**Propane**

# Programming Distributed Control Planes

Ryan Beckett (Princeton & MSR)

Ratul Mahajan (MSR)

Todd Millstein (UCLA)

Jitu Padhye (MSR)

David Walker (Princeton)

# Configuring Networks is Error-Prone

# Objectives ⟷ Mechanisms

GAP

## Objectives: Network-wide

- Prefer traffic to send traffic through customers over providers
- Don't use our network as transit between A and B
- Traffic must stay within national boundaries

## Mechanisms: Device-by-Device



aggregate
12.34.0.0/16

add community
8075:400

failure

filter
path(_100_)

# Propane:
# Programming a Distributed Control Plane

**1) Language for expressing high-level objectives with:**

- Path constraints and relative preferences with fall-backs in case of failures

- Uniform abstractions for intra- and inter-domain routing



**2) Compiler that generates low-level, distributed configs:**

- Efficient algorithms to synthesize a set of policy-compliant BGP configs

- Failure analysis guarantees *policy compliance* under all failures

# Example: A Data Center Network

# A Data Center Network

**Goals:**

- P1: Announce global services externally as the aggregate PG

- P2: Do not announce local services externally

- P3: Prefer Backbone1 to Backbone2

# A Data Center Network

**Goals:**

- P1: Announce global services externally as the aggregate PG

- P2: Do not announce local services externally

- P3: Prefer Backbone1 to Backbone2

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG



Global Services    Local Services

# A Data Center Network

**Goals:**

- P1: Announce global services externally as the aggregate PG
- P2: Do not announce local services externally
- P3: Prefer Backbone1 to Backbone2

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside
- do *not* export announce's from G,H outside
  - appeal: X,Y do not need to know which prefixes are local vs global
- aggregate to PG if announce is subset of PG

**Consider X-G, X-H Failure:**



8

# A Data Center Network

**Goals:**

- P1: Announce global services externally as the aggregate PG

- P2: Do not announce local services externally

- P3: Prefer Backbone1 to Backbone2

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG

**Consider X-G, X-H Failure:**

- PL* announcements travel H-Y-D-X

- PL* announcements are then leaked



Global Services          Local Services

# A Data Center Network

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG

- **disallow "valley" paths**

# A Data Center Network

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X,Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG

- <span style="color:darkred">**disallow "valley" paths**</span>

**Consider D-A, X-C Failure:**

# A Data Center Network

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG

- **disallow "valley" paths**
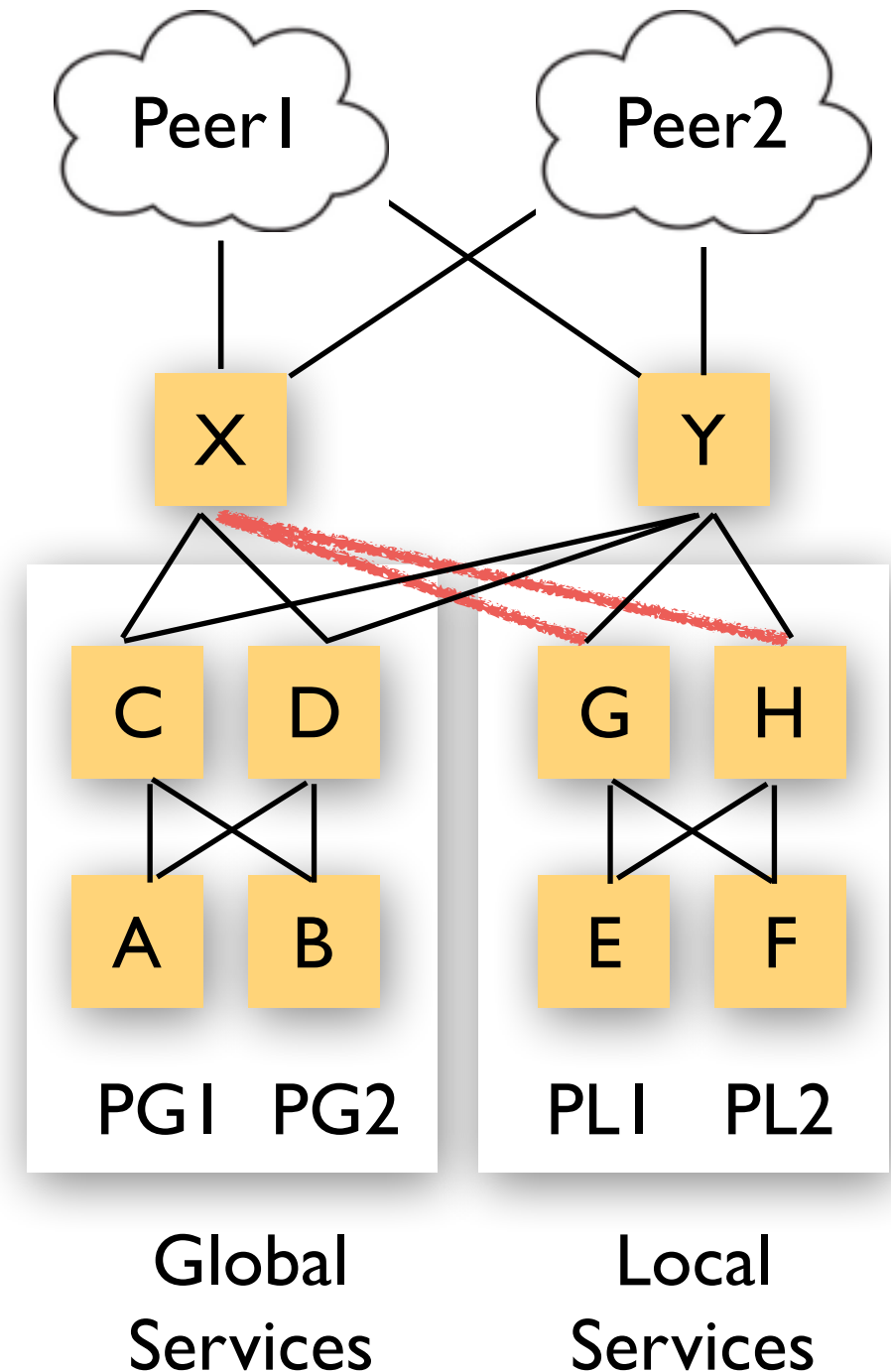
**Consider D-A, X-C Failure:**

