

开放搜索服务：从系统、平台 到服务架构演变

阿里巴巴搜索事业部

郭瑞杰 (花名：国泊)

2015年4月24日

促进软件开发领域知识与创新的传播



ArchSummit
全 球 架 构 师 峰 会

【深圳】2015年7月17日-18日

QCon
全 球 软 件 开 发 大 会

【上海】2015年10月15-17日



关注InfoQ官方微信
及时获取QCon演讲视频信息



阿里搜索系统架构演变



系统、平台、服务

搜索引擎基础数据结构

Inverted Index

Attribute
(FieldData)

Summary
(Stored Field)

Doc	Term
doc0	a, b, c, a
doc1	b, c, d
doc2	a, c, d
doc3	c, e



Term	
a	0<0,3>, 2<0>
b	0<1>, 1<0>
c	0<2>, 1<1>, 2<1>, 3<0>
d	1<2>, 2<2>
e	3<1>

Per Doc:

Field	attribute
field1	1, 2, 4, 9
field2	abc, cde, a
field3	1.02
field4	true

Per Doc:

Field	attribute
field1	abc cde a
field2	aa bb cc
field3	a b c
field4	...

使用Inverted Index加速查询

使用Attribute过滤、分组统计、聚合

使用Summary做展示和摘要飘红

淘宝网 Taobao.com

更多市场 宝贝 羽绒服 搜索

所有宝贝 天猫 二手 我的搜索

所有分类 共 981.26万 件宝贝

品牌: 波司登 yaloo/雅鹿 艾莱依 韩国SZ La Chapelle/拉夏... 优衣库 多选 更多

欧时力 海澜之家 Afs Jeep/战地吉普 Goldfarm/高梵 金羽杰 秋水伊人

衣长: 中长款 短款 常规 长款 超短 多选 更多

服装款... : 拉链 口袋 带毛领 拼接 纽扣 蝴蝶结 系带 立体装饰 多选 更多

女表: 羽绒服 棉衣/棉服 大码女装 皮衣 中老年服装 毛呢外套 短外套 更多

筛选条件: 流行男装 选购热点 尺码 服装版型 衣门襟 上市年份季节

您是不是想找: 短款羽绒服 羽绒服男 欧时力羽绒服 孕妇羽绒服 女士羽绒服 羽绒服女中长款 棉衣 羽绒衣

综合排序 人气 销量 信用 价格 ￥ - 发货地 1/100

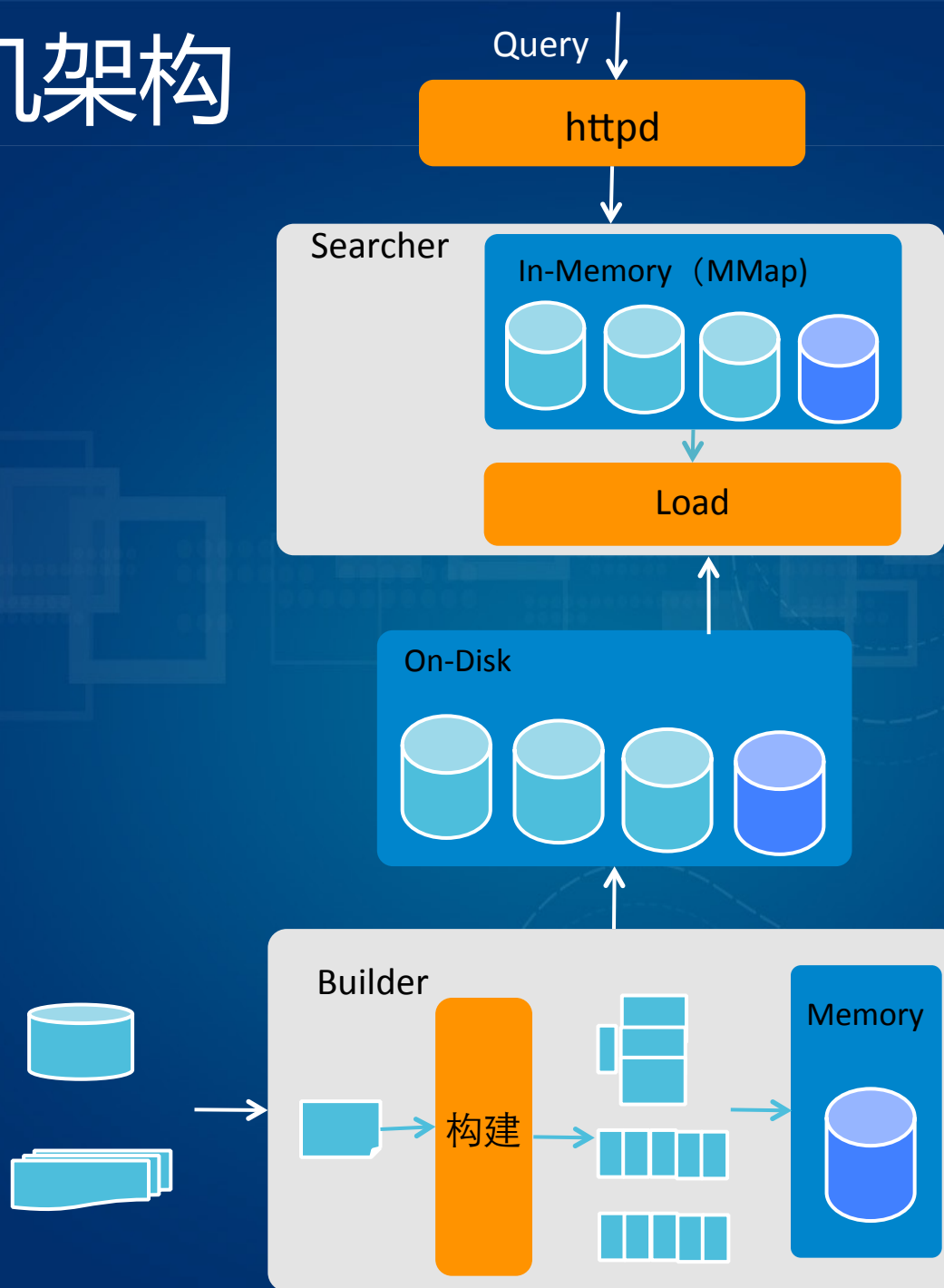
☐ 包邮 ☐ 赠送退货运费险 ☐ 货到付款 ☐ 新品 ☐ 海外商品 ☐ 二手 ☐ 天猫 更多 合并卖家

