

Declarative Routing: Extensible Routing with Declarative Queries

Boon Thau Loo¹

Joseph M. Hellerstein^{1,2}, Ion Stoica¹, Raghu Ramakrishnan³,

¹University of California at Berkeley, ²Intel Research Berkeley,

³University of Wisconsin-Madison

Motivation

- ◆ Lack of extensibility and flexibility in today's Internet routing
- ◆ Hard to add/improve/update routing protocols

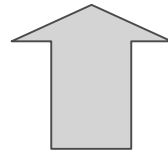
Two “Extremes”:

“Hard-coded” protocols:

- ⊕ Efficiency, safety
- ⊖ Flexibility, evolvability

Active Networks

- ⊕ Flexibility, evolvability
- ⊖ Safety, efficiency



Declarative Routing:

- ⊕ Flexibility, evolvability, safety

Restricted instantiation of Active Networks for the control plane

Key Idea

- ◆ Recursive query language for expressing routing protocols:
 - Datalog: a declarative recursive query language
 - Well-researched in the database community
 - Well-suited for querying properties of graphs

Advantages

- ◆ Expressiveness: Compact and clean representation of protocols
- ◆ Safety: Datalog has desirable safety properties on termination
- ◆ Efficiency: No fundamental overhead when executing standard protocols.

Usage Scenarios

◆ ISP administrators

- Run different protocols for different nodes
- Modify existing protocols in routers

◆ End-hosts

- Set up customized routes for different quality-of-service and policy requirements of applications

