

ArchSummit全球架构师峰会 深圳站2016

七牛自定义数据处理



促进软件开发领域知识与创新的传播



关注InfoQ官方微信
及时获取ArchSummit
大会演讲信息



全球软件开发大会

[上海站] 2016年10月20-22日

咨询热线: 010-64738142



全球架构师峰会 2016

[北京站] 2016年12月2-3日

咨询热线: 010-89880682

自我介绍

- 袁晓沛
- 经历：盛大、七牛、EMC、七牛
- 领域：分布式存储，容器、微服务，大规模数据处理



大纲

- 业务、产品介绍
- 官方数据处理
 - 业务特点、挑战
 - 架构演化
 - 解决方案
- 自定义数据处理
 - 业务特点、挑战
 - 注册、开发、构建
 - 启动、升级、伸缩

业务定义

- 针对海量数据
- 提供零运维、高可用、高性能的数据处理服务
- 日处理数近百亿次
- 让用户轻松应对图片、音视频以及其他各类数据的实时、同步处理场景



处理方式

- 官方数据处理
 - 提供基础的数据处理服务，包括但不限于图片转码、水印、原图保护、防盗链等，及音视频的转码、切片和拼接等。
- 自定义数据处理
 - 允许用户构建、上传自定义的私有数据处理服务，并无缝对接存储在七牛的数据及其他数据处理服务。
- 第三方数据处理
 - 开放应用平台，提供大量功能丰富的第三方数据处理服务，如图片鉴黄、人脸识别、广告过滤、语言翻译、TTS等。