- X and Y will hear PG2

# A Data Center Network

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global

- aggregate to PG if announce is subset of PG

- **disallow "valley" paths**

**Consider D-A, X-C Failure:**

- X and Y will hear PG2

- X and Y will announce aggregate PG



Peer1    Peer2

PG    PG

X    Y

C  D    G  H

A  B    E  F

PG1  PG2    PL1  PL2

Global Services    Local Services

# A Data Center Network

**Implementation Techniques for X, Y:**

- do export announce's from C, D outside

- do *not* export announce's from G,H outside

  - appeal: X, Y do not need to know which prefixes are local vs global
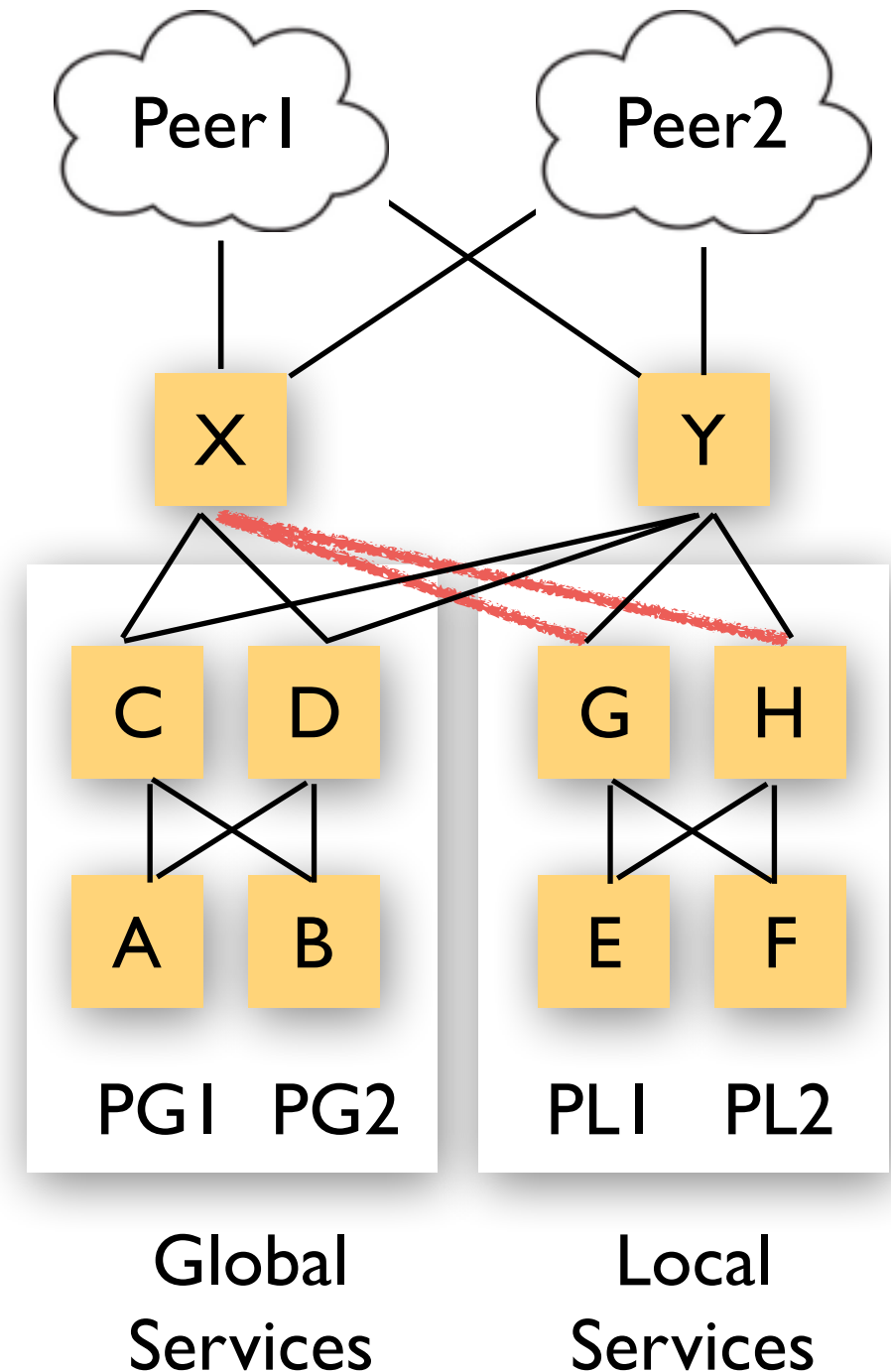
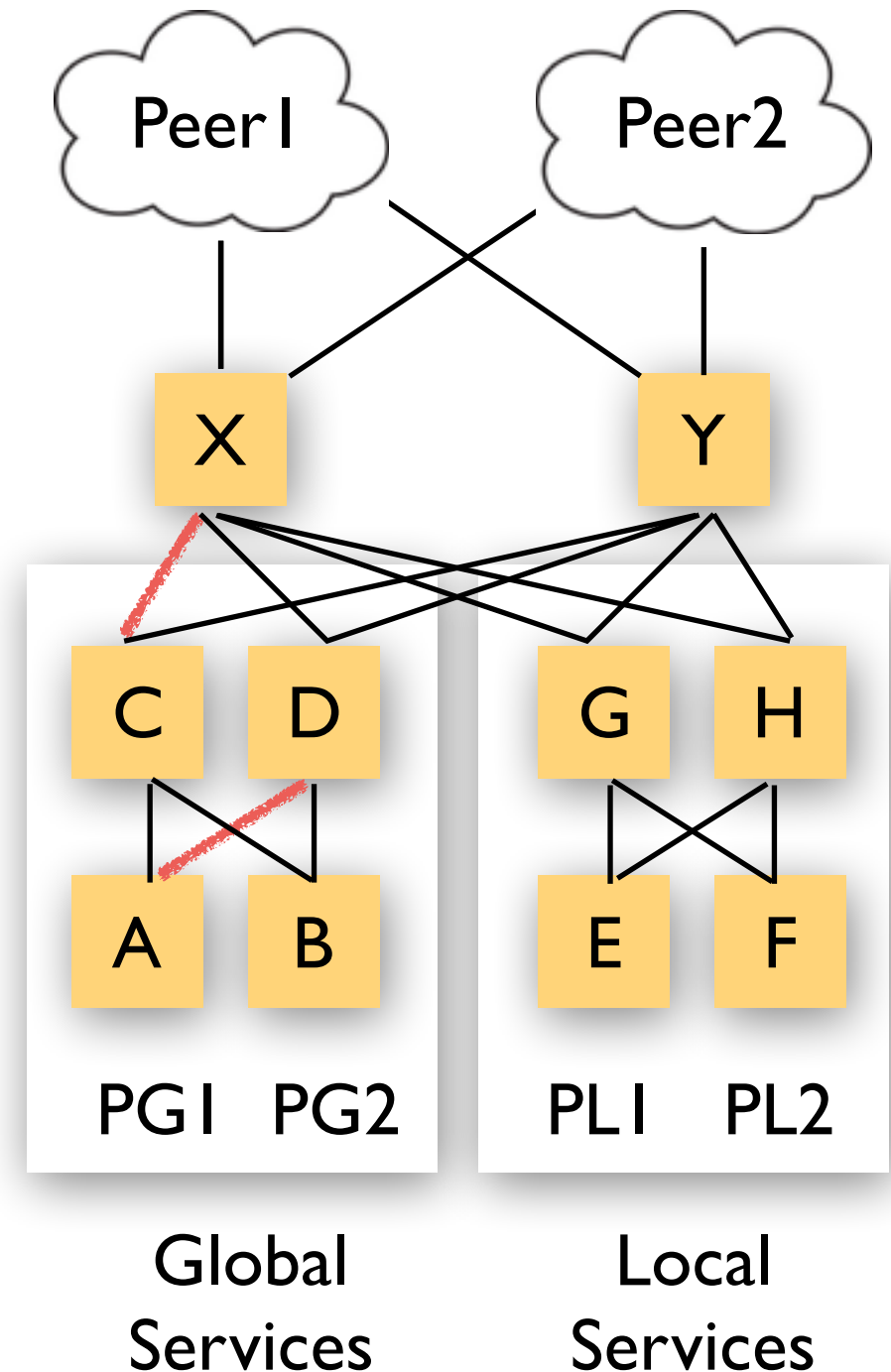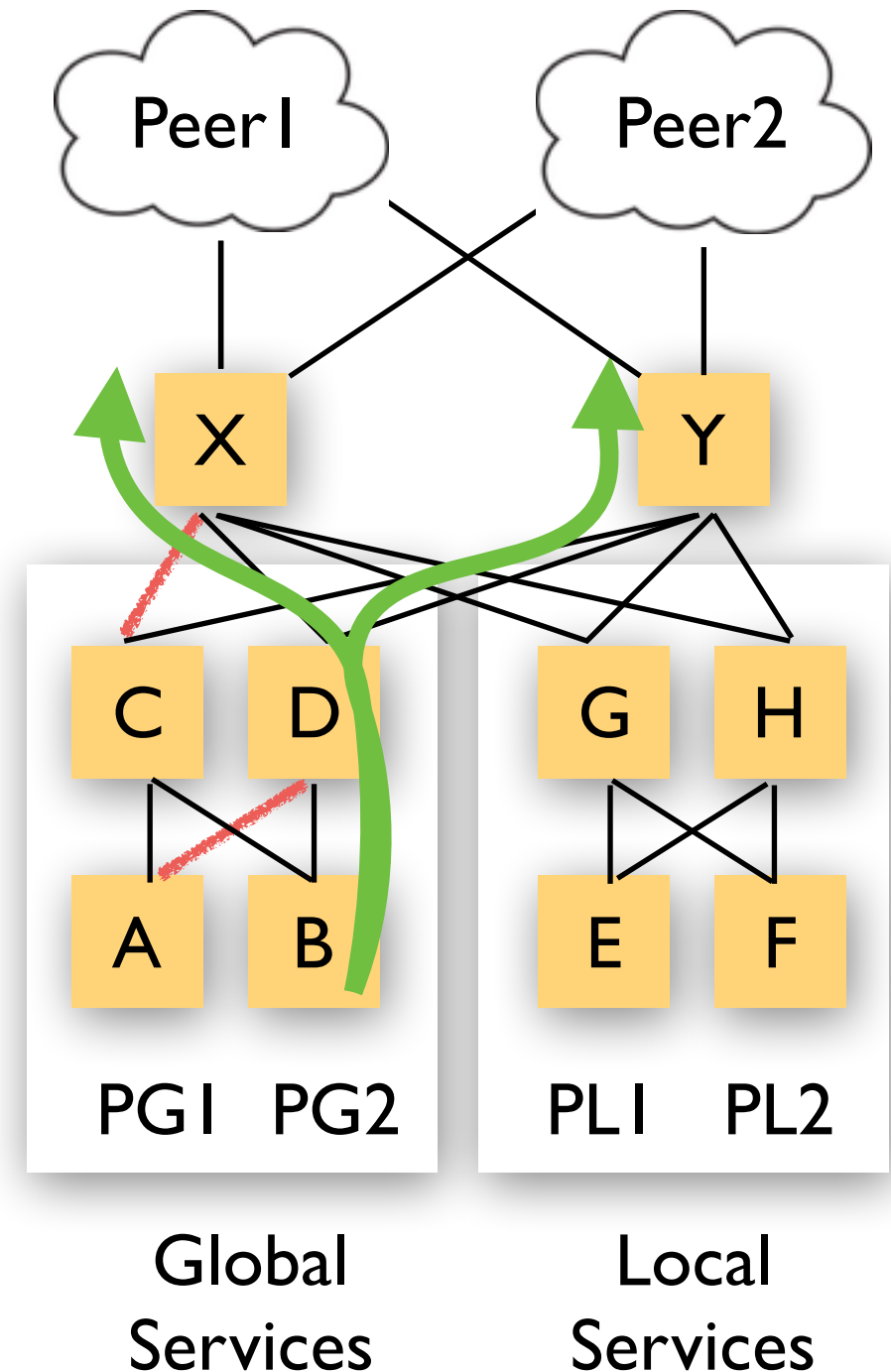- aggregate to PG if announce is subset of PG

