# AQM

## Drop-tail Queues
Good: **easy to implement**
Bad: Filled buffers (buffer delay), synchronized window



(RTT Flow)
half the slope → half $W_{max}$
(2×RTT, Flow) ↑
$t$    2×RTT → ¼ throughput

## : RED
Idea: Drop probabilistically to prevent cong. early and desynchronize flows

Implementation: Drop based on avg. queue len. $x(t)$



every T:
$x(t) \leftarrow w_{q_k} \cdot q_k(t) + (1-w_{q_k}) x(t-T)$
$q_k(t)$ instant queue len.

Problems:
① Params hard to tune
② Queue depend on RTT and number of flows

## * PI (Proportional Integral)
Ideas: ① Remove EWMA ⇒ responds faster than RED
② Integral Control → Decouple queue len. & num. flow
③ Use derivative of queue → more stable

$P(t) \leftarrow P(t-T) + \alpha(q_k(t) - q_{ref}) + \beta(q_k(t) - q_k(t-T))$

Integral control to drive $q$ to $q_{ref}$
derivative, we should respond faster when $\frac{dq}{dt}$ is big

## * PI Enhanced
Ideas: ① Control delay instead of length
② Auto tune params $(\alpha, \beta)$ based on $P$

## XCP
Sender reports RTT and CWND, router specify $\Delta$CWND
No per-flow state at router. $CWND \mathrel{+}= \Delta CWND$

**✷ Efficiency Controller**
Match input to output capacity to keep queue short

---

S: Spare Bandwidth
d: average RTT
→ match $B_{in}$ and $B_{out}$
→ drain the queue
$\Delta = \boxed{\alpha \cdot d \cdot S} - \boxed{\beta \cdot Q}$
$\alpha$: queue len.

## * Fairness Controller
Divide $\Delta$ between flows (look at info reported by TX)
$\Delta > 0$: distribute evenly → same throughput increase

Per-flow $\Delta_{CWND_i} \propto RTT_i$

Per-packet $P_i \propto \dfrac{RTT_i}{\frac{CWND_i}{RTT_i}}$ ← Rate of incoming packets

$\sum \dfrac{P_i}{RTT_i} = \Delta/d$ ← need to change $\Delta$ by time d

↑ this change takes effect after time $RTT_i$
for all packets during control interval $d$

$\Rightarrow k \sum \dfrac{RTT_i}{CWND_i} = \frac{\Delta}{d}$ , $k = \dfrac{\Delta}{d} / \sum \dfrac{RTT_i}{CWND_i}$

$\Delta < 0$: distribute prop. to throughput
↓
faster flow decrease faster ($\propto$ throughput)

Per-flow $\Delta_{CWND_i} \propto CWND_i$
Per-packet $P_i \propto \dfrac{CWND_i}{\frac{CWND_i}{RTT_i}} = RTT_i$

## Fair Queueing & CSFQ

Work conserving: Never idle if has packet in Q
Max-min Fairness: converge to $\alpha$ that
request $r_i < \alpha$: give $r_i$
$r_i > \alpha$: give $\alpha$
Scheduling: Which pkt to send next?
How to fair? Bit-by-bit Round-robin
Real world: Emulate bit-by-bit RR

## * Fair Queueing
Round: each queue sends one bit
$\mu$: output rate ; N: num of flows
$\dfrac{dR}{dt} = \dfrac{\mu}{N}$  (Num of rounds to finish a pkt is independent of N)

$P_i^\alpha$: $i$-th pkt of queue $\alpha$ (the size of)
$S_i^\alpha$: when it reaches head of queue

---

$F_i^\alpha$: when it is finished transmitting
$S_i^\alpha = \max(R(t), F_{i-1}^\alpha)$ ; $F_i^\alpha = S_i^\alpha + P_i^\alpha$

Send packet with smallest $F_i^\alpha$
Deficit Round Robin: for each queue, credit increases at rate of fair rate, decrease by pkt sz.

