

# ScalableApplicationLayerMulticast

Suman Banerjee

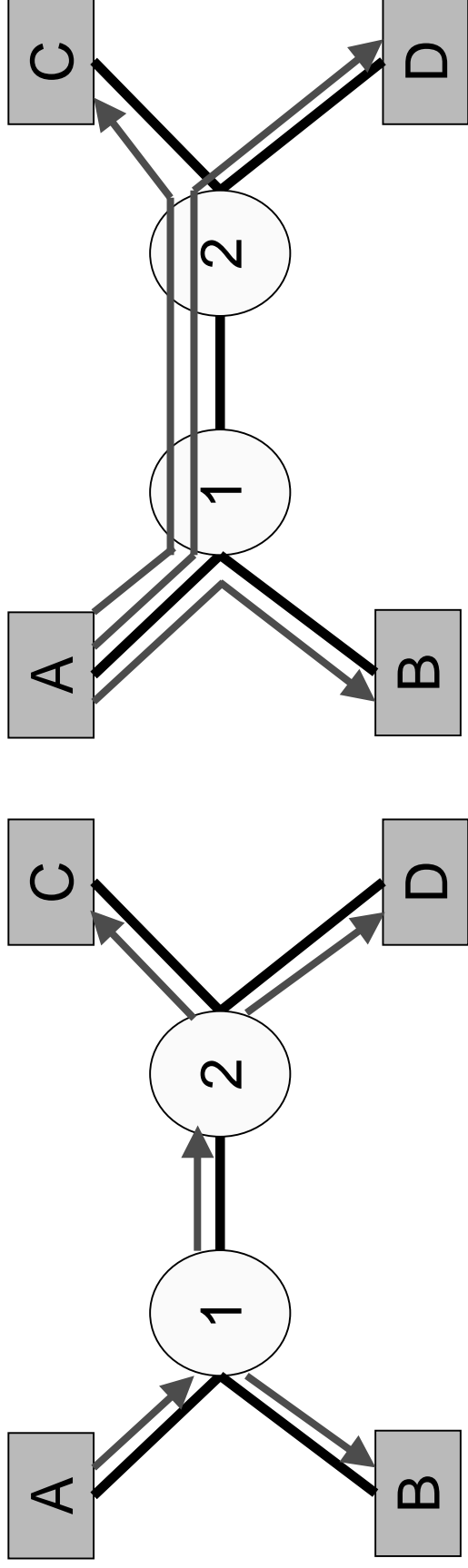
Bobby Bhattacharjee

Christopher Kommareddy



<http://www.cs.umd.edu/projects/nice>

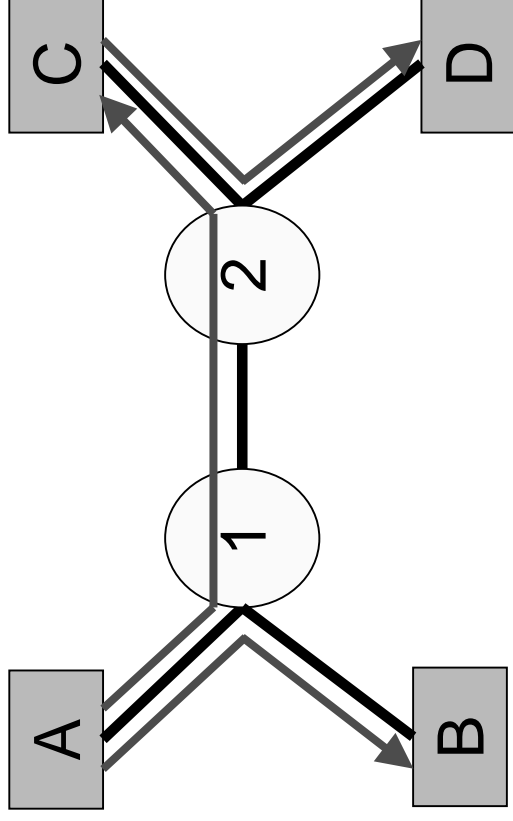
# Group Communication



Network-layer Multicast  
Replication at routers

Sequence of Direct Unicasts  
Replication only at source

# Application-layer Multicast

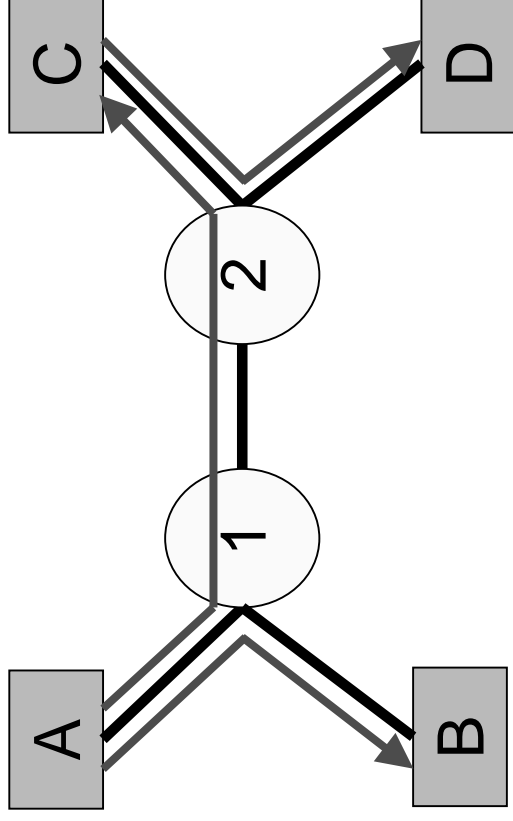


Replication at end-hosts

Examples:

