



大数据云和应用生态的构建之路

陈夏明 博士

星环信息科技（上海）有限公司
www.transwarp.io



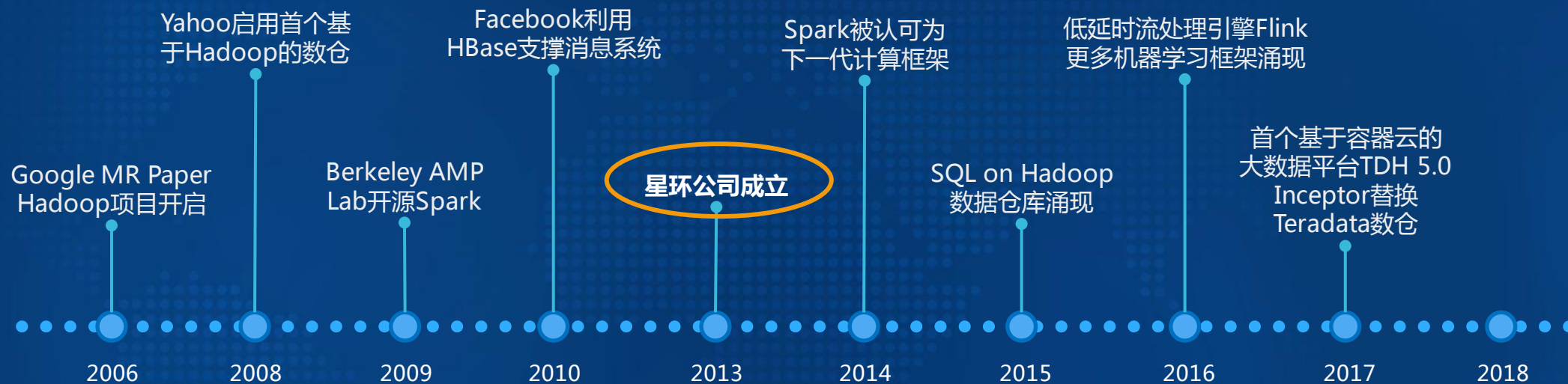
内容

- ◆背景
- ◆复杂应用编排
- ◆数据云架构演进
- ◆Q&A



大数据和云技术发展历程

大数据



云计算





大数据市场格局

Gartner Visionary (Global), Feb. 2016

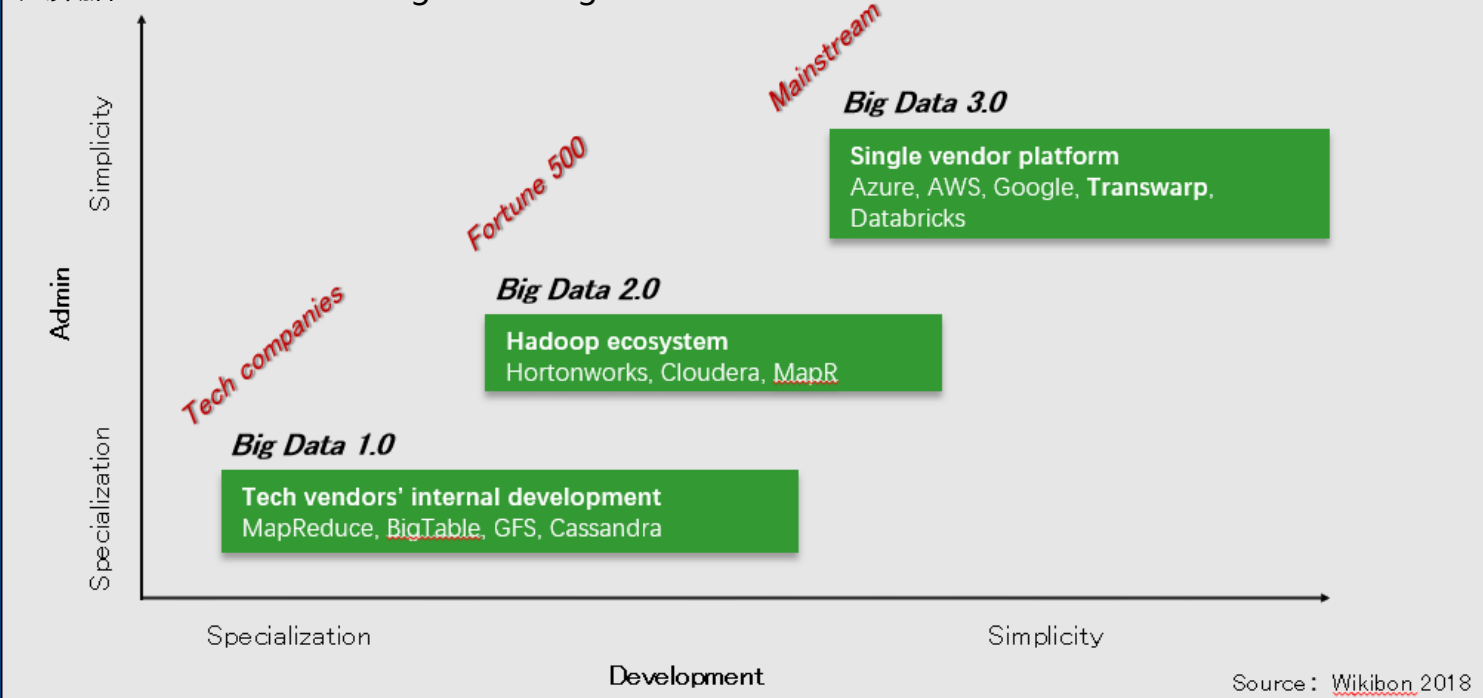


IDC Market Scape (China) , Oct. 2017



大数据行业发展趋势

大数据3.0 - Artificial Intelligence + Big Data + Cloud



BigData 3.0：提供“ABC”技术融合的综合平台，满足客户多元化、复杂化的需求，同时降低开发和管理难度，在整个产业链占据更重要地位。

数据生态化

