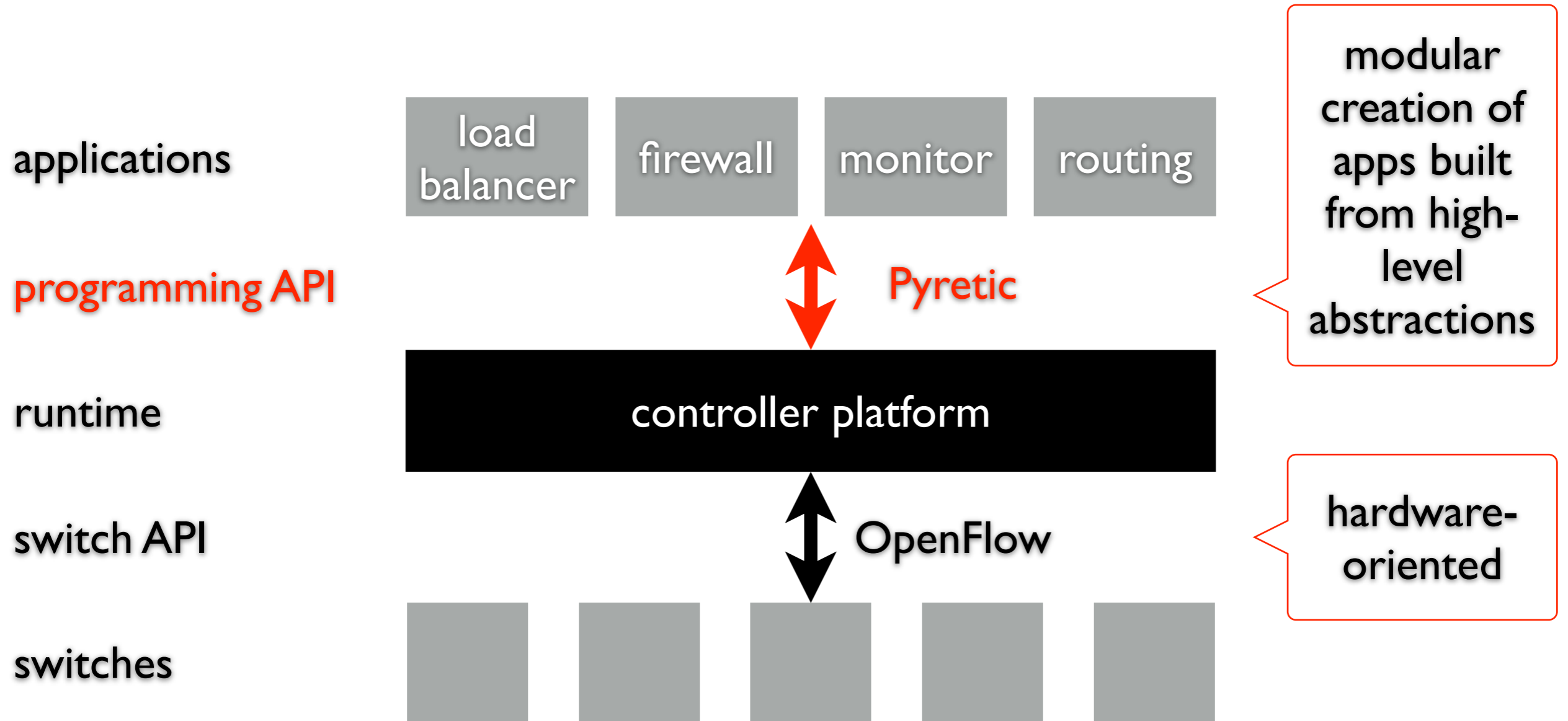# state management

## 5590: software defined networking

anduo wang, Temple University
T 17:30-20:00

# OpenFlow, Pyretic

applications

**load balancer**    **firewall**    **monitor**    **routing**

programming API     Pyretic

runtime     **controller platform**

switch API     OpenFlow

switches

modular creation of apps built from high-level abstractions
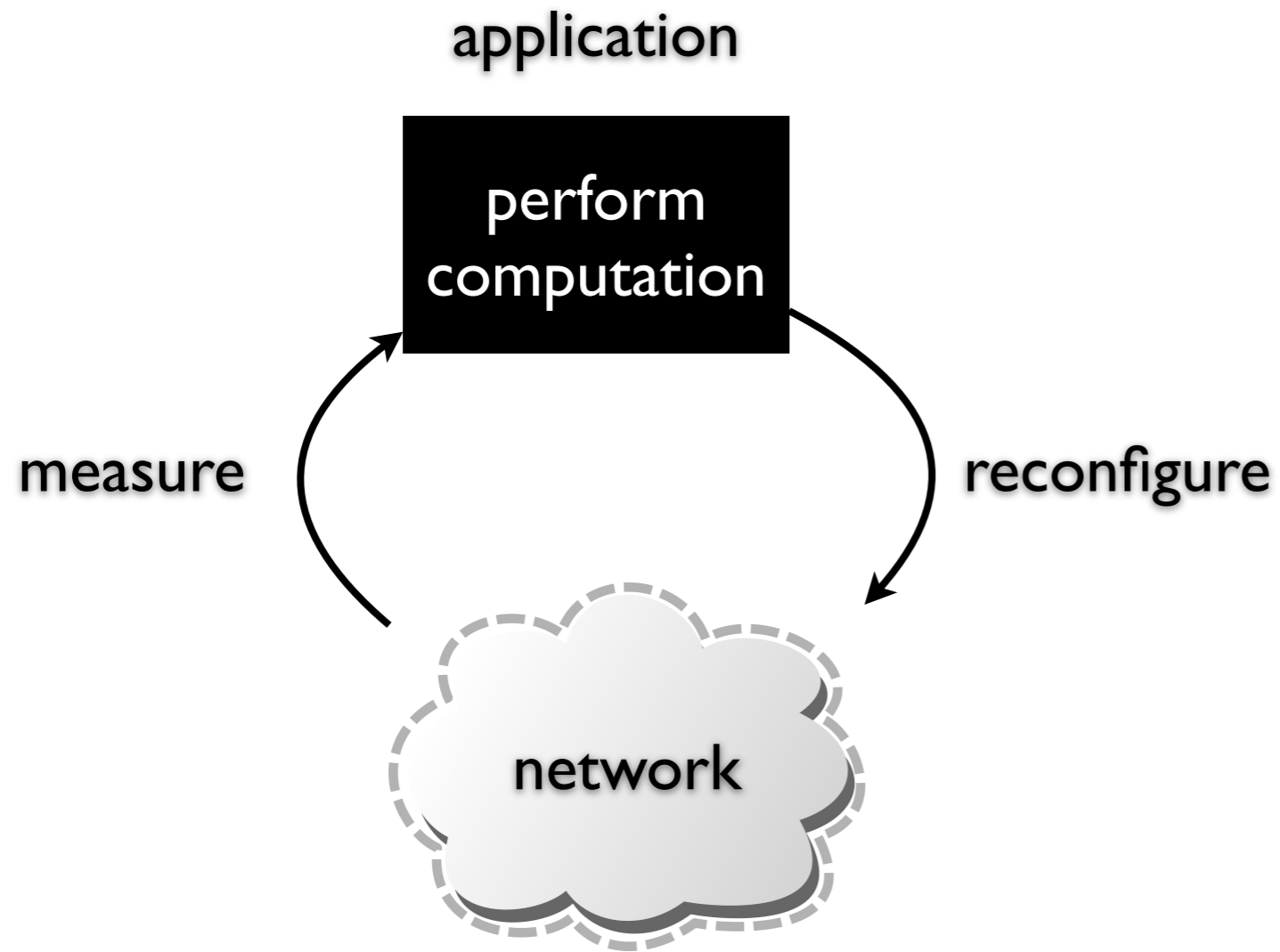
hardware-oriented

# datacenter network (DCN)

runs multiple management applications

- traffic engineering
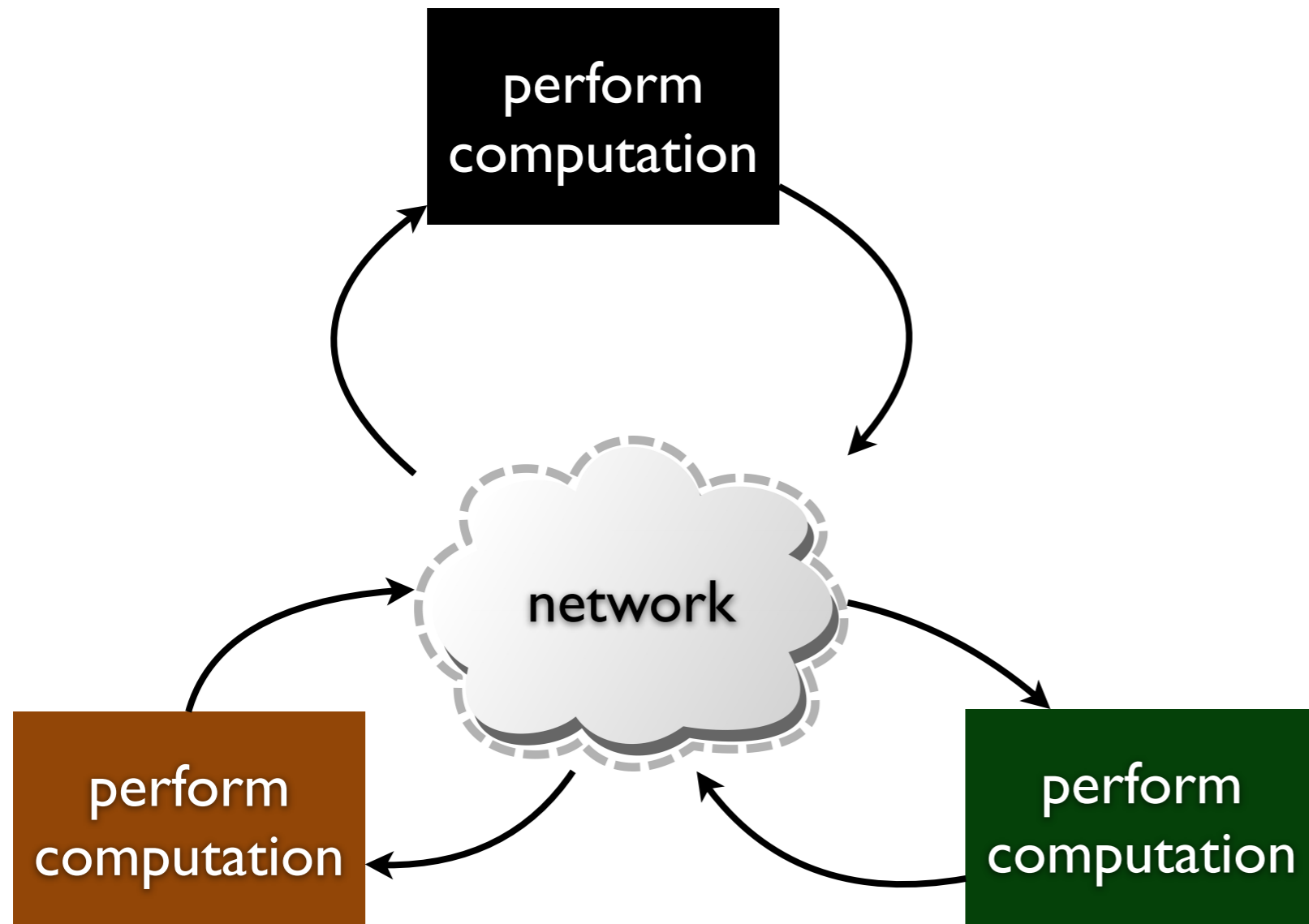- server load balancing
- network virtualization

