

# 目 录

1 前言 .....	1
2 云计算 .....	1
2.1 集中云.....	3
2.2 分散云.....	5
2.3 云计算小结.....	8
3 数据中心.....	9
3.1 Client 与 Server .....	9
3.2 层次化与扁平化.....	11
3.3 三层结构与两层结构.....	13
3.4 Server 与 Storage .....	14
3.5 数据中心多站点.....	16
3.6 多站点选择.....	18
3.7 数据中心小结.....	20
4 网络 .....	20
4.1 路由与交换.....	20
4.2 EOR 与 TOR .....	21
4.3 控制平面与转发平面.....	22
4.4 Box 与集中式转发.....	22
4.5 Chassis 与分布式转发.....	26
4.6 Clos 与 VOQ .....	32
4.7 网络小结.....	34
5 技术 .....	34
5.1 技术结构.....	35
5.2 网络虚拟化.....	36
5.2.1 网络多虚一技术.....	37
5.2.2 网络一虚多技术.....	39
5.3 技术理解.....	40
5.4 VM 本地互访网络技术.....	42
5.4.1 Cisco 接入层网络虚拟化 .....	43
5.4.2 802.1Qbg EVB.....	53
5.4.3 小结.....	59
5.5 Ethernet 与 FC 网络融合技术-FCoE.....	60
5.5.1 FC.....	60
5.5.2 FCoE.....	64
5.5.3 NPV.....	68
5.5.4 小结.....	70
5.6 跨核心层服务器二层互访 .....	71
5.6.1 控制平面多虚一技术.....	72
5.6.2 数据平面多虚一技术.....	75
5.6.3 控制平面一虚多技术.....	82
5.6.4 小结.....	85

5.7 数据中心跨站点二层网络.....	85
5.7.1 光纤直连.....	85
5.7.2 MPLS 核心网.....	87
5.7.3 IP 核心网.....	89
5.7.4 小结.....	94
5.8 数据中心多站点选择.....	95
5.8.1 GSLB .....	95
5.8.2 LISP .....	97
5.9 技术总结.....	102
6 终章 .....	102
7 感言 .....	104
8 文章原目录（来源 51CTO） .....	104

# 云计算数据中心网络技术全面剖析（图）

**摘要：**本文希望能够帮读者对云计算的数据中心的网络的技术建立起全面的结构性认识，因此除了总体思路的描述外，在介绍过程中也会力争用三言两语对前面部分中涉及的每个技术点都有所说明。

## 1 前言

题目并不吸引人，主要是作者犯懒，罗列了一下关键词而已，当然好处是一看就知道文章要说啥。

简单说下结构，首先讲讲云计算，其次是数据中心，再然后是网络，重点还是技术。内容是循序渐进的，可以理解前面每个词都是后面词的定语。

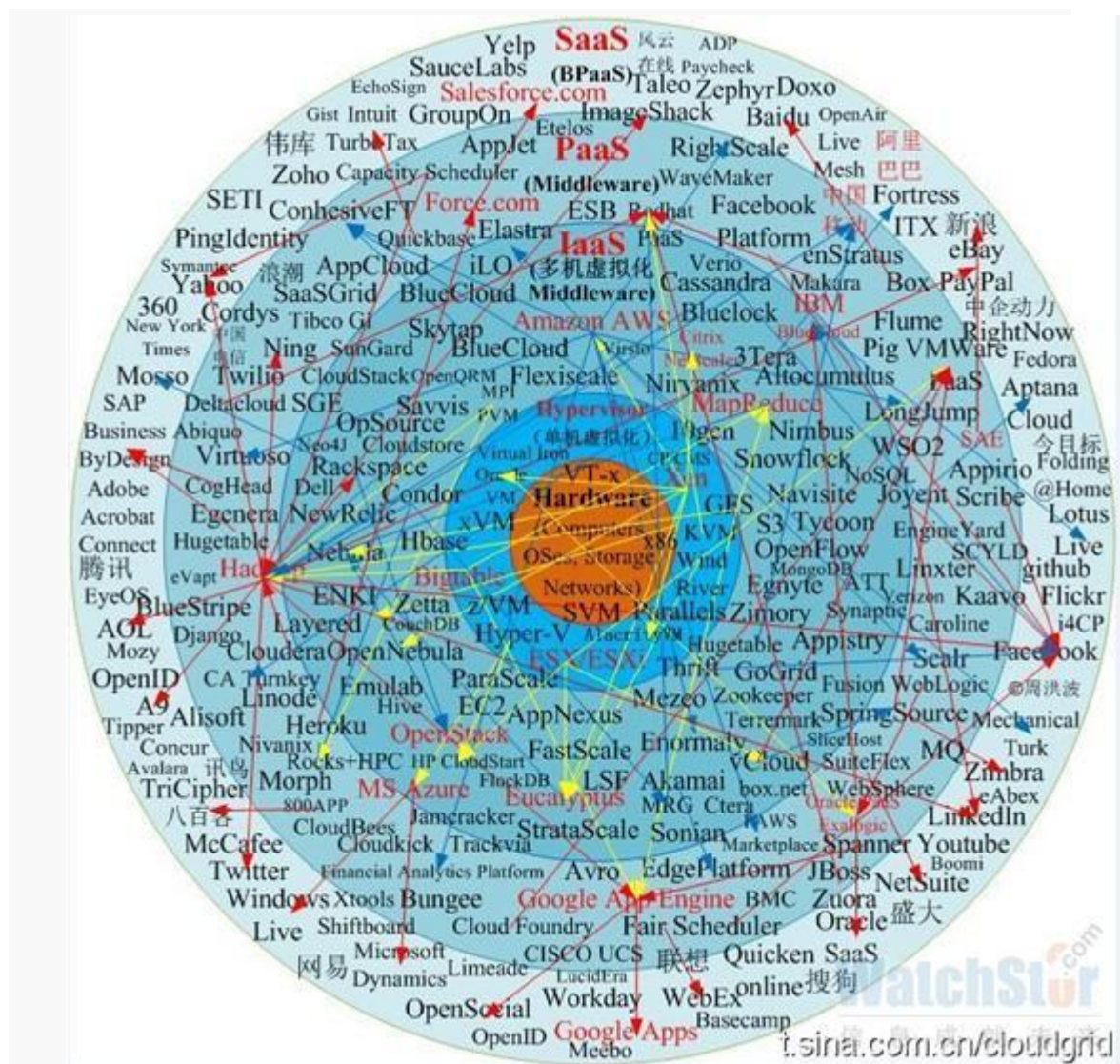
本文希望能够帮读者对云计算的数据中心的网络的技术建立起全面的结构性认识，因此除了总体思路的描述外，在介绍过程中也会力争用三言两语对前面部分中涉及的每个技术点都有所说明，至少让人明白这个东东怎么来的，要干啥和怎么干。但由于受篇幅所限，无法做到很详细，大家如果对某个技术点真感兴趣时，还是去网上找些更细节的资料来理解，本文是打算没有写成一本书的。

力争做到让文档读起来不感到枯燥吧，对作者来说那是相当有挑战的。

## 2 云计算

最早接触这个词好像是 06 年了，当时也是刚刚开始接触数据中心不久，这几年眼睁睁看着它被炒作得一塌糊涂，现在已经成为非常给力的一个概念。和别人谈数据中心要是不提云计算，你还真不好意思张这个嘴。

服务器厂商在喊云计算，网络、操作系统、应用软件甚至存储厂商都在喊。大家各喊各的，让我们感觉听上去都有那么点儿味道，但下来仔细一琢磨大都还在云里雾里。看看这张网上截取的云计算产业全景图，估计没有几个能够不头晕的。



云计算的各方面定义很多，基于用户的视角来看，目的就是让使用者在不需了解资源的具体情况下做到按需分配，将计算资源虚拟化为一片云。站在高处看，当前的主流云计算更贴切于云服务，个人认为可理解为早先运营商提供数据中心服务器租用服务的延伸。以前用户租用的是一台台物理服务器，现在租用的是虚拟机，是软件平台甚至是应用程序。公认的三个云计算服务层次是 IaaS（Infrastructure as a Service）、PaaS（Platform as a Service）和 SaaS（Software as a Service），分别对应硬件资源、平台资源和应用资源。对于用户来说：

- 1、当提供商给你的是一套 a 个核 CPU、b G 大小内存的主机、c M 带宽网络以及 d G 大小存储空间，需要你自己去装系统和搞定应用程序，那么这就是 IaaS，举例如 Amazon EC2；
- 2、当提供的是包含基本数据库和中间件程序的一套完整系统，但你还需要根据接口编写自己的应用程序时，那么就是 PaaS，举例如 Google AppEngine、Microsoft Azure 和 Amazon SimpleDB, SQS；

3、最傻瓜的方式自然是连应用程序都写好了，例如你只需要告诉服务提供商想要的是个 500 人的薪酬管理系统，返回的服务就是个 HTTPS 的地址，设定好帐号密码就可以访问过去直接使用，这就是 SaaS 了，如 SalesForce、Yahoo Hadoop 和 Cisco Webex: Collaboration SaaS 等。

服务属性	Amazon EC2	Google App Engine	Microsoft Azure	Yahoo Hadoop
架构	IaaS/PaaS	PaaS	PaaS	SaaS
服务形态	Compute/Storage	Web application	Web and non-web	Software
管理技术	OS on Xen hypervisor	Application container	OS through Fabric controller	Map / Reduce Architecture
使用者界面	EC2 Command-line tools	Web-based Administration console	Windows Azure portal	Command line and web
APIs	yes	yes	yes	yes
收费	yes	yes	yes	no
编程语言	AMI (Amazon Machine Image)	Python	.NET framework	Java, 成就未来

为啥举例都是国外的呢，因为国内目前的云服务状况是，能提供的都处于 IaaS 阶段，有喊着要做 PaaS 的，但还没听说有 SaaS 的。

说完公共的，该讲些私货了。

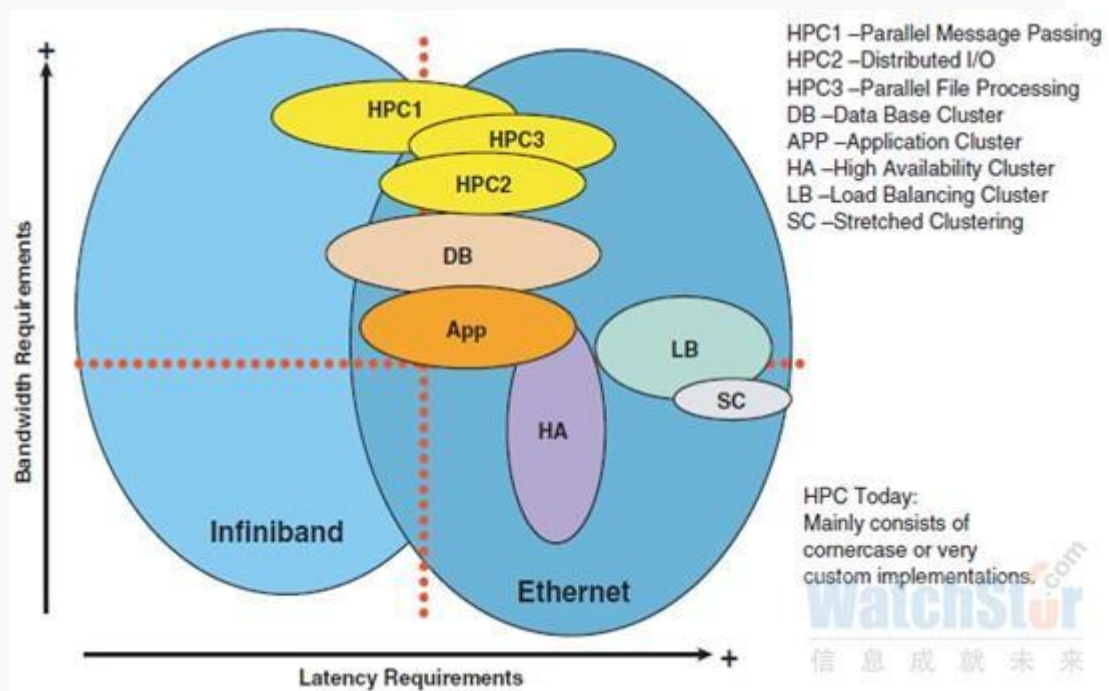
个人理解云计算的核心首先是计算，什么网络、存储、安全等等都是外延，从技术上讲云计算就是计算虚拟化。最早的云计算来自于网格计算，通过一堆性能较差的服务器完成一台超级计算机才能完成的计算任务，简单的说就是计算多虚一。但是现如今一虚多（VM/XEN 等）也被一些厂商扯着大旗给忽悠进来，并且成为主流。但是单从技术角度来看，这两者是南辕北辙的。因此云计算技术在下面被作者主观的分为集中云与分散云两个概念来阐述。

## 2.1 集中云

首先是集中云，根正苗红的多虚一，最早期的也是目前最大的一个典型实际用户就是 Google 了（注意这里说的不是现在 Google 云服务）。搜索引擎是超级消耗资源的典型应用，从你在网页上一个关键词的搜索点击，到搜索结果的产生，后台是经过了数百上千台服务器的统一计算。至于搜索引擎的工作模型本文就不多说了，网上很多资料的。随着互联网的发展，现在的开心、淘宝、新浪微博等等（好孩子不翻墙），虽然使用者看到的只是在简单的



页面进行点击输入，但是后台的工作量已经远远不是少量几台大型服务器能够胜任的了，即使天河一号也不见得能搞定。集中云的应用主力就是这些大型的互联网内容提供商们，当然还有一些传统应用如地震、气象和科研项目的计算也会存在此类需求。



了解了需求，下面简单谈下技术，上图是 Cluster 集群多虚一技术的简单分布，除了按照承载网络类型可分成 Infiniband 和 Ethernet 外，根据技术分，还可分为 Active-Standby 主备与 LoadBalance 负载均衡两类。

主备模式好理解，所有的 Server 里面只有一台干活，其他都是候着的，只有侦听到干活的歇菜了，才开始接管处理任务。主备模式大部分就二虚一提供服务，多了如三虚一什么的其实意义都不太大，无非是为了再多增加些可靠性。主备模式以各类 HA 集群技术为代表。

而负载均衡模式复杂一些，在所有的 LB 技术中都存在两个角色，协调者与执行者，协调者一般是一个或多个（需要主备冗余时），主要工作就是接活儿和分活儿（有点儿像包工头）；而执行者就只处理计算了，分到啥就完成啥，典型的苦力。从流量模型上来说，LB 集群技术有来回路径一致和三角传输两种，来回路径一致指流量都是客户发起连接，请求协调者进行处理，协调者分配任务给执行者进行计算，计算完成后结果会都返回到协调者，再由协调者应答客户。

这种结构简单，计算者不需要了解外界情况，由协调者统一作为内外接口，安全性最高。此模型主要应用于搜索和地震气象科研计算等业务处理中。三角传输模型指计算者完成计算后直接将结果反馈给客户，此时由于计算者会和客户直接通信，造成安全性降低，但返回流量减少了协调者这个处理节点，性能得到很大提升。此模型主要应用于腾讯新浪的新闻页面和阿里淘宝的电子商务等 WEB 访问业务。

集中云在云服务中属于富人俱乐部的范围，不是给中小企业和个人玩的，实际上都是各大互联网服务提供商自行搭建集中云以提供自己的业务给用户，不会说哪天雅虎去租用个 Google 的云来向用户提供自己的新闻页面访问。集中云服务可能的租用对象是那些高度科研项目，因而也导致当前集中云建设上升到国家宏观战略层面的地位。你能想象哪天百度的云服务提供给总装研究院去计算个导弹轨迹，核裂变什么嘛，完全不可能的事。

最后是多虚一对网络的需求。在集中云计算中，服务器之间的交互流量多了，而外部访问的流量相对减少，数据中心网络内部通信的压力增大，对带宽和延迟有了更高的要求，自然而然就催生后面会讲到的一些新技术（L2MP/TRILL/SPB 等）。

题外话，当前的多虚一技术个人认为不够给力，现在把 10 台 4 核 CPU 的服务器虚拟合一后，虚拟的服务器远远达不到一个 40 核 CPU 的计算能力。准确的说现在的多虚一只能基于物理服务器的粒度进行合并，理想的情况应该是能够精细到 CPU 核以及每台设备的内存缓存等等物理构件虚拟合一。这块应该就涉及到超算了，不熟不深谈。总的来说认为技术进步空间巨大，有些搞头。

## 2.2 分散云

再讲分散云，这块是目前的主流，也是前面提到的云服务的关键底层技术。由于有 VMware 和 Citrix 等厂家在大力推广，而且应用内容较集中云更加平民化，随便找台 PC 或服务器，装几个虚拟机大家都能玩一玩，想干点儿啥都成，也就使其的认知度更加广泛。

一虚多的最主要目的是为了提高效率，力争让所有的 CPU 都跑到 100%，力争让所有的内存和带宽都占满。以前 10 台 Server 干的事，我整两台 Server 每台跑 5 个虚拟机 VM (Virtual Machine) 就搞定了，省电省空间省制冷省网线，总之省钱是第一位的（用高级词儿就是绿色环保）。技术方面从实现方案来看，目前大致可分为三类：

### 操作系统虚拟化 OS-Level

在操作系统中模拟出一个个跑应用程序的容器，所有虚拟机共享内核空间，性能最好，耗费资源最少，一个 CPU 号称可最多模拟 500 个 VPS (Virtual Private Server) 或 VE (Virtual Environment)。缺点是操作系统唯一，如底层操作系统跑的 Windows，VPS/VE 就都得跑 Windows。代表是 Parallels 公司（以前叫 SWsoft）的 Virtuozzo（商用产品）和 OpenVZ（开源项目）。Cisco 的 Nexus 7000 猜测也是采用这种方案运行的 VDC 技术，但不太清楚为什么会有最多 4 个 VDC 的数量限制，也许是基于当前应用场景进行规格控制的一种商业手段。

### 主机虚拟化 Hosted

先说下 Hypervisor 或叫做 Virtual Machine Monitor (VMM)，它是管理虚拟机 VM 的软件平台。在主机虚拟化中，Hypervisor 就是跑在基础操作系统上的应用软件，与 OS-Level 中 VE 的主要区别在于：

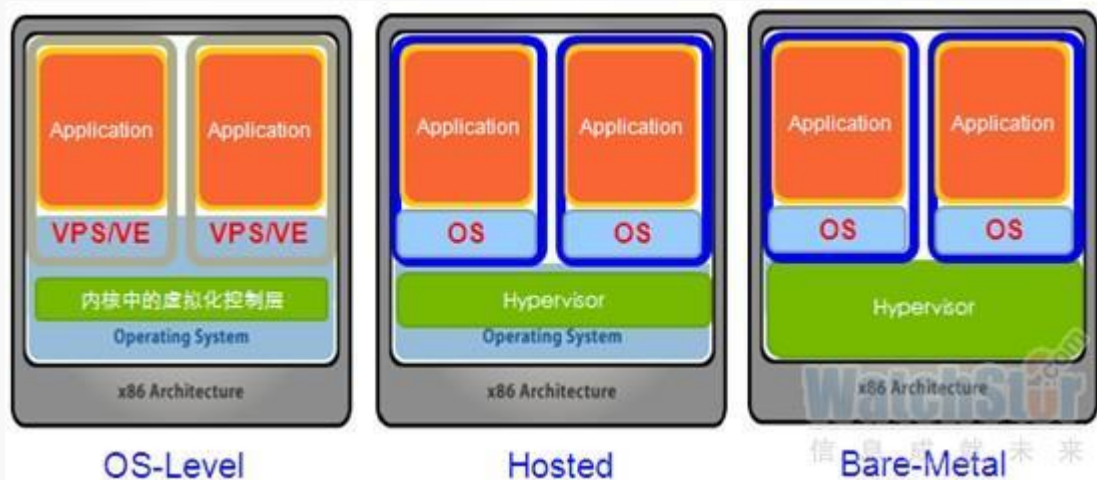
Hypervisor 构建出一整套虚拟硬件平台（CPU/Memory/Storage/Adapter），上面需要你再去安装新的操作系统和需要的应用软件，这样底层和上层的 OS 就可以完全无关化，诸如 Windows 上跑 Linux 一点儿问题没有；

VE 则可以理解为盗用了底层基础操作系统的资源去欺骗装在 VE 上的应用程序，每新创建一个 VE，其操作系统都是已经安装好了的，和底层操作系统完全一样，所以 VE 比较 VM（包括主机虚拟化和后面的裸金属虚拟化）运行在更高的层次上，相对消耗资源也少很多。

主机虚拟化中 VM 的应用程序调用硬件资源时需要经过:VM 内核->Hypervisor->主机内核，导致性能是三种虚拟化技术中最差的。主机虚拟化技术代表是 VMware Server（GSX）、Workstation 和 Microsoft Virtual PC、Virtual Server 等。

### 裸金属虚拟化 Bare-metal

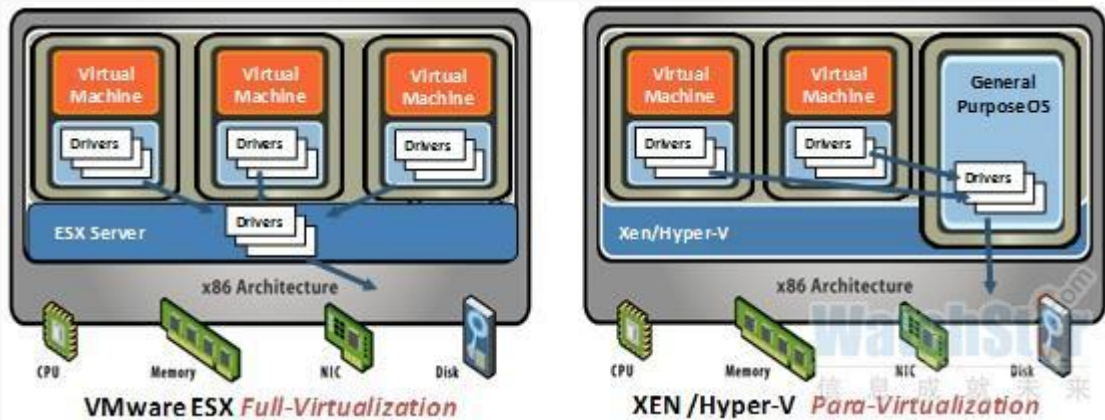
裸金属虚拟化中 Hypervisor 直接管理调用硬件资源，不需要底层操作系统，也可以理解为 Hypervisor 被做成了一个很薄的操作系统。这种方案的性能处于主机虚拟化与操作系统虚拟化之间。代表是 VMware ESX Server、Citrix XenServer 和 Microsoft Hyper-V。



上图描述了三种虚拟化方案的形态区别。当前分散云数据中心服务器虚拟化使用的主要是 Bare-Metal 方案。分散云给数据中心网络带来了新的挑战，虚拟机之间的数据通信管理需求促使了一系列网络新技术的发展。在 OS-Level 与 Hosted 方案中，虚拟机都是架设于操作系统之上的，因此 VM/VE 之间的通信主要由同样运行于基础操作系统之上的网络交换应用程序来完成。而在最主流的 Bare-Metal 结构中，由于 Hypervisor 薄操作系统的引入，性能、管理、安全和可靠性等多维度的考虑，造成 VM 间网络通信管理发展出不同的技术道路（EVB 与 BPE），后文会对这些技术方向加以详述。

VMware ESX 与 Xen/Hyper-V 的 Bare-Metal 方案实现结构有所不同，简单如下图所示。





分散云除了给网络带来上述的 VM 通信问题，同样由于其对服务器硬件能力的极端榨取，造成网络中的流量压力增大，与集中云一样存在着带宽扩展的需求。原本一台服务器一个操作系统跑一个应用只需要 10M 流量带宽就够了，现在装了 10 个 VM 跑 10 个应用，带宽可能就需要 100M 了。

大型机与小型机的一虚多技术早在 30 年前 IBM 就做出来了，现在 RISC 平台上已经相当完善了，相比较而言 X86 架构的虚拟化才处于起步阶段，但 X86 架构由于性价比更高成为了分散云计算的首选。

X86 架构最早期是纯软件层面的 Hypervisor 提供虚拟化服务，缺陷很多，性能也不够，直到 2006 年 Intel 推出了实现硬件辅助虚拟化的 VT 技术 CPU 产品后才开始迅猛发展（AMD 也跟着出了 VM 技术）。硬件辅助虚拟化技术主要包括 CPU/Chipset/Network Adapter 等几个方面，和网络技术紧密相关的就是网卡虚拟化，后文会对如 SR-IOV 等网卡虚拟化技术应用进行更具体分析。随着 2007 年 Intel VT FlexMigration 技术的推出，虚拟机迁移成为可能，2009 年 Intel 支持异构 CPU 间动态迁移再次向前迈进。

### vMotion

这里再多唠叨几句 vMotion 技术。vMotion 是 VMware 公司提出的虚拟机动态迁移技术名称（XEN 也有相应的 XENMotion 技术），由于此名称被喊得最早，范围最广，认知度最高，因此下文提到虚拟机迁移技术时大都会使用 vMotion 来代称。

先要明确 vMotion 是一项资源管理技术，不是高可靠性技术，如果你的某台服务器或 VM 突然宕机了，vMotion 是有助于应用访问进行故障切换和快速恢复的。vMotion 是将一个正常的处于服务提供中的 VM 从一台物理服务器搬家到另一台物理服务器的技术，vMotion 的目的是尽可能方便的为服务管理人员提供资源调度转移手段，当物理服务器需要更换配件关机重启啦，当数据中心需要扩容重新安排资源啦，这种时候 vMotion 就会有用武之地了。

设想一下没有 vMotion 上述迁移工作是怎么完成的，首先需要将原始物理服务器上的 VM 关机，再将 VM 文件拷贝到新的物理服务器上，最后将 VM 启动，整个过程 VM 对外提供的服务中断会达到几分钟甚至几小时的级别。而且需要来回操作两台物理服务器上的 VM，对管理人员来说也很忙叨。

使用 vMotion 后，两台物理服务器使用共享存储来保存 VM 文件，这样就节省了上述步骤 2 中的时间，vMotion 只需在两台物理服务器间传递当前的服务状态信息，包括内存和 TCP 等上层连接表项，状态同步的拷贝时间相对较短，而且同步时原始 VM 还可以提供服务使其不会中断。同步时间跟 VM 当前负载情况及迁移网络带宽有关，负载大了或带宽较低使同步时间较长时，有可能会使 vMotion 出现概率性失败。当状态同步完成后，原始物理服务器上的 VM 会关闭，而新服务器上的 VM 激活（系统已经在状态同步前启动完毕，始终处于等待状态），此时会有个较短的业务中断时间，可以达到秒级。再者 vMotion 是通过 VMware 的 vCenter 管理平台一键化完成的，管理人员处理起来轻松了许多。

这里要注意 vMotion 也一定会出现业务中断，只是时间长短区别，不要容易被一些宣传所忽悠。想想原理，不论怎么同步状态，只要始终有新建发生，在同步过程中原始服务器上新建立的客户连接，新服务器上都是没有的，切换后这部分连接势必被断开重建，零丢包只能是理想值。VMware 也同样建议将 vMotion 动作安排在业务量最少的时候进行。

vMotion 什么场景适用呢？首先肯定得是一虚多的 VM 应用场景，然后是对外业务中断恢复的可靠性要求极高，一般都是 7\*24 小时不间断应用服务才用得上，最后是计算节点规模始终在不断增长，资源调度频繁，管理维护工作量大的数据中心。

另外共享存储这个强制要求会给数据中心带来了整体部署上的限制，尤其是下面提到的跨数据中心站点 vMotion 时，跨站点共享存储的问题解决起来是很麻烦的，由于这部分内容和网络关系不大，属于存储厂商的地盘，对跨站点共享存储技术有兴趣的读者可以参考 EMC/IBM 等厂商的资料看看，本文就不过多介绍了。

vMotion 的出现推动了数据中心站点间大二层互联和多站点动态选路的网络需求，从而导致 OTV 和 LISP 等一系列新技术的出现。

## 2.3 云计算小结

通过前面的描述，希望大家能对云计算有个较为清晰的概念。云计算还有一大块内容是平台管理资源调度方面（目前很多厂家吆喝的云计算都是云平台）。这部分主要针对客户如何更便捷的创建与获取虚拟化服务资源，实际过程就是用户向平台管理软件提出服务请求，管理平台通过应用程序接口 API（Application Program Interface）将请求转化为指令配置下发给服务器、网络、存储和操作系统、数据库等，自动生成服务资源。需要网络做的就是设备能够识别管理平台下发的配置，从技术创新的角度讲，没有啥新鲜东西，就不多说了。当

前的云平台多以 IaaS/PaaS 为主，能做到提供 SaaS 的极少。但在今后看来，SaaS 将会成为云服务租用主流，中小企业和个人可以节省出来 IT 建设和维护的费用，更专注于自身的业务发展。

**总结一下云计算给数据中心网络带来的主要变化：**

- 1、 更高的带宽和更低的延迟
- 2、 服务器节点（VM）规模的增加
- 3、 VM 间通信管理
- 4、 跨数据中心站点间的二层互联以承载 vMotion

题外再多说两句，计算虚拟化中一虚多与多虚一结合使用才是王道。但目前云计算服务提供商能够提供的只是先将物理服务器一虚多成多台 VM，再通过 LB/集群计算等技术将这些 VM 对外多虚一成一个可用的资源提供服务。个人感觉，如果能做到先将一堆物理服务器虚拟成一台几万个核 Super Computer，用户再根据自己的需要几个几十个核的自取资源，这样才更有云计算的样子， Super Computer 就是那朵云。当然计算虚拟化的时候不光是核的调配，还要包括 IO/Memory 等一起进行调度，这里只是简单举例。

## 3 数据中心

数据中心的产生有多早？从人类开始将信息记录到介质上传递开始就有了数据中心，那个记载信息的介质（石头或树皮）就是数据中心，不过那时的网络是靠手手相传而已。如果更甚一些，可以理解人类产生语言开始，知识最多的人（酋长/祭祀）就是数据中心，口口相传就相当于现如今的网络传输。有人该说，夸张了哈，写作手法而已，只是想突出一下数据中心的重要性。

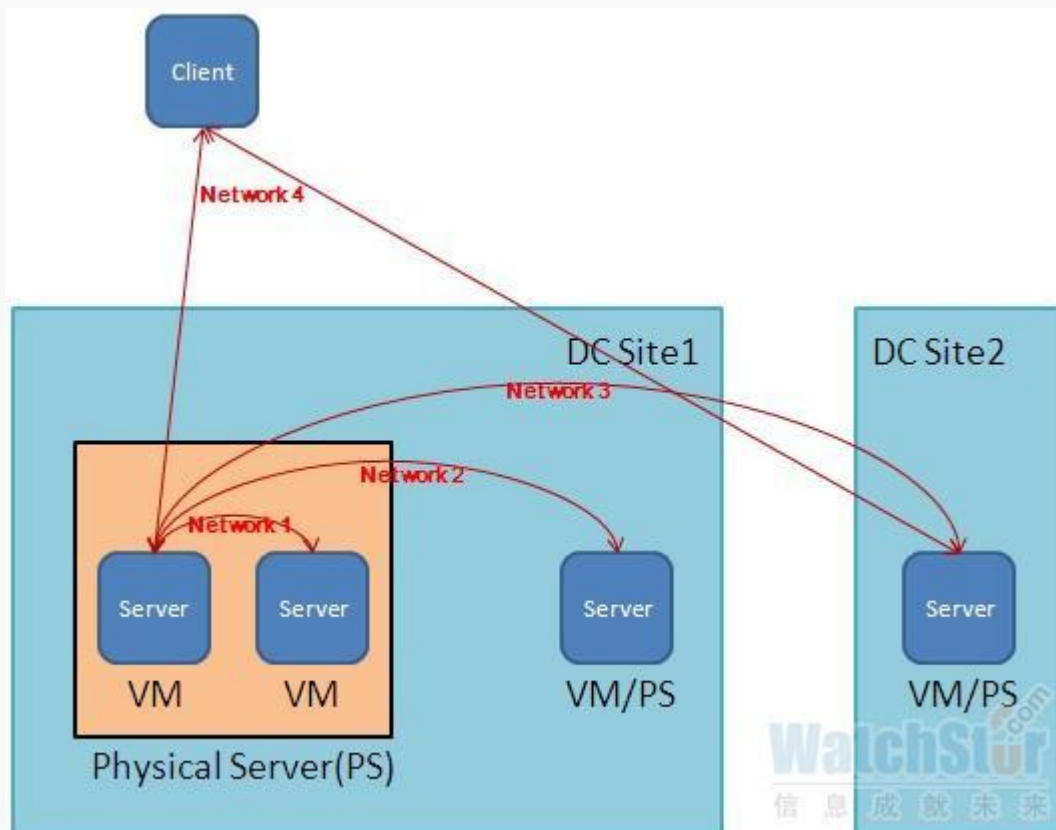
当计算机网络连接到 Server 的那一刻起，整个世界的网络就从网状变成了树状，一个个数据中心就是网络世界的根。

### 3.1 Client 与 Server

在所有的数据通信会话中，只有两个永恒的角色，Client 与 Server。为了下文叙述简便，作者把数据中心内部的终端统一称之为 Server，数据中心外部的为 Client。这样网络间的流量通信就只剩下 Client-Server（CS）与 Server-Server（SS）两种了。其实更准确说还是只有 CS 一种，SS 通信也是有个发起方和响应方的。QQ/MSN 等即时通信软件的流量模型实际可理解为 CSC；唯有 P2P 对 CS 结构有所颠覆，但不管怎么处理也必定会存在 Server 角色进行最初的调度。

所有数据中心需要处理的业务就是 CS 和 SS 两种,CS 肯定是基于 IP 进行 L3 转发的了,SS 则分为基于 IP 的 L3 和基于 MAC 的 L2 两种转发方式。基于 IP 的 SS 通信主要是不同业务间的数据调用,如 WEB/APP 服务器去调用 DB 服务器上的数据,再如有个员工离职,职工管理系统会同步通知薪酬管理、考勤管理、绩效管理等一系列系统进行删除信息的相关操作。基于 MAC 的 SS 通信则是同一类服务器间的数据同步计算,比如使用 WEB 集群分流用户访问时,需要对修改或增删的数据进行集群同步;再比如多虚一中集群一起计算任务时协调者和执行者之间的大量通信进行任务调度。

可以看出云计算数据中心给网络带来的挑战主要是基于 MAC 的二层 (OSI 模型) SS 通信。在一虚多技术影响下,Server 的概念已经扩展到以单台 VM 为基础单元,因此可以引出下面这个图,看看新网络技术是如何划分的。



Network1: VM 到 VM 之间的 SS 二层互连网络

Network2: DC 站点内部 SS 二层互连网络

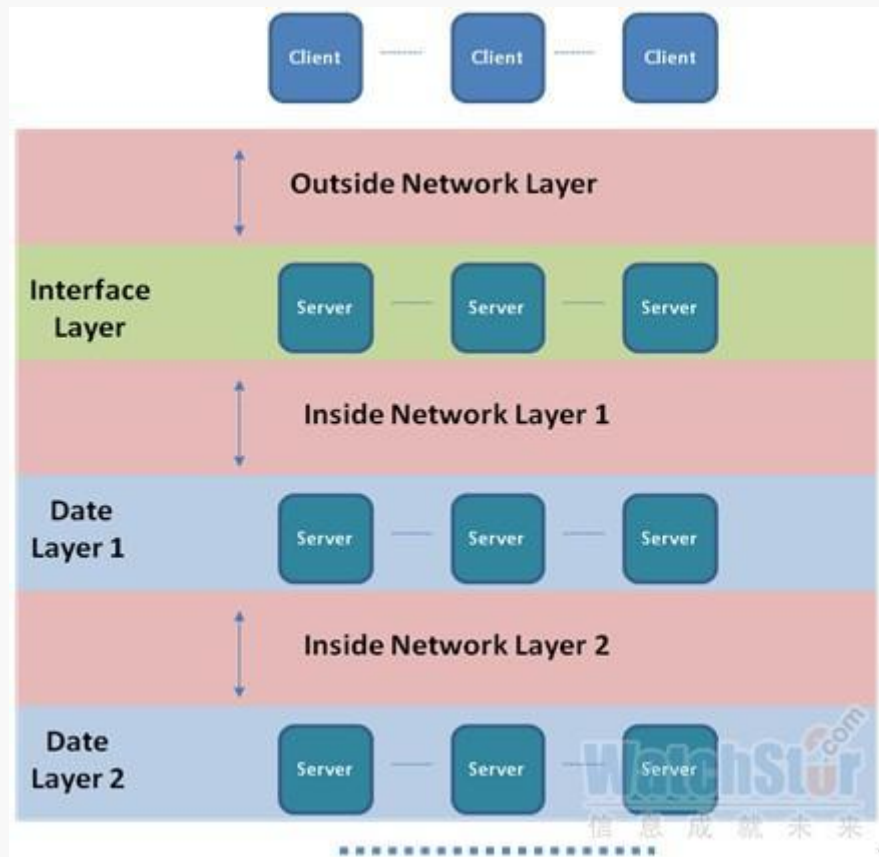
Network3: 跨 DC 站点间的 SS 二层互连网络

Network4: DC 到 Client 之间的 CS 三层互连网络

后文的技术章节就会针对这些部分进行展开,详细说下都有哪些技术分别对应在这四段网络中,这些技术的特点是什么。

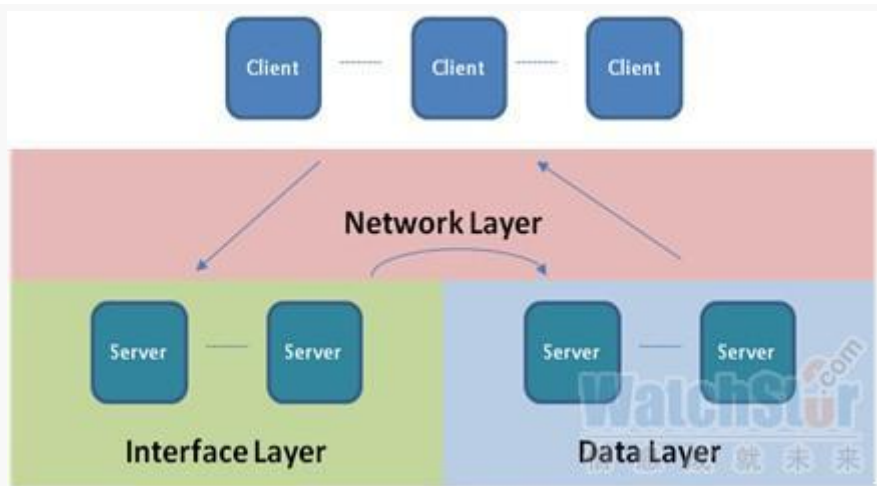
## 3.2 层次化与扁平化

数据中心的网络结构取决于应用计算模型，计算模型主要分为层次化与扁平化两种结构。层次化结构如下图所示，典型的应用如 WEB-APP-DB、搜索引擎或高性能计算（地震、科研）等。特点是客户请求计算结果必须逐层访问，返回数据也要逐层原路返回。

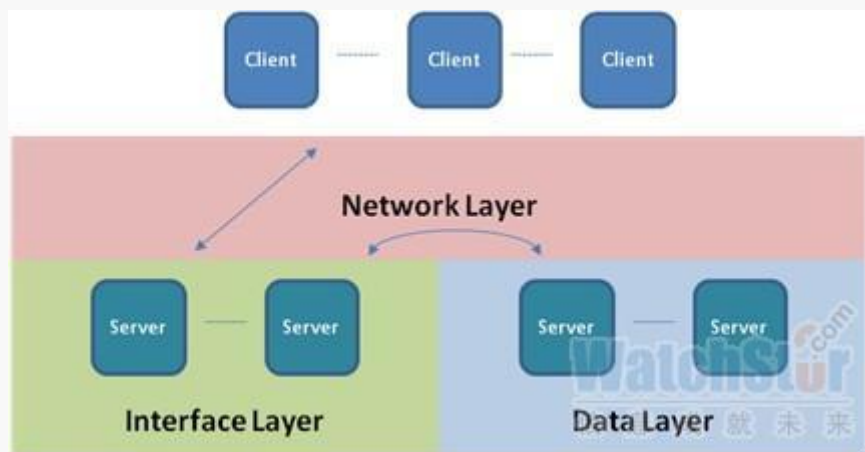


计算模型扁平化结构如下图所示，特点是数据层服务器会将结果直接返回给客户，不需要再由接口层服务器进行处理，也有管这种模型叫做三角传输的。典型的应用如一些 Internet 网站服务采用的 LB 结构，LB 服务器就是只做调度，WEB 服务器会直接将请求结果返回给用户。





注意，上面说的是计算模型，和网络模型并不是一一对应，采用层次化结构计算模型一样可以进行扁平化组网，如下图所示。

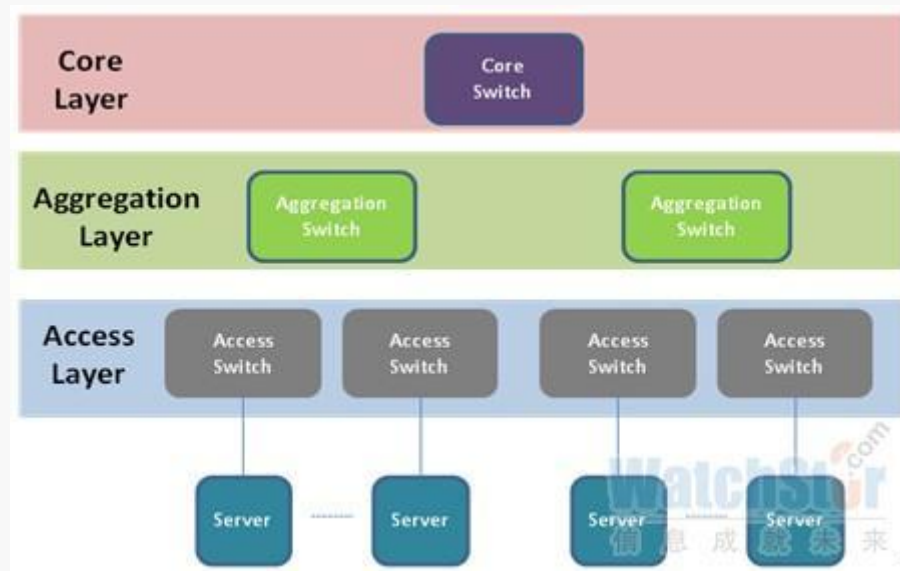


从网络角度讲，扁平化相比较层次化结构最大的好处是可以减少服务器的网卡接口数量（省钱），然而缺点是没有清晰的层次，部署维护的复杂度就会相应提升。总体来说，当前数据中心实际组网建设中，这两种方式谁都没占据到绝对优势，上哪种结构完全看规划者的考量重点是在哪个方面。

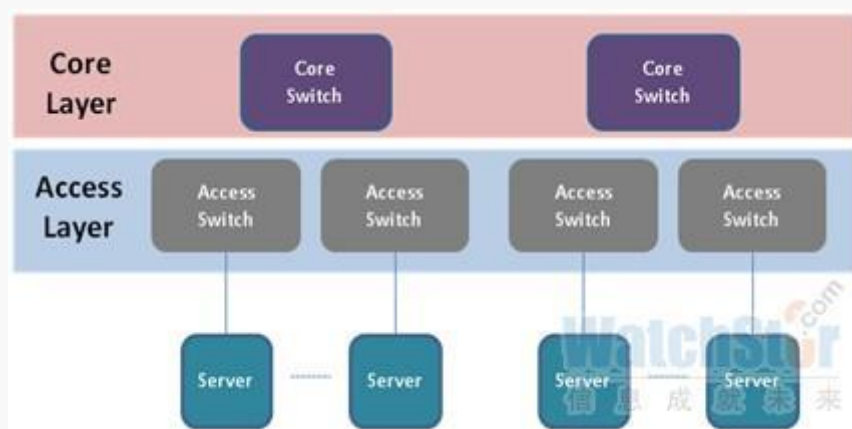
前面说过，云计算主要分为多虚一与一虚多两种虚拟化结构。一虚多对传统计算模型没有太大影响，只是将其服务器从物理机到虚拟机数量规模扩大了  $N$  倍，网络规模也随之进行扩大。而多虚一中，协调者角色对应了接口层服务器，执行者角色则对应数据层服务器，由于此时大量的通信交互是在不同执行者之间或执行者与协调者之间，需要重点关注的大规模网络就由原来的接口层服务器之前，转移到了接口层服务器与数据层服务器之间。

### 3.3 三层结构与两层结构

在以往的数据中心网络建设时，关注的重点都是指接口层服务器前的网络，传统的三层网络结构如下图所示。其中的汇聚层作为服务器网关，可以增加防火墙、负载均衡和应用加速等应用优化设备。

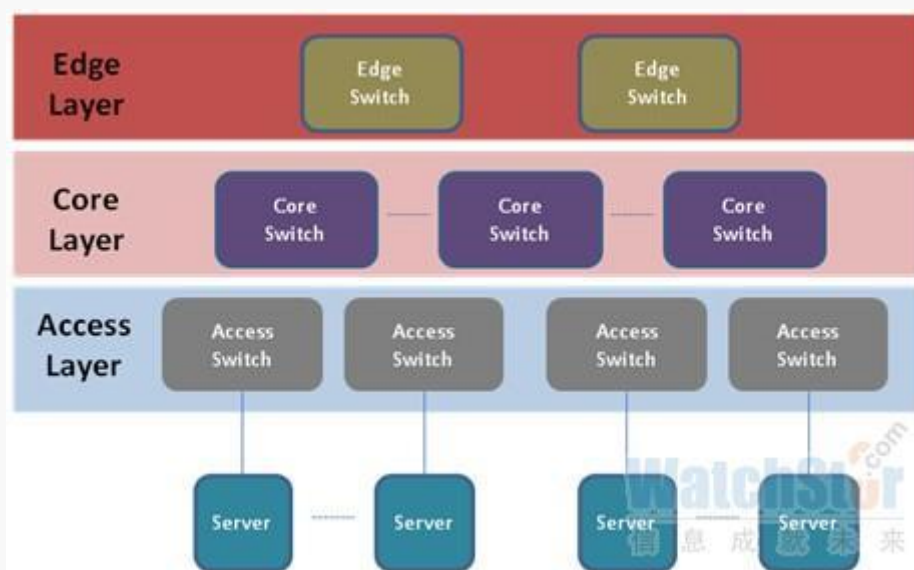


但在云计算数据中心里面 Ethernet 网络规模扩大，流量带宽需求增加，因此不会在网络中间位置再插入安全和优化设备了，转发性能太低，上去就是瓶颈，汇聚层的位置也就可有可无了。再加上带宽收敛比的问题，短期内大型云计算数据中心网络里面不会出现汇聚层的概念。以前是百兆接入、千兆汇聚、万兆核心，现在服务器接入已经普及千兆向着万兆迈进了，除非在框式交换机上 40G/100G 接口真的开始大规模部署，还有可能在云计算数据中心里面再见到超过两层的级联结构网络。现如今的云计算数据中心流行的都是如下图所示的千兆接入，万兆核心的两层网络结构。

