# EWFQ:一种新的高速网络分组调度算法\*<sup>?</sup>

EWFQ: A Novel Packer Scheduling Algorithms in High Speed Networks

# 任立勇 卢显良

(电子科技大学计算机学院 成都610054)

Abstract Packet scheduling algorithm is one of crucial technologies of routers in high speed networks. In this paper we first discuss the limitation of some existing packet scheduling algorithms, then show the quantitative relationships between the GPS system and its corresponding packet WFQ system. A novel packet scheduling algorithm is proposed. It is proven to have following properties: (1) it ensures fair allocation of bandwidth among all sessions; (2) it provides deterministic delay upper bounds to a session whose traffic is constrained by a leaky bucket; (3) it has a relatively low asymptotic complexity of O(logN); (4) it has a relatively low Worst-case Fair Index (WFI). So it can be deployed in routers of high speed networks.

Keywords Fair queueing, Scheduling algorithms, Quality of service, High speed networks

## 1 引官

宽带综合业务网要求能给不同的应用提供不同的服务质量(QoS),其中分组调度算法作为网络路由器中的一个重要组件起着相当关键的作用,传统的Internet 是基于尽力而为(best-effort)模型实现的,该模型采取先来先服务(FCFS)的分组调度算法,这种模型具有实现简单的特点,它在假定所有应用互相协作的情况下工作得非常好。但当网络发生拥塞时,实时应用的服务质量往往得不到保证。同时,连接间的隔离性能也非常差,吞吐量大的连接得到更多的服务,某些不良行为的连接可能造成其他连接的服务质量急剧下降。

A. K. Parekh 等提出的广义处理器共享(GPS)<sup>[1]</sup> 能较好地解决上述问题:1)当在数据源端实施漏桶算 法的流量整形时,GPS 能提供端到端的延迟界限,2) GPS 能为有分组积压的连接提供公平的带宽分配。由于 GPS 是基于流体模型的理想化调度算法,不能用于实际系统,于是各种各样的基于分组调度的近似 GPS 算法相继提出<sup>[2~5]</sup>,这些算法均是在连接带宽公平性分配与计算复杂度间作出平衡,如 WFQ 尽管提供了可与 GPS 相当的性能,但由于其计算复杂度较高(O(N))而不能在高速网络中运行,SCFQ 算法重新定义了系统虚拟时间计算函数。将计算复杂度降低为 O(I),但不足之处在于当网络连接增多时,连接的隔离和保护性能变差,同时 SCFQ 提供的端到端延迟较大。

为此,本文在研究 WFQ 与 SCFQ 的工作机制与不足的基础上,提出了一种新的分组调度算法 EWFQ (Extended Weighted Fair Queueing),理论分析与实验证明 EWFQ 不仅具有计算复杂度低的特点,同时也更好地接近了 GPS 性能水平。

\*)本文得到国家九五重点攻关项目基金资助。信息产业部生产发展基金资助。任立勇 博士研究生。主要研究方向为网络资源管理、网络协议等。卢显良 教授,博士生导师,主要研究方向为操作系统与网络应用技术、

由器的体系结构正朝着速度更快、服务质量更好和更 易于综合化管理三个方向发展。

#### 参考文献

- 1 Asthana A.Delhp S. Jagdish H. et al. Towards a gigabit IP router. Journal of High-Speed Networks, 1992, 1(4):281 ~288
- 2 Partidge C. Carvey P. Burgess E. et al. A 50Gb/s IP router. IEEE Trans. on Networking, 1998.6(3):237~

248

- 3 Chen J S. Guerin R. Performance study of an input queueing packet switch witch with two priority classes. IEEE Trans. Commun[J].1991.39(1):117~126
- 4 Mckeown N. Fast Switched Backplane for a Gigabit Switched Router. White Paper, http://www.cosco.com
- 5 Semeria C. Internet Backbone Routers and Evolving Internet Design. White Paper http://www.jumper.com
- 5 李津生,等编著。下一代 Internet 网络技术, 人民邮电出版社, 2001

