

互联网环境下海量数据的处理与应用

周 宇

DTCC2014 长沙 湖南金恩科技有限公司

关于我

简历

- 周宇，湖南长沙人
- 06~09 北京航空航天大学 计算机专业 硕士研究生
- 09~10 微软亚洲工程院 SDE
- 10~12 百度 SE -> Team Leader
- 12~Now 湖南金恩科技(Keen Tech)

领域

- 专注于系统/网络级解决方案
- 专注于Web相关技术

目录：

从互联网的数据说起...

- 数据的价值与面临的问题
- 思路和解决方案

目录：

从互联网的数据说起...

- 数据的价值与面临的问题
- 思路和解决方案

数据的价值与面临的问题

- 互联网有哪些数据



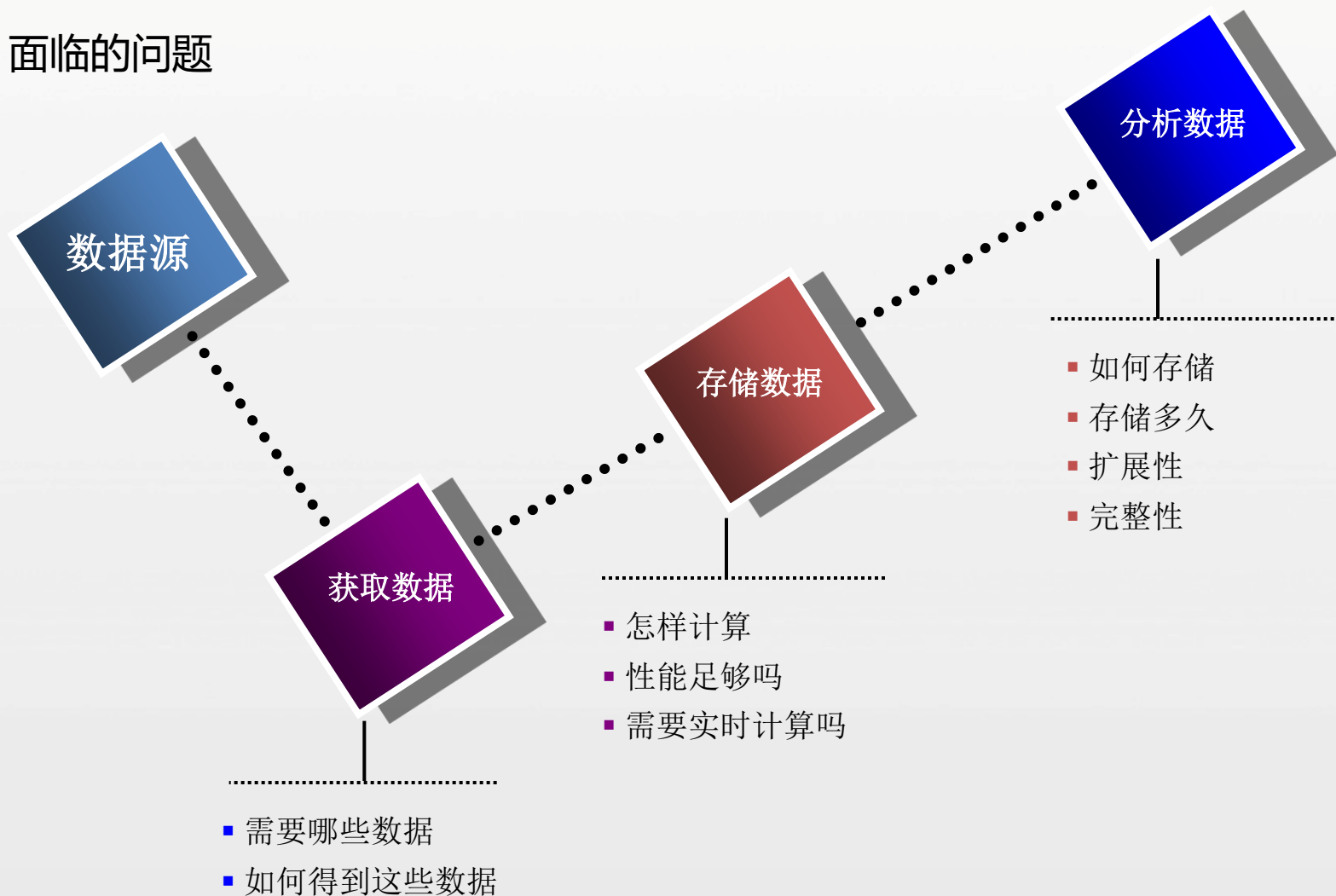
1. 数据是事实

2. 数据支撑决策

3. 数据决定成败

数据的价值与面临的问题

- 面临的问题



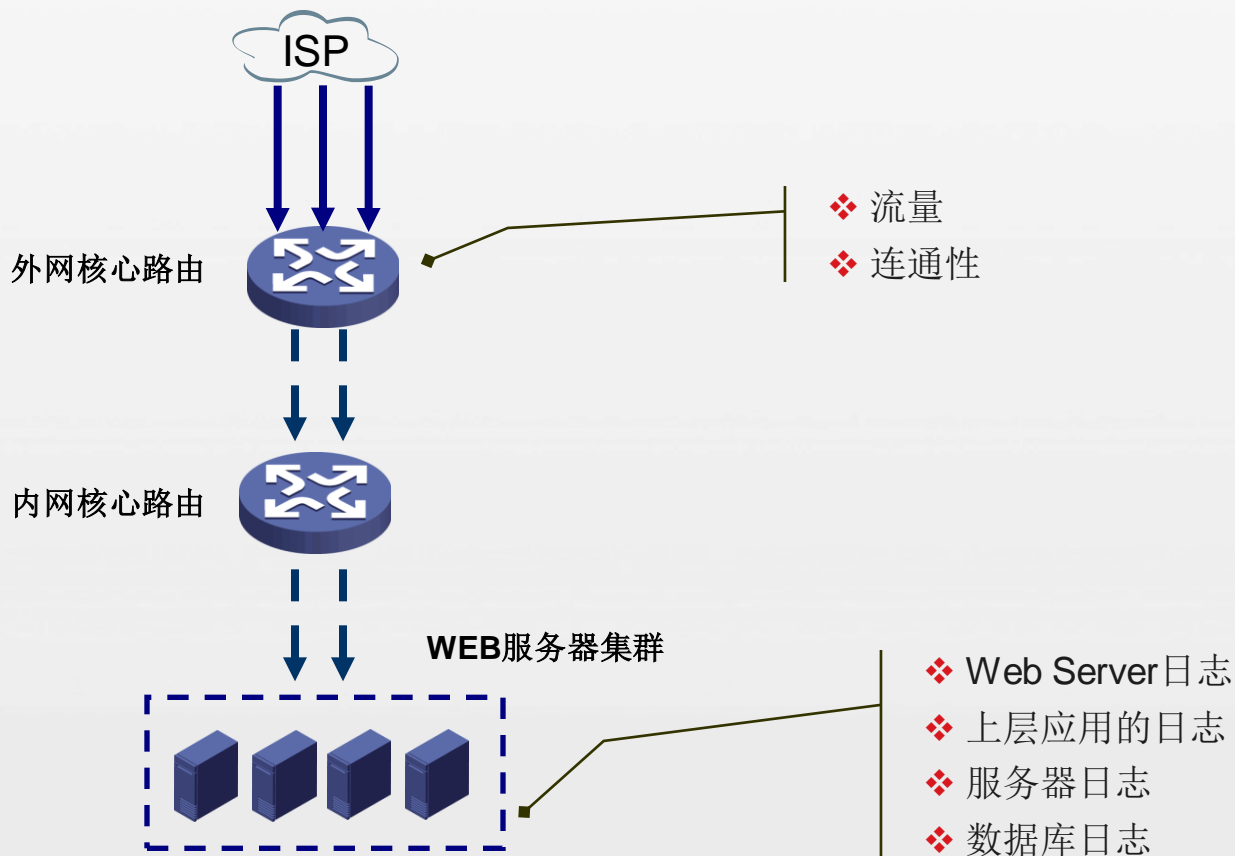
目录：

从互联网的数据说起...

