GETTING THE WORD OUT:

# MEMBERSHIP, DISSEMINATION & POPULATION PROTOCOLS

SEAN CRIBBS
SENIOR PRINCIPAL ENGINEER
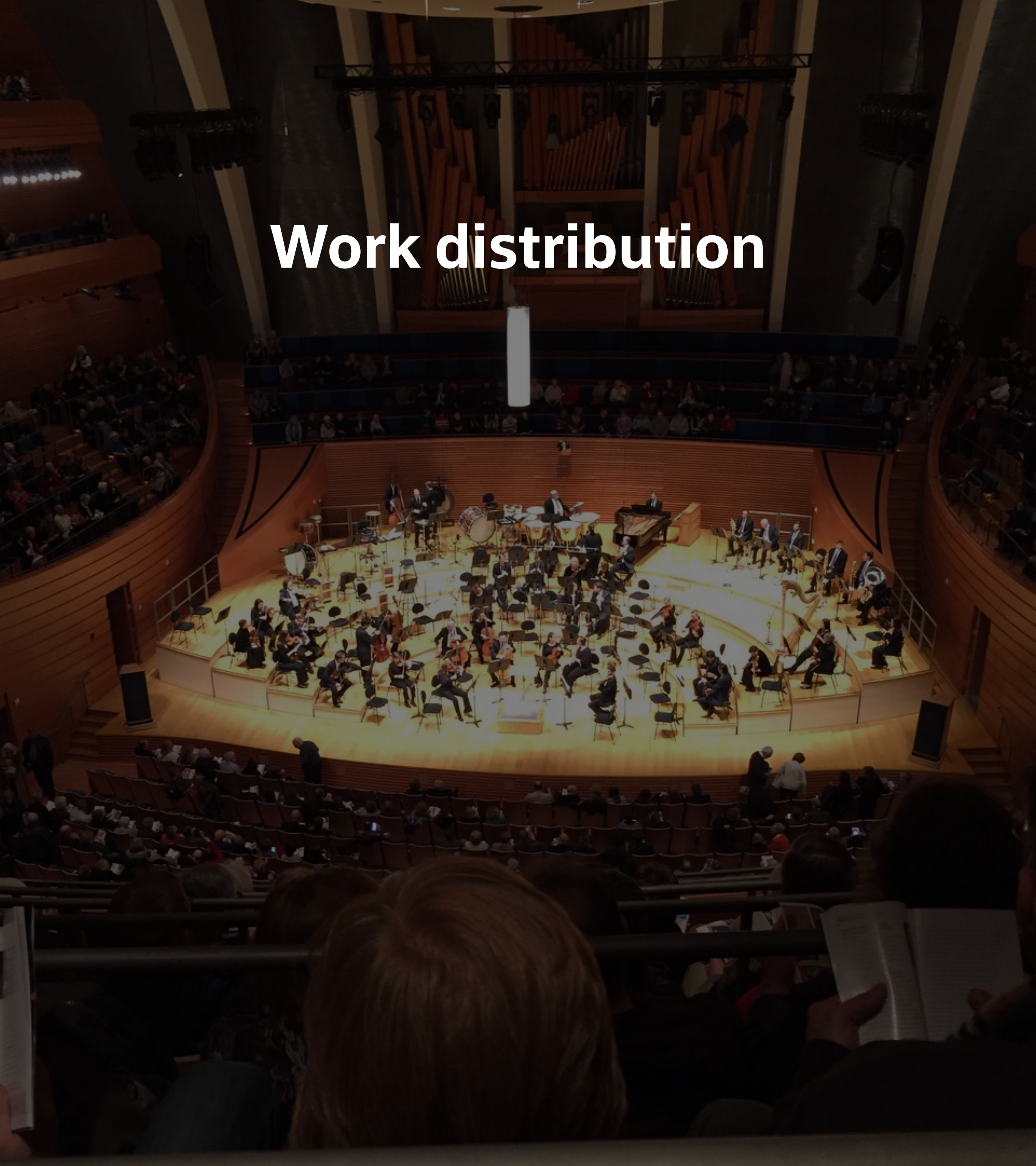
COMCAST

All photos are my own unless attributed

# WHY BUILD PEER-TO-PEER SYSTEMS?

**Work distribution**

**Work distribution**

parallelism

**Work distribution**

parallelism

concurrency

**Work distribution**

parallelism
concurrency
independence

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

detection

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

detection

recovery

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

detection
recovery
redundancy

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

detection
recovery
redundancy

**Locality**

**Work distribution**

parallelism
concurrency
independence

**Fault Tolerance**

detection
recovery
redundancy

**Locality**

regional presence

**Work distribution**

parallelism

concurrency

independence

**Fault Tolerance**

detection

recovery

redundancy

**Locality**

regional presence

code-to-data
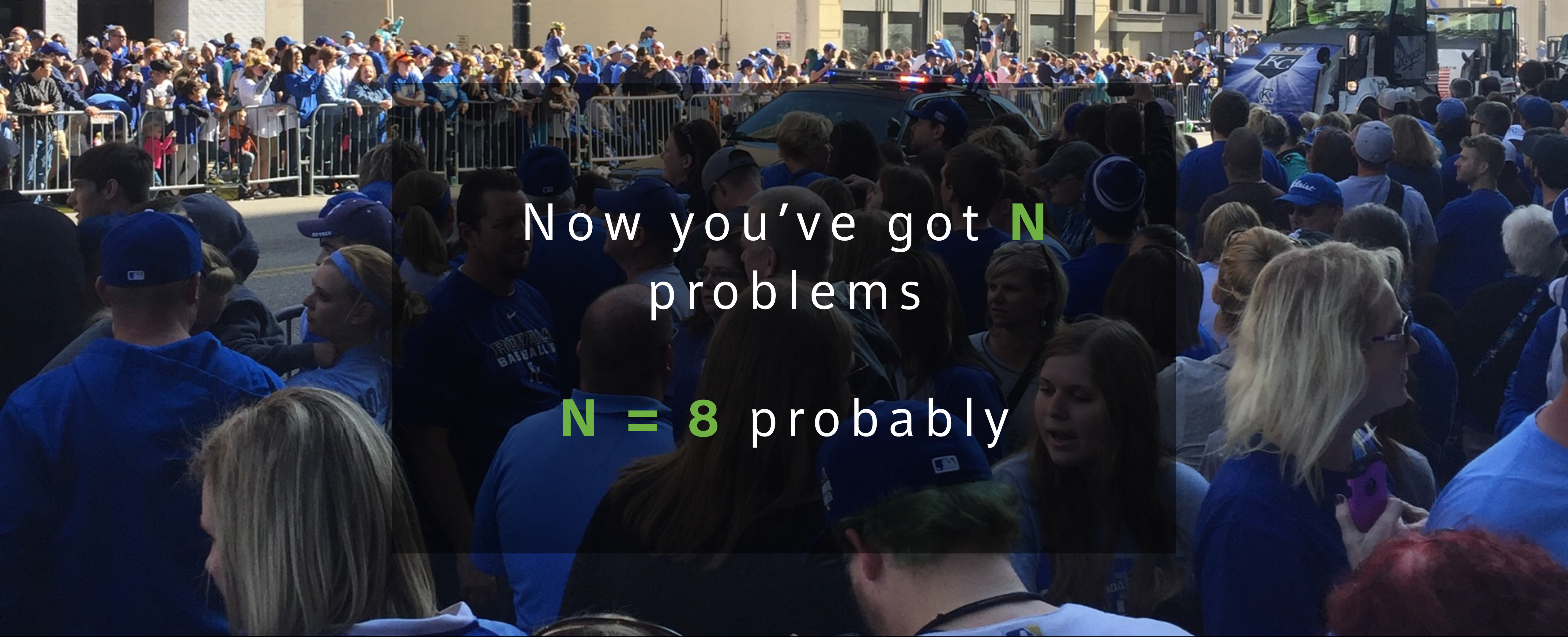
# WHY **NOT** PEER-TO-PEER SYSTEMS?

# WHY **NOT** PEER-TO-PEER SYSTEMS?

Now you've got **N** problems

# WHY **NOT** PEER-TO-PEER SYSTEMS?

Now you've got **N** problems

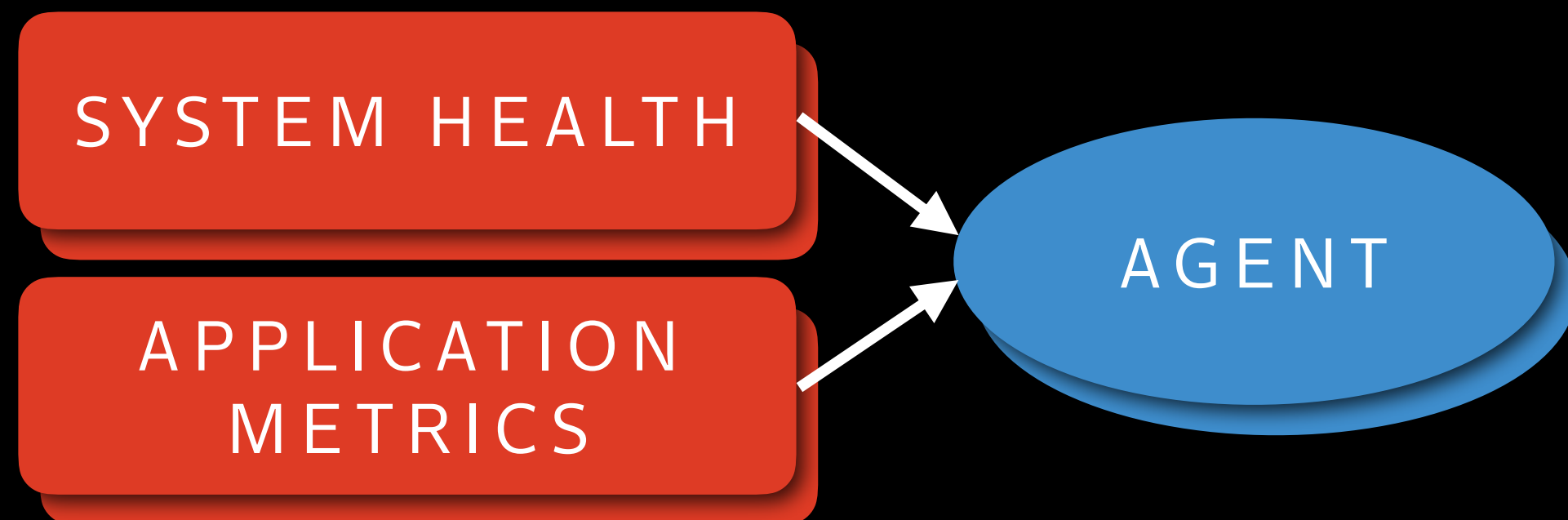**N = 8** probably

# WHAT ARE WE BUILDING?

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

EXTERNAL SERVICES

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

EXTERNAL SERVICES

COMCAST

ARGUS OPERATIONAL VISIBILITY PROJECT

"MULTI-TENANT"

SYSTEM HEALTH

APPLICATION METRICS

AGENT

EXTERNAL SERVICES

COMCAST

ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

EXTERNAL SERVICES

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

EXTERNAL SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

SYSTEM HEALTH

APPLICATION METRICS

AGENT

EXTERNAL SERVICES

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

EXTERNAL SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT



SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

EXTERNAL SERVICES

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

➡ Work distribution

EXTERNAL SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

➡ Work distribution

➡ Fault-tolerance

EXTERNAL SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

➡ Work distribution

➡ Fault-tolerance

➡ Locality

EXTERNAL SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

SYSTEM HEALTH

APPLICATION METRICS

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

EXTERNAL SERVICES

➡ Work distribution

➡ Fault-tolerance

➡ Locality

✓ Peer to Peer!

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

AGENT

EXTERNAL SERVICES

How do cluster nodes find each other?

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

AGENT

EXTERNAL SERVICES

How do cluster nodes find each other?

Distribute code and configuration?

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

AGENT

EXTERNAL SERVICES

How do cluster nodes find each other?

Distribute code and configuration?

Know what happened when?

Emoji provided free by Emoji One

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

Where do agents send data?

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

EXTERNAL SERVICES

How do cluster nodes find each other?

Distribute code and configuration?

Know what happened when?

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

Where do agents
send data?

How do cluster nodes
find each other?

AGENT

"MULTI-TENANT"

"MULTI-REGION"

"HIGHLY-AVAILABLE"

"REAL-TIME"

"STREAMING"

"PLATFORM"

Distribute code and
configuration?

How to get fault-tolerance
without spam?

Know what happened
when?