- disallow "valley" paths
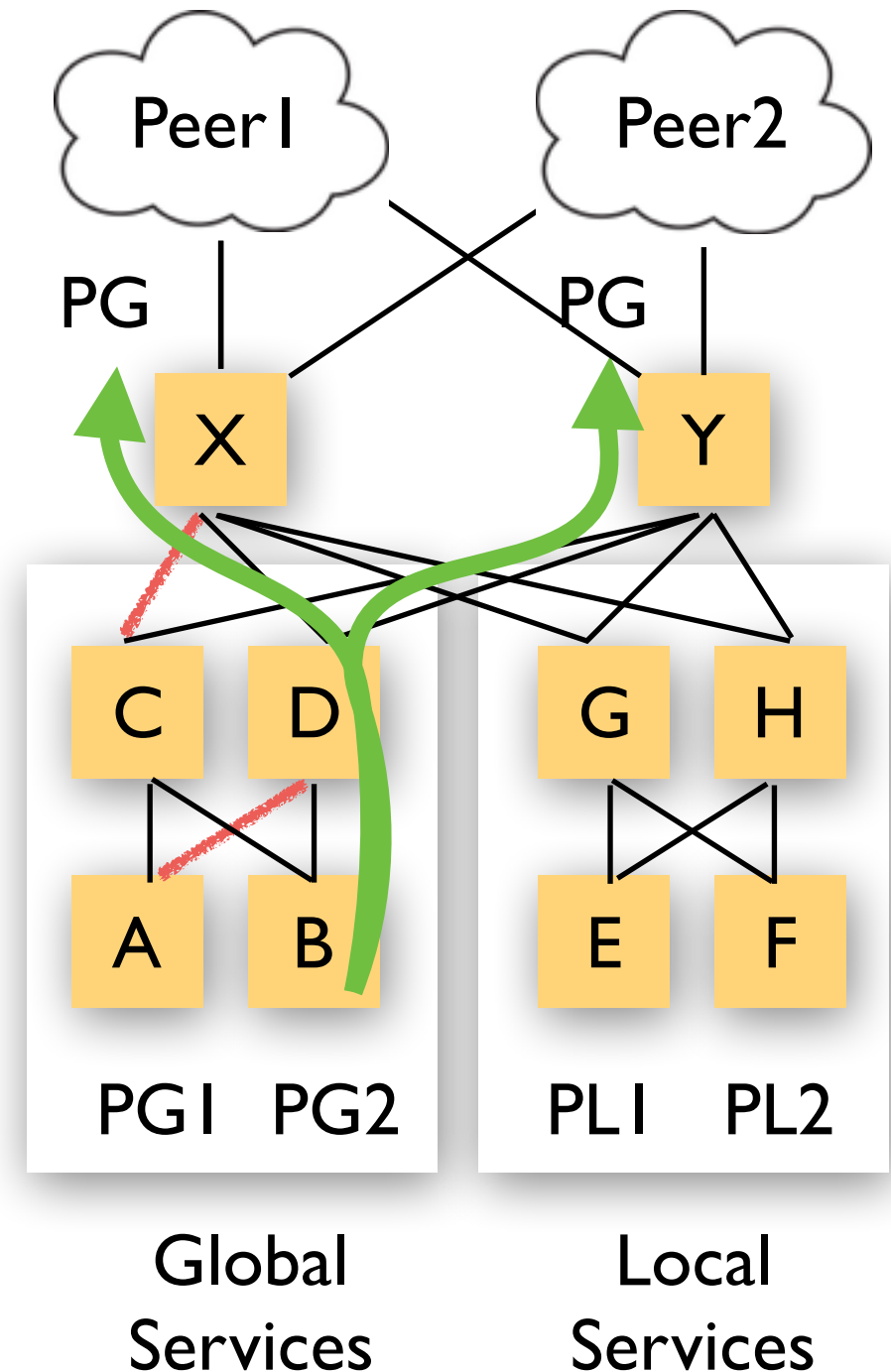
**Consider D-A, X-C Failure:**