infrastructure

- failure recovery *NetPilot [SIGCOMM'12]*
- energy saving *Elastic tree [NSDI'10]*
- switch configuration
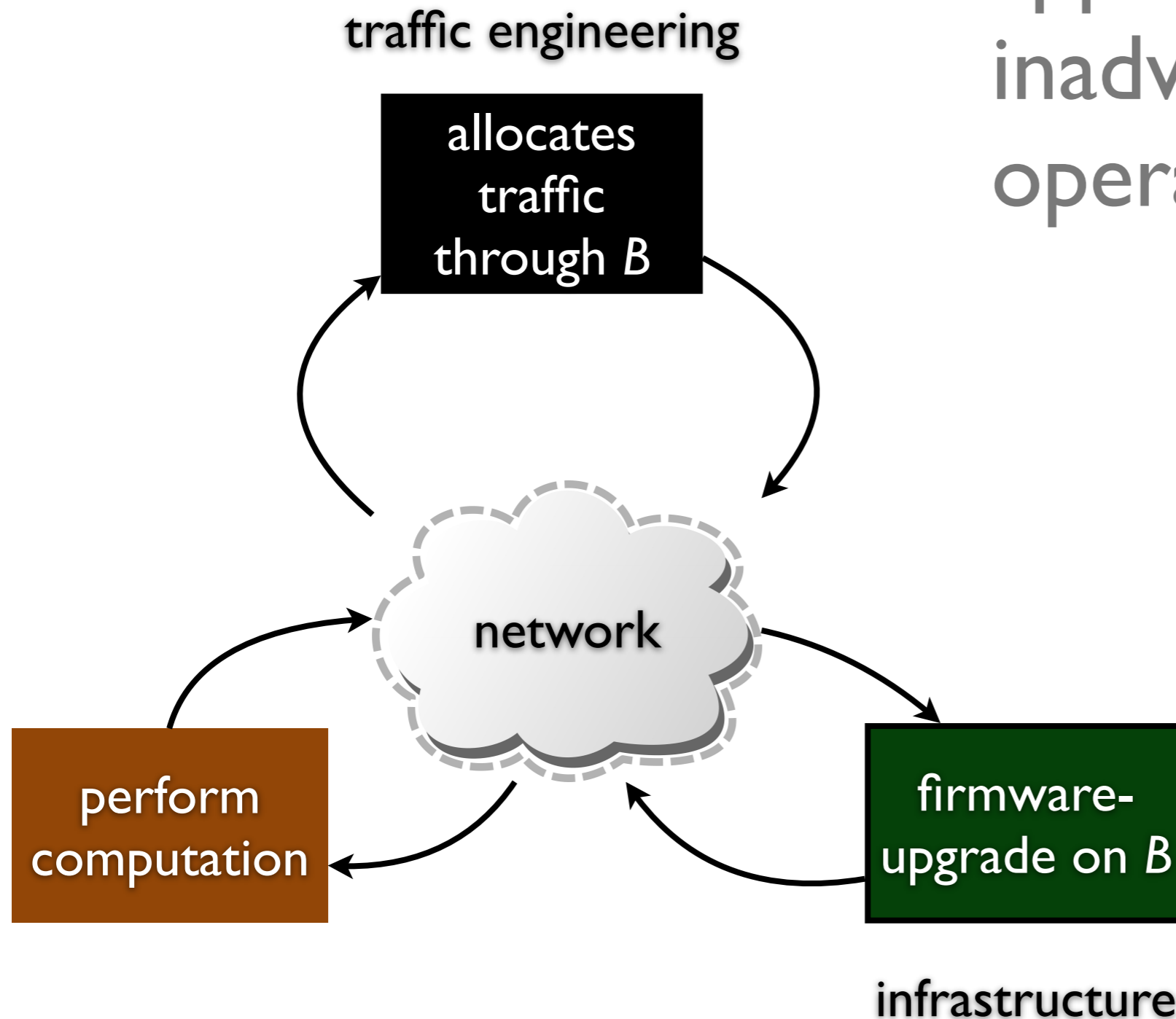
# management applications

# running multiple applications
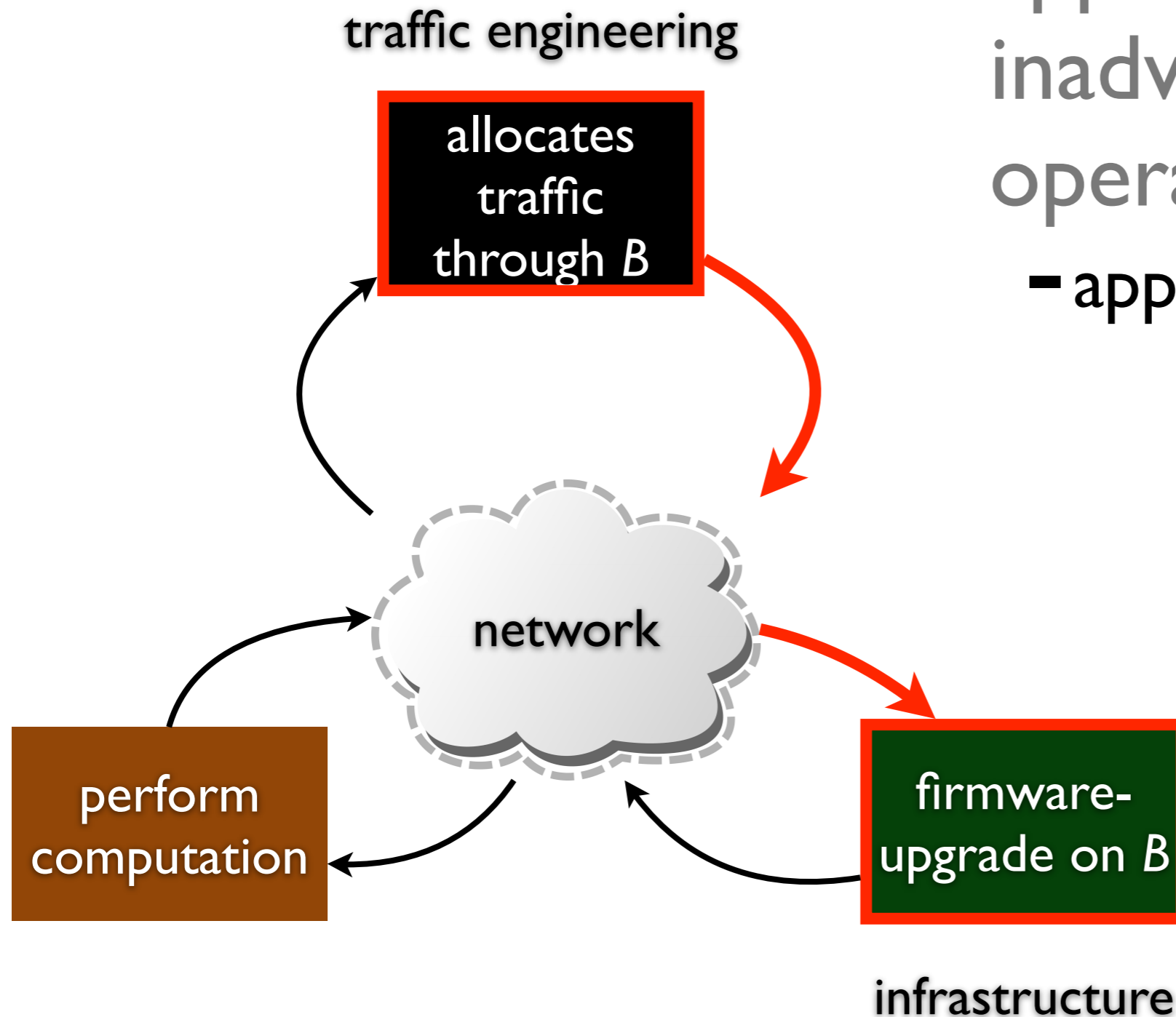
# running multiple applications

applications can inadvertently affect the operations of another

traffic engineering

allocates traffic through *B*

network

perform computation

firmware-upgrade on *B*

infrastructure

# running multiple applications



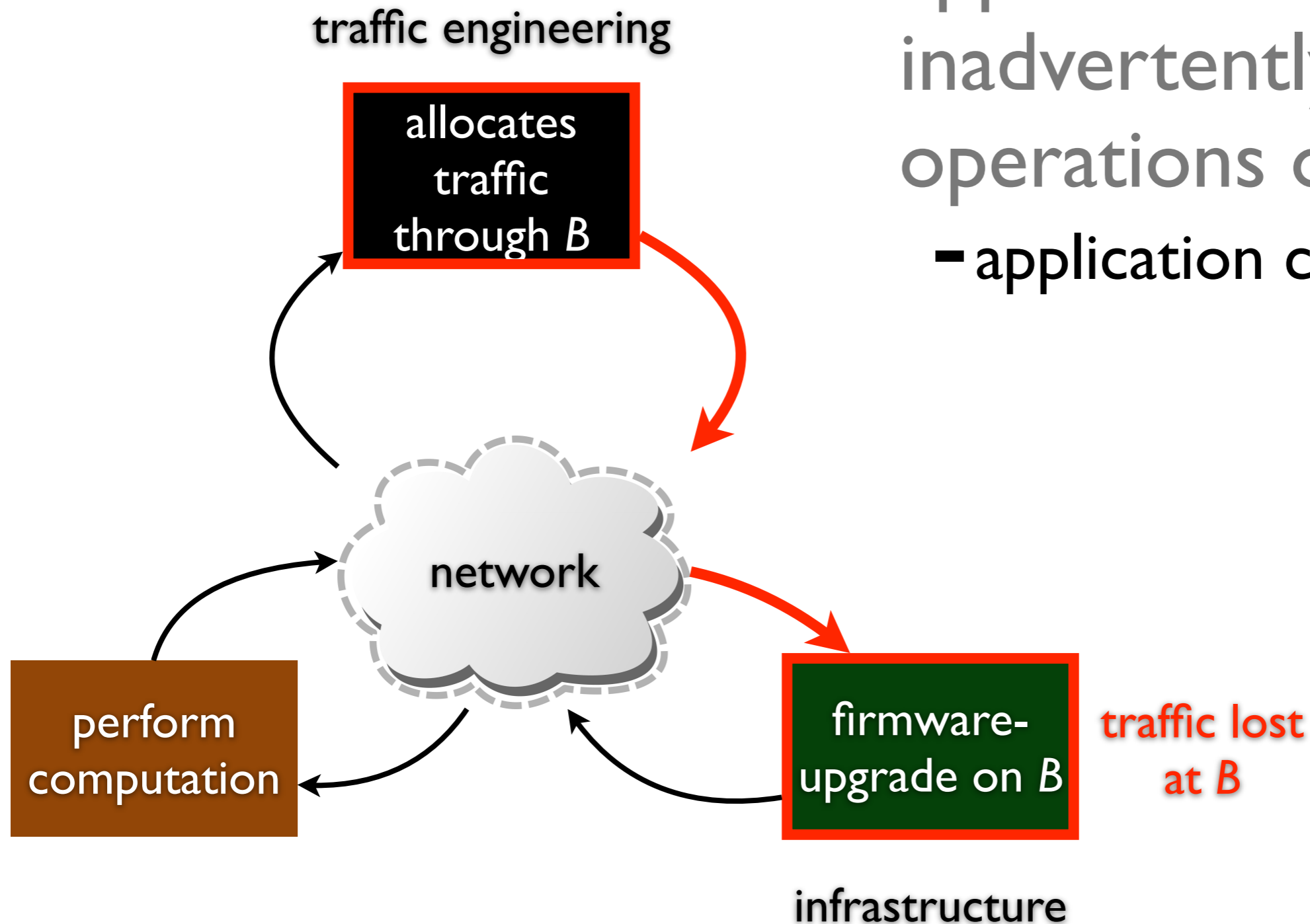applications can inadvertently affect the operations of another

- application conflict
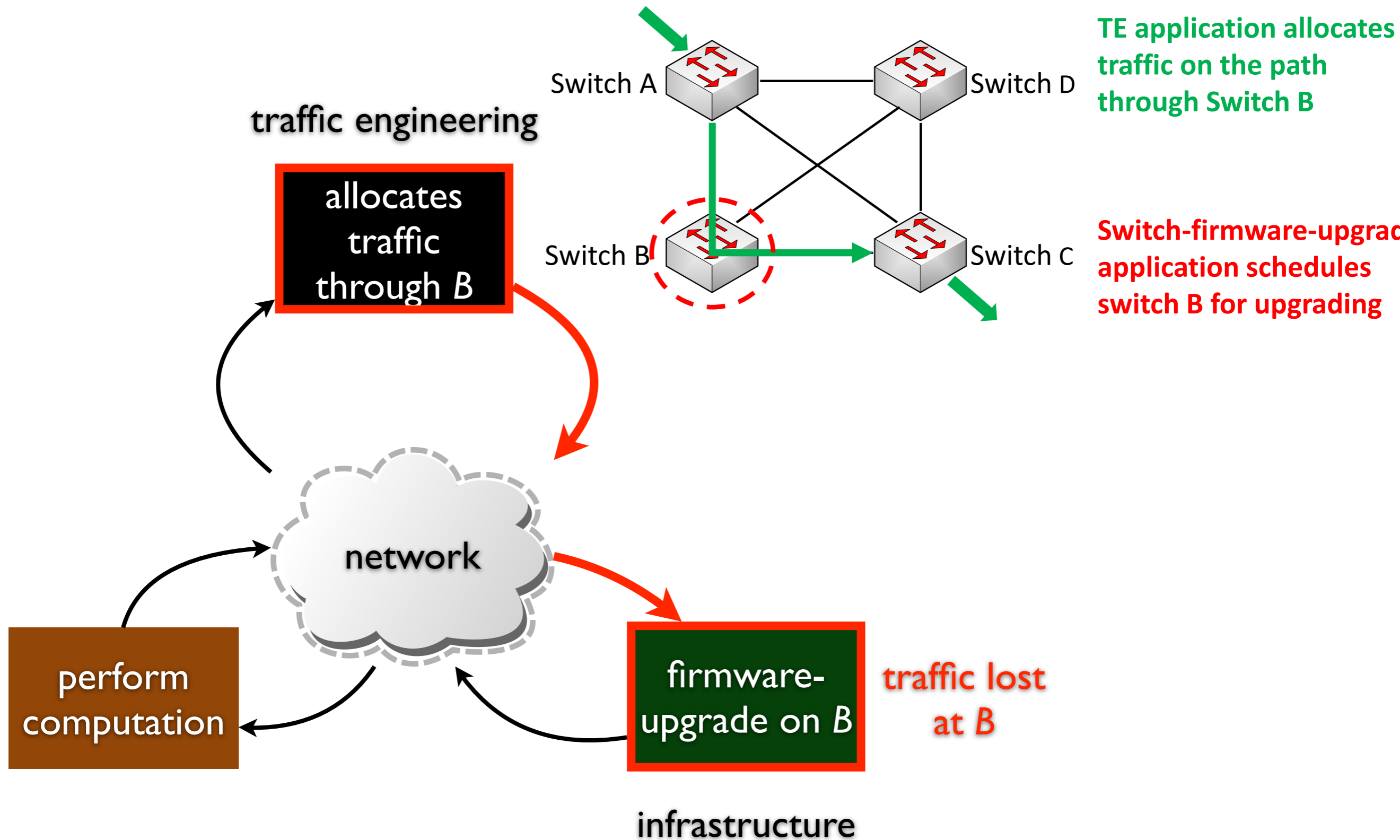
# running multiple applications

