

# Docker swarm

# Docker Swarm



- Поставляется из коробки с Docker (начиная с версии 1.12)
- Разворачивается 1-й командой (начиная с версии 1.12)
- Не использует внешние распределенные хранилища
- Позволяет описывать конфигурацию в формате Docker Compose
- Есть встроенный Service Discovery
- Есть встроенный Load Balancer

# Swarm Cluster



# Service

- Описывает что должно быть запущено и как (количество задач, необходимые ресурсы и т.д.)
- При запуске Service, *лидер* назначает *Worker*-ам необходимые задачи (*tasks*)
- Сервисы и их зависимости объединяем в Stack

## Управление с помощью

- `docker service create/update/rm/...`
- `docker stack deploy -c compose-file.yml`

# Stack

version: '3.3'

services:

mongo:

image: mongo:\${MONGO\_VERSION}

volumes:

- mongo\_data:/data/db

networks:

back\_net:

aliases:

- post\_db
- comment\_db

post:

image: \${USER\_NAME}/post:\${POST\_VERSION}

networks:

- front\_net
- back\_net

comment:

image: \${USER\_NAME}/comment:\${COMMENT\_VERSION}

networks:

- front\_net
- back\_net

ui:

image: \${USER\_NAME}/ui:\${UI\_VERSION}

ports:

- "\${UI\_PORT}:9292/tcp"

networks:

- front\_net

volumes:

mongo\_data: {}

networks:

back\_net: {}

front\_net: {}

# Stack

```
>> docker stack deploy --compose-file docker-compose.yml ENV
```

не поддерживает переменные окружения и .env файлы

## Workaround

```
>> docker stack deploy --compose-file=<(docker-compose -f docker-compose.yml  
config 2>/dev/null) ENV
```

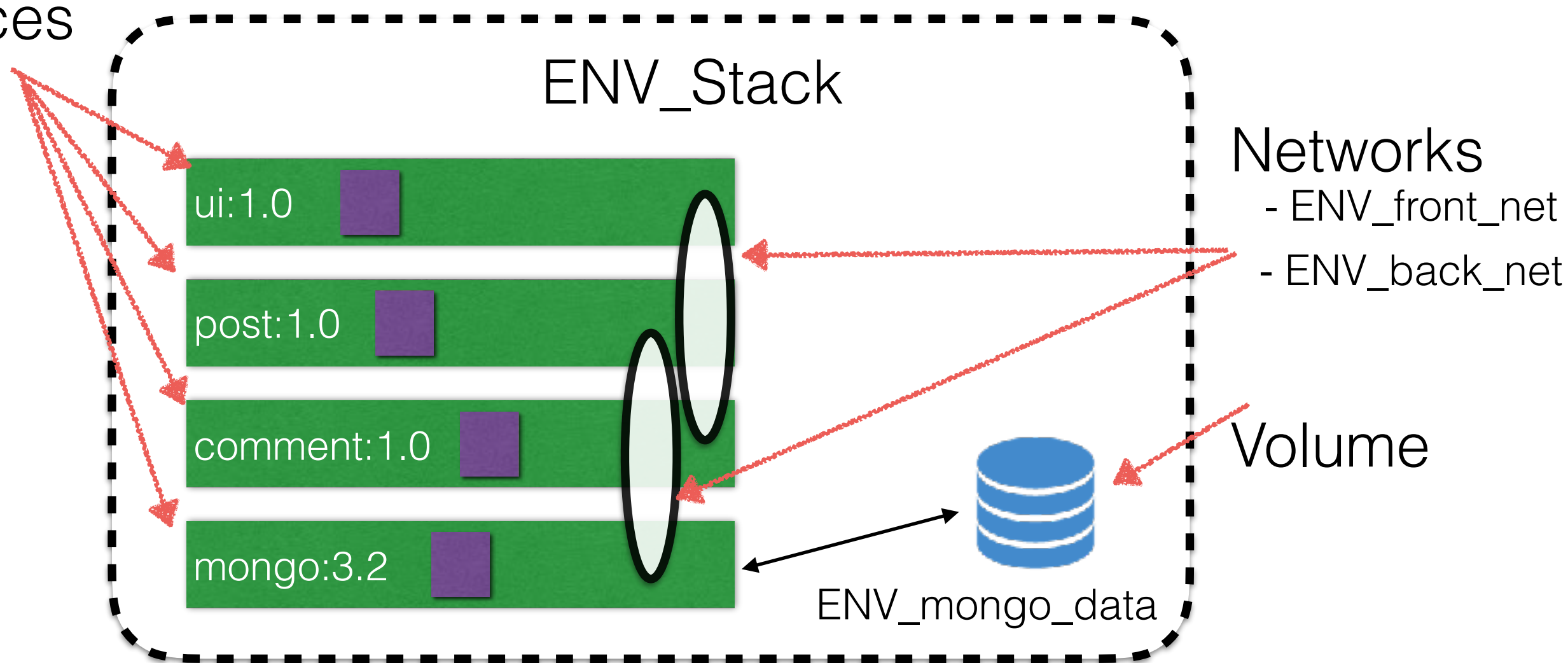
Имя стека



# Stack

>> docker stack deploy

Services



# Директивы планировщика

- Placement - **где разместить** задачу
- Deploy Mode - **как распределять** задачи
- Resources - **сколько ресурсов выделить** на задачу
- Update-policy - **как обновлять** задачи
- Restart-policy - **что делать**, если задача завершилась
- Labels - **какую метку присвоить** задаче
- Endpoint Mode - **как задачам взаимодействовать** между собой



# Описание задачи

*docker-compose.yml*

```
post:
  image: ${USER_NAME}/post:${POST_VERSION}
  deploy:
    replicas: 2
    placement:
      constraints:
        - node.labels.worker == true
  ...
```

ИЛИ

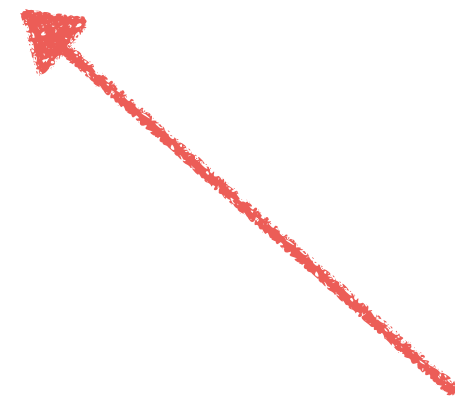
```
>> docker service create ...
```

# Где разместить задачу?

```
$ docker service create \  
  --name mongo \  
  --constraint 'node.labels.reliability == high' \  
  --constraint 'engine.labels.provider == google' \  
  mongo:3.2
```



```
mongo:  
  image: mongo:3.2  
  deploy:  
    placement:  
    constraints:  
      - node.labels.reliability == high  
      - engine.labels.provider == google
```



2 типа меток:

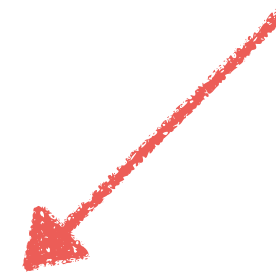
- node
- engine

...



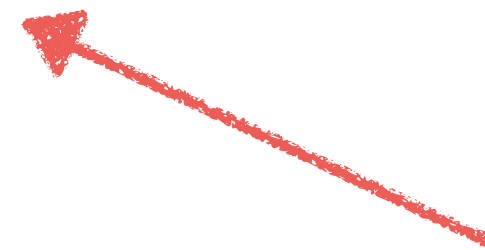
# Labels

добавить node label



```
>> docker node update --label-add reliability=high master-1
```