- 数据驱动业务闭环
- 服务和应用共享
- 数据对外赋能

数据服务化

- 数据化运营
- 智能应用
- 在线数据服务

数据资产化

- 多源数据融合
- 数据质量管理
- 资产化与计量

数据集约化

- 数据集中采集
- 统一的元数据
- 统一的计算平台



大数据3.0应用生态





复杂应用编排





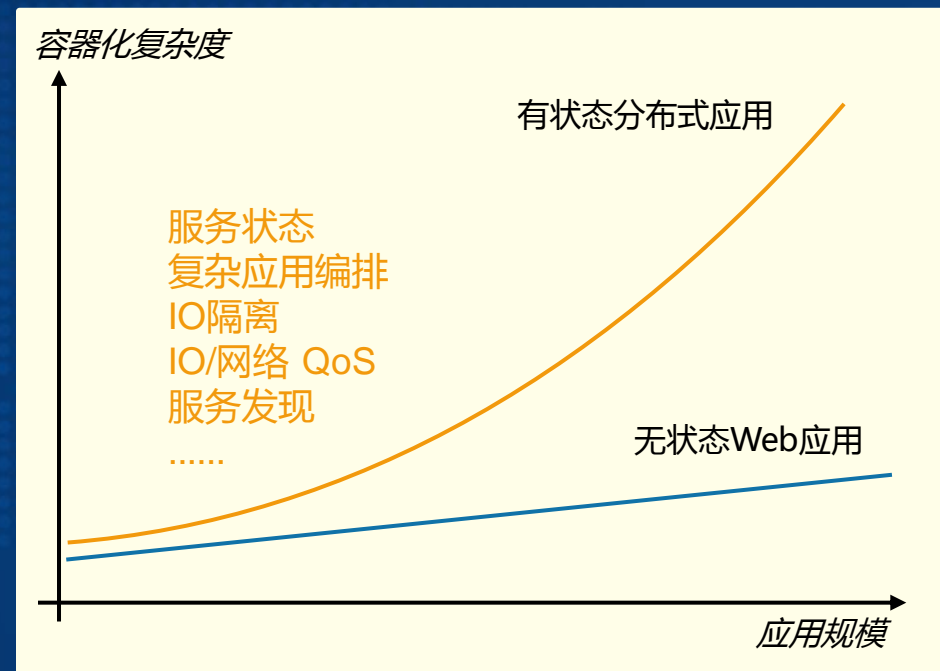
应用部署技术变迁





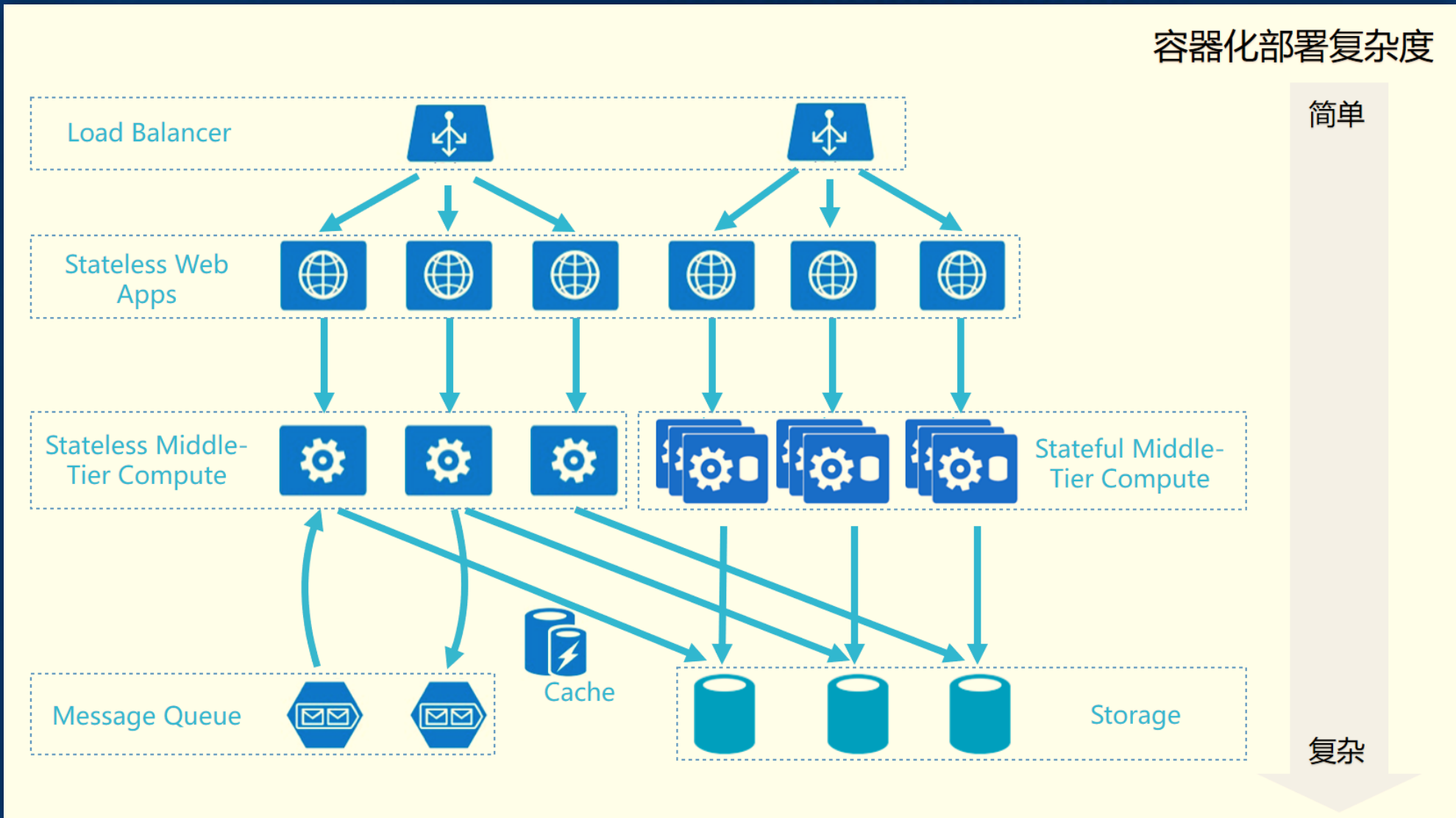
应用编排

- ◆ 协调多个服务，从而对外提供更强大的功能
- ◆ 编排的问题范畴
 - ◆ 配置
 - ◆ 系统配置 (cpu, memory, storage)
 - ◆ 服务配置 (容器内应用的配置)
 - ◆ 调度
 - ◆ 网络
 - ◆ 存储
- ◆ 编排工具
 - ◆ docker-compose, swarm, mesos, kubernetes etc.

