

Algorithms and Distributed Systems 2019/2020 (Lab Two)

**MIEI - Integrated Master in Computer Science and
Informatics**

Specialization block

João Leitão (jc.leitao@fct.unl.pt)



**FACULDADE DE
CIÊNCIAS E TECNOLOGIA
UNIVERSIDADE NOVA DE LISBOA**

Class structure:

- One Unstructured Overlay: HyParView.
- One Partially Unstructured Overlay: Plumtree.

Class structure:

- **One Unstructured Overlay: HyParView.**
- One Partially Unstructured Overlay: Plumtree.

Unstructured Overlay Management

Reactive strategy:

Partial views are updated only as a reaction to some external event.

e.g. Scamp.

Cyclic strategy:

Partial views are updated as a result of some periodically operation.

e.g. Cyclon.

HyParView

HyParView: a membership protocol for reliable gossip-based broadcast*

João Leitão
University of Lisbon
jleitao@lasige.di.fc.ul.pt

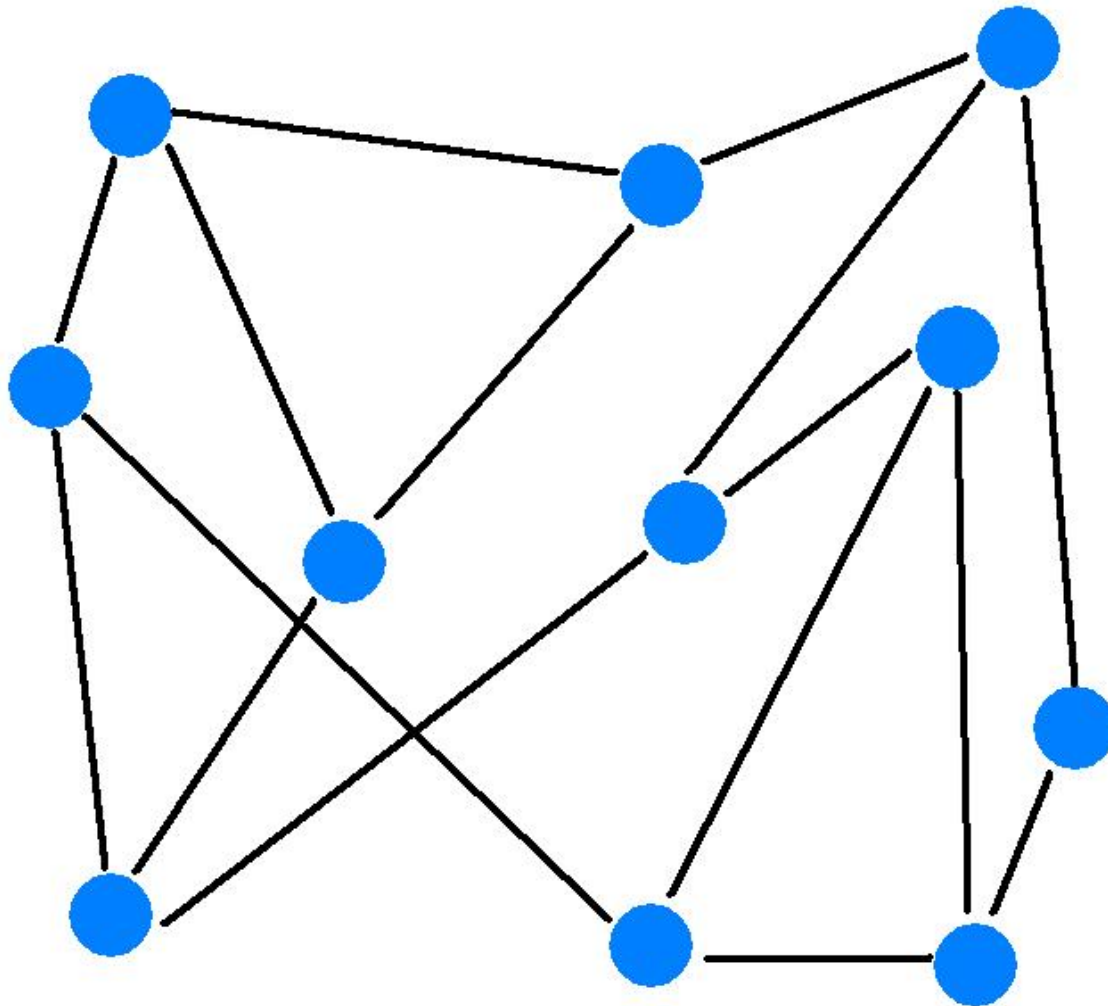
José Pereira
University of Minho
jop@di.uminho.pt

Luís Rodrigues
University of Lisbon
ler@di.fc.ul.pt

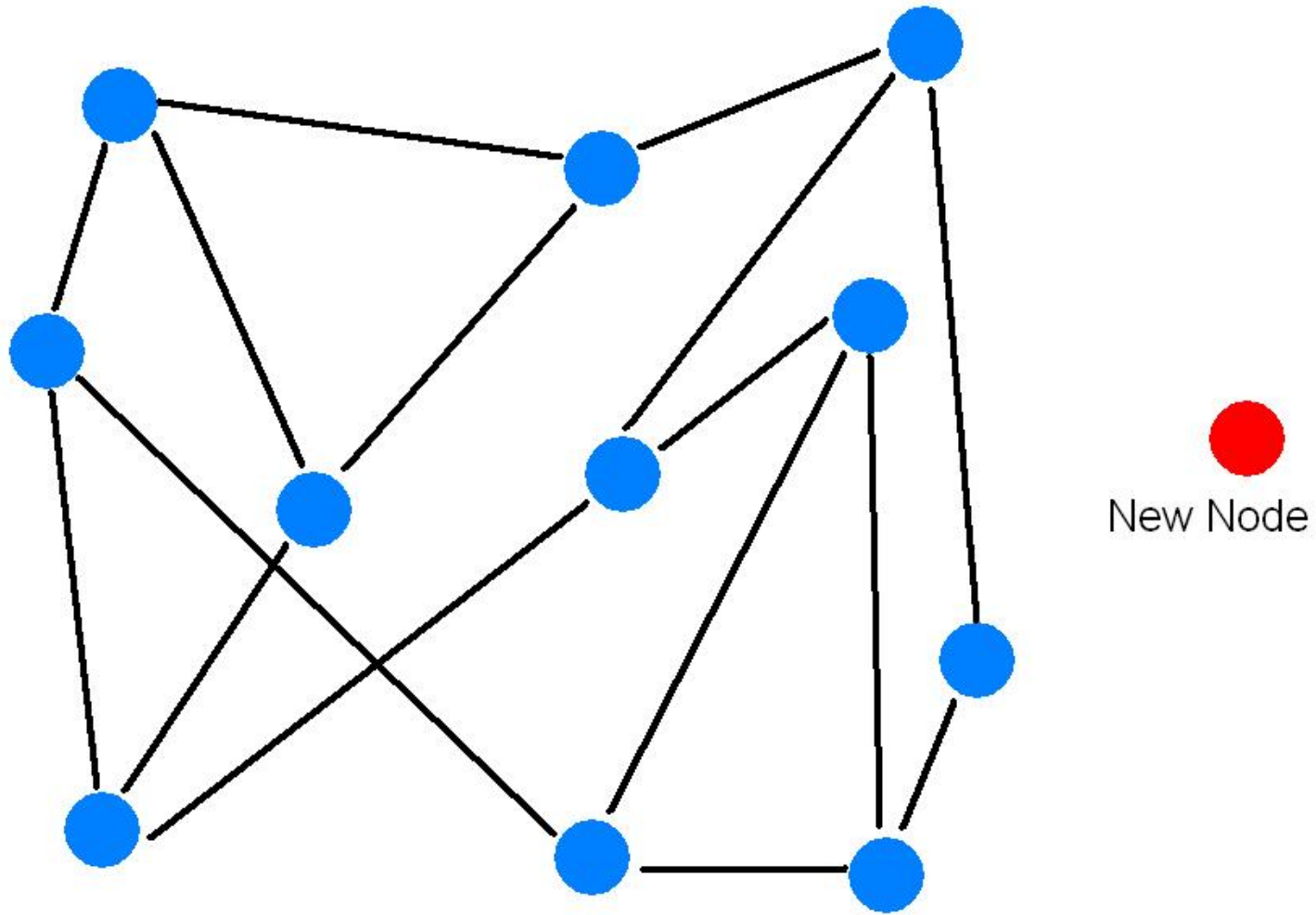
HyParView

- Name stands for **Hybrid Partial View**
 - It maintains two distinct partial views.
 - Active view:
 - Small sized ($fanout+1$).
 - Symmetric.
 - TCP connections are maintained with nodes in this views.
 - Used to disseminate messages.
 - Passive view:
 - Larger size in order of: $k \times \#(\text{active View})$
 - Used for fault tolerance.

