

BRISA:

Combining Efficiency and Reliability in Epidemic Data Dissemination

Miguel Matos^{*}, Valerio Schiavoni⁺, Pascal Felber⁺
Rui Oliveira^{*} and Etienne Rivière⁺

^{*}INESC TEC & University of Minho, Portugal
⁺Université de Neuchâtel, Switzerland

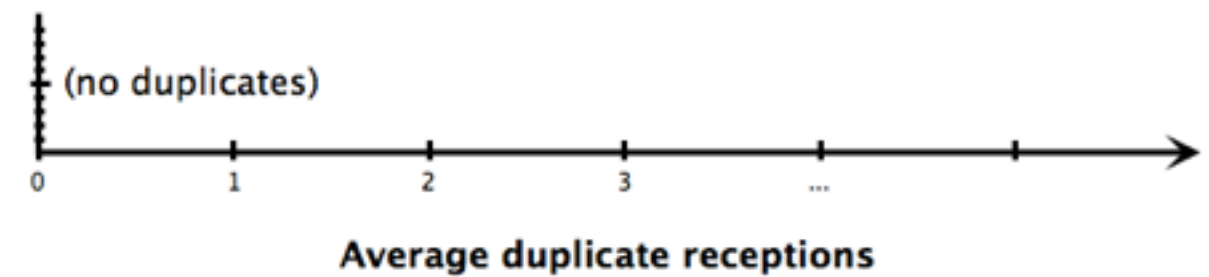
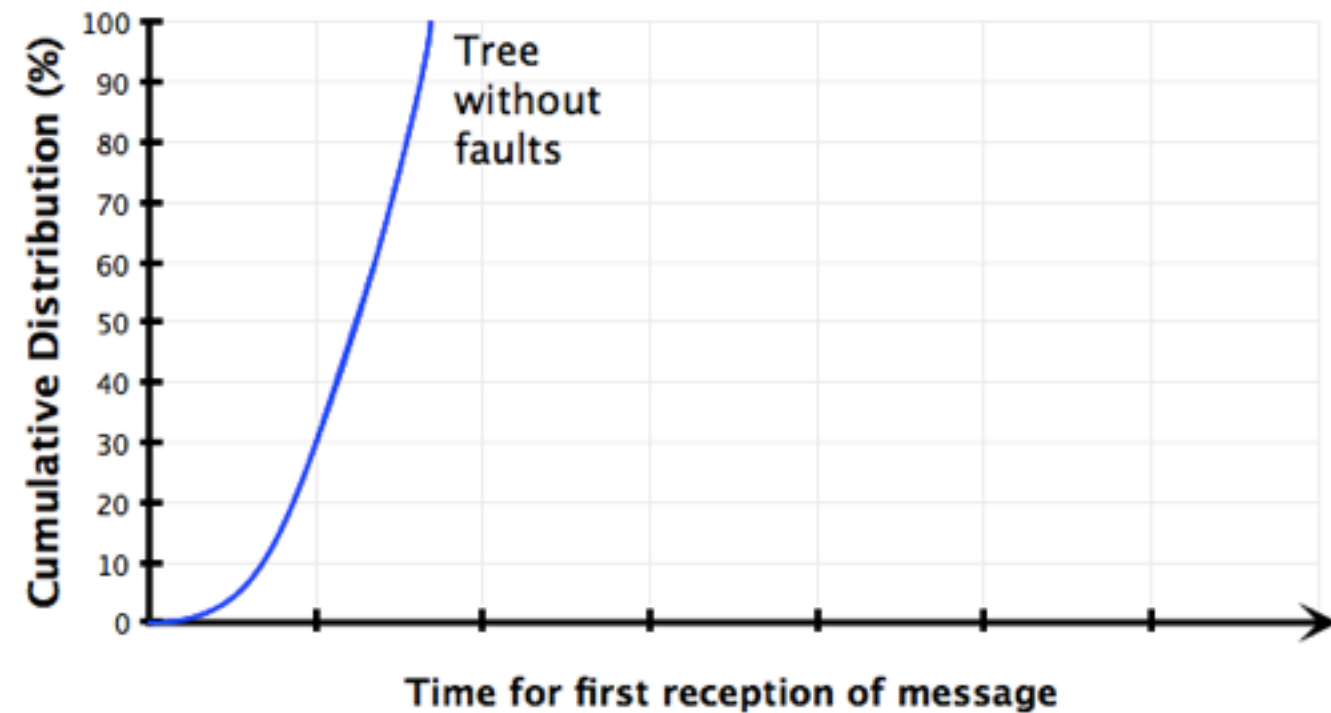
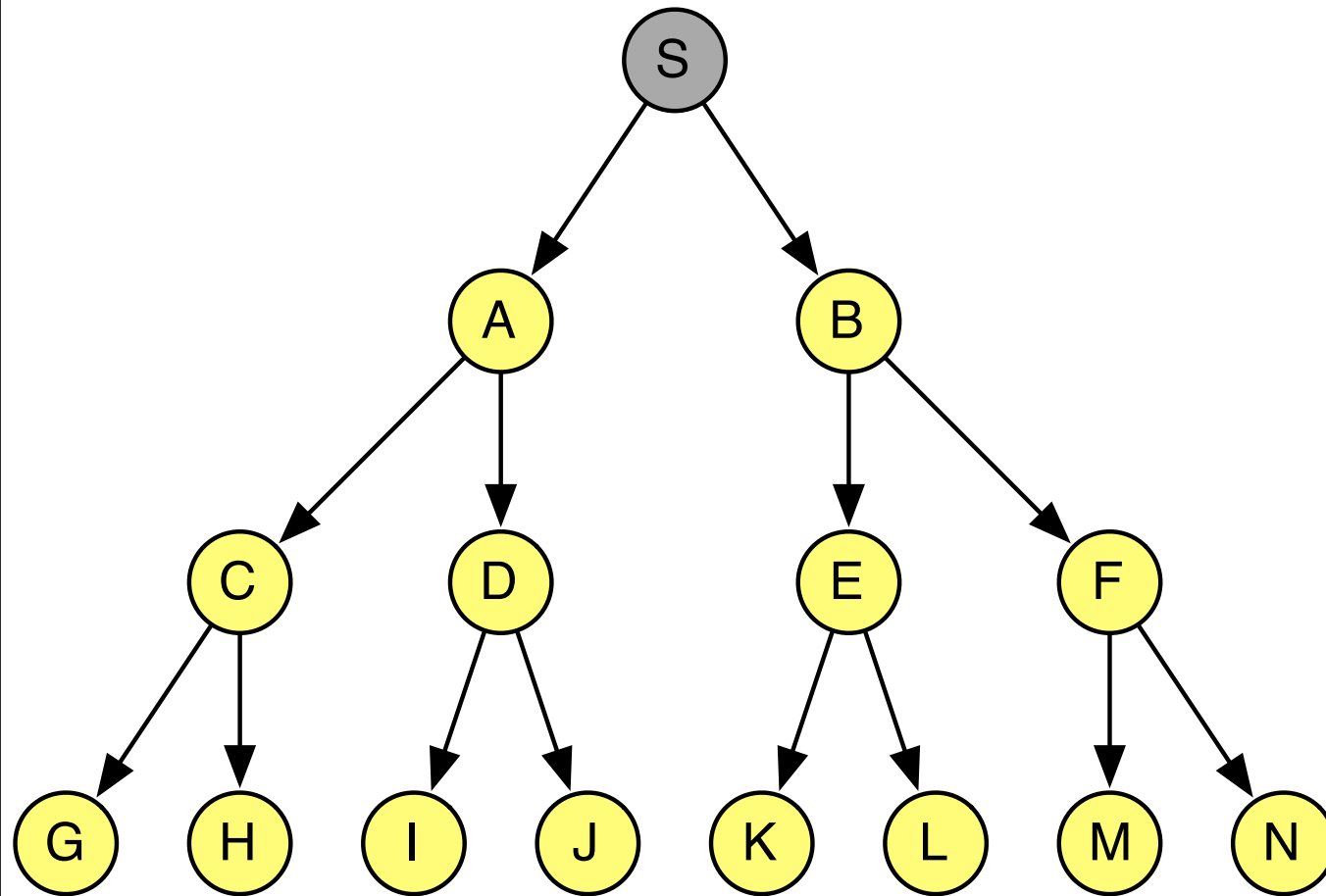
miguelmatos@di.uminho.pt

International Parallel & Distributed Processing Symposium
24 May, 2012

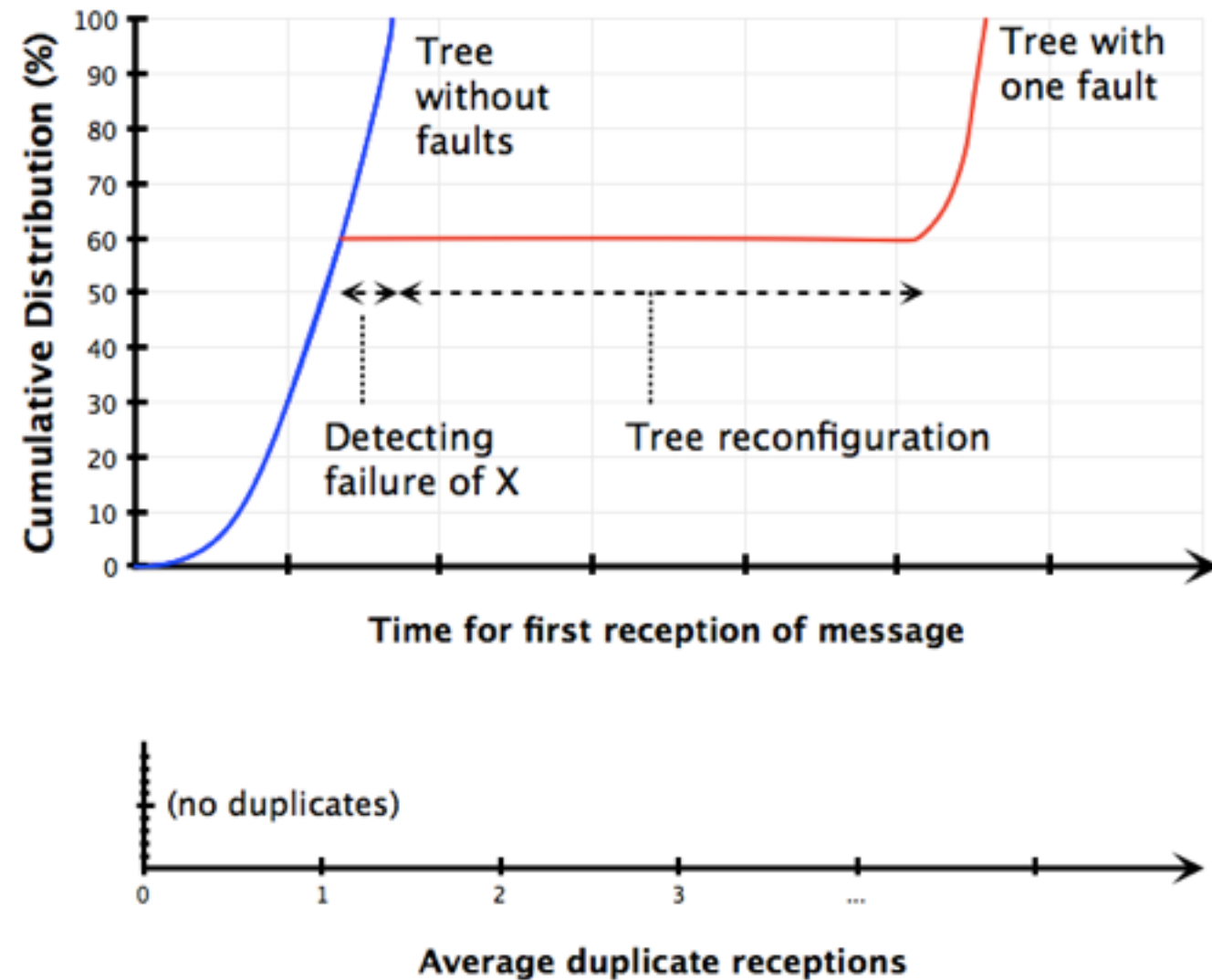
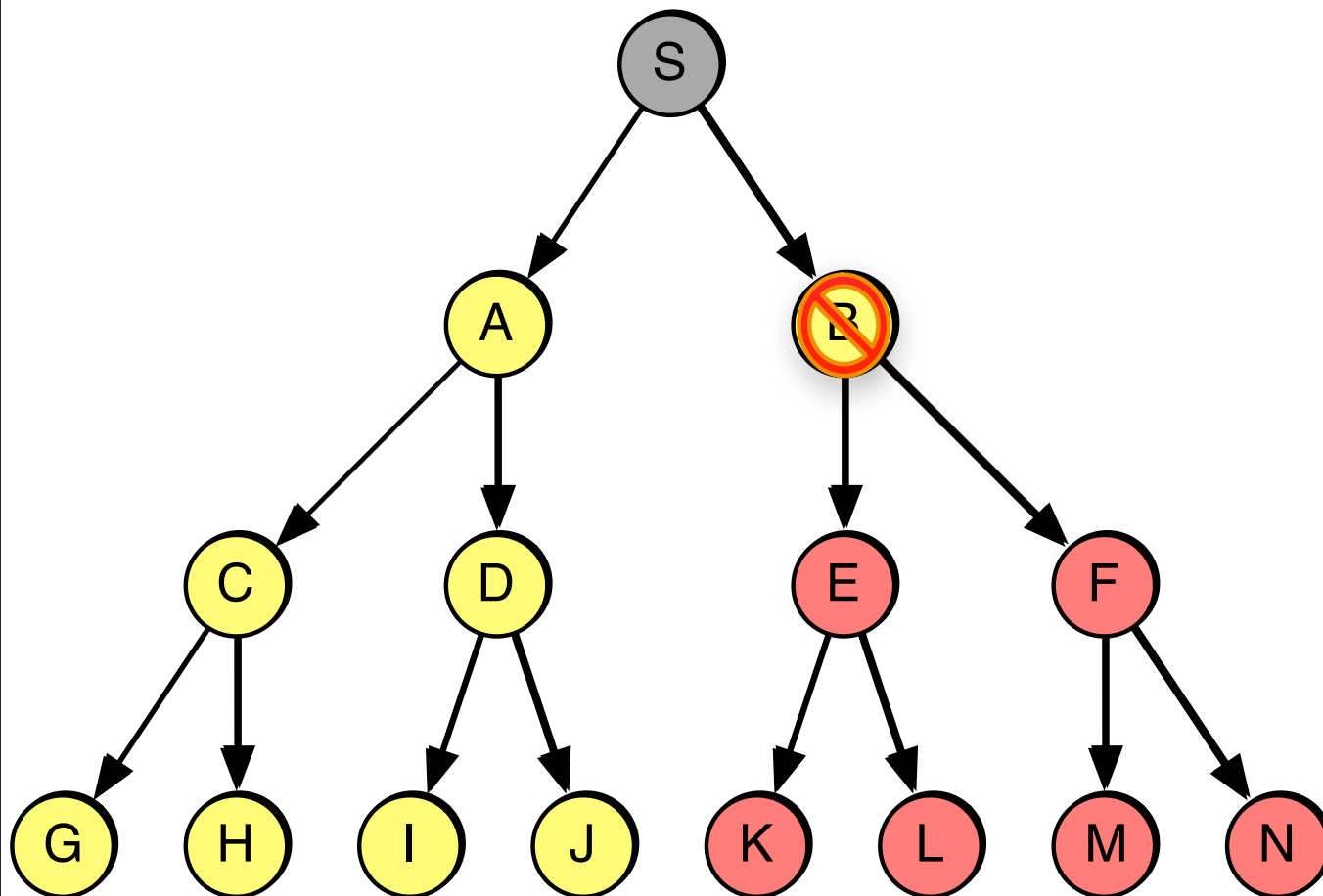
Motivations

- **Broadcast** of small or medium sized messages
- **Challenging** environment
 - (very) **Large** scale
 - Highly **dynamic**
- Performance criteria
 - Low **delays** for all
 - Distribution of first reception time
 - Low **overhead** (overlay construction)
 - Low **message** count (for disseminations)
- + **Simplicity!**