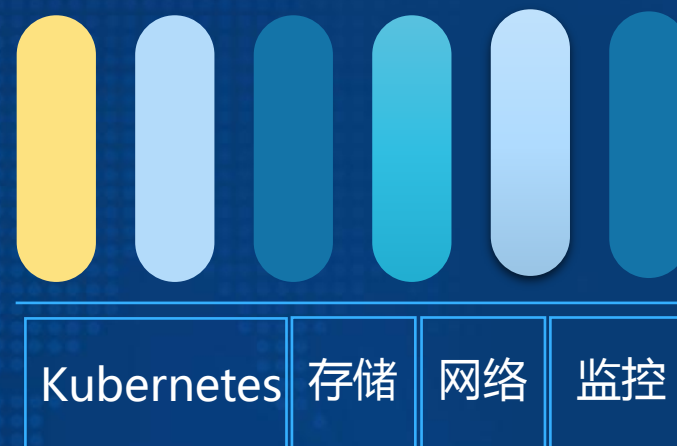
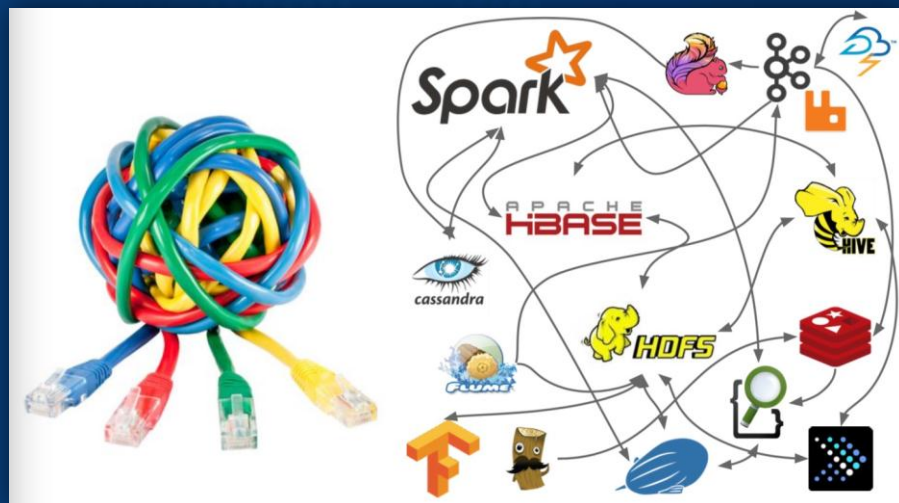




不同应用的容器化难度



大数据应用的服务拆分



- ◆ 提前规划节点的服务
- ◆ 一个物理集群很难安装多个Hadoop集群
- ◆ 容易与第三方应用有兼容性问题
- ◆ 升级操作复杂

传统部署

- ◆ 每个服务按需灵活扩容/缩容、上线/下线、升级/降级
- ◆ 拆分完之后，解决服务间依赖和服务发现的问题
 - ◆ 集群内使用Headless service
 - ◆ 服务之间通过DNS和配置注入进行服务发现
 - ◆ 集群外的访问
 - ◆ NodePort
 - ◆ REST服务 + LoadBalancer