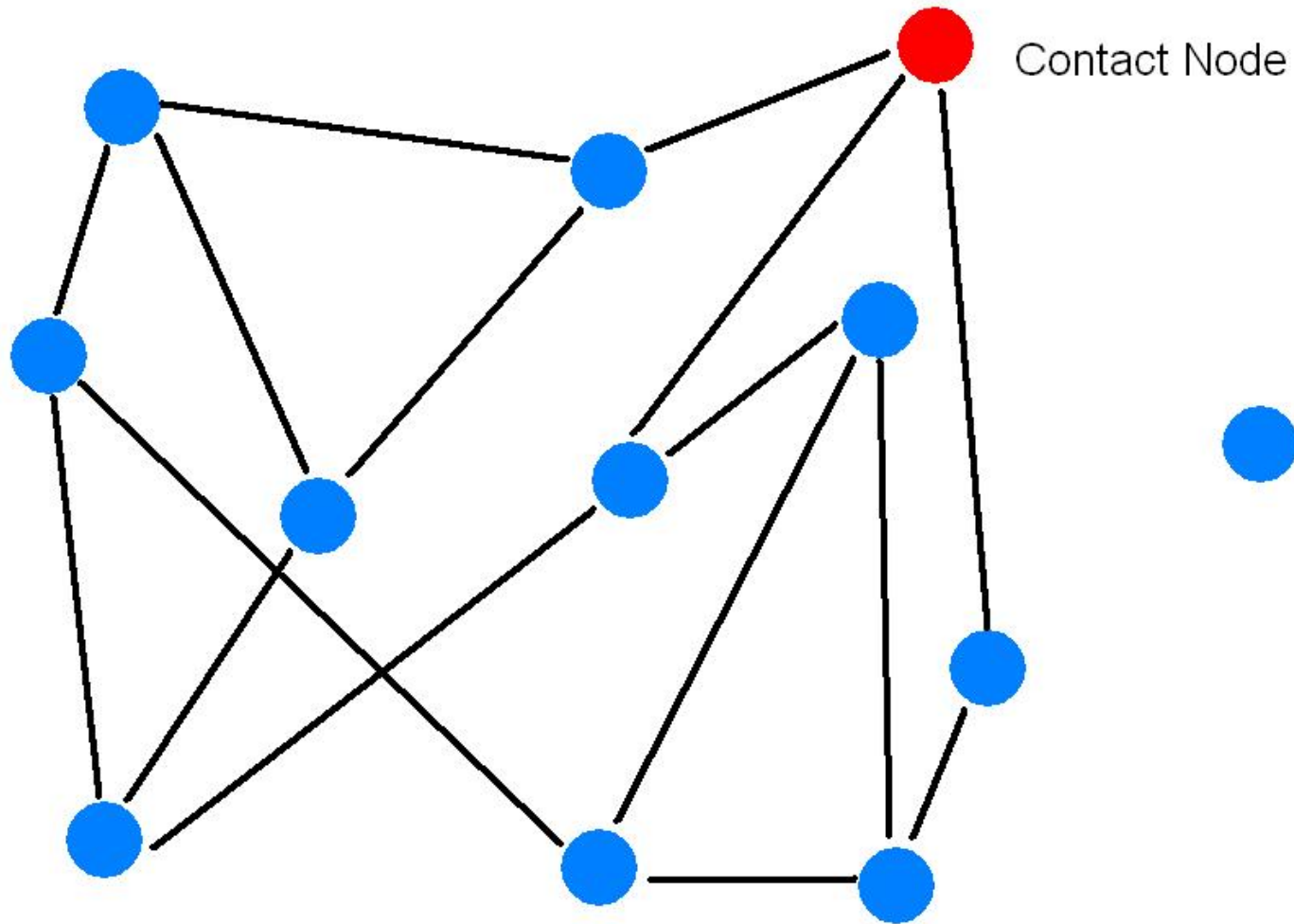
HyParView Join Mechanism



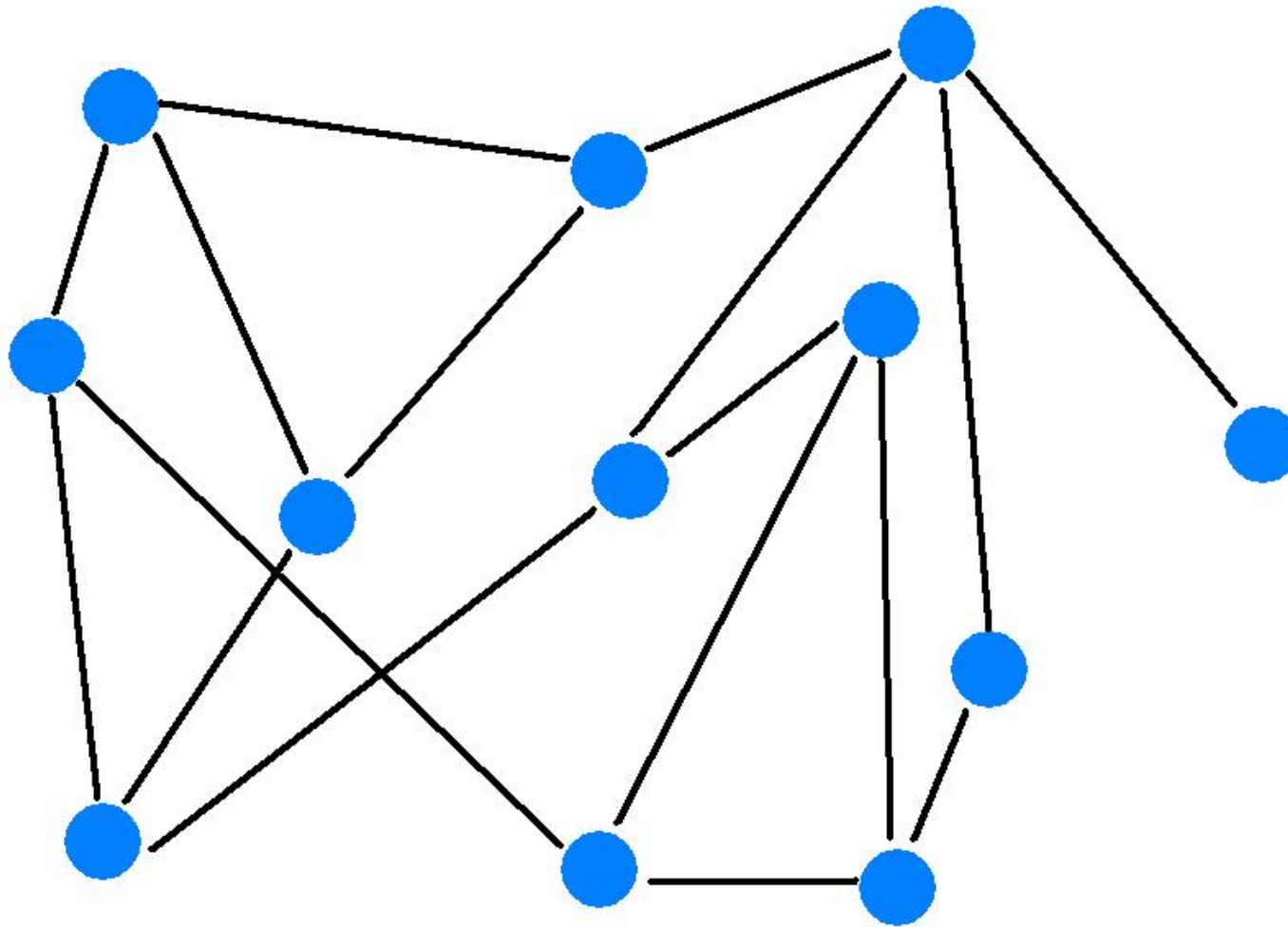
HyParView Join Mechanism



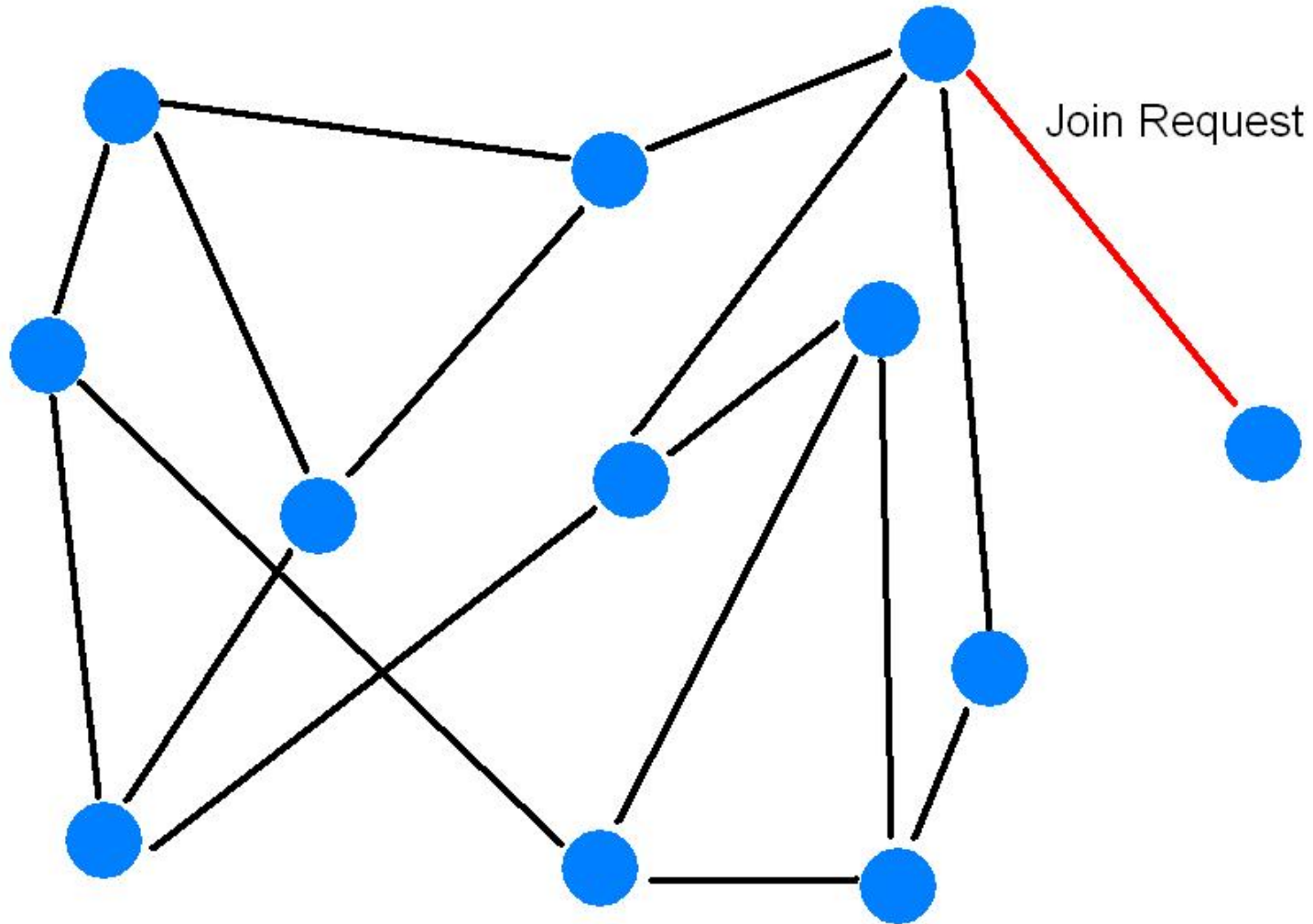
HyParView Join Mechanism



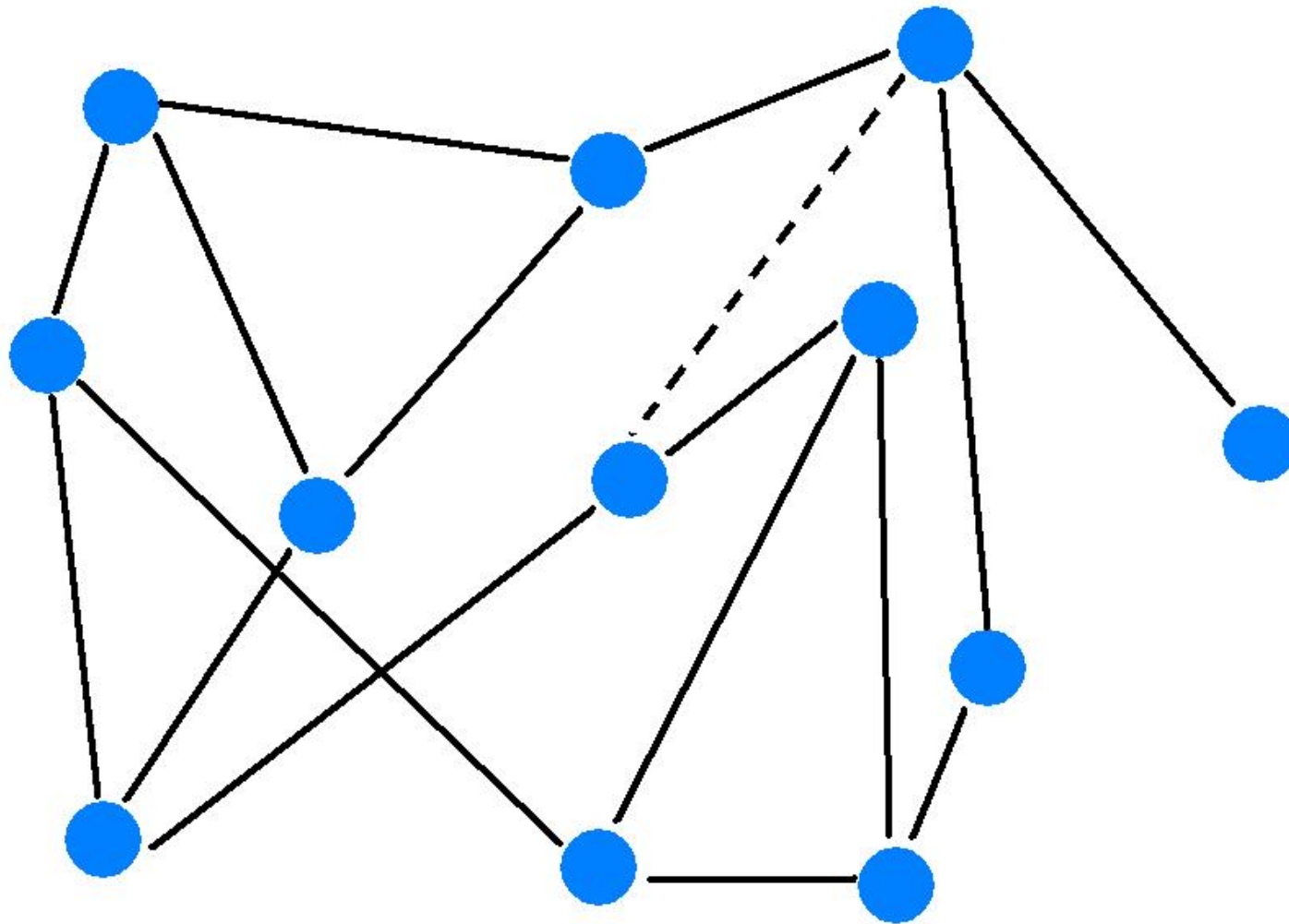
HyParView Join Mechanism



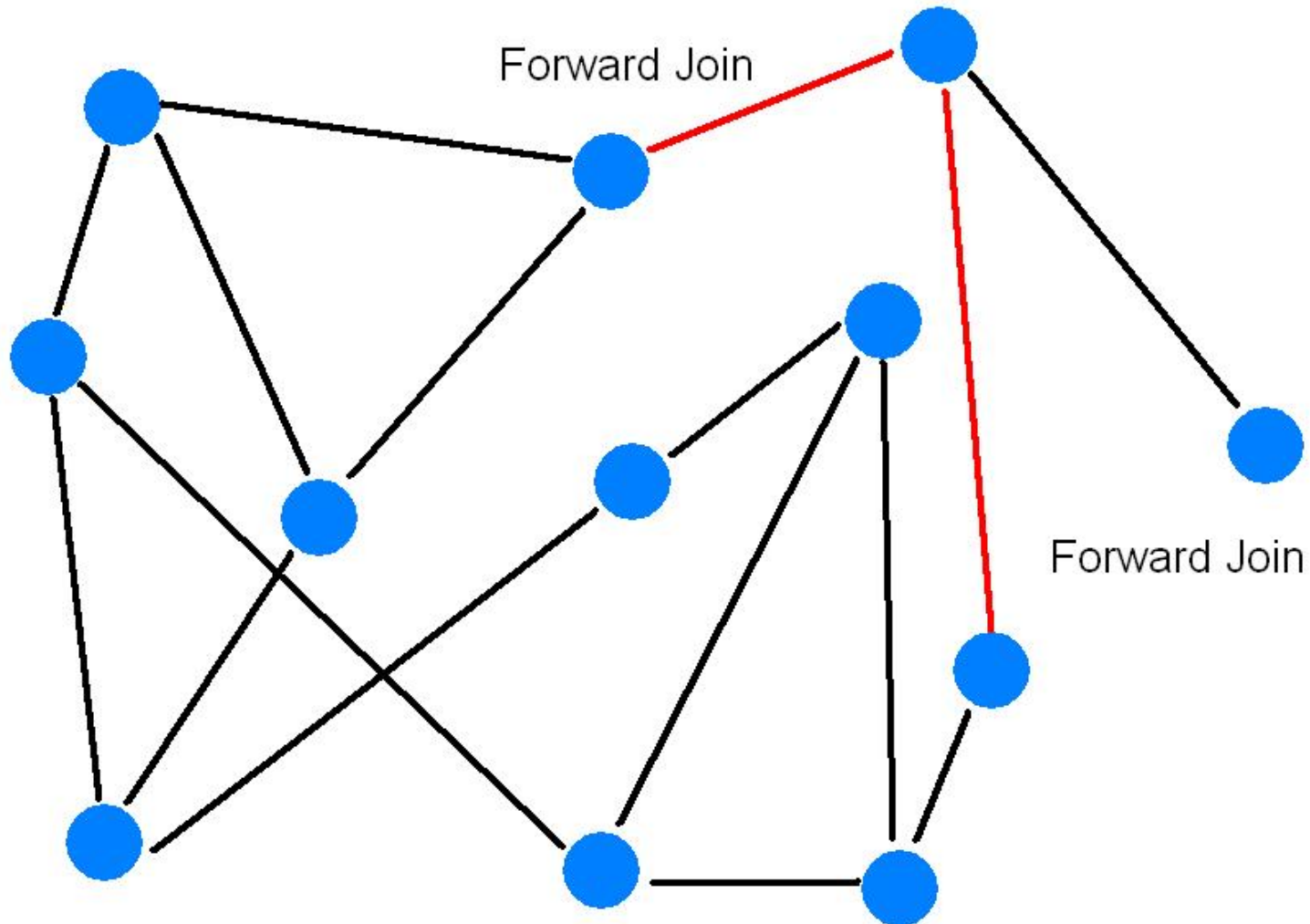
HyParView Join Mechanism



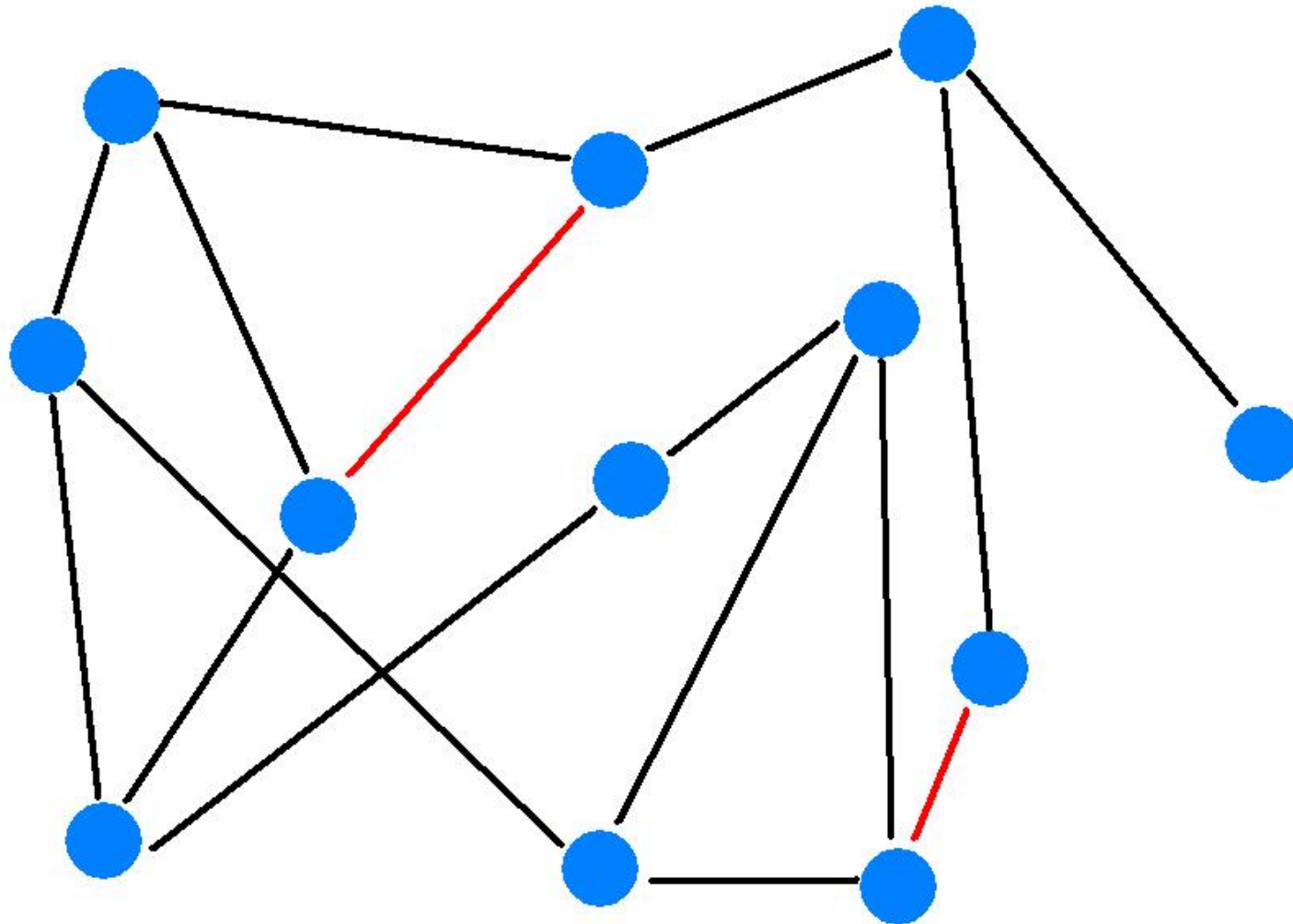
HyParView Join Mechanism



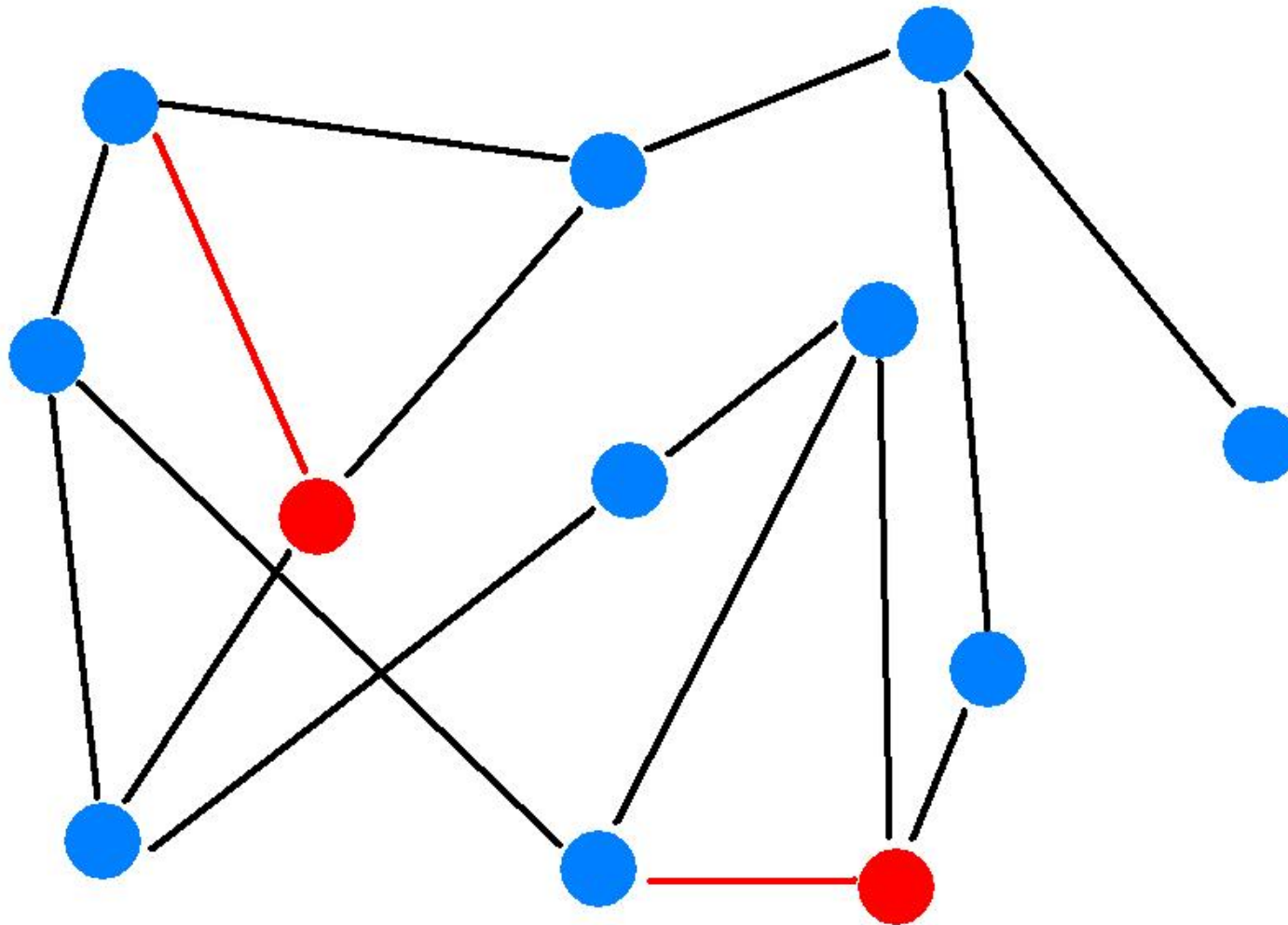
HyParView Join Mechanism



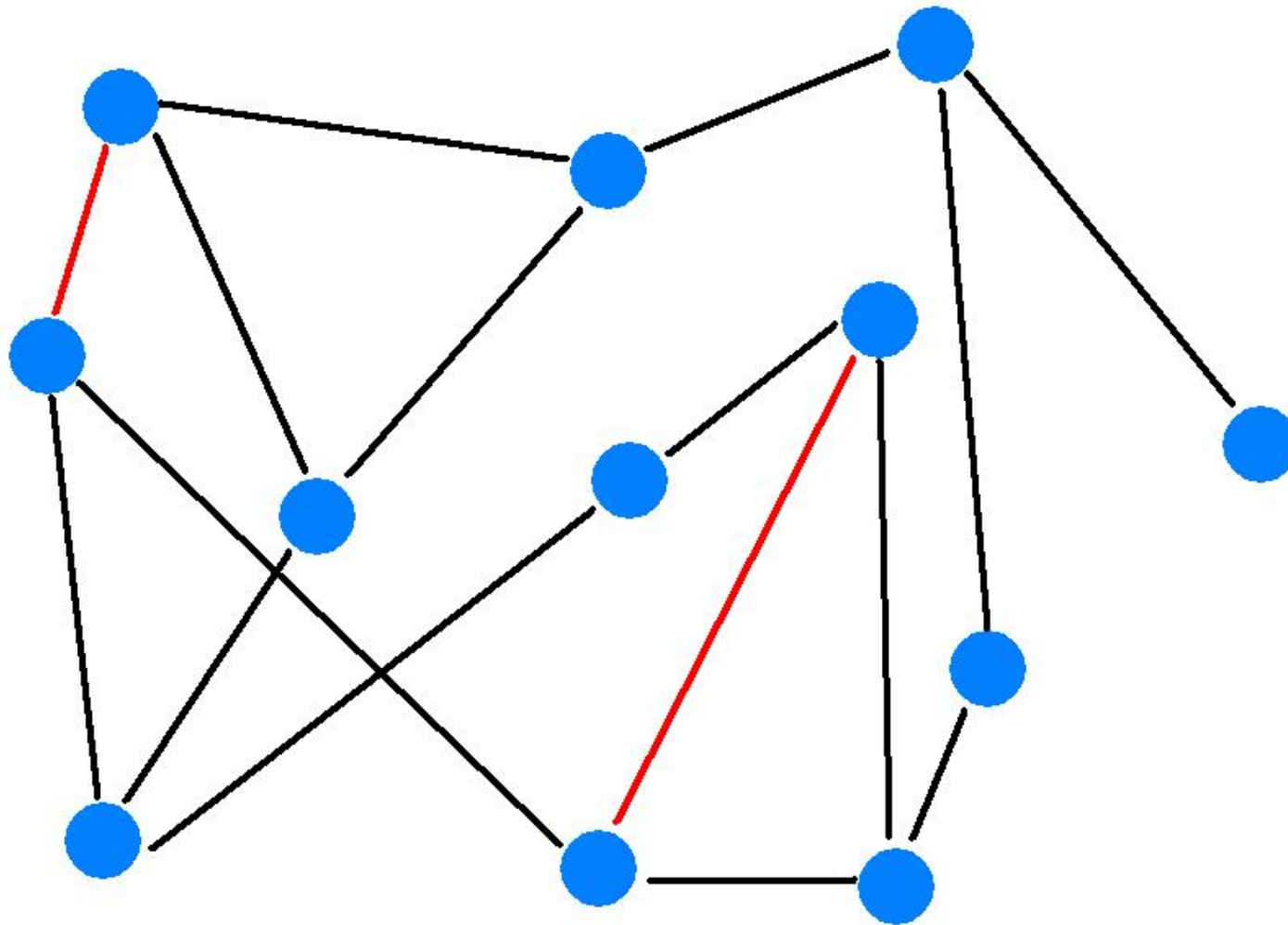
HyParView Join Mechanism



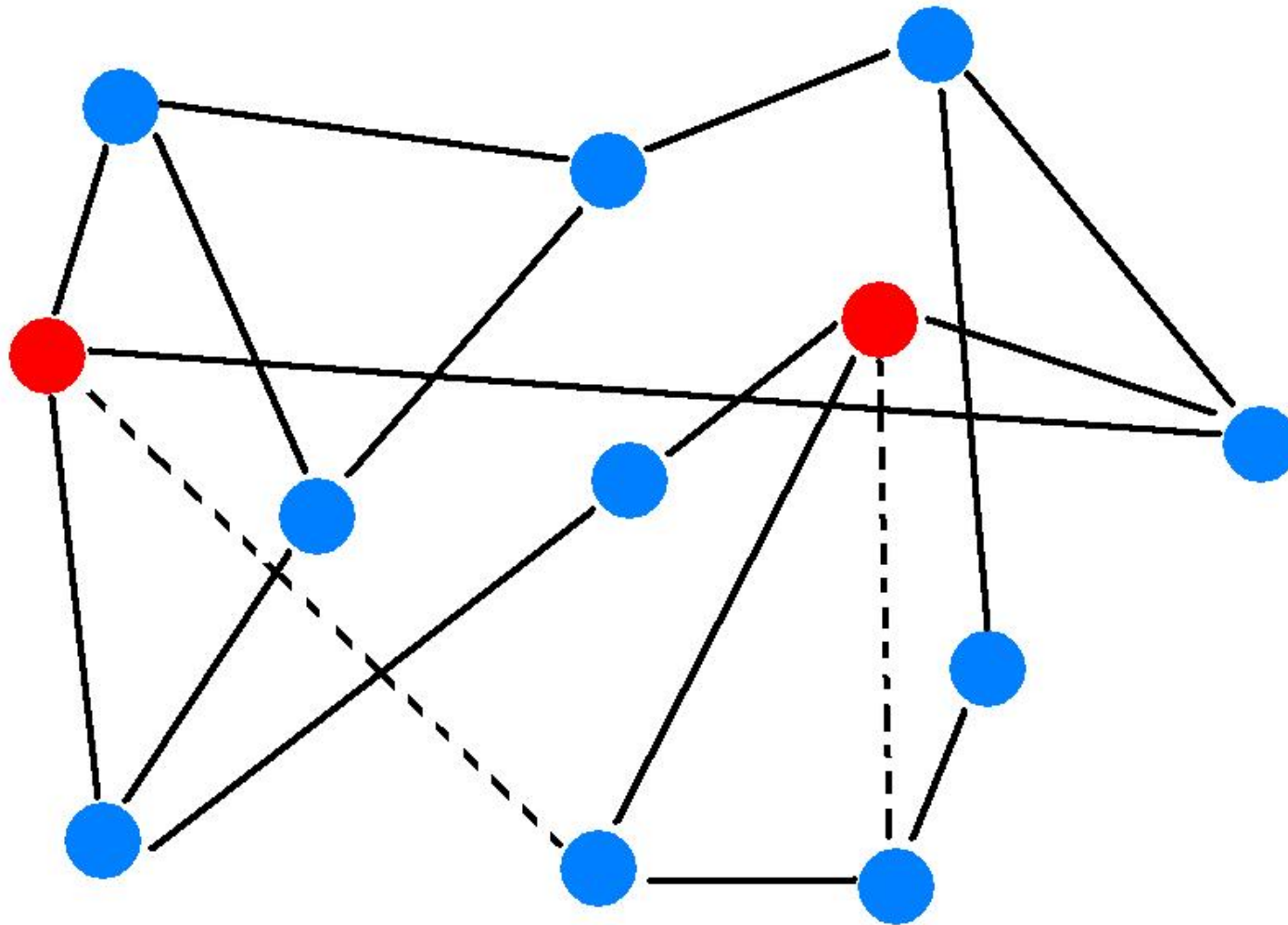
HyParView Join Mechanism



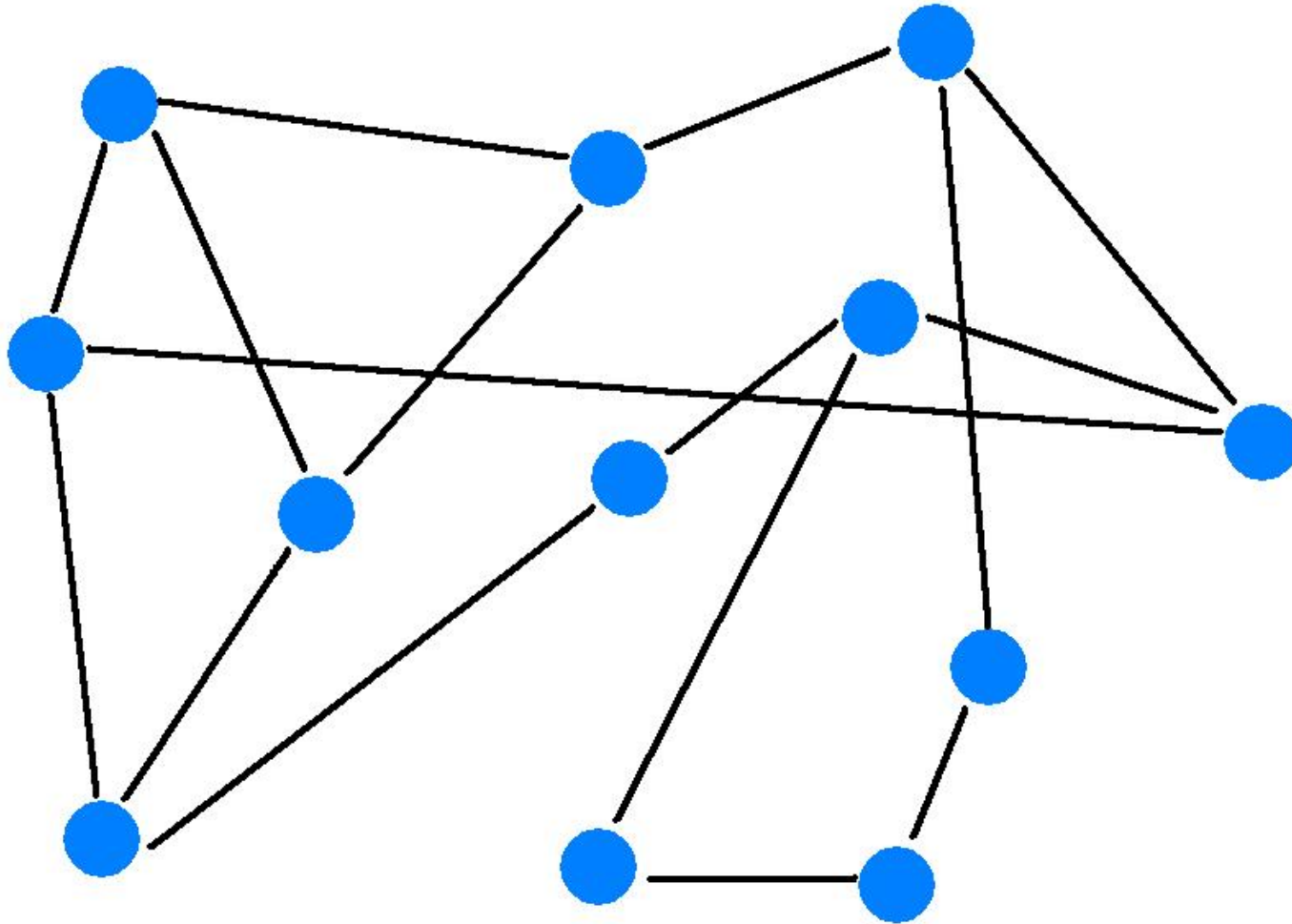
HyParView Join Mechanism



HyParView Join Mechanism



HyParView Join Mechanism



HyParView: Active View

- Maintained through a reactive strategy.
- TCP is used as an unreliable failure detector.
- When a node suspects a neighbor has failed:
 - The suspected neighbor is removed from the active view.
 - Random nodes from the passive view are contacted in order to promote them to the active view.

HyParView: Passive View

- Maintained through a cyclic strategy.
- Periodically each node starts a shuffle process with a random node.
- They exchange messages with samples from their partial views:
 - That node that starts the process send both himself and some elements from his active view.
 - His peer only send elements from his passive view.
- Both nodes update their passive views with the information received in the process.

HyParView: The Algorithm

Algorithm 1: Join mechanism


Data:

myself: the identifier of the local node
activeView: a node active partial view
passiveView: a node passive view
contactNode: a node already present in the overlay
newNode: the node joining the overlay
ARWL: Active random walk length
PRWL: Passive random walk length

```
1  upon init do
2      Send(JOIN, contactNode, myself);

3  upon Receive(JOIN, newNode) do
4      call addNodeActiveView(newNode)
5      foreach  $n \in \text{activeView}$  and  $n \neq \text{newNode}$  do
6          Send(FORWARDJOIN, n, newNode, ARWL, myself)

7  upon Receive(FORWARDJOIN, newNode, timeToLive, sender) do
8      if timeToLive==0 || #activeView==1 then
