



FTF 2016
TECHNOLOGY FORUM CHINA

基于ARM的QorIQ SoC上的KVM可 铸就卓越的虚拟化性能

FTF-DES-N1887

DIANA CRĂCIUN
软件工程师
2016年5月16日



软件产品和服务

开发工具

- CodeWarrior

运行时产品

- VortiQa软件解决方案

CodeWarrior
QorIQ

VortiQa



集成服务

- 安全咨询
- 强化Linux

解决方案参考

- 物联网网关
- OpenWRT+

Linux®服务

- 商业支持

- 性能调谐



加快客户产品上市时间



提供商用软件、支持、服务和解决方案



简化与恩智浦的软件合作



创造成功!



引领网络中的64位ARM浪潮



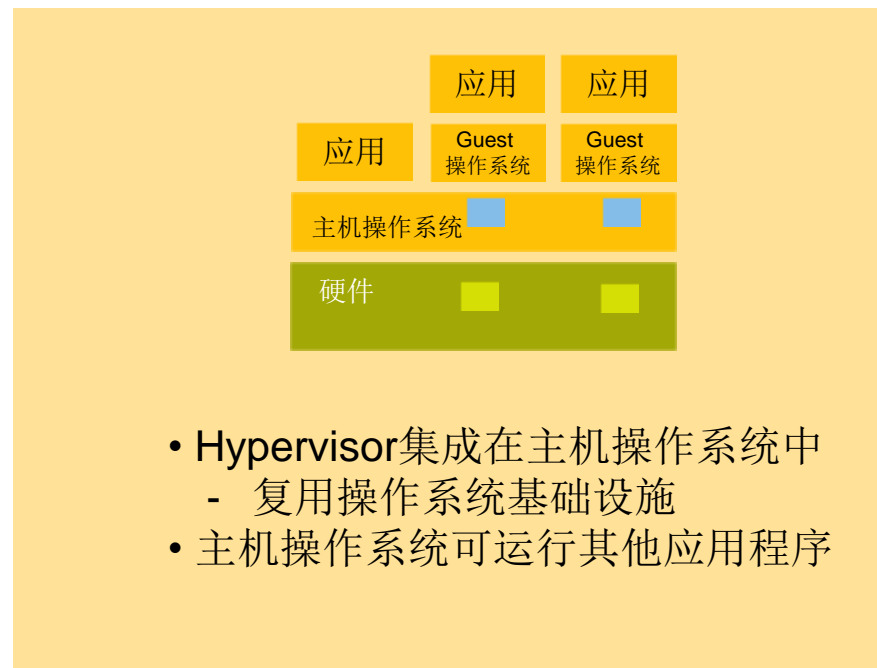
议程

- 简介
- KVM/QEMU
- ARM虚拟化扩展
- ARM上的KVM
- 结果
- 结论

简介

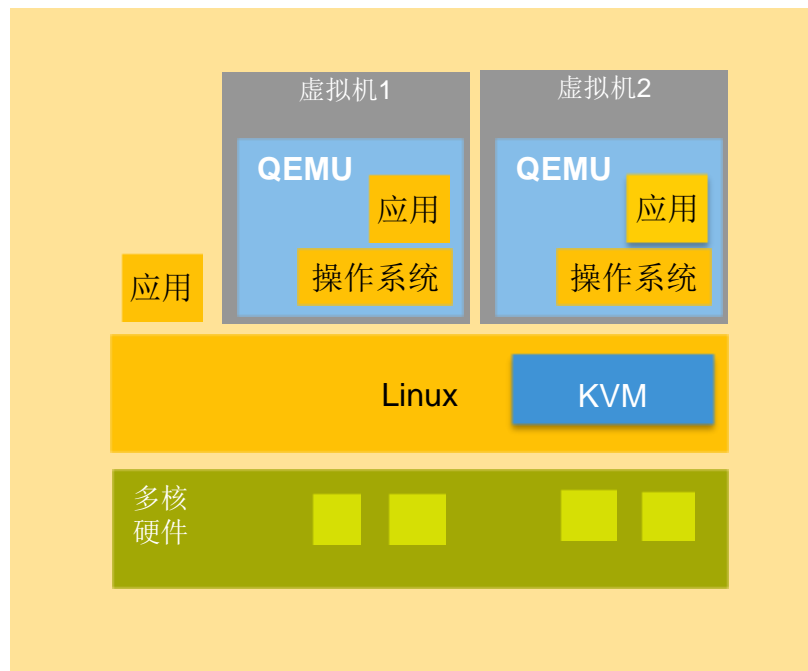
虚拟化和Hypervisor

- 虚拟化 – 能够提供一个抽象层以在单个计算机系统上运行多个操作系统的硬件和软件技术
- **Hypervisor** 是一种可创建和管理能够运行Guest操作系统的虚拟机的软件组件



KVM/QEMU

KVM/QEMU

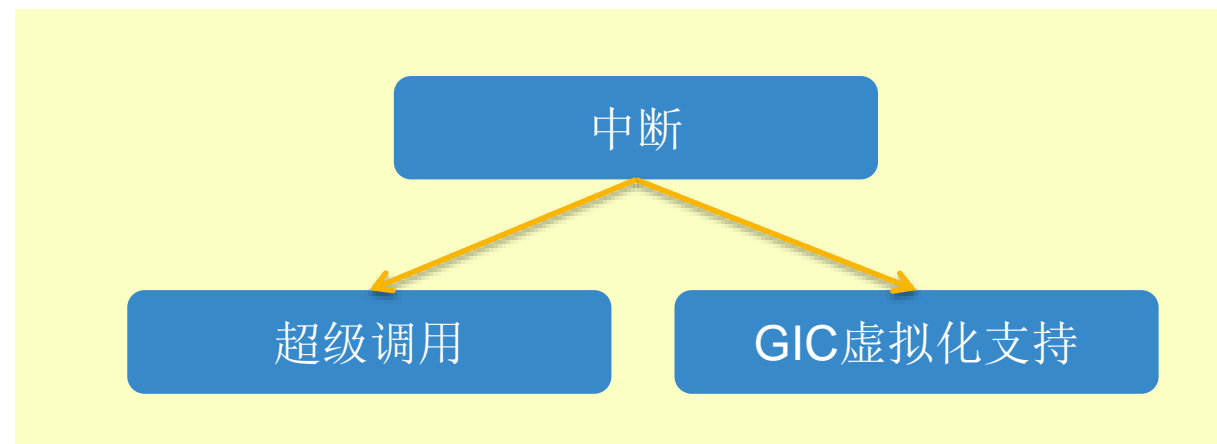
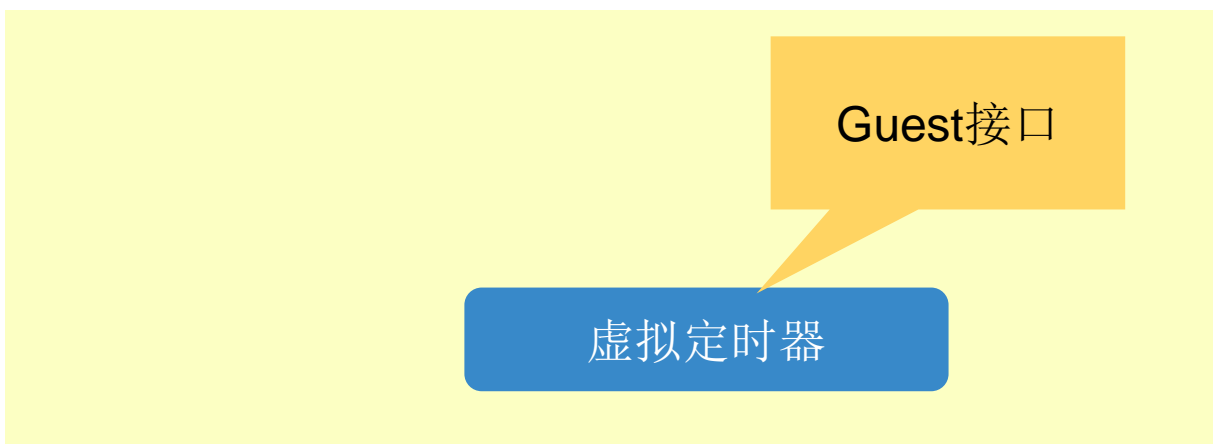
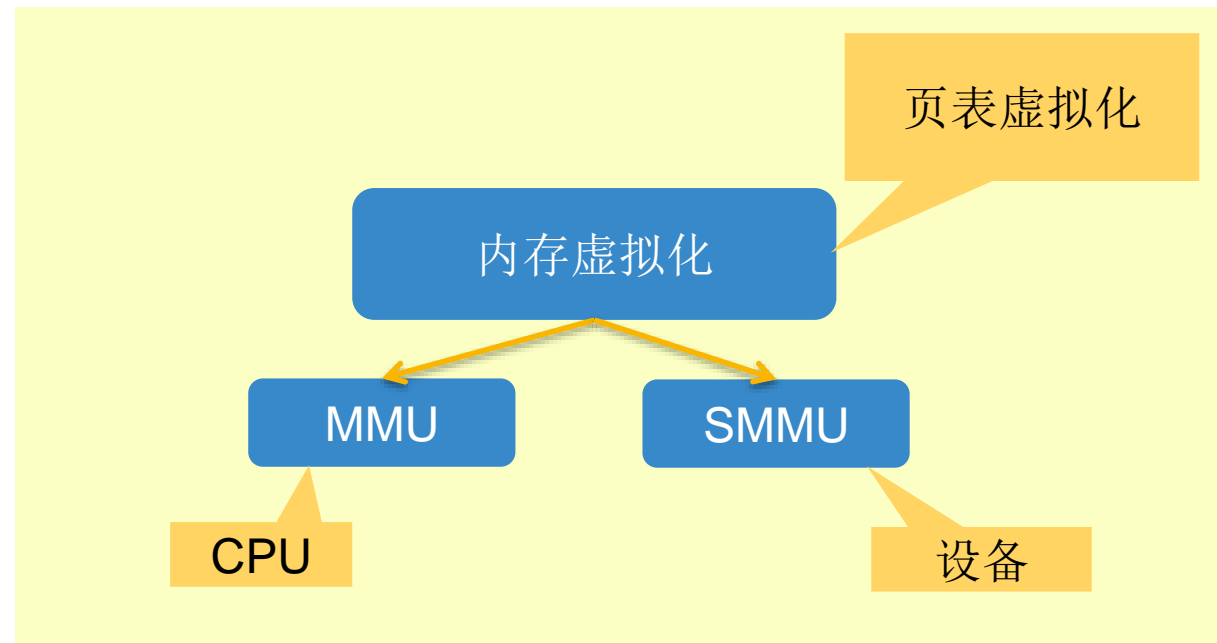
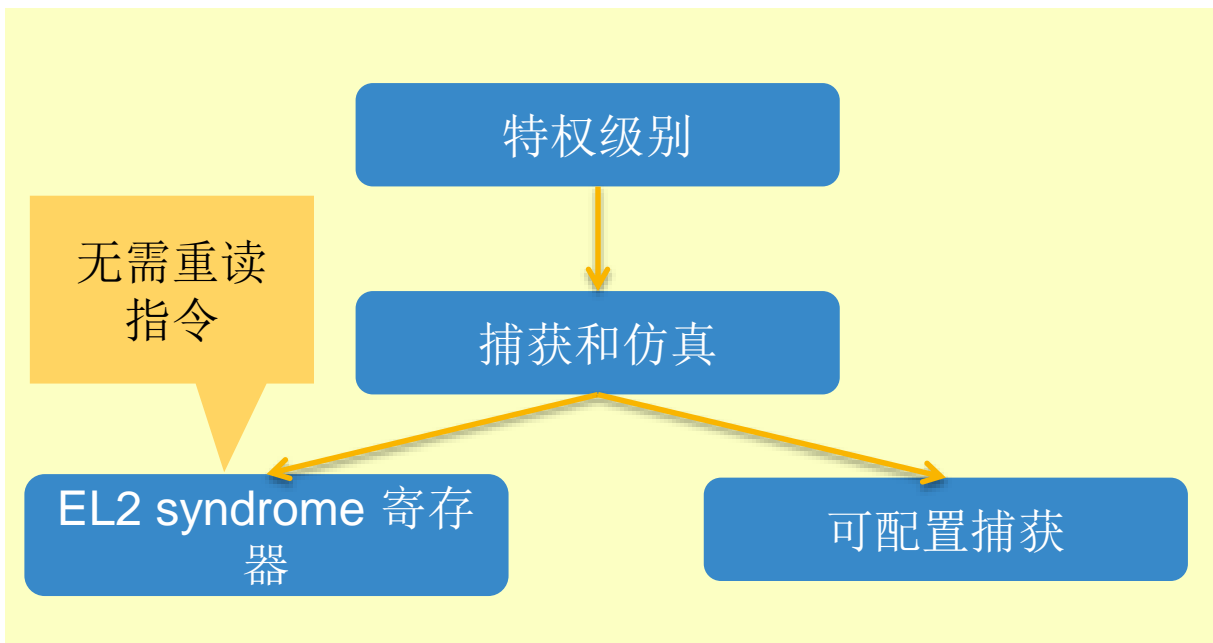