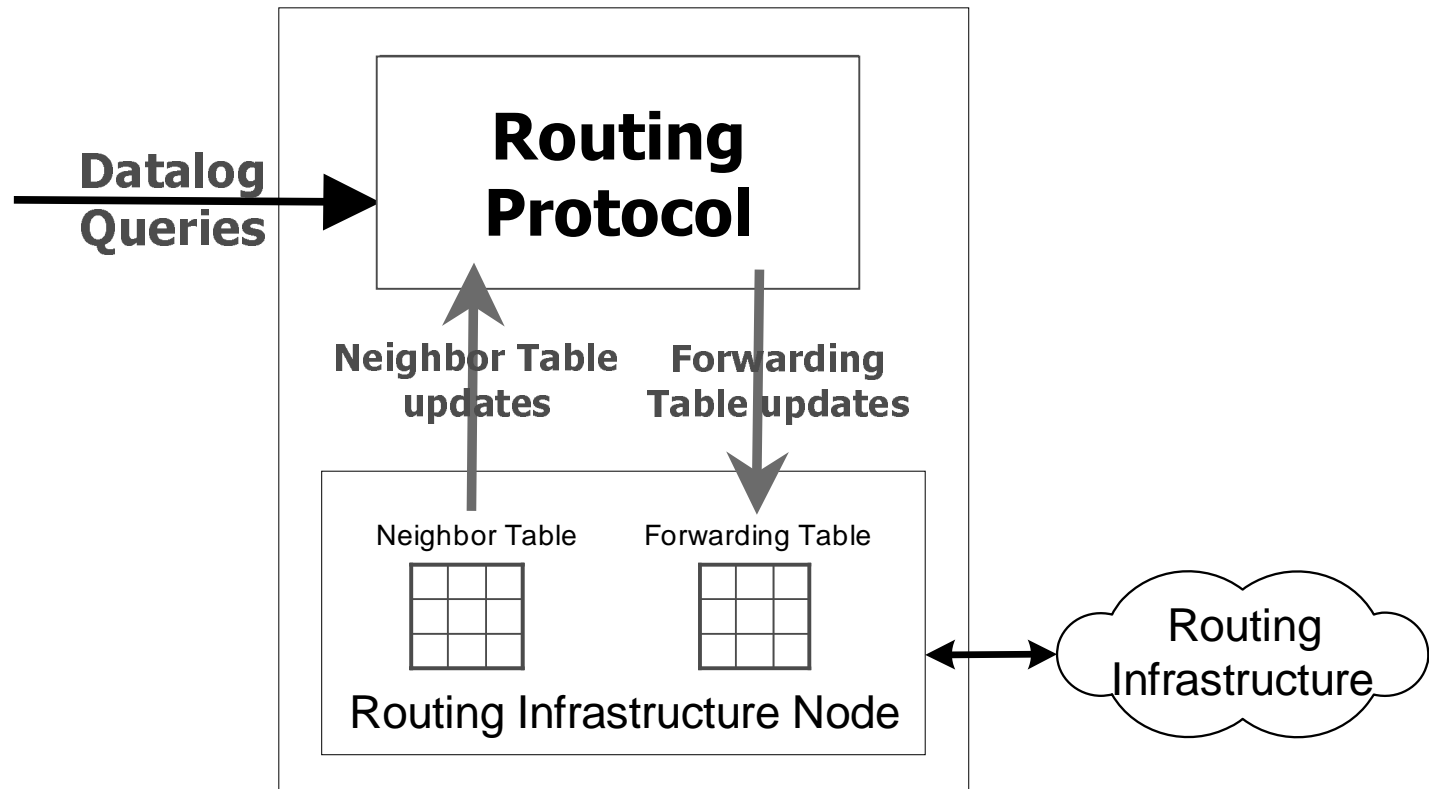
Roadmap

- ◆ Execution Model
- ◆ Introduction to Datalog
- ◆ Path-Vector Protocol Example
- ◆ Advantages:
 - Expressiveness
 - Safety
 - Efficiency
- ◆ Evaluation

Centralized Execution Model

- ◆ Store entire network state into a centralized database
- ◆ Issue Datalog queries on the centralized database for customized routes

Distributed Execution Model



Distributed Routing

Introduction to Datalog

Datalog rule syntax:

<head> ← <precondition1>, <precondition2>, ... , <preconditionN>.

All-Pairs Reachability

➔ R1: $\text{reachable}(S,D) \leftarrow \text{link}(S,D)$

R2: $\text{reachable}(S,D) \leftarrow \text{link}(S,Z), \text{reachable}(Z,D)$

$\forall \text{link}(a,b) \in E$, "there is a link from node a to node b "

$\text{reachable}(a,b)$ – "node a can reach node b "

"For all nodes S,D ,

If there is a link from S to D , then S can reach D ".

◆ Input: $\text{link}(\text{source}, \text{destination})$

◆ Output: $\text{reachable}(\text{source}, \text{destination})$

All-Pairs Reachability

R1: $\text{reachable}(S,D) \leftarrow \text{link}(S,D)$

➔ R2: $\text{reachable}(S,D) \leftarrow \text{link}(S,Z), \text{reachable}(Z,D)$

“For all S, D and Z ,
If $\text{link}(S,Z)$ exists AND $\text{reachable}(Z,D)$ exists, generate
 $\text{reachable}(S,D)$.”

“For all nodes S, D and Z ,
If there is a link from S to Z , AND Z can reach D , then S
can reach D ”.

◆ Input: $\text{link}(\text{source}, \text{destination})$

◆ Output: $\text{reachable}(\text{source}, \text{destination})$

All-Pairs Reachability

➔ R1: $\text{reachable}(S,D) \leftarrow \text{link}(S,D)$

R2: $\text{reachable}(S,D) \leftarrow \text{link}(S,Z), \text{reachable}(Z,D)$

Query: $\text{reachable}(M,N)$ ← **All-Pairs**

Input table:

link

S	D
a	b

link

S	D
b	c

link

S	D
c	d



Output table (Round 1):

reachable

S	D
a	b

reachable

S	D
b	c

reachable

S	D
c	d

E.g. R1: $\text{reachable}(b,c) \leftarrow \text{link}(b,c)$

All-Pairs Reachability

R1: $\text{reachable}(S,D) \leftarrow \text{link}(S,D)$

➔ R2: $\text{reachable}(S,D) \leftarrow \text{link}(S,Z), \text{reachable}(Z,D)$

Query: $\text{reachable}(M,N)$

Input table:

link

S	D
a	b

link

S	D
b	c

link

S	D
c	d



**Output table
(Round 2):**

reachable

S	D
a	b
a	c

reachable

S	D
b	c
b	d

reachable

S	D
c	d

R2: $\text{reachable}(b,d) \leftarrow \text{link}(b,c), \text{reachable}(c,d)$

All-Pairs Reachability

R1: $\text{reachable}(S,D) \leftarrow \text{link}(S,D)$

➔ R2: $\text{reachable}(S,D) \leftarrow \text{link}(S,Z), \text{reachable}(Z,D)$

Query: $\text{reachable}(M,N)$

Input table:

link	
S	D
a	b

link	
S	D
b	c

link	
S	D
c	d



reachable reachable reachable

Output table
(Round 3)

Recursive queries are natural for querying graph topologies

a	d
---	---

Roadmap

- ◆ Execution Model
- ◆ Introduction to Datalog
- ◆ Path-Vector Protocol Example
 - Distributed Datalog → Execution Plan → Protocol
- ◆ Advantages:
 - Expressiveness
 - Safety
 - Efficiency
- ◆ Evaluation

Distributed Datalog

R1: $\text{reachable}(\underline{S}, D) \leftarrow \text{link}(\underline{S}, D)$

R2: $\text{reachable}(\underline{S}, D) \leftarrow \text{link}(\underline{S}, D), \text{reachable}(\underline{Z}, D)$

Query: $\text{reachable}(\underline{M}, N)$

Input table:

link

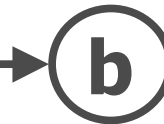
<u>S</u>	D
a	b

link

<u>S</u>	D
b	c

link

<u>S</u>	D
c	d



reachable

reachable

reachable

Output table:

<u>S</u>	D
a	b
a	c
a	d

<u>S</u>	D
b	c
b	d

<u>S</u>	D
c	d

Path Vector Protocol Example

R1: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, D), P=(S, D).$

R2: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{Z}, S), \text{path}(\underline{Z}, D, P_2), P=S+P_2.$

Query: $\text{path}(\underline{S}, D, P)$

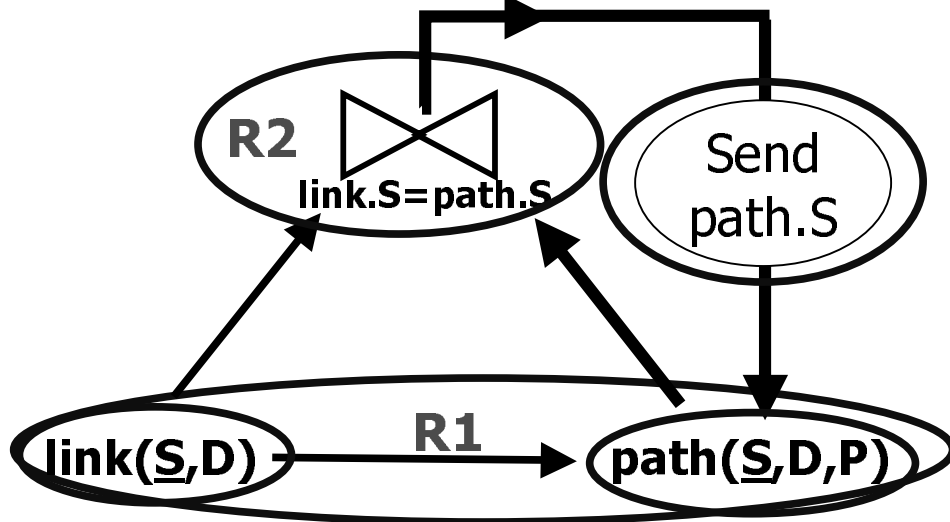
- ◆ Input: $\text{link}(\text{source}, \text{destination})$
- ◆ Query output: $\text{path}(\text{source}, \text{destination}, \text{pathVector})$