applications can inadvertently affect the operations of another
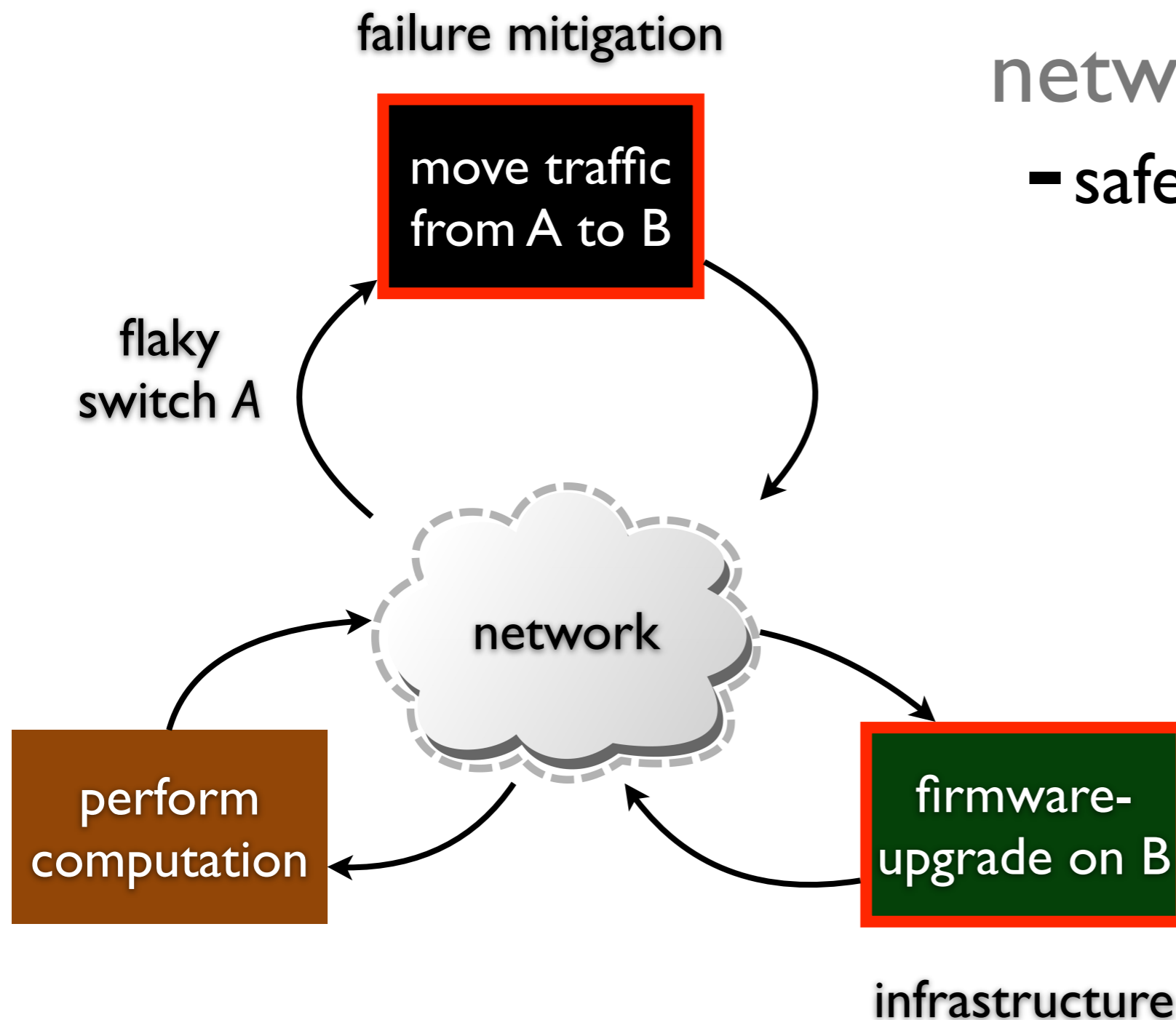
- application conflict

traffic engineering

allocates traffic through *B*

network

perform computation

firmware-upgrade on *B*

traffic lost at *B*

infrastructure

# running multiple applications



traffic engineering

allocates traffic through *B*

network

perform computation

firmware-upgrade on *B*

infrastructure

TE application allocates traffic on the path through Switch B

Switch-firmware-upgrade application schedules switch B for upgrading

traffic lost at *B*

Switch A
Switch D
Switch B
Switch C

8

# running multiple applications



combined effects lead to network-wide failures
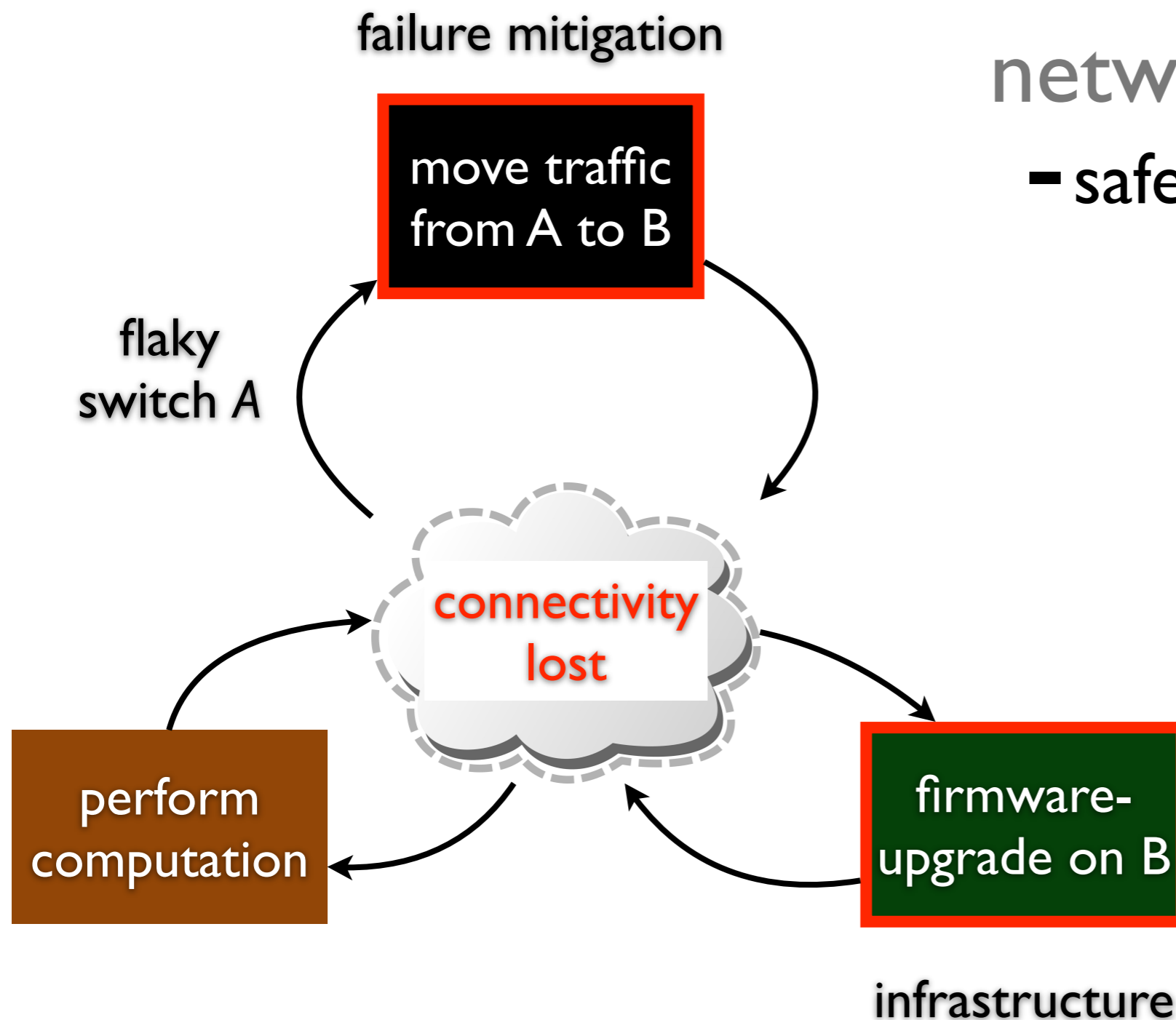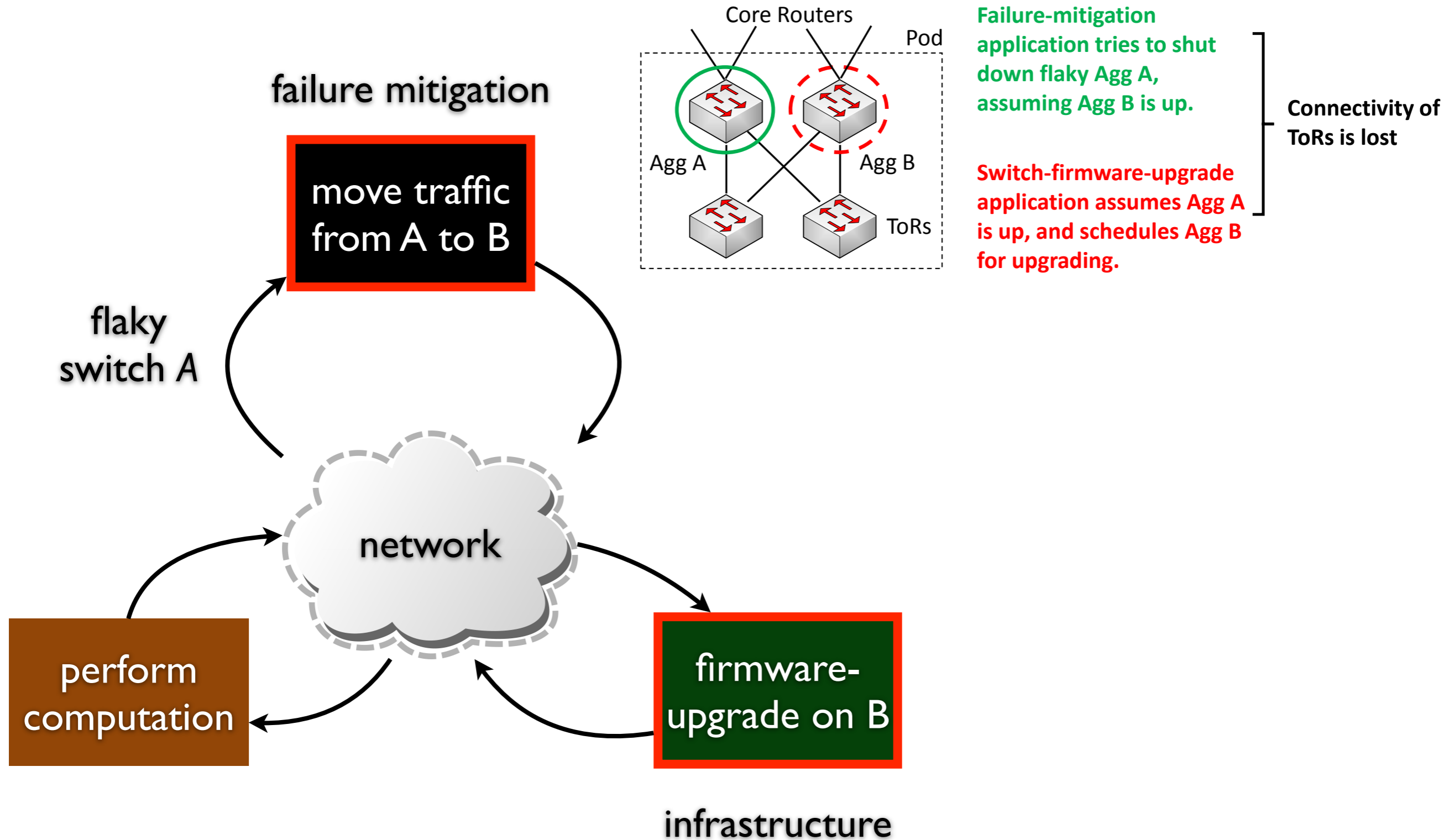- safety failure