EXTERNAL
SERVICES

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

➡ Cluster membership and discovery

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

➡ Cluster membership and discovery

➡ Code and configuration dissemination

COMCAST

# ARGUS OPERATIONAL VISIBILITY PROJECT

➡ Cluster membership and discovery

➡ Code and configuration dissemination

➡ Relative and convergent time

COMCAST

# MEMBERSHIP PROTOCOLS

JUST RUB SOME CONSENSUS ON IT

# WHY NOT ZOOKEEPER/CONSUL/ETCD?

COMCAST

➡ Connectedness

COMCAST

# MEMBERSHIP: DESIRABLE PROPERTIES

➡ Connectedness

➡ Balance

# MEMBERSHIP: DESIRABLE PROPERTIES

➡ Connectedness

➡ Balance

➡ Short path-length

COMCAST

# MEMBERSHIP: DESIRABLE PROPERTIES

➡ Connectedness          ➡ Low clustering

➡ Balance

➡ Short path-length

COMCAST

# MEMBERSHIP: DESIRABLE PROPERTIES

➡ Connectedness

➡ Balance

➡ Short path-length

➡ Low clustering

➡ Scalability
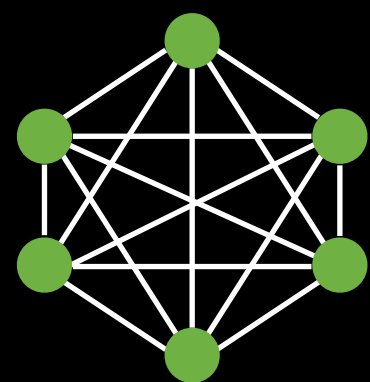
COMCAST

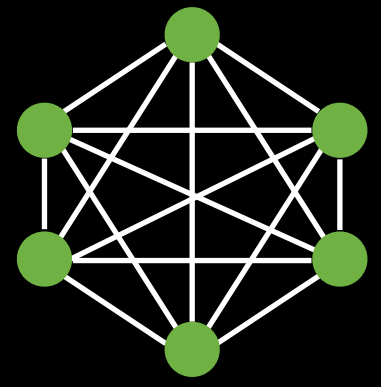# MEMBERSHIP: "VIEW" FLAVORS

# Partial

# SWIM - 2002

**SWIM:** *Scalable Weakly-consistent Infection-style Process Group Membership* **Protocol**

Abhinandan Das, Indranil Gupta, Ashish Motivala[*]
Dept. of Computer Science, Cornell University
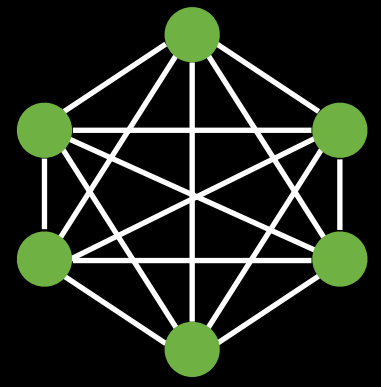Ithaca NY 14853 USA
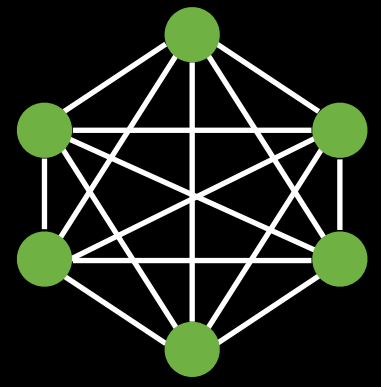{asdas,gupta,ashish}@cs.cornell.edu

COMCAST

# SWIM - 2002

COMCAST

# SWIM - 2002

# Heartbeat protocols

COMCAST

# SWIM - 2002

💕 Heartbeat protocols

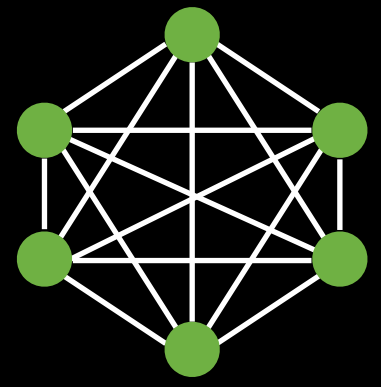◉ Quadratic load

COMCAST

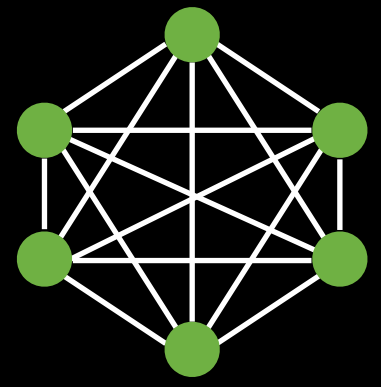# SWIM - 2002

Heartbeat protocols

◉ Quadratic load

◉ Failure detection

COMCAST

# SWIM - 2002

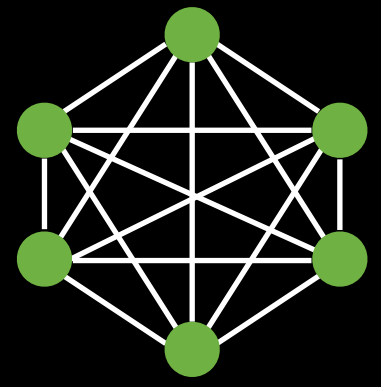Heartbeat protocols

◉ Quadratic load

◉ Failure detection

◉ Response times

COMCAST

# SWIM - 2002

Heartbeat protocols

◉ Quadratic load

◉ Failure detection

◉ Response times

◉ False positives

COMCAST

# SWIM - 2002
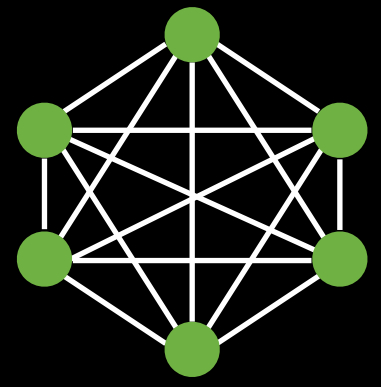
## Heartbeat protocols

- Quadratic load

- Failure detection

  - Response times

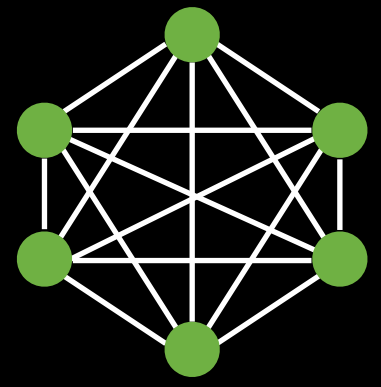  - False positives

## SWIM solutions

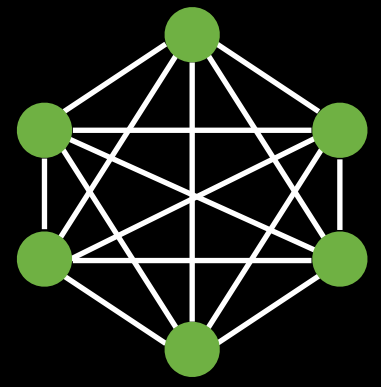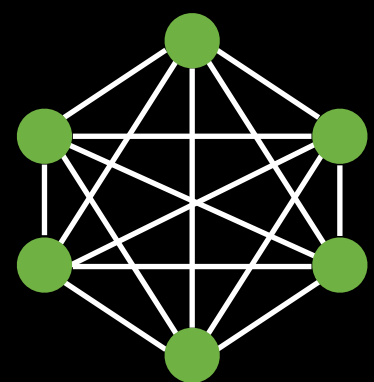COMCAST

# SWIM - 2002

### Heartbeat protocols

- ◉ Quadratic load
- ◉ Failure detection
  - ◉ Response times
  - ◉ False positives

### SWIM solutions

➡ Separate membership and failure detection

COMCAST

# SWIM - 2002

## 💕 Heartbeat protocols

- ◉ Quadratic load

- ◉ Failure detection

  - ◉ Response times

  - ◉ False positives

## 🏊 SWIM solutions

➡ Separate membership and failure detection

➡ Randomized probing

COMCAST

# SWIM - 2002

## ❤️ Heartbeat protocols

- ◉ Quadratic load

- ◉ Failure detection

  - ◉ Response times

  - ◉ False positives

## 🏊 SWIM solutions

➡️ Separate membership and failure detection

➡️ Randomized probing

➡️ Piggyback membership on probes
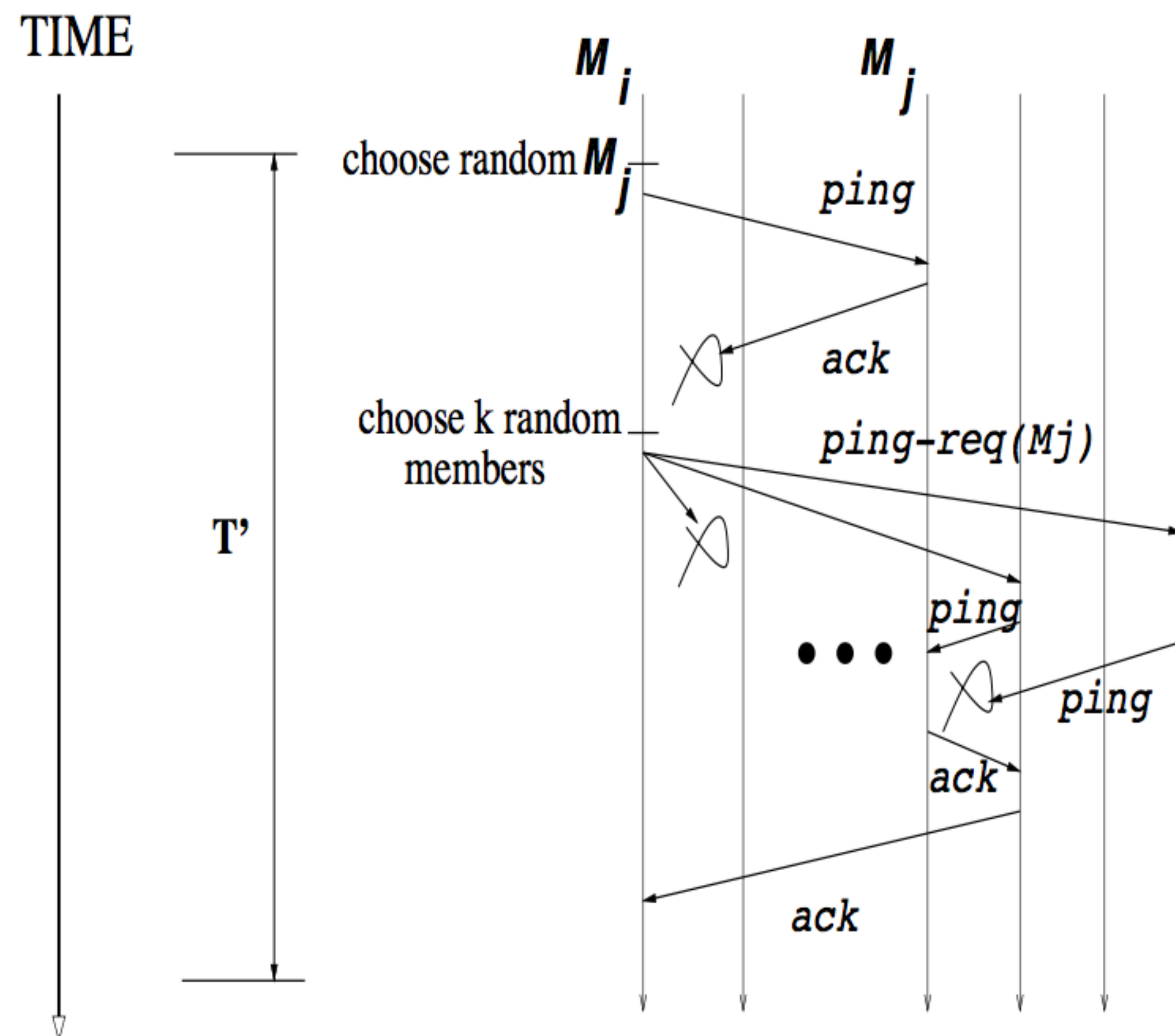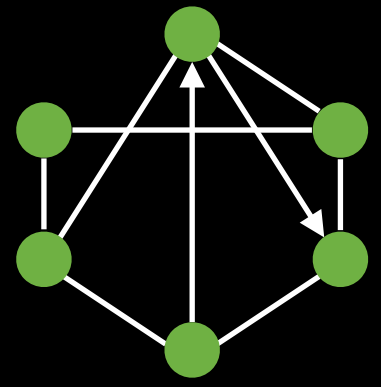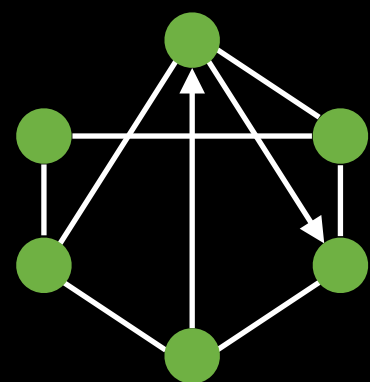
COMCAST

# SWIM - 2002



**Figure 1.** SWIM failure detection: Example protocol period at $M_i$. This shows all the possible messages that a protocol period may initiate. Some message contents excluded for simplicity.
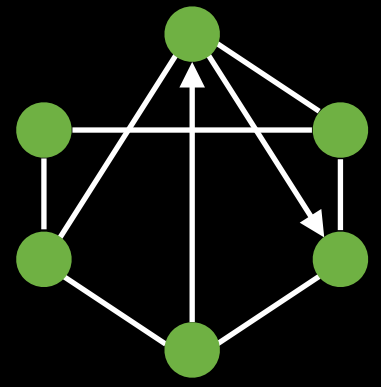
# SCAMP - 2003

Peer-to-Peer Membership Management
for Gossip-Based Protocols

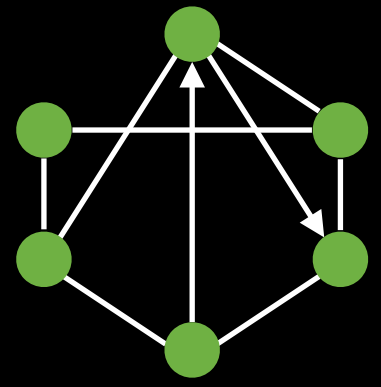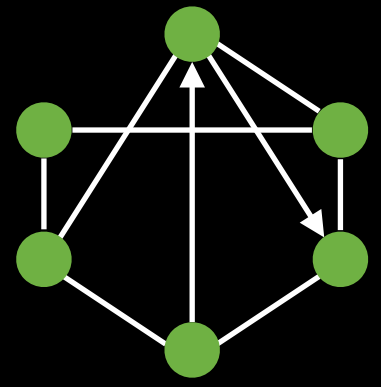Ayalvadi J. Ganesh, Anne-Marie Kermarrec, and Laurent Massoulié

COMCAST

SCAMP - 2003

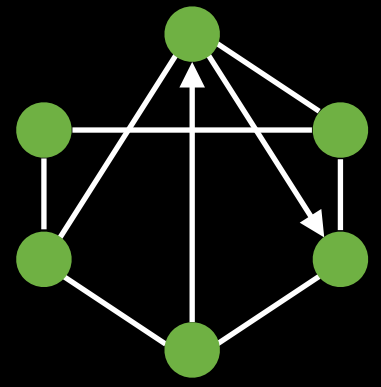COMCAST

# SCAMP - 2003

- Full views limit scalability

# SCAMP - 2003

⦿ Full views limit scalability
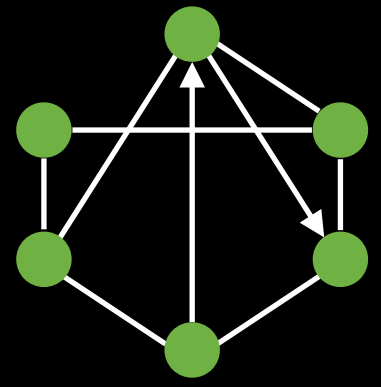
➡ Flexible partial-view size, asymmetric

