

Edge-to-edge Control: Congestion Avoidance and Service Differentiation for the Internet

David Harrison

Rensselaer Polytechnic Institute
 harrisod@cs.rpi.edu
<http://networks.ecse.rpi.edu/~harrisod>

Single Vs. Multi-Provider Solutions

- ❑ ATM and frame relay operate on single datalink layer.
 - ❑ All intermediate providers must agree on a common infrastructure. Requires upgrades throughout the network. *Coordination* to eliminate *heterogeneity*.
 - ❑ Or operate at lowest common denominator.
- ❑ Overprovision:
 - ❑ Operate at single digit utilization.
 - ❑ More bandwidth than sum of access points. 1700 DSL (at 1.5 Mbps) or 60 T3 (at 45 Mbps) DDoS swamps an OC-48 (2.4 Gbps).
 - ❑ Peering points often last upgraded in each upgrade cycle. Performance between MY customers more important.
 - ❑ *Hard for multi-provider scenarios.*

Outline

- ❑ QoS for Multi-Provider Private Networks
- ❑ Edge-to-Edge Control Architecture
- ❑ Riviera Congestion Avoidance
- ❑ Trunk Service Building Blocks
 - ❑ Weighted Sharing
 - ❑ Guaranteed Bandwidth
 - ❑ Assured Bandwidth

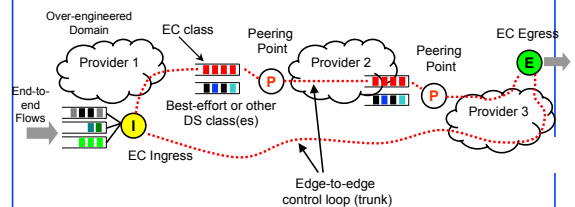
Scalability Issues

- ❑ Traditional solutions:
 - ❑ Use QoS:
 - ❑ ATM, IntServ: per-flow/per-VC scheduling at *every* hop.
 - ❑ Frame Relay: Drop preference, per-VC routing at *every* hop.
 - ❑ DiffServ: per-class (eg: high, low priority) scheduling, drop preference at *every* hop. Per-flow QoS done only at network boundaries (edges).

QoS for Multi-Provider Private Networks

- ❑ Principle Problems
 - ❑ **Coordination:** scheduled upgrades, cross-provider agreements
 - ❑ **Scale:** thousands-millions connections, Gbps.
 - ❑ **Heterogeneity:** many datalink layers, 48kbps to >10Gbps

Edge-to-Edge Control (EC)



Use Edge-to-edge congestion Control to push queuing, packet loss and per-flow bandwidth sharing issues to edges (e.g. access router) of the network

QoS via Edge-to-Edge Congestion Control

- Benefits:
 - Conquers *scale* and *heterogeneity* in same sense as TCP.
 - Allows QoS *without upgrades* to either end-systems or intermediate networks.
 - Only *incremental* upgrade of edges (e.g., customer premise access point).
 - Bottleneck is CoS FIFO.
 - Edge knows congestion state and can apply stateful QoS mechanisms.
- Drawbacks:
 - Congestion control cannot react faster than propagation delay. → Loose control of delay and delay variance.
 - Only appropriate for data and streaming (non-live) multimedia.
 - Must configure edges and potential bottlenecks.

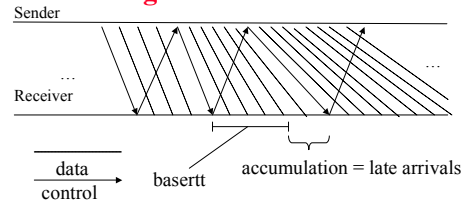
Mechanisms for Fairness and Bounded Queue

- Estimate this control loop's backlog in path.
If backlog > max_thresh
Congestion = true
Else if backlog <= min_thresh
Congestion = false
- All control loops try to maintain between min_thresh and max_thresh backlog in path.
→ bounded queue (**Goal 4**)
- Each control loop has roughly equal backlog in path → proportional fairness [Low] (**Goal 5**)
- Well come back to goal 5.

