# PRINCIPLES OF COMMUNICATIONS

## Routing

Finding path from source to destination in a network. Issues include:
- what route?
- is that the shortest route
- what happens if link goes down?
- Multipath routing (for mobile wireless link)
- latency, bandwidth etc

① Interdomain Routing: socioeconomic issue — comms between different stakeholders
② Multicast Routing
③ Routing in telephone networks → often just random routing
 - useful also for load balancing
④ Routing for mobile hosts.

## Packet Switching

(1) Routing protocol sets up a (2) routing table in routers and switch controller
- node makes a local choice depending on global topology and state
- Separation between control and data view of network (and group addressing). We may have a centralised system
  - inherently large — can't topologically sort — if we can sort hierarchically then aggregate route entries
  - Dynamic (traffic conditions)
  - hard to collect (failures, etc) → restore fast → do we wait or find a new path.
- → Must intelligently summarize relevant information

### Requirements
① Minimize routing table space
- fast to look up ⎤ expensive
- less to exchange ⎦
② Minimize number and frequency of control messages — dealing with lots of changes
③ Robustness, avoiding:
- Black Holes: Islamabad example where YouTube ISP was made by to redirect to them — hence, best route was to go via this hence entire world followed this path

- Loops: TTL is a last resort; especially with multiclass traffic
- Oscillations:
④ Use optimal path

Features: Goal: maximise throughput, subject to min delay and cost.
① Packets
② Topology is complicated
③ Many provides
④ Traffic sources are bursty
⑤ Traffic matrix is unpredictable and large — traffic is more important than latencies

## Routing Model

(1) Dynamic routing and (2) Intra and Inter AS — AS = each ISP is an AS locus of admin control
Interior Gateway Protocols → IGPs inside AS's: RIP, OSPE, HELLO → autonomous systems
  - oldest
  - bay → distribute cost of paths

↳ Exterior Gateway Protocols (EGPs) for inter AS routing.
  ↳ EGP, BGP-4

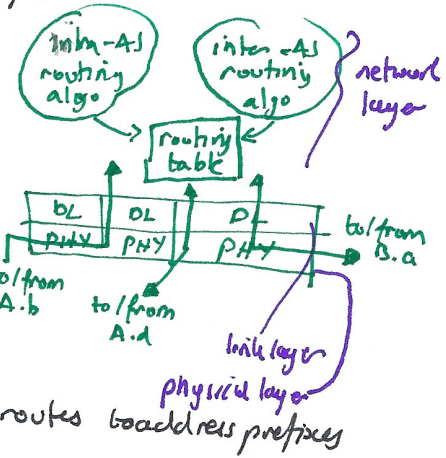## Intra-AS Routing Requirements

① Scale for size of AS

② Different requirements on routing convergence after topology changes

③ Operational / Admin / Management (OAM) complexity

④ Traffic Engineering Capabilities

## Inter-AS Routing Requirements

① should scale for global internet size
  ↳ reachability not |optimality| — lots of possible goals
     ↳ latency, throughput, to/from
  ↳ address aggregation techniques
  ↳ Allow flexibility in topological structure.

② Allow policy-based routing between autonomous systems.
  ↳ In case of routing, options include advertised AS-level routes & address prefixes
  ↳ Fully distributed routing is the only possibility
  ↳ Extensible

**Gateways**
① perform inter-AS routing amongst themselves
② perform intra-AS routing.



## Basic Dynamic Routing Methods

① Source Based: source gets a map of the network — IP allows that to occur
  ↳ source finds route and (1) signals route-setup
                          or (2) encode route into packets.
  ↳ But, can have attacks on the way

② Link State Routing (per link info)  → centralized
  ↳ get map of network at all nodes and find next-hops locally

③ Distance Vector (per node info)
  ↳ at every node, set up distance signposts to destination nodes
  ↳ look at neighbours signposts ⇒ BGP

## Choices

① Centralized?: simpler but prone to failure and congestion
  ↳ need very reliable central system

② Source-based vs hop-by-hop.
  ↳ more expensive + packet header size
  ↳ Intermediate : loose source route

③ Stochastic vs deterministic
  ↳ spreads load, avoiding oscillations but misorders

④ Single vs multiple path

⑤ State-dependent vs state-independent ⇒ routes depend on current network state?

## Central Control over Distributed Routing (fibbing)

SDNs (software defined networks) :- latency can be very low using cut-through switches, also single managment area. Effectively using remote control switching. change the table to additional entries to do what they want for entries ⇒ if they define breaks, it returns to default. Effectively a hybrid of link state and centralised controller
  ↳ central system adds security

* thanks to partial distribution - fallback system

| | Traditional | SDN | fibbing |
|---|---|---|---|
| Manageability | low | high | high ~ same as SDN |
| Flexibility | low | highest | high fibbing is a hybrid SDN architecture |
| Scalability | by design | ad hoc | by design per destination full control |
| Robustness | high | low | *high some functions are distributed |

SDN does three things: (1) computes paths, (2) derives FIB entries, (3) installs FIB entries
└→ effectively allows you to control the switches
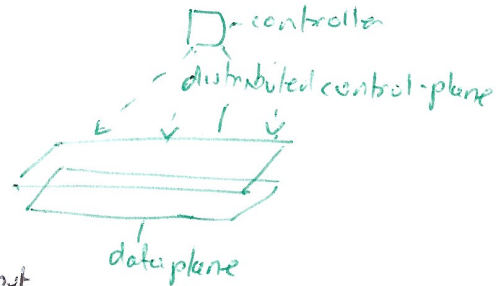
Fibbing centralizes only high-level decisions
    └→ controller:
        └→(1) computes paths
    └→ distributed control-plane:
        └→(1) computes FIB entries
        └→(2) install FIB entries

To control the IGP output, the Fibbing controller inverts the shortest-path function.

- controller
- distributed control-plane

data plane

input
weighted
topology  →  shortest-path  →  output
                function / computation    forwarding plans

output format determined by operators

## Flexibility

Operators can modify shortest path outputs by injecting information as fake nodes and links to the IGP control-plane
    └→ any set of forwarding DAGs can be enforced by Fibbing through:
        └→① fine-grained traffic steering – ⬚ middleboxing
        └→② per-destination load balancing – ⬚ traffic engineering
        └→③ backup paths provisioning – ⬚ failure recovery

## Scalability

Fibbing controller says algorithms to be run in sequence
    └→ merger iteratively merges fake nodes – programs multiple next-hop changes with a single fake node
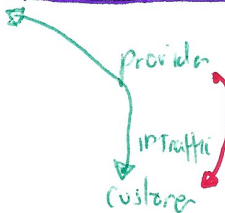Does not impact IGP convergence – achieves fast forwarding changes

## Robustness

Improves it by relying on underlying IGP which: → re-converge quickly,
        └→ provides fast failure detection and control-plane sync
        └→ supports fail-open and fail-close semantics
        └→ survives replica failures with no impact on forwarding

## INTERDOMAIN ROUTING

Provider

IP Traffic

Customer

customer pays provider for access to the internet.

peers provide transit between their respective customers. They do not provide transit between peers. Allows connectivity between customers of Tier-1 providers

| peer | don't peer |
|---|---|
| reduces upstream transit costs | would rather have customers |
| Increase upstream performance | Peers usually competition |
| Only way to connect customers to Tier-1 | Relationships require periodic renegotiation ↓ can be messy negotiations |