- KVM/QEMU – 基于Linux®内核的开源虚拟化技术
- KVM是一个Linux内核模块
- QEMU是一个使用KVM进行加速的用户空间仿真器
- 同时运行虚拟机和Linux应用程序
- 操作系统无需更改或仅需小幅更改
- 虚拟I/O功能
- Direct / pass through I/O – 为VM分配I/O设备

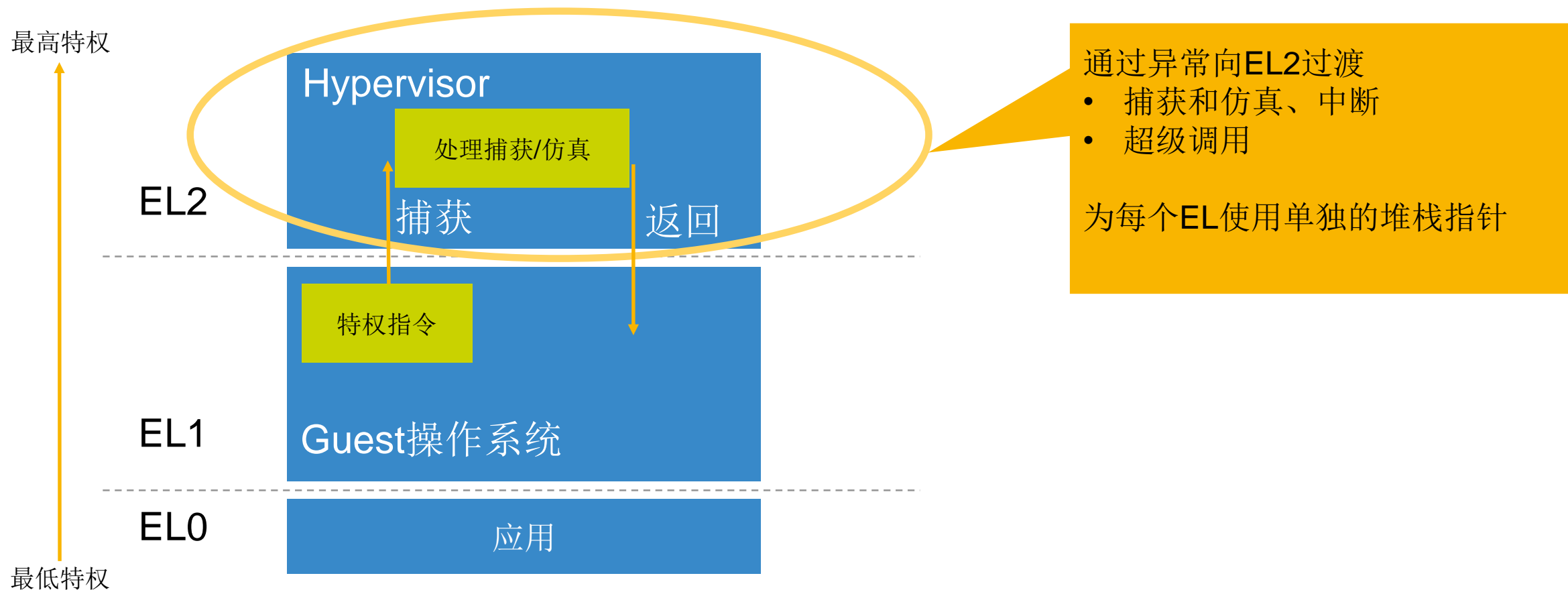
ARM虚拟化扩展



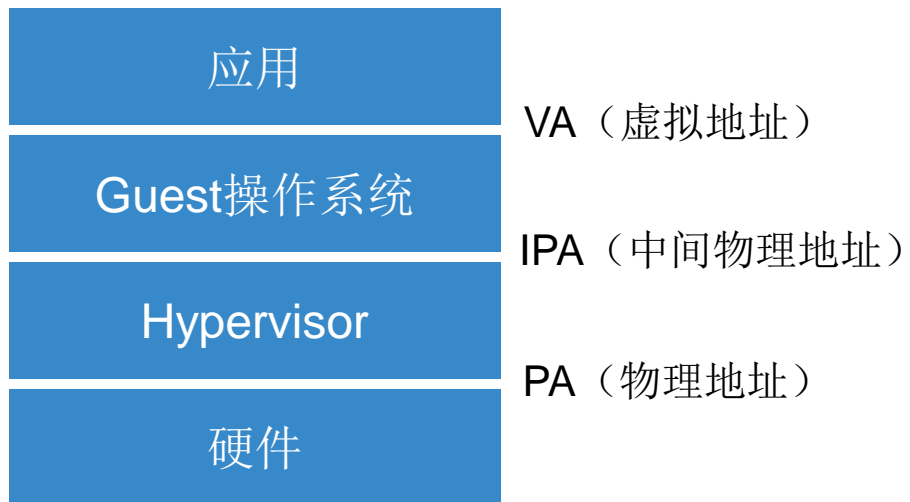
ARM虚拟化扩展



特权级别



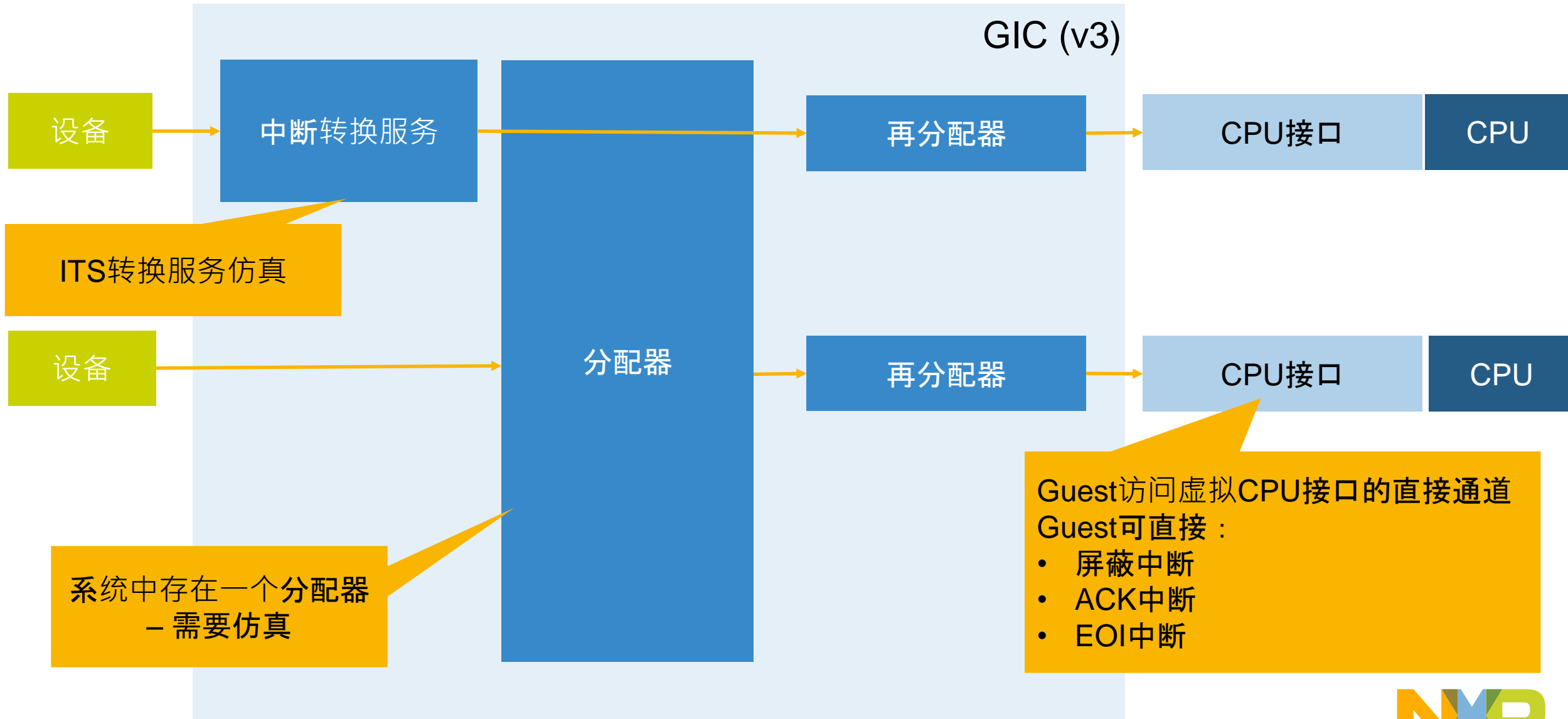
内存虚拟化



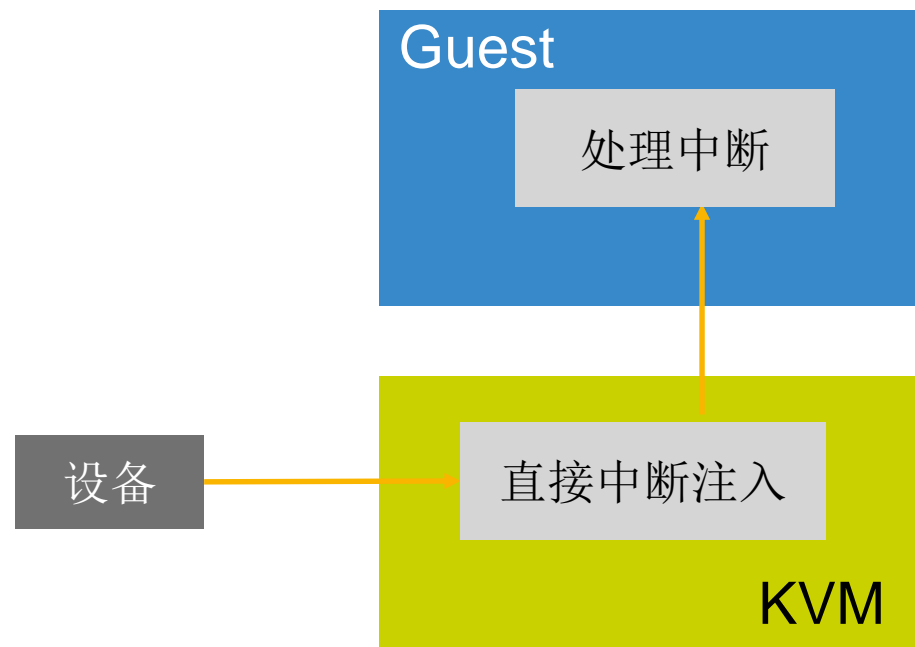
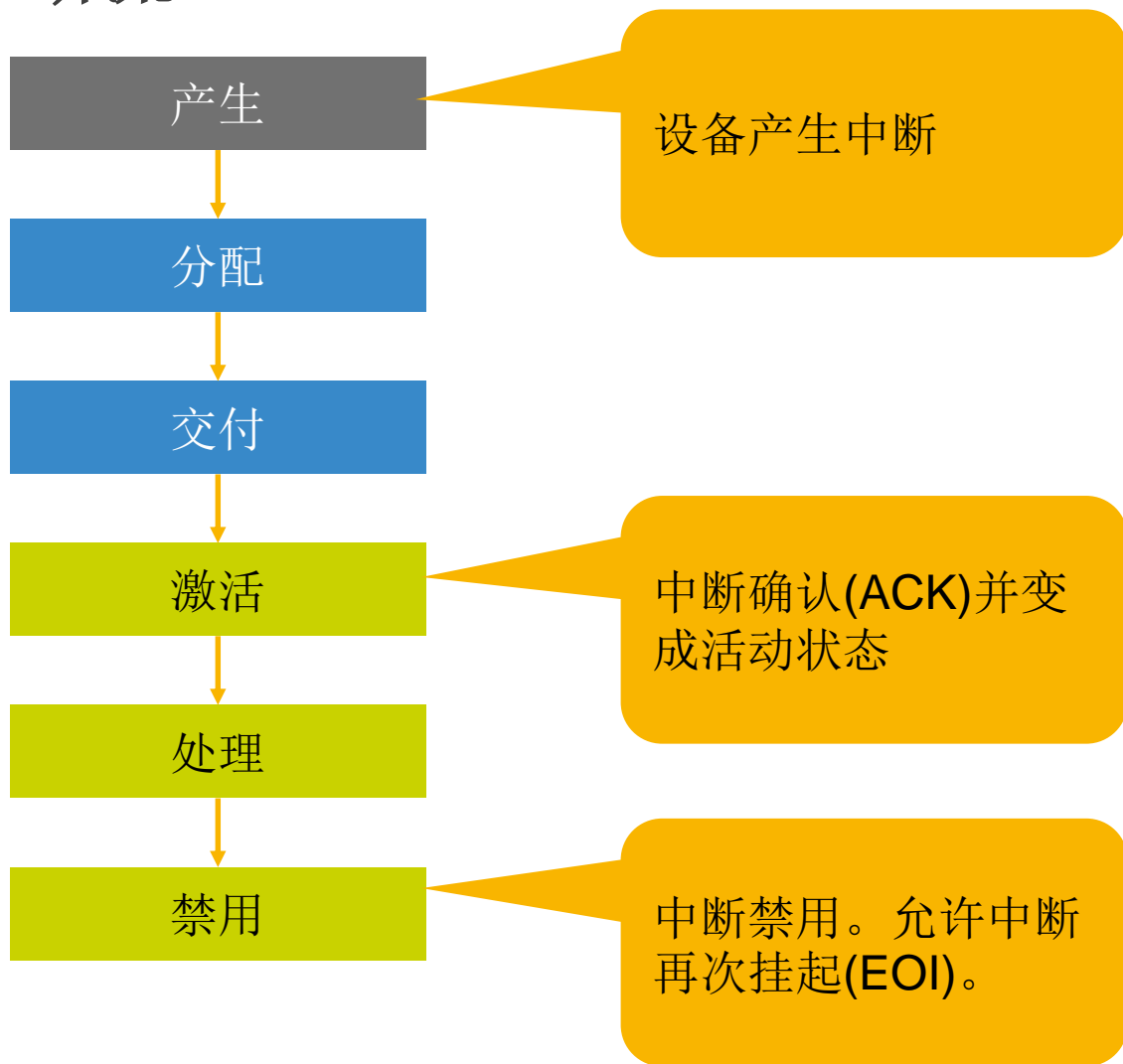
- 2级内存转换
- 首次转换将内存从VA转换至IPA
 - 由Guest所有
- 二级转换从IPA转换至PA
 - 表项由Hypervisor维护