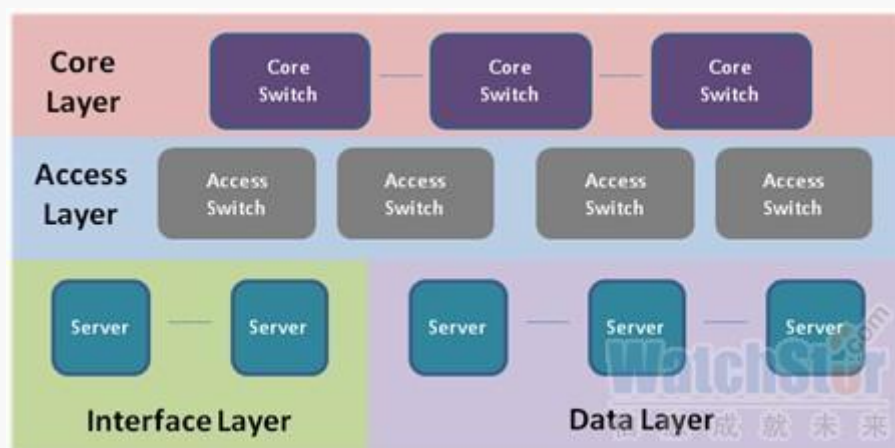


此两层网络结构部署在接口层服务器之前，则一般会将服务器网关部署在 Core Switch 上，但前提是网络规模不会太大，Core 不会太多（2 个就差不多了），否则 VRRP/HSRP 等

多网关冗余协议只能走到一个活动网关，会导致网络效率很低。还有一种方式是将服务器网关部署在 Access Switch 上，Access SW 与 Core SW 之间通过 OSPF 等动态路由协议达到全互联，使用等价路由达到多 Core SW 的负载均担。但此方式的缺点是 L3 路由交互转发效率较低，部署复杂且占用大量 IP 地址。在未来的 TRILL/SPB 等二层 Ethernet 技术结构中，可能会出现专门作为网关与外部进行 IP 层面通信的边缘交换机（由于出口规模有限，2-4 台足够处理），内部的 Core SW 只做二层转发，可以大规模部署以满足内部服务器交互的需求，如下图所示。



当遇到多虚一高性能计算的模型，则二层网络结构会被部署在接口服务器与数据服务器之间，为二者构建纯二层的大规模交互网络，结构如下图所示。



### 3.4 Server 与 Storage

前面说的 CS/SS 网络可以统称为数据中心前端网络，目前和以后基本上都是 IP+Ethernet 一统天下（IB Infiniband 只能吃到高性能计算的一小口）。有前端当然就有后端，在数据中

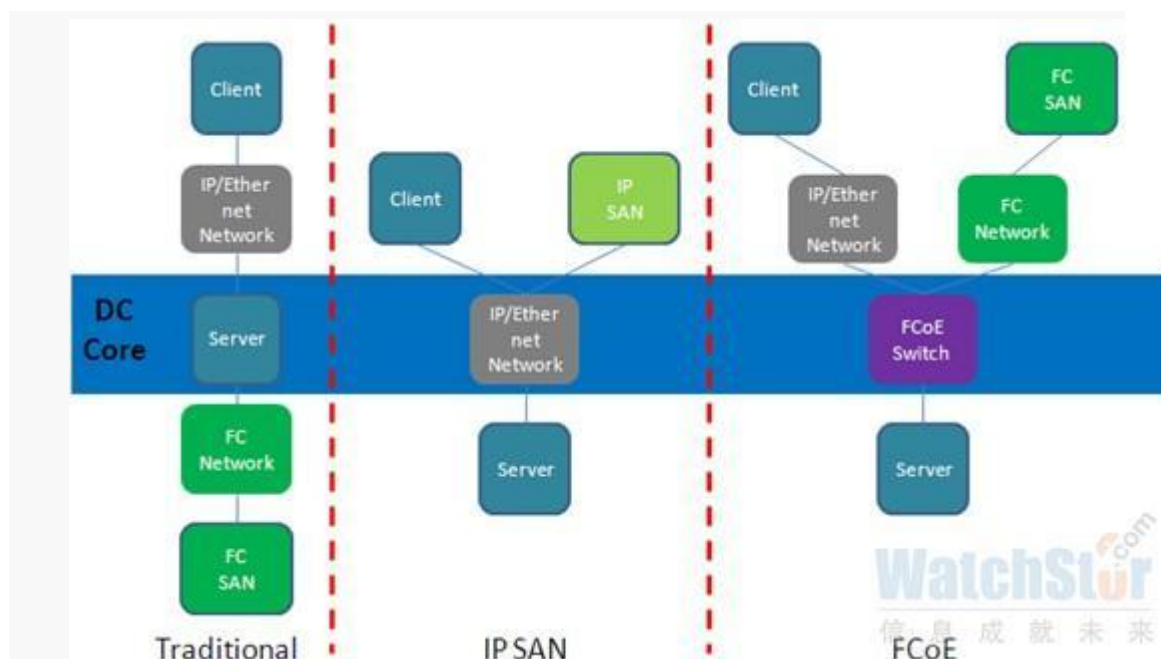
心里面，服务器与存储设备连接的网络部分统称为数据中心后端网络。就目前和短期的未来来看，这块儿都是 FC 的天下。

简单说两句存储，DAS (Direct Attached Storage) 直连存储就是服务器里面直接挂磁盘，NAS (Network Attached Storage) 则是网络中的共享文件服务器，此二者大多与数据中心级别存储没什么关系。只有 SAN (Storage Area Network) 才是数据中心存储领域的霸主，磁盘阵列会通过 FC 或 TCP/IP 网络注册到服务器上模拟成直连的磁盘空间。而目前 FC SAN 是主流中的主流，基于 TCP/IP 的 IP SAN 等都是配太子读书的角色。

在服务器到存储的后端网络中，涉及到 IO 同步问题，高速、低延迟与无丢包是对网络的基本需求，而 Ethernet 技术拥有冲突丢包的天然缺陷，FC 的无丢包设计使其领先一步，加上早期 Ethernet 还挣扎在 100M 带宽时，FC 已经可以轻松达到 2G，所以在后端网络中从开始到现在都是 FC 独占鳌头。但是从发展的眼光看，Ethernet 目前已经向着 40G/100G 迈进，而 FC 的演进并不理想，无论是 BASE10 的 10/20/40G 路线（主要用在 FC 交换机之间，目前基本已经被淘汰）还是 BASE2 的 2/4/8/16/32G 路线（当前主流 FC 应用）都已经落后，加上各种以太网零丢包技术（CEE/DCE/DCB）的出现，以后鹿死谁手还真不好说。

在目前阶段，为了兼容数据中心已有的主流 FC 网络和存储设备，在基于 iSCSI 技术的 IP SAN 技术没能开花结果的情况下，众多 Ethernet 网络厂商又推出了 FCoE 来蚕食服务器到存储这块蛋糕。下文技术章节会专门介绍 FCoE 的内容。

先简单说下，FCoE 没有惦着像 IP SAN 那样一下子完全取代 FC 去承载后端网络，而是走前后端网络融合，逐步蚕食的路线，是网络厂商们将数据中心的核​​心由服务器向网络设备转移的重要武器。如下图，就是看谁做太阳，谁做星星。相比较 IP SAN 的壮烈牺牲，FCoE 采用了一条更为迂回的兼容道路，目前已经出现了支持 FCoE 的存储设备，也许 Ethernet 完全替代 FC 的时代真的能够到来。



如果 FCoE 能成功，虽然短期内交换机、服务器和存储的价格对比不会有太大的变化，但是占据了核心位置，对未来的技术发展就有了更大的话语权，附加值会很高。又如当前的 EVB（Edge Virtual Bridging）和 BPE（Bridging Port Extend）技术结构之争也同样是话语权之争。

顺便一提，当一项完全不能向前兼容的全新技术出现时，除非是有相当于一个国家的力量去推动普及，而且原理简单到 8-80 岁都一听就明白，否则注定会夭折，与技术本身优劣无太大关系。老话说得好，一口吃不成胖子。

### 3.5 数据中心多站点

这是个有钱人的话题，且符合 2-8 原则，能够建得起多个数据中心站点的在所有数据中心项目中数量也许只能占到 20%，但他们占的市场份额肯定能达到 80%。

建多个数据中心站点主要有两个目的，一是扩容，二是灾备。

#### 扩容

首先说扩容，一个数据中心的服务器容量不是无限的，建设数据中心时需要考虑的主要因素是空间、电力、制冷和互联。数据中心购买设备场地建设只是占总体耗费的一部分，使用过程中的耗能维护开销同样巨大，以前就闹过建得起用不起的笑话。当然现在建设时要规范得多，考虑也会更多，往往做预算时都要考虑到 10 年甚至以上的应用损耗。

再讲个故事，以前曾有某大型 ISP 打算找个雪山峡谷啥的建数据中心，荒郊野外空间本来就大，融雪制冷，水力发电，听上去一切都很美，但是就忘了一件事，互联。光纤从哪里拉过去，那么远的距离中间怎么维护，至少从目前阶段来说这个问题无解。也许等到高速通



信发展到可以使用类似铱星的无线技术搞定时，数据中心就真的都会建到渺无人烟的地方吧，现在还只能在城市周边徘徊。貌似听说过国外有建得比较偏远的大型数据中心，但个人觉得应该还是人家通信行业发达，光纤资源丰富，四处都能接入。但至少目前国内的运营商们不见得会支持，大城市周边搞搞就算了，远了没人会陪你玩。

有些扯远，回到正题。现在国内已经有超过 10k 台物理服务器在一个数据中心站点的项目了，再多我还没有听说过。只有几百上千的物理服务器就敢喊搞云计算是需要一定勇气的，既然是云，规模就应永无止境。所以建多个数据中心站点来扩容就成了必然之举。这时就可能遇到 Cluster 集群计算任务被分配在多个站点的物理服务器或虚拟机来完成的情况，从而提出了跨多个数据中心站点的 Ethernet 大二层互联需求。

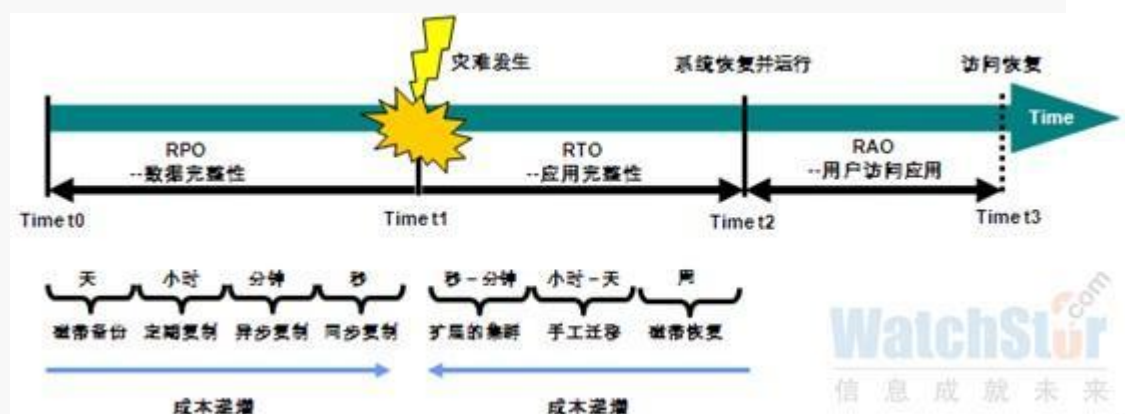
在扩容时，就可以充分利用 vMotion 等虚拟机迁移技术来进行新数据中心站点的建设部署，同样需要站点间的大二层互通。支持 IP 层的 vMotion 目前虽然已经出现，但由于技术不够成熟，限制很多，实用性不强，还是以 Ethernet 二层迁移技术为主。

## 灾备

再说说灾备，最近几年天灾人祸着实不少，数据中心容灾就越来越受到重视。扩容和灾备的主要区别就是：扩容的多个站点针对同一应用都要提供服务；而灾备则只有主站点提供服务，备份站点当主站点挂掉的时候才对外服务，平时都处于不运行或者空运行的状态。

参考国标《信息系统灾难恢复规范》GB/T 20988—2007，灾备级别大致可划分为数据级别（存储备份），应用级别（服务器备份），网络级别（网络备份），和最高的业务级别（包括电话、人员等所有与业务相关资源）。

国内外统一的容灾衡量标准是 RPO（Recovery Point Objective）、RTO（Recovery Time Objective）和 RAO（Recovery Access Objective）了，通过下图形象一些来体现他们的关系。



简单来说 RPO 衡量存储数据恢复，RTO 衡量服务器应用恢复，RAO 衡量网络访问恢复。一般来说 RPO 设计都应小于 RTO。国外按照 RTO/RPO 的时间长短对灾难恢复分级参考由高到低为：

Class 1/A RTO=0-4 hrs; RPO=0-4 hrs

Class 2/B RTO=8-24 hrs; RPO=4 hrs

Class 3/C RTO=3 day; RPO=1 day

Class 4/D RTO=5+ days; RPO=1 day

标准归标准，真正建设时候最重要的参考条件还是应用的需求，像银行可以直接去调研储户能容忍多长时间取不出来钱，腾讯去问问 QQ 用户能容忍多长时间上不了线，就都知道该怎么设计容灾恢复时间了。

真正在玩多中心灾备的行业，国内集中在金融系统（尤其是银行），政府和能源电力等公字头产业，国外的不太清楚，但我想以盈利为主要目的企业不会有太强烈意愿去建设这种纯备份的低效益站点，更多的是在不同站点内建设一些应用服务级别的备份，所有站点都会对外提供服务。

### 小结

在云计算规模的数据中心中，对于 LB 类型的多虚一集群技术，执行者（概念参见文档前面集中云部分）少上几个不会影响全局任务处理的，只要在扩容时做到数据中心间大二层互通，所有站点内都有计算任务的执行者，并且配合 HA 技术将协调者在不同站点做几个备份，就已经达到了应用容灾的效果。针对一虚多的 VM 备份，VMware/XEN 等都提出了虚拟机集群 HA 技术，此时同样需要在主中心站点与备份中心站点的服务器间提供二层通道以完成 HA 监控管理流量互通，可以达到基于应用层面的备份。

云计算数据中心多站点主要涉及的还是扩容，会部署部分针对 VM 做 HA 的后备服务器，但是不会搞纯灾备站点。针对多站点间网络互联的主要需求就是能够做而二层互联，当站点数量超过两个时所有站点需要二层可达，并部署相关技术提供冗余避免环路。

## 3.6 多站点选择

数据中心建设多站点后，由于同一应用服务可以跑在多个站点内部，对 Client 来说就面临着选择的问题。

首先要记住的是一个 Client 去往一个应用服务的流量必须被指向一台物理或虚拟的 Server。你可以想象一个 TCP 请求的 SYN 到 ServerA，而 ACK 到了 ServerB 时，ServerA 和 B 为了同步会话信息都会疯掉。想办法维持一对 Client-Server 通信时的持续专一是必须的。

Client 到 Server 的访问过程一般分为如下两步：

- 1、 Client 访问域名服务器得到 Server IP 地址（很少人会去背 IP 地址，都是靠域名查找）
- 2、 Client 访问 Server IP，建立会话，传递数据。

当前的站点选择技术也可以对应上面两个步骤分为两大类。

第一类是在域名解析时做文章，原理简单来说就是域名服务器去探测多个站点内 IP 地址不同的服务器状态，再根据探测结果将同一域名对应不同 IP 返回给不同的 Client。这样一是可以在多个 Client 访问同一应用时，对不同站点的服务器进行负载均衡，二是可以当域名服务器探测到主站点服务器故障时，解析其他站点的服务器 IP 地址给 Client 达到故障冗余目的。这时要求不同站点的服务地址必须在不同的三层网段，否则核心网没法提供路由。缺点很明显，对域名解析服务器的计算压力太大，需要经常去跟踪所有服务器状态并 Hash 分配 Client 请求的地址。此类解决方案的代表是 F5/Radware/Cisco 等厂商的 3DNS/GSLB/GSS 等技术。

第二类就是把多个站点的服务 IP 地址配置成一样，而各个站点向外发布路由时聚合成不同位数的掩码（如主中心发布/25 位路由，备中心发布/24 位路由），或通过相同路由部署不同路由协议 Cost 值以达到主备路由目的。使用掩码的问题是太细则核心网转发设备上的路由数量压力大，太粗则地址使用不好规划很浪费。使用 Cost 则需要全网 IP 路由协议统一，节点规模受到很大限制。另外这种方式只能将所有 Client 访问同一服务 IP 的流量指向同一个站点，负载分担只能针对不同的服务。好处则是这种站点选择技术谁都能用，不需要专门设备支持，部署成本低成为其存活的根据。

在云计算大二层数据中心部署下，各个站点提供同一服务的 Server 都处于一个二层网络内，且不能地址冲突，与前面描述的两种站点选择技术对服务器 IP 设置要求都不匹配，因此需要配合 SLB 设备一起使用。可以理解其为一种基于 IP 粒度的多虚一技术，使用专门 LB 硬件设备作为协调者，基于 IP 地址来分配任务给服务组中不同的 Server 执行成员。LB 设备通常将多个 Server 对应到一个 NAT 组中，外部访问到一个 NAT Server 虚拟 IP 地址，由 LB 设备按照一定算法分担给各个成员。LB 设备同时会探测维护所有 Server 成员状态。当各个站点内 LB 设备将同一服务对外映射为不同的虚拟 IP 地址时，可以配合域名解析方式提供 Client 选路；而配置为相同时则可以配合路由发布方式使用。

现有的站点选择技术都不尽如人意，即使是下文介绍的 Cisco 新技术 LISP 也只是部分的解决了路由发布技术中，发布服务器地址掩码粒度过细时，给核心网带来较大压力的问题，目前还不算是一套完整的站点选择解决方案。个人感觉，最好的路还是得想法改造 DNS 的处理流程，目前的 DNS 机制并不完备，在攻击面前脆弱不堪，后面的安全附加章节中会对此再深入讨论。

## 3.7 数据中心小结

又到了小结部分，云计算数据中心相比较传统数据中心对网络的要求有以下变化：

- 1、 Server-Server 流量成为主流，而且要求二层流量为主。
- 2、 站点内部物理服务器和虚拟机数量增大，导致二层拓扑变大。
- 3、 扩容、灾备和 VM 迁移要求数据中心多站点间大二层互通。
- 4、 数据中心多站点的选路问题受大二层互通影响更加复杂。

题内话，FCoE 并不是云计算的需求，而是数据中心以网络为核心演进的需求，至于云计算里面是不是一定要实现以网络为核心，就看你是站在哪个设备商的角度来看了。

## 4 网络

说到网络，这里关注的重点是前文提到的数据中心内部服务器前后端网络，对于广泛意义上的数据中心，如园区网、广域网和接入网等内容，不做过多扩散。

### 4.1 路由与交换

网络世界永远的主题，至少目前看来还没有出现能取代这二者技术的影子，扩展开足够写好几本书的了。

数据中心的网络以交换以太网为主，只有传统意义的汇聚层往上才是 IP 的天下。参考前文的需求可以看出，数据中心的以太网网络会逐步扩大，IP 转发的层次也会被越推越高。

数据中心网络从设计伊始，主要着眼点就是转发性能，因此基于 CPU/NP 转发的路由器自然会被基于 ASIC 转发的三层交换机所淘汰。传统的 Ethernet 交换技术中，只有 MAC 一张表要维护，地址学习靠广播，没有什么太复杂的网络变化需要关注，所以速率可以很快。而在 IP 路由转发时，路由表、FIB 表、ARP 表一个都不能少，效率自然也低了很多。

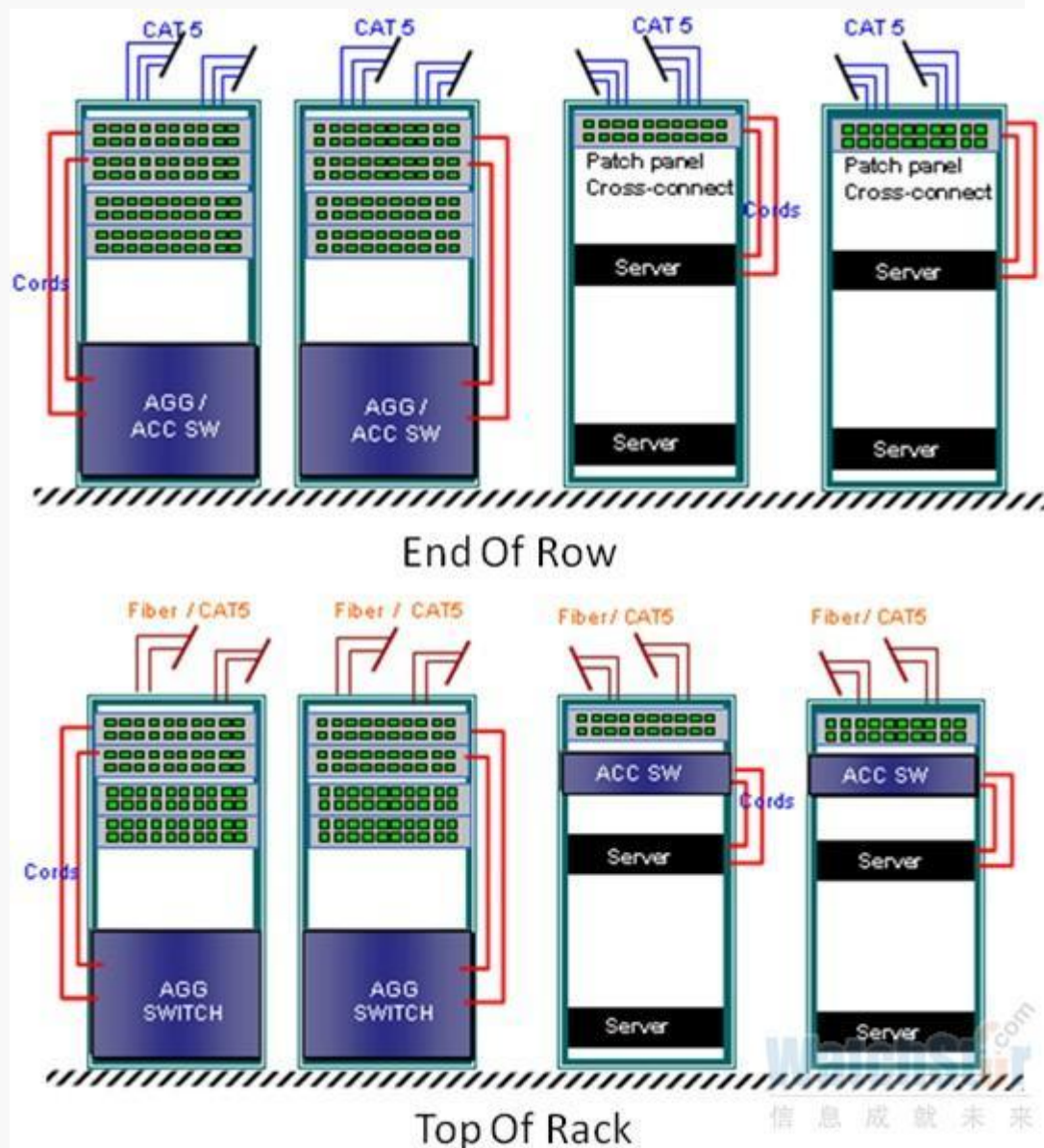
云计算数据中心对转发带宽的需求更是永无止境，因此会以部署核心-接入二层网络结构为主。层次越多，故障点越多、延迟越高、转发瓶颈也会越多。目前在一些 ISP（Internet Service Provider）的二层结构大型数据中心里，由于传统的 STP 需要阻塞链路浪费带宽，而新的二层多路径 L2MP 技术还不够成熟，因此会采用全三层 IP 转发来暂时作为过渡技术，如前面提到过的接入层与核心层之间跑 OSPF 动态路由协议的方式。这样做的缺点显而易见，组网复杂，路由计算繁多，以后势必会被 Ethernet L2MP 技术所取代。

新的二层多路径技术会在下文做更详细的介绍，不管是 TRILL 还是 SPB 都引入了二层 ISIS 控制平面协议来作为转发路径计算依据，这样虽然可以避免当前以太网单路径转发和广

播环路的问题,但毕竟是增加了控制平面拓扑选路计算的工作,能否使其依然如以往 Ethernet 般高效还有待观察。MPLS 就是一个尴尬的前车之鉴,本想着帮 IP 提高转发效率而生根发芽,没想到却在 VPN 路由隔离方面开花结果了,世事难料啊。

## 4.2 EOR 与 TOR

前面说了,数据中心网络设备就是交换机,而交换机就分为框式与盒式两种。当前云计算以大量 X86 架构服务器替代小型机和大型机,导致单独机架 Rack 上的服务器数量增多。受工程布线的困扰,在大型数据中心内 EOR(End Of Row)结构已经逐步被 TOR(Top Of Rack)结构所取代。盒式交换机作为数据中心服务器第一接入设备的地位变得愈发不可动摇。而为了确保大量盒式设备的接入,汇聚和核心层次的设备首要解决的问题就是高密度接口数量,由此框式交换机也就占据了数据中心转发核心的位置。





## 4.3 控制平面与转发平面

对交换机来说，数据报文转发都是通过 ASIC（Application Specific Integrated Circuit）芯片完成，而协议报文会上送到 CPU 处理，因此可以将其分为控制平面与转发平面两大部分。

控制平面主体是 CPU，处理目的 MAC/IP 为设备自身地址和设备自身发给其他节点的报文，同时下发表项给转发 ASIC 芯片，安排数据报文的转发路径。控制平面在三层交换机中尤为重要，需要依靠其学习路由转发表项并下发到 ASIC 芯片进行基于 IP 的转发处理。而二层交换机中数据报文的转发路径都是靠 MAC 地址广播来直接学习，和控制平面 CPU 关系不大。

转发平面则是以 ASIC 芯片为核心，对过路报文进行查表转发处理，对交换机来说，ASIC 转发芯片是其核心，一款交换机的能力多少和性能大小完全视其转发芯片而定。而控制平面 CPU 虽然也是不可或缺的部分，但不是本文介绍的重点，下文将以分析各类型交换机的转发处理为主。

## 4.4 Box 与集中式转发

经常看到设备商们今天推出个“高性能”，明天推出个“无阻塞”，后天又搞“新一代”的网络交换产品，各种概念层出不穷，你方唱罢我登台，搞得大家跟着学都学不过来，总有一种是不是被忽悠了的感觉。其实很多时候真的是在被忽悠。

先说说 Box 盒式设备。盒式交换机从产生到现在，以转发芯片为核心的集中式转发结构就没有过大的变化。集中式转发盒子的所有接口间流量都是走转发芯片来传输，转发芯片就是盒子的心脏。

而这个心脏的叫法多种多样，神马 Port ASIC、Switch Chip、Fabric ASIC、Unified Port Controller 等等都是各个厂家自行其说罢了，关键就看各个物理接口的 PHY（将 0/1 信号与数据互相转换用的器件）连接到哪里，哪里就是核心转发芯片。一般的中小型交换机设备厂商（H3C/中兴/锐捷/Foundry/Force10 等，Juniper 目前的数据中心 Switch 不提也罢，下文会简单说说未来的 QFabric）都会直接采购 Broadcom 和 Marvell 等芯片生产厂商的产品，只有 Cisco 和 Alcatel 等寥寥几家大厂商有能力自己生产转发芯片。但说实话，从转发能力来看这些自产的还真不见得能赶上公用的，人家专业啊。自产的最大好处其实在于方便搞些私有协议封包解包啥的，我的芯片我做主。

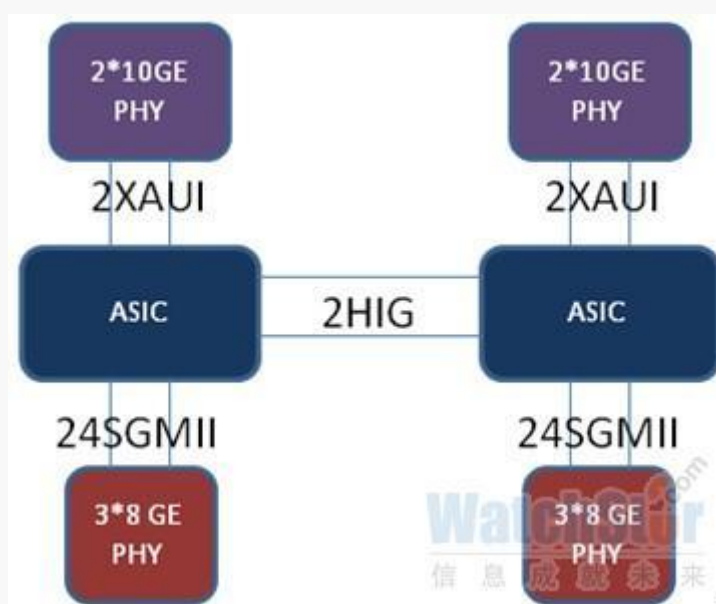
下面来说说集中式转发能力的计算，假设一个盒子喊自己的转发能力是  $x$  Gbps/ $y$  Mpps， $x$  是依靠所有外部接口带宽总和算出来的，如 48GE+2\*10GE 的盒子，转发能力就是单向 68GE，双向 136GE，一般  $x$  都会取双向的值；而  $y$  则是整机的包转发能力，

$y = x * 1000 / 2 / 8 / (64 + 20)$ , 1000 是 G 到 M 的转换, 2 是双向, 8 是每字节 8 比特, 64 是报文最小载荷, 20 是 IP 头长。要注意下面的机框式转发就不是这么算的了。大部分盒子的包转发能力还是能够很接近这个理论值的, 毕竟能选的转发芯片就那么多, 设备厂商在这里自己搞不出太多猫腻来。唯一有可能用来混淆客户的就是用芯片转发能力替代设备接口转发能力作为  $x$  值来宣传, 绝大部分交换机使用的芯片转发能力是大于所有接口带宽总和的。这时  $x$  与  $y$  都会比实际的要大一些, 但是很明显, 芯片再强, 没有接口引出来也没用的。所以这里的防忽悠技巧就是数接口数自己加一下。

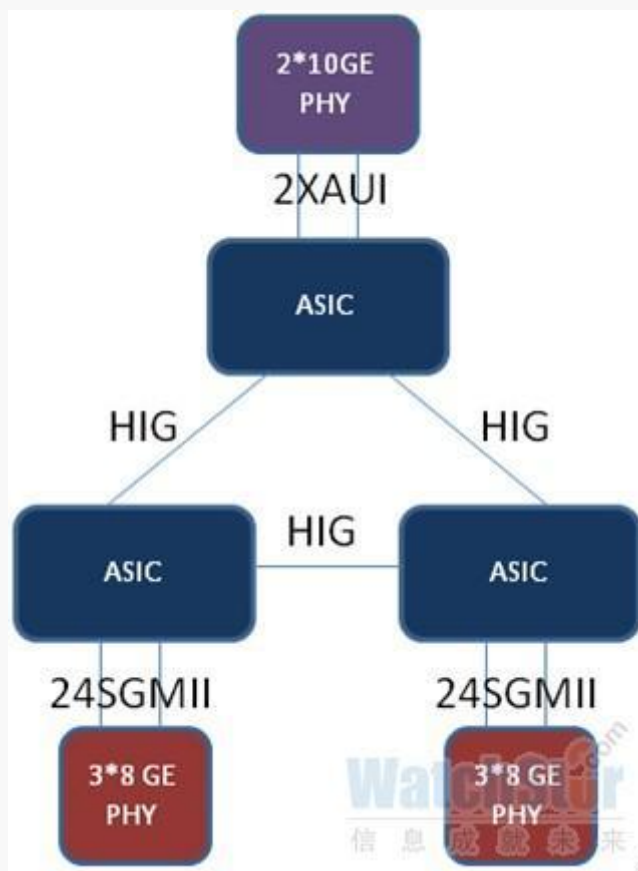
再说结构, 决定一款盒式交换机的接口转发容量的是转发芯片, 反之你看一款盒子的接口排布情况大概能反推出其使用的芯片能力。转发芯片的接口多种多样(如 SGMII、XAUI、HIG、Senders 等), 但通常每个芯片只连接 24 个 GE 接口(8 个口一个 PHY, 3 个 PHY 为一组), 因此遇到 24GE 口交换机, 通常都是单芯片的, 而 48GE 或更多就肯定是多芯片的了。而 10GE 接口的多少要看芯片的能力, 个人了解 Broadcom 有支持 24 个 10GE 的转发新片, 应该还有能力更强的。现在作者知道的 10GE 接口密度最高的盒子是 Arista 的 7148SX 和 Juniper 的 QFX3500, 都支持 48 个 10GE 接口, 具体布局有待拆机检查。

多芯片交换机还是很考验设备厂商架构设计人员的, 既要保证芯片间足够带宽互联互通, 又要考虑出接口不能浪费, 需拿捏好平衡。所以现在的多芯片盒式交换机设备大多是 2-3 个转发芯片的, 再多的就极少了, 芯片间互联设计起来太麻烦了。举两个例子, 大家可以看看下面这两种结构有没有问题。

首先是我們能不能用两块 6 个 HIG 接口级别转发能力的 ASIC(HIG 接口带宽 12.5GE), 设计一款 48GE+4\*10GE 的交换机呢? 答案是可以做, 但存在结构性拥塞, 芯片间至少需要 4 条 HIG 才能满足完全无阻塞。

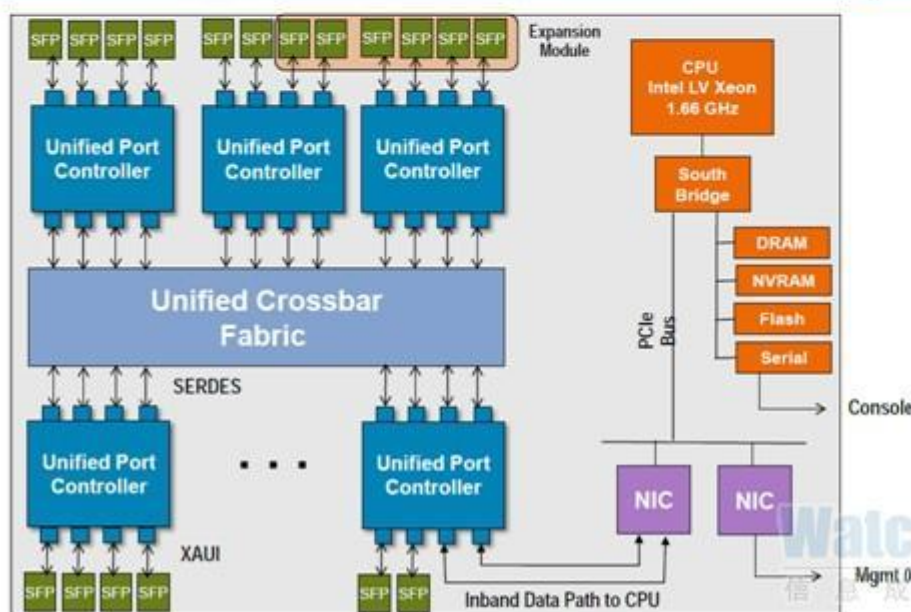
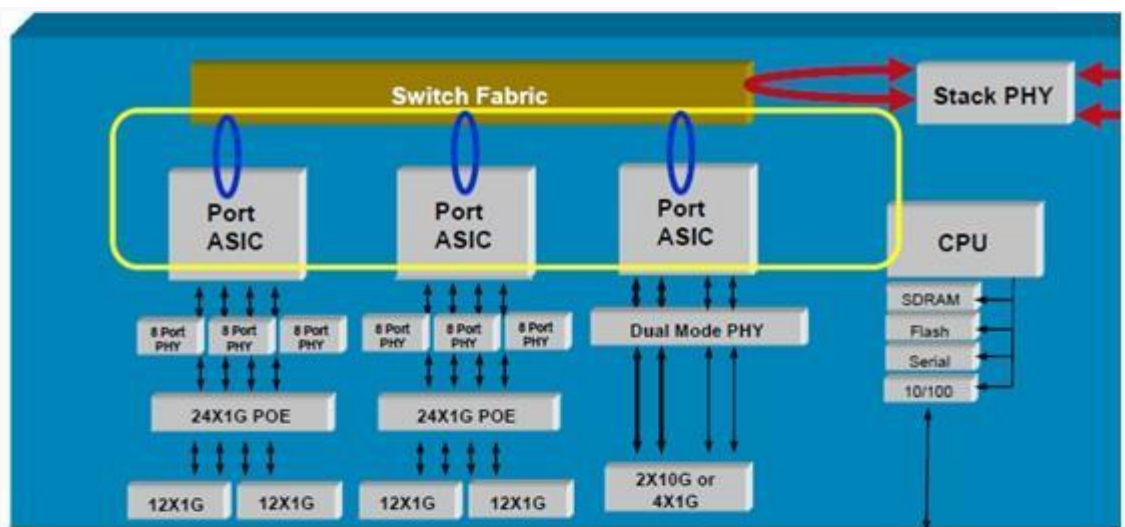


再来看一个，能不能用 3 块 4 个 **HIG** 接口级别转发能力的 **ASIC** 搭建出一款 48GE+2\*10GE 的交换机呢？没有问题，如下图所示内部结构是完全无阻塞的，缺点是部分流量会多绕经 1 个 **ASIC** 转发。



看完了前面这部分，相信大家对盒式交换机都能有个大致了解了，这里只讲讲结构，更详细的转发功能流程就需要有兴趣的童鞋自行去查看下各种芯片手册。另外上述两个例子只为讲解，请勿将当前市场产品对号入座。

刚刚说了，当盒子里面芯片较多时候连接起来很麻烦，于是出现了新的转发单元 **Switch Fabric**（Cisco N5000 上新的名词叫做 **Unified Crossbar Fabric**）。其实这个东东在框式交换机里面很常见，下面会有更详细的介绍。而在盒式交换机里面，目前看到的发布资料使用此种架构的就是 Cisco 的 3750X 和 N5000 了，连接方式如下图所示，这已经接近分布式转发的范围了。



作者将这个 Fabric 单元叫做交换芯片，便于和前面的 ASIC 转发芯片区分，二者的主要区别是，交换芯片只处理报文在设备内部的转发，类似 Cut-Through，为不同转发芯片间搭建路径，不做过滤和修改。而转发芯片要对报文进行各种查表、过滤和修改等动作，包括缓存都在其中调用，大多是基于 Store-Forward 方式进行报文处理，是交换机处理数据报文的核心部件。

3750X 目前还没有看到进一步的发展需要，而 N5000 其实是为了 Cisco 的网络虚拟化架构而服务，不再单单属于传统意义上的 Ethernet 交换机了。Juniper 为 QFabric 设计的 QFX3500 接入盒子（48\*10GE+4\*40GE）估计也是类似于 N5000 这种带交换芯片的分布式架构。另外怀疑 Arista 的 7148SX 也是分布式架构的，应该是 6 个 8\*10G 的转发芯片通过交换芯片连接，和它的机框式交换机中 48\*10G 接口板布局相同。

总的来说盒子里面搞分布式的最主要原因就是希望提高接口密度,尤其是万兆接口密度,后面相信还会有其他厂商陆续跟进,但是其接口数量需求是与部署位置息息相关的,盲目的扩充接口数并不一定符合数据中心的需要。

再唠叨几句数据中心 Box 交换机的选型,前面说了 Top Of Rack 是 Box 的主要归宿,一个标准 Rack 目前最高的 42U,考虑冗余怎么也得搞 2 台 Box,剩下最多装 40 台 1U 的 Server,那么上 48GE+4\*10GE 的 Box 就是最适合的。依此类推,接口数量多的 box 不见得真有太大作用,位置会很尴尬。考虑选择 Box 的最大转发容量时,直接根据服务器接口数来计算接口即可。目前随着 FCoE 的推进,服务器提供 10GE CNA 接口上行到接入交换机越来越常见,那么对 Box 的要求也随之提升到 10GE 接入 40G/100G 上行的趋势,像 Juniper 的 QFX3500 (48\*10GE+4\*40GE) 明显就是上下行带宽 1:3 收敛的交换机,估计下一代 Top Of Rack 的数据中心交换机怎么也得要 40\*10GE+4\*100GE 的接口才能彻底搞定 42U 机架,如果全部署 2U 的服务器,则最少也需要 16\*10GE+4\*40GE 接口的 Box 才靠谱一些。

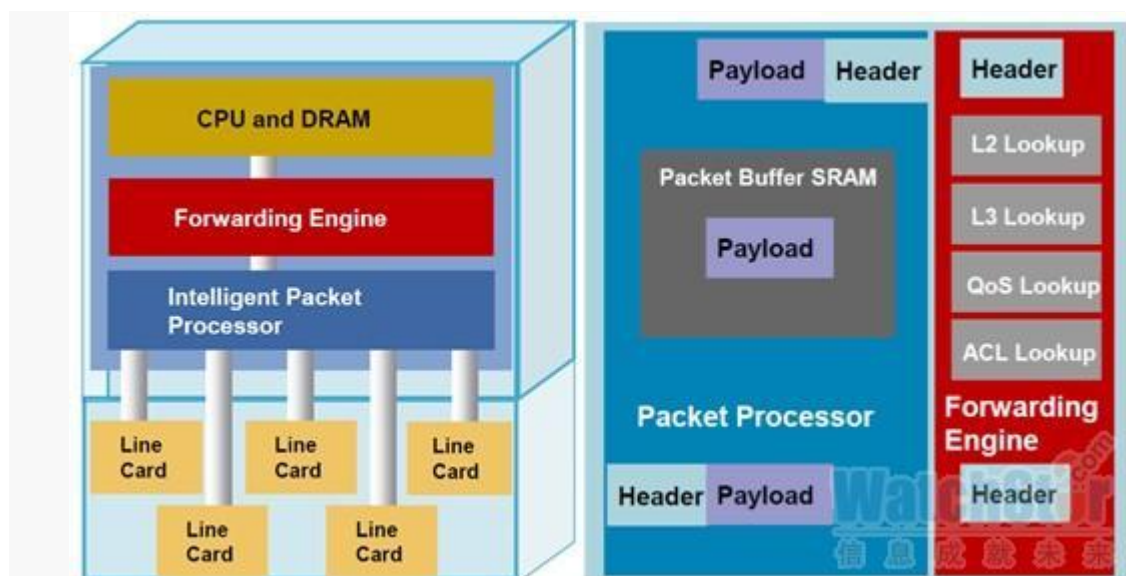
## 4.5 Chassis 与分布式转发

本章节涉及转发能力的举例计算量较大,对数字不感兴趣的同学可以直接略过相关内容。

盒子说完了讲讲框,盒式设备发展到一定程度,接口密度就成了天花板,必须要搞成机框式才能继续扩展了。可以把机框里面的板卡理解为一个独立的盒子,然后通过交换网络将其连接起来形成整体。

罗马不是一天建成的,机框式交换机最初也是按照集中式转发架构来进行设计。例如 Cisco4500 系列(又是 Cisco,没办法,就他家产品最全,开放出来的资料最多,而且确实是数通领域的无冕之王,下文很多技术也都跟其相关),其接口板(LineCard)上面都没有转发芯片的(XGStub ASIC 只做接口缓存和报文排队的动作),所有的数据报文都需要通过背板通道(Fabric),上送到主控板(Supervisor)的转发芯片(Forwarding Engine)上进行处理。结构如下图所示,其中 PP(Packet Processor)是做封包解包的,FE(Forwarding Engine)是做查表处理的。