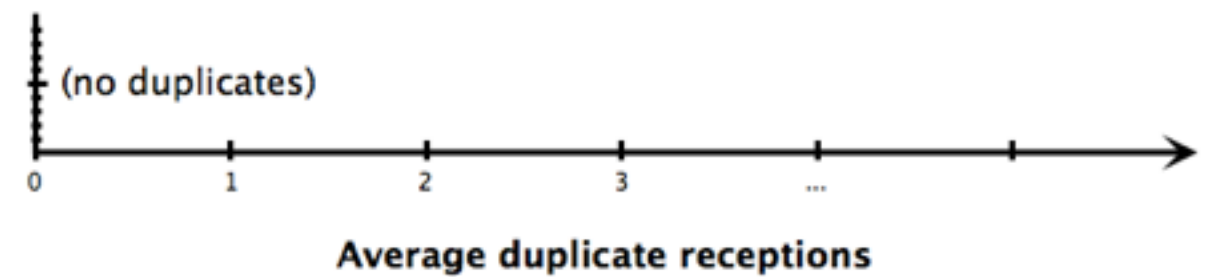
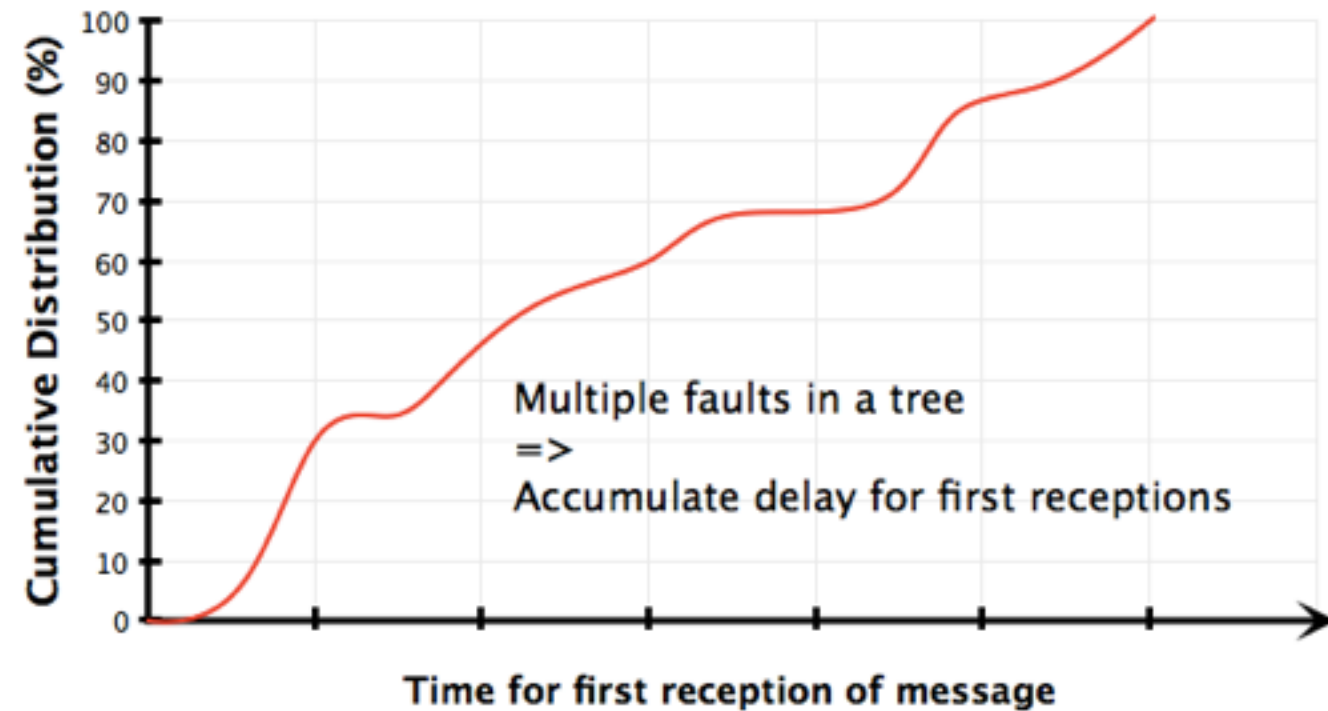
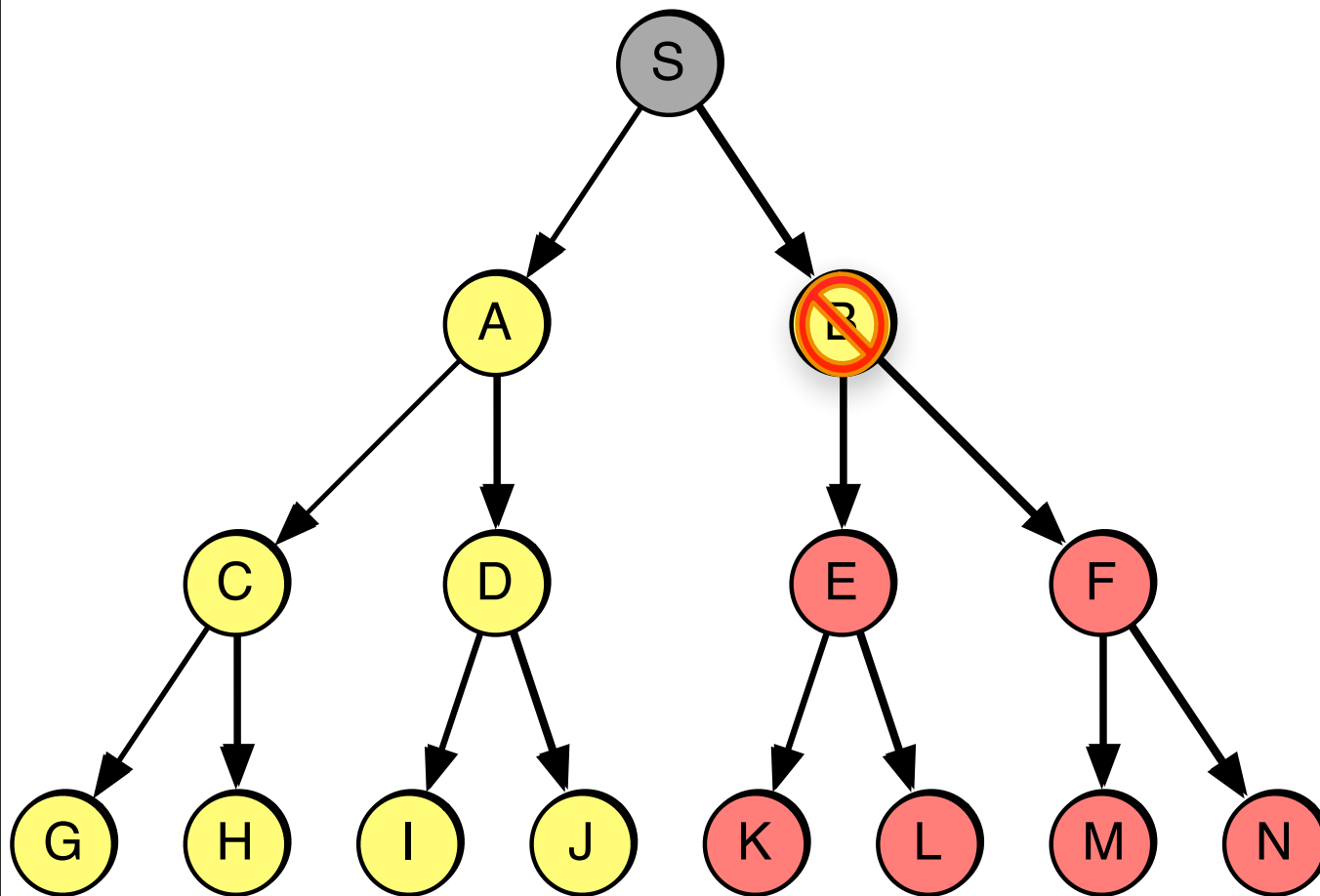
Using a tree?



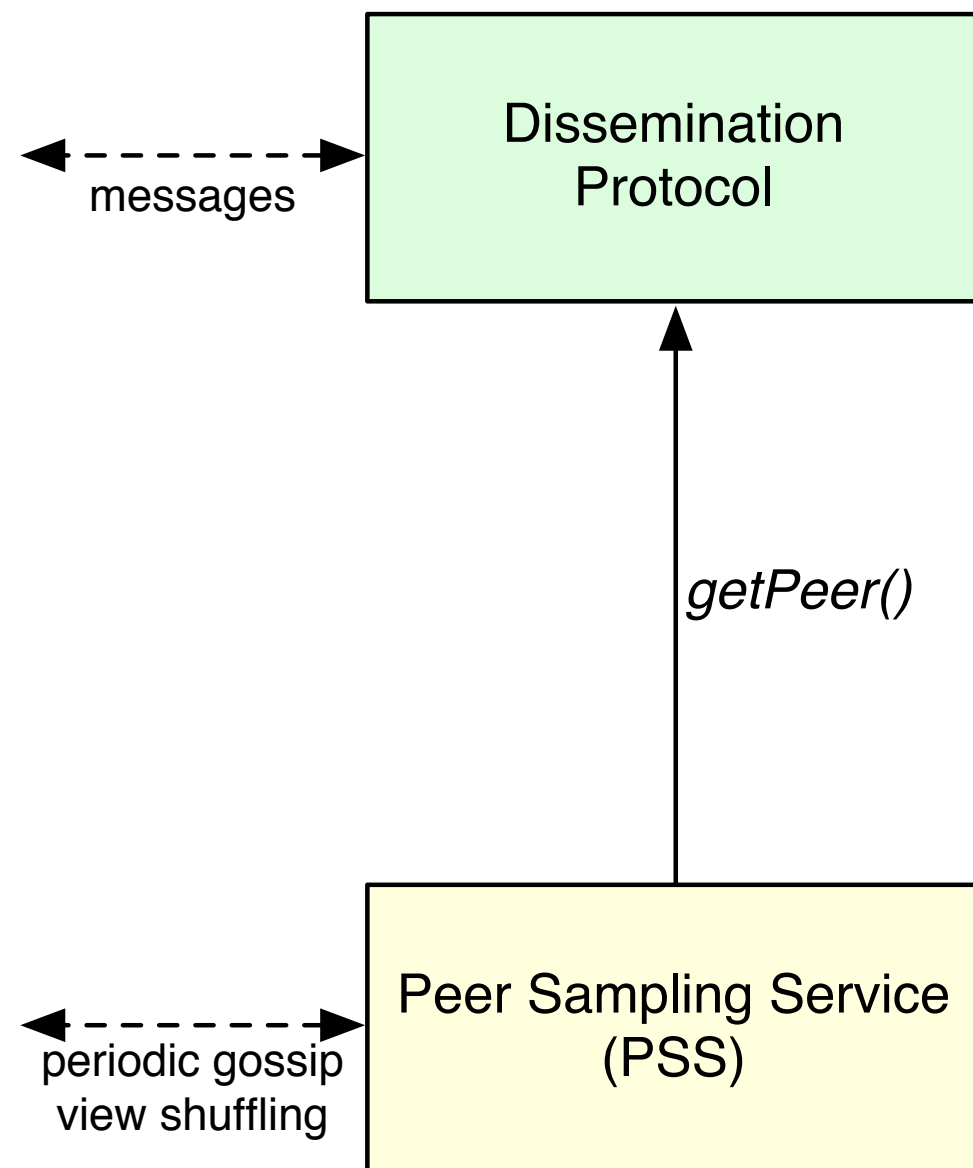
Faults happen



Faults cumulate



Classic epidemic-based dissemination



Flooding: send message to all nodes upon first reception

=> **Probabilistic Broadcast**

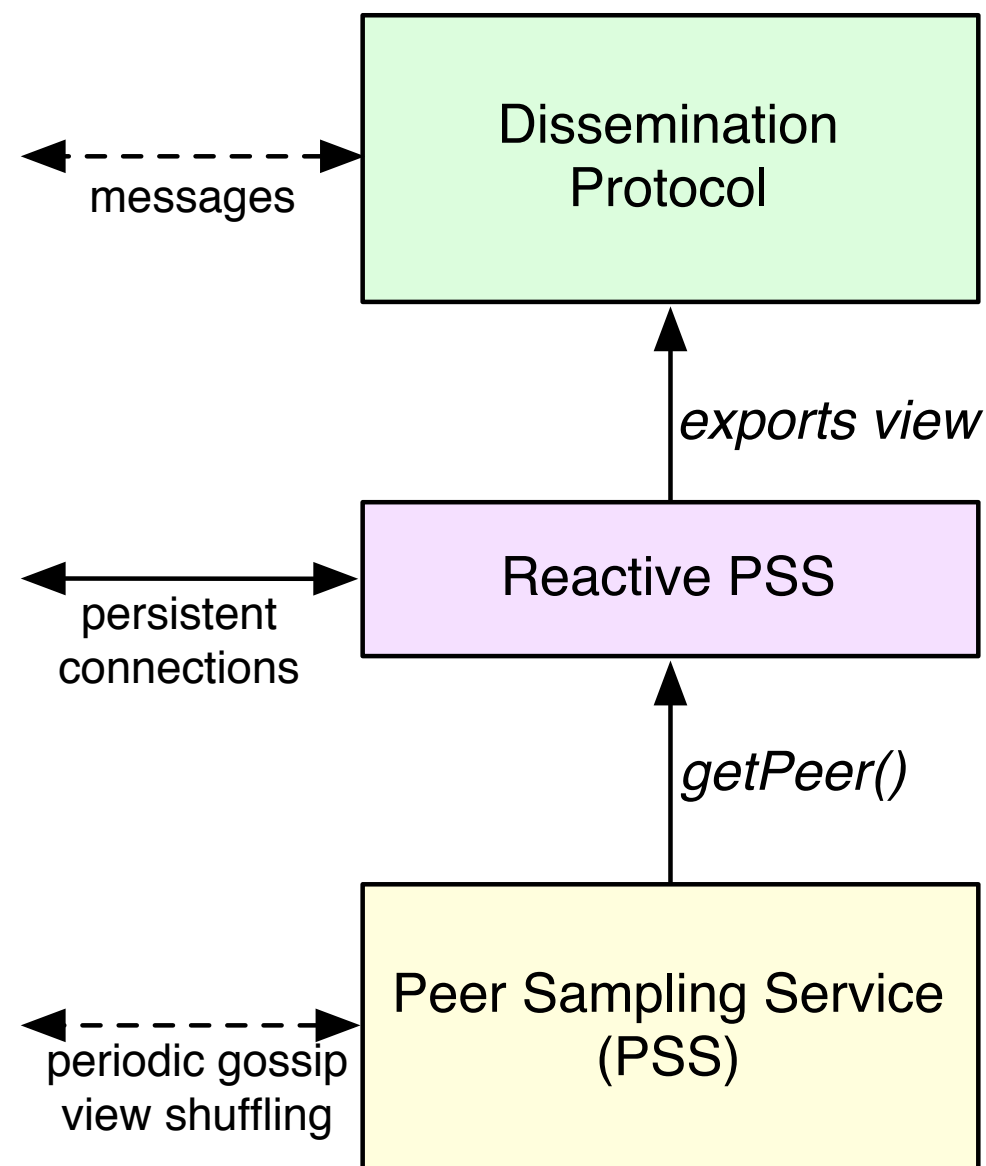
Provides constant *stream* of random peers