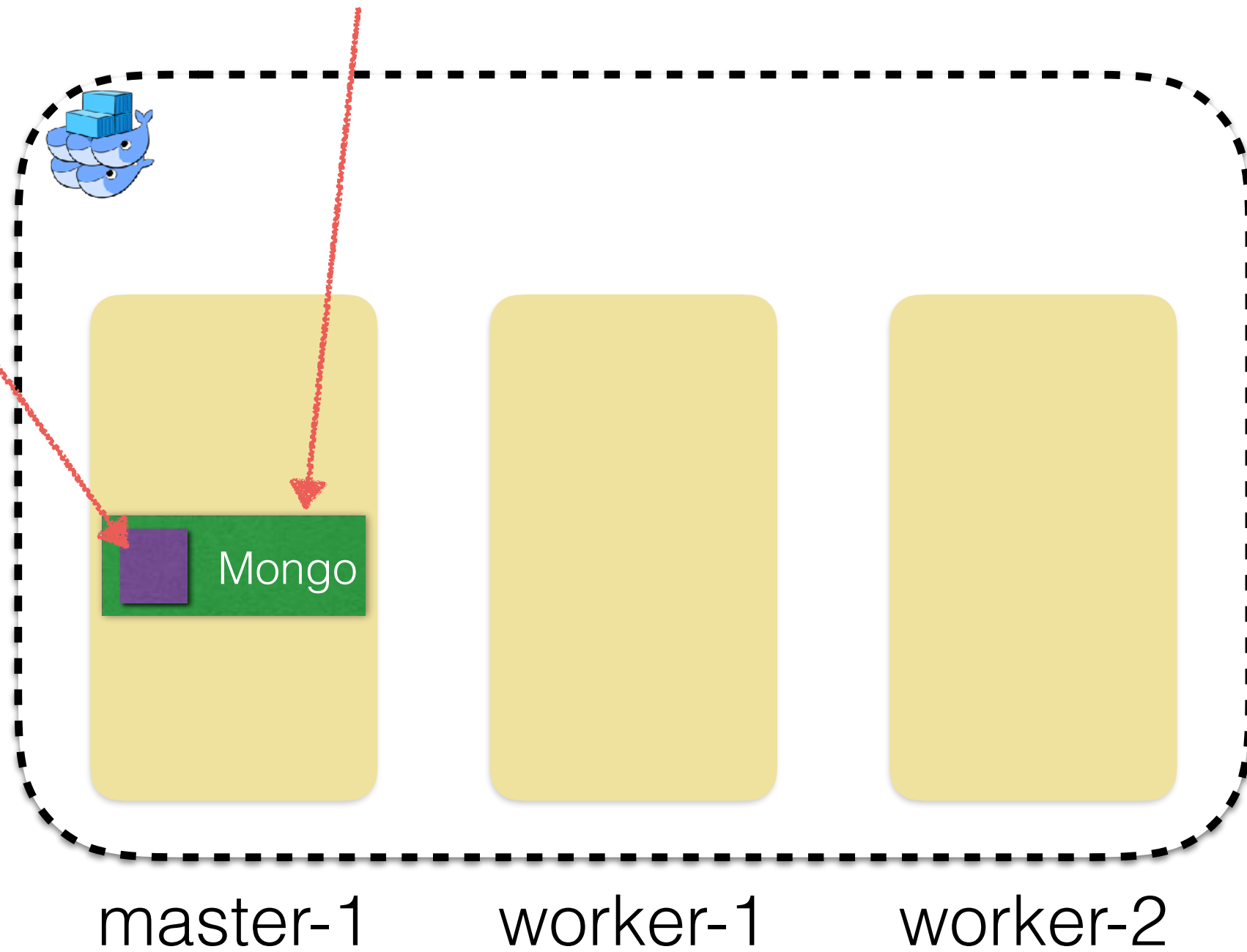
```
>> /usr/bin/dockerd --label provider=google ...
```



добавить engine label

Сервис

Контейнер



reliability=high

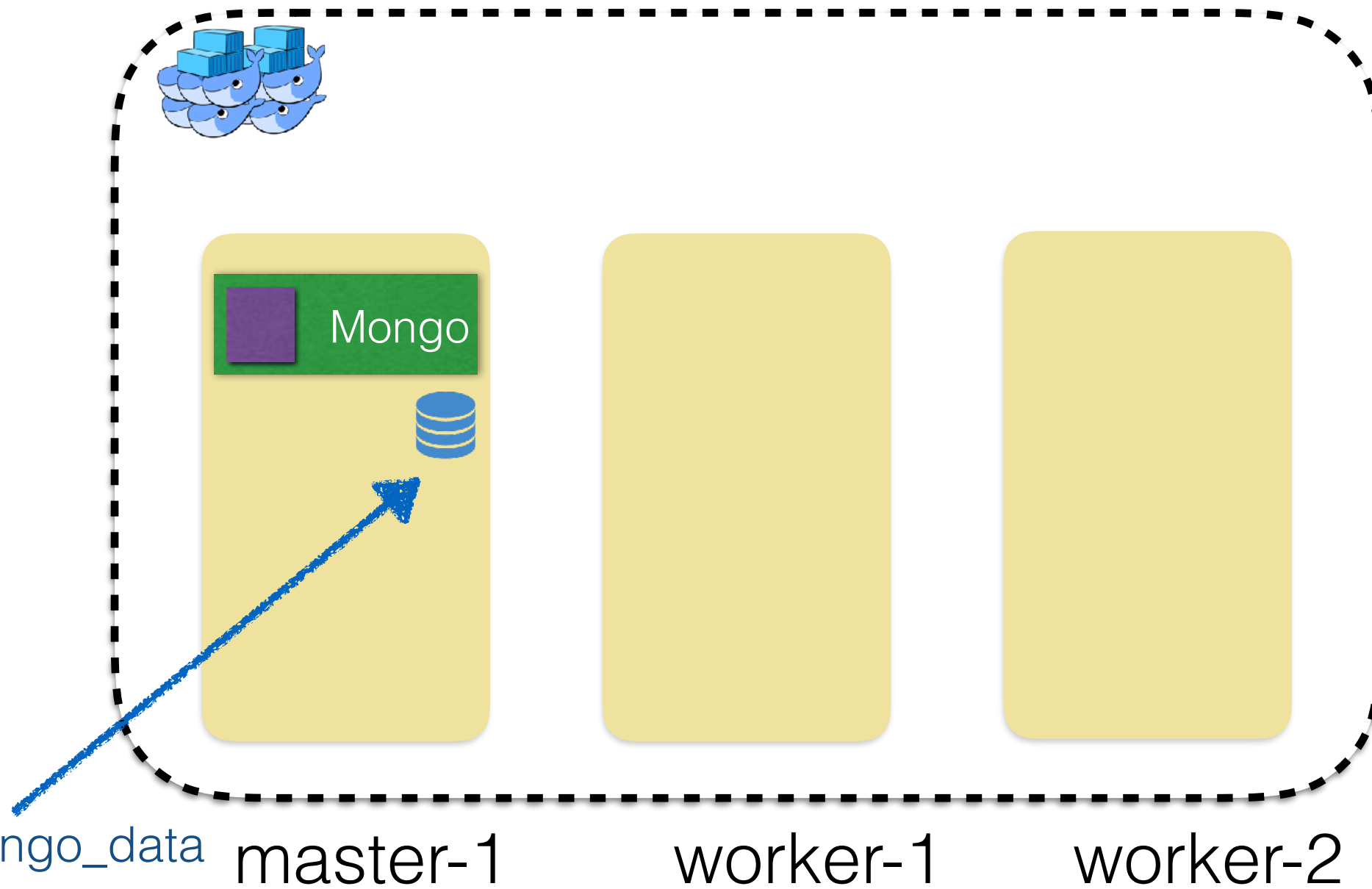
# А что с данными?

```
mongo:
  image: mongo:${MONGO_VERSION}
  deploy:
    placement:
      constraints:
        - node.labels.reliability == high
        - engine.labels.provider == google
  volumes:
    - mongo_data:/data/db
```

Экземпляр данных будет доступен только на одной машине

**Как думаете чем это может грозить?**

# Volume



reliability=high

# Deploy Mode

**Как распределять задачи (в соответствии с placement) ?**

- replicated - запустить определенное число задач (default)
- global - запустить задачу на каждой ноде