# SCAMP - 2003

◉ Full views limit scalability

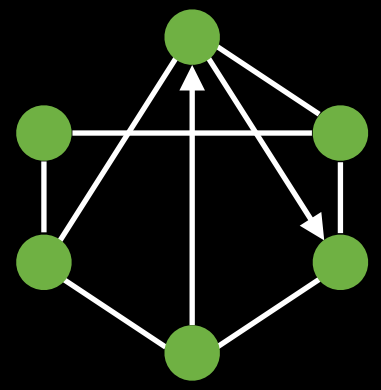➡ Flexible partial-view size, asymmetric

➡ Reactive view management

# SCAMP - 2003

- Full views limit scalability

➡ Flexible partial-view size, asymmetric

➡ Reactive view management

➡ Join ("subscribe") via random walk

COMCAST

# SCAMP - 2003

◉ Full views limit scalability

➡ Flexible partial-view size, asymmetric

➡ Reactive view management

➡ Join ("subscribe") via random walk

➡ Automatic balancing via indirection and leases
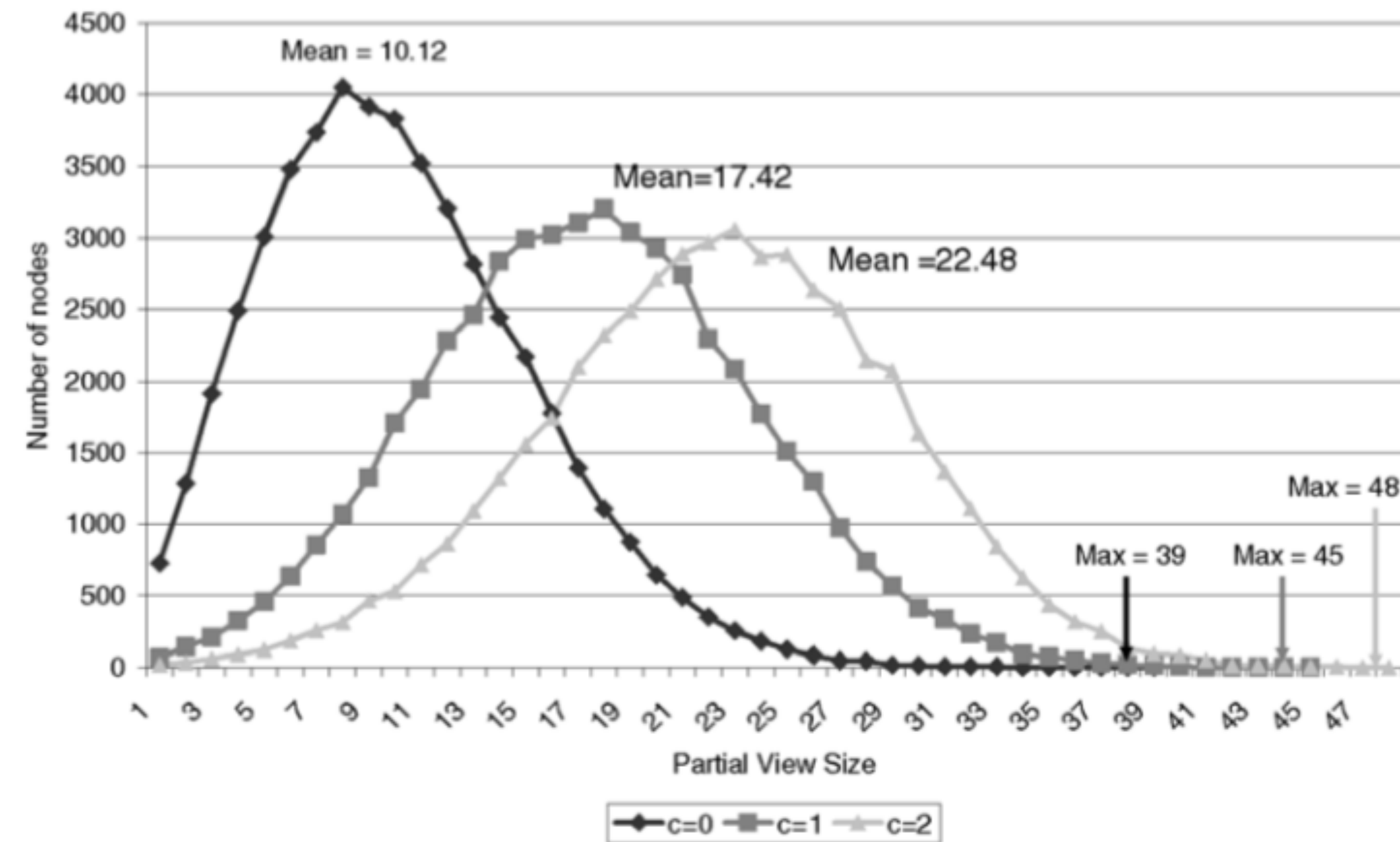
COMCAST

# SCAMP - 2003

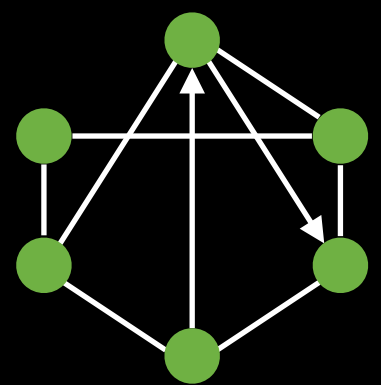

Fig. 5. Impact of c on the partial view size distribution in a 50,000 node group.
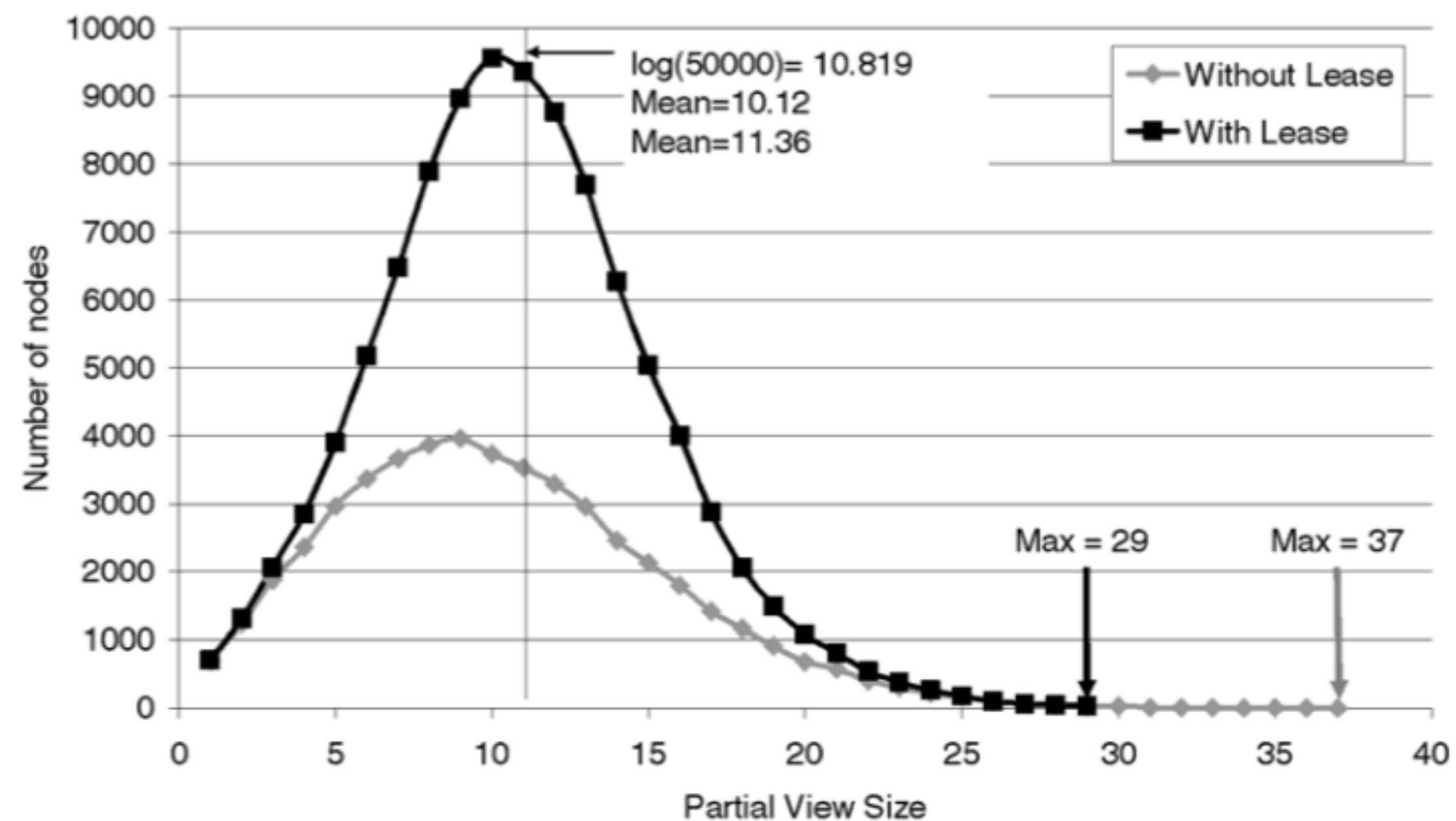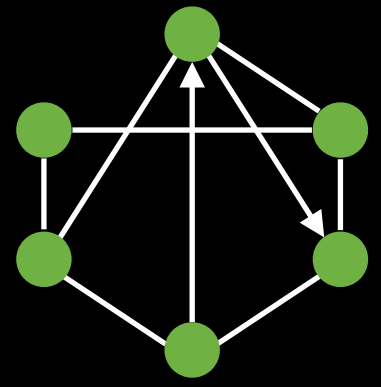
# SCAMP - 2003



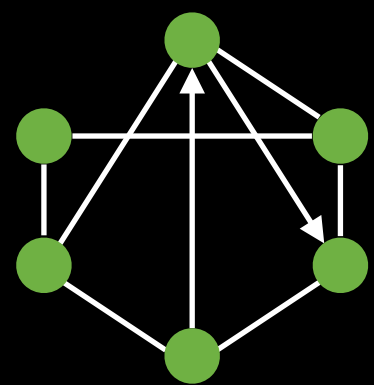Fig. 8. Distribution of the partial views by size in a 50,000 node group with and without the lease mechanism.

# CYCLON - 2005

# CYCLON - 2005
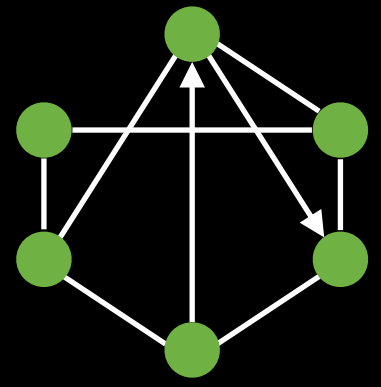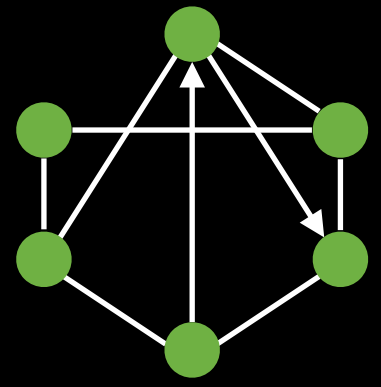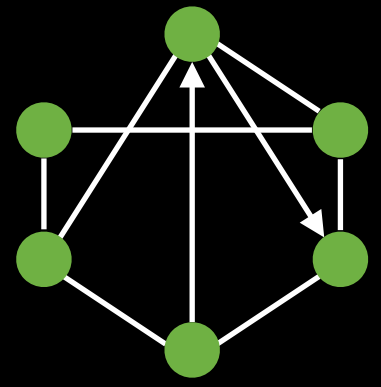
- Random shuffling doesn't create good balance

COMCAST
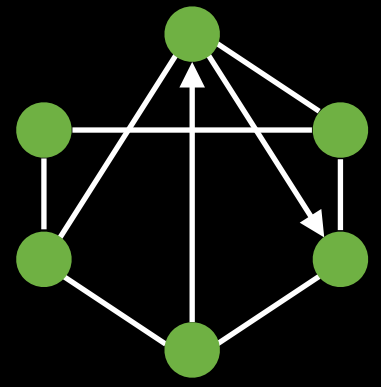
◉ Random shuffling doesn't create good balance
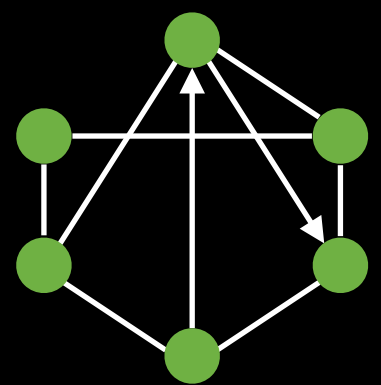
➡ Fixed partial-view size, symmetric