=> **Membership Management**

=> **Avoids Partition**

=> **Insert new peers, forget old ones**

Epidemic-based dissemination on reactive PSS



Deterministic flooding: send message to all nodes on first reception

=> **Total Broadcast**

Reactive view management

Maintains **persistent** connections to peers

Bi-directional links

Fault Detection

Replaces failed peers from PSS `getPeer()`

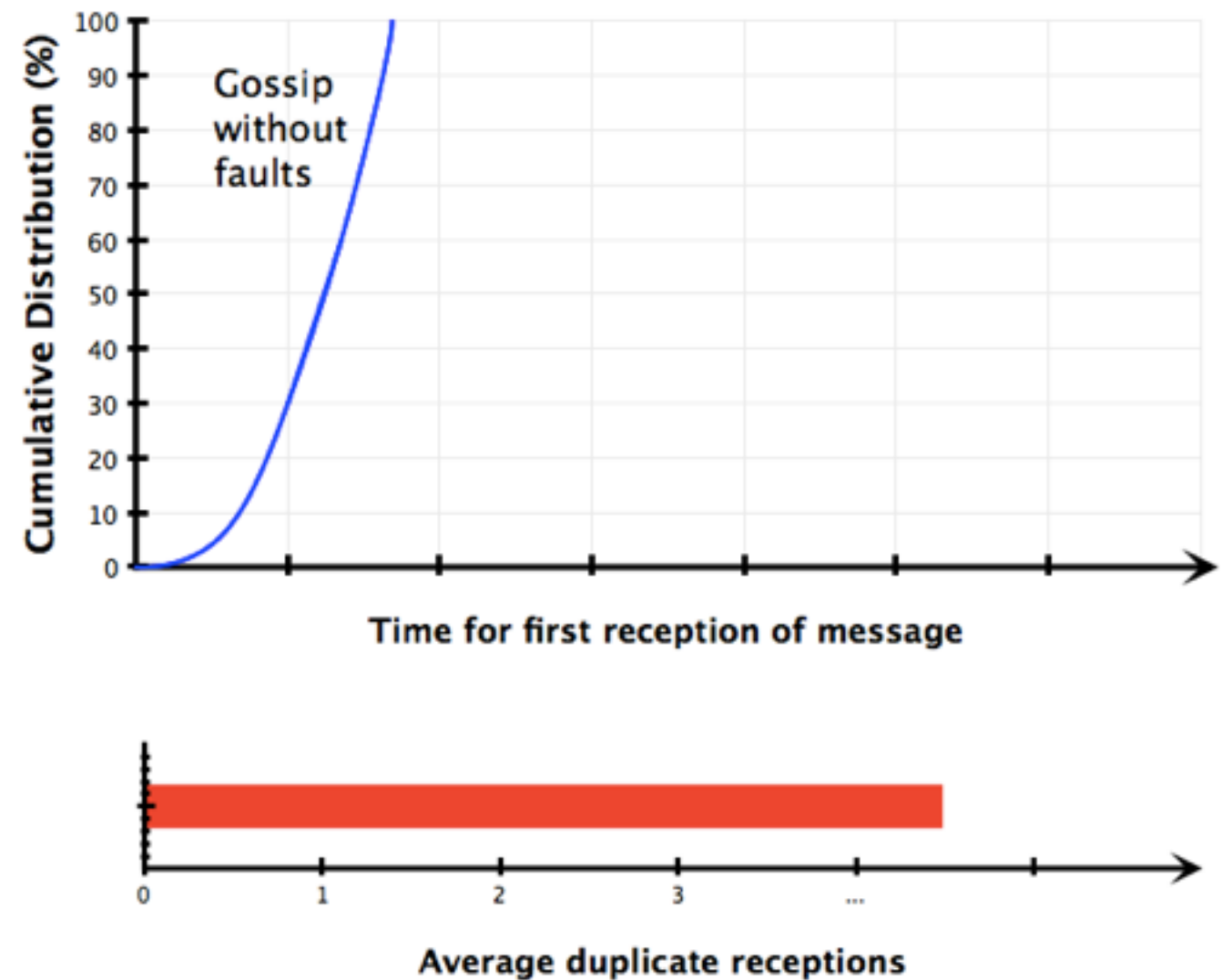
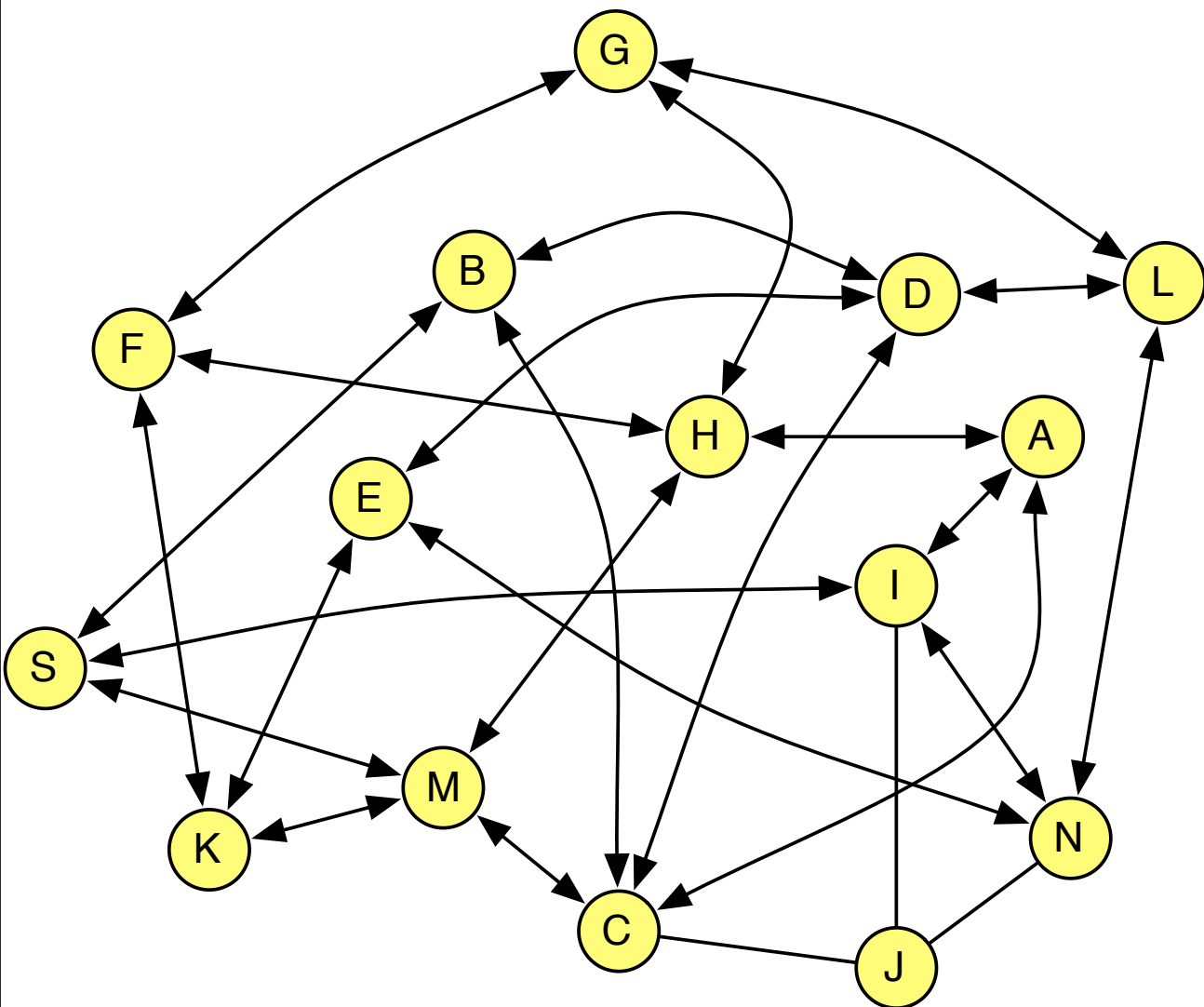
Provides constant *stream* of random peers

=> **Membership Management**

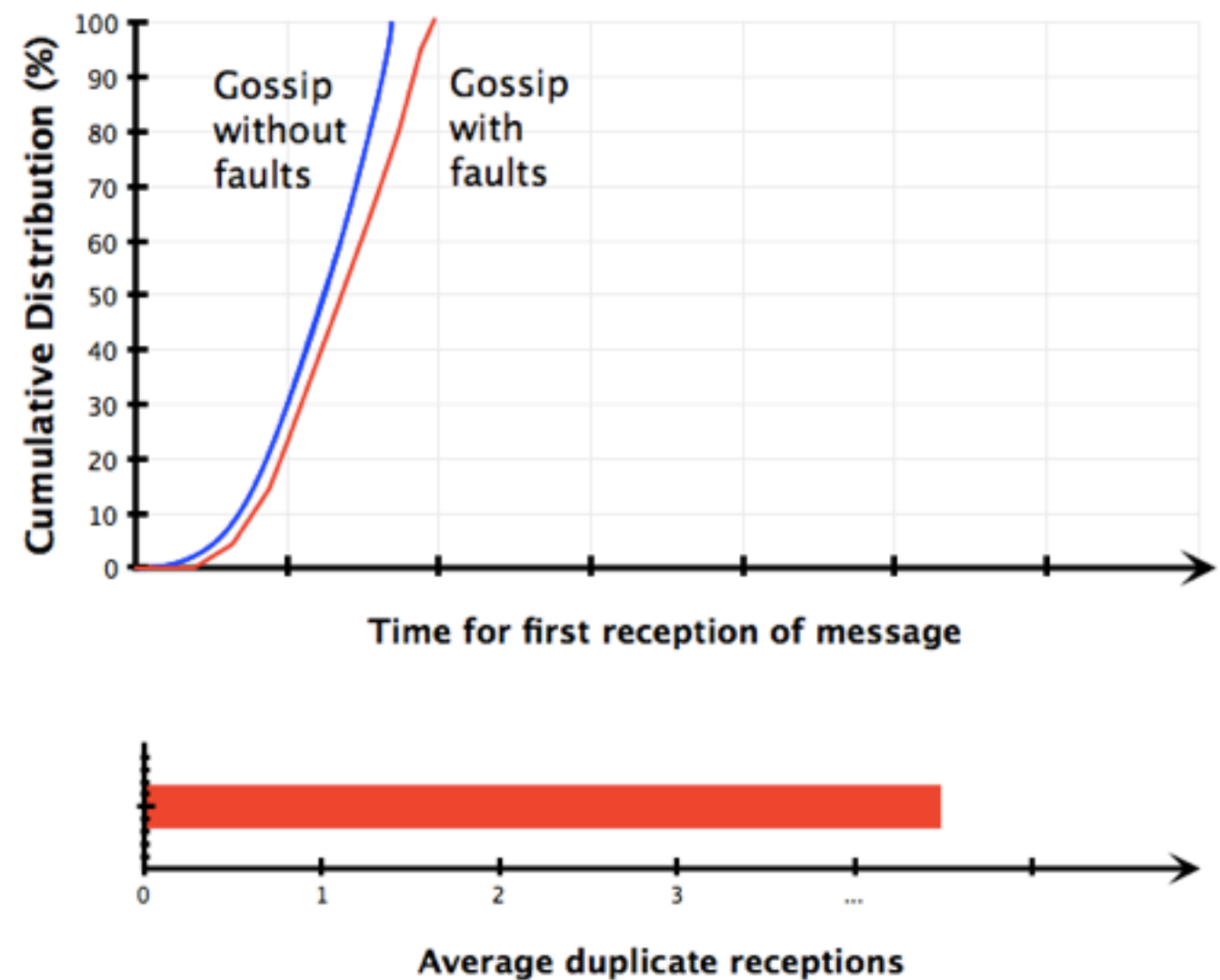
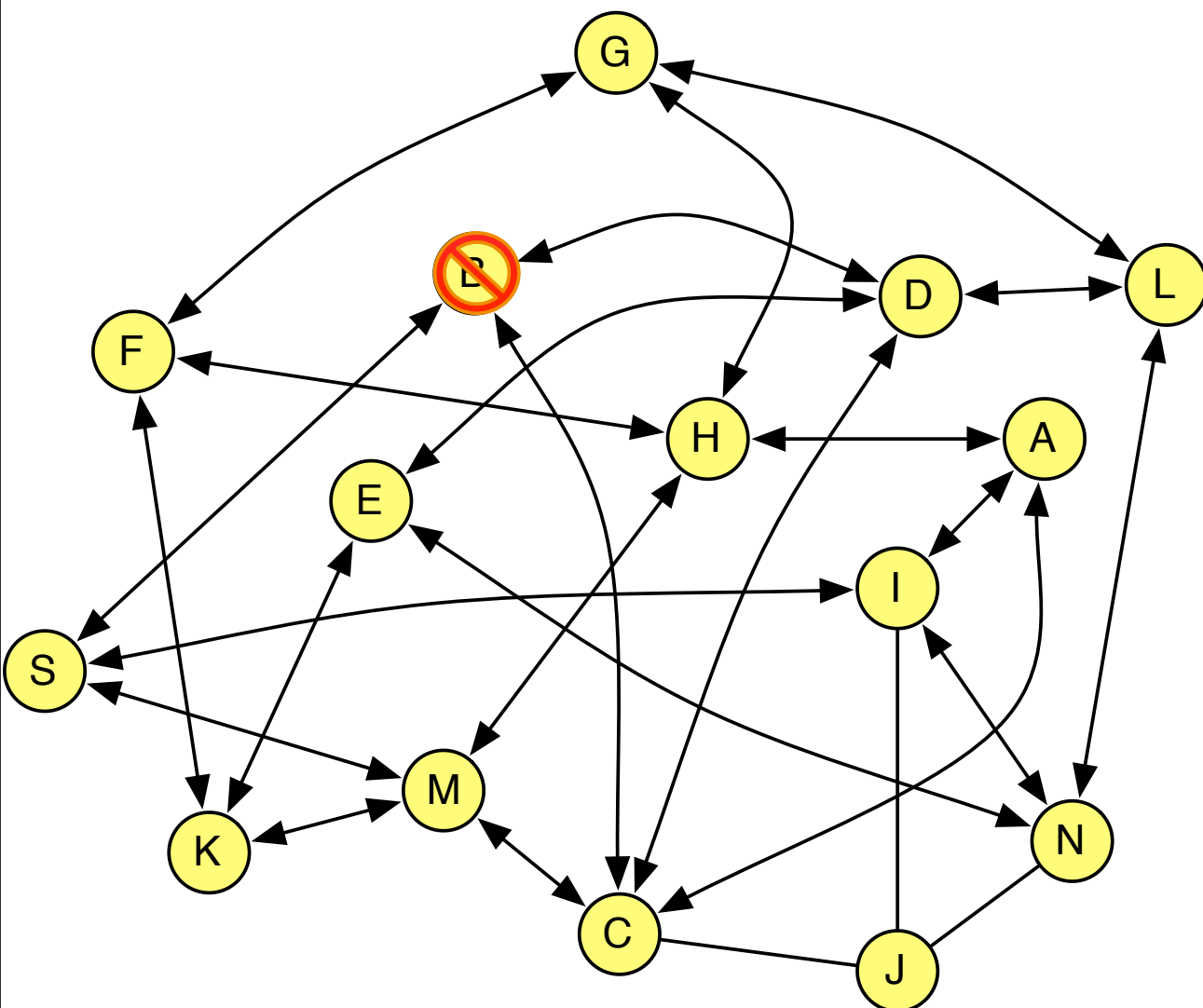
=> **Avoids Partition**