Forwarding: determine next hop → always works
Routing: establish end-to-end paths
    └→ can be badly broken
    └→ may not always find a path

④ Forwarding tables used to implement routing:

↳ Statically :- world manually configures routes : Admins manually configures
   forwarding table entries  — mostly at the edges
↳ Dynamically : routers exchange reachability information to compute best routes
   ↳ Interior Gateway Protocol (IGP)
      ↳ metric based : OSPF, IS-IS, RIP, EIGRP
   ↳ Exterior Gateway Protocol (EGP)
      ↳ policy based : BGP

| Link State | Vectoring |
|---|---|
| · Topology information flooded within the routing domain. | · Router knows little about network topology |
| · Best end-to-end paths computed locally at each router | · Best next-hops chosen by each router |
| · Best paths determine next hop | · Best paths result from composition of all next-hop choices |
| · Based on minimizing some notion of distance | · No notion of distance |
| · Works only if policy is shared and uniform | · No uniform policies |

IGP = OSPF /
IS-IS

IGP = RIP
EGP = BGP

Routing computation distributed among routers within routing domain

Autonomous Routing Domains (ARD): collection of physical networks glued together using IP that have a unified administrative routing policy.

LAS is an ARD assigned an Autonomous System Number (ASN) — 32 bit values ⇒ these can be shared

↳ ARDs use static routing

BGP-4 : de-facto EGP today.
   ↳① Establish TCP session
      ↳ authentication
      ↳ certification
   ↳② Exchange active routes
      ↳ advertise what's reachable + border router
   ↳③ Exchange incremental updates
      ↺

① eBGP (external) in a different Autonomous system
② iBGP (internal) in same Autonomous system.
   ↳ iBGP mesh does not scale

Route Reflectors

Can pass it on iBGP updates to client — passes through only best routes, avoiding loops
BGP Confederations : has multiple internal ASs that present one external AS

Messages
   ↳① Open
   ↳② Keep alive
   ↳③ Notification
   ↳④ Update  ⇒ announcement = IP_prefix + attributes_values
Attributes including :
   ↳ AS_PATH        ↳ COMMUNITY
   ↳ NEXT_HOP       ↳ ORIGIN
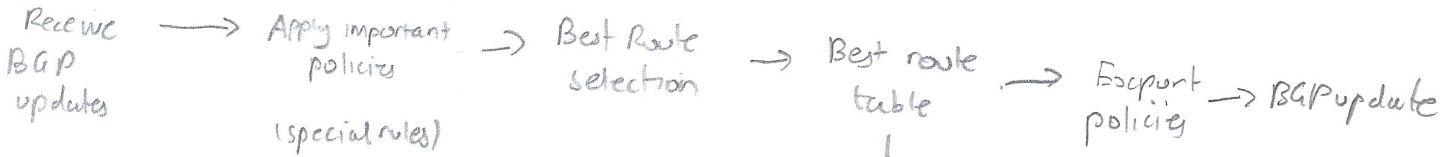   ↳ MULTI EXIT DISC  ↳ CLUSTER LIST
   ↳ LOCAL PREF

(5)

## Route Selection
↳① Highest Local Preference —— used only in IBGP
↳② Shortest ASPATH
↳③ Lowest MED  multi-exit descrimination ⎤
↳④ iBGP < eBGP                          ⎟ traffic engineering
↳⑤ Lowest IGP cost to BGP egress        ⎦
↳⑥ Lowest router ID ⎤ break all ties

## BGP Route Processing

The process is as
follows:

Receive ——→ Apply important ⟶ Best Route ⟶ Best route ⟶ Export ⟶ BGP update
BGP           policies          selection       table       policies
updates      (special rules)

Next hop attribute: if crosses AS boundary, NH
attribute is changed to IP address of border router
that announced the route

EGP joins to IGP by connecting forwarding tables and
EGP tables.

~~RSA~~.

BGP Communities: has a community value of 32 b.b
↳ import: customer routes, peer routes, provider routes
↳ export: to customers, to peers, to providers

⌐→ Next 16 is community number ⎤
⌐→ First 16 b.b is ASN          ⎦ used for signalling within and between ASes
Community Attribute = list of communites
values — one route can belong to
multiple communites

IP forwarding
table
↓
propagate
to clients

### Enforcing Customer/Provider and Peer/Peer relationships
① Enforce transit relationships -
output route filtering
② Enforce order of route preference-
to provider < peer < customer

## Traffic Engineering with BGP

| inbound | Outbound |
|---|---|
| ① Filter outbound routes | ① Filter inbound routes |
| ② tweak attributes on these routes to influence best neighbour's best route selection. | ② Tweak attributes on inbound routes to influence best route selection |

↳ don't accept BGP AS PATHS that
would lead to a loop
↳ traffic often follows the ASPATH but
sometimes doesn't
↳ we can implement backup links with
local preference for outbound traffic
and inbound traffic
↳ use ASPATH padding to make the course
less good by making it look longer
but this does not always work

AS has more control over outbound traffic

↳ customer routes +
local preference > ASPATH
length
↳ can fix this with a
community attribute

## Hot potato routing means we go for the
closest egress point
↳ this can have many issues though ⇒ can have
Multi-Exit Discriminator Attribute (MEDs), hence
provider network as far as possible

## Route Pinning Example

original goal was just reachability and loop free
→(i) Couldn't optimise for a metric
→ (ii) couldn't scale to large scale
→ traffic remains pinned to backup if disaster strikes
N.B BGP not guaranteed to converge on stable routing ⇒ leads to livelock
→ not guaranteed to recover from network failure

For static semantics, BGP policies solves the Stable Paths Problem.
→ dynamic semantics, BGP solves SPVP
→ Simple Path Vector Protocol = distributed algorithm for solving SPP

## Stable Paths Problem

→ graph of nodes and edges
→ origin nodes
→ For non-origin node, set of permitted paths to origin
→ ranking of permitted paths
→ origin is n
→ solution does not represent a shortest path tree or spanning tree
→ can have multiple solutions ⇒ results in route triggering
→ bad gadget is where there is no solution

## Internet Growth Trends

Large BGP tables is an issue — causes a large problem and slows things down.
→ goals are ① fast convergence
② minimal updates
③ path redundancy

Can have bad apples ⇒ leads to local traffic meaning people have to recalculate paths — reduce effect of this by squashing updates.
→① Rate limiting on sending updates - 30 seconds
→② Route Flap Dampening - punishment system

→① rate limiting dampens some of oscillations inherent in a vectoring protocol.
→② Routes given penalty for changing - route is dampened which decays exponentially.
If penalty goes below reuse limit, then announced again
→ applied only on eBGP inbound only

→ however, lots and lots of BGP updates:
→① Misconfiguration
→② Route flap dampening not widely used
→③ Software bugs
→④ BGP exploring alternate paths

| MEDs export internal instability

## Multicast Routing

→ packet sent by any member of a group are received by all.
Multicast Group : associates sender and receiver
→ sender does not know receivers identities
→ has its own class D address which send to and receives request packet from that address.
→ dynamic directory service associates the two

→ Issues
→① Currently active groups
→② How to express interest in joining
③ Discover receivers in a group
④ Delivery data to members in group

# Ring Search (Expanding)

Method of using multicast groups for resource discovery
↳ reaches all receivers - discovering local resources first

## Multicast Flavours

↳ Unipoint : point to point .
↳ Multicast : point to multipoint ⟶ simulate by set of point to point unicast
multipoint to multipoint ⟶ simulate by set of point to multipoint

(shortest-path tree)
Form multicast tree, routed at the sender

Issue in Wide-Area Multicast - exploit LAN's broadcast capability

↳① sources join and leave dynamically
↳② leaves of tree members of broadcast LAN
↳③ want receivers to join and leave without notifying sender
↳ Goal: distribute packets coming from sender to all routes on path from group member

*each endpoint is a router ⟶ uses IGMP to get members in LAN*

↳ does LAN contain members for a given group
↳ what MAC address corresponds to IP address
↳ there's a well known translation ~~table~~ table ⟹ no need algorithm for a translation table

## Group Management Protocol : detects if LAN has any members
for a particular group. (If no members, prune shortest path tree for that group by telling parent).