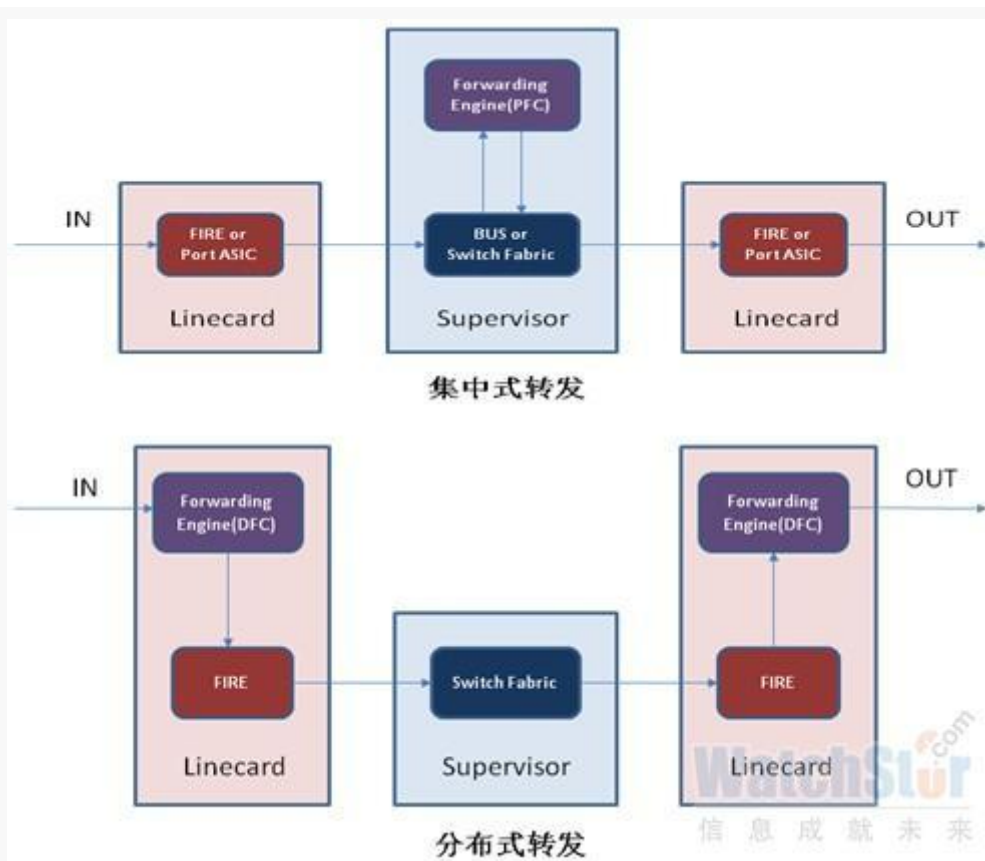




在早期 Cisco6500 系列交换机设备上同样是基于总线（BUS）的集中式转发结构。如 Classic 类型接口板（Module）就只有 Port ASIC 做缓存和排队，所有的报文同样要走到主控板（Supervisor32 或 720）上的转发芯片（PFC3）来处理。普通的 CEF256 和 CEF720 系列接口板虽然以 Switch Fabric 替代 BUS 总线通道来处理接口板到主控板的流量转发，但仍然是靠主控板上的 PFC3 对流量进行集中处理，因此还是集中式转发。直到 CEF256 和 CEF720 的 DFC（Distributed Forwarding Card）扣板出来，才能在板卡上进行转发，称得上是真正的分布式架构。而最新的第四代接口板 dCEF720 Linecards 已经直接将 DFC 变成了一个非可选组件直接集成在接口板上。

分布式架构指所有的接口板都有自己的转发芯片，并能独立完成查表转发和对报文的 L2/L3 等处理动作，接口板间通过交换芯片进行报文传递，机框的主控板只通过 CPU 提供协议计算等整机控制平面功能。分布式架构接口板上都会专门增加一个 Fabric 连接芯片（Fabric Interface 或 Fabric Adapter Process 等），用以处理报文在框内接口板间转发时的内部报头封装解封装动作。当报文从入接口板向交换芯片转发时，连接芯片为报文封装一个内部交换报头，主要内容字段就是目的出接口板的 Slot ID 和出接口 Port ID，交换芯片收到报文后根据 Slot ID 查找接口转发，出接口板的连接芯片收到后根据 Slot ID 确认，并将此内部交换报头去掉，根据 Port ID 将报文从对应出接口转出交换机。很显然分布式对比集中式的区别主要是芯片更多，成本更高，转发能力也更高。目前各厂商最新一代的主流数据中心交换机都已经是完全的分布式转发架构（如 Cisco 的 N7000，H3C 的 12500 等）。

下面说下 Chassis 的转发能力，这个可比盒子要复杂多了，各个厂家多如繁星的机框、主控和接口板种类足以使用户眼花缭乱。还是以 Cisco6500 系列交换机举例，一法通万法通，搞明白这个其他的也不过尔尔了。选择 Cisco6500 还有一个主要原因就是其结构从集中式跨越到分布式，从 BUS 总线通道跨越到 Crossbar 转发，堪称传统机框交换机百科全书。



FIRE (Fabric Interface & Replication Engine) 为 Cisco 的接口板连接芯片，除了作为连接 Switch Fabric 的接口对报文进行内部报头的封包解包动作外，还能提供本地镜像和组播复制功能。图中举例了报文在 65 机框式交换机中跨接口板转发的主要节点。集中式转发时板内接口间流量转发同样适用此图，而分布式转发时板内转发流量不需要走到 Switch Fabric。

另外报文走到出方向接口板时是否经过转发芯片处理各个厂家的设备实现并不一致，最简单的一个方法就是看交换机接口板支持不支持出方向的报文 ACL (Access Control List) 过滤，就知道其有没有上出口板转发芯片处理了。

从上图可以看出接口板的转发能力都受限于板卡连接 BUS 或 Switch Fabric 的接口带宽，而衡量整机转发能力时，集中式转发受限于转发芯片 FE 的转发能力，分布式转发受限于交换芯片 Switch Fabric 的转发能力。先说接口板转发能力，大家以前可能经常会听到接口板存在非线性速和收敛比的概念，看到这里就很好明白了，例如 CEF256 类型接口板的 Switch Fabric 接口带宽是 8G，那最多就支持 8 个 GE 口和其他接口板进行流量转发，其 WS-X6516-GBIC 接口板的面板上有 16 个 GE 口，明显就是一块 2:1 的收敛比的非线性速板。

再如 CEF720 类型接口板的 Switch Fabric 接口是 2\*20G (单板上有两个 FIRE)，那 48GE 口的单板也明显不可能是线速的了。即使是号称第四代的 dCEF720 接口板，其 Switch Fabric 接口和 CEF720 一样都是 2\*20G 接口，那么 X6708-10G 接口板 (提供 8\*10GE 接口) 和 X6716-10G 接口板 (提供 16\*10GE 接口) 只能是 2:1 和 4:1 收敛的非线性速板了。

背板通道预留不足，Switch Fabric 交换能力不够，6500 系列的这些架构缺陷促使 Cisco 狠下心来为数据中心重新搞出一套 Nexus7000，而其他交换机厂商也都几乎同时期推出了新架构的机框式交换机，都是被逼的啊，谁让 1000M 接入这么快就替代了 100M 接入呢，核心更得开始拼万兆了。

再说说整机转发能力。在集中式转发时，Cisco6500 不论使用 Supervisor32 还是 Supervisor720 主控，FE 转发芯片都是走 BUS 的，带宽都是 16G（双向 32G），因此只要用的接口板没有 DFC，整机最大也就双向 32G 了。而其中 Supervisor32 不支持 Switch Fabric，也就支持不了 DFC 的分布式转发，名称里的 32 就代表了其双向 32G 的最大整机转发能力。

Supervisor720 主控支持 18\*20 的 Switch Fabric 交换，名称中的 720 是指整个 Switch Fabric 的双向交换能力  $18*20*2=720G$ 。但其中 1 个通道用于连接 FE 转发芯片，1 个通道暂留未用，只有 16 个通道留给了接口板，意味着整机实际最大能够支持的双向转发能力是  $16*20*2=640G$ 。Supervisor720-10GE 支持 20\*20 的 Switch Fabric，多出来的 2 个 10G 通道给了 Supervisor 上的 2 个 10GE 接口，实际提供给接口板的交换通道仍然是 16\*20G。

刚刚说了，目前最新的 CEF720 系列接口板每块有 2\*20G 的出口，简单做个除法， $16/2=8$ ，主控板的交换芯片最多能够承载 8 块 CEF720 接口板，熟悉 Cisco6500 产品的同学这时候就会想到 6513 机框怎么办呢。6513 除去 7-8 的主控槽位外，一共有 11 个接口板槽位，1-6 槽位背板只提供 1 个 Switch Fabric 通道，9-13 才能提供 2 个通道，正好是  $6+2*5=16$  个通道满足主控板的 Switch Fabric 交换能力。而 6513E 虽在 1-6 槽位背板提供了 2 个通道，但实际上 1-6 槽位也同样只能支持 1 个 Switch Fabric 通道，否则 Supervisor720 的 Switch Fabric 也搞不定的。如果想 6513E 的接口板通道全用起来，只能等 Cisco6500 出下一代引擎了，至少是 Supervisor880 才能搞定 6513E 的全线速转发，不过从交换芯片的发展来看，Supervisor960 的可能性更大一些，1280 就有些拗口了。由上看出即使将 CEF720 接口板插到 6513/6513E 的 1-6 槽，也只能跑 20G 的流量，这下连 24GE 接口板都无法线速了。

前面算了好多数，好在都是加减乘除，只要搞明白了，完全可以避免选型时再被设备厂商忽悠。题外话，很多厂商的机框千兆接口板（24 或 48 个光/电口）都可以在其同时代盒式交换机中找到相似的影子，假如看到支持相同接口数量类型的接口板和盒子，相信里面的转发芯片十之八九也用的一样。万兆接口板不做成盒式是因为接口密度太低，价格上不去；而高密万兆的盒子做不成接口板则是因为框式交换机交换芯片和背板通道结构限制导致跨板转发能力上不去。

框式交换机架构从集中式发展到分布式后，整机的转发能力迎来了一次跳跃性发展，从 Cisco6500 的 Supervisor32 到 Supervisor720 就可见一斑。那么下一步路在何方呢，各个厂家都有着不同的看法。看回到前面分布式转发的结构图，可以想到要继续提升转发能力有两个

主要方向，一个是将单芯片处理能力提升，交换芯片只处理一次查表转发，工作简单相对更容易提升，而转发芯片要干的事情太多就不是那么好换代的了。

而另一条路就是增加芯片的数量，转发芯片由于要排布在接口板上，毕竟地方就那么大，发展有限，现在的工艺来说，一块单板放 4 个转发芯片基本上已经到极限了，6 个的也只看到 Arista 的 7548 接口板上有，再多的还没有见过，因此转发芯片的发展还是要看芯片厂商的能力了。而像 Cisco6500 的 Supervisor720 一样将交换芯片布在主控板上的话，同样面临空间的限制，上面还得放些 CPU/TCAM 什么的，最多每块主控上面放 2 个交换芯片就顶天了，双主控能支撑 4 个，但是全用做转发的话就做不到冗余了。最新的思路是将交换芯片拿出来单独成板，这样只要新机框设计得足够大，交换芯片的数量就不再是限制。例如 Cisco 的 N7000 可以插 5 块交换网板，而 H3C 的 12500 能够插 9 块交换网板。当然转发能力并不是交换芯片的数量越多就越好，还要看具体其单体转发能力和整机背板通道布局。

以 Cisco 的 N7000 举例分析，其交换网板 Fabric Modules 上的 CFA (Crossbar Fabric ASICs) 宣称是支持每槽位 (Slot) 2\*23G 的通道交换，整机最大支持 2\*23\*5=230G 的每槽位单向转发能力。这样能看出来啥呢？

1、N7010 上 8 个板卡槽位，2 个主控槽位 (主控槽位支持 1 条 23G 通路)，一共是  $8*2+2=18$  条通道，可以看出 7010 的交换网板上就一块 Crossbar Fabric ASIC，这个交换芯片和以前 Cisco6500 Supervisor720 上的 18\*20 交换芯片除了每通道带宽从 20G 提升到了 23G 以外，通路数都是 18 条没有变化，应该属于同一代交换芯片产品。7018 可以算出是  $16*2+2=34$  条通道，那么其每块交换网板上应该是 2 个与 7010 相同的 CFA 交换芯片，而且还空了 2 条通道暂时没用上。

2、其接口板上的数据通道同样应该与交换网板通道相匹配，升级到 23G 的容量。看下 48GE 接口板的图，上面只有一块 2 通道的转发芯片 Forwarding Engine，于是了解为啥其只能提供 46G 的全线速转发，而且使用一块交换网板就可以达到最大转发能力了。

3、再看其 10GE 接口板，8 口万兆板上上面有两块 2 通道的转发芯片，这样 80G 流量完全够处理，那么算算需要 2 块交换网板才能线速跨板转发，1 块就只能转 40G 了。而 32 口万兆板上上面就一块 4 通道的转发芯片，只能搞定 80G 流量转发，是收敛比 4:1 的非线速板，同样需要两块交换网板才能达到最大的跨板转发能力。

4、由上面 3 点可以看出，只使用目前 Cisco N7000 的接口板的话，交换网板 2+1 冗余就完全足够用了。Cisco 的下一步换代目标肯定是要想办法提升接口板转发芯片的能力了。首先应该搞定两块 4 通道转发芯片 FE 的工艺布局 (VOQ 和 Replication Engine 芯片的数量都要翻倍)，这样能把 16 口线速万兆板先搞出来，然后是否研究 20\*10GE 接口板就看其市场战略了。再下一步由于目前交换网板支持每接口板 230G 的总带宽限制，24/32 口万兆线速板肯定是搞不定的。只能先想法将交换网板升级一下，至少得让交换能力再翻一翻才好拿

出来搞定 32/40 口万兆板的线速转发，至于交换芯片是换代还是数量翻番就都有可能了。不过无论走哪条路都不是可以一蹴而就的事情，一两年内应该没戏。

再简单说说 H3C 的 12500，由于其公布资料太少，说多了会有问题。还是从网站公布的宣传值来看。12508 背板 7.65T，交换容量 3.06T/6.12T，包转发率 960Mpps/2400Mpps；12518 背板 16.65T，交换容量 6.66T/13.32T，包转发率 2160Mpps/5400Mpps。12508 与 12518 都是最大支持 9 块交换网板，当前主要接口板与 N7000 相似，含 48GE 和 4 万兆、8 万兆的线速接口板，32 万兆非线速接口板。

1、从背板算起，首先  $16.65 - 7.65 = 9T$  就是 10 个槽位的容量，考虑到厂商的宣传值都是双向，那么每接口板槽位应该是预留了  $9000 / 2 / 10 = 450G$  的最大出口带宽。根据 12508 推算，双主控每主控板槽位应该是预留  $(7650 / 2 - 450 * 8) / 2 = 112.5G$  的最大出口带宽。由此背板预留通道数接口板与主控板为  $450:112.5 = 4:1$  的关系。基于省钱原则，主控板上肯定只有一条通道，那么接口板都是 4 条通道，12508 背板槽位一共给接口和主控板留了  $4 * 8 + 2 = 34$  条通道。

2、12508 交换网板总的交换容量 3.06T，则每条通道的带宽应该是  $3060 / 34 = 90G$ ，由此可以推算出实际每块接口板的出口带宽为  $90 * 4 = 360G$ ，同样由于 3.06T 肯定是个双向值，则每接口板最大交换即可偶带宽理论值为 180G（比较 Cisco N7000 的 230G 理论值要低一些）。“交换容量 3.06T/6.12T”的写法应该指新一代的交换网板芯片能力翻倍或者是数量翻倍，那时其接口板理论带宽就可以达到 360G 了，还是小于前面计算的背板预留 450G 的最大带宽，说明背板设计还是考虑不错的。

3、再来算算接口板，从 8 万兆接口板支持线速转发看来，首先 4 个通道应该对应到 4 块转发芯片，每转发芯片对应 2 个万兆接口，处理 20G 的流量。而  $32 * 10G$  非线速板应该是同样使用 4 块转发芯片，所以也是 4:1 的收敛比。而其 48GE 和  $4 * 10GE$  的接口板应该是只用了 2 块同样的转发芯片，转发芯片的接口应该是使用类似于前面盒式交换机中的 12.5G 带宽线路类型，每块转发芯片对应 2 组 12GE 接口或 2 个 10GE 接口。考虑其所有接口板采用完全相同转发芯片是因为大量采购时存在价格优势，不像 Cisco 自己做芯片。

4、返回来说下 12518，总的通道数应该是  $18 * 4 + 2 = 74$ ，则总的交换容量应该为  $90 * 74 = 6.66T$  与其宣传值相同。有个小问题，这里的通道数计算是按照接口板与主控板来统计的，以交换板的角度来看时，12508 每块交换板一个交换芯片要连 8 块接口板，每接口板最大 4 条通道，既需要  $8 * 4 = 32$  个出口（主控板通道不见得会连接到交换芯片上，也可能是连接到交换网板的 CPU）；而 12518 每块交换板肯定是两个交换芯片，每芯片需要  $18 * 4 / 2 = 36$  个出口。这说明 12500 系列交换机网板上的交换芯片要不就都是 32 出口的，那么 12518 有 2 个槽位只有一半的转发能力；要不就都是 36+出口的，12508 存在部分出口空余用不上。

5、最后说下包转发率的计算，机框式交换机的包转发率应该是所有转发芯片转发能力的总和。如每个转发芯片 20G 处理带宽（单向），则转发率应该为  $20 / 8 / (64 + 20) = 29.76Mpps$ ，

取整为 30Mpps。按每接口板最大 4 个转发芯片计算，则 12508 整机为  $30 \times 4 \times 8 = 960\text{M}$ ，12518 为  $30 \times 4 \times 18 = 2160\text{M}$ ，符合其宣传值。至于其后面的 2400M 和 5400M 两个值，反向推算，每接口板转发能力为  $2400/8 = 5400/18 = 300\text{Mpps}$ ，带宽则为  $300 \times 8 \times (64 + 20) = 201600$  约 200G，难道是预示着其下一代接口板能够使用 2 个 100G 的转发芯片支持 2 个 100G 接口，拭目以待。

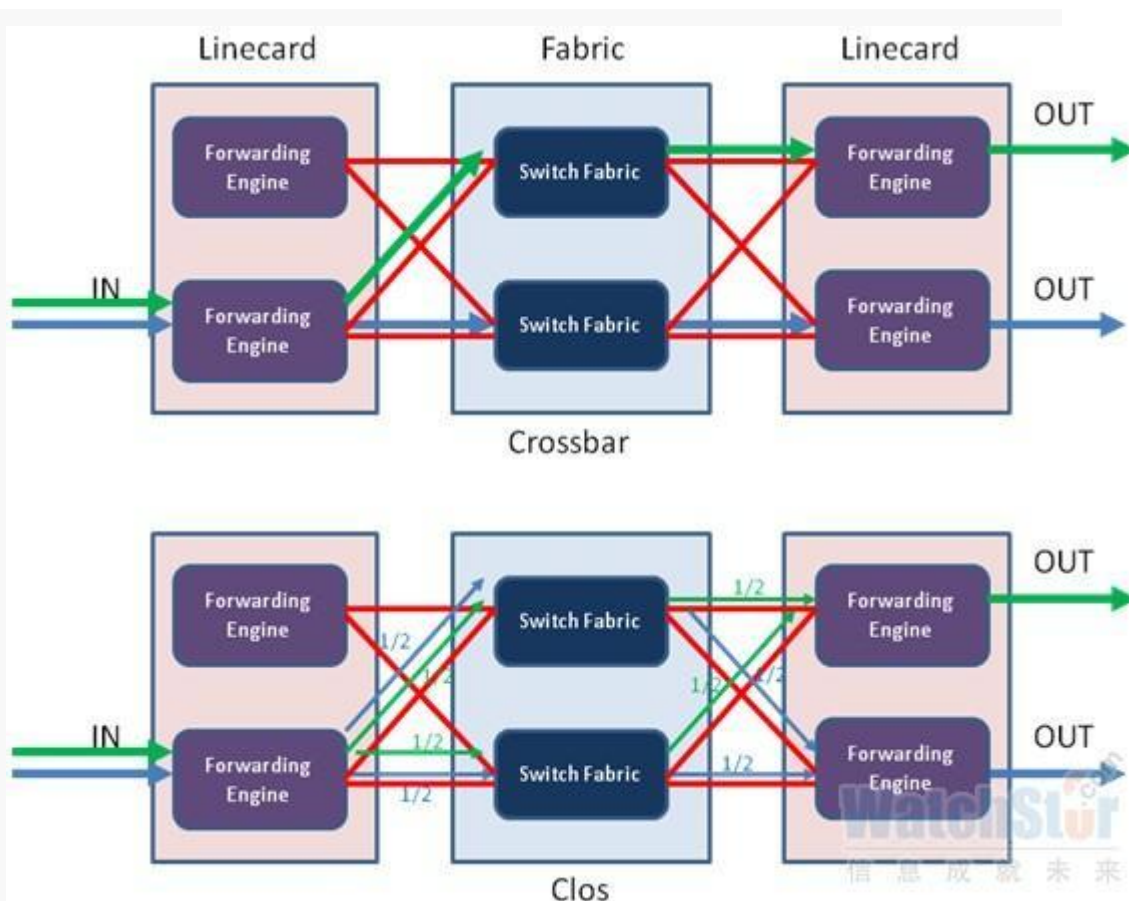
前面算了这么多，希望不会导致头晕吧。

## 4.6 Clos 与 VOQ

在 Crossbar 里面，任何两个转发芯片之间的只会经过一块交换芯片，路径是根据背板固定死的。这就导致了两个结构性问题的产生：一是多交换芯片时同一对转发芯片之间的流量不能被负载分担，如 Cisco N7000 就是如此；二是当多块入方向接口板往一块出方向接口板打流量的时候，流量可能都走到一块交换芯片上，导致本来应该在出接口板发生的拥塞，提前发生到交换芯片上，产生结构性拥塞，影响其他经过此芯片转发的流量。而且交换芯片采用 Cut-Through 方式转发是没有缓存的，报文都会直接丢弃，对突发流量的处理不理想。

为解决这两个问题，H3C 的 12500、Force10 的 E 系列和 Foundry BigIron RX 等设备都引入了 Clos 架构的概念（Cisco 的 CRS 系列高端路由器也是 Clos 结构，但 Nexus7000 不是）。Clos 架构是 1953 年贝尔实验室研究员 Charles Clos 设计的一种多级交换结构，最早应用在电话网络中。主要是两个特点，一是可以多级交换，二是每个交换单元都连接到下一级的所有交换单元上。上述厂商设备中基本都是入接口板-交换网板-出接口板的 3 级交换结构，而根据 Clos 设计，后续交换网可以扩展成多层结构。Crossbar 与 Clos 的主要区别如下图所示。



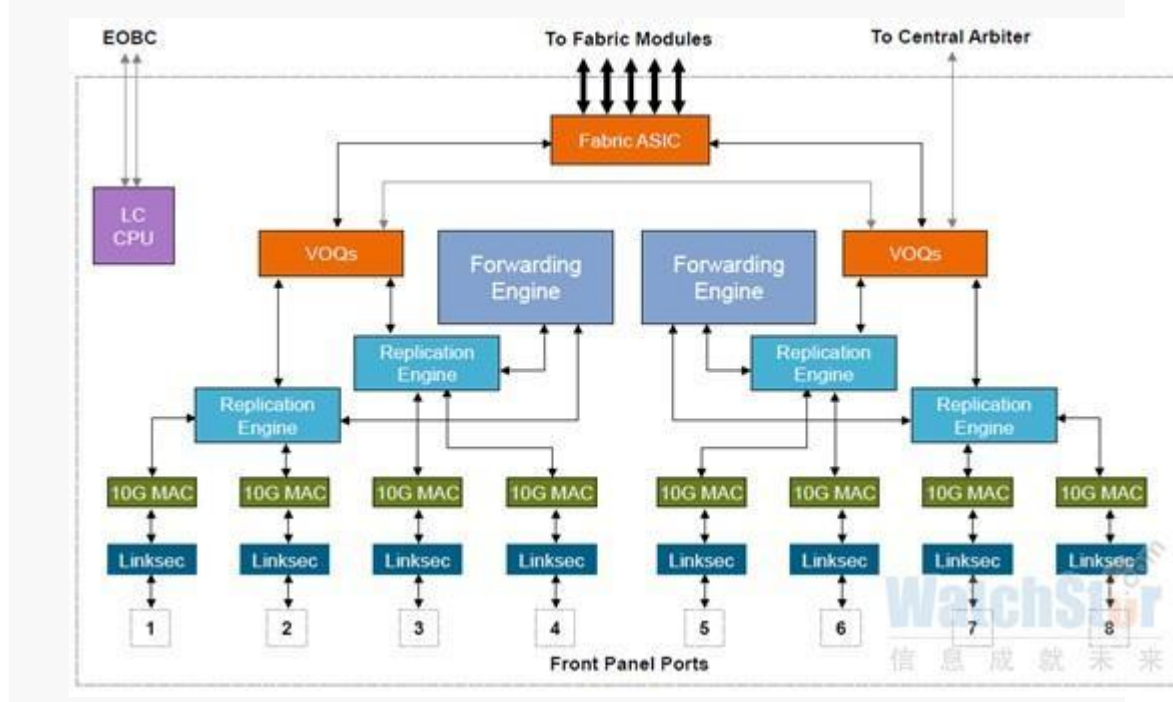


全连接的方式满足了对中间交换芯片的负载均衡需求，同时可以避免单交换芯片的结构性阻塞。不过话说回来，目前机框式交换机的转发能力提升瓶颈还是在转发芯片上，像前面举例的 N7000 和 125 都是结构转发能力远远大于实际接口板处理能力。所以暂时还不好说 Clos 就一定是趋势或代表啥下一代结构，就好像我现在一顿能吃 2 碗米饭，你给 10 碗还是 20 碗对我没有啥区别，都得等我胃口先练起来再说，但当我胃口真练起来那天，说不定又改吃馒头了呢。

多说一句，H3C 的 12500 在交换芯片转发流量时，报文是在入接口板先被切成等长信元再交给交换芯片的，到出接口板再组合，有些类似 ATM 转发，号称效率更高。而 Force10 的 E 系列则是按报文逐包转发，号称是为了避免乱序等问题。又是各有道理，管他呢，不出问题就什么都好。

目前新的分布式转发交换机另一项重要的技术就是 VOQ（Virtual Output Queues）。刚才说的 Crossbar 第二个阻塞问题在 Clos 架构中，虽然流量不会在 Switch Fabric 阻塞，但是多打一的情况下仍然会在出接口板阻塞。VOQ 就是在入接口板将报文发给 Switch Fabric 之前，先用 VOQ 缓存一下，然后通过中央裁决线路，发一个问询给出接口板，看看那边还有没有空间接收，有的话就发，没有先缓存一会儿，和 FC 网络中实现零丢包的 Bufer to Bufer Credit 机制很相似（BB Credit 机制详见下文 FCoE 技术部分）。这样就使出接口板的缓存能

力扩充到多块入接口板上，容量翻倍提升，可以有效的缓解突发拥塞导致的丢包问题。看下图 Cisco N7000 的 8 口万兆板结构图可以较好理解 VOQ 在接口板中的位置。



## 4.7 网络小结

数据中心网络看交换，交换机发展看芯片，分布式转发是必然，Clos 架构有得盼。

本章内容是下文数据中心内部服务器通信网络发展技术的重要铺垫。充分了解机框式交换机，可以对后面提出的新一代数据中心网络虚拟化技术，（如 Cisco 的 VN-Tag、Fabric Extend 和 Fabric Path/E-TRILL 等）在理解时起到巨大的帮助。

题外话，目前很多企业规模大了以后，网络部门负责网络，业务部门负责应用和服务器，很多时候互不搭界，于是设计网络和应用的时候就各搞各的，等数据中心建起来之后发现这也是问题那也是问题，各个都变身救火队员，不是啥好现象。有一本书建议所有的网络规划设计人员翻看，《自顶向下的网络设计》，即使找不到或没时间看也请一定要记住这个书名，终身受益的。对应用业务设计人员，也请稍微了解下网络，最少也得能估算出业务上线后理论上的平均带宽和峰值带宽，好向网络设计人员提出需求，免得出事时焦头烂额互相推诿。

# 5 技术

终于到本文的根本了，前面 balabala 的说了那么多，都是本章的铺垫，就是希望大家明白下面这些技术是为何而来，要解决什么样的需求和问题。再次对前面的需求进行个汇总。

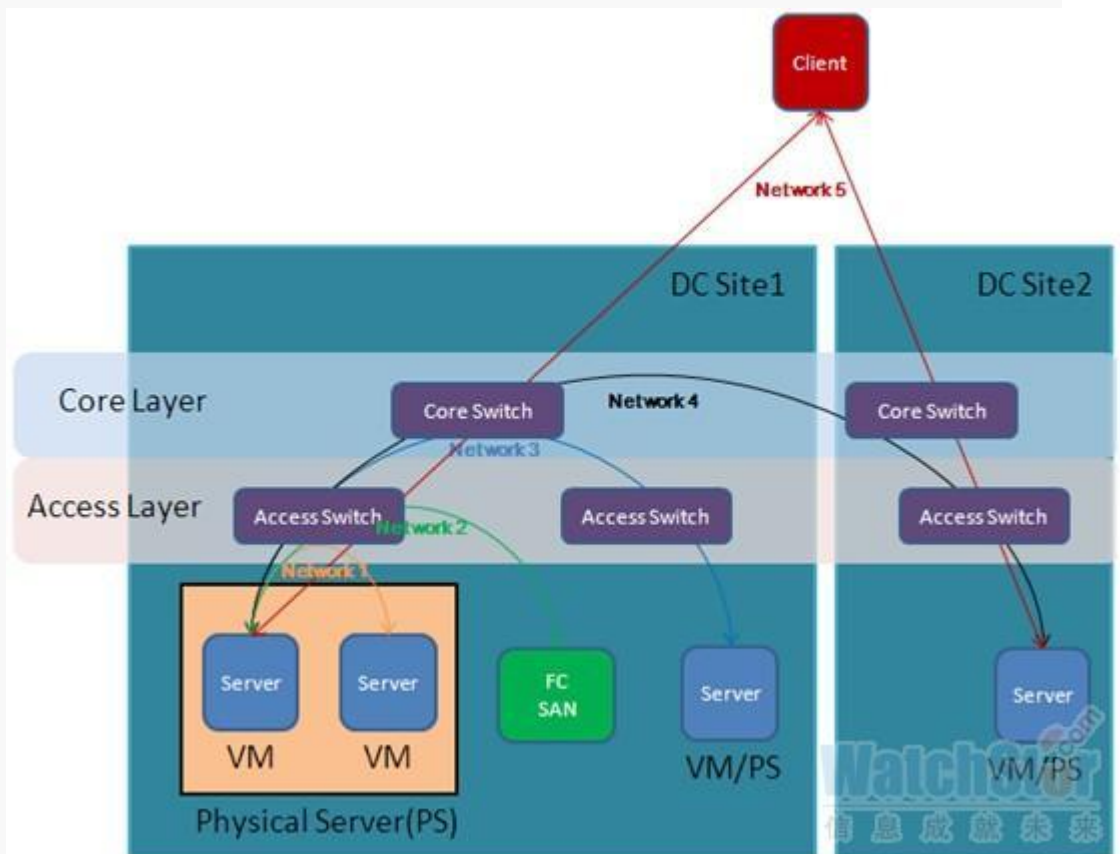
## 1、VM 之间的互通

- 2、更多的接口，更多的带宽
- 3、二层网络规模扩大
- 4、数据中心站点间二层互联
- 5、VM 跨站点迁移与多站点选路
- 6、服务器前后端网络融合（这个属于厂家引导还是用户需求真不好说）

下面就来看看下面这些网络技术是如何解决上述需求问题的。

## 5.1 技术结构

前面说了，数据中心网络流量的根本出发点是 Server，结合云计算最适合的核心-接入-二层网络结构，可以把下面要介绍的各种技术分类如下图所示。此处只做结构上的介绍，具体技术细节将在下文展开。



Network1-VM 本地互访网络，边界是 Access Switch，包括物理服务器本机 VM 互访和跨 Access Switch 的不同物理服务器 VM 互访两个层面。原有技术以服务器内部安装软件虚拟交换机 VSwitch 为主，新技术则分为以服务器为主体的 802.1Qbg EVB

（VEPA/Multi-channel）和 Cisco 以网络交换机为主体的 802.1Qbh BPE（Port Extend/VN-Tag/VN-Link）两大 IEEE 标准体系。

Network2-Ethernet 与 FC 融合，就是 FCoE，边界仍然是 Access Switch。在服务器物理网卡到 Access Switch 这段，将 FC 数据承载在 Ethernet 的某个 VLAN 中传输。但实际上各个厂商当前实现都是做 NPV 交换机，并不是真正的 FCoE，只有很少的产品如 Cisco 的 Nexus5000 系列和 Brocade 的 8000 系列等能够支持做 FCF。

Network3-跨核心层服务器互访网络，边界是 Access Switch 与 Core Switch。可理解为服务器互访流量从进入 Access Switch，经过 Core Switch，再从另一个 Access Switch 转出过程的网络处理技术。原有技术就是 STP 了，新技术分为设备控制平面虚拟化（VSS/vPC/IRF）和整网数据平面虚拟化（SPB/TRILL/Fabric Path）两大体系。这两个体系都是网络虚拟化中的多虚一方向，在一虚多方向除去传统的 VLAN/VRF 外，Cisco 的 N7000 系列还依照 X86 架构虚拟化整出了个 VDC。

Network4-数据中心跨站点二层网络，边界是 Core Switch。目标是跨越核心网为多个数据中心站点的 Core Switch 之间建立一条二层通道。根据站点间互联核心网的区别，分为以下三类技术：

光纤直连（SDH/DWDM 等）对应 Ethernet（RPR）

MPLS 核心网对以 L2VPN（VLL/VPLS）

IP 核心网对应 IP 隧道技术（VLLoGRE/VPLSoGRE/L2TPv3/OTV）

Cisco 的 OTV 虽然主要应用在 IP 核心网中，但实际前面两种方式下同样可以使用，只要多个数据中心站点的 Core Switch 设备间能够建立可达的 IP 路径即可部署。使用 VLL/VPLS 相关技术时必须增加专门的 PE 设备为站点间的 Core Switch 建立二层隧道，而 OTV 可以直接部署在 Core Switch 上。

Network5-数据中心多站点选择，技术边界在数据中心与广域网相连的边缘。在云计算中，VM 跨站点迁移后，业务服务器 IP 地址不变，网络指向需要随之变化。这块前面也提到现有技术就是 DNS 域名解析与 ServerLB 的 NAT 配合，以及主机 IP 路由发布等方式。新技术则是 Cisco 提出 LISP 以 IPinIP 技术结构绕开 DNS，由网络设备单独处理 Client 在广域网中选择站点的情况。

## 5.2 网络虚拟化

云计算就是计算虚拟化，而存储虚拟化已经在 SAN 上实现得很好了，那么数据中心三大件也就剩下网络虚拟化。那么为什么要搞网络虚拟化呢？还是被计算逼的。云计算多虚一时，所有的服务资源都成为了一个对外的虚拟资源，那么网络不管是从路径提供还是管理维护的角度来说，都得跟着把一堆的机框盒子进行多虚一统一规划。而云计算一虚多的时候，物理服务器都变成了一堆的 VM，网络怎么也要想办法搞个一虚多对通路建立和管理更精细化一些不是。

### 5.2.1 网络多虚一技术

先说网络多虚一技术。最早的网络多虚一技术代表是交换机集群 Cluster 技术，多以盒式小交换机为主，较为古老，当前数据中心里面已经很少见了。而新的技术则主要分为两个方向，控制平面虚拟化与数据平面虚拟化。

#### 控制平面虚拟化

顾名思义，控制平面虚拟化是将所有设备的控制平面合而为一，只有一个主体去处理整个虚拟交换机的协议处理，表项同步等工作。从结构上来说，控制平面虚拟化又可以分为纵向与横向虚拟化两种方向。

纵向虚拟化指不同层次设备之间通过虚拟化合多为一，代表技术就是 Cisco 的 Fabric Extender，相当于将下游交换机设备作为上游设备的接口扩展而存在，虚拟化后的交换机控制平面和转发平面都在上游设备上，下游设备只有一些简单的同步处理特性，报文转发也都需要上送到上游设备进行。可以理解为集中式转发的虚拟交换机

横向虚拟化多是将同一层次上的同类型交换机设备虚拟合一，Cisco 的 VSS/vPC 和 H3C 的 IRF 都是比较成熟的技术代表，控制平面工作如纵向一般，都由一个主体去完成，但转发平面上所有的机框和盒子都可以对流量进行本地转发和处理，是典型分布式转发结构的虚拟交换机。Juniper 的 QFabric 也属于此列，区别是单独弄了个 Director 盒子只作为控制平面存在，而所有的 Node QFX3500 交换机同样都有自己的转发平面可以处理报文进行本地转发。

控制平面虚拟化从一定意义上来说是真正的虚拟交换机，能够同时解决统一管理 with 接口扩展的需求。但是有一个很严重的问题制约了其技术的发展。在前面的云计算多虚一的时候也提到过，服务器多虚一技术目前无法做到所有资源的灵活虚拟调配，而只能基于主机级别，当多机运行时，协调者的角色（等同于框式交换机的主控板控制平面）对同一应用来说，只能主备，无法做到负载均衡。

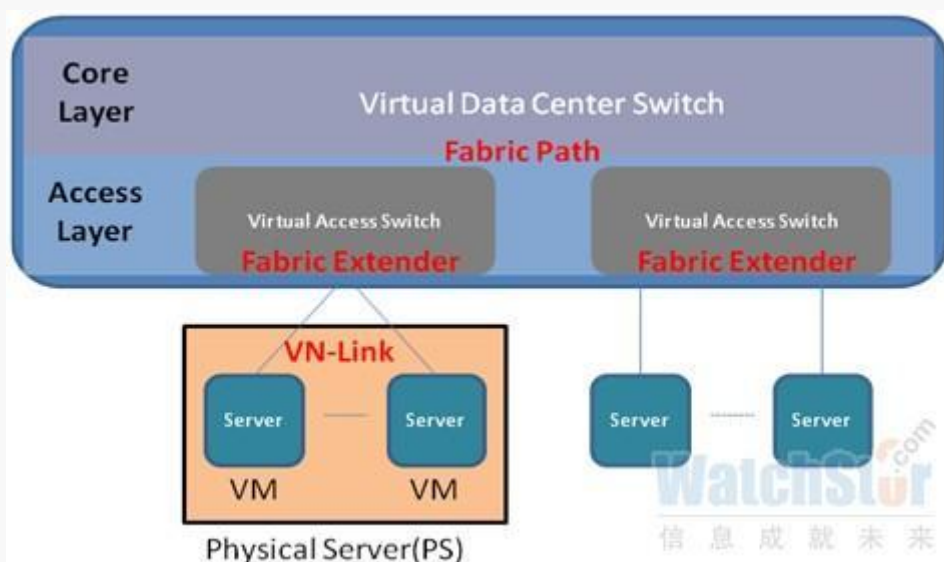
网络设备虚拟化也同样如此，以框式设备举例，不管以后能够支持多少台设备虚拟合一，只要不能解决上述问题，从控制平面处理整个虚拟交换机运行的物理控制节点主控板都只能有一块为主，其他都是备份角色（类似于服务器多虚一中的 HA Cluster 结构）。总而言之，虚拟交换机支持的物理节点规模永远会受限于此控制节点的处理能力。这也是 Cisco 在 6500 系列交换机的 VSS 技术在更新换代到 Nexus7000 后被砍掉，只基于链路聚合做了个 vPC 的主要原因。三层 IP 网络多路径已经有等价路由可以用了，二层 Ethernet 网络的多路径技术在 TRILL/SPB 实用之前只有一个链路聚合，所以只做个 vPC 就足矣了。另外从 Cisco 的 FEX 技术只应用于数据中心接入层的产品设计，也能看出其对这种控制平面虚拟化后带来的规模限制以及技术应用位置是非常清晰的。

#### 数据平面虚拟化

前面说了控制平面虚拟化带来的规模限制问题，而且短时间内也没有办法解决，那么就想法子躲过去。能不能只做数据平面的虚拟化呢，于是有了 TRILL 和 SPB。关于两个协议的具体细节下文会进行展开，这里先简单说一下，他们都是用 L2 ISIS 作为控制协议在所有设备上上进行拓扑路径计算，转发的时候会对原始报文进行外层封装，以不同的目的 Tag 在 TRILL/SPB 区域内部进行转发。对外界来说，可以认为 TRILL/SPB 区域网络就是一个大的虚拟交换机，Ethernet 报文从入口进去后，完整的从出口吐出来，内部的转发过程对外是不可见且无意义的。

这种数据平面虚拟化多合一已经是广泛意义上的多虚一了，相信看了下文技术理解一节会对此种技术思路有更深入的了解。此方式在二层 Ethernet 转发时可以有效的扩展规模范围，作为网络节点的 N 虚一来说，控制平面虚拟化目前 N 还在个位到十位数上晃悠，数据平面虚拟化的 N 已经可以轻松达到百位的范畴。但其缺点也很明显，引入了控制协议报文处理，增加了网络的复杂度，同时由于转发时对数据报文多了外层头的封包解包动作，降低了 Ethernet 的转发效率。

从数据中心当前发展来看，规模扩充是首位的，带宽增长也是不可动摇的，因此在网络多虚一方面，控制平面多虚一的各种技术除非能够突破控制层多机协调工作的技术枷锁，否则只有在中小型数据中心里面刨食的份儿了，后期真正的大型云计算数据中心势必是属于 TRILL/SPB 此类数据平面多虚一技术的天地。当然 Cisco 的 FEX 这类定位于接入层以下的技术还是可以与部署在接入到核心层的 TRILL/SPB 相结合，拥有一定的生存空间。估计 Cisco 的云计算数据中心内部网络技术野望如下图所示：（Fabric Path 是 Cisco 对其 TRILL 扩展后技术的最新称呼）





## 5.2.2 网络一虚多技术

再说网络一虚多，这个可是根源久远，从 Ethernet 的 VLAN 到 IP 的 VPN 都是大家耳熟能详的成熟技术，FC 里面也有对应的 VSAN 技术。此类技术特点就是给转发报文里面多插入一个 Tag，供不同设备统一进行识别，然后对报文进行分类转发。代表如只能手工配置的 VLAN ID 和可以自协商的 MPLS Label。传统技术都是基于转发层面的，虽然在管理上也可以根据 VPN 进行区分，但是 CPU/转发芯片/内存这些基础部件都是只能共享的。目前最新的一虚多技术就是 Cisco 在 X86 架构的 Nexus7000 上实现的 VDC，和 VM 一样可以建立多个 VDC 并将物理资源独立分配，目前的实现是最多可建立 4 个 VDC，其中还有一个是做管理的，推测有可能是通过前面讲到过的 OS-Level 虚拟化实现的。

从现有阶段来看，VDC 应该是 Cisco 推出的一项实验性技术，因为目前看不到大规模应用的场景需求。首先转发层面的流量隔离（VLAN/VPN 等）已经做得很好了，没有必要搞个 VDC 专门做业务隔离，况且从当前 VDC 的实现数量（4 个）上也肯定不是打算向这个方向使劲。如果不搞隔离的话，一机多用也没有看出什么实用性，虚拟成多个数据中心核心设备后，一个物理节点故障导致多个逻辑节点歇菜，整体网络可靠性明显降低。另外服务器建 VM 是为了把物理服务器空余的计算能力都用上，而在云计算数据中心里面网络设备的接口数应该始终是供不应求的，哪里有多少富裕的还给你搞什么虚拟化呢。作者个人对类似 VDC 技术在云计算数据中心里面的发展前景是存疑的。

### SR-IOV

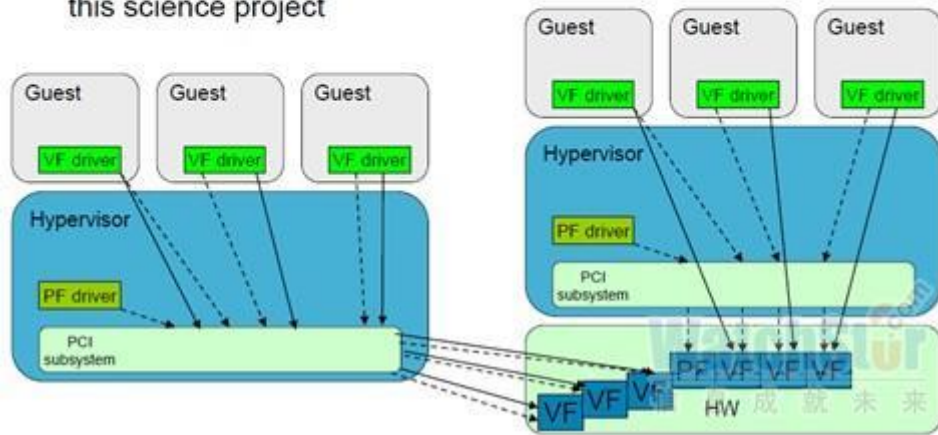
对网络一虚多这里还有个东西要补充一下，就是服务器网卡的 IO 虚拟化技术。单根虚拟化 SR-IOV 是由 PCI SIG Work Group 提出的标准，Intel 已经在多款网卡上提供了对此技术的支持，Cisco 也推出了支持 IO 虚拟化的网卡硬件 Palo。Palo 网卡同时能够封装 VN-Tag（VN 的意思都是 Virtual Network），用于支撑其 FEX+VN-Link 技术体系。现阶段 Cisco 还是以 UCS 系列刀片服务器集成网卡为主，后续计划向盒式服务器网卡推进，但估计会受到传统服务器和网卡厂商们的联手狙击。

SR-IOV 就是要在物理网卡上建立多个虚拟 IO 通道，并使其能够直接一一对应到多个 VM 的虚拟网卡上，用以提高虚拟服务器的转发效率。具体说是对进入服务器的报文，通过网卡的硬件查表取代服务器中间 Hypervisor 层的 VSwitch 软件查表进行转发。另外 SR-IOV 物理网卡理论上加块转发芯片，应该可以支持 VM 本地交换（其实就是个小交换机啦），但个人目前还没有看到实际产品。SR(Single Root)里面的 Root 是指服务器中间的 Hypervisor，单根就是说目前一块硬件网卡只能支持一个 Hypervisor。有单根就有多根，多根指可以支持多个 Hypervisor，但貌似目前单物理服务器里面跑多个 Hypervisor 还很遥远，所以多根 IO 虚拟化 MR-IOV 也是个未来未来时。摘录 Cisco 胶片对 MR-IOV 描述如下：(HW 为 Hardware，PF 为 Physical Function，VF 为 Virtual Functions)

## Multi-Root IO Virtualization

- HW Shared at the PCI Level
- PCI translation and routing requirements – new protocols
- PCI not designed to leave the enclosure – Gartner advocated this science project

**This is not happening**



SR-IOV 只定义了物理网卡到 VM 之间的联系，而对外层网络设备来说，如果想识别具体的 VM 上面的虚拟网卡 vNIC，则还要定义一个 Tag 在物理网卡到接入层交换机之间区分不同 vNIC。此时物理网卡提供的就是一个通道作用，可以帮助交换机将虚拟网络接口延伸至服务器内部对应到每个 vNIC。Cisco UCS 服务器中的 VIC (Virtual Interface Card) M81-KR 网卡 (Palo)，就是通过封装 VN-Tag 使接入交换机 (UCS6100) 识别 vNIC 的对应虚拟网络接口。