=> **Insert new peers, forget old ones**

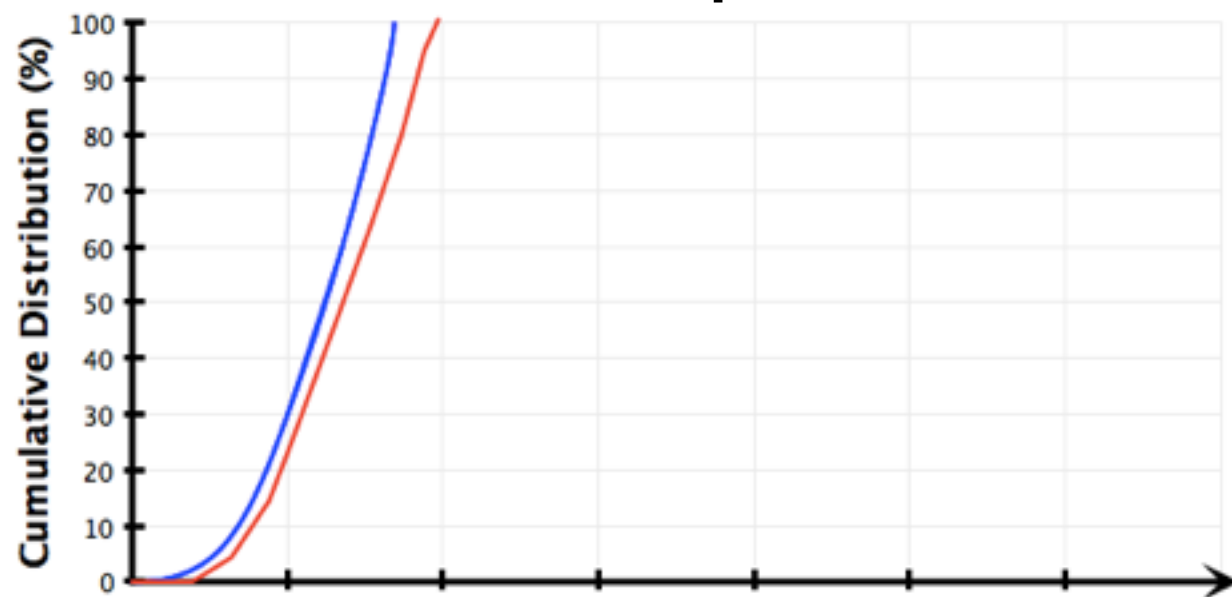
Epidemic-based dissemination



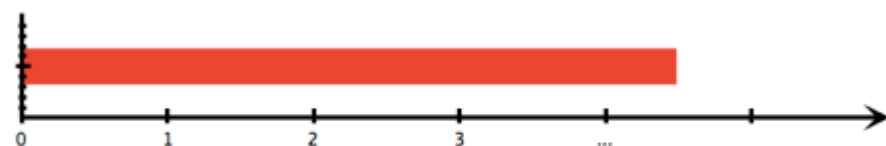
Epidemic-based dissemination



Gossip

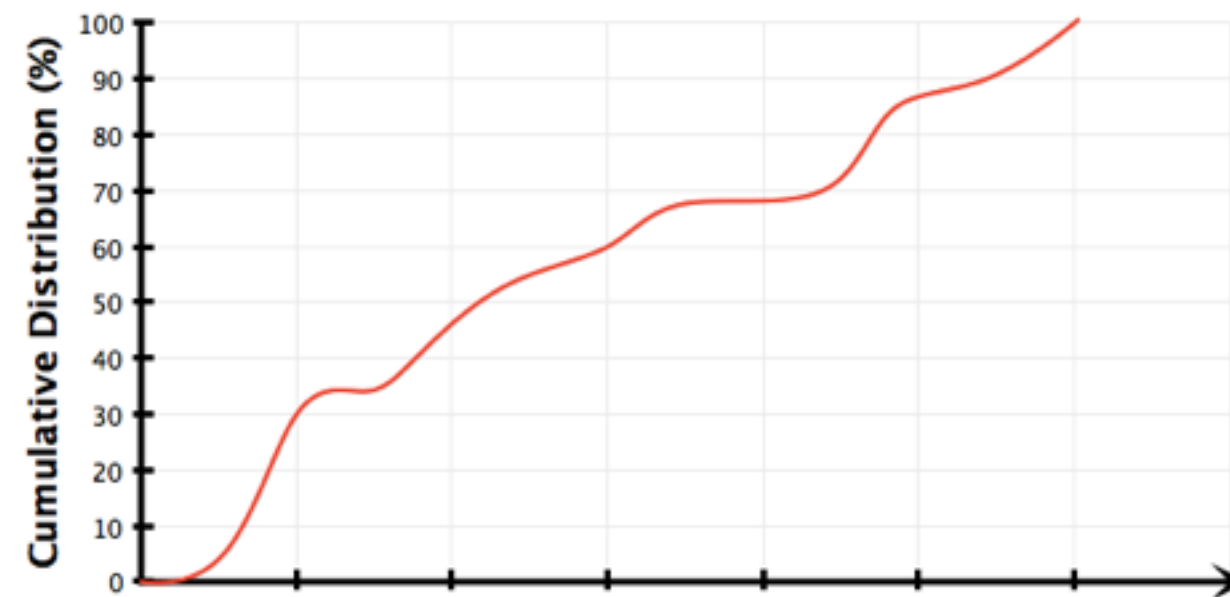


Time for first reception of message

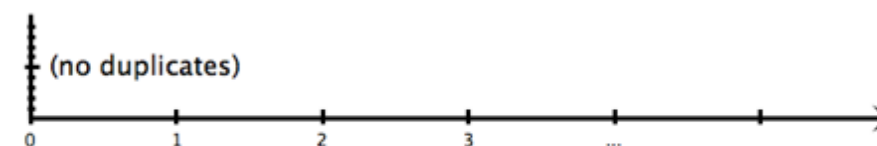


Average duplicate receptions

Tree

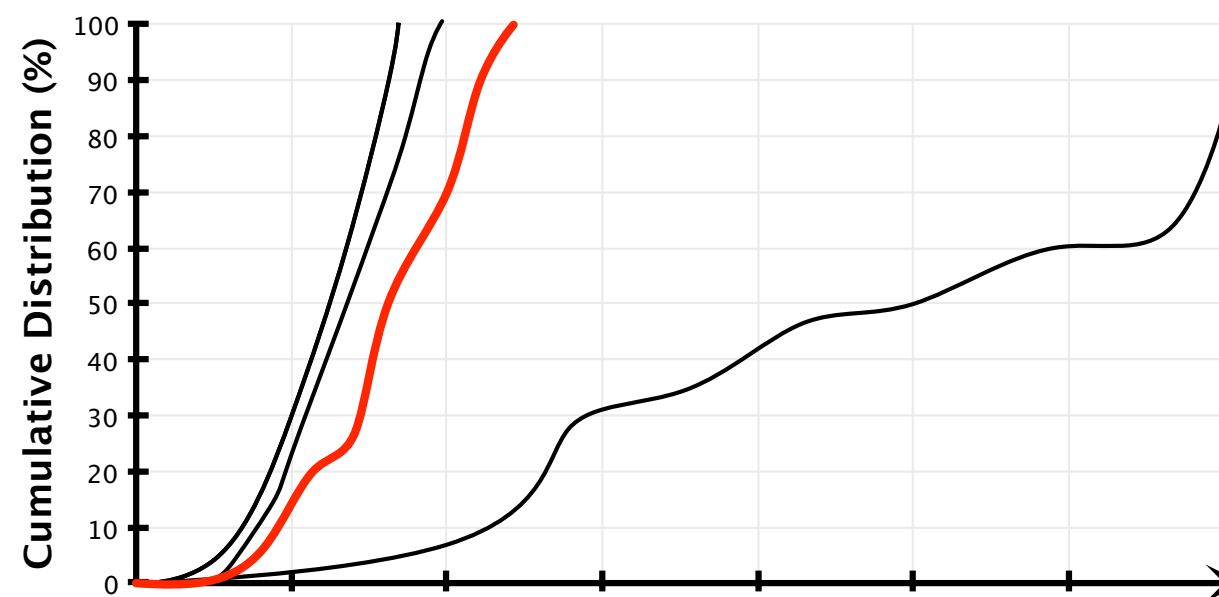


Time for first reception of message

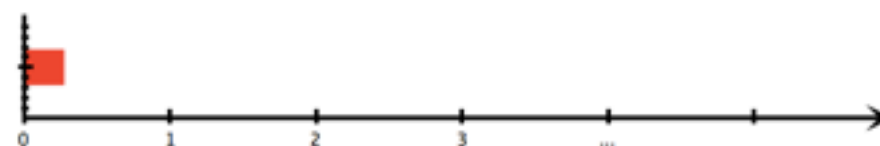


Average duplicate receptions

Objective



Time for first reception of message

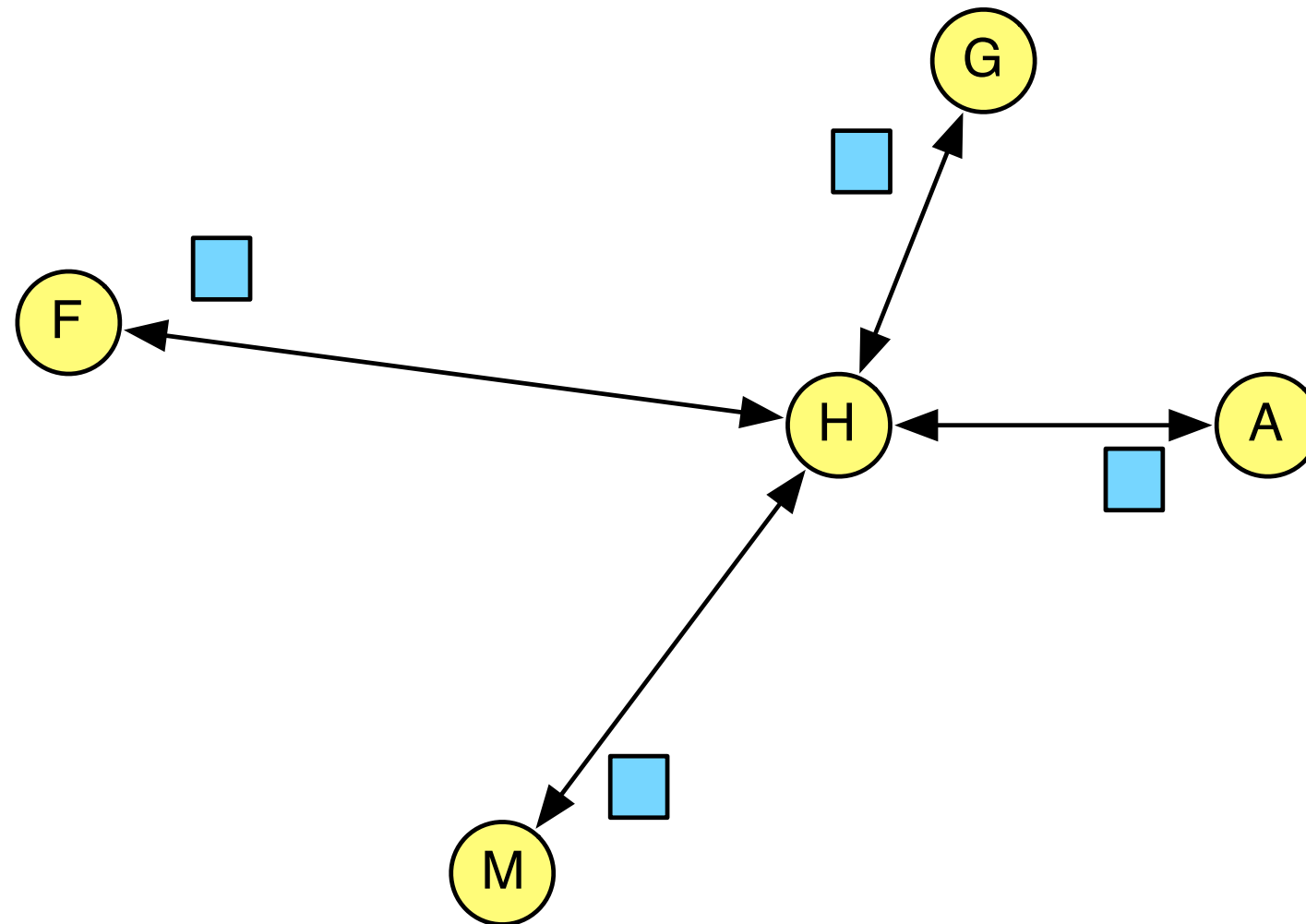


Average duplicate receptions

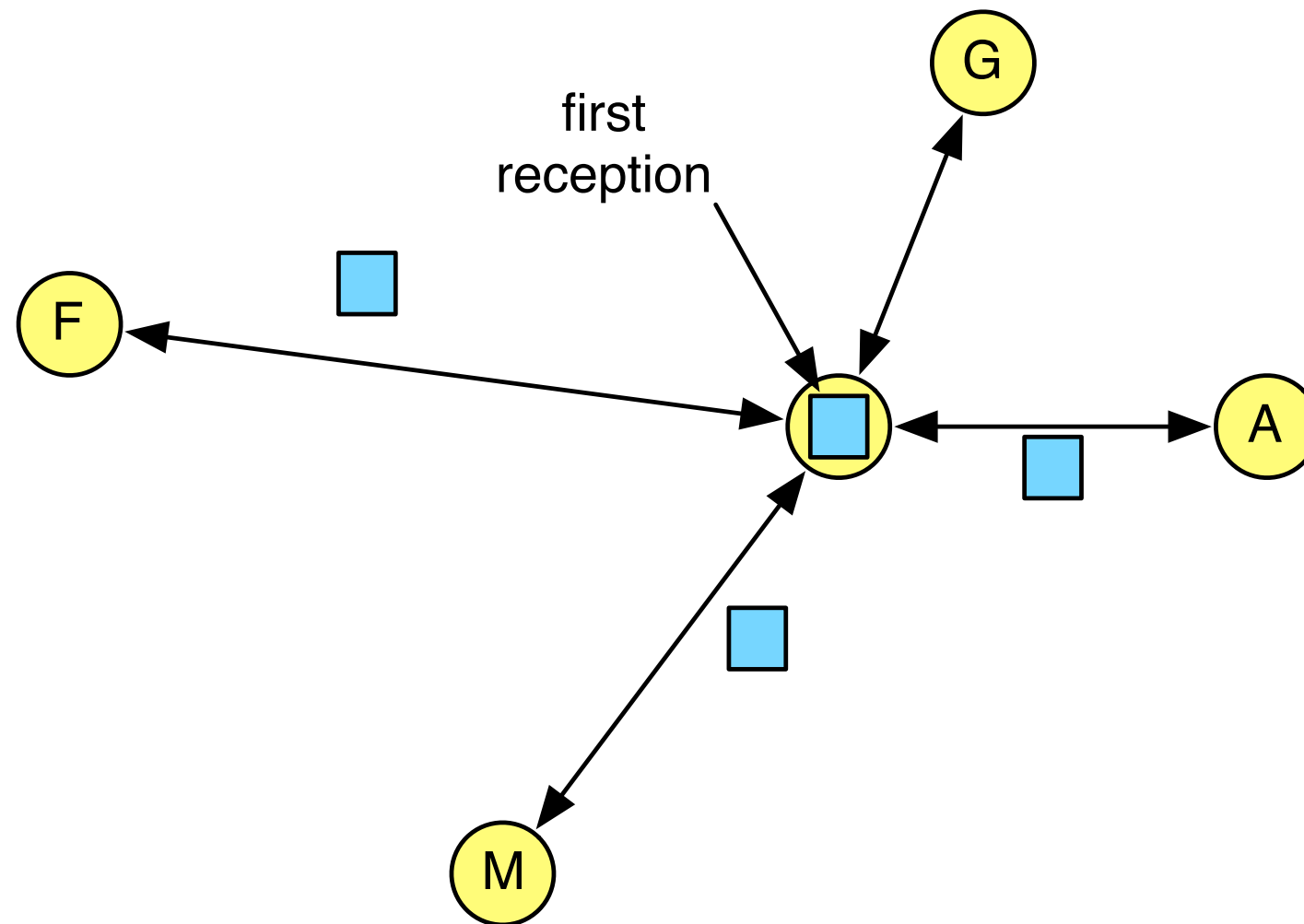
BRISA *in a nutshell*

- Start from existing *Reactive Peer Sampling Service* (RPSS)
 - First dissemination = flooding over RPSS
- Emerge an **embedded tree**
 - Select between the links maintained by the RPSS
 - One of the receptions = **parent** link
 - Others **deactivated**
- Allow only a **small number of duplicates**
 - Upon join
 - Upon failure
- Quick **recovery** based on RPSS persistent connections and quick fault detection