## 2 公平排队调度算法

公平排队调度算法是近似 GPS(Generalized Processor Sharing)服务器[日的分组调度算法。GPS 服务器持续工作(Work-Conserving)并以固定的速率 r 发送数据,它是一个基于流体模型的理想化调度算法,即 GPS 服务器假定分组无限可分并能同时为多个连接服务。一个有 N 个队列的 GPS 系统表示为 N 个正实数  $(x_1,x_2)$ 为连接。在时间间隔 $(t_1,t_2)$ 得到的服务量、 $W(t_1,t_2)$ 为服务器总的服务量,则 GPS 服务器可定义为。

$$\frac{W_{i}(t_{1}, t_{2})}{\varphi_{i}} = \frac{W(t_{1}, t_{2})}{\sum_{i \in B(t_{1}, t_{2})}^{N} \varphi_{i}} \ge \frac{W(t_{1}, t_{2})}{\sum_{i=1}^{N} \varphi_{i}}$$
(1)

其中  $B(t_1,t_2)$ 为时间间隔 $(t_1,t_2)$ 内所有有分组积压的连接的集合。令 g.=r\*g.由(1)可推导出  $W.(t_1,t_2)$   $\geqslant$   $g.*(t_2-t_1)$ (假定  $\sum_i g_i=1$ ),则连接 i 在有分组积压时的实际发送速率  $r.\geqslant_{g.}$ 。可见、GPS 服务器可保证给连接 i 分配最小的服务速率 g.=r\*g.。因此,通过调整连接权值 g. GPS 服务器可灵活地为连接分配带宽。

GPS 良好的工作特性促使许多研究者研究各种近似 GPS 的分组调度算法,其中加权公平排队(WFQ)被认为是最理想的近似GPS 算法。

定理1 GPS 调度算法与相应的分组调度算法 WFQ(或 PGRS)有如下关系:

(1)设 d'wra, d', cps 分别为 WFQ 与 GPS 系统中连接 i 的第 k 个分组的离开时间,则

$$d_{i,w_{PQ}}^{t} - d_{i,GPS}^{t} \leqslant \frac{L_{\max}}{r} \tag{2}$$

(2) 设 W<sub>1,GPS</sub> (0,t), W<sub>1,WFQ</sub> (0,t) 分别为 GPS 与 WFQ 系统中第 t 个连接在时间间隔[0,t)内的服务量, 测

$$W_{t,GPS}(0,t) - W_{t,WFQ}(0,t) \leq L_{max}$$
 (3)

(3)当在信源端实施参数为(σ, r,)的漏桶控制时, WFQ 系统中连接 i 的延迟界限为:

$$\frac{\sigma_{t}}{r_{t}} + \frac{L_{max}}{r} \tag{4}$$

(4)令 WFQ 系统中有 N 个连接,连接 1 的速率为  $r_1$ ,则它的最坏情况公平指数为:

$$C_{t,weq} = N \cdot \frac{r_t}{r} \cdot \frac{L_{war}}{r}$$
 (5)  
文[1]给出了定理1的证明。

WFQ 系统在实现时,设置了一个系统虚拟时间函数 V(t),每个分组到达时计算其虚拟发送时间 S:与虚拟完成时间 F,在选择下一个发送分组时,WFQ 系统采用的是 SFF(Smallest Finish First)服务机制。尽管 WFQ 调度算法提供了与 GPS 相当的特性、但它存在两个致命的弱点:(1)在最坏情况下,服务器的 N个连接全部活跃,此时该算法的计算复杂度为 O(N),因此,几乎不可能在高速主干网中的路由器中配置 WFQ 调度算法;(2) WFQ 服务器的最坏情况公平指数(WFI,Worst-case Fair Index)与系统中的连接数 N 成正比,因此造成时延抖动增加。

为克服 WFQ 的上述缺点,S. Golestan: 提出了另 一种算法:SCFQ(Self-Clock Fair Queueing)[3.,SCFQ 算法取当前发送分组的完成时间为系统虚拟时间 V (t),这使算法的计算复杂度降为 O(l),但 SCFQ 算法 不足之处在于当网络连接增多时,连接的隔离和保护 特性将变差,同时 SCFQ 算法的延迟(Delay)与延迟抖 动(Delay-jitter)要比 WFQ 大得多",这主要是因为 SCFQ 在计算系统虚拟时间时的不准确性产生的。 J. C. R. Bennett 提出了一种叫 WF<sup>2</sup>Q(Worst-case Fair Weighted Fair Queueing)的算法[5],由于 WF2Q 在选 择下一个发送分组时采用的是 SEFF(Smallest Eligible Finish First)服务机制,因此其最坏情况公平指数 与服务器中的连接数无关,同时,在任一时间段[0,t) 内、WF3Q 服务器提供给连接;的服务量与相应 GPS 服务器提供给连接:的服务量相比不会超过该连接的 最大分组长度的一个百分比,即 $W_{i,u,p^2Q}(0,t)-W_{i,cp}$ ,  $(0,t) \leq (1-\frac{r_t}{\mu})L_{t,\max}$ ,但遗憾的是,WF'Q 与 WFQ 具 有相同的时间复杂度(O(N))。

