



面向数据中心的高性能 SDN交换技术

深圳市风云实业有限公司

2016年10月

1. 云数据中心网络的发展趋势
2. Nebula Flex Connect 技术介绍

云数据中心网络的发展趋势

应用层

物理层



智能化的应用软件开发逻辑
具备可编程能力的硬件
自动化的配置和管理方法

来自客户的疑问

- 敏捷性
- 安全性
- 自动化
- 移动性
- 解耦合

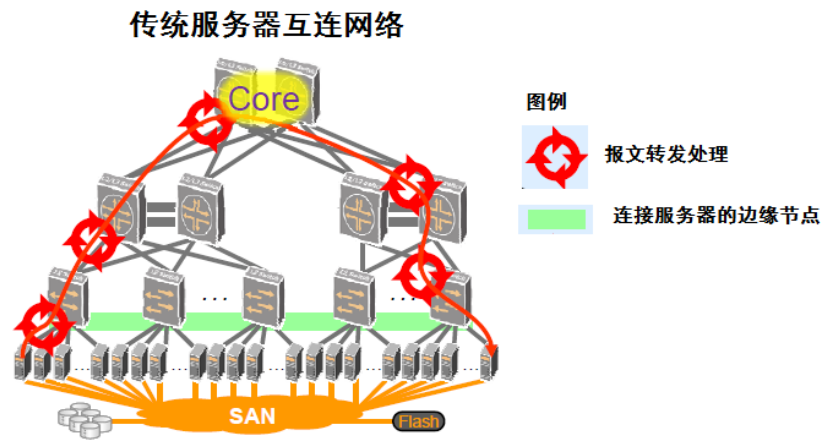
• 数据交换 →

大规模
高密度
低延迟
高吞吐
集约化
弹性化



传统通信网络互连服务器带来如下问题：

- 通过交换机级联构建网络
- 各个服务器的性能与其网络位置相关, 服务器节点之间带宽与通信路径强相关, 带宽受限
- 每一步交换都要进行报文转发处理, 解析、查表、封装报文, 增大服务器节点之间跳步, 增加了通信延迟大, 通信性能随跳步数降低
- 静态负载均衡, 不能感知网络状态, 避免拥塞, 不可预知的通信热点, 造成路径拥塞, 加剧延迟的抖动, 链路利用率低, 网络性能在重负载下急剧下降
- 交换机的配置复杂繁琐, 难以维护管理
- 仅支持以太网、FC或者Infiniband 单一网络接口



