

程 (traffic engineering) [RFC 3346; RFC3272; RFC 2702; Xiao 2000], 其中网络运行者能够超越普通的 IP 路由选择, 迫使某些流量沿着一条路径朝着某给定的目的地引导, 并且朝着相同目的地的其他流量沿着另一条路径流动 (无论是由于策略、性能或某些其他原因)。

将 MPLS 用于其他目的也是可能的。能用于执行 MPLS 转发路径的快速恢复, 例如, 经过一条预计算的无故障路径重路由由流量来对链路故障作出反应 [Kar 2000; Huang 2002; RFC 3469]。最后, 我们注意到 MPLS 能够并且已经被用于实现所谓虚拟专用网 (Virtual Private Network, VPN)。在为用户实现一个 VPNR 的过程中, ISP 使用它的 MPLS 使能网络将用户的各种网络连接在一起。MPLS 能被用于将资源和由用户的 VPN 使用的寻址方式相隔离, 其他用户利用该 VPN 跨越该 ISP 网络, 详情参见 [DeClercq 2002]。

这里有关 MPLS 的讨论是简要的, 我们鼓励读者参阅我们提到的这些文献。我们注意到对 MPLS 有许多可能的用途, 看起来它将迅速成为因特网流量工程的瑞士军刀!

## 5.6 数据中心网络

近年来, 因特网公司如谷歌、微软、脸谱 (Facebook) 和亚马逊 (以及它们在亚洲和欧洲的同行) 已经构建了大量的数据中心。每个数据中心都容纳了数万至数十万台主机, 并且同时支持着很多不同的云应用 (例如搜索、电子邮件、社交网络和电子商务)。每个数据中心都有自己的**数据中心网络** (data center network), 这些数据中心网络将其内部主机彼此互联并与因特网中的数据中心互联。在本节中, 我们简要介绍用于云应用的数据中心网络。

大型数据中心的投资巨大, 一个有 100 000 台主机的数据中心每个月的费用超过 1200 万美元 [Greenberg 2009a]。在该费用中, 用于主机自身的开销占 45% (每 3~4 年需要更新一次); 变压器、不间断电源系统、长时间断电时使用的发电机以及冷却系统等基础设施的开销占 25%; 用于功耗的电力设施的开销占 15%; 用于联网的开销占 15%, 这包括了网络设备 (交换机、路由器和负载均衡设备)、外部链路以及传输流量的开销。(在这些比例中, 设备费用是分期偿还的, 因此费用通常是由一次性购买和持续开销 (如能耗) 构成的。) 尽管联网不是最大的费用, 但是网络创新是减少整体成本和性能最大化的关键 [Greenberg 2009a]。

主机就像是数据中心的工蜂: 它们负责提供内容 (例如, 网页和视频), 存储邮件和文档, 并共同执行大规模分布式计算 (例如, 为搜索引擎提供分布式索引计算)。数据中心中的主机称为**刀片** (blade), 与比萨饼盒类似, 一般是包括 CPU、内存和磁盘存储的商用主机。主机被堆叠在机架上, 每个机架一般堆放 20~40 台刀片。在每一个机架顶部有一台交换机, 这台交换机被形象地称为**机架顶部** (Top of Rack, TOR) **交换机**, 它们与机架上的主机互联, 并与数据中心中的其他交换机互联。具体来说, 机架上的每台主机都有一块与 TOR 交换机连接的网卡, 每台 TOR 交换机有额外的端口能够与其他 TOR 交换机连接。尽管目前主机通常有 1Gbps 的以太网与其 TOR 交换机连接, 但 10Gbps 的连接也许成为标准。每台主机也会分配一个自己的数据中心内部的 IP 地址。

数据中心网络支持两种类型的流量: 在外部客户与内部主机之间流动的流量, 以及内部主机之间流动的流量。为了处理外部客户与内部主机之间流动的流量, 数据中心网络包括了一台或者多台**边界路由器** (border router), 它们将数据中心网络与公共因特网相连。数据中心网络因此需要将所有机架彼此互联, 并将机架与边界路由器连接。图 5-30 显示了一个数据中心网络的例子。**数据中心网络设计** (data center network design) 是互联网络和协议设计的