#### 3 改进的分组调度算法

如第2节所述、WFQ 与 WF<sup>2</sup>Q 算法虽然具有较为满意的分组调度特性、公平性、确定的分组延迟上界、小的 WFI(只针对 WF<sup>2</sup>Q),但由于它们采用相同的系统虚拟时间计算函数  $V_{CPS}(t)$ ,因此具有相同计算复杂度  $O(N)^{[1]}$ ,很难在高速网络中实现。为此,我们设计并实现了一种新的分组调度算法 EWFQ(Extended Weight Fair Queueing),EWFQ 采用与 WF<sup>2</sup>Q 相同的调度机制 SEFF,但重新定义了系统虚拟时间函数,将计算复杂度降为  $O(\log N)$ ,该函数具体如下:

$$V(t) = \max\{V(\tau) + W(\tau, t), \min_{\tau \in \mathbb{R}^n} \{S_{\tau}(\tau)\}\}$$
 (6)

<sup>:</sup>D当连接 ι 在信源实施漏桶控制 (σ, ι, r, )时 ι WFQ 提供的端到端延迟及延迟抖动分别为  $\frac{\sigma_i + nL_{\max}}{r_i} + \Sigma_{j=1}^{m} \frac{L_{\max}}{C_j}$  ,  $\frac{\sigma_i + nL_{\max}}{r_i}$  ,  $\frac{\sigma_i + nL_{\max}}{C_i}$  ,  $\frac{\sigma_i + nL_{\max}}{C_i} + \Sigma_{j=1}^{m} (K_j - 1) \frac{L_{\max}}{C_j}$  . 其中 n 为从源到目的跳数 n , n 分别为第 n 个路由器的链路速率与共享该链路的连接数。

式中 r≤t,W(r,t)为时间间隔(r,t)内系统的服务量,B (r)为时刻 r 时有分组积压的连接的集合 ,S<sub>c</sub>(r)为连接 2的第一个分组的虚拟开始时间。

如图1所示,EWFQ 算法的具体实现可分为两部 分:11当有新的分组到达时,ARRIVE(x,k,P)函数计 算该分组的虚拟开始时间与虚拟完成时间。并给分组 打上时间标记,然后将该分组送入相应连接的队列;3) 分组发送完成后、FINISH(1.k.P)函数重新计算系统 虚拟时间,然后选择下一个合格的具有最小完成时间 (SEFF)分组进行发送。

定理2 EWFQ的计算复杂度为 O(log N)。

证明·由图1可知,EWFQ的计算时间主要花费在 两部分:计算系统虚拟时间 V(t)与选择下一个最早完 成的分组进行发送。其中第一部分需要对所有有分组 积压连接的队头分组的虚拟开始时间进行排序以选择 在 GPS 系统中最早开始服务的分组(即合格分组),第 二部分需要对所有合格分组按虚拟完成时间排序以选 择在实际系统中最早服务分组。最坏情况下,N个连 接全部活跃,此时两部分的计算复杂度均为 () (logN)[8],因此 EWFQ 的计算复杂度为上述两部分计 算复杂度之和,即 $O(\log N)$ 。

> $ARRIVE(\iota,k,p)$  $1 \ IF \ Q_i(t) = 0$ 2 THEN  $S_i^k = \max\{F_i^{(k+1)}, V(t)\}$ ELSE  $S_i^i = F_i^{i-1}$ 6  $F_{i}^{k} = S_{i}^{k} + L_{n}/r_{i}$ 7 Stamp(P.St.F!) B enqueue(Q.P) 1 dequeue(Q,P) 2  $(emp \leftarrow mun_{i \in B}(a_i^*)(S_i)$ 3  $V(t) \leftarrow V(a_t^2) + L_t$  $4|V(t) \leftarrow max\{V(t), temp\}$ 5 select(b)