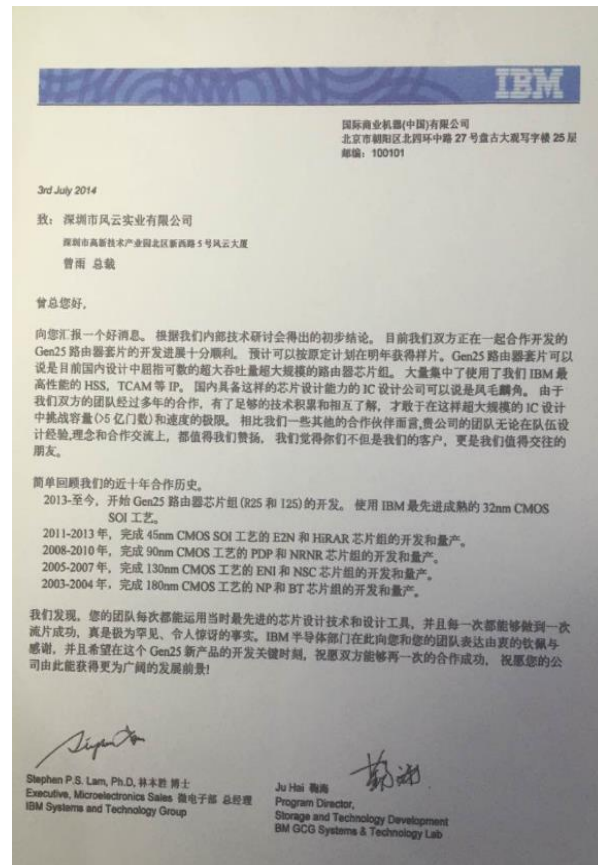
1. 云数据中心网络的发展趋势

2. Nebula Flex Connect 技术介绍



Nebula 互联方案概述

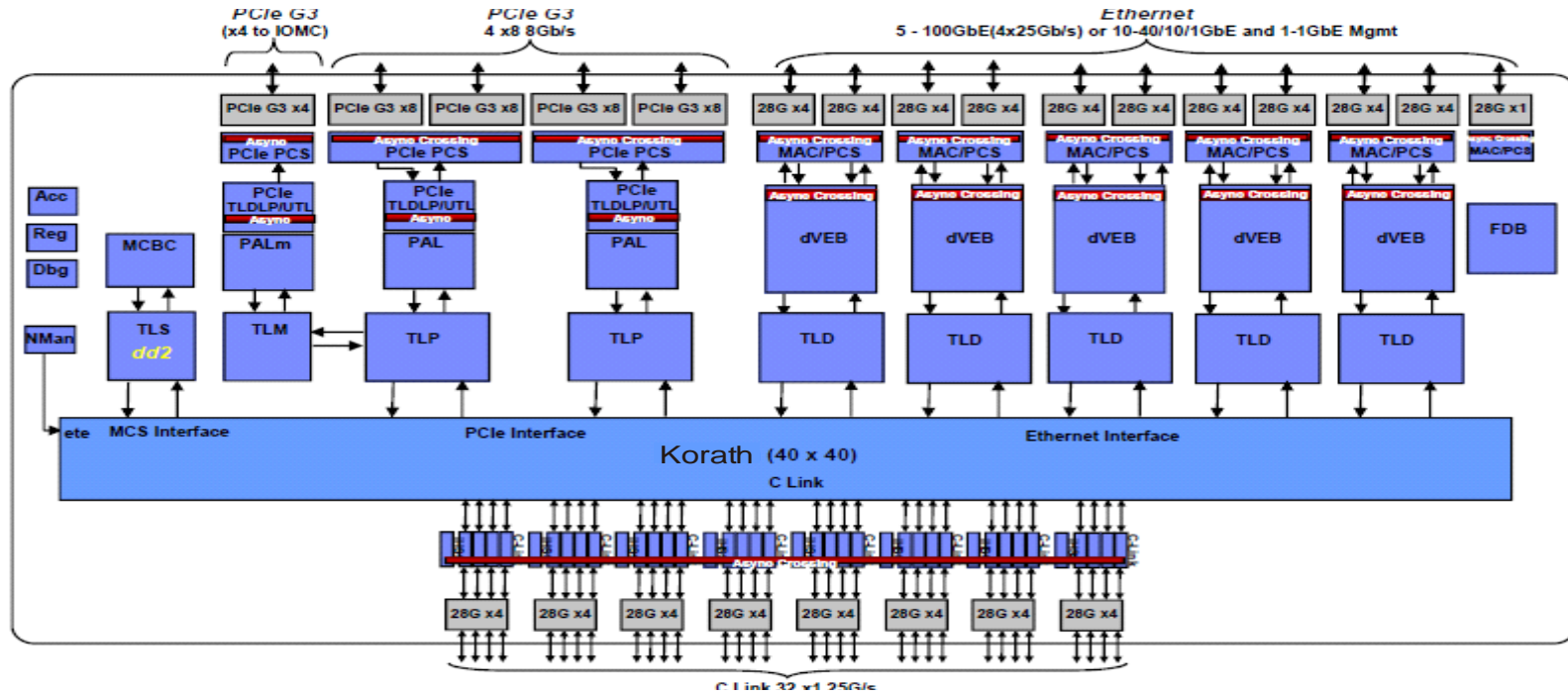
- **定位：**面向大数据与云计算中心市场，消除传统互连网络的弊端，设计研发基于标准以太网和PCIE协议接口的高性能、多协议、跨平台的高性能互连系统核心技术，支持**SDN应用场景**，高效的支持**虚拟化和融合网络**。
- **核心技术：**高带宽、低延迟、高可用、可编程的互连芯片技术为核心。
两款核心的高端互连芯片：通用可编程处理芯片**Korath**和高密度交换芯片**Collector**。
配套的开发和管理软件**IOMC**，以及相关原型系统的机框、机架、板卡的参考设计。
- **知识产权情况：**风云公司和**IBM**拥有所有相关知识产权，相互授权对方不受约束的研发、生产、销售和改进技术。





Korath芯片架构

SDN的尝试和思考





Korath芯片技术规格

(1) 性能指标

- 5路100GE 以太网接口/10路40GE/10GE、5路 PCIE GEN3接口、32路动态均衡的25Gbps C-link接口，提供共计2Tbps 带宽
- 最小转发延迟250ns、30MB 片上报文缓冲
- 支持HASH+TCAM结合的多级流表、内嵌多核可编程协处理器
- 支持136个Korath芯片互连扩展能力(同时支持680个100G 端口和544个PCIE端口)