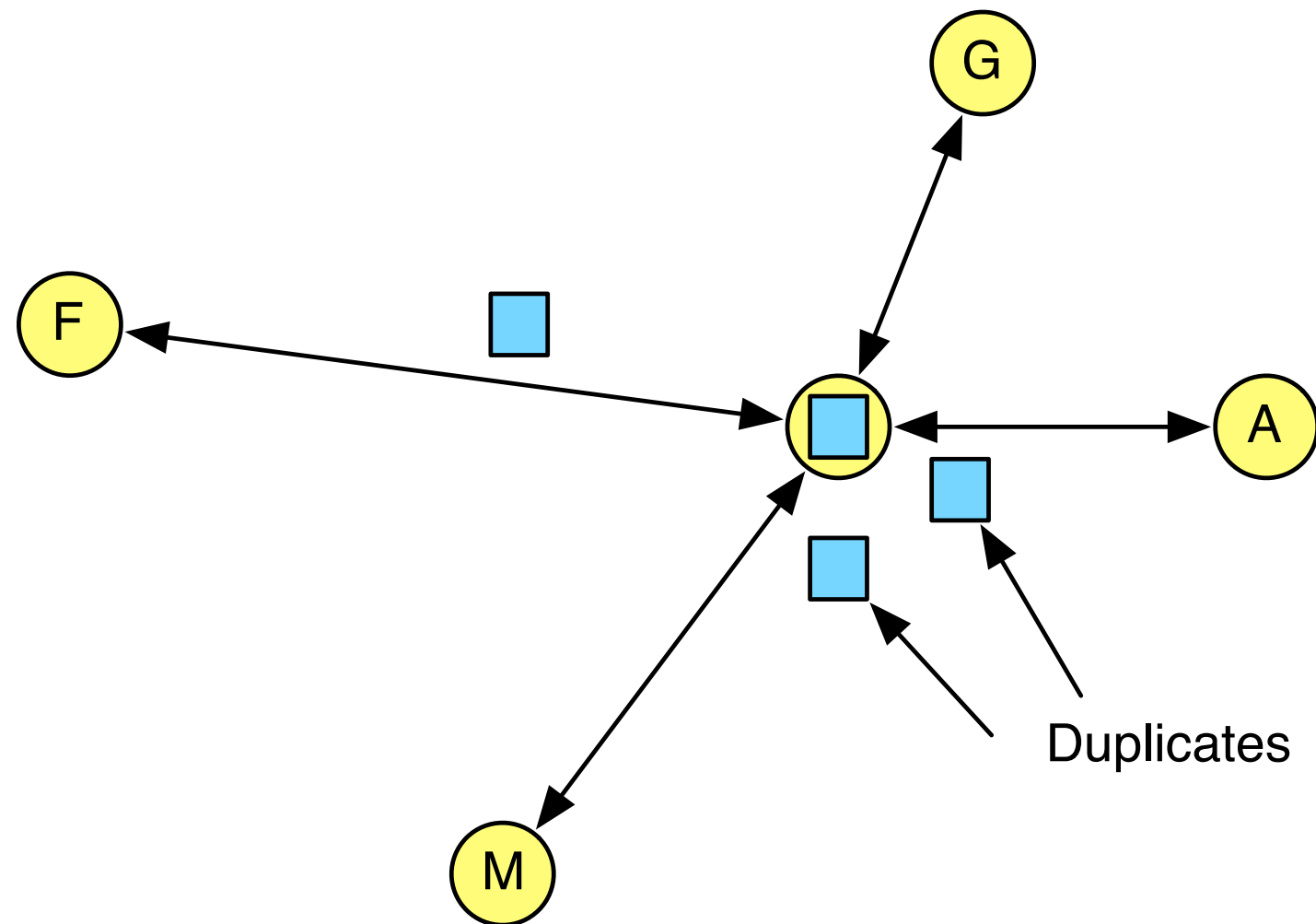
Emerging a tree through deactivations



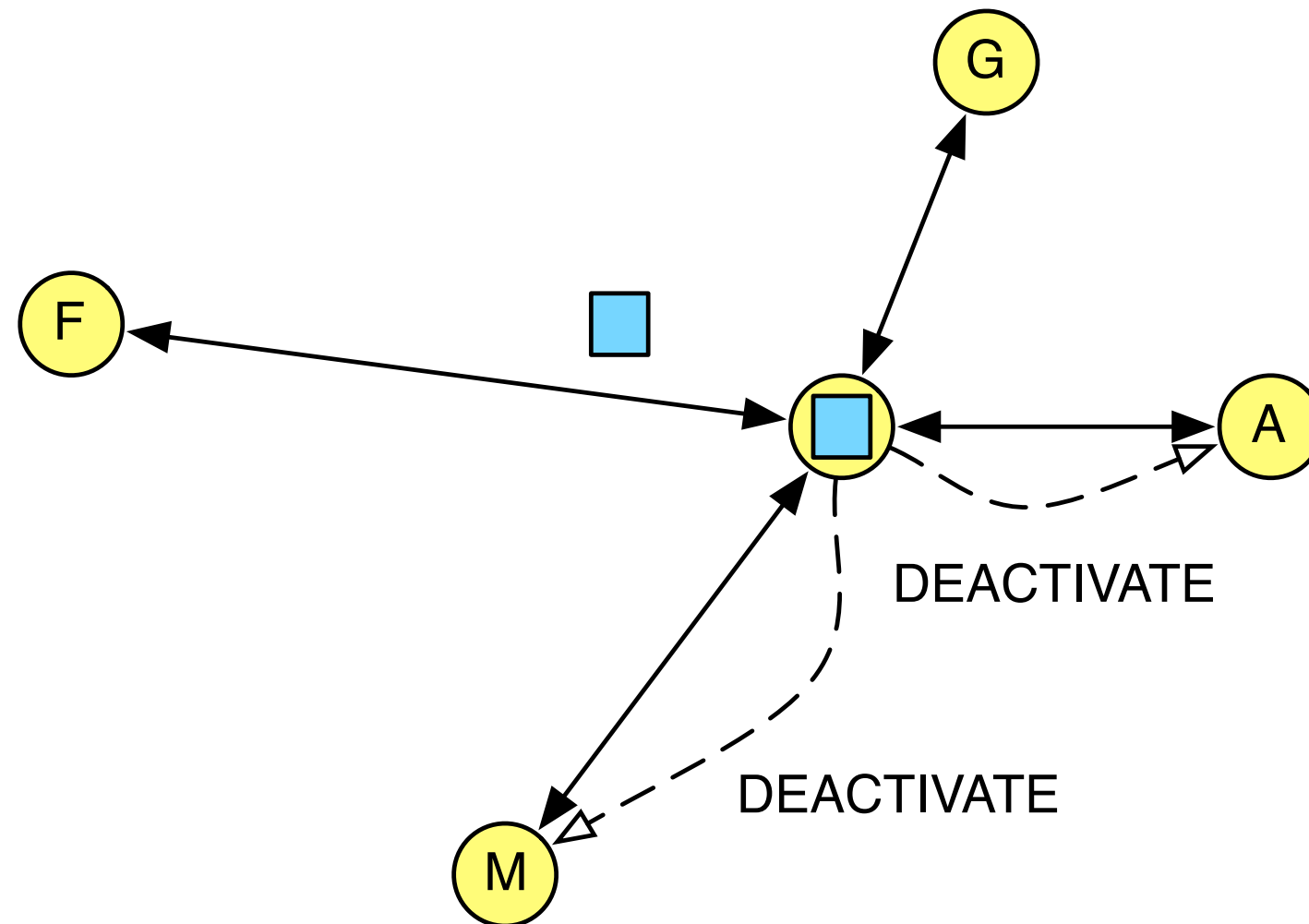
Emerging a tree through deactivations



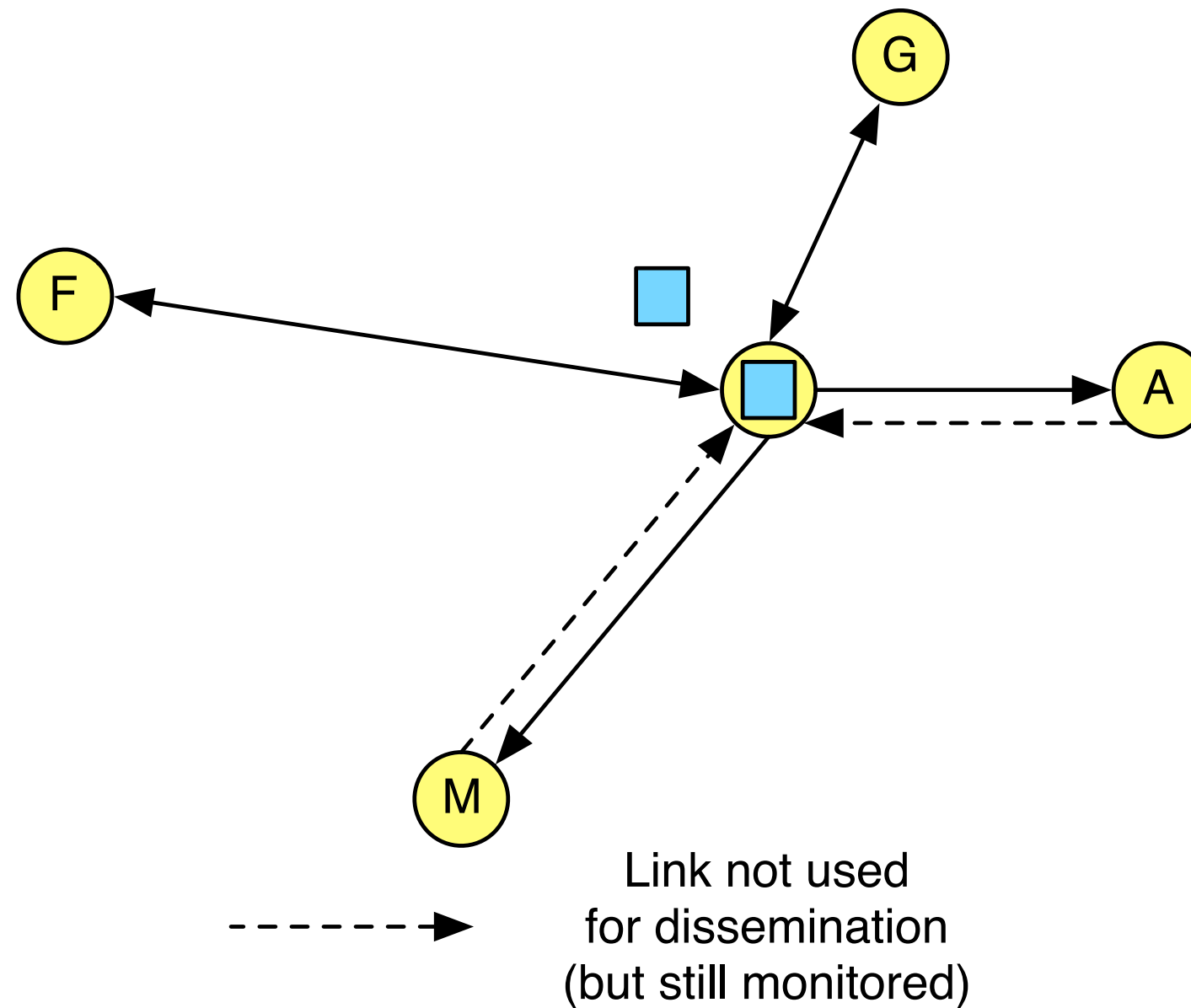
Emerging a tree through deactivations



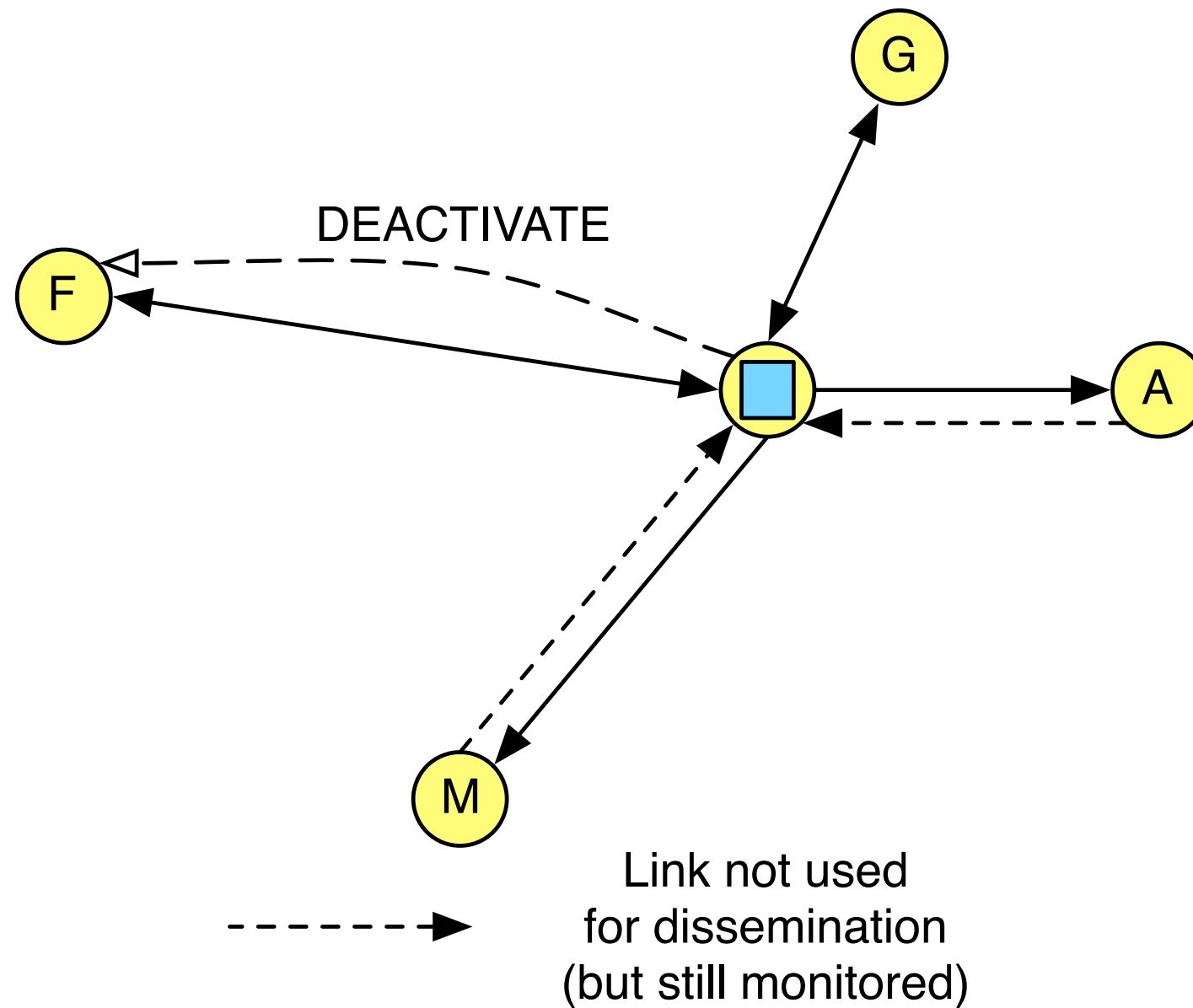
Emerging a tree through deactivations



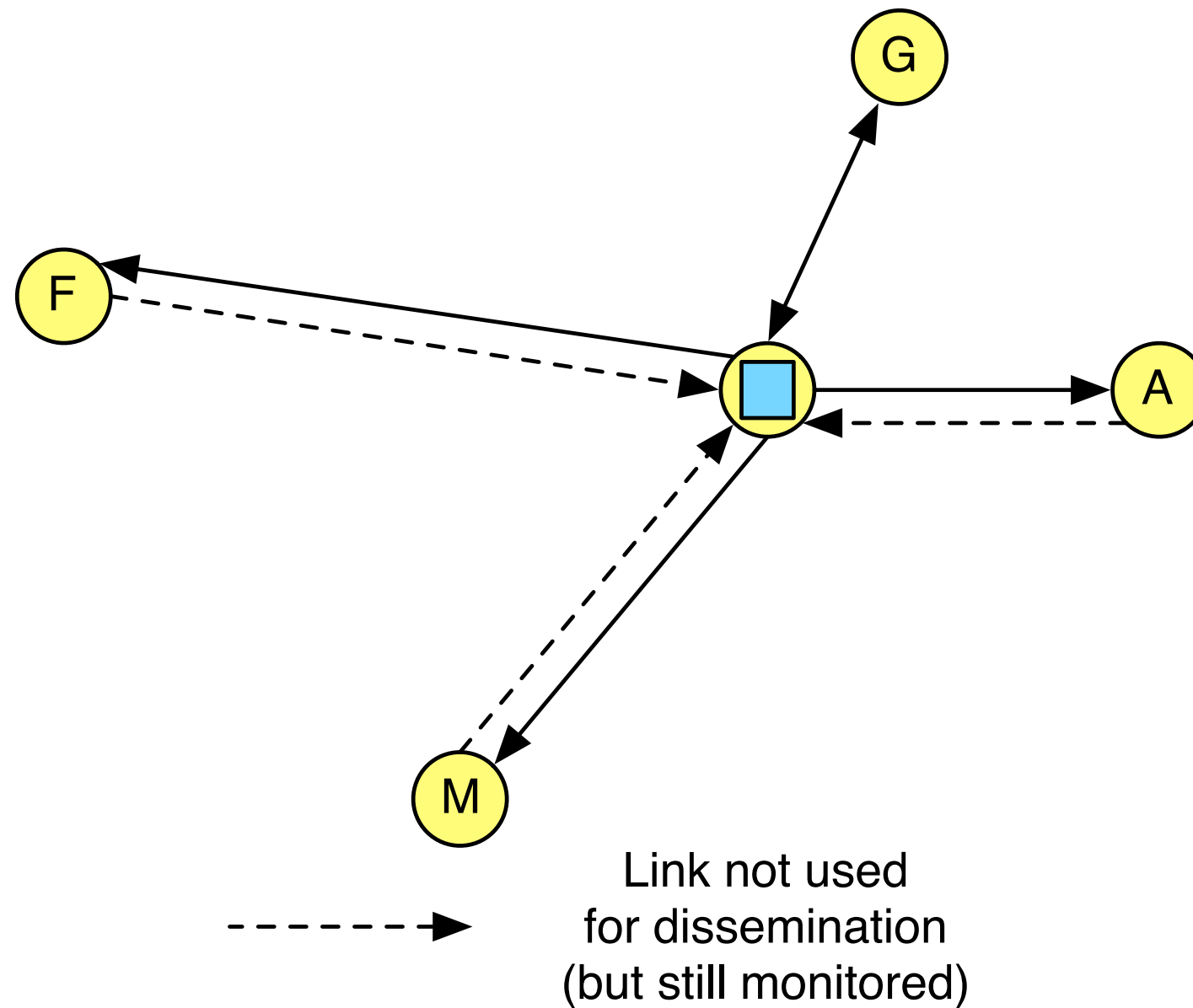
Emerging a tree through deactivations



Emerging a tree through deactivations



Emerging a tree through deactivations

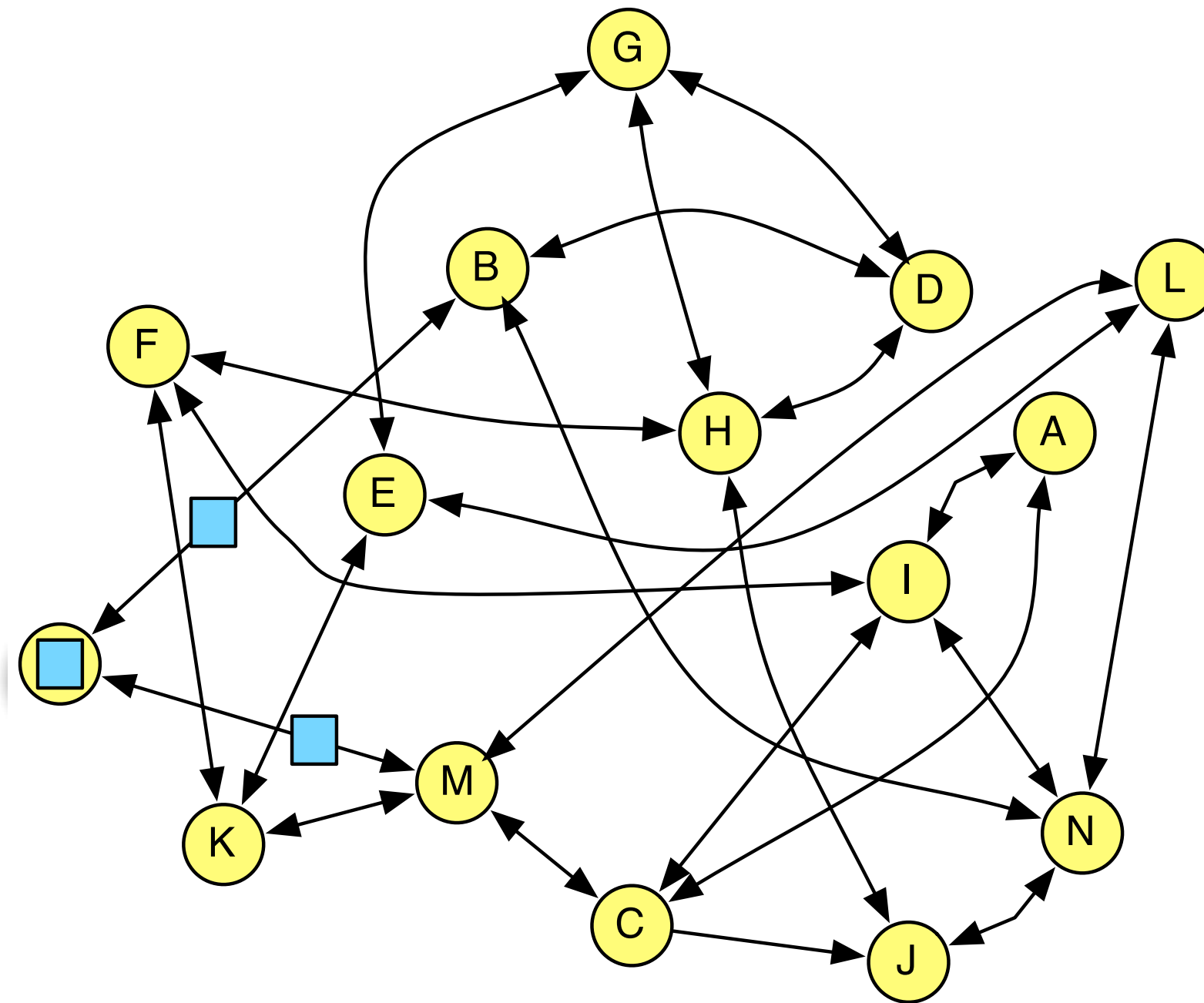


The graph consists of 14 nodes labeled A through L. Node A is a yellow circle with a blue square inside. All other nodes are yellow circles. The graph shows a complex network of directed edges between the nodes. The edges are as follows:

- A → B, A → C, A → D, A → E, A → F, A → G, A → H, A → I, A → J, A → K, A → L, A → M, A → N
- B → C, B → D, B → E, B → F, B → G, B → H, B → I, B → J, B → K, B → L, B → M, B → N
- C → D, C → E, C → F, C → G, C → H, C → I, C → J, C → K, C → L, C → M, C → N
- D → E, D → F, D → G, D → H, D → I, D → J, D → K, D → L, D → M, D → N
- E → F, E → G, E → H, E → I, E → J, E → K, E → L, E → M, E → N
- F → G, F → H, F → I, F → J, F → K, F → L, F → M, F → N
- G → H, G → I, G → J, G → K, G → L, G → M, G → N
- H → I, H → J, H → K, H → L, H → M, H → N
- I → J, I → K, I → L, I → M, I → N
- J → K, J → L, J → M, J → N
- K → L, K → M, K → N
- L → M, L → N
- M → N

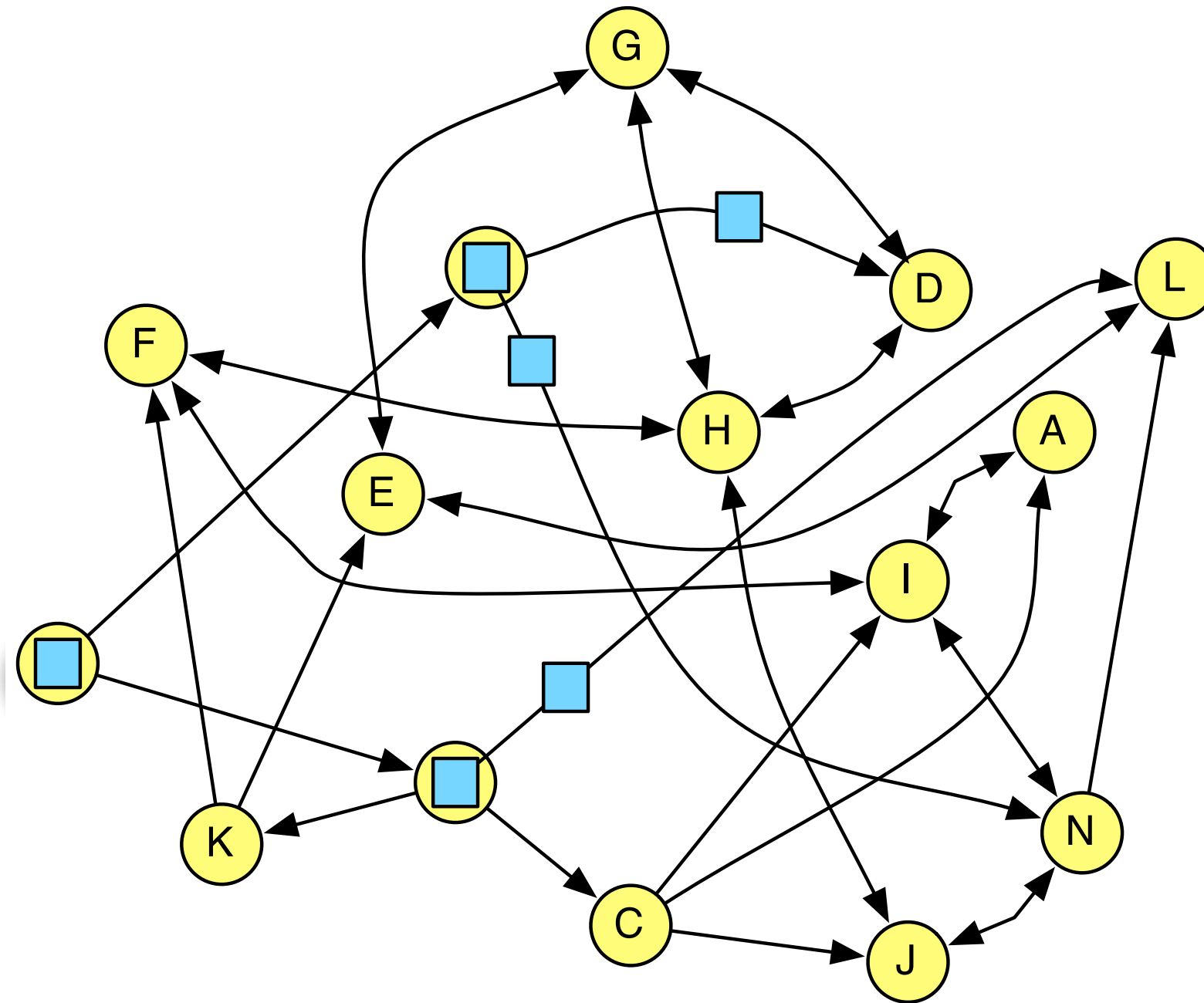
BRISA: Combining Efficiency and Reliability in Epidemic Data Dissemination

BRISA in action



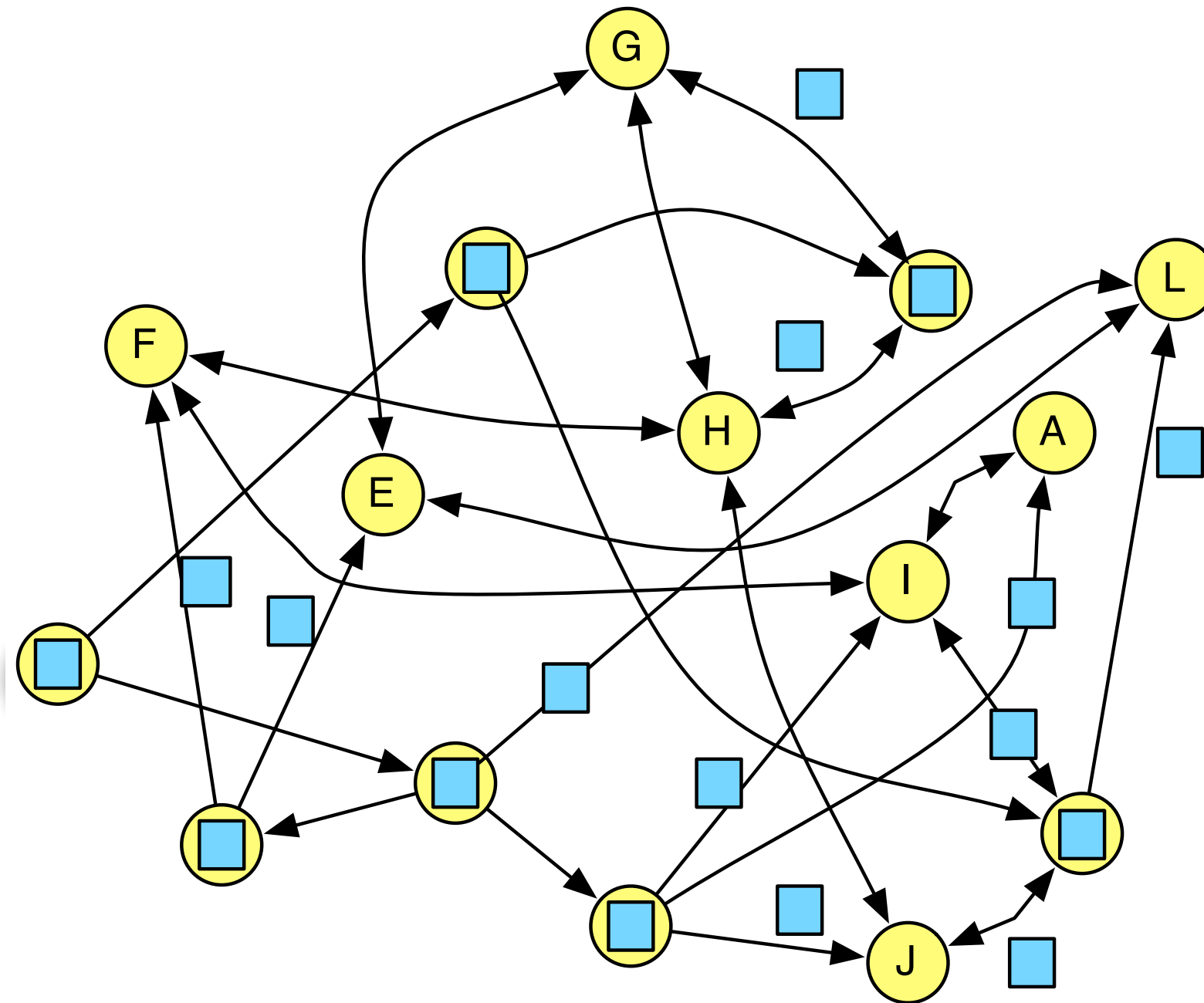
(Inactive Links not shown)