中断虚拟化

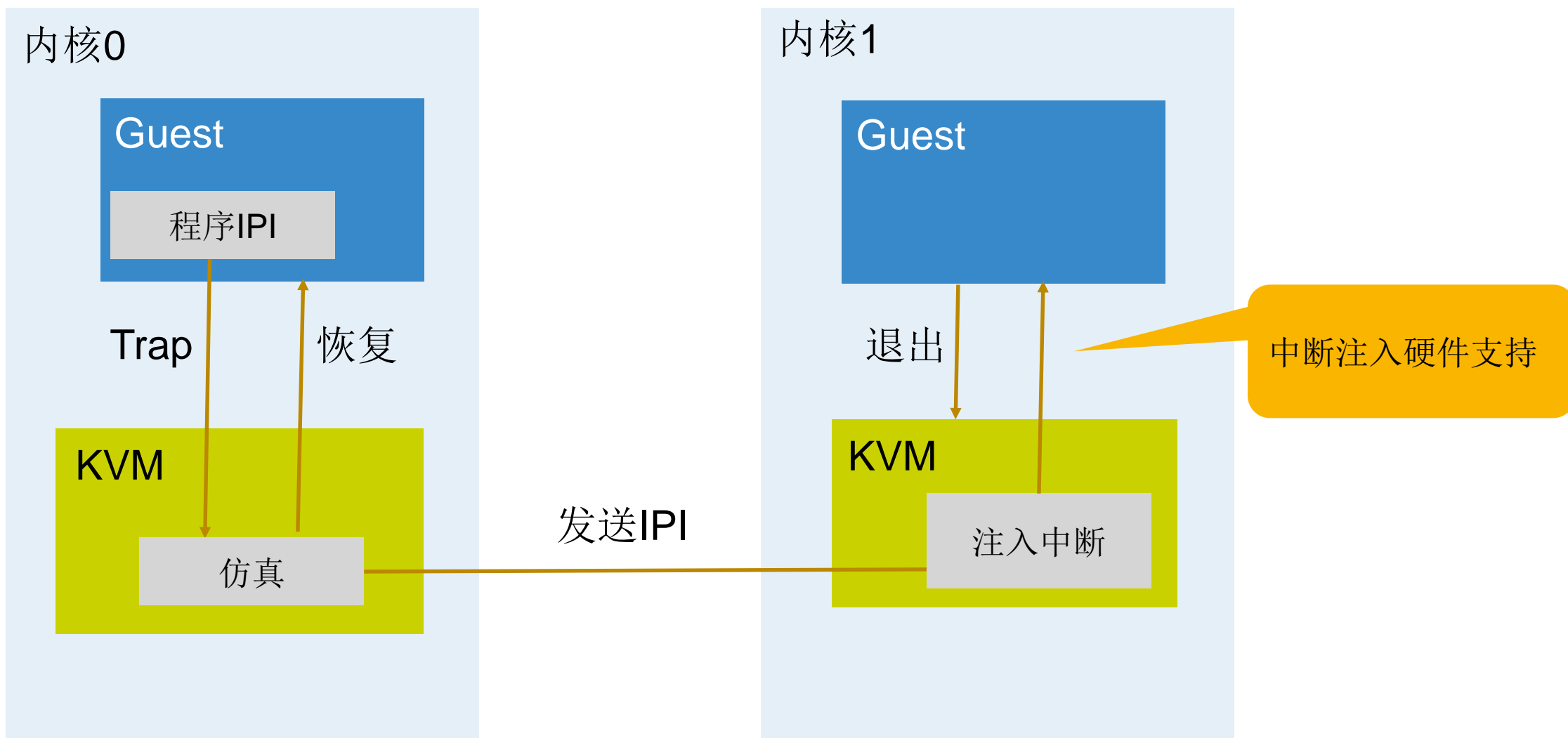


中断流



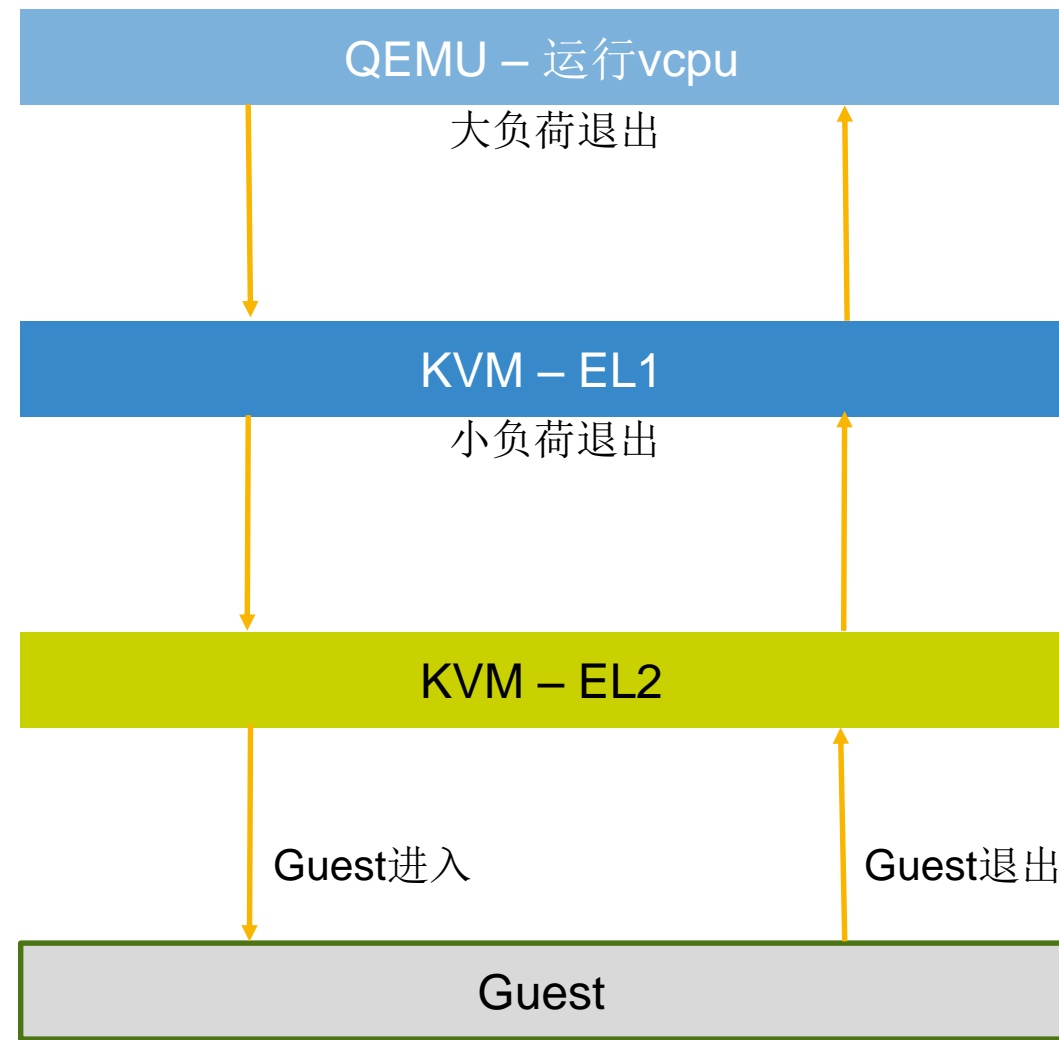
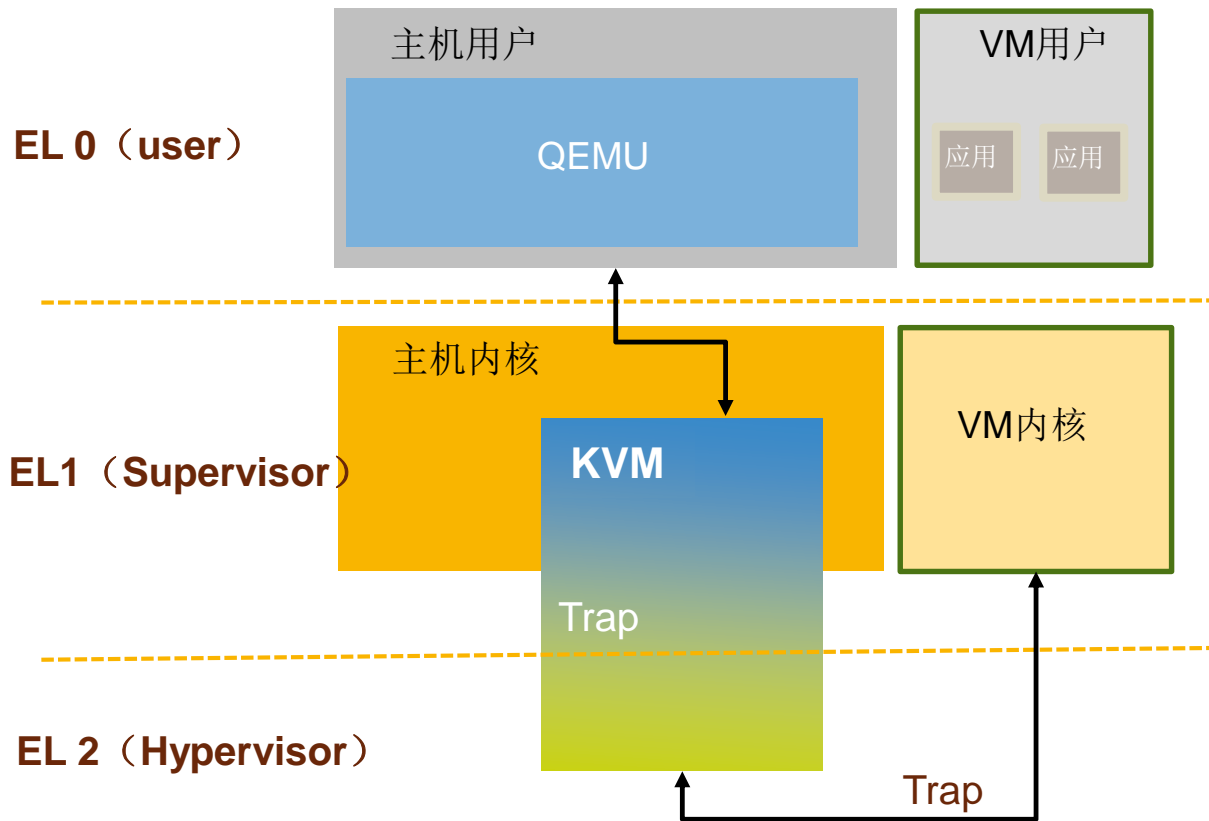
- 接收到中断后**Guest**被中断
- 硬件支持中断注入
- 某些情况下可能存在额外退出（但应该很少）。

IPI流



ARM上的KVM

ARM上的KVM/QEMU



结果

开销来源

- 虚拟化可能附带成本：开销
- 但是在具有硬件扩展时是什么导致开销呢？
- 因Guest退出导致开销
 - 陷阱 (Trap)、中断
- Guest速度
 - 内存转换步骤更多
 - TLB/缓存污染/竞争
 - 锁竞争
- 应用延迟
 - 延迟敏感型应用在虚拟化环境中的工作特性可能不同

Guest退出 – 示例

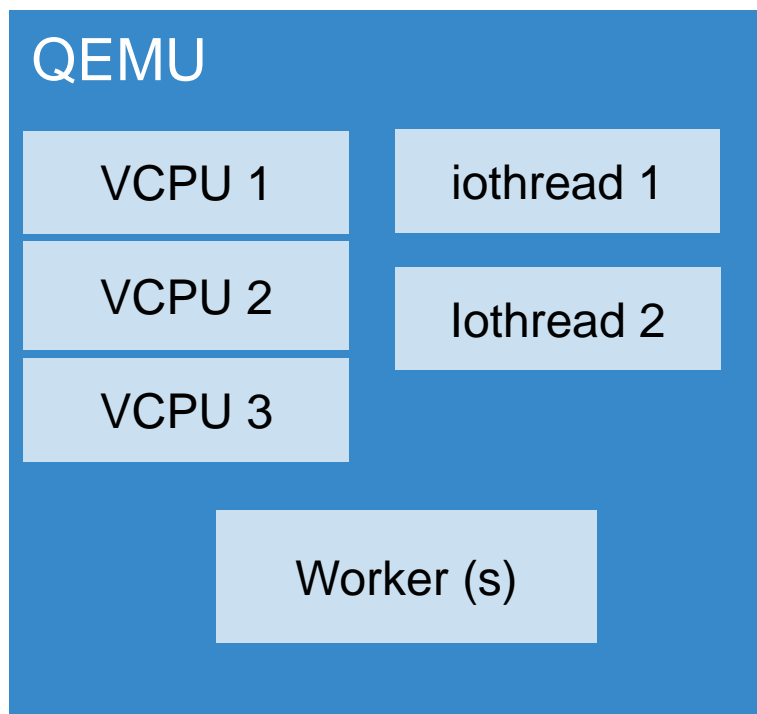
- 退出定时框架

- 为各种退出类型报告了在Hypervisor中花费的时间

类型	计数	最小值 (cycle)	最大值 (cycle)	平均值 (cycle)	合计 (cycle)	标准差 (cycle)	计数	合计
WFX	0	0	0	0	0	0	0	0
CPC15_32	0	0	0	0	0	0	0	0
CPC15_64	0	0	0	0	0	0	0	0
CP14_MR	0	0	0	0	0	0	0	0
CP14_LS	0	0	0	0	0	0	0	0
CP14_64	0	0	0	0	0	0	0	0
HVC32	0	0	0	0	0	0	0	0
SMC32	0	0	0	0	0	0	0	0
HVC64	0	0	0	0	0	0	0	0
SMC64	0	0	0	0	0	0	0	0
SYS64	0	0	0	0	0	0	0	0
IABT_LOW	0	0	0	0	0	0	0	0
DABT_LOW	0	0	0	0	0	0	0	0
DABT_IO_MEM	0	0	0	0	0	0	0	0
DABT_USER_MEM	0	0	0	0	0	0	0	0
DABT_IO_MEM_IPI	157225	10090	57218	13385	2.104.548.054	172	15722,5	210.454.805,40
INTERRUPT	159395	4963	39654	6792	1.082.746.418	226	15939,5	108.274.641,80
TIMEINGUEST	316620	163	356454	47376	15.000.203.563	1230	31662	1.500.020.356,30
DESCHEDULED	2	7036	7036	7036	14072	0	0,2	1.407,20



KVM Benchmark考虑因素



?