- X and Y will hear PG2

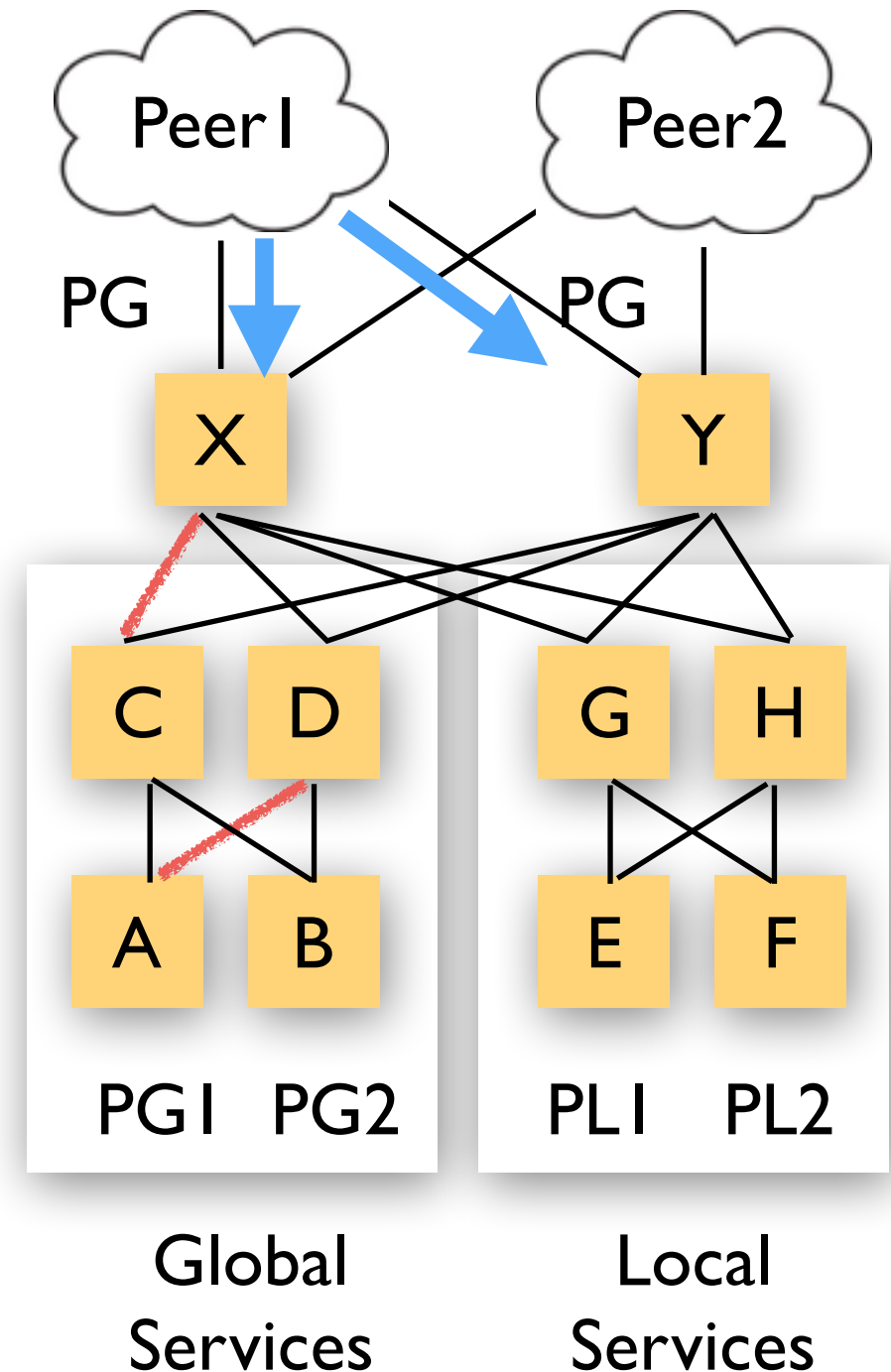- X and Y will announce aggregate PG

- But PG1 is inaccessible through X because there is no valley routing

- An aggregation-induced black hole is created [See Le et al, CoNext '11]

14

# Quick Demo

- Originate prefixes for each TOR router
- Do not announce local services externally
- Aggregation on global prefixes
- Prefer Peer1 over Peer2
- Prevent transit between Back1 and Back2



Peer1  Peer2

X  Y

C  D    G  H

A  B    E  F

PG1  PG2    PL1  PL2

Global Services    Local Services

# Compiling Propane

**Propane Compiler**

| | |
|---|---|
| Propane | Front End Constraint Language |
| Regular IR | Regular Expression-based IR |
| Topology → Product Graph IR ← Failure Analyses | |
| Abstract BGP | Vendor-independent BGP |
| CISCO    Juniper | Vendor-specific configurations |

17

# Propane Regular IR


Propane
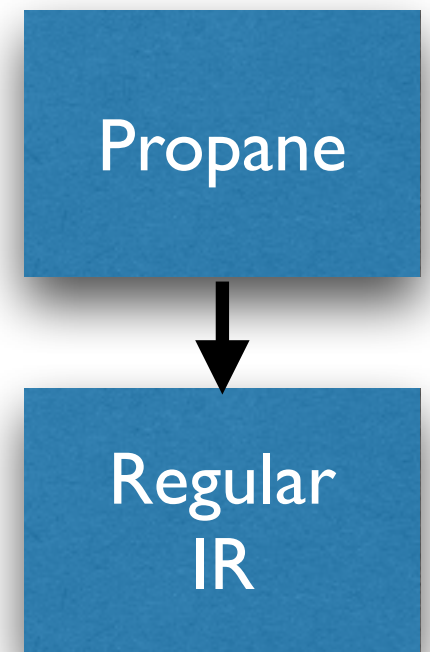
Regular IR

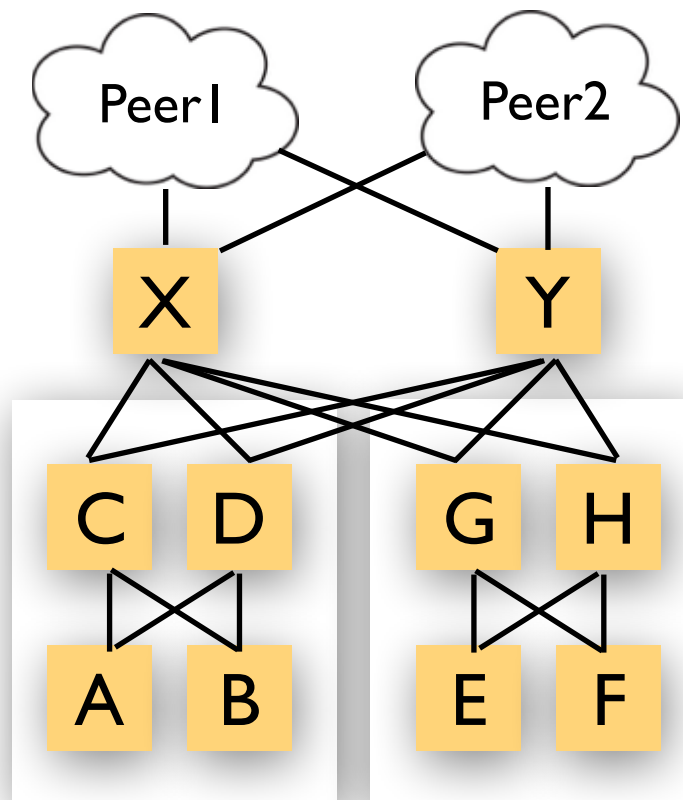**Expand constraints in to regular expressions.  EG:**