使用方式

`www.clouddn.com/beauty.jpg?facecrop/200x200`

图片URL

UFOP命令

请求参数

原图



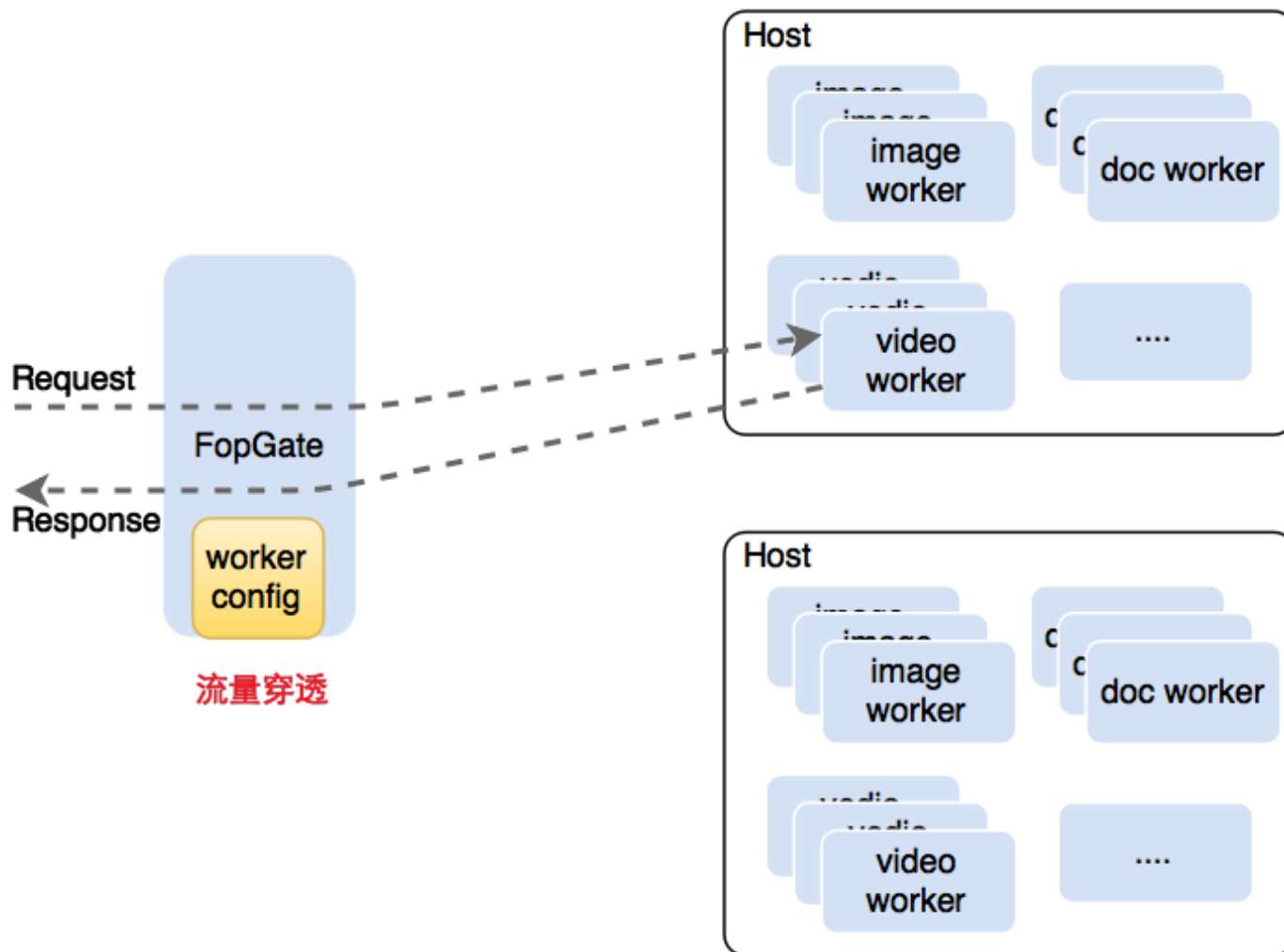
结果



官方数据处理 挑战

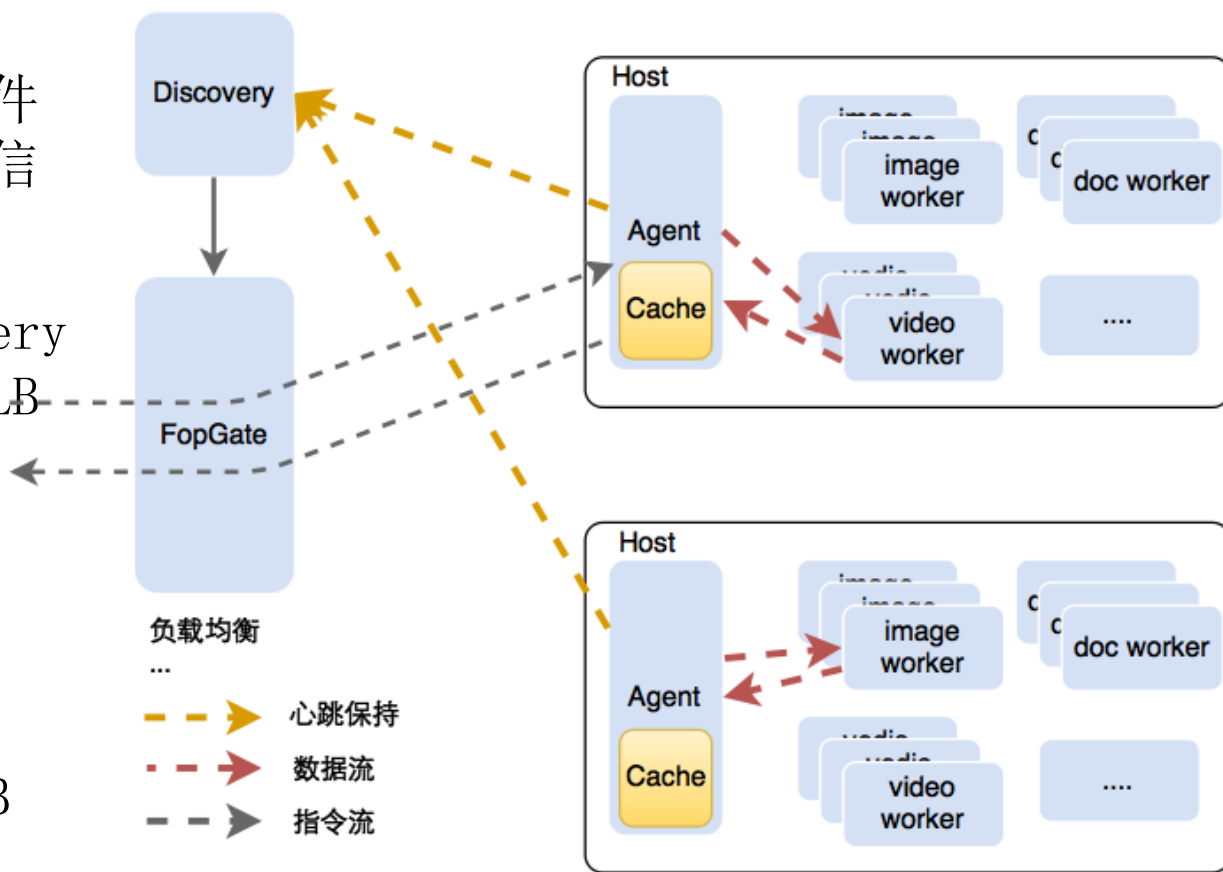
- 请求量非常大
- 突发流量频繁
- CPU密集型计算
- IO操作频繁

官方数据处理 - v1



官方数据处理 - v2

- 增加 Discovery 组件，收集 Agent 上报信息
- FopGate 从 Discovery 获取集群信息，做 LB
- 增加业务 Agent
 - 上报后端信息
 - 上报保活信息
 - 单机内 worker LB