高梵 GOLD FARM 高梵2014冬加厚中长款羽绒服 2691人付款 ¥569.00 包邮 高梵旗舰店 安徽 合肥

茵曼 茵曼女神的新衣NANA同款 2394人付款 ¥659.00 包邮 茵曼旗舰店 广东 广州

千仞岗 千仞岗正品女装新款羽绒服 1533人付款 ¥649.00 包邮 千仞岗服饰旗舰店

单机架构



全量索引构建

- 对所有文档重新构建索引，并替换旧索引数据

增量索引构建

- 对新增文档批量构建增量索引，合并至索引库中

实时索引更新

- Builder通过定时flush内存索引至磁盘，Searcher从磁盘Load新索引，实现分钟级或小时级更新时效性
- 性能较差

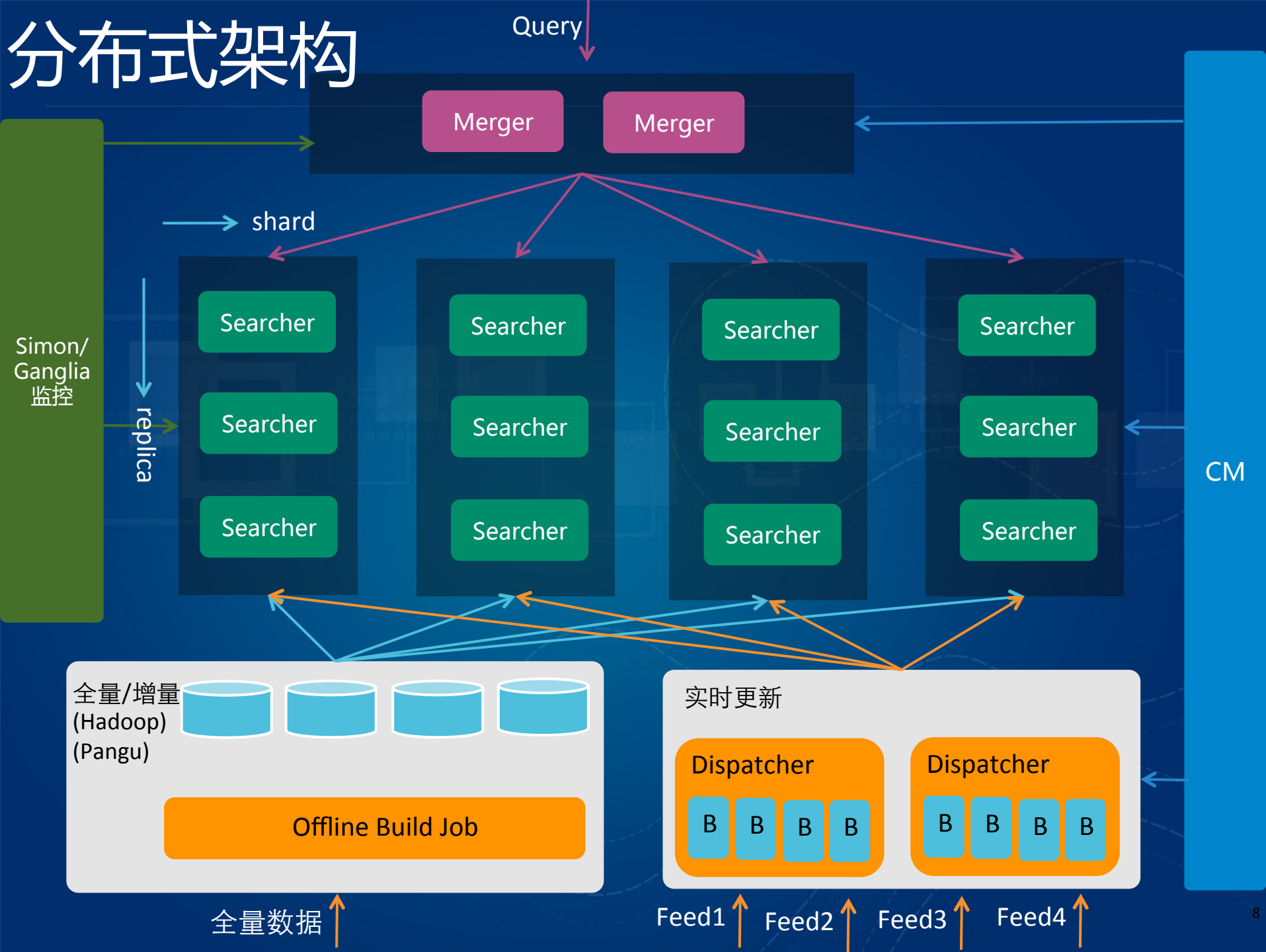
数据



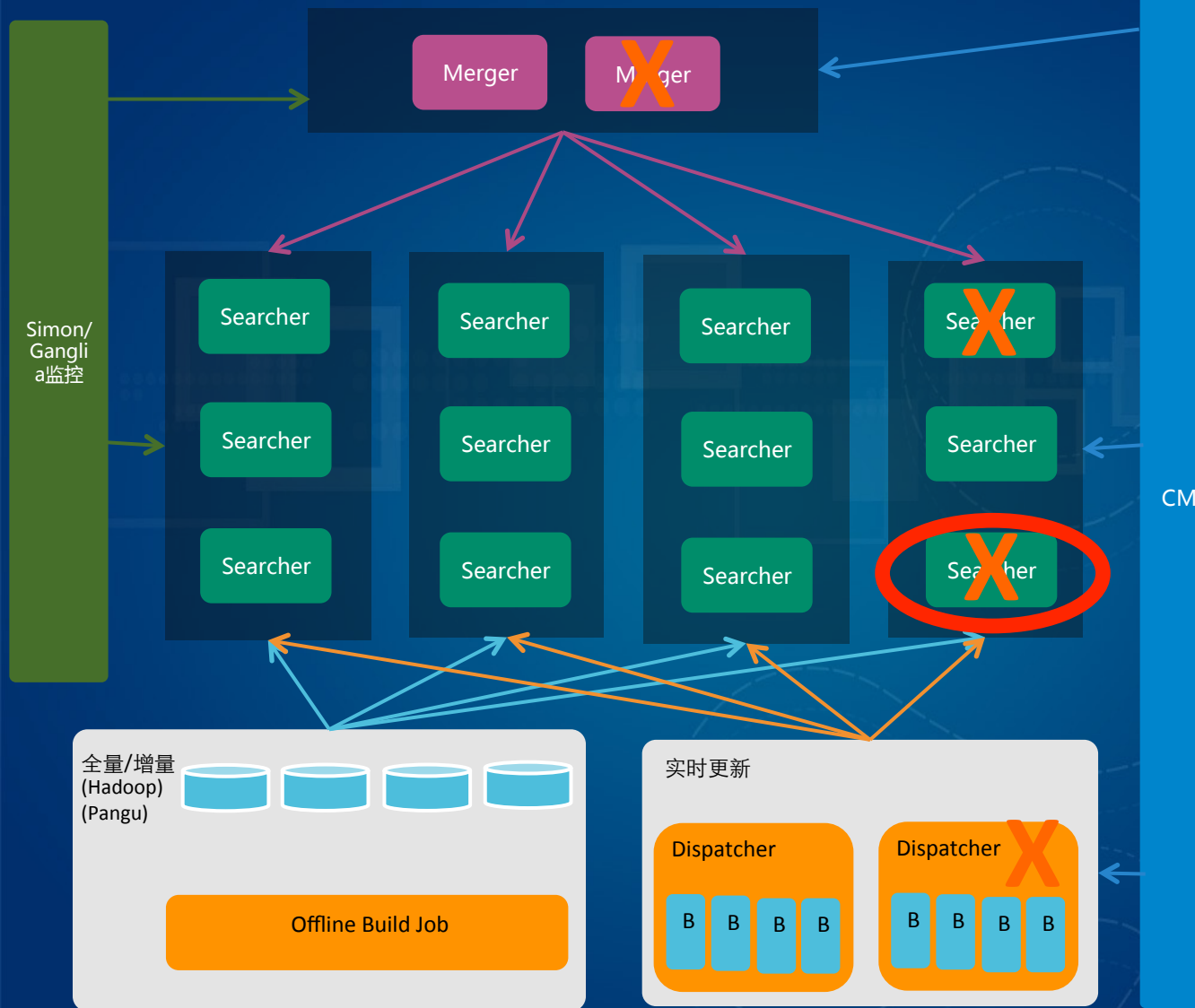