#### 图1 改进的加权公平排队(EWFQ)

定理3 EWFQ是持续工作的(Work-Conserving),

证明:由(5)式可以看出,系统虚拟时间至少等于 所有队头分组中有最小虚拟开始时间分组的开始时 间,反过来讲就是在任意时刻t时,只要系统中有分组 积压,就至少有一个分组,其虚拟开始时间小于或等于 系统虚拟时间,即至少有一个合格分组,这就保证了以 SEFF 为调度机制的 EWFQ 是持续工作的。

定理4 当在信源端实施参数为(a, r,)的漏桶控 制时,EWFQ系统中连接,的延迟界限为:

$$\frac{\sigma_r}{r_s} + \frac{L_{\text{max}}}{r}$$

由于本定理的证明相当复杂,因此这里只给出相 关说明,EWFQ采用了SEFF服务机制,运用文[7]中

相同的方法,可以证明 EWFQ 实现的系统为分组速率 比例服务器,同时,文气气证明了对于具有 SEFF 调度 机制的分组速率比例服务器,如果在连接,的信源端 实施参数为(σ, ε,)的漏桶控制时,连接:的延迟界限 为 $\frac{\sigma_i}{T} + \frac{L_{\text{max}}}{T}$ .

为精确量化分组调度算法与流体 GPS 的差别,文 [5]中给出了一个叫最坏情况公平指数(WFI, Worstcase Fair Index)的定义。

定义1 任意时间上时,若连接上的分组在服务器 s中的延迟小于等于 $\frac{1}{r}Q_{\alpha}(t)+C_{\alpha}$ ,则称服务器。对 连接,是最坏情况公平的,即

$$d_{i,j}^{k} \leqslant a_{i}^{k} + \frac{Q_{i,j}(a_{i}^{k})}{r_{i}} + C_{i,j}$$
 (8)

式中, d, d, 分别表示连接: 的第 & 个分组到达时间与 窗开时间,Q,(a)表示连接,的第点下到达时,连接, 在服务器中分组的队列长度,, 为连接,的服务速率, C.,,为连接:的最坏情况公平指数 WFI.

从定义中可以看出为保证服务器对每个连接都是 最坏情况公平的,其WFI必须尽量小,并与服务器中 连接数无关。文[5]中证明了 Cure = 0.同时举例证明了  $C_{1,WFQ} = N^{\frac{L_{max}}{n}},$ 

定理5 EWFQ 系统中连接:的最坏情况公平指

数 WFI 为 
$$c_{r,EWEQ} = \frac{L_{t,max}}{r_t} - \frac{L_{t,max}}{r} + \frac{L_{max}}{r}$$
,即,
$$d_{t,EWEQ}^k = a_t^k \leqslant \frac{Q_{t,EWEQ}(a_t^k)}{r_t} + \frac{L_{t,max}}{r_t} - \frac{L_{t,max}}{r} - \frac{L_{max}}{r} \quad \forall t. k$$
(9)

式中:ai.di, ewea分别表示连接,的第 k 个分组到达时 间和离开时间、Queusa为时刻或时连接。的队列长度、  $L_{\text{t.max}}$ ,  $L_{\text{max}}$ 分别表示连接:与系统中最大分组长度。

证明 由文[5]可知,WF'Q系统具有以下特件。

$$d_{r,WF^2Q}^k - d_{r,GFS}^k \leqslant \frac{L_{max}}{r}$$

$$W_{\text{total}}(0,t) = W_{\text{total}}(0,t) \leqslant (1 + \frac{r_t}{r}) L_{\text{total}}$$
 (11)

由于 EWFQ 采用了 SEFF(Smallest Eligible Finish First)调度机制。因此可采用与文[5]中定理1相同 的方法证明 EWFQ 也具有上述特性,即。

$$d_{1,\text{EWFQ}}^{4} - d_{1,\text{GPS}}^{4} \leqslant \frac{L_{\text{max}}}{r} \tag{12}$$

$$W_{i,SWFQ}(0,t) - W_{i,GPS}(0,t) \le (1 - \frac{r_i}{r}) L_{i,mix}$$
 (13)

由定理2可知,EWFQ是持续工作的,因此由(13) 式立即可得

$$Q_{cGFS}(\tau) - Q_{cFM,PQ}(\tau) \leq (1 - \frac{r_c}{\tau}) L_{c,max}$$
 (14)

同时由 WFI 的定义可知、GPS 的 WFI 为u,即

$$d_{r,GFS}^{*} - a_{i}^{*} \leqslant \frac{Q_{r,GFS}(a_{i}^{*})}{r}$$
 (15)