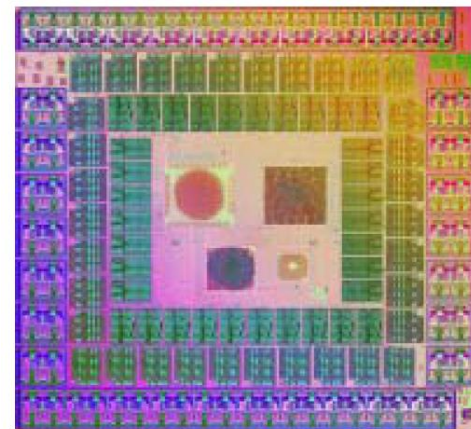
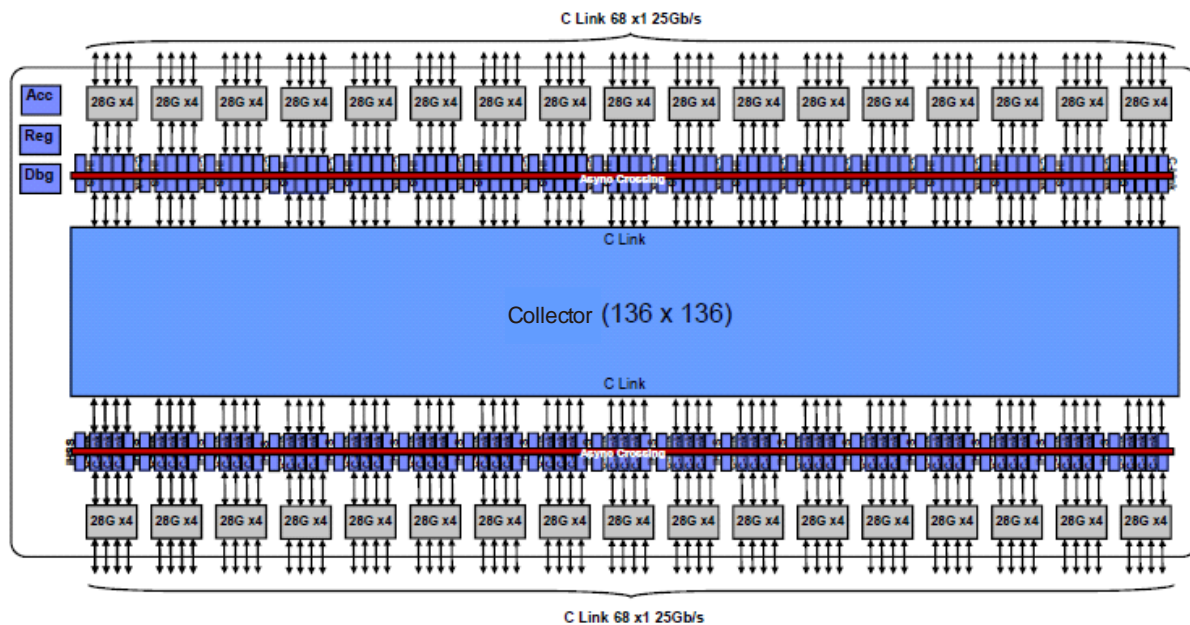
(2) 支持完整的以太网特性

- 支持传统以太网协议
- 支持数据中心扁平交换特性
- 支持数据中心无损以太网协议

(3) 丰富的PCIE交换特性

- 支持每个PCIE端口可以独立配制成上行或下行端口，上行端口接服务器的PCIE Root Port，下行端口接PCIE Adapter
- 支持灵活配置的分布式PCIE Switch，每个上行端口可以连接最大32个下行端口，每个下行MR-IOV端口支持64个上行端口
- 支持单根虚拟化 (SR_IOV)，支持多根虚拟化 (MR_IOV),支持最大64个虚拟层次(VH)、
- 支持通过MATCH-ACTION的机制实现PCIE接口数据交换

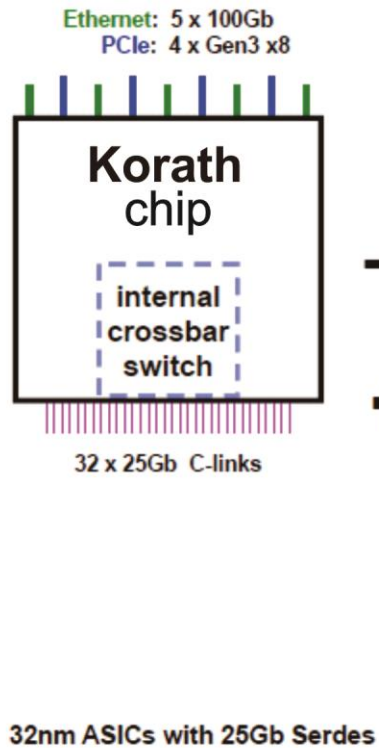
- 136个25G C-link端口，交换带宽6.8Tbps
- 无阻塞交换，交换延迟175ns
- 支持端到端可靠性（自动CRC检查与重传机制）
- 支持芯片级自适应路由、支持芯片级的多路径负载均衡和等价路由