PV



分布式架构



错误恢复



Merger进程挂掉

- CM发现后屏蔽
- 进程重启，向CM注册，恢复服务

Searcher进程挂掉

- CM发现后屏蔽
- 进程重启，向CM注册
- Dispatcher获取断点信息，续传实时更新的文档

Dispatcher进程挂掉

- CM发现后屏蔽
- 进程重启，向CM注册，恢复服务

机器挂掉

- CM发现后自动屏蔽
- 每天晚上全量索引完成后换机器
- 不支持动态切换机器

索引损坏

- 同上

机器

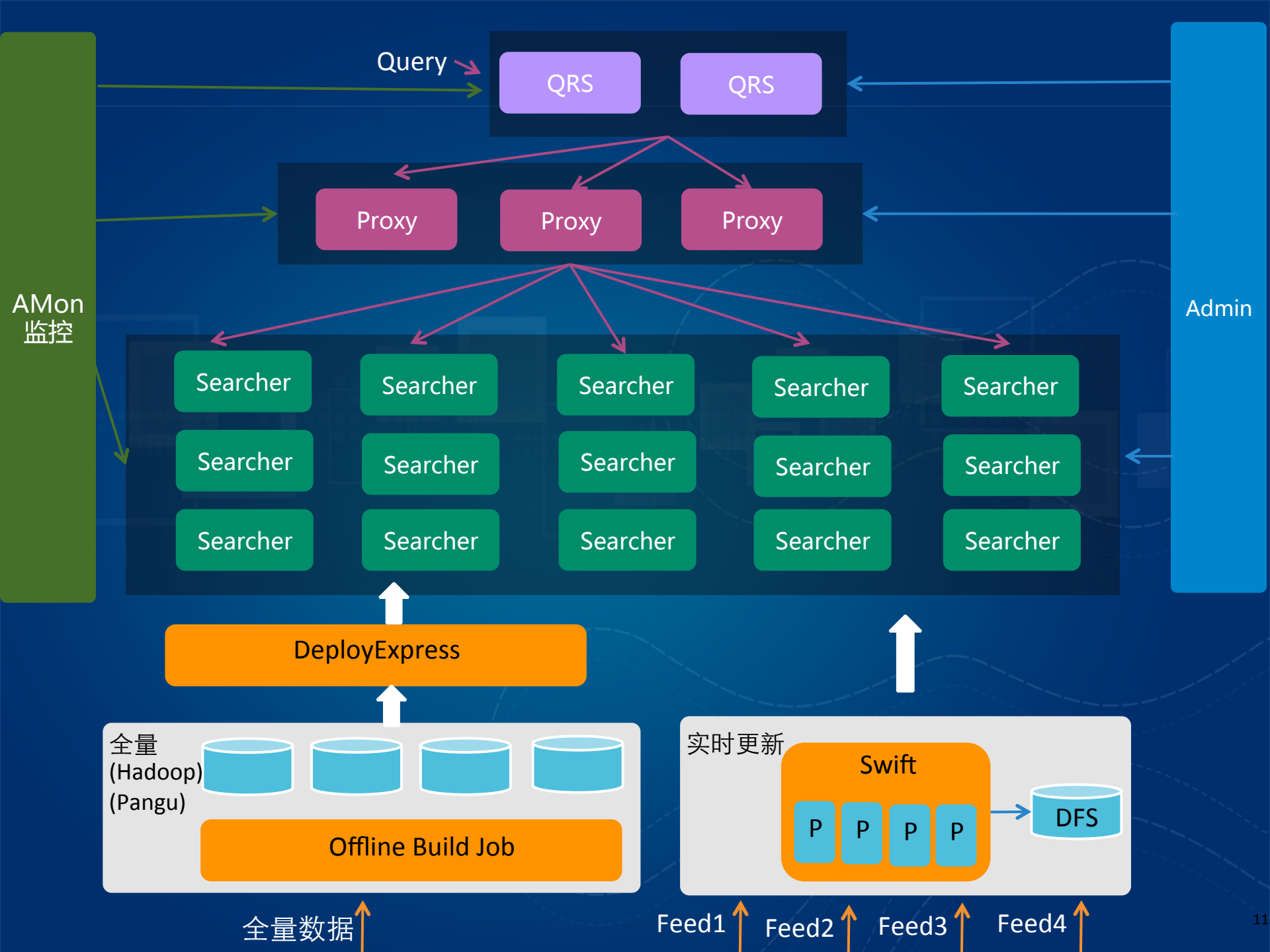
50



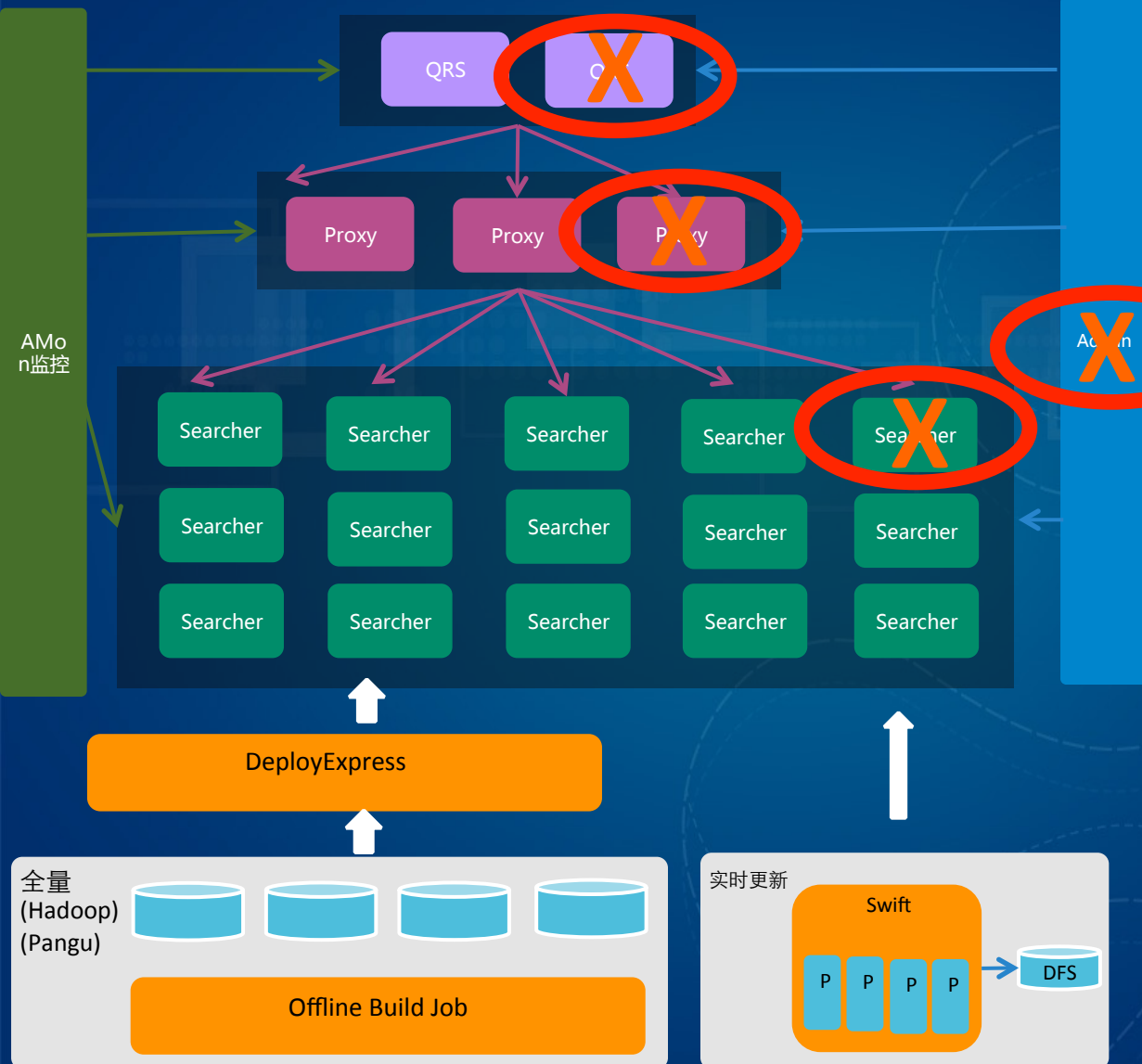
500+

PE





错误恢复



QRS/Proxy机器挂掉