由(12)式可得

$$(d_{i,EVEQ}^{b} - u_{i}^{b}) - (d_{i,DES}^{b} - u_{i}^{b}) \leqslant \frac{L_{\max}}{r}$$

$$(16)$$

$$d_{i,EWEQ}^{b} - u_{i}^{b} \leqslant (d_{i,DES}^{b} - u_{i}^{b}) + \frac{L_{\max}}{r} \leqslant \frac{Q_{i,DES}(u_{i}^{b})}{r} + \frac{L_{\max}}{r}$$

$$\leqslant \frac{Q_{i,EWEQ}(u_{i}^{b}) + (1 - \frac{r_{i}}{r})L_{i,\max}}{r} - \frac{L_{\max}}{r}$$

$$= \frac{Q_{i,EWEQ}(u_{i}^{b})}{r} + \frac{L_{i,\max}}{r} - \frac{L_{\max}}{r} + \frac{L_{\max}}{r}$$

$$\Leftrightarrow \frac{Q_{i,EWEQ}(u_{i}^{b})}{r} + \frac{L_{i,\max}}{r} - \frac{L_{\max}}{r} + \frac{L_{\max}}{r}$$

$$\Leftrightarrow \frac{Q_{i,EWEQ}(u_{i}^{b})}{r} + \frac{L_{i,\max}}{r} + \frac{L_{\max}}{r} + \frac{L_{\max}}{r}$$

由于 EWFQ 是在 WFQ 算法的基础上设计而成,因此它除了具有上述特件之外,同时还可参照文[5]相同的方法证明 EWFQ 在带宽分配方面具有公平特性,各连接间的隔离特性也相当好。

# 4 仿真实验

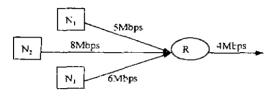
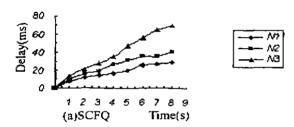
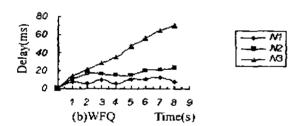


图2 模拟实验网络模型





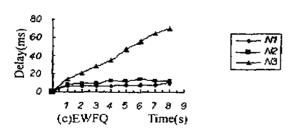
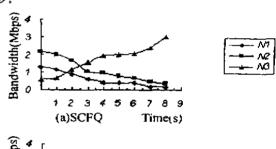
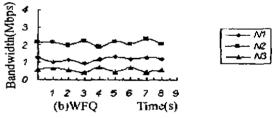


图3 三种算法的延迟比较 下面通过仿真实验进一步说明 EWFQ 算法的性

能。为简单起见,我们采用图2所示网络模型,设该仿真系统有3个连接共享输出链路,其中 N<sub>1</sub>,N<sub>2</sub>发出的信源符合 ON-OFF 模型,两种状态持续时间分别服从负指数分布,这两个连接信源流都经过漏桶算法整形,整形后流量符合(a,p)模型,假定 N<sub>2</sub>发出的信源为不良行为流,实验中,该流持续以超过其分配带宽的速率发送数据。如图3、图4分别为路由器采用 SCFQ、WFQ 与EWFQ 算法时各连接的时延及实际占用带宽比较。

如图3所示,仿真实验表明,由于 SCFQ 算法的隔离性和保护性较差,尽管 N<sub>1</sub>,N<sub>2</sub>发出的流经过漏桶算法整形,但由于 N<sub>2</sub>发出的流来经任何限制,持续以超过其分配带宽的速率发送数据,从而导致 N<sub>1</sub>,N<sub>2</sub>流延迟特性变差。同时,尽管 WFQ 系统中漏桶算法整形后的连接有确定的时延上界(如(4)式),但由于 WFQ 算法计算复杂度较高,WFI 也比较大,因此其时延特性虽比 SCFQ 算法好,但时延抖动相对较大。比 EWFQ 算法提供的时延特性差。同样道理,如图4所示,随着时间的变化,N<sub>2</sub>发出的流将占用 SCFQ 系统的绝大部分带宽,而 WFQ 算法与 EWFQ 算法由于采用相同的带宽分配策略,因此在带宽分配方面就公平得多,平滑得多。