↳① Router periodically broadcasts query message
↳② Hosts reply with list of groups

## Solution

*Pruning
router tells tree parent to stop forwarding
↳ can be associated with multicast group or with source and group*

Simplest: flood packets from source to network ⟹ if router not seen it forward to all interfaces except incoming one

Clever: reverse path forwarding :- ① forward packet from S to all interfaces iff packet arrives on the interface that corresponds to shortest path to S.
↳ need to know shortest path ⟶ need routing table.
↳ doesn't work if routers do not support multicast
② No need to remember past packets
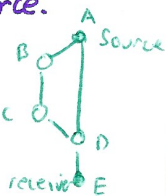③ Don't always need to forward all packets

Better: Don't send packet downstream if not on shortest path from downstream route to the source.

↳ can rejoin though (if C wants to rejoin)
↳ IGMP lets C know of host's interest
↳ ~~When~~ C sends join message to B which propagates from A



## Multicast Routing Protocol

Interface on shortest path to source depends on whether the path is real or virtual

But need to discover shortest paths only taking into account multicast capable routes
↳ DVMRP - Distance Vector Multicast Routing Protocol
Used in conjunction with (1) flood and prune, (2) reverse path forwarding,
(3) Explicit join messages to reduce join latency

MOSPF : routers flood group membership information with LSA
↳ each router independently computes shortest path tree with only multicast-capable routes
(no need to flood and prune
↳ Complex: → ① Must interact with external and summary records
↳ ② Need storage per group per link
↳ ③ Need to compute shortest path tree

## Core Based Trees : coordinate multicast with core router : (1) host sends join request to core router, (2) routers along path mark incoming interface for forwarding. Named as rendezvous point - periodically sends 'I am alive' messages. Leaf routers set timer on receipt - if timer goes off, send request to another rendezvous point.

## Routing for mobile hosts: issues

① Location
② Routing

### Cellular routing

↗ if you reduce size of cell, bandwidth increases

↳ System knows which cell you are seen in ⇒ send message in the backbone
↳ Each cell phone has a global ID that it tells MTSO when turned on
↳ mobile routing in the internet is very similar to mobile telephony but outgoing traffic does not go through home
↳ remote MTSO tells home MTSO then to closest base ⇒ new MTSOs added as load increases
↳ registration packets instead of slotted ALOHA – passed onto home address agent.
↳ Old care-of-agent forwards packet to new care-of-agent
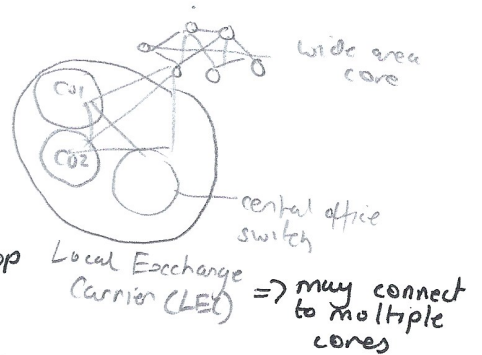
### Problems
① Security
② Loops

## Routing in telephone networks

Old was lines to every town and would ring every town

Algorithm: ① If endpoints within same CO, directly connect
② If call is between COs in same LEC, use one-hop path between COs
③ Else send call to one of the cores
④ Only decision at toll switch
↳ max two-path hop  ] – need to decide which two-path hop path

wide area core

central office switch

Local Exchange Carrier (LEC) ⇒ may connect to multiple cores
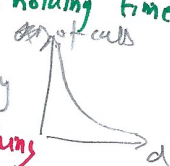
### Features of telephone network routing

① Stable Load
② Reliable Switches
③ Single organization controls entire core
④ Very highly connected network
⑤ Connections require resources

### Statistics
people make calls independently
↳ Poisson call arrival – independence assumption
↳ Exponential call holding time ⇒ long calls less likely
↳ not true sadly

Goal: minimise call blocking

Simplicity – historical necessity, but requires:
↳ (1) reliability in every component
↳ (2) logically fully-connected core.

## Dynamic Nonhierarchical Routing
– leads to metastability : simplest core routing protocol – accept call if one-hop path is available, else drop.

↳ divides day into around 10-periods
↳ in each period, each toll switch assigned primary one-hop path + alternatives
↳ overflow to alternative if necessary
↳ drop only if alternate paths are busy
↳ rely on accurate predictions

every ~~week~~ week

↳ burst of activity can cause network to enter metastable state
↳ high blocking prob even with a low load

### Real-time Network Routing

↳ Each toll switch maintains list of lightly loaded links
↳ Intersection of source and destination lists gives set of lightly loaded paths

↳ Removed by trunk reservation
↳ Trunk Status Map Routing – updates measurements once an hour

no centralized control

### Dynamic Alternative Routing:
↳ whenever link is saturated, use alternative node (tandem)

### Fixed Tandem
For any pair of nodes, assign fixed node k as tandem
↳ needs careful traffic analysis
↳ inflexible during breakdowns and unexpected bch at tandem.

## Sticky Random Tandem

If no free circuit along $(i,j)$, a new call is routed through a randomly chosen tandem $k$ → as long as it does not fail

    ↳ if fails, call is lost and new tandem is selected

Decentralized and flexible with no pre-analysis of traffic required.

$$p_k(i,j) = \text{proportion of calls between } i \text{ and } j \text{ which go through } k$$

$$q_k(i,j) = \text{proportion of calls that are blocked}$$

$$p_a(i,j)\, q_a(i,j) = p_b(i,j)\, q_b(i,j)$$

Assign different frequencies to different tandems

Trunk Reservation: tandem



accepts to forward calls if it has
free capacity more than $R$.     → penalize two link calls, at least
                                      when lines are busy

Erlang's Bound: node connected to $C$ circuits
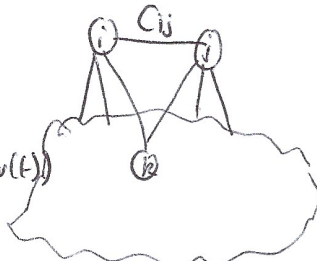         Arrival: poisson with mean $v$

      ↗ don't memorise this.