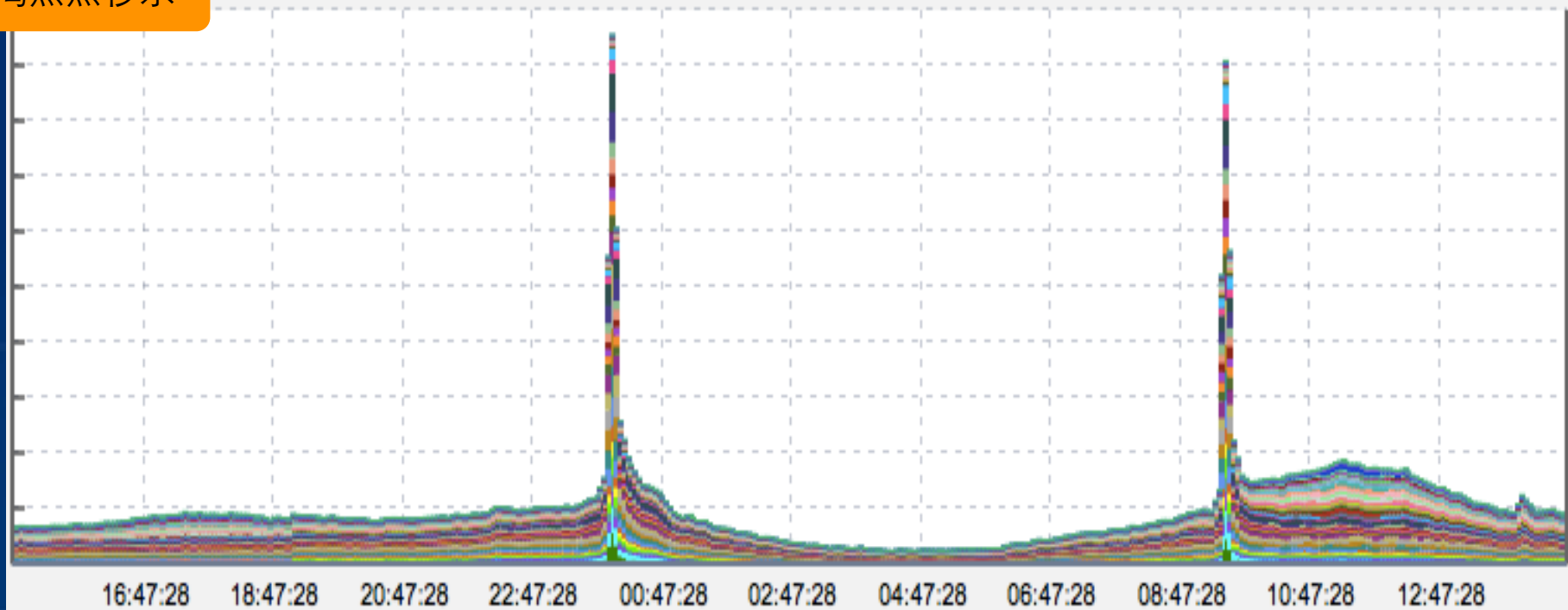
- Admin从Buffer机器列表中选择1台机器
- 启动进程，替代挂掉机器提供服务

Searcher机器挂掉

- Admin从Buffer机器列表中选择1台机器
- 启动进程
- Admin发送命令给DeployExpress，拉取增量索引数据
- 从Swift拉取实时更新数据，数据完整后提供服务

Admin机器挂掉

- 启动备份Admin
- 恢复配置，接管服务

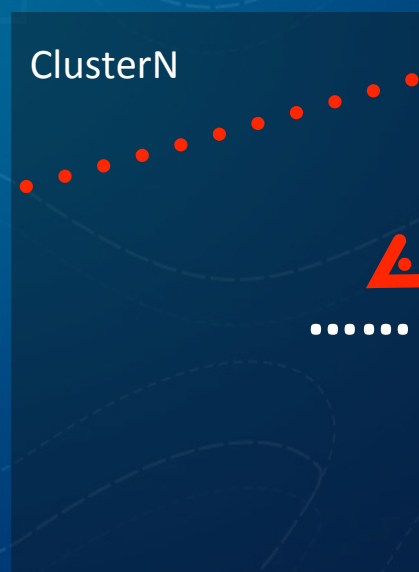
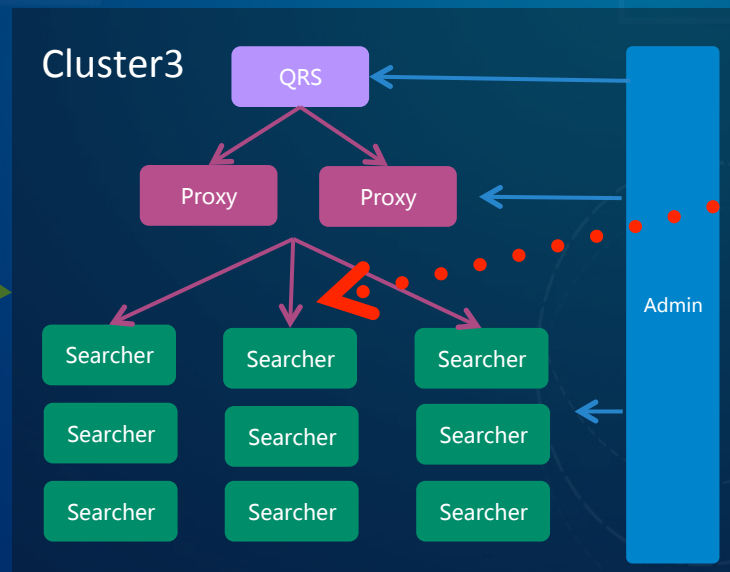
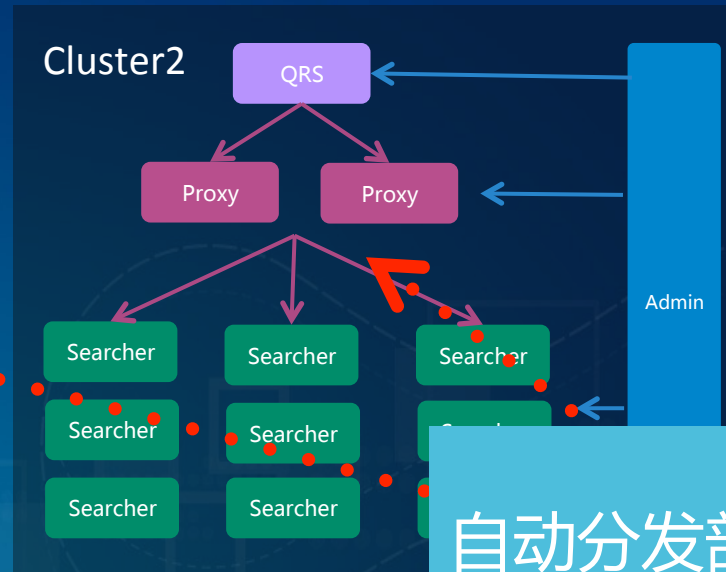
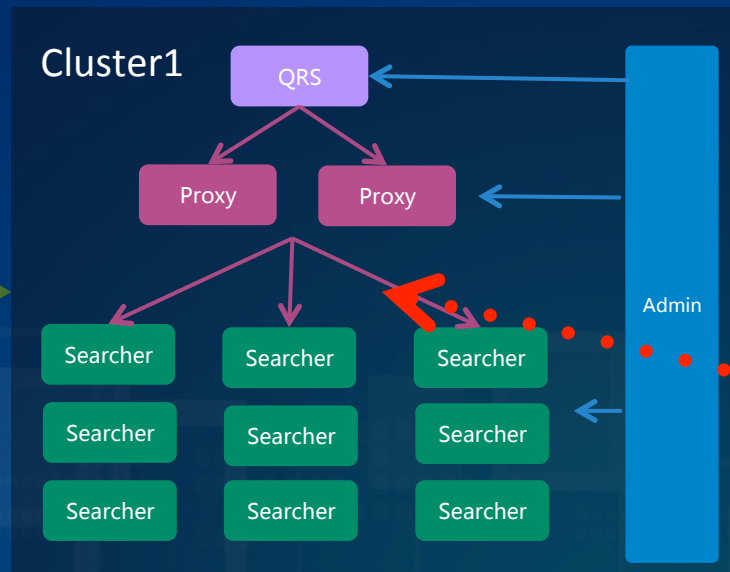