## * CSFQ
Ideas: ① Edge routers mark estimate of arrival time
② Core routers use ① and internal measure of fair share to compute prob. of drop.
③ Estimation of fair share converges quickly
(Arrival rate)
$r_i^{new} = (1 - e^{-T_i^k/k}) \dfrac{l_i^k}{T_i^k} + e^{-T_i^k/k} r_i^{old}$
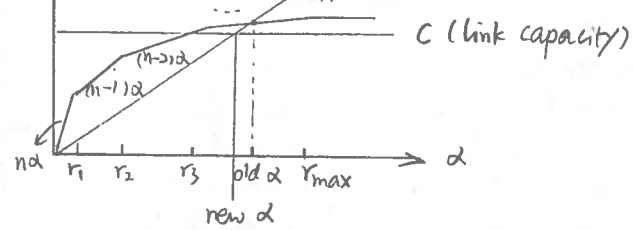
$t_i^k$: arrival time of $k$th packet in flow $i$
$l_i^k$: length of $k$th packet in flow $i$
$T_i^k = t_i^k - t_i^{k-1}$   $k$: constant

$P(drop) = \max(1 - \dfrac{\alpha}{r_i(t)}, 0)$ , $\alpha$ is fair rate

Acceptance rate $F(\alpha) = \sum_i \min(r_i(t), \alpha)$



## BGP
Route: link/next-hop to a dest (local info)
Path: Sequence of edges  (the whole path)
Why scalable: ① Nearby in topology ⇒ similar IP
② Route announced as prefixes
Longest Prefix Match: use the longest match of prefix
AS: unit of who announces a route
Transit: provider provides access for customer
Peering: mutual access to subset of each other
B advertise P to A: B will forward pkts to P from A
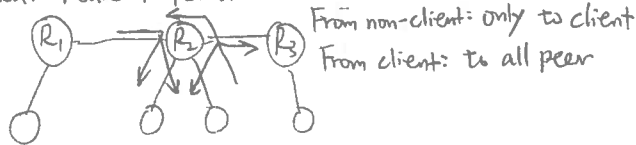Pricing: 95% of 5-min moving avg. throughput
Pick route: Customer > Peer > Provider (LOCAL PREF)
Announce Route (see reverse side)

| To \ from | Customer | Peer | Provider |
|---|---|---|---|
| Customer | ✓ | ✓ | ✓ |
| Peer | ✓ | X | X |
| Provider | ✓ | X | X |

BGP Goals: ① Scalability ② Policy ③ Cooperative Competition

iBGP Route Reflector



From non-client: only to client
From client: to all peer

BGP Attributes:
NEXT HOP: IP of next-hop router
ASPATH: all AS that announcement goes thru
MED: should use entry point with smallest MED

BGP Path Selection:
LOCAL PREF > len(ASPATH) > MED > eBGP or iBGP (learned from?) > IGP Path cost > Smaller Router IP

Measurement

What? Transport: performance, congestion control
Network: routing fail, topology, performance
AS-level Topo: find ASes and BGP links
Router Alias Resolution: map IP to Router
Methods: ① Probe and see reply IP
② Increasing IPID field
ISP Topology Inference: combine tracert info
Challenge: how to choose targe IP to trace?
Rocket fuel: ① sufficient num of vantage points
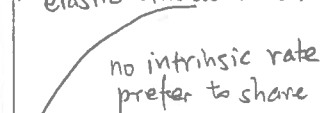② select target IP
③ Deal with tracert issues

* ZMap
Idea: ① bypass TCP/IP and craft Eth frames
② Encode dest info in packet → Stateless

NUM (Network Utility Max.)

$x$) Utility Function: benefit of sending at rate $X$
elastic (file download)    inelastic (real-time)



no intrinsic rate prefer to share

intrinsic rate prefer to randomize

---

Max. $\sum_{i=1}^{N} U_i(x_i)$   N flows, L links

routing matrix

When $R_{L \times N} \times [x_i] \leq [c_i]$  (capacity)

If $U(x)$ concave, have unique solution.

