Route or path: criteria of goodness

- hop count
- delay
- bandwidth
- loss rate
- policy
- manual configuration

Composition of performance metric:

- \longrightarrow quality of end-to-end path
- additive: hop count, delay
- min: bandwidth
- multiplicative: loss rate

Goodness of routing:

- \longrightarrow assume N users or sessions
- \longrightarrow suppose path metric is delay

Two approaches:

- system optimal routing
 - \rightarrow choose paths to minimize $\frac{1}{N} \sum_{i=1}^{N} D_i$
 - \rightarrow good for the system as a whole
- user optimal routing
 - \rightarrow each user i chooses path to minimize D_i
 - \rightarrow selfish route selections by each user
 - \rightarrow end result may not be good for system as a whole

Pros/cons:

- system optimal routing:
 - good: minimizes delay for the system as a whole
 - bad: complex and difficult to scale up
- user optimal routing:
 - good: simple
 - bad: may not make efficient use of resources
 - \rightarrow low utilization
 - \rightarrow tragedy of commons

Two pitfalls of user optimal routing:

- fluttering or ping pong effect
 - \rightarrow induced synchronization
- Braess paradox
 - \rightarrow adding more resources (extra link) can make things worse

Increasing resource should improve things but has the opposite effect

- \longrightarrow D. Braess (1969)
- → paradox possible due to user optimal routing
- → cannot arise in system optimal routing

Modus operandi of the Internet: user optimal routing

- \longrightarrow simplicity wins the day
- ... conceptually related problem in operating systems
 - → page fault replacement algorithms
 - \longrightarrow Belady's anomaly

Algorithms:

Find short, in particular, shortest paths from source to destination.

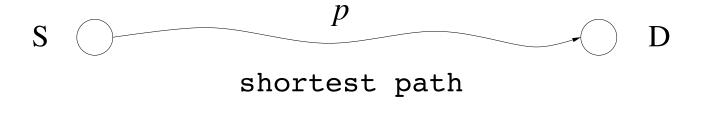
Key observation on shortest paths:

- Assume p is a shortest path from S to D $\to S \stackrel{p}{\leadsto} D$
- \bullet Pick any intermediate node X on the path
- \bullet Consider the two segments p_1 and p_2

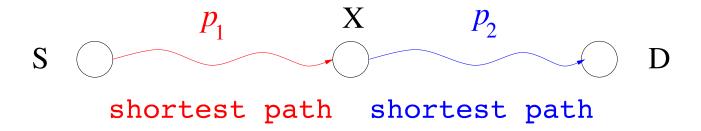
$$\rightarrow S \stackrel{p_1}{\leadsto} X \stackrel{p_2}{\leadsto} D$$

- The path p_1 from S to X is a shortest path, and so is the path p_2 from X to D
 - → leads to Dijkstra's algorithm

Illustration:







 \rightarrow suggests algorithm for finding shortest path

Dijkstra: single-source all-destination

Features:

- running time: $O(n^2)$ time complexity
 - $\rightarrow n$: number of nodes
- if heap is used: $O(|E| \log |V|)$
 - $\rightarrow O(n \log n)$ if |E| = O(n)
- can also be run "backwards"
 - \rightarrow start from destination D and go to all sources
 - \rightarrow a variant used in inter-domain routing
 - \rightarrow forward version: used in intra-domain routing
- ullet source S requires global link distance knowledge
 - \rightarrow centralized algorithm (center: source S)
 - \rightarrow every router runs Dijkstra with itself as source
 - \rightarrow lots of broadcast management packets

- Internet protocol implementation
 - \rightarrow OSPF (Open Shortest Path First)
 - \rightarrow also called link state algorithm
 - \rightarrow broadcast protocol
- \bullet builds minimum spanning tree rooted at S:
 - \rightarrow to all destinations
 - \rightarrow if select destination: called multicasting
 - \rightarrow multicast group
 - \rightarrow complexity including group membership management

Distributed/decentralized shortest path algorithm:

 \longrightarrow Bellman-Ford algorithm

Key procedure:

- Each node X maintains current shortest distance to all other nodes
 - \rightarrow a distance vector
- Each node X advertises to neighbors its current best distance estimates
 - \rightarrow i.e., neighbors exchange distance vectors
- ullet Each node X updates shortest paths based on neighbors' advertised information

$$d(X,Z) \leftarrow \min\{d(X,Z), d(Y,Z) + \ell(X,Y)\}$$

→ same update criterion as Dijkstra's algorithm

Features:

- running time: $O(n^3)$, i.e., O(n|E|)
 - \rightarrow parallel/distributed: O(|E|)
- each source or router only talks to neighbors
 - \rightarrow local interaction
 - \rightarrow no need to send update if no change
 - \rightarrow if change, entire distance vector must be sent
- knows shortest distance but not path
 - \rightarrow just the next hop is known
- elegant but additional issues compared to Dijkstra's algorithm
 - \rightarrow e.g., stability
- Internet protocol implementation
 - \rightarrow RIP (Routing Information Protocol)

Data center networks: leaf-spine connectivity

- \rightarrow each leaf switch connected to all spine switches
- \rightarrow use equal-cost multipath routing (ECMP)

Suppose n leaf switches, m spine switches. At leaf switch $i \in [1, n]$:

- IP packet p forwarded to spine switch $j \in [1, m]$ if h(p) = j where hash function
 - \rightarrow input of h: IP source/destination addresses and port numbers
 - → prevent packet reordering issue
- facilitates load balancing
 - \rightarrow TCP incast problem remains

QoS routing:

Given two or more performance metrics—e.g., delay and bandwidth—find path with delay less than target delay D (e.g., 100 ms) and bandwidth greater than target bandwidth B (e.g., 10 Mbps)

- → from shortest path to best QoS path
- → multi-dimensional QoS metric
- \longrightarrow other: jitter, hop count, etc.

How to find best QoS path that satisfies all requirements?

Brute-force

- enumerate all possible paths
- rank them

Is there a poly-time algorithm?

- \rightarrow as of November 2025: unknown
- \rightarrow specifically: QoS routing is NP-hard
- \rightarrow strong belief no fast algorithm

In networking: several problems turn out to be NP-complete

- \rightarrow e.g., scheduling, crypto, . . .
- \rightarrow "P = NP" problem
- \rightarrow one of the hardest problems in science

In practice: doesn't matter too much for QoS routing

- \rightarrow no pressing demand (yet) for very good algorithm
- \rightarrow intra-domain: short paths
- \rightarrow inter-domain: policy routing

Policy routing:

- → meaning of "policy" is not precisely defined
- \longrightarrow almost anything goes

Criteria include:

- Performance
 - \rightarrow e.g., short paths
- Trust
 - \rightarrow what is "trust"?
- Economics
 - \rightarrow pricing
- Geo-politics, etc.

Implementation:

Internet routing protocols:

- RIP: intra-domain, Bellman-Ford
 - \rightarrow also called distance vector routing
 - \rightarrow metric: hop count
 - \rightarrow UDP
 - \rightarrow nearest neighbor advertisement
 - \rightarrow popular in small intra-domain networks
 - \rightarrow e.g., used at Purdue for many years
- OSPF: intra-domain, Dijkstra
 - \rightarrow also called link state
 - \rightarrow metric: average delay
 - \rightarrow directly over IP: protocol number 89
 - \rightarrow broadcasting via flooding
 - \rightarrow popular in larger intra-domain networks

- IS-IS: intra-domain, Dijkstra
 - → directly over link layer (e.g., Ethernet)
 - \rightarrow less complex than OSPF
 - \rightarrow popular in larger intra-domain networks
 - → intra-domain routing: performance based

iBGP (interior BGP): intra-domain protocol that bridges gap with BGP inter-domain routing

→ eBPG (exterior BGP)

BGP (Border Gateway Protocol):

 \longrightarrow inter-domain routing

Autonomous System B Peering Border Routers

- \rightarrow peering between two domains
- \rightarrow typical: customer-provider relationship
- \rightarrow in some cases: A and B are equals (true peers)

