

构建多公有云系统部署实践

上海曙安数据服务有限公司

叶向宇

QCon

2016.10.20~22

上海·宝华万豪酒店

全球软件开发大会 2016

[上海站]



购票热线: 010-64738142

会务咨询: qcon@cn.infoq.com

赞助咨询: sponsor@cn.infoq.com

议题提交: speakers@cn.infoq.com

在线咨询 (QQ): 1173834688

团 · 购 · 享 · 受 · 更 · 多 · 优 · 惠

7折

优惠 (截至06月21日)
现在报名, 立省2040元/张

关于我

- 过去：
 - 联想集团服务器事业部
 - 微软中国技术中心
 - 惠普中国技术中心
 - 微软云计算事业部
- 现在：
 - 上海曙安数据服务有限公司
 - VC3多云管理平台架构师

今天的话题

- 我们的目标是什么？
- 实现业务目标过程中遇到了什么问题？
- 我们是如何思考的？
- 我们是如何实践的？
- 我们下一步的计划是什么？

- 单数据中心，VMWare环境
- 宕机4小时

1

- 切换云供应商
- 再次宕机

3

2

- 单云供应商
- 宕机6小时

4

- 再找一家云供应商???

我们如何走到这一步？

从宕机中学到的几件事（1）

SLA < 99.95% → \$\$

- 供应商SLA不是保证不宕机，而是索赔的依据

从宕机中学到的几件事（2）



- 小范围宕机几乎不可避免

从宕机中学到的几件事（3）

Region	Status	30 Day Availability	1 block = 1 mins	Outages	▼ Downtime
ams2	↑	86.7087%		<u>1</u>	95.7 hours
de-germany	↑	93.8392%		<u>2</u>	44.36 hours
us-virginia	↑	95.6044%		<u>1</u>	31.65 hours
perth	↑	98.5208%		<u>4</u>	10.65 hours
paris	↑	99.7083%		<u>1</u>	2.1 hours
frankfurt	↑	99.8611%		<u>2</u>	60 mins
us-nevada	↓	99.877%		<u>2</u>	53.15 mins
ZRH	↑	99.882%		<u>1</u>	50.98 mins
UC1	↑	99.8954%		<u>2</u>	45.2 mins
sweden-south	↑	99.9126%		<u>4</u>	37.75 mins

<https://cloudharmony.com/status-of-compute>

- 大范围宕机发生可能性依然存在

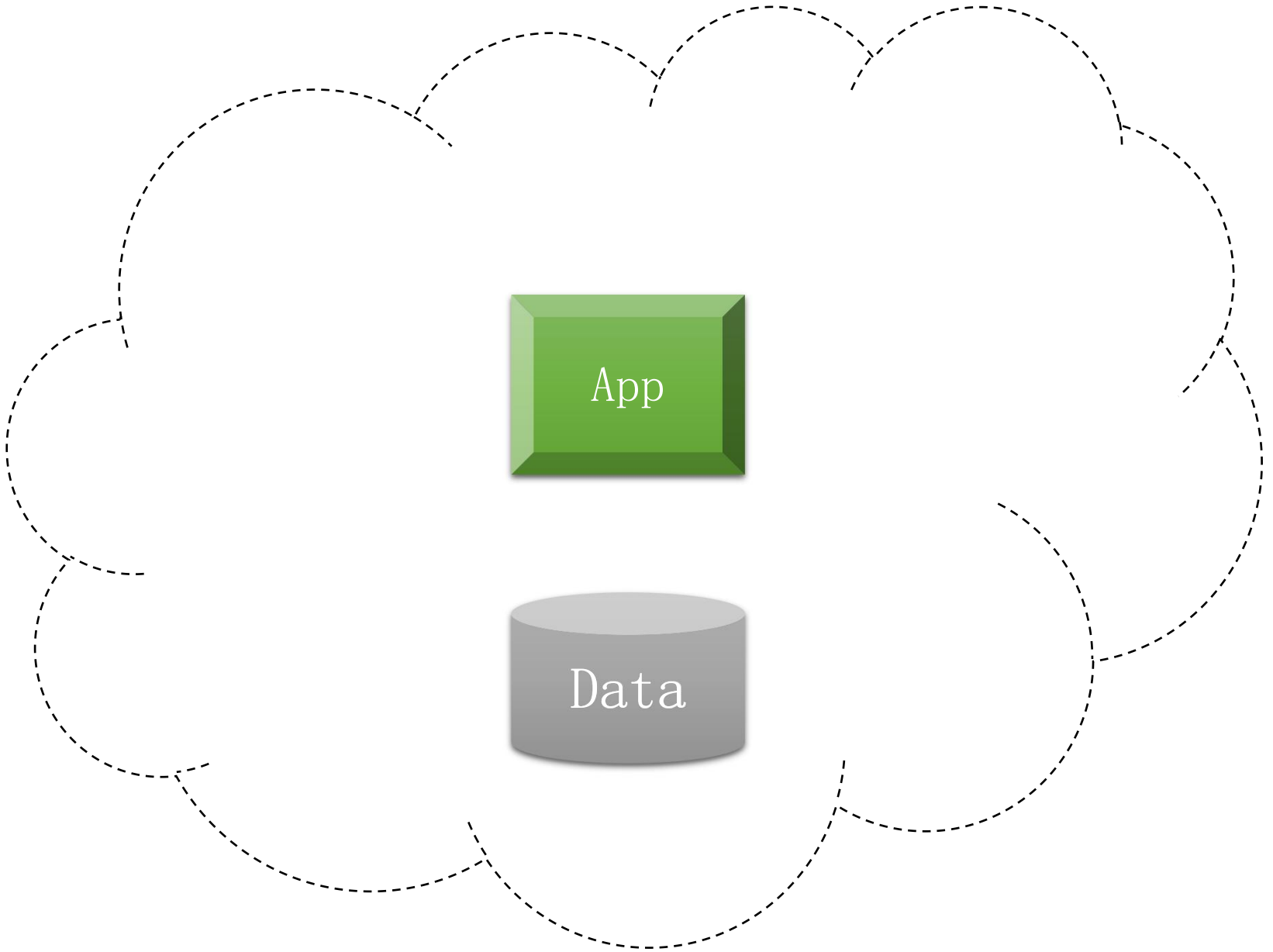
“不宕机是核心需求”

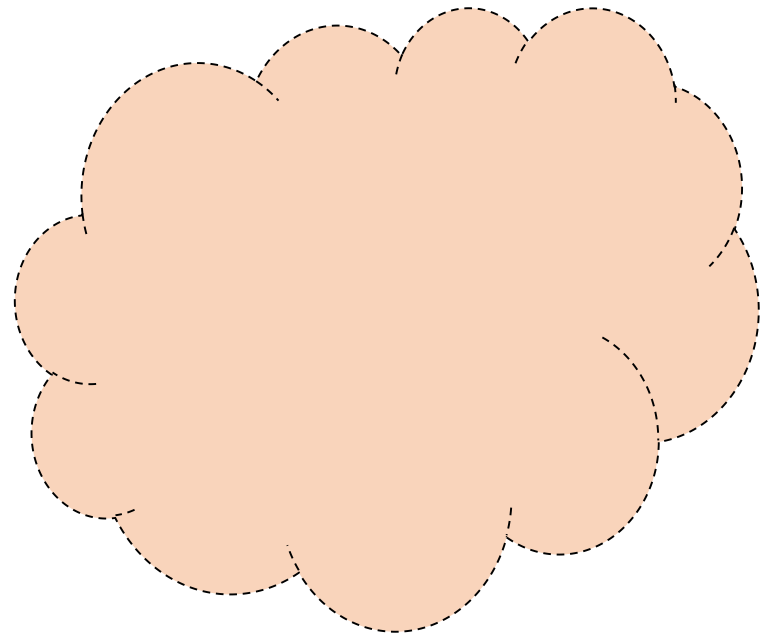
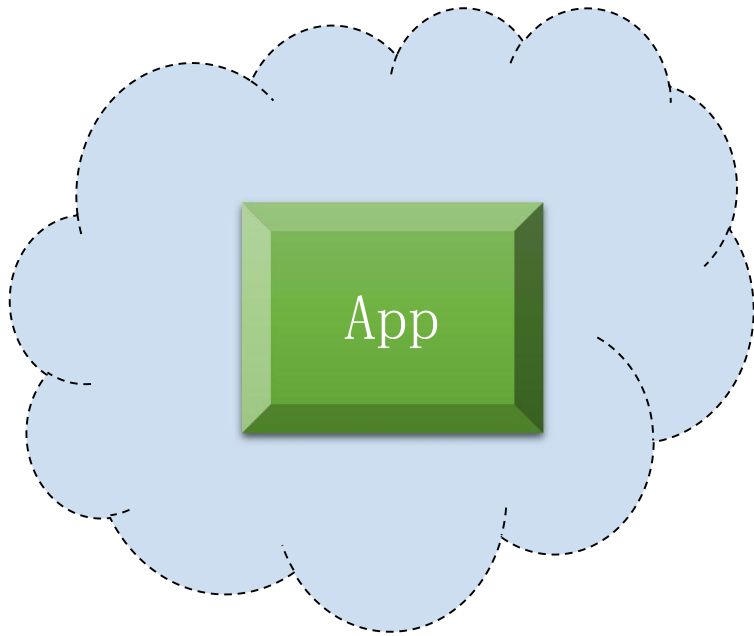