Datalog \rightarrow Execution Plan

R1: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, D), P=(S, D).$

R2: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{Z}, S), \text{path}(\underline{Z}, D, P_2), P=S+P_2.$

Matching variable Z = "Join" \bowtie
Recursion



Pseudocode at node Z:

```

while (receive < path(Z, D, P2) > ) {
  for each neighbor S {
    for each neighbor S {
      newpath = path(S, D, S + P2)
      newpath = path(S, D, S + P2)
    }
    send newpath to neighbor S
  }
}

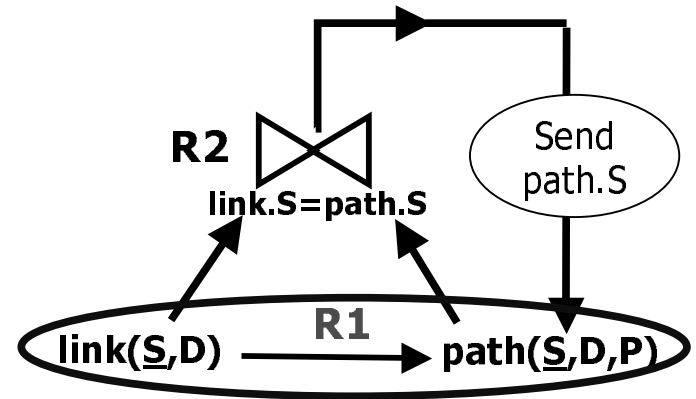
```

Query Execution

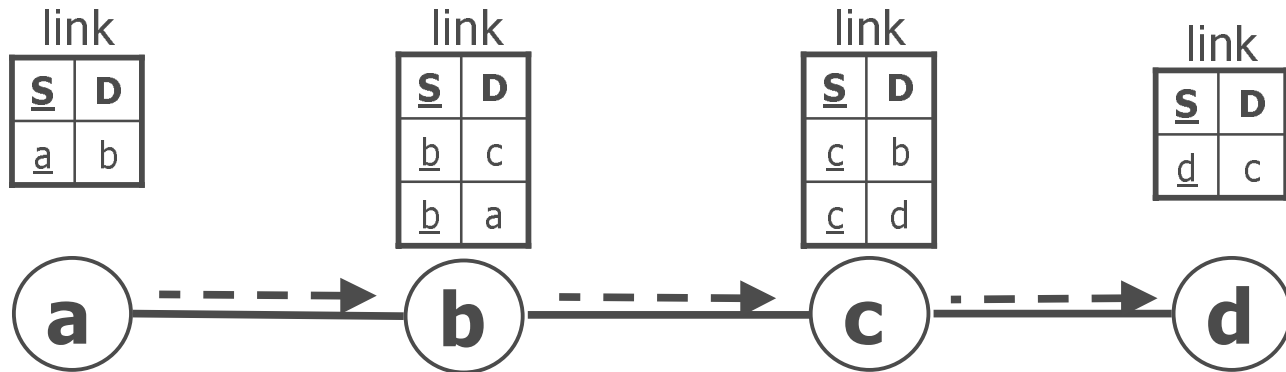
R1: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, D), P=(S, D).$

R2: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{Z}, S), \text{path}(\underline{Z}, D, P_2), P=S+P_2.$

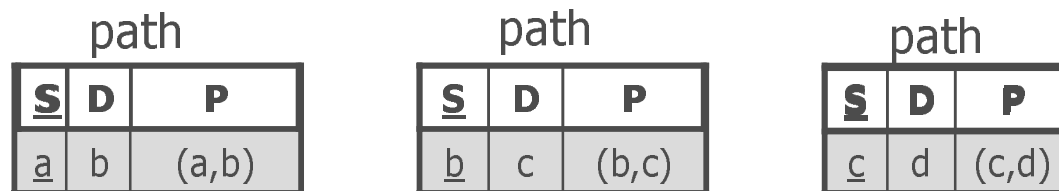
Query: $\text{path}(\underline{S}, D, P, C)$