* $\alpha$-fairness

maximize $\sum_{i=1}^{N} \frac{x_i^{1-\alpha}}{1-\alpha}$  ($\alpha \geq 0$)

$\alpha = 0$: throughput max.
$\alpha \to 1$: prop. fairness ($\log(x_1) + \log(x_2) + \cdots$)
$\alpha \to \infty$: max-min fair

* Solve NUM

$p_\ell$: cong. price for link $\ell$ (price per unit bw)
$q_i$: total price for source $i$ ($\sum p_\ell$ along path)
profit = $U_i(x_i) - q_i x_i$ (max. done at source indep)
$p_\ell(t+1) = \max(p_\ell(t) + k(y_\ell(t) - c_\ell), 0)$
(done at link)
total traffic   capacity

* TCP+PI:
$P(\ell)$ follow the same update rule as prob. drop
solution: $x_i = \frac{1}{RTT_i \sqrt{q_i}}$  (TCP formula)

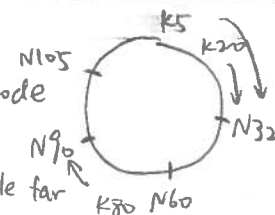$U_i(x_i) = \frac{-1}{RTT_i^2 x_i}$  ($\alpha$-fair, $\alpha = 2$)

DHT / Chord



Node and key share a ring
key stored on next higher node
Consistent hashing
Each node point to ½, ¼, cycle far

* Fault tolerance:
Each node keep next $r$ nodes
Each key stored by $r$ nodes after owner

* Joining
1. look up itself to find the successor
2. set self. successor
3. copy keys from successor
4. call stabilize
① find succ. pred. and set self. succ
② notify self. succ about itself

---

* Lookup
Each node stores $succ(n+2^{i-1})$, i.e. node that have n+2
Get k: search j in table that j is closest smaller to k

VL2 / Datacenter

Agility: any service at any server
Layer 2 Good: Auto-config, failover
Bad: Broadcast (ARP), no multipath (STP)
Layer 3 Good: Scalable, multipath
Bad: Hard to migrate (change IP), config

* Conventional Problem    * Goal of VL2



1. L2 semantics
2. Uniform high capacity
3. Performance isolation

* ECMP LB
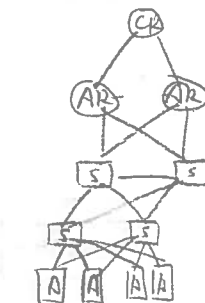→ Pick equal-cost paths by hash of 5-tuples (flow-level LB)
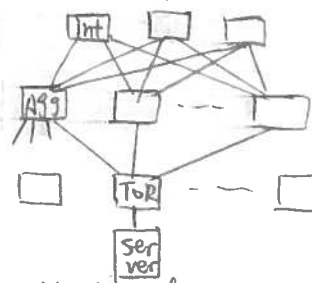Problem: hash collisions
Not Problem:
① Flows are many and small
② Switch-switch link thick than flow (NIC) size

* Clos Topo



* Name-Loc Separation

Number of servers: (24-port switch)
Top: all 24 ports connect to next level
Other layers: half talk to upper level

Congestion Ctrl



* Cubic:
$W(t) = \alpha(t-k)^3 + W_{max}$
Loss: $t=0$, $W(0) = (1-\beta)W_{max}$
$P \propto \frac{1}{W_m^{4/3}}$, Throughput $\propto \frac{1}{p^{3/4}}$

* Reno:
$BDP + Q = W_{max}$, $\frac{W_{max}}{2} \geq BDP \Rightarrow Q \geq BDP$ for max util
$P \propto \frac{1}{W_m^2}$, $W_m \propto \frac{1}{\sqrt{P}}$, Throughput $\propto \frac{1}{\sqrt{P} \cdot RTT}$