# running multiple applications
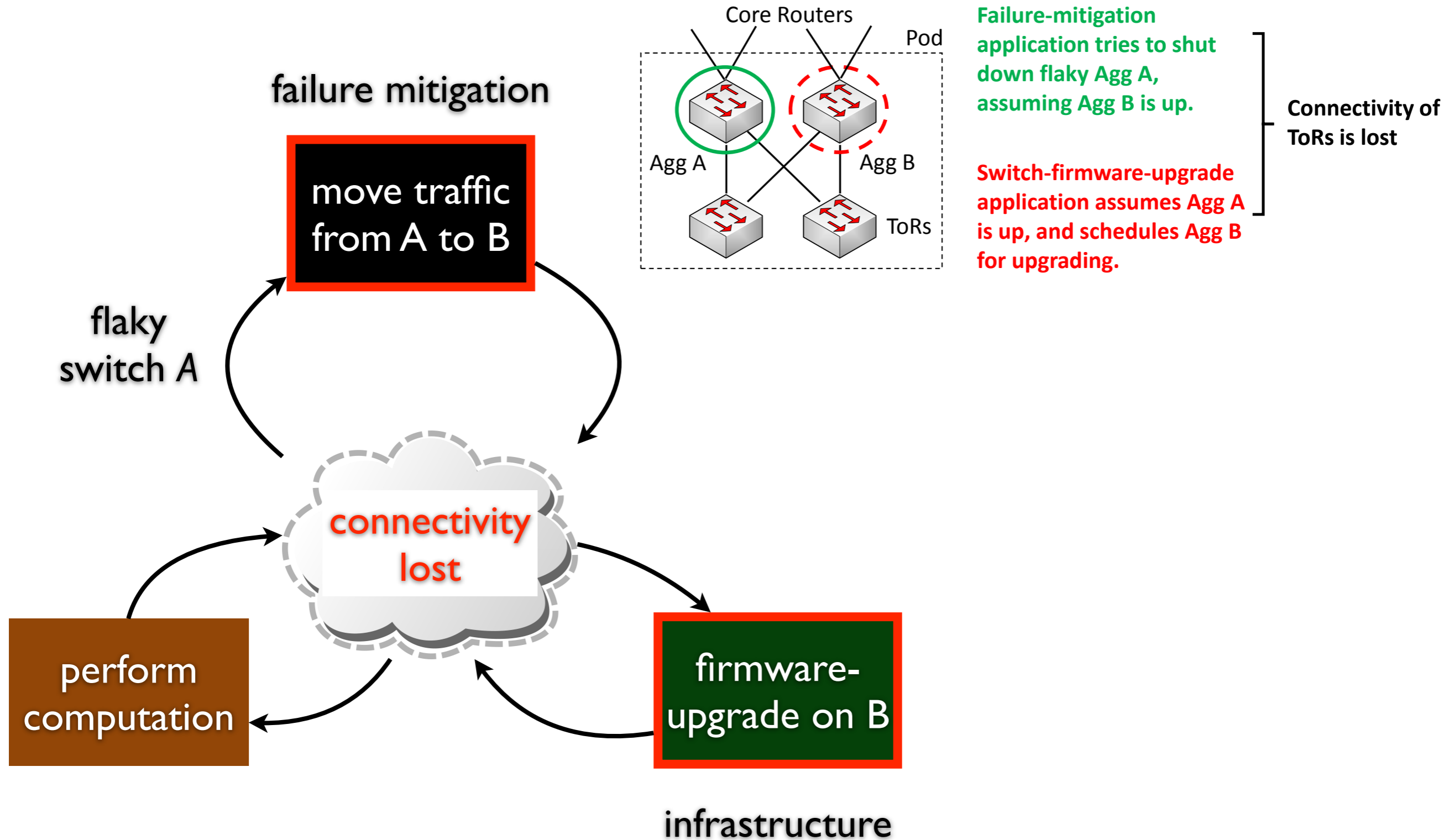


combined effects lead to network-wide failures
- safety failure

# running multiple applications

failure mitigation

**move traffic from A to B**

flaky switch *A*

network

**perform computation**

**firmware-upgrade on B**

infrastructure

Core Routers

Pod

Agg A          Agg B

ToRs

**Connectivity of ToRs is lost**

10

# running multiple applications

failure mitigation



move traffic from A to B

flaky switch A

connectivity lost

perform computation

firmware-upgrade on B

infrastructure

Core Routers

Pod

Agg A

Agg B

ToRs

**Failure-mitigation application tries to shut down flaky Agg A, assuming Agg B is up.**

**Connectivity of ToRs is lost**

**Switch-firmware-upgrade application assumes Agg A is up, and schedules Agg B for upgrading.**

# alternative to running multiple applications

## one single monolithic application

- complex
- explicit coordination
- high overhead on applications

# alternative to running multiple applications

## one single monolithic application

- complex
- explicit coordination
- high overhead on applications

*tightly coupled, repeated extension*

# statesman approach

# statesman approach

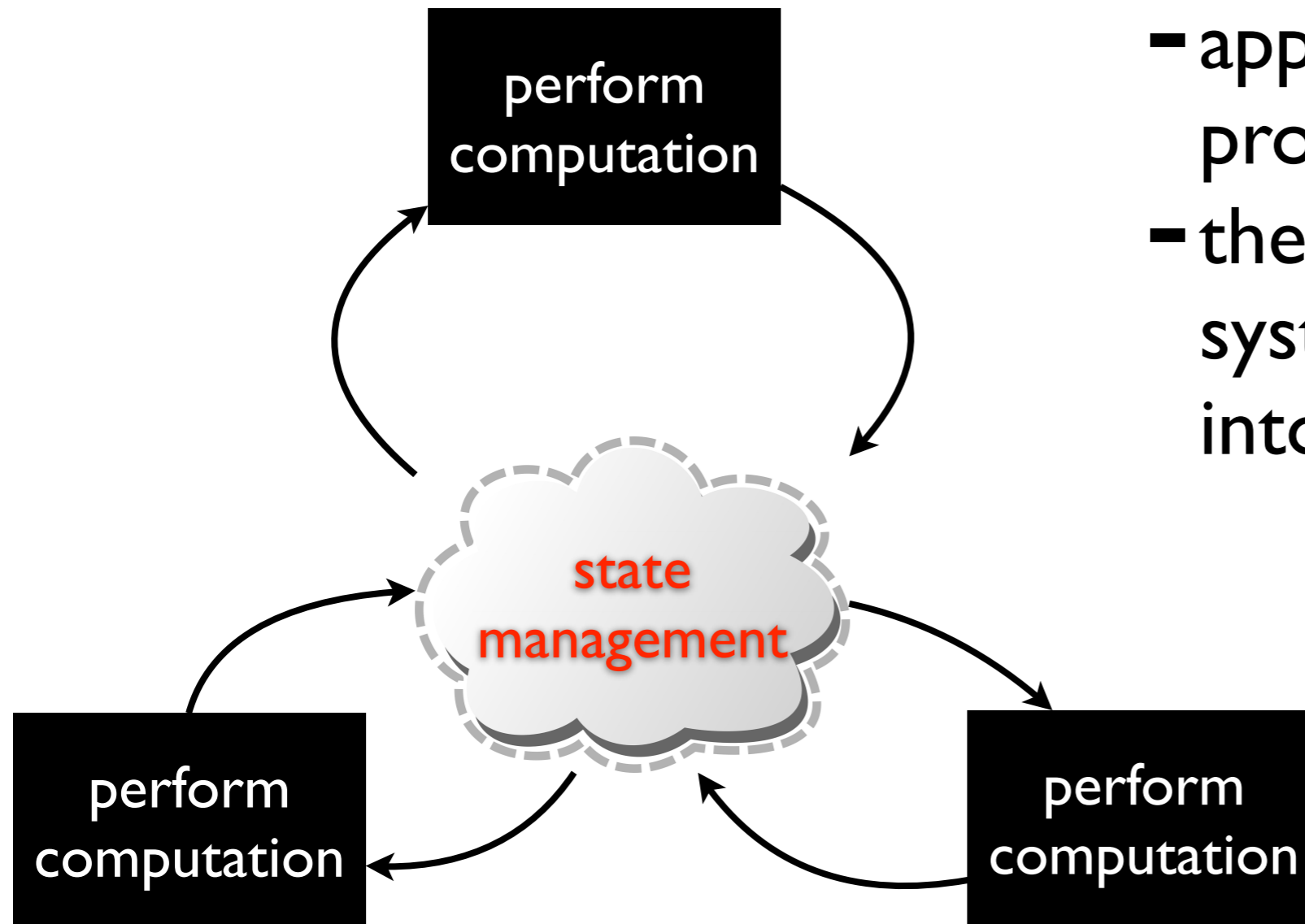build and run applications in a loosely coupled manner

# statesman approach

build and run applications in a loosely coupled manner

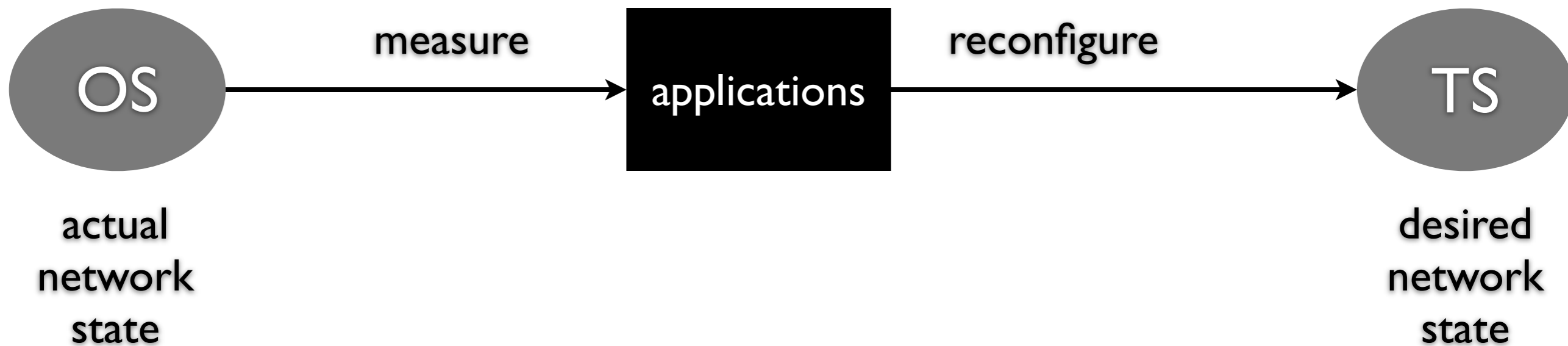introduce a separate (state) management system

- conflict resolution
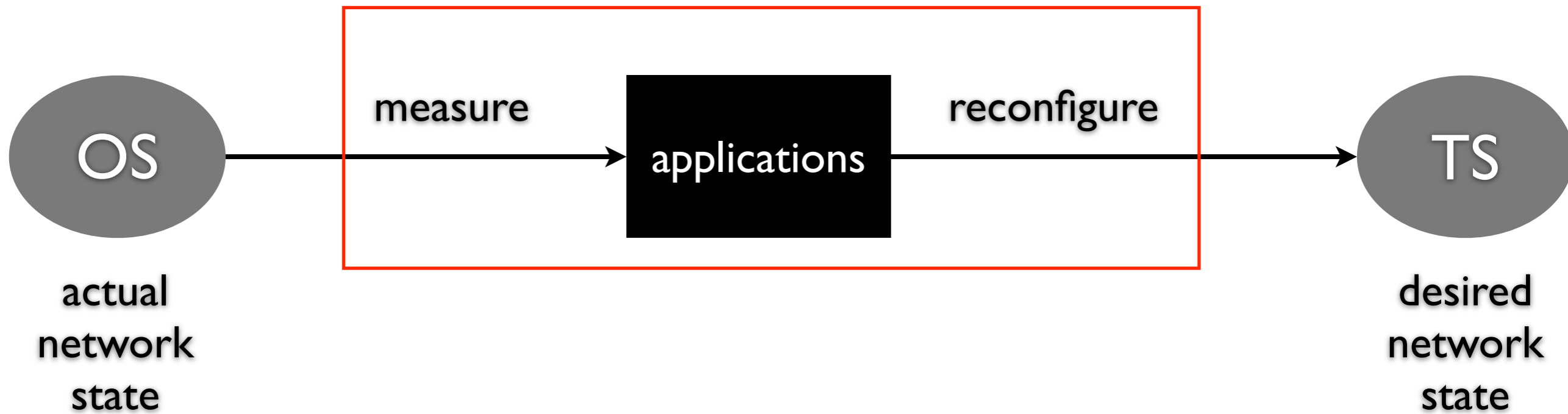- invariant enforcement

