

Route or path: criteria of goodness

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

Composition of performance metric:

→ quality of end-to-end path

- hop count, delay: additive
- bandwidth: minimum
- loss rate: multiplicative

Goodness of routing:

→ assume  $N$  users or sessions

→ suppose path metric is delay

Two approaches:

- system optimal routing

→ choose paths to minimize  $\frac{1}{N} \sum_{i=1}^N D_i$

→ good for the system as a whole

- user optimal routing

→ each user  $i$  chooses path to minimize  $D_i$

→ selfish route selections by each user

→ 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
    - also, fluttering/ping pong effect, Braess paradox

Internet: user optimal routing

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  $D$

$$\rightarrow S \overset{p}{\rightsquigarrow} D$$

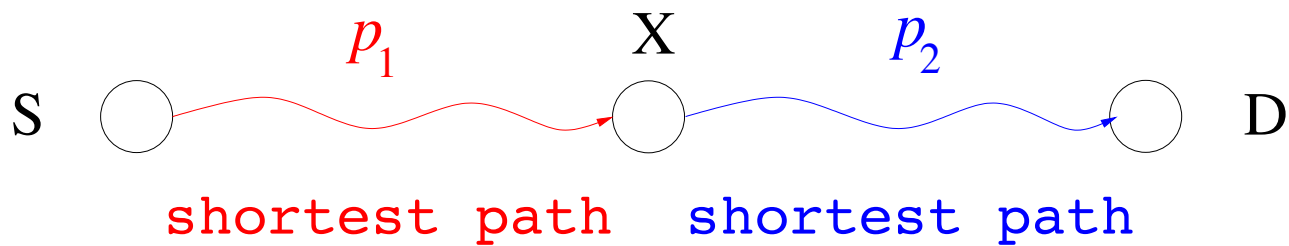
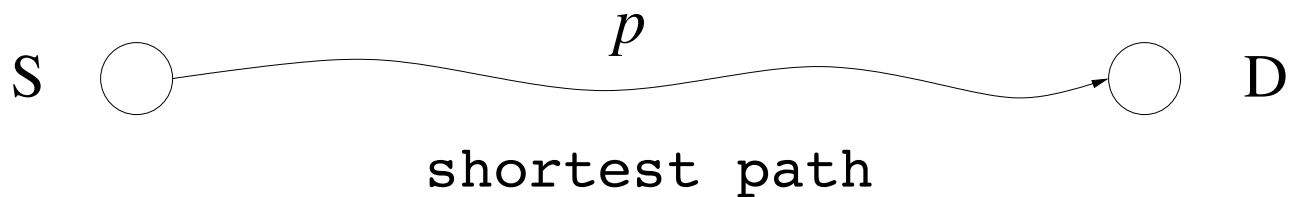
- Pick any intermediate node  $X$  on the path
- Consider the two segments  $p_1$  and  $p_2$

$$\rightarrow S \overset{p_1}{\rightsquigarrow} X \overset{p_2}{\rightsquigarrow} D$$

- 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:



→ suggests algorithm for finding shortest path

Dijkstra: single-source all-destination

Features:

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

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

Distributed/decentralized shortest path algorithm:

→ Bellman-Ford algorithm

Key procedure:

- Each node  $X$  maintains current shortest distance to all other nodes

→ a distance vector

- Each node  $X$  advertises to neighbors its current best distance estimates

→ i.e., neighbors exchange distance vectors

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

Data center networks: leaf-spine connectivity

→ each leaf switch connected to all spine switches

→ 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  $h$  is hash function
  - input of  $h$ : IP source/destination addresses and port numbers
- facilitates load balancing

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
- other: jitter, hop count, etc.

How to find best QoS path that satisfies all requirements?

Brute-force

- enumerate all possible paths
- rank them
- super-exponential

Is there a poly-time algorithm?

→ as of April 2026: unknown

→ technically: QoS routing is NP-hard

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

→ scheduling related, crypto, etc.

→ “ $P \neq NP$ ” problem

→ one of the hardest problems in science

In practice: doesn't matter much for QoS routing

→ no pressing demand for good QoS routing algorithm

→ intra-domain: short paths

→ inter-domain: policy routing

Policy routing:

- meaning of “policy” is not precisely defined
- almost anything goes

Criteria include:

- Performance
  - e.g., short paths
- Trust
  - what is “trust”?
- Economics
  - pricing
- Geo-politics, etc.