BRISA in action



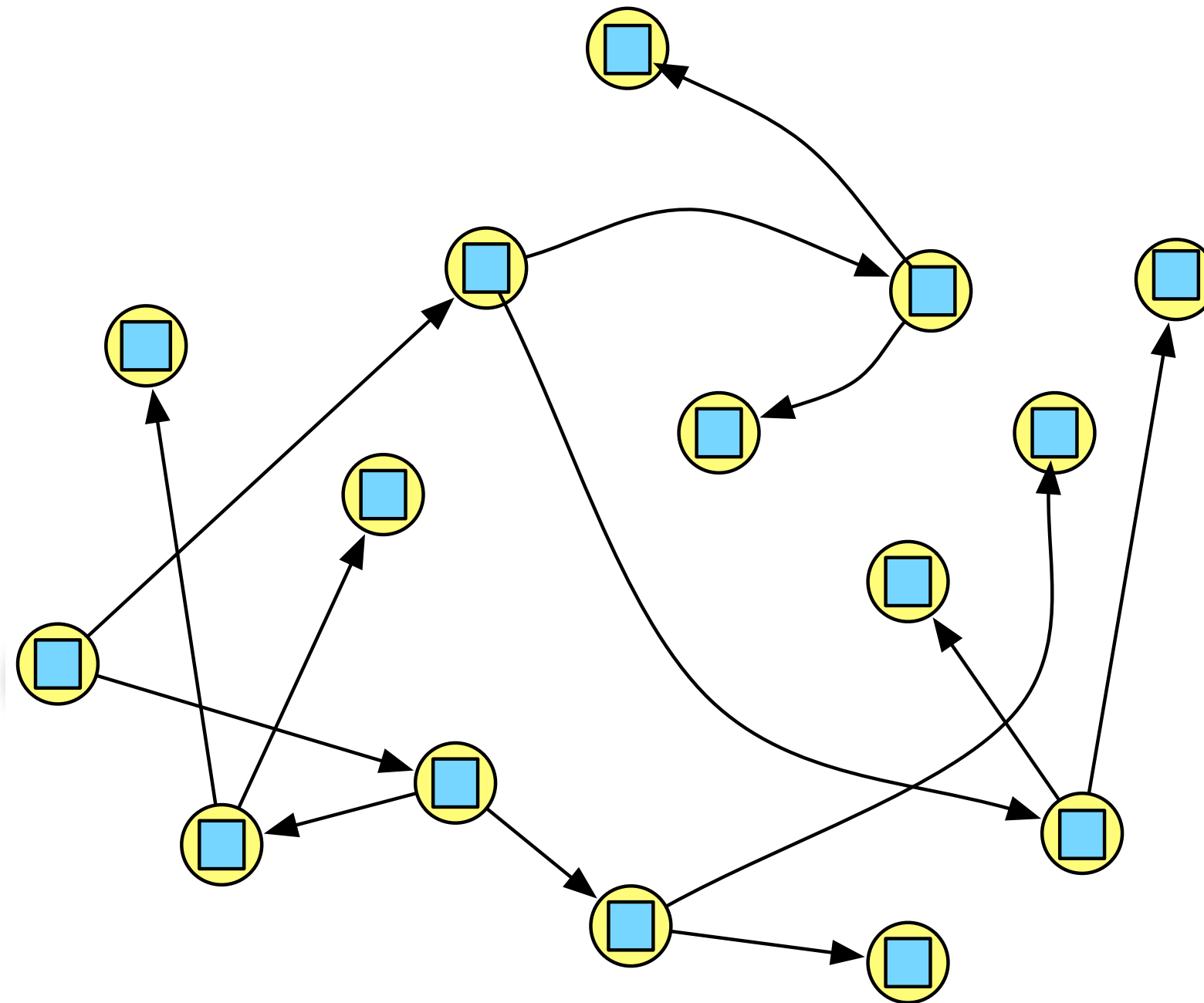
(Inactive Links not shown)

BRISA in action



(Inactive Links not shown)

BRISA in action

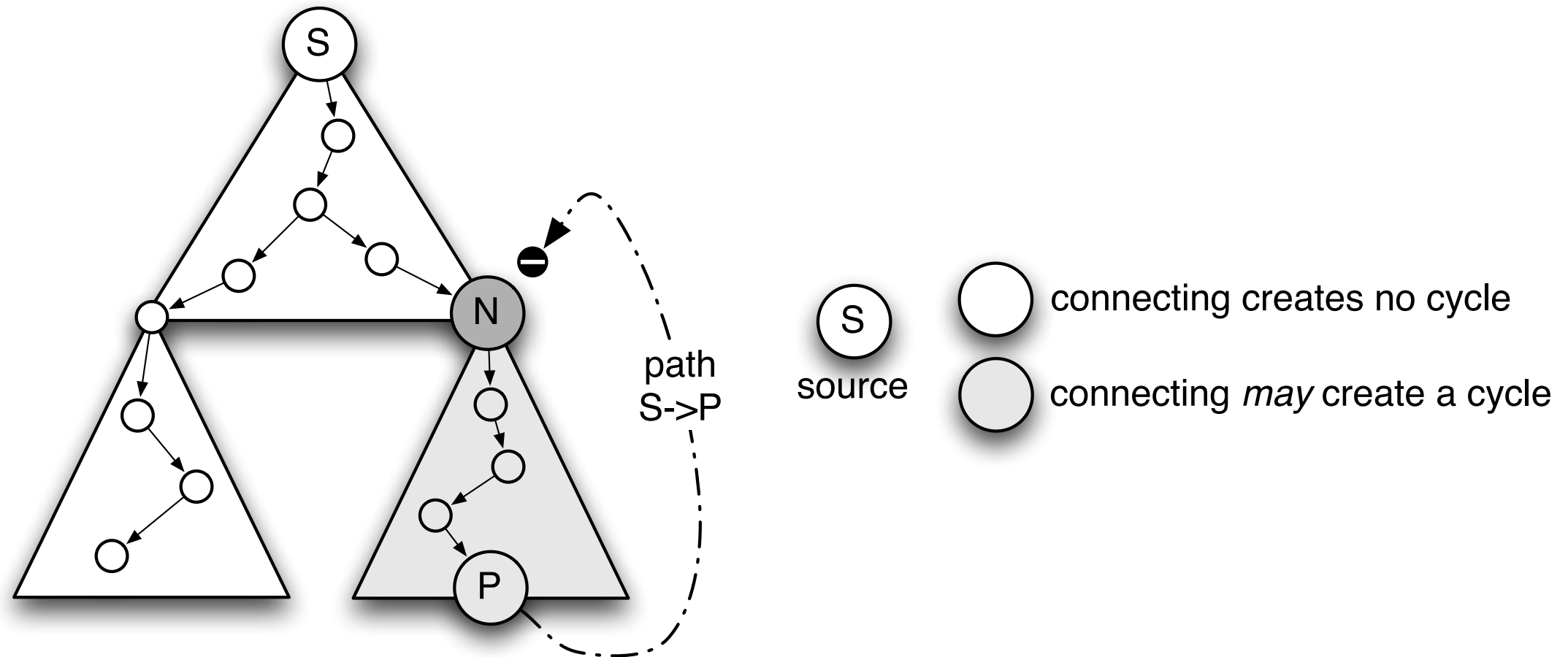


(Inactive Links not shown)

Failure Detection: RPSS level

- Parent failure implies **reactivating** a link
 - Send *activate* to one peer
 - Remember that connections are **persistent**
- Problem:
 - **Which peer** to reactivate ?
 - Avoid creating cycles
- **Cycle detection** mechanism
- Selection between eligible peers
 - **Performance** criteria
 - Criteria on tree **structure**

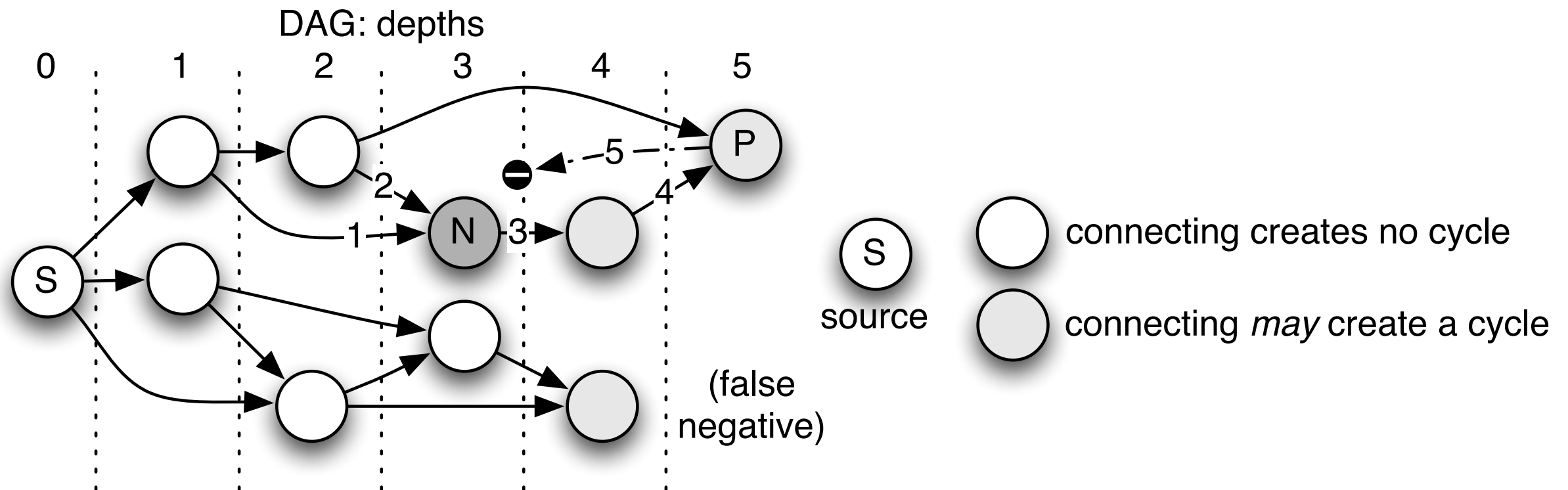
Avoiding cycles



Parent selection strategies

- Selection of peer to reactivate among **eligible** peers
- **Simplest**: first-come first picked
- **Performance** criteria
 - Delay-aware
- Criteria on tree **structure**
 - **Heterogeneity**-aware
 - Favors peer with most available bandwidth
 - **Load-balancing**
 - Seeks to reduce out-degree variance
 - **Gerontocratic**
 - Favors peers with highest uptime
 - Observations show that these peers are more stable

Extension to DAGs



- **DAG** = more than one parent
 - **Control** the number of received **duplicates**
 - Allows supporting **higher churn** levels
 - Lazy parent replacement
 - **Service continuity** if at least one parent remains
- Path-embedding is impossible for DAGs
 - Use node depth instead

Evaluation Setup

- Evaluation using a **prototype**
 - Supported by the **Splay** distributed system evaluation framework
- Two testbeds
 - 15 nodes **cluster** with Splay lightweight virtualization
 - 512 nodes
 - Uses Splay's churn replay mechanism
 - 128 nodes on **PlanetLab**

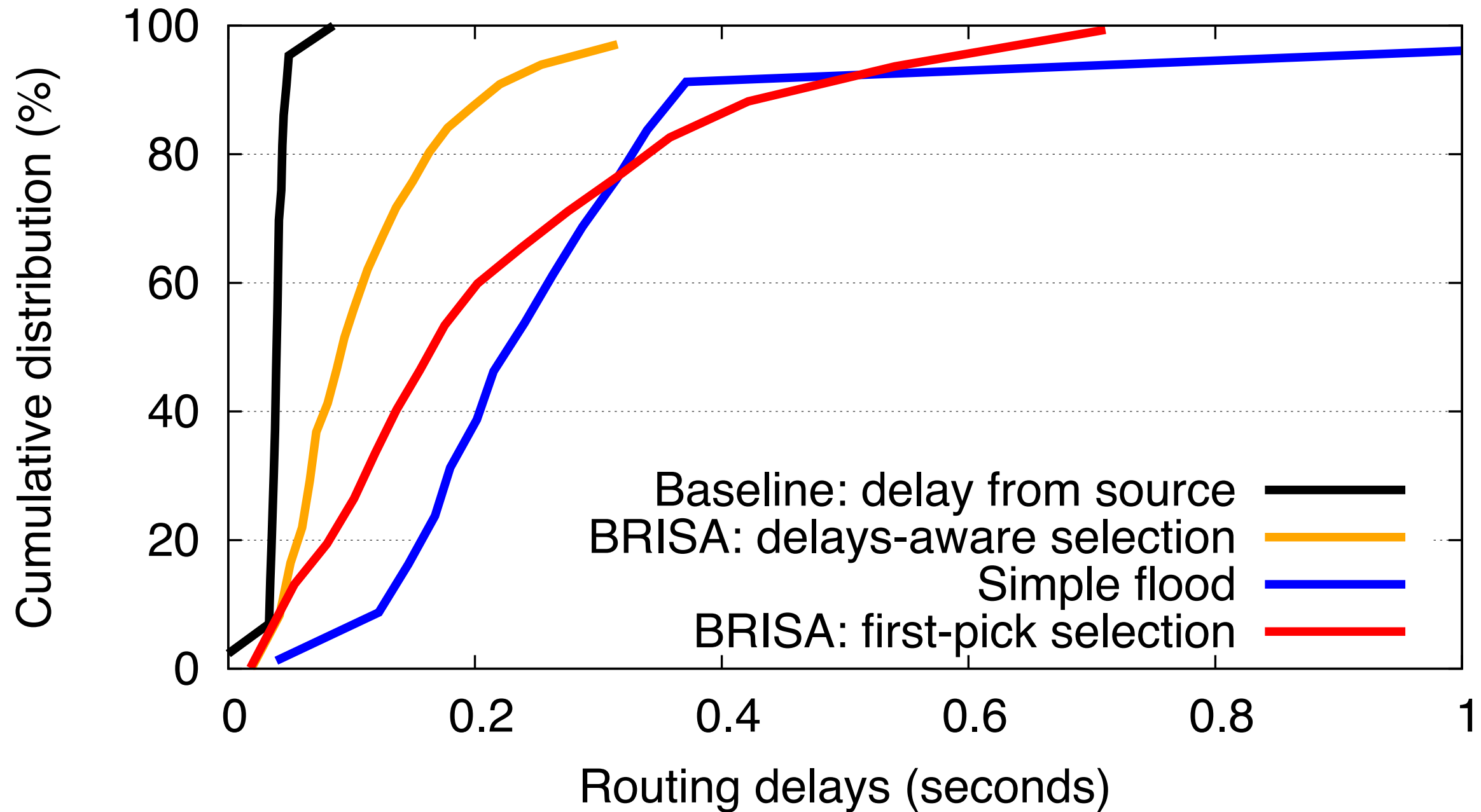
<http://www.splay-project.org/>

Evaluation Setup

- SimpleTree: statically constructed tree-based protocol
- SimpleGossip: traditional gossip-based protocol
- **TAG** [Liu & Zhou]
 - Closest related work
 - Another protocol that combines trees and epidemic overlays
 - Overlay used to build and repair the tree
 - No cycle detection: data is pulled from the overlay and from the tree
 - Tree repairing done by traversing the overlay