- VM scaling
- 集群 (Cluster)
- QEMU线程亲和性
- CPU scaling
- 空闲/繁忙主机
- 复现性
- 中断亲和性

测试方法和分析工具

- Benchmark
 - Coremark
 - Lmbench
- 分析工具
 - 退出定时测量
 - Perf计数器（硬件计数器）
- 平台
 - LS2080 QorIQ硬件

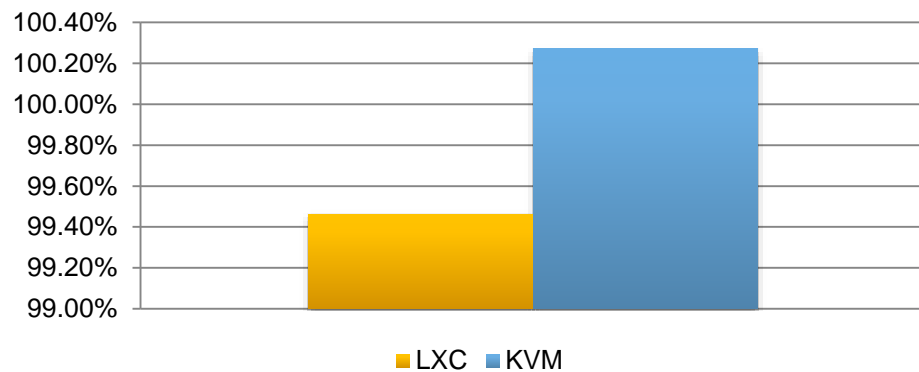
CoreMark

- Microbenchmark
- Core centric

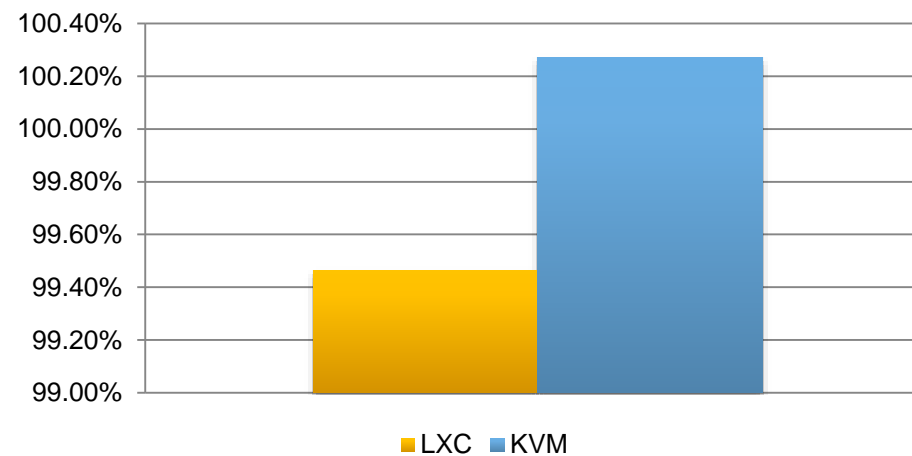
CoreMark – 结果

guest/native[%]

64b CoreMark /MHz - 虚拟化与 本机



64b CoreMark /MHz - 2VM与1VM



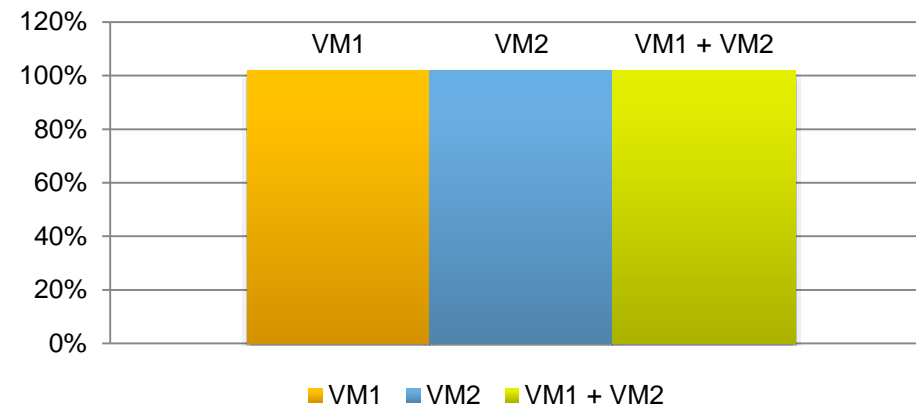
CoreMark – VM的并行性

过量预订 (Oversubscription)

- Host: 2个CoreMark进程在同一CPU上运行
- Guest: 2个VM (VCPU)在同一CPU上运行, 每个均在同一CPU上运行一个CoreMark实例

guest/native[%]

64b CoreMark /MHz – guest vs native



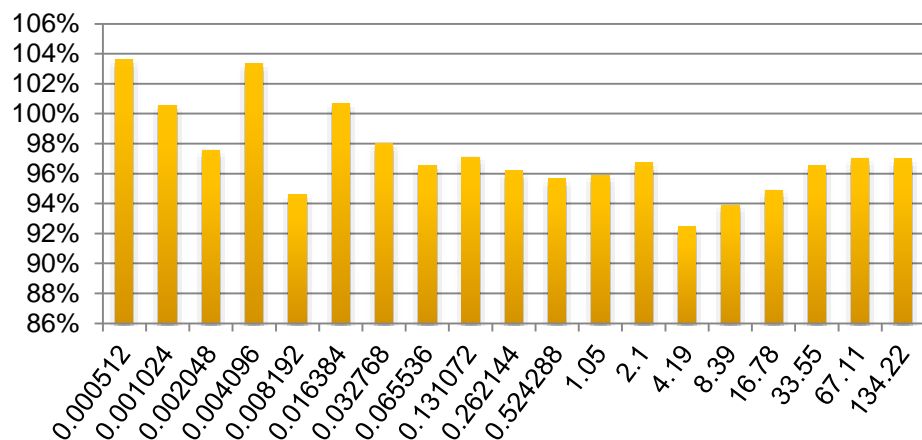
LMbench

- Synthetic microbenchmark
 - 带宽Benchmark
 - 存储器带宽
 - IPC带宽
 - 缓存I/O带宽
- Latency Benchmark
 - 存储器读取
 - 信号处理
 - 进程创建
 - 上下文切换
 - 进程间通信
 - 文件系统

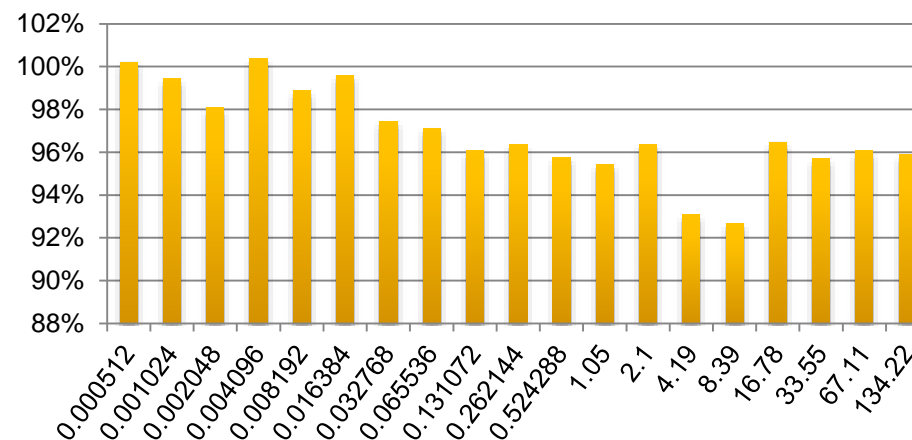
LMBench – 通信带宽

guest/native[%]