系统测量

- FopGate
 - 单机最大请求数、句柄数
 - 根据实际的业务量，确定机器数
- Image/Audio/Video Worker
 - 找到资源使用最佳范式
 - 根据最佳范式，合理分配资源、配置实例
- 意外发现
 - 大实例、高并发，不如多实例、限制并发
 - 操作系统对CPU调度，比进程好

增加队列

- 服务质量
 - 请求排队，不争抢资源
 - 保证运行速度最快
- 运营角度
 - 根据节点个数、队列长度，
 - 区分免费、付费客户
 - 免费用户，确保高可用
 - 付费用户，确保高质量

限流

- 为什么限流？
 - 大量长链接影响FopGate性能
 - 突发流量，导致队列过长
- 限流手段
 - 并发HTTP请求限制
 - 单用户请求数限制
 - 但Cmd数限制

合理协调IO、CPU

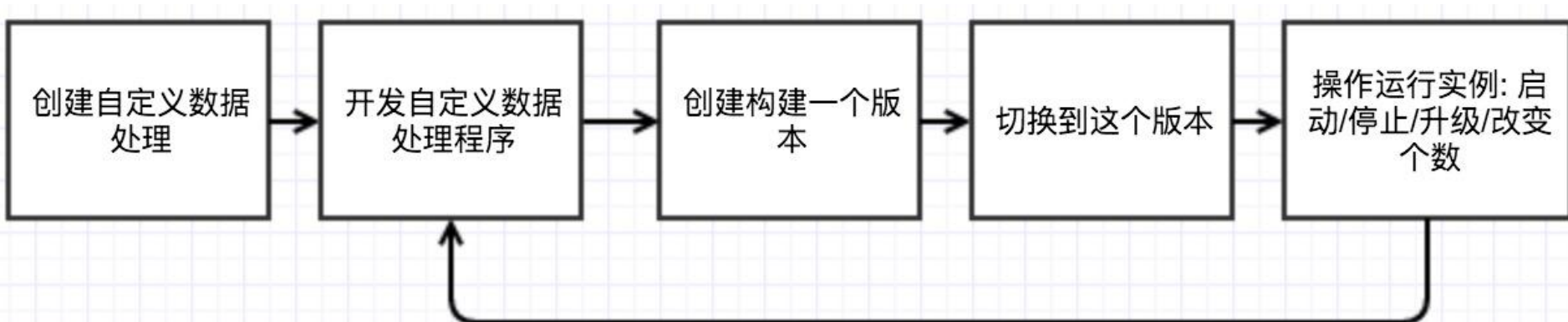
- 为什么？
 - 下载、写盘、处理、写盘、返回
- 协调方式
 - 总原则：就近计算
 - FopAgent、Worker混布(1:N)
 - 缩减网络IO的路由次数
 - 挂载ramfs，将内存当磁盘使用
 - 跳过磁盘IO