9          call addNodeActiveView(newNode)
10     else
11         if timeToLive==PRWL then
12             call addNodePassiveView(newNode)
13          $n \leftarrow n \in \text{activeView}$  and  $n \neq \text{sender}$ 
14         Send(FORWARDJOIN, n, newNode, timeToLive-1, myself)
```



HyParView: The Algorithm

Algorithm 2: View manipulation primitives

Data:

activeView: a node active partial view

passiveView: a node passive view

```
1  procedure dropRandomElementFromActiveView do
2       $n \leftarrow n \in \text{activeView}$ 
3      Send(DISCONNECT,  $n$ , myself)
4       $\text{activeView} \leftarrow \text{activeView} \setminus \{n\}$ 
5       $\text{passiveView} \leftarrow \text{passiveView} \cup \{n\}$ 

6  procedure addNodeActiveView(node) do
7      if node  $\neq$  myself and node  $\notin$  activeView then
8          if isfull(activeView) then
9              call dropRandomElementFromActiveView
10          $\text{activeView} \leftarrow \text{activeView} \cup \text{node}$ 

11 procedure addNodePassiveView(node) do
12     if node  $\neq$  myself and node  $\notin$  activeView and node  $\notin$  passiveView then
13         if isfull(passiveView) then
14              $n \leftarrow n \in \text{passiveView}$ 
15              $\text{passiveView} \leftarrow \text{passiveView} \setminus \{n\}$ 
16          $\text{passiveView} \leftarrow \text{passiveView} \cup \text{node}$ 

17 upon Receive(DISCONNECT, peer) do
18     if peer  $\in$  activeView then
19          $\text{activeView} \leftarrow \text{activeView} \setminus \{\text{peer}\}$ 
20         call addNodePassiveView(peer)
```

HyParView: The Algorithm

Algorithm 2: View manipulation primitives

Data:

activeView: a node active partial view

passiveView: a node passive view

```
1  procedure dropRandomElementFromActiveView do
2       $n \leftarrow n \in \text{activeView}$ 
3      Send(DISCONNECT,  $n$ , myself)
4       $\text{activeView} \leftarrow \text{activeView} \setminus \{n\}$ 
5       $\text{passiveView} \leftarrow \text{passiveView} \cup \{n\}$ 

6  procedure addNodeActiveView(node) do
7      if node  $\neq$  myself and node  $\notin$  activeView then
8          if isfull(activeView) then