28nm 150W



Nebula应用模式



+

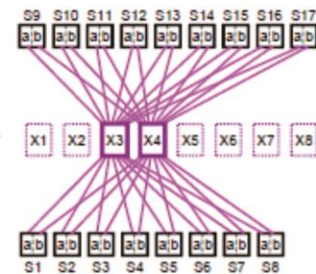


System Configurations

FULL MESH
Best cost/power/density
for Smaller configs

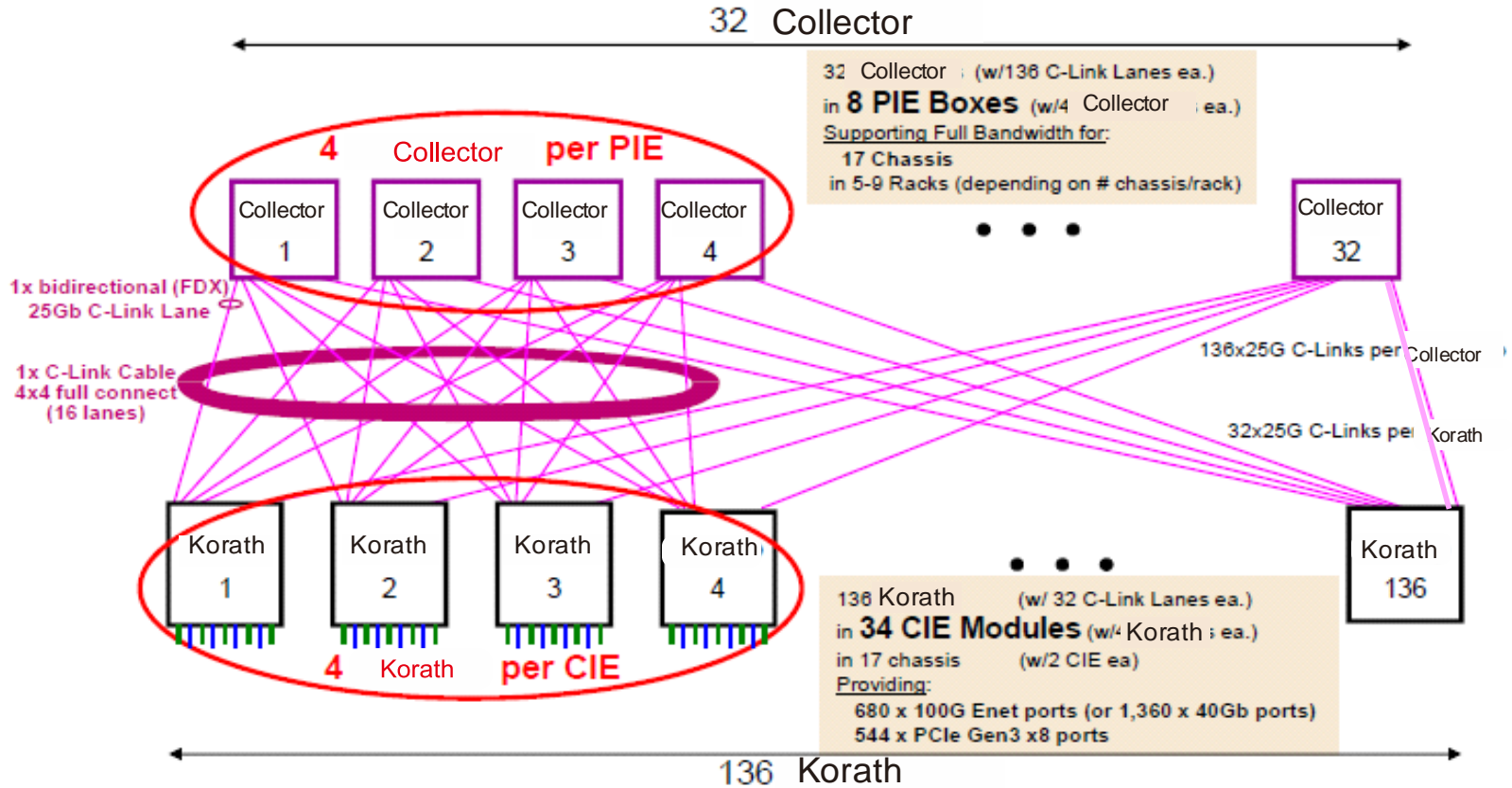


SPINE-LEAF
Best scalability
for Larger configs





Nebula机柜互联



■扁平化的单层分布式交换架构

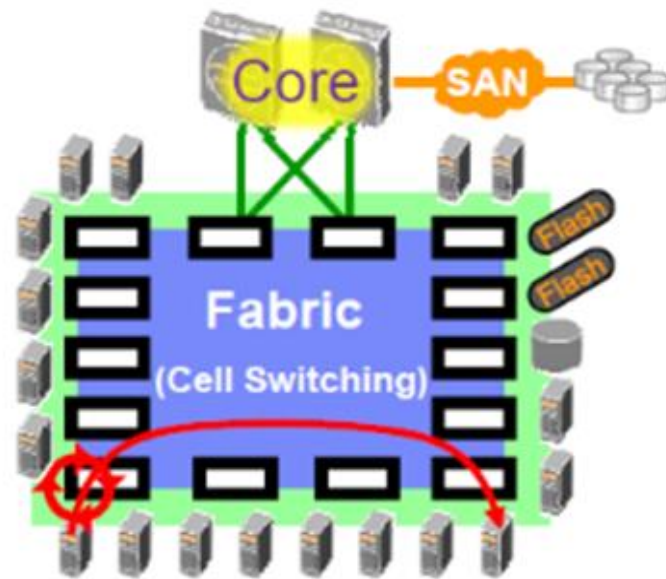
- 高带宽（40Gbps~100Gbps/通道）
- 低延迟（任意两网络节点间一跳可达，且延迟小于1us）
- 性能与服务器、存储节点在网络中的位置无关