Expected value of $i$-th call blocking $E(v,C) = \dfrac{v^C}{C!} \left[ \displaystyle\sum_{i=0}^{C} \dfrac{v^i}{i!} \right]^{-1}$

### Max-Flow Bound

mean load $= v_{ij}$



$$E\left[ \sum_{i<j} n_{ij}(t) \right] \le f(v(t))$$

$f$ is solution to LP: $\max \displaystyle\sum_{i<j} \left( x_{ij} + \sum_{j, \, i \neq k} x_{ikj} \right)$

Extension to DAR:-
    ↳① $n$-link paths
    ↳② Multiple alternatives
    ↳③ Least-busy alternative
    ↳④ Rerouting - call in progress can be rerouted

↳ traffic > capacity for some links
↳ not always feasible to find set of tandems   ↳ for $p < \frac{1}{3}$, successful with $p$ approaching 1
    ↳ Greedy Algorithm
        ↳(i) no saturated links → DONE
        ↳(ii) saturated link and good triangle (one saturated, two not)
            ↳ DONE
            ↳ Add good triangle to list
        ↳(iii) saturated link and no good triangle - fail

tradeoffs between overhead for ① stability and ② simplicity for unfairness

### Flow Control Problem
↳ sender needs to match receiver's rate
↳ occurs at transport or datalink layer

Source                     sink



### Open Loop Flow Control
    ↳ Open Loop: ① source describes desired flow rate, ② network admits call, ③ source sends at this rate

① Call Setup
    ↳(1) network prescribes parameters
    ↳(2) user chooses parameter values
    ↳(3) network admits or denies call

② Data transmission
    ↳(1) User send with parameter range
    ↳(2) Network polices the user
    ↳(3) Scheduling policies gives the user QoS

Problem 1: choosing descriptor at source
    ↳ envelope that constrains worst case behaviour, used as:
        ↳(1) basis for traffic contract
        ↳(2) input to regulator and policer

Descriptor requirements:

↳ (1) Representativity - adequately describes flow
↳ (2) Verifiability
↳ (3) Preservability
↳ (4) Usability

examples: time series of intersarrival times (1, 2, !4), ② peak rate (!1, 2, 3, 4), ③ average rate, ④ Linear Bounded Arrival Process ①

Regulator → timer set on packet transmission and if expires, send packet

### → PEAK RATE

② Highest rate at which source can send data

↳ two ways to compute:

↳ min interpacket spacing for networks with [fixed-size packets]

↳ highest rate over all packets of particular interval   for networks with variable size packets

↳ NB it is sensitive to extremes

### ③ AVERAGE RATE

Rate over some time period hence less susceptible to outliers (window)

↳ JUMPING: over ~~consecutive~~ consecutive intervals of length $\Delta t$, only a bit sent
↳ regulator reinitialized every interval

↳ MOVING: over all intervals of length $t$, only a bit sent
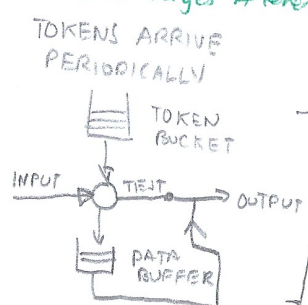↳ regulator forgets packet sent more than $t$ seconds ago

### ④ Linear Bounded Arrival Process : (1) Source bounds - # of bits sent is any time interval by a linear function of time, (2) # of bits transmitted in any active interval of length

↳ Leaky bucket is a regulator for LBAP, token bucket fills up at a rate and stores # tokens < s

$t < rt + s$
↳ $r$ = long term rate
↳ $s$ = burst limit
insensitive to outliers

TOKENS ARRIVE PERIODICALLY

TOKEN BUCKET

INPUT → TEST → OUTPUT

DATA BUFFER

Sum of token and data bucket is what matters

→ Choosing parameters: tradeoff between $r$ and $s$    and $s$ (data / bucket size)
↳ minimal descriptor costs less and doesn't simultaneously have smaller $r$ and $s$
↳ Choosing : → ① Keep loss rate same:
↳ as $s \uparrow$, $r \downarrow$. For reach $r$ we have at least $s$
↳ Then choose the knee of the curve

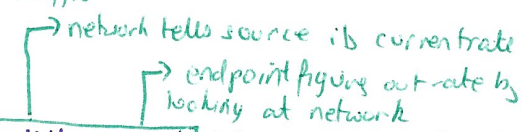LBAP is popular → sort of representative, verifiable, sort of preservable, sort of usable.
↳ NB struggles with multiple time scale traffic

Closed Loop : monitor available bandwidth ⇒ Generalized Processor Sharing and adapt to it. Can lead to loss / delay - normally it can't describe traffic

### Taxonomy

↳ First Generation : ignore network state and match receiver
↳ Second Generation : responsive to state with three choices : (1) Explicit or implicit state measurement, (2) Control (flow control window size or rate), (3) Point of control (endpoint or within network)

→ network tells source its current rate
→ endpoint figures out rate by looking at network

↳ largest number of packets outstanding
↳ if endpoint has sent all packets, it must wait ⇒ slows down its rate
↳ Hence have transmission window (error control and flow control window)
↳ Window vs rate: window has no need for timer and is self-limiting. rate has better control
↳ simpler   → better control
↳ distributes overflow

### Hop-by-hop vs end-to-end

↳ easy to implement with 1st gen flow control at every link
↳ sender matches all servers on path.

## On-Off

Receiver gives ON and OFF signals — used in serial lines and LANs



Source O    Router    O Dest
wait    ack

**Stop and Wait:** (1) Send packet and (2) wait for ack before sending the next packet.
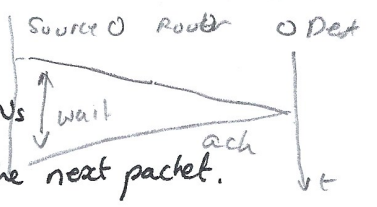
**Static window:** send at most one packet per RTT — works well if $b$ and $R$ are fixed
    ↳ but need to adapt window

**Set window size:** ⟹ bottleneck service rate = $b$ pkts/sec

$$RTT = R$$

flow control window = $w$

sending rate = $w/R$

$$w/R > b \Rightarrow w > bR \Rightarrow \text{OPTIMAL WINDOW SIZE}$$

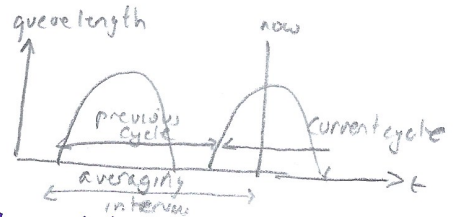## DEC bit Flow Control

**Intuition:** → every packet has a bit in header
   ↳ intermediate routers set bit if queue has built up ⟹ source window is too large
   ↳ sink copies bit to ack
   ↳ if bit set, source reduces window size
   ↳ in steady state, oscillate around optimal size

## Router Actions

① Measure demand and mean queue length of each source ⎤ balance between sensitivity
② Computed over queue regeneration cycles    ⎦ and stability

If mean queue length > 1.0, set bit on sources whose demand exceeds fair share

if > 2.0, set bit on everyone


queue length     now
previous cycle     current cycle
averaging interval → t

## Source Actions

↳ Keep track of bits and can't take control actions too fast ⎤
↳ wait for past change to take effect
↳ measure bits of past + present window size
   ↳ if > 50% set, decrease window, else: increase
   ↳ (additive increase, multiplicative decrease)

**Evaluation:** works with FIFO but requires per-connection state (demand). Works with software but assumes cooperation.

## TCP Flow Control: implicit with dynamic window. Very similar to DEC bit:

   ↳ ① no support from routers
   ↳ ② increase if no loss
   ↳ ③ window decrease on timeout
   ↳ ④ additive increase, multiplicative decrease

