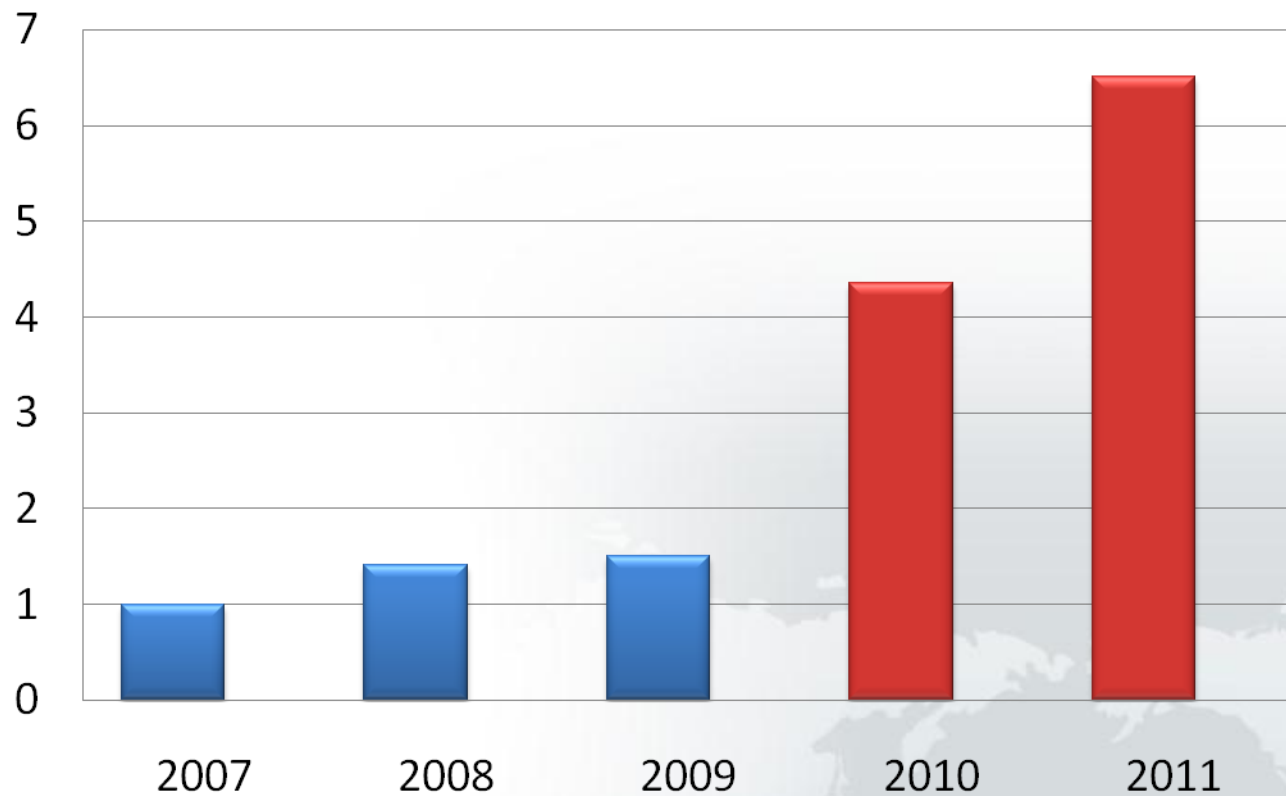


百度的脚印



今天的百度搜索离线系统

- ✓ 百亿网页需要在一天内完成分析
- ✓ 每天千亿级链接分析
- ✓ 数个异地机房
- ✓ 数十个异构集群
- ✓ 数万台机器





百度离线集群整合之道

内 容



背景



目标



架构



效果



QA



目 标

“简单、可依赖” 的CPU/MEM池

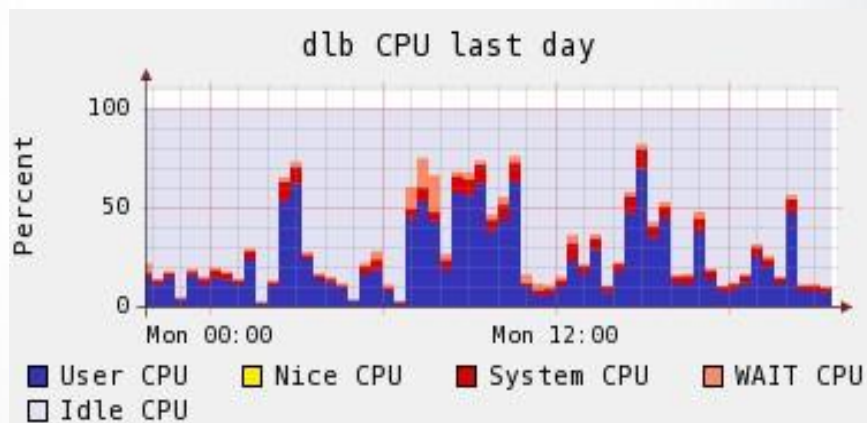


目标--- 整合的资源

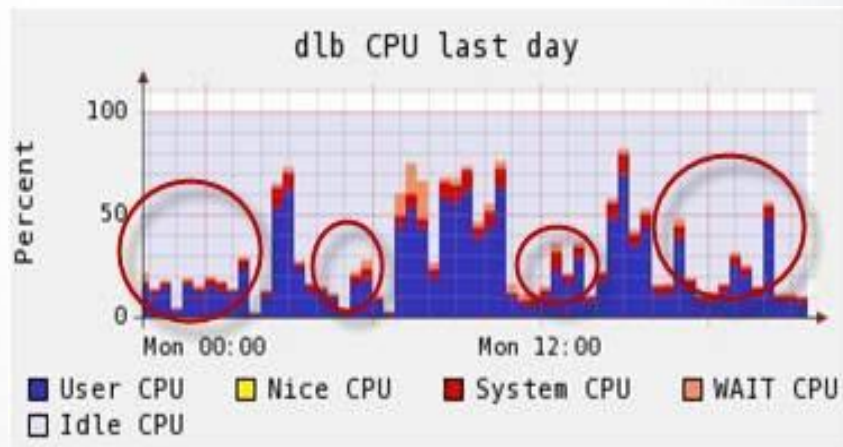
System	硬件偏好	OS特性	分布式系统	使用规律	运维要求
网页库	IO吞吐	block access	自有系统	稳定	<1分钟
链接库	IO吞吐	block access	开源+自主系统	无规律	数分钟
基础检索	CPU, MEM, IO	mem, cache, tcp/ip, cgroup	无	稳定、有规律	<<秒
调研	MEM, IOPS	mem, cache, tcp/ip, cgroup	无	稳定	数分钟



目标--- CPU使用规律



目标--- CPU使用规律



Volunteer Computing system

project	start	Affiliation	Area	Peak#Hosts	Current status	Computing power
Bitcoin	2009	Bitcoin	Digital currency	6,500 [40]	active	3500 TFLOPS[41]
Ibercivis	2008	National science agencies of Portugal and Spain	Biomedicine, physics, other	14,710[39]	active	6 TFLOPS
Rosetta@home	2005	University of Washington	Biology	100,000	active	100 TFLOPS
Einstein@home	2005	LIGO	astrophysics	280,000	active	370 TFLOPS
Climateprediction.net	2003	University of Oxford	Climate change	150,000	active	90 TFLOPS
BOINC	2002	University of California, Berkeley	Biomedicine, other	527,000	active	5430 TFLOPS
Folding@home	2000	Stanford University	biology	406,000	active	7870 TFLOPS
SETI@home	1999	University of California, Berkeley	SETI	362,000	active	684 TFLOPS
GIMPS	1996	?	mathematics	10,000	active	40 TFLOPS