网络虚拟化技术在下一个十年中必定会成为网络技术发展的重中之重，谁能占领制高点谁就能引领数据中心网络的前进。从现在能看到的技术信息分析，Cisco 在下一个十年中的地位仍然不可动摇。

### 5.3 技术理解

进入正式介绍之前再多说两句如何快速理解技术的思路。搞网络的一般最头疼最不愿意的就是去读 RFC 等标准技术文档了，至少心底里有抵触。各种各样的报文、状态机、数据库、链表充斥于字里行间，再加上标准文档为了避免歧义，一句话能说清楚的也得分成三四句解释来解释去。也许是眼界不够开阔，反正我还真不识别能把草案标准当成小说看的牛人。

这里只是简单介绍一下协议入门的经验，如果想要深入甚至精通，那还得去一字一字的考究了，做学问来不得半点马虎。

学习一门技术前，首先要了解的是由何而来与从何而去，既技术产生的背景和应用的地方，这样对其要解决什么样的问题大致能有个印象。举例来说 PBB 是运营商城域以太网技

术，运营商技术的特点就是组网规模大，节点众多，路径众多。而传统以太网只能使用 STP 避免环路，阻塞了一堆链路，这个在运营商里面也是不可想象的，那一条条链路都是钱啊。因此 PBB 肯定是要在避免环路的同时，能够增大以太网组网规模 and 将所有路径都利用起来的技术。

再来就是看技术的类型， Routed Protocol 和 Routing Protocol 两个词很好的对技术组成部分进行了分类。这里的 Route 可以进行广义理解，不要只限于 IP，作者倾向于将 Routed 解释成封装，Routing 解释成寻址。任何一段数据信息从起点 A 发送到终点 B 的过程中，中间网络做的事情就是封装与寻址两件事。

由于中间网络只做传输，是不需要了解数据信息的，因此要封装一个可以识别的目的地址 Tag，这个 Tag 可以理解为目的 IP/目的 MAC/MPLS 标签等等，所有中间设备只要能识别这个 Tag 即可，这就是封装。

再说寻址，网络设备能够识别目的 Tag 后，还需要知道对应的本地出接口在哪才能将报文转发出去。最傻瓜的处理方式有两种，一个是通过手工配置的方式将 Tag 静态对应到本地出接口上（如静态路由、静态 MAC 等），再有就是在所有接口广播了（Ethernet）。高级的方式则是使用一种寻址用的动态协议，自动的进行邻居发现、拓扑计算和 Tag 传递等动作，如使用 RIP/OSPF/BGP/ISIS/LDP/PIM/MSDP 等等。这里需要注意的是传统 Ethernet 是通过广播来寻址的，注定规模不能太大。STP 的唯一作用就是防止环路，通过拓扑计算将多余的路径阻塞掉，与寻址无关。而前面提到的那些寻址协议主要任务都是传递 Tag 计算转发路径，大部分协议会通过计算拓扑来防止环路，但也有如 RIP 这种不计算拓扑的协议，搞些水平分割、毒性逆转和最大跳数等机制来避免环路。

封装解封装技术是网络入口与出口节点在原始数据信息前将 Tag 进行加载剥离动作，寻址技术则是在网络节点之间运行的交互动作。在很多协议技术中提到的数据平面其实就是封装转发，而控制平面就是标识寻址。



图是不是很眼熟，大部分的网络协议够可以照着这个模型去套的。

对于 IP 来说，Sender 和 Receiver 就是 TCP 协议栈，Edge 就是 IP 协议栈，Core 就是 Router，Payload 就是 TCP 数据，Tag 就是 IP 头中的目的 IP；

对于 Ethernet 来说，Sender 和 Receiver 就是 IP 协议栈，Edge 就是网卡接口，Core 就是 Switch，Payload 就是 IP 数据，Tag 就是 Ethernet 头中的目的 MAC；

对于 MPLS 来说，Sender 和 Receiver 就是 CE，Edge 就是 PE，Core 就是 P；Payload 就是 Ethernet/IP 数据，Tag 就是 MPLS 标签；

甚至对于分布式结构机框交换机来说，Sender 和 Receiver 就是接口板转发芯片，Edge 就是接口板上的交换接口芯片，Core 就是交换芯片，Payload 就是 Ethernet 数据报文，Tag 就是目的 Slot ID 和 Port ID（交换芯片转发时只看 Slot ID，目的接口板查看 Port ID）。

传统的 FC/IP/Ethernet 技术体系上面已经玩不出来花了，现在新的技术大都是在 FC/IP/Ethernet 等数据载荷外面增加个新的 Tag 并设计一套对应的寻址协议机制（如 MPLS 和下文的 FEX/ TRILL 等），或者干脆就还使用原有的 IP/MAC 作为外层封装 Tag，只对寻址进行变化。对于后者，作者喜欢称呼其为嫁接技术，神马 MACinMAC，IPinIP，MACinIP 等等都属于此列。此类技术的好处是兼容，缺点是继承，缝缝补补肯定没有全新设计来得自由。

封装比较好明白，协议理解的难点其实在于寻址。前面说了，静态寻址要手工一条条配置，规模大了能累死人。动态寻址技术配置工作量小了很多，但复杂度就上升了好几个台阶。不劳力就劳心，目前看来大家还是更喜欢劳心一些。回来说动态寻址，除了 RIP 这种早期的靠广播来传递路由 Tag 的寻址协议外，后面出来的都是先建邻接，后画拓扑，再传 Tag 的三步走了，从 OSPF/BGP/ISIS 到下面要讲到的 TRILL/SPB/OTV 皆是如此。对寻址技术主要内容简单归纳如下，细的就要看各协议具体实现了，希望有助于读者在学习寻址协议时能够少死些脑细胞。（文学素养有限，合辙押韵就算了吧）

建立邻居靠 Hello（Advertise），拆除邻接等超时。各自为根绘周边，主根扩散画整网。Tag 同步传更新，本地过期发删除。

技术理解部分就说这些，希望对读者认识新技术时能够有所帮助。下面开始进入技术主题。

## 5.4 VM 本地互访网络技术

本章节重点技术名词：EVB/VEPA/Multichannel/ SR-IOV/VN-Link/FEX/VN-Tag/ UCS/ 802.1Qbh/802.1Qbg

题目中的本地包含了两个层面，一个是从服务器角度来物理服务器本地 VM 互访，一个是从交换机角度来接入层交换机本地 VM 互访。这两个看问题的角度造成了下文中 EVB 与 BPE 两个最新技术体系出发点上的不同。

在 VM 出现伊始，VMware 等虚拟机厂商就提出了 VSwitch 的概念，通过软件交换机解决同一台物理服务器内部的 VM 二层网络互访，跨物理服务器的 VM 二层互访丢给传统的 Ethernet 接入层交换机去处理。这时有两个大的问题产生了，一是对于 VSwitch 的管理问题，前面说过大公司网络和服务一般是由两拨人负责的，这个东西是由谁来管理不好界定；二是性能问题，交换机在处理报文时候可以通过转发芯片完成 ACL packet-filter、Port Security (802.1X)、Netflow 和 QoS 等功能，如果都在 VSwitch 上实现，还是由服务器的 CPU 来处理，太消耗性能了，与使用 VM 提高服务器 CPU 使用效率的初衷不符。

Cisco 首先提出了 Nexus1000V 技术结构来解决前面的问题一，也只解决了问题一。为了解决问题二，IEEE (Institute of Electrical and Electronics Engineers) 标准组织提出了 802.1Qbg EVB (Edge Virtual Bridging) 和 802.1Qbh BPE (Bridge Port Extension) 两条标准路线了，Cisco 由 802.1Qbh 标准体系结构实现出来的具体技术就是 FEX+VN-Link。

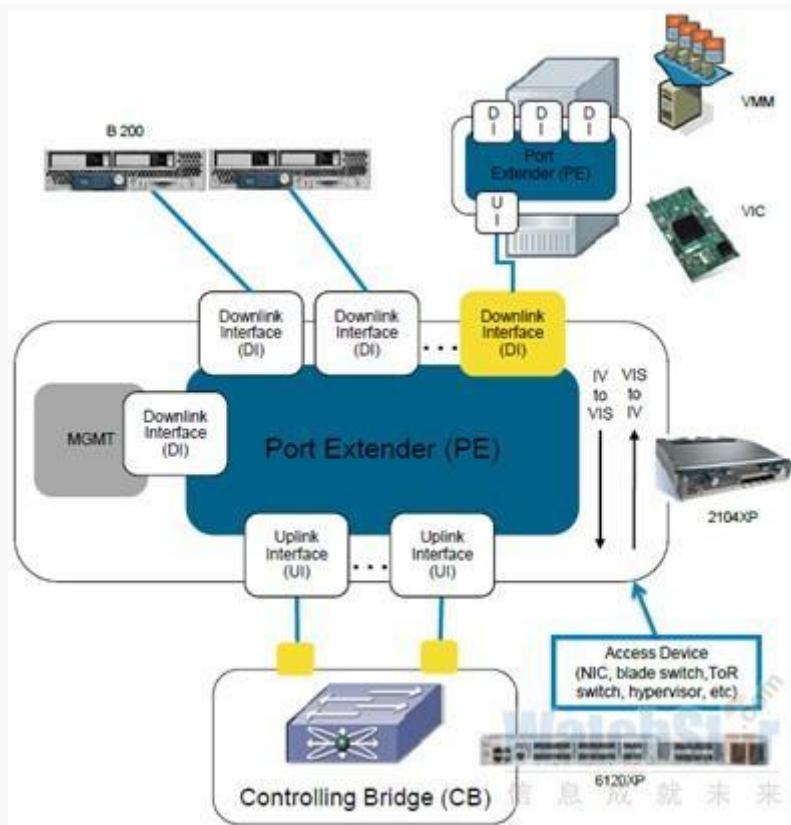
在数据通信世界只有两个阵营：Cisco 和非 Cisco。而就目前和能预见到的未来而言，非 Cisco 们都仍是 Cisco 的跟随者和挑战者，从数据中心新技术发展就可见一斑。在 VM 本地互访网络技术章节中会先介绍 Cisco 的相关技术与产品，再讲讲挑战者们的 EVB。

### 5.4.1 Cisco 接入层网络虚拟化

Cisco 在其所有的 VM 接入技术中都有两个主要思路：一是将网络相关内容都虚拟化为一个逻辑的框式交换机来集中由网络进行管理，二是给每个 VM 提供一个虚拟交换机接口 (vETH/VIF)。目的都是以网络为根，将枝叶一步步伸到服务器里面去。

#### 802.1Qbh

先来看下 802.1Qbh BPE (Bridge Port Extension)，下图是 Cisco 以 UCS 系列产品对应的结构图。



802.1Qbh 定义的是 VM 与接入层交换机之间的数据平面转发结构，不包括控制平面。这里可以将其看为一台虚拟的集中式框式交换机，其中 CB 可以理解为带转发芯片的主控板，PE 就是接口板。PE 进入服务器内部是通过硬件网卡来实现的，后续可能在 Hypervisor 层面会做软件 PE 来实现。Cisco 通过 FEX 来定义 CB 到 PE 以及 PE 到 PE 的关系，其数据平面是通过封装私有的 VN-Tag 头来进行寻址转发；通过 VN-Link 来定义 PE 的最终点 DI 到 VM 的 vNIC 之间的关系，提出了 Port Profile 来定制 DI 的配置内容。

在 802.1Qbh 结构中，整个网络是树状连接，每个 PE 只能上行连接到一个逻辑的 PE/CB，因此不存在环路，也就不需要类似于 STP 这种环路协议。所有的 VM 之间通信流量都要上送到 CB 进行查表转发，PE 不提供本地交换功能。PE 对从 DI 收到的单播报文只会封装 Tag 通过 UI 上送，UI 收到来的单播报文根据 Tag 找到对应的 DI 发送出去。对组播/广播报文根据 Tag 里面的组播标志位，CB 和 PE 均可以进行本地复制泛洪。更具体的转发处理流程请参考下文 Nexus5000+Nexus2000 的技术介绍。

Cisco 根据 802.1Qbh 结构在接入层一共虚拟出三台框式交换机，Nexus1000V (VSM+VEM)、Nexus5000+Nexus2000 和 UCS。其中 1000V 还是基于 Ethernet 传统交换技术的服务器内部软件交换机，没有 FEX，主要体现 VN-Link；而 Nexus5000+Nexus2000 则是工作于物理服务器之外的硬件交换机盒子，以 FEX 为主，VN-Link 基本没有；只有到 UCS 才通过服务器网卡+交换机盒子，完美的将 FEX+VN-Link 结合在一起。下面来逐台介绍。



## Cisco Nexus1000V

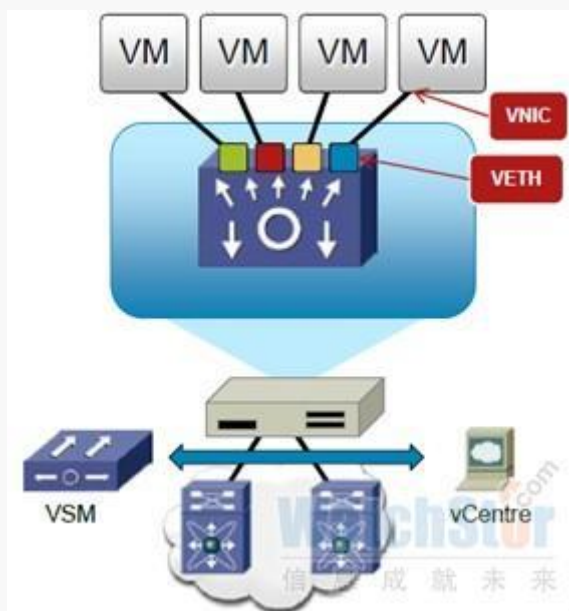
Nexus1000V 包含两个组件 VSM (Virtual Supervisor Module) 与 VEM (Virtual Ethernet Module)。看名字就能瞧出 VSM 对应机框交换机的主控板 Supervisor，而 VEM 对应其接口板。

VEM 就是一台安装运行在采用裸金属虚拟化结构的物理服务器中 Hypervisor 层次的软件交换机，其虚拟接口 vETH 分为连接 VM 虚拟网卡 vNIC 的下行接口和连接到每个物理网卡接口的上行接口，使用 Ethernet 基于 MAC 方式进行报文转发。由于其处于网络的末端，不需要运行 STP，通过不允许上行接口收到的报文从其他上行接口转发的规则来避免环路产生。与早期的 VSwitch 相比多了很多交换机相关功能。

VSM 则有两种形态，可以是独立的盒子，也可以是装在某个 OS 上的应用软件。要求 VSM 和 VEM 之间二层或三层可达，二层情况下 VSM 与 VEM 之间占用一个 VLAN 通过组播建立连接，三层情况下通过配置指定 IP 地址单播建立连接。VSM 是一个控制平台，对 VEM 上的 vETH 进行配置管理。通过 VSM 可以直接配置每台 VEM 的每个 vETH。

VSM 在管理 vETH 的时候引入了 Port Profile 的概念，简单理解就是个配置好的模板，好处是可以一次配置，四处关联。在 VM 跨物理服务器迁移时，VSM 就可以通过 vCenter 的通知了解到迁移发生，随之将 Port Profile 下发到 VM 迁移后对应的 vETH 上，使网络能够随 VM 迁移自适应变化。

VN-Link 是 CISCO 在虚拟接入层的关键技术，VN-Link=vNIC+vETH+Port Profile。

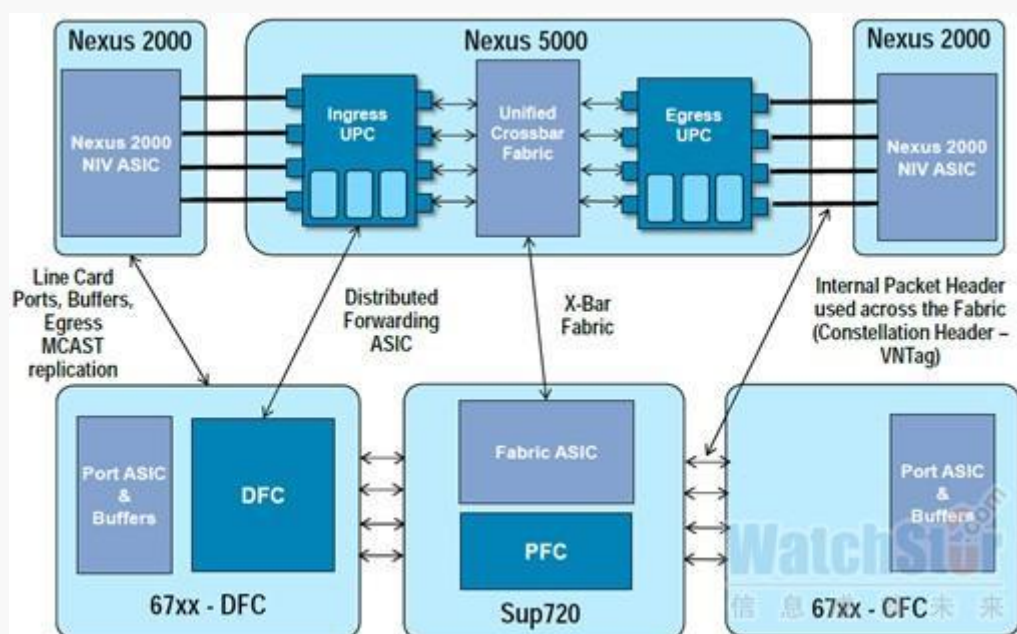


Nexus1000V 中的 vETH 是建立在软件交换机上的，而下文 UCS 系统里面的 vETH 就建立在 Cisco 的网卡硬件上了，对应到 UCS 虚拟交换机上就是 VIF (Virtual Interface)，同时 UCS 通过硬件实现可以把 FEX 里面要介绍的 VN-Tag 网络封装标识引入到物理服务器里面。

VEM 之间通过普通 Ethernet 交换机相连，跨 VEM 转发的流量也是传统以太网报文，因此 Nexus1000V 虽然可以理解为一台虚拟交换机，但不是集中式或分布式结构，也不存在交换芯片单元，仅仅是配置管理层面的虚拟化，属于对传统 VSwitch 的功能扩展，只解决了最开始提到的管理边界问题，但对服务器性能仍然存在极大耗费。

从产品与标准的发布时间上看，Nexus1000V 是先于 802.1Qbh 推出的，因此推测 Cisco 是先做了增强型的 VSwitch-Nexus1000V，然后才逐步理清思路去搞 802.1Qbh 的 BPE 架构。1000V 属于过渡性质的兼容产品，后续应该会对其做些大的改动，如改进成可支持 VN-Tag 封装的软件 PE，帮助 N5000+N2000 进入物理服务器内部，构造 FEX+VN-Link 的完整 802.1Qbh 结构。

### Nexus5000+Nexus2000



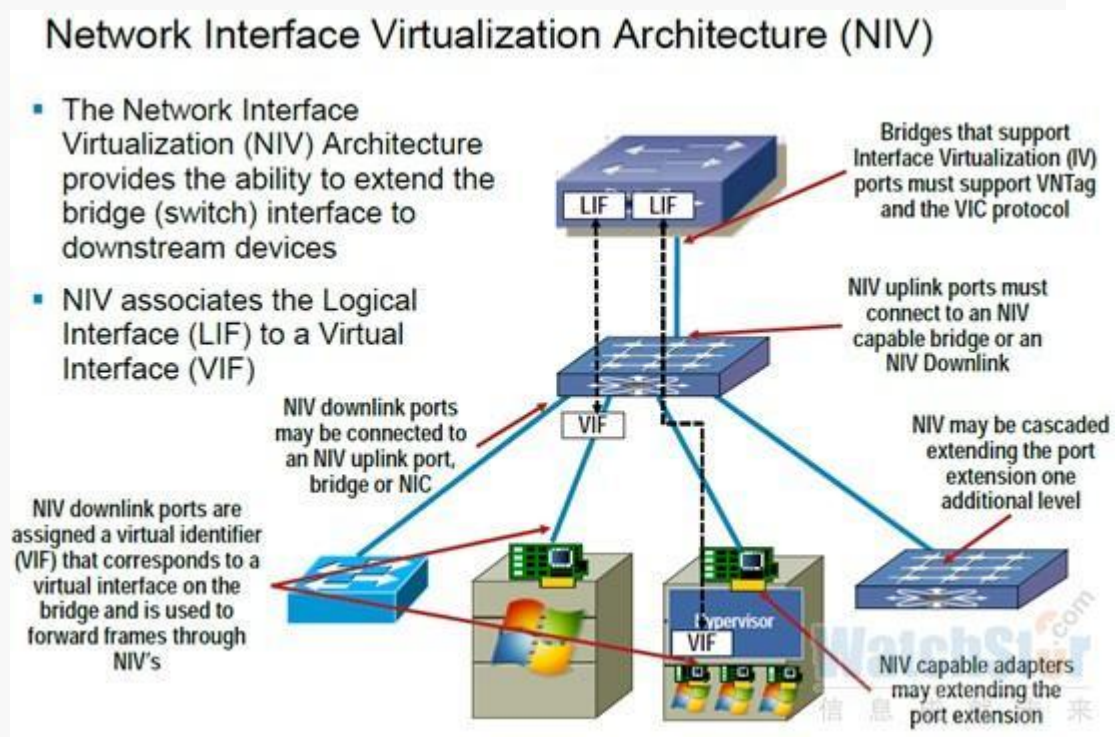
N5000+N2000 组成了一台集中式结构的虚拟交换机，集中式是指所有的流量都要经过 N5000 交互，N2000 不提供本地交换能力，只是作为 N5000 的接口扩展。对应 802.1Qbh 结构，N5000 就是 CB，而 N2000 就是 PE。组合出来的虚拟交换机中，N5000 就是带转发芯片和交换芯片的主控板，而 N2000 则是接口板，整体更像 Cisco 早期的 4500 系列机框或使用主控板 PFC 进行转发的 6500 系列机框，但是在 N5000 盒子内部又是以分布式结构处理转发芯片与交换芯片连接布局的，可参考如下的 N5000 和 6500 结构比较图。整了半天其实数据平面转发报文还是那几个步骤。

N5000+N2000 实现了 Cisco 的 FEX 典型结构（Fabric Extend，等同于 Port Extend）。在 N5000 上看到每台 N2000 就是以 FEX 节点形式出现的接口板。N2000 拥有两种物理

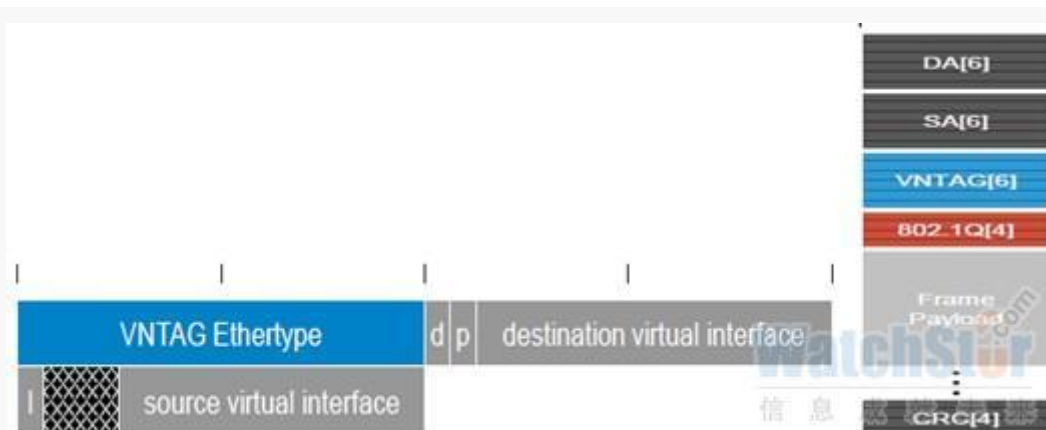
接口类型，连接下游设备（可以是服务器或 N2000，FEX 结构支持级联扩展）的 HIF（Host Interface）和连接上游 N5000 和 N2000 的 NIF（Network Interface），此两种接口是固定于面板上的，且角色不可变更。以 2248T 举例，右侧黄色接口为 NIF，其他为 HIF。



Cisco 将 FEX 结构又称为 Network Interface Virtualization Architecture (NIV)，在 NIV 中将 N2000 上的 HIF 称为 Virtual Interface (VIF)，将 N5000 上对应 HIF 的逻辑接口称为 Logical Interface (LIF)。截取 Cisco 胶片如下描述 NIV 的内容。



在 NIV 模型中所有的数据报文进入 VIF/LIF 时均会被封装 VN-Tag 传递，在从 VIF/LIF 离开系统前会剥离 VN-Tag。VN-Tag 就是在 FEX 内部寻址转发使用的标识，类似于前面框式交换机内部在转发芯片与交换芯片传输报文时定义槽位信息与接口信息的标识。VN-Tag 格式与封装位置如下：



d 位标识报文的走向，0 代表是由 N2000 发往 N5000 的上行流量，1 代表由 N5000 发往 N2000 的下行流量。

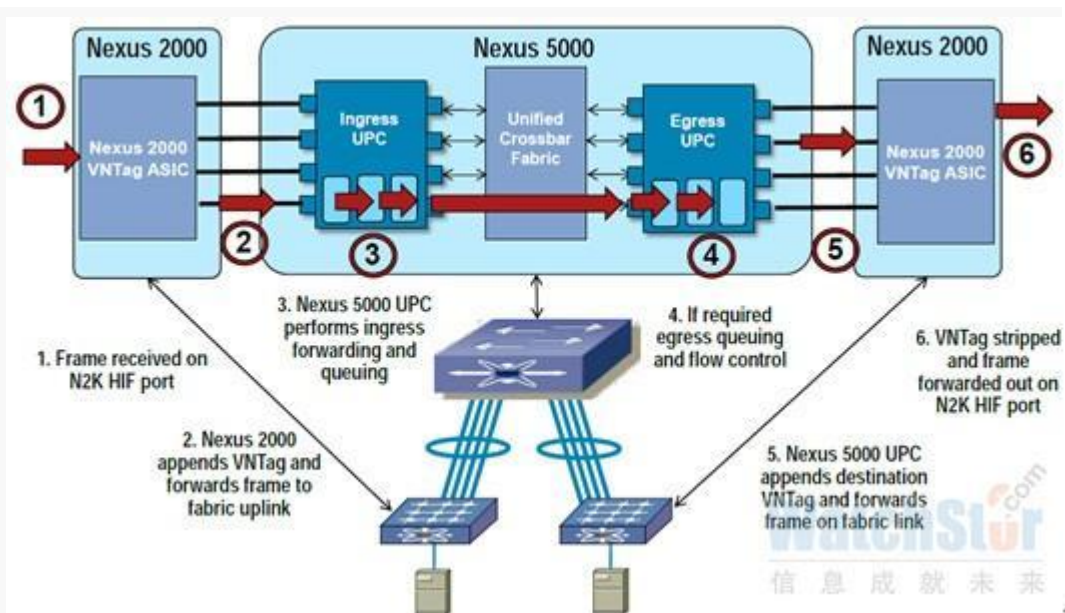
p 位标识报文复制，0 代表不需要复制，1 代表 N2000 收到此报文后需要向同 VLAN 的所有本地下行接口复制。此位只有 N5000 可以置位。

1 位标识报文是否返回给源 N2000，既 0 代表源和目的接口不在同一个 N2000 上，1 代表目的与源接口都在同一个 N2000 设备上。

DVI（Destination Virtual Interface）标识目的 HIF 接口，SVI（Source Virtual Interface）标识源 HIF 接口。每个 HIF 接口 ID 在一组 FEX 系统中都是唯一的。

流量转发时，N2000 首先从源 HIF 收到报文，只需要标识 SVI 的对应 HIF 信息，其他位都置 0 不用管，直接从 NIF 转发到 N5000 上即可。N5000 收到报文，记录源 HIF 接口与源 MAC 信息到转发表中，查 MAC 转发表，如果目的 MAC 对应非 LIF 接口则剥离 VN-Tag 按正常 Ethernet 转发处理；如果目的接口为 LIF 接口，则重新封装 VN-Tag。其中 DVI 对应目的 HIF，SVI 使用原始 SVI 信息（如果是从非 LIF 源接口来的报文则 SVI 置 0），d 位置 1，如果是组播报文则 p 位置 1，如果目的接口与源接口在一台 N2000 上则 1 位置 1。N2000 收到此报文后根据 DVI 标识查找本地目的出接口 HIF，剥离 VN-Tag 后进行转发，如果 p 位置 1 则本地复制转发给所有相关 HIF。每个 FEX 的组播转发表在 5000 上建立，所有 2000 上通过 IGMP-Snooping 同步。转发过程截取 Cisco 胶片介绍如下：



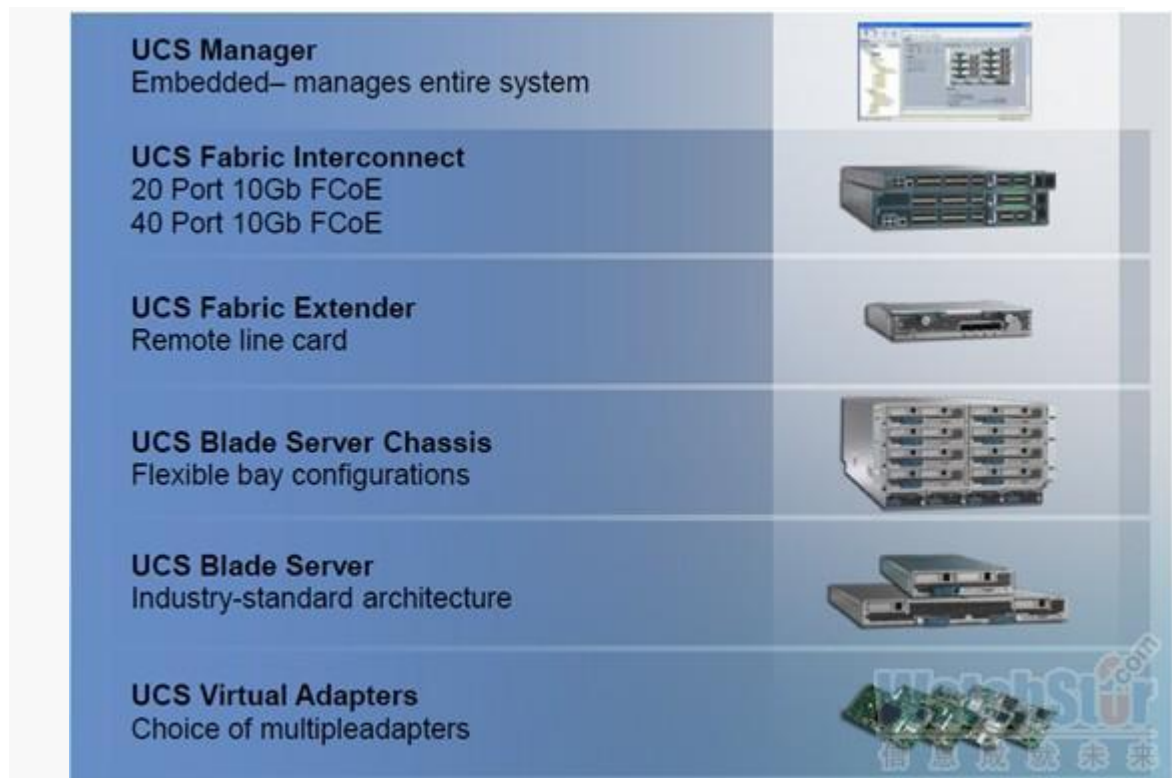


从前面 NIV 的结构图中可以看到 Cisco 希望将 FEX 通过网卡推进到服务器内部，但实际上目前阶段由于 Cisco 在服务器网卡方面的市场地位，这个还只是一个梦想。N2000 还是只能基于物理服务器的物理网卡为基本单元进行报文处理，搞不定 VM 的 vNIC，因此前面说 N5000+N2000 这台虚拟交换机只实现了 FEX，但没有 VN-Link。

好吧，搞不定服务器就搞不定网卡，更没有办法推行 FEX+VN-Link 的 802.1Qbh 理念。于是 Cisco 一狠心，先搞了套 UCS 出来。

## UCS

UCS (Unified Computing System) 是包括刀箱、服务器、网卡、接口扩展模块、接入交换机与管理软件集合的系统总称。这里的各个单元独立存在时虽然也能用，但就没有太大的价值了，与其他同档产品相比没有任何优势，只有和在一起才是 Cisco 征战天下的利器。UCS 产品结构如下图所示：

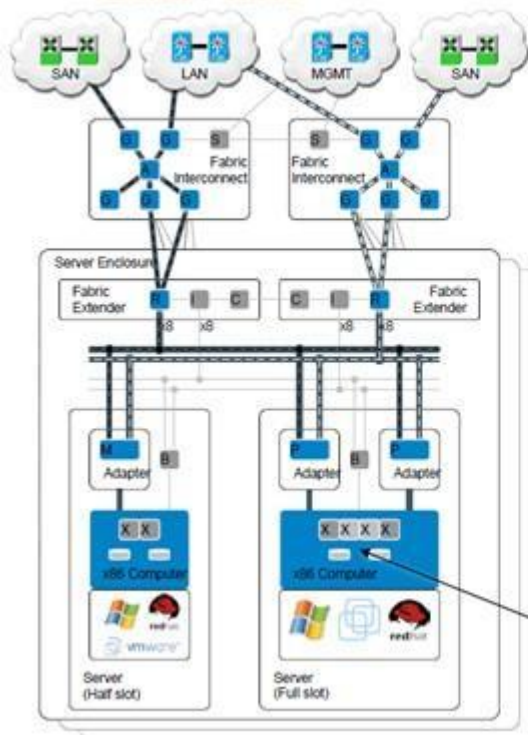


其中服务器、刀箱机框和管理软件都是标准的东西，没啥可多说的。关键部件就是网卡、交换机和刀箱的扩展线卡（Fabric Extender，这个名字个人觉得不好，容易与 FEX 结构混淆，可以叫个 FE Blade 或者 Fabric Line Card 什么的）。Interconnect 交换机对应 N5000，Fabric Extender 对应 N2000，这样加上 Virtual Adapters 就可以实现前面 NIV 结构中将 VIF（HIF）直接连到 VM 前的期望了，从而也就能完美实现 802.1Qbh BPE（Cisco FEX+VN-Link）的技术体系结构。

整个 UCS 系统结构就在下面这三张图中体现，分别对应数据平面（转发平面）、控制平面、管理平面。由于技术实现上和前面将的 FEX 和 VN-Link 没有大的区别，不再做重复赘述，有兴趣的同学可自行细琢磨。



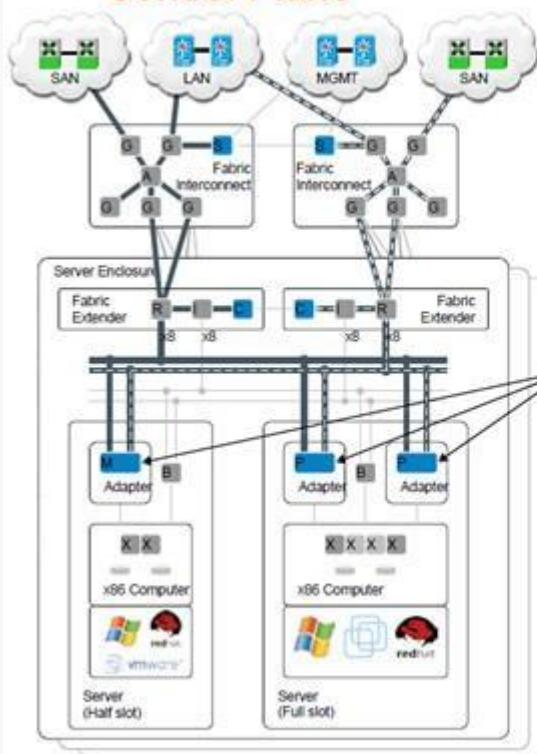
## Data Plane



- A) X-bar ASIC
  - 1.12 Tbps Xbar
  - 3 Unicast and 1 Multicast crosspoints
- G) Forwarding ASIC
  - XE/FC/GE Media Access Controllers
  - Forwarding - Ethernet, Fibre Channel, Multipath
  - Policy Engine
  - Packet Buffering
- R) VNTag ASIC
  - Host to uplink traffic engineering
  - Connectivity detection & management portal
- M) DCB ASIC
  - Couple Industry standard NICs/HBAs to ServerArray
- P) Virtualization (SRIOV/VNTAG) ASIC
  - Virtualized adapter for single OS and hypervisor systems
  - Ethernet and Fibre Channel vNICs
  - Direct Data Placement for Fibre Channel
- Memory Controller ASICs
  - Large memory configurations

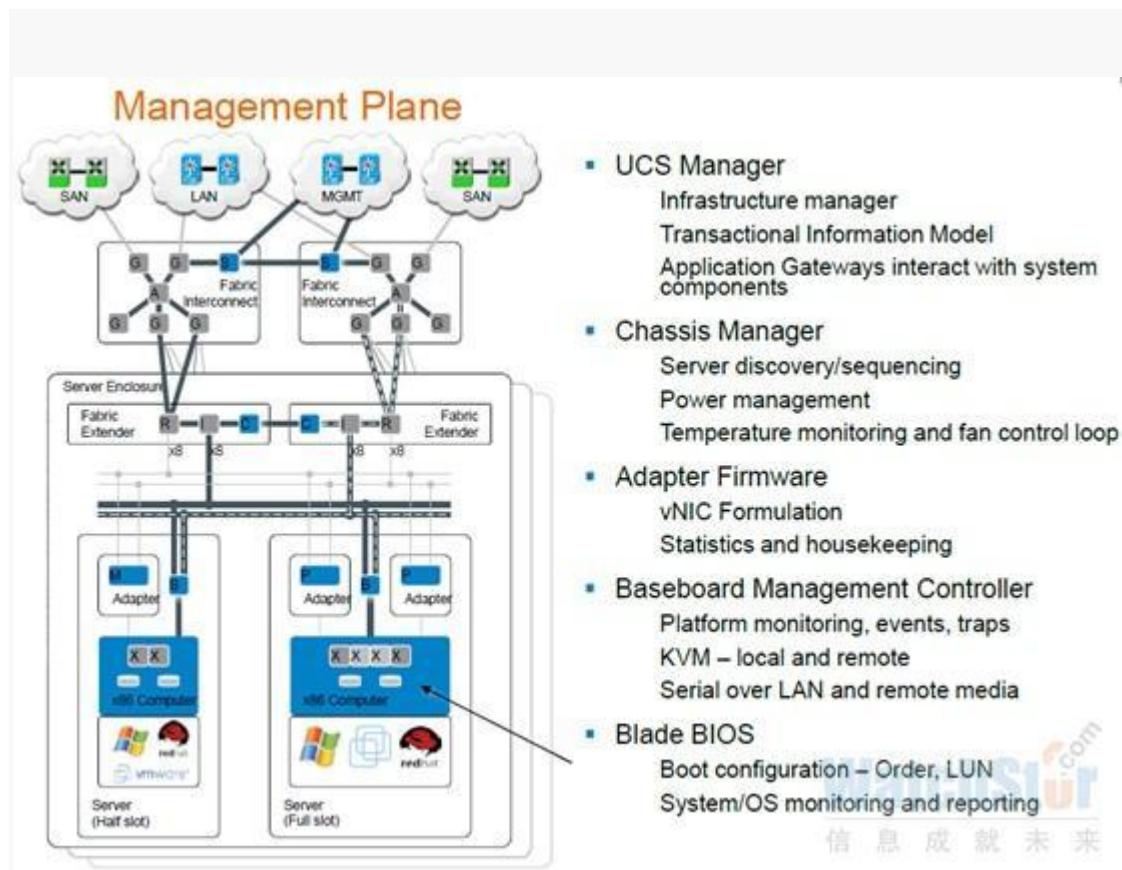
WatchStar.com  
信息成就未来

## Control Plane

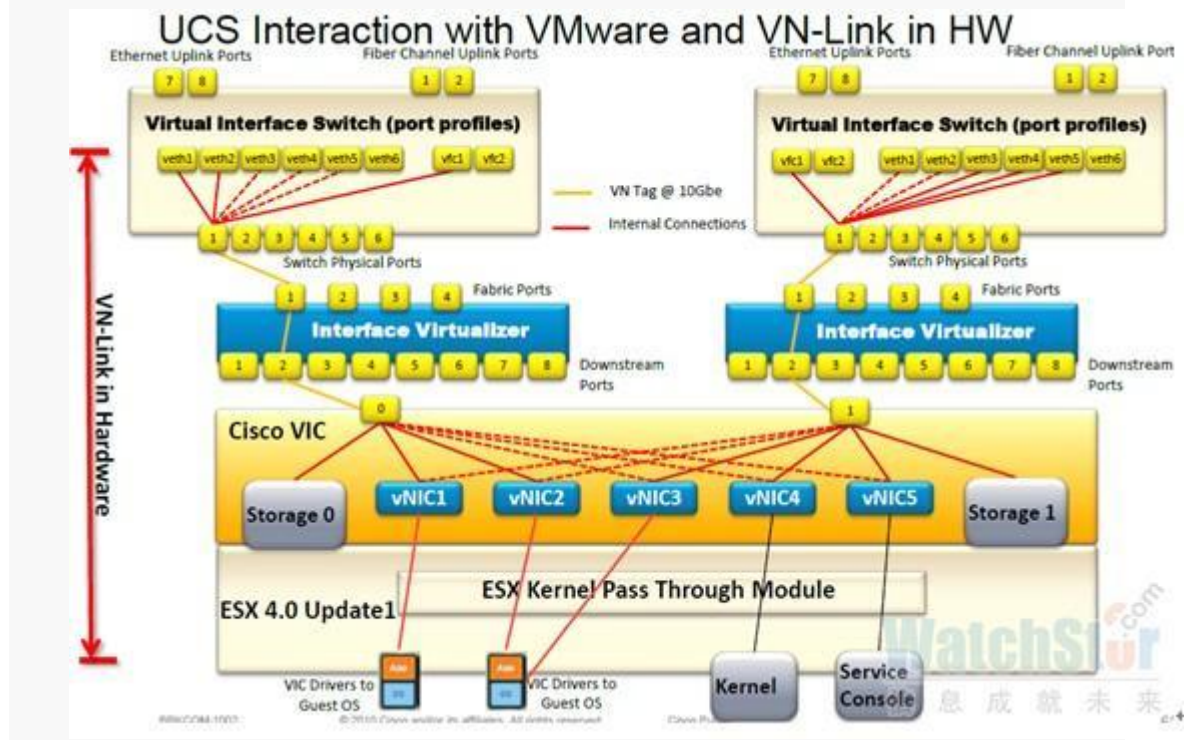


- Interconnect Supervisor
  - Infrastructure and Ethernet
  - Consolidated Ethernet/Fibre Channel
  - Network Interface Virtualization
  - Distributed Interconnect Fabric
- Fabric Extender
  - Fabric Connectivity
  - Satellite Interconnect ports and vNIC channels
- Adapter Firmware
  - Network controlled
  - Inaccessible from the host
  - vNIC instantiation
  - Fabric based balancing and failover
  - Fibre Channel/SCSI control suite (M81KR)

WatchStar.com  
信息成就未来



再补充一张 UCS VN-Link 的示意图，可以对应看一下三要素 vNIC、vEth 和 Port Profile 的位置连接关系。



UCS 就说这么多了，前面介绍的 Cisco 三台虚拟交换机如果铺开了讲每个都有上百页的内容，本文只是希望能够从结构上帮助大家理解，将作者的知识大厦框架拿出来与读者共享参考，至于每个人的楼要怎么盖还需自己去添砖加瓦。包括下文的技术点讲解也是如此，作者会将自己认为最重要的关键部分讲出来，细节不会过于展开。

### 5.4.2 802.1Qbg EVB

说完了 Cisco 再说说非 Cisco 阵营，在如 802.1Qbg EVB 和 802.1aq SPB 等所谓挑战技术的参与编纂者中，都会看到 Cisco 的身影。如下图为 2009 年 IEEE Atlanta, GA 时发出的 EVB 所有撰稿相关人名单。



The image shows a slide titled "IEEE 802 Contributors and Supporters". It lists 24 individuals and their affiliations, organized into two columns. The affiliations include various companies and organizations like 3Com, Dell, HP, Fulcrum, Cisco, Broadcom, Intel, Brocade, Neterion, PMC-Sierra, IBM, Marvell, Extreme, BNT, InMon, Solarflare, and QLogic. A watermark "Watchdog 信息成就未来" is visible in the bottom right corner of the slide.