SSthresh is called slow start threshold — contains window size in case of (loss)
loss detected by duplicate ACK and timeout

   ↳ DEPENDS ON VERSION OF TCP ⟹ Tahoe vs Reno

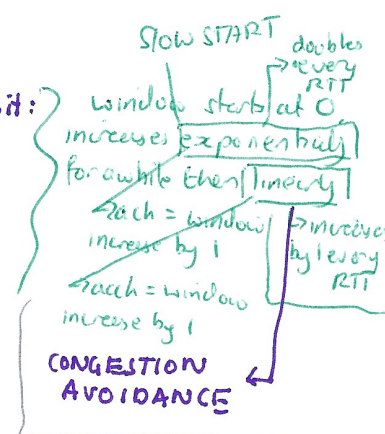   ↳ NB TCP uses implicit measurement of congestion → operates at the cliff

**SLOW START** doubles every RTT
window starts at 0 increases exponentially for a while then linearly
ₐₐₖ = window → increases
increase by 1    by 1 every RTT
ₐₐₖ = window increase by 1

**CONGESTION AVOIDANCE** ←

## TCP Vegas

$$E_{throughput} = \text{transmission\_window\_size} / \text{prop. delay}$$
known     measure smallest RTT
   ↳ can also measure actual throughput ⟹ measure difference and adjust

## NETBLT

→ losses + retransmissions don't affect flow rate

Rate-based flow control scheme, separating error control and flow control. Application data sent as series of buffers (at particular rate) ⟹ Rate = (burst size + burst rate)
   ↳ if received rate > sending rate, increase sending rate

## Packet Pair

Improves basic ideas of NETBLT : → better measurement of bottleneck
  ↳ control based on prediction
  ↳ finer granularity

Spacing between packets at receive = 1 / slowest server    (Assuming bottleneck serve packets in round robin order)

ACKs give time series of service rates in the past
  ↳ use this to predict next rate

  ↳ Exp averager with fuzzy rules to change averaging factor

Predicted rate feeds into flow control equation

$I$ = source rate

$X$ = #packets in bottleneck buffer
$S$ = #outstanding packets
$R$ = RTT
$b$ = bottleneck rate

$$X = S - Rb$$
$$I(k+1) = b(k+1) + (\text{setpoint} - x)/R$$

### Hybrid Flow Control
  ↳ source gets minimum rate, but can use more
  ↳ BUT have problems of both open loop and closed loop

**Design Goals**
① Stability guarantee it
② Transient Response → short settling time, small overshoot
③ Steady-(small) state error
④ Robustness
  ↳(i) Disturbance rejection
  ↳(ii) Sensitivity Low

## FEEDBACK CONTROL THEORY

Why? Allows for interaction with physical environment w.r.t QoS guarantees especially in open, unpredictable environments

Control: applying input to cause system variables to conform to desired values called the reference
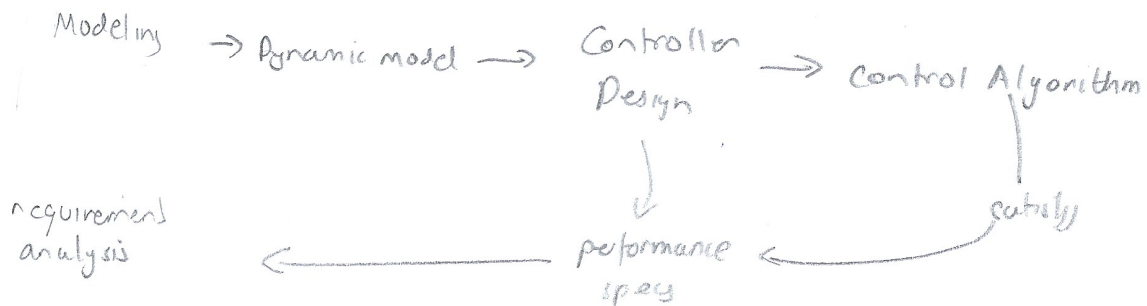
Open-Loop Control : compute control input without continuous variable measurement (but need to know everything accurately & and things not to change)

Close-loop Feedback Control : measure variables and use to control input → more complicated
  ↳ continuously measure and correct. Don't need to know everything and things can change.

### Advantages of Feedback Control Theory
① Adaptive Resource Management Heuristics : laborious iterations for everything
② Ability to handle feedback → it's good at bad conditions, less good at okay conditions
③ Feedback Control Theory ⇒ systematic approach

### Control Design Methodology

Modeling → Dynamic model → Controller Design → Control Algorithm

Requirement analysis ← performance spec ← satisfy

### System Models
  ↳① Linear vs non-linear
  ↳② Deterministic vs Stochastic
  ↳③ Time-invariant vs Time-varying
  ↳④ Continuous time vs Discrete time
  ↳⑤ System ID vs First Principle

### Dynamic Model.
Characterize relationships using differential equations in either time domain or frequency domain ⇒ Fourier Transform to get between

Controller definition :
$$a_2 \ddot{y}(t) + a_1 \dot{y}(t) + a_0 y(t) = b_1 \dot{u}(t) + b_0 u(t)$$

Transfer function (f-domain)
$$G(s) = \frac{b_1 s + b_0}{a_2 s^2 + a_1 s + a_0} = \frac{c_1}{s - p_1} + \frac{c_2}{s - p_2}$$

## Model Differential Equation

$U =$ utilisation

Error : $E(t) = U_s - U(t)$

$$U(t) = \int_{s=0}^{t} (R_a(s) - R_c(s)) \, ds$$

N.B. Convolution is very simple, just multiplication in the frequency domain

## Laplace Transform — similar to Fourier Transform

↳ translation from time domain to $s$ (frequency)

$$F(s) = L[f(t)] = \int_{0^-}^{\infty} f(t) e^{-sp} \, dt \qquad \text{where } s = \sigma + i\omega \text{ (complex variable)}$$

$$a_2 \ddot{y}(t) + a_1 \dot{y}(t) + a_0 y(t) = b_1 \dot{u}(t) + b_0 u(t)$$

$$\iff Y(s) = \frac{b_1 s + b_0}{a_2 s^2 + a_1 s + a_0} \cdot U(s)$$

### Examples

Impulse : $f(t) = \delta(t) \iff F(s) = 1$

Step signal: $f(t) = a_0 \mathbb{1}(t) \iff F(s) = 1/s$

Ramp signal: $f(t) = a_0 t \iff F(s) = a/s^2$

Exp signal : $f(t) = e^{at} \iff F(s) = 1/(s-a)$

Sine signal : $f(t) = \sin(at) \iff F(s) = a/(s^2 + a^2)$

LINEARITY :
$$L(af(t) + bg(t)) = aLf(t) + bLg(t)$$

DIFFERENTIATION
$$L \frac{df(t)}{dt} = sF(s) - f(0_-)$$

INTEGRATION
$$L\left(\int_t f(s)\,ds\right) = F(s)/s$$

### Transfer Function : models a linear time-invariant system

$$G(s) = Y(s)/U(s) \implies Y(s) = G(s)U(s)$$

$$U(s) \longrightarrow G(s) \longrightarrow Y(s)$$

### Poles and Zeros

$$F(s) = \frac{b_m s^m + b_{m-1} s^{m-1} + \dots + b_0}{a_n s^n + a_{n-1} s^{n-1} + \dots + a_0} = K \frac{\prod_{i=1}^{m}(s - z_i)}{\prod_{i=1}^{n}(s - p_i)} = \frac{C_1}{s - p_1} + \dots + \frac{C_n}{s - p_n}$$