BVC	2010	BAIDU	Search engine	12,000	active	740 TFLOPS(4th)
-----	------	-------	---------------	--------	--------	-----------------

ref : http://en.wikipedia.org/wiki/List_of_distributed_computing_projects

BVC : Baidu Volunteer Computing system



VC横向对比

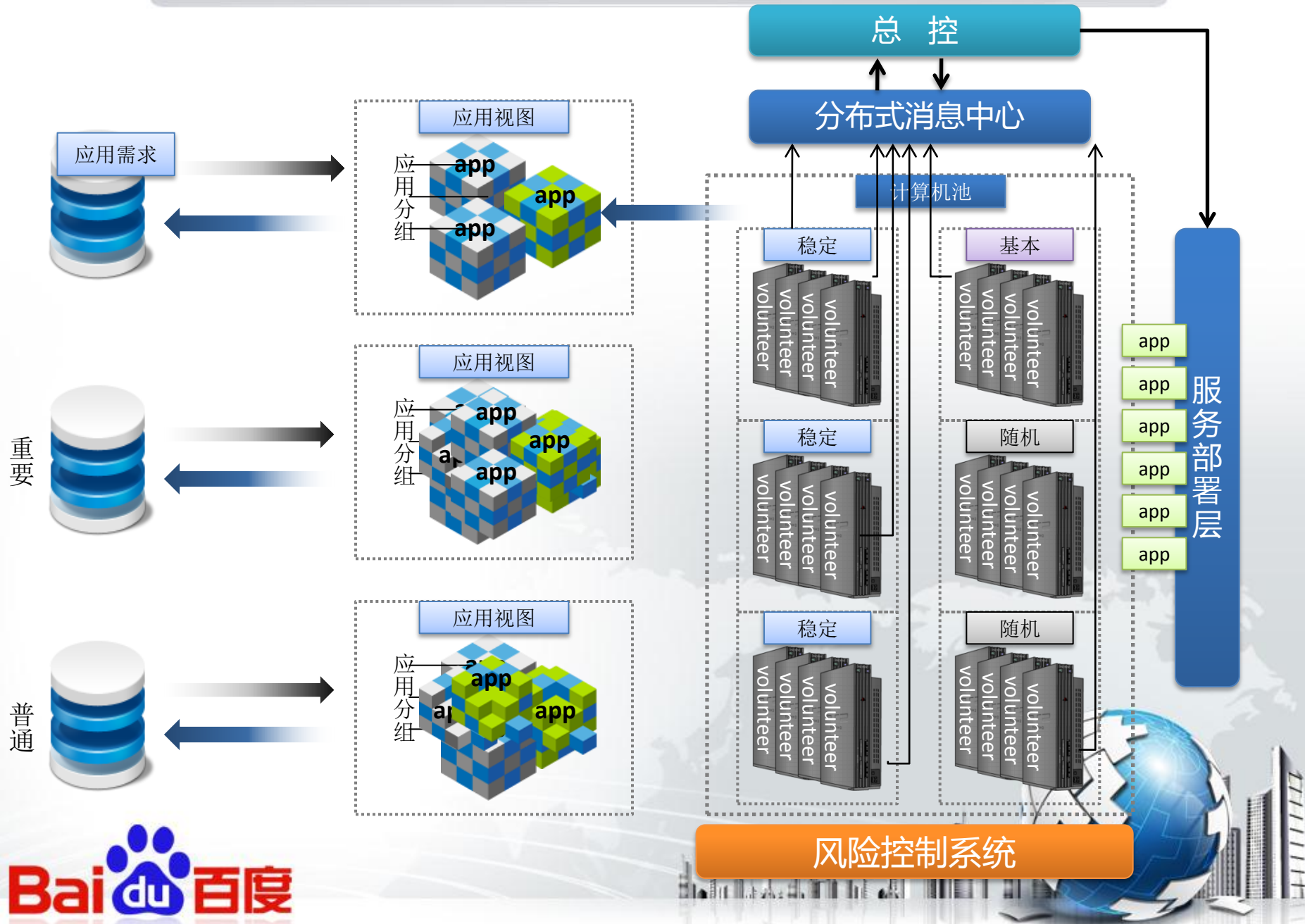
- ✓ 复杂时效性需求
- ✓ 灾难处理的代价极大
- ✓ 网络整体的交互量大
- ✓ BVC不用解决的问题
 - 存储的整合
 - 计算可靠性