$$\mathbf{end}(X) = (\Sigma*.X)$$

$$\mathbf{exit}(X) = (\mathbf{out}*.\mathbf{in}*.(X \cap \mathbf{in}).\mathbf{out}+)\cup$$
$$(\mathbf{out}*.\mathbf{in}+.(X \cap \mathbf{out}).\mathbf{out*})$$
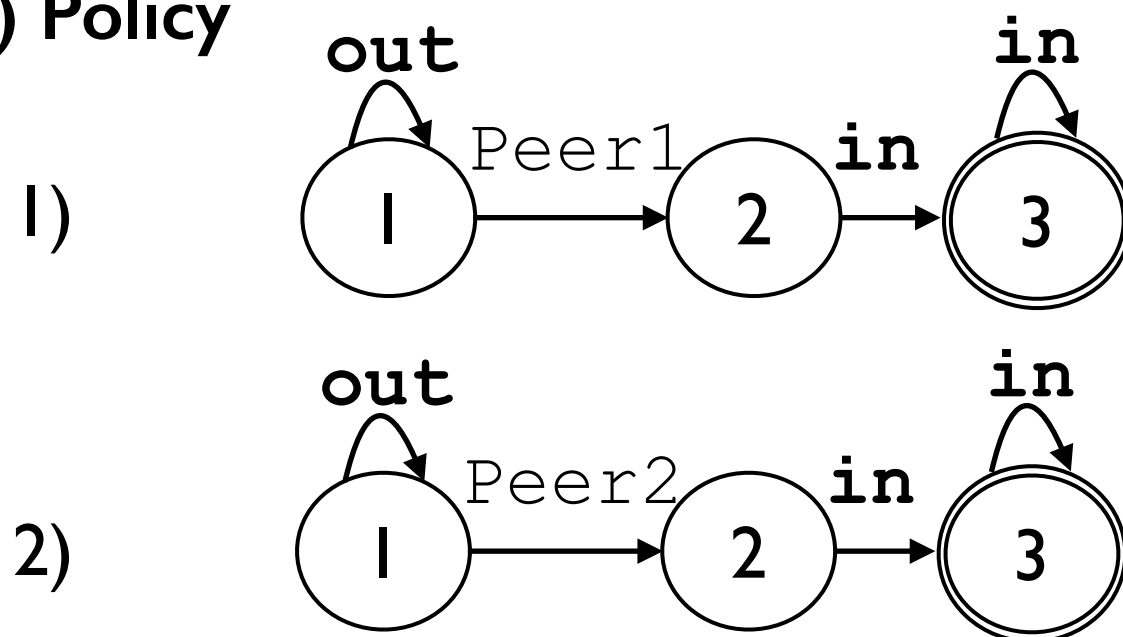
**A few other simple transformations:**

- conjunction of constraints ==> intersection of regular expressions

- conjunction of policies ==> prefix-by-prefix intersection

- nested preferences lifted: (x >> y) . z ==> (x.z) >> (y.z)

# Constructing the Product Graph (PG)

Regular IR

Product Graph IR

## (a) Topology



## (b) Policy
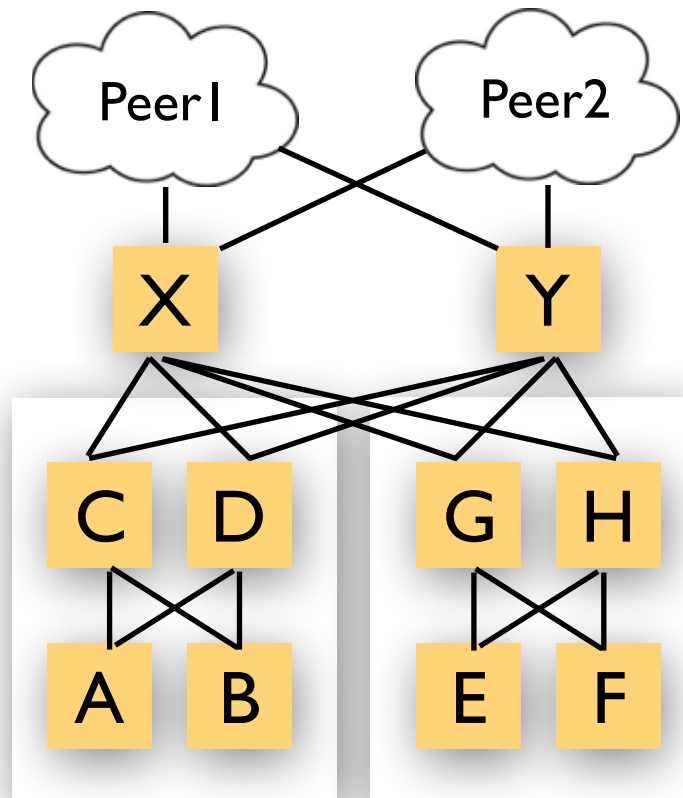
1)



2)

# Constructing the Product Graph (PG)
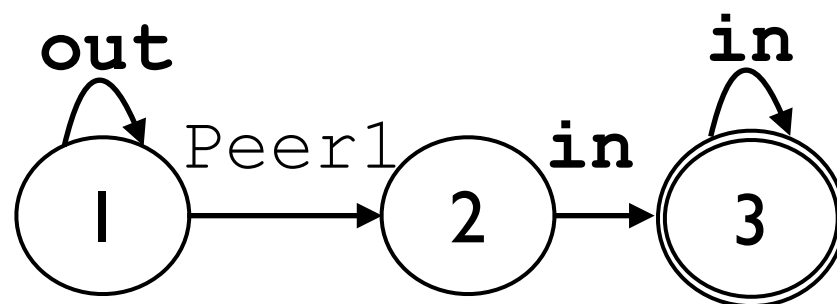
## (a) Topology



General Idea:
PG represents locations reachable in the topology while following the policy
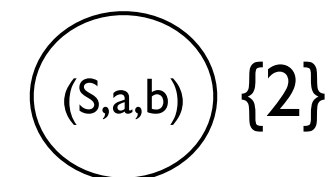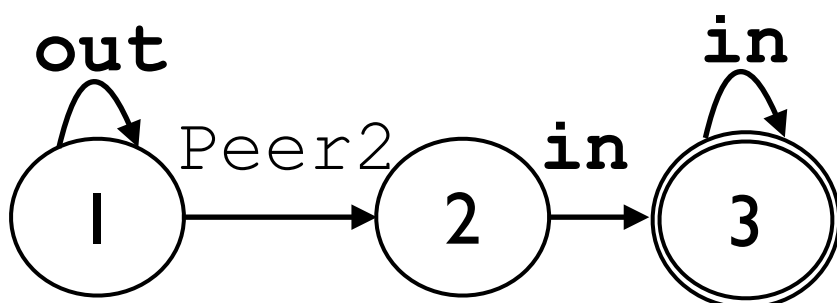
Each PG node contains:
- topology node (S)
- state of automaton 1 (a)
- state of automaton 2 (b)
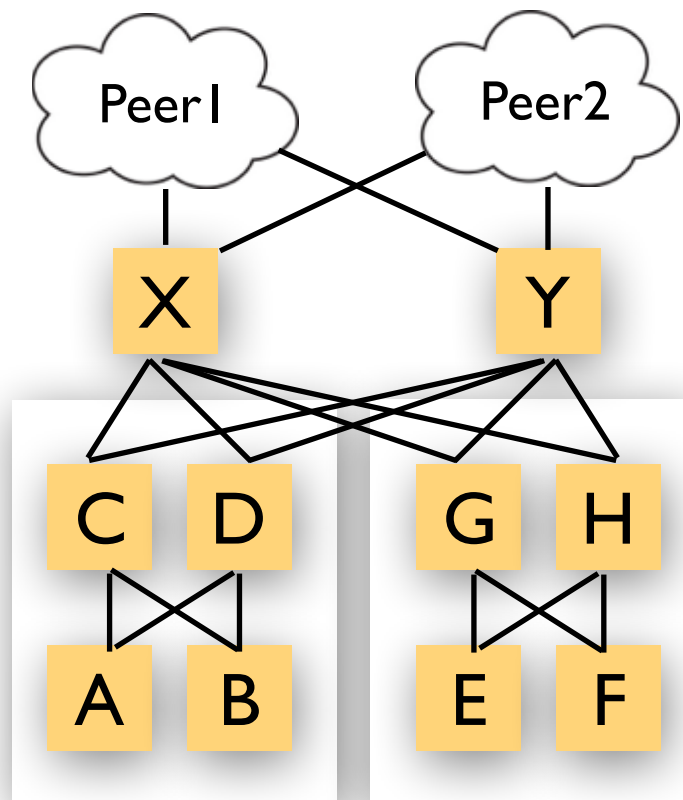- set of preferences achieved
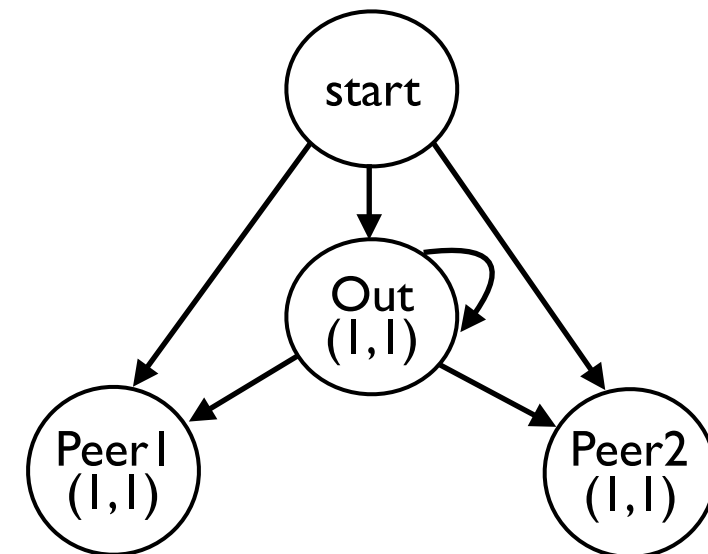
(S,a,b)  {2}