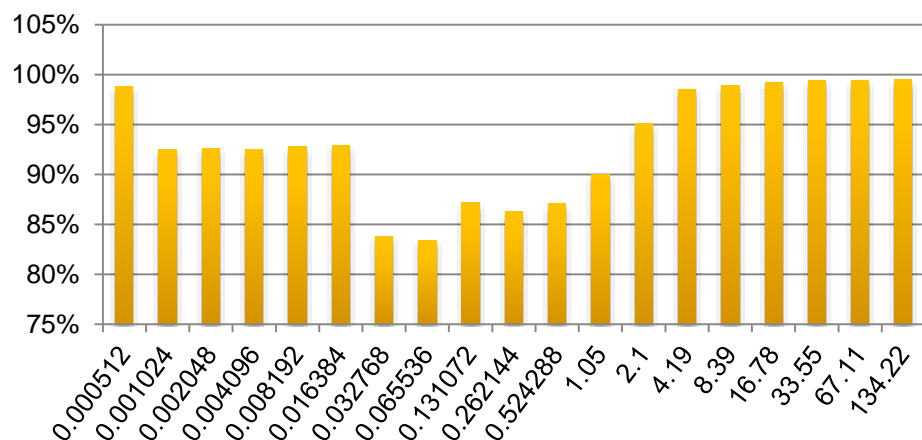
文件读取带宽



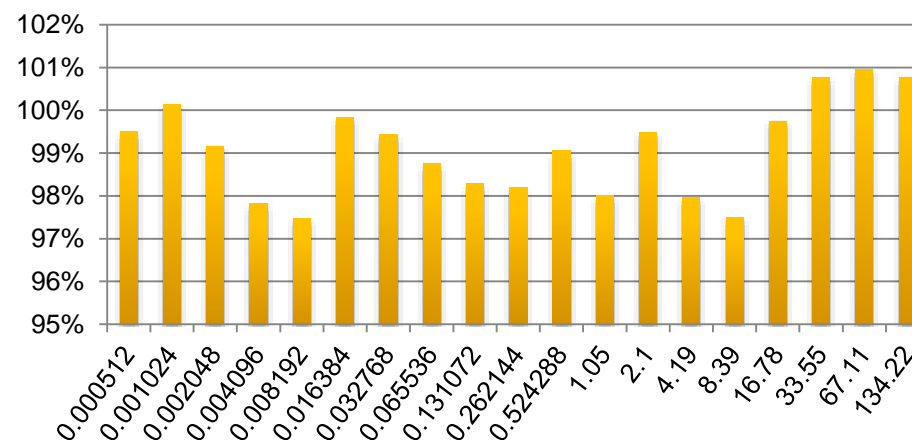
读取open2close带宽



Mmap读取带宽



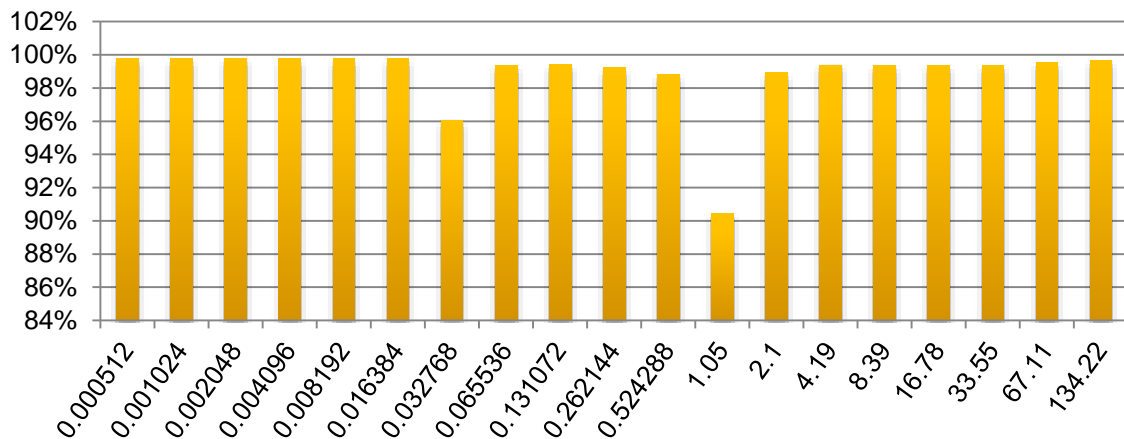
Mmap读取open2close带宽



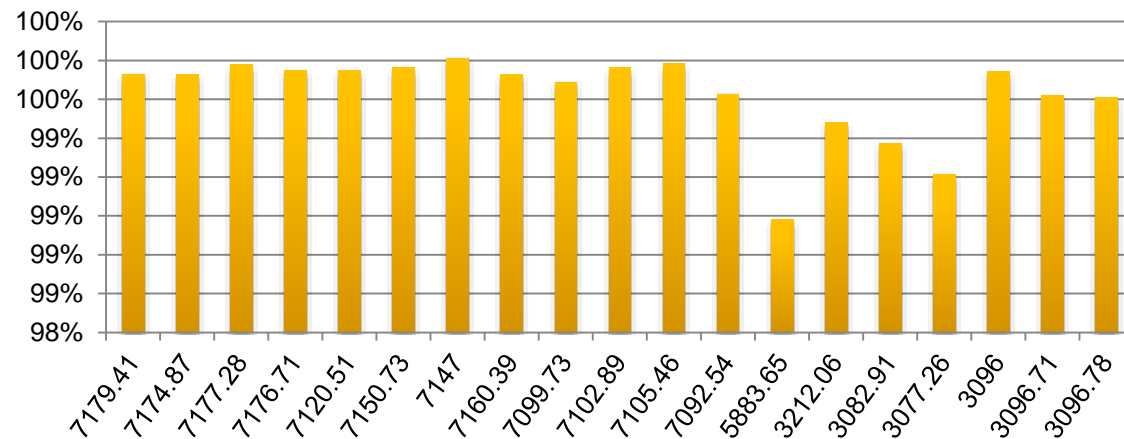
LMBench – 内存带宽

guest/native[%]

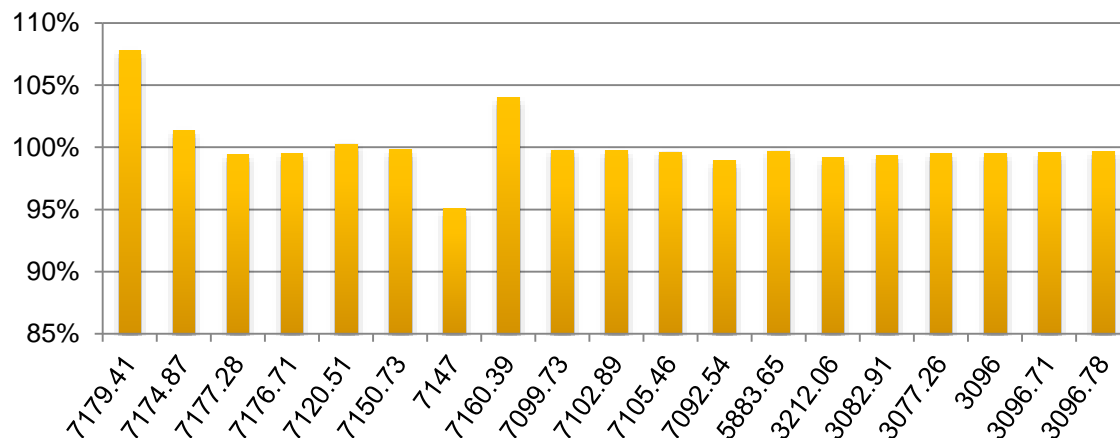