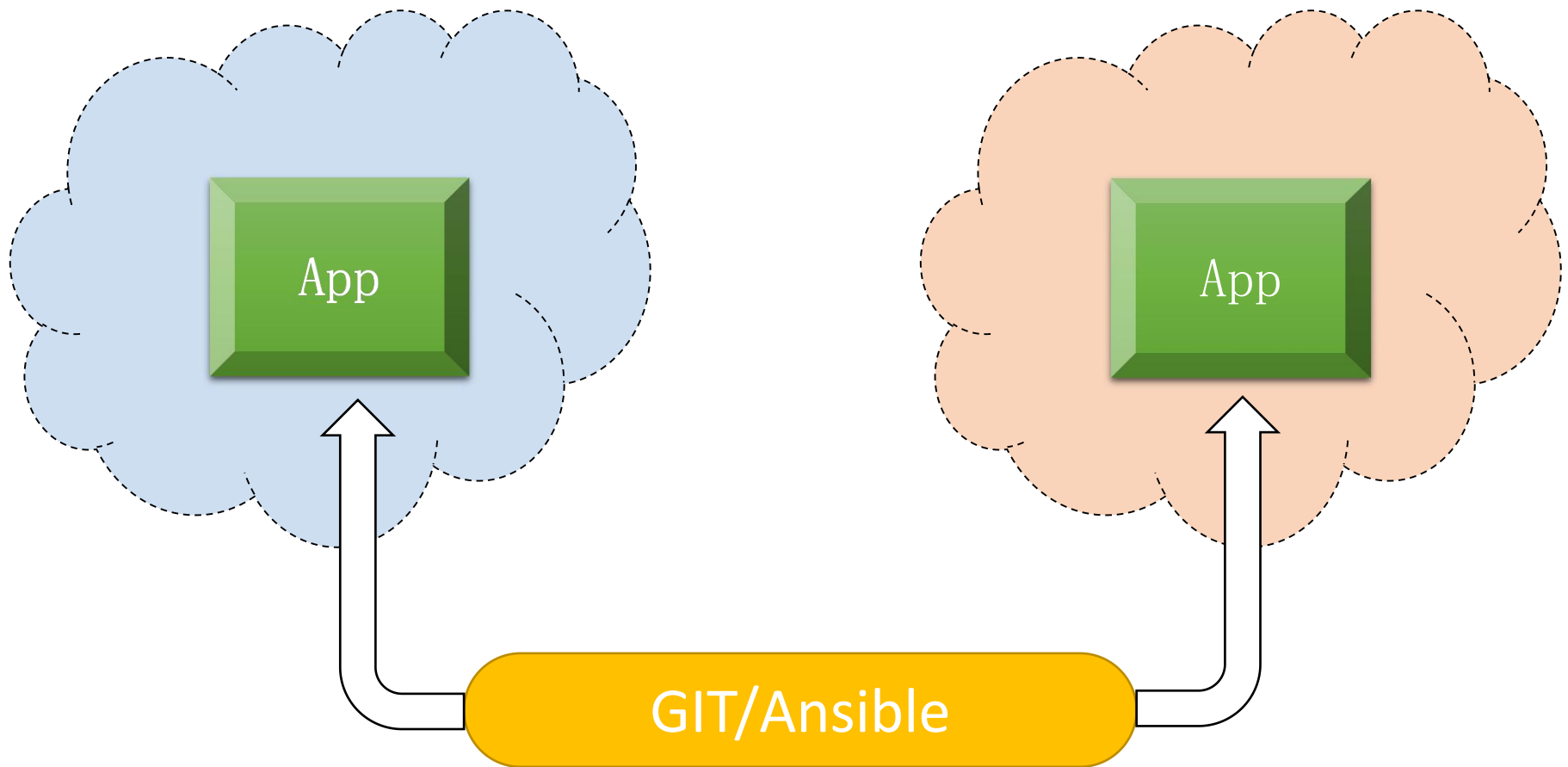
“如果一朵云宕机不可避免
那就把应用部署到多个云上”