自定义数据处理挑战

- 处理程序由客户提供
 - 安全性
 - 隔离性
- 业务规模不确定性
 - 可伸缩性

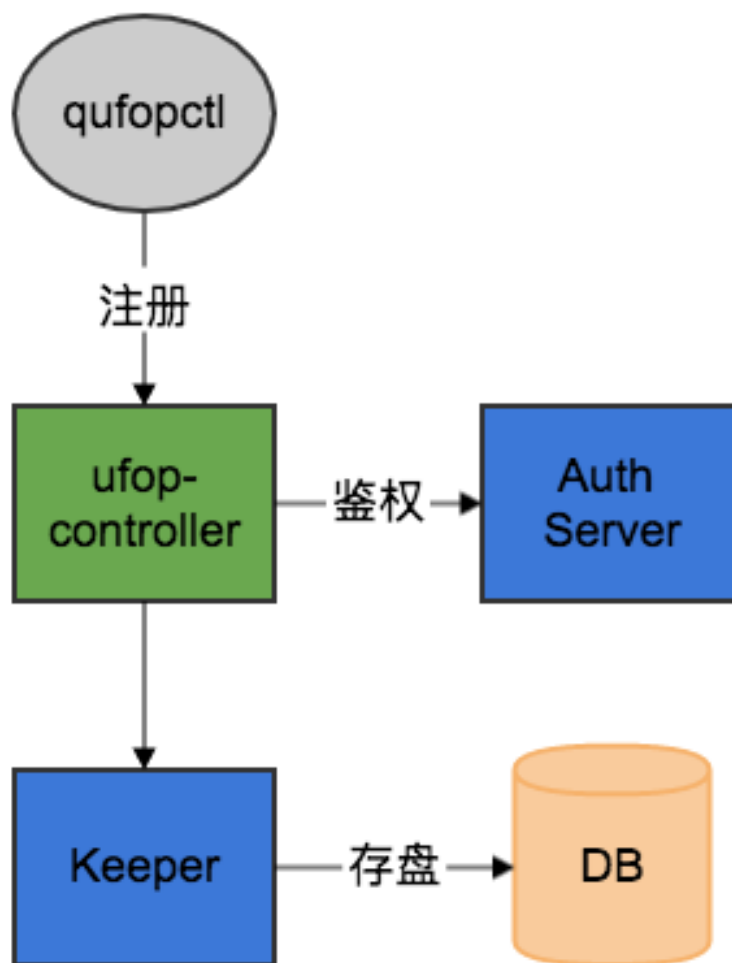
自定义数据处理 - 业务流程



注册

```
$ ./qufopctl reg ufop-demo -m 2
```

注册后端



```

12 type ReqArgs struct {
13     Cmd  string `json:"cmd"`
14     Mode uint32  `json:"mode"`
15     Src  struct {
16         Url      string `json:"url"`
17         Mimetype string `json:"mimetype"`
18         Fsize    int32  `json:"fsize"`
19         Bucket  string `json:"bucket"`
20         Key     string `json:"key"`
21     } `json:"src"`
22 }
23
24 func demoHandler(w http.ResponseWriter, req *http.Request) {
25     body, _ := ioutil.ReadAll(req.Body)
26
27     var args ReqArgs
28     json.Unmarshal(body, &args)
29
30     resp, _ := http.Get(args.Src.Url)
31     defer resp.Body.Close()
32
33     buf := make([]byte, 512)
34     io.ReadFull(resp.Body, buf)
35     contentType := http.DetectContentType(buf)
36     lengthStr := strconv.Itoa(int(resp.ContentLength))
37     w.Write([]byte("Hello World!\n"))
38     w.Write([]byte("The file's mime type is: " + contentType))
39     w.Write([]byte("The file's length is: " + lengthStr))
40 }
41
42 func main() {
43     http.HandleFunc("/uop", demoHandler)
44     err := http.ListenAndServe(":9100", nil)
45     if err != nil {
46         log.Fatal("Demo server failed to start:", err)
47     }
48 }

```

ufop. yml

```

1 image: ubuntu
2 build_script:
3   - chmod a+x ufop-bin
4   run: ./ufop-bin

```

ufop. tar

```

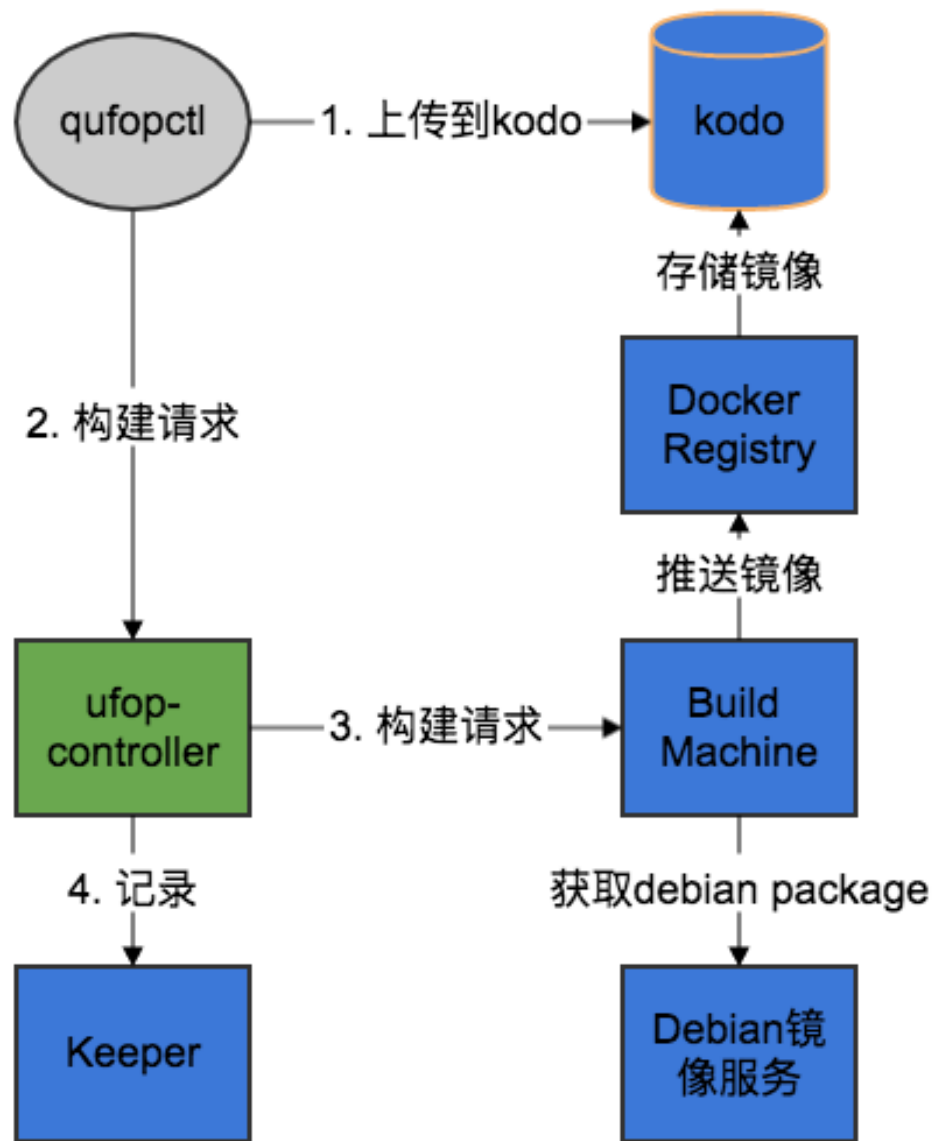
berty@udev:~$ tree ufop-demo-folder/
ufop-demo-folder/
├─ ufop-bin
└─ ufop.yaml