Riviera Congestion Avoidance

- Implements EC Traffic Trunks.
- EC Constraints:
 - Cannot assume access to TCP headers.
 - No new fields in IP headers (no sequence numbers)
 - Cannot assume existence of end-to-end ACKs (e.g., UDP)
 - Cannot impose edge-to-edge ACKs (doubles packets on network)
- → *No window-based control.*
- Solution: *rate-based control.*

Backlog Estimation and Goal 2



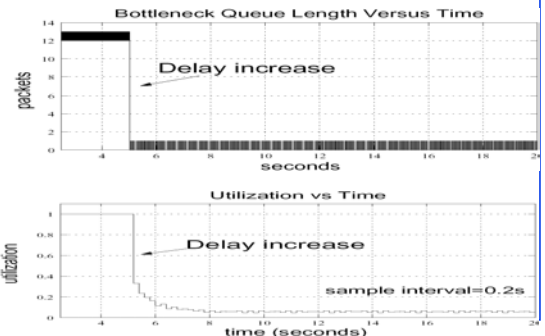
- Use basertt like Vegas backlog estimation.
- As with Vegas, when basertt is wrong → **gross unfairness** (violates Goal 2).
- Sol'n: *ensure good basertt estimate.*

Congestion Avoidance Goals

- 1. Avoid of congestion collapse or persistent loss.
 - Behave like TCP Reno in response to loss.
- 2. Avoid starvation and gross unfairness.
 - Isolate from best effort traffic.
 - Solve Vegas RTPD estimation errors.
- 3. High utilization when demand.
- 4. **Bounded queue.**
 - **Zero** loss with sufficient buffer.
 - Accumulation.
- 5. **Proportional fairness.**
- ...
- **Attack goals 2,4, and 5 in reverse order.**

Vegas & Delay Increase (Goal 2)

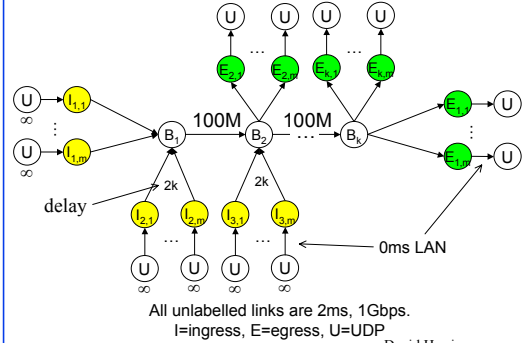
- Vegas sets basertt to the minimum RTT seen so far.
- → **GROSS UNFAIRNESS!**



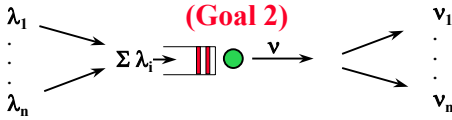
Riviera Round-trip Propagation Delay (RTPD) Estimation (Goal 2)

- Reduce gross unfairness w/ good RTPD estimation.
- Minimum of last $k=30$ control packet RTTs.
- Drain queues in path so RTT in last k RTTs likely reflects RTPD.
 - Set max_thresh high enough to avoid excessive false positives.
 - Set min_thresh low enough to ensure queue drain.
- Provision drain capacity with each decrease step

Proportional Fairness Topology (Goal 5)



Increase/Decrease Policy to Drain Queue (Goal 2)



r_i = rate limit on leaky bucket (σ, ρ) shaper. $\lambda_i \leq r_i$

Increase/decrease Policy

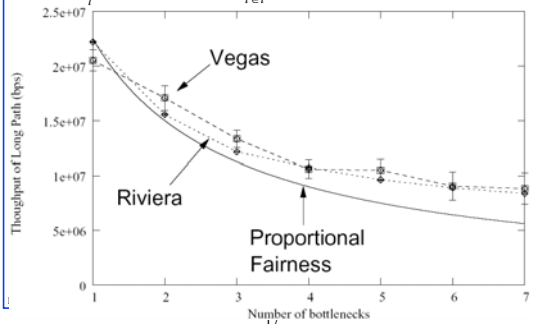
$$r_i = \begin{cases} v_i + \text{MTU}/\text{RTT} & \text{if no congestion} \\ \beta v_i & \text{if congestion} \end{cases}$$

$1 > \beta \gg 0$

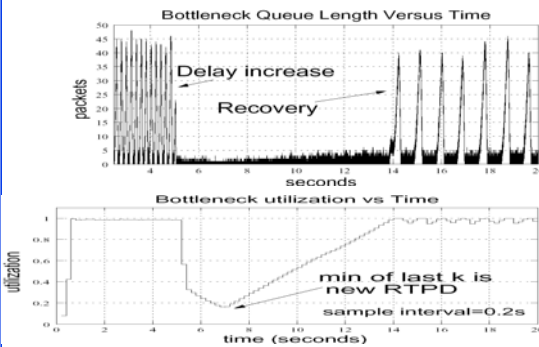
Lower β improves probability queues drain at cost to utilization.