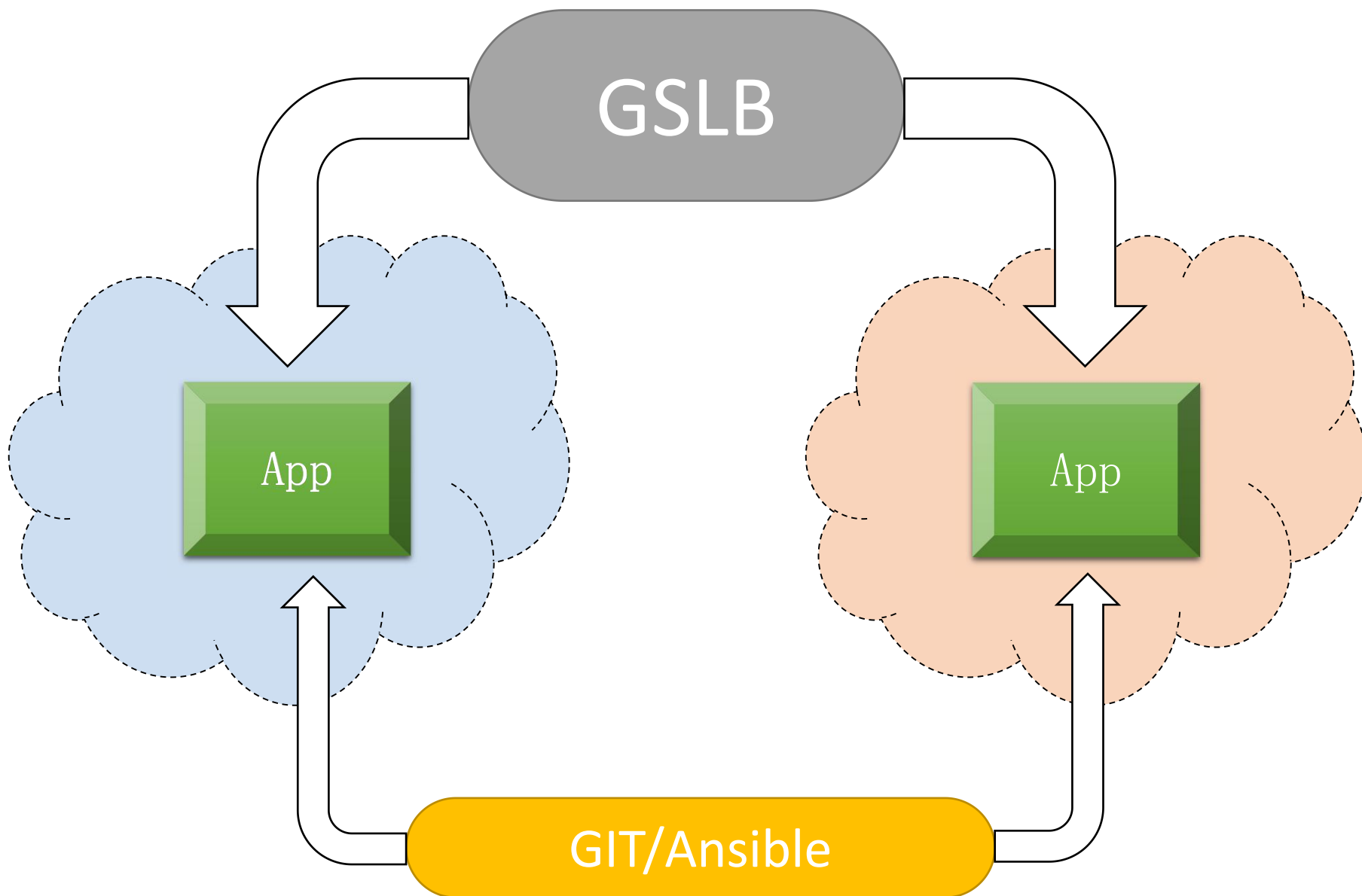
目标

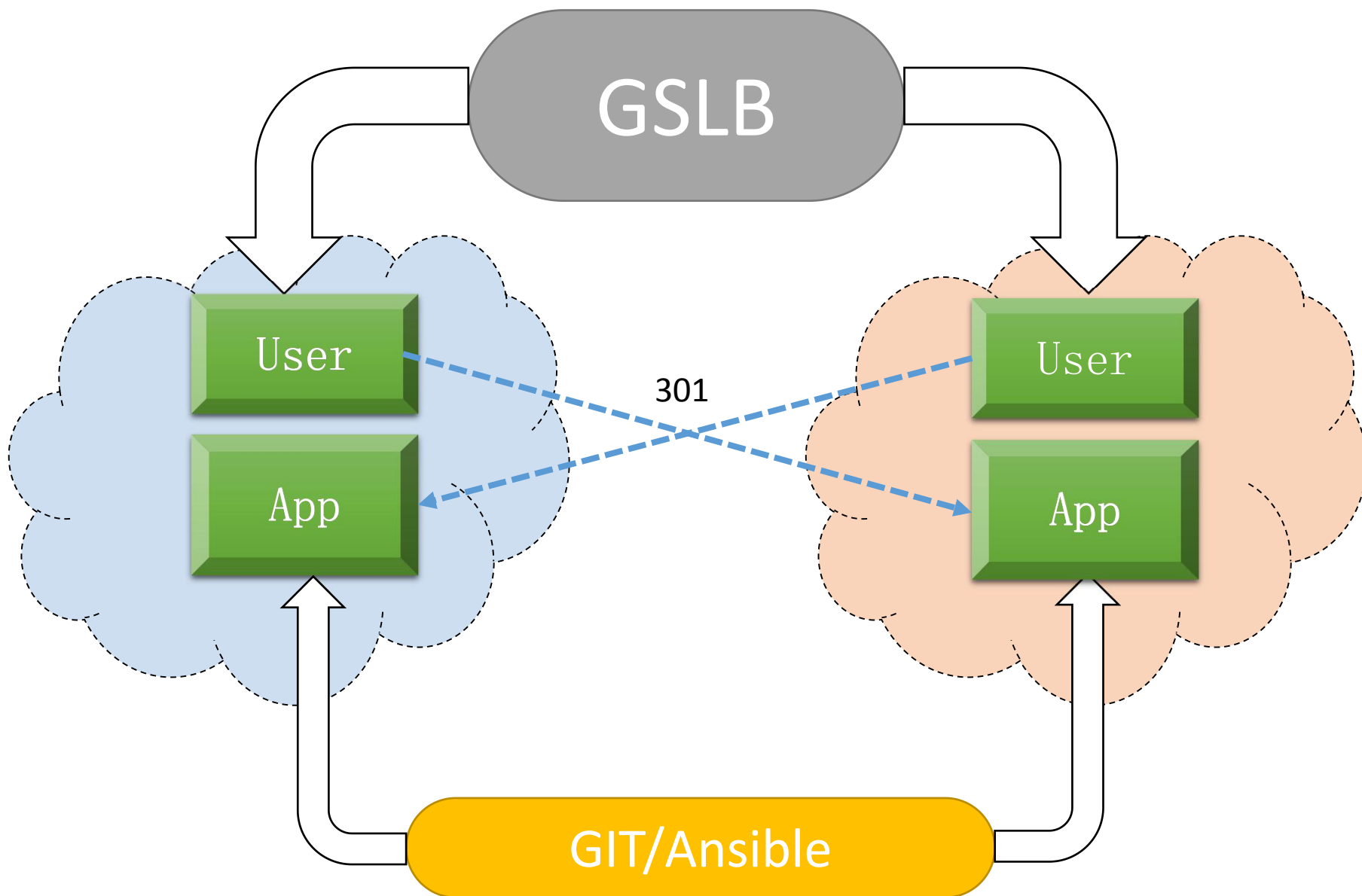
- 多数据中心多活
- 节省成本、可负担的解决方案
- 在灾难发生的过程中，如果无法达到完全可用，则至少应该保证部分可用：
 - 部分业务功能可用
 - 部分客户业务可用
 - 部分数据可用
- 尽量少的人工干预





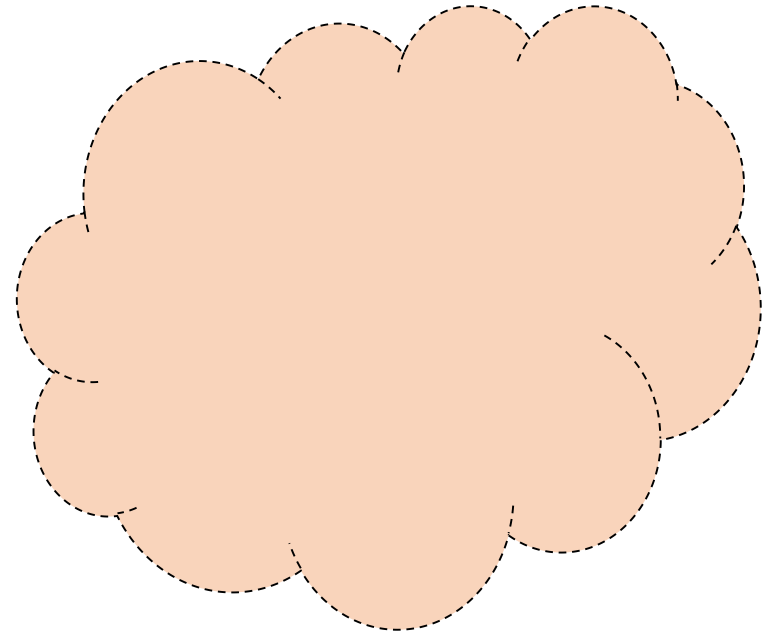
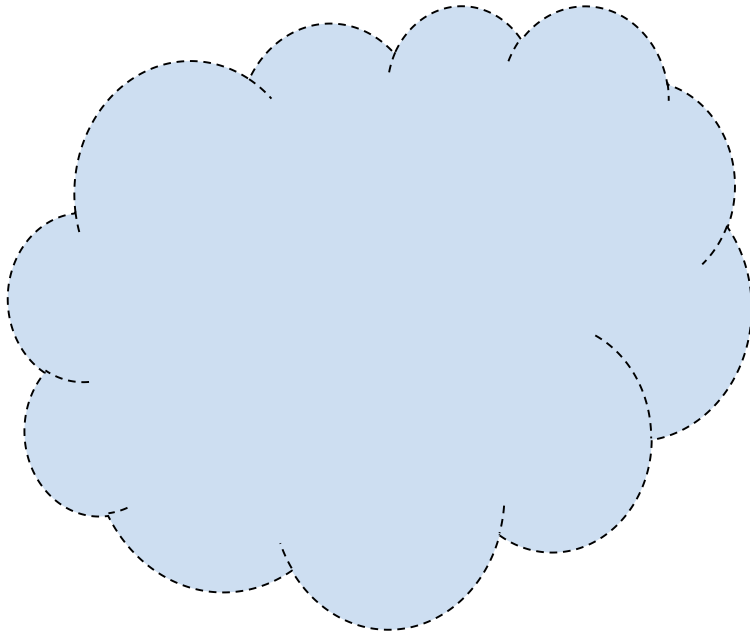








?