容器部署

容器镜像标准化，e.g. Hadoop



启动方式

方案一：配置打包到容器，通过变量替换修改

- 优点：配置与image匹配，本地测试比较容易
- 缺点：修改/扩展配置比较麻烦



方案二：配置通过ConfigMap的方式挂入

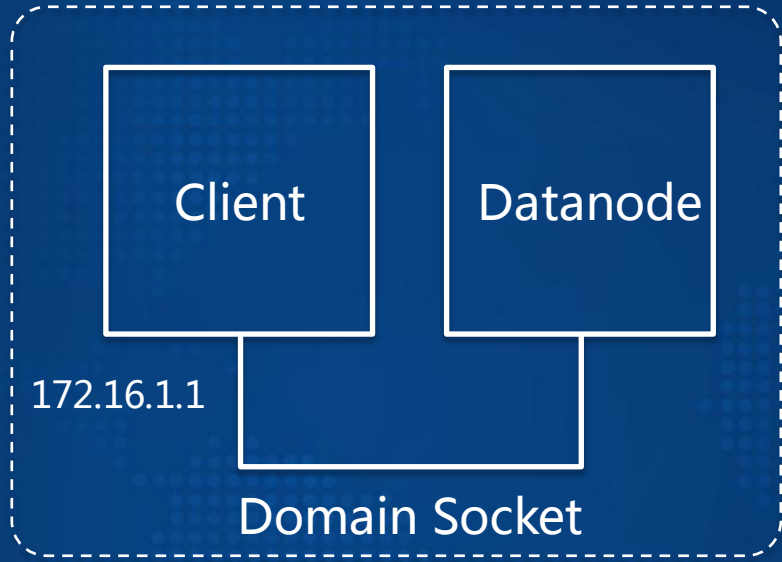
- 优点：修改/扩展配置容易
- 缺点：需要配置管理模块，本地测试比较麻烦



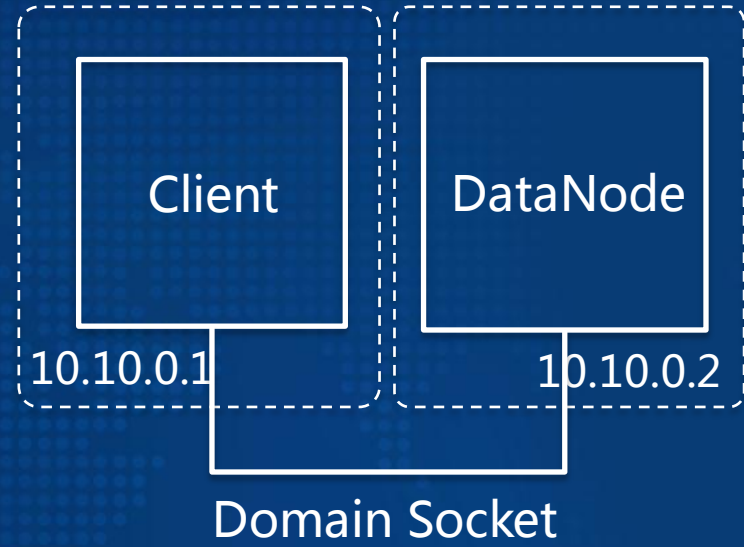
配置文件



Data Locality



- 通过判断IP是否一致决定是否进行本地读写
- Domain socket优化，免去本地IO走tcp协议栈



- 每个容器有一个独立IP
- 每个容器独立的文件系统

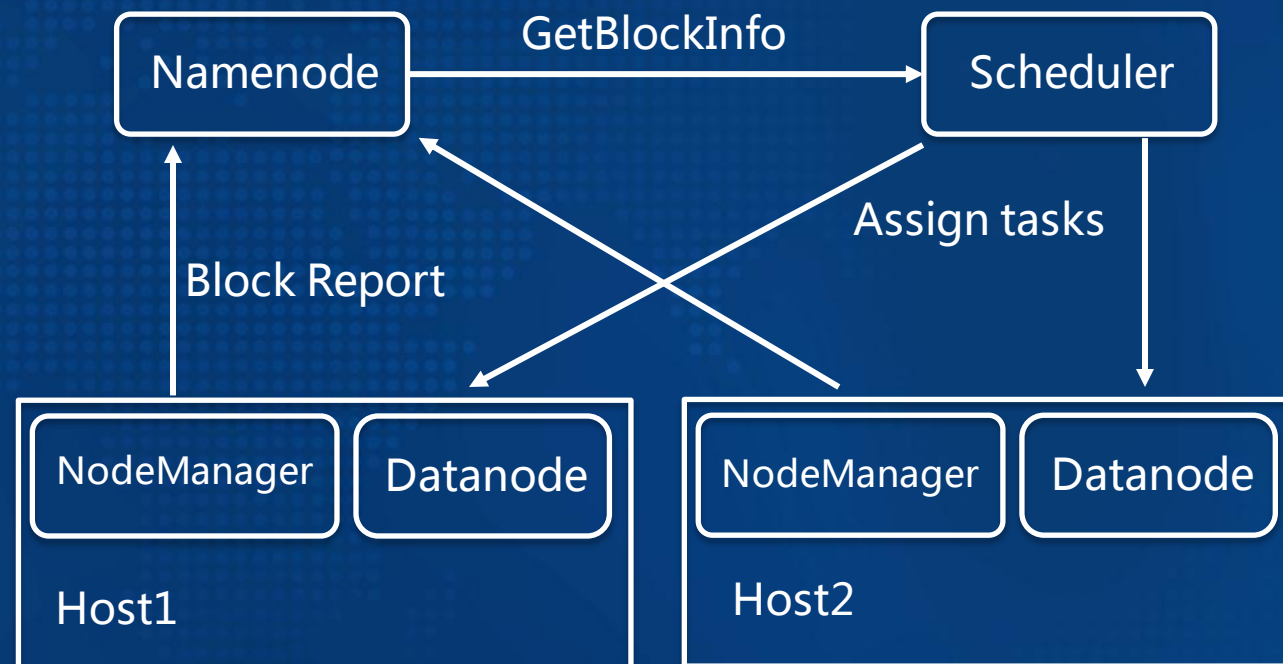
需要的改进：

- 更改Locality 判断逻辑
- 多个pod共享Domain socket，而且每个租户的hdfs的domain socket互相独立



Yarn/Spark 调度逻辑

- ◆更改Yarn/Spark调度逻辑，
判断对应的计算切片是否在
同一台host上





复杂应用编排策略

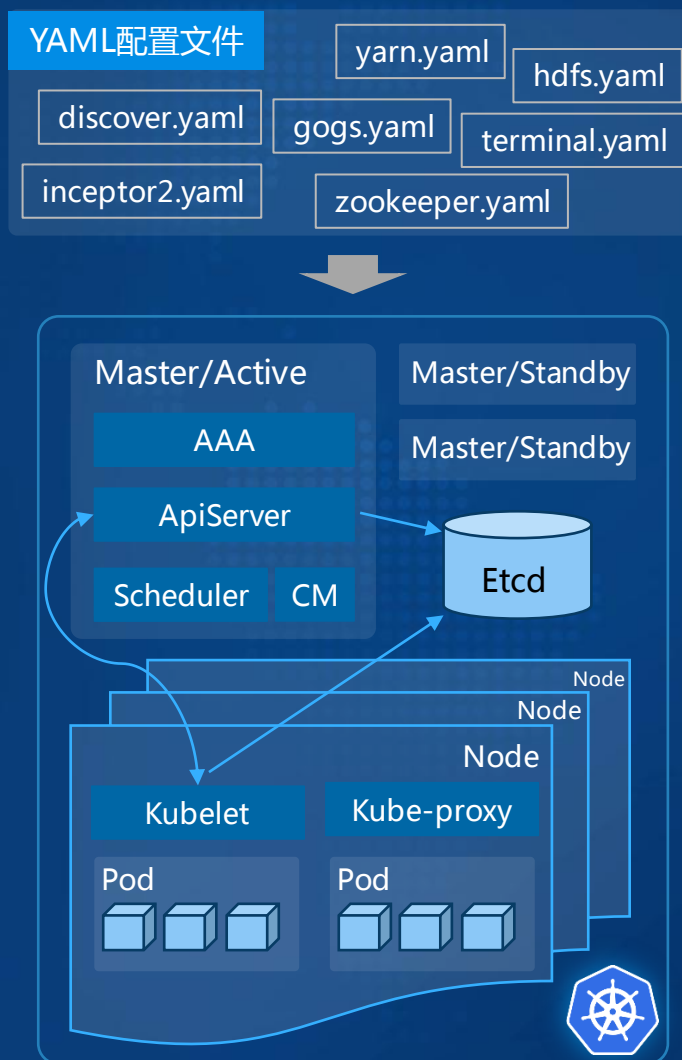
- ◆ 服务拆分（细分应用程序为不同的微服务，并包装在各自的容器中）
- ◆ 镜像标准化（镜像启动入口、配置文件挂载方式、CI/CD流程化）
- ◆ 动态编排（优化这些容器的资源利用率）
- ◆ 调度策略（Disk/Network Throughput, DiskSize, DiskType, CPU, Memory, Priority）
- ◆ 亲和性（计算容器和存储容器在同一个节点，降低网络IO开销）
- ◆ 反亲和性（多副本在不同节点，保证可用性）
- ◆ 服务优先级（实时计算任务和批处理任务拥有不同的等级）
- ◆ 资源超售（在不影响应用正常运行的情况下，降低成本）
- ◆ P2P镜像分发



数据云架构演进

Image source: <https://www.wired.com/story/container-ships-use-super-dirty-fuel-that-needs-to-change/>

数据云架构 1.0



◆Kubernetes, v1.2

- ◆ 早期版本，框架很灵活，功能相对较弱
- ◆ 产品化程度低，需要用户手动配置大量YAML参数
- ◆ 需求多实例、多租户
- ◆ 需求依赖注入
- ◆ 需求Image更换
- ◆ 需求特定服务开关，如Kerberos安全
- ◆