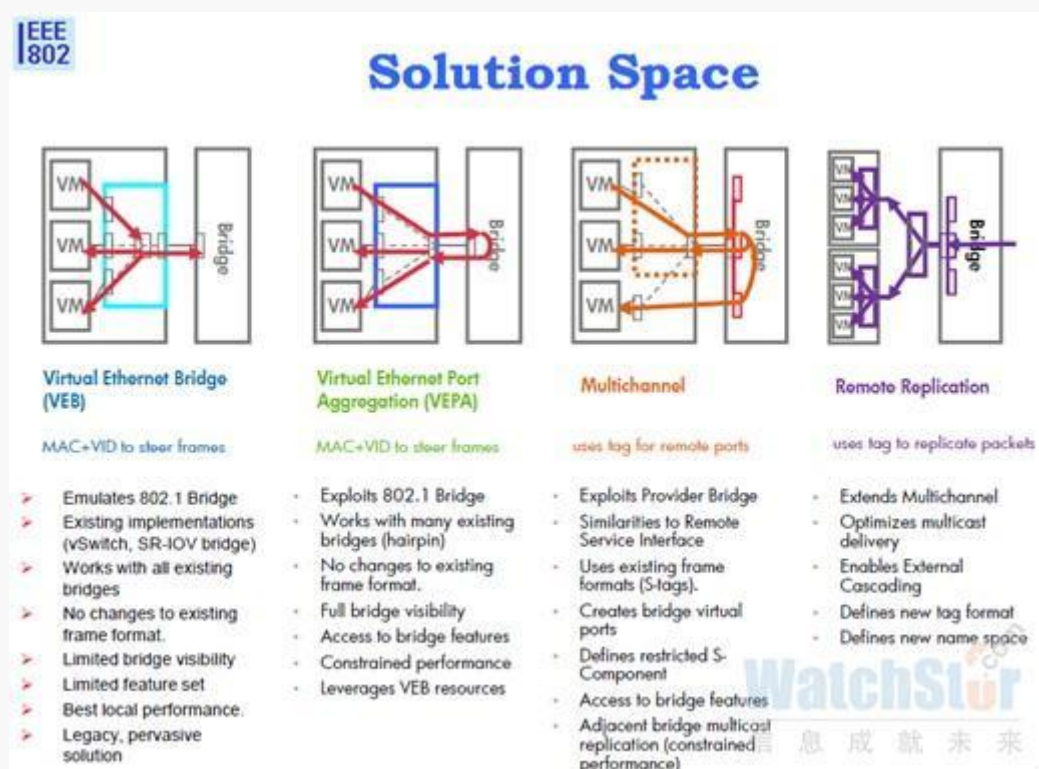
Contributors and Supporters	
Siamack Ayandeh	(3Com)
Guarav Chawla	(Dell)
Paul Congdon	(HP)
Dan Daly	(Fulcrum)
Claudio DeSanti	(Cisco)
Uri Elzur	(Broadcom)
Norm Finn	(Cisco)
Ilango Ganga	(Intel)
Anoop Ghanwani	(Brocade)
Leonid Grossman	(Neterion)
Chuck Hudson	(HP)
Brian L'Ecuyer	(PMC-Sierra)
Pankaj K Jha	(Brocade)
Jeffry Lynch	(IBM)
David Koenen	(HP)
Charles R. (Rick) Maule	(consultant)
Menu Menuchehry	(Marvell)
Shehzad Merchant	(Extreme)
Vijoy Pandey	(BNT)
Joe Pelissier	(Cisco)
Peter Phaal	(InMon)
Renato Recio	(IBM)
Rakesh Sharma	(IBM)
Jeelani Syed	(Juniper)
Patricia Thaler	(Broadcom)
Neil Turton	(Solarflare)
Manoj Wadekar	(QLogic)
Martin White	(Marvell)
Robert Winter	(Dell)

802.1Qbg 当时的主要撰写人是 HP 的 Paul Congdon，不过最近几稿主要 Draft 已经由 Paul Bottorff 取代。其中 Cisco 的 Joe Pelissier 也是 802.1Qbh 的主要撰写人，而 Bottorff 也同样参与了 802.1Qbh 的撰写工作。单从技术上讲，这二者并不是对立的，而是可以互补的，上述两位 HP 和 Cisco 的达人都正在为两种技术结构融合共存而努力。具体可以访问 <http://www.ieee802.org/1/pages/dcbridges.html> 对这两个处于 Active 阶段的 Draft 进行学习。

802.1Qbh 通过定义新的 Tag (VN-Tag) 来进行接口扩展，这样就需要交换机使用新的转发芯片能够识别并基于此新定义 Tag 进行转发，因此目前除了 Cisco 自己做的芯片外，其他厂商都无法支持。只有等 Broadcom 和 Marvel 等芯片厂商的公共转发芯片也支持了，大家才能跟进做产品，这就是设备厂商有没有芯片开发能力的区别。而 802.1Qbg 就走了另外一条路，搞不定交换机转发芯片就先想办法搞定服务器吧。下面从 IEEE 截取的图中可以看



到 EVB 的四个主要组成部分，也可以看做四个发展阶段。当前处于 VEPA 的成长期，已经出现部分转化完成的产品，而 Multichannel 还在产品转化前的研究状态。



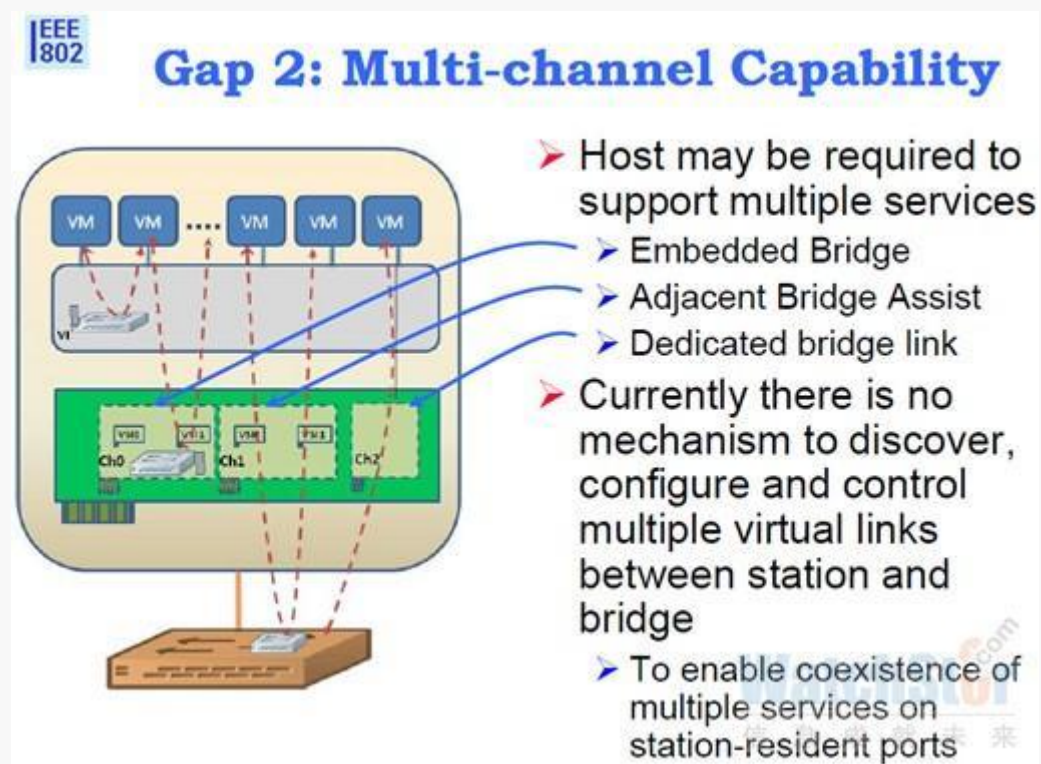
先说 VEB，这个最好理解，就是定义了物理服务器内部的软、硬件交换机。软件交换机就是前面提过的 VSwitch，硬件交换机就是从 SR-IOV 演进来的网卡交换机。SR-IOV 已经可以使 VM 的 vNIC 在物理网卡上一一对应通道化了，那么再加个转发芯片基本就可以做成最简单的交换机了，当然这只是原理上可行，实际中作者还没有见过成熟产品。VEB 与普通 Ethernet 交换机的最大区别是定义了连接交换机的上行口与连接 VM 的下行口，而 VEB 的上行口间是不允许相互转发报文的，这样可以在不支持 STP 的情况下保证无环路产生。Cisco 的 N1000V 就可以认为是个 VEB。VEB 的优点是好实现，在 Hypervisor 层面开发软件或者改造网卡就可以出成品，缺点是不管软件的还是硬件的相比较传统交换机来说能力和性能都偏弱，网卡上就那么点儿大的地方，能放多少 CPU、TCAM 和 ASIC 啊。

于是有了 VEPA，VEPA 比 VEB 更简单，不提供 VM 间的交换功能，只要是 VM 来的报文都直接扔到接入交换机上去，只有接入交换机来的报文才查表进行内部转发，同样不允许上行接口间的报文互转。这样首先是性能提升了，去掉了 VM 访问外部网络的流量查表动作。其次是将网络方面的功能都扔回给接入层交换机去干了，包括 VM 间互访的流量。

这样不但对整体转发的能力和性能有所提升，而且还解决了前面最开始 VSwitch 所提出的网络与服务器管理边界的问题。相比 Cisco 将网络管理推到 VM 的 vNIC 前的思路，这种做法更传统一些，将网络管理边界仍然阻拦在服务器外面，明显是出于服务器厂商的思路。在传统 Ethernet 中，要求交换机对从某接口收到的流量不能再从这个接口发出去，以避免环

路风暴的发生。而 VEPA 的使用要求对此方式做出改变，否则 VM 之间互访流量无法通过。对交换机厂商来说，这个改变是轻而易举的，只要变动一下 ASIC 的处理规则即可，不需要像 VN-Tag 那样更新整个转发芯片。从理论上讲，如 VEB 一样，VEPA 同样可以由支持 SR-IOV 的网卡来硬件实现，而且由于需要实现的功能更少，因此也更好做一些。个人认为 VEPA 的网卡可能会先于 VEB 的网卡流行起来。

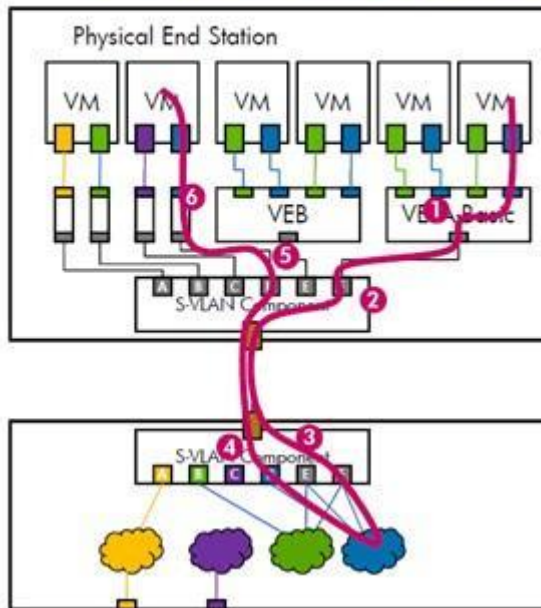
下面说说 Multichannel，这个东西就有点儿意思了。802.1Qbg 的说法是在混杂场景中，物理服务器中同时有 VEB、VEPA 和需要直接通过 SR-IOV 连到交换机的 vNIC，而当前对这多种流量在网卡到交换机这条链路上是无法区分识别的，于是整出个 Multichannel。参见 IEEE 的胶片原文如下：



想在一条通道内对相同类型流量进行更细的分类，看了前面技术理解一节大家应该有个思路了，加 Tag 呗。Multichannel 借用了 QinQ 中的 S-VLAN Tag（就是个 VLAN 标签）。在数据报文从网卡或交换机接口发出时封装，从对端接口收到后剥离。简单的转发过程如下：

## MultiChannel Approach

Example: VM through VEPA to Directly Accessible VSI



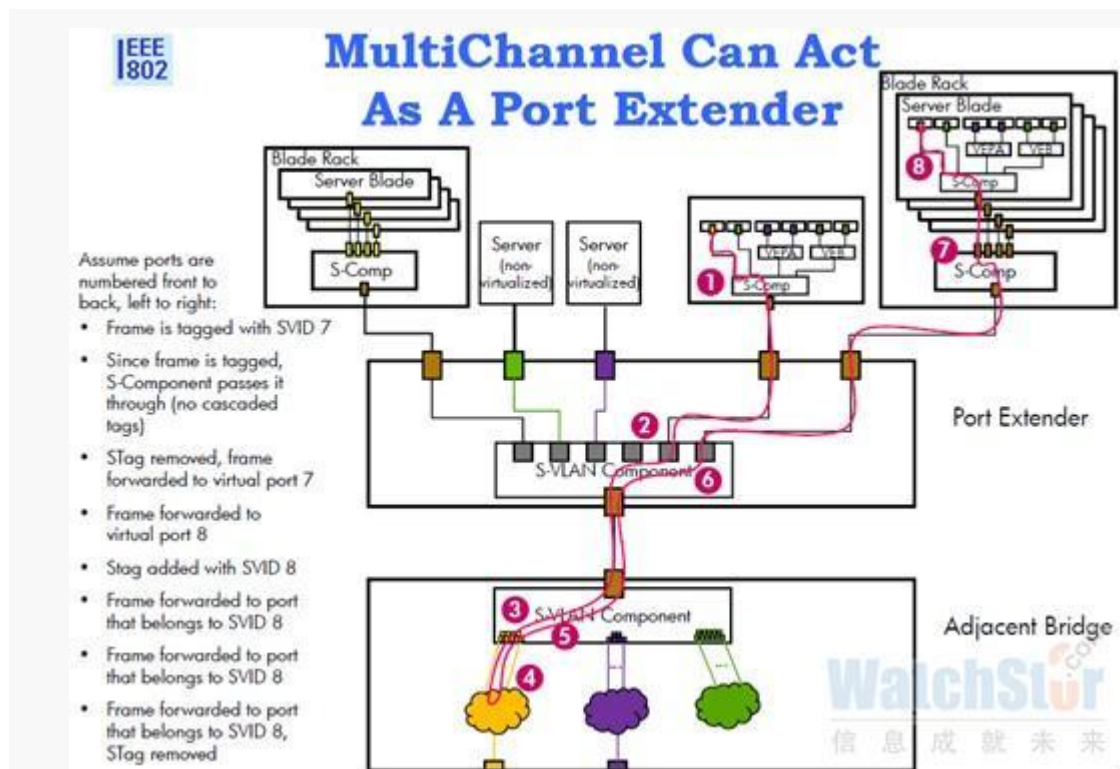
1. VEPA ingress frame from VM forwarded out VEPA uplink to S-Component
2. Station S-Component adds SVID (F)
3. Bridge S-Component removes SVID and forwards to port F
4. Frame is forward back to port D, S-Component adds SVID D
5. Station S-Component removes SVID D
6. S-Component forwards frame on Port D on Blue VLAN.

WatchStar.com  
信息成就未来

诸位看官看到这里可能会产生疑问，这个和 Cisco 的 BPE 很像啊，无非是用 S-VLAN 取代了 VN-Tag 作用在网卡和交换机之间。作者个人觉得 Multichannel 真正瞄准的目标也不是什么多 VEB 和 VEPA 之间的混杂组网，至少目前做虚拟化的 X86 服务器上还没有看到这种混杂应用的需求场景。

真正的目标应该就是通过 S-VLAN Tag 建立一条 VM 上 vNIC 到交换机虚拟接口的通道，和 Cisco FEX+VN-Link 的目标是等同的，只是没有考虑网络接入层上面的 FEX 扩展而已。Cisco 的达人 Joe Pelissier 目前在 EVB 工作组中做的事情也是将其与 802.1Qbh 在 Port Extend 方面做的尽量规则一致。可参考如下胶片内容：

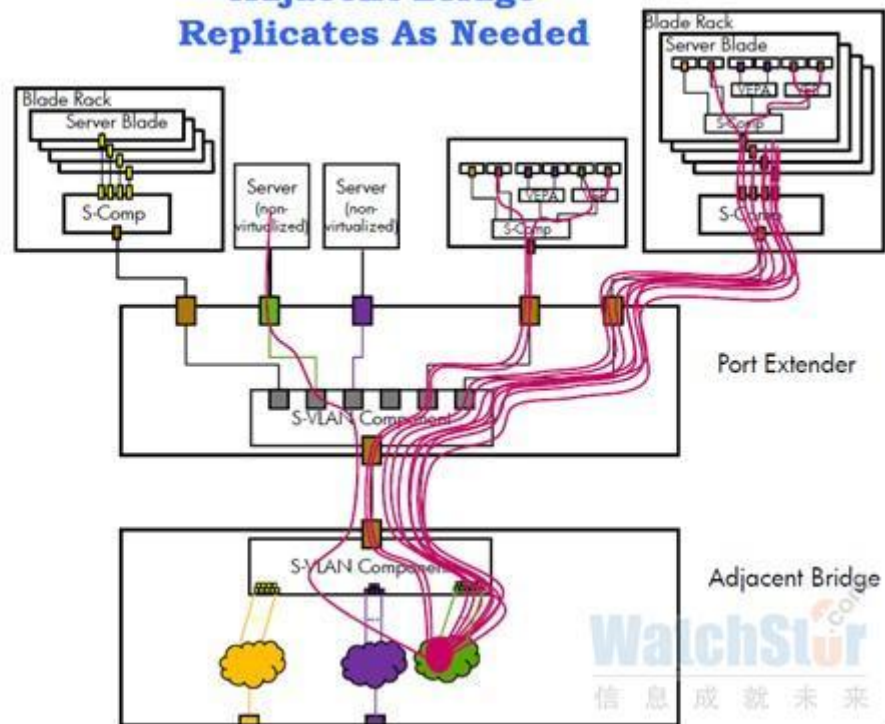




Multichannel 相比 VN-Tag 的优势是交换机目前大部分的转发芯片就已经支持多层 VLAN 标签封装的 QinQ 技术了，而网卡封装 VLAN Tag 也是现成的，只要从处理规则上进行一些改动就可以完全实现。但由于其未考虑网络方面的扩展，S-VLAN 还不能进行交换机透传，只能在第一跳交换机终结，所以从接入层网络部署规模上很难与 FEX 抗衡。

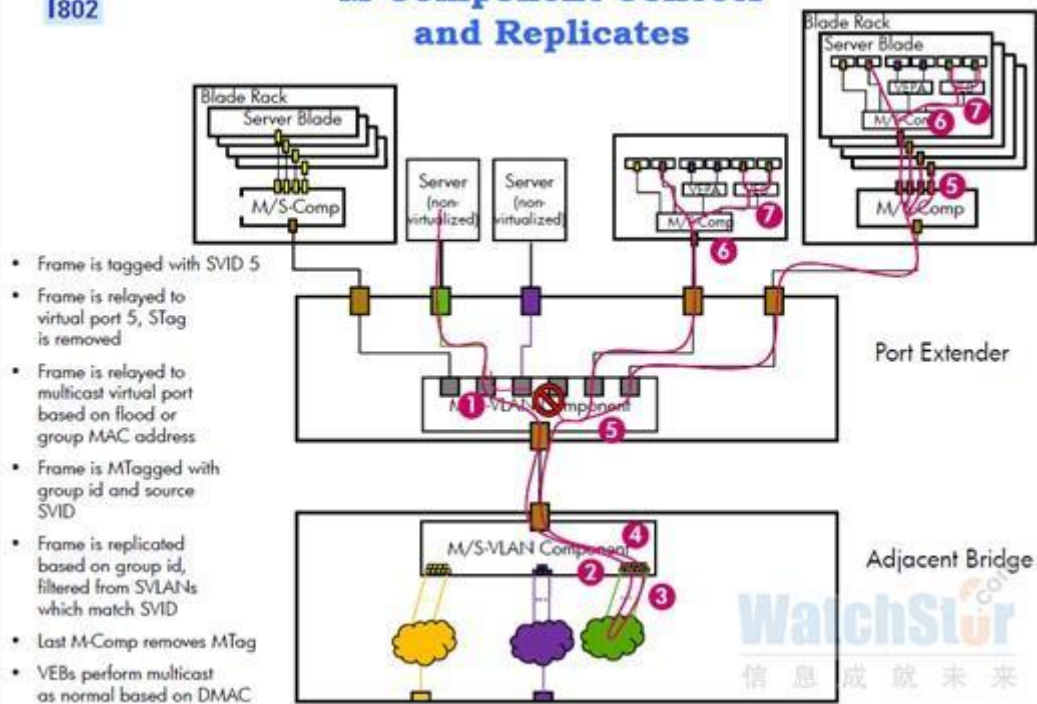
最后一个是 Remote Replication 复制问题。Ethernet 网络当中广播、组播和未知单播报文都需要复制，而前面的 Multichannel 结构所有的复制工作都会在交换机完成，于是会造成带宽和资源的极大浪费，如下图所示：

## Adjacent Bridge Replicates As Needed



此时需要定制一个标志位，以通知每个 S-VLAN 组件进行本地复制。当前在 802.1Qbg 中此标志位叫做 M 标识，其实和 VN-Tag 字段中的 p 标志位一个作用，因此这块也是由 Cisco 的达人 Joe Pelissier 在完善。M 标识的位置和作用如下图所示：

## M-Component Collects and Replicates



整个 EVB 就这些内容了，虽然目前还是在不断改动，但都是些实现细节方面的东西，大的框架结构应该就如上所述不会再变化了。从工作内容上就可以看出，HP 等服务器厂商的思维范畴就到 VEPA 了，Multichannel 只是提出个概念，实际上后续的东东还都要是借鉴 Cisco 的网络思路来实现，个人甚至怀疑连 Multichannel 都可能是 Pelissier 提出的，谁专业谁知道。从大体上来说，EVB 总的思路就是希望在尽量对现有设备最小变更的情况下解决 VM 接入互访的问题，从改变 Ethernet 交换机接口转发规则到增加 S-VLAN 标签都是如此，但到 M 标识就不见得还能控制得住了。不过就目前协议技术完善和进行产品转化的速度来看，还有得时间进行考虑变化。

### 5.4.3 小结

又到了小结的时间。该有人问了，讲了半天 802.1Qbh 和 802.1Qbg 这两个技术体系谁优谁劣，谁胜谁败啊？作者不是半仙也不是裁判，只能推测无法判断。从技术角度讲，尤其从网络技术角度讲，802.1Qbh BPE 提出了一整套的网络虚拟化解决方案，而 802.1Qbg EVB 则只是提出了解决几个 VM 网络接入问题的办法，二者的技术深度不可同日而语。然而对于市场上来说，一时的技术优势并不能完全左右胜负，各方博弈会使结局充满不可预测性，Nortel 和 3Com 的没落就是例子。从一名技术至上者的角度出发，作者更倾向于 Cisco，不过一切都有待于市场的检验。当然在这个世界上还有些地方存在可以代表市场直接做出裁决的裁判们，他们的裁决结果看看各公司在当地的财务贡献就可以简单预测，和真正的市场选择有没有关系，你懂的。

顺便说一句，802.1Qbh 需要变更交换机的转发芯片以适应 VN-Tag 转发，而 802.1Qbg 的 VEPA 和 Multichannel 目前则只需要交换机做做软件驱动方面的变动即可支持。不论将来谁成为了市场技术主导，大家觉得 Cisco 设备同时支持两套标准会有什么难度。从市场发展上大胆预测一下，Cisco 会在 802.1Qbg 标准成熟后，在新一代 N2000 位置产品上实现对 S-VLAN 组件和 M-Tag 的支持，以后的主流结构就是服务器内部用 VEPA+SR-IOV，网卡和 N2000 之间使用 S-VLAN 区分通道，N2000 再往上到 N5000 还是封装 VN-Tag 的 FEX。BPE+EVB 才是王道。

YY 一下，还有木有啥其他的技术可以起到类似的作用呢？目前 802.1Qbg 和 802.1Qbh 都是通过定义一个 Tag 来为接入层交换机标识 VM（VN-Tag/S-VLAN），那么实际上还有个现成的可用 Tag，就是 MAC，每个 VM 的 vNIC 都拥有独一无二的 MAC，那么是否可以让交换机根据源和目的 MAC 来建立 VIF 去对应处理每个 VM 的 vNIC 流量呢。当然细节上还要设计很多机制来保障各种情况的正常处理，但是暂时从方向上感觉应该有些搞头，回头有时间细琢磨一下吧。能不能成标准无所谓，当做头脑游戏锻炼思维了。

短期来看，上述厂家标准有得一争，但从长远来看，其实硬件交换机进入服务器内部才是王道。毕竟转发芯片会越来越便宜，性能会越来越高，再发展几年，不管是放在网卡上还是集成在主板上都没有太大难度。

## 5.5 Ethernet 与 FC 网络融合技术-FCoE

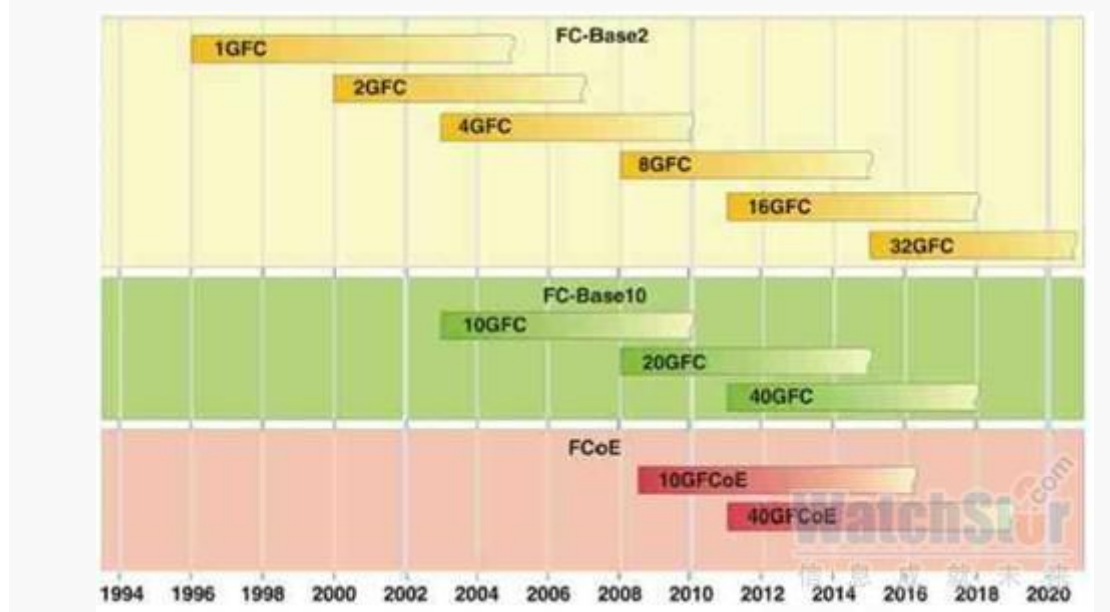
本章节重点技术名词：FC/FC ID/FCoE/FIP/FCF/DCB/NPV

服务器后端连接存储设备的 FC 网络与前端 Ethernet 网络融合是目前传统以太网交换机厂商进军后端存储网络的阳谋，Cisco 称其为统一 IO。下面会先介绍下 FC，然后是 FCoE，最后说说 NPV。

### 5.5.1 FC

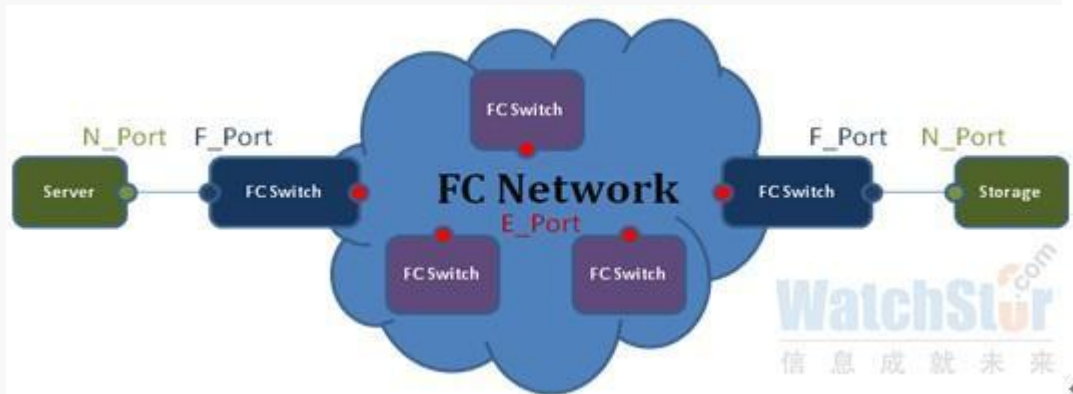
FC（Fibre Channel）在 1994 年由 ANSI T11 制定首个标准，从开始就达到 1G 带宽，大家可以想想那时候 Ethernet 还是什么年代。同时由于其无丢包的协议特性，受到了存储网络的青睐，成为 Server 到 Storage 之间 SAN 网络的霸主，相信由于 Ethernet 40G/100G 的缓慢发展速度，FC 的地位至少在下个 5 年内还是无可动摇的。在 FC 通信领域是典型的寡头独占市场，前三位中 Brocade 占据了 70% 的市场份额，Cisco 占 20%，Qlogic 占个位数。2010 年统计整个 FC Switch 市场销售额约 \$2B，对众多的传统交换机厂商来说，谁不想通过 FCoE 进去分杯羹呢，即使是 Cisco 也忍不住想翻身把歌唱。

FC 包含 Base2 与 Base10 两套演进道路，Base10 主要应用在 FC 交换机之间，使用较少，已经快消亡了。下面截图可以看出其演进情况。



FC 拥有自己的独立层次结构，FC-0 到 FC-4 对应 OSI 模型的 1-5 层，但也并非一一对应，完整协议内容请大家自行查阅标准文档。其中 FC-2 定义了数据通信的内容，是与网络方面息息相关的，下面介绍的内容也都是以 FC-2 为主。

在 FC 网络中一共有三种主要的接口角色，NPort，FPort 和 EPort，其中 N 是服务器或存储等终端节点连接 FC 网络的接口，F 是 FC 交换机设备连接服务器或存储等终端节点的接口，E 是 FC 交换机互联接口。还记得前面技术理解里面的典型结构么？



### FC 设备都拥有 2 个重要标识：

**WWN（World Wide Name）：**64bit，节点和每个接口都有各自固定的 WWN 且所有的 WWN 均是唯一的，WWN 的作用是为了身份识别和安全控制，有些类似于 MAC，但不做转发寻址使用。

**FC ID：**24bit，由 8 个 bit 的 Domain ID，8bit 的 Area ID 和 8bit 的 Port ID 组成，每个 Domain ID 代表一台 FC Switch（由此可以算出每个 FC 网络最多支持 256 个 Switch 节点，减去部分保留 ID，实际能够支持最多 239 个 Switch）。终端节点的 FC ID 是基于接口的，每个 NPort 的 FC ID 是由直连的 FC Switch 动态分配。FC ID 的主要作用就是供数据报文在 FC 网络中寻址转发。

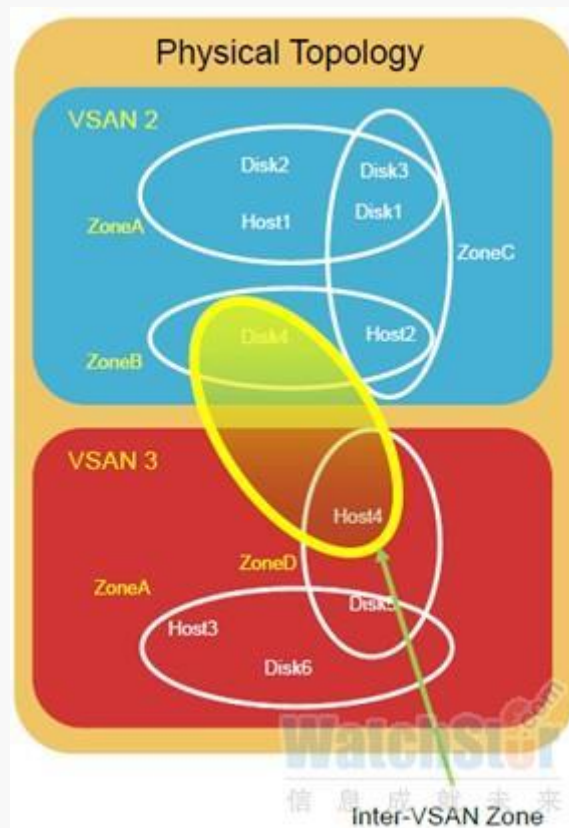
有了标识的 Tag，那么就需要个动态协议供 FC Switch 互相学习了，FC 网络使用 FSPF（Fabric Shortest Path First）进行 FC ID 的寻址学习，看名字就知道其协议机制和 OSPF 没有什么大的区别，不多说了。

FC 网络中的另外两个重要概念就是 VSAN 和 Zone，VSAN 和 VLAN 很类似，都是手工配置的，不同 VSAN 的流量相互隔离，这样不同 VSAN 中可以分配相同的 FC ID。而且由于 VSAN 是非公有定义协议字段，各个厂家实现并不见得一致，因此实际的 FC 组网中很难见到不同厂商设备的混合组网。Zone 则是类似于 ACL 的安全特性，配置为同一个 Zone 的成员可以互访，不同 Zone 的就会被隔离。Zone 是作用于 VSAN 内部的，可以理解 VSAN 是底层物理隔离，Zone 是上层逻辑分隔。同一个设备节点可以属于不同的 Zone，Zone 成员以 WWN 进行标识，可以简单类比为 ACL 中的同一个源 IP 地址可以配置在不同的 Rule 中



对应不同的目的 IP，以匹配不同的流量。Zone 的控制可以是软件实现，也有相应的 ASIC 可以做硬件处理。

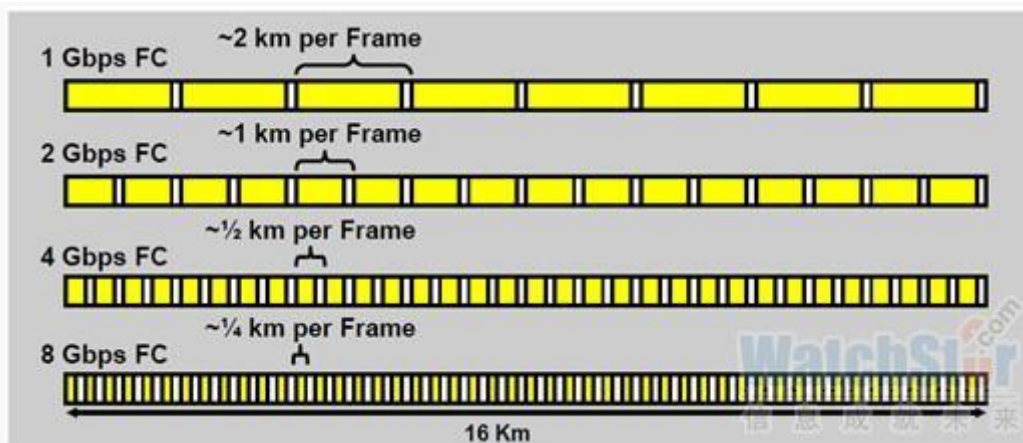
当然有隔离还得能互通，就好像做了 VPN 后还惦记着跨 VPN 互访，有了 FW 还搞个 HTTPS 翻墙(目前又在研究怎么控制 HTTPS 内容了),FC 中也有 IVR(Inter-VSAN Routing) Zone 的概念，就是通过一些静态配置的手段翻越已有的 VSAN 隔离。人们总是处在不断的制造藩篱和打破藩篱的循环中。Zone、VSAN 和 IVR Zone 的关系如下 Cisco 资料截图所示：



FC 技术体系还有最重要的一个关键流控技术 Buffer to Buffer Credits 用来确保无丢包转发。BB Credits 和 TCP 滑动窗口相似，规则很简单，两个相邻 FC 节点在连接初始化的时候先协商一个度量收包设备 Buffer 大小的数值 N 出来，发包设备每发一个数据报文就做 N-1，收包设备每收一个报文就回一个 R\_RDY 报文回来，发包设备每收到一个 R\_RDY 就做 N+1，当 N=0 时，发包设备就停止发包。

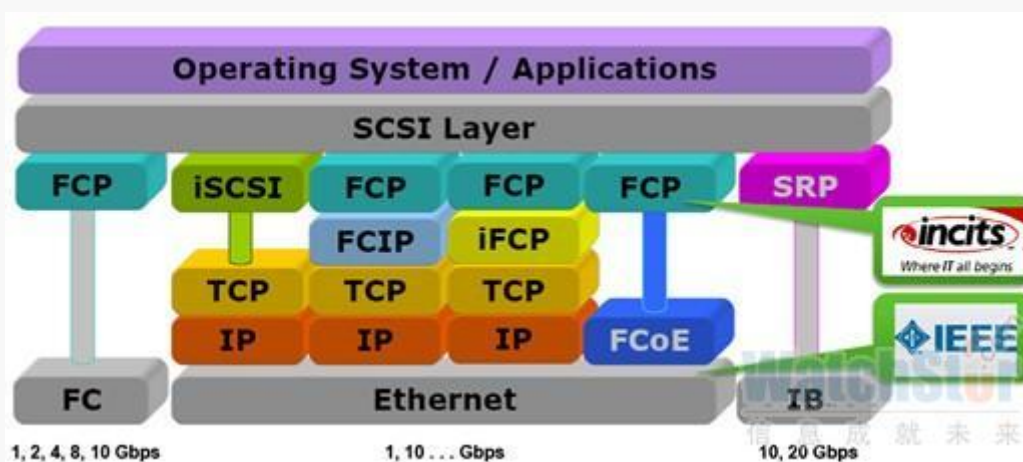
这样当突发拥塞时，上游设备们都把报文存在本地缓存中等着，下游有空间时再发，可以最简单的避免丢包。BB Credits 是以报文数目衡量 buffer 能力，与报文长度无关（FC 报文最大长度 2112Byte）。另外 Credits 协商数目大小与带宽和距离存在比率关系，可参考如下图示的 Cisco 建议：





FC 设备（一般指服务器，称为 Initiator）在传输数据之前需要进行两步注册动作，NPort 先通过 FLOGI（Fabric Login）注册到最近的 Fabric 交换机上，获取 FC ID 及其他一些服务参数并初始化 BB Credits。然后再通过 PLOGI（Port Login）注册到远端的目的设备（一般指存储，称为 Target）的 NPort 上建立连接，并在 P2P 直连的拓扑下初始化 BB Credits。

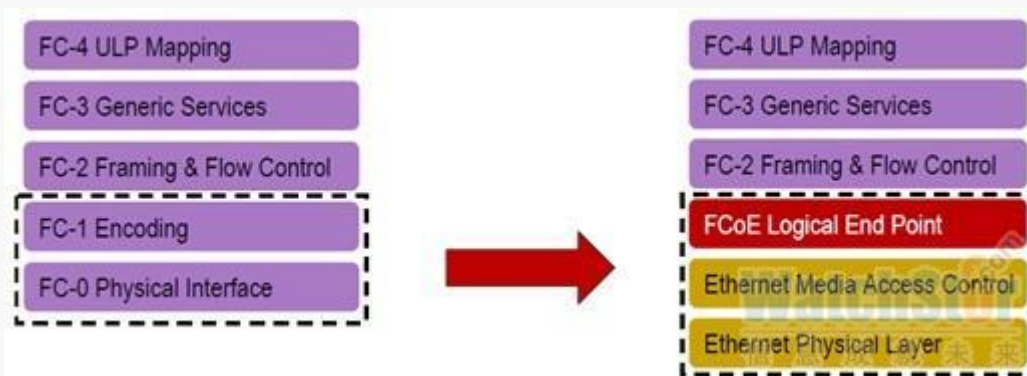
FC 从标准建立伊始就开始被研究跨传统 TCP/IP/Ethernet 网络传播，目前主要有 iSCSI（IP SAN）、FCIP、iFCP 和 FCoE 四条道路。其中 FCIP 和 iFCP 应用最少，iSCSI 缓慢增长，FCoE 后来居上。



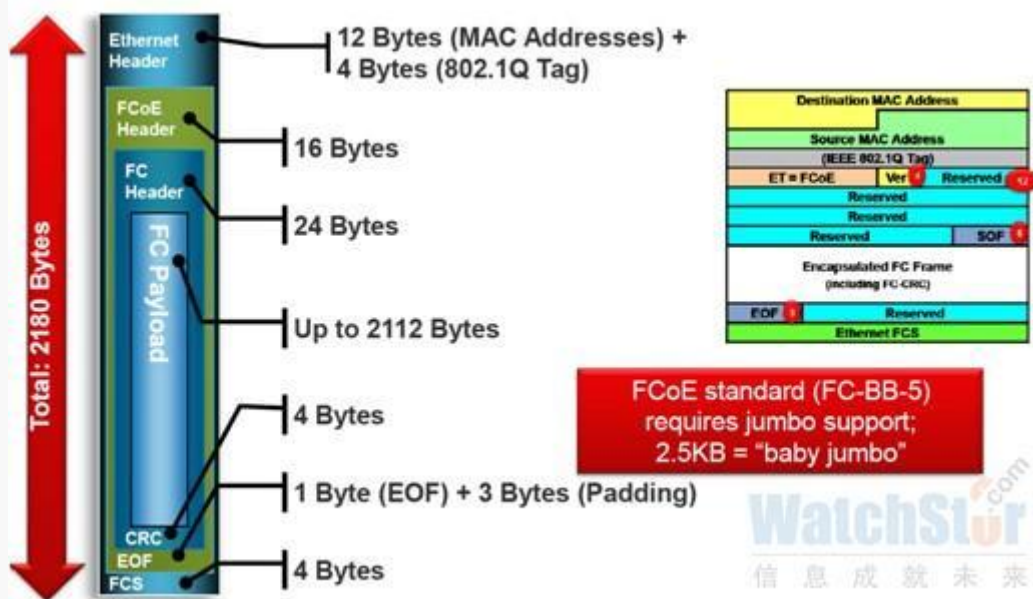
SCSI 不熟，这里不多说。FCP（Fibre Channel Protocol）是用来协助 SCSI 进行寻址的协议。iSCSI、FCIP 和 iFCP 都是依靠 TCP 的可靠连接确保无丢包，但封装的报头多了开销很大。iSCSI 由于需要全新的存储设备支持，过于激进，目前虽然有发展，但是受传统存储设备厂商制约始终很缓慢。FCIP 和 iFCP 都是支持 FC 网络跨 IP 核心网传输时用到的网络协议，由于目前 SAN 还是本地组网或使用光纤直连方式的远程组网较多，此场景并不多见，因此也应用很少，其中 FCIP 已经成为 RFC，而 iFCP 止步于 Draft。FCoE 相比较来说对上层协议改动较少，开销较低，且有利于减少服务器网络接口数量，在传统交换机厂商的大力鼓吹下当前发展最为迅猛，数据中心网络毕竟会是交换机的天下。

## .5.2 FCoE

FCoE 是在 2007 年 INCITS(国际信息技术标准委员会)的 T11 委员会（和 FC 标准制定是同一组织）开始制定的标准，2009 年 6 月标准完成（FC-BB-5）。FCoE 基于 FC 模型而来，仍然使用 FSPF 和 WWN/FC ID 等 FC 的寻址与封装技术，只是在外层新增加了 FCoE 报头和 Ethernet 报头封装和相应的寻址动作，可以理解为类似 IP 和 Ethernet 的关系。



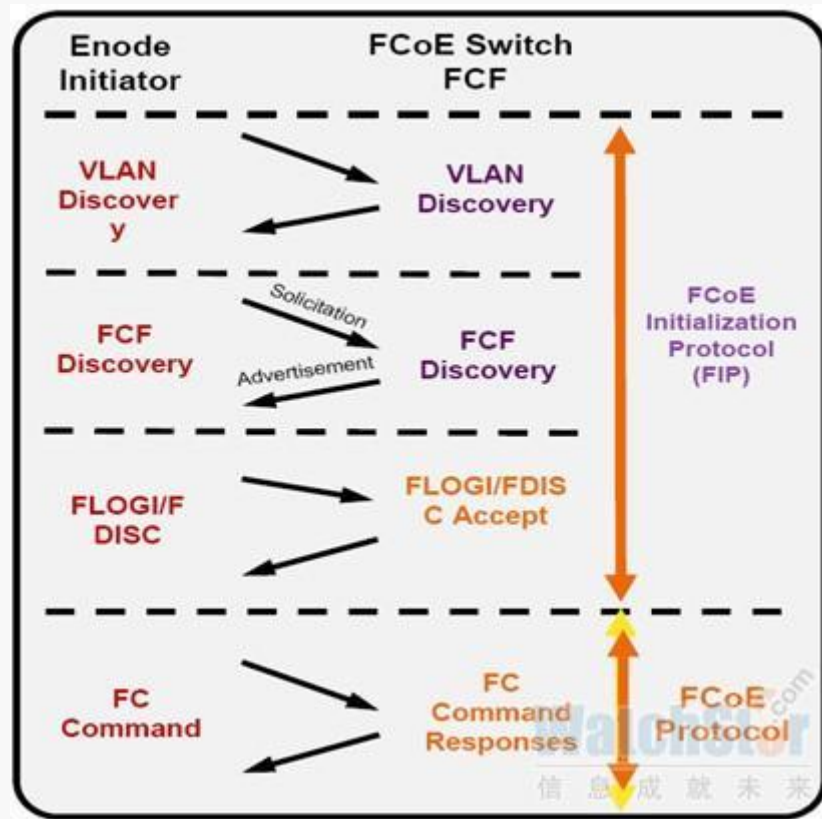
FCoE 标准定义了数据平面封装与控制平面寻址两个部分。封装很好理解，大家看看下面这张图就了然了。



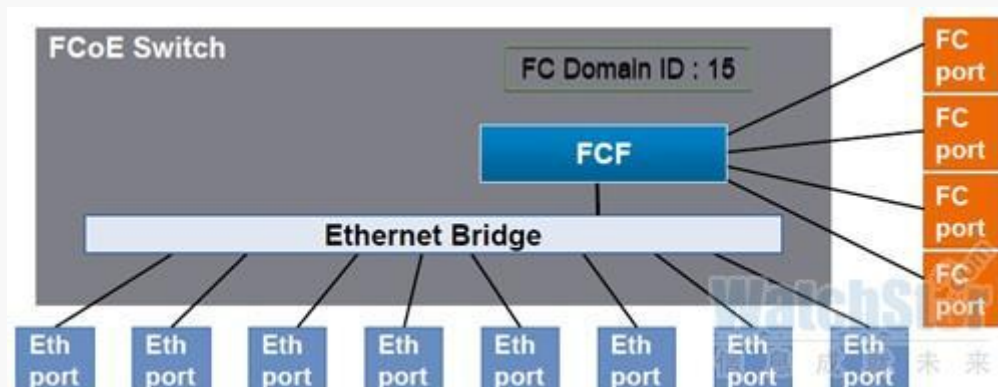
寻址稍微说一下，FCoE 使用 FIP（FCoE Initialization Protocol）进行初始化连接，FIP 运行于 VFPort 和 VNPort 之间或 VEPort 之间，所谓的 V 就是前面介绍 FC 的接口角色中的名称前面加了个 Virtual。FIP 在接口使能后一共做了三件事：

- 1、使用本地 VLAN（如 VLAN1）确认 FCoE 数据报文将要使用的 VLAN ID。
- 2、和 FCF 建立连接。

3、FLOGI/FDISC（Discover Fabric Service Parameters，FC 节点设备第一次向 FC 交换机注册请求 FC ID 时使用 FLOGI，后面再续约或请求其他 FC ID 时都使用 FDISC）