# Replicated

**Replicated** - запустить определенное число задач (default)

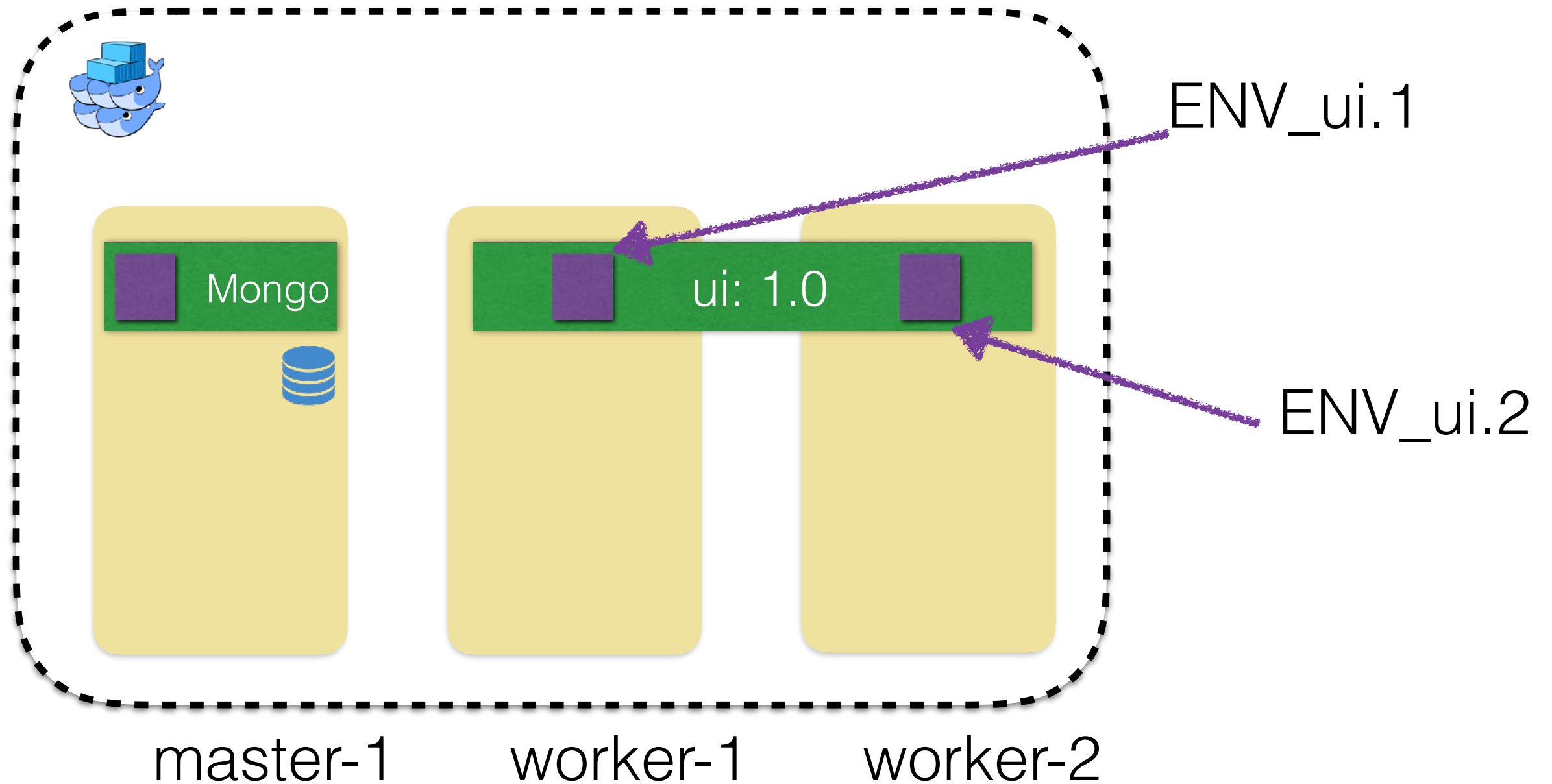
```
ui:  
  image: ${USER_NAME}/ui:${UI_VERSION}  
  deploy:  
    mode: replicated  
    replicas: 2  
    placement:  
      constraints:  
        - node.role == worker  
...
```



Не будем размещать на manager



# Replicated



reliability=high

# Replicated

Запустим остальные микросервисы

post:

image: \${USER\_NAME}/post:\${POST\_VERSION}

deploy:

mode: replicated

replicas: 3

placement:

...

comment:

image: \${USER\_NAME}/comment:\${COMMENT\_VERSION}

deploy:

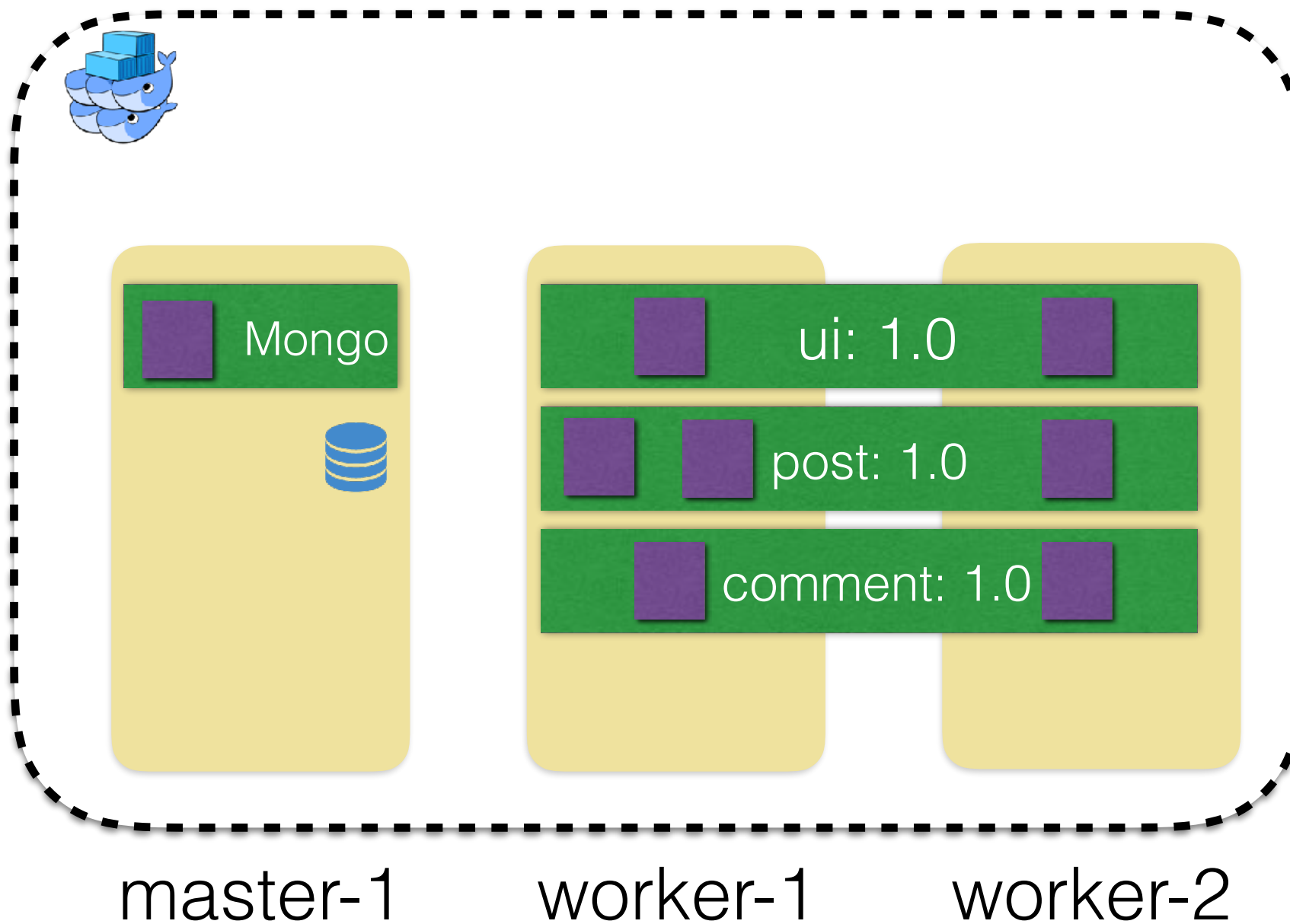
mode: replicated

replicas: 2

placement:

...