# statesman approach



- applications "pull" observed states (OS)
- applications "push" proposed states (PS)
- the separate statesman system "merges" the states into target states (TS)
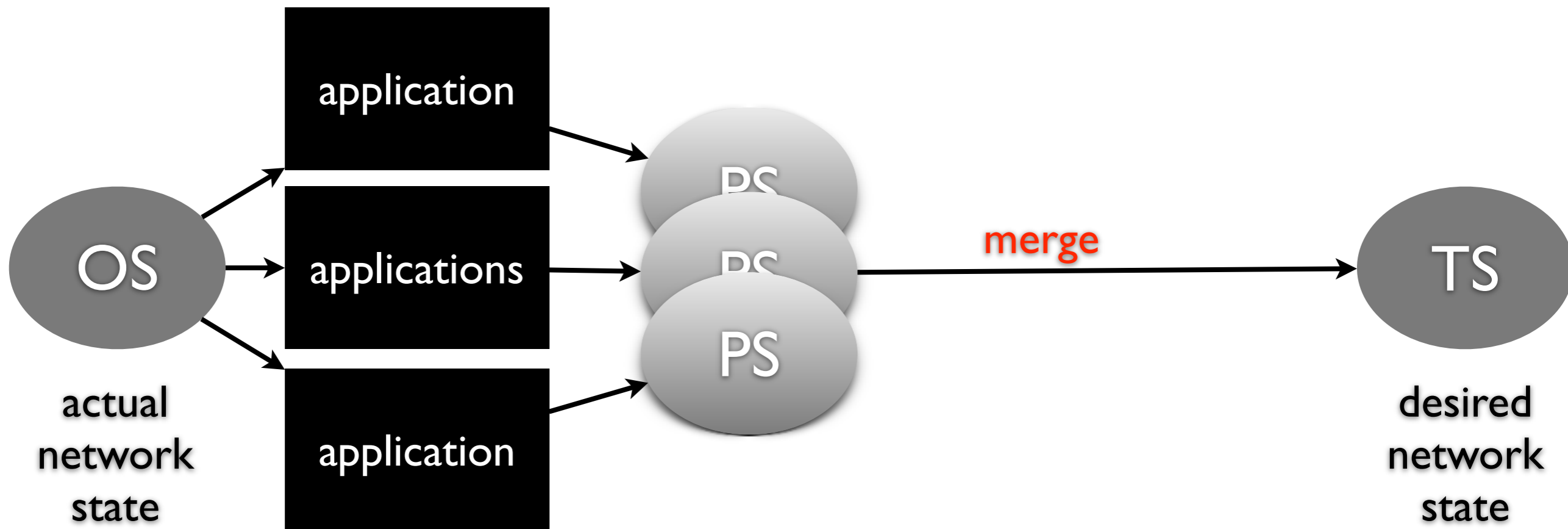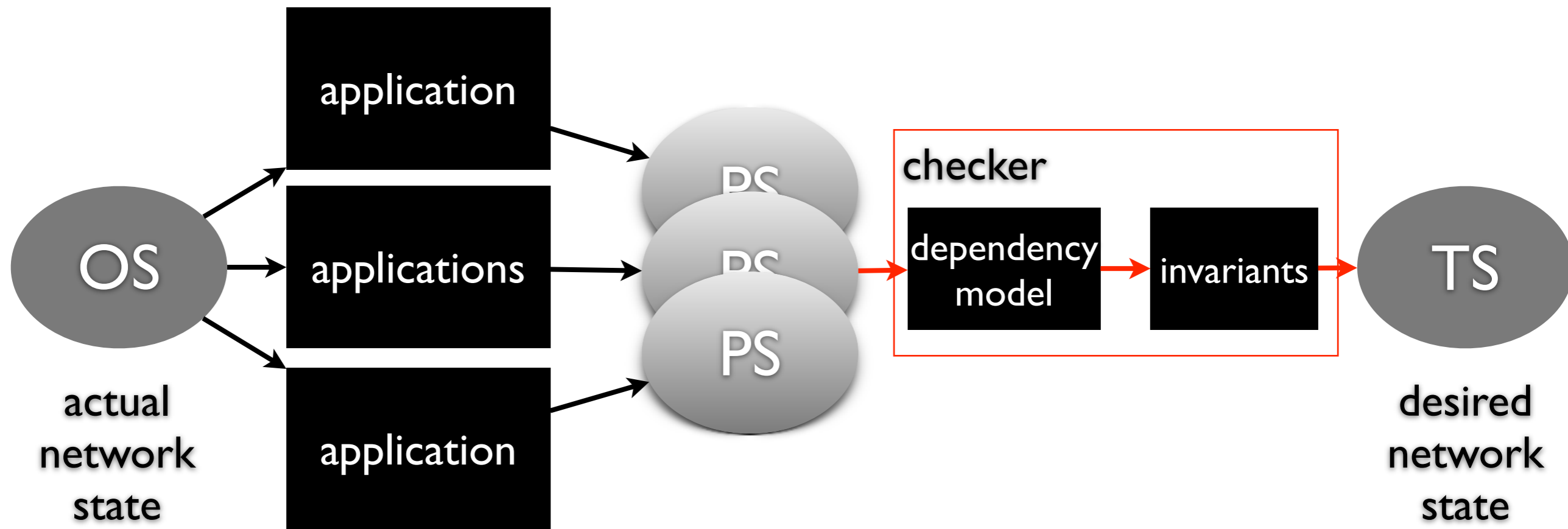
# statesman approach

OS

*measure*

applications

*reconfigure*

TS

actual
network
state

desired
network
state

# statesman approach



OS
actual network state

measure

applications

reconfigure

TS
desired network state

# statesman approach

# statesman approach

checker

# checker

## use state dependency model

- determine whether PSes applicable to existing OSes
- detect and resolve conflicts among PSes
- form TSes

# checker

## use state dependency model

- determine whether PSes applicable to existing OSes
- detect and resolve conflicts among PSes
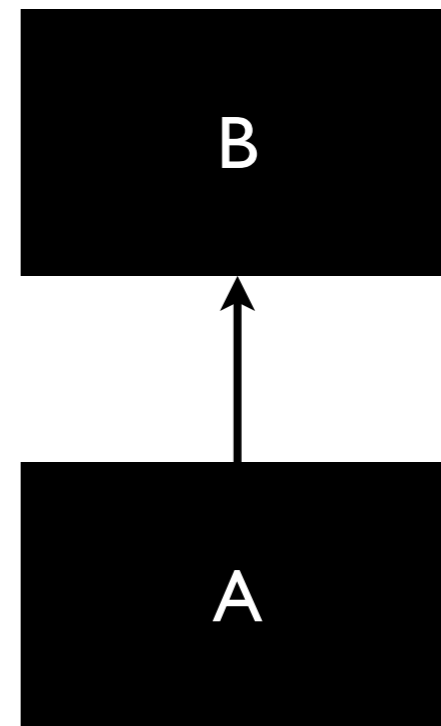- form TSes

## use operator-specified invariants

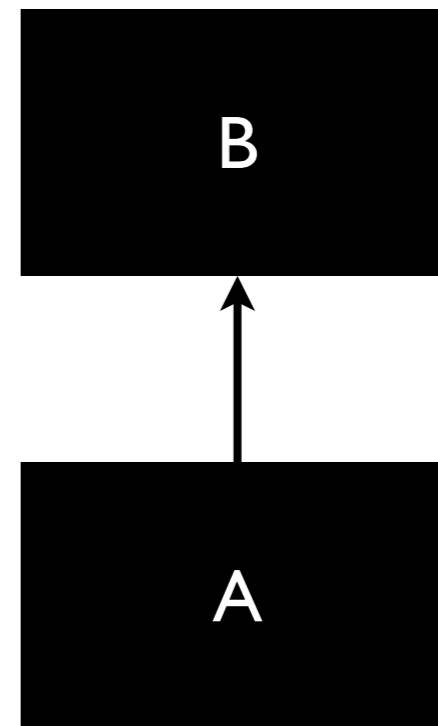- examine the TSes against the invariants

# state dependency model

*state variables in one application's PS can depend on state variables in another application's PS*

## B depends on A

- A is a prerequisite for <span style="color:red">writing</span> B states
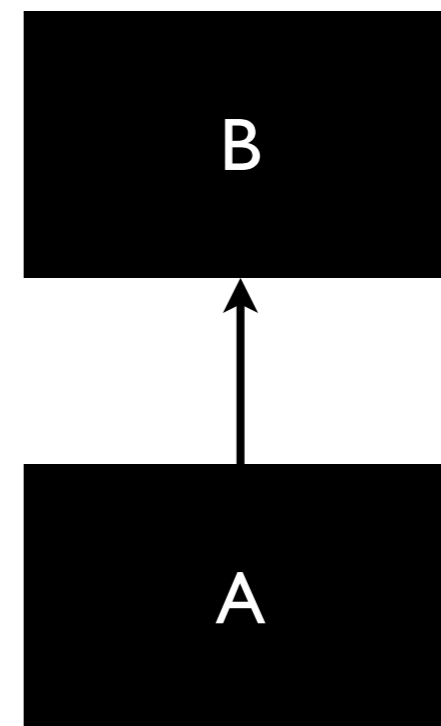
# state dependency model

*state variables in one application's PS can depend on state variables in another application's PS*

## B depends on A

- A is a prerequisite for writing B states
- B is controllable only if A value is appropriate

# state dependency model

*state variables in one application's PS can depend on state variables in another application's PS*
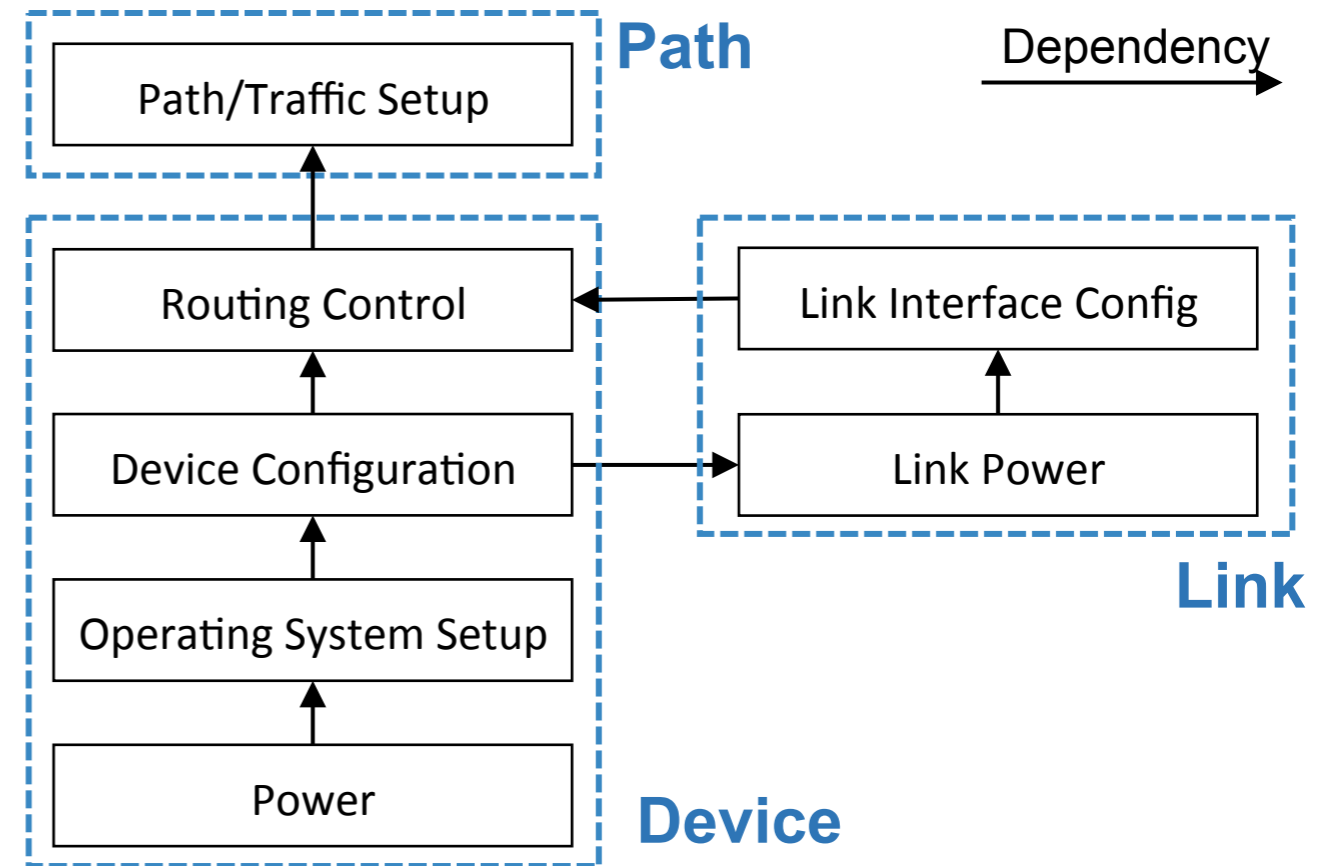
## B depends on A

- A is a prerequisite for writing B states
- B is controllable only if A value is appropriate

## conflicts

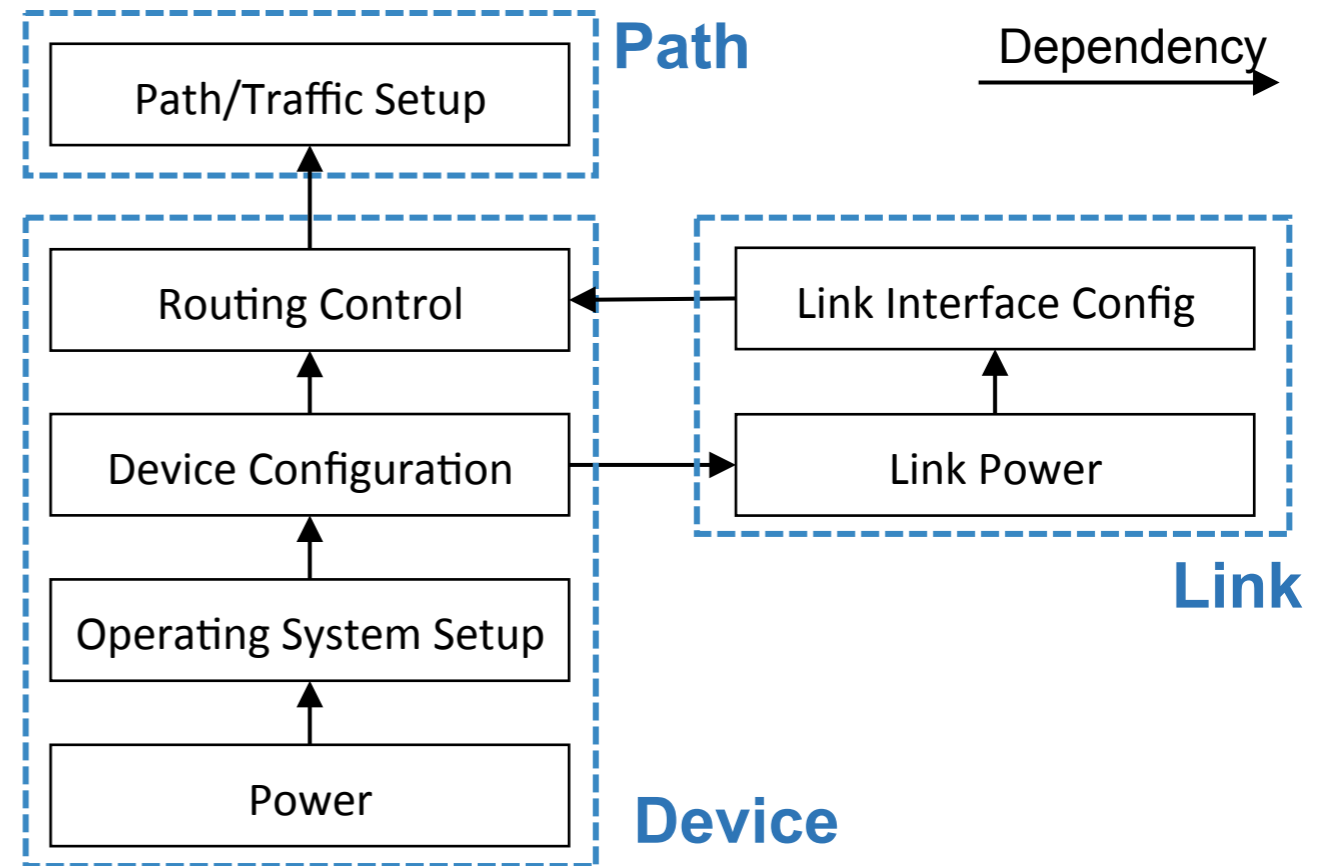- B is uncontrollable due to state (or state change) in A