9              call dropRandomElementFromActiveView
10          $\text{activeView} \leftarrow \text{activeView} \cup \text{node}$ 

11 procedure addNodePassiveView(node) do
12     if node  $\neq$  myself and node  $\notin$  activeView and node  $\notin$  passiveView then
13         if isfull(passiveView) then
14              $n \leftarrow n \in \text{passiveView}$ 
15              $\text{passiveView} \leftarrow \text{passiveView} \setminus \{n\}$ 
16          $\text{passiveView} \leftarrow \text{passiveView} \cup \text{node}$ 

17 upon Receive(DISCONNECT, peer) do
18     if peer  $\in$  activeView then
19          $\text{activeView} \leftarrow \text{activeView} \setminus \{\text{peer}\}$ 
20         call addNodePassiveView(peer)
```

HyParView: Message Dissemination

- The size of the active view is $fanout+1$.
- Whenever a node receives a message for the first time:
 - It forwards it to all nodes in the active view.
 - Except the one from which he received the message.
- Message dissemination is performed:
 - Deterministically in a random overlay...
 - ...by flooding an overlay with small degree.
 - Each node tests all its active view each time it forwards a message.

HyParView: Message Dissemination

Algorithm 3: Eager push protocol

Data:

myself: the identifier of the local node

receivedMsgs: a list of received messages identifiers

f: the fanout value

```
1  upon event Broadcast( $m$ ) do
2       $mID \leftarrow \text{hash}(m + \text{myself})$ 
3       $\text{peerList} \leftarrow \text{getPeer}(f, \text{null})$ 
4      foreach  $p \in \text{peerList}$  do
5          trigger Send(GOSSIP,  $p$ ,  $m$ ,  $mID$ ,  $\text{myself}$ )
6      trigger Deliver( $m$ )
7       $\text{receivedMsgs} \leftarrow \text{receivedMsgs} \cup \{mID\}$ 

8  upon event Receive(GOSSIP,  $m$ ,  $mID$ , sender) do
9      if  $mID \notin \text{receivedMsgs}$  then
10          $\text{receivedMsgs} \leftarrow \text{receivedMsgs} \cup \{mID\}$ 
11         trigger Deliver( $m$ )
12          $\text{peerList} \leftarrow \text{getPeer}(f, \text{sender})$ 
13         foreach  $p \in \text{peerList}$  do
14             trigger Send(GOSSIP,  $p$ ,  $m$ ,  $mID$ ,  $\text{myself}$ )
```

HyParView: Is this really Good?