# Replicated



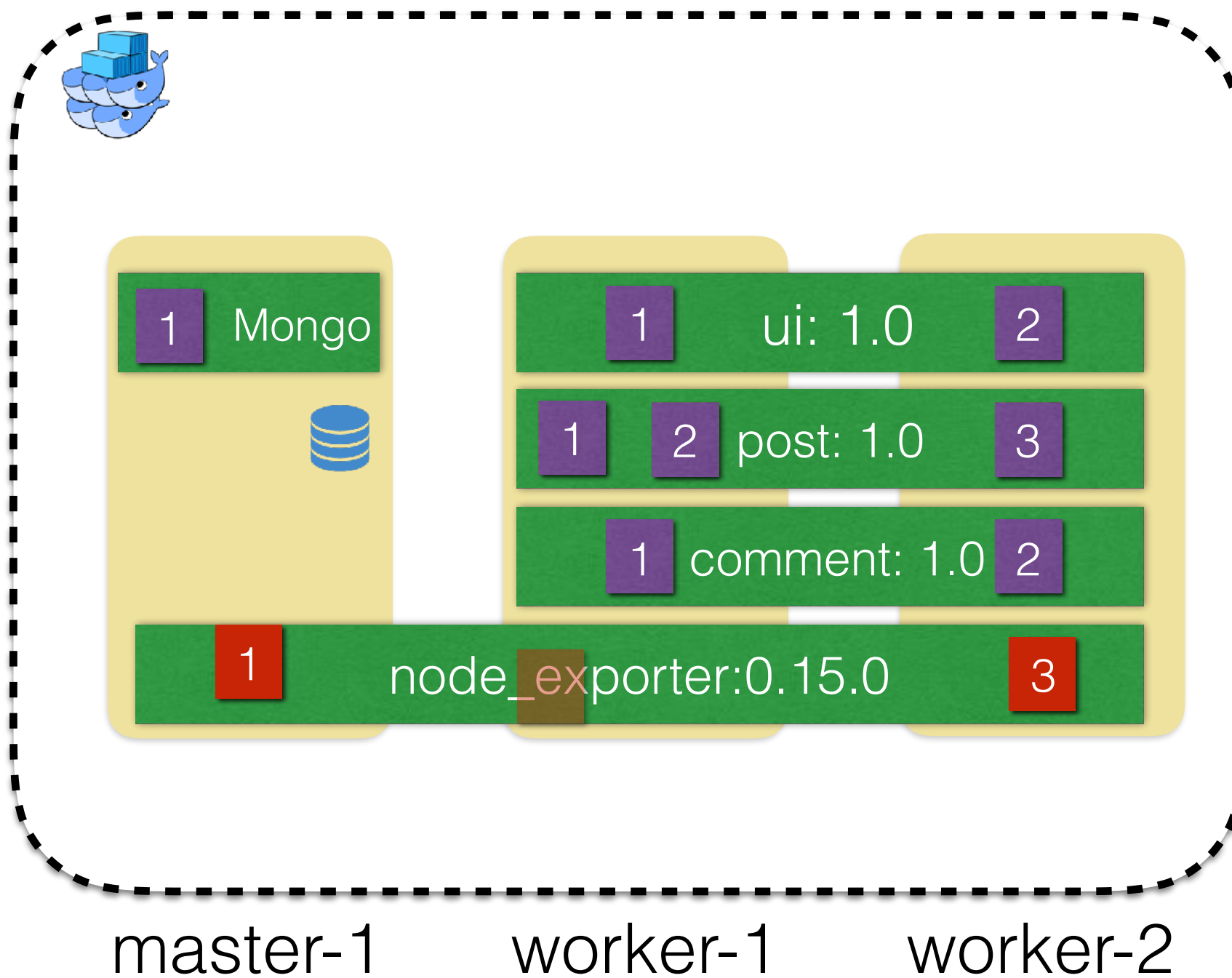
reliability=high

# Global

**Global** - запустить на задачу на каждой ноде

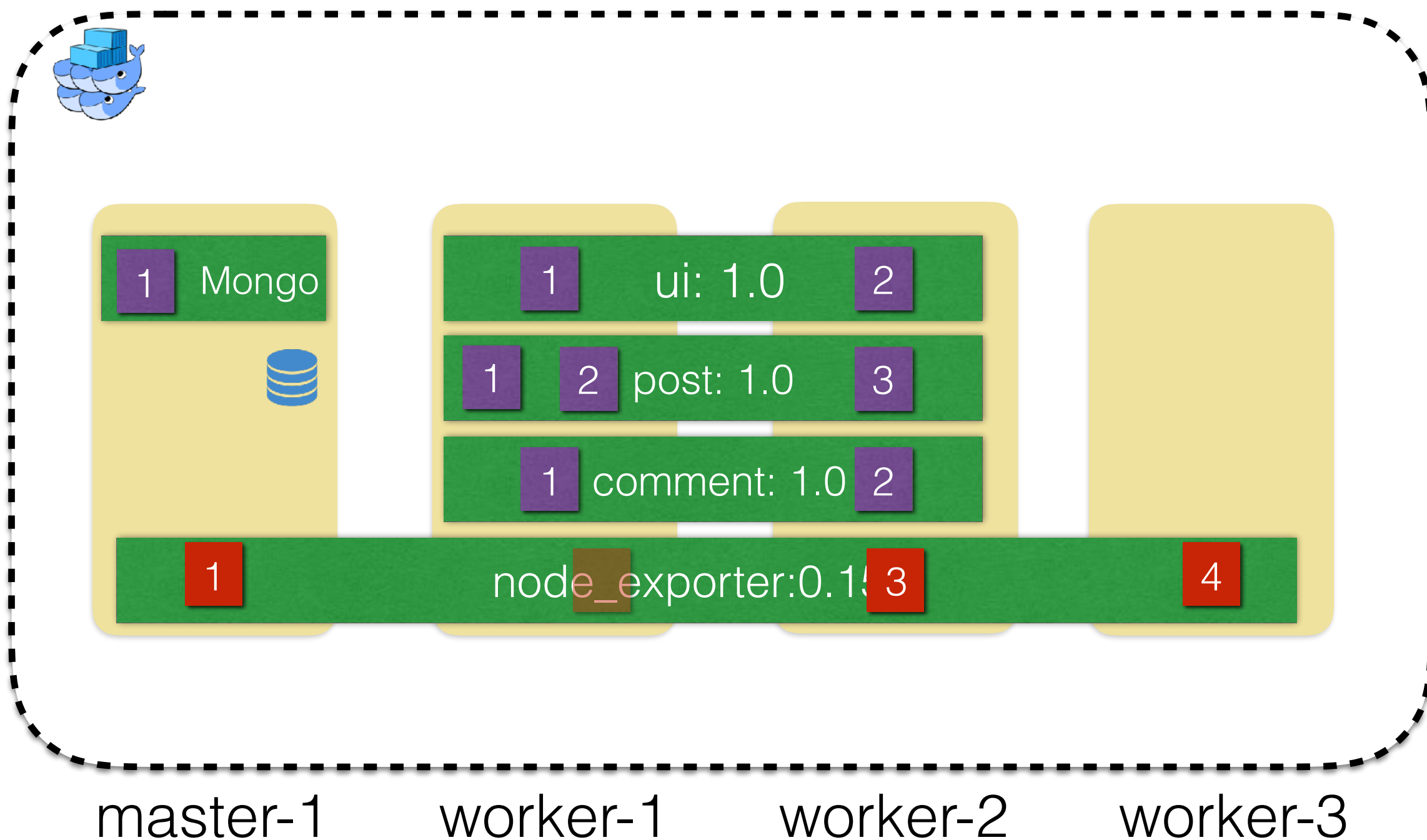
```
node-exporter:  
  image: prom/node-exporter:v0.15.0  
  deploy:  
    mode: global  
  user: root  
  volumes:  
    - /proc:/host/proc:ro  
  ...
```