```

构建

```
$. /qufopctl build ufop-demo -d ./ufop-demo-  
folder
```

构建后端



使用Debian镜像服务

AppRox

- 经常下载超时
- 下载出错后，需要手动清除

Debian Pkg Mirror

- 首次全量下载
- 定时增量更新

避免Docker构建缓存

Wrong

```
-RUN curl -o jdk.tar.gz https://dn-qcos.qbox.me/jdk-7u15-linux-x64.tar.gz
```

```
-RUN mkdir -p opt && tar -xf jdk.tar.gz -C /opt && rm -rf jdk.tar.gz
```

Correct

```
+RUN curl -o jdk.tar.gz https://dn-qcos.qbox.me/jdk-7u15-linux-x64.tar.gz \
```

```
+      && mkdir -p opt \
```

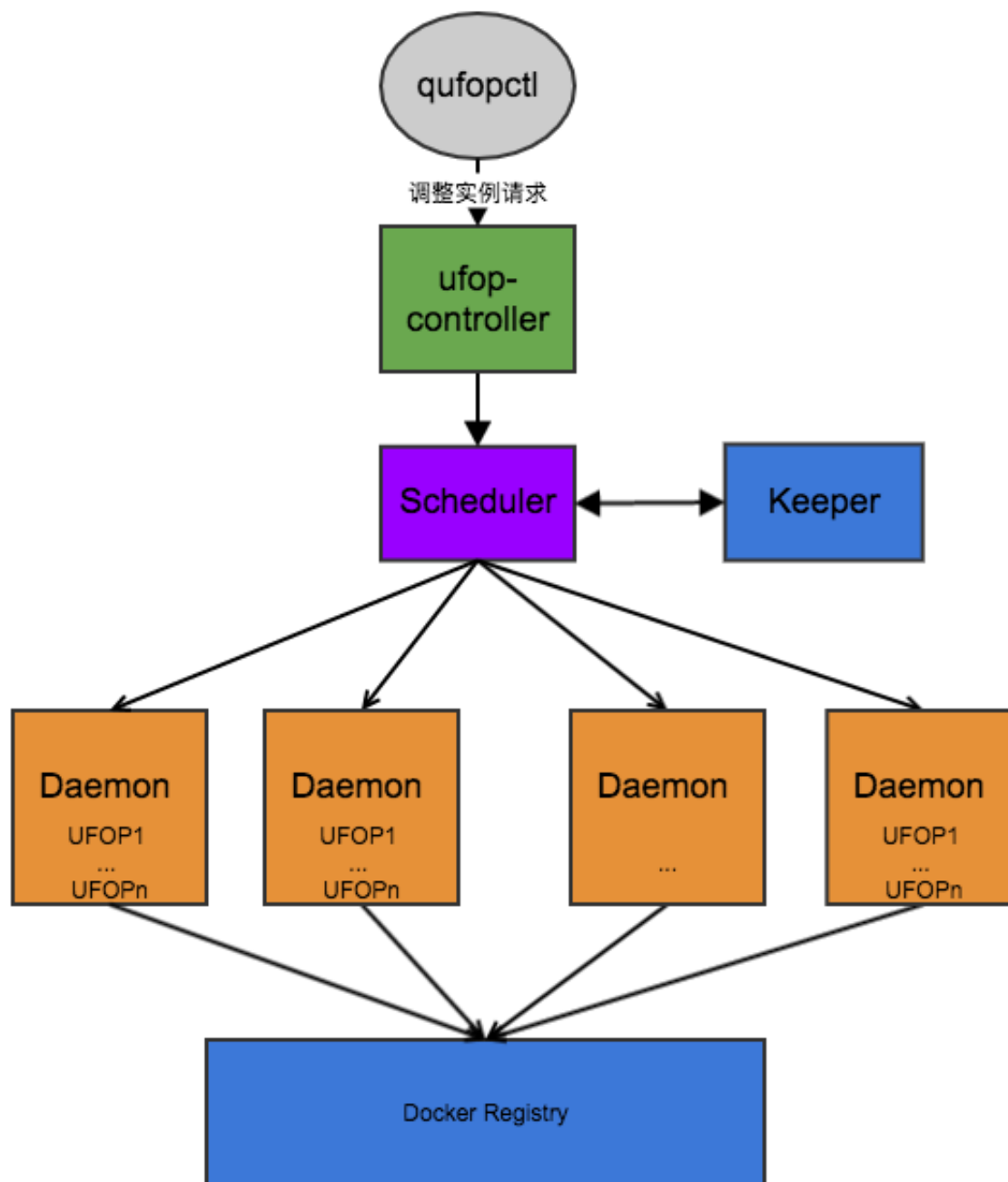
```
+      && tar -xf jdk.tar.gz -C /opt \
```

```
+      && rm -rf jdk.tar.gz
```