- 数据的价值与面临的问题
- 思路 and 解决方案
 - 数据获取
 - 数据存储
 - 数据分析

思路和解决方案—数据获取

- 互联网企业的典型架构

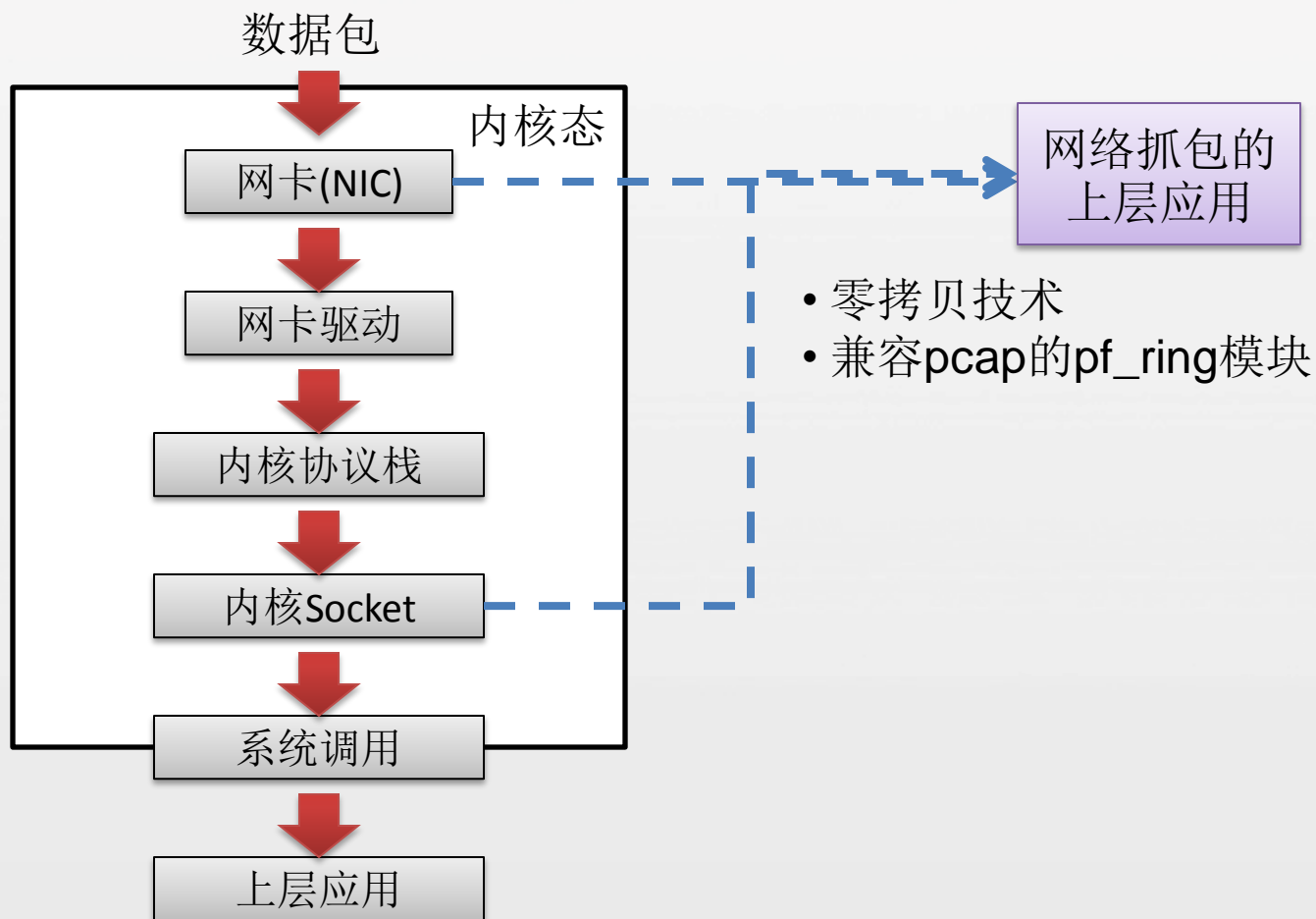


思路和解决方案—数据获取

- 两个问题
 - 日志占用了紧缺的服务器资源
 - 日志类型繁杂，容易重复造轮子

思路和解决方案—数据获取

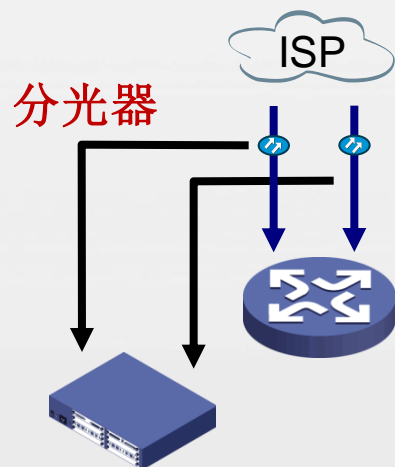
- 将资源压榨到底——网络抓包的例子



- 结束了吗
 - 1Gbps Vs 10Gbps