⦿ Random shuffling doesn't create good balance

➡ Fixed partial-view size, symmetric

➡ Cyclic view management

COMCAST

# CYCLON - 2005

⦿ Random shuffling doesn't create good balance

➡ Fixed partial-view size, symmetric

➡ Cyclic view management

➡ Join via random walk

# CYCLON - 2005



$$2 \rightarrow 9 : \{2, 0, 6\}$$
$$2 \leftarrow 9 : \{0, 5, 7\}$$
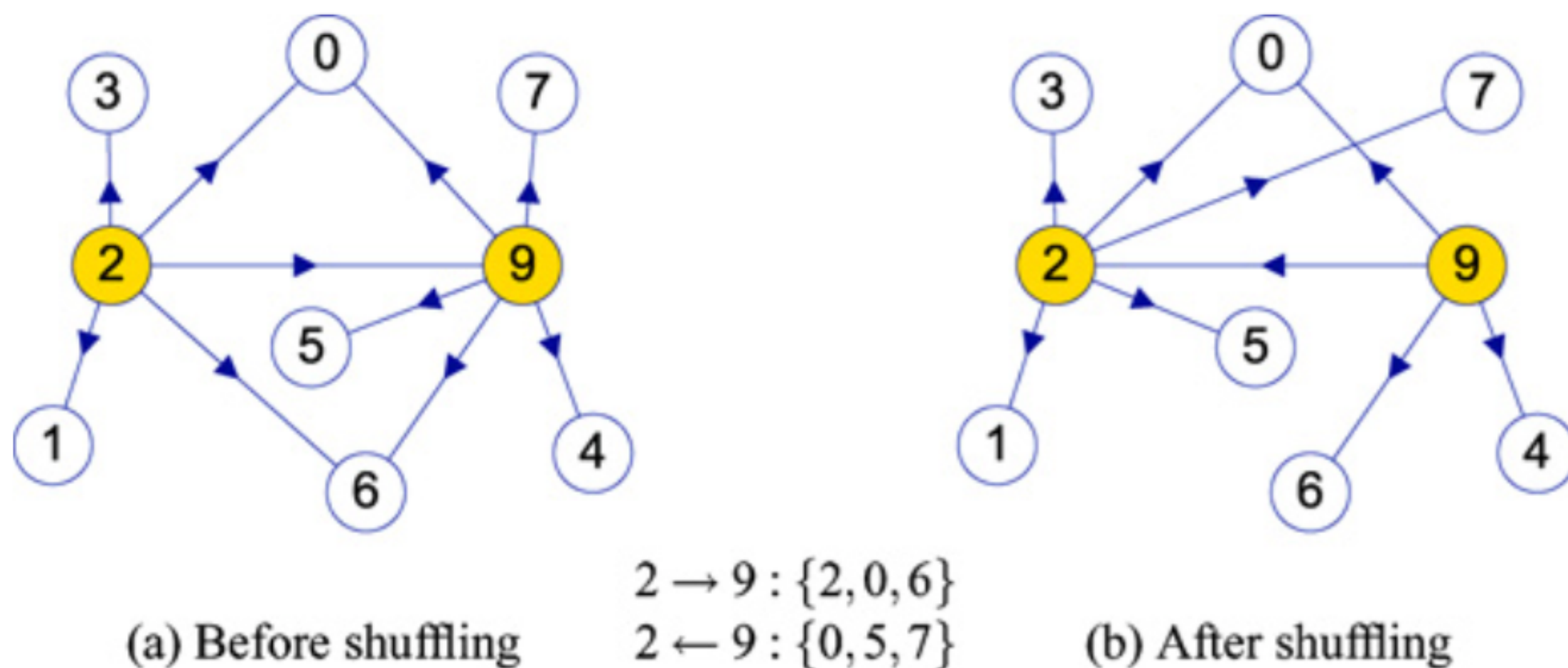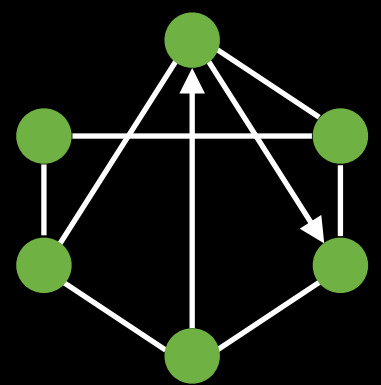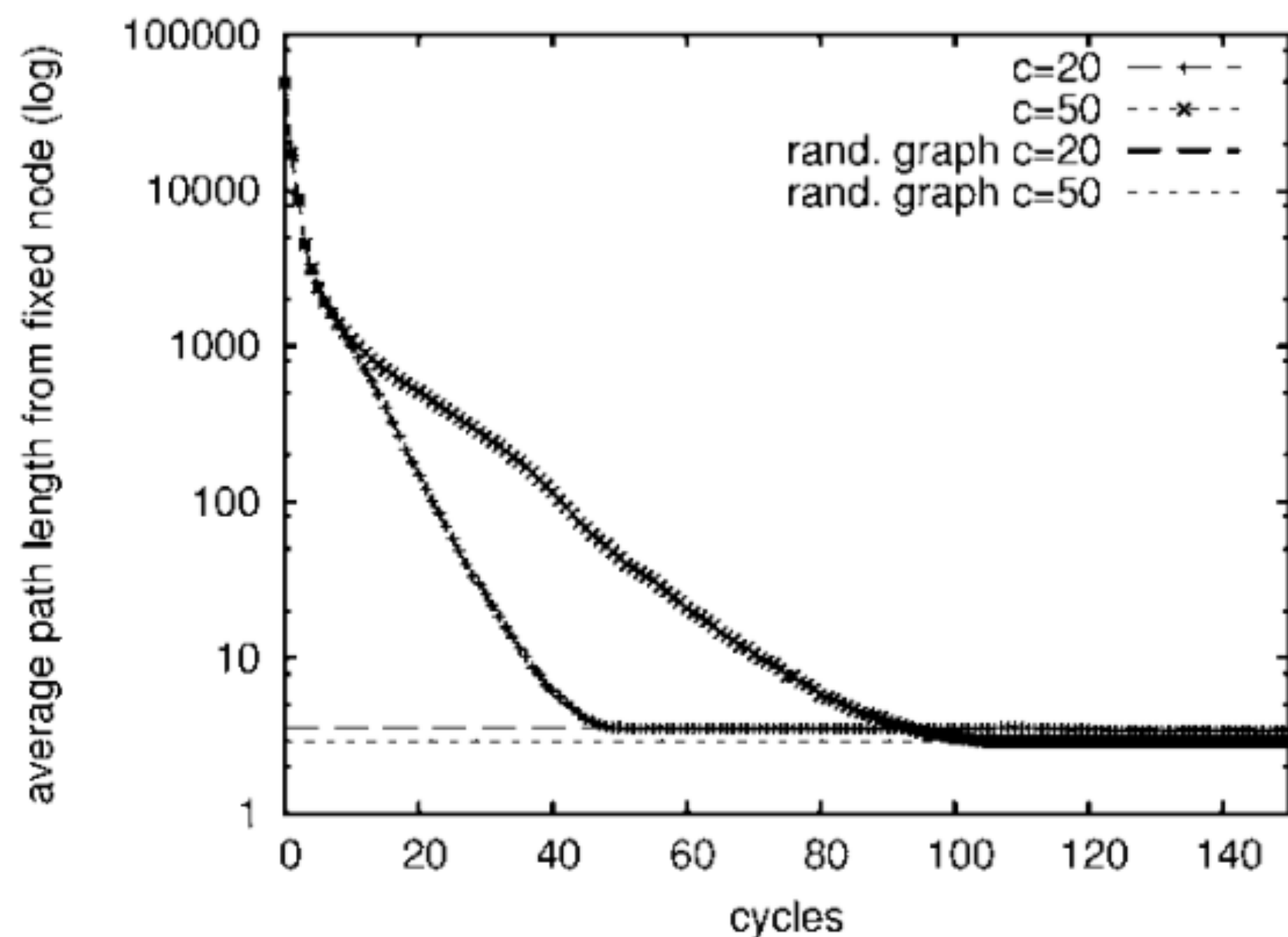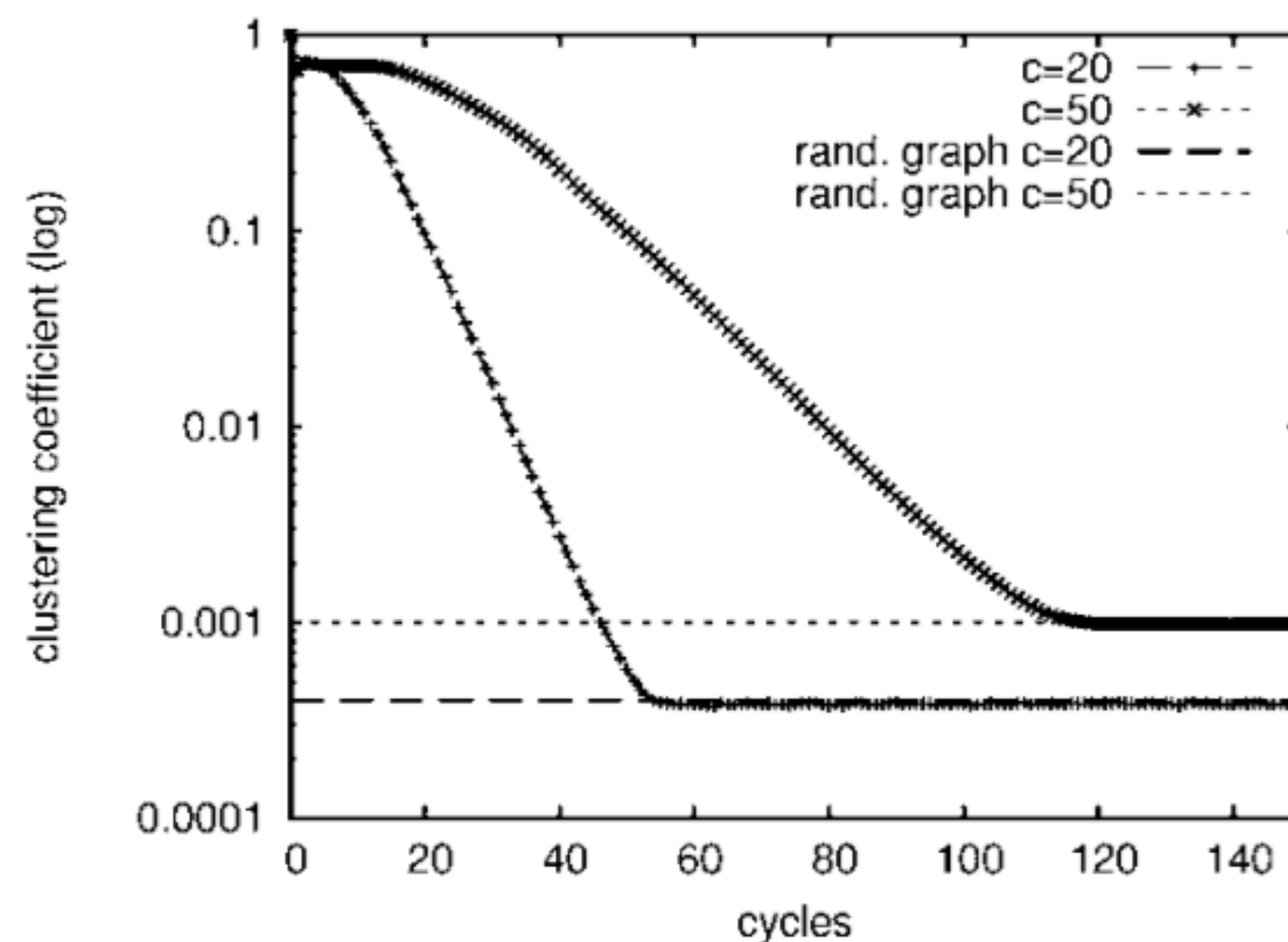
(a) Before shuffling    (b) After shuffling

**Fig. 1.** An example of shuffling between nodes 2 and 9. Note that, among other changes, the link between 2 and 9 reverses direction.
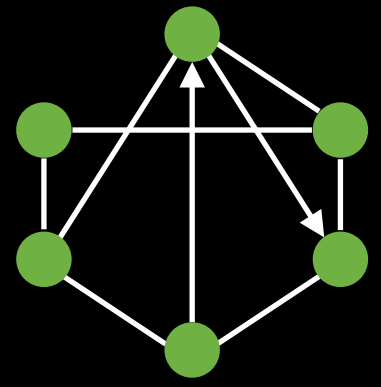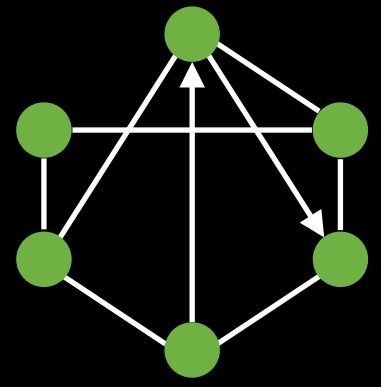
**Fig. 2.** (a) Average shortest path length between two nodes for different cache sizes. (b) Average clustering coefficient taken over all nodes.

# PROBLEMS WITH SCAMP & CYCLON

- No failure detectors

- SCAMP: asymmetric views $\implies$ disconnection

- SCAMP: unbounded view size $\implies$ imbalance
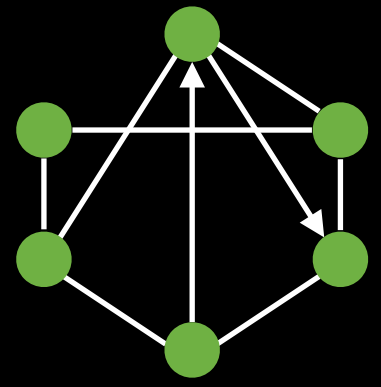
COMCAST

# HYPARVIEW - 2007

HyParView: a membership protocol
for reliable gossip-based broadcast

João Leitão
José Pereira
Luís Rodrigues

COMCAST

HYPARVIEW - 2007
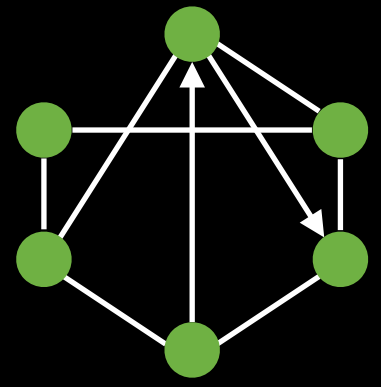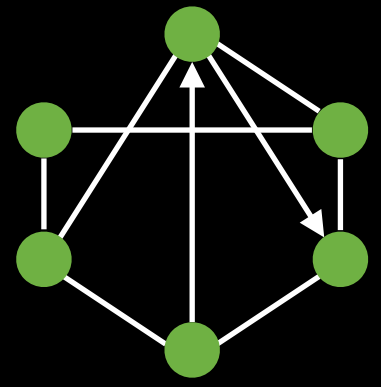
COMCAST

- Fanout is related to reliability

# HYPARVIEW - 2007

◉ Fanout is related to reliability

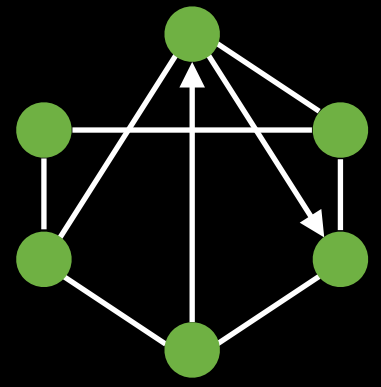◉ High failure rates decrease quality

COMCAST

# HYPARVIEW - 2007

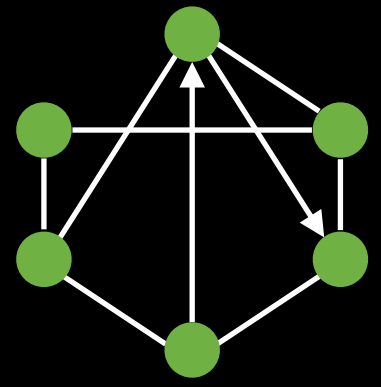- Fanout is related to reliability

- High failure rates decrease quality

# HYPARVIEW - 2007

- Fanout is related to reliability
- High failure rates decrease quality
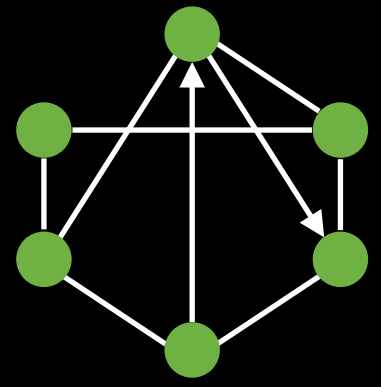
➡ TCP for transport and failure detector