调整实例数

```
$. /qufopctl resize ufop-demo -n 3
```



升级实例

```
$. /qufopctl upgrade ufop-demo -r 1:2
```

灰度升级阶段

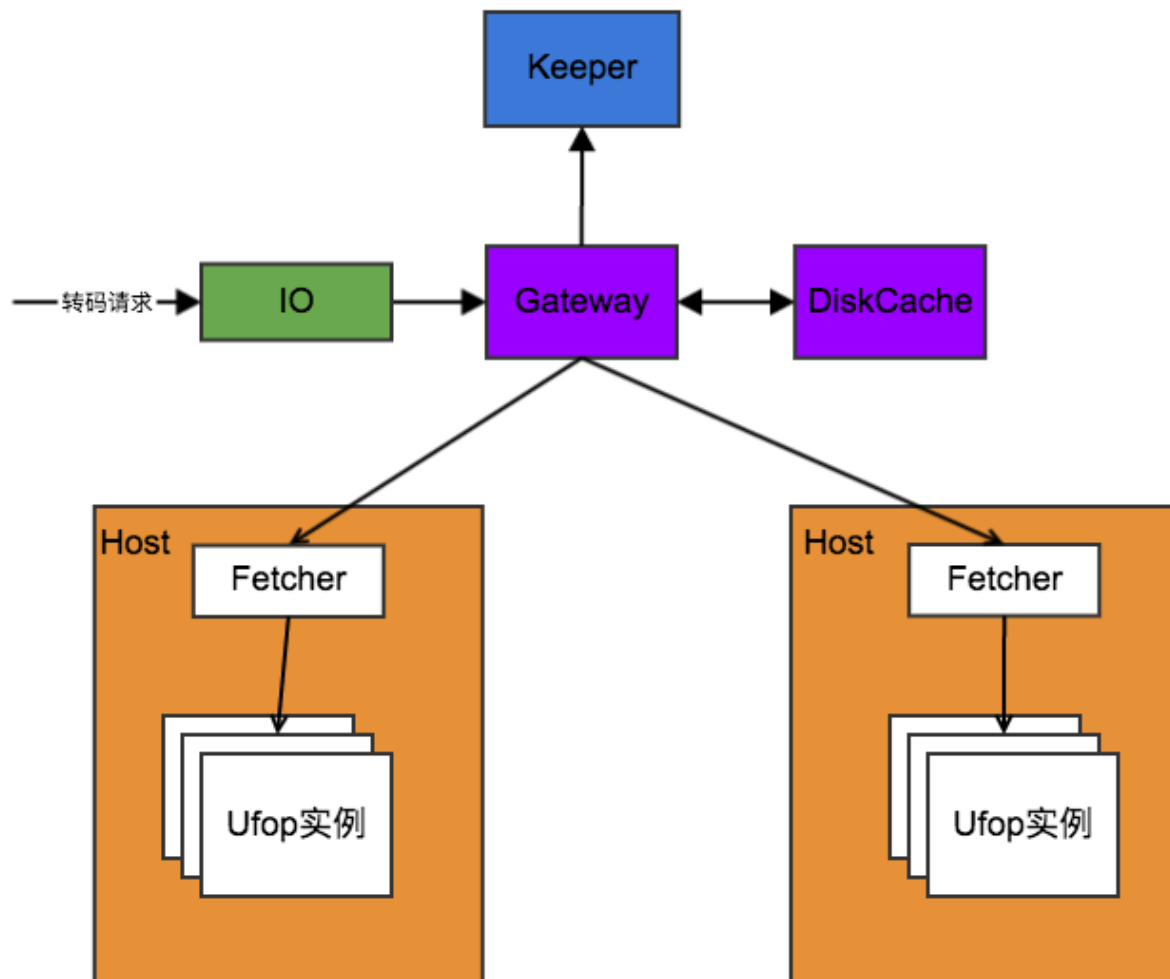


升级的细化