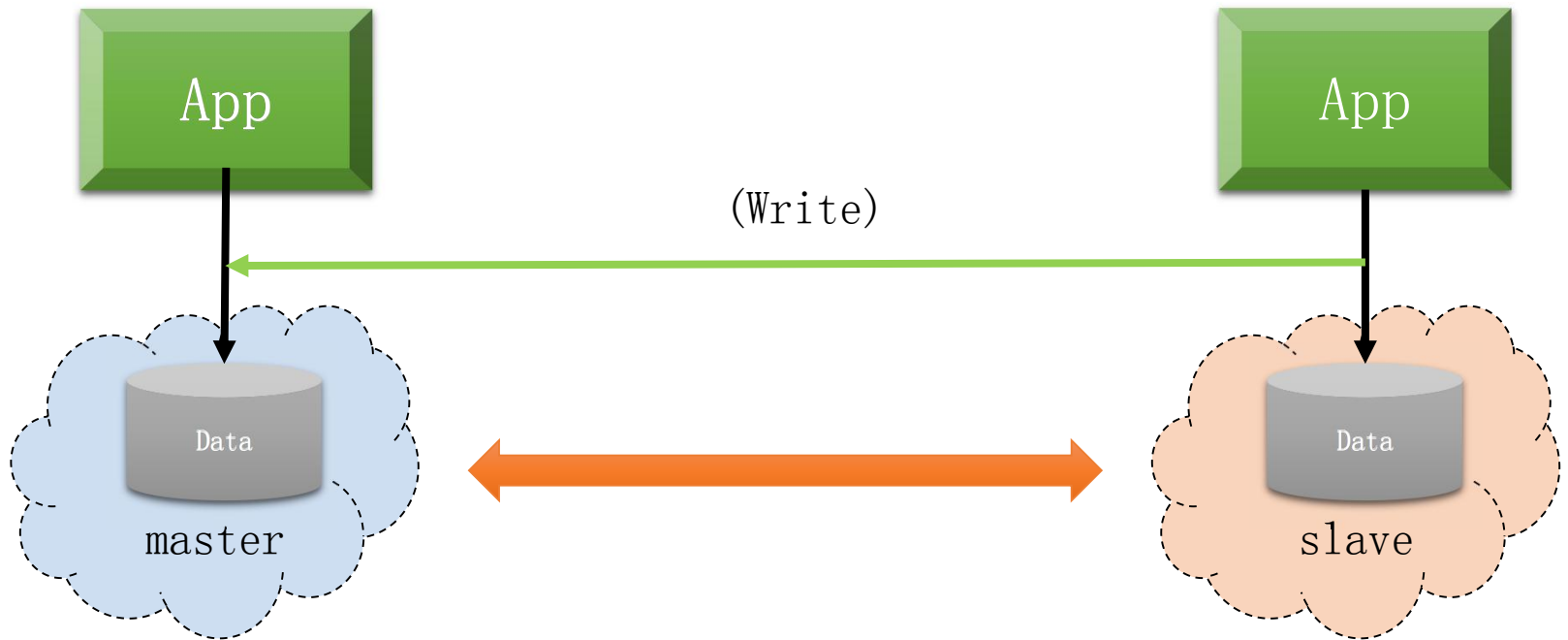
Master / Slave



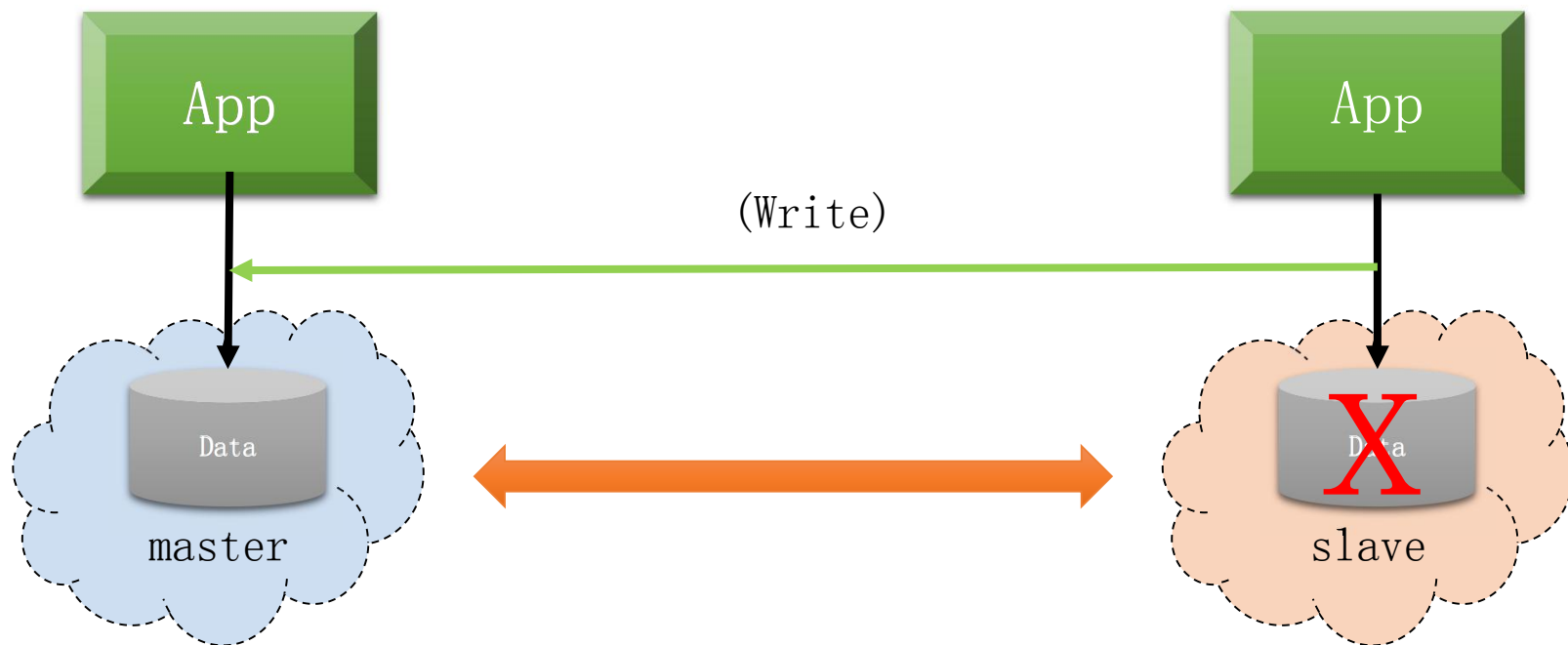
Master / Slave



Master / Slave



Master / Slave



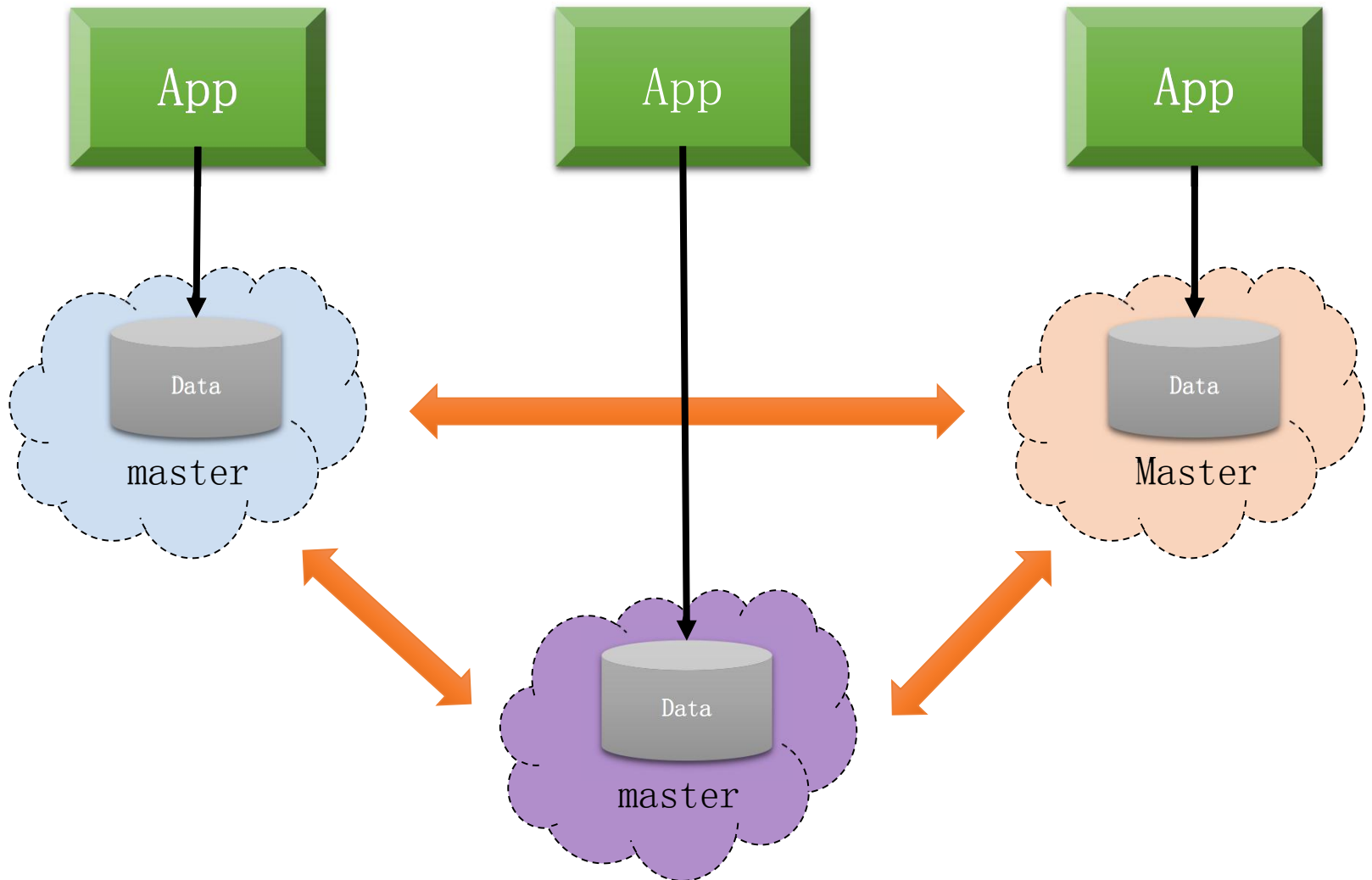
- 如果设计目标是随时保持2份数据拷贝，那么slave宕机的情况下，master应不能写入