# Global



reliability=high

# Global

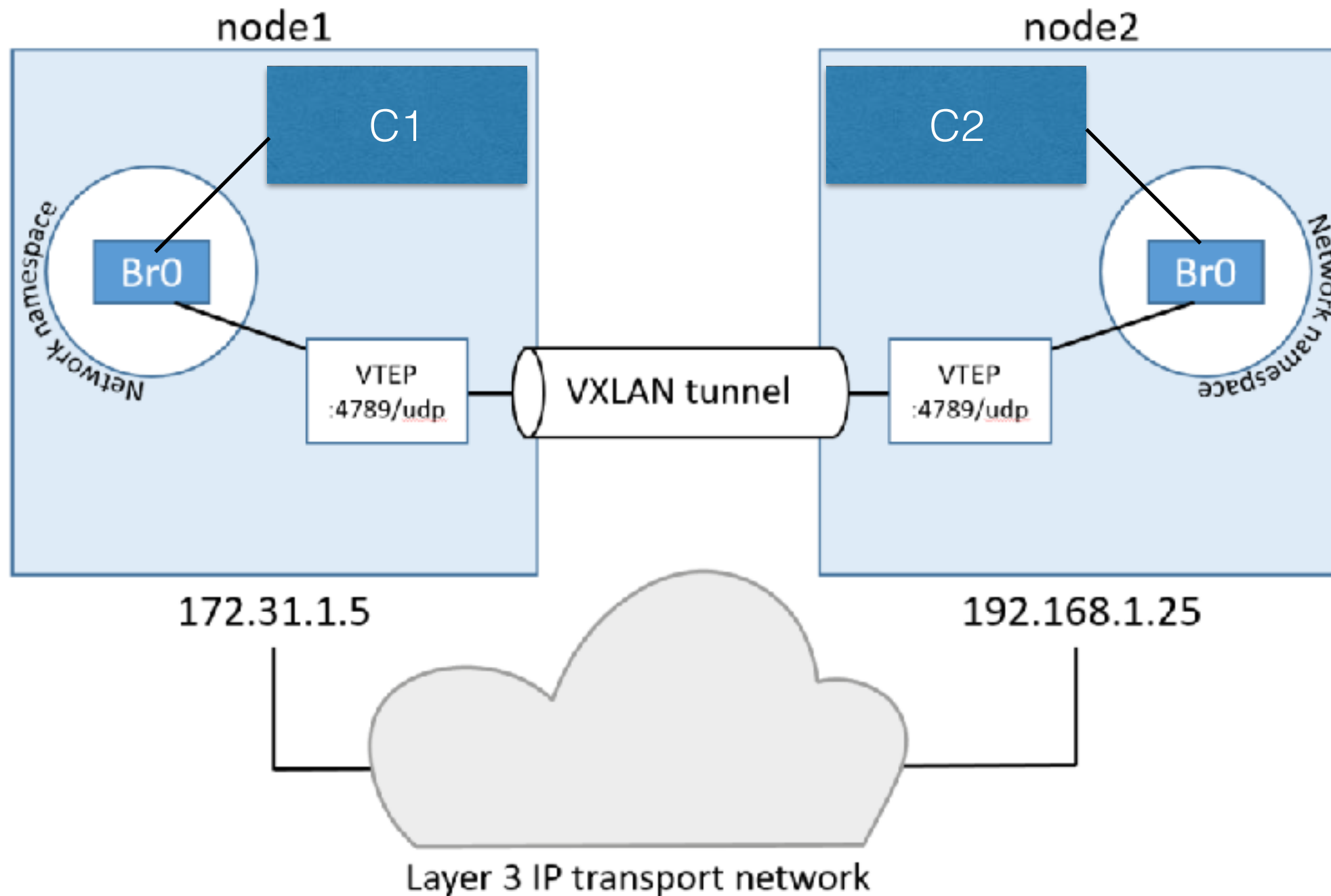


reliability=high

# Overlay network

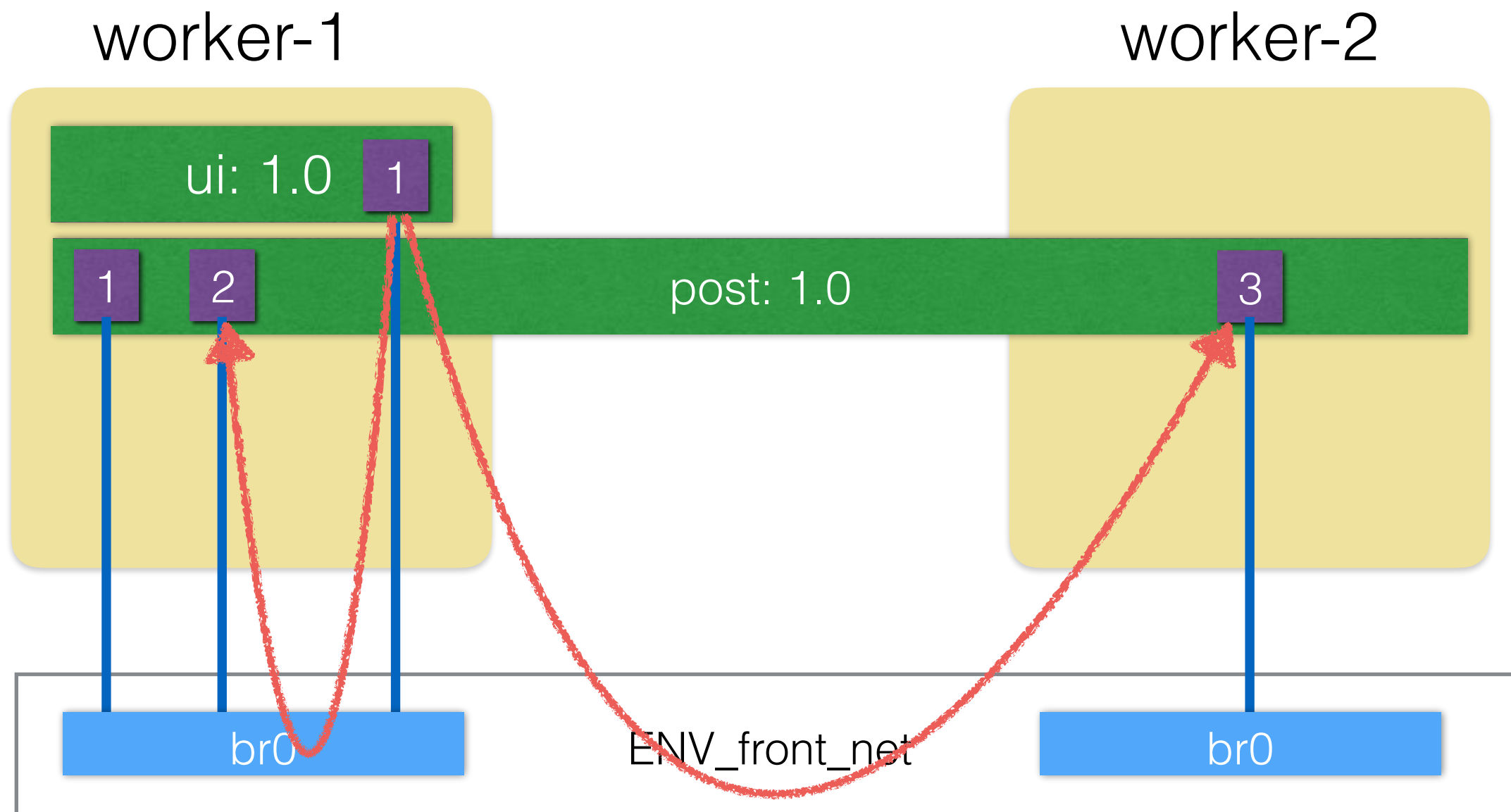
- Overlay network – позволяет объединить в одну сеть контейнеры нескольких докер хостов.
- работает поверх vxlan
- Нужно хранить состояние распределенной сети

# Overlay network





# Как им общаться между собой?



# IP Virtual Server (IPVS)

- Реализация балансировщика нагрузки в ядре Linux
- Включен по-умолчанию (Round Robin LB)
- Сервису назначается 1 VIP (Виртуальный IP-адрес)
- Клиенты не знают, сколько за этим адресом участников

post:

image: \${USER\_NAME}/post:\${POST\_VERSION}

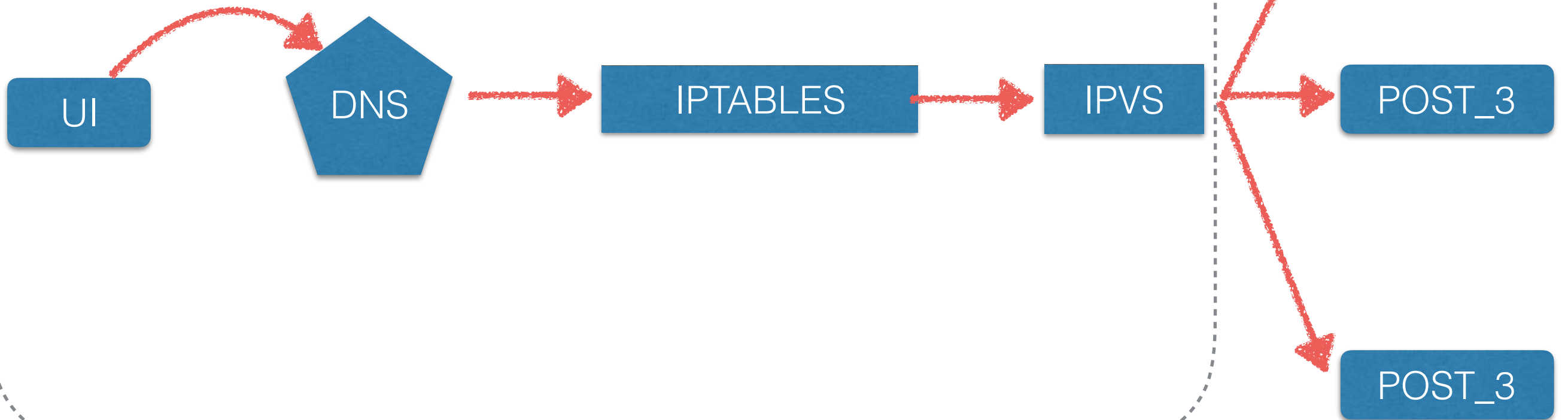
deploy:

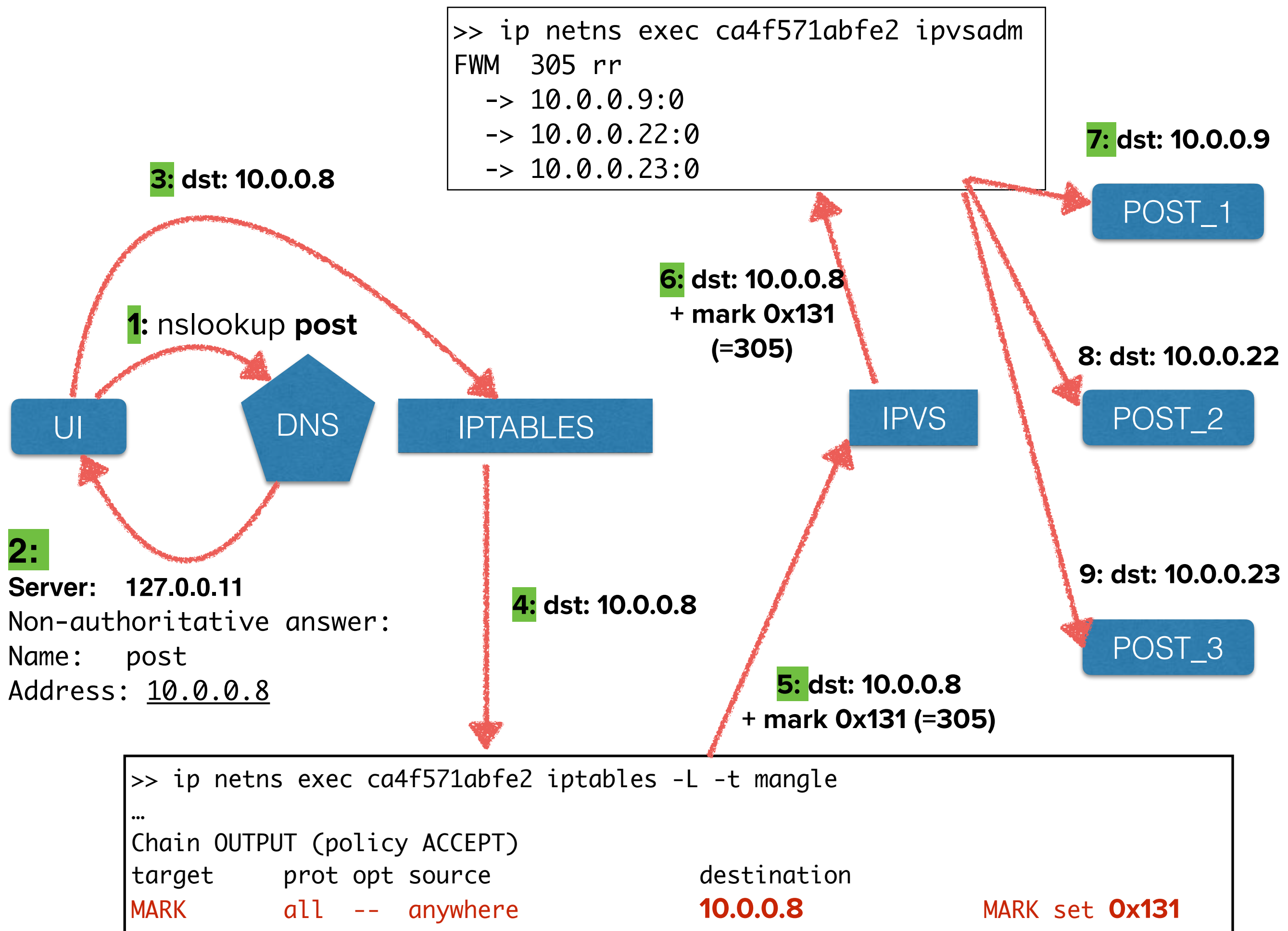
endpoint\_mode: vip

replicas: 3

# IP Virtual Server (IPVS)

ENV\_UI.1\_Network\_namespace





# DNS Round Robin

Включается отдельно.

DNS возвращает адреса всех участников

```
post:
```

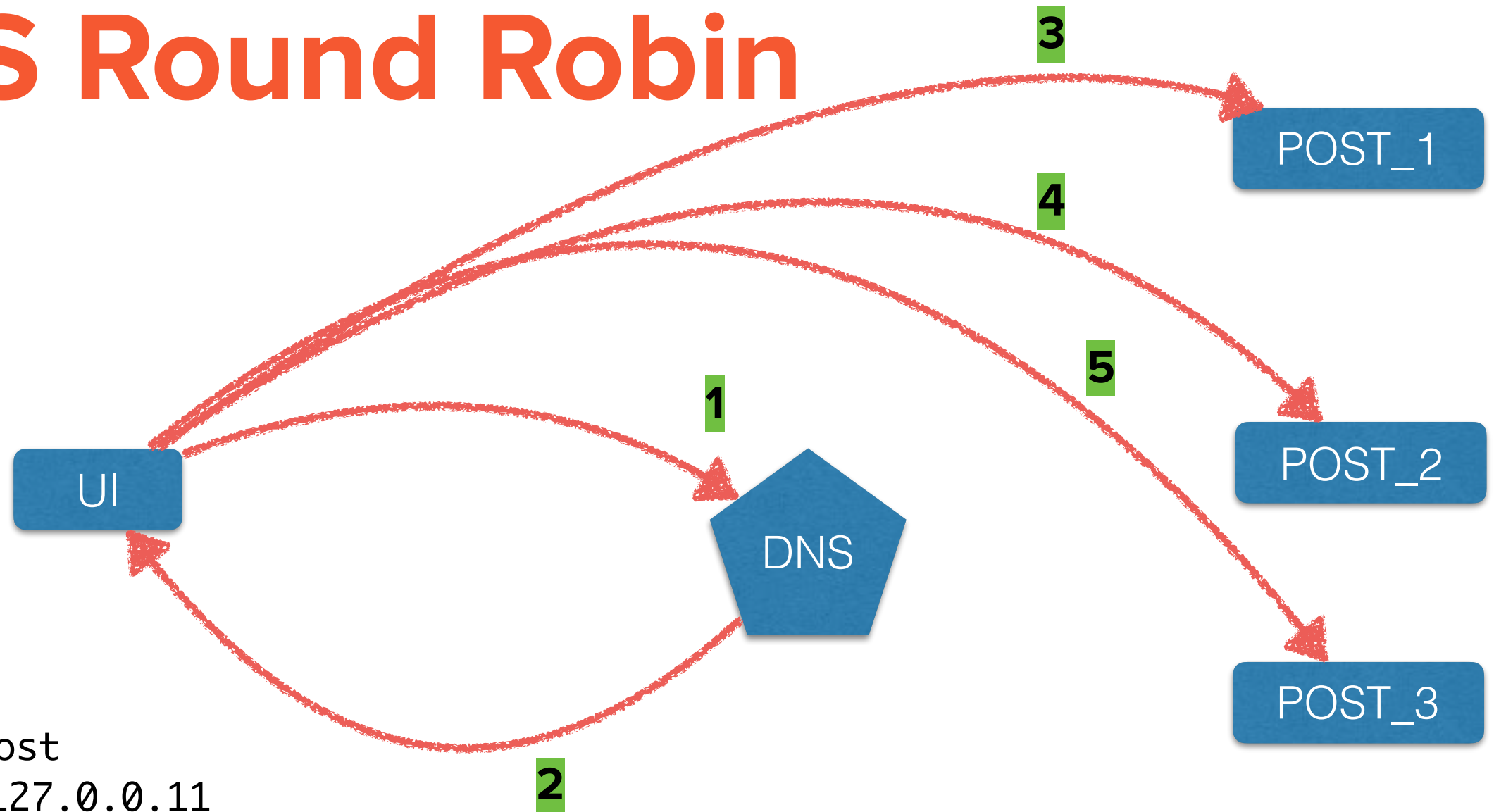
```
  image: ${USER_NAME}/post:${POST_VERSION}
```

```
  deploy:
```

```
    endpoint_mode: dnsrr
```

```
    replicas: 3
```