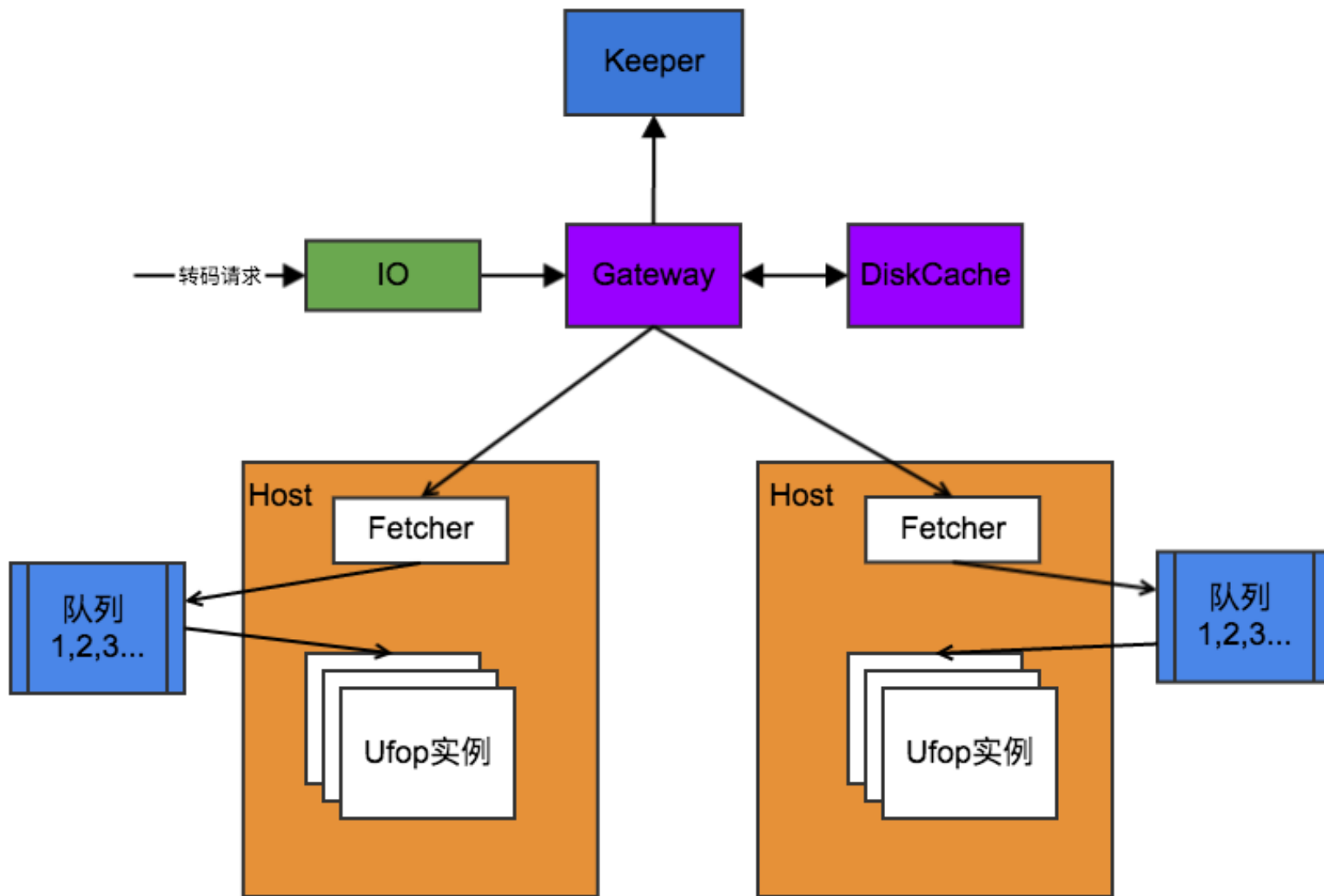
- 新实例WarmUp
 - 内存池、线程池、连接池初始化，初始请求太慢
 - 设定预热时间段，期间请求权重比正常小一点
- 老实例CoolDown
 - 老的请求正在处理，直接停掉影响可用性
 - 应用Docker StopWait
- 计算冗余
 - 预留足够的计算冗余
 - 升级步长 $<$ 冗余实例数