资源问题

-流量波峰波谷、多业务间、在线离线资源无法复用

运维问题

-业务多、机器多、峰值突增，扩容和机器迁移耗时耗力



自动分发部署

- 依赖数据、Binary分发部署

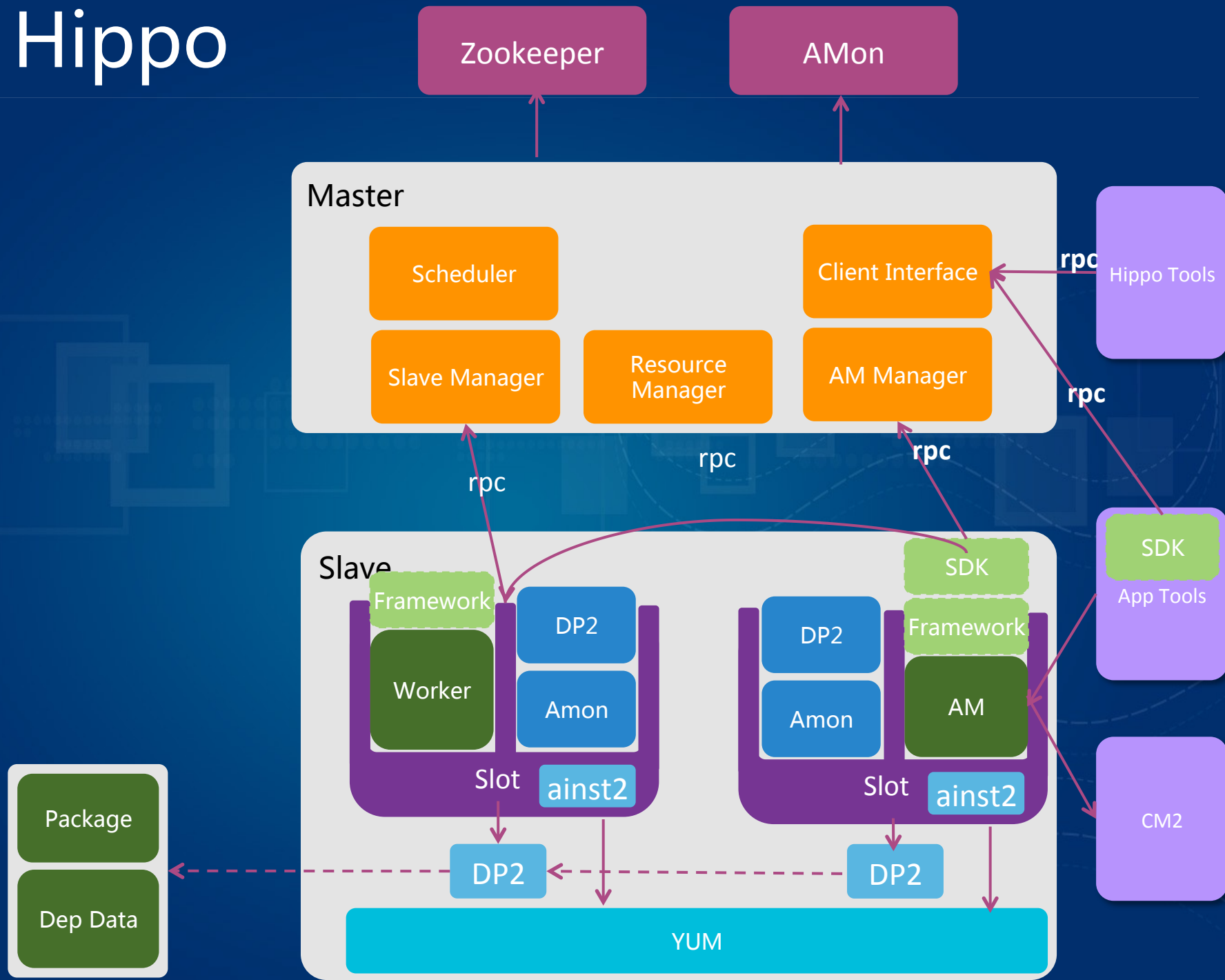
集群间资源复用

- 共用Buffer机、机器分配、资源迁移

提升资源利用率

- 低谷复用(跑测试算法任务)、
不同服务间复用、在线离线
复用

Hippo



两层调度实例

Application

app_name: test

user: test

packages: { ha3, java }

resource : { cpu : xx, mem: xx }

process: hippoadmin_worker

parameters: -p xx -c xx

envs: xxxx

....

1. Submit

Hipoo Master

qrs resource X 2
searcher resource X 4

2. Select

3. Start

4. assign

127.0.0.[2-7]

QRS

qrs_0

Hippo Slave

QRS

qrs_1

Hippo Slave

5. Start

Searcher
search_0

Hippo Slave

Searcher
search_1

Hippo Slave

Searcher
search_2

Hippo Slave

Searcher
search_3

Hippo Slave

Admin

QRS

qrs_0

QRS

qrs_1

Searcher

search_0

Searcher

search_1

Searcher

search_2

Searcher

search_3

Hippo Slave

错误恢复

QRS/Proxy机器挂掉

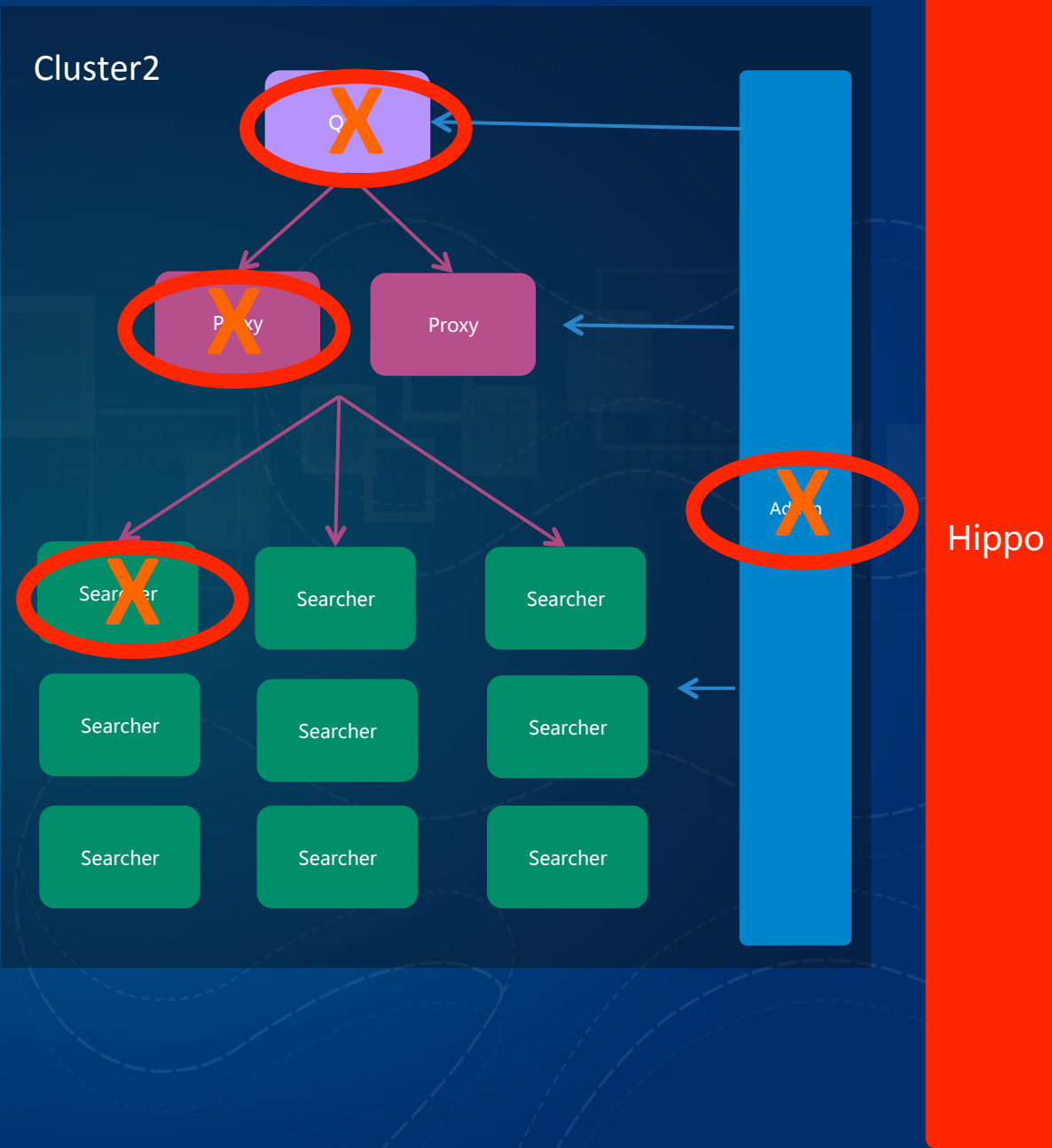
- Admin向Hippo申请机器资源
- Hippo分发Binary包和数据，启动进程，替代挂掉机器提供服务

Searcher机器挂掉

- Admin向Hippo申请机器资源
- Hippo分发Binary包和数据，启动进程
- Admin发送命令给DeployExpress，拉取全量索引数据
- 从Swift拉取实时更新数据，数据完整后提供服务