内存读取带宽



内存写入带宽

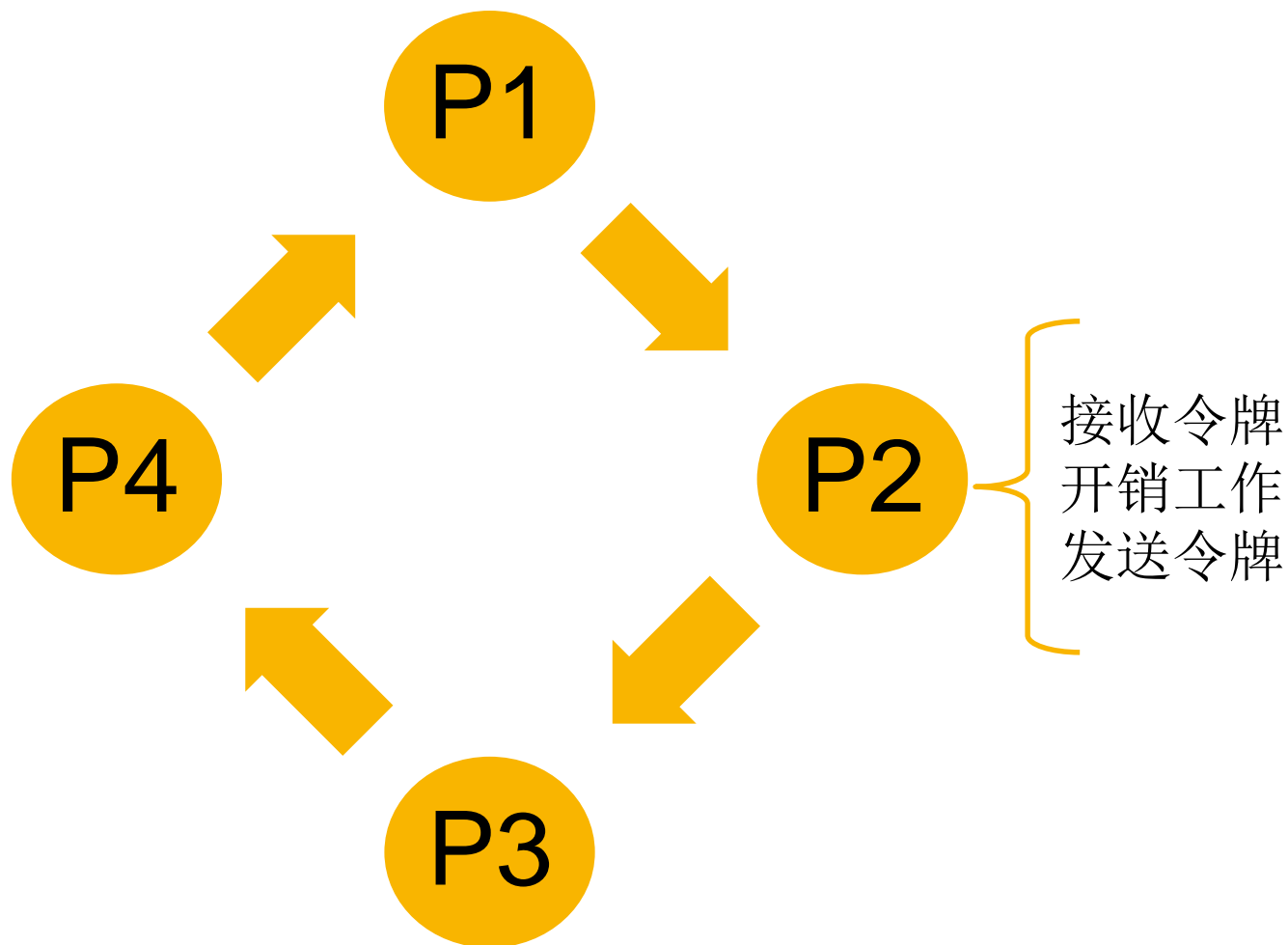


内存部分读取/写入带宽



LMBench – 上下文切换 sub-benchmark

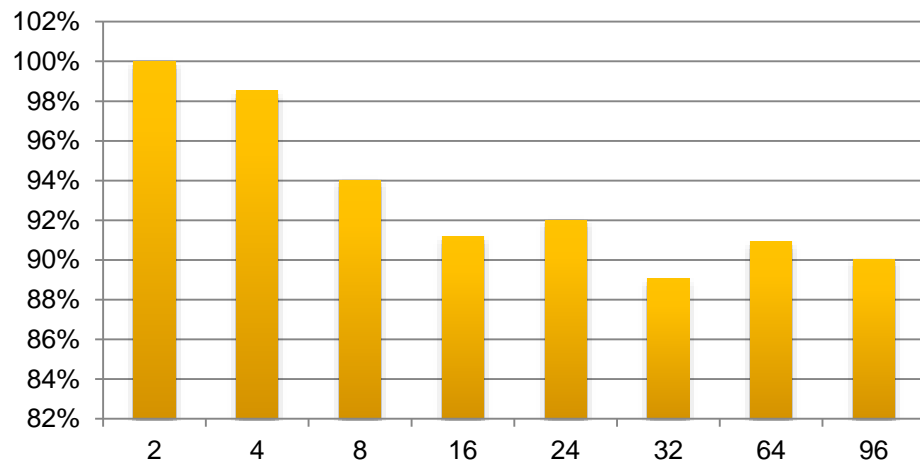
- 平衡
 - 所有进程在同一内核上
- 唯一
 - 每个进程都在不同内核上



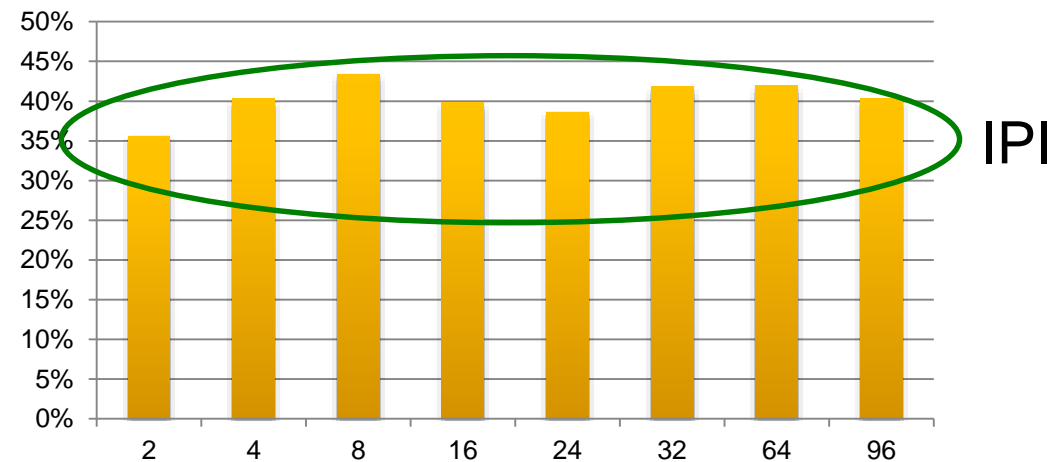
LMBench – 上下文切换延迟

guest/native[%]