- Narada, Yoid, Gossamer, HMTP, Scribe, Bayeux, CAN -multicast, DT, ...
- **NICE**

# Application-layer Multicast



Replication at end-hosts

## Metrics

Tree Quality

State/Control Overheads

Robustness

# TalkOutline

- Introduction
- NICEApplication -layerMulticastProtocol
- Results
- Conclusions

# NICE Application -layerMulticast

- Scale to large group sizes
  - Low average and worst case control overheads
  - Does not compromise tree equality or robustness
- Even low -bandwidth applications are efficient
  - Web tickers

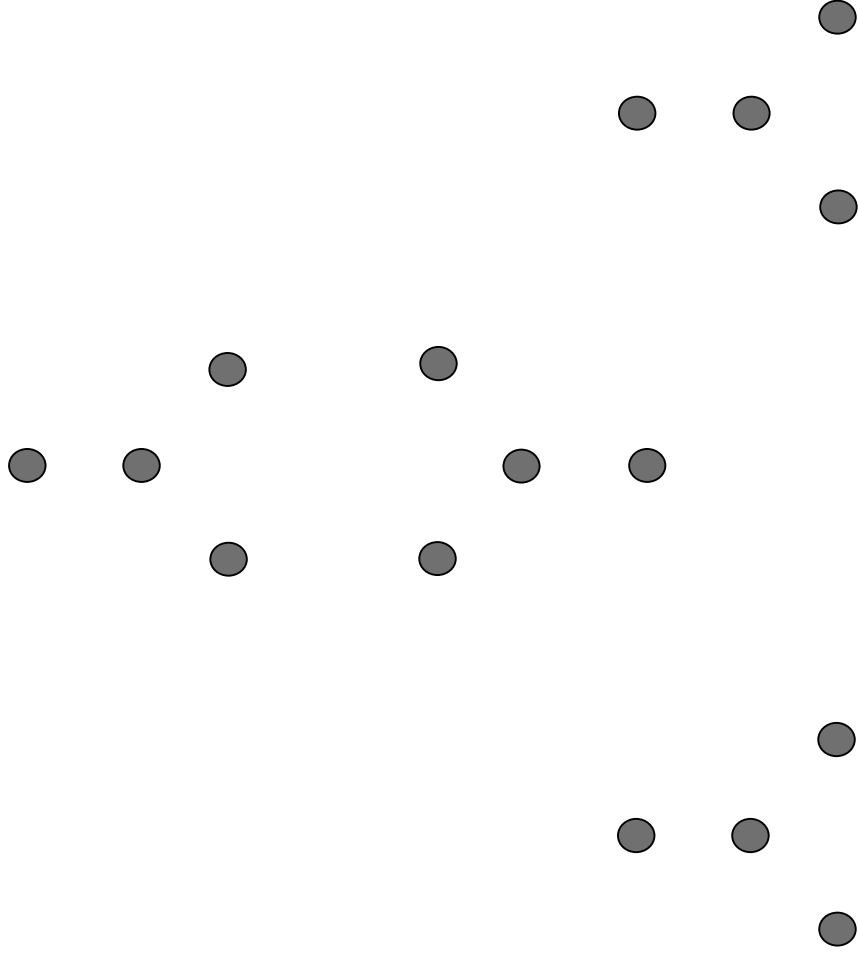
Uses a hierarchy



# NICETopologies

- Controltopology
  - Detectshostfailuresandre overlay -structurethe
- Datadeliverytopology
  - Basicpath:Implicitlydefinedbythehierarchy
  - Canbeindependentofthecontrolpath

# NICE Hierarchy



A Set of Members

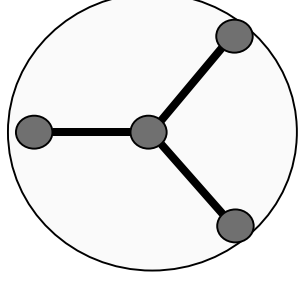
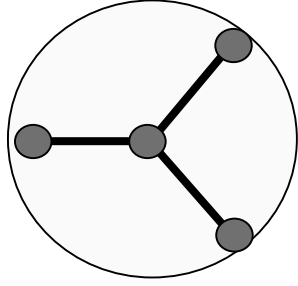
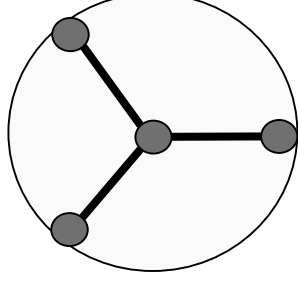
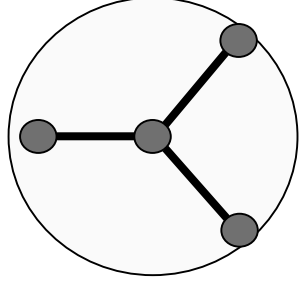




# NICE Hierarchy

## Clusters