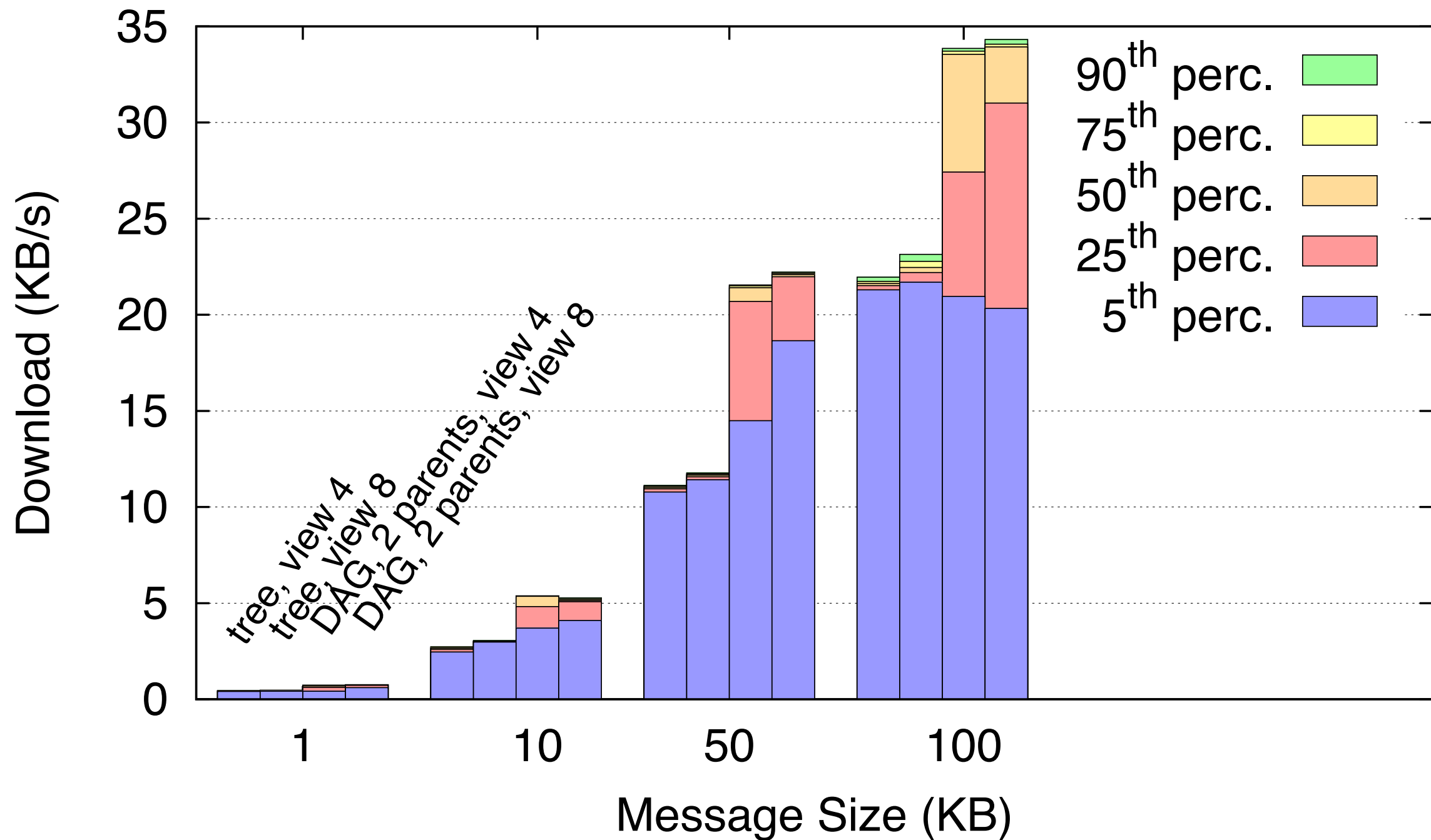
BRISA Dissemination Delay

150 PlanetLab nodes, 1KB



BRISA Bandwidth usage: download

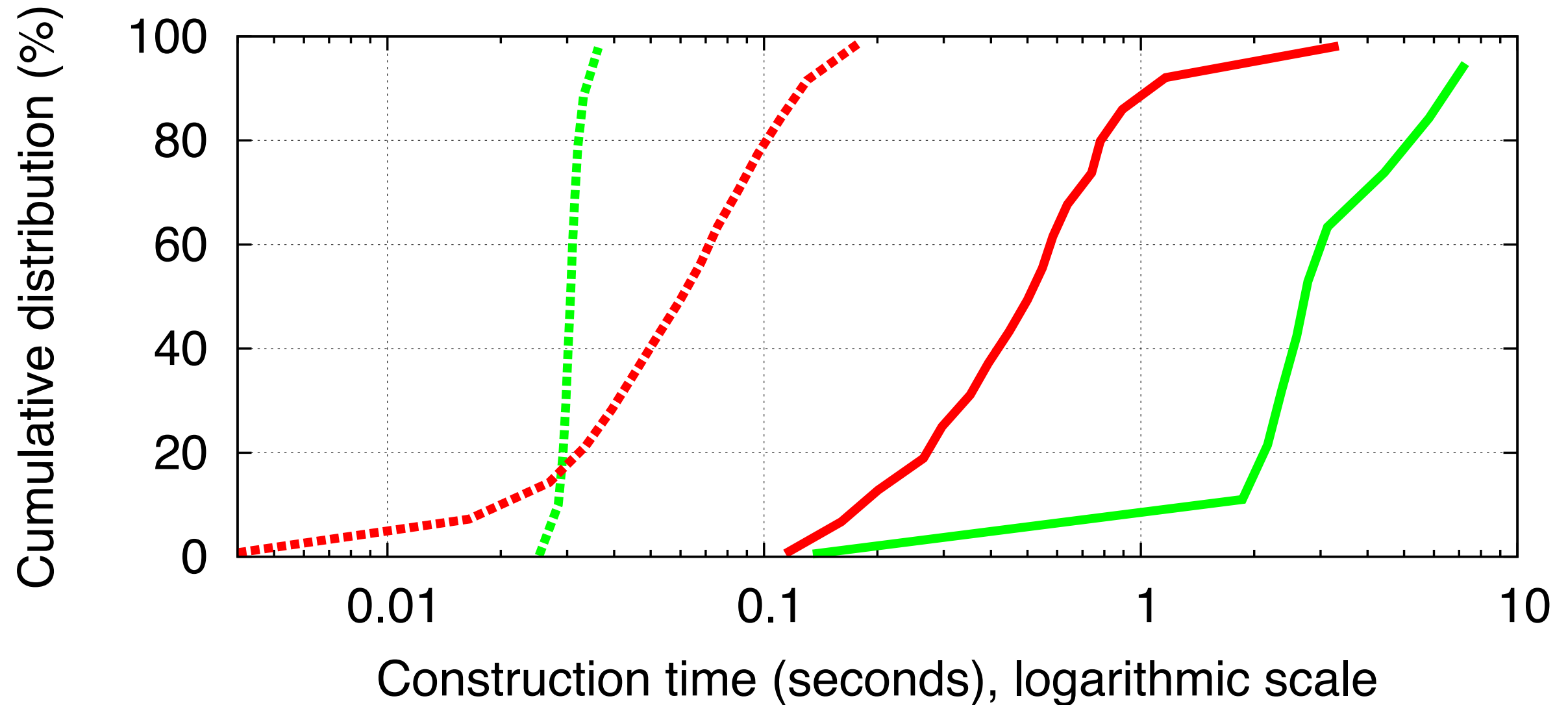
512 cluster nodes



Structure construction time

512 nodes

TAG (cluster) BRISA (PlanetLab) ———
BRISA (cluster) TAG (PlanetLab) ———



Dissemination Latency

512 cluster nodes, 500 x 1KB message, 5 msg / sec.

| <i>Protocol</i> | <i>Latency (seconds)</i> | <i>Overhead</i> |
|---------------------|--------------------------|-----------------|
| <i>SimpleTree</i> | 100.02 | - |
| <i>BRISA</i> | 106.59 | +6% |
| <i>SimpleGossip</i> | 128.23 | +28% |
| <i>TAG</i> | 200.48 | +100% |

Churn

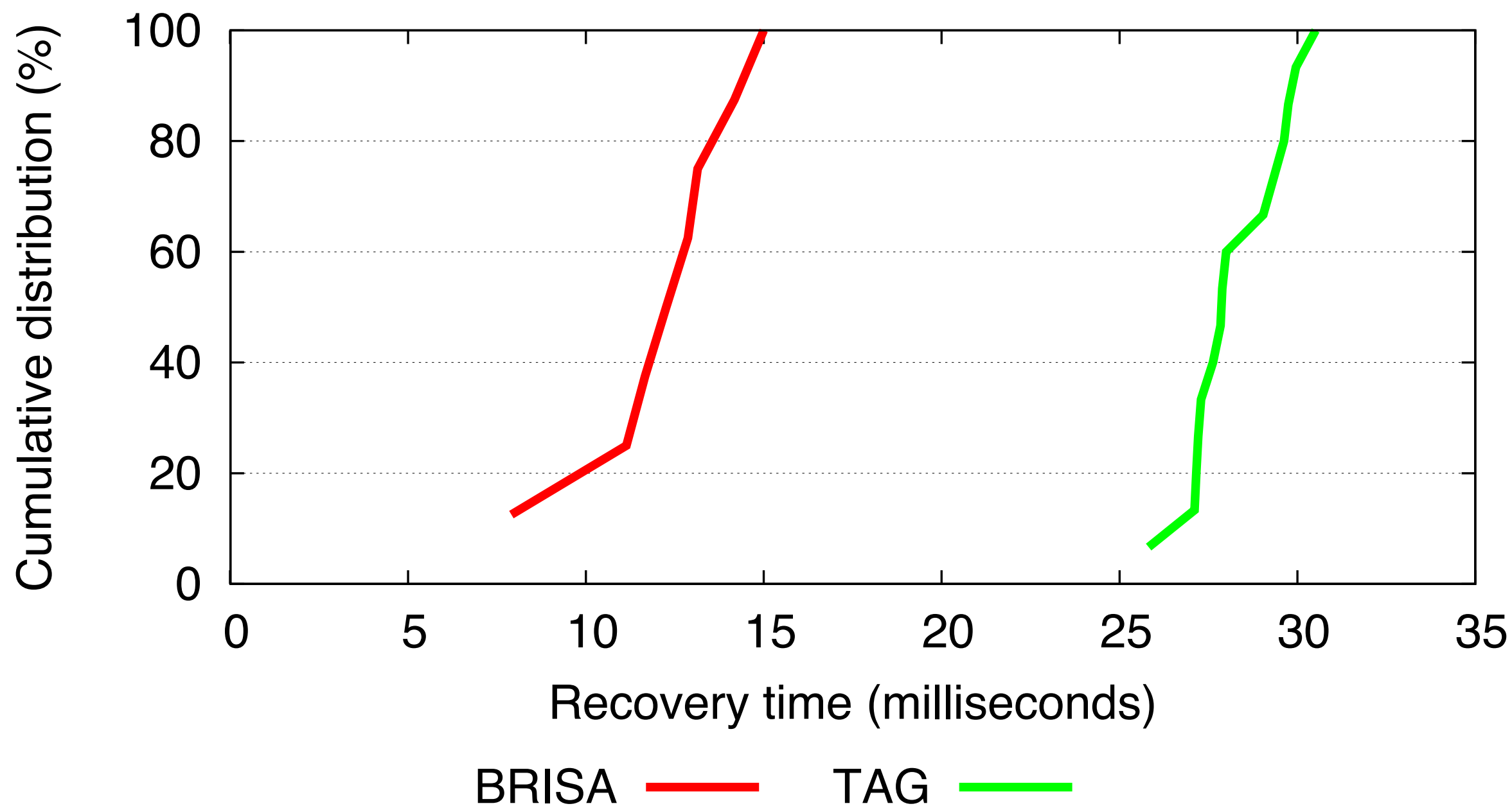
128 cluster nodes, 10 minutes

| | <i>Churn rate</i> | <i>Parents lost/min</i> | <i>Orphans/min</i> | <i>Direct fallback on RPSS link</i> | <i>Tree rejoin</i> |
|-----------------------|----------------------------------|-------------------------|--------------------|-------------------------------------|--------------------|
| Tree | 3% replacement per minute | 2.3 | 2.3 | 87% | 13% |
| DAG, 2 parents | | 4 | 0.2 | 92.5% | 7.5% |

(average per minute)

Parent recovery delay (tree rejoins)

128 cluster nodes; churn rate 3%



Conclusion & perspectives

- BRISA *combines* advantages of tree-based dissemination with advantages of epidemic dissemination
 - Robust to faults and churn
 - Quick failure recovery
 - Low overhead
 - Low latency
- Perspectives
 - Use DAGs as base for multiple trees
 - Low latency and overhead makes BRISA suitable for video streaming

BRISA:

Combining Efficiency and Reliability in Epidemic Data Dissemination

Miguel Matos^{*}, Valerio Schiavoni⁺, Pascal Felber⁺
Rui Oliveira^{*} and Etienne Rivière⁺

^{*}INESC TEC & University of Minho, Portugal
⁺Université de Neuchâtel, Switzerland

miguelmatos@di.uminho.pt

International Parallel & Distributed Processing Symposium
24 May, 2012

Backup slides

Future Perspectives

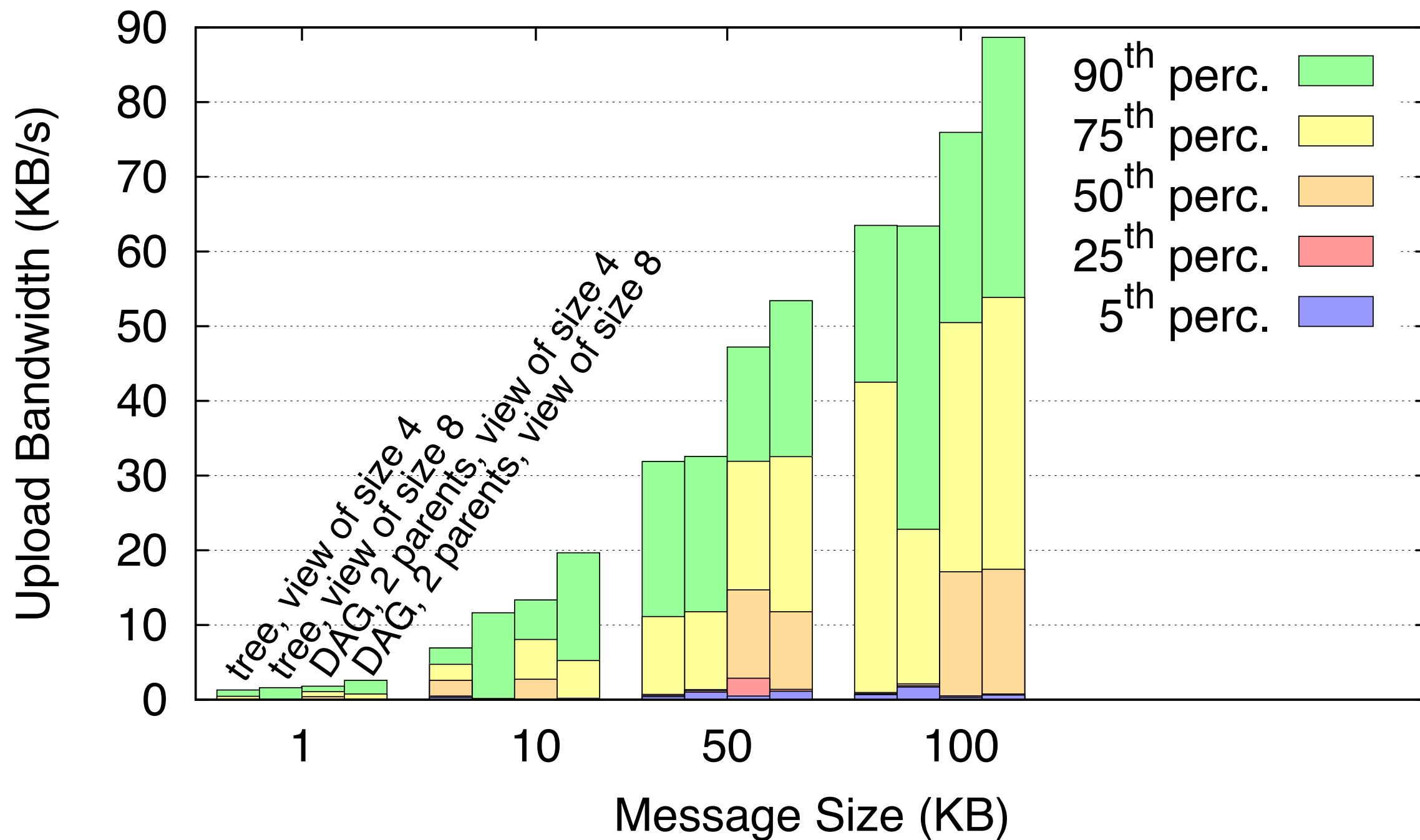
- More elaborate parent selection strategies
- Use DAGs as base for multiple trees
- Low latency and overhead makes it suitable to video streaming

Failure Detection

- At the *Reactive Peer Sampling Service* (RPSS) level
 - Constant neighbors **monitoring**
 - Takes place on persistent connection
 - Also for links not used for dissemination
 - Classical failure detector
 - “**Informed**” failure detector
 - Use **hint** from dissemination layer
 - Use deactivated links to send notifications of received message *ids*
 - Reception of a *id* for an (yet) unknown message
=> triggers instantaneous monitoring of parent
 - Allows reducing monitoring frequency while keeping **reactiveness**
- ➡ Detection of failed parent triggers link **reactivation**

BRISA Bandwidth usage: upload

512 cluster nodes



Goals

- Efficient: management overhead is low with respect to application data; low latency
- Robust: service continuity even under faults and churn
- Scalable: able to handle system growth sublinearly