# using state dependency model



**Path**

Path/Traffic Setup

**Dependency**

Routing Control ← Link Interface Config

Device Configuration → Link Power

Operating System Setup
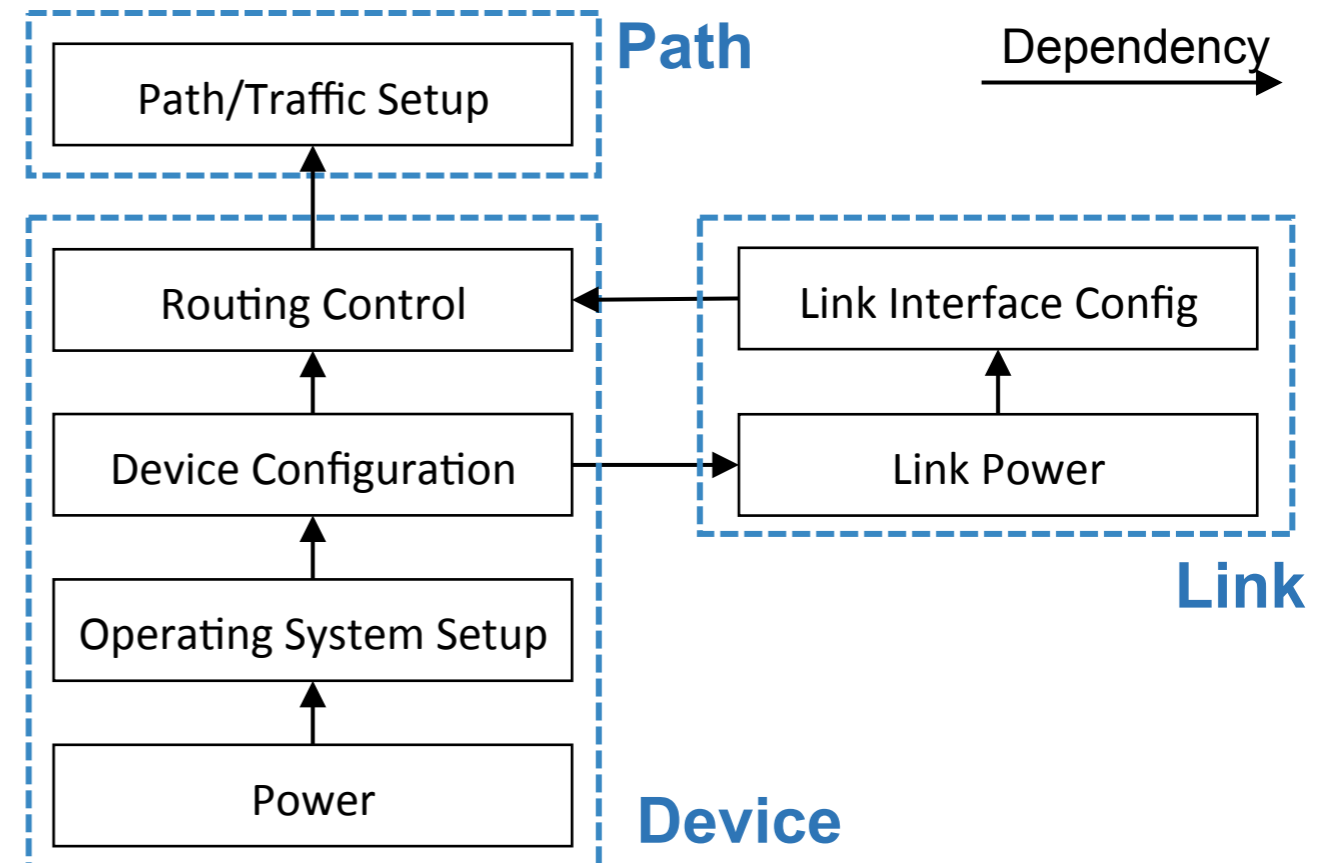
Power

**Device**

**Link**

# using state dependency model

## Statesman exposes

- B's (latest) value
- together with a logical controllability variable
  - "1" only if all B's parents are controllable

# using state dependency model

## Statesman exposes

- B's (latest) value
- together with a logical controllability variable
  - "1" only if all B's parents are controllable

## question

- how to <span style="color:red">extend</span> the dependency model?
- <span style="color:red">advantage</span> of having an explicit separate model?



**Path**

Path/Traffic Setup

Dependency →

Routing Control ← Link Interface Config

Device Configuration → Link Power

**Link**

Operating System Setup

Power

**Device**

# resolving conflicts

# resolving conflicts

## TS-OS, PS-OS

- conflicts due to the changing OS
  - makes some variables in TS/PS uncontrollable
- solution: simply reject

# resolving conflicts

## TS-OS, PS-OS

- conflicts due to the changing OS
  - makes some variables in TS/PS uncontrollable
- solution: simply reject

## PS-TS

- a PS can conflict with the TS due to an accepted PS from another application
  - TS is really just the accumulation of all accepted in the past
- solution
  - accept with last-write-wins / priority-based locking
  - at the level of individual switches and links

# maintaining invariants

## what

- invariants: infrastructure's operational stability, independent of apps
  - suffice to safeguard the network & not too stringent with app goals
- examples: connectivity, capacity

## how

- checking TS against invariants
  - maintain a base network state graph using values from the OS
  - compute difference between TS and OS
  - check invariants on the new network state

# discussion

## making multiple applications coexist

- ONIX, NOX: no support
- Pyretic, Pane, Maple: compose target traffic management applications
- Corybantics: hosting multiple applications on isolated slices

# statesman: use cases, evaluations …

# statesman deployment

10 geographically-distributed datacenters (DCs)

- cover switches, links within each DC and across DC (WAN)

three applications

- switch-upgrade
- failure-mitigation
- inter-DC TE

# challenges—maintaining globally available and distributed states

- inter-DC
  - due to WAN failures, DCs may be disconnected
- within-DC
  - huge volume of state data: hundreds of thousands of  switches and links
  - millions of state variables

# challenges—updating DCN states

- heterogeneity: diverse range of network elements expose heterogenous interfaces for updates
- device can fail during an update
- device respond slow, dominating the application control loop

solution—maintaining globally available and distributed states

solution—maintaining globally available and distributed states

partitioning checker's responsibility into impact groups

- one impact group per DC
- one additional impact group with border routers of all DCs and the WAN links

solution—maintaining globally available and distributed states

partitioning checker's responsibility into impact groups

- one impact group per DC
- one additional impact group with border routers of all DCs and the WAN links

partitioning monitor

- split monitor's responsibility into many instances
  - each covers 1k switches

# solution—updating DCN states

# solution—updating DCN states

heterogeneity

- OpenFlow and command templates

# solution—updating DCN states

heterogeneity

- OpenFlow and command templates

dynamic failures

- stateless updates
- simply push to the devices the latest OS-TS difference

# use case: maintaining invariants



**CORE** — 1 ... 4

**AGG** — [1 ... 4] ... ✗[1 ... 4] ... [1 ... 4]

**ToR** — [1 ... n] [1 ... n] [1 ... n]

**Pod 1**     **Pod 4**     **Pod 10**

✗ Link with FCS error

switch_upgrade and failure_mitigation coexist

statesman goal: maintaining capacity invariant

- 99% ToR pairs have at least 50% capacity

# use case: maintaining invariants



CORE
AGG
ToR

Pod 1    Pod 4    Pod 10

✕ Link with FCS error

## one DC with 10 pods

- each pod has 4 AGGs and a number of ToRs

## switch_upgrade

- upgrade all 40 AGGs
- (sequentially) pod by pod
- attempt parallel upgrades within each pod

# use case: maintaining invariants



## 90 ToR pairs

- one ToR from each pod
- put the 9 ToR pairs from the same pods together
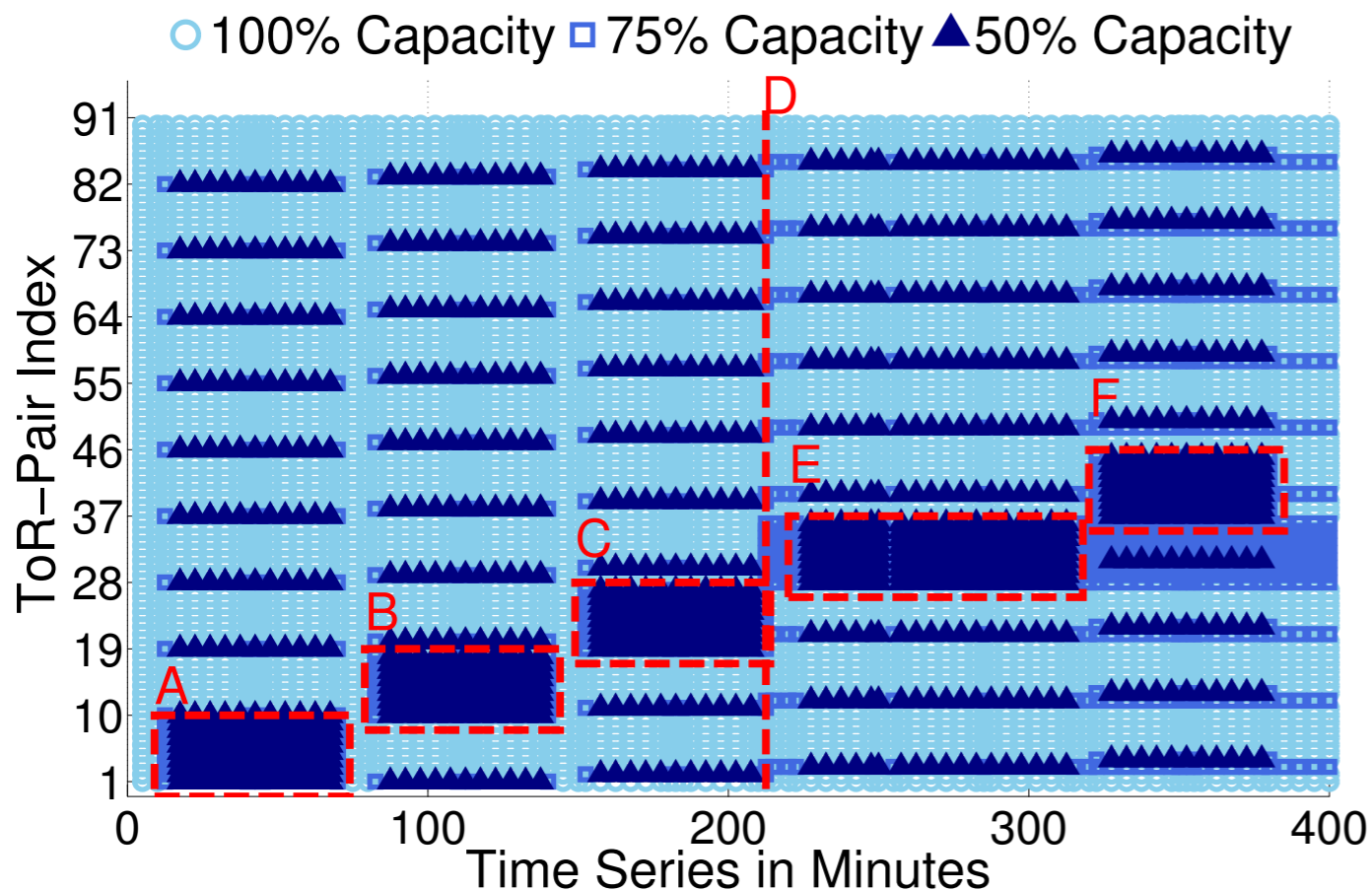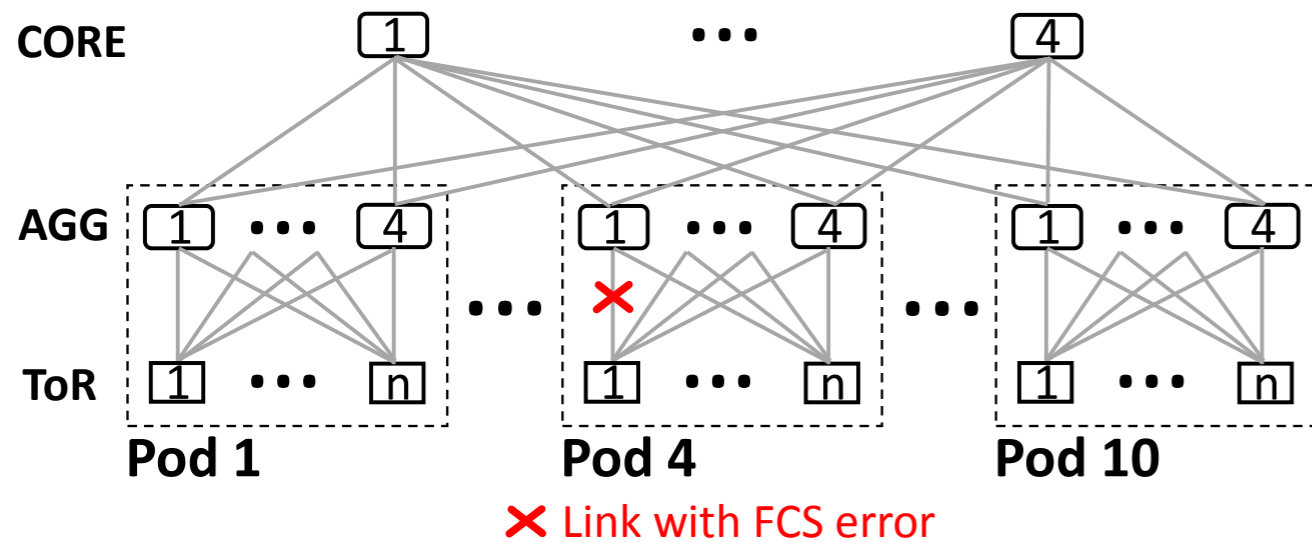
# use case: maintaining invariants