# DNS Round Robin



```
>>nslookup post
Server:      127.0.0.11
Address: 127.0.0.11#53
```

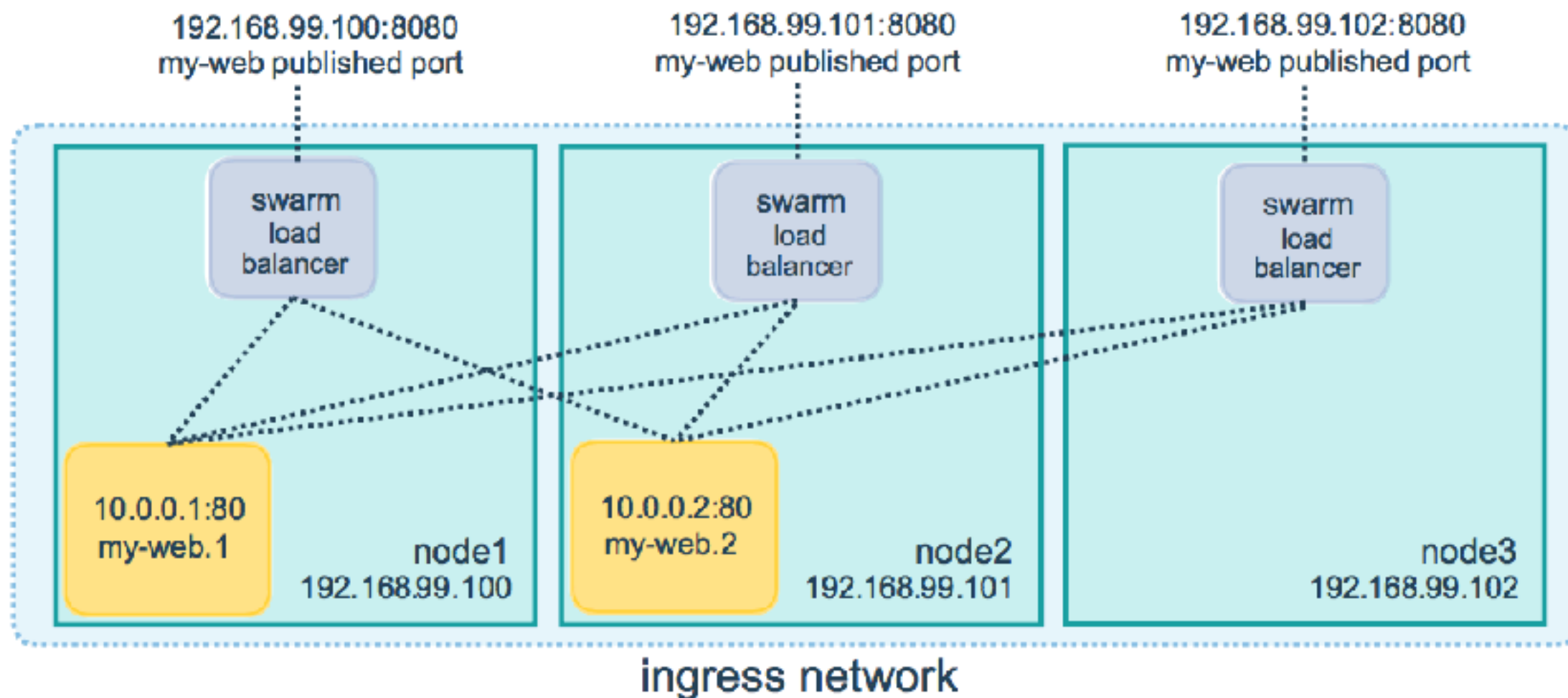
Non-authoritative answer:

```
Name:  post
Address: 10.0.0.29
Name:  post
Address: 10.0.0.28
Name:  post
Address: 10.0.0.27
```

# Как мы общаемся с приложением?

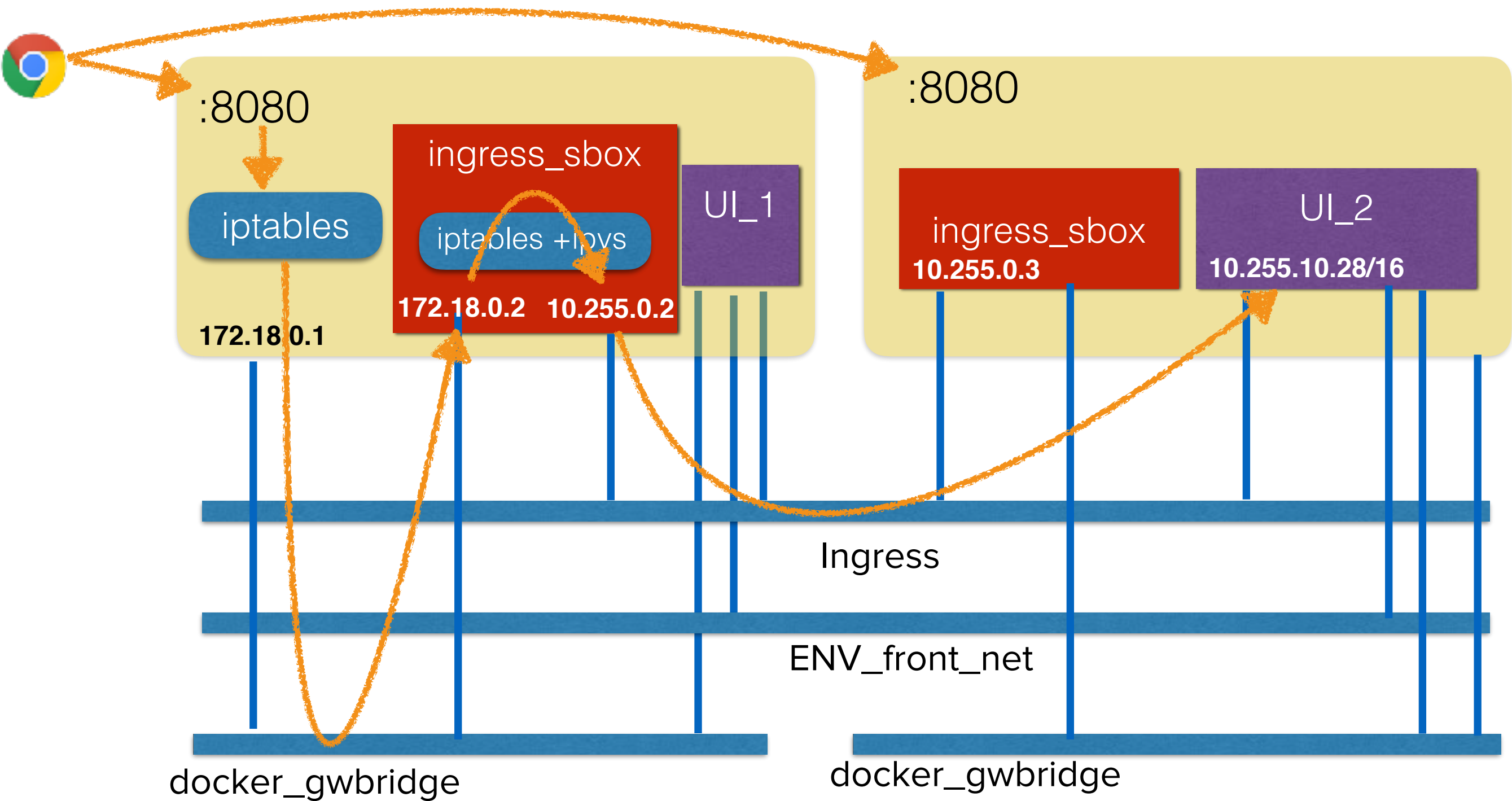
```
ui:  
  image: ${USER_NAME}/ui:${UI_VERSION}  
  deploy:  
    ...  
  ports:  
    - "8080:9292/tcp"
```