## (b) Policy



Two PG nodes are connected if topology nodes are connected and the automata make the specified transition
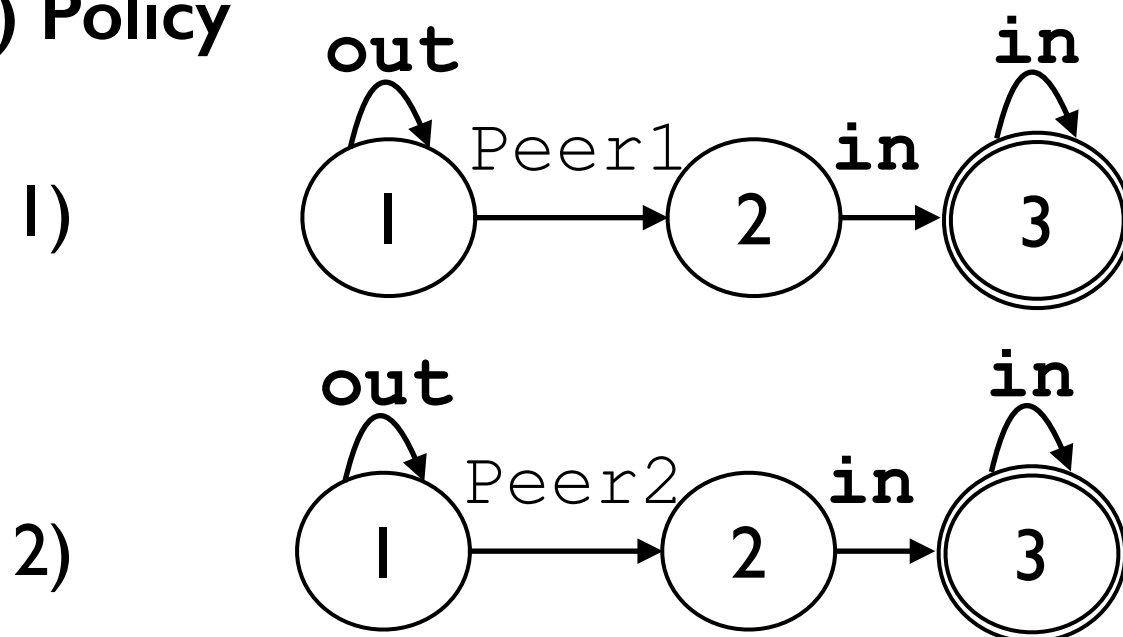
# Constructing the Product Graph (PG)