Master / Master

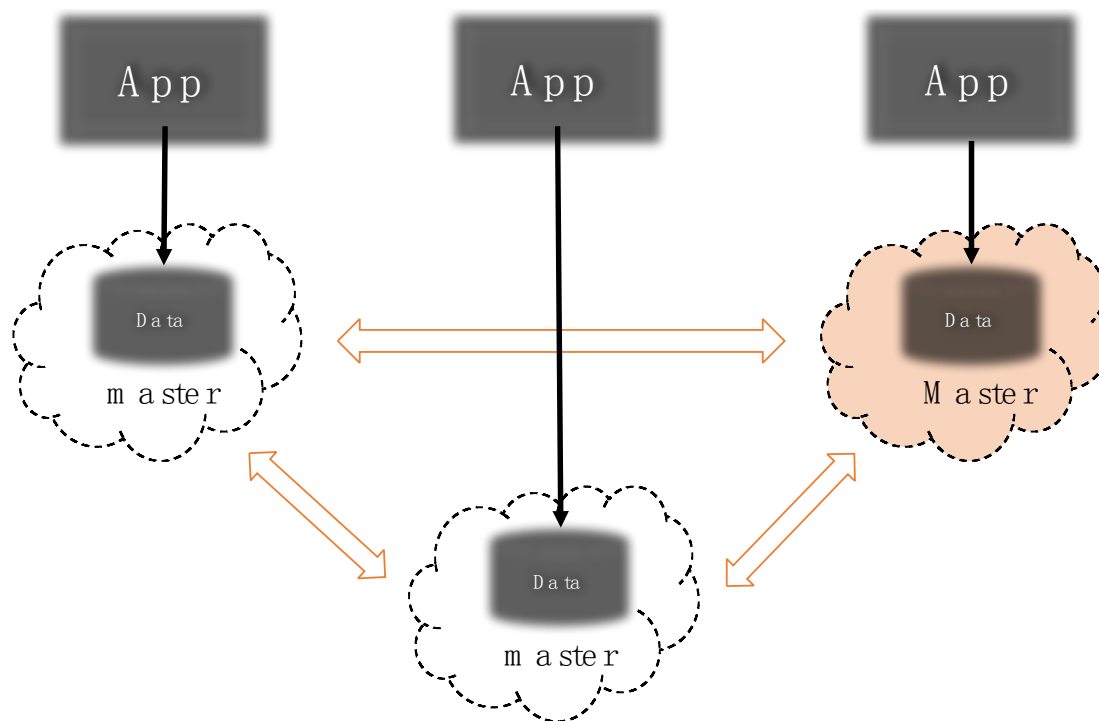


- 需要假定网络可靠（拜占庭将军问题）
- Master越多越慢，代价越高，不可扩展
- 适合单数据中心内部，可以用来解决局部故障
- 跨数据中心则可能需要投入专线

Master / Master / Master



Master / Master / Master

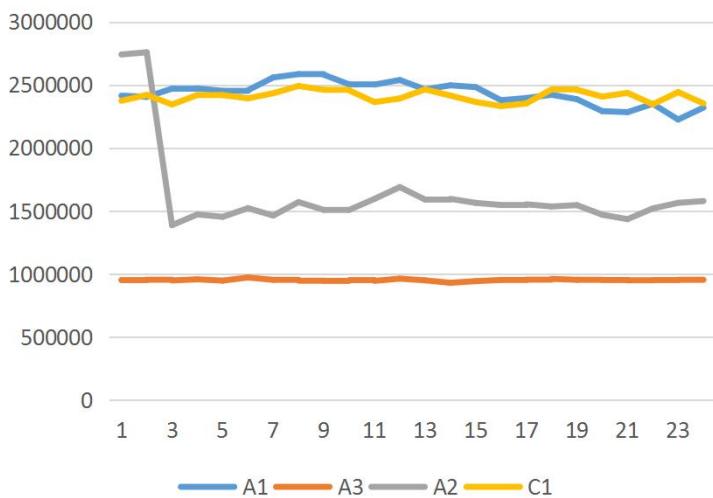


1. 允许任意一个站点宕机、断网的情况下保证依然有2个可用站点
2. 可以形成2:1的多数派解决数据不一致的问题
3. 超高的可用性

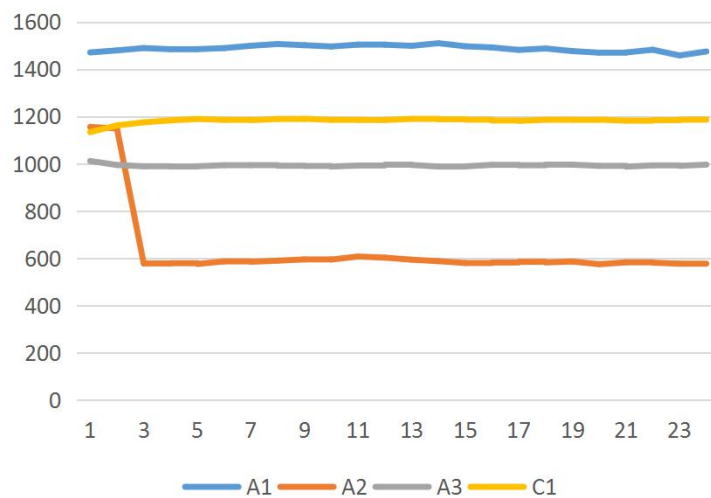
实战之：找到合适的云

不同云的性能

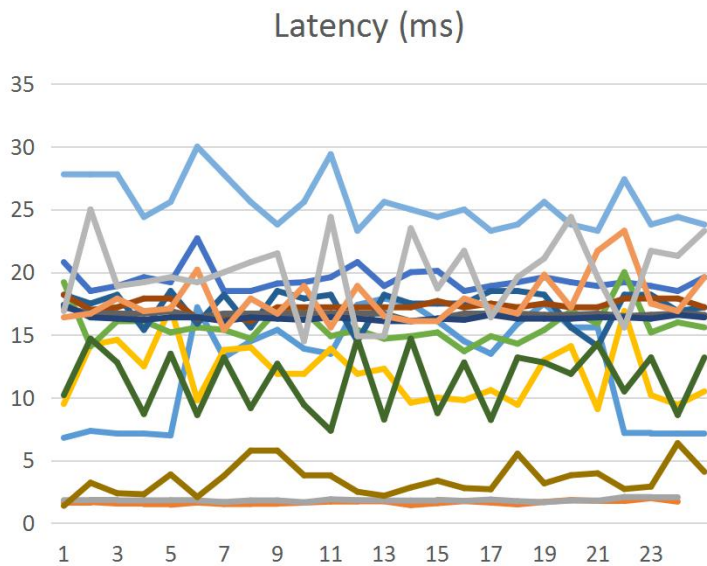
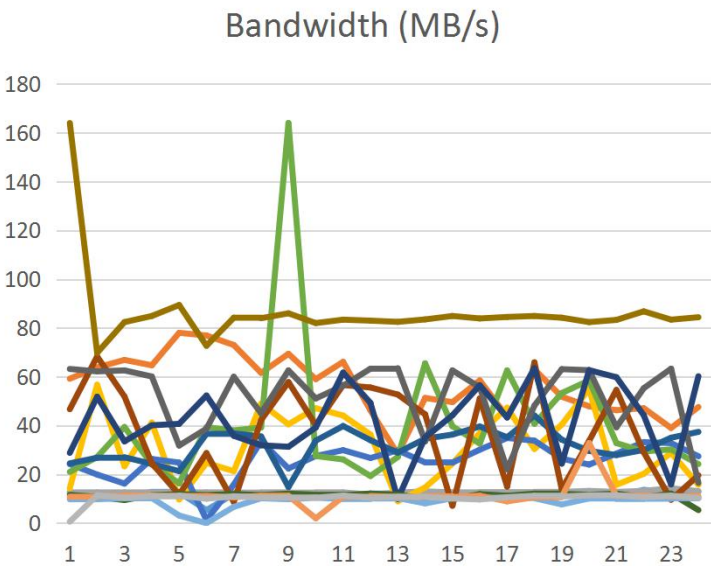
IO Performace