# Routing mesh

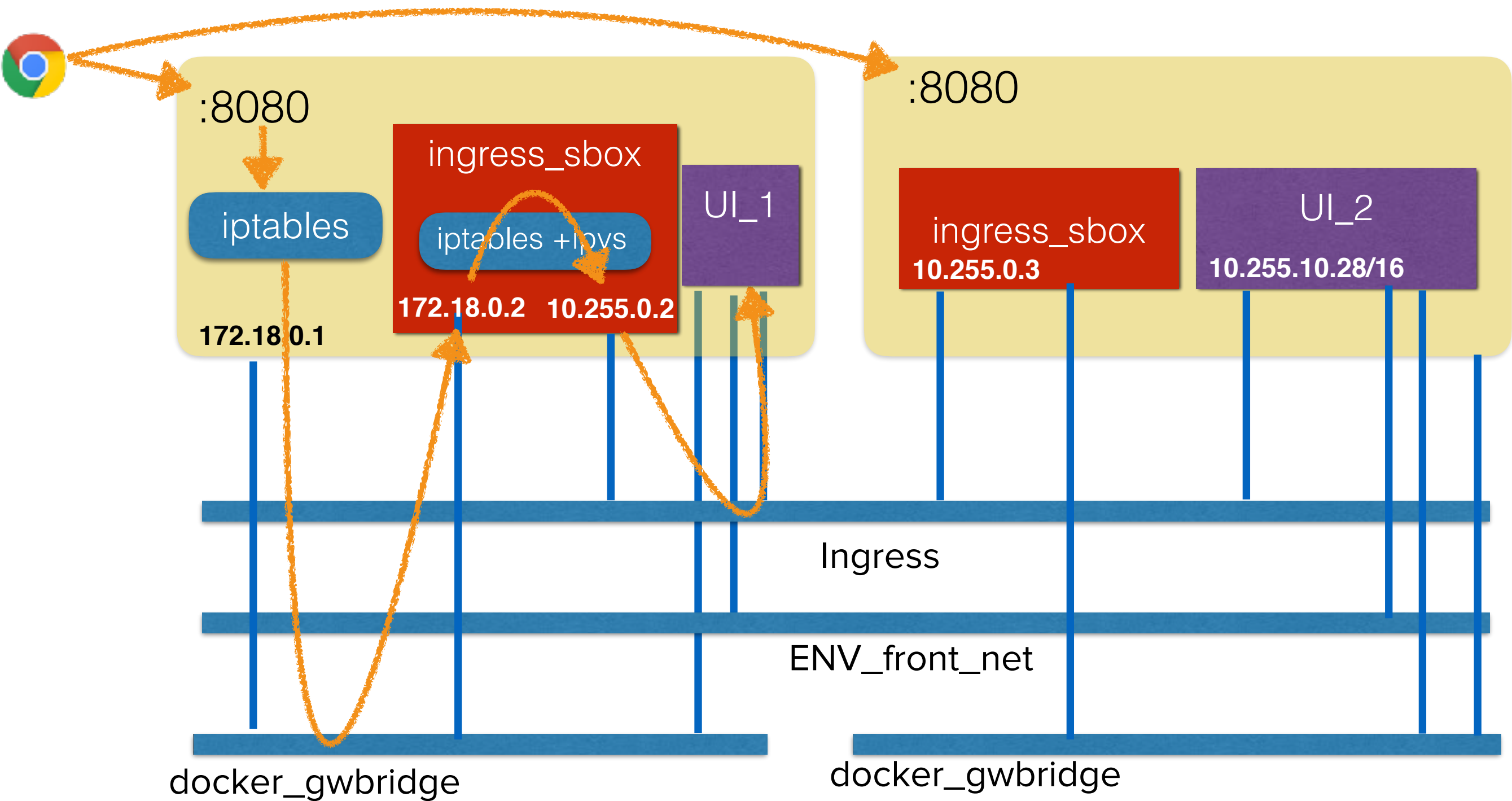




# Routing mesh



# Routing mesh



# Update Service

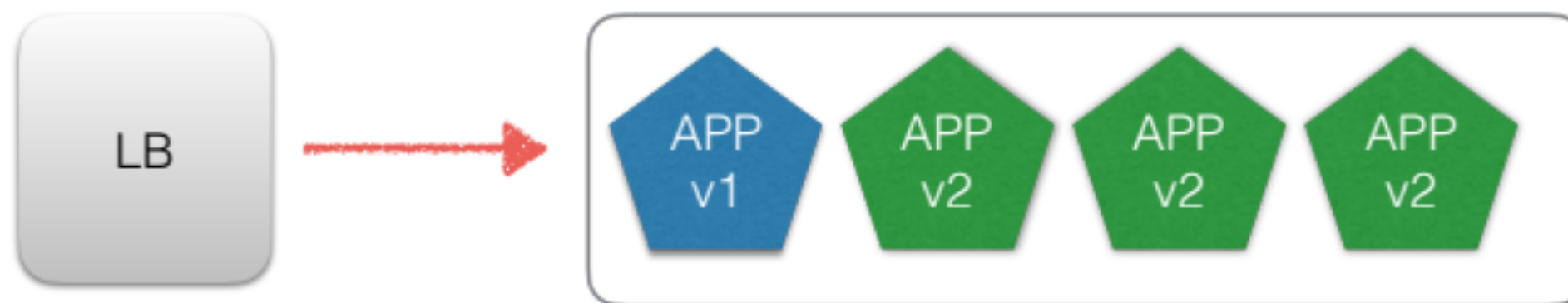
Обновить конфигурацию сервиса:

- `docker stack deploy -c docker-compose.yml ...`
- `docker service update ...`

# Rolling Update



# Rolling Update



# Rolling Update

**как обновлять задачи ?**

```
ui:  
  image: ${USER_NAME}/ui:${UI_VERSION}  
  deploy:  
    update_config:  
      parallelism: 2  
      delay: 5s  
      failure_action: rollback  
      monitor: 5s  
      max_failure_ratio: 2  
      order: start-first
```

# parallelism

Сколько контейнеров обновить одновременно?

```
ui:  
  image: ${USER_NAME}/ui:${UI_VERSION}  
  deploy:  
    update_config:  
      parallelism: 2
```

# Delay

Сколько ждать между обновлениями групп контейнеров?

```
ui:  
  image: ${USER_NAME}/ui:${UI_VERSION}  
  deploy:  
    update_config:  
      parallelism: 2  
      delay: 5s
```



# failure\_action

Что делать, если обновление прошло с ошибкой?

- rollback - откатить все задачи на предыдущую версию
- pause (default) - приостановить обновление
- continue - продолжить обновление

ui:

image: \${USER\_NAME}/ui:\${UI\_VERSION}

deploy:

update\_config:

parallelism: 2

delay: 5s

failure\_action: rollback

# order

В каком порядке обновлять?

- stop-first (default) - остановить старые задачи, перед стартом новых
- start-first - запустить новые задачи, затем остановить старые

ui:

image: \${USER\_NAME}/ui:\${UI\_VERSION}

deploy:

update\_config:

parallelism: 2

delay: 5s

failure\_action: rollback

order: start-first

**Примечание!**

Работает только в compose **3.4**

# Restart policy

Что делать, если задача завершилась ? **Перезапустить!**

Когда ?

- any - **всегда**
- on-failure - перезапускать, если приложение вернуло **код != 0**
- none - никогда

ui:

image: \${USER\_NAME}/ui:\${UI\_VERSION}

deploy:

restart\_policy:

condition: on-failure

delay: 5s

max\_attempts: 3

window: 120s

# Limits

**сколько ресурсов выделить** на задачу (CPU, Memory) ?

ui:

image: \${USER\_NAME}/ui:\${UI\_VERSION}

deploy:

resources:

limits: 

cpu: '0.50'

memory: 50M

reservations: 

cpu: '0.25'

memory: 20M

Сколько **максимум** можно  
на один контейнер

Сколько **минимум** нужно  
на один контейнер

# Healthchecks

Как проверить, что задача работает правильно?

```
ui:
```

```
  healthcheck:
```

```
    test: ["CMD", "curl", "-f", "http://localhost"]
```

```
    interval: 1m30s
```

```
    timeout: 10s
```

```
    retries: 3
```

Либо **HEALTHCHECK CMD** в Dockerfile

# Healthchecks

Что делать если Healthcheck не прошел?