**(a) Topology**
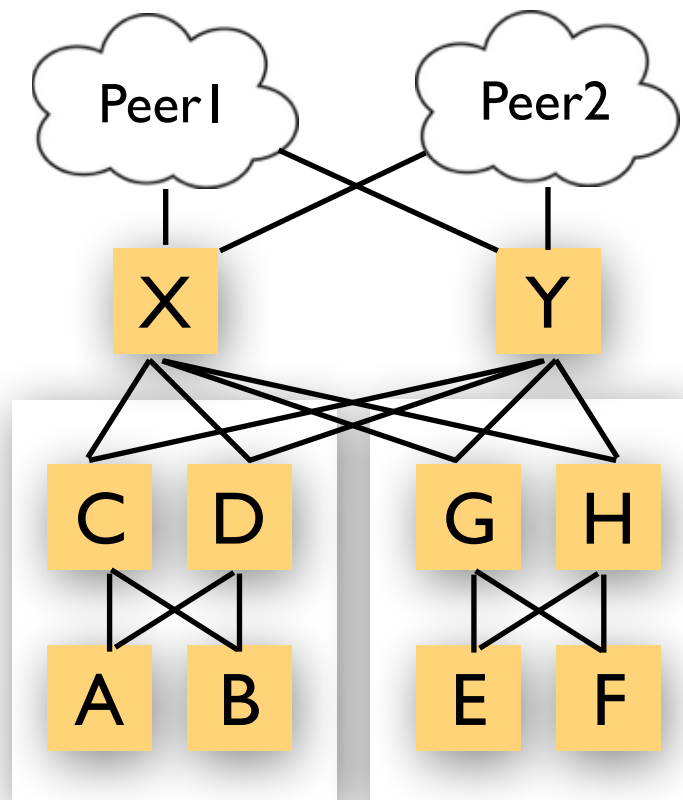


**(c) Product Graph**



**(b) Policy**

# Constructing the Product Graph (PG)

**(a) Topology**



**(c) Product Graph**



**(b) Policy**

# Constructing the Product Graph (PG)

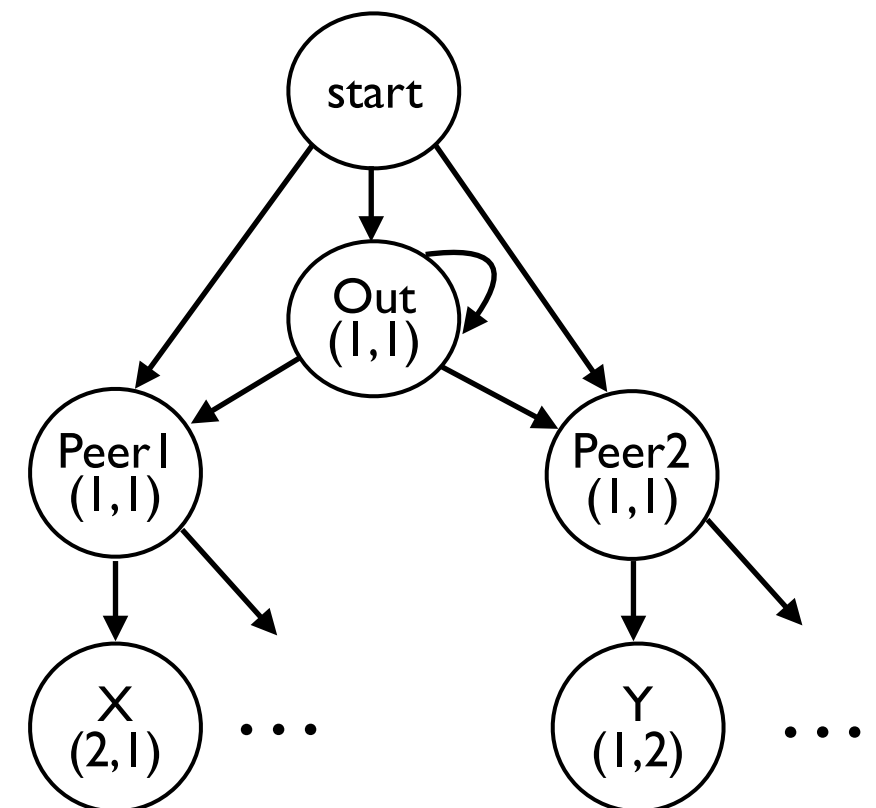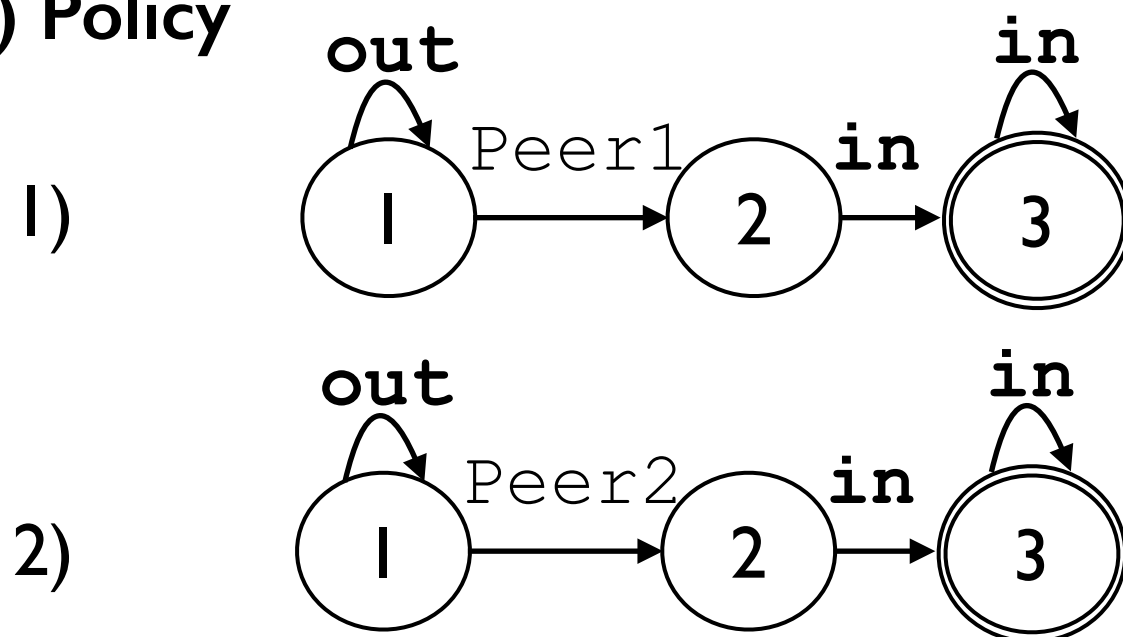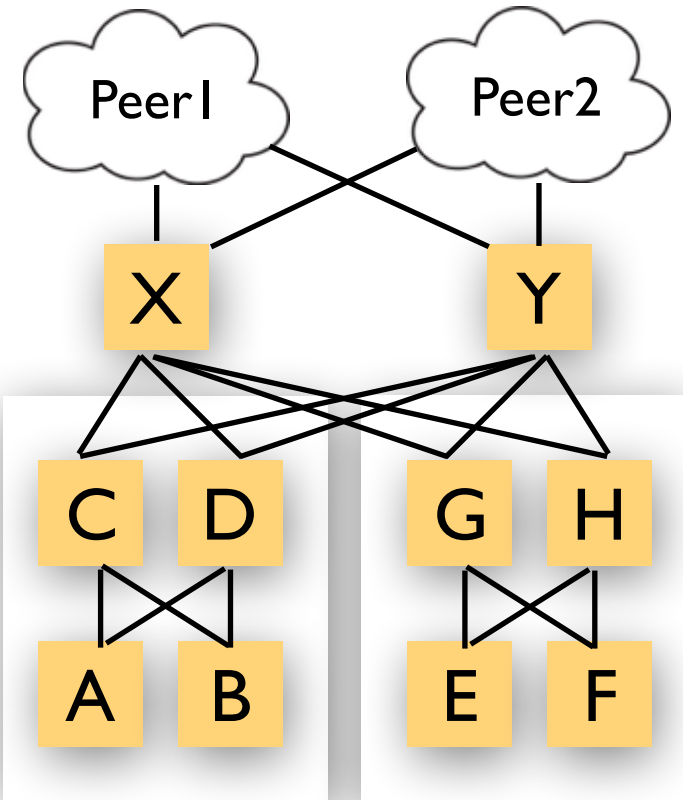**(a) Topology**