$$\implies f(t) = \sum_{i=0}^{n} C_i e^{p_i t}$$

where $p$ are the poles of the function and decide the system behavior
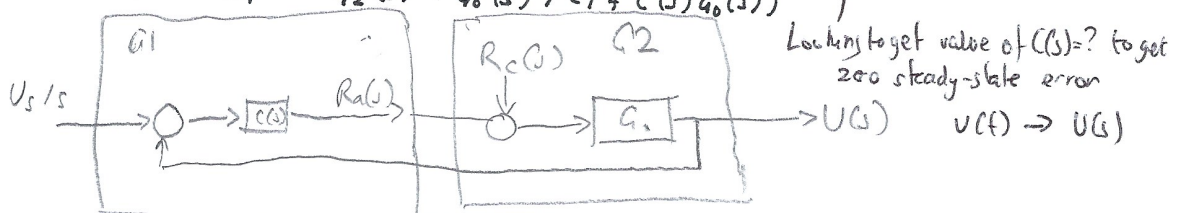
### Transfer Function & Block Diagram

CPU model: $U(t) = \int_{s=0}^{t} (R_a(s) - R_c(s))\,ds \iff U(s) = \frac{R_a(s) - R_c(s)}{s} \iff G_0(s) = \frac{1}{s}$

(set-point)    task

Inputs : reference $U_s(s) = U_s/s$ ; completion rate $R_c(s)$

let $K = C_s G_0(s)$

Pole $P_0 = -K < 0 \iff$ system is BIBO stable iff $K > 0$

Close-loop system transfer functions have two possible inputs :

   ↳ $U_s(s)$ as input : $G_1(s) = C(s)G_0(s)/(1 + C(s)G_0(s))$

   ↳ $R_c(s)$ as input : $G_2(s) = G_0(s)/(1 + C(s)G_0(s))$

Looking to get value of $C(s) = ?$ to get zero steady-state error



$v(t) \to U(s)$

Output: $U(s) = G_1(s) U_s/s + G_2(s) R_c(s)$

## BIBO STABILITY : bounded input results in bounded output. LTI system is BIBO stable if poles of transfer function in LHP ($\forall p_i, Re[p_j] < 0$)

**Steady-State Error**   $e_{ss} = \lim_{t \to \infty} e(t) = \lim_{t \to \infty} (r(t) - y(t))$

ref input ↑    ↑ system output

↳ final value theorem:

$$\lim_{t \to \infty} f(t) = \lim_{s \to 0} sF(s)$$

$$e_{ss} = \lim_{t \to \infty} e(t) = \lim_{s \to 0} sE(s)$$

## Robustness

① Disturbance Rejection: steady state error caused by external disturbances
↳ reg Dos

② Sensitivity : relative change in steady-state output / relative change of system parameter

## Proportional - Integral - Derivative Control

↳ may have non-zero steady state error

↳ Proportional Control: $x(t) = Ke(t) \iff C(s) = k$

↳ Integral Control : $x(t) = KK_i \int_0^t e(\tau) d\tau \iff C(s) = \frac{KK_i}{s}$
   ↳ improves steady state tracking

↳ Derivative Control : $x(t) = KK_d \dot{e}(t) \iff C(s) = KK_d s$
   ↳ may improve stability and transient response

## PI Controller: stability

↳ $r_a(t) = K(e(t) + K_i \cdot \int_t e(t) dt)$ ; $C(s) = K(1 + K_i/s)$

Transfer Functions :

↳ $U_s/s$ as input : $G_1(s) = (Ks + KK_i)/(s^2 + Ks + KK_i)$

↳ $R_c$ as $C(s)$ as input: $G_2(s) = s/(s^2 + Ks + KK_i)$

Stability : poles $Re[P_0] < 0 \iff$ system is BIBO stable iff $K > 0$ & $K_i > 0$

↳ Steady-State Error

Completion rate $R_c(s) = R_c/s$

SystemResponse = $U(s) = \frac{U_s G_1(s)}{s} + \frac{R_c G_2(s)}{s} = \frac{(KU_s + R_c)s + KK_iU_s}{s(s^2 + Ks + KK_i)}$

$e_{ss} = 0$ hence PI control accurately achieves desired response

## Discrete Control & Modelling

z-transform : $f(k) \to F(z)$

$$F(z) = Z[f(k)] = \sum_{k=0}^{\infty} f(k) z^{-k}$$

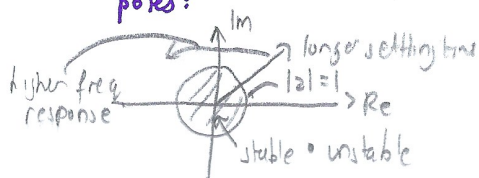Output in $m^{th}$ sampling window $= V(m) = a_1 V(m-1) + a_2 V(m-2) + b_1 U(m-1) + b_2(m-2)$

$U(m)$ is input in $m^{th}$ sampling window

$V(z) = a_1 z^{-1} V(z) + a_2 z^{-2} V(z) + b_1 z^{-1} V(z) + b_2 z^{-2} U(z)$

Transfer function $G(z) = (b_1 z + b_2)/(z^2 - a_1 z - a_2)$

## Root Locus Analysis
① stability boundary : $|z| = 1$, (2) settling time = dist from origin, (3) speed =
↳ Effect of discrete poles : locations relative to Im axis



higher freq response
longer settling time
$|z| = 1$
→ Re
stable · unstable
Im

**Advanced Control**

① Robust noise → tolerance to noise

② Adaptive control

③ MIMO control

④ Stochastic control - minimise variance

⑤ Optimal control - minimize cost function of energy and error

⑥ Non-Linear System

## Issues

① Systems are non-linear
② First-principles modelling is difficult
③ Tough to map control objectives to feedback control loops
④ Deeply embedded networking → decentralized

## Router Queuing Behaviour in packet switched networks

↳ traditionally scheduling is very simple — better to increase speed of link and router.

↳① DATA TRANSFER : individual packets, no recognition of flows and no signalling
↳② FORWARDING: based on per-datagram using forwarding tables
    ↳no priority ster system
↳③ TRAFFIC PATTERNS    ↳add priority system to support QoS

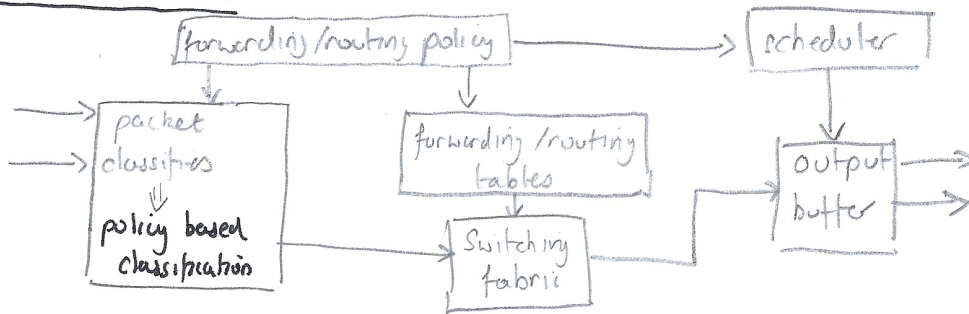↳ deal with congestion by tail drop where packets are lost abd as you put onto output queue

↳ Can use FIFO queue not FCFS    → local/global, protect flow from another
                                               MAX-MIN

## SCHEDULING

→ requirements: ↳① Ease of implementation
                →② Fairness and protection → per flow bounds, deterministic,
                ↳③ Performance bounds    probabilistic, data rate, delay,
                ↳④ Admission control             jitter loss

Scheduler ① decides service order based on policy and (2) manages service (output) queues.