FCF（Fibre Channel Forwarder）是 FCoE 里面重要的角色，可以是软件或者芯片硬件实现，需要占用 Domain ID，处理 FCoE 交换机中所有与 FC 相关的工作，如封装解封装和 FLOGI 等。



Enode 是指网络中所有以 FCoE 形式转发报文的节点设备，可以是服务器 CAN 网卡、FCoE 交换机和支持 FCoE 的存储设备。FCoE 外层封装的 Ethernet 报头中 MAC 地址在 Enode 间是逐跳的，而 FC ID 才是端到端的。（不好理解就琢磨下 IP/Ethernet 转发模型，将 Enode 想成路由器和主机，一样一样滴）

与三层交换机中的 VLAN 接口一样，每个 FCF 都会有自己的 MAC，由于 FC ID 是 FCF 分配给 Enode 的，继承下来的终端 Enode MAC 也是由 FCF 分配的并具有唯一性，这个地址叫做 FPMA（Fabric Provided MAC Address）。FPMA 由两部分组成，FC-MAP 与 FC ID，结构如下所示，这样当 FCoE 交换机收到此报文后可以根据 FC-MAP 判断出是 FC 报文，直接送给 FCF，FCF 再根据 FC ID 查表转发，处理起来更简单。



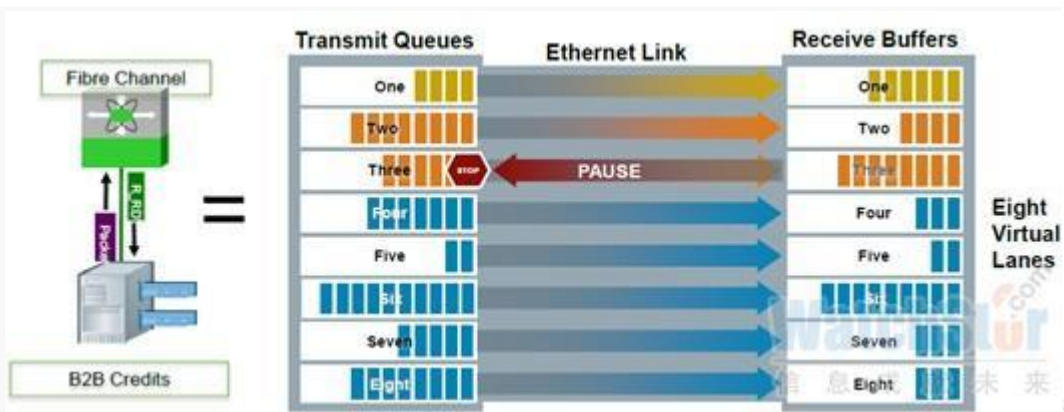
由上面 FC-MAP 的定义也可以看出，每个 FCF 下联的 Enode 终端最多也就 255 个(00-FF)。

## DCB

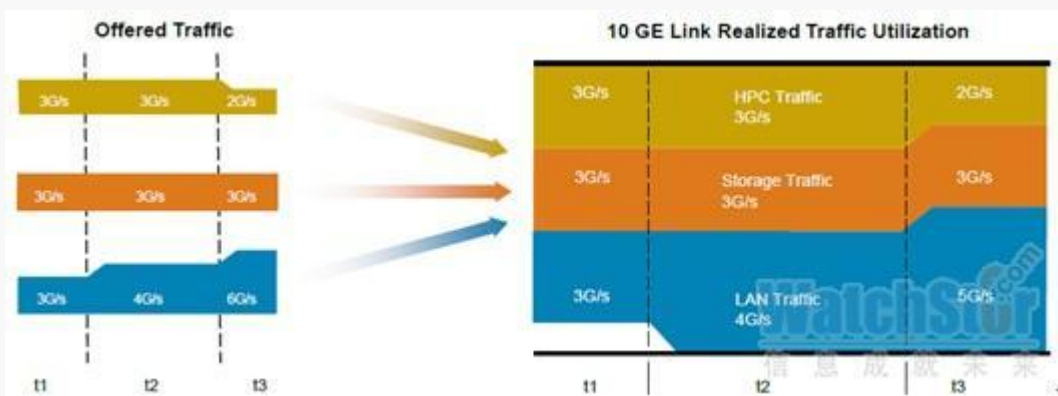
由于 Ethernet 是冲突丢包的，为了保证 FCoE 的无丢包，IEEE 引入了一系列的无丢包以太网技术（Lossless Ethernet），都定义在 802.1Q DCB（Data Centre Bridging）标准系列中。DCB 等同于 DCE（Data Centre Ethernet）和 CEE（Converged Enhanced Ethernet）的含义，就是不同厂商和工作组的不同称谓，内容都是一致的。DCB 是 IEEE 为了在数据中心对传统以太网技术进行扩展而制定的系列标准，前面说过的 VM 接入技术标准中 802.1Qbg 和 802.1Qbh 都是 DCB 中的一部分，另外还有 802.1Qau CN（Congestion Notification），802.1Qaz ETS（Enhanced Transmission Selection）和 802.1Qbb PFC（Priority-based flow control）。其中 802.1Qau CN 定义了拥塞通知过程，只能缓解拥塞情况下的丢包，加上其必须要全局统一部署与 FCoE 逐跳转发的结构不符，因此不被算成无丢包以太网技术的必要组成部分。常见的无丢包技术主要是 PFC 和 ETS，另外还有个 DCBX（Data Center Bridging Exchange Protocol）技术，DCBX 也是一起定义在 802.1Qaz ETS 标准中。

PFC 对 802.3 中规定的以太网 Pause 机制进行了增强，提供一种基于队列的无丢包技术，实际达到的效果和 FC 的 BB Credits 一样。简单理解如下图所示。





ETS 是带宽管理技术,可以在多种以太网流量共存情况下进行共享带宽的处理,对 FCoE 的流量报文进行带宽保障。简单理解如下图所示。

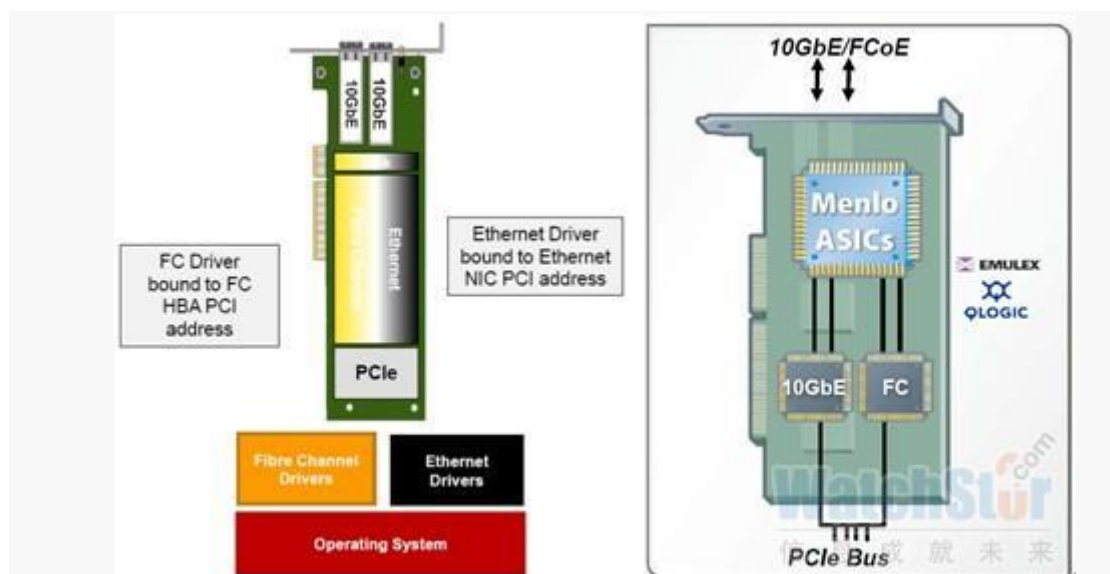


DCBX 定义了通过 LLDP 在两个相邻 Enode 之间进行 PFC, ETS 等参数自协商交互的过程。

DCB 的几个标准目前都还处于 Draft 阶段,其中 PFC 是由 Cisco 的 Claudio DeSanti 主编,ETS 由 Qlogic 的 Craig Carlson 主编。

## CNA

再补充一句服务器上的 FCoE 网卡 CNA (Converged Network Adapter),这个东西就是万兆 Ethernet 和 FC HBA (Host Bus Adapter) 网卡的合体,里面包含两个独立芯片处理 Ethernet 和 FC 各自的流量,在操作系统上看到的就是两个独立的 Ethernet 和 FC 网络接口,其上再增加第三个芯片进行流量混合封包处理。可参考下面两张图示。

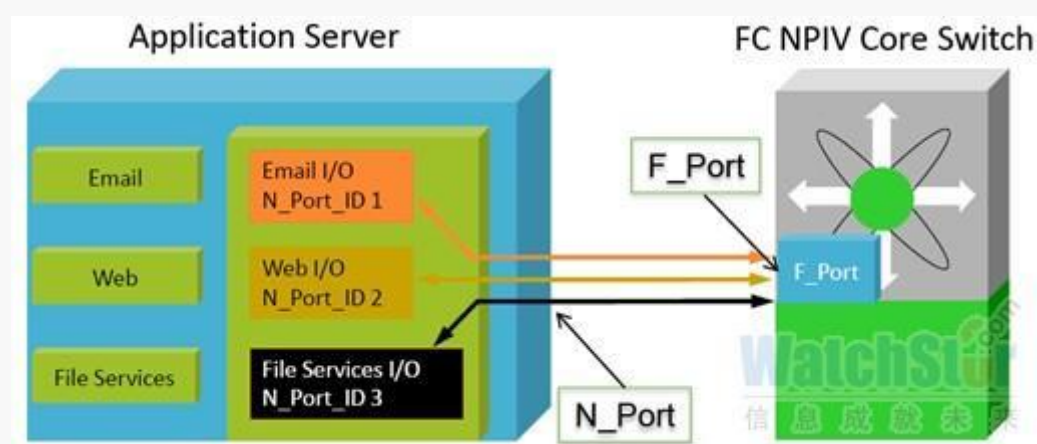


FCoE 的技术要点就这么多了，需要记住的关键是，FCoE 标准提供的是服务器到存储设备端到端的网络连接模型，FCF 是 FCoE 交换机的关键特性。目前已经有支持 FCoE 的存储设备出现，估计服务器到存储的全 FCoE 商用项目组网也很快会在市场上出现。

### 5.5.3 NPV

目前市面上 80% 以上的标榜自己实现了 FCoE 的交换机产品其实都是只实现了 NPV 功能，和前面描述的 FCoE 标准内容沾不上多少边儿。那么 NPV 是啥呢？

先说 NPIV（NPort ID Virtualization），这个还是 FC 里面的概念。前面说了 Server 的 NPort 需要向 FC Switch 进行 FLOGI 注册获取 FC ID 进行路由，那么如果一台物理服务器里面搞了好多虚拟机后，每个 VM 都打算弄个 FC ID 独立通信，但只有一块 FC HBA 网卡咋办呢。FC 中通过 NPIV 解决了这种使用场景需求，可以给一个 NPort 分配多个 FC ID，配合多个 pWWN（private WWN）来进行区分安全控制。



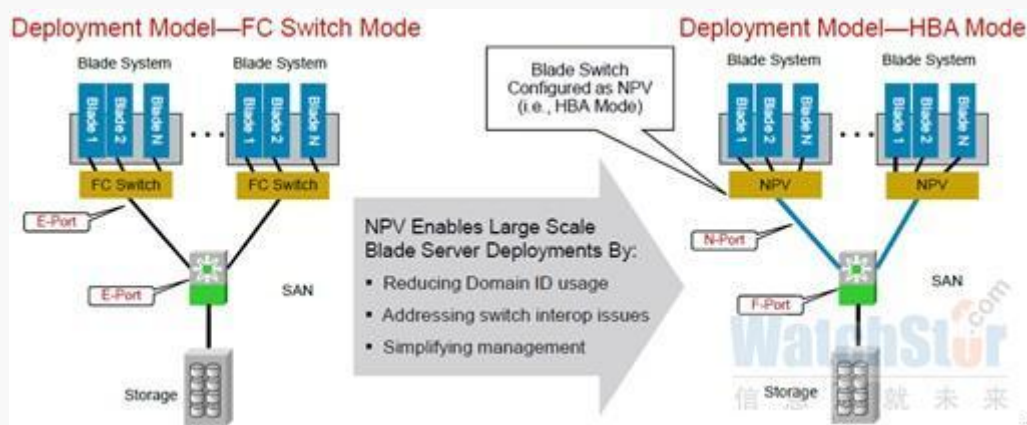
理解了 NPIV 就好说 NPV 了，我们把上图中的 NPort 拿出来作为个独立设备给后面服务器代理进行 FC ID 注册就是 NPV（NPort Virtualization）了。NPV 要做的两件事：



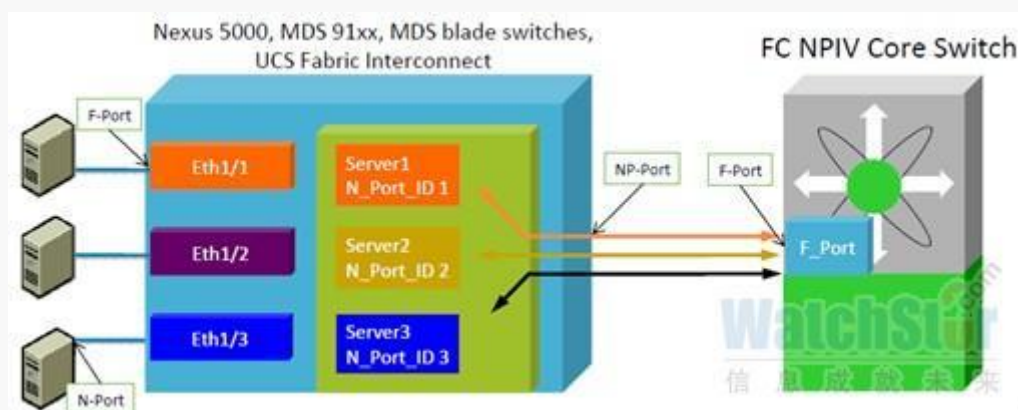
1、自己先通过 FLOGI 向 FC Switch 注册去要个 FC ID

2、将后续 Server 过来的 FLOGI 请求代理成 FDISC 请求，向 FC Switch 再去申请更多的 FC ID

NPV 的好处是可以不需要 Domain ID（每个 FC 区域最多只有 255 个），同时能将 FC 交换机下联服务器规模扩大。NPV 在 FC 网络中最常见的应用是在刀片交换机上。



随之有人将 FCoE 的脑筋动到了 NPV 与服务器之间的网络上，如下图所示：



在 FCoE 中的 NPV 相比较 FC 中要多做三件事，参考前面 FIP 流程：

1、回应节点设备关于 FCoE 承载 VLAN 的请求

2、回应节点设备的 FCF 查找请求，根据自己初始化时从 FC Switch 得到的 FC ID 生成仿冒 FCF 使用的 MAC 地址

3、在 CNA 网卡和 FC Switch 之间对转发的数据报文进行 FCoE 头的封包解包。

NPV 不是 FCoE 标准中定义的元素，因此各个厂家在一些细节上实现起来都各玩各的。比如都是将连接服务器的 Ethernet 接口和连接 FC Switch 的 FC 接口绑定起来使用，但是对应的绑定规则就可能不同。再有如 FC 接口故障时，如何将服务器对应的通道切换到其他 FC 接口去，是否通知服务器变化重新进行 FLOGI 注册，及通知等待时长等设定都会有所区别。

说说 NPV 的好处，首先是实现容易，之前描述的那几件主要的任务现在都已经有了公共芯片可以直接搞定，所以包装盒子就是了。其次是部署简单，不需要实现 FCF，不用管 FC 转发，不计算 FSPF，不占 Domain ID。最后是扩展方便，使用 FC Switch 的少量接口就可以连接大量的服务器。

由于 NPV 与服务器之间网络为传统以太网，因此 NPV 交换机也必须支持 DCB 标准中相关的无丢包以太网技术。

严格来讲，NPV 交换机不是 FCoE 标准中定义的 FCoE 交换机，但可以在接入层交换机上实现与服务器之间的 Ethernet 网络复用，减少了服务器的物理网卡数量（并未减少操作系统层面的网络通道数量），扩展了 FC 网络接入服务器节点的规模，适用于云计算大规模服务器部署应用。

补充一个 ENPV (Ethernet NPV) 的概念，这个东东是 Cisco 提的，就是在服务器与 FCoE 交换机 (FCF) 之间串个 NPV 进去，还是做些代理的工作，可以对 FIP 进行 Snooping，监控 FIP 注册过程，获取 VLAN/FC ID/WWN 等信息，对过路流量做些安全控制啥的。这种东东存在既合理，但以后有没有搞头就不好说了，市场是检验技术的唯一标准。

#### 5.5.4 小结

FCoE 是端到端的，FCF 是不可少的，NPV 是干代理的，目前是最合适的。

个人觉得 NPV 这个东西真的很彪悍，设计一套 FCoE 标准虽然是技术含量很高的活儿，但第一个搞出来 NPV 的才是真正的人精。如果没有 NPV，FCoE 想从 FC 口里夺食难度至少会增加上百倍，没准儿就跟 iSCSI 一样落得个鸡肋的地步。强烈建议 FCoE 标准尽快将 NPV 搞进来，要不单独出个 FC-BB-7/8 啥的独立标准体系也不错。

随着互联网的发展，对网络最大的需求就是带宽增长，云计算更是如此，因此如果 FC 的带宽演进继续这么不紧不慢的话，势必会被 100G Ethernet 取代，至于时间点就要看带宽需求增速了，个人估计不会超过 10 年，到时就有 FCoE 的用武之地了，至于之前这段时间应该都还是 NPV 在前面冲锋陷阵。

继续大胆 YY，可否搞个什么技术冲击下 FCoE 呢？由于 IP 是无连接的，Ethernet 是冲突丢包的，因此想保证数据传输的可靠性就只能像 iSCSI 一样在 TCP 上做文章，但是层次做高了，报头一多开销又太大，矛盾啊矛盾。如果完全替代 Ethernet，那就类似于重建一套 FC 协议了，但是目前也看不到带宽速率发展能超过 Ethernet 的替代技术。那么中间手段就是搞下 IP 这个层面，像 FCoE 的思路可以理解为用 SCSI/FCP 替代 TCP/IP 在 Ethernet 上传输，由于 SCSI/FCP 这套协议和 Ethernet 已经很成熟了，只是搞个接口 (FCoE 报头) 在中间承上启下就够了。不过 FCoE 需要引入无丢包以太网的设计，Pause 帧会不会降低 Ethernet 的转发效率还不好说。作者思路是设计一套带连接状态的传输机制（类似 TCP 带重传确认，

而且可以像 IP/FC 一样能够寻址) 替代 IP/FCP 这个层面, 上面还是承载 SCSI, 下面跑着传统以太网。不见得靠谱, 仅供拓展一下思路, 有兴趣的同学可深思。

## 5.6 跨核心层服务器二层互访

本章节重点技术名词: L2MP/ VSS/IRF/ vPC/ TRILL/SPB/FabricPath/QFabric/VDC/VPN

在服务器跨核心层二层互访模型中, 核心层与接入层设备有两个问题是必须要解决的, 一是拓扑无环路, 二是多路径转发。但在传统 Ethernet 转发中只有使用 STP 才能确保无环, 但 STP 导致了多路径冗余中部分路径被阻塞浪费带宽, 给整网转发能力带来了瓶颈。因此云计算中需要新的技术在避免环路的基础上提升多路径带宽利用率, 这是推动下面这些新技术产生的根本原因。

前面网络虚拟化章节部分提到了两个解决上述需求的思路。

首先是控制平面多虚一, 将核心层虚拟为一个逻辑设备, 通过链路聚合使此逻辑设备与每个接入层物理或逻辑节点设备均只有一条逻辑链路连接, 将整个网络逻辑拓扑形成无环的树状连接结构, 从而满足无环与多路径转发的需求。这种思路的代表技术就是 VSS/IRF/vPC, 前两者都是控制平面全功能同步的整机虚拟化技术, vPC 则是精简后只处理控制平面与跨设备链路聚合使用相关功能的多虚一技术。此类技术必定都是私有技术, 谁家的控制平面都不可能拿出来完全开放。

另一个思路是数据平面多虚一, 在接入层与核心层交换机引入外层封装标识和动态寻址协议来解决 L2MP (Layer2 MultiPath) 需求, 可以理解这个思路相当于在 Ethernet 外面搞出一套类似 IP+OSPF 的协议机制来。对接入层以下设备来说, 整个接入层与核心层交换机虚拟成了一台逻辑的框式交换机, Ethernet 报文进 Ethernet 报文出, 中间系统就是个黑盒, 就好像 IP 层面用不着了解到 Ethernet 是怎么转发处理的一样。这种思路的代表技术是 IETF (Internet Engineering Task Force) 标准组织提出的 TRILL 和 IEEE 提出的 802.1aq SPB 两套标准, 以及一些厂商的私有技术。如 FabricPath 是 Cisco 对 TRILL 做了一些变更扩展后的私有技术称谓 (以前也有叫 E-TRILL), QFabric 则是 Juniper 的私有技术, 推测是基于 MACinMAC 封装和自己搞的私有寻址协议来做的。

这里唠叨两句私有协议, 从有网络那天开始, 私有协议就始终相生相随。从早期的 EIGRP 和 HSRP 到现在的 VSS/IRF/vPC/OTV/QFabric 都是各厂家的私货。即使是标准还有 IETF/IEEE 等不同的标准化组织呢, 哪会有厂家就大公无私到研究出个啥东西都恨不得全人类共享, 逐利才是企业发展的根本。私有协议还有两个有趣的现象, 首先单从技术角度来说, 私有协议基本都代表了同时代同类技术的最先进生产力, 如果是个落后的技术, 只有脑袋进水了才会砸钱进去研发。还有就是只有在市场上占了一定地位的重量级厂商才会经常推自己的私有协议, 而且推得越多也代表着其地位越高。

从市场上看，最早由于没得选，大家基本上都能接受使用私有协议，后来出于不想将所有鸡蛋放在一个篮子里的心理，开始对私有协议有了抵触情绪，尤其是运营商级别用户明确要求不能使用私有协议。但就目前随着云计算的爆炸式增长，数据中心网络技术面临着一次新的飞跃，传统技术已经无法满足需求，因此私有协议再次进入了人们的视线。预计随后几年中，新的云计算数据中心会以站点 Site 为单位来部署单一厂商设备，如可以看到全 Cisco 设备运行 FabricPath/vPC 的 Site，全 Juniper 设备 QFabric 的 Site，或全 H3C 设备运行 IRF 的 Site 等等，站点之间或对外再采用一些公有协议如 MSTP、RPR、BGP 等进行连接。

下面会分别介绍几项主要技术，顺序为控制平面多虚一技术 VSS/IRF/vPC，数据平面多虚一技术 TRILL/SPB/Fabric Path/QFabric 和控制平面一虚多技术 VDC。

### 5.6.1 控制平面多虚一技术

#### VSS/IRF

先说下 VSS/IRF，这两个技术基本上没啥差别。VSS（Virtual Switching System）是只在 Cisco 6500 系列交换机上实现的私有技术，IRF（Intelligent Resilient Framework）是在 H3C 所有数据中心交换机中实现的私有技术。二者的关键技术点如下：

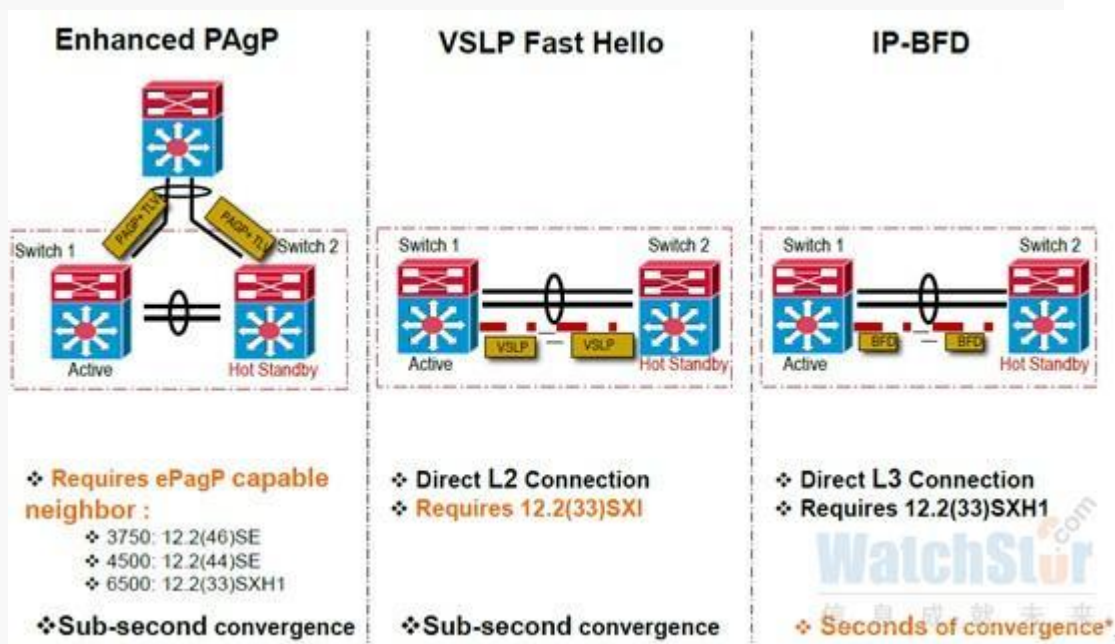
1、专用链路跑私有协议：VSS 使用 VSL（Virtual Switch Link），IRF 使用 IRF link 来承载各自的控制平面私有交互协议 VSLP 和 IRF。专用链路使用私有协议来初始化建立邻接、协商主备（描绘拓扑）、同步协议状态，同时会在虚拟化完成后，传输跨机框转发的数据流量。二者都推荐使用 10GE 链路捆绑来做专用链路，说明私有协议交互和跨框传输的流量会很大。

2、基于引擎的主备模式：二者的控制平面都是只有一块主控引擎做为虚拟交换机的主控制引擎，其他的引擎都是备份。所有的协议学习，表项同步等工作都是由这一块引擎独立完成。好在这些设备大都是分布式交换，数据转发的工作由交换板自己完成了，只要不是类似 OSPF 邻居太多，拓扑太大等应用情况，一块主控大部分也都能搞定了。注意 Cisco 6500 必须使用 Supervisor720 主控配合带转发芯片的接口板才能支持 VSS。

3、跨设备链路聚合：前面说了网络虚拟化主要是应对二层多路径环境下防止环路，因此跨设备链路聚合就是必须的了。Cisco 配合 VSS 的专用技术名词是 MEC（Multichassis EtherChannel），IRF 倒是没有啥专门的名词，其链路聚合也和单设备上配置没有区别。

4、双活检测处理：当 VSL 或 IRF link 故障后，组成虚拟化的两个物理设备由于配置完全相同会在网络中出现双活节点，对上下游设备造成 IP 网关混乱。因此 VSS/IRF 都设计了一些双活处理机制以应对专用链路故障。

1) 首先网络中如果有跨设备链路聚合时, VSS 使用 PAgP、IRF 使用 LACP 扩展报文来互相检测通知; 2) 如果有富裕接口在虚拟化的两台物理设备间可以单独再拉根直连线路专门用做监控, VSS 使用 VSLP Fast Hello、IRF 使用 BFD 机制进行检测通知; 3) 另外 VSS 还可以使用 IP BFD 通过互联的三层链路进行监控, IRF 则支持使用免费 ARP 通过二层链路进行监控。上述几种方式都是监控报文传输的链路或者外层承载协议不同。当发现专用链路故障时, VSS/IRF 操作结果目前都是会将处于备份状态的物理机框设备的所有接口全部关闭, 直到专用链路恢复时再重新协商。需要注意这两种虚拟化技术在进行初始协商时都需要将角色为备份的机框设备进行重启才能完成虚拟化部署。下面以 Cisco VSS 的三种故障检测方式举例, IRF 也差不多。



除了上述 4 个关键技术点外, VSS/IRF 还有一些小的相似技术设定, 如 Domain 的设定、版本一致性检查、三层虚接口 MAC 协商等等, 都是基于方方面面的细节需求来的。由于应用环境相似, 因此实现的东西也区别不大。想想 RFC 中的 OSPF 和 BGP 到现在都还在不断的推出新的补充标准和 Draft 来查漏补缺, 就知道细节的重要性了。

VSS 特色一些的地方是结合了 Cisco 6500 的 NSF (None Stop Forwarding) 和 SSO 进行了主控板故障冗余和版本升级方面的可靠性增强。而 IRF 则是将虚拟化延伸到了 H3C 接入层设备 5800 系列盒式交换机上 (最大支持 8 或 9 台物理设备虚拟化为一台逻辑设备), 可以打造逐层虚拟化的数据中心网络。

VSS 和 IRF 都是当前较为成熟的虚拟化技术, 其优点是可以简化组网, 便捷管理, 缺点则是扩展性有限, 大量的协议状态同步工作消耗系统资源, 而且纯主备的工作方式也导致了主控引擎的资源浪费。







用推广。估计得等到能够支撑上百块接口板的 Super Supervisor 出现，或者更完善的算法以提供多主控负载均担，才能打破当下控制平面的多虚一规模限制了。

从 Cisco Nexus 系列产品的技术发展来看，在网络虚拟化的路线上 Cisco 已经开始偏向于数据平面虚拟化的 TRILL 等新兴网络技术，VSS/vPC 等技术受主控引擎性能影响的部署规模局限性和协议私有化特征是制约其发展的硬伤，终将逐渐淡出大家的视线。

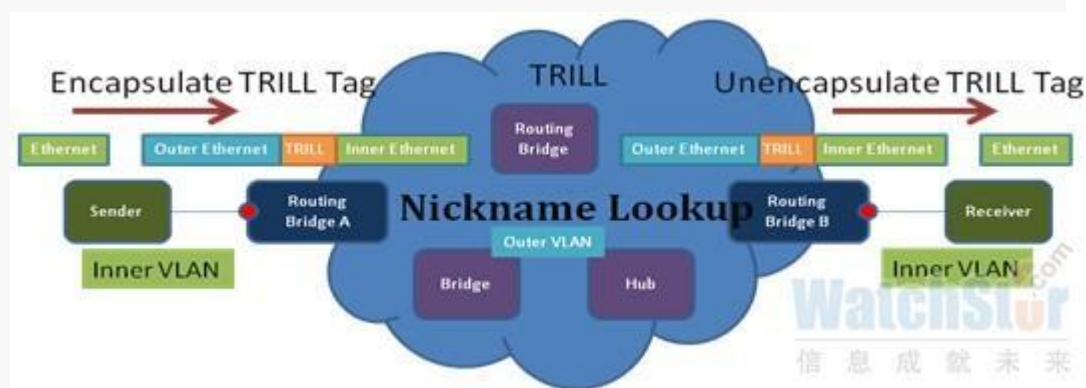
## 5.6.2 数据平面多虚一技术

数据平面多虚一技术的统一特征就是在二层 Ethernet 报文外面再封装一层标识用于寻址转发，这样基于外层标识就可以做些多路径负载均衡和环路避免等处理工作了。TRILL/SPB 都是属于此列，QFabric 目前开放出来的资料较少，猜测其应该也使用了类似 MACinMAC 之类的方式在其网络内部传输报文。

### TRILL

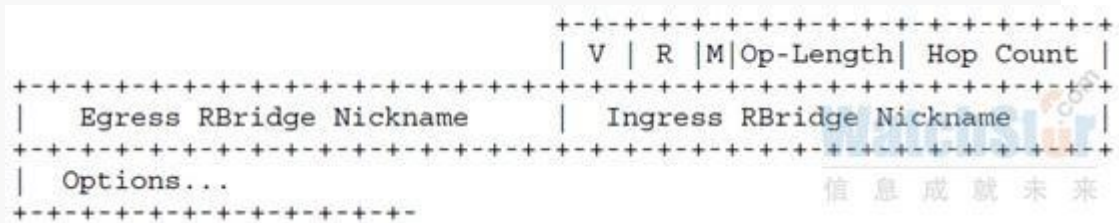
先说 TRILL (TRAnsparent Interconnect of Lots of Links) 和 FabricPath。2010 年 3 月时 TRILL 已经提交了 IETF RFC 5556 规范(Transparent Interconnection of Lots of Links (TRILL): Problem and Applicability Statement)，此 RFC 只是描述了 TRILL 要解决的问题以及应用范围，定义协议细节的文档目前都还处于 Draft 阶段，形成完整的协议标准体系应该还得 1-2 年。

TRILL 并不是专门为数据中心开发的技术，其定义的是在大型 Ethernet 网络中解决多路径问题的方案。FabricPath 是 Cisco 在 TRILL 标准之上加入了很多私货的专门为数据中心而设计的一个超集，基本的控制平面与数据平面二者没有明显区别。



控制平面上 TRILL 引入了 L2 ISIS 做为寻址协议，运行在所有的 TRILL RB (Routing Bridge) 之间，部署于一个可自定义的独立协议 VLAN 内，做的还是建立邻接、绘制拓扑和传递 Tag 那几件事。数据平面在内外层 Ethernet 报头之间引入了 TRILL 报头，使用 NickName 作为转发标识，用于报文在 TRILL 网络中的寻址转发（可理解为类似 IP 地址在 IP 网络里面转发的作用）。每个 RB 都具有唯一的 Nickname，同时维护其他 RB 的 TRILL 公共区

域 MAC 地址、Nickname 和私有区域内部 MAC 地址的对应关系表。因为 TRILL 封装是 MACinMAC 方式，因此在 TRILL 公共区域数据报文可以经过传统 Bridge 和 Hub 依靠外部 Ethernet 报头转发。TRILL 报头格式如下图所示：



V (Version): 2 bit，当前 Draft 定义为 0。

R (Reserved): 2 bits，预留。

M (Multi-destination): 1 bit, 0 为已知单播, 1 为未知单播/组播/广播, 此时 Egress RBridge Nickname 意味着当前转发使用多播树的根。

Op-Length (Options Length): 5 bit，Option 字段长度。

Hop Count: 6 bit，最大跳数，逐跳减一，为 0 丢弃，防止环路风暴。

Egress RBridge Nickname: 16 bit，已知单播标示目的私网 MAC 对应的 RB，多播则标示多播树根 RB。中间传输 RB 节点不能改变此字段值。

Ingress RBridge Nickname: 16 bit，标示报文进入 TRILL 区域的初始边缘 RB，中间传输 RB 节点不能改变此字段值。

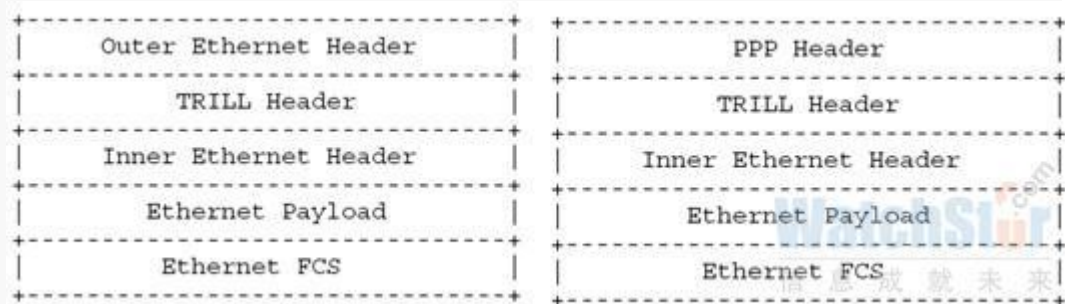
Options: 目前只定义了 CHbH (Critical Hop by Hop) 和 CltE (Critical Ingress to Egress) 两个 1bit 的标志位，用于说明后面的 Option 预留内容是需要逐跳设备识别处理的或是首末端设备必须识别处理的。至于真正的 Option 目前都还没有定义。下图为 Option 字段内容：



普通 Ethernet 报文在首次从 TRILL 边缘 RB 设备进入 TRILL 区域时，作为未知单播还是依照传统以太网传播方式，广播给所有其他的 RB 节点。但是除了边缘 RB 外，TRILL 区域中间的 RB 和传统 Bridge 都不会学习此数据报文中私有区域内部 MAC 地址信息，有效的降低了中间设备的 MAC 地址表压力。为了防止环路同时做到多路径负载均衡，TRILL 的每个 RB 在初始建立邻接绘制拓扑时，都会构造出多个多播树，分别以不同的 Nickname 为根，将不同的未知单播/组播/广播流量 Hash 到不同的树，分发给其他所有 RB。

由于全网拓扑唯一且构造树时采用的算法一致，可保证全网 RB 的组播/广播树一致。在 RB 发送报文时，通过将报文 TRILL 头中的 M 标志位置 1 来标识此报文为多播，并填充树根 Nickname 到目的 Nickname 字段，来确保沿途所有 RB 采用同一颗树进行广播。组播与广播报文的转发方式与未知单播相同。已知单播报文再发送的时候，会根据目的 RB 的 Nickname 进行寻路，如果 RB 间存在多条路径时，会逐流进行 Hash 发送，以确保多路径负载均衡分担。

另外 TRILL 除了支持外层 Ethernet 封装在传统以太网中传输外，还规定了一种外层 PPP 封装方式可以跨广域网技术传输。以下是两种典型的 TRILL 报文封装方式：



TRILL 的主要技术结构就是上面这些了，对更细节内容感兴趣的同学可以自行去 IETF 翻翻相关 Draft。目前各个芯片厂商都已经进入 TRILL Ready 的阶段，只要技术标准完善发布并被广泛客户所接受，相关产品商用是 So 快的。

### FabricPath

FabricPath 是 Cisco 2010 年 6 月底正式发布的专门针对数据中心设计的私有技术，以前也叫做 L2MP/E-TRILL，Cisco 资料上的原话是：

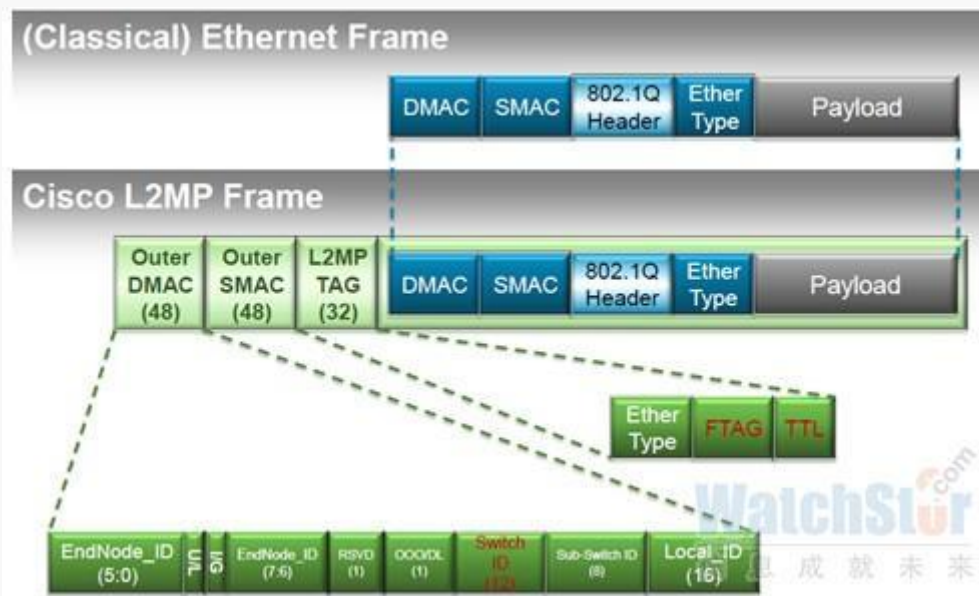
The Cisco engineers who developed L2MP only pushed part of it to IETF TRILL.