 - 多队列技术
 - 硬件Hash
 - 延迟中断技术
- 可以做得更好吗

思路和解决方案—数据获取

- 通过交换机的端口镜像获取数据
 - 释放了服务器资源
 - 转发能力有限
- 我们都是光纤接入



1Gbps



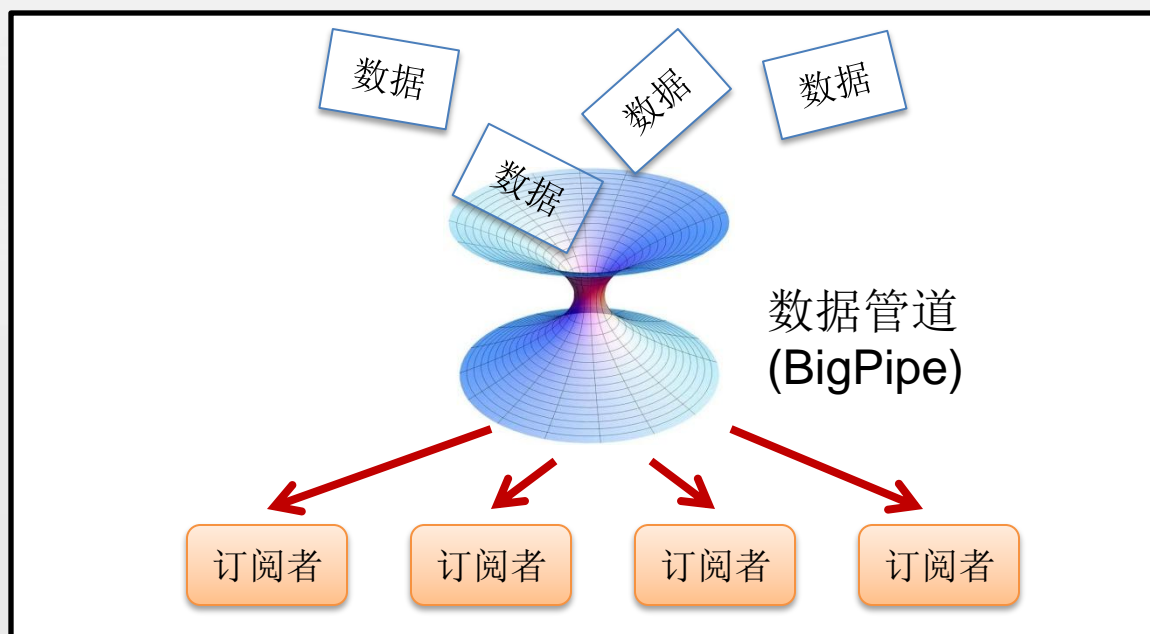
10Gbps



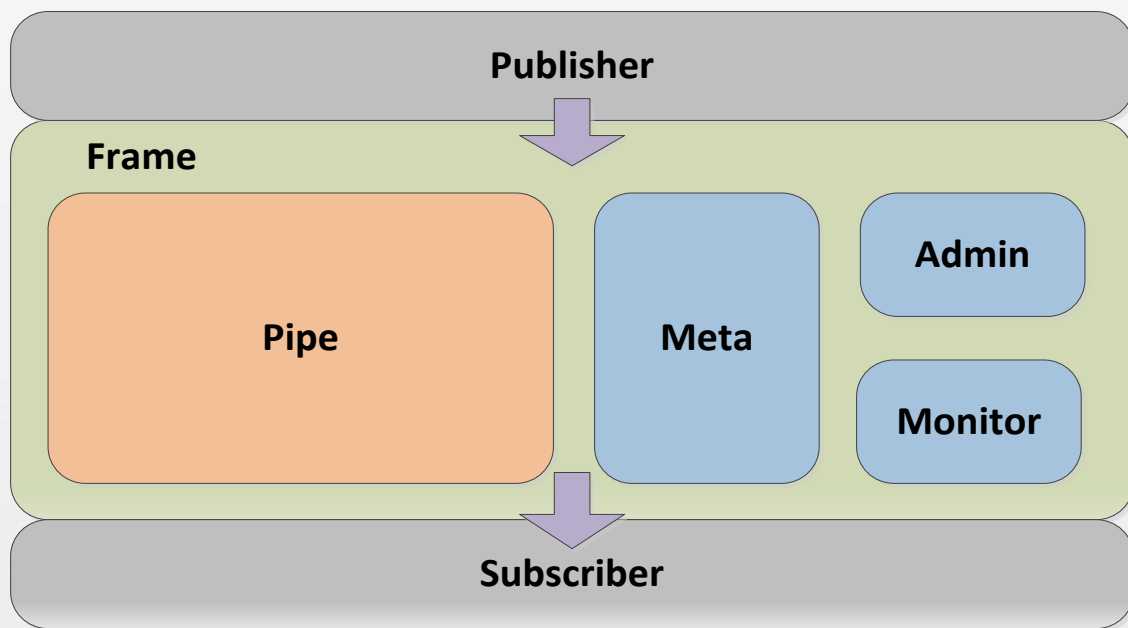
100Gbps

思路和解决方案—数据获取

- 谁来分析？
 - 每个团队都需要
 - 每个团队的想象力都很丰富