UNIXBench Index

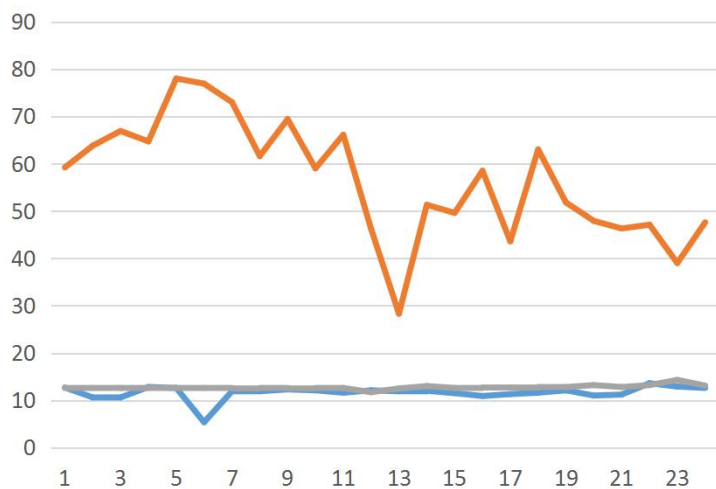


云间的网络

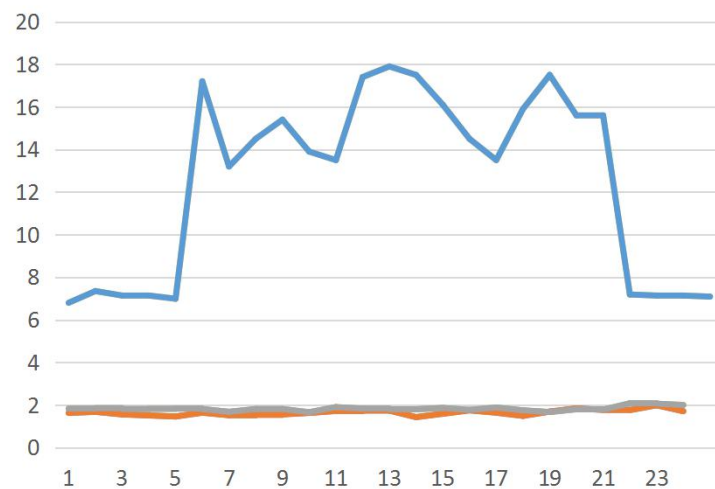


最后的选择

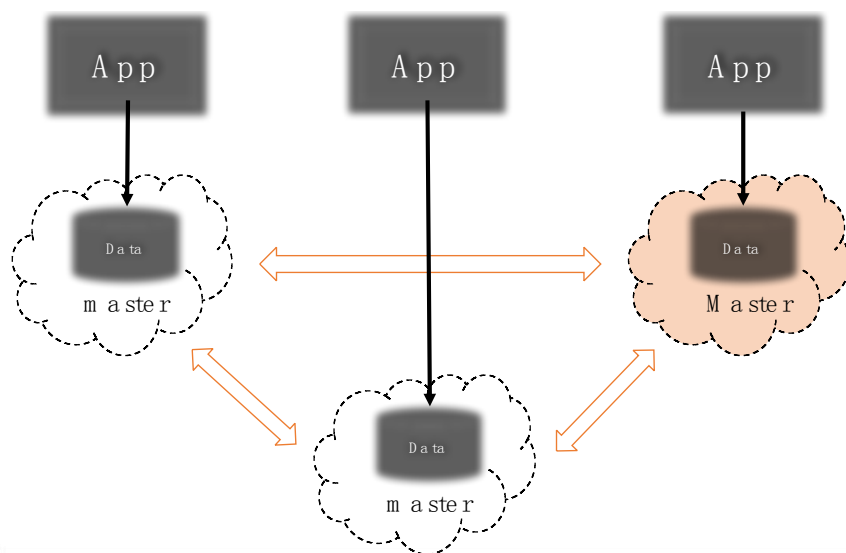
Bandwidth (MB/s)



Latency (ms)



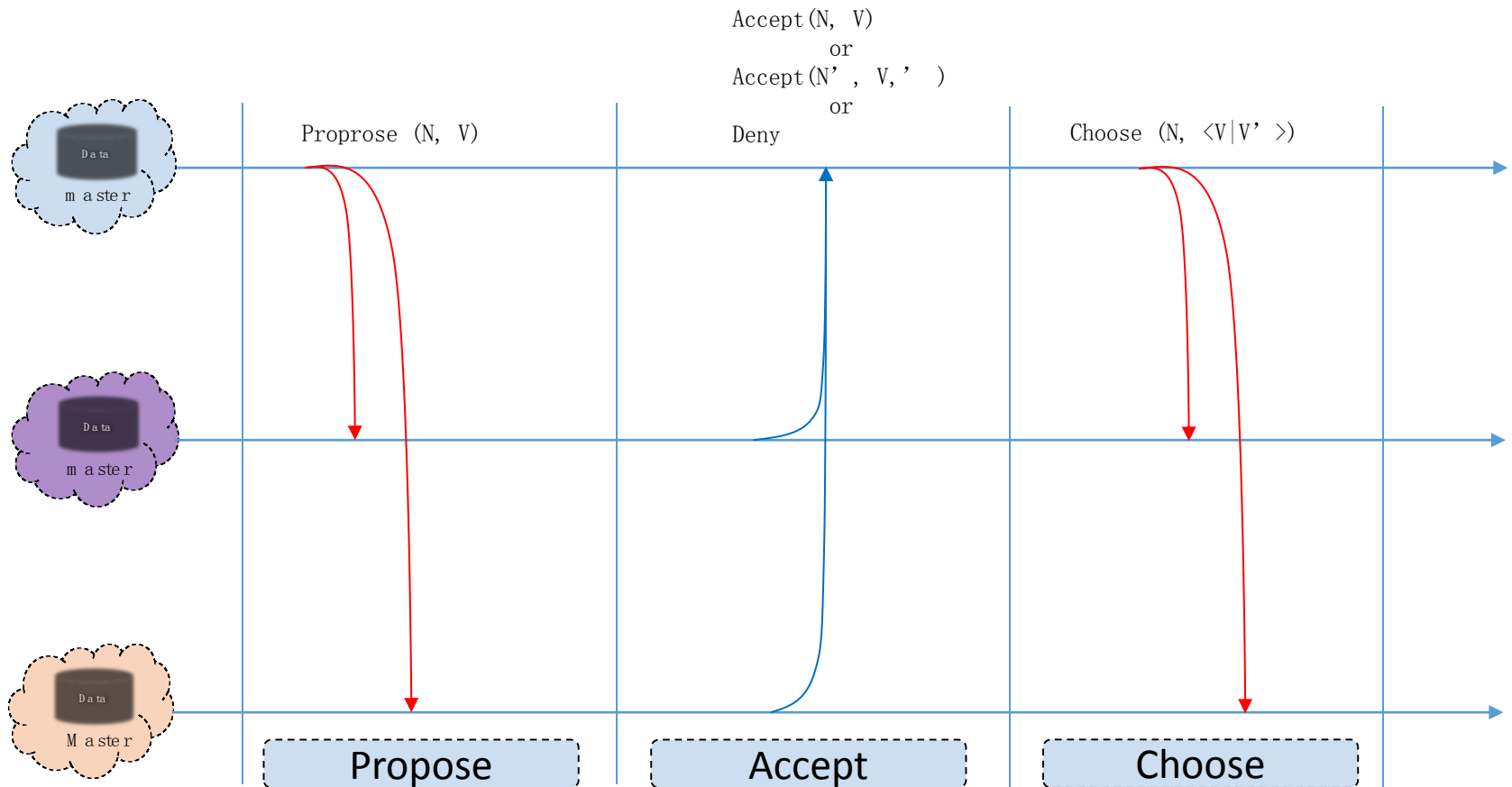
Master / Master / Master



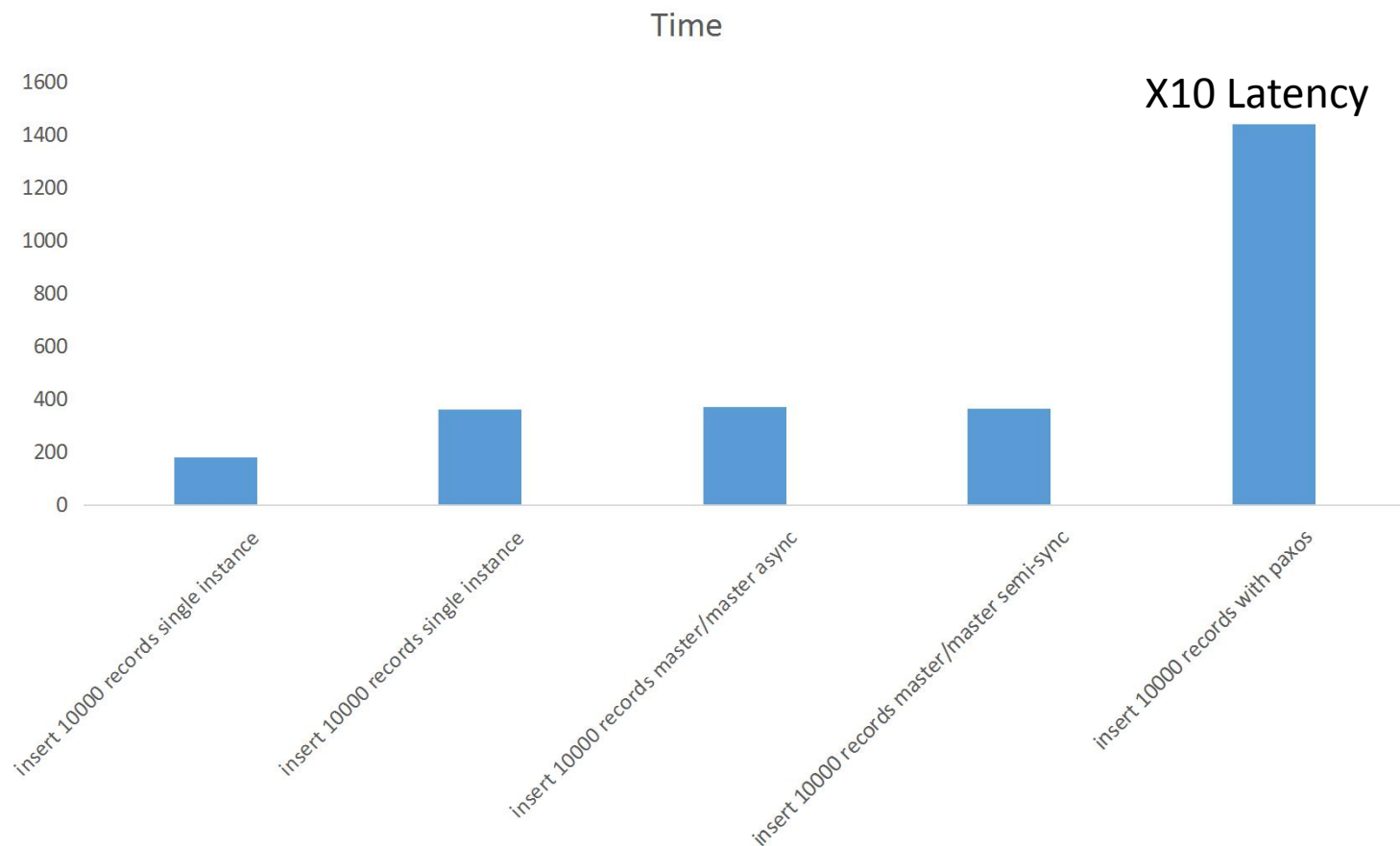
- 公网传输数据（专线成本过高）
- 允许1个云计算数据中心宕机
- 允许网络传输不稳定
- 允许时钟不同步
- 当数据差异发生时可以做到多数票