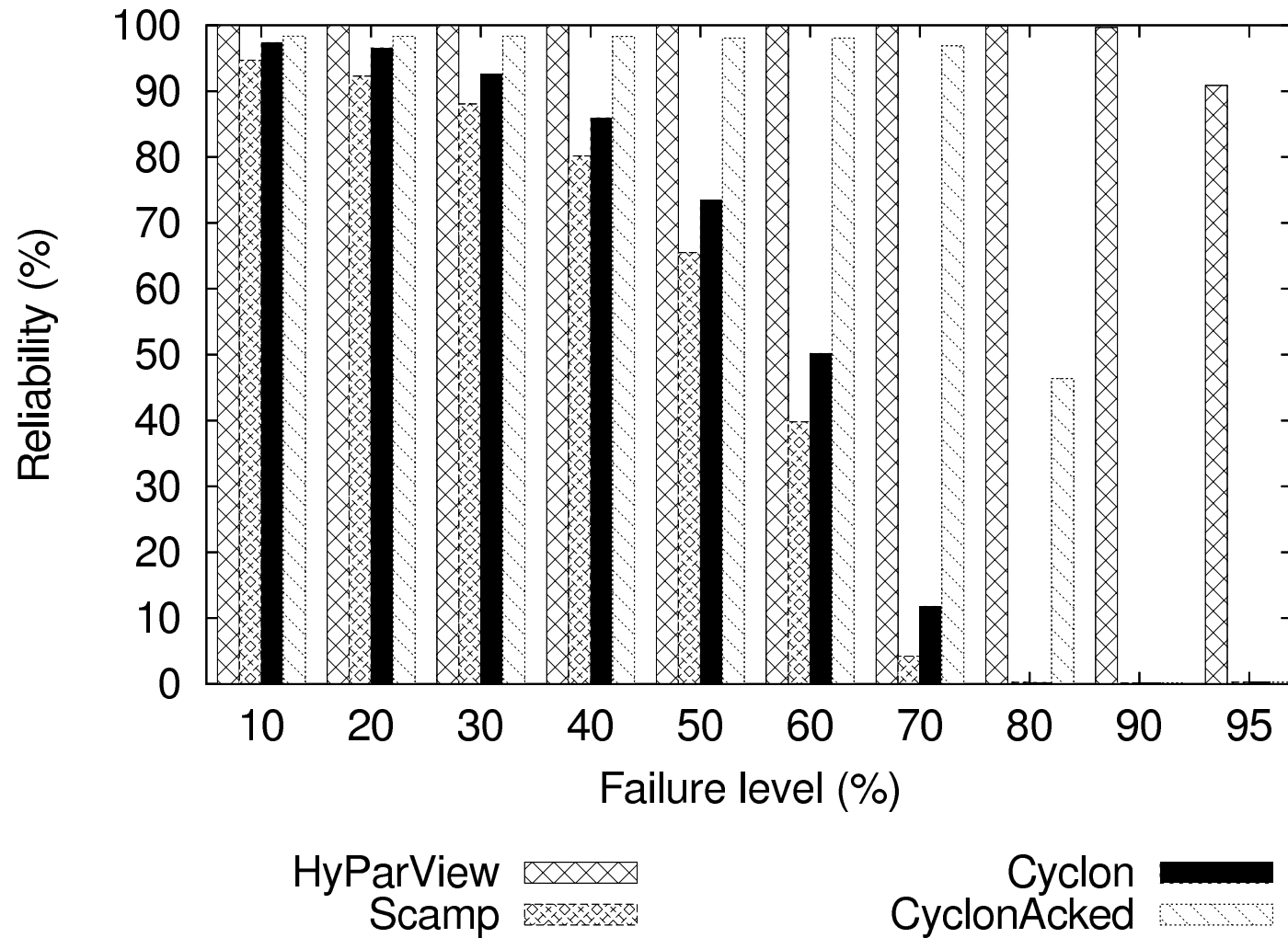


Figure 4.2: Average reliability for 1000 messages

Class structure:

- One Unstructured Overlay: HyParView.
- **One Partially Unstructured Overlay: Plumtree.**

Class structure:

- One Unstructured Overlay: HyParView.
- **One Partially Unstructured Overlay: Plumtree.**

Plumtree is not required for the Phase 1 of the course project, but you might want to think about it on the second phase.

Plumtree Protocol

26th IEEE International Symposium on Reliable Distributed Systems

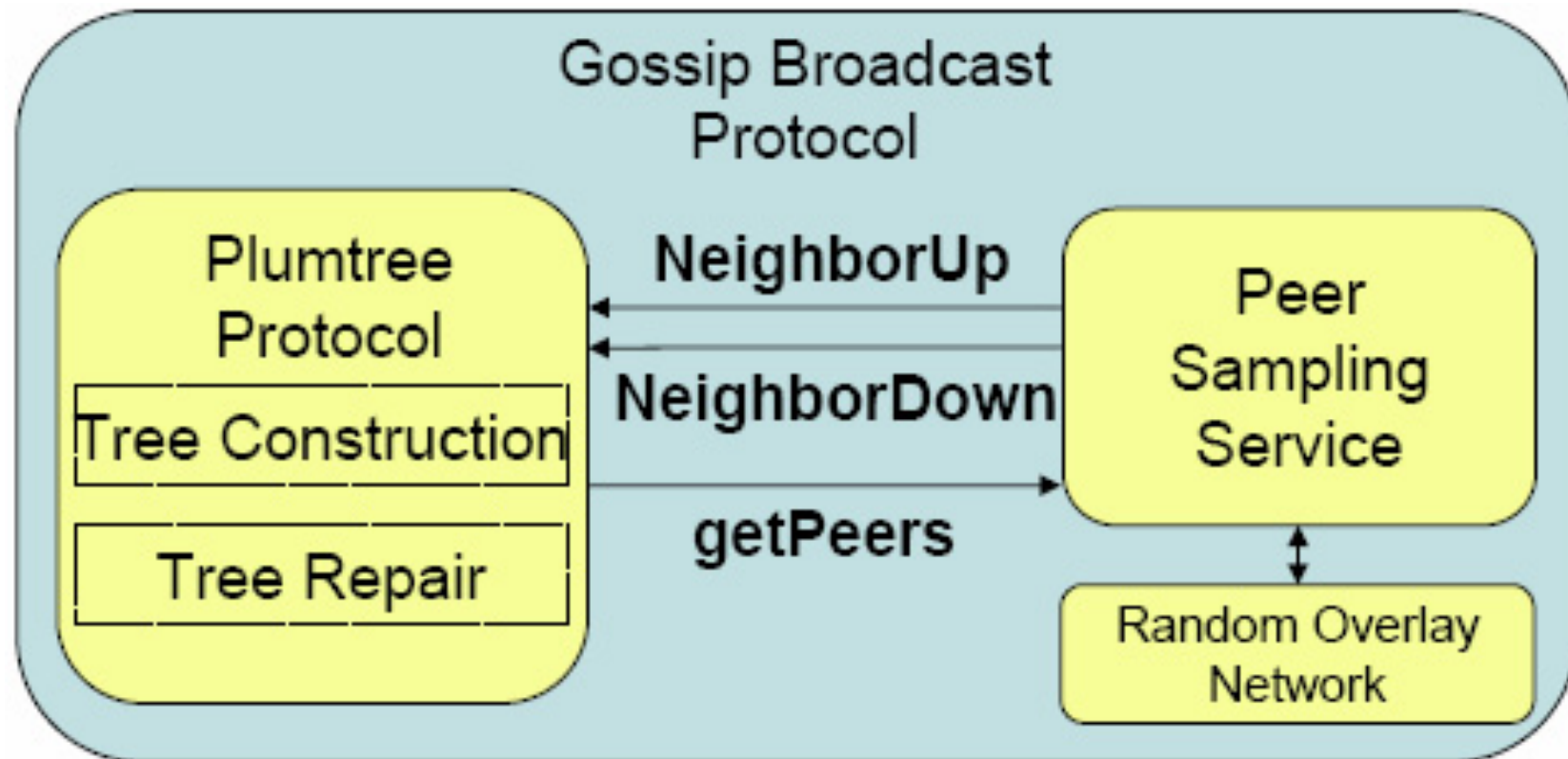
Epidemic Broadcast Trees *

João Leitão
University of Lisbon
jleitao@lasige.di.fc.ul.pt

José Pereira
University of Minho
jop@di.uminho.pt

Luís Rodrigues
University of Lisbon
ler@di.fc.ul.pt

Plumtree Architecture



Plumtree Requirements of the Unstructured Overlay Network

- Ensure connectivity
- Scalable
- Reactive membership

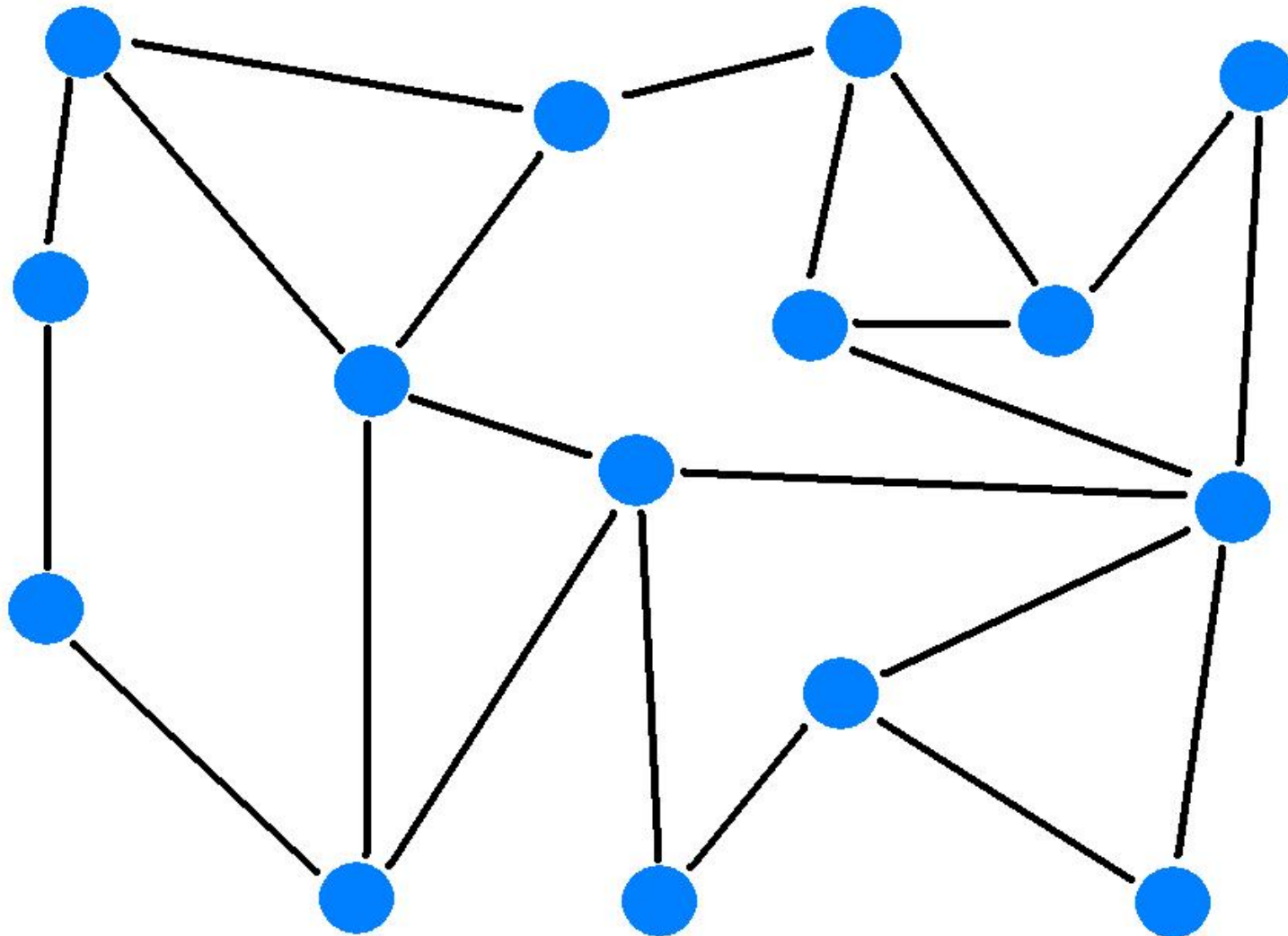
Optionally:

- Symmetric partial views

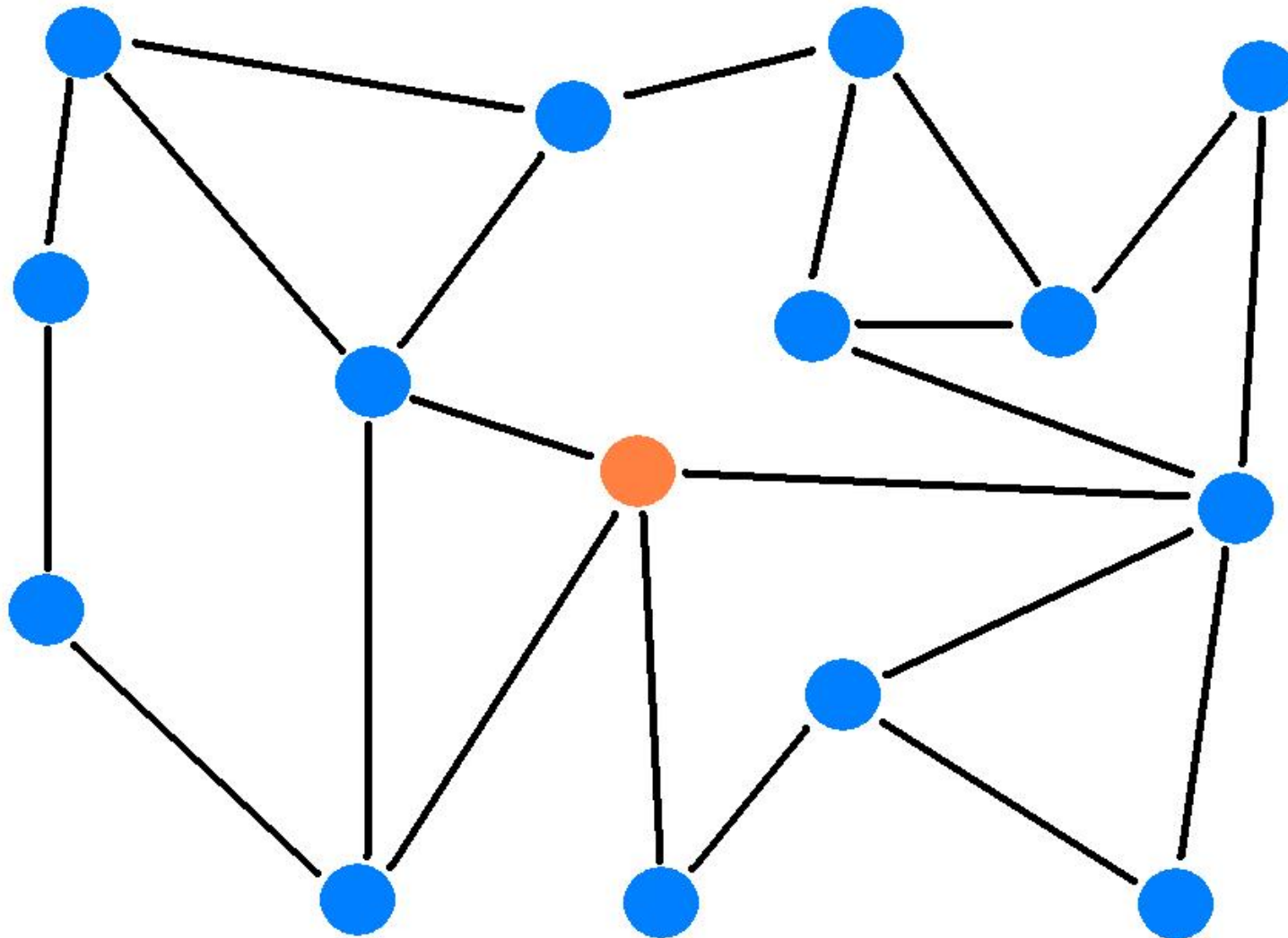
Plumtree: The Intuition

- Think about a gossip protocol (eager push) operating on top of a static overlay network, where you actually use flooding.
- What is going to happen here?