**(b) Policy**



**(c) Product Graph**

# Constructing the Product Graph (PG)

**(a) Topology**



**(c) Product Graph**



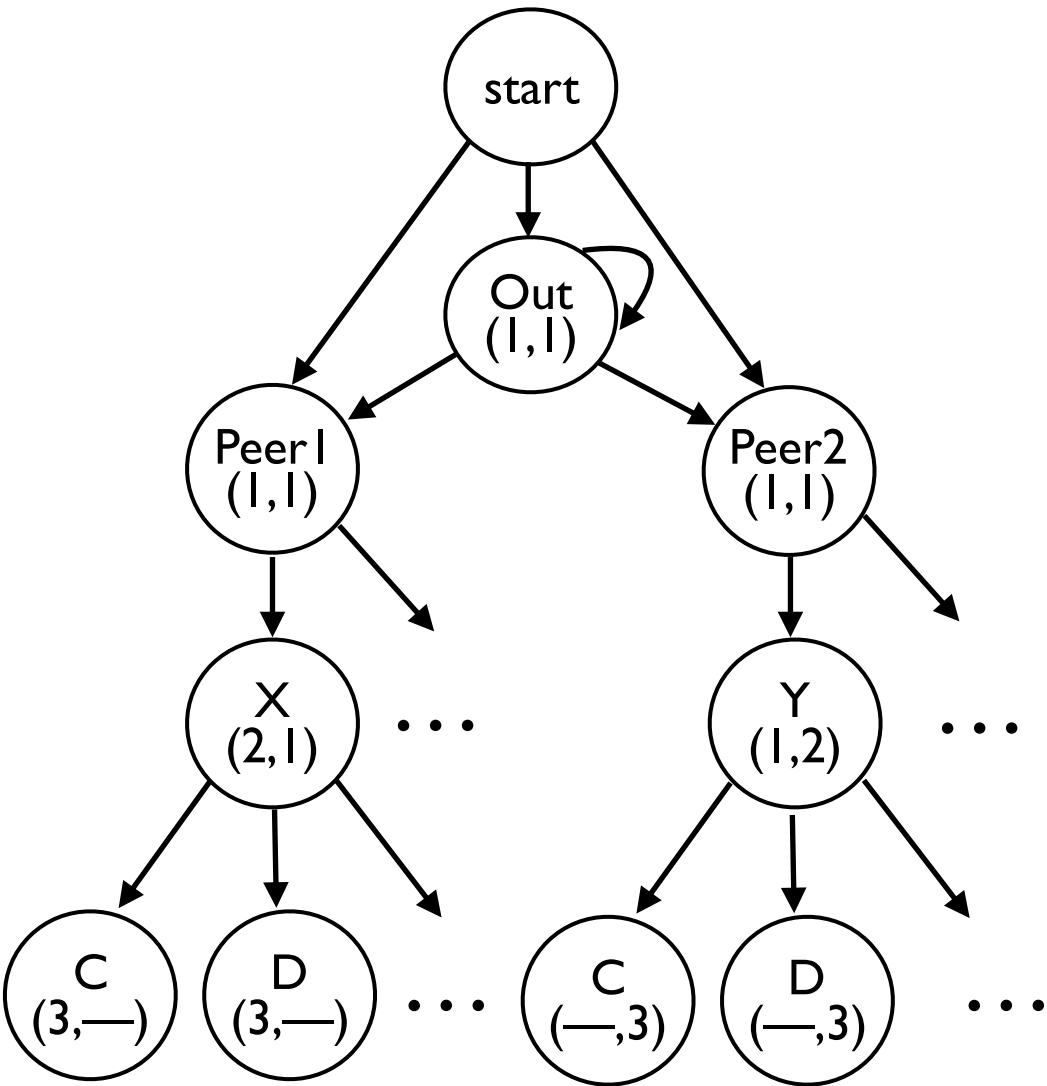**(b) Policy**

# Compilation to BGP:



**Idea:**

- Filter import messages according to incoming PG edges.

- For each internal location, decide which announcements to prefer, forward messages along PG edges

- Use a community value to tag the state of the automata

- For each external location, do nothing

# Compilation to BGP:



Router C
    **allow** peer=X comm=(2,1)
        peer ←{A,B} comm ← (3, —)
    **allow** peer=Y comm=(1,2)
        peer ←{A,B} comm ← (—, 3)

# Compilation to BGP:



Router C
    **allow** peer=X comm=(2,1)
        peer ←{A,B} comm ← (3,—)
    **allow** peer=Y comm=(1,2)
        peer ←{A,B} comm ← (—, 3)

Router X
    **allow** regex(Peer1 . out*)
        peer ←{C,D,G,H} comm ← (2, 1)
    **allow** regex(Peer2 . out*)
        peer ←{C,D,G,H} comm ← (1, 2)

# Compilation to BGP:



Router C
    **allow** peer=X comm=(2,1)
        peer ←{A,B} comm ← (3,—)
    **allow** peer=Y comm=(1,2)
        peer ←{A,B} comm ← (—, 3)

**Graph Analysis**

Router X
    **allow** regex(Peer1 . out*) **with** lp = 100
        peer ←{C,D,G,H} comm ← (2, 1)
    **allow** regex(Peer2 . out*) **with** lp = 99
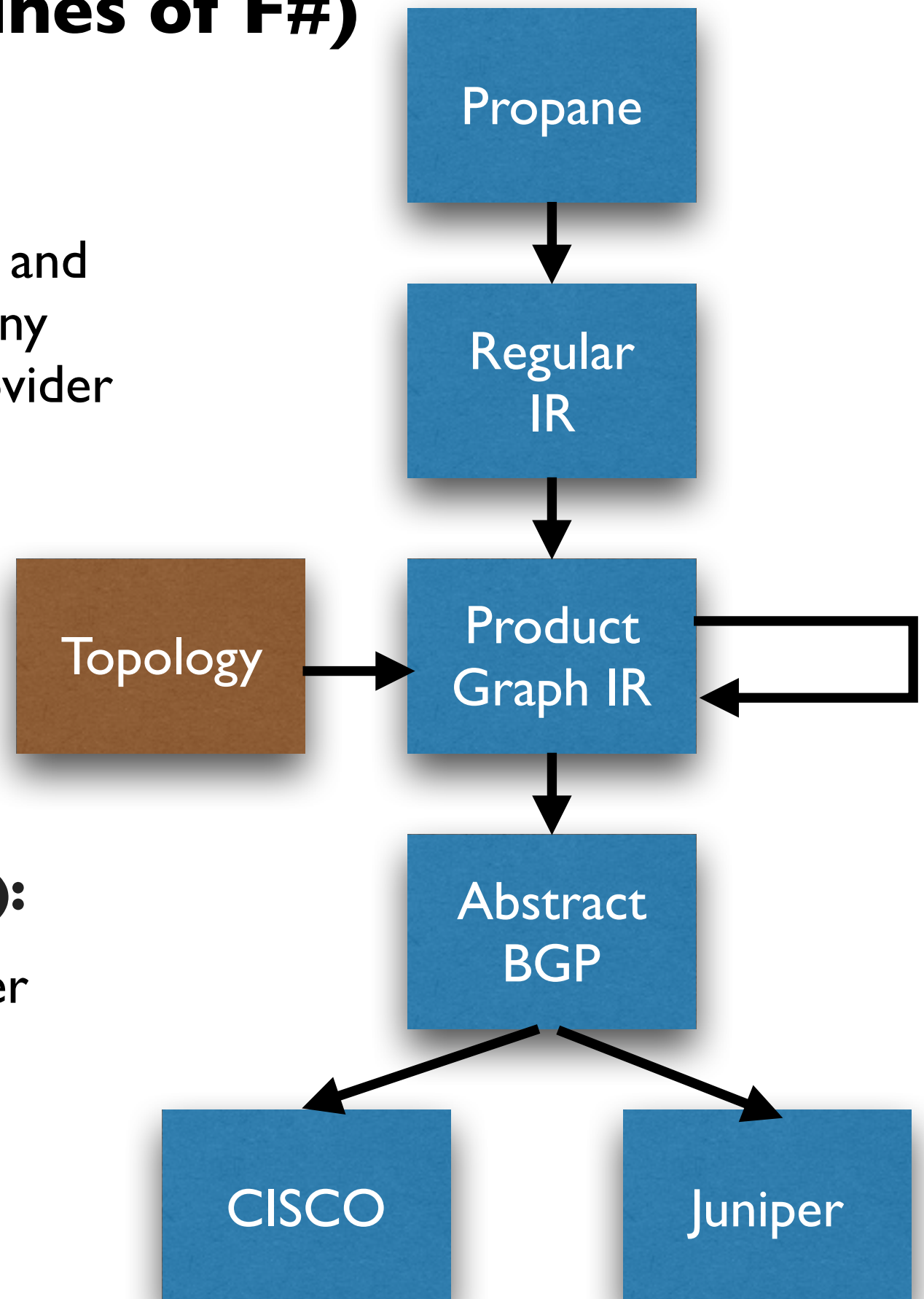        peer ←{C,D,G,H} comm ← (1, 2)

# Implementation (5,500 lines of F#)

**Benchmarks:**

- data center policies (~1600 routers) and backbone policies (~200 routers, many peers/router) from a large cloud provider

- policy from English docs

- Ignoring prefix, customer group and ownership definitions:
  - 31 lines for data center
  - 43 lines for backbone

**Scaling (8 core Windows machine):**

- 10s/pfx (mean) for largest data center

- 45s/pfx (mean) for largest backbone

- 3 minutes total for the backbone

- 9 minutes total for the data center

Propane

→

Regular IR

→

Topology →

Product Graph IR

→

Abstract BGP

CISCO          Juniper

# Thanks