## Router Schematic



Scheduler decides which output queue is serviced next

## FCFS Scheduling : packets queued to outputs in order they arrive with no differentiation and no notions of flows

work-conserving ↴
scheduler   → not idle if packets waiting

## Conservation Law

Capacity is sum of utilisation and delay: $C = \sum_{n=1}^{N} \rho_n q_n$    → mean packet rate
                                                → mean per packet service
Tradeoff between throughput and latency    ↳ $\rho_n = \lambda_n \mu_n$           rate
for different flows
                                ↳ mean link utilisation

Non-work-conserving schedules - receiver doesn't buffer receiving things
↳ Allows smoothing of packet flows (can be idle even if waiting packets)
                             ↳ wait until packet is eligible for transmission
  ↳ Less jitter               ↳ higher end-to-end delay ↳ fixed time per router or
  ↳ Downstream traffic more predictable   ↳ Complex in practise        fixed time across network
  ↳ Less buffer space

**Max-Min Fair Share Criteria**: flows allocated resource in order of increasing demand
↳ they get no more than they need

↳ weighted max-min fair share possible. If fair, it provides protection

$M_n$ = actual resource allocation to flow $n$
$= \min(x_n, M_n)$
↳ resource demand by flow $n$
↳ resource available to flow $n$

**Dimensions**
① Priority Levels
② Work conserving or not
③ Degree of aggregation
↳ flow granularity
↳ per application flow
↳ per user?
↳ per end-system?
↳ cost vs control
④ Servicing within a queue

Capacity of resource (max resource) $= M_n = \dfrac{C - \sum_{i=1}^{n-1} m_i}{N - n + 1}$

$1 \leq n \leq N$
$\phi(n)$ - weight given to flow $n$
$S(i,j,t)$ $1 \leq i \leq N$
↳ service to flow $i$ in interval $[t_j, t]$ flow $i$ has non empty queue

$\dfrac{S(i,j,t)}{S(j,j,t)} \geq \dfrac{\phi(i)}{\phi(j)}$

**Simple Priority Queuing**: higher priority queues are serviced first
↳ not max-min fair - starvation

~~Garanteed~~ **Generalised Process Sharing**

Work conserving with max-min fairness ⇒ can provided weighted max-min fair share.
**NOT IMPLEMENTABLE** ↳ round robin with infinitessimally low quantum
↳ used as a comparison (others emulate GPS).

protection ⇒ $(n-1)$ priority max

$RFB = \left| \dfrac{S(i,j,t)}{g(i)} - \dfrac{S(j,j,t)}{g(j)} \right| \Rightarrow$ ↳ ① Relative Fairness Bound - fairness of scheduler with respect to other flow it is servicing

$AFB = \left| \dfrac{S(i,j,t)}{S(i)} - \dfrac{g(i,j,t)}{g(i)} \right|$ ↳ ② Absolute Fairness Bound - fairness of scheduler compared to GPS for same flow

↳ GPS service for flow $i$ in $[s,t]$

$g(i) = \min \{g(i,1), \ldots, g(i,K)\}$

$g(i,k) = \dfrac{\phi(i,k) \, r(k)}{\sum_{j=1}^{N} \phi(j,k)}$ — service rate of router $k$

$i$ = flow number
$k$ = router number

**Weighted Round Robin**: Queues visited round robin in proportion to weights assigned
↳ different mean packet sizes ⇒ this is unpredictable and may cause unfairness
↳ service is fair over long timescales

If we instead compute packet size on the fly, we have **Deficit Round-Robin**
↳ each queue has deficit counter initially at zero
Scheduler attempts to serve one quantum of data from non-empty queue
↳ packet at head served if size ≤ quantum + dc
↳ else: dc += quantum
↳ set to max expected packet size

$RFB = 3T/r$
max packet service time ↗ ↳ link rate

## Weighted Fair Queuing

: calculate idealised for each round for each packet size → calculated per flow. Need data store per flow and need info of destination. Hence, lookup necessary every time

**Buffer Drop Policy:** if full queue, drop in order of decreasing finish time.

↳ GPS emulation to get finish-numbers for packets in queue
  ↳ serves packets bit-by-bit round-robin
  ↳ time packet would have completed service under (bit by bit) GPS. Tags finish-number for each packet. Smallest finish number served first

↳ ROUND NUMBER → NB need to be computed every time a packet arrives or leaves
  ↳ Execution of round by bit-by-bit round-robin server. Finish number calculated from round number
    ↳ if queue empty: $fN = n(\text{bit in packet}) + \text{round-number}$
    ↳ else : $fN = \max(fN\text{s in queue}) + n(\text{bit in packet})$

$F(i,k,t)$ = Finish-number for packet $k$ on flow $i$ at time $t$

$P(i,k,t)$ = size of packet $k$ on flow $i$ arriving at time $t$

$R(t)$ = round-number at time $t$ — depends on number of active flows and their weights

$\phi(i)$ = weight given to flow $i$

$F(i,k,t) = \max\{F(i,k-1,t), R(t)\} + P(i,k,t)$

$F_\phi(i,k,t) = \max\{F_\phi(i,k-1,t), R(t)\} + \frac{1}{\phi(i)} P(i,k,t)$

## Class-Based Queuing

Assign a capacity and priority to each node, which can borrow spare capacity from a parent hence meaning fine-grained flows are possible.

## Queue Management

↳① Ensuring buffers are available; ie memory management
↳② Organising packets within a queue
↳③ Packet dropping when a queue is full
↳④ Congestion control
  ⇒① Dealing with misbehaving sources
  ↳② Source synchronisation
  ↳③ Routing instability

↳) Packet Dropping Policies
  ① Drop from tail
  ② Drop from head
  ③ Random drop
  ④ Flush queue
  ⑤ Intelligent drop

**End system reaction to packet drops**
①TCP - works well
②UDP - hurts real time adaptive flows

Can adapt for real time flows - use ECN or add multicast to get over issue of packet drop for TCP

## Random Early Detection

Idea is to spot congestion before it happens, some can drop packets simplying preemptive — can mark offending packets which are more likely to be dropped

congestion signal, stopping real congestion

$p(\text{packet drop}) \propto \text{queue length}$   exp average: smooths reaction to short bursts

(18)

# Datacenter Networks (QJump)

Use: ① commodity hardware ⟹ recent
② Static Network topology
③ Policies are under single admin domain
④ Cooperation
⑤ ~~Stati~~ Statistically multiplexed from internet

constraints
① unmodified apps
② Unmodified kernel code

NB, no guarantees about latency
↳ Hadoop (one app) causes problems - queues caused by this slow down other applications

Solve problems by applying queuing concepts in datacenter

## Delays
↳ Queuing Delay $D_q$
↳ Servicing Delay $D_s$  } causes queuing delay
↳ if we can bound servicing delay, rate limit hosts so we don't get query delay

n of sending hosts
$\tau = \delta \times \dfrac{P}{R}$ — packet size
— edge speed