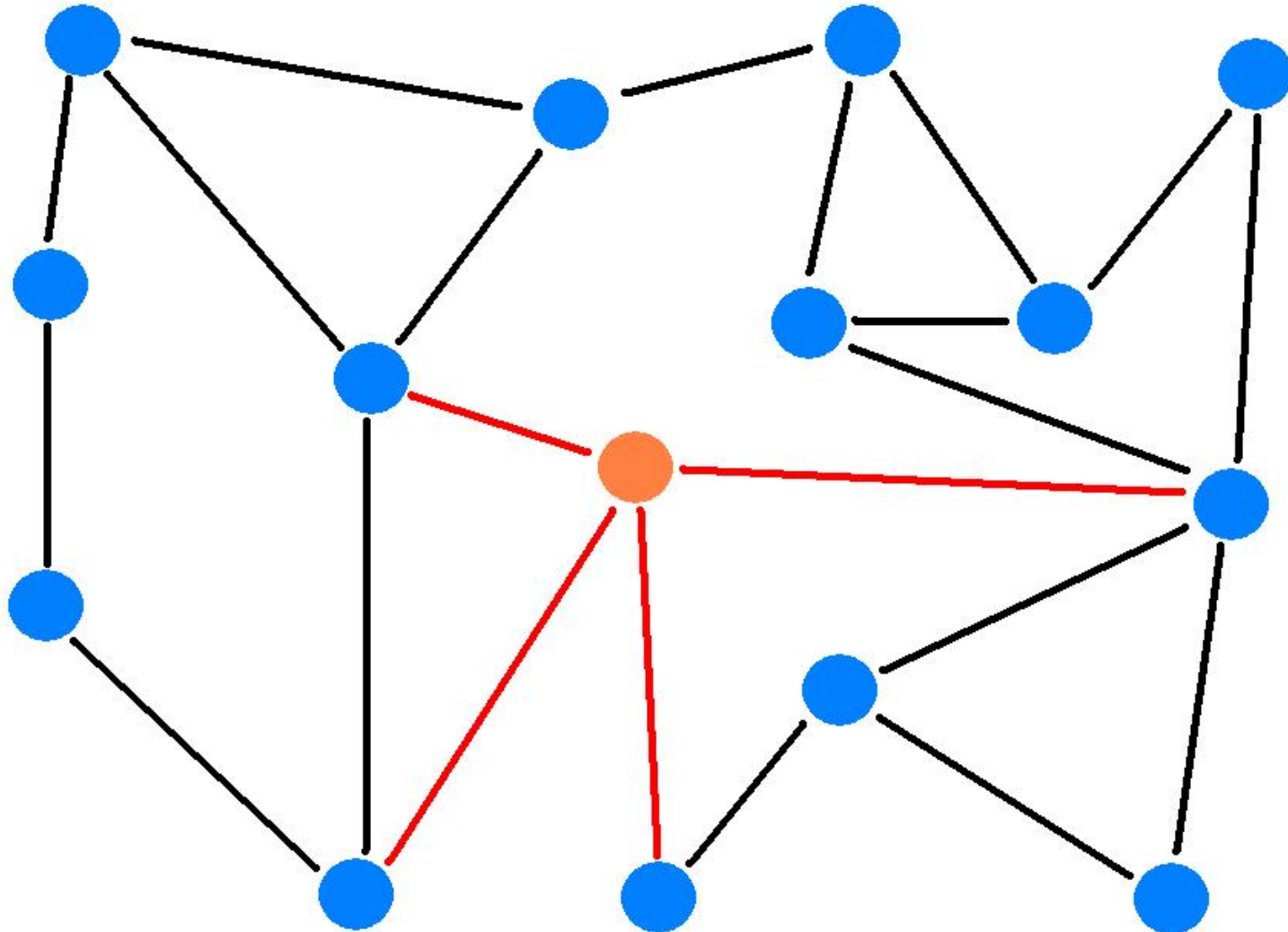
Plumtree: The Intuition



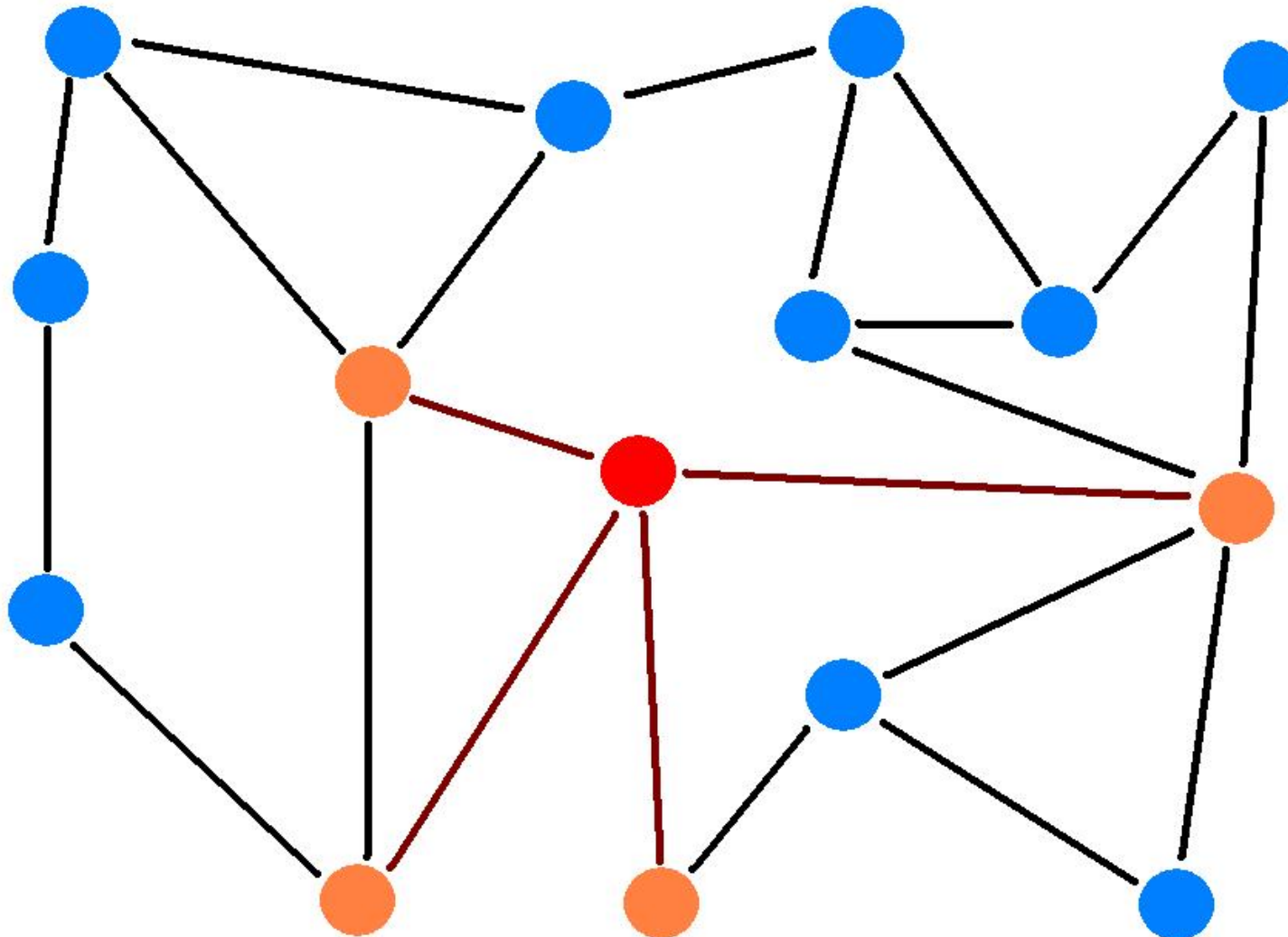
Plumtree: The Intuition



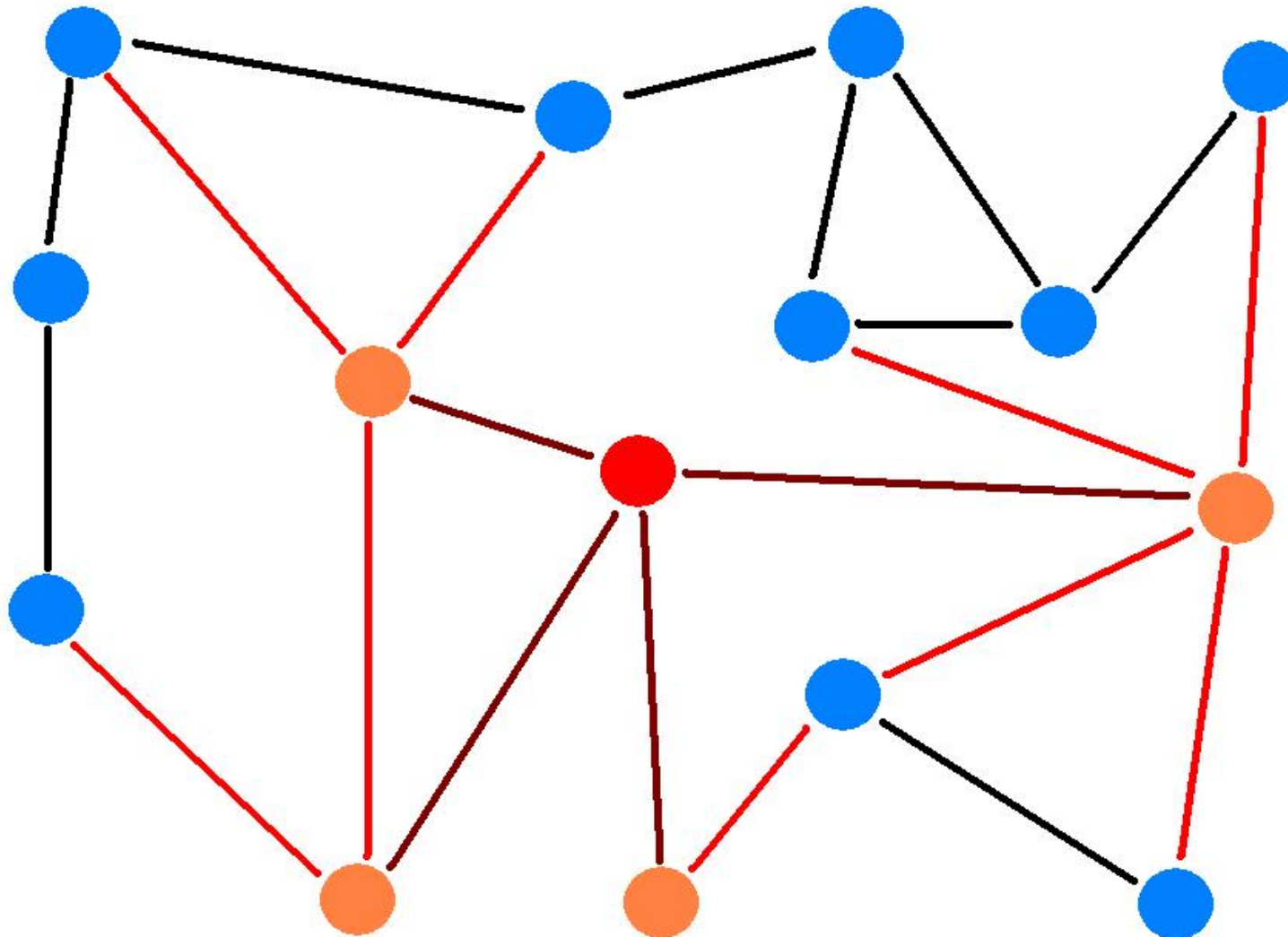
Plumtree: The Intuition



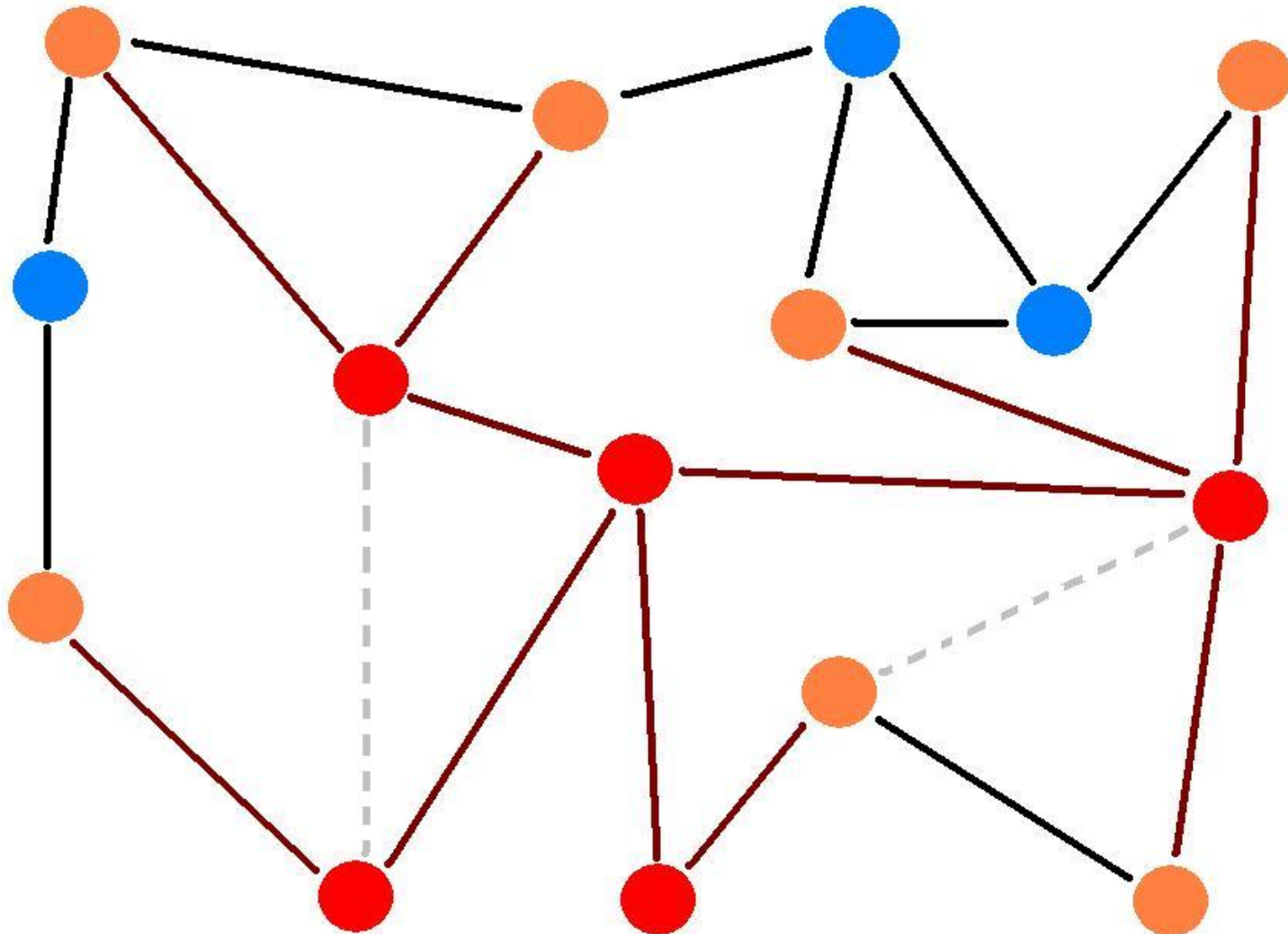
Plumtree: The Intuition



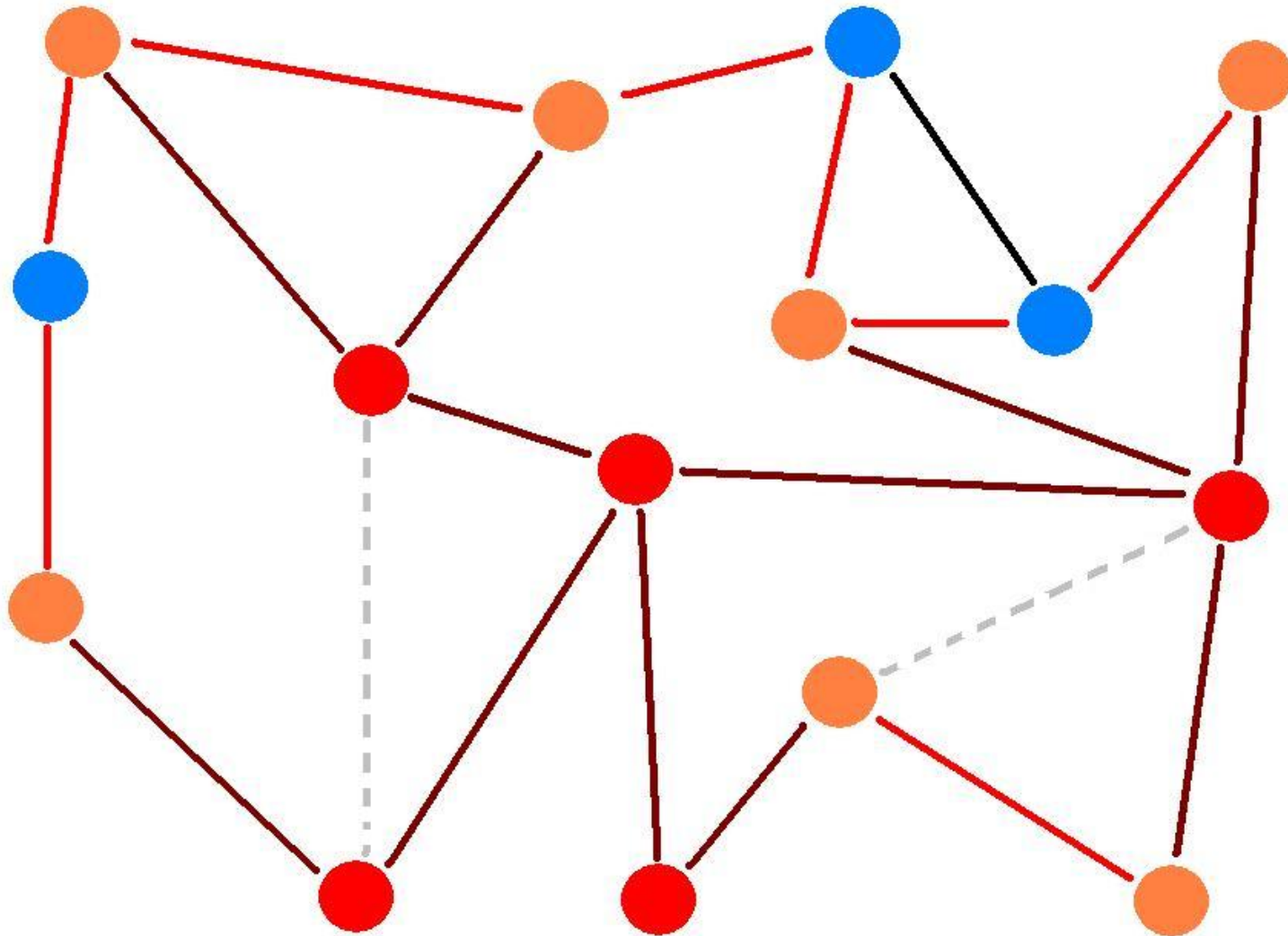
Plumtree: The Intuition



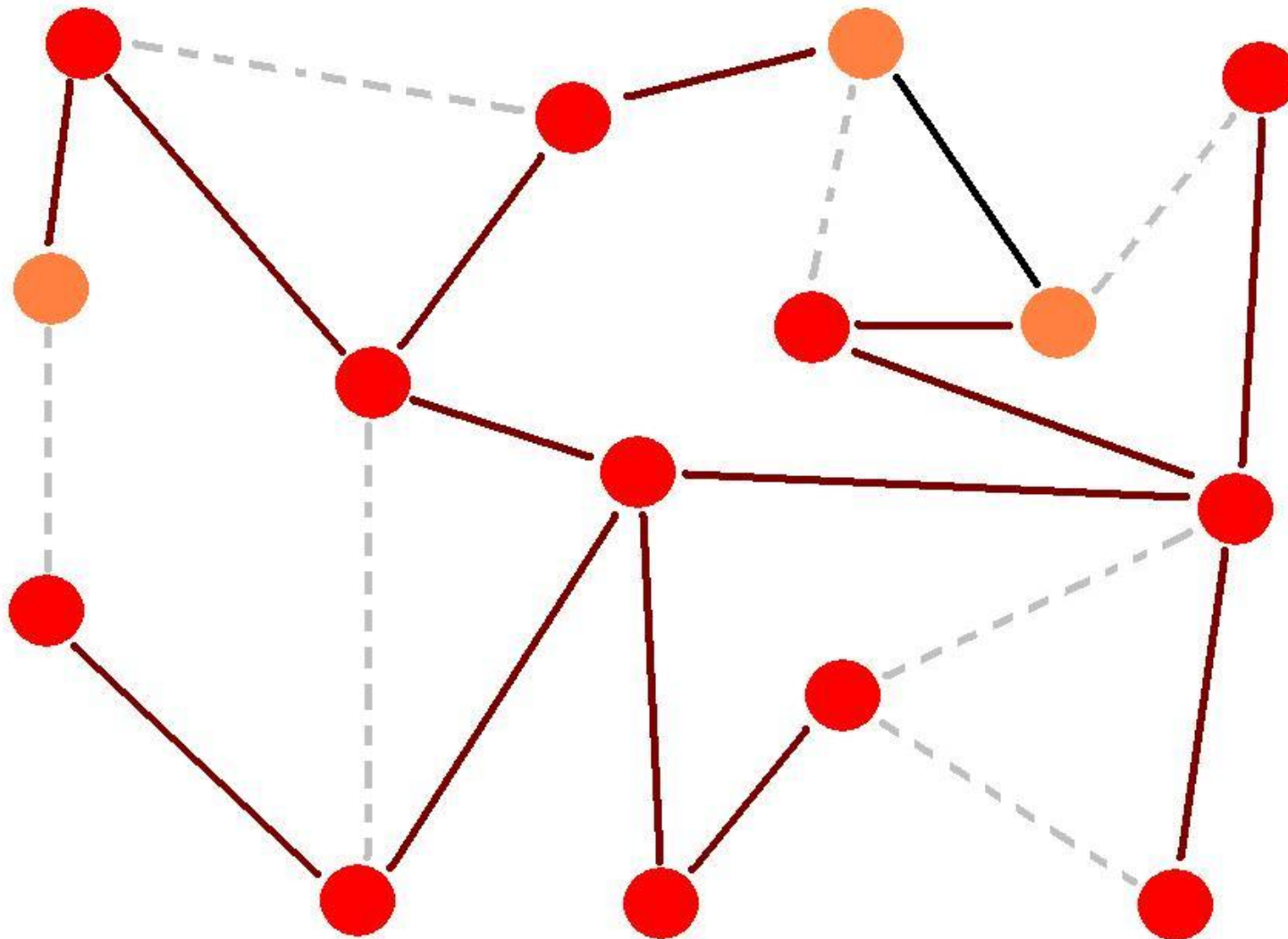
Plumtree: The Intuition



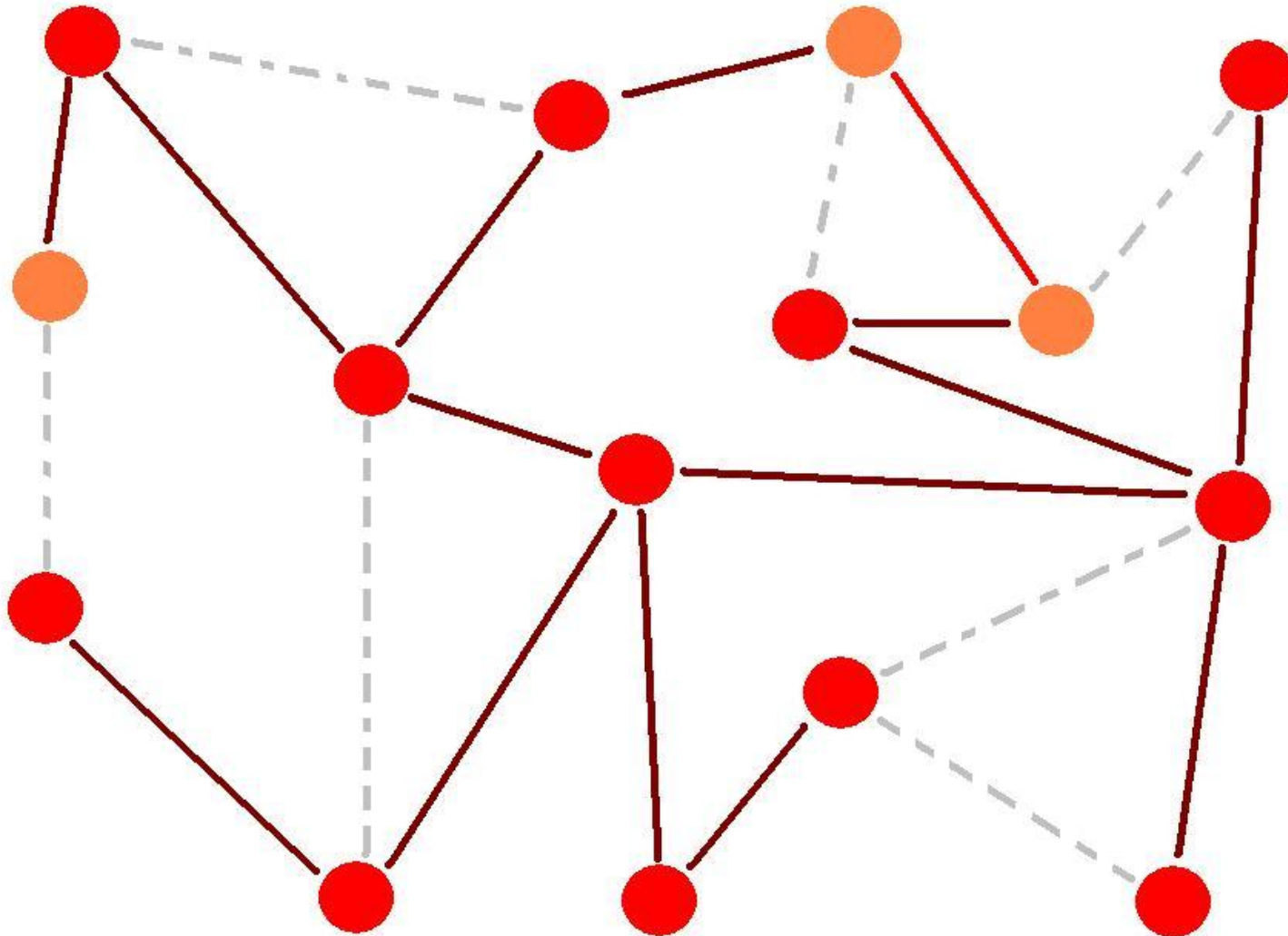
Plumtree: The Intuition



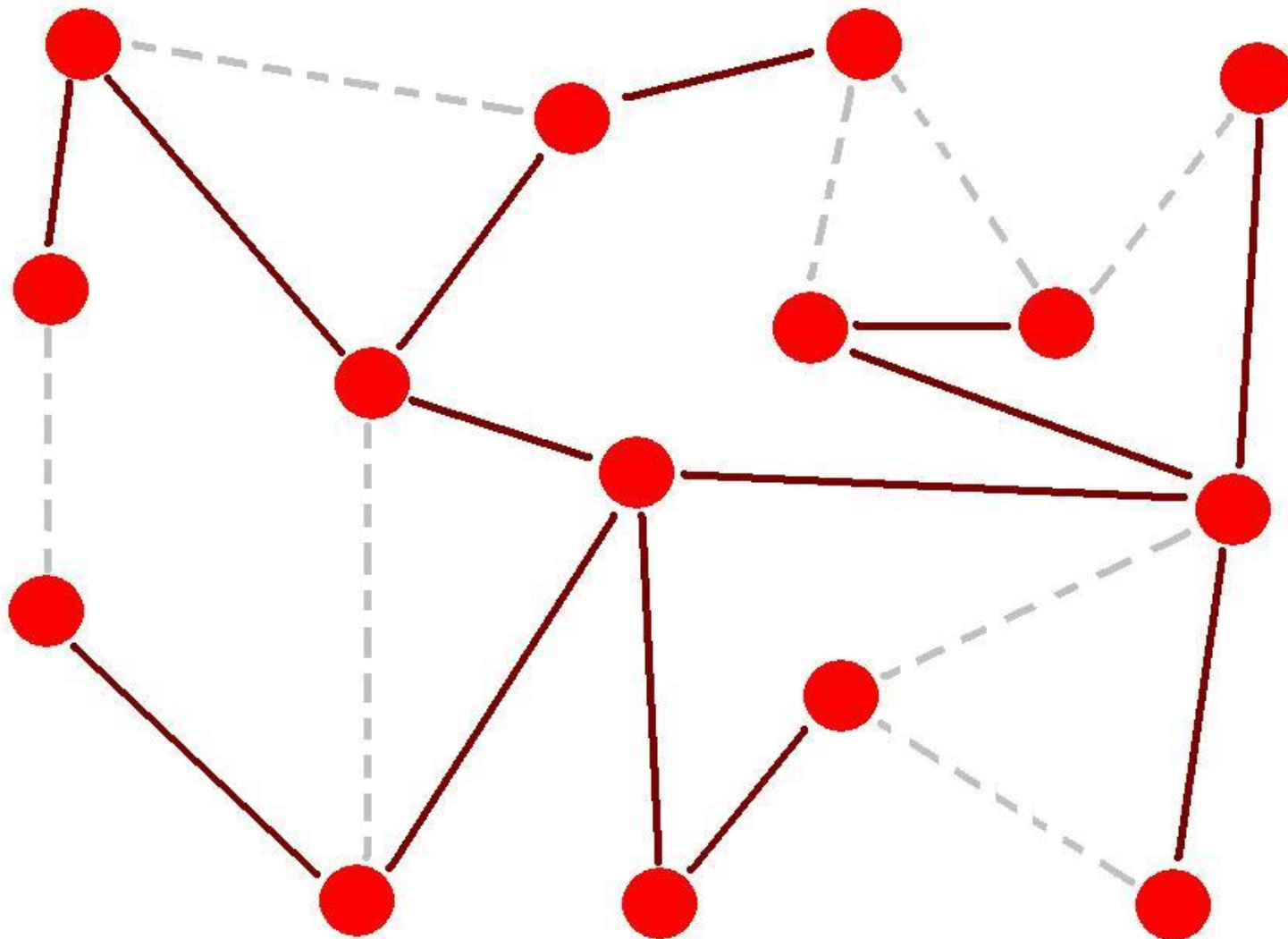
Plumtree: The Intuition



Plumtree: The Intuition



Plumtree: The Intuition



Plumtree Overview

Plumtree: **push-lazy-push multicast tree**.

- Protocol operates as any gossip protocol, each node gossips with t neighbors provided by a peer sampling service by combining eager push and lazy push.
- To support the operation of the protocol, two distinct sets of neighbors are maintained:
 - eagerPushPeers
 - lazyPushPeers
- During initialization t nodes are obtained from the peer sampling service and inserted in the eagerPushPeers.
- TCP is used between nodes as it offers a better reliability and an additional fault tolerance mechanism.

Plumtree: Tree Construction

- Intuition: Use the links in the overlay that generated a message delivery.
- After initialization all links in the overlay are candidates do belong to the broadcast tree.
- When a message is broadcasted, all links used to disseminate a redundant message are pruned from the tree.
- If the random overlay is connected, after the dissemination of a message the tree will cover all nodes.

Plumtree: Tree Repair

- If a node fails some nodes might become disconnected from the embedded tree.
- When a node receives a announcement for a payload message it did not received it starts a timer.
- If the timer expires before receiving the payload message the node explicitly request the payload.
- Moves the node to whom it makes the request to its eagerPushPeers list.
- The node that receives the request sends the payload.
- Moves the requester to its eagerPushPeers list.
- At the end of this process the embedded tree is repaired.

Plumtree: The Algorithm

- State

Algorithm 4: Internal data structure

Data:

myself: the identifier of the local node

receivedMsgs: a list of received messages identifiers

f: the push fanout value

eagerPushPeers: a list of neighbors whom links form the spanning tree

lazyPushPeers: a list of neighbors whom links does not belong to the spanning tree

lazyQueue: a list of tuples {mID,node,round}

Plumtree: The Algorithm

- Tree Construction

Algorithm 5: Spanning tree construction algorithm

```
1  procedure dispatch do
2    announcements  $\leftarrow$  policy (lazyQueue) //set of IHave messages
3    trigger Send(announcements)
4    lazyQueue  $\leftarrow$  lazyQueue  $\setminus$  announcements

5  procedure EagerPush (m, mID, round, sender) do
6    foreach  $p \in$  eagerPushPeers:  $p \neq$ sender do
7      trigger Send(GOSSIP,  $p$ ,  $m$ ,  $mID$ , round, myself)

8  procedure LazyPush (m, mID, round, sender) do