- 解耦数据提供方与数据获取方



- BigPipe架构



- 发布/订阅模式

- 解决数据传输/存储/压缩的问题

目录：

从互联网的数据说起...

- 数据的价值与面临的问题
- 思路 and 解决方案
 - 数据获取
 - 数据存储
 - 数据分析

- 存储时间
 - 计算需要
 - 政策需要
- 存储需要考虑什么
 - 完整性
 - 读写性能

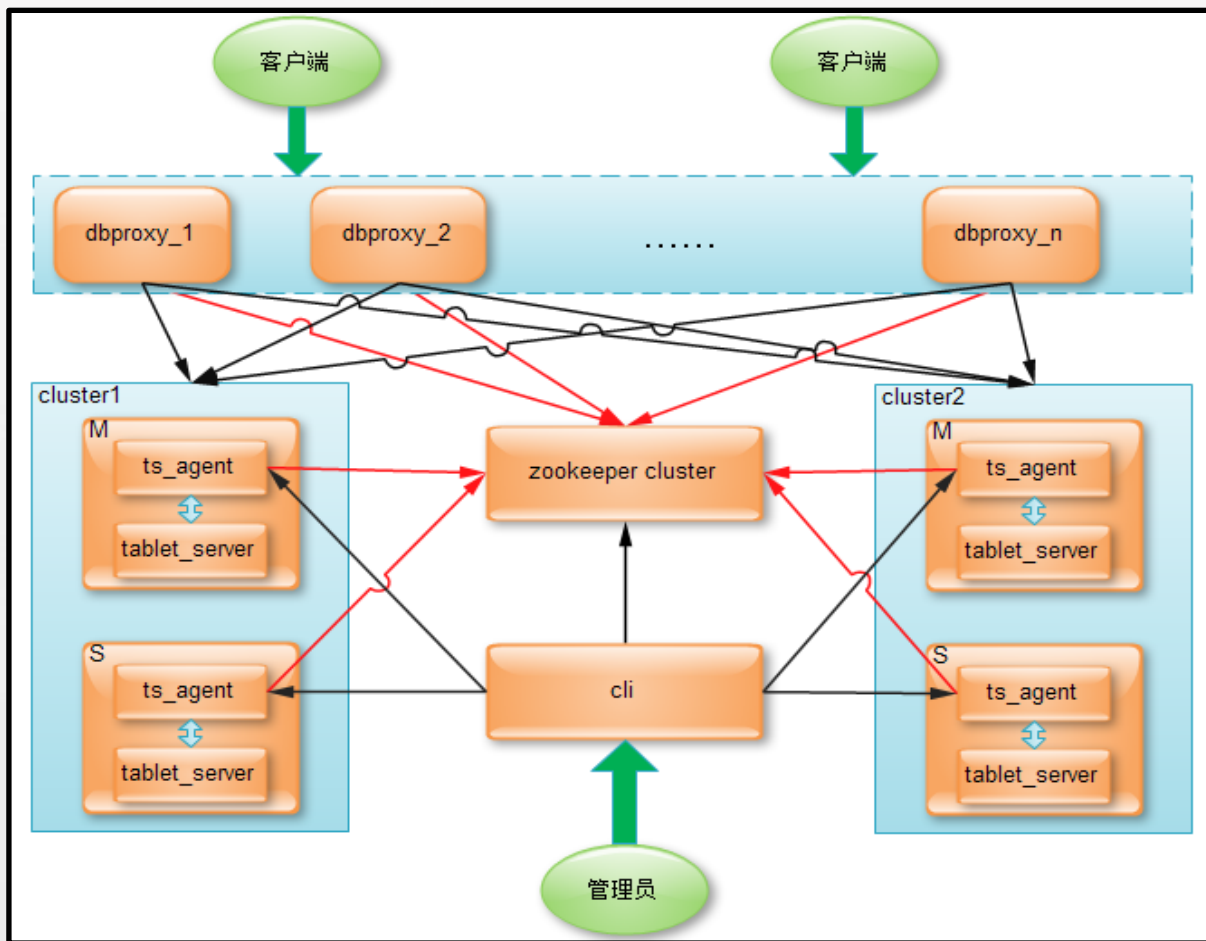
- HADOOP
 - Map Reduce
 - HDFS
- 模型简洁
- 互联网数据存储的事实标准