Neighbor table:



Forwarding table:

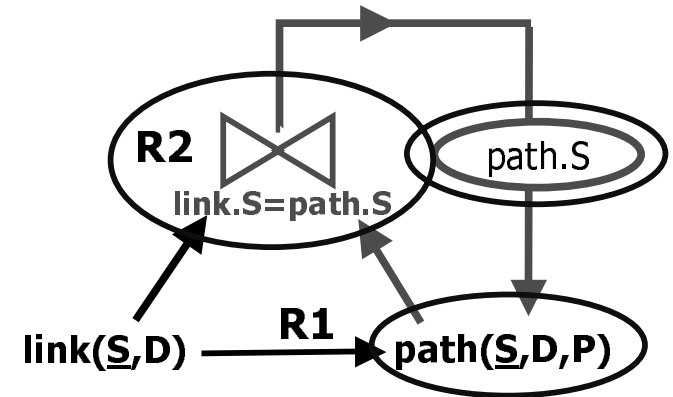


Query Execution

R1: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, D), P=(S, D).$

R2: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, Z), \text{path}(\underline{Z}, D, P_2), P=S+P_2.$

Query: $\text{path}(\underline{S}, D, P, C)$



Neighbor table:

link

<u>S</u>	D
a	b

link

<u>S</u>	D
b	c
b	a

link

<u>S</u>	D
c	b
c	d

link

<u>S</u>	D
d	c



$p(a, c, [a, b, c])$

$p(b, d, [b, c, d])$

Forwarding table:

path

<u>S</u>	D	P
a	b	(a, b)
a	c	(a, b, c)

path

<u>S</u>	D	P
b	c	(b, c)
b	d	(b, c, d)

path

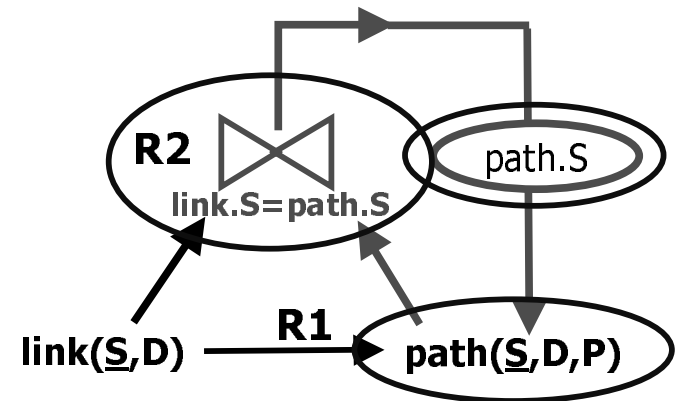
<u>S</u>	D	P
c	d	(c, d)

Query Execution

R1: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, D), P=(S, D).$

R2: $\text{path}(\underline{S}, D, P) \leftarrow \text{link}(\underline{S}, Z), \text{path}(\underline{Z}, D, P_2), P=S+P_2.$

Query: $\text{path}(\underline{S}, D, P, C)$



Neighbor table:

link	
<u>S</u>	D
a	b

link	
<u>S</u>	D
b	c
b	a

link	
<u>S</u>	D
c	b
c	d

link	
<u>S</u>	D
d	c



$p(a, d, [a, b, c, d])$

path

path

path

Forward table

Communication patterns are identical to those in the actual path vector protocol

<u>a</u>	d	(a, b, c, d)
----------	---	--------------

Roadmap

- ◆ Execution Model
- ◆ Introduction to Datalog
- ◆ Path-Vector Protocol Example
 - Distributed Datalog → Execution Plan → Protocol
- ◆ Advantages:
 - Expressiveness
 - Safety
 - Efficiency
- ◆ Evaluation

Expressiveness

- ◆ Best-Path Routing
 - ◆ Distance Vector
 - ◆ Dynamic Source Routing
 - ◆ Policy Decisions
 - ◆ QoS-based Routing
 - ◆ Link-state
 - ◆ Multicast Overlays (Single-Source & CBT)
- } Minor variants give many options!