Riviera Achieves Proportional Fairness? (Goal 5)

$$\max \sum_i \log \lambda_i \quad \text{with} \quad \sum_{i \in L} \lambda_i \leq C_i, \quad i \in L, \quad \text{and} \quad \lambda_i \geq 0$$

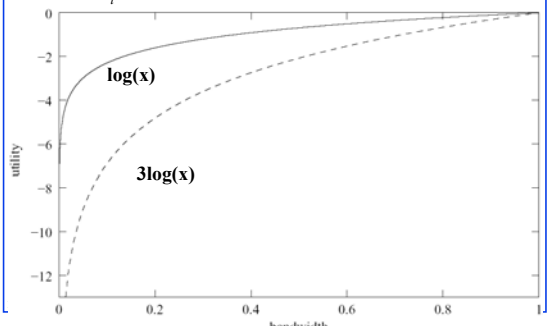


Riviera & Propagation Delay Increase (Goal 2)



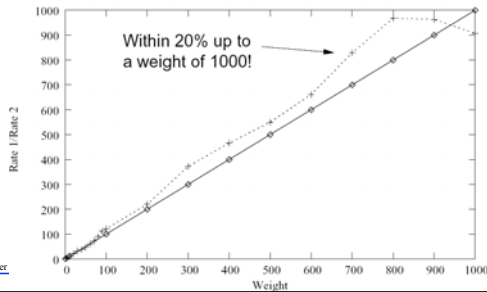
Weighted Proportional Fairness

$$\max \sum_i w_i \log \lambda_i$$



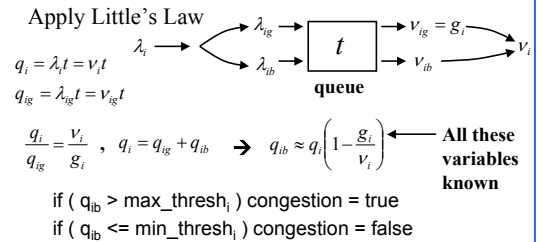
Weighted Service Building Block

- Modify accumulation thresholds:
 $\max_thresh_i = w_i * \max_thresh_1$
 $\min_thresh_i = w_i * \min_thresh_1$



Quasi-Leased Line (QLL)

- Converges on guaranteed bandwidth allocation.
- Accumulation Modification:**

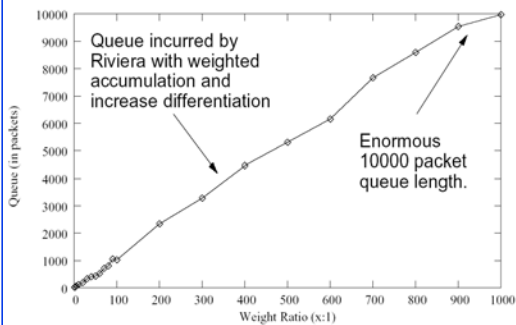


Rensselaer Polytechnic Institute

David Harrison

22

Weighted Service Building Block (2)



20

QLL Increase/Decrease Policy

- Increase/decrease policy:**

$$r_i = \begin{cases} \max(g_i, v_i + \text{MTU}/\text{RTT}) & \text{if no congestion} \\ \max(g_i, \beta(v_i - g_i) + g_i) & \text{if congestion} \end{cases}$$

$1 > \beta \gg 0$

Go immediately to guarantee and refuse to go below.

Decrease based only on the rate that is above the guarantee

- No admission control → unbounded queue.