CORE

AGG

ToR

Pod 1    Pod 4    Pod 10

✗ Link with FCS error

## 90 ToR pairs

- one ToR from each pod
- put the 9 ToR pairs from the same pods together



○ 100% Capacity  □ 75% Capacity  ▲ 50% Capacity



BR = Border Router

DC 2
BR 3    BR 4

BR 1
DC 1
BR 2

BR 6    BR 5
DC 3

# use case: maintaining invariants



CORE

AGG

ToR

Pod 1     Pod 4     Pod 10

✗ Link with FCS error

## 90 ToR pairs

- one ToR from each pod
- put the 9 ToR pairs from the same pods together



○ 100% Capacity  □ 75% Capacity  ▲ 50% Capacity

ToR–Pair Index

Time Series in Minutes



BR = Border Router

DC 2
BR 3    BR 4

BR 1

DC 1

BR 2

DC 3
BR 6    BR 5

# use case: maintaining invariants



CORE, AGG, ToR diagram with Pod 1, Pod 4, Pod 10

✖ Link with FCS error

## 90 ToR pairs

- one ToR from each pod
- put the 9 ToR pairs from the same pods together



○ 100% Capacity □ 75% Capacity ▲ 50% Capacity

ToR–Pair Index vs Time Series in Minutes



BR = Border Router

DC 1, DC 2, DC 3, BR 1, BR 2, BR 3, BR 4, BR 5, BR 6

# use case: resolving conflicts

BR = Border Router



## setup

- 8 border routers (BRs)
- 24 (12 physical links x 2 directions) inter-DC links
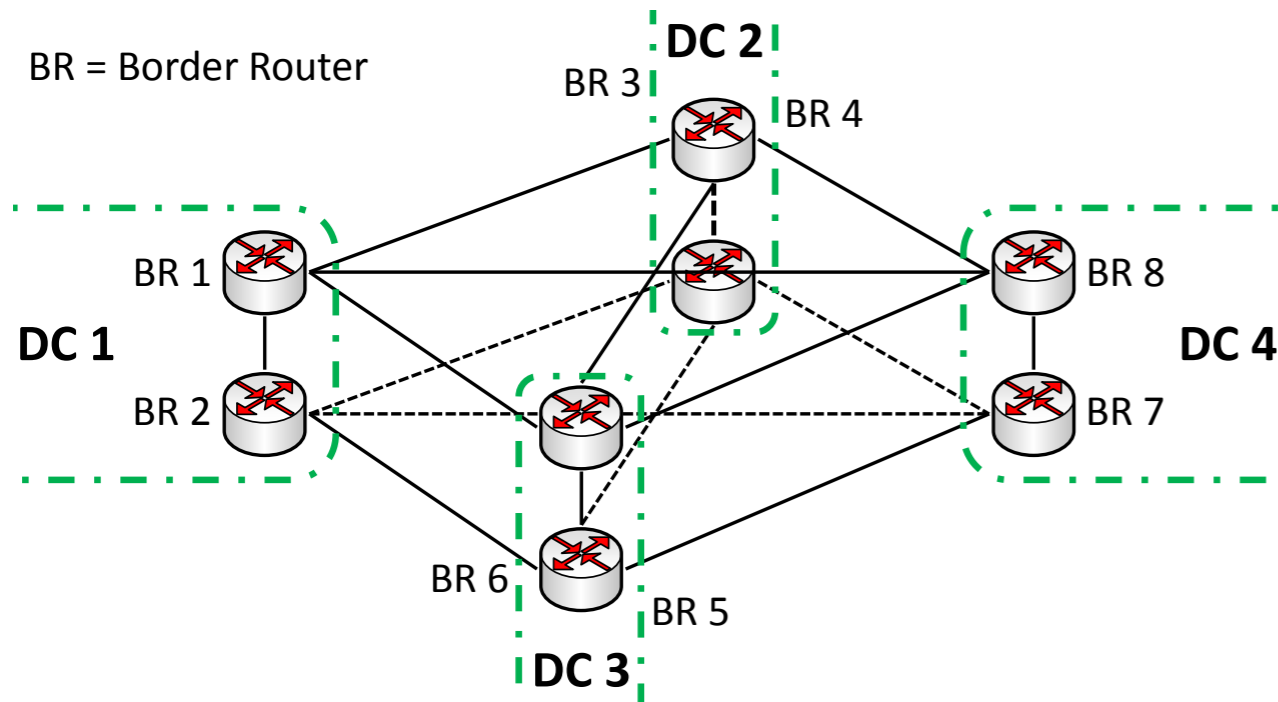
## goal

- upgrade BRs while inter-DC TE is on

# use case: resolving conflicts



BR = Border Router

DC 2
BR 3
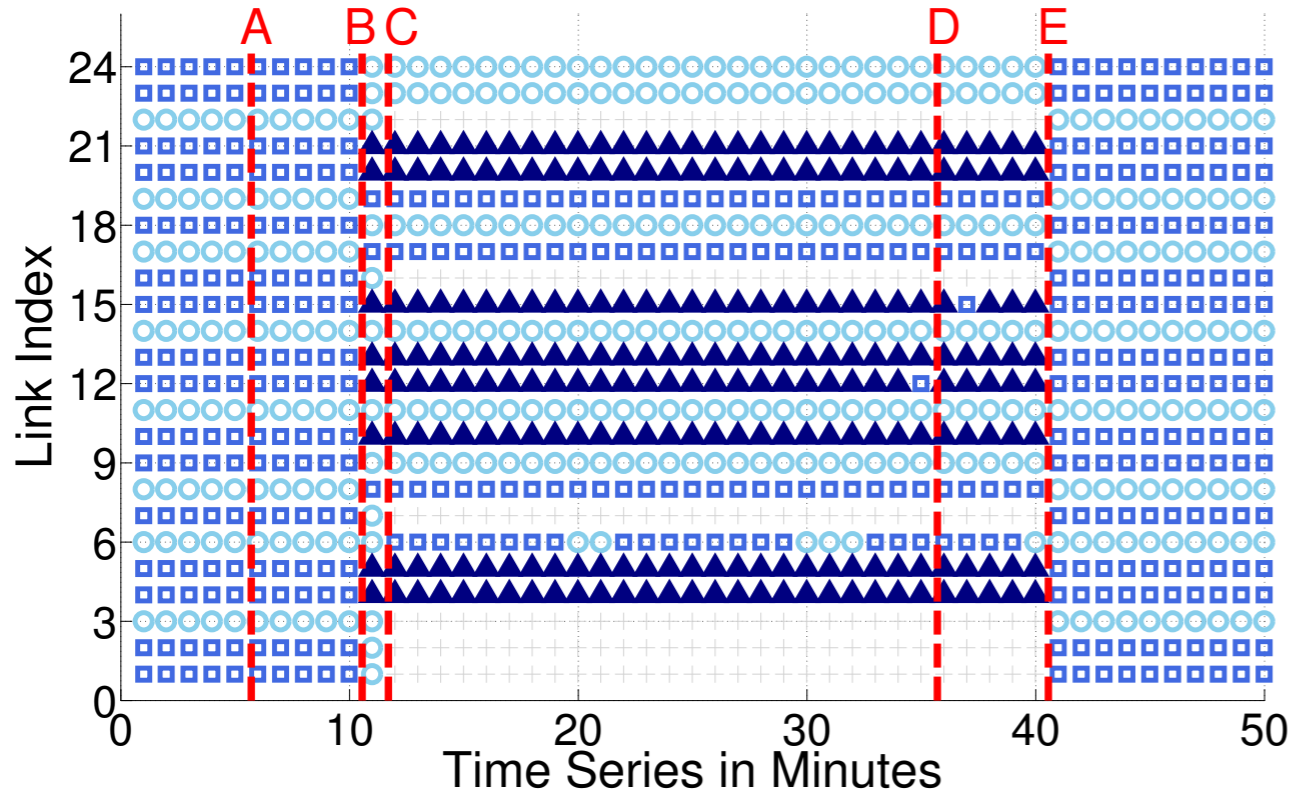BR 4
BR 1
DC 1
BR 2
BR 8
DC 4
BR 7
BR 6
BR 5
DC 3

solution: statesman coordinates, by locks, swtich_upgrade, TE

- assign TE low-level lock
- switch_upgrade high-level lock


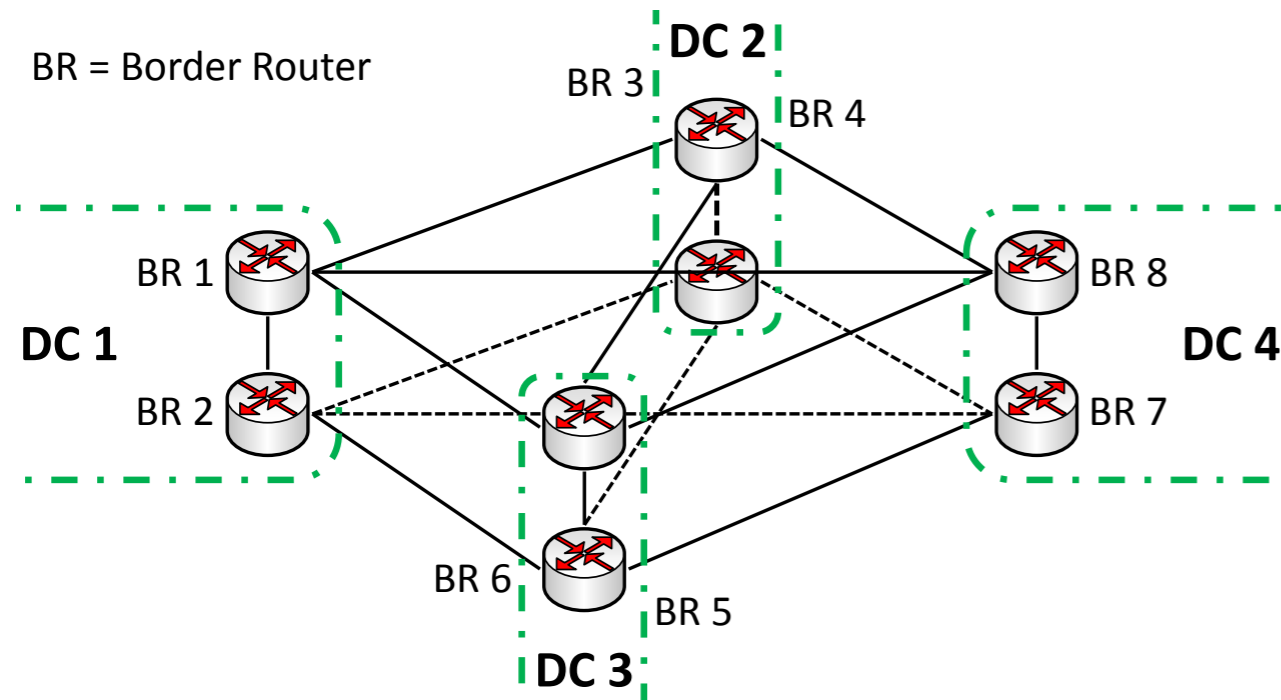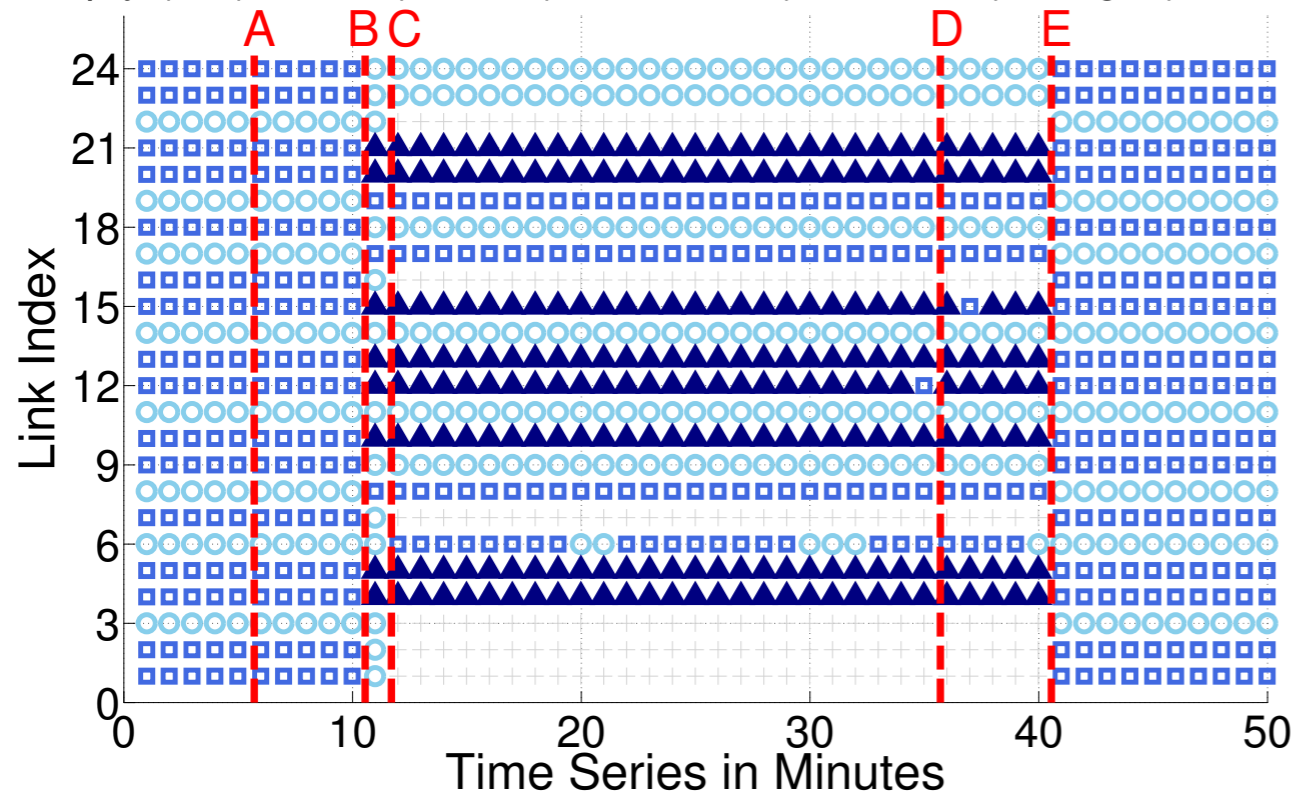
Empty (0%)  Low (1~40%)  Medium (40%~80%)  High (80%~100%)