Expressiveness

◆ All-pairs all-paths:

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C), P = (\underline{S}, D).$

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, P_2, C_2), C = C_1 + C_2,$
 $P = S + P_2.$

Query: $\text{path}(\underline{S}, D, P, C)$

Expressiveness

◆ Best-Path Routing:

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C), P = (\underline{S}, D)$.

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, P_2, C_2), C = C_1 + C_2, P = \underline{S} + P_2$.

R3: $\text{bestPathCost}(\underline{S}, D, \min\langle C \rangle) \leftarrow \text{path}(\underline{S}, D, Z, C)$

R4: $\text{bestPath}(\underline{S}, D, Z, C) \leftarrow \text{bestPathCost}(\underline{S}, D, C), \text{path}(\underline{S}, D, P, C)$

Query: $\text{bestPath}(\underline{S}, D, P, C)$

Expressiveness

◆ Best-Path Routing:

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C)$, $P = (\underline{S}, D)$.

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, Z, C_1)$, $\text{path}(\underline{Z}, D, P_2, C_2)$, $C = \text{FN}(C_1, C_2)$, $P = S + P_2$.

R3: $\text{bestPathCost}(\underline{S}, D, \text{AGG}\langle C \rangle) \leftarrow \text{path}(\underline{S}, D, Z, C)$

R4: $\text{bestPath}(\underline{S}, D, Z, C) \leftarrow \text{bestPathCost}(\underline{S}, D, C)$, $\text{path}(\underline{S}, D, P, C)$

Query: $\text{bestPath}(\underline{S}, D, P, C)$

Customizing C, AGG and FN: lowest RTT, lowest loss rate, highest available bandwidth, *best-k*

Expressiveness

◆ All-pairs all-paths:

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C), P = (\underline{S}, D).$

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, P_2, C_2), C = C_1 + C_2,$
 $P = S + P_2.$

Query: $\text{path}(\underline{S}, D, P, C)$

Expressiveness

◆ Distance Vector:

R1: $\text{path}(\underline{S}, D, D, C) \leftarrow \text{link}(\underline{S}, D, C)$

R2: $\text{path}(\underline{S}, D, Z, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, W, C_2), C = C_1 + C_2$

R3: $\text{shortestLength}(\underline{S}, D, \min\langle C \rangle) \leftarrow \text{path}(\underline{S}, D, Z, C)$

R4: $\text{nextHop}(\underline{S}, D, Z, C) \leftarrow \text{nextHop}(\underline{S}, D, Z, C), \text{shortestLength}(\underline{S}, D, C)$

Query: $\text{nextHop}(\underline{S}, D, Z, C)$

Count to Infinity problem?

Expressiveness

◆ Distance Vector with Split Horizon:

R1: $\text{path}(\underline{S}, D, D, C) \leftarrow \text{link}(\underline{S}, D, C)$

R2: $\text{path}(\underline{S}, D, Z, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, W, C_2), C = C_1 + C_2, W \neq S$

R3: $\text{shortestLength}(\underline{S}, D, \min\langle C \rangle) \leftarrow \text{path}(\underline{S}, D, Z, C)$

R4: $\text{nextHop}(\underline{S}, D, Z, C) \leftarrow \text{nextHop}(\underline{S}, D, Z, C), \text{shortestLength}(\underline{S}, D, C)$

Query: $\text{nextHop}(\underline{S}, D, Z, C)$

Expressiveness

◆ Distance Vector with Poisoned Reverse:

R1: $\text{path}(\underline{S}, D, D, C) \leftarrow \text{link}(\underline{S}, D, C)$

R2: $\text{path}(\underline{S}, D, Z, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, W, C_2), C = C_1 + C_2, W \neq S$

R3: $\text{path}(\underline{S}, D, Z, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, W, C_2), C = \infty, W = S$

R4: $\text{shortestLength}(\underline{S}, D, \min\langle C \rangle) \leftarrow \text{path}(\underline{S}, D, Z, C)$