■ 更好的支持虚拟机和存储

- 任意位置的虚拟机和存储获得的带宽和延迟一致
- 支持虚拟机灵活迁移和镜像
- 支持超过10万个虚拟机的应用环境
- 可支持RDMA协议用于高效的快速的虚拟机迁移
- 为两个镜像节点提供的高速、低延迟、可保证的通信链路

■ 更好支持SDN

- 大容量多级流表和可编程以太网协议处理引擎
- 将传统的以太网SDN技术引入到PCIE交换中
- 支持目前的主流SDN架构



Nebula数据中心网络

■ 芯片级的自适应路由、动态负载均衡以及Qos技术

- 硬件自动感知网络状态，避免拥塞，
- 链路利用率高
- 网络性能在重负载下保持一致
- 为存储、通信、计算、IO流量提供可预知的带宽和QoS保证
- 在设计上，保证了大型的“elephant”流不影响“mice”流。使得虚拟机迁移、虚拟机镜像、数据库备份、磁盘阵列-闪存阵列交互等应用产生的大数据流不影响延迟敏感的通信应用。
- 管理便捷，无需专家级人员规划和维护网络负载均衡

■ 支持Scale-out可扩展性和灵活性

- 每一个机柜就是一个独立系统，多个机柜可以互连成更大规模系统，资源可以跨机柜共享。目前可支持17个机柜，680个100G以太端口和544个PCIE GEM3X8端口，远超传统高端SMP服务器规模。并且还可具备进一步扩大到数百个机柜规模的潜力。
- 每个机柜内部可以根据需求灵活的配置计算模块、存储模块、PCIE IO模块、闪存模块、网络模块的数量和位置，并且支持动态调整。
- 达到同等带宽所需的交换机和线缆的数量为传统网络的1/4

■ 丰富的应用协议和高扩展性

- 面向网络应用支持以太网协议，面向存储应用支持FCoE，面向高性能集群支持ROCE (RDMA over Ethernet)，面向I/O和协处理器应用支持PCIE。针对FCoE有专门的优化拥塞避免算法。
- 利用内部的低延迟C-link信元交换技术，将来可方便的引入更多的开放标准接口，如SAS、SCSI、FC等等
- 支持端到端的重传，硬件支持链路上CRC校验与重传

■ 数据中心融合网络

- 针对数据库集群、高性能计算集群，同时满足高带宽的大数据流存储需求和低延迟的消息传递需求
- 一次投入，存储、I/O、计算、共享内存、网络等都享有当前最高速的互连带宽，更低的投资、使用和维护成本
- 二次开发成本低，用户可自定义专用协议和接口，一次转换可接入整个系统，例如SAS、SCSI、FC可通过PCIE适配器直接连入系统。