Link Index

Time Series in Minutes

BR = Border Router



DC 2
BR 3    BR 4

BR 1

DC 1

BR 2

BR 6    BR 5

DC 3

DC 4

BR 8

BR 7

licts

ordinates

ade, TE

-level lock

de high-level



+ Empty (0%)  ○ Low (1~40%)  □ Medium (40%~80%)  ▲ High (80%~100%)

A    B C            D    E

A

= switch_upgrade acquires high-level lock on BR₁

Link Index

Time Series in Minutes

40

BR = Border Router

DC 2

BR 3    BR 4

BR 1

DC 1

BR 2

BR 8

DC 4

BR 7

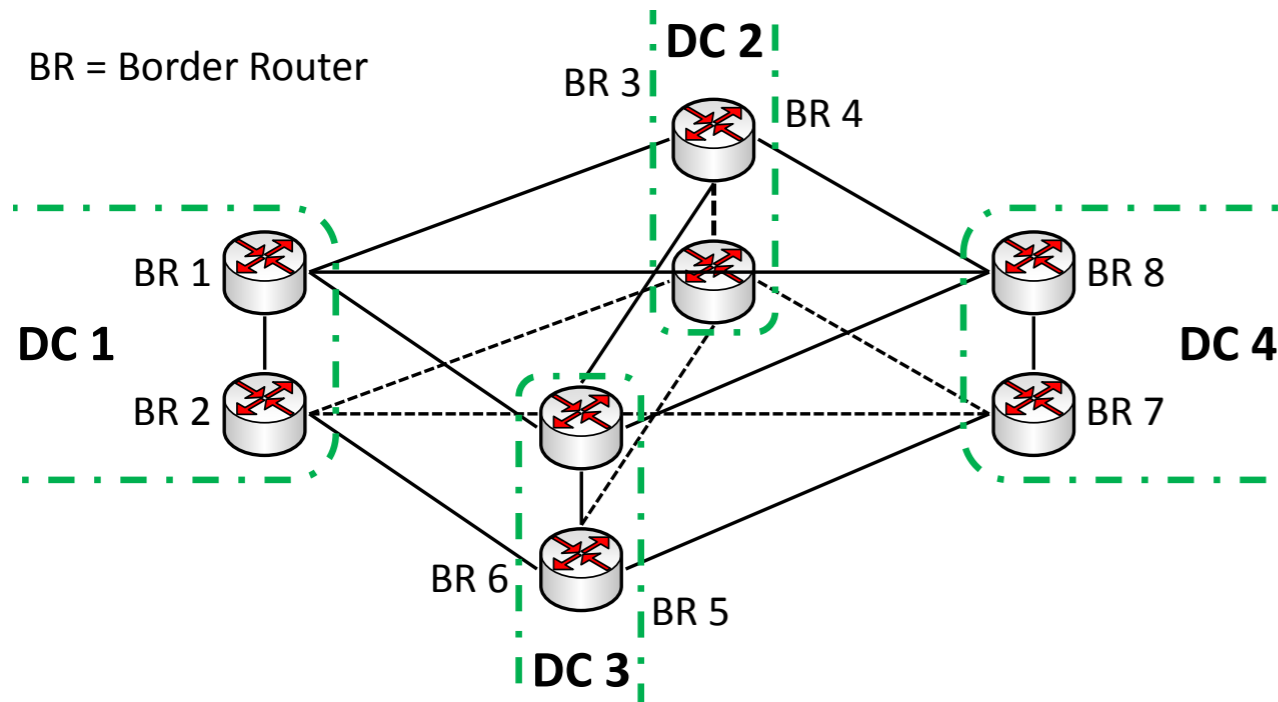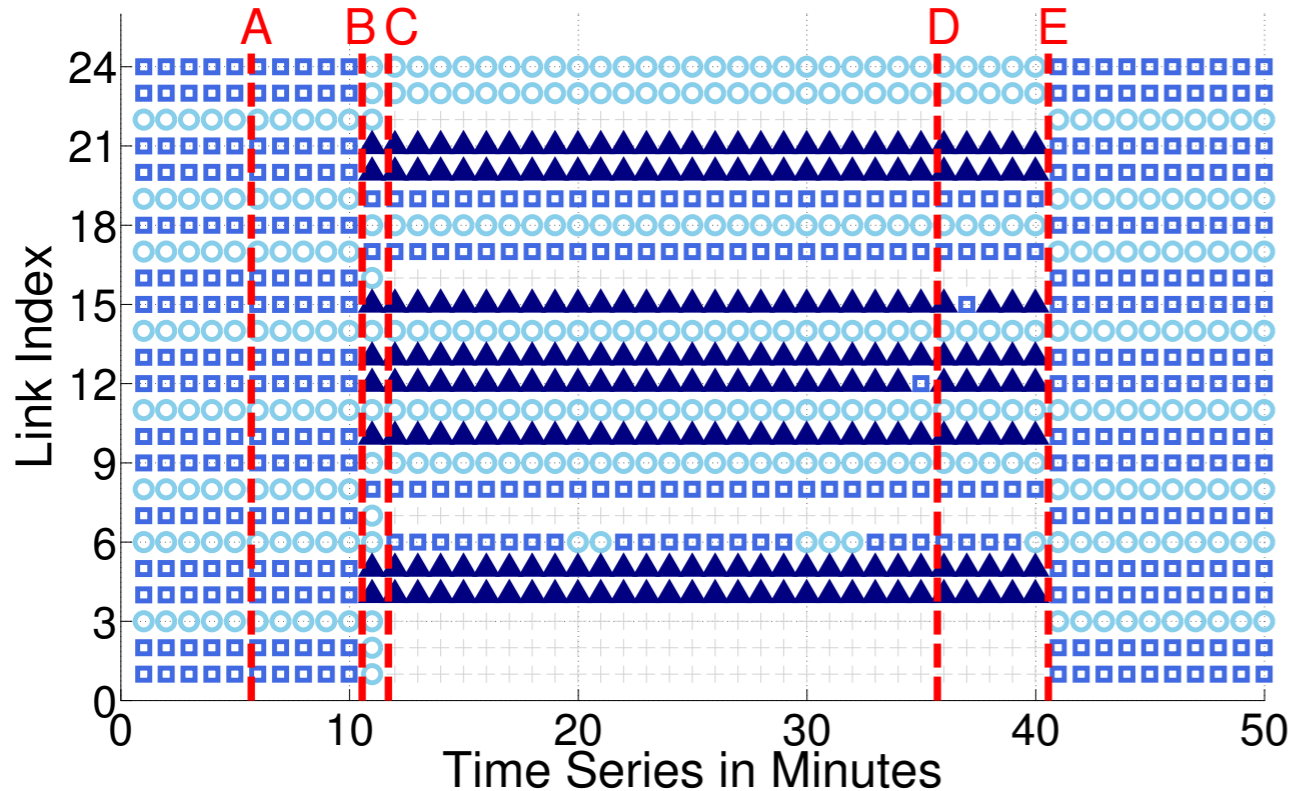BR 6    BR 5

DC 3

licts

ordinates

ade, TE

-level lock

de high-level

+ Empty (0%) ○ Low (1~40%) □ Medium (40%~80%) ▲ High (80%~100%)



B

- TE fails to hold low-level lock, moving traffic away from BR₁

BR = Border Router

DC 2
BR 3    BR 4
BR 1
DC 1
BR 2    BR 8
DC 4
BR 7
BR 6    BR 5
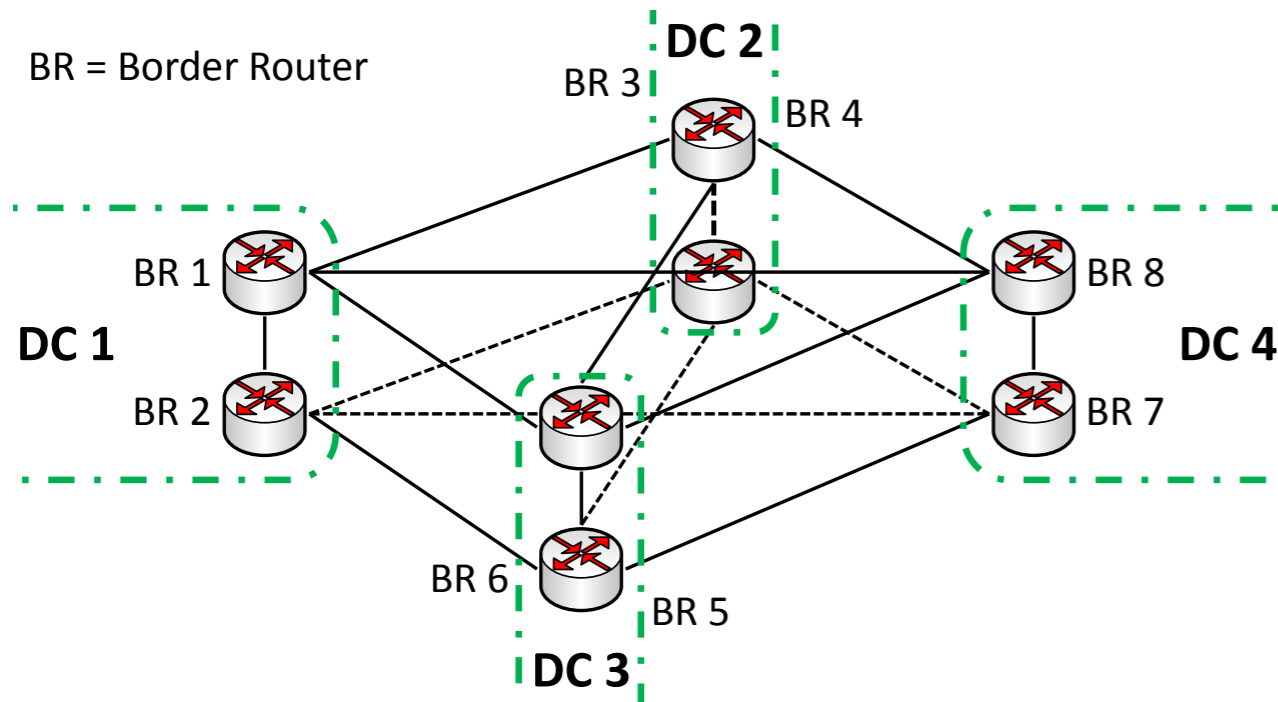DC 3

licts

ordinates

ade, TE

-level lock

de high-level



+ Empty (0%)  ○ Low (1~40%)  □ Medium (40%~80%)  ▲ High (80%~100%)

C,D

- C upgrading BR₁ in progress
- D upgrading done at BR₁, releasing high-level lock

BR = Border Router

DC 2

BR 3          BR 4

BR 1

DC 1

BR 8

DC 4

BR 2                                    BR 7

BR 6          BR 5

DC 3

ordinates

ade, TE

-level lock

de high-level

+ Empty (0%)  ○ Low (1~40%)  □ Medium (40%~80%)  ▲ High (80%~100%)



Link Index

Time Series in Minutes

E

- TE grabs low-level lock, in operation

BR = Border Router

DC 2

BR 3    BR 4

BR 1

DC 1

BR 8

DC 4

BR 2    BR 7

BR 6    BR 5

DC 3

licts

ordinates

ade, TE

-level lock

de high-level

+ Empty (0%) ○ Low (1~40%) □ Medium (40%~80%) ▲ High (80%~100%)



question
- what next?

# statesman performance

evaluating latency

- application: (<10ms) negligible
- checker: seconds
- updater: (>50%) dominating