大数据应用编排 1.0

模板字符串替换

- Python模板Jinja/Django
- Java模板FreeMarker

带来的挑战

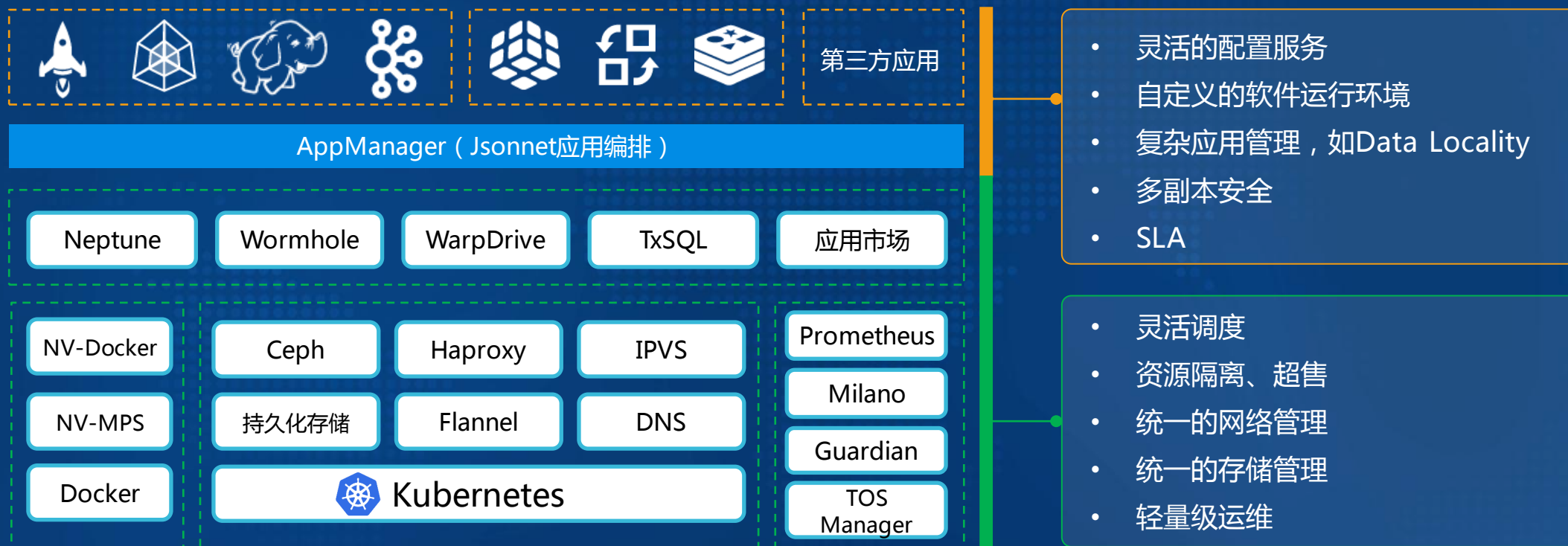
- ◆ 通用功能的维护代价大，复用度底
 - ◆ k8s模板中相似的部分需要复用
 - ◆ 如node selector, network name, 定制存储卷等
- ◆ 测试复杂
 - ◆ 模板渲染是否符合预期
 - ◆ 模板渲染出的YAML是否符合k8s语法
- ◆ 复杂的逻辑需要借助编程语言，如Python/Java
- ◆ 版本衍生代价大

```
Please edit this file and save it to the same location, then the file will be
and an empty file. The file will be reopened with the same permissions.

resourcequota
* : Invalid value for 'resourceVersion': invalid character '0' in
invalid character '0' in '66247062'

apiVersion: v1
kind: ResourceQuota
metadata:
  creationTimestamp: 2018-11-23T02:50:05Z
  name: billing.quota
  namespace: billing
  resourceVersion: "66247062"
  selfLink: /api/v1/namespaces/billing/resourcequotas/billing.quota
  uid: 7b269430-eeca-11e8-ab82-0cc47ae2fc6e
spec: {
  hard:
    limits.cpu: 1k
    limits.memory: 1028991300Ki
    requests.cpu: 1k
    requests.memory: 1028991300Ki
  prioritizedHards:
    - priority: 100
      resources:
        limits.cpu: "100"
        limits.memory: 100Gi
        requests.cpu: "100"
        requests.memory: 100Gi
  status: {}
```

数据云架构 2.0



版本总结：配置灵活，但是普通用户上手难度较大；面向业务程度低



大数据应用编排 2.0

◆ Jsonnet模板

- ◆ Google开源的一门动态配置语言，增强并兼容JSON
- ◆ 在JSON基础上增加新特性：
 - ◆ 注释
 - ◆ 引用
 - ◆ 算术与逻辑运算
 - ◆ 条件操作符
 - ◆ 数据与对象深入
 - ◆ 引入、函数
 - ◆ 局部变量
 - ◆ 继承
 - ◆

```
6 # Copyright 2016 Transwarp Inc. All rights reserved.
7
8 local kube = import "../../applib/kube.libsonnet";
9 local app = import "../../applib/app.libsonnet";
10 local hdfs = import "../hdfs.libsonnet";
11 local hdfs_datanode_common = import "../hdfs-datanode-common.jsonnet";
12
13 local default_config = {
14 };
15
16 function(user_config={})
17   local config = default_config + user_config;
18   kube.tos.Deployment(
19     kube["extensions/v1beta1"].Deployment
20     (
21       generateName="hdfsdatanode",
22       moduleName="hdfsdatanode",
23       config=config
24     ) {
25       spec+: {
26         replicas: config.hdfs_data_replicas,
27         strategy: kube["extensions/v1beta1"].DeploymentStrategy(config.hdfs_update_strategy_configs),
28         template+: hdfs_datanode_common(config) {
29           spec+: {
30             volumes: [
31               kube.v1.PersistentDirVolume("datanodedir", {
32                 datanode: "1",
33                 ["datanode.install." + std.toString(config.Transwarp_Install_ID)]: "true",
34               }),
35               kube.v1.HostShareDirVolume("socketdir", "/root/docker/common"),
36               kube.v1.EmptyVolVolume("transwarp", "100000000"),
37             ],
38             priority: app.parseInt(config.hdfs_data_priority),
39           },
40         },
41       },
42     },
43     config
44   )
45 }
```

HDFS datanode deployment via jsonnet

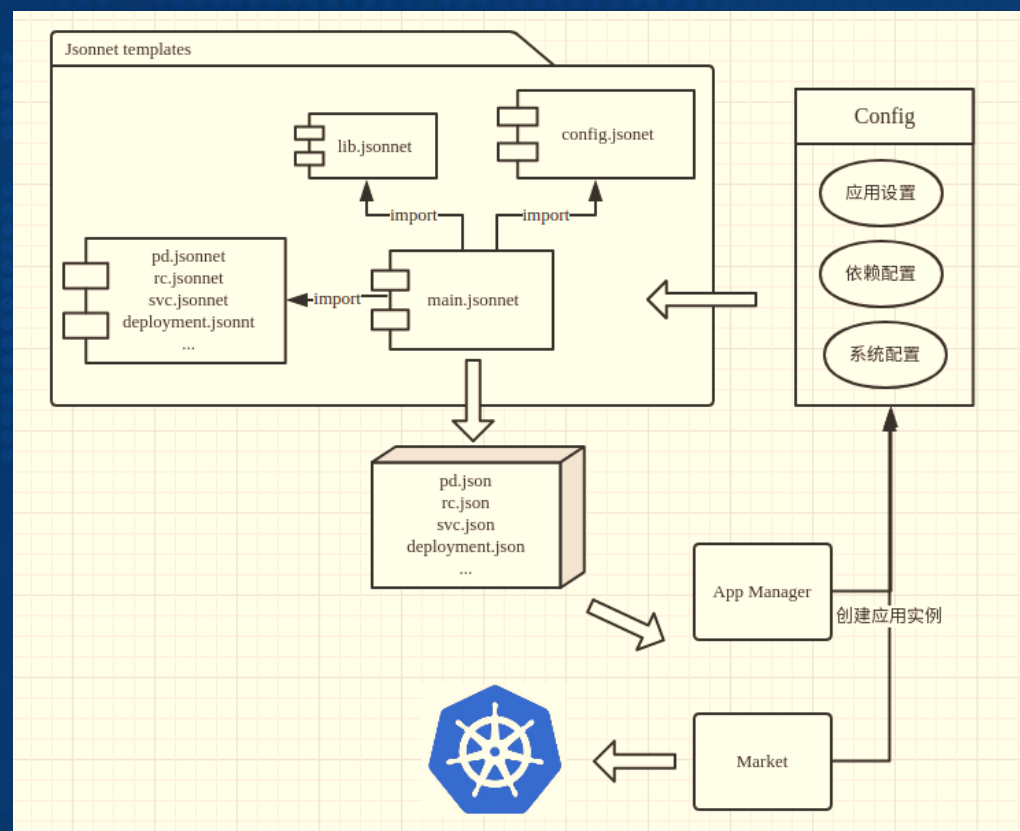
大数据应用编排 2.0 (cond.)

◆ 主要模块

- ◆ app.yaml 描述依赖关系，依赖变量注入
- ◆ config.jsonnet 描述配置参数，定义输入
- ◆ main.jsonnet 模板入口，定义输出

◆ 模板约定

- ◆ 每个应用
 - ◆ 通过transwarp.app=name区分
- ◆ 每个应用实例
 - ◆ 通过transwarp.install=id区分
 - ◆ 不同模块通过transwarp.name区分
 - ◆ 都有一个dummy service实现服务暴露，且transwarp.meta=true
 - ◆ 在annotation中transwarp.meta=“json string”，记录实例需要暴露的信息





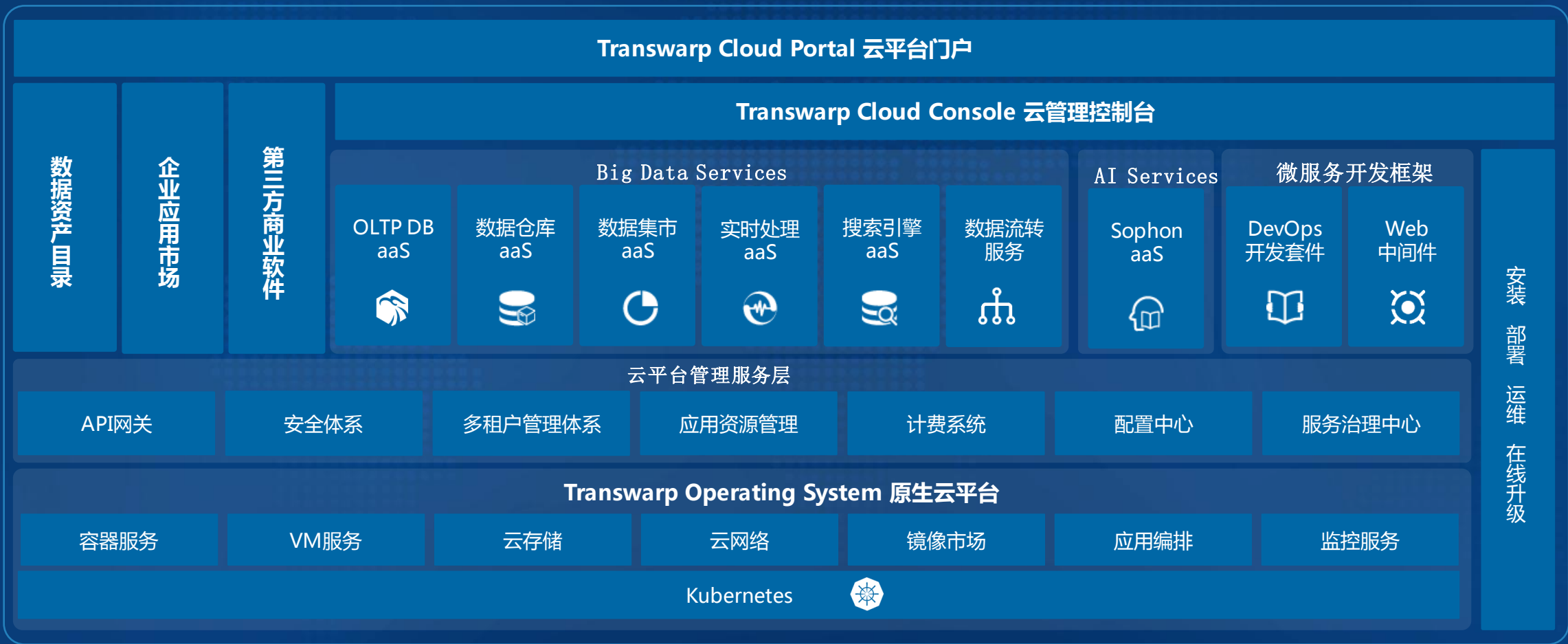
大数据应用编排 2.0 (cond.)

◆面临的挑战

- ◆ 调试麻烦: jsonnet -> svc/pod/deploy
- ◆ 语法复杂, 普通开发人员学习曲线陡
 - ◆ 继承语法比较复杂
 - ◆ 复杂的if-else
- ◆ 应用集中管理, 无法独立分发
- ◆ 社区不兼容

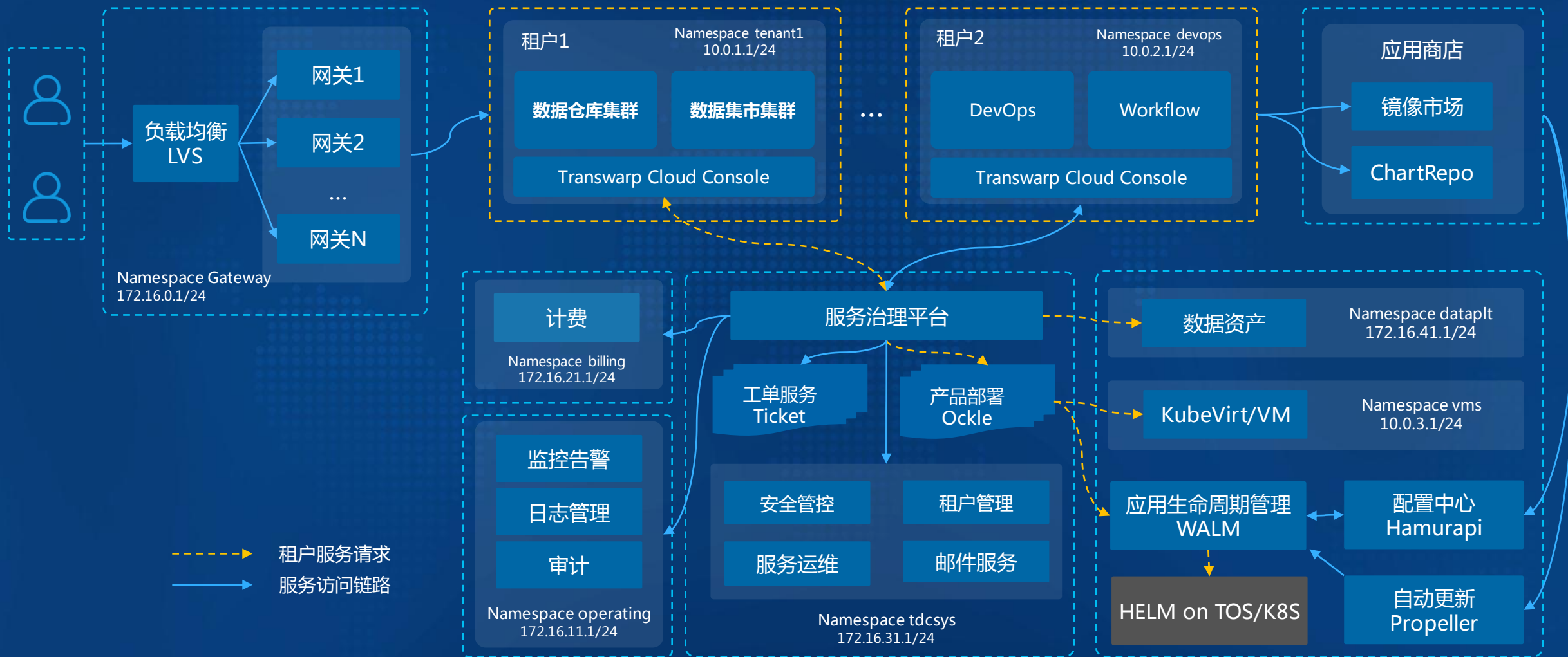


数据云架构 3.0 – Transwarp Data Cloud





TDC拓扑结构





TDC多租户实践





大数据应用编排 3.0

◆设计思路

- ◆ 开发、使用、运维工作内容解耦
- ◆ 完善的DevOps流程
- ◆ 应用包管理
- ◆ 提高社区兼容性
- ◆ 完善的服务动态依赖



大数据应用编排 3.0 (cond.)

◆ 社区版Helm

- ◆ K8S包管理服务，类似于CentOS的yum
- ◆ **Chart**: 应用配置打包格式，类似于rpm
- ◆ **Templates**: k8s manifest配置模板
- ◆ **Release**: chart部署实例
- ◆ **Repository**: chart仓库

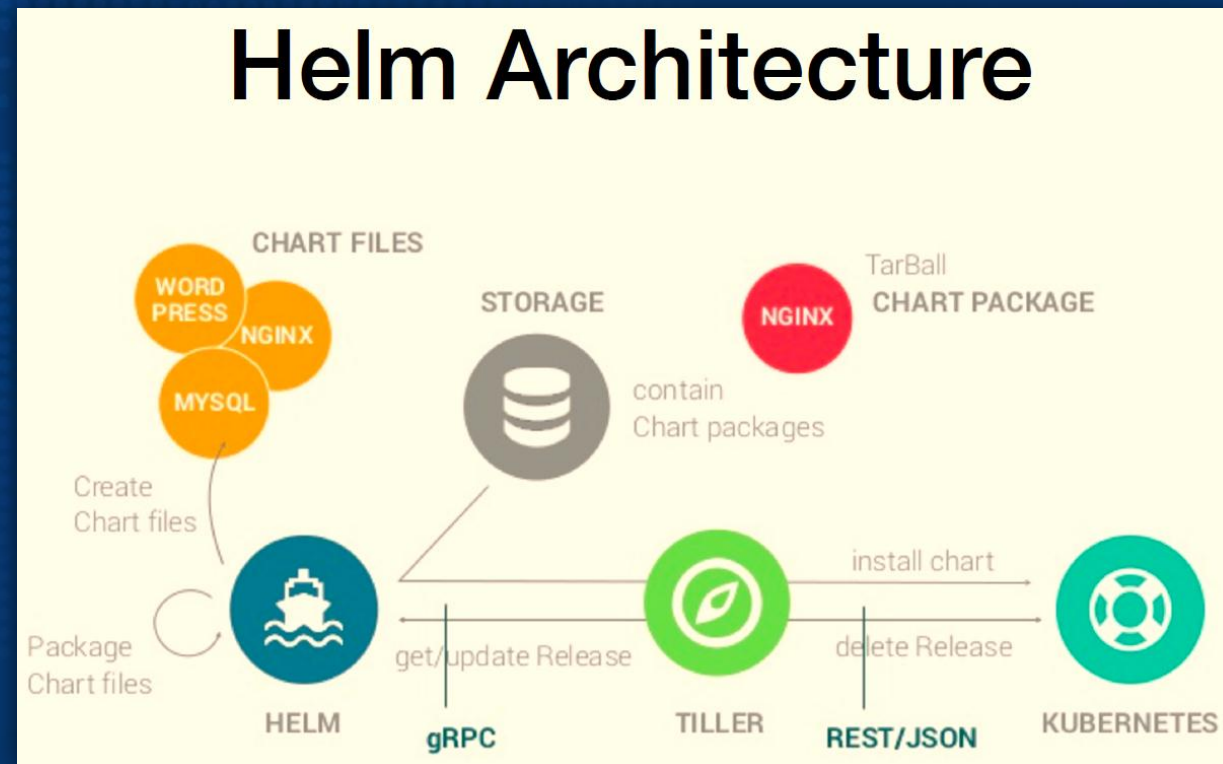


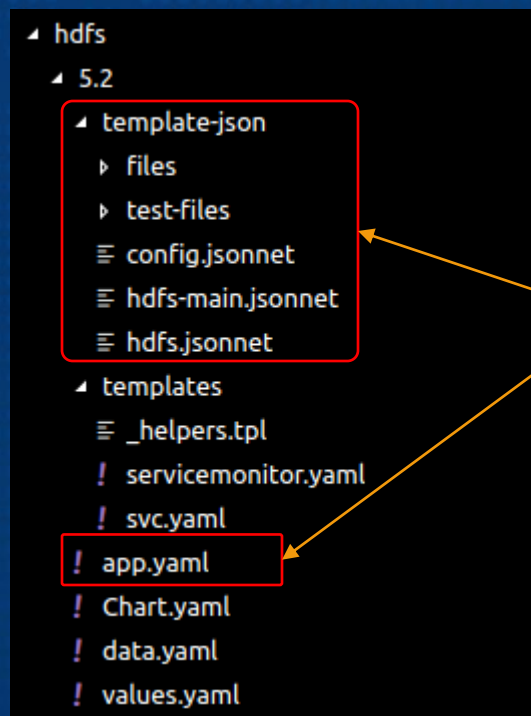
Image source: <https://www.slideshare.net/alexLM/helm-application-deployment-management-for-kubernetes>

大数据应用编排 3.0 (cond.)

◆ TDC Helm+

- ◆ 完全兼容社区Helm Apps
- ◆ 多渲染引擎支持
 - ◆ 社区go-template/Ksonnet
 - ◆ 星环Jsonnet
- ◆ 通过参数控制扩容/缩容
- ◆ 支持服务动态依赖，自定义导出/引用变量
- ◆ 多种依赖模式
 - ◆ 依赖现有应用
 - ◆ 重新创建
 - ◆ 依赖变更
 - ◆ 一对多/多对一的依赖

第三代Chart

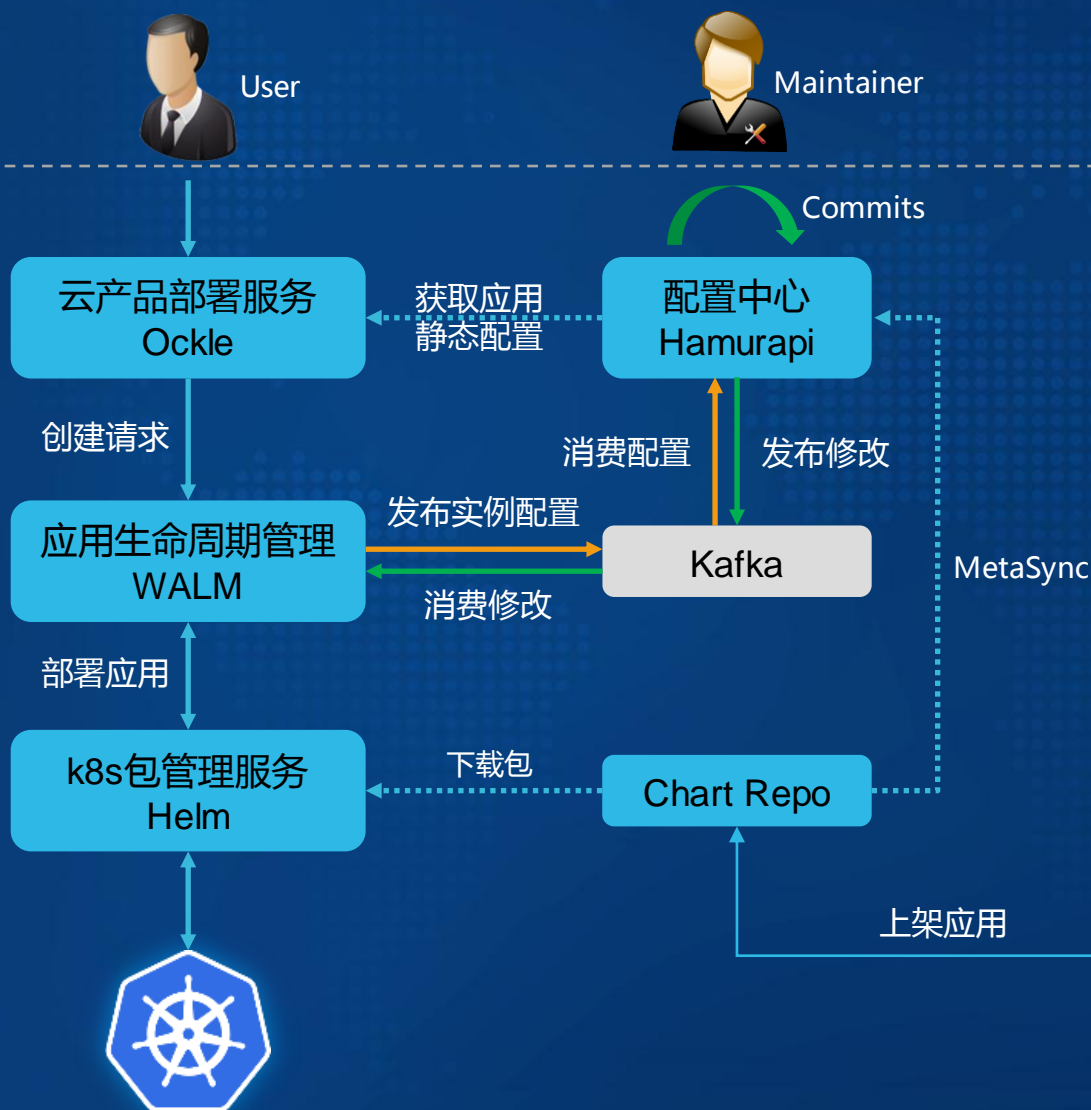


第二代Jsonnet





大数据应用编排 3.0 (cond.)



The screenshot shows the '应用开发平台' (Application Development Platform) interface, specifically the '2 服务信息' (2 Service Information) tab. The interface includes the following fields and controls:

- *微服务名称** (Microservice Name): A text input field with a placeholder '请输入' (Please enter) and a note '必须为英文数字组合, 长度小于10' (Must be a combination of English letters and numbers, length less than 10).
- *副本数** (Replica Count): A text input field with the value '3'.
- *容器名称** (Container Name): A text input field with the value '云和镜像处理' (Cloud and Image Processing). There are buttons for '清除内容' (Clear Content) and '删除容器' (Delete Container).
- *镜像地址** (Image Address): A text input field with a placeholder '请输入镜像地址' (Please enter image address). There is a checkbox for '登录信息' (Login Information).
- *端口** (Port): A text input field with the value '8080'. There are checkboxes for '服务入口' (Service Entry) and '默认路径' (Default Path).
- 服务名称** (Service Name): A text input field with the value 'tomcat'.
- 路径** (Path): A text input field with the value '/'. There is a trash icon next to it.
- + 添加端口** (Add Port): A button to add more ports.
- *启动方式** (Startup Method): Radio buttons for '启动命令' (Startup Command) and '启动脚本' (Startup Script).
- 请输入英文字符串** (Please enter English string): A text input field.
- 资源配置** (Resource Configuration): A section with a plus icon.
- 高级选项** (Advanced Options): A section with a plus icon.
- 健康检查** (Health Check): A section with a plus icon.



We Are Hiring

感谢您的聆听

www.transwarp.io