So: →① Network idle
② Hosts send $\leq P$ bytes
③ Wait $(n \times P/R)$ secs
④ Goto 1
} Network epoch

Network epoch $= (2)n \times \dfrac{P}{R}$

mesochronous compensation

Can also add in priorities then which allows for queue jumping
↳ hardware

Throughput $= \dfrac{R}{2n}$

## Optimization Based Routing

Framework:
↳① W set of source-destination pairs
↳② rw : rate of sd pair w
↳③ Pw : set of paths between sd pair w
↳④ $x_p$ : flowrate on path p.
↳⑤ $c_{ij}$ : capacity of link i,j

$P_{ij}(F_{ij})$

$c_{ij}$ on x, $F_i$

In routing problem, rw given
In rate control problem, rw variable

Question is how to set rates on various paths
↳ they can by set ① centrally, at ② edges or ③ at routers
↳ minimizing either (i) system wide delay or maximize system wide utility

$F_{ij} = \sum x_p$   **ROUTING OPTIMIZATION PROBLEM**
↳ all paths p crossing link i,j   choose $\{x_p\}$

Aim: minimize $\sum\limits_{\substack{all\ links\\ ij}} D_{ij}\left( F_{ij} = \sum\limits_{all\ link} \dfrac{F_i}{c_{ij} - F_{ij}} \right)$ subject to $\sum\limits_{\substack{all\ paths\\ in\ Pw}} x_p = \forall w \in W,$

~~$x \geq 0$~~ all incoming flow routed to destination

At optimum: $\dfrac{\partial P(x^*)}{\partial x_p} = \dfrac{\partial P(x^*)}{\partial x_p}$  ∀ nonzero $x_{p1}, x_{p2}$ in $P_w$

↳ Algo: ① ∀ sd pairs w : evaluate $\dfrac{\partial P(x)}{\partial x_p} \ldots \dfrac{\partial D(x)}{\partial x_p}$ for all hw pairs for w

② Move small amount of flow to paths with min marginal increase from other paths and repeat until all equal.

Can treat resource allocation as an optimization problem

Model

$L(s)$ = link used by source $s$
$U_s(x)$ = utility if source rate $= x_s$

**USER PROBLEM**

$$\max \left(\frac{w_s}{p_s}\right) - w_s \quad w_s \geq 0$$

System problem maximise $\sum_s U_s(x_s)$ subject to $\sum_{s \in S(l)} x_s \leq c \quad \forall l \in L$
$x_s \geq 0$

bandwidth less than capacity

$\quad \hookrightarrow x_s = \dfrac{w_s}{p_s} \Rightarrow$ unit time

$p_s \Rightarrow$ charge per unit flow for source
$\hookrightarrow$ COST

$\hookrightarrow$ network wants to maximise : $\sum_s w_s \log x_s$ (NETWORK PROBLEM)
$\hookrightarrow$ converges to relaxation of network problem

Idea is $\exists$ prices $p_s, x_s, w_s$ , $w_s = p_s x_s$ st $\{w_s\}$ solves user problem
$\{x_s\}$ solves network problem
$\{x_s\}$ solves system problem

$x_s$ is fair if feasible and for $\{x_s^*\}$
$\hookrightarrow$ any other feasible :

$$\sum_{s \in S} \frac{x_s^* - x_s}{x_s} \leq 0$$

System properties
① Convergence
② Achieve objective
③ Benchmarks $\rightarrow$ ① $\max \sum_i U_i(x_i)$ st $Rx \leq c$
② 
③ Var $x, R$
④ Utility gap between joint system and benchmark

System Design

Putting together network resources to extract max usage

$\rightarrow$ computation, storage and transmission resources.

$\hookrightarrow$ in any system, some resources are more freely available than others

some combination of:
① time ② space,
③ computation, ④ money,
⑤ labour

Aim. To maximise performance metrics given resource constraints

Metrics

① Time : response time, throughput, (degree of parallelism = response time × throughput)
② Space : limit of available space (kB) , and bandwidth (kb/s)
③ Computation : processing / unit time
④ Money : £

Social constraints

① Standards
② Market Requirement

Scaling is a design constraint (hard to measure) but very necessary for success. Also important for economies of scale.

**Bottleneck:** most constrained element in a system.
- removing bottleneck improves performance - but creates another bottleneck.
- Aim is generate balanced system where all resources are simultaneously bottlenecked

**Techniques to trade off one resource for another**
- ① Multiplexing - trades time and space for money
  - increases response time and less space but costs less → economies of scale
  - examples: ① Multiplexed links
             ② Shared memory
  - server controls access to shared resource using schedule to resolve contention

**Statistical Multiplexing**
Resource has capacity C, shared by N identical tasks, each task requires capacity c. If $Nc \leq C$ then resource underloaded.
- if at most 10% tasks active $C \geq Nc/10$ enough
- statistical multiplexing gain

- **Spatial multiplexing:** expect only a fraction of tasks to be simultaneously active.
- **Temporal multiplexing:** expect task to be active only part of the time

② **Pipelining**
Break up task into independent subtasks - optimal if all subtasks take the same amount of time. But, ∃ dependencies. Pipeline is a special case of serially dependent subtasks
- best decomposition - degree of parallelism = $\frac{R}{S}$, maximised when $N = \frac{R}{S}$.
  hence $S = R/N$

③ **Batching**
Group tasks together to amortize overhead - only works when overhead for N tasks < N times overhead for one task. Time taken to accumulate batch shouldn't be too long.
- tradeoff of reduced overhead for longer worst case response time and increased throughput

- Can also exploit locality, spatial and temporal through caching ④
- 80/20 rule: 80% of time spent in 20% of code
  - this is part that we should optimize
  - Measure using [Amdahl's Law]

Execution time after improvement = $\frac{(\text{execution affected by improvement})}{\text{amount of improvement}}$

$+$ execution unaffected.

⑤ Hierarchy: decomposition of system into smaller pieces that depend only on parent ~~for~~ for proper execution

    ↳ No single point of control

    ↳ Higly scaleable

    ↳ But, leaf to leaf communication can be expensive

⑥ Binding and indirection

    ↳ translation from an abstraction to an instance

    ↳ <u>Indirection</u> : can bind automatically if translation table is stored in a well known place

⑦ Virtualization : combination of indirection and multiplexing

    ↳ refer to virtual instance that gets matched to instance at run time.

    ↳ Can cleanly and dynamically reconfigure - build as if real resource is available

⑧ Randomization

    ↳ allows us to break a tie fairly

Soft State:

    ↳ memory in a system that influences future behaviour.

    ↳ delete the state on a timer - refresh if you want to keep it.

    ↳ Automatically cleans up after a failure but increases bandwidth requirement

    ↳ Important to use explicit state exchange where network elements need to exchange state.

⑨ Hysteresis: when need to detect if value above or below threshold where the variable fluctuates near the threshold, we can use state-dependent threshold (hysteresis)

Data vs Control

    ↱ data path vs control path

    ↳ divide actions that happen once per data transfer and once per packet. - increase throughput by minimizing actions in data path.

    ↳ But keeping control information in data element has it's advantages

        ↳ per packet QoS

⑩ Extensibility : good idea to leave hooks to allow for future growth

<u>Tuning Existing Systems</u>

① Measure

② Characterise workload

③ Build system model

④ Analyse

⑤ Implement