Implementation:

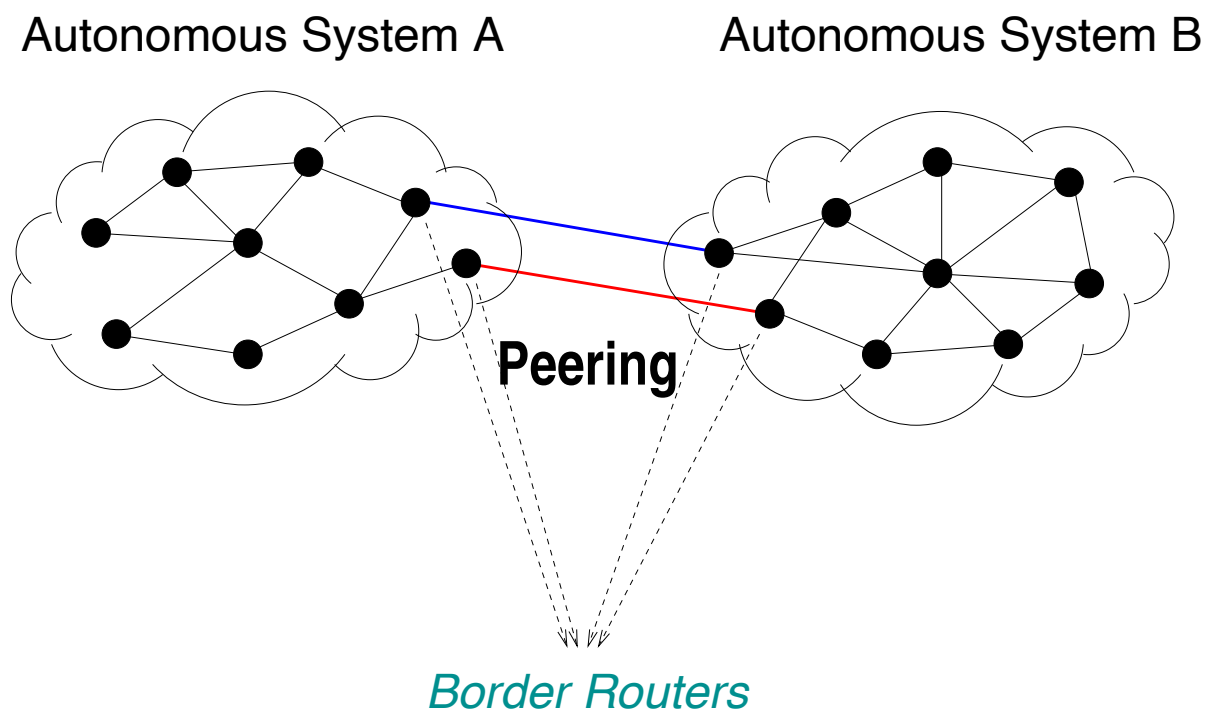
Internet routing protocols:

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

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

BGP (Border Gateway Protocol):

→ inter-domain routing



→ peering between two domains

→ typical: customer-provider relationship

→ in some cases:  $A$  and  $B$  are equals (true peers)

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

BGP route update:

→ BGP update message propagation

BGP update message format:

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

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

→ called path vector

→ also AS-PATH

Some AS numbers:

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

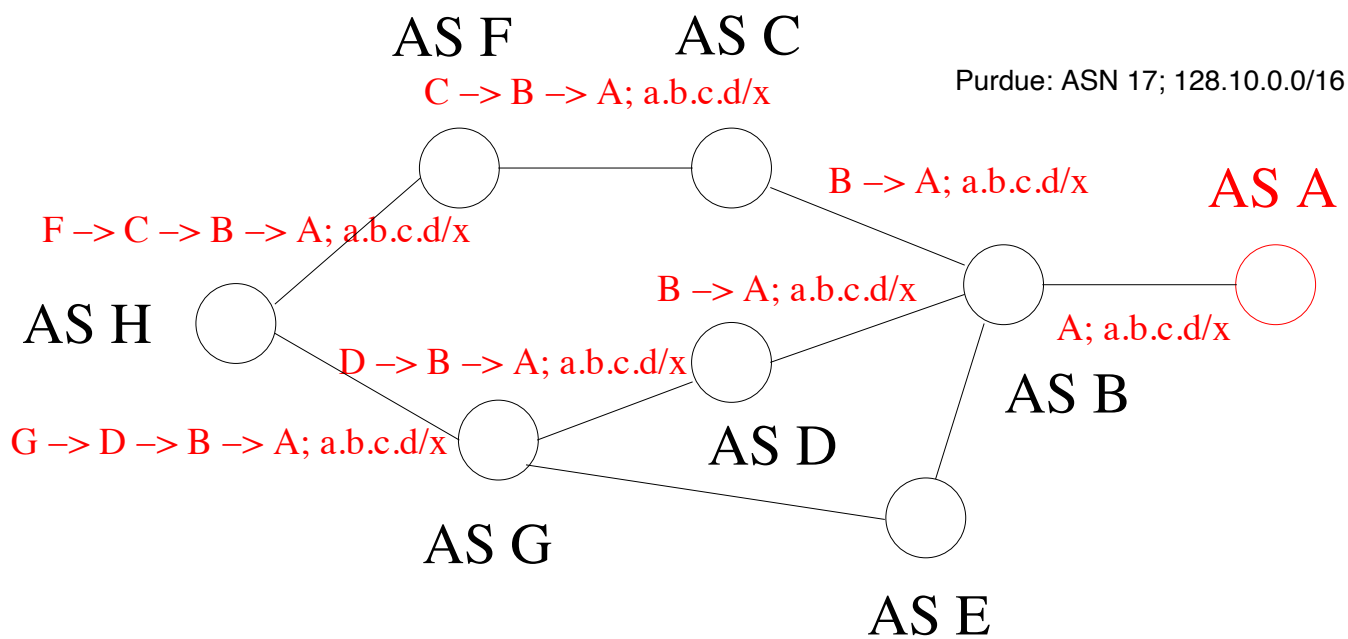
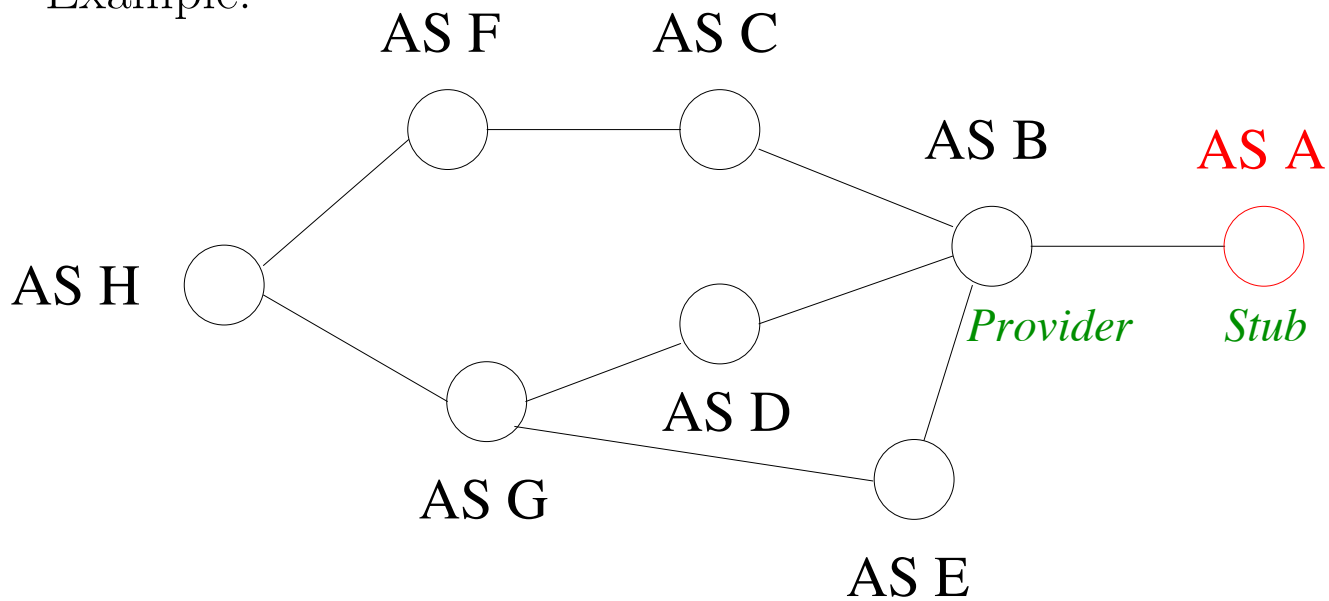
Policy:

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

SLA (service level agreement):

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

Example:



Performance:

Route update frequency:

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

Other factors for route instability:

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

Route amplification:

- shortest AS path  $\neq$  shortest router path
- 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
- BGP routing problems
- further aggravated by DNS