

- 1、资源规划层理想方面
 - a. 自身负责产品的核心业务指标
 - b. 熟悉并着力建设资源平台的计算规划
 - c. 对业务资源数据合理性逻辑, 如千核设备能承载多少pv&bps
 - d. 对全链路各模块负载情况的把控
 - e. 验证产品线资源使用流转情况, 预算->资源规划->资源交付
 - a. 定义业务核心SLA、SLO指标, 并和技术测试达成共识, 双方共同承担考核
- 2、质量保障方面
 - b. 对负责的业务功能逻辑核心链路梳理清晰, 调用关系明确掌握
 - c. 具备排查从网络层、操作系统层、云组件、编译发布工具异常的诊断能力
 - d. 具备度量产品架构存在的短板和可能存在隐患及退场的能力
 - e. 具备故障、事故完成溯源复盘和落地改善能力
 - f. 具备质疑、改良内部运维平台逻辑或问题的基本功, 并建设度统一的产品集群
- 3、研发能力
 - a. 具备java, golang, python任意两种语言的开发能力
 - b. 具备SQL、Hadoop、Spark、Flink等开源产品的使用及开发能力
 - c. 能够熟练使用CI/CD工具, 自身能够高度自律, 严格按照代码规范进行开发
 - e. 具备参与研发架构设计能力, 提出建设性且对业界顶级标准的设计方案
- 4、云、中台应用运维能力
 - a. 掌握不限于阿里云、腾讯云、AWS等云厂商产品的使用逻辑, 以及迁移、容灾、融合业务场景的架构方案
 - b. 掌握云原生社区各类成熟的建设方案, 如K8s, Prometheus等
 - c. 具备中台架构的思维, 包括但不限于打造中台能力、成本精细化拆分、链路透明化展示、资源快速交付等能力