COMCAST

# HYPARVIEW - 2007

- Fanout is related to reliability

- High failure rates decrease quality

➡ TCP for transport and failure detector

➡ Small reactive view ("active")

COMCAST

# HYPARVIEW - 2007

- Fanout is related to reliability

- High failure rates decrease quality

➡ TCP for transport and failure detector

➡ Small reactive view ("active")

➡ Larger cyclic view ("passive")

COMCAST

# HYPARVIEW - 2007

- Fanout is related to reliability

- High failure rates decrease quality

➡ TCP for transport and failure detector

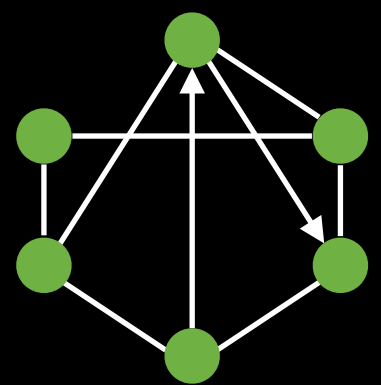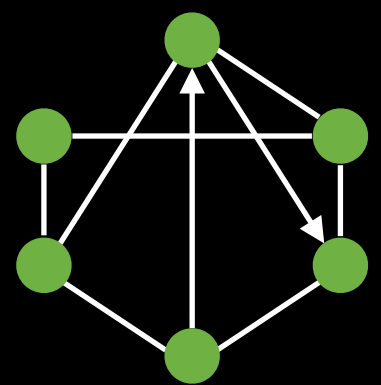➡ Small reactive view ("active")

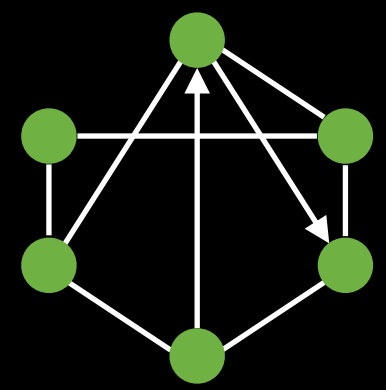➡ Larger cyclic view ("passive")

➡ Join and shuffle via random walk

COMCAST

## Algorithm 1: Membership Operations

**upon init do**
    Send(JOIN, contactNode, myself);

**upon** *Receive*(JOIN, newNode) **do**
    **if** isfull(activeView) **then**
        **trigger** dropRandomElementFromActiveView
    activeView ⟵ activeView ∪ newNode
        **foreach** $n \in$ activeView and $n \neq$ newNode **do**
            Send(FORWARDJOIN, $n$, newNode, ARWL, myself)
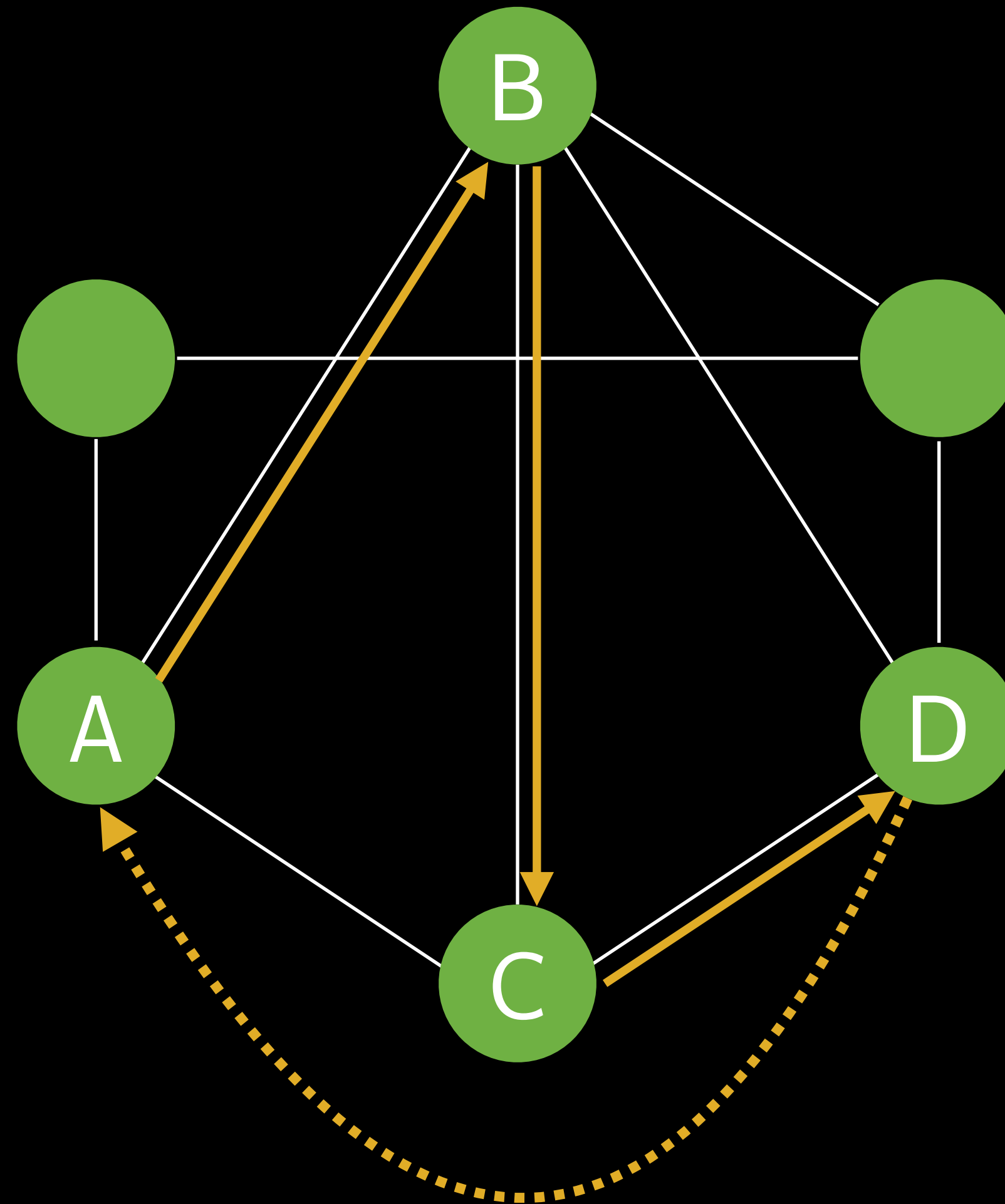
**upon** *Receive*(FORWARDJOIN, newNode, timeToLive, sender) **do**
    **if** timeToLive== 0‖#activeView== 0 **then**
        **trigger** addNodeActiveView(newNode)
    **else**
        **if** timeToLive==PRWL **then**
            **trigger** addNodePassiveView(newNode)
        $n \longleftarrow n \in$ activeView and $n \neq$ sender
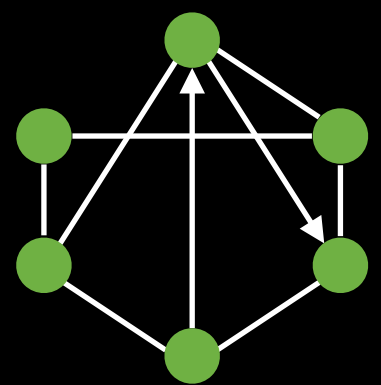        Send(FORWARDJOIN, $n$, newNode, timeToLive-1, myself)

**upon** dropRandomElementFromActiveView **do**
    $n \longleftarrow n \in$ activeView
    Send(DISCONNECT, $n$, myself)
    activeView ⟵ activeView ∖ $\{n\}$
    passiveView ⟵ passiveView ∪ $\{n\}$

HYPARVIEW - 2007

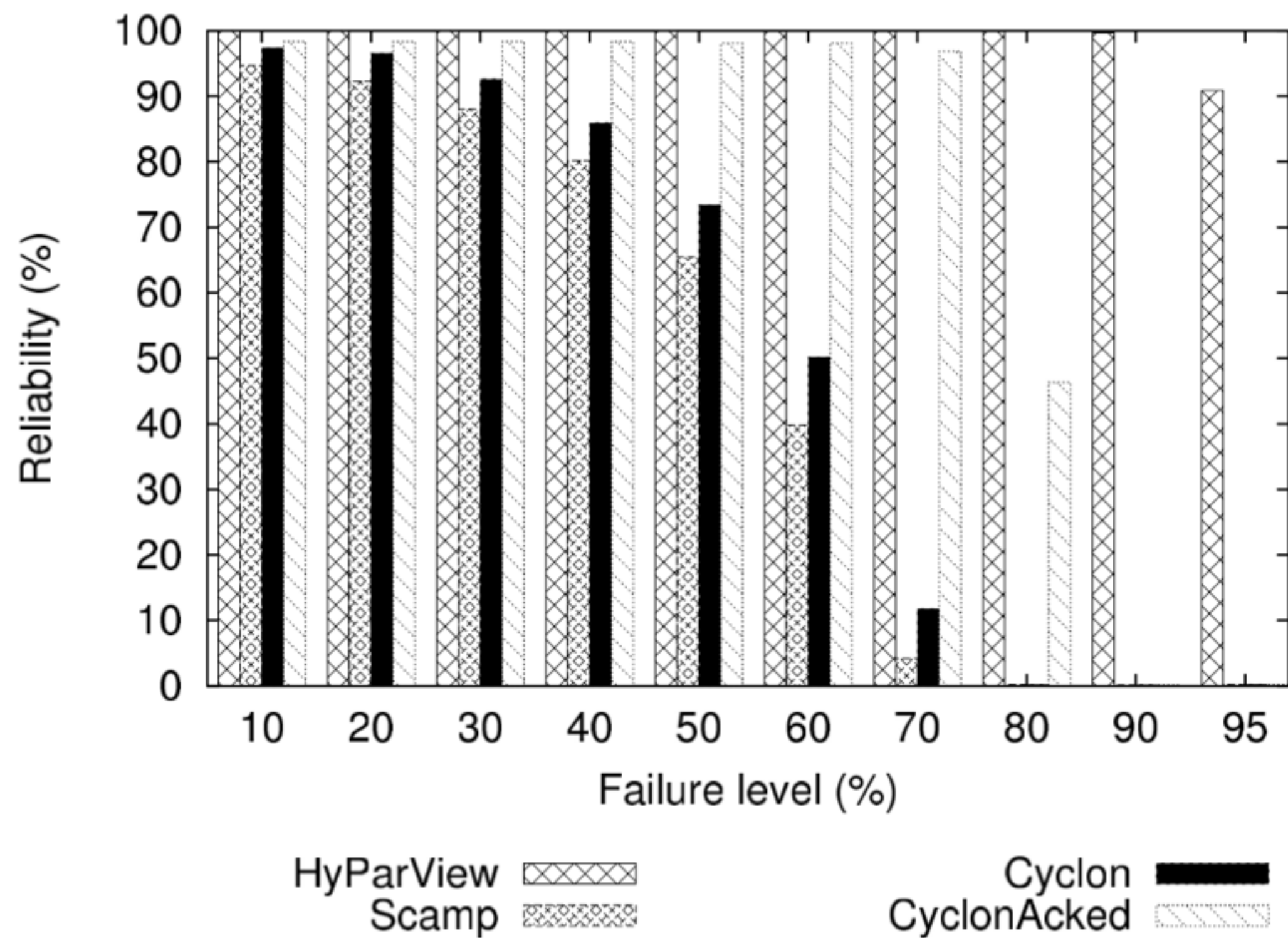Passive view maintenance

COMCAST
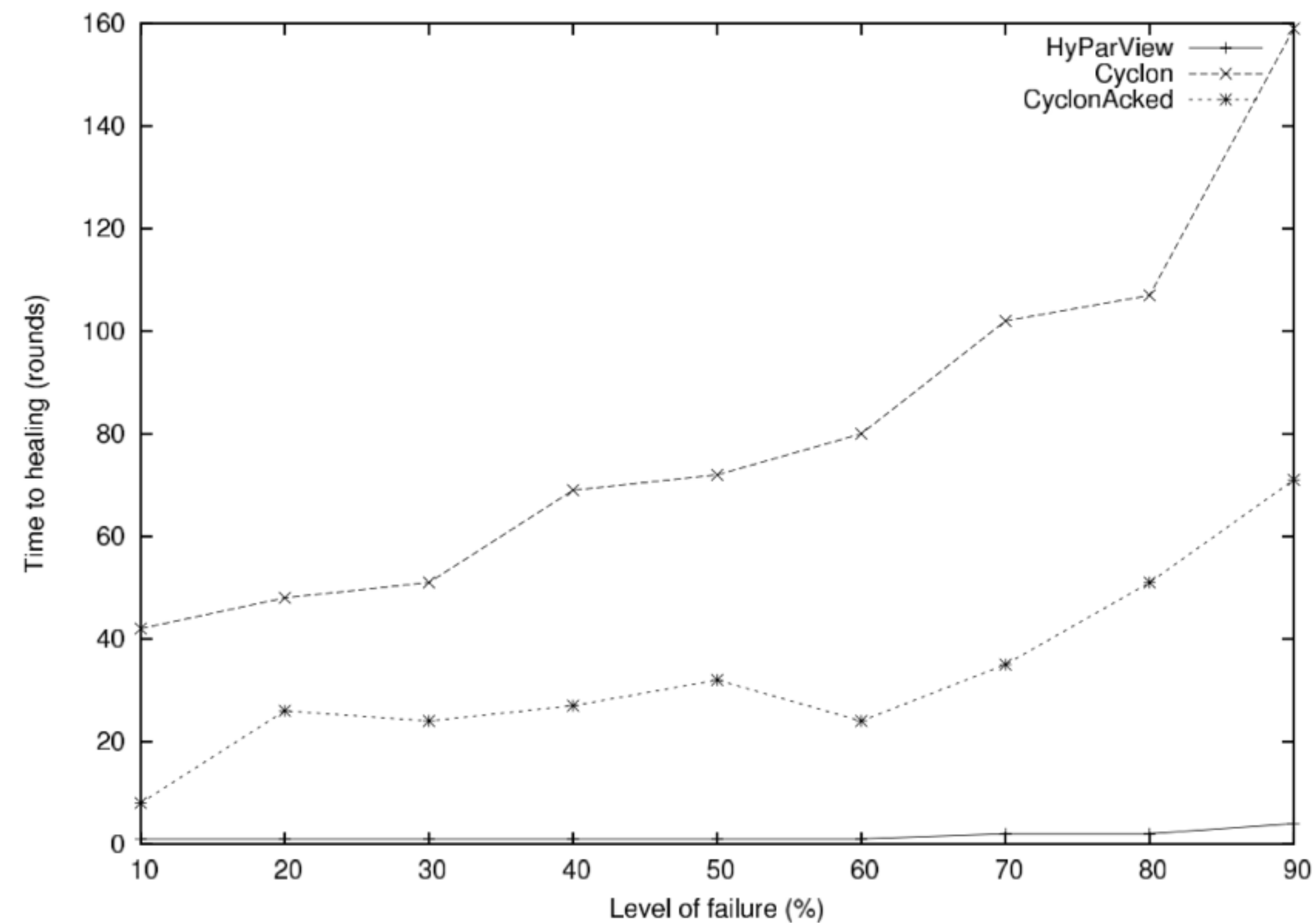
Figure 2: Reliability for 1000 messages