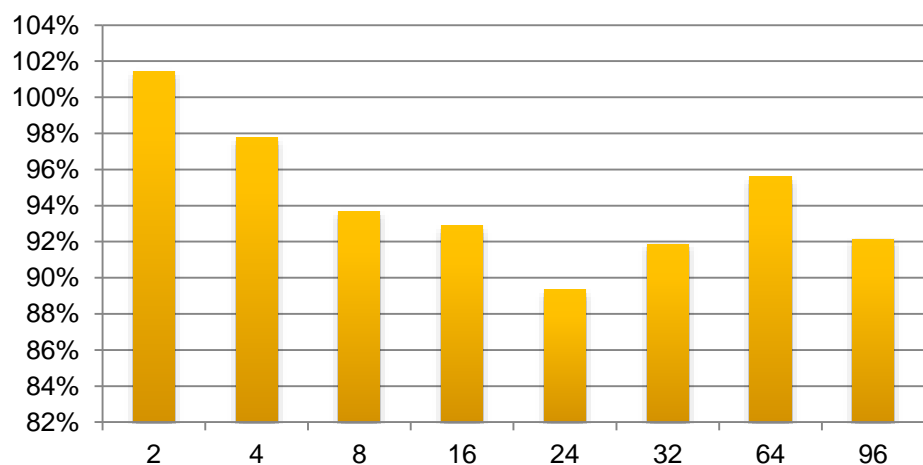
大小 = 0K (平衡)



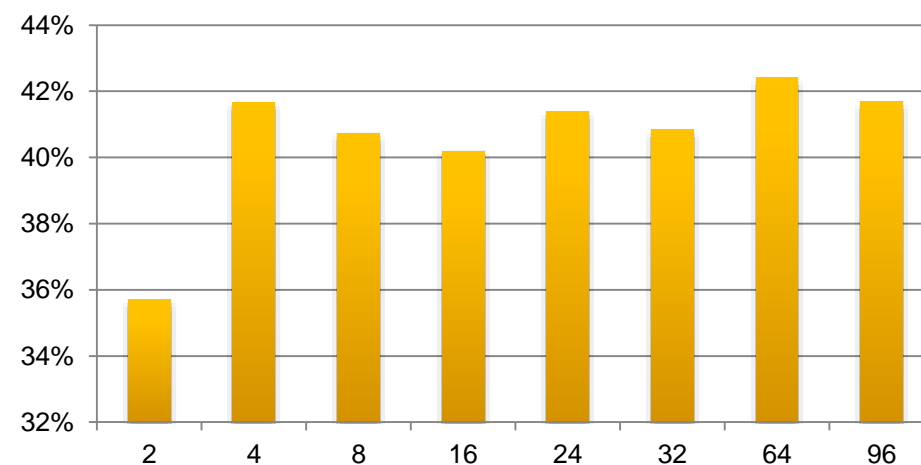
大小 = 0K (唯一)



大小 = 4K (平衡)



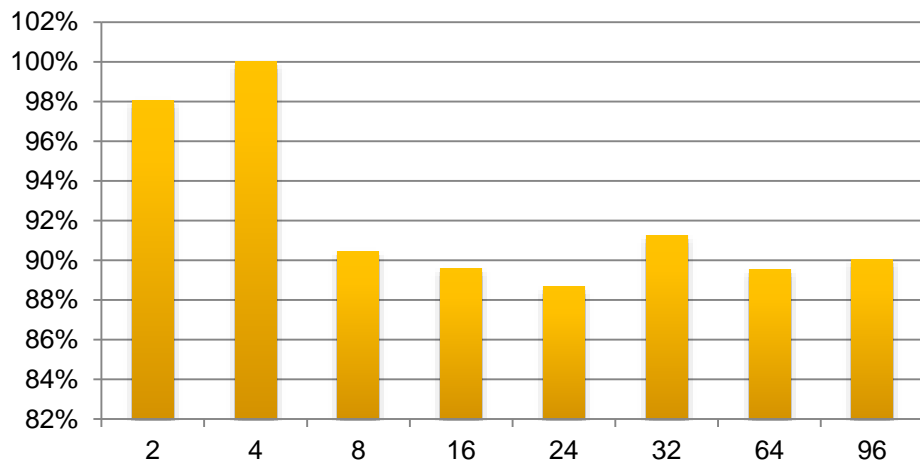
大小 = 4K (唯一)



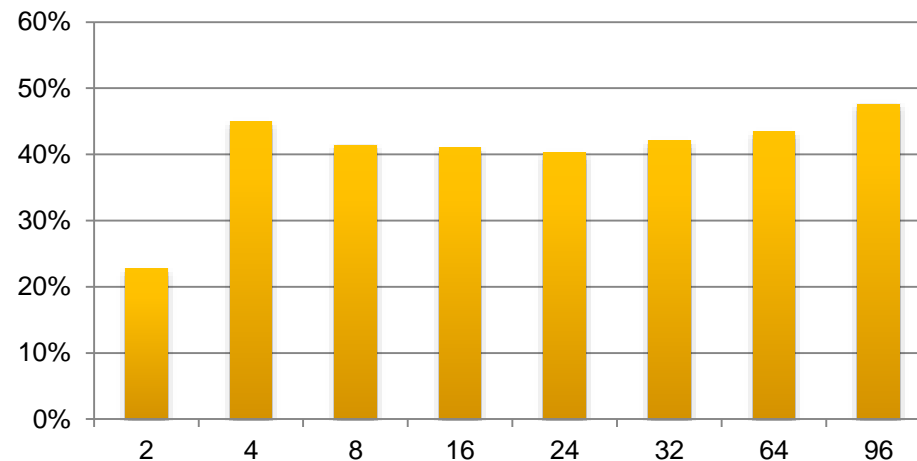
LMBench – 上下文切换延迟 – Scaling

guest/native[%]

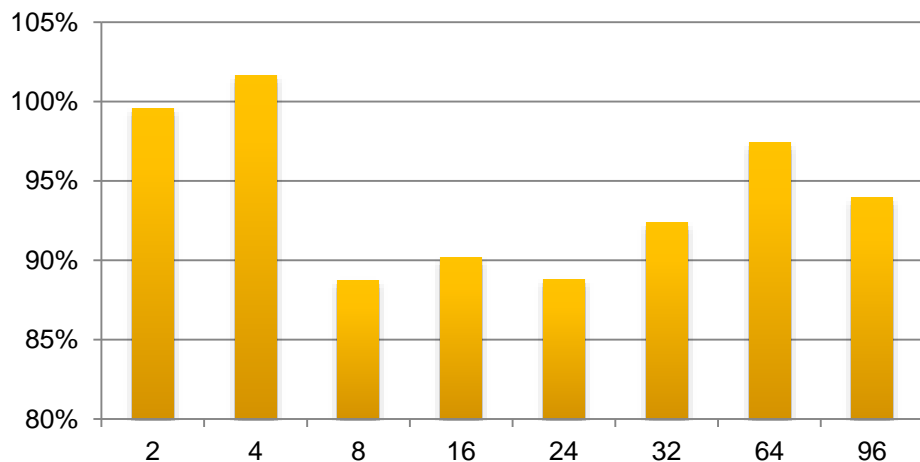
大小 = 0K (平衡)



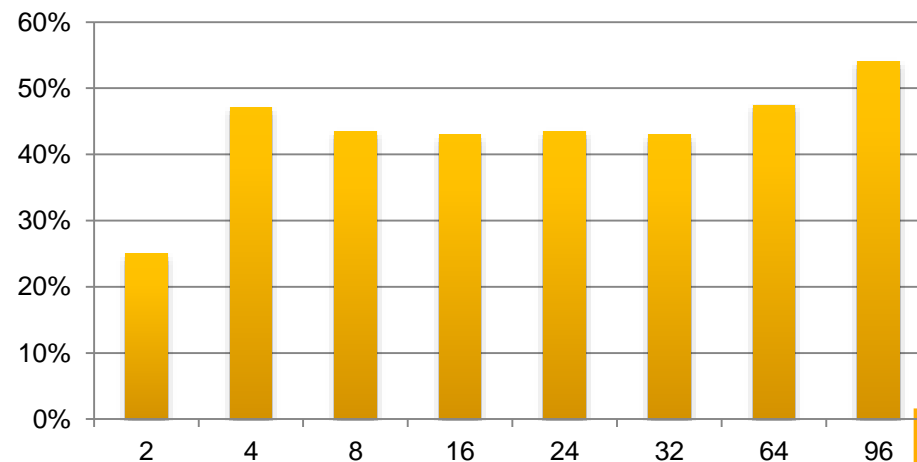
大小 = 0K (唯一)