- CIDR addressing
 - \rightarrow i.e., a.b.c.d/x
 - \rightarrow Purdue: 128.10.0.0/16, 128.210.0.0/16, 204.52.32.0/20
 - → check at www.iana.org (e.g., ARIN for US)
- Metric: policy
 - \rightarrow e.g., shortest-path, trust, pricing
 - → meaning of "shortest": delay, router hop, AS hop
 - → mechanism: path vector routing
 - \rightarrow BPG update message

BGP route update:

→ BGP update message propagation

BGP update message format:

$$ASNA_k \rightarrow \cdots \rightarrow ASNA_2 \rightarrow ASNA_1$$
; a.b.c.d/x

Meaning: ASN A_1 (with CIDR address a.b.c.d/x) can be reached through indicated path

- \longrightarrow called path vector
- \longrightarrow also AS-PATH

Some AS numbers:

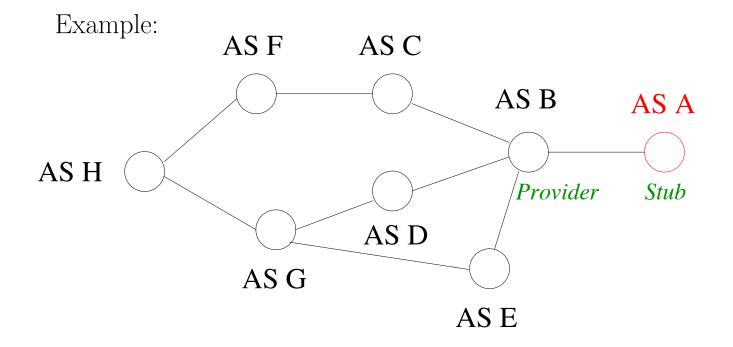
- BBN: 1
- UUNET: 701
- Level3: 3356
- Abilene (aka "Internet2"): 11537
- AT&T: 7018
- Purdue: 17

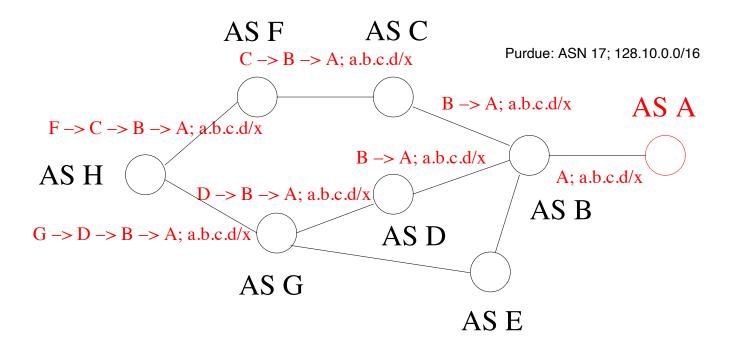
Policy:

- if multiple AS-PATHs to target AS are known, choose one based on policy
 - \rightarrow e.g., shortest AS path length, cheapest, least worrisome
- advertise to neighbors target AS's reachability
 - \rightarrow also subject to policy
 - \rightarrow no obligation to advertise!
 - \rightarrow specifics depend on bilateral contract (SLA)

SLA (service level agreement):

- \longrightarrow bandwidth (e.g., 40 Gbps)
- \longrightarrow delay (e.g., avrg. 20ms US), loss (e.g., 0.01%)
- \longrightarrow peak vs. average
- \longrightarrow pricing
- → availability, among others





Performance:

Route update frequency:

- → routing table stability vs. responsiveness
- \longrightarrow rule: not too frequently
- \longrightarrow 30 seconds
- \longrightarrow stability wins
- → hard lesson learned from the past (sub-second)
- \longrightarrow legacy: TTL

Other factors for route instability:

- \longrightarrow selfishness (e.g., fluttering)
- → BGP's vector path routing: can be unstable
- \longrightarrow more common: slow convergence
- → target of denial-of-service (DoS) attack

Route amplification:

- \longrightarrow shortest AS path \neq shortest router path
- \longrightarrow e.g., may be several router hops longer
- → AS graph vs. router graph
- → policy: company in Denmark

Route asymmetry:

- → routes are not symmetric
- → mainly artifact of inter-domain policy routing
- → various performance implications
- → source traceback

Black holes:

→ persistent unreachable destination prefixes

 \longrightarrow BGP routing problems

→ further aggrevated by DNS