Figure 4: Healing time

# WE CHOSE HYPARVIEW

- **Only** active view maintenance

- Passive view maintains **full membership** (unbounded)

- **Later**: switch to complete passive maintenance

COMCAST

# DISSEMINATION PROTOCOLS

COMCAST

➡ Reliability

COMCAST

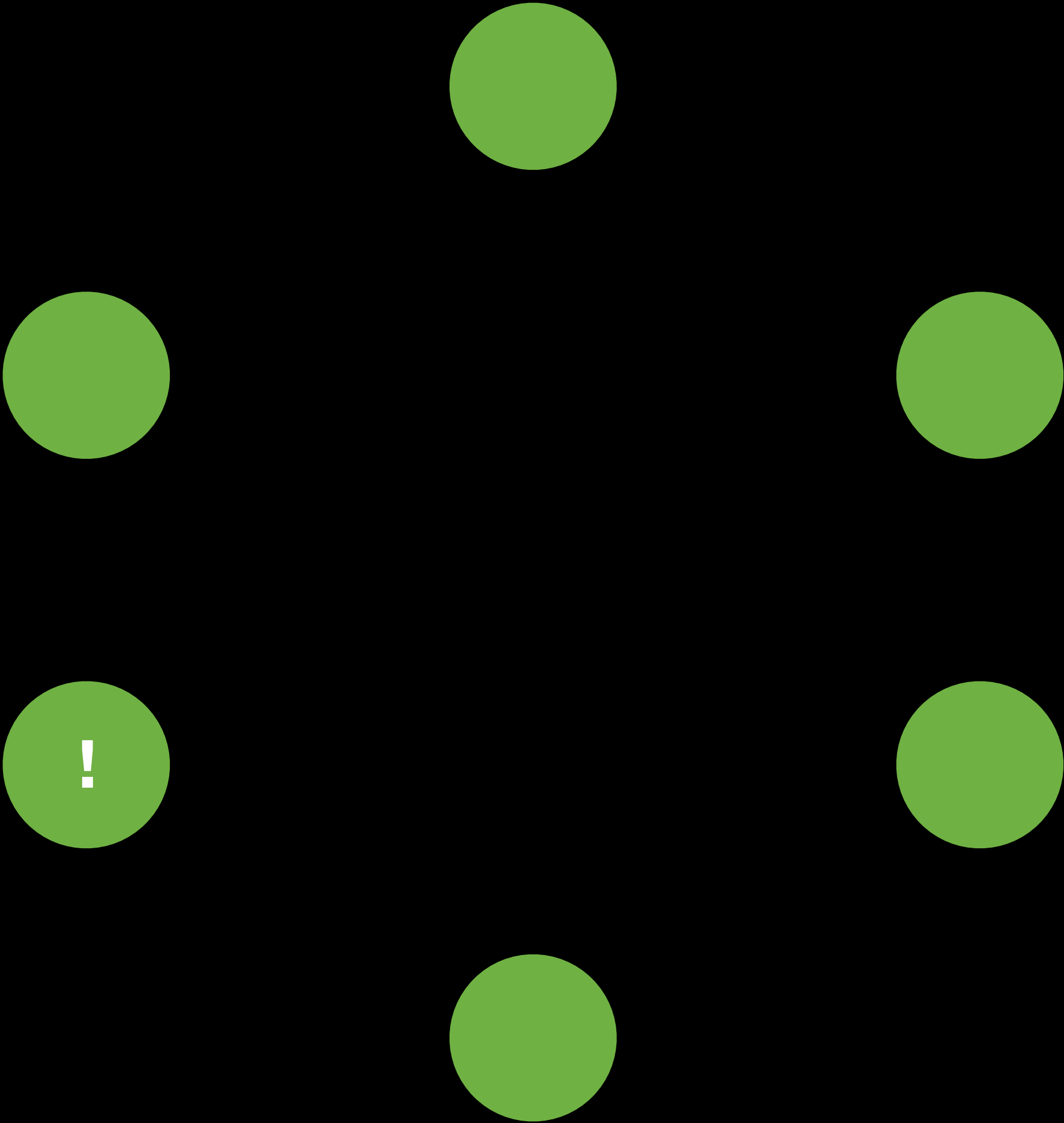# DISSEMINATION: DESIRABLE PROPERTIES

➡ Reliability

➡ Scalability

COMCAST

# DISSEMINATION: DESIRABLE PROPERTIES

➡ Reliability

➡ Scalability

➡ Efficiency

COMCAST

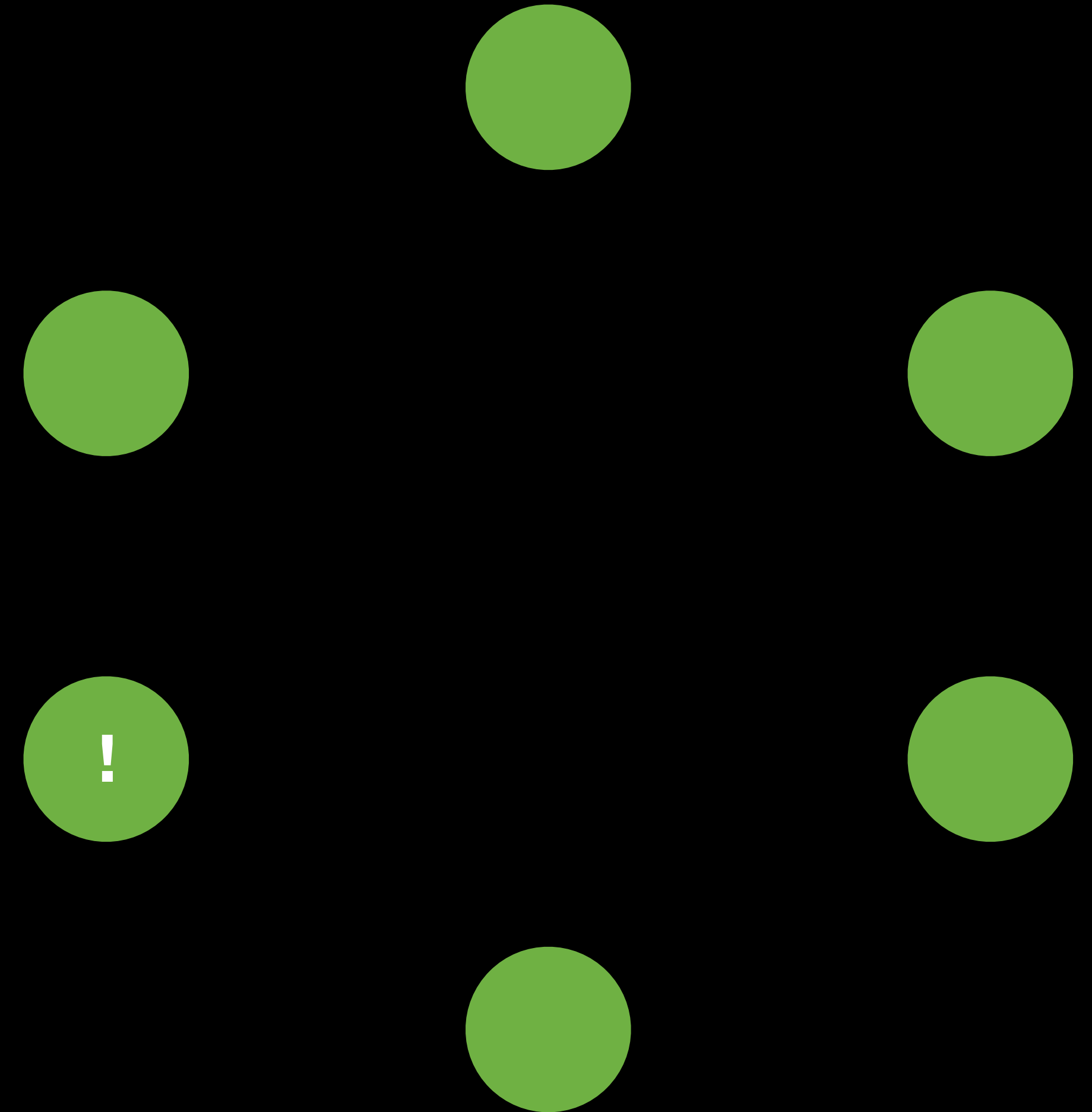EPIDEMIC
BROADCAST (GOSSIP)

COMCAST

# EPIDEMIC BROADCAST (GOSSIP)

➡ Send to random peers

COMCAST

# EPIDEMIC BROADCAST (GOSSIP)

➡ Send to random peers
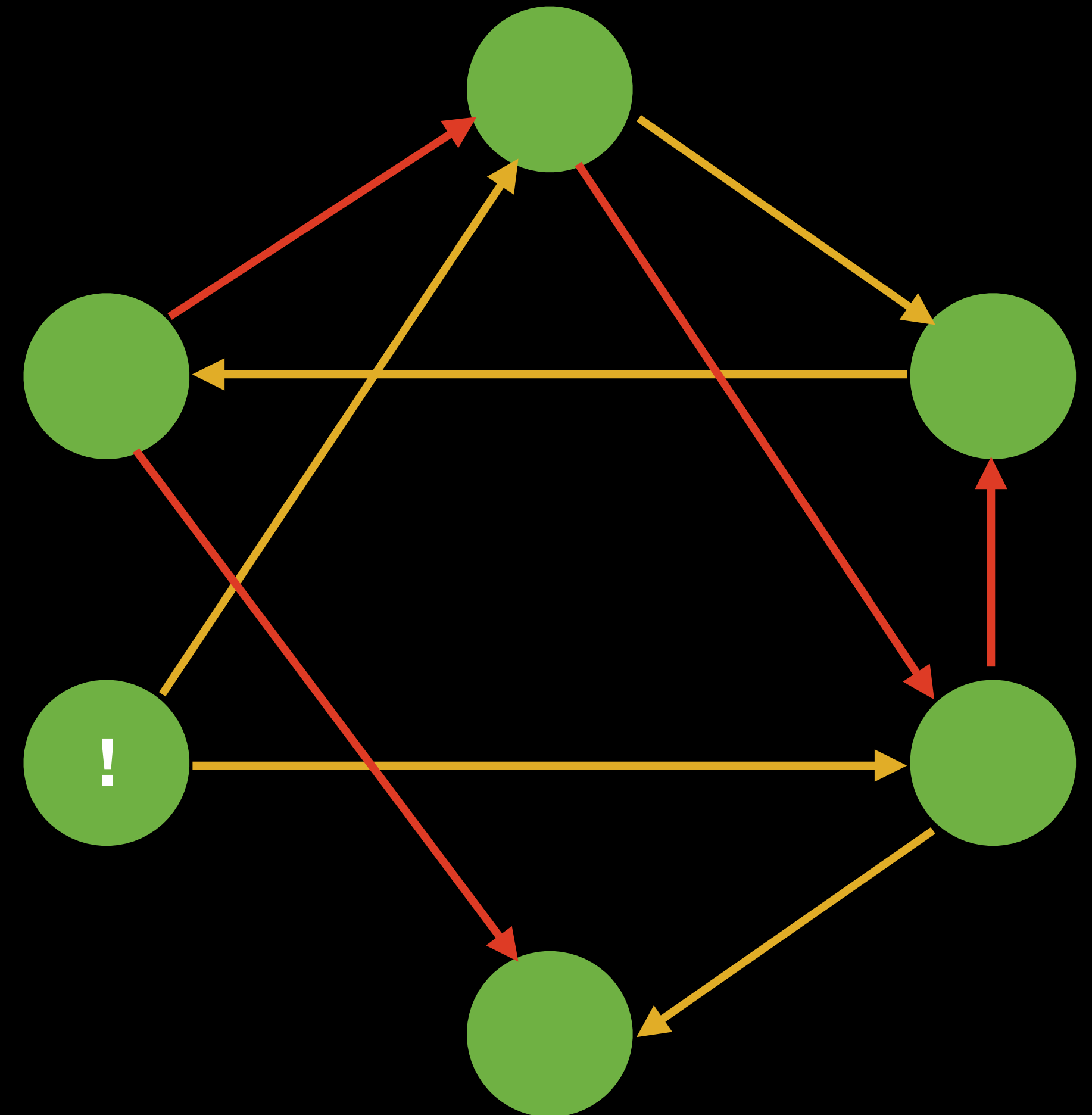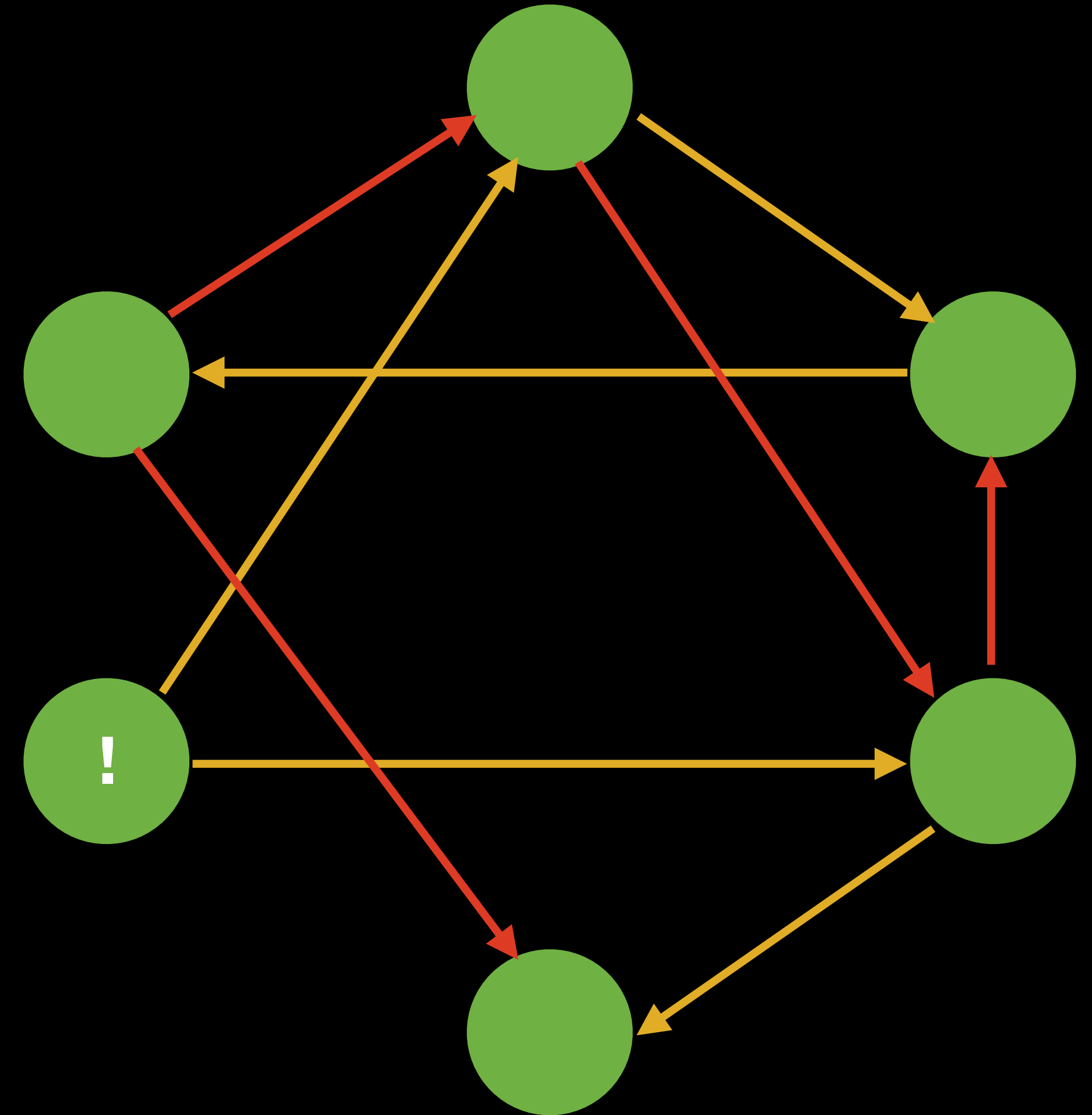
➡ Messages rebroadcast by recipients