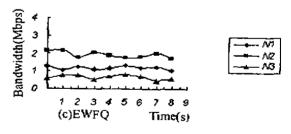


图4 三种算法的带宽比较

结束语 本文提出了一种新的分组调度算法

# 基于宽带的无线接入技术 LMDS

Wireless Access Technology Based on Broadband

## 马燕

(重庆师范学院物理学与信息技术系 重庆400047)

LMDS (Local Multipoint Distribution Service) is a broadband wireless point-to-multipoint communication system operating above 20 GHz. It occupies the spectrum frequency at around 28 GHz and 31 GHz. It can handle telephony and video programming as well as data services such as Internet access. A big advantage of LMDS is that cable or copper lines for access to the home are not needed. LMDS provides an effective last-mile solution for the incumbent service provider and can be used by competitive service providers to deliver services directly to end users. In this paper. The author analyses the performance, working principle, structure of system of LMDS and its key technology. Also, the paper discusses its foreground-

Keywords LMDS, Base station, User station, Protocol

随着信息化社会的到来,作为连接全球信息的 Internet 的业务量日益巨增,正呈现出爆炸性的增长,因 此近年各国都高度重视 Internet 骨干网络的建设、核 心网络实现了光纤化,其带宽基本可满足当前对信息 获取的需求。而网络的瓶颈越来越集中体现在接人网 方面、即用户与核心网络的连接部分。

从理论上说,全光纤接人网络将是比较完美的解 决方案,但实现上面临很多的困难,即使在发达国家也 还远远未能实现,因此,人们在网络的接人技术方面做 了大量的研究,现已提出了多种解决方案,如 cablemodem,即广电的宽带多媒体传送系统,它利用了有线 电视系统可用频谱中的一小部分来传送数字信号: DSL,即"数字用户线路",其最具代表性的 ADSL 技术 是利用现有的电话双绞铜线对用户提供高速 Internet 接入等业务, ISDN 技术, 是传送话音、数据、图像等综

EWFQ,理论分析与实验证明该算法不仅具有连接独 立的特性,能更公平,灵活地为连接分配带宽,而且它 的计算复杂度比 WFQ 算法低得多。同时、EWFQ 还具 有较好的时延特性和较低的 WFI,可以为连接提供确 定的时延上界和稳定的时延抖动。因此,该算法能更好 地满足高速网络的要求。

#### 参考文献

- i Parekh A K. Gallager R G. A generalized processor sharing approach to flow control in integrated services networks the single-node IEEE/ACM Trans. Networking. 1993.1(3) 344~357
  2 Demers A. Keshav S. Shenker S. Analysis and simulation

台业务的主要手段之一;还有其它如 HFC、PON(无源 光网络)、APONC(基于 ATNK 的无源光网络)、SDH 等方式的光纤宽带接入等,这些接入技术都有各自的 特点和针对性。

另一方面,无线接入技术特别是宽带无线技术正 迅速发展。运用该技术可以将数据、Internet、话音、视 频和多媒体应用传送到商业和家庭用户。无线网的组 建方便快捷、对人群和环境影响最少,而且无需巨额基 础设施和场地投资,作为具有较强竞争性的宽带无线 接入技术具有更美好的前景。特别是近年来发展起来 的 LMDS 技术,是无线宽带技术的一个典范。

#### 一、LMDS 综述

LMDS(Local Multipoint Distribution Service)是 近年来逐渐发展起来的一种工作于24GH2~38GH2 频

- of a fair queueing algorithm. In Proc. ACM SIGC OMM'89,  $3{\sim}12$
- Golestant S. A self-clock fair queueing scheme for broad-band applications. Proc. IEEE INFOCOM'94, Toronto. CA, 1994,4.636~646
- CA, 1994,4, 636~646
  4 王宏宇, 顾冠群. 集成服务网络中的分组调度算法研究综选. 计算机学报,1999,22(10):1090~1099
  5 Bennett J C R, Zhang H. WF<sup>2</sup>Q: Worst-case fair weighted lair queueing. IEEE INFOCOM'96, San Francisco. CA, 1996,3,120~128
- Stoica I. Wahab H A. Earliest eligible virtual deadline first. A flexible and accurate mechanism for proportional share resource allocation [Tech. Rep. TR-95-22]. Old Dominion Univ., 1995, 11
- Stiliadis D. Varma A. Rate-proportional servermethodology for fair queueing algorithms IEEE/ACM Trans. Networking, 1998.6(2):164~174