需要具备的能力

SRE

基础方案

- 容量
 - 容灾基本条件: 业务模块需要具备多节点部署能力, 并且节点均为无状态。
 - 多节点容灾: 底层调度系统在节点或主机故障或者网络设备不通时, 可将业务节点迁移至新的资源上以正常提供服务。
 - 跨机房容灾: 单一机房不可用情况下, 名字服务可将该机房流量导入就近机房, 并且就近机房能够根据流量现象进行相应的资源自动扩缩容。
 - 跨地域容灾: 跨地域容量完全冗余, 具备自动切换能力, 共同承载服务。
- 容量
 - 各模块已占用资源的容量水位情况
 - 资源功能能够补充的容量情况
 - 弹性能力
 - 自动扩容: 根据合理预算, 保证资源池有充足设备 (适用于新闻APP后台策略, 因无法准确预估下一次大流量峰值时间)
 - 自动缩容: 结合业务模块能力范围, 统计低负载实例, 制定合理的最小保留实例数, 以保证平台自动缩容不会影响到业务
 - 自动迁移: 云环境故障自愈, 业务容器需满足无状态, 不写死IP, 允许平台在任意时间动态调度迁移
- 性能指标
 - QPS
 - 成功率
 - 耗时
 - DAU、MAU
 - 日志、月活
- 核心业务指标
 - UV、PV
 - UV指向一次登录多次访问记录为1个访客、PV指页面累计点击量
 - 次留: 次留人数是指新增的用户在第二天登录的用户数量。
 - PCU: 最高同时在线数
 - 业务属性指标: 如类似抖音日视频上传量、成功率、视频过审率等
- 资源指标
 - CPU、GPU、LOAD、内存、网卡、磁盘、连接数等
- 监控
 - 监控覆盖范围
 - 1. 人工构建负责业务模块从公司流量接入、到最后的存储设备全链路的监控链
 - 2. 绘制模块间调用关系的大盘视图, 落地全链路监控
 - 有效性、时效性
 - 测试告警性能, 新建一个告警, 将告警接收人写进自己, 降低阈值触发告警 (笨办法)
 - 在测试相关业务并允许操作情况下, 通过演习、注入故障、人工触发节点不可用、流量陡增故障等手段, 验证告警的准确性、时效性
 - 告警分级
 - 重要告警电话通知, 一般告警聊天工具通知, 可第二天处理的告警, 邮件通知
 - 告警收敛
 - 将相关告警人工与自动关联, 致力于做到一条告警即可屏蔽业务大致问题所在, 而不是发多条告警, 人为判断。
 - 监控日报、周报
 - 工聚该业务单位天内, 触发了多少次告警, 多少条未处理, 收件人可反馈超预期无效告警
 - 通过收件人反馈告警有效性, 监控平台可根据结果进一步收敛告警
 - 同时也可以反馈相应SRE在问题发生时的响应及时性
- 运维链
 - 业务本身各模块间调用关系
 - 云上资源: 例如: CDB、Kafka、云Redis、ES等
 - 业务本身所依赖中间件的调用关系
 - 自研中间件: 例如: 自研kafka, hadoop等
 - 业务提供给上游调用的接口情况
 - 可通过业务自行上报日志至ES集群用以可用性监控
 - 亦可通过拨测接口验证接口可用性
 - 业务依赖的下游服务SLA
 - 数据链路, 把打好变量链接, 明确操作负责人, 并且变更需上级审批
- 变更发布
 - 业务模块是否需快速发布部署能力
 - 业务模块快速发布部署能力
 - 业务模块支持灰度, 灰度等包括节点级别、用户级别等
 - 业务模块具备从Git仓库到可执行程序流水线的自动化构建能力
 - 发布配置文件等权限收敛到最低, 通过平台化的标准规范限制研发可操作空间, 避免人为疏忽导致发版失败, 影响线上业务
- 资源规划
 - 资源展示Portal: 展示各业务所使用的资源概数、CVM、云资源配置及使用
 - 资源申请: 每季度核算自身业务当前资源使用情况, 取研产品确定下一季度是否有新功能、活动上线、促销申请资源
 - 成本分析: 比如资源管理平台资源, 提供给各上游业务的能力, 根据上游调用量算成硬件成本, 同步给运维能力, 形成制衡
 - 资源交付效率
 - 公司资源交付给业务效率
 - 交付设备后, 业务场景利用率是否达标
- 演习 (蓝红工程)
 - 制定完整可行的演习方案
 - 每周固定一个月进行业务演习
 - 范围上下游
 - 在测试环境, 共同验证演习方案可执行性
 - 挑选业务合适时间, 进行线上演习切换
 - 演习过程是否满足要求 (演习是否成功)
 - 演习平台, 负载均衡, 监控系统等基础设施是否在演习过程中满足预期
 - 故障演练结果形成总结
 - 参与人员是否安装事先准备好的故障脚本能解决问题
 - 故障预案是否详细可执行, 是否达到预期, 查漏补缺
 - 形成文档落地改进措施
 - 需求: 本风险范围及感知方向的灵感
- 故障操作
 - 故障面
 - 对于可预知的故障, 需提前做好故障预案
 - 能自动化的, 通过故障自愈平台解决
 - 不能自动化的形成文档, 在演习过程中主动触发, 反复训练此类故障的应急处置, 并尽量往平台一键操作方向靠拢, 防止真正业务故障发生时, 人为操作复杂, 容易出错
 - 建设对监控链大盘, 在故障时能快速获取定位问题作用
 - 故障中
 - 对于不可预知的故障, 优先选择回滚、降额、切流量至可用节点等操作, 以恢复系统可用性为最高优先级, 事件再去详细定位问题
 - 故障后
 - 每一次故障都是难得的事发业务系统鲁棒性宝贵经验
 - 以日志的形式详细记录故障发生前后的操作 (如右图示例)
 - 故障复盘
 - 通过分析故障, 验证基础平台是否正常运行期间预期, 例如监控是否及时准确告警, 名字服务是否及时删除故障节点等
 - 故障能否通过技术手段避免, 或者发生时是否能够通过规范处理流程避免人为失误
 - 是否为已知故障或者曾经发生过的故障, 总结已知故障再次发生的原因