!

COMCAST

# EPIDEMIC BROADCAST (GOSSIP)

➡ Send to random peers

➡ Messages rebroadcast by recipients

# EPIDEMIC BROADCAST (GOSSIP)

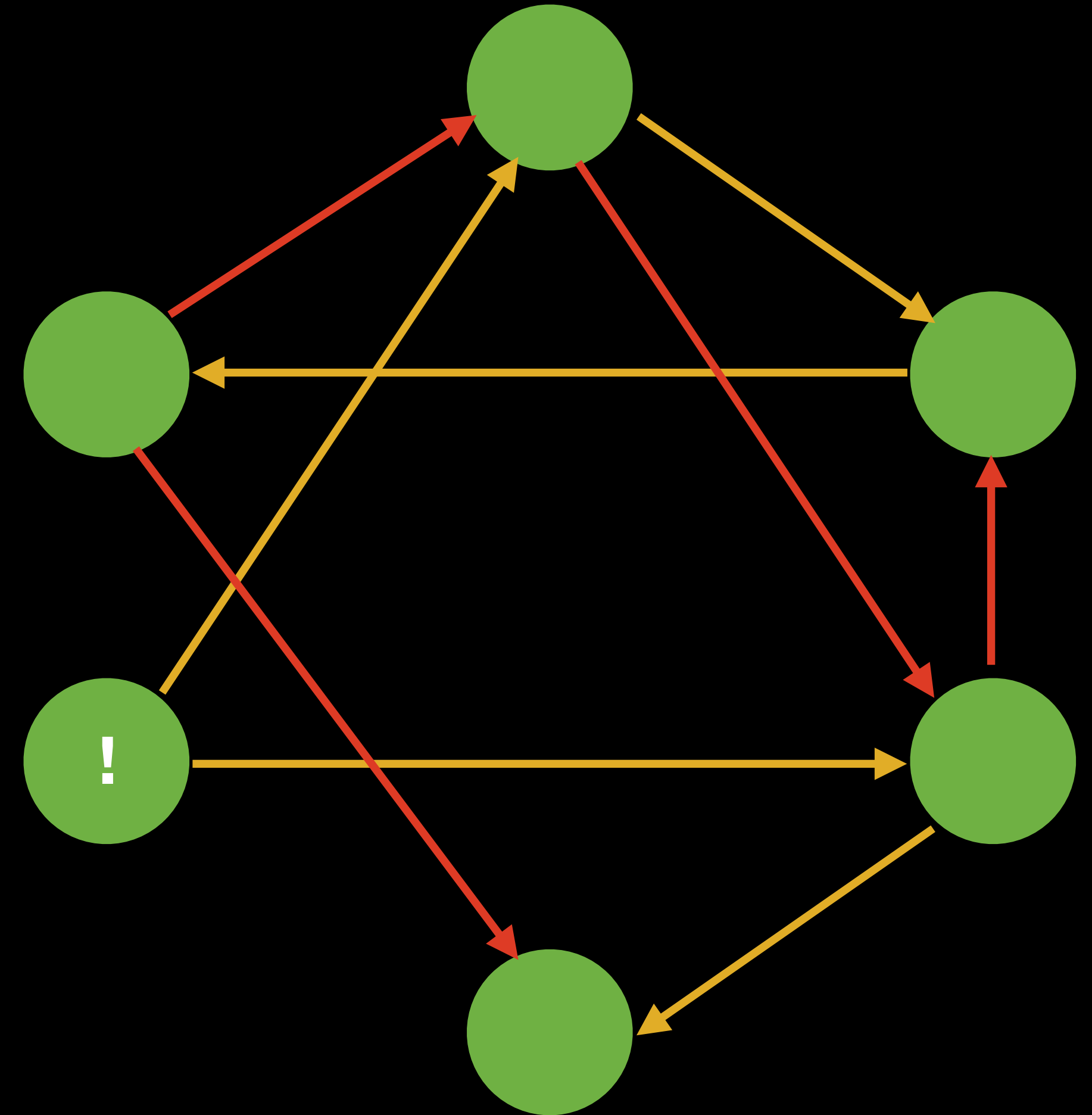➡ Send to random peers

➡ Messages rebroadcast by recipients

◉ High redundancy

# EPIDEMIC BROADCAST (GOSSIP)

➡ Send to random peers

➡ Messages rebroadcast by recipients

◉ High redundancy

◉ Low scalability

WITHOUT REDUCING DELIVERY GUARANTEES, WE NEED

# INCREASED EFFICIENCY

# PLUMTREE - 2009 CONSTRUCTION
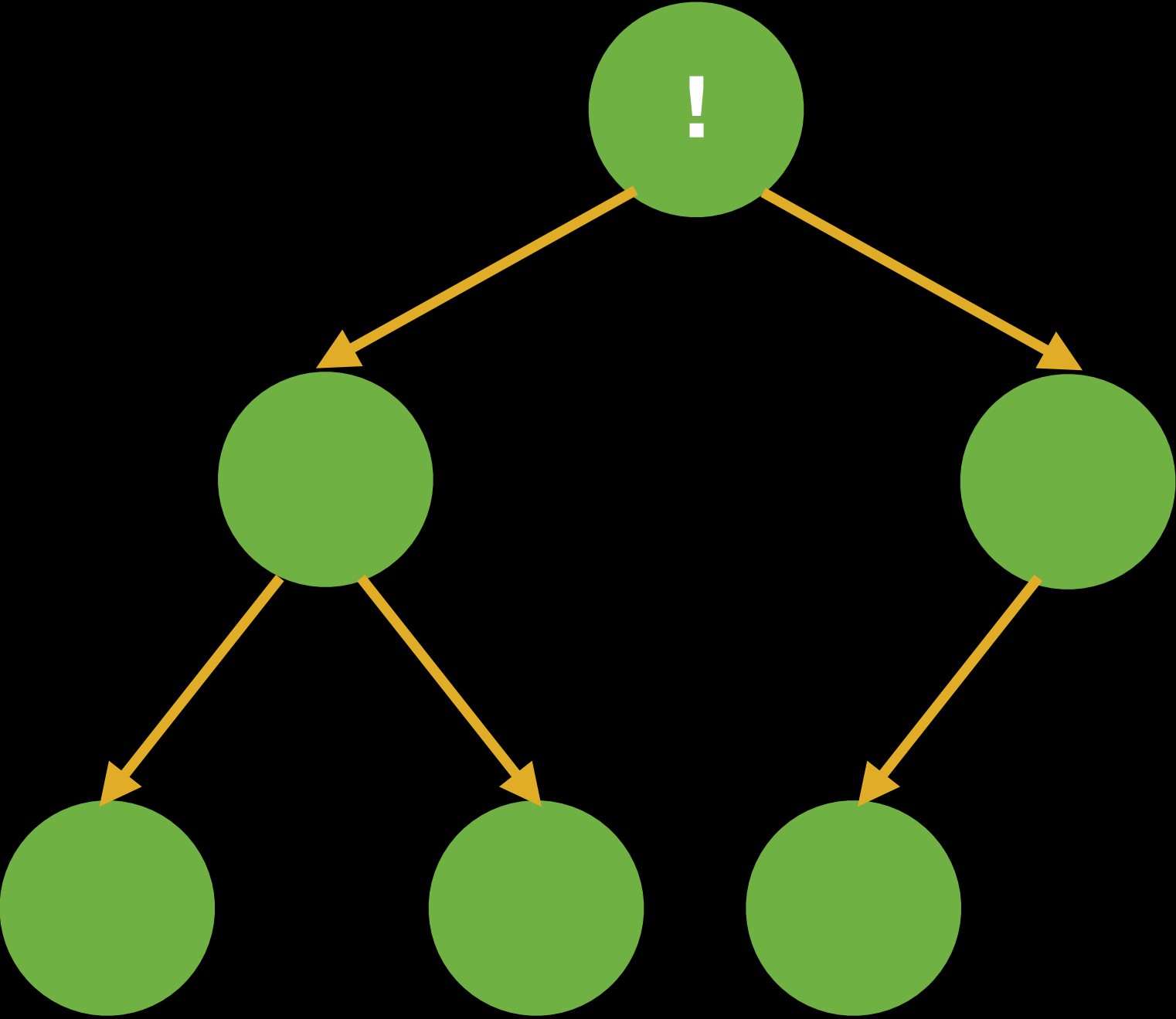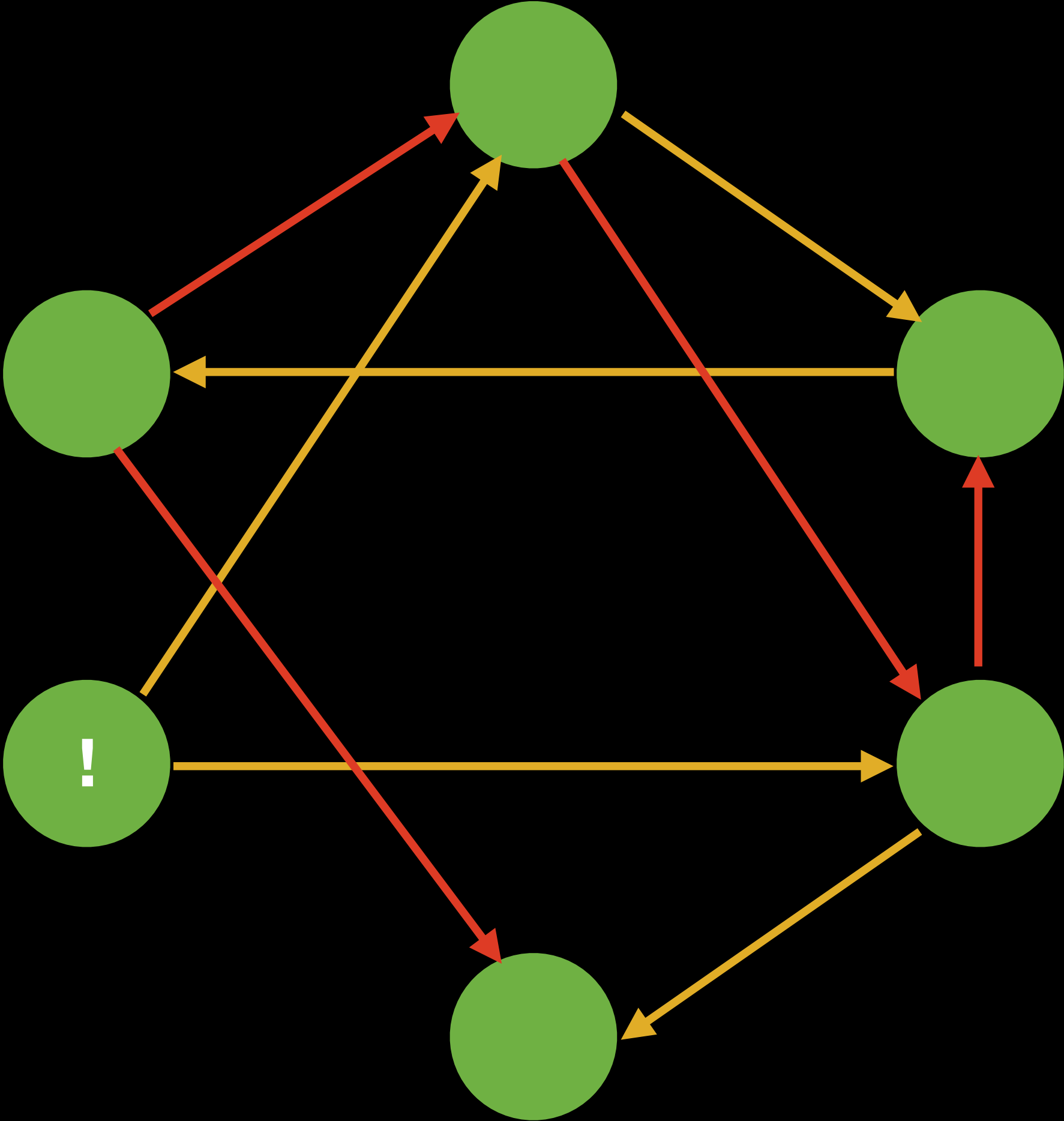
- All nodes start with full "eager" set