■ 数据中心虚拟化

- 支持SR-IOV、MR-IOV，任意节点的PCIE适配器，全系统共享可见，物理共享、虚拟隔离

基于传统互联网的多网络平台

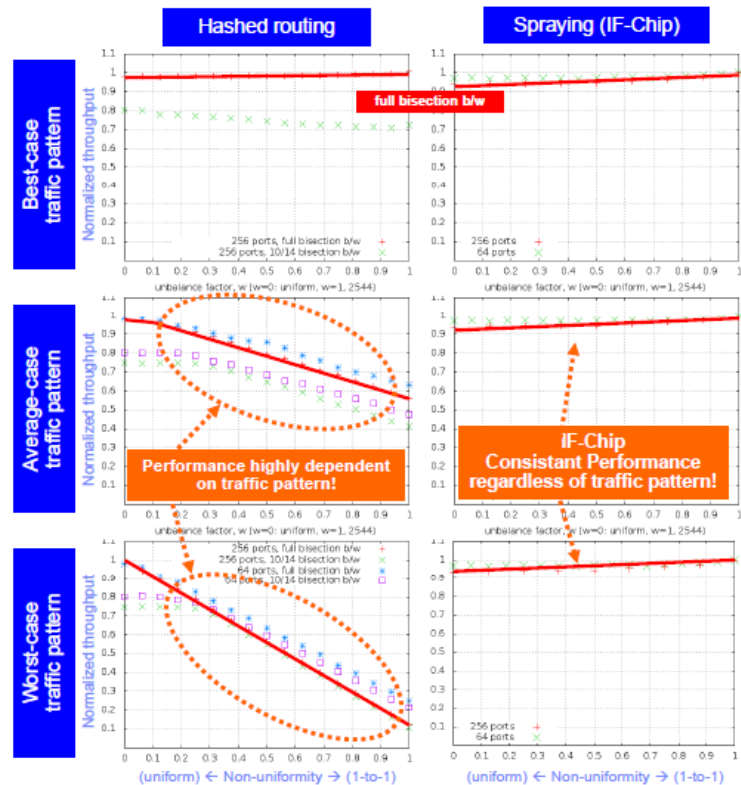
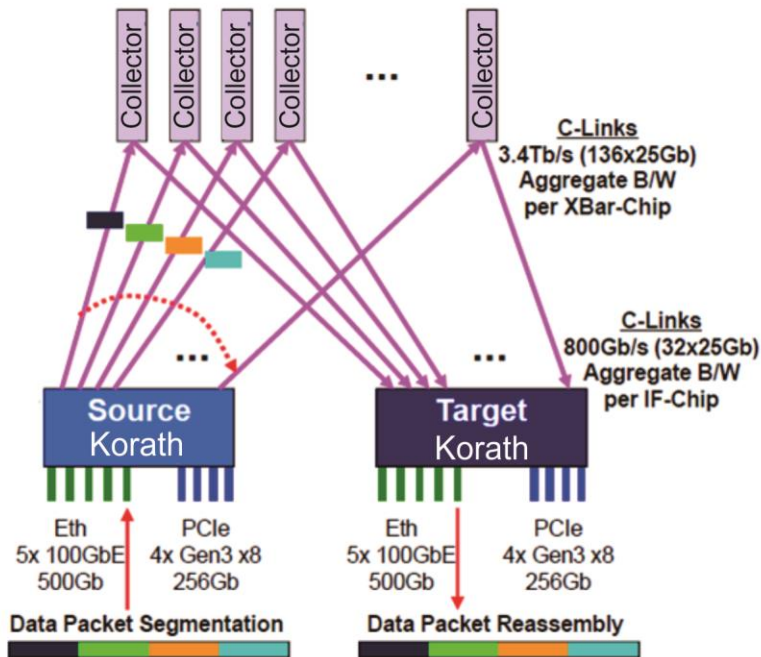


基于Nebula的融合网络平台





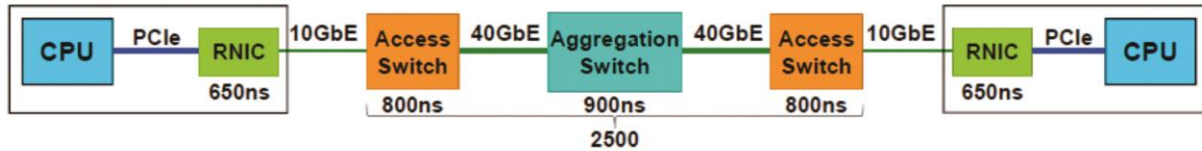
Nebula信元交换技术





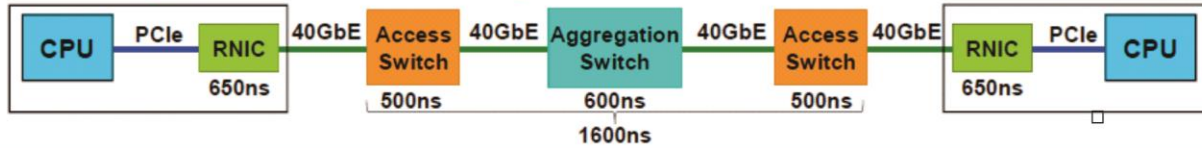
Nebula低延迟转发

3800ns : Industry Multi-tier Ethernet Switching (representative of: 2012-14 10GbE Embedded, ToR, Stacked Switches)



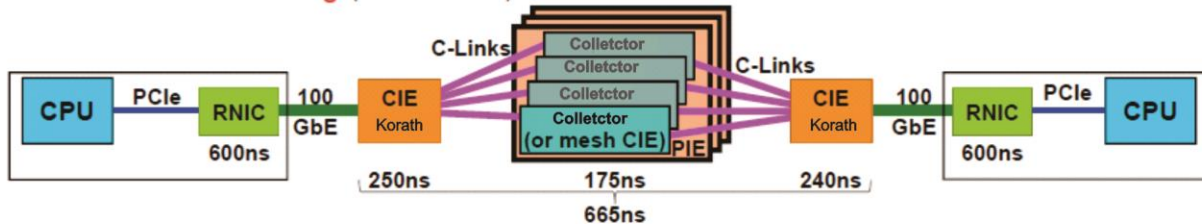
Updated Switch ASICs
Saves 900ns → 23% lower CPU-CPU
36% lower latency Enet-Enet

2900ns : Industry Multi-tier Ethernet Switching (representative of: 2014-16 40GbE)