Functionality-wise, L2MP is a superset of TRILL

L2MP = TRILL + Cisco extensions

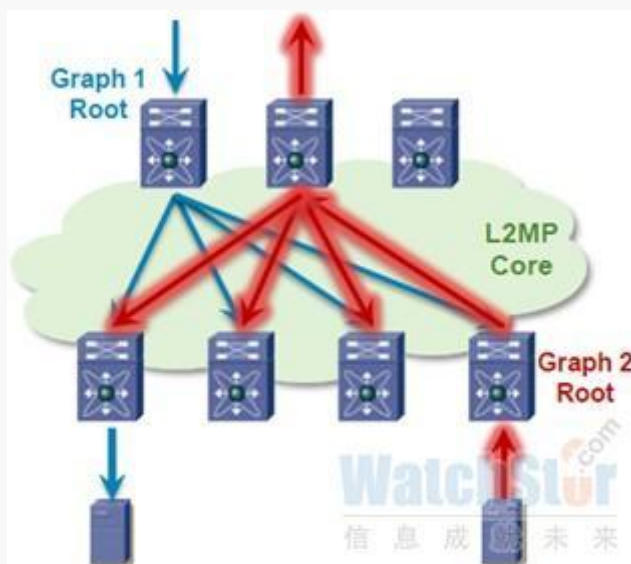
控制平面和转发规则上 FabricPath 和 TRILL 没有大的区别，一些主要变化如下，也可以说这些是 Cisco 专门针对数据中心网络对 TRILL 做出的变更：

1、FabricPath 只支持 RB 间的点到点直连，不能加入传统 Bridge 等设备。因此 FabricPath 的报文格式相比较 TRILL 就更加简化，不再需要依靠外层的目的 MAC 进行 Ethernet 转发，数据报文封装有较大不同，FabricPath 的报文格式如下图所示。



其中的 Switch ID 就是 TRILL 里面的 Nickname。TTL 字段用于避免环路风暴

2、采用 FTAG（Forwarding TAG）标示不同的多播树 Graph，用于多拓扑中未知单播/组播/广播报文的转发。多拓扑指可以在同一套 FabricPath 网络中支持不同的拓扑转发，而目前 TRILL 的多拓扑还未明确定义。每套拓扑缺省使用 2 个 Graph，每个 Graph 可以用一个 FTAG 标示。目前 NOS 发布版本只开放支持 2 棵树 Graph，既一套拓扑，据称最多可扩展至 64 棵树。多拓扑结构如下图所示。



3、MAC 基于会话进行学习。当 FabricPath 的边缘设备从 FabricPath 区域中收到报文进行 MAC 地址学习时，会进行目的 MAC 地址检查，只有目的 MAC 在本地私有区域存在的，



才会学习报文的源 MAC 到地址表中，这样可以避免 MAC 的不必要扩散。TRILL 中 RB 设备还是传统的 Ethernet 方式，收到报文就会学习源 MAC，不做判断。

#### 4、FabricPath 支持基于 vPC、FHRP 等 Cisco 私有协议的组合应用扩展。

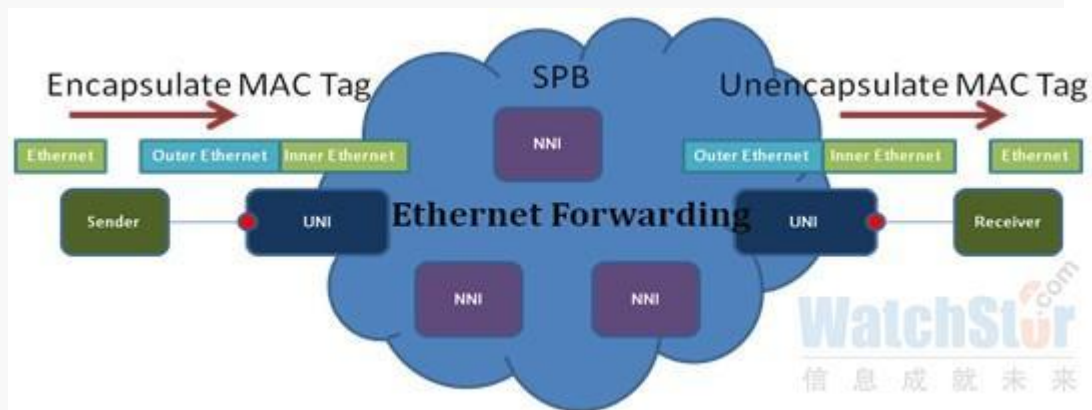
在当前 TRILL 还未有完整明确的标准出台情况下，Cisco 已经用 FabricPath 走在了所有人前面，可以支持云计算大规模节点二层通信的数据中心建设，当然其主要的被攻击点就是私有协议。另外在 Cisco 的发布资料中也指出，其产品均已进入 TRILL Ready 状态，以后只需要命令变更就可以切换设备分别运行于纯粹的 TRILL 和扩展的 FabricPath 模式下。

### SPB

要说 SPB 需要先谈谈 PBB。PBB (Provider Backbone Bridging) 是 IEEE 于 2008 年完成的 802.1ah 标准，为运营商城域以太网定义了一整套 MACinMAC 的转发机制。但 PBB 只定义了转发平面的封装内容，当报文封装上外层 Ethernet 报头在运营商骨干区域二层网络中时，仍然需要依靠传统的 STP 进行环路避免和转发控制。于是 IEEE 在 2009 年又定义了 802.1Qay PBB-TE (Provider Backbone Bridge Traffic Engineering)，用于在运营商的骨干区域中进行拓扑管理与环路保护，说白了就是通过手工方式配置一堆指定路径取代 STP 的自动收敛。目前 IEEE 还有个相关的标准 P802.1Qbf, PBB-TE infrastructure protection 处于草案阶段，预计 2011 年发布。

PBB-TE 静态规划转发路径，明显无法适用于大型二层网络扩展，于是 IEEE 再搞出个 P802.1aq SPB (Shortest Path Bridging) 来，当前也还处于草案阶段。从 IEEE 的资料上看 SPB 主要是为了解决 STP 阻塞链路浪费带宽的问题而研究出来的。从实现上来看，同样是采用了 L2 ISIS 作为其控制平面协议进行拓扑学习计算，用 MACinMAC 封装方式在 SPB 区域内部进行报文传输。和 TRILL 很像吧，好在 IEEE 和 IETF 都是开放的标准化组织，不存在专利之争，不然肯定要掐架了。

SPB 可细分为 SPBV (VLAN QinQ) 和 SPBM (MACinMAC) 两个部分，目前看主要用到的是 SPBM。





SPBM 是标准的 MACinMAC 封装，在 SPB 区域中数据报文也都是依靠外层 MAC 做传统 Ethernet 转发。外层 Ethernet 报头中的源目的 MAC 就代表了 SPB 区域边缘的 UNI 设备，此设备 MAC 是由 L2 ISIS 在 SPB 区域中传递的。

由于在 SPB 网络中还是采用传统 Ethernet 进行转发，因此需要定义一系列的软件算法以保证多路径的广播无环和单播负载均衡。下面介绍几个主要的部分：

1、首先 SPB 定义了 I-SID 来区分多个拓扑，I-SID 信息在数据报文中以 BVID（外层 Ethernet 报头中的 VLAN Tag）形式携带，这样可以解决不同业务多拓扑转发的问题。

2、每个 SPB 节点都会为每个 I-SID 计算三棵树：到达所有相关 UNI 节点的 SPT(Shortest Path Tree) 用于单播与组播报文的转发；ECT（Equal Cost Tree）以处理两个 UNI 间存在多条等价路径时负载均衡转发；自己为根的多播树 MT（Multicast Tree）用于未知单播与广播报文转发。

3、任意两点间的 Shortest Path 一定是对称的；ECT 的负载均衡是基于不同 I-SID 分担的；

总的来说，SPB 和 TRILL/FabricPath 相比主要有以下不同：

	SPB	TRILL	FabricPath
多拓扑	支持	研究中	支持
外层封装	标准Ethernet	Ethernet+TRILL	Ethernet+L2MP
转发标识	目的MAC	Nickname	Switch ID
多播树	各自为根	全局统一	全局统一
多路径负载分担	ECT端到端分担	逐跳Hash	逐跳Hash
环路避免	RPFC（反向路径检测）	Hop Count/RPFC	TTL/RPFC
标准组织	IEEE	IETF	Cisco
转发芯片支持	现有芯片	新一代芯片	Cisco自有芯片

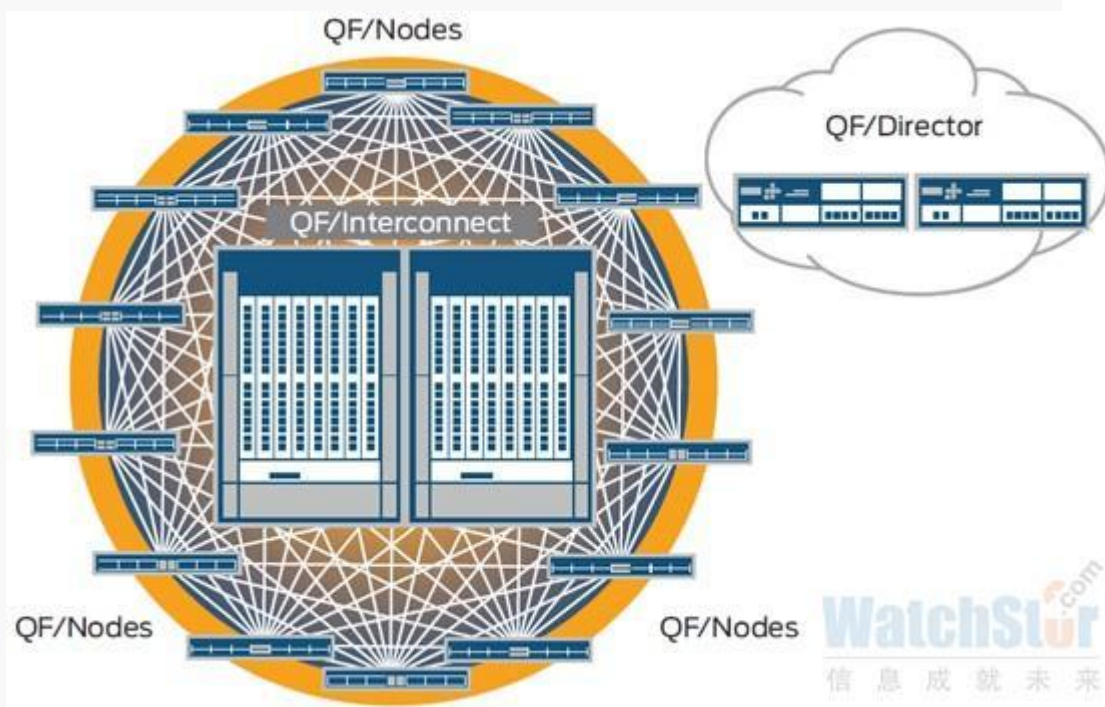
SPB 目前的最大困扰是转发路径靠软件算法保障，尤其在多路径负载分担时，对 CPU 计算压力远远超过 TRILL 和 FabricPath，因此实际转发效率令人存疑。而且 SPB 的出发点是运营商的城域以太网环境应用，是否能适用于数据中心网络还有待观察。当前 802.1aq SPB 已经进入到 Draft4.0，对其细节有兴趣的同学可以去 IEEE 网站注册下载学习。

多说一句，SPB 是纯软件的解决方案，不需要更新转发芯片去支持，因此只要其标准化后，任何厂家都可以很快推出支持的版本，包括 Cisco。

**QFabric**

Juniper 的 QFabric 也是目前喊得很大声的数据中心下一代网络技术,但由于还没有正式发布,开放的技术原理性文档基本没有,大都是些市场方面的资料。个人理解有以下几个要点:

- 1、首先控制平面一定是 Juniper 的私有协议,肯定要全 J 设备建设才成
- 2、由于 J 是自己没有芯片研发能力的,因此采购的基本只能是 Broadcom/Marvel 等几家通用芯片厂商的片子,再因此其转发肯定是基于 MACinMAC 的标准报文封装方式。
- 3、由 1 和 2 可以推断其实现方式应该是转发平面公共化,控制平面私有化。
- 4、从其如下的结构图中可以看出,对于虚拟后的逻辑交换机,可以理解 Director、Interconnect 和 Nodes 应该分别对应一台框式交换机的主控引擎、交换网板和接口板。



目前 Juniper 只发布了 Nodes 节点的 QFX3500 设备,等 Interconnect 和 Director 都出来估计怎么也得 2012 了。

### 小结

TRILL/FabricPath/SPB/QFabric 都引入了控制平面协议来处理拓扑管理和转发路径判定的工作,都肯定会导致转发效率上,相比较传统 Ethernet 的下降,同时引入拓扑变化影响流量路径变更收敛速度的问题。但是这些技术毕竟比以前的 STP 在带宽上多了一倍的扩充,组网规模上也得到扩展,更适用于云计算数据中心的网络需求,总体来讲得大于失,属于更先进的生产力。

从开放标准上讲，个人倾向于 IETF 的 TRILL。毕竟 SPB 出身不正，定位于运营商城域互联应用，而且就连 IEEE 都没有将其放入 DCB（DataCenter Bridging）的大技术体系中。

从私有技术来看，VSS/IRF 受到组网规模有限的硬伤限制，随着云计算网络规模的增大会逐渐退出大型数据中心的舞台，但用户服务器规模也不是说上就能上来的，至少还能有 2、3 年的赚头。而 FabricPath/QFabric 会较前者有更宽广的舞台和更长久的生存期，但由于其私有化的特征，当开放标准成熟铺开，部署规模也只会日渐萎缩。可以想想 OSPF 和 EIGRP 的昨天和今天。

下面说些个人极度主观的预测（每写到这里都有些神棍的感觉）：

1、未来 2-3 年投入使用的云计算数据中心将是 FabricPath/IRF/QFabric 这些私有技术乱战的天下。在旧有开放技术不能满足使用，而新的标准仍未完善的情况下，私有技术成了人们唯一的选择。（VSS 由于基于 Cisco6500 系列产品，受设备性能影响 1-2 年内会完全退出数据中心的历史舞台）

2、未来 5 年左右大型数据中心网络将会是 TRILL 一统天下，SPB 半死不活，而各种私有技术也会延续之前的一部分市场，但很难得到进一步普及和推广。

最后完全依据个人喜好将上述技术做个排位，仅供参考，请勿拍砖。

FabricPath > TRILL > QFabric > IRF/VSS > SPB

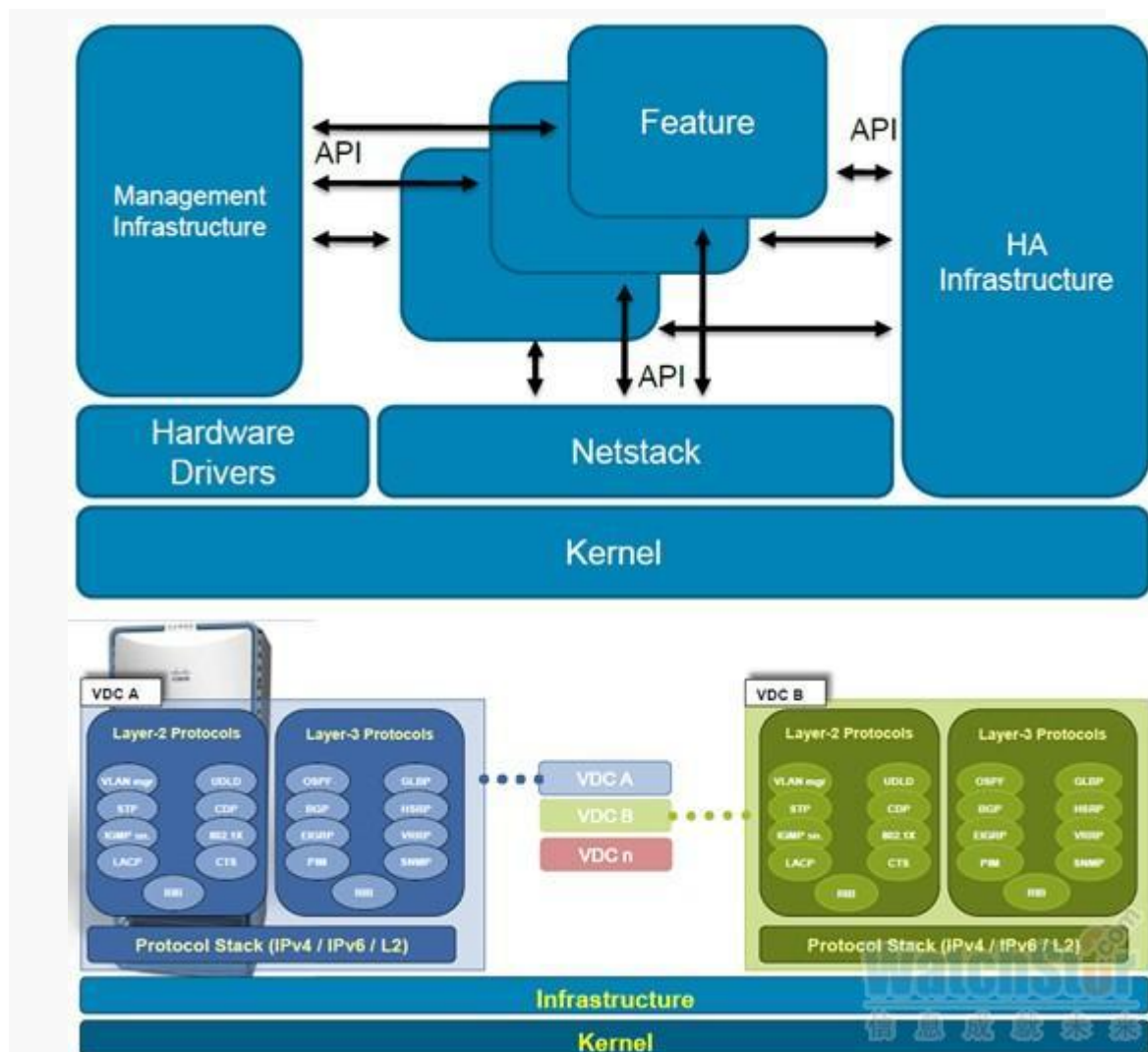
### 5.6.3 控制平面一虚多技术

这个目前就是 VDC 了，没有看到任何厂家有类似的技术出来，另外由于此技术完全是本地有效的使用范围，也不会存在什么标准化的互通问题，就算有人去提个标准大家也不会理睬的，可以参考 VMware ESX/XEN/HyperV 之间的关系。

VDC（Virtual Device Contexts）是 Cisco 基于操作系统级别的一虚多网络虚拟化技术。下面列表用于展示几项主要网络一虚多技术的区别。

技术	逻辑虚拟化对象
VLAN	数据平面
VRF	数据平面+少部分控制平面（路由协议）
虚拟防火墙	数据平面+管理平面
VDC	数据平面+控制平面+管理平面+系统资源

从下图的 NXOS 模型和 VDC 模型，可以看出 VDC 是建立在底层 OS 之上的，因此推测其采用的应该是前文中提到的 OS-Level 虚拟化技术。



现阶段 Cisco 公布的软件规格为每物理设备最多虚拟 4 个 VDC，并支持每 VDC 4k VLANs 和 200 VRFs 进行分层次虚拟化部署。按照云计算的需求来看，一般给每个用户分配的都是以带宽为度量的网络资源，不会以交换机为单位进行虚拟网络资源非配。即使是给的话，一台 N7000 只能虚拟 4 个 VDC 也不够用户分的。如果纯粹的流量隔离需求，最多使用到 VRF 也就够了，再多层次的虚拟化目前还看不清使用场景需求。

举个类似的例子，分层 QoS 一段时间以来也很火爆，但个人认为分 2 层也好，分 10 层也好，业务用不上都是白搭，当然在不考虑具体应用情况，纯粹拼指标时还是很有用的。又回到了前面的老话，正确的网络设计原则永远是自顶向下的。

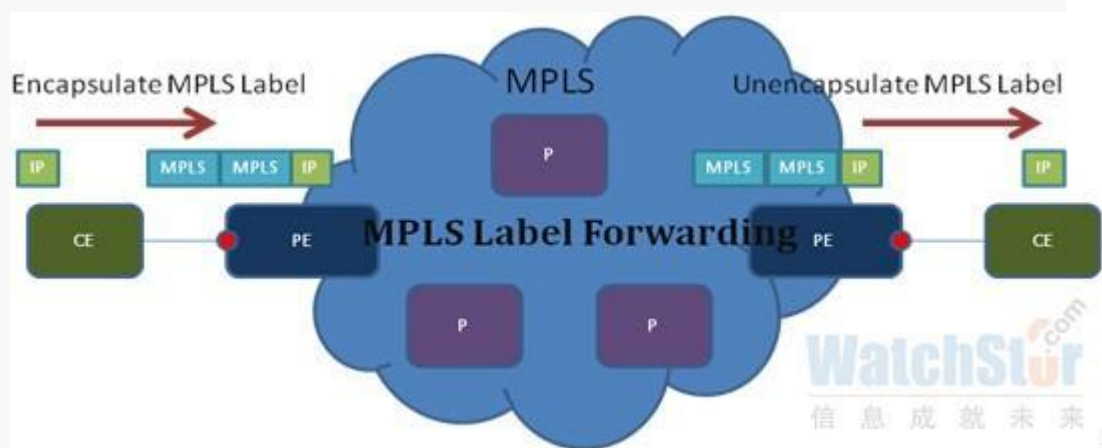
VDC 可说的东西就这么多了，更多的就涉及到 NXOS 的核心设计了，Cisco 也不太可能会向外公布。从 Cisco 的行业地位来看，未来的 1-2 年左右，其他各个设备厂商相应的私有技术都会随之跟进。也许实现手段和市场宣传各有春秋，但就技术使用和用户体验来说不会有有多大差别。



再简单聊两句 VLAN 和 VPN 这两个熟透了的一虚多技术。VLAN 好理解，就是在 Ethernet 报头中加个 Tag 字段，多了个层次区分，将 Ethernet 层面数据报文划分从一维转成二维，瞬间规模就大了 N 倍，随之也使 Ethernet 的复杂度大大增加。VPN 就稍微复杂些了，谁让当初 IP 设计时候没有考虑预留其他 Tag 标识字段这块呢。

VPN（Virtual Private Network）分为本地有效的 VRF（Virtual Routing and Forwarding，也有说是 VPN Routing and Forwarding）和用于跨设备的 MPLS VPN 两个部分。VRF 就是将本地的路由表、转发表和三层接口都给个 VPN Tag，统一做 IP 层面的路由转发处理。例如从属于 VPN A 的接口进入设备的报文，只能查 VPN A 的路由表和转发表，从 VPN A 的其他接口转发出设备。当然后来又根据需求设计了一些跨 VPN 互访的技术，就不多说了。由于每个 VPN 都要维护自己的路由转发表，因此需要将各个路由协议 RIP/OSPF/ISIS/BGP 等都通过 VPN 标识隔离出多个数据库进程用于分开构造各自的路由表。这个也没啥难的，不管是数据报文还是协议报文都是基于接口出入设备的，因此只要将不同接口绑定到不同的 VPN 中，就可以做到 IP 路由层面的隔离了。而且这个 VRF 都是本地有效，各个厂家做的小有区别也不会相互影响。

跨设备转发时就麻烦了。首先得设计个 Tag 来让所有设备统一 VPN Tag，于是有了 MPLS Label；再就得让数据报文传输过程中带着这个 Label 四处游走，于是有了 MPLS 报头；还得让全体设备能够统一 MPLS 报头中 Label 对应本地 VPN 及 IP 路由的关系，于是有了 LDP/BGP VPNv4 等专用和扩展协议用于传递 Label。继续通过万能图解释 VPN 跨设备转发。



大的体系有了，还得补充细节。MPLS 报头在 IP 报头外面，而且字段又少，只能多裹层头，分层标识不同 PE 设备和 PE 设备上的不同 VPN（公网 Label 与私网 Label）；规模大了，不同区域的 Label 无法互相识别，于是要想办法 VPN 跨域；三层 IP 报文搞定了又想在 VPN 里面传二层 Ethernet 报文，于是有了 VLL/VPLS。。。

有兴趣的同学可以去 IETF 统计下 VPN 跨设备转发相关的 RFC，个人是实在数不过来了，到今天都还有各种相关 Draft 不断地推陈出新，排队审核呢。



## 5.6.4 小结

数据中心内部的服务器互访技术介绍到这里就告一段落了，无论是前面的服务器跨接入层互访还是后面的跨核心层互访模型，都对网络虚拟化提出了严峻的要求。可以说在后面的云计算数据中心建设中，非虚拟化的网络将很快的被挤出市场舞台。未来 10 年的大型数据中心网络是属于网络虚拟化的，从技术层面，目前只能看到一个独领风骚的前行者。

PS：声明一下，作者不是唯 Cisco 派的，是唯技术派的。

## 5.7 数据中心跨站点二层网络

本章节重点技术名词：RPR/VLL/VPLS/A-VPLS/GRE/L2TPv3/OTV

前面说了，数据中心跨站点二层网络的需求来源主要是集中云时的多站点服务器集群计算和分散云时的虚拟机 vMotion 迁移变更。搭建二层网络时，依据中间网络的承载方式不同，主要分为光纤直连、MPLS 核心和 IP 核心三类。这里的多站点一般指 3 个及 3 个以上，只搞两个站点的云计算喊出去会比较掉价。

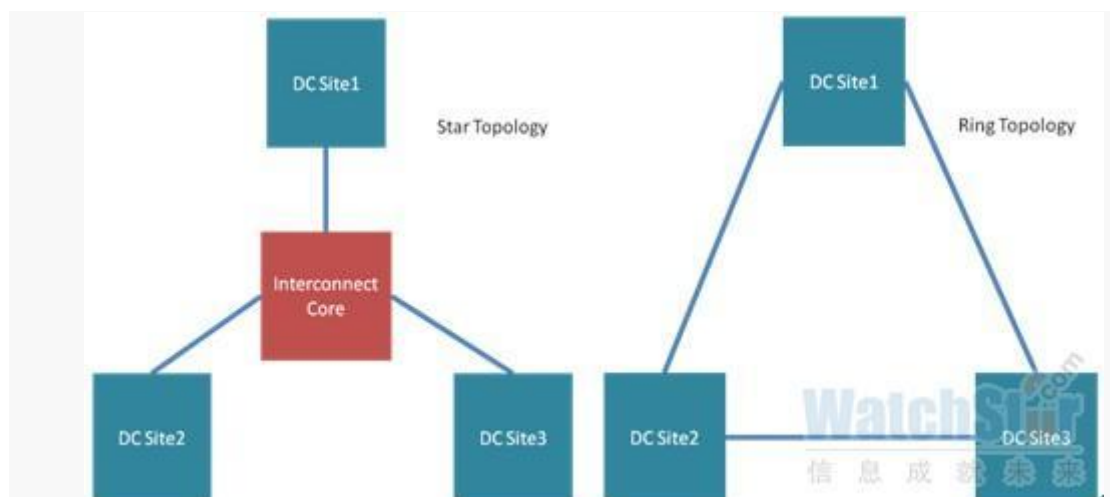
### 5.7.1 光纤直连

两个站点就不多说了，直接在两个站点的核心或汇聚设备之间拉两根光纤就 OK 了，也用不到什么特别的技术。唯一需要注意的是在两个站点之间的链路上做些报文控制，对广播和 STP 等报文限制一下发送速率和发送范围，避免一个站点的广播风暴或拓扑收敛影响到其他站点的转发。

当站点较多时，理论上有两种结构可用：

星形结构：专门找几台设备作为交换核心，所有站点都通过光纤直连到此组交换核心设备上，缺点是可靠性较低，核心挂掉就都连不通了，而且交换核心放置的位置也不易规划。这种结构不是值得推荐的模型。

环形结构：推荐模型，尤其在云计算这种多站点等同地位互联的大型数据中心组网下，环形结构既省设备省钱，又能提供故障保护，以后肯定会成为建设趋势。



从技术上讲星形拓扑不需要额外的二层互联技术，只部署一些报文过滤即可，可以通过链路捆绑增强站点到核心间链路故障保护和链路带宽扩展。而环形拓扑必须增加专门的协议用于防止环路风暴，同样可以部署链路捆绑以增加带宽冗余。

环形拓扑的公共标准控制协议主要是 STP 和 RPR (Resilient Packet Ring IEEE802.17)，STP 的缺点前面说了很多，RPR 更适合数据中心多站点连接的环形拓扑。另外很多厂商开发了私有协议用于环路拓扑的控制，如 EAPS (Ethernet Automatic Protection Switching, IETF RFC 3619, Extreme Networks)，RRPP (Rapid Ring Protection Protocol, H3C)，MRP (Metro Ring Protocol, Foundry Networks)，MMRP (Multi Mater Ring Protocol, Hitachi Cable)，ERP (Ethernet Ring Protection, Siemens AG) 等。

这里简单介绍一下 RPR。从控制平面看，环路拓扑组网相对简单，控制协议交互规则制定也比较前面的 TRILL/SPB 更加简化，了解了全网各节点位置后，确定内外环两条通路即可。在数据平面上，RPR 通过 MACinMAC 方式在环上封装外层节点 MAC 信息方式确认已知单播传递节点对象，非目标节点会将数据报文直接转给环上的下一跳，只有当目标节点收到此报文后根据外层目的 MAC 信息确认本地为终点，将报文下环转发。环上每个节点都会对未知单播/组播/广播报文着做下环复制和逐跳转发处理，直到转了一圈后，源节点再次收到此报文丢弃终止转发。

由于 RPR 在环路传输数据报文封装时增加了 1 个 Byte 的基本环控制和 1 个 Byte 的扩展环控制用于环路信息识别，因此也必须使用专用硬件处理环路接口的报文收发封装工作。RPR 虽然很早就确立了标准内容，但由于其初始应用针对运营商城域以太网，且只能支持环路拓扑，因此各个厂商并没有花太大力气去开发产品进行支撑推广，当前使用不多。

就作者看来，未来几年的云计算数据中心建设时，除非在所有站点采用相同厂家的设备还有可能使用一些私有协议组环（可能性比较低），前文提到预测会以站点为单位选择不同厂家进行建设，这时就需要公共标准用于多站点互联了。在光纤直连方式下成熟技术中最好

的选择就是 RPR，但如果 TRILL 能够将多拓扑这块内容定义好，未来是能够将其取而代之的。

### 5.7.2 MPLS 核心网

一些大型的行业企业（如政府军工）自建内部网络时，会使用 MPLS 技术搭建各个地方的互联核心网。此时可以将各地的数据中心站点复用 MPLS 核心网进行跨地域连接，省钱才是王道。在自建的 MPLS 核心网中，需要在各个站点的 PE 设备间搭建 VPLS 隧道用于传输 Ethernet 报文。如果是租用运营商的 VPLS 隧道则不需要考虑这么多，那时 PE 是由运营商提供的，对用户来说组网部署和前面的光纤直连没有区别。

#### VLL

如果是只有两个站点互联的情况，可以使用 VLL（Virtual Leased Line）。VLL 是一种点到点的虚拟逻辑链路技术，数据报文从隧道入口入，只能从定义好的另外一端出口出，不存在多个隧道终点一说。数据平面没啥可说的，A 点收到的二层报文进隧道直接封装上 MPLS 报头发给 B 点就 OK 了，整个过程框架可参考前面的 MPLS 转发图。控制平面由于隧道都是点到点连接方式，不需要复杂寻址，只要在数据流量传输时，给 VPN 分配外层封装的对应 Label 即可。分配方式有以下四种：

**CCC（Circuit Cross Connect）：**全网静态为 VPN 分配一个 Label，包括所有路径的 PE 和 P 设备都需要手工配置。此 Draft 已经处于 Dead 状态，目前基本也没人用了。

**SVC（Static Virtual Circuit）：**只在 PE 上静态配置私网 VPN 的 Label，公网标签不管。有用的但也不多，静态配置这种方式对故障处理总是心有余而力不足的。

**Martini：**RFC4762，使用 LDP 协议在 PE 间建立连接，为 VPN 动态分配 Label，省事好用。

**Kompella：**RFC4761，使用 BGP 协议在 PE 间建立连接，使用 BGP VPNv4 扩展字段携带 VPN 对应 Label 信息进行传递，这个实现起来比 Martini 复杂一点点，用得也就少了一些些。

后面两种 Martini 和 Kompella 方式在 MPLS L3 VPN 和 VPLS 里面也都有应用，都是作为控制协议来为 VPN 分配和传递 Label 用的。

#### VPLS

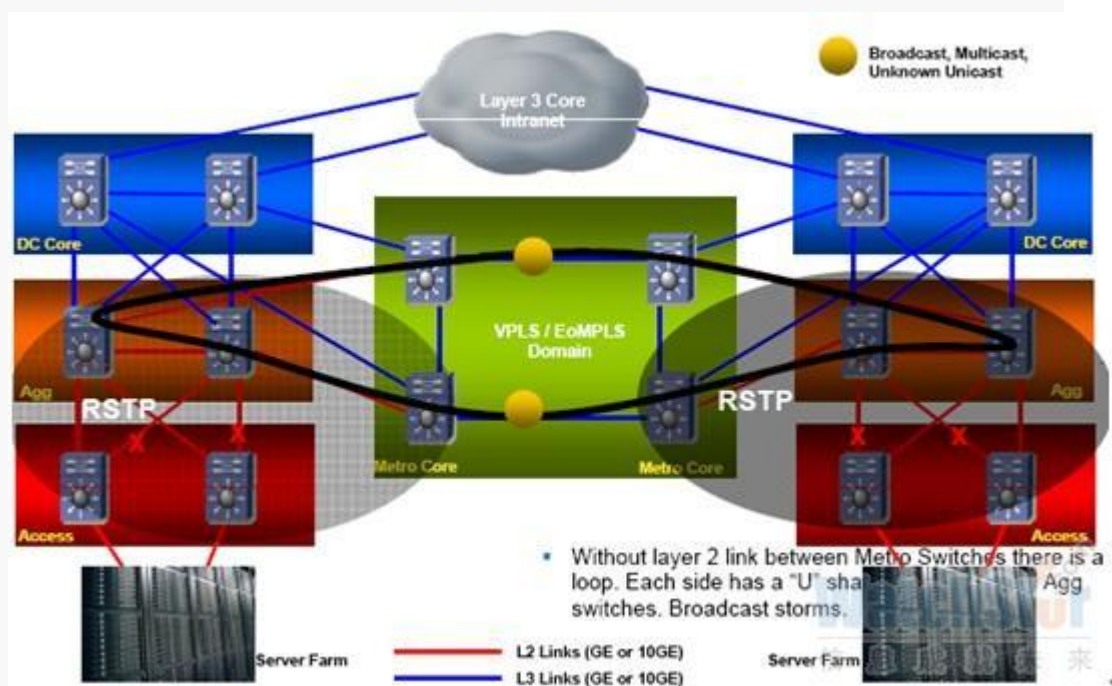
当存在多个站点时，A 站点收到的二层报文就有个选 B 还是选 C 进行转发的问題。于是有了 VPLS（Virtual Private Lan Service）。VPLS 是支持点到多点的虚拟链路技术，从隧道入口进入后，可以根据 VPLS MAC 地址表从多个隧道出口中去选择正确的出口，或者广播给所有出口。控制平面还是通过 Martini 和 Kompella 两种方式分配与传递 VPN 对应的

Label。数据平面则要多维护一张 VPN 的 MAC 对应 VC（Virtual Circuit）转发表，既前面提到的 VPLS MAC 地址表，本地接口收到的报文，MAC 地址学习方式还是和传统 Ethernet 一样；只有当报文从远端 PE 过来时，记录的源 MAC 需对应远端 PE 的 VC ID。

由于 VPLS 透传的是二层 Ethernet 报文，就涉及到 VLAN 标识处理的问题。VPLS 可以配合 QinQ 技术，将用户侧发来的带 VLAN 标签报文打上外层 VLAN 标签，以扩展 VLAN 数量规模。当然现在的交换机一般都是最大支持 4k 的 VLAN，大部分场景都是够用的了，还没有听说谁家的数据中心 VLAN 部署超过 4k。但云计算服务器节点数量规模成倍增加以后就不好说了，留出冗余总是好的。

为了防止广播风暴，VPLS 做了水平分割特性，PE 设备从远端 PE 收到的广播/未知单播报文只能发给本地的 CE，不能转发给其他 PE。还有其他的分层 PE 和 Hub-Spoke 等技术在数据中心多站点互联环境中一时还应用不上，这里也就不过多介绍了。

VPLS 技术已经很完善了，喜欢细节的同学可以去查下 RFC 相关文档。这里再说下其在数据中心多站点互联应用中的不足之处。数据中心要求的是全冗余，无单点故障点或单点故障链路，而 VPLS 在双 PE 冗余方面没有专门的定义，因此造成技术上的使用不便，一是会形成如下图所示的跨站点二层环路，二是本端 CE 无法感知对端 CE-PE 间链路状态情况，故障时导致流量黑洞问题。



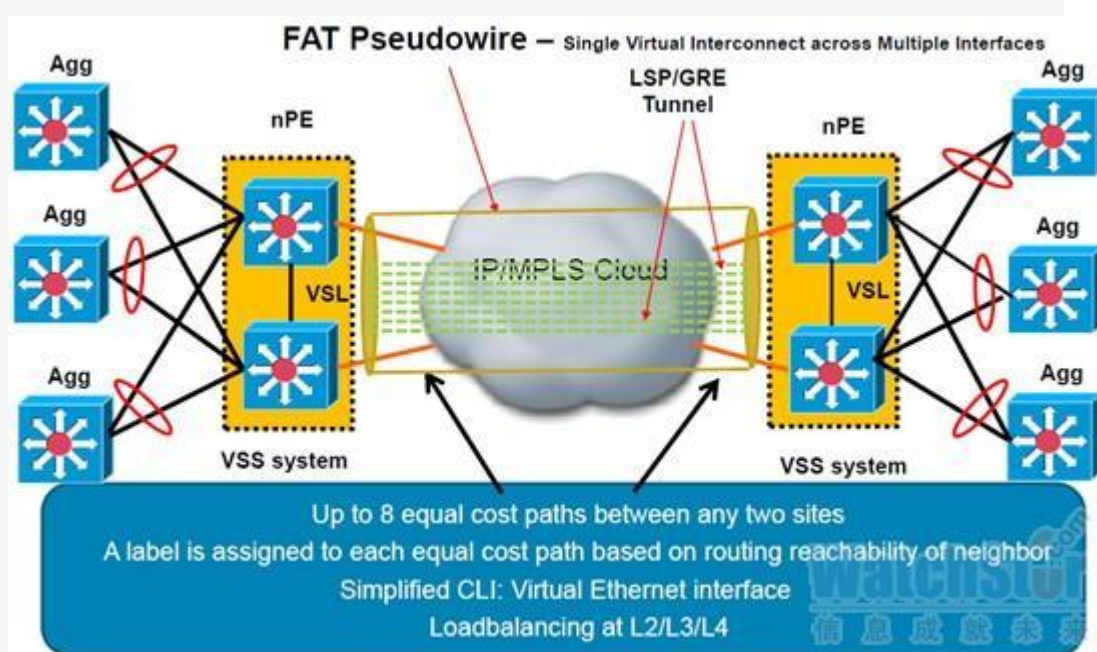
解决上述问题有以下两个思路：

首先是使用万能的 STP 构建出整网拓扑，即可避免环路，还可检测故障切换。缺点和前面在其他地方使用时一样，浪费带宽与收敛速度慢，另外就是要让 STP 跨站点组网，会

导致一个站点出现问题，其他站点全部受影响。此方案的好处就是公共标准大家都能做，而且不存在互通问题。

其次是使用控制平面多虚一技术，如 VSS/IRF 和 vPC，使多个物理节点变为唯一的逻辑节点将整个拓扑由环状变为链状，以避免环路。同时通过链路检测监控联动路径切换动作以避免流量黑洞问题，如 Cisco 的 EEM。这些小技术组合起来可以解决上述问题，但缺点是都是私有技术，没有统一标准，无法支持不同厂商产品混合组网。

另外可以一提的是 Cisco 的私有技术 A-VPLS（Advanced VPLS），此技术配合其 VSS，可以将多条 VPLS 的 PW（Pseudo Wire，可理解等同于 VC）虚拟化为一条逻辑的 Fat PW，达到多 PW 路径负载分担的效果，和链路聚合很类似。如下图所示。



此技术由于需要往 MPLS 报头中添加一个 Flow Label 的标签字段，用于处理多 PW 的流量路径 Hash，因此别的厂家设备肯定无法识别，只能在全 Cisco 设备环境下部署。其他厂商也有开发出类似的技术，对多 PW 进行流量负载分担，但也都是私有的小特性，无法互通组网。估计再有 1-2 年 IETF 可以搞出个多 PW 负载分担的标准来，到时大家就好做了。

### 5.7.3 IP 核心网

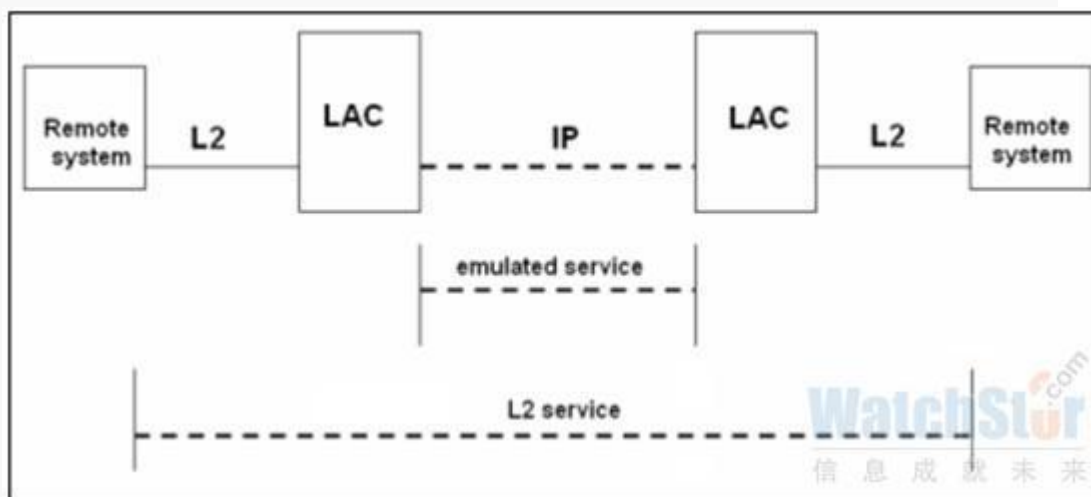
全球最大的公共 IP 核心网就是 Internet 了，只要解决了报文加密的安全问题，且 Internet 出口带宽足够大，谁说以后的数据中心站点间二层互联不能走 Internet 呢。另外也有很多大企业的核心网采用 IP 建网，国内的如金融电力，国外则遍地开花。

从技术上来看，公共的技术标准主要有 VLLoGRE/VPLSoGRE 和 L2TPv3，私有技术就是 Cisco 的 OTV 了。VLLoGRE/VPLSoGRE 没啥好说的，就是在 IP 层打通个 GRE 隧道，再把 Ethernet 报文扔到隧道里面传。下面主要说说 L2TPv3 和 OTV。



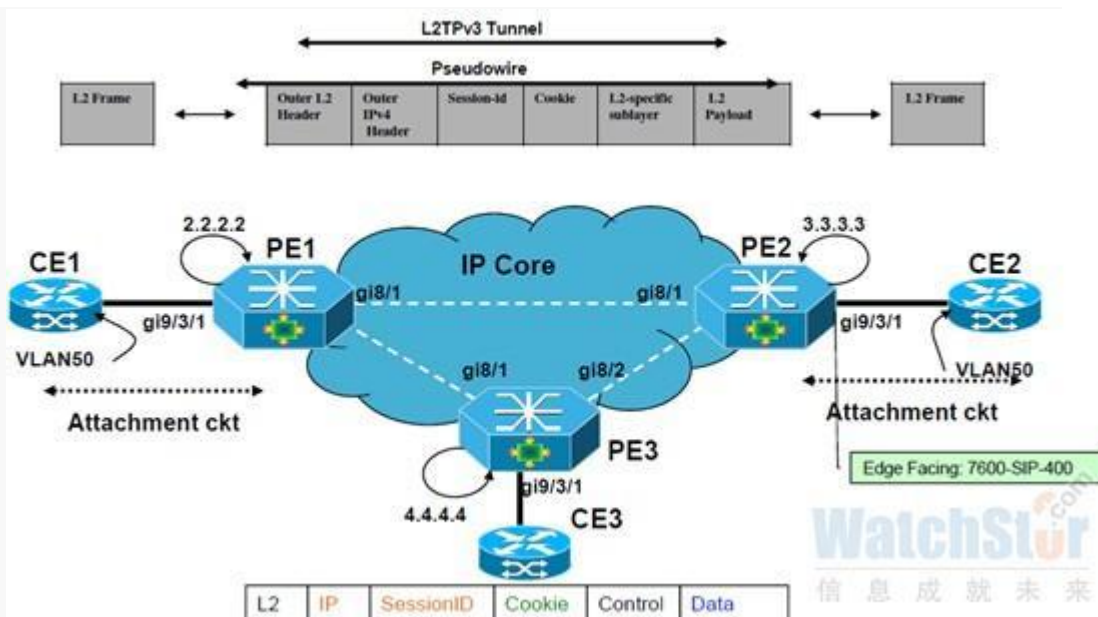
### L2TPv3

L2TP (Layer 2 Tunneling Protocol) 是 IETF RFC2661 (L2TPv2) 定义的, 已经有不少年头了, 主要应用于移动办公通过 Internet 进行 VPN 接入等场景。L2TPv2 的标准封装是基于 PPP 格式的, 后来为了进行应用扩展, 推出了 L2TPv3, RFC3931。L2TPv3 可以封装在 IP/UDP 层之上, 摆脱了 PPP 的束缚。



LAC (L2TP Access Concentrator) 是 L2TP 的角色名称, 作为 IP 隧道的起终点, 另外还有个 LNS (L2TP Network Server) 角色, 但在数据中心多站点互联场景中应用不到, 这里就不做介绍了。可以简单理解 LAC 等同于 VPLS 的 PE。

从控制平面看, L2TPv3 使用自己的控制协议报文建立 IP 隧道, 由于隧道都是点到点的, 也不存在什么拓扑学习的问题, 这点和 VLL 相同。数据平面就是进隧道封包转发了, 封装时只需在原始 L2 报文和外层 IP 头之间插入一个 L2TPv3 报头即可, 里面包含 SessionID 和 Cookie 字段, 其中 SessionID 字段 32Bit, Cookie 字段可选, 最长 64Bit, 整个 L2TPv3 报头最大 12Byte, 要大于前面那两个 oGRE 的了。整体结构可参考下面 Cisco 胶片的截图。但是注意里面虽然画了 3 个 PE, 但是实际上 L2TPv3 和 VLL 一样只能支持点到点的传输, 如果 CE3 上也有 VLAN50 想和 CE1/CE2 一起组成二层网络 L2TPv3 是搞不定的。这个图有那么点儿混淆概念的意思。其中的 Control 和 L2-specialfic sublayer 字段在数据中心互联场景中就是指 Ethernet 报头, 而在其他场景中也可以使用 PPP 等其他二层链路协议报头, 毕竟 L2TP 定义的是要承载所有 L2 数据报文。

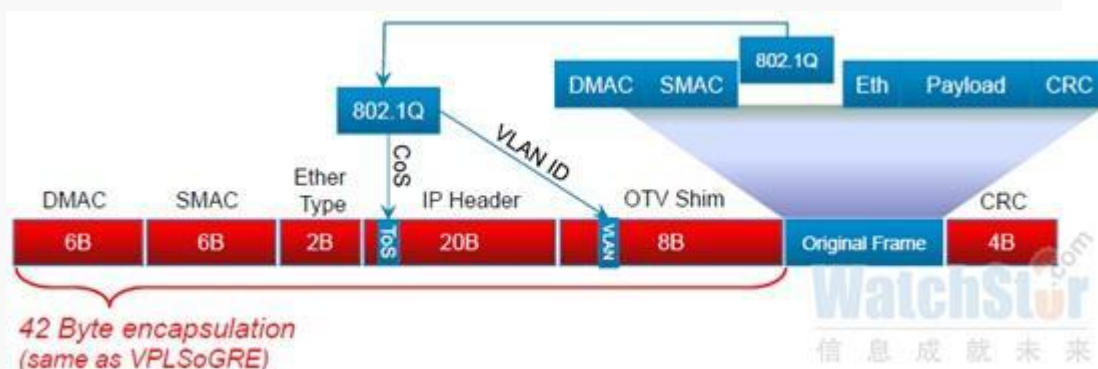


L2TPv3 是 RFC 标准，谁家都能用，但只能像 VLL 一样支持两站点互通的场景，因此使用并不广泛，当前也就看到 Cisco 有在小范围推广。VPLSoGRE 还是大部分厂商在多中心互联中主推的技术。顺便一提，L2TPv3 只解决跨 IP 封装传输的问题，对于前面提到的多 PE 和流量黑洞的问题并没有对应方案，还是得配合 VSS/IRF 和探测处理等私有技术才能适用于数据中心。

## OTV

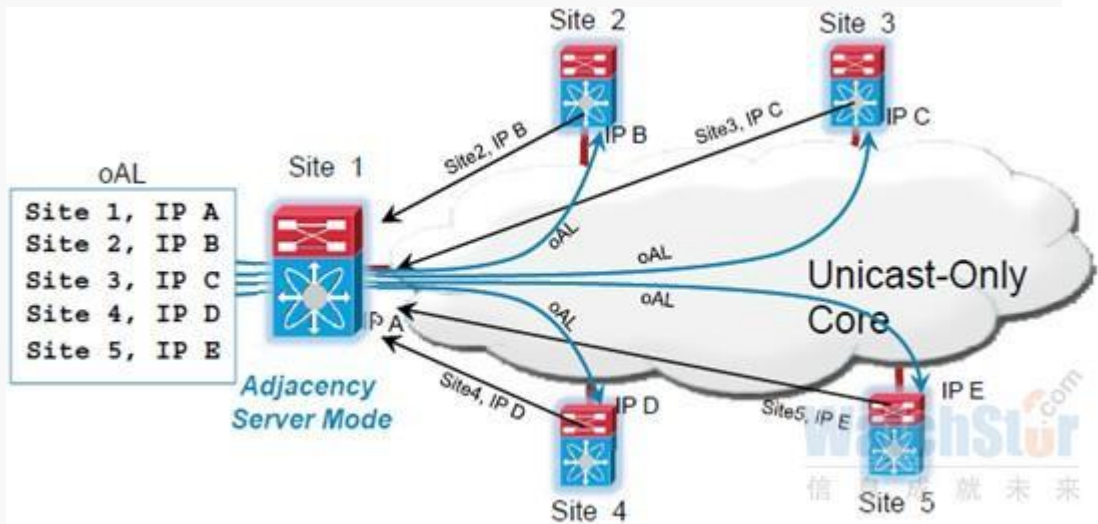
前面说了，在 IP 核心网情况下，公共标准只有使用 VPLSoGRE 才能支持多站点的二层互联，而 VPLSoGRE 同样存在前面 VPLS 组网中的多 PE 连接和流量黑洞问题，需要配合其他一系列的私有技术才能一并解决，部署起来相当繁琐，而且真出了问题定位也很困难。公共标准在这种场景下不是不能用，而是不好用，只能看私有技术了。

Cisco 在其新一代数据中心交换机 Nexus 7000 中推出了 OTV（Overlay Transport Virtualization）私有技术来专门处理数据中心多站点二层互联使用场景。数据平面 OTV 以 MACinIP 方式封装原始 Ethernet 报文，报文结构如下：



可以看到对比 VPLSoGRE 只是将 8 字节 GRE 头替换成 OTV 标识，长度没有变化。因此转发效率估计和 VPLSoGRE 相当。

控制平面上，OTV 有组播与单播两种方式建立邻接拓扑，组播方式适用于支持组播的 IP 核心网，个人觉得还是很少见的。单播方式需要设置一台 AS（Adjacency Server），保存所有的邻接设备信息列表 oAL（overlay Adjacency List）。所有 OTV 节点需要手工设定此 AS 地址，上线时去取得其他邻居节点信息以建立邻接。当然还可以配置备份 AS 节点避免单点故障。OTV 单播邻接建立方式如下截图所示。