- All nodes start with full "eager" set

- Broadcast triggers eager-push

- All nodes start with full "eager" set

- Broadcast triggers eager-push

- Duplicate messages cause "pruning" (move to "lazy")

- All nodes start with full "eager" set

- Broadcast triggers eager-push

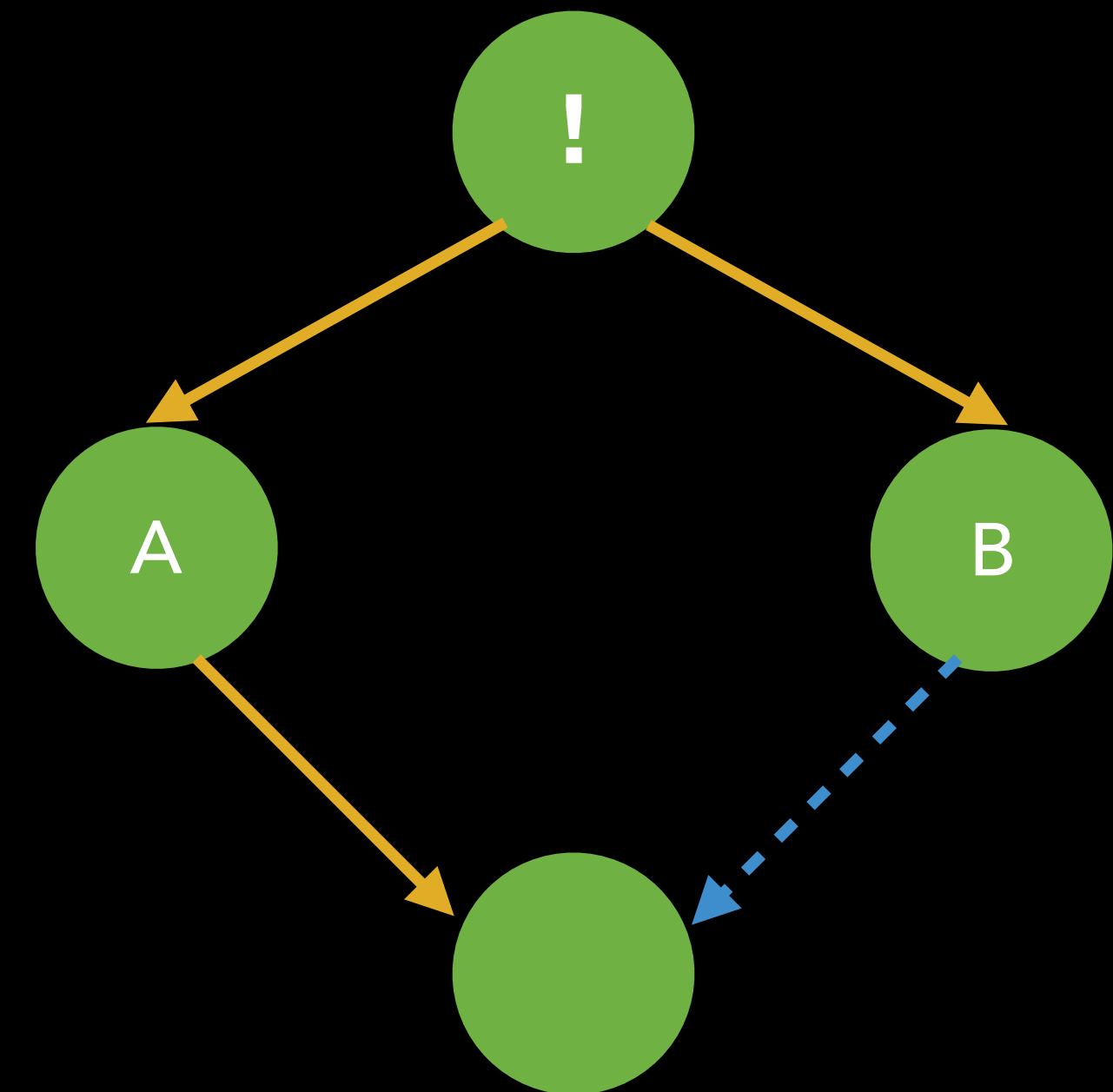- Duplicate messages cause "pruning" (move to "lazy")
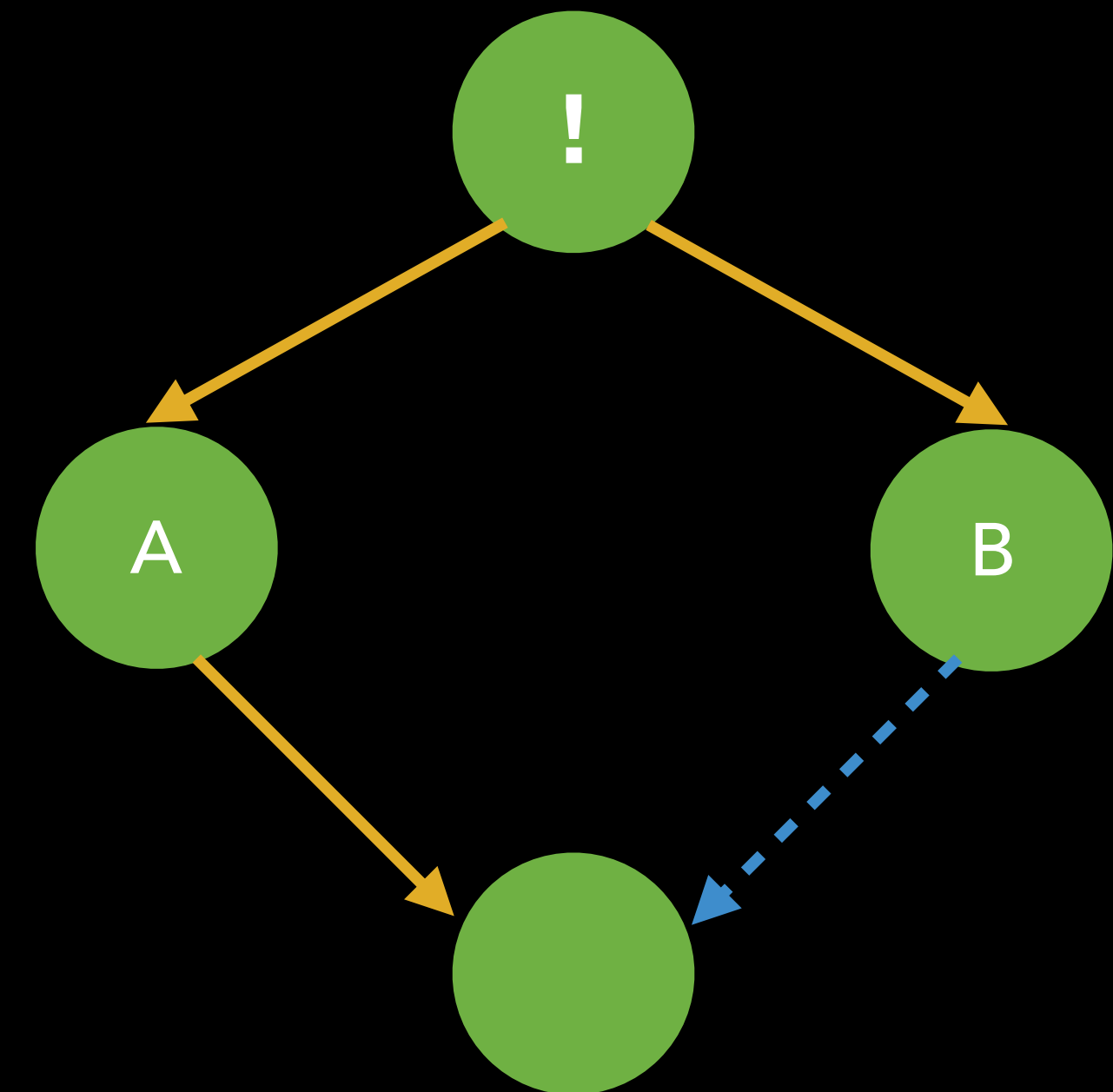
# PLUMTREE - 2009 CONSTRUCTION

- All nodes start with full "eager" set

- Broadcast triggers eager-push

- Duplicate messages cause "pruning" (move to "lazy")

- Regular broadcasts proceed with new "eager" sets

- Lazy-push sends "I Have" messages

# PLUMTREE - 2009 REPAIR

- Lazy-push sends "I Have" messages

- Timeout triggers "grafting" (move to "eager")



COMCAST

# PLUMTREE - 2009 REPAIR

- Lazy-push sends "I Have" messages

- Timeout triggers "grafting" (move to "eager")

# PLUMTREE - 2009 REPAIR

- Lazy-push sends "I Have" messages

- Timeout triggers "grafting" (move to "eager")

- Lazy-push batched to reduce overhead

WE CHOSE PLUMTREE

COMCAST

- **Good tradeoff** between reliability and redundancy

COMCAST

# WE CHOSE PLUMTREE

- **Good tradeoff** between reliability and redundancy

- Optimizes for **lowest-latency paths**

COMCAST

# WE CHOSE PLUMTREE

- **Good tradeoff** between reliability and redundancy

- Optimizes for **lowest-latency paths**

- **Existing** open-source implementations

COMCAST

# WE CHOSE PLUMTREE

- **Good tradeoff** between reliability and redundancy

- Optimizes for **lowest-latency paths**

- **Existing** open-source implementations

- **Excellent fit** with HyParView

COMCAST

# POPULATION PROTOCOLS

POPULATION PROTOCOLS USE

# RANDOMIZED INTERACTIONS

# Logical Physical Clocks
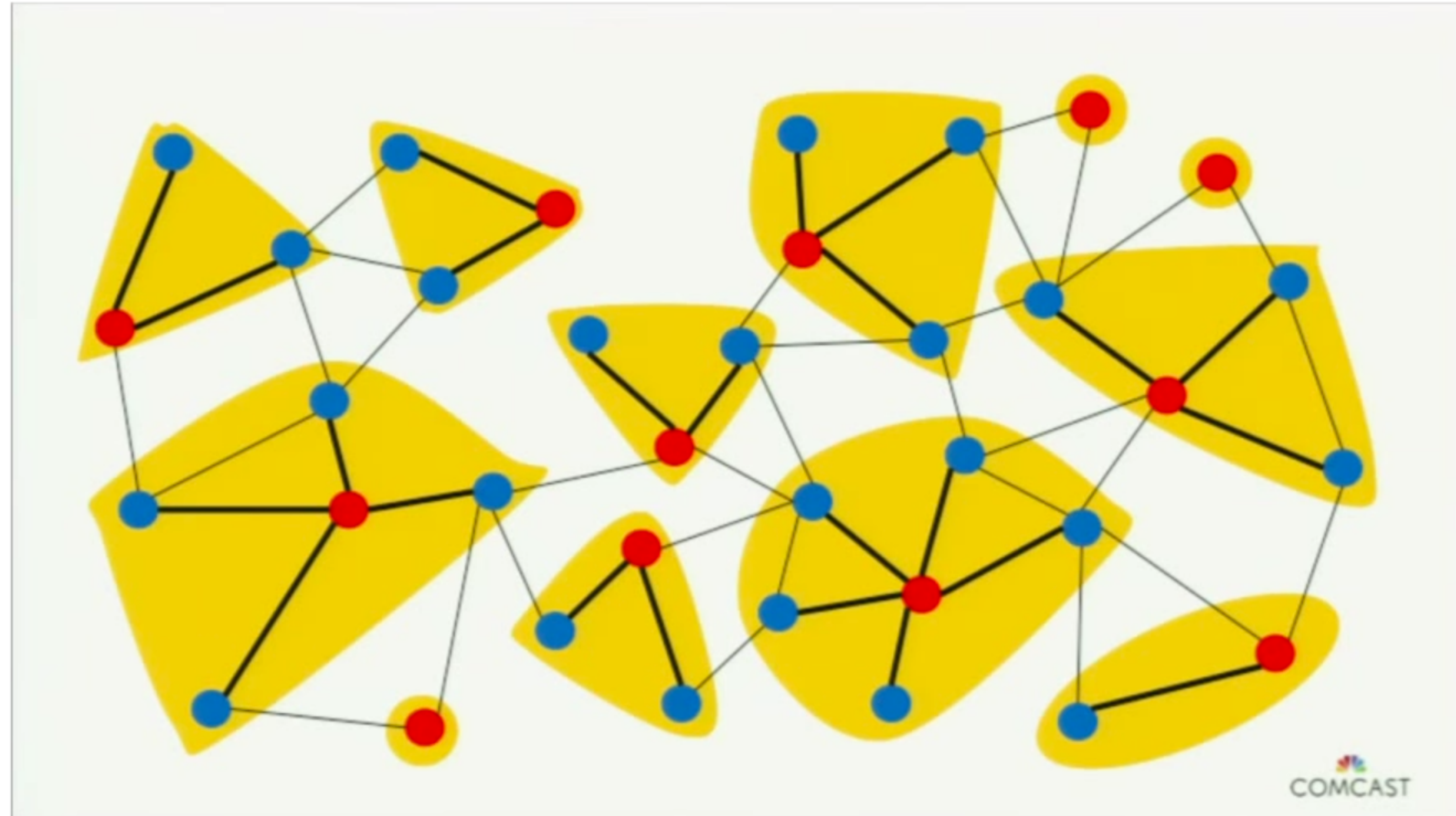# and Consistent Snapshots in Globally Distributed Databases

Sandeep Kulkarni[*], Murat Demirbas[**], Deepak Madeppa[**], Bharadwaj Avva[**], and Marcelo Leone[*]

[*]Michigan State University
[**]University at Buffalo, SUNY

COMCAST

JON MOORE

# DISTRIBUTED MONOTONIC CLOCKS

Vidcap from StrangeLoop 2015: https://youtu.be/YqNGbvFHoKM

# DMC PROBLEMS

- "Wacky clock mode"

- Hierarchy imbalances load

- Long-lived partitions

- No convergence proof

# Bridging the Gap between Population and Gossip-based Protocols

Yann Busnel, Marin Bertier, Anne-Marie Kermarrec

- Use **existing dissemination** with DMC

- Transmit clocks **along with other messages**

- Use monotonic clocks as a **drift-detection mechanism**

COMCAST

LESSONS LEARNED

Vidcap from RICON West: https://youtu.be/s4cCUTPU8GI

# THANK YOU!

## @SEANCRIBBS