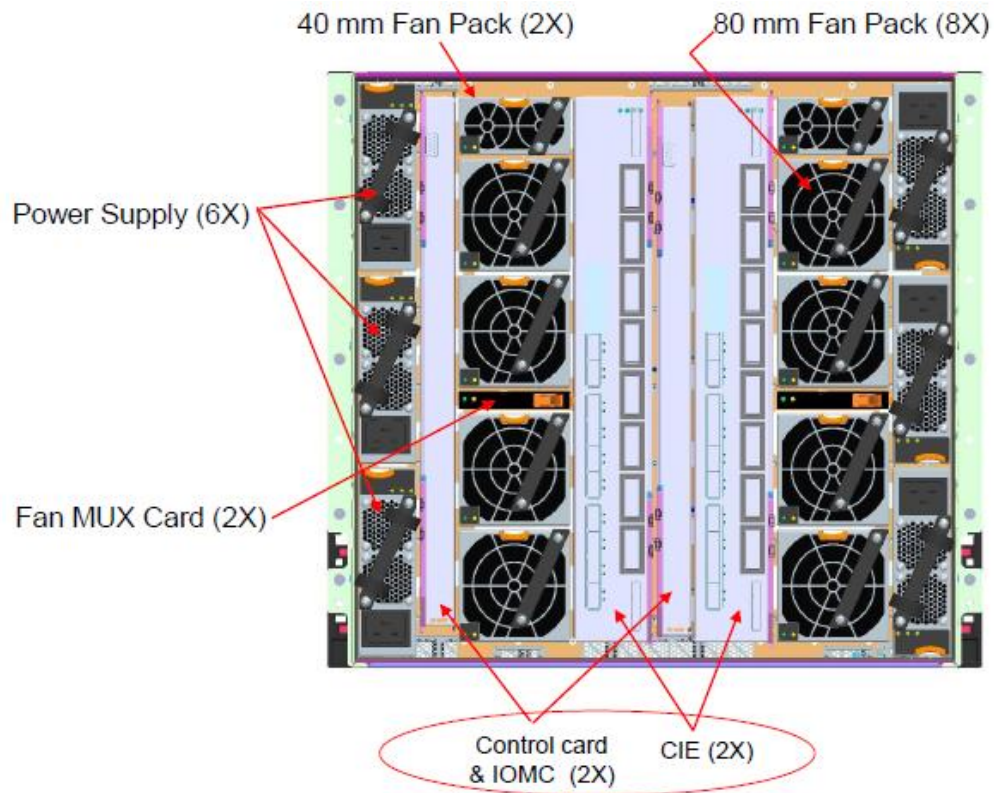
Remove Packet Processing
Saves 1035ns → 35% lower latency CPU-CPU
935ns → 58% lower Enet-Enet

1865ns : Harrier Cell Switching (2015+ 100GbE)



