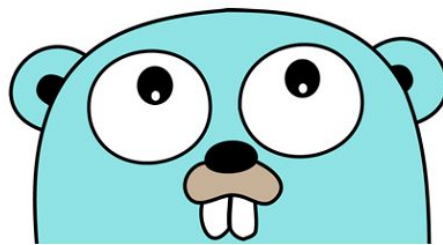


Срок разработки - 2 мес

Пропускная способность - 10 000  
msg/sec с машины

Что не доделали:

- UI еще в разработке, нужен plugin для grafana

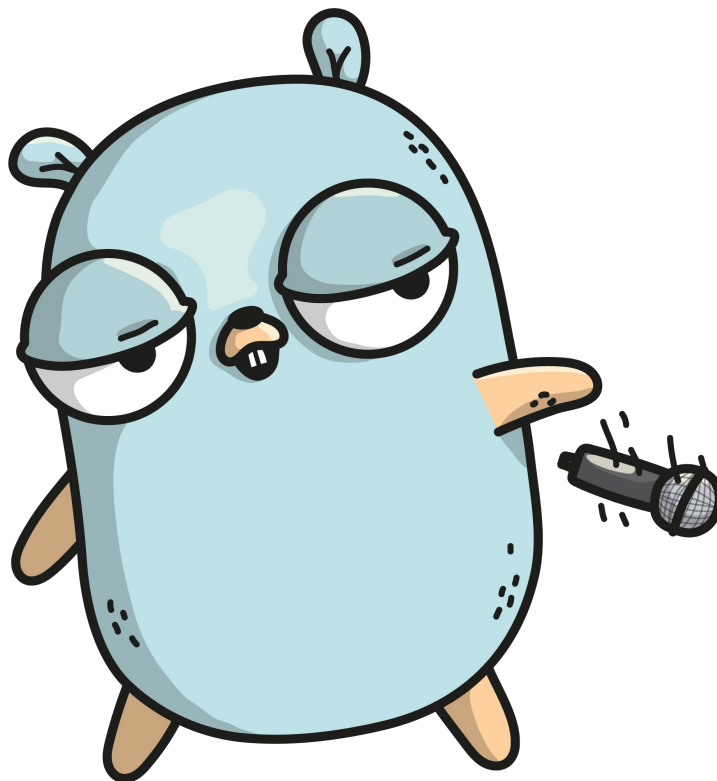


# Когда рассказываешь об этом проекте





# Тинькофф



Tinkoff.ru