



智能网络构建高效云计算平台

张辉， 亚太及中国区解决方案营销总监

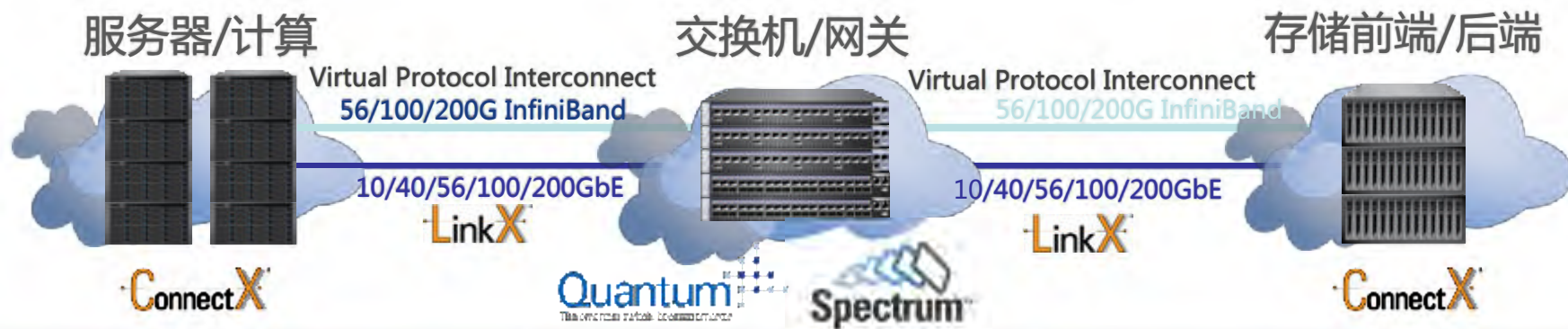
April 19th . 2017



领先的端到端互联解决方案提供商



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT



完整的端到端互联产品家族

芯片



网卡



NPU & 多核



交换机/网关



软件



Metro / WAN



网线/光模块





虚拟化

不牺牲性能的高效、弹性



实现终极云性能的三个障碍



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球新开源
GLOBAL CLOUD-COMPUTING OPEN SOURCE SUMMIT

低效的网络协议

有状态的传输协议(如TCP使卸载更复杂,造成额外的CPU开销)