设计目标

- ✓ 优先级控制，全局负载均衡
- ✓ 风险控制，最大利用率
- ✓ 插头、插座





BVC总控

- ✓ 任务启停，权限控制，合法性校验
- ✓ 优先级控制
- ✓ 负载均衡
- ✓ BVC 总控不是系统中的单点



BVC优先级控制

✓ 时效性需求

- 秒（分钟），小时，半天，天，周，月，季度
- 优先级动态调整（时效性，相关性，覆盖率）

✓ 基于优先级分级控制

- 优先级越高，越早调度
- 同一级别内部基本均匀
- 被延迟越久的任务，“级内优先级”越高
- 优先级动态调整，人工确认



BVC负载均衡

- ✓ 通过volunteer分组，控制组间均衡
- ✓ 分组内部通过计算单元轮询控制均衡
- ✓ 任务在集群之间迁移



BVC 风险控制 最大利用率

✓ 风险控制

- 内核控制 ($<<s$)
- 智能分析 (10s, 30s)

✓ 最大利用率

- 有规律：可配置，准确率很高
- 无规律：行为分析，动态调整

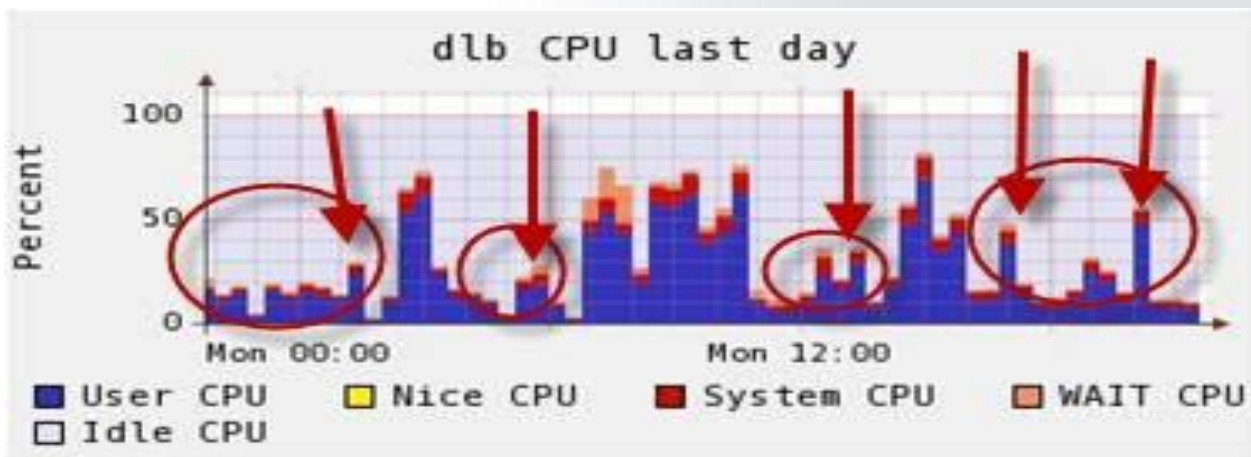


BVC 风险控制体系



BVC 风险控制 vs 利用率

暂停接受计算单元

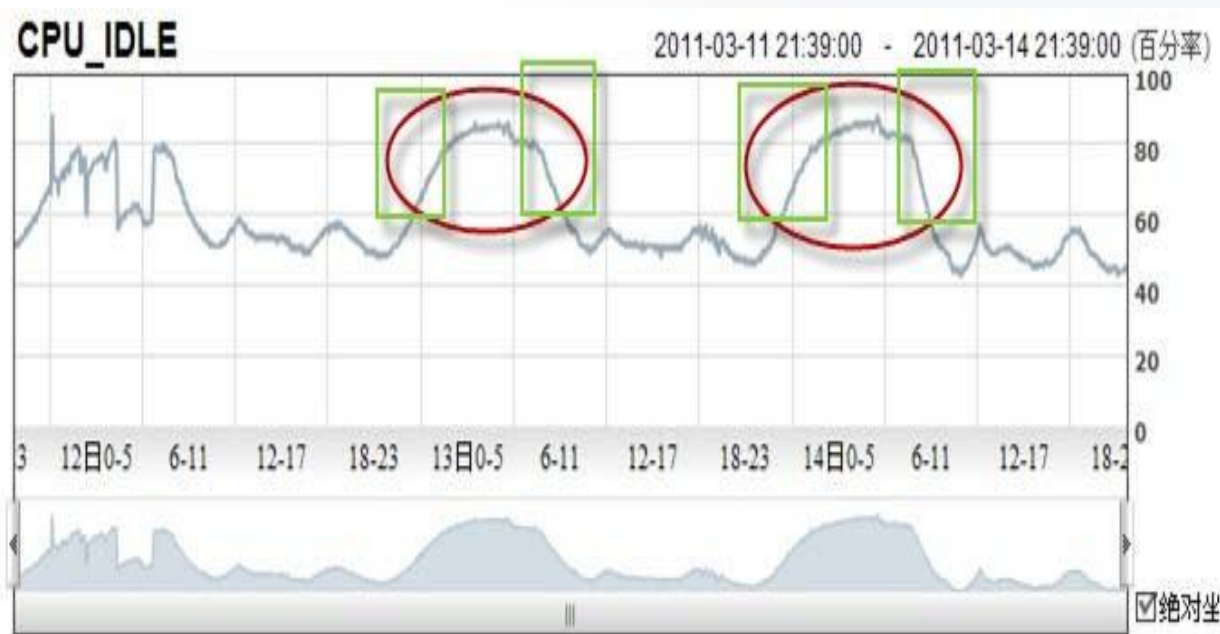


BVC 风险控制 vs 利用率

任务线程、内存控制

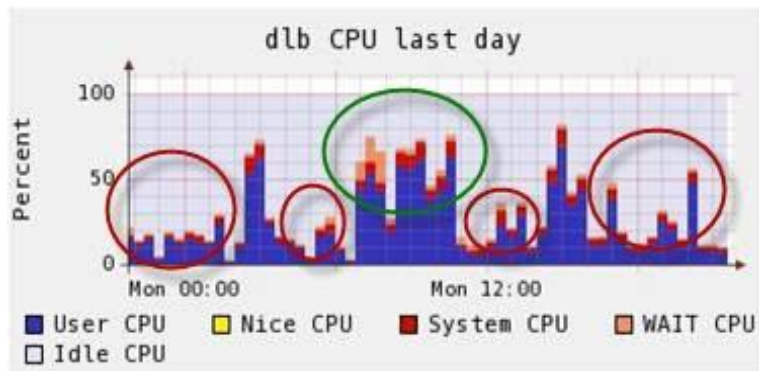
$\text{total mem} = \text{default} + \text{thr\#} * \text{thr_mem}$