R5: $\text{nextHop}(\underline{S}, D, Z, C) \leftarrow \text{nextHop}(\underline{S}, D, Z, C), \text{shortestLength}(\underline{S}, D, C)$

Query: $\text{nextHop}(\underline{S}, D, Z, C)$

Expressiveness

◆ All-pairs all-paths:

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C), P = (S, D).$

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, Z, C_1), \text{path}(\underline{Z}, D, P_2, C_2), C = C_1 + C_2, P = S + P_2.$

Query: $\text{path}(\underline{S}, D, P, C)$

Expressiveness

◆ Dynamic Source Routing (DSR):

R1: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{S}, D, C), P = (S, D).$

R2: $\text{path}(\underline{S}, D, P, C) \leftarrow \text{link}(\underline{Z}, D, C_2), \text{path}(\underline{S}, Z, P_1, C_1), C = C_1 + C_2,$
 $P = P_1 + D.$

Query: $\text{path}(\underline{S}, D, P, C)$

Switching *Right-recursion* to *Left-recursion* execution
=> Path vector protocol to DSR.

Expressiveness

- ◆ Best-Path routing
- ◆ Distance Vector
- ◆ Dynamic Source Routing
- ◆ Policy-based routing
- ◆ QoS-based routing
- ◆ Link-state
- ◆ Multicast Overlays (Single-Source & CBT)

Safety

- ◆ Queries are sand-boxed within query engine
 - Queries use input tables to produce output tables
 - No side-effects on existing input tables
- ◆ Pure Datalog guarantees termination:
 - Natural bound on resource consumption of queries
- ◆ Static termination checks for our extended Datalog:
 - Identify recursive definitions and check for termination
 - E.g., monotonically increasing/decreasing cost fields whose values are upper/lower bounded
- ◆ Orthogonal security issues:
 - Denial-of-service attacks, compromised routers

Efficiency

- ◆ Explore well-known database techniques
 - Aggregate selections: avoid sending unnecessary paths to neighbors
 - Limit computation to portion of network
 - ◆ Few sources and destinations
 - ◆ Magic sets + left-right recursion rewrite
 - Multi-query sharing:
 - ◆ Identify “similar” queries, share their computations
 - ◆ Reuse previously computed paths

Queries under Churn

- ◆ Long-running continuous queries
- ◆ Maintain all intermediate derived tuples for query duration
- ◆ Incremental updates:
 - Link failures are treated as link updates with $\text{cost}=\text{infinity}$.
 - Paths invalidated ($\text{cost}=\text{infinity}$), and new paths are incrementally recomputed.

Evaluation Setup

- ◆ PIER: Distributed relational query processor
 - Each node runs the query engine of PIER
 - Initialized neighbor table directly accessible by PIER.
- ◆ Simulation:
 - Bandwidth and latency bottlenecks
 - Transit-stub topologies
- ◆ PlanetLab
 - 72 PIER nodes
 - Random, location-aware topologies
 - Long-running queries

Summary of Results

◆ Simulations:

- When all nodes issue the same query,
 - ◆ Scalability properties show similar trends as traditional DV/PV protocols
- When few nodes issue the same query,
 - ◆ Overhead is reduced using standard query optimizations

◆ PlanetLab experiments:

- Long-running all-pairs shortest RTT paths query
- Ability to handle link RTT changes
- Induced failures (up to 20% of nodes)
 - ◆ Recovery time: <1s (median), <3.6s (average)

Conclusion

◆ Declarative routing:

- Express routing protocols using a recursive query language
- Better balance between routing extensibility and safety

◆ Future work:

- Expressing policies as declarative rules
- Run-time query optimizations:
 - ◆ Cost-based decisions on query rewrites
 - ◆ Multi-query sharing
- Static checker for routing protocols
- Run-time monitoring of routing protocols

◆ Declarative Networks

- Research agenda: Specify and construct networks declaratively
- P2: "Implementing Declarative Overlays" (SOSP 2005)

Thank You