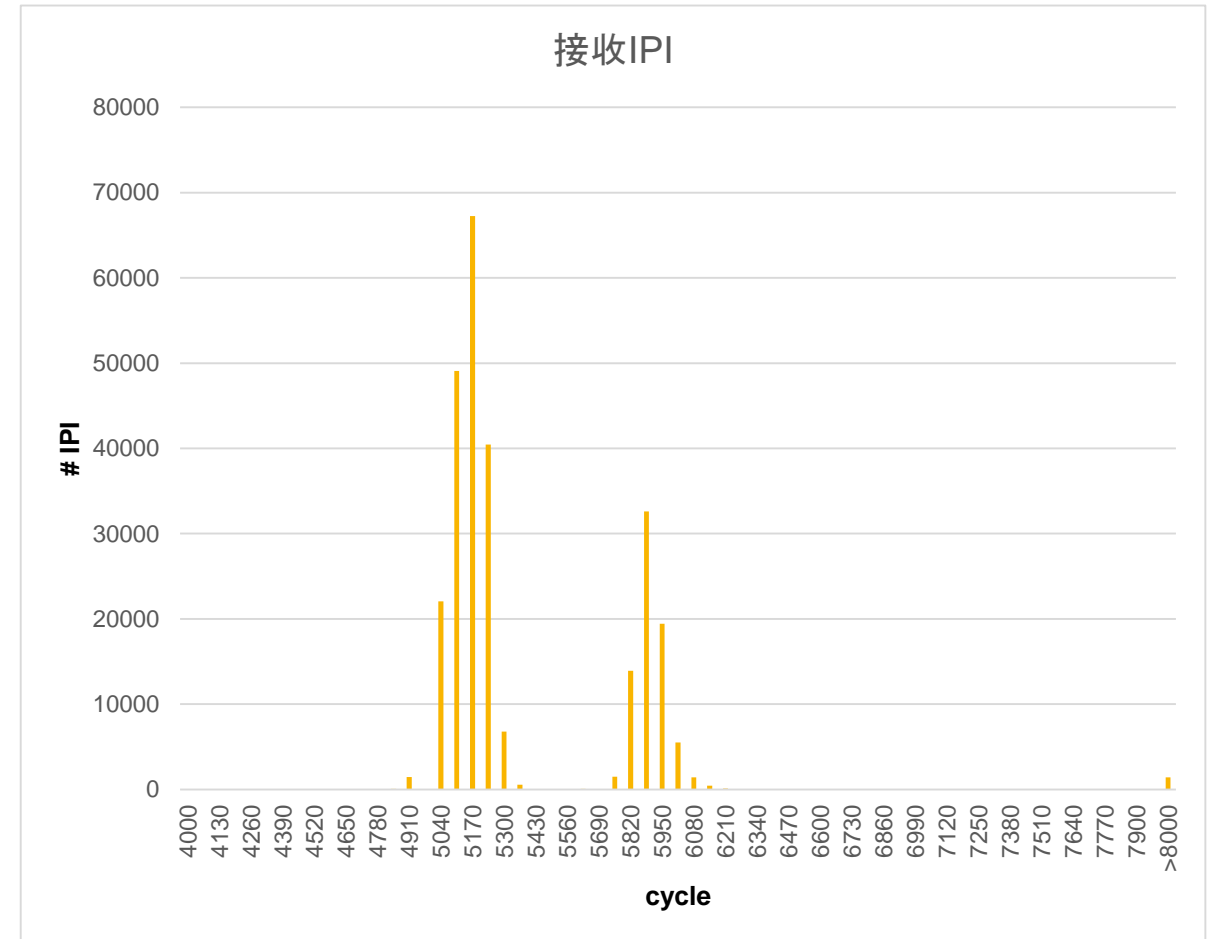
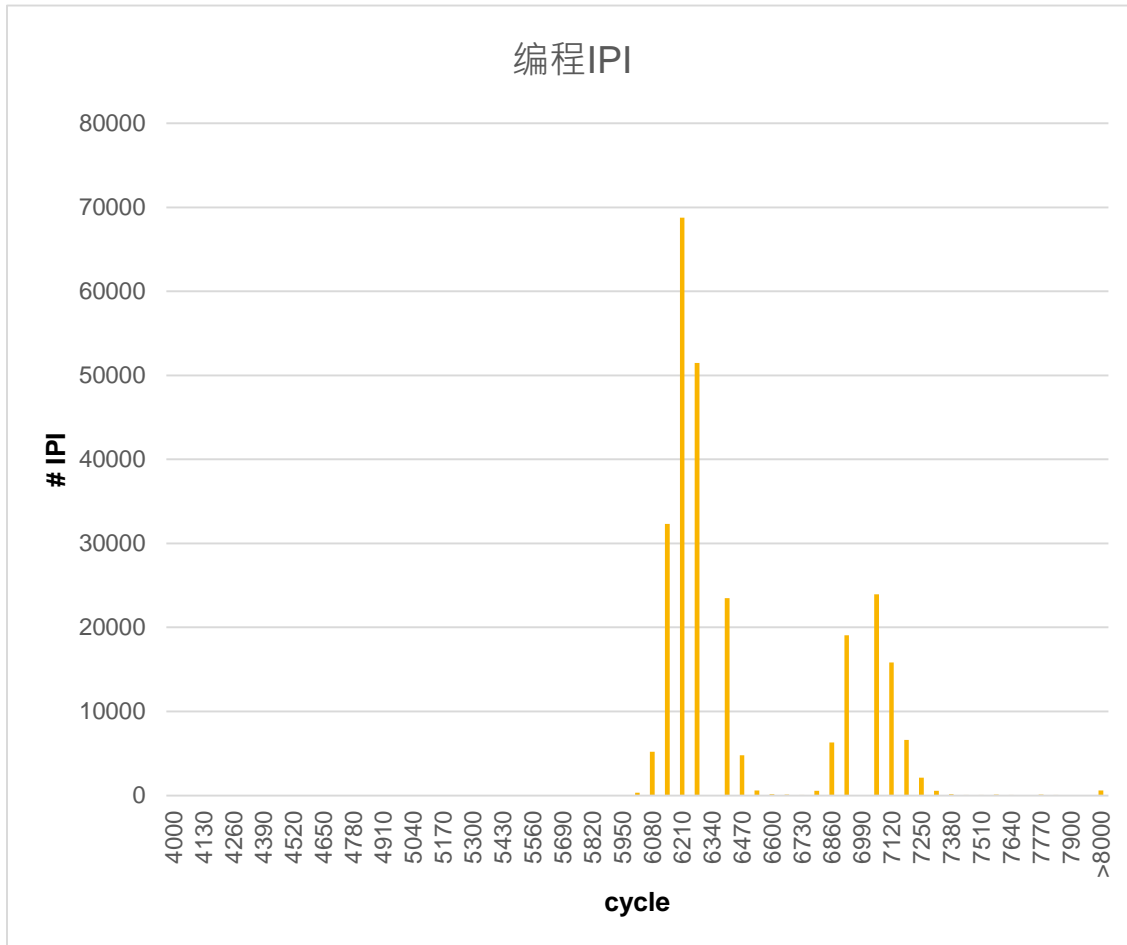
大小 = 4K (平衡)



大小 = 4K (唯一)

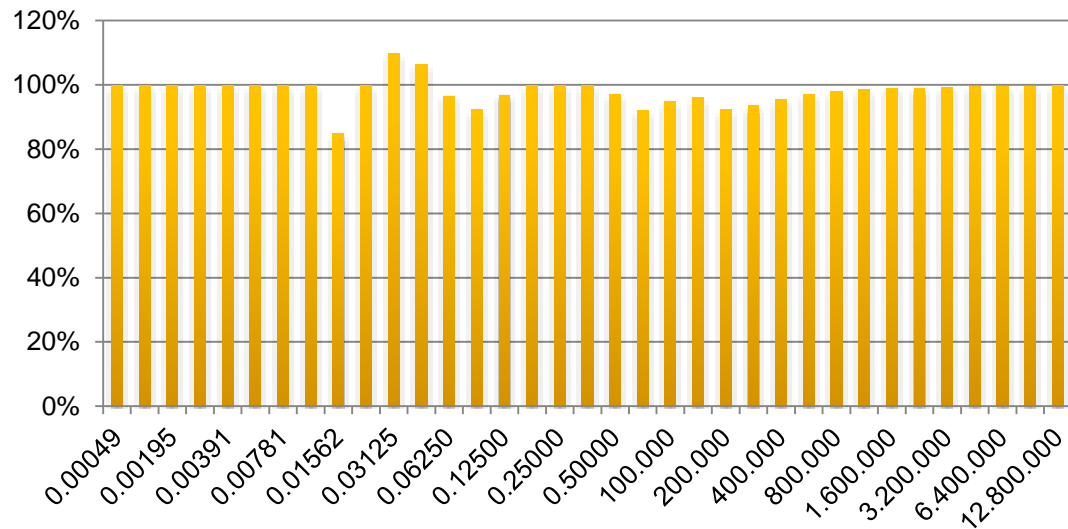


上下文切换延迟分配

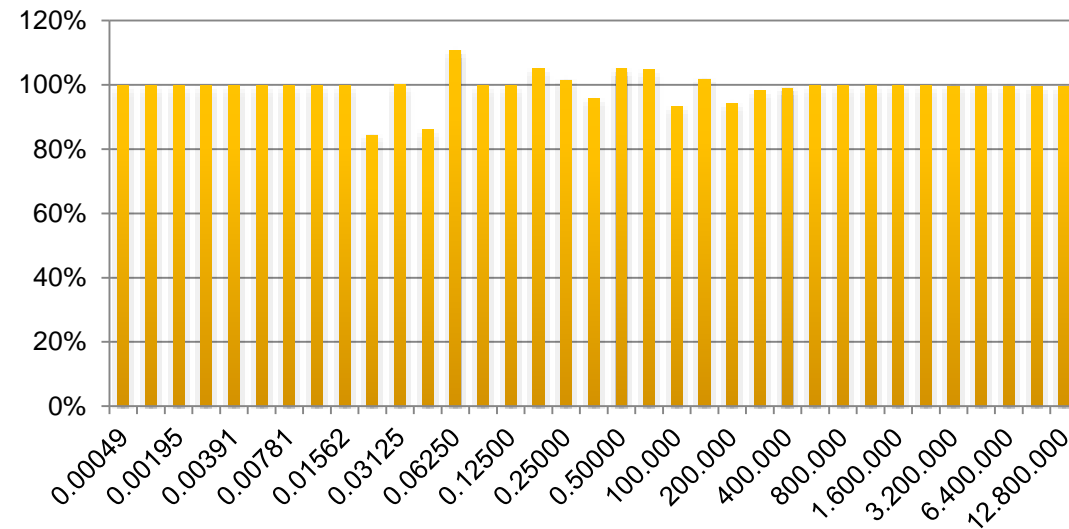


存储器加载延迟 – VM与本机

存储器加载延迟 – 线性



存储器加载延迟 – 随机



结论

结论

- 对于与内核相关的 **benchmark**, 由于退出很少, 因此性能良好
- 开销来源为: 特别是由**IPI**和中断仿真造成的**Guest**退出
 - 通过重新设计**GIC**分配器仿真而提升性能
- 与内存相关的**benchmark**在虚拟化环境中没有显示出重要的开销
 - 页表级数不会产生重大影响



SECURE CONNECTIONS
FOR A SMARTER WORLD

版权声明

恩智浦、恩智浦徽标、恩智浦“智慧生活，安全连结”、CoolFlux、EMBRACE、GREENCHIP、HITAG、I2C BUS、ICODE、JCOP、LIFE VIBES、MIFARE、MIFARE Classic、MIFARE DESFire、MIFARE Plus、MIFARE Flex、MANTIS、MIFARE ULTRALIGHT、MIFARE4MOBILE、MIGLO、NTAG、ROADLINK、SMARTLX、SMARTMX、STARPLUG、TOPFET、TrenchMOS、UCODE、飞思卡尔、飞思卡尔徽标、AltiVec、C 5、CodeTEST、CodeWarrior、ColdFire、ColdFire+、C Ware、高效解决方案徽标、Kinetis、Layerscape、MagniV、mobileGT、PEG、PowerQUICC、Processor Expert、QorIQ、QorIQ Qonverge、Ready Play、SafeAssure、SafeAssure徽标、StarCore、Symphony、VortiQa、Vybrid、Airfast、BeeKit、BeeStack、CoreNet、Flexis、MXC、Platform in a Package、QUICC Engine、SMARTMOS、Tower、TurboLink和UMEMS是NXP B.V.的商标。所有其他产品或服务名称均为其各自所有者的财产。ARM、AMBA、ARM Powered、Artisan、Cortex、Jazelle、Keil、SecurCore、Thumb、TrustZone和 μ Vision是ARM Limited（或其子公司）在欧盟和/或其他地区的注册商标。ARM7、ARM9、ARM11、big.LITTLE、CoreLink、CoreSight、DesignStart、Mali、mbed、NEON、POP、Sensinode、Socrates、ULINK和Versatile是ARM Limited（或其子公司）在欧盟和/或其他地区的商标。保留所有权利。Oracle和Java是Oracle和/或其关联公司的注册商标。Power Architecture和Power.org文字标记、Power和Power.org徽标及相关标记是Power.org的授权商标和服务标记。© 2015–2016 NXP B.V.