数据流 - v1



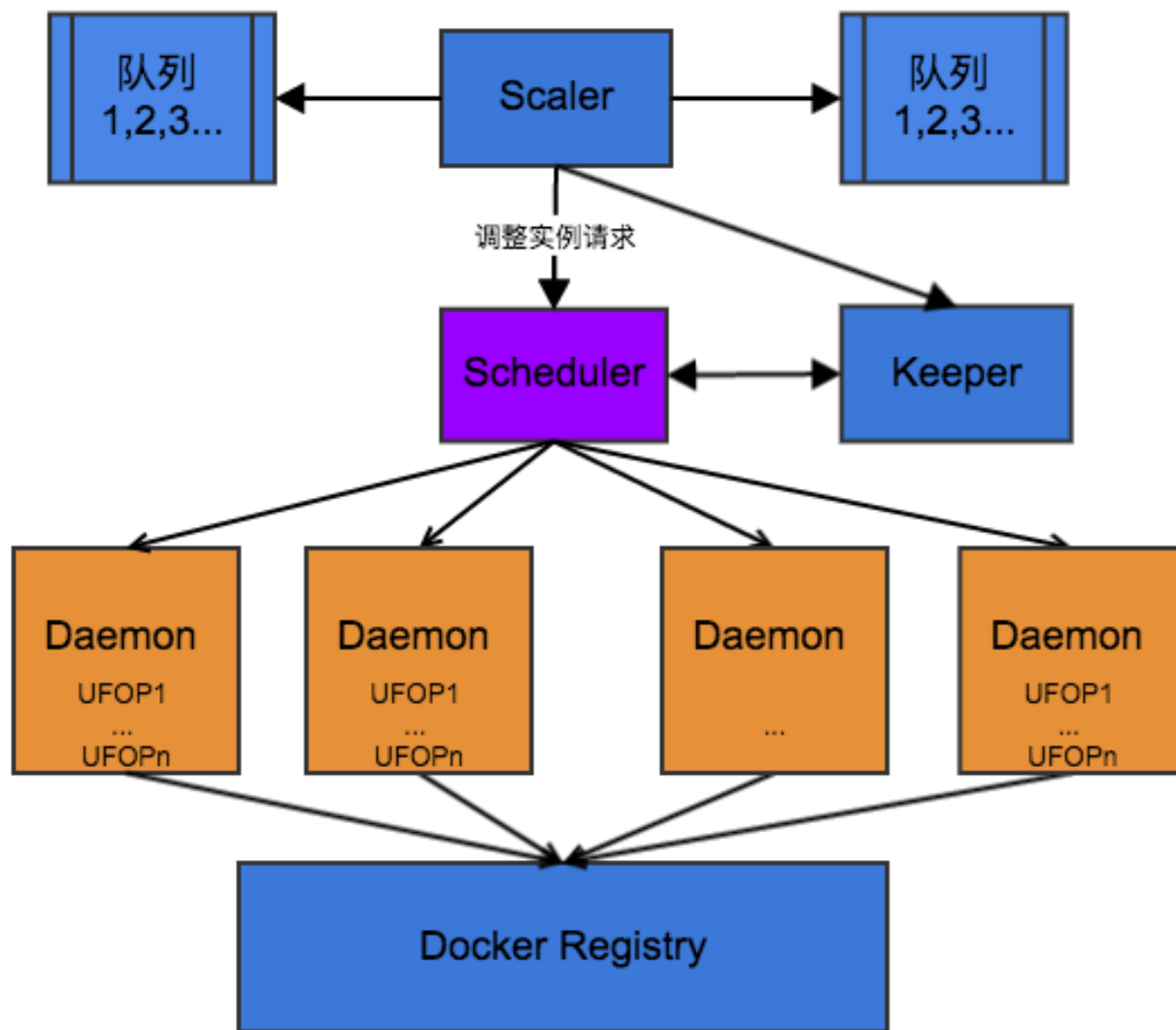
数据流 - v2



自动伸缩设置

- 用户配置
 - 默认实例数
 - 平均单实例待处理任务数
 - 是否自动伸缩
- 自动伸缩
 - 增大、或者缩小实例数，以保证：
 - 平均单实例待处理任务数

自动伸缩后端



解决方案

- 安全性
 - 借助iptables
- 隔离性
 - 借助容器的cgroup
- 可伸缩性
 - 实现容器调度系统，支持秒级伸缩
 - 暴露伸缩API，手动伸缩
 - 利用队列长度，自动伸缩



Thanks!