9    foreach  $p \in$  lazyPushPeers:  $p \neq$ sender do
10     lazyQueue  $\leftarrow$  (textscIHave( $p$ ,  $m$ ,  $mID$ , round, myself)
11     call dispatch()

12 upon event Init do
13   eagerPushPeers  $\leftarrow$  getPeer(f)
14   lazyPushPeers  $\leftarrow$   $\emptyset$ 
15   lazyQueue  $\leftarrow$   $\emptyset$ 
16   missing  $\leftarrow$   $\emptyset$ 
17   receivedMsgs  $\leftarrow$   $\emptyset$ 

18 upon event Broadcast( $m$ ) do
19   mID  $\leftarrow$  hash( $m$ +myself)
20   call EagerPush (m, mID, 0, myself)
21   call LazyPush (m, mID, 0, myself)
22   trigger Deliver( $m$ )
23   receivedMsgs  $\leftarrow$  receivedMsgs  $\cup$  {mID}

24 upon event Receive(GOSSIP,  $m$ , mID, round, sender) do
25   if  $mID \notin$  receivedMsgs then
26     trigger Deliver( $m$ )
27     receivedMsgs  $\leftarrow$  receivedMsgs  $\cup$  {mID}
28     if  $\exists (id, node, r) \in$  missing : $id=mID$  then
29       cancel Timer(mID)
30     call EagerPush (m, mID, round+1, myself)
31     call LazyPush (m, mID, round+1, myself)
32     eagerPushPeers  $\leftarrow$  eagerPushPeers  $\cup$  {sender}
33     lazyPushPeers  $\leftarrow$  lazyPushPeers  $\setminus$  {sender}
34     call Optimize ( $m$ , mID, round, sender) // optional
35   else
36     eagerPushPeers  $\leftarrow$  eagerPushPeers  $\setminus$  {sender}
37     lazyPushPeers  $\leftarrow$  lazyPushPeers  $\cup$  {sender}
38     trigger Send(PRUNE, sender, myself)

39 upon event Receive(PRUNE, sender) do
40   eagerPushPeers  $\leftarrow$  eagerPushPeers  $\setminus$  {sender}
41   lazyPushPeers  $\leftarrow$  lazyPushPeers  $\cup$  {sender}
```

Plumtree: The Algorithm

- Tree Repair

Algorithm 6: Spanning tree repair algorithm

```
1  upon event Receive(IHAVE, mID, round, sender) do
2      if  $mID \notin receivedMsgs$  do
3          if  $\nexists Timer(id): id=mID$  do
4              setup Timer(mID, timeout1)
5              missing  $\leftarrow$  missing  $\cup \{(mID, sender, round)\}$ 

6  upon event Timer(mID) do
7      setup Timer(mID, timeout2)
8      (mID, node, round)  $\leftarrow$  removeFirstAnnouncement(missing, mID)
9      eagerPushPeers  $\leftarrow$  eagerPushPeers  $\cup \{node\}$ 
10     lazyPushPeers  $\leftarrow$  lazyPushPeers  $\setminus \{node\}$ 
11     trigger Send(GRAFT, node, mID, round, myself)

12 upon event Receive(GRAFT, mID, round, sender) do
13     eagerPushPeers  $\leftarrow$  eagerPushPeers  $\cup \{sender\}$ 
14     lazyPushPeers  $\leftarrow$  lazyPushPeers  $\setminus \{sender\}$ 
15     if  $mID \in receivedMsgs$  do
16         trigger Send(GOSSIP, sender, m, mID, round, myself)
```

Plumtree: Multiple Senders

Small number of senders

- One Plumtree instance for each sender.
- Lower latency.

Large number of senders

- One single Plumtree instance shared by all senders.
- Higher latency.

Summary

- You know everything you need to know to do the entire phase 1 of the project.
- We already identified some (expected) limitations that we will find in this phase, start thinking how we can address those limitations (that is going to be one of the focus on phase 2).