$\text{thr \#} = \min(\text{cpu_idle}, \text{free_mem})$

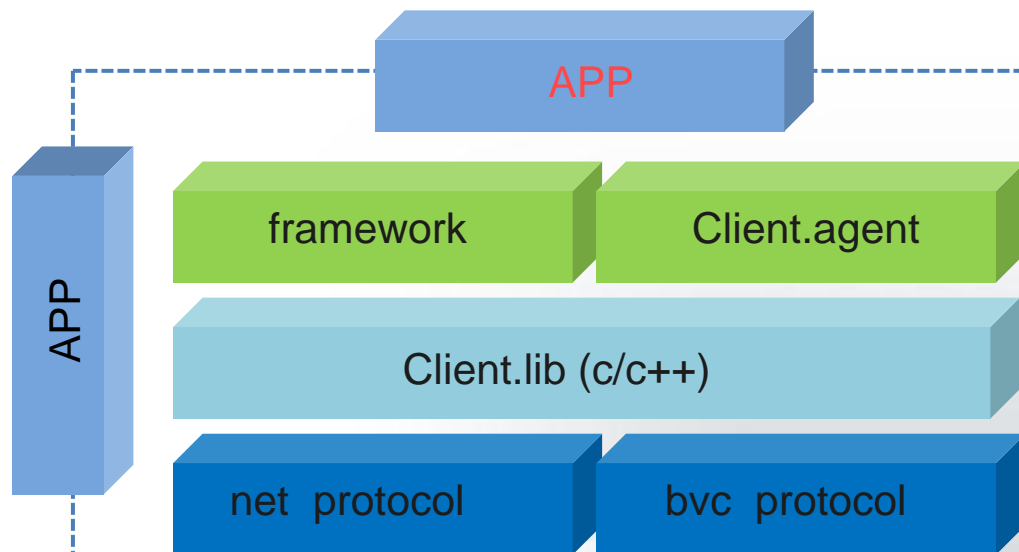


BVC 风险控制 vs 利用率

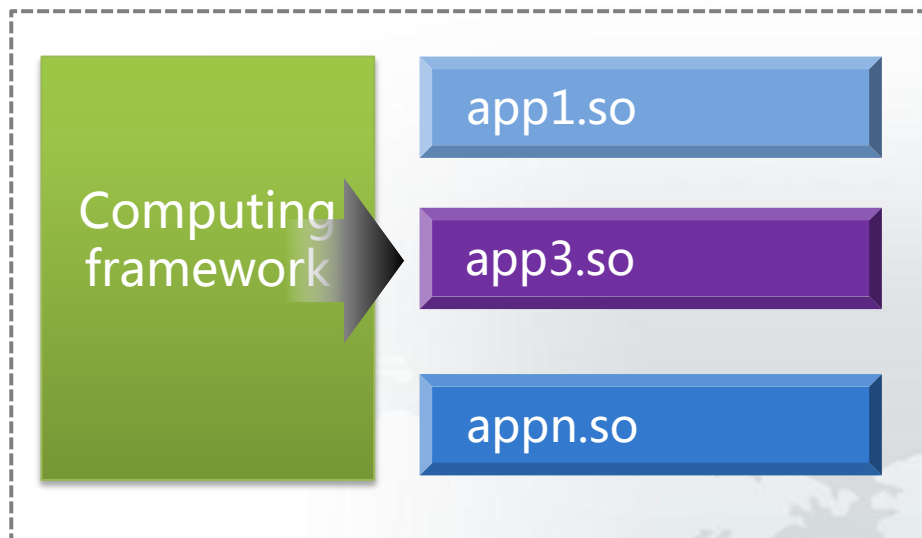
任务启停



BVC 应用接入体系



BVC 应用CF接入



BVC的实际应用

- ✓ 部分小规模调研完全由BVC支持
- ✓ 主要业务之一的预算的 20%由BVC支撑
- ✓ 2011全年6月的计算能力峰值将达到1PFLOPS
- ✓ 2012 :
 - 80%的网页搜索离线集群会加入bvc系统
 - 提供离线预算的30~40%



经验 教训

- ✓ 集群整合是成长过程中的痛苦
- ✓ Volunteer computing System 是很好的选择
- ✓ 优先级控制，风险控制，简明易用的接口



Q&A



Thank you !



The QCon logo features a large, stylized letter 'Q' in a vibrant green color, followed by the letters 'Con' in a bold, blue sans-serif font. The background of the entire slide is a photograph of a traditional Chinese stone pagoda with multiple tiers and circular openings, situated in a body of water with a hazy cityscape in the distance.

QCon

杭州站 · 2011年10月20日~22日

www.qconhangzhou.com (6月启动)

QCon北京站官方网站和资料下载

www.qconbeijing.com

全球企业开发大会

THE ANNUAL
INTERNATIONAL
SOFTWARE DEVELOPMENT
CONFERENCE