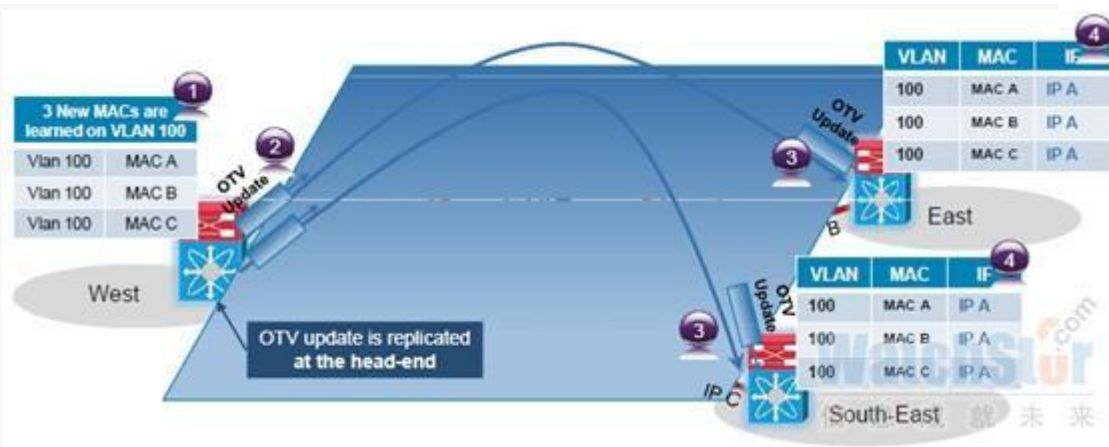
针对数据中心多站点互联的场景，OTV 设计了一系列的机制处理。

1、STP 隔离。将 STP BPDU 报文在 OTV 边缘设备 ED（Edge Device）上进行阻塞，禁止其跨站点传播。

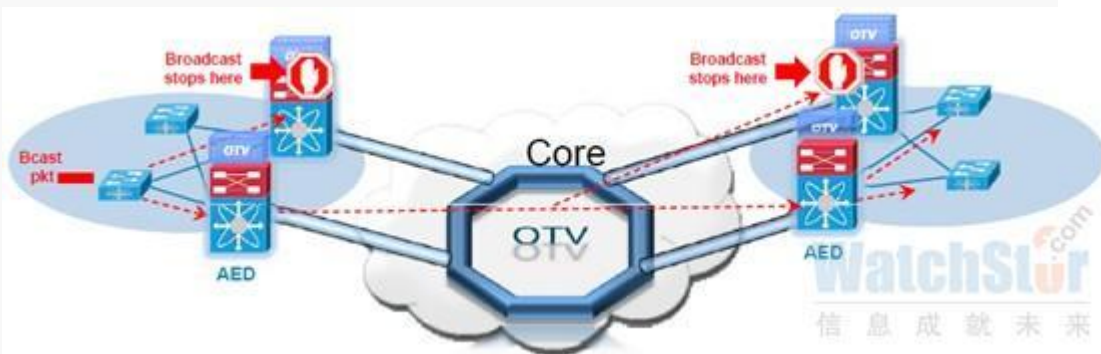
2、未知单播隔离。将未知单播数据报文在 ED 上进行阻塞，禁止跨站点广播。同时可以手工配置静态 MAC 地址对应远端 OTV 接口的表项，以应对部分静默主机应用场景。

3、ARP 控制。对远端站点返回的 ARP Reply 报文进行 Snoop 和 Cache，当再收到本地查询同样目的 IP 的 ARP Request 时直接代答，不向其他 OTV 站点扩散，减少跨站点的 ARP Request 广播。

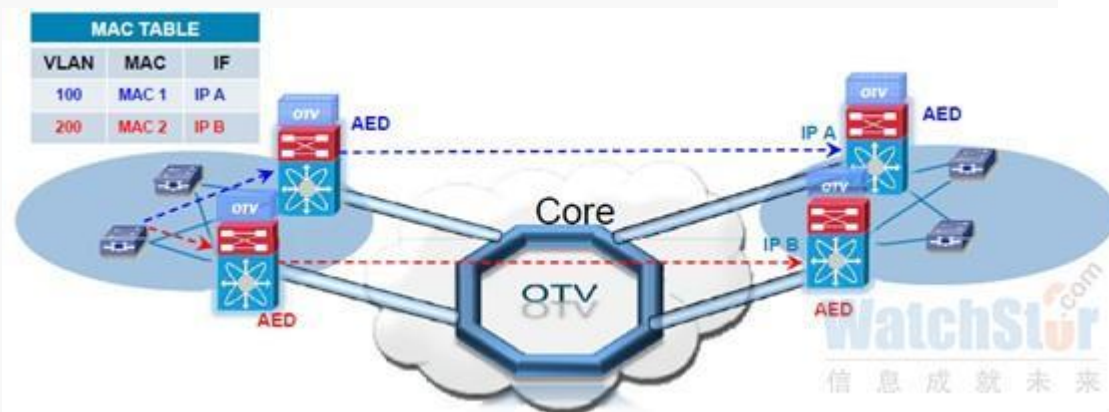
4、MAC 地址学习控制。由于未知单播报文被隔离了，因此需要通过 OTV 协议报文进行站点间的 MAC 地址学习同步。过程如下图所示。跨站点的广播报文的 MAC 地址学习规则仍然与传统 Ethernet 相同，而且 OTV 不会对其做特殊控制。广播报文限速这种功能现在基本是个交换机就能支持了，算不上特色技术。



5、站点 ED 双机冗余。可以在一个站点使用多台 ED 接入 OTV 核心网，前提是要保证多 ED 在站点内部可以二层互通。运行控制协议进行 AED（Authoritative Edge Device）设备的选举，只有 AED 可以转发和接收广播/组播报文，以避免环路风暴。如下截图所示。



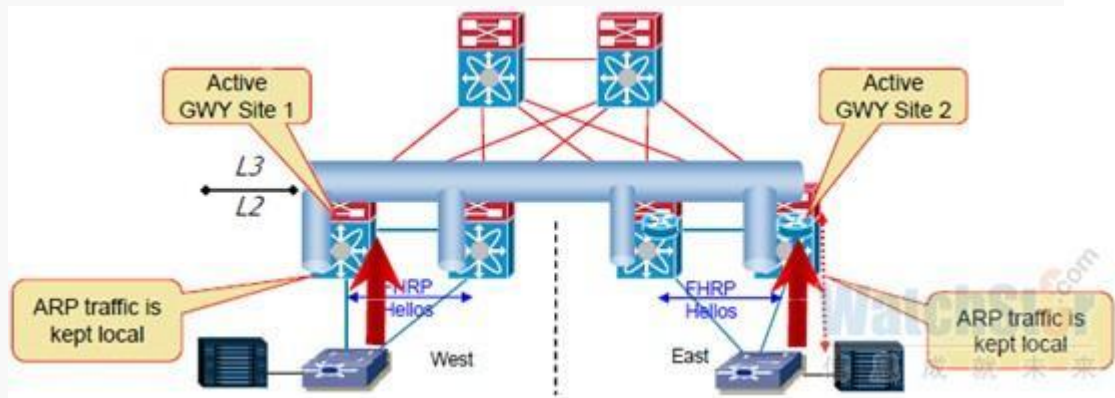
另外可以在多个 ED 间基于 VLAN 进行不同的 AED 选举，达到全局意义上的流量转发负载均衡。如下截图所示。



6、HSRP 隔离。不管二层技术再怎么搞，除非采用 VSS/IRF 等私有技术去做控制平面整体虚拟化，否则网关冗余在数据中心里面都是必不可少的。而当多个站点二层通信时，也必然存在网关部署位置的选择问题。如果站点 A 的主机每次和外界通信都走站点 B 的网关，则会导致大量的非必要流量途径站点间二层互联链路，浪费带宽且路途绕远。Cisco 提出了



更好的处理方式，OTV 在 ED 上通过过滤 HSRP HELLO 报文，将 HSRP 进行站点间隔离，这样各个站点的网关相同，但各自为政，上下行流量路径最优。



OTV 的主要内容差不多就这些了。从技术上来讲，其细节考虑最全面，而且瞄准的市场是所有站点间 IP 可达的应用场景，既 OTV 同样可以部署在光纤直连和 MPLS 核心网的场景中，远远领先于其他的数据中心多站点互联技术方案，但由于其私有协议的地位和可能的性能瓶颈，在市场上最终能占多大地盘还有待观察。这里有两点猜测：

1、OTV 如果今年仍在米国拿不下专利，估计明年就该去 IETF 提 Draft 了。米国的专利还是很严格的，这种继承性居多的技术审查通过不易。再想想国内的一些相关领域技术专利，都是神马的浮云。

2、未来的一两年间，其他各个瞄准数据中心的设备厂商类似（用仿字不好听）OTV 的私有技术将会如雨后春笋般发布出来。

## 5.7.4 小结

数据中心跨站点二层互联目前还是个比较新的需求，实际上马的项目并不多，以前大多是满足存储需求的光纤直通。但大潮已经涨起来了，国内的各大运营商纷纷启动测试，国外估计也有一些项目应用。各个厂商最好尽快做好弄潮的准备，不然很快就会被淹没的。

从各个角度分析，个人认为上述三种互联方式中还是光纤直连最靠谱。

首先从省钱角度来看，建得起数据中心多站点的就不会差那几个钱拉不起光纤，如果使用 CWDM/DWDM 等波分复用设备可以在一对光纤上建多个通道供存储和不同业务共享使用，总带宽最高可以上 T，性价比不见得比租用和复用 MPLS/IP 核心网络要低。

其次从重要性角度来说，需要跨站点运算的应用程序肯定是以企业关键业务为主，给关键业务拉根专线，不去和其他业务搅合，也是可靠性设计的必要需求。

最后从当前可使用技术的角度考虑，VLL/VPLS 的数据报文查表封装处理工作，新一些的转发芯片可以搞定，但 L2TPv3/GRE/OTV 一般的转发芯片肯定是搞不定的，这样就需要



额外的 CPU 或 NP 去处理数据报文，性能上自然也不会有什么太好的期待，万兆线速都是奢求。而且不管是怎么封装，传输的时候外面一堆报头都很损耗带宽，列举个极限的例子，OTV 和 VPLSoGRE 外层报头要封 42Byte，对 64Byte 载荷的数据小包，报头带宽损耗达到了  $42/(42+64)=40\%$ ，一条 10G 链路，跑满了才能用 6G 传数据，效率那是相当的低啊。

而光纤直连方式目前主要技术瓶颈在多站点互联时，缺少专业高效的公共标准技术，RPR 如果想在数据中心有所作为，还需要进行一些协议上的标准改进，多考虑一些如未知单播/广播数据报文和 STP/ARP/VRRP 等协议报文的数据中心场景专门处理，以及站点多边缘设备冗余接入应用场景处理。另外各个厂商对相关产品的高调发布和大力市场推广也是必不可少的。个人觉得在数据中心多站点互联这块，私有技术早期也许能忽悠住一些用户，但最终肯定是站不住脚的。

## 5.8 数据中心多站点选择

本章节重点技术名词：SLB/GSLB/LISP

数据中心多站点建设时考虑应用服务冗余，则必然面临着以下问题：1) Client 访问 Server 的时候选 A 站点还是 B 站点；2) A 站点的 Server 故障或服务迁移到 B 站点后，Client 访问如何能随之快速切换。

前面已经提到，在选路技术中主要有两个解决方案思路，一是 DNS，二是路由。

DNS 方案的缺点首先是应用扩展性不强。DNS 协议以处理 HTTP/HTTPS 的应用为主，其他类型应用较少。不过话说回来，目前的应用中基于 WEB 的 BS (Browser-Server) 结构快打遍天下无敌手了，这个问题倒也影响不大。还有个问题就是 DNS 自身协议设计时，没有考虑多 IP 选择问题，所以此类解决方案都要和主机探测、迁移同步等私有技术配合起来使用，一般都得好几个盒子联动起来才行。DNS 的好处是简单，使用的技术都是成熟技术，准备好接口写两行代码谁都能分上一杯羹，只要操心优化处理性能方面即可。目前的典型技术代表就是 GSLB+SLB (可选)+vMotion 通知 (可选)，下面会进行详细介绍。

路由的方案是完全基于 IP 技术出发，网络设备厂商自己就能搞定，不需要找 DNS 或 vMotion 去联动。其中的主备中心的路由掩码比明细和主机路由发布两种方案都不是啥好招，属于拆了东墙补西墙，会造成其他的问题。而稍微完整一点儿的 LISP 目前也还是试验探索阶段，vMotion 后主动路由刷新这个最主要需求，也没能得到太好的解决。下文会简单介绍 LISP，并探讨更优的处理方案。

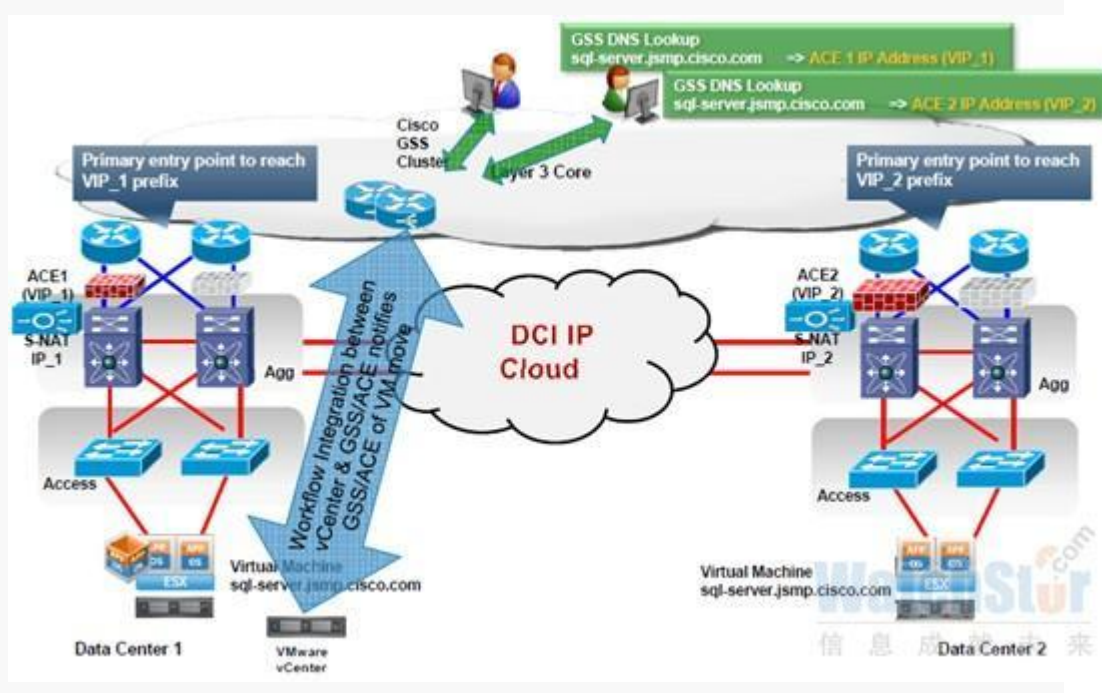
### 5.8.1 GSLB

GSLB (Global Server Load Balance) 全局负载均衡技术，这个貌似是 Redware 先提出来的叫法，其他各个厂家叫法都有区别，如 F5 的 3DNS 和 Cisco 的 GSS 等。作者觉得这个词

更贴切而且技术性普适一些，就按这个名词进行技术介绍，从原理上讲和 3DNS/GSS 什么的没有啥大的区别。

GSLB 就是一台 DNS 解析服务器。最主要功能就是对不同 Client 发往相同域名服务的请求，以一定算法规则进行 Hash，回应不同 Server IP 地址。常见的算法包括轮询（90%都用的）、最少连接数和服务器最快响应速度等。增强功能是可以对 Server IP 进行探测，如果探测到某台 Server 故障，则会使用其他正常的 Server IP 进行 Client 的 DNS 响应。常见的探测方式有 ICMP（90%都用的）、TCP 以及上层应用如 FTP/HTTP 等。当然也可以再搞些 HTTP 重定向等特性，将 GSLB 设备放在数据中心站点的入口便于故障快速切换。

在 vMotion 应用场景中，由于 VM 在迁移前后的 IP 地址不变，因此两个数据中心站点的 VM 对外提供服务的 Server IP 地址相同，GSLB 此时就需要服务器前面的 SLB（Server Load Balance）设备进行配合了。SLB 是个 NAT（Network Address Translation）服务器，主要作用就是将后端真实服务器的 IP 和 TCP/UDP Port 等映射为对外提供服务的虚拟 IP 和 TCP/UDP Port，然后将不同 Client 访问虚拟 IP 的流量修改目的 IP 后，分别发到后端的不同真实服务器上，以达到对后端多台真实服务器的流量负载均衡效果。SLB 同样需要对真实服务器进行探测，以及根据不同的算法规则将 Client 流量均匀 Hash 到不同的真实服务器上。使用不同的 SLB 可以将后端相同的真实服务器 IP 映射为对外的不同虚 IP 地址，此时配合 GSLB 就可以解决 vMotion 前后 VM 服务 IP 相同的迁移切换问题。另外如 VMware 的 VM 管理控制平台 vCenter 可以将 vMotion 动作通知 GSLB 设备，达到快速切换效果。下面以 Cisco 的技术结构截图举例，其中的 GSS 就是 GSLB，ACE 为 SLB。其他厂家的方案在技术结构上也没有啥根本区别。

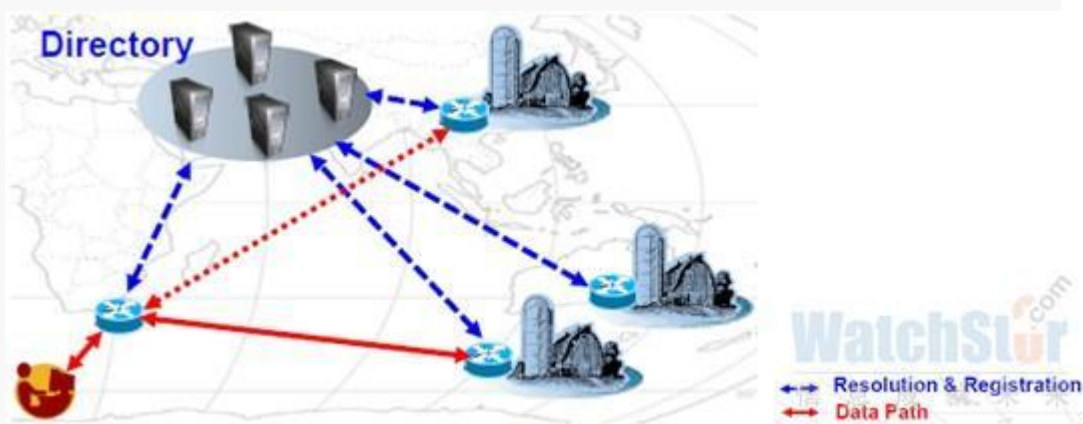


说实话个人觉得 GSLB 技术没啥难度可言，DNS 解析和调度算法都是现成的东东，稍微有点儿技术实力的厂商都能做个盒子出来。这个东西关键在于性能，由于 DNS 解析等上述功能行为都得靠 CPU 实现，没啥公用芯片，那么就看谁家的算法实现效率高，谁家支持的新建并发规格大。如果性能规格差不多，都能满足需求，再要考虑的就是可靠性和性价比了。关于衡量数据中心性能和可靠性的问题，会在本文的相关外篇中再深入讨论。

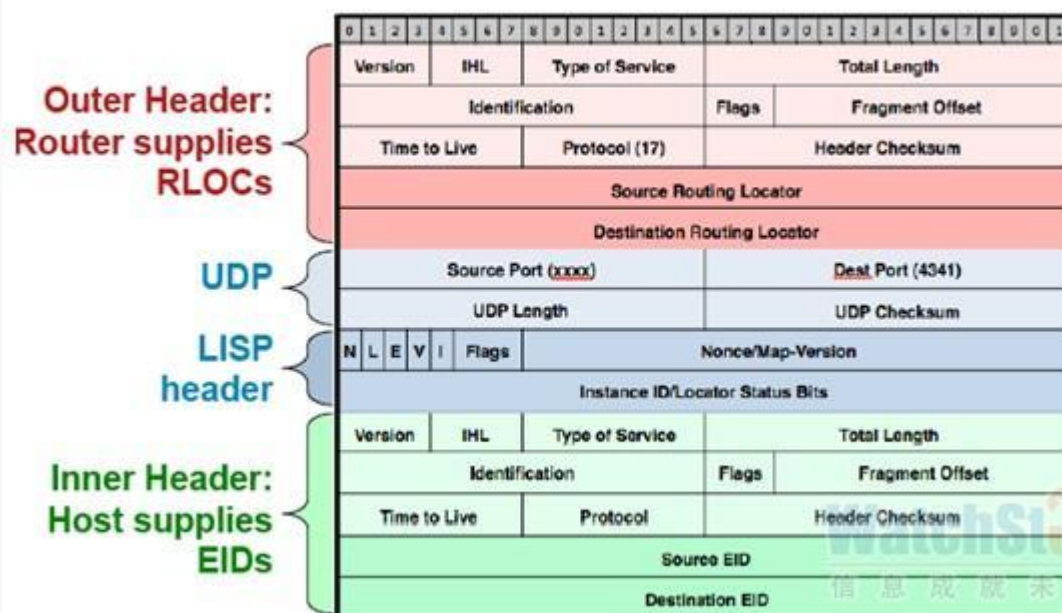
## 5.8.2 LISP

LISP（Locator/ID Separation Protocol）实质是个 IPinIP 的协议，其主要思想早在 15 年前就已经被人提出来进行研究，然而一直没有太具体的东西产出。直到 2006 年，Cisco 重新开始投入资源进行研究，目前已经提交了很多的 IETF Draft，最新版是今年 4 月份的 Draft Version12。但就应用来说，Cisco 的 LISP 目前也只处于试验阶段，距离能够推广商用还有不短的时间，很多技术细节方面问题需要解决。

LISP 的应用结构如下截图所示：

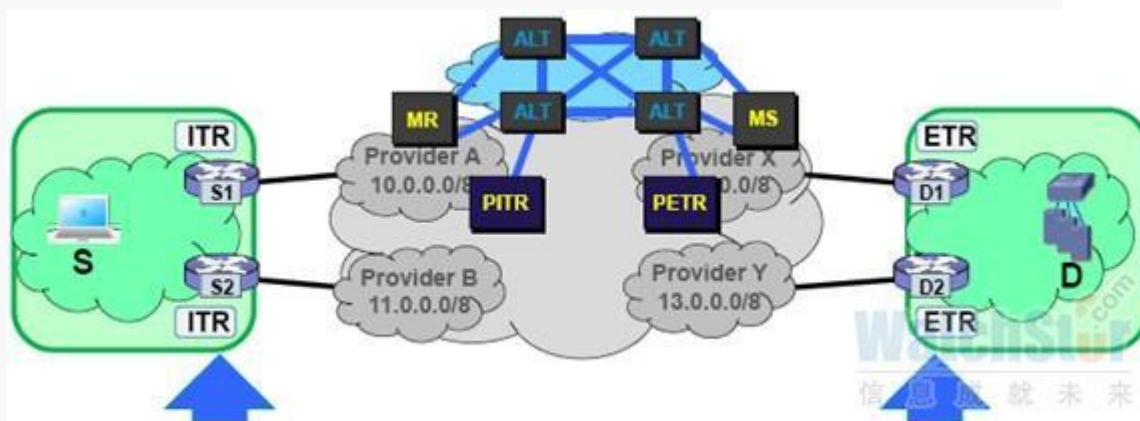


LISP 提出将标识 Locator 的 IP（RLOC）和标识目的节点 ID 的 IP（EID）进行区分和叠加封装，在公网传输时只根据 Locator IP 转发，只有到达站点边缘时才会剥离外层 IP，使用内层标识 EID 的 IP 进行转发。从下面的报文封装截图就可以看到其 IPinIP 的思想。



LISP 有两个主要的目标：一是公网设备不需要学习站点内部明细 IP 路由项，可以有效减少公网的路由数目；二是当访问的目标服务在站点间迁移时，可以只变更 Locator 的外层 IP，不需要对服务节点的内部 IP 地址进行变更，可以避免 TCP 等上层应用的中断重建，此点主要是应用于数据中心 vMotion 和手机上网漫游的场景。

LISP 的技术结构如下截图所示：



上图的名词很多，通过简单描述整个数据转发过程来帮助大家进行理解：

- 1、由源主机发往目的主机的数据报文第一次到达客户区域的 LISP 边界设备 ITR。
- 2、ITR 会根据报文的目的 IP 地址 EID，向 Directory 区域的本地查询服务器 MR 请求 EID 对应的目的站点边缘设备公网 Locator IP。

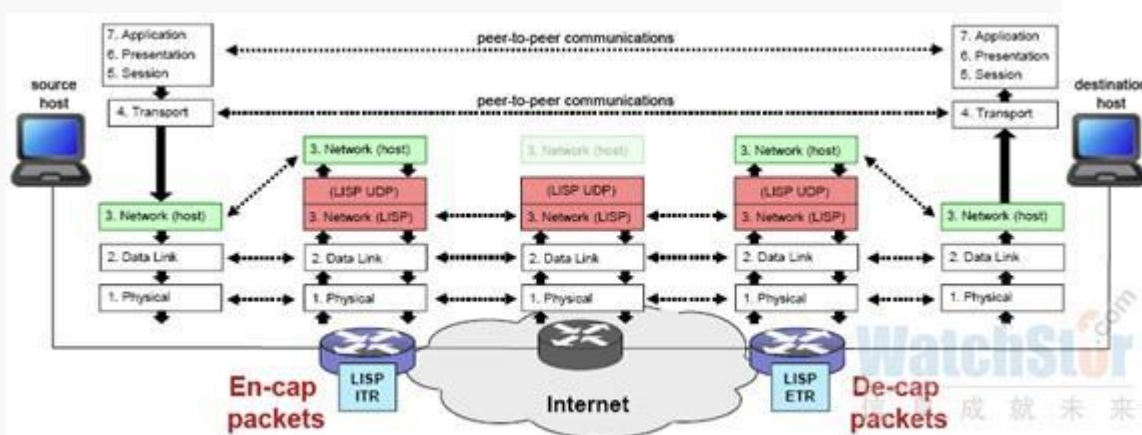


3、MR 会根据本地 EID 所属路由网段与 MS 的对应表项将查询请求提交给数据服务器 MS。

4、MS 上拥有 EID 对应目的站点边缘设备 ETR 的对应表项信息，MS 会查表将此请求转发给 ETR。

5、ETR 会根据自身设置的规则（如优先级等）选择站点的某个公网 IP 作为 Locator IP 反馈给 ITR。

6、ITR 会根据 ETR 回应报文中的 Locator IP 封装外层报头将数据报文发到公网上，同时记录此 EID 与 Locator IP 的临时对应表项，当再有去往此 EID 的数据报文流经时直接封装转发。封装转发的过程如下截图所示，也有些眼熟吧。



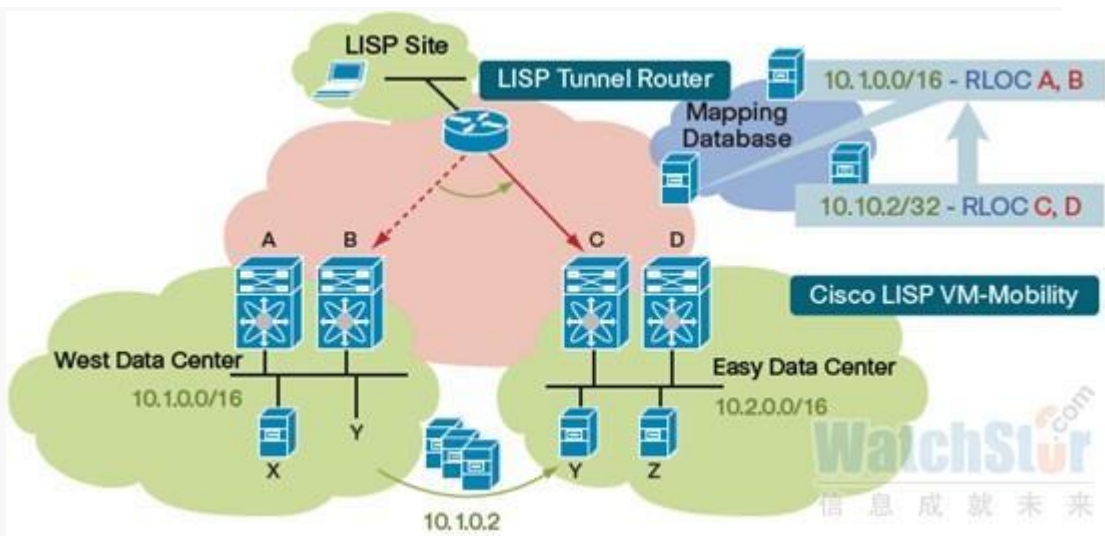
如果觉得 1-5 步的过程复杂不易理解，请回想一下 DNS 整套的域名解析过程，都是相通的。注意上述名称都是技术角色，一台设备可以实现多个 LISP 的角色功能，如同时实现 ITR/ETR 功能，或同时实现 MR/MS 功能等。

另外 ALT 是用于搭建 MR 和 MS 之间 Directory 区域互联用的中间角色，通过 BGP 扩展报文为 MR 和 MS 之间传递路由信息；PITR 和 PETR 是用于 LISP 与不支持 LISP 的网络对接时做 ITR 和 ETR 代理用的角色。由于目前 LISP 也还没有定稿，此部分设备功能没有完整定义，有兴趣的同学请自行深入研究。

LISP 中各角色之间大都通过手工指定的方式建立连接关系，如 ITR 上需要指定 MR 地址，ETR 上需要指定 MS 地址，只有 MR 和 MS 之间可通过 BGP 来建立邻接关系并通过扩展报文传递 EID 表项信息，但目前实现出来的还是以手工指定方式为主。而且 ETR 上要将哪些 EID 信息发送到 MS 上，也同样需要通过配置网段掩码的方式手工进行指定。

LISP 并不是专门为数据中心开发的技术，因而 Cisco 如果想将其在数据中心场景进行研究推广，估计会进行一些协议改造使其更加适用于数据中心的场景需求。目前 Cisco 给出的 LISP 数据中心实现 vMotion 过程如下截图所示：





上图是能找到的里面相对描述最清楚的了，但说实话感觉还是很糙。个人理解如果希望 LISP 应用于数据中心多站点选路，还需要解决以下一些技术问题：（下面这几段读起来可能会有些费劲，珍惜脑细胞的同学慎入）

1、迁移后服务器 EID 在新站点的 ETR 注册问题。既 VM 迁移后，新站点的 ETR 如何知道它此时需要向外发布对应 EID。Cisco 的当前做法是使用 IP 报文侦听，先配置个侦听范围，既可能会迁移过来的 IP 地址段，当监听到本地出现此地址段为源 IP 地址的 IP 报文时，会激活此 EID 表项并发给 MS 进行注册。个人感觉侦听免费 ARP 会更方便一些，vMotion 后 VM 肯定会发免费 ARP 报文的，但发 IP 报文就得看服务器的应用层协议设置了。不过此方案需要在站点间过滤免费 ARP，不能使其跨站点传输，否则二层隧道会将免费 ARP 扩散到所有站点，侦听就没意义了，而过滤后会不会有其他问题还需细琢磨。这里只随口提个思路，什么方案都是有利就有弊的，需权衡清楚再实现。另外也可以让 vCenter 等管理平台去通知 ETR，类似于前面的 GSLB 方案，这样由必须和 VMware 等虚拟机厂商做强联动，有些违背使用 LISP 的初衷。

2、迁移后通知 ITR 快速切换新的 Locator IP 问题。这个就更加复杂了，VM 刚由站点 A 迁移到站点 B，ITR 不知道啊，还是在用旧的 Locator IP 封包发送，此时应用业务肯定就断了，直到 ITR 获取到新的 Locator IP 后才能再建立起连接恢复业务。如果是时间敏感型的业务，中断个几分钟，上下几百万就没了不是。当前 LISP 提了很多解决方向出来，但还没有什么确定的技术方案。个人思路如下，此问题需分解为两个小问题各自解决：

首先是要解决 ITR 如何感知迁移发生的问题：1）管理平台通知，需要联动，而且 ITR 那么多，也不会都注册到管理平台上。不太靠谱。2）由 ITR 探测 EID 存活状态，探测范围太大，EID 可能以主机路由居多，对 ITR 负担太重，可行但不是很好。3）由原始站点 ETR 探测 EID 状态，当迁移后探测到 EID 不在本地站点，则通知 ITR 删除临时表项。ETR 上是保存了所有 ITR 表项的，可以很方便知道都要通知谁。但由于各站点间服务器前端网络二

层互通，因此还要想办法将此探测报文在站点间隔离，否则迁移前后始终都是通的，和前面的 ARP 侦听方案存在相同的问题。4) 在问题一已经解决的情况下，当服务器 EID 在新站点的 ETR 注册完成后，由新站点 ETR 向所有原始站点 ETR 发通知，再由原始站点 ETR 通知 ITR 删除临时表项。这个方案感觉相对更靠谱一些，可以将所有可能运行同组业务的数据中心站点 ETR 都设定到一个组里面，大家有事没事互通有无一下。

再有要解决 ITR 感知迁移后如何切换 EID 对应的 Locator IP 问题：1) 使用上面几种感知迁移发生的方法后，都可以将 ITR 的 EID 对应 Locator IP 临时表项删除，由 ITR 重新发起一套 EID 的寻址流程获取新的 Locator IP，缺点是稍慢。2) 在上面第 4 种解决方案中，原始 ETR 收到新 ETR 的通知后，向 ITR 发个 EID 变更信息告知新的 ETR 地址，由 ITR 向新的 ETR 直接发请求获取新的 Locator IP，不经 MR/MS 倒腾一遍手了，这样需要在 LISP 中多定义一个 EID 变更报文和相关处理流程。切换速度能提升些，但也稍微复杂了些。

提问题->找多个解决方案->比较不同方案的利弊，协议设计就是这么个过程。

LISP 即使能够成事，至少也得 2 年以后了，所以大家可以先看个热闹，等 RFC 标准立起来再介入也不迟。也许 Cisco 回头觉着整这个太费劲，说不定哪天就偃旗息鼓了。别的厂商又真不见得有这么大号召力能把 LISP 忽悠起来。全当学着玩了，目前别拿 LISP 太当真。

另外多说一句 LISP 在 Mobile 里面的应用场景，Cisco 已经今年 5 月份已经向 IETF 提出了 draft-meyer-lisp-mn-05，其中 mn 就是 mobile node 缩写。简单来说就是在手机上支持 ITR/ETR 功能，可参考下图：



技术上很有想法，市场发展上不咋看好。感觉 Cisco 把手伸到手机里面，还不如前文提到的伸进服务器虚拟化里面有搞头。

### 小结

小结一下，数据中心多站点选路目前网络厂商单从路由角度来看没有什么好的方案，还是用 DNS 搞定更靠谱一些，只要应用程序的 BS 结构始终领先前行，暂时就还不用考虑其他的解决方案。让那些有钱有势的大厂商去试验吧，大家在后面跟进就是了。搞预研是有一定风险的行为，资源消耗了，万一路没选对，大厂商还能壮士断腕，小厂商就得折腰而终了，须慎研慎行。

凡事都有两面，换言之，混乱的局面也是崛起的机会，如果谁能想到个啥招，搞个盒子出来，从路由或者其他层面独立解决多站点选路的问题，还是可以一试的。大赚不好说，但搞到像 F5/Redware/Citrix 这种规模并不是没得奔头。

## 5.9 技术总结

对数据中心网络而言，当前的技术发展正处于一个关键时期，虽然 Cisco 暂时占先，但只是通过技术的先进性增大了其市场话语引导能力，为其在今后 10 年的数据中心市场的竞争中增加了一些砝码，绝不是说就一定会所向披靡。未来充满变数，哪怕是什么时候 Google 或 Tencent 推出了应用于云计算数据中心的网络设备或技术，作者都完全不会感到吃惊，一切皆有可能。就像 Arista 这两年横空出世，彻底的取代了 Force10 在作者心中“傻快”的霸者地位，预计其在今后几年的数据中心市场中势必会有所斩获，天下武功，唯快不破。

新技术是层出不穷的，前文介绍的这些技术只是作者知道的内容，眼界有限，相信还有更多更先进的技术在此章介绍之外。期望通过本文能够对读者学习其他技术也有所帮助，万事万物有果必有因，透过纷乱的报头字段、状态机变化和报文交互等协议机制设计表象，去了解清楚技术产生的背景原因和要解决的问题，势必在学习时可以达到事半功半的效果。

用以下言语与技术同好共勉：在技术之路上，了解得越多，敬畏之心越重，但仍需不断前行，即使无法成为引领者，也必将超越原地踏步者。

## 6 终章

完了就是完了，其实没啥好多说的，想说的要说的该说的前面都已经说过了。文中做了不少预测，根据作者的恶趣味，最后在这里对那些神棍内容再总结一下凑凑字数。

在未来 5-10 年间作者认为：

市场方面

1、云计算市场出现不了如 Microsoft 于操作系统、Google 于网络搜索或 Cisco 于数据通信的一家独大局面，但对多虚一的集中云与一虚多的分散云，市场划分会更加清晰，客户抗忽悠能力也将得到大幅度提升。

2、解决了安全问题后，基于服务的 SaaS 会占据更多的业务租用市场，中小企业自身 IT 资源消耗进而降低，业务能力反而提升。例如同时租用 Google 云的数据管理，Amazon 云的人力资源，Microsoft 云的 ERP，Cisco 云的统一通信和作者云的客户关系管理等系统来综合搭建企业 IT 平台，会成为很时髦很常见的思路。（想创业的抓紧，SaaS 机会贼多的，而且初始投入规模并不需要太大）

3、提供云服务（以 SaaS 为主）的产业将如雨后春笋般出现，这些服务提供商将会搭建大量的数据中心为客户提供云租用服务，也会成为网络设备厂商们的衣食父母。需求较小的企业都去直接租用服务了，因此数据中心步入了大型与巨型为主的时代，动辄成千上万的服务器节点绝对是小 Case。数据中心产品的销售也将随之进入规模化采购阶段，搞定几个大客户，厂商一年下来就吃穿不愁了。

## 技术方面

1、VM 之间互通技术之争会以硬件交换机进入服务器内部为最终结局，有可能是在网卡上实现，也有可能直接在主板上加转发芯片。毕竟从现在发展情况来看，芯片价格会越来越便宜，集成度会越来越高。

2、存储方面 FCoE 基于 Ethernet 带宽发展方面的优势，必将取代 FC，当然过程会比较漫长的，估计 10 年之后 FC 也还能占有一定的空间。

3、数据中心站点内部 TRILL 将会一统天下。巨型数据中心内，基于 IP 层面的交互会导致传输效率降低和部署复杂度提升，因此仍然会以 Ethernet 技术为主，而 TRILL 是目前看得到的最有希望胜出的公共标准。各个厂商的私有技术会将在规模稍小一些的大中型数据中心内有所应用，比如前面说的 SaaS 创业企业，其可能会更看重网络的高可靠性、高性能、易管理和易维护等私有技术强项的地方。而且网络规模较小，搞一家厂商的设备就差不多了，不需要考虑互通。

4、数据中心跨站点二层互联方面，RPR 由于是公共标准可以成为种子选手，但其成长空间目前并不充分，还要看技术发展演进和各个厂家的态度。当然如果有哪家厂商愿意把自己的私有技术拿出来推成标准，也还是很有希望在市场上占据高点的。

5、在多站点选路方面，应该会有些新的技术标准出来，DNS 方案一统天下的局面不会长久。这块谁都有机会，就看投入与机遇了。

## 7 感言

沥沥拉拉写了小两个月，长度和时间都远远超出了最初的计划，也耗费了不少的热情和精力，以后是不敢随便写这种大文章了。但整个写作过程对作者来说受益匪浅，不断总结是自我提升的重要动力。后面休息休息还会再整理一些关于云计算数据中心安全、存储、性能和可靠性等方面的外篇，先在这里立个目标好做自我督促。

套用一些书中常看到的话，谨以此文献给我的家人朋友和同事，并纪念作者步入而立之年。顺便感谢每一位能从头读到这儿的读者，你们的存在是我写作乐趣的源泉。

## 8 文章原目录（来源 51CTO）

### 【内容导航】

- |  |  |
|--|--|
| 第 1 页: <a href="#">云计算</a>                     | 第 2 页: <a href="#">集中云</a>                       |
| 第 3 页: <a href="#">分散云</a>                     | 第 4 页: <a href="#">Bare-Metal 方案</a>             |
| 第 5 页: <a href="#">vMotion</a>                 | 第 6 页: <a href="#">云计算小结</a>                     |
| 第 7 页: <a href="#">Client 与 Server</a>         | 第 8 页: <a href="#">层次化与扁平化</a>                   |
| 第 9 页: <a href="#">三层结构与两层结构</a>               | 第 10 页: <a href="#">Server 与 Storage</a>         |
| 第 11 页: <a href="#">数据中心多站点</a>                | 第 12 页: <a href="#">多站点选择</a>                    |
| 第 13 页: <a href="#">数据中心小结</a>                 | 第 14 页: <a href="#">EOR 与 TOR</a>                |
| 第 15 页: <a href="#">控制平面与转发平面</a>              | 第 16 页: <a href="#">48GE+4*10GE 交换机设计</a>        |
| 第 17 页: <a href="#">Switch Fabric</a>          | 第 18 页: <a href="#">Chassis 与分布式转发</a>           |
| 第 19 页: <a href="#">Chassis 转发能力</a>           | 第 20 页: <a href="#">框式交换机架构</a>                  |
| 第 21 页: <a href="#">H3C 的 12500</a>            | 第 22 页: <a href="#">Clos 与 VOQ</a>               |
| 第 23 页: <a href="#">Virtual Output Queues</a>  | 第 24 页: <a href="#">网络小结</a>                     |
| 第 25 页: <a href="#">技术结构</a>                   | 第 26 页: <a href="#">网络虚拟化</a>                    |
| 第 27 页: <a href="#">数据平面虚拟化</a>                | 第 28 页: <a href="#">网络一虚多技术</a>                  |
| 第 29 页: <a href="#">技术理解</a>                   | 第 30 页: <a href="#">VM 本地互访网络技术</a>              |
| 第 31 页: <a href="#">Cisco Nexus1000V</a>       | 第 32 页: <a href="#">Nexus5000+Nexus2000</a>      |
| 第 33 页: <a href="#">VN-Tag 格式与封装位置</a>         | 第 34 页: <a href="#">Unified Computing System</a> |
| 第 35 页: <a href="#">UCS 系统结构</a>               | 第 36 页: <a href="#">UCS VN-Link</a>              |
| 第 37 页: <a href="#">802.1Qbg EVB</a>           | 第 38 页: <a href="#">802.1Qbh 接口扩展</a>            |
| 第 39 页: <a href="#">Multichannel</a>           | 第 40 页: <a href="#">Multichannel 转发过程</a>        |
| 第 41 页: <a href="#">Multichannel 对比 VN-Tag</a> | 第 42 页: <a href="#">Remote Replication 复制</a>    |
| 第 43 页: <a href="#">S-VLAN 组件本地复制</a>          | 第 44 页: <a href="#">802.1Qbh 和 802.1Qbg</a>      |
| 第 45 页: <a href="#">Ethernet 与 FC 网络融合技术</a>   | 第 46 页: <a href="#">FC 网络三种主要接口</a>              |
| 第 47 页: <a href="#">FC IVR Zone 概念</a>         | 第 48 页: <a href="#">FC 设备进行两步注册</a>              |
| 第 49 页: <a href="#">FCoE 技术详解</a>              | 第 50 页: <a href="#">FCoE 初始化连接</a>               |



- 第 51 页: [Fibre Channel Forwarder](#)
- 第 53 页: [802.1Q DCB](#)
- 第 55 页: [NPort ID Virtualization](#)
- 第 57 页: [FCoE、FCF 与 NPV](#)
- 第 59 页: [控制平面多虚一技术](#)
- 第 61 页: [virtual Port-Channel](#)
- 第 63 页: [FabricPath 技术](#)
- 第 65 页: [P802.1aq SPB](#)
- 第 67 页: [Juniper QFabric 网络技术](#)
- 第 69 页: [控制平面一虚多技术](#)
- 第 71 页: [数据中心跨站点二层网络](#)
- 第 73 页: [Virtual Private Lan Service](#)
- 第 75 页: [IP 核心网](#)
- 第 77 页: [OTV 技术的机制处理](#)
- 第 79 页: [HSRP 隔离](#)
- 第 81 页: [数据中心多站点选择](#)
- 第 83 页: [LISP IPinIP 协议](#)
- 第 85 页: [LISP 中各角色之间的连接](#)
- 第 87 页: [如何切换 EID 对应的 Locator IP](#)
- 第 52 页: [FCF 拥有自己的 MAC](#)
- 第 54 页: [FCoE 网卡 CNA](#)
- 第 56 页: [NPV 与服务器之间的网络](#)
- 第 58 页: [跨核心层服务器二层互访](#)
- 第 60 页: [Cisco VSS 三种故障检测方式](#)
- 第 62 页: [数据平面多虚一技术](#)
- 第 64 页: [Forwarding TAG 技术](#)
- 第 66 页: [SPB 网络定义软件算法](#)
- 第 68 页: [云计算数据中心未来](#)
- 第 70 页: [VLAN 和 VPN 一虚多技术](#)
- 第 72 页: [MPLS 核心网](#)
- 第 74 页: [Advanced VPLS 技术](#)
- 第 76 页: [Cisco OTV 私有技术](#)
- 第 78 页: [站点 ED 双机冗余](#)
- 第 80 页: [数据中心跨站点二层互联](#)
- 第 82 页: [全局负载均衡技术](#)
- 第 84 页: [LISP 的技术结构](#)
- 第 86 页: [数据中心多站点选路](#)
- 第 88 页: [技术总结](#)