Admin机器挂掉

- 启动备份Admin
- 恢复配置，接管服务





系统、平台、服务

搜索业务和系统



一套系统、多个代码分支、多套部署

Search For
神马搜索

Search For
淘宝

Search For
B2B

Search For
天猫

Search For
聚划算

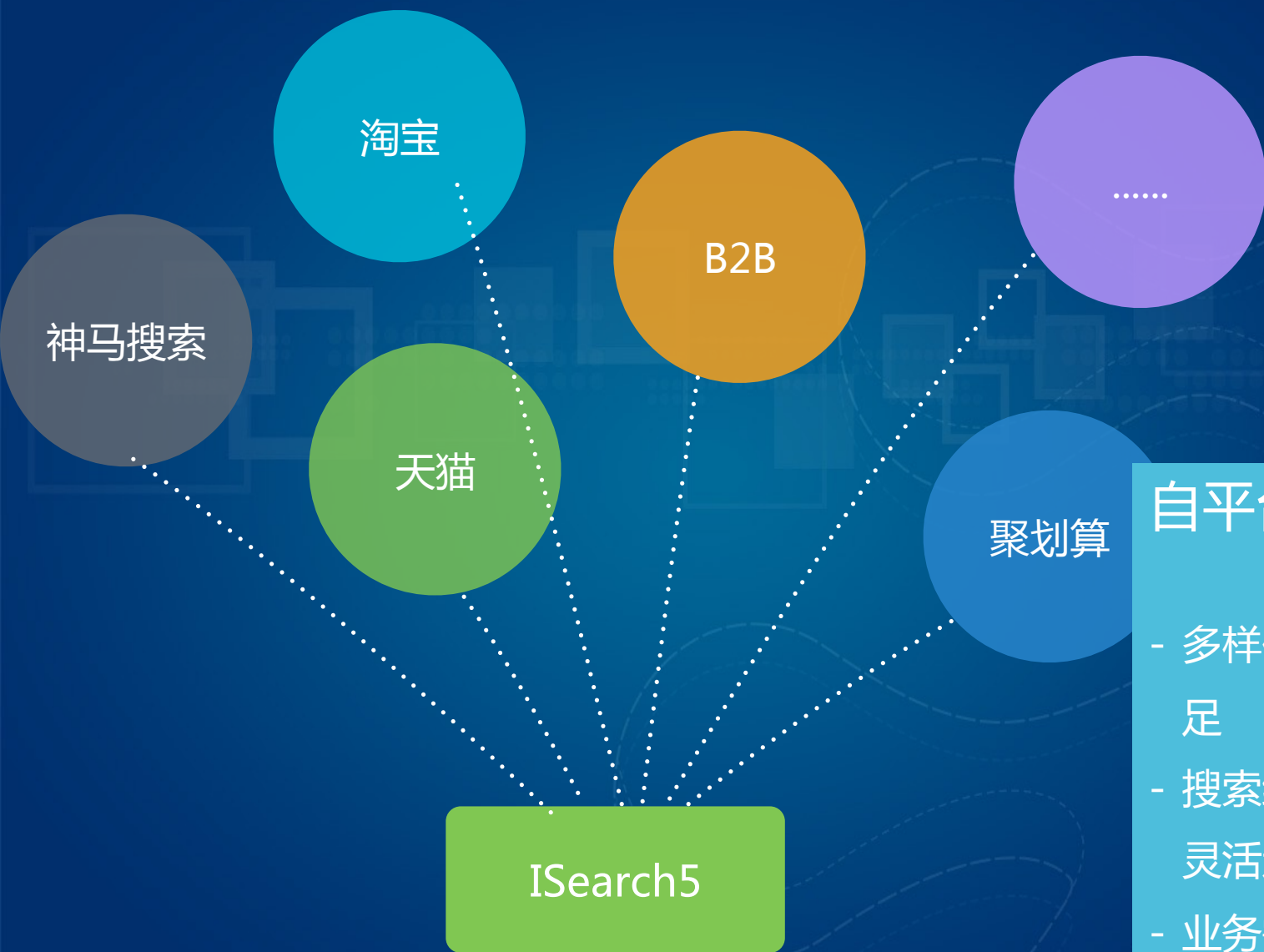
.....