思路和解决方案—数据存储

- MySQL(RDBMS)
 - dbproxy
 - 分布式集群

思路和解决方案—数据存储

- DDBS架构



目录：

从互联网的数据说起...

- 数据的价值与面临的问题
- 思路 and 解决方案
 - 数据获取
 - 数据存储
 - 数据分析

- Map Reduce
 - 每个人都可以计算
 - 人人都是科学家
 - 实时性问题
- 实时计算

- 程序员喜欢简单
 - 单机能满足吗
 - 大型机？
- 后Moore时代到来
 - 高主频
 - 多核
 - 定制化

- 04年开始，多核是处理器发展的主要方向
 - Intel sandy bridge 8 core
 - Tiler 64 core
 - Cavium 12 core
- 定制化
 - Nehalem : SSE4.2
 - Sandybridge : 视频加速 (not only graphic)
 - Haswell : 硬件随机数 , transactional memory

- 基于多核处理器的单机处理能力
 - 10Gbps线速转发
 - 接近10Gbps的7层分析能力
- 多核处理器的编程有何不同
 - 熟悉体系结构
 - 设计Lock Free的数据结构
 - 关注Cache Miss和TLB Miss
 - 善用Profiling

- 服务器还是x86的天下吗
 - 应用和机器规模的高速膨胀
 - 软件不再去适配通用硬件架构，而是硬件通过定制化来适应软件
- ARM的崛起
 - 低功耗
 - A9，A15具备高端处理器特征
- 其他
 - SSD
 - 硬件压缩/解压缩