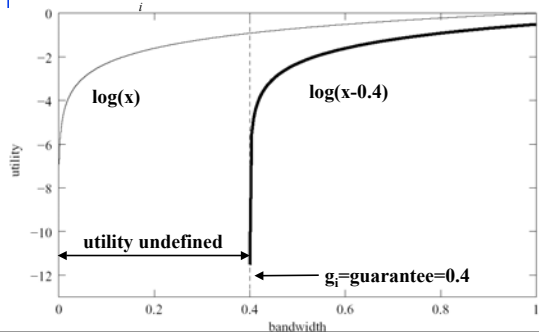
Rensselaer Polytechnic Institute

David Harrison

23

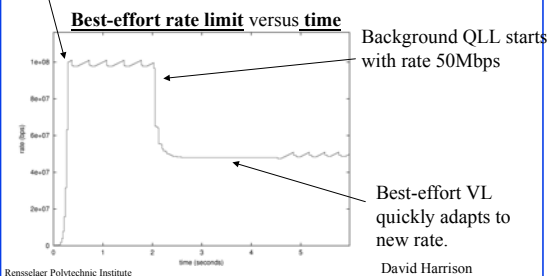
Guaranteed Bandwidth Allocation

$$\text{maximize } \sum w_i \log(\lambda_i - g_i)$$



Quasi-Leased Line Example

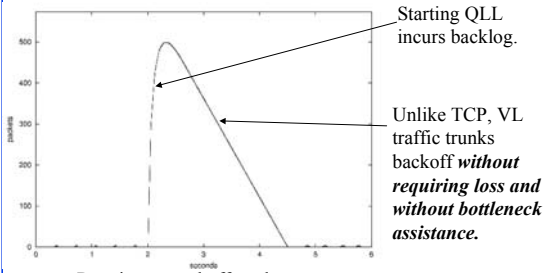
Best-effort VL starts at $t=0$ and fully utilizes 100 Mbps bottleneck.



24

Quasi-Leased Line Example (cont.)

Bottleneck queue versus time



Requires more buffers: larger max queue

Assured Building Block

- Accumulation: $q_{ib} \approx q_i \left(1 - \frac{a_i}{v_i}\right)$
 - if $(q_{ib} > \max_thresh \ || \ q_i > w_i * \max_thresh)$
congestion = true
 - else if $(q_{ib} \leq \min_thresh \ \&\& \ q_i \leq w_i * \max_thresh)$
congestion = false

Increase/Decrease Policy:

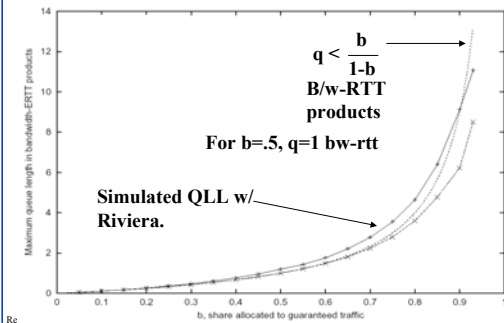
$$r_i = \begin{cases} v_i + MTU/RTT & \text{if no congestion} \\ \min(\beta_{AS} v_i, \beta_{BE}(v_i - a_i) + a_i) & \text{if congestion} \end{cases}$$

$1 > \beta_{AS} > \beta_{BE} >> 0$

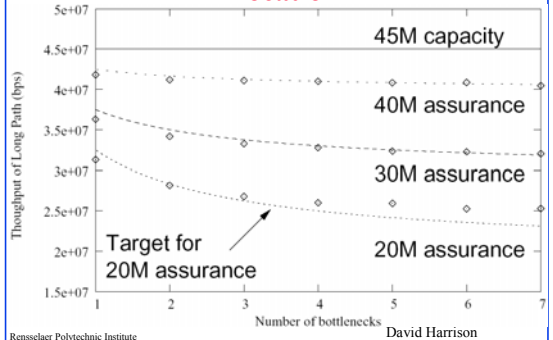
- Backoff little (β_{AS}) when below assurance (a),
- Backoff (β_{BE}) same as best effort when above assurance (a)

Quasi-Leased Line (cont.)

Single bottleneck queue length analysis:

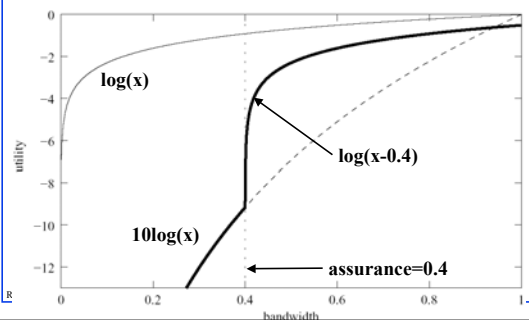


Assured Building Block Vs. Assured Allocation

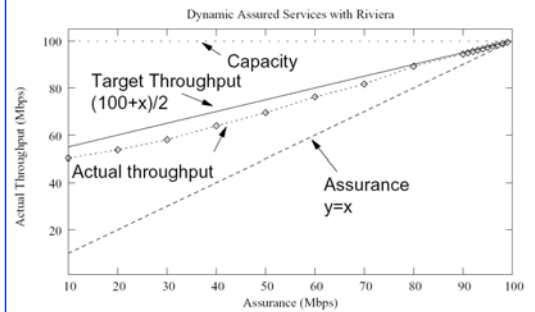


Assured Bandwidth Allocation

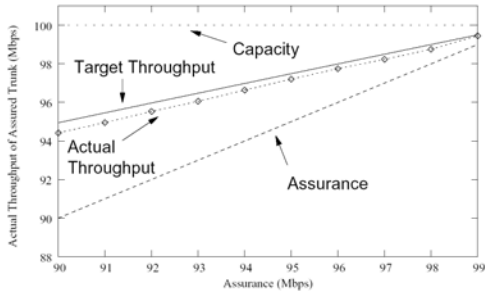
$$\text{maximize } \sum_i \begin{cases} \log(\lambda_i - a_i) & \text{if } \lambda_i > C_i \\ w_i \log \lambda_i & \text{if } \lambda_i \leq C_i \end{cases} \quad \text{with } \log(C_i - a_i) = w_i \log C_i$$



Wide Range of Assurances



Large Assurances

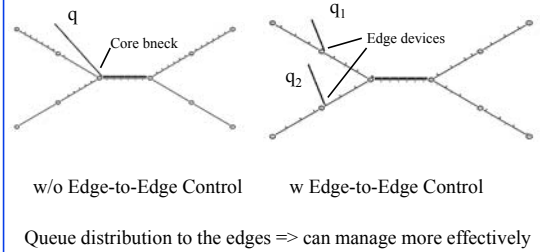


Rensselaer Polytechnic Institute

David Harrison

31

Edge-to-Edge Queue Management



Rensselaer Polytechnic Institute

David Harrison

34

Summary



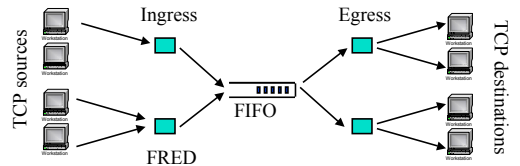
- **Issues:**
 - Simplified overlay QoS architecture
 - **Intangibles:** deployment, configuration advantages
 - **Edge-based Building Blocks & Overlay services:**
 - A **closed-loop** QoS building block
 - Weighted services, **Assured** services, **Quasi**-leased lines

Rensselaer Polytechnic Institute

David Harrison

32

Distributed Buffer Management (1)



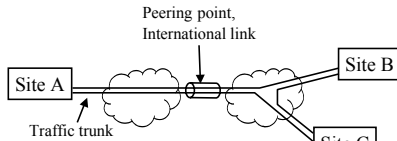
- Implement FRED AQM at edge rather than at bottleneck. Bottleneck remains FIFO.
- Versus FRED at bottleneck and NO edge-to-edge control.

Rensselaer Polytechnic Institute

David Harrison

35

Private Versus Public Networks



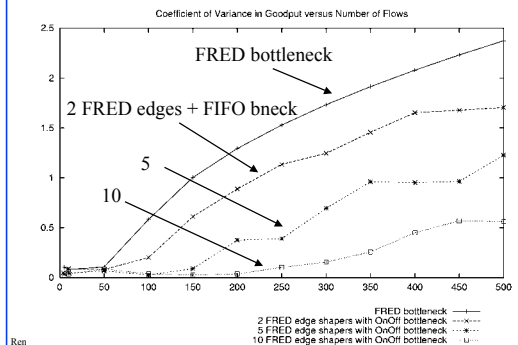
- **Private Networks:**
 - < Tens or hundreds of communicating parties
 - limits per-loop state and configuration at edges.
 - Parties remain in VPN for moderate to long term
 - Amortize configuration over lifespan of VPN.
 - **Aggregation**
 - Amortize control packet overhead.
- **Public Internet:**
 - None of these assumptions hold.