计算虚拟化的惩罚

由于传输额外的软件程序开销引起的I/O性能下降

网络虚拟化的惩罚

Overlay 网络虚拟化导致新的数据包报头格式的I/O，系统可能无法处理



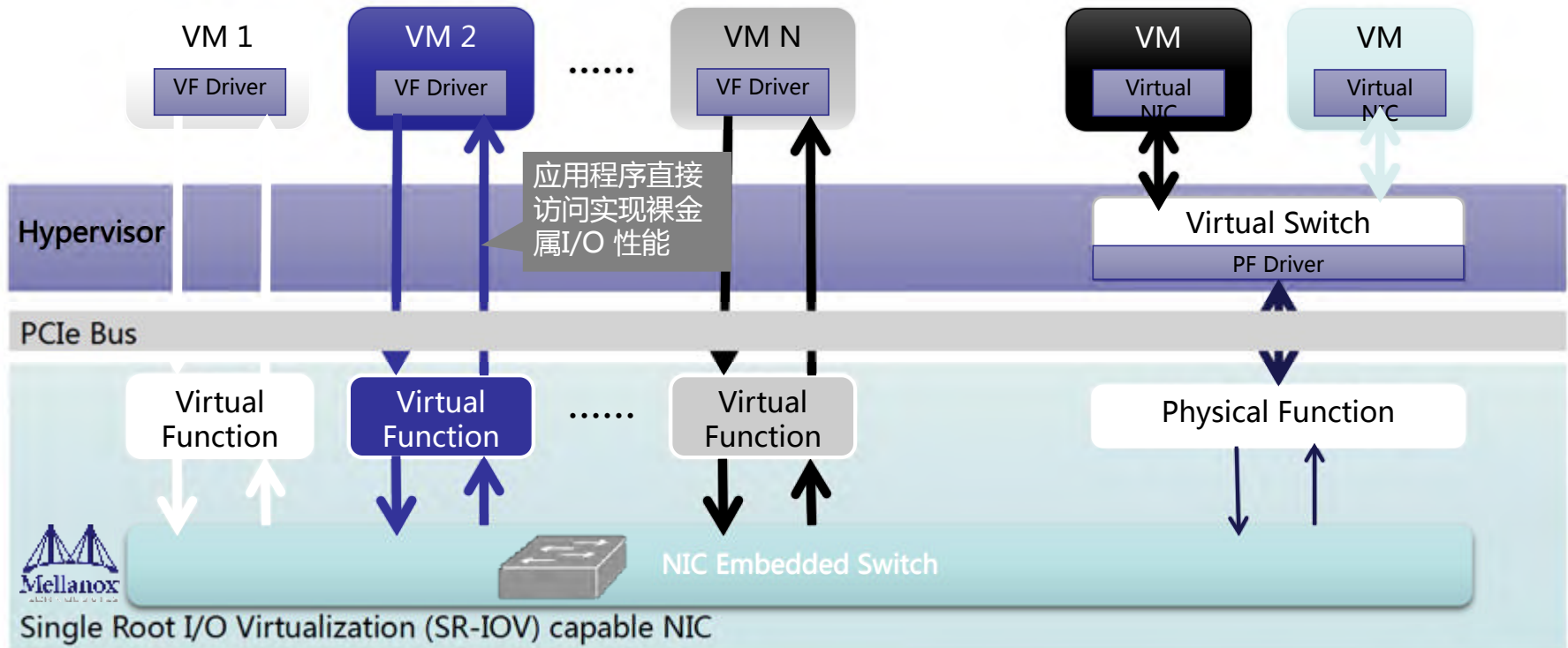
SR-IOV – 克服计算虚拟化的惩罚



全球云计算开源峰会2017
联合云计算新势力，拥抱全球新开源
GLOBAL CLOUD-COMPUTING OPEN-SOURCE SUMMIT

虚拟机利用SR-IOV和Mellanox eSwitch
实现近线性的性能而无CPU消耗

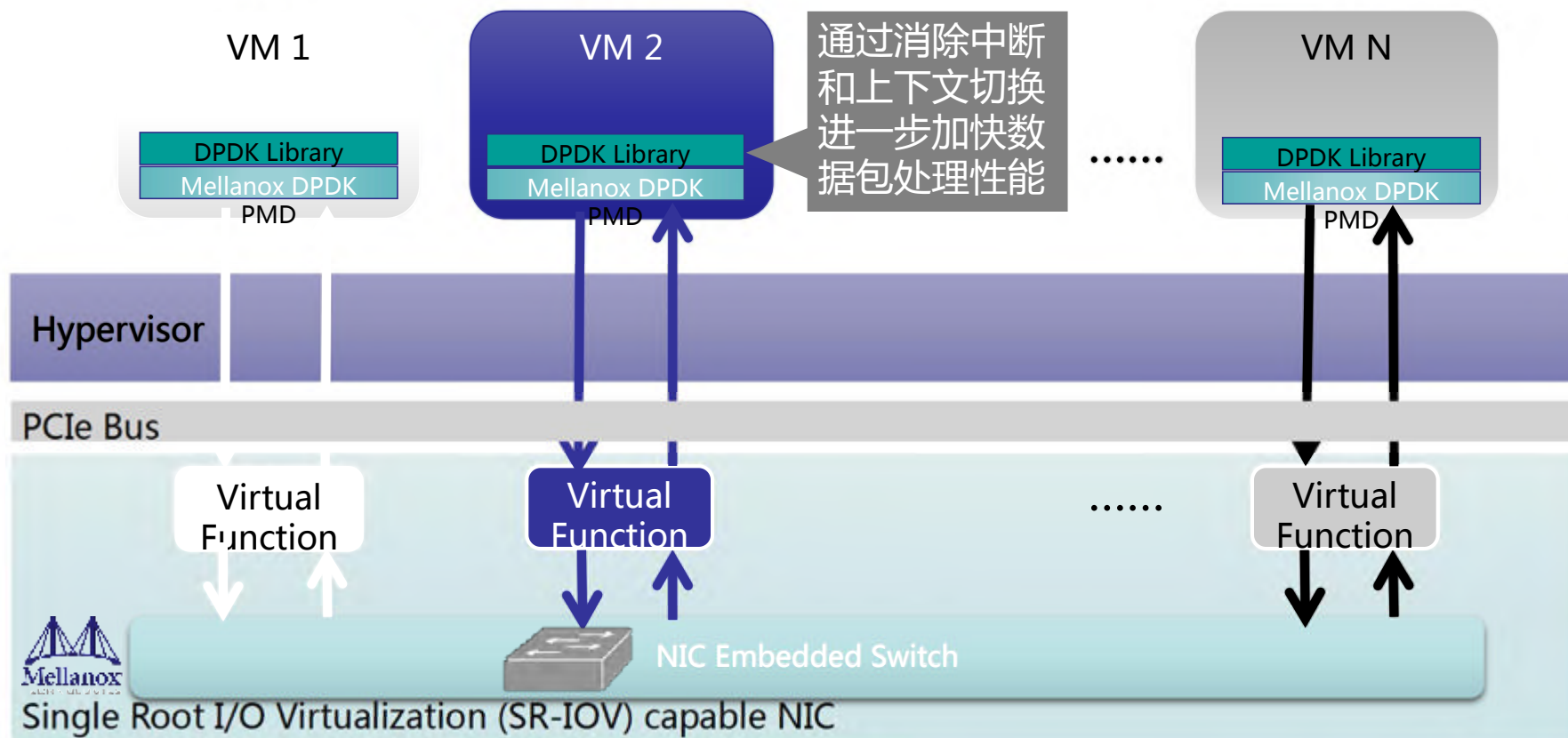
基于软交换的虚拟机饱受计算虚拟化的
惩罚



SR-IOV + DPDK: 与Mellanox PMD协同效果更佳



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球新开源
GLOBAL CLOUD-COMPUTING OPEN SOURCE SUMMIT

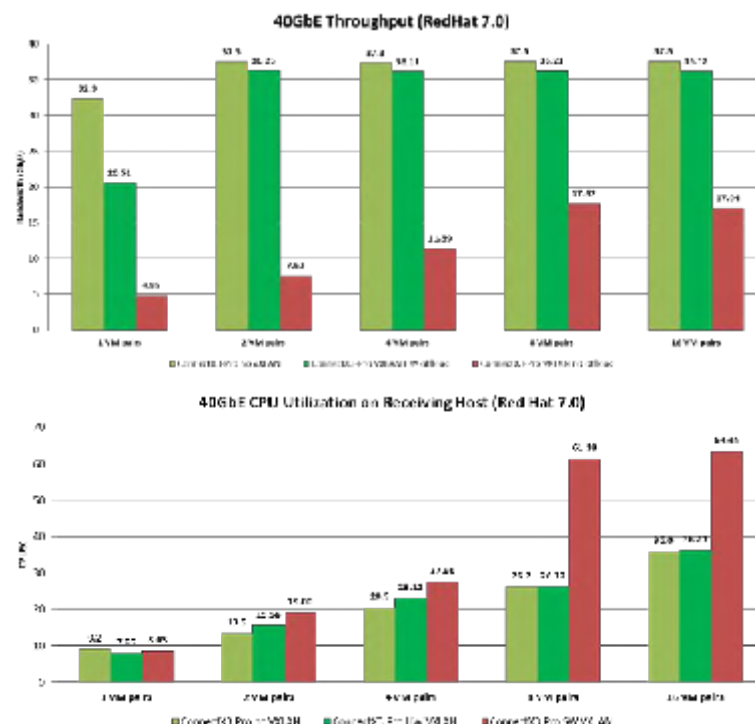


VXLAN卸载 – 克服网络虚拟化的惩罚



全球云计算开源峰会2017
聚合云计算新势力，拥抱金世界新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

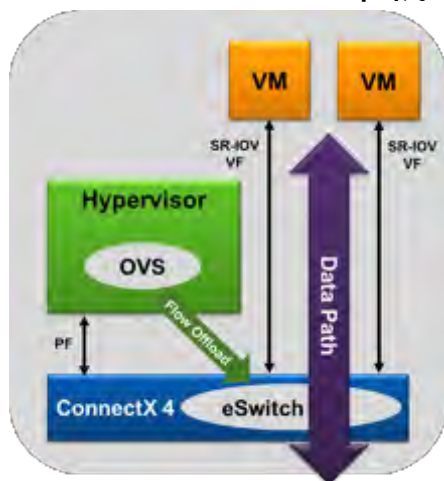
- 解决方案：
 - 网卡中的Overlay网络加速器
 - 以裸金属速度运行无性能牺牲
 - 主流的SDN厂商已验证和集成
- 收益：
 - 40Gb链路上达到**37.5Gb/s**，**2倍**的性能提升（相对没有使用VXLAN卸载）
 - 在一个20个Core的系统中，释放出7个Core 用于运行虚拟机，节省**35%**的内核，同时实现吞吐率**两倍**提升！



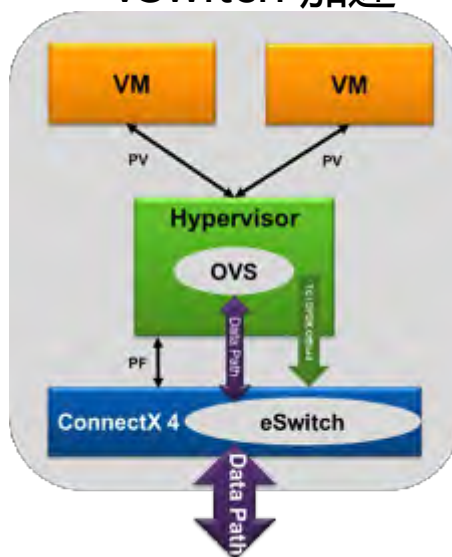
Accelerated Switch And Packet Processing (ASAP²)

- ASAP² 利用ConnectX-4能力加速或卸载 “主机” 网络堆栈
- 三种主要使用方式

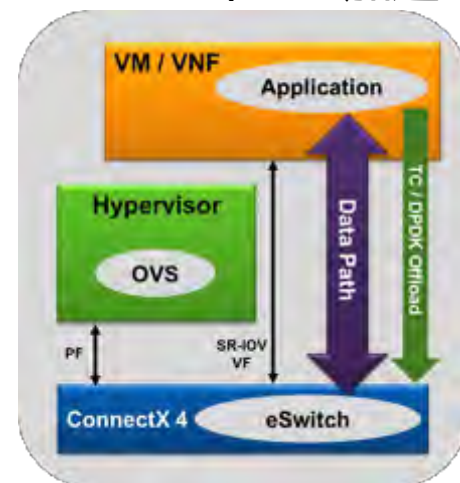
ASAP² Direct
全vSwitch卸载



ASAP² Flex
vSwitch 加速



ASAP² Flex
VNF/VM 加速

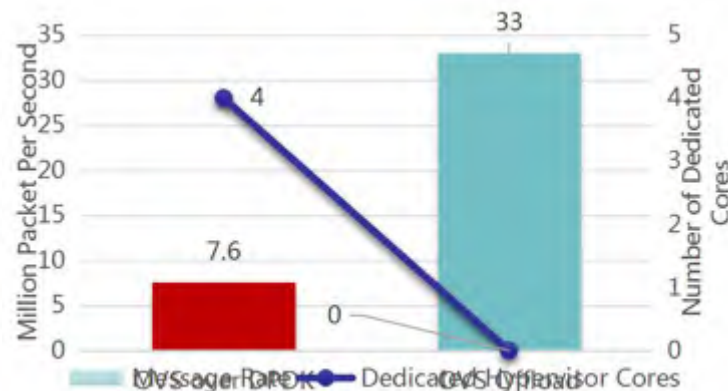
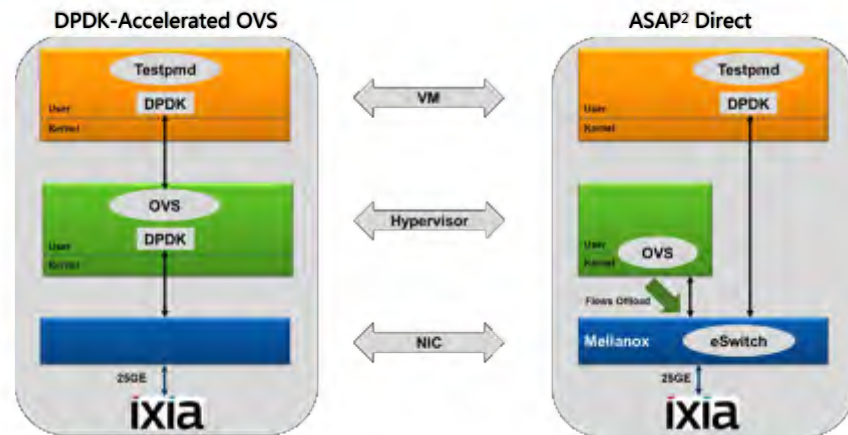


DPDK加速的OVS vs Mellanox ASAP² Direct



全球云计算开源峰会2017
聚力云计算新势力，拥抱全球新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 1 flow, 没有VXLAN
 - 消息传输率是OVS over DPDK的3.3倍
 - CPU消耗率为零！对比OVS over DPDK却使用4个CPU内核
 - 释放出来的CPU可以用于运行虚拟机
- 2000 flows, VXLAN硬件封装/解封装
 - OVS offload 实现 ~25MPPS
 - CPU消耗依然为零，DPDK却使用4个CPU内核





加速

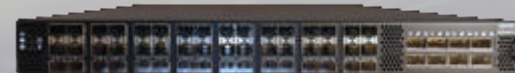
为横向扩展的应用程序启用快速网络和存储访问



Spectrum: 25/100GbE交换机中的旗舰



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球新开源
9th Global Cloud Computing Open Source Summit



简单易用

- 智能可视
- 云敏捷及QoS
- 自动化

性能

- 10倍延迟
- 包零丢失
- 智能缓冲

扩展性

- 更低的功耗
- 2倍的转发表
- 更低的成本\$/Gbps

- 最高的包转发率
- 更好的缓冲：10-15倍的 microburst性能提升
- 可预测的运行状态
- 更低的延迟：降低50%
- 更低的能耗



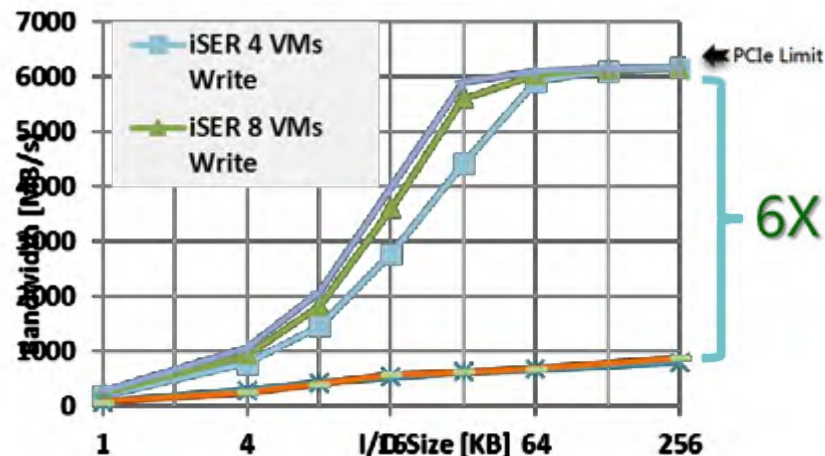
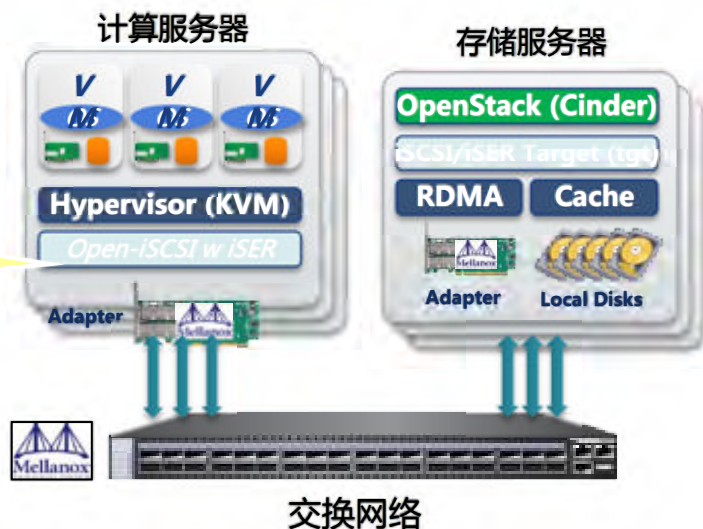
RDMA加速访问OpenStack存储



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT



使用RDMA加速iSCSI存储



- 内置的OpenStack组件和管理
 - 不需要额外的软件
 - RDMA已经内置到OpenStack并被用户广泛使用

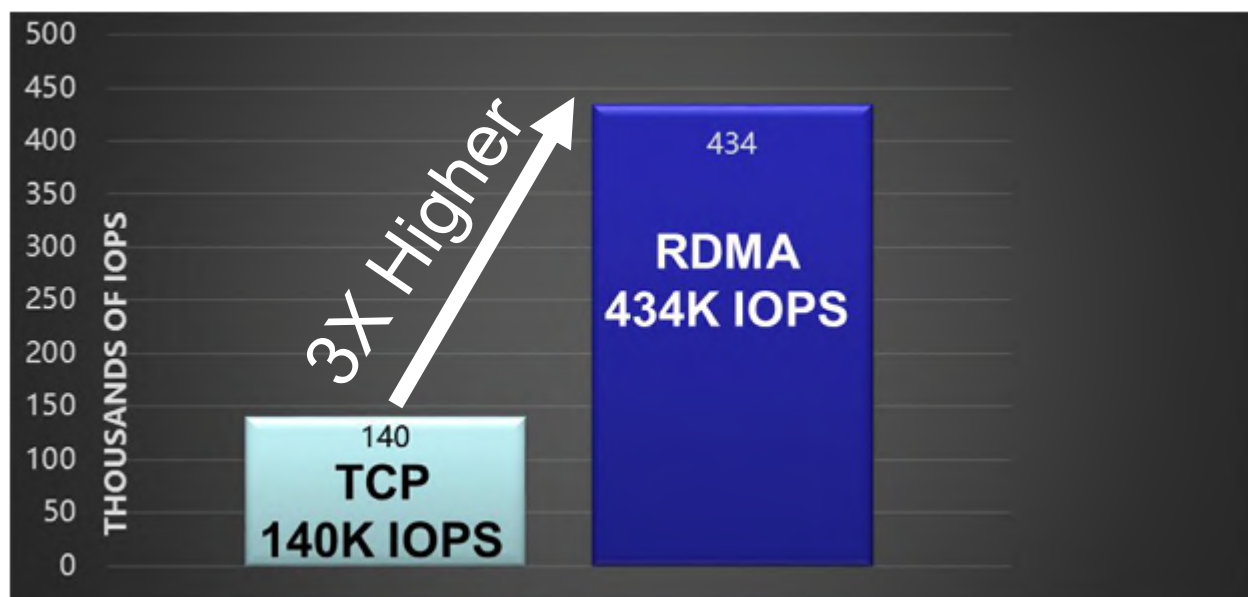
RDMA实现带宽提升6倍，延迟降低5倍，并且降低CPU利用率

使用RDMA加速Ceph



全球云计算开源峰会2017
聚合云计算新势力，拥抱新世界新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

- 更快的RDMA集成至应用
- 最大化消息和CPU并发
- 单节点 > 10GB/s
- 延迟 < 10usec



Ceph Read IOPS: TCP vs. RDMA

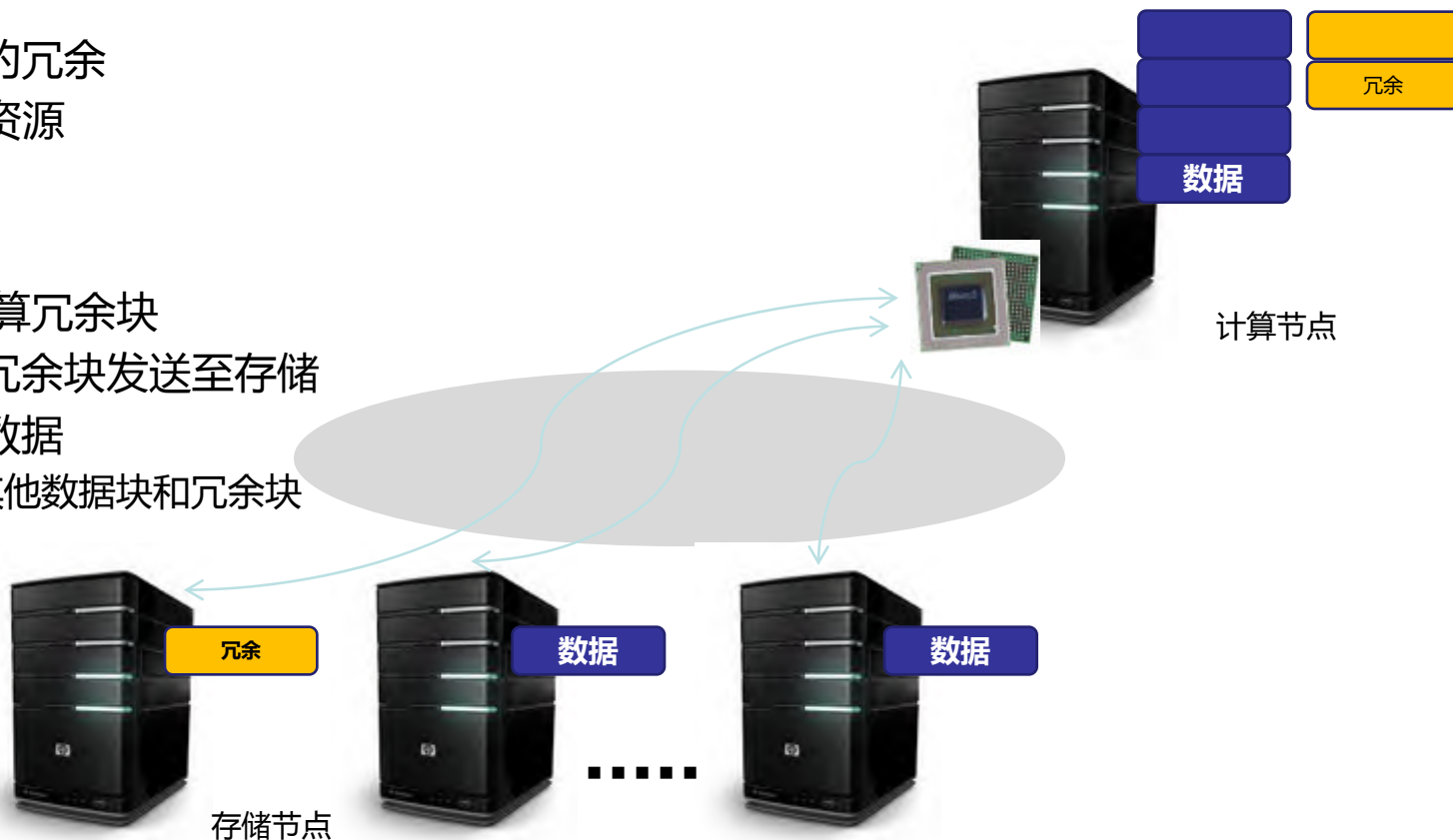
纠删码卸载 (Erasure Coding)

- 价值

- 集群级别的冗余
- 释放CPU资源

- 硬件卸载

- RS编码计算冗余块
- 数据块和冗余块发送至存储
- 重算丢失数据
 - 给予其他数据块和冗余块





自动化

操作简单、高效



开放组合网络 (OCN)



精选的网络
操作系统



开放的API



自动化



端到端互连



操作简单、高效



全球云计算开源峰会2017
聚合云计算新势力，拥抱全球开源生态
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

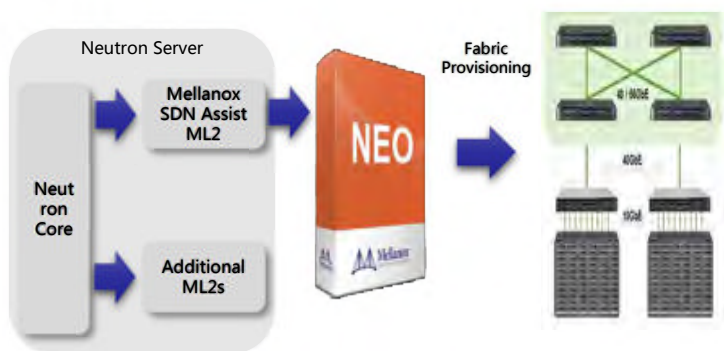


MLNX Switches in
vSphere inventory

NEO与vSphere集成



全面支持以太网和InfiniBand



零接触网络配置



网络感知调度



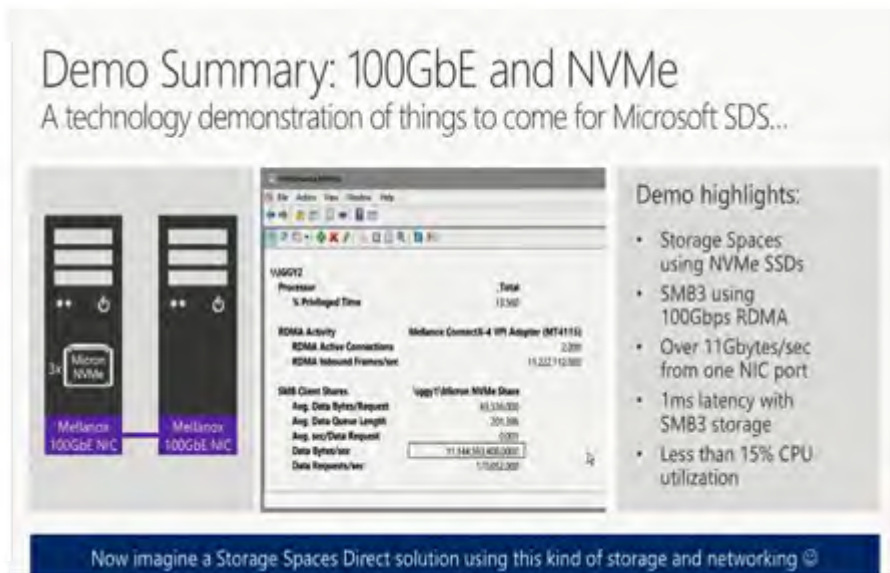
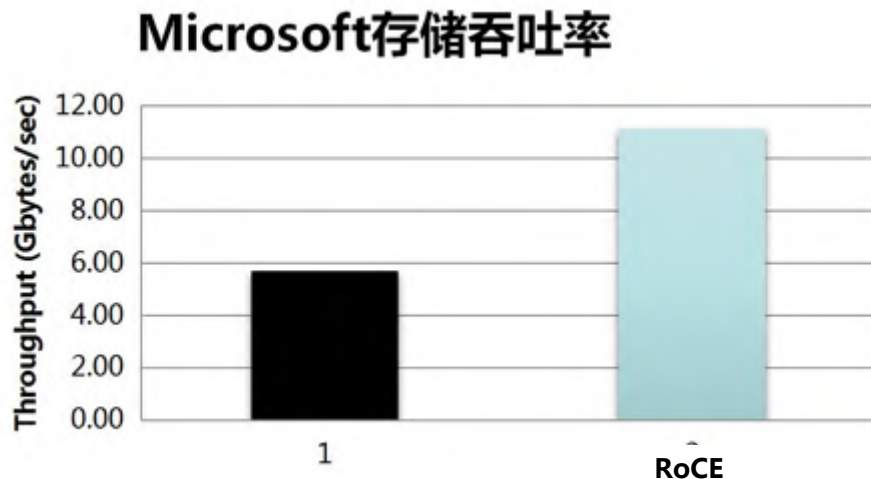
案例分享



Mellanox 100G方案使云效率最大化



全球云计算开源峰会2017
联合云计算新势力，拥抱全球新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT



- RDMA在Windows 2016及 AzureStack被广泛支持
- 使用 RoCE 可以提升2倍带宽及2倍CPU利用率
- RoCE 充分利用的闪存存储带宽
 - 远程存储再不妥协



Microsoft

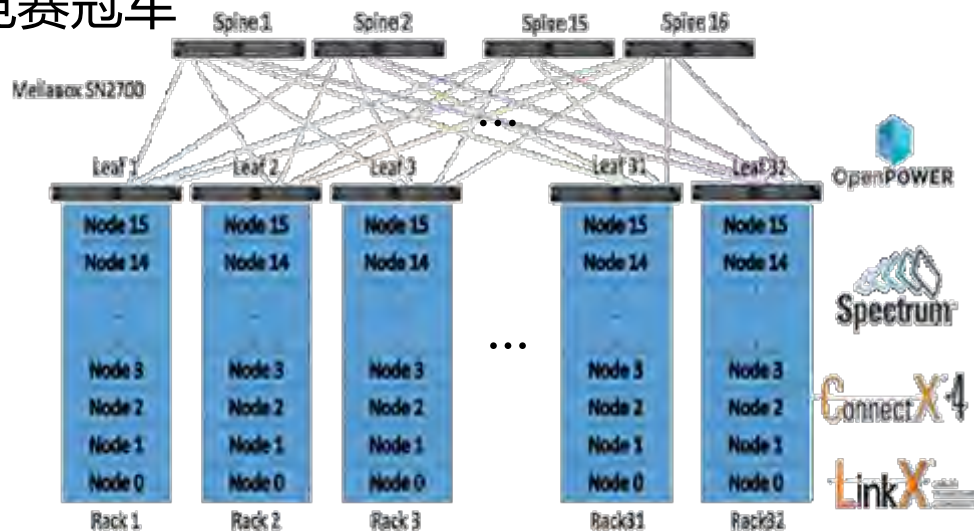
腾讯云创造 Terasort 世界纪录



全球云计算开源峰会2017
聚合云计算新势力，拥抱新世界新开源
GLOBAL CLOUD COMPUTING OPEN SOURCE SUMMIT

TeraSort 年度排序基准全球计算竞赛冠军

- Mellanox 和 IBM 携手 Tencent 实现超过 1TB/秒的 TeraSort 性能新纪录
- 打破了 GraySort 和 MinuteSort 类别的纪录，将去年的成绩提高了**5.8**倍
- 每节点的性能比上一届提升了**33**倍



* <http://sortbenchmark.org/>



非常感谢