- 一个例子——实时异常URL分析

- 20Gbps的流量
- 实时分析
- 什么是异常URL：由于人为，机器或者程序原因，造成同一个URL在短时间内被大量访问，这样的URL被认为是异常的

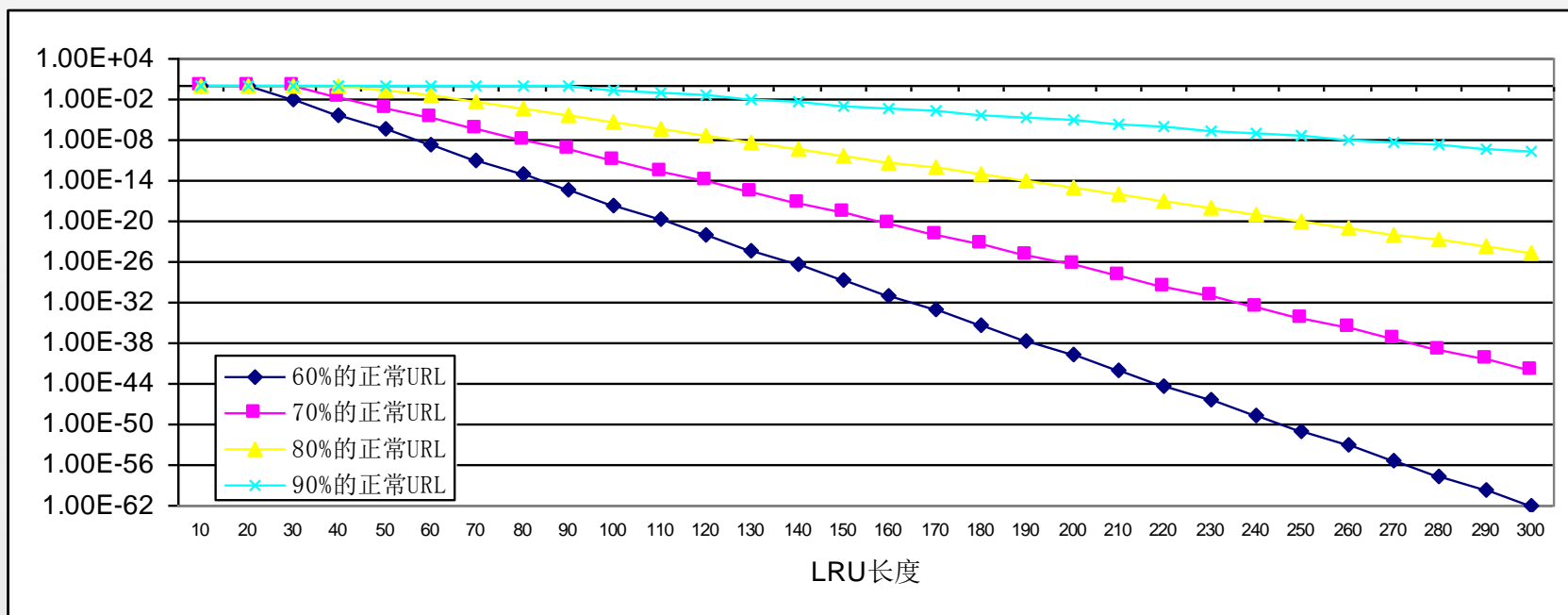
- 一个例子——实时异常URL分析

- 如何高效提取URL(Cycle级别的优化)
- 如何保持好的局部性
- 如何控制内存使用量
- 基于LRU的统计算法可行吗

思路和解决方案—数据分析

• 一个例子——实时异常URL分析

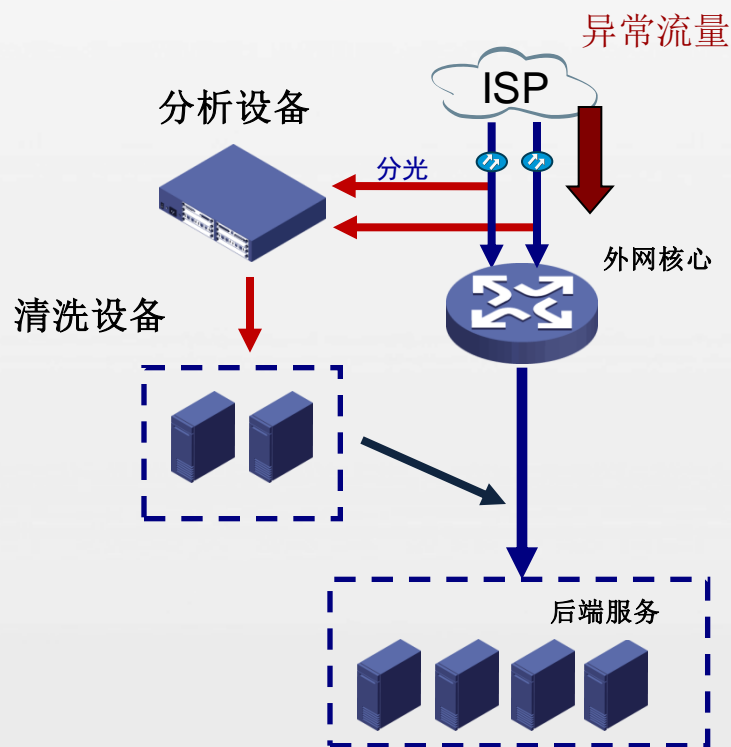
— 是时候使用一些数学工具了



思路和解决方案—数据分析

• 一个例子——实时异常URL分析

- 发现了怎么办
- 通知服务器做处理
- 网络层面做处理



谢谢