Rensselaer Polytechnic Institute

David Harrison

33

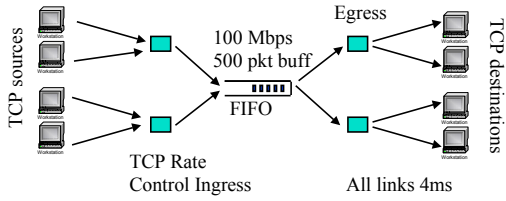
Distirbuted Buffer Management (2)



Ren

36

TCP Rate Control (Near Zero Loss)



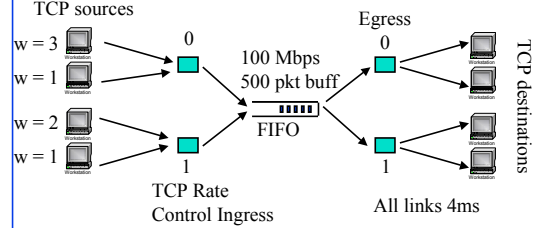
- Use Edge-to-edge Control to push bottleneck back to edge.
- Implement TCP rate control at edge rather than at bottleneck. Bottleneck remains FIFO.

Rensselaer Polytechnic Institute

David Harrison

37

Remote Bottleneck Bandwidth Management



- Edge redistributes VL's fair share between end-to-end flows.

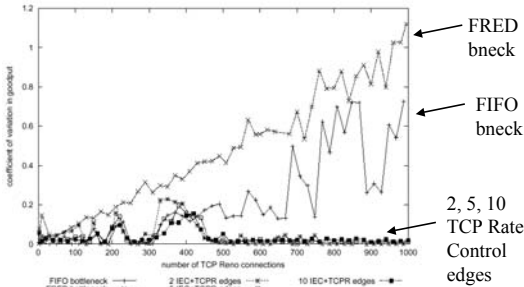
Rensselaer Polytechnic Institute

David Harrison

40

TCP Rate Control (2)

Coefficient of Variation in Goodput vs. 10 to 1000 TCP flows

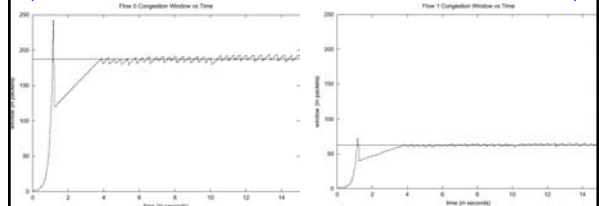


Rensselaer Polytechnic Institute

David Harrison

38

Remote Bandwidth Management (2)



TCP 0 with weight 3.
obtains 3/4 of VL 0

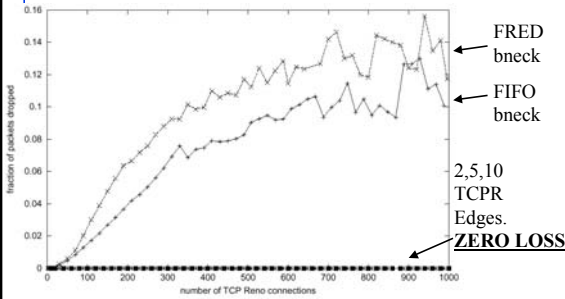
TCP 1 with weight 1
obtains 1/4 of VL 0

Rensselaer Polytechnic Institute

David Harrison

41

TCP Rate Control (3)

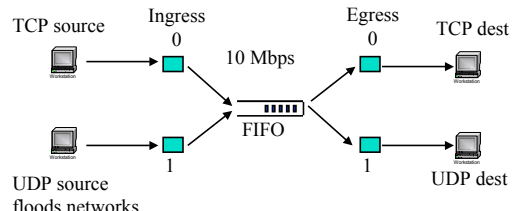


Rensselaer Polytechnic Institute

David Harrison

39

UDP Congestion Control, Isolate Denial of Service

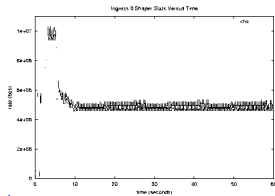


Rensselaer Polytechnic Institute

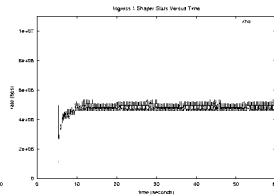
David Harrison

42

UDP Congestion Control, Isolate Denial of Service

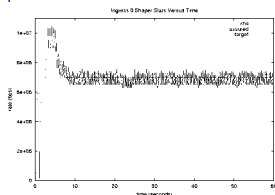


Trunk 0 carries TCP starting at 0.0s

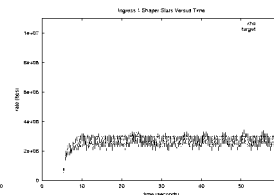


Trunk 1 carries UDP flood starting at 5.0s

Effects: Bandwidth Assurances



TCP with 4 Mbps assured + 3 Mbps best effort



UDP with 3 Mbps best effort