PAXOS

Paxos PG数据库

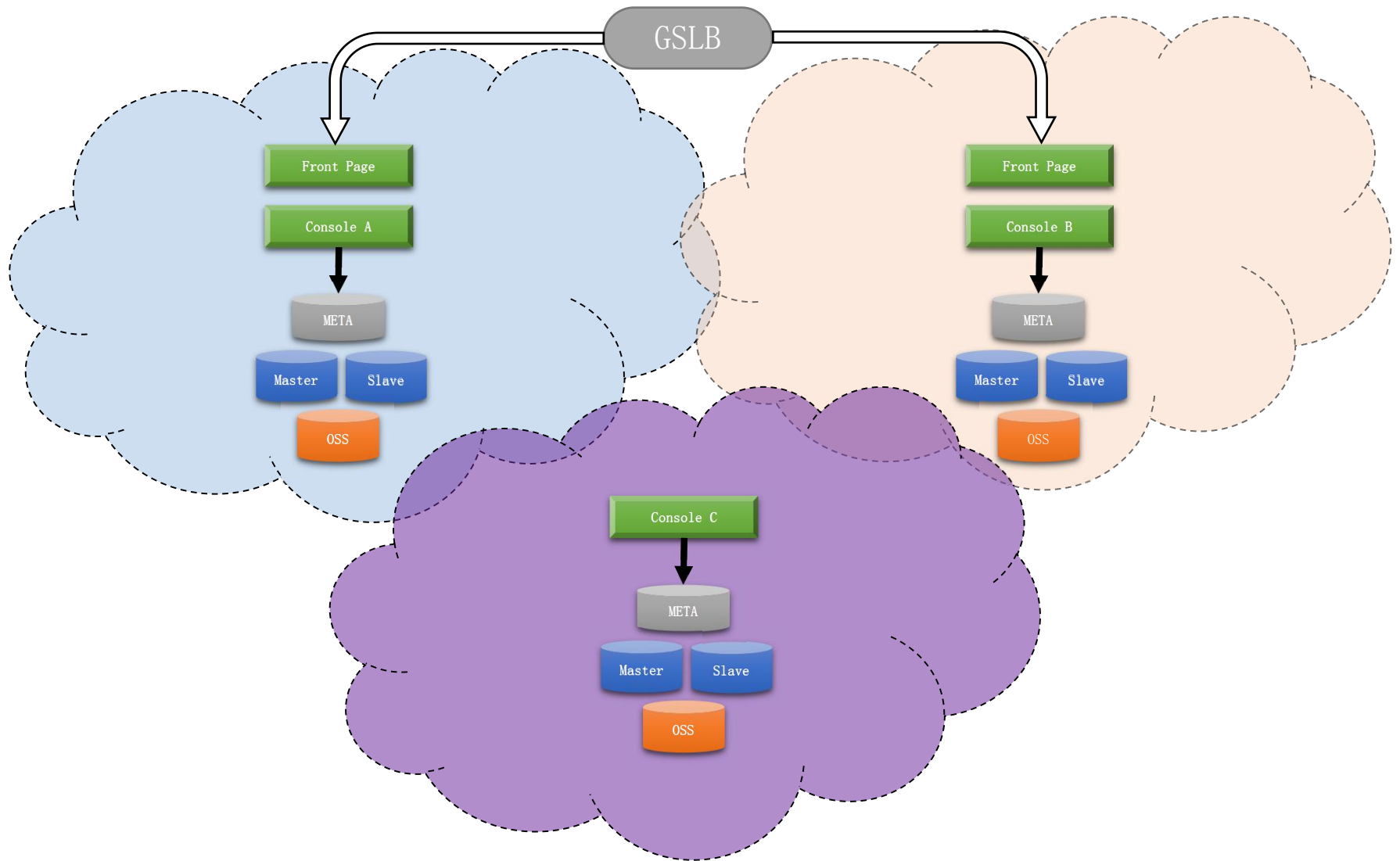


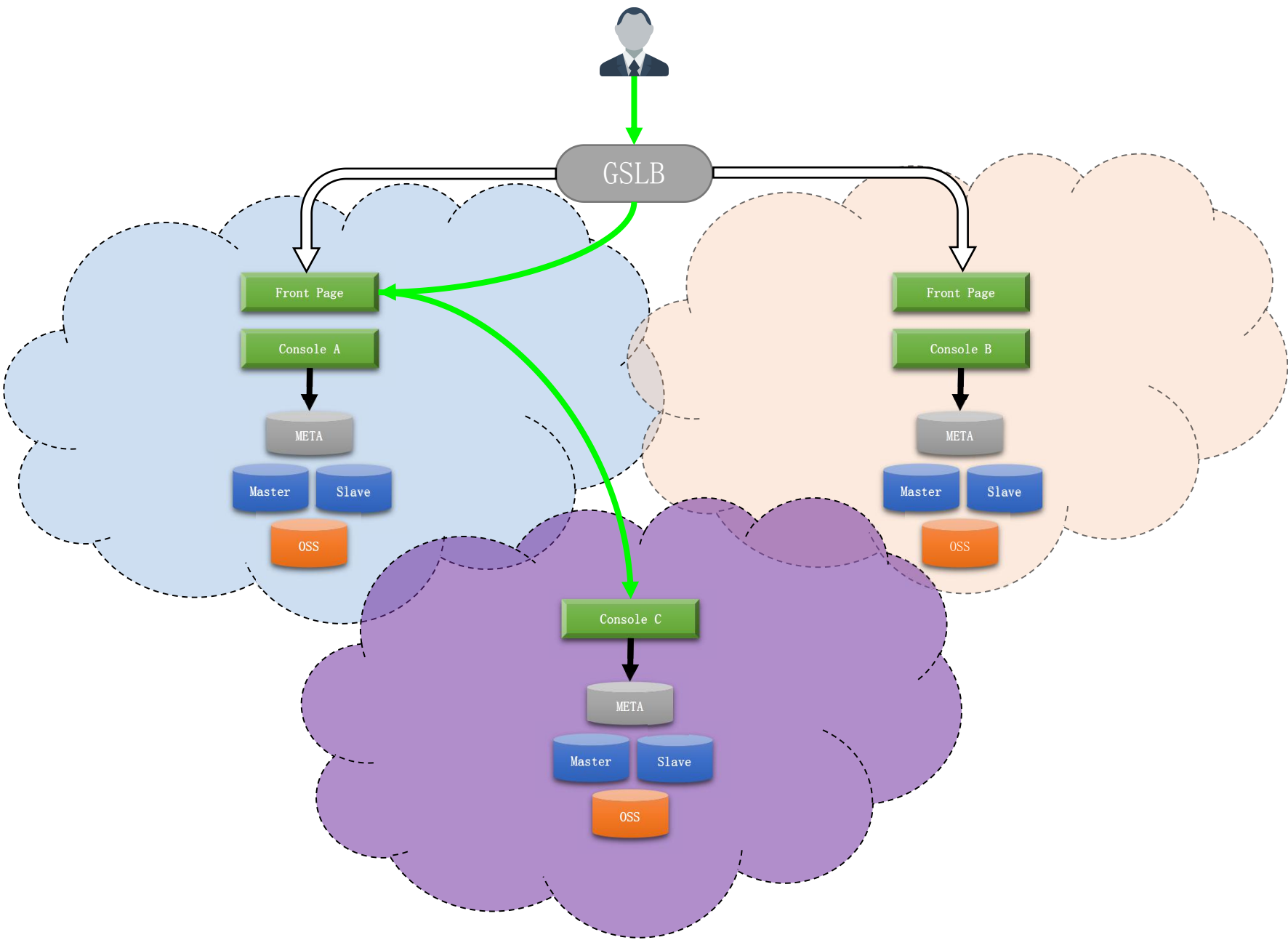
Paxos PG 性能问题

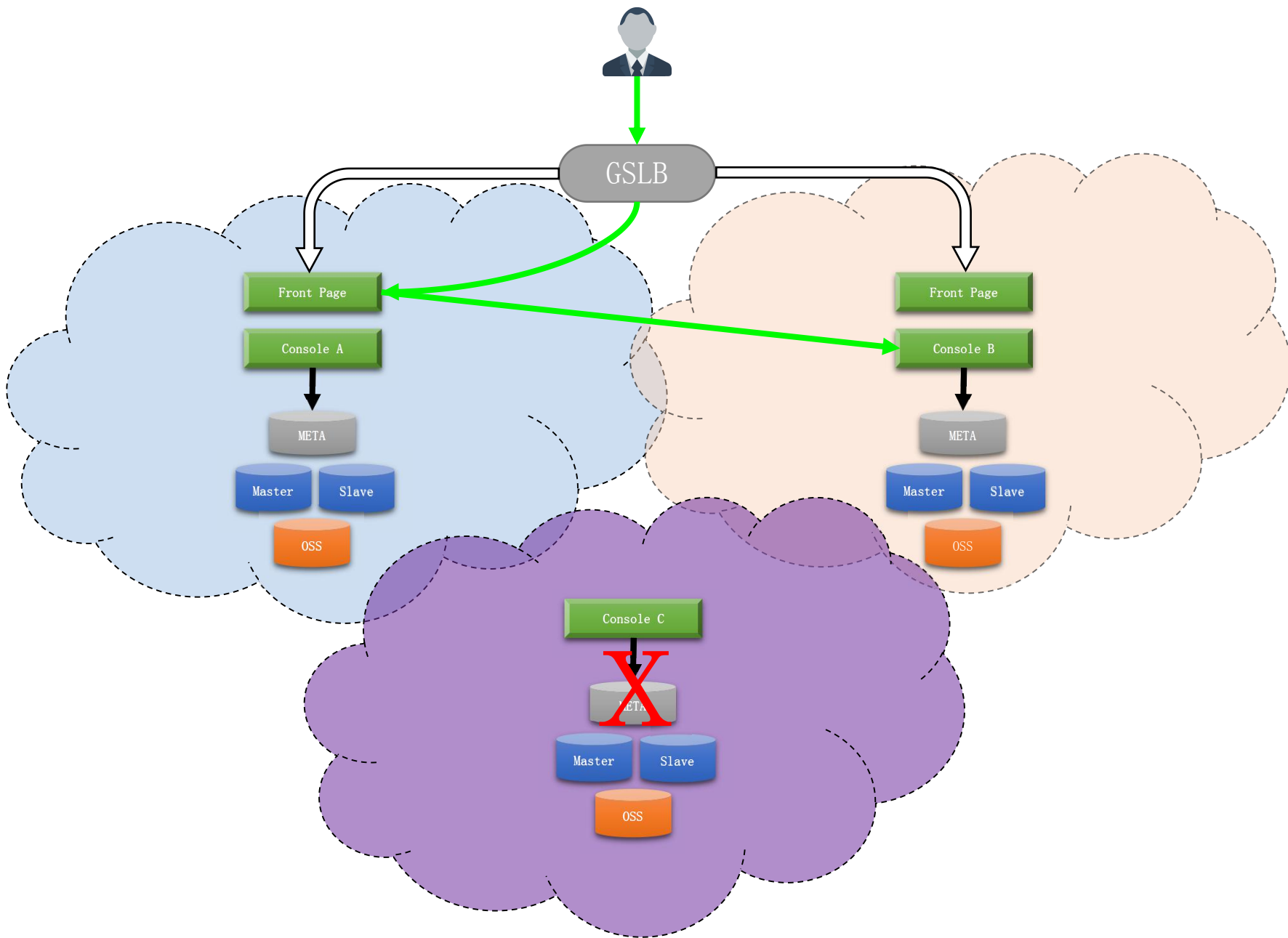


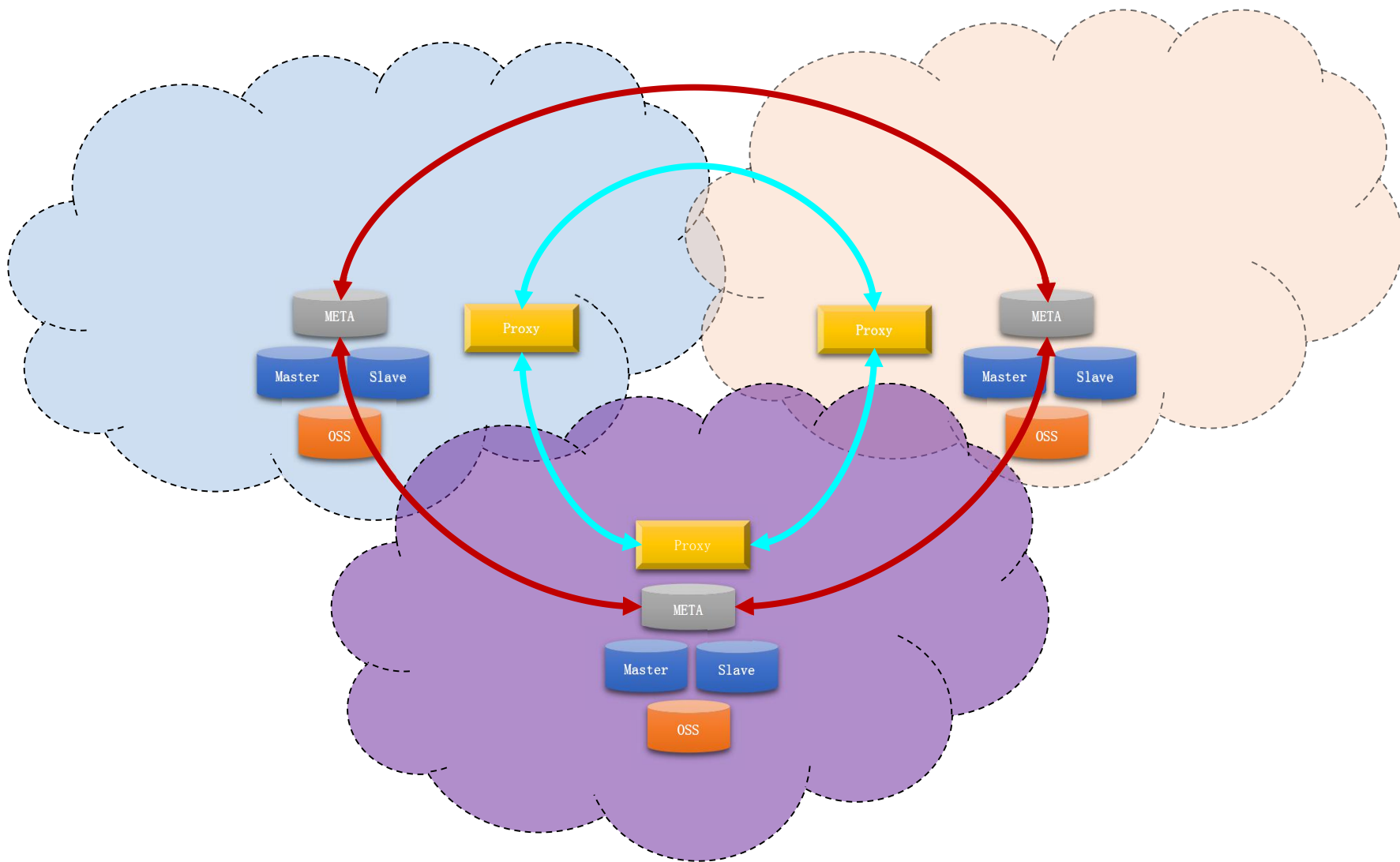
数据的分类处理

数据分类	常见操作	存储选择
元数据 (用户信息、权限、记账)	Create Read Update Delete	Paxos PostgreSQL
资产数据 (云上资源资产信息)	Create Read Update Delete	双实例数据库 + 定期批量复制到从站 + 资源实际状态定期更新
操作数据 (云上资源的操作日志)	Create Append Read	OSS和异步复制









尚需解决的问题

- paxos_pg 事务性的问题
- paxos的算法的数据代理？？？
- 异地数据中心部署的可能性
- 10~100万量级云资产管理

总结

- 故障理所当然发生
- 打破Dev | Ops的边界，双方共同构建可用性
- 理解业务目标，按需架构设计
- 学会妥协

VC³云管家

跨平台、智能化管理运维
让所有云轻松用、易管理

Make Cloud Work for You

WWW.VC3MARKET.COM