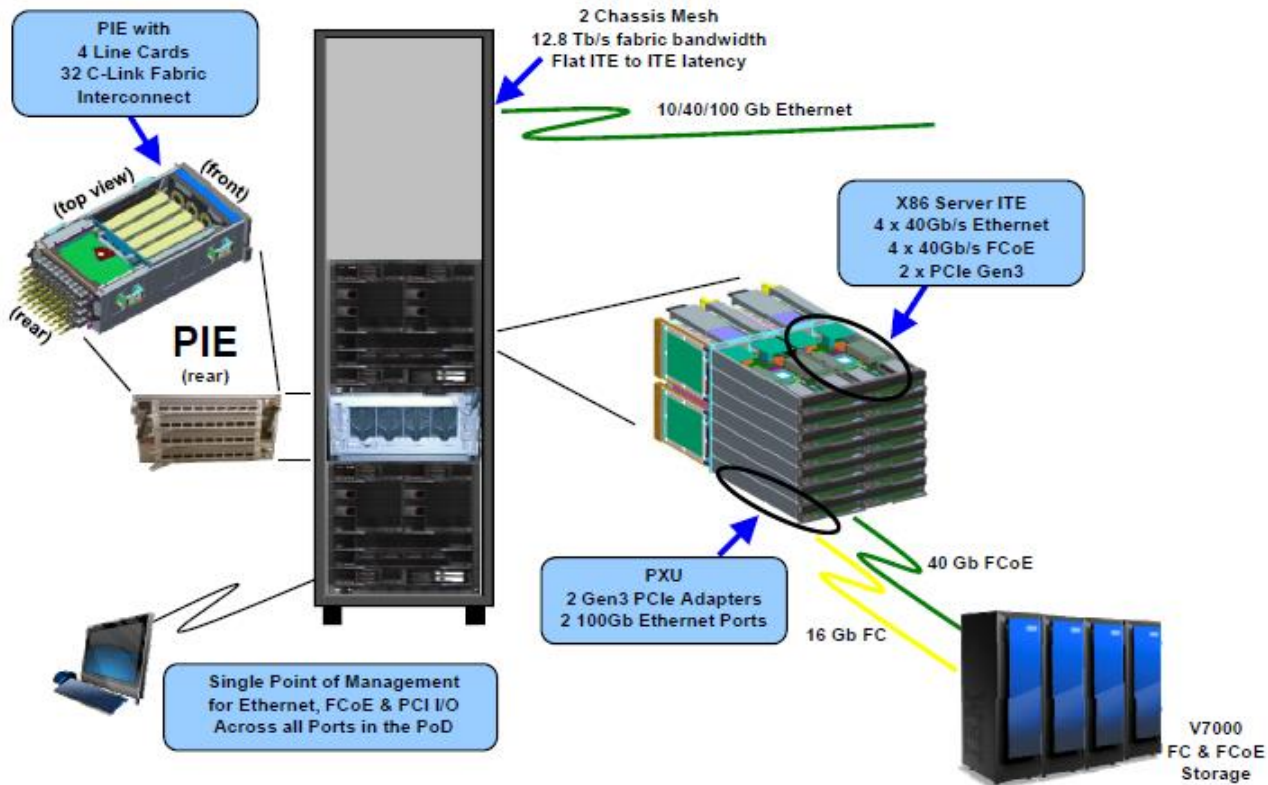


Nebula系统机框



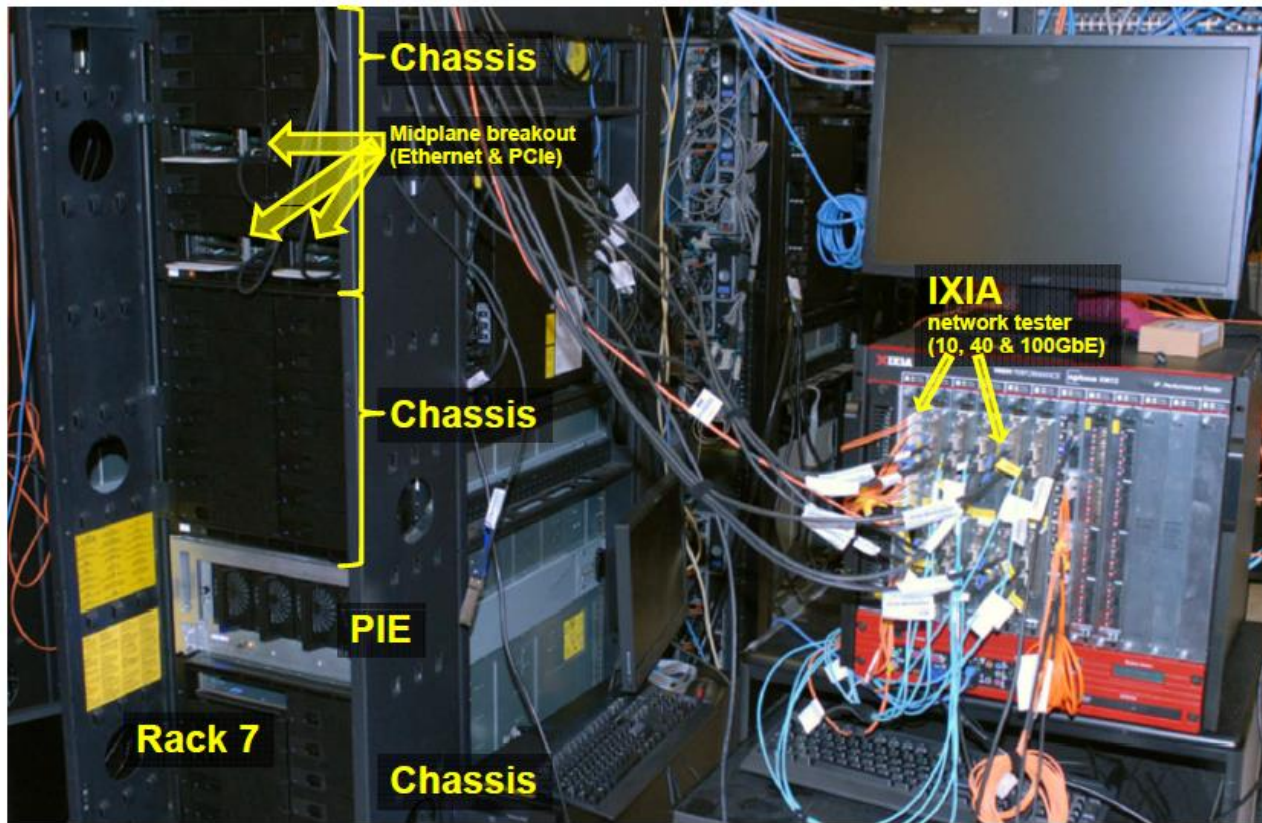


Nebula系统形态





Nebula原型测试系统





Thanks !