代码不能复用维护困难

系统提升无法共享

随着业务数量增长效率急速下降

解决方案：平台化



自平台化的挑战

- 多样化功能需求如何满足
- 搜索结果排序规则如何灵活定制
- 业务需求如何快速响应

系统插件化

算分插件

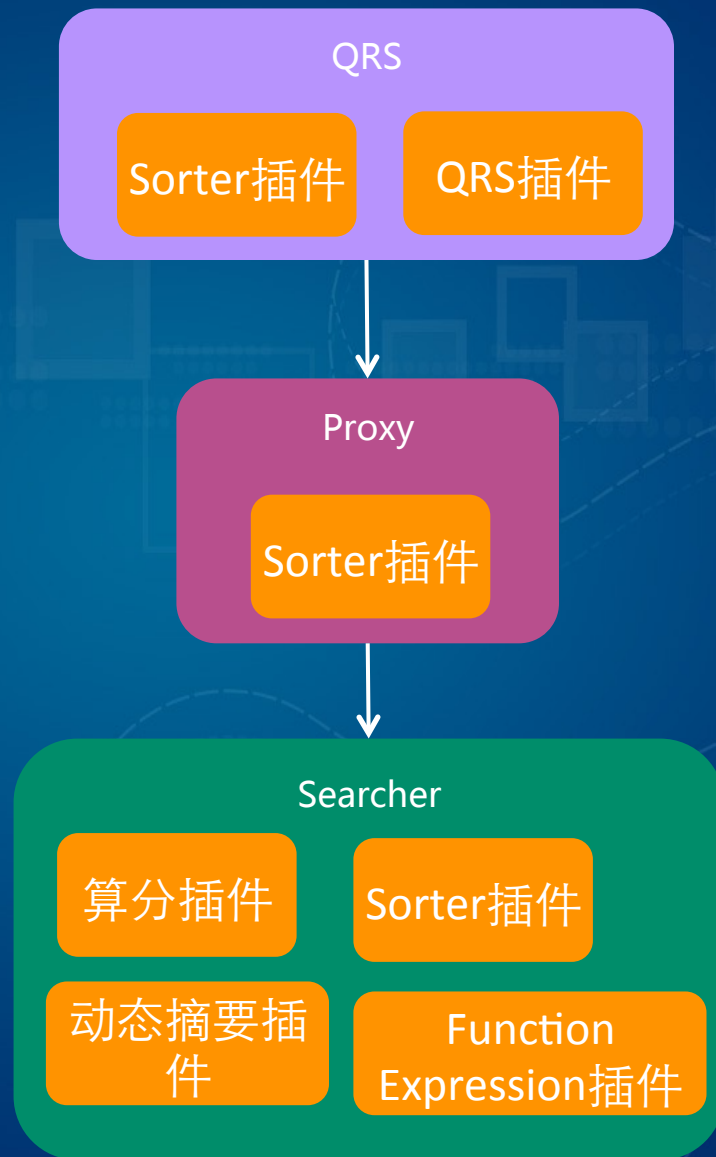
QRS插件

文档处理插件

动态摘要插件

Function
Expression插件

Sorter插件



团队组织架构保障

业务



业务平台团队



算法团队



引擎平台团队

业务数量

100+



1000+



服务化

不让用户写插件

屏蔽系统细节

自助式定制

OpenSearch：自助式云搜索服务

UI & API

- 所有功能可以在控制台中操作完成
- 所有主要功能可通过API操作实现
- 按需使用服务，按需调整配额

索引结构定制

- 文档字段、字段类型、索引方式可定制
- 运行中的服务实例，索引结构可动态修改

相关性定制

- 通过表达式定制搜索结果排序逻辑（后续将支持Lua）
- 内置文本、地理位置相关等算分特征函数
- 使用查询分析服务更好理解用户query

UI

Console

应用 模板 数据源 相关性 数据统计 调试

API

Search Push Config SolrAPI ESAPI

AngularJS

Tengine/PHP

在线计算

Hippo
集群部署

Suggest
下拉提示

QP
查询分析

Aggregator

Free Schema
多租户数据模型

相关性排序

简单脚本 战马 Feature Lua

QuotaServer
配额管理

Swift
分布式消息队列

HA3
搜索引擎平台

Amon
集群监控

离线计算

统计

Join
主附表Join

Processor
数据处理插件

Adapter

RDS ODPS OSS MySQL 云梯1 DRDS

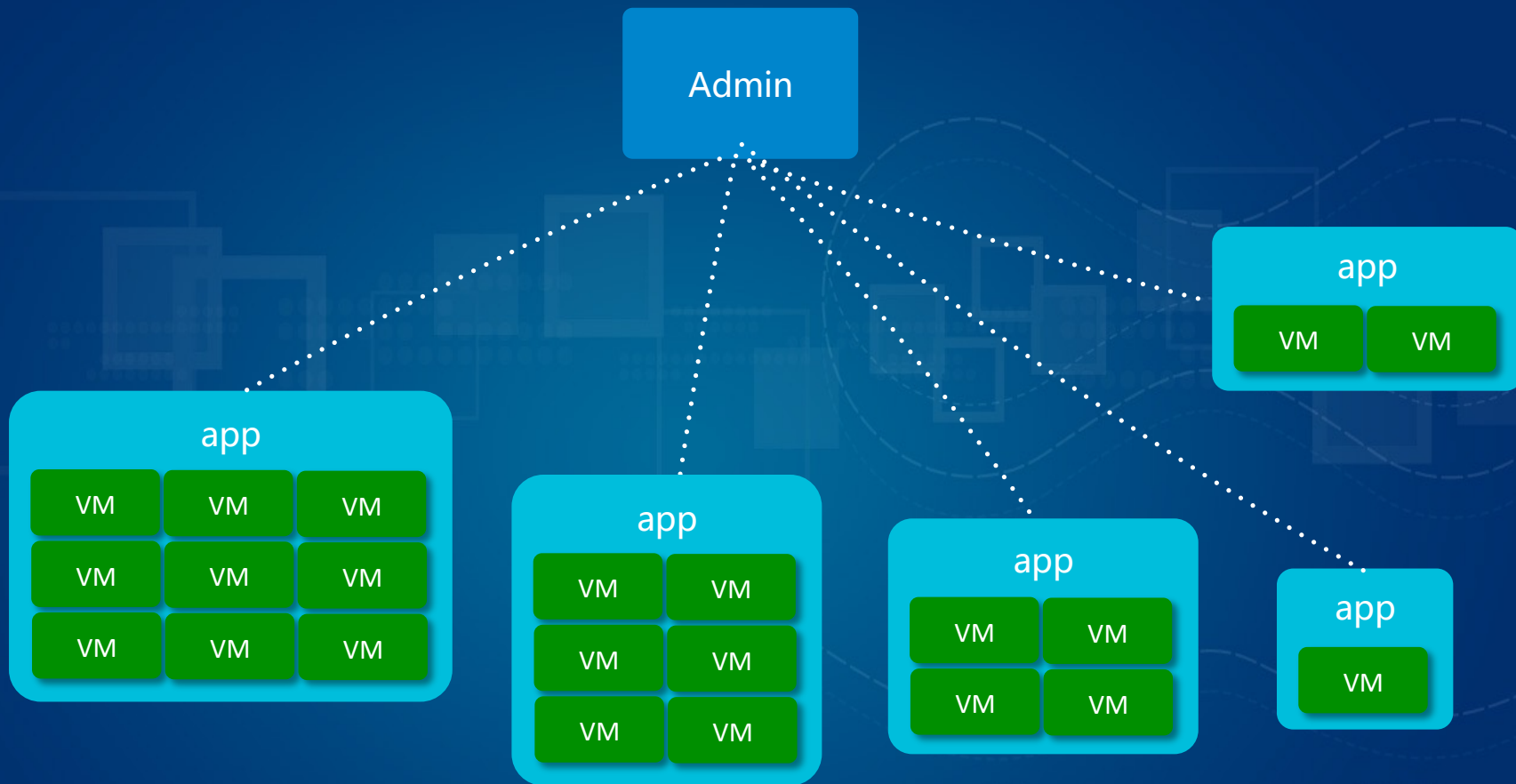
MapReduce

HQueue

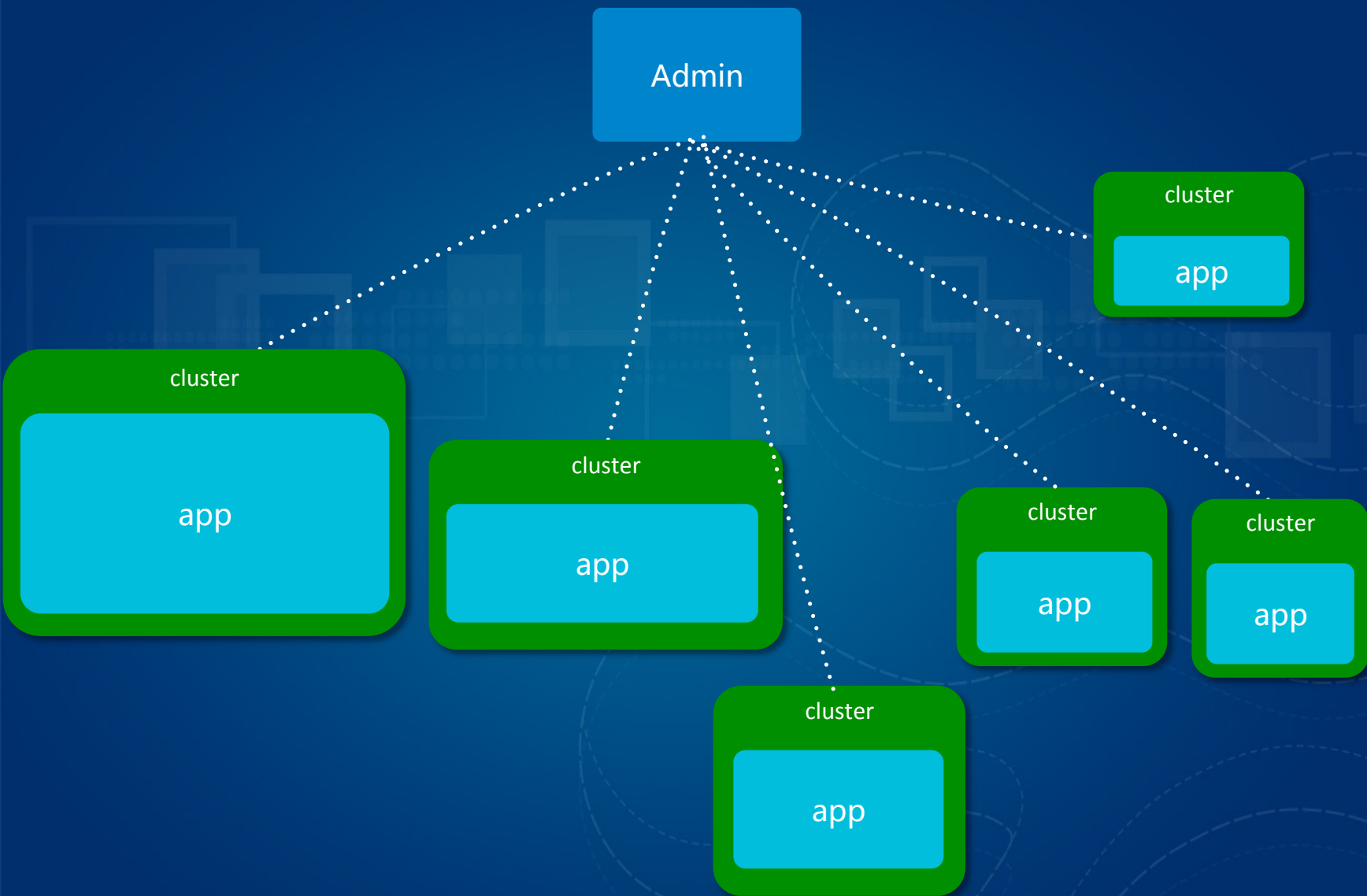
iStream

HBase/Hadoop

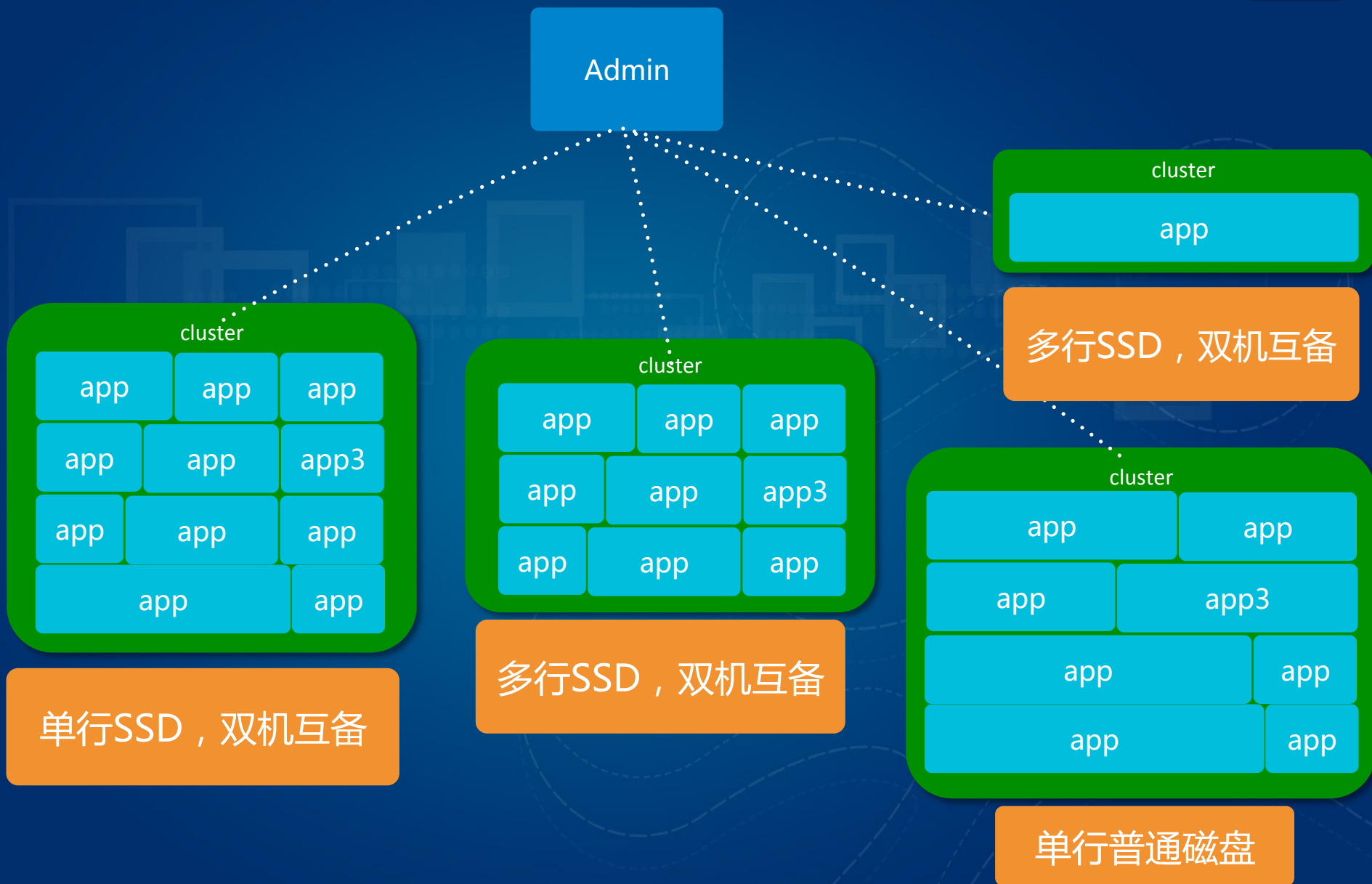
多租户数据模型：虚拟机



多租户数据模型：海量cluster



多租户数据模型：混合



按需使用，自动迁移

流量高峰

流量低谷

Admin

Cluster(ssd)

app

多行，双机互备

Cluster(sata)

app

app

app

app

app

app

app

app

多行，双机互备

单行

cluster (ssd)

app

app

app

app

app

app

app

app

app

app

app

单行，双机互备

Cluster (ssd)

app

app

app

app

app

app

app

app

app

多行，双机互备



开放搜索服务OpenSearch

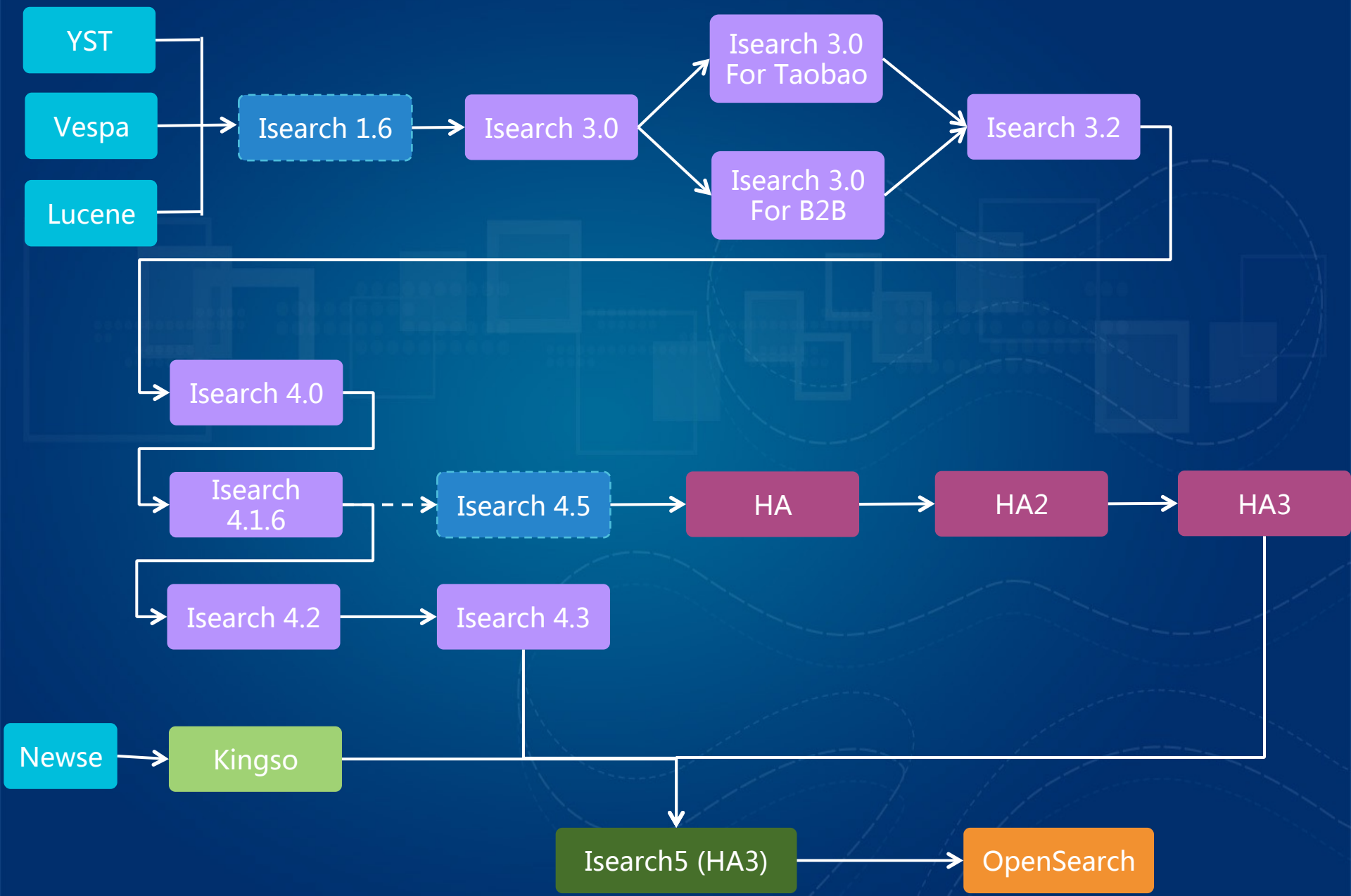
www.aliyun.com/product/opensearch

微博@阿里云开放搜索服务

旺旺群@1318169830

QQ群@370015616

阿里搜索引擎平台演变历史



开源: 2015



Thanks

微博@ruijieguo

微信@jared_guo

InfoQ^{ueue}

专注中高端技术人员的
社区媒体



EGO^{ueue} EXTRA GEEKS' ORGANIZATION
NETWORKS

高端技术人员
学习型社交网络



StuQ^{ueue}

实践驱动的
IT职业学习和服务平台



极客邦科技

InfoQ | EGO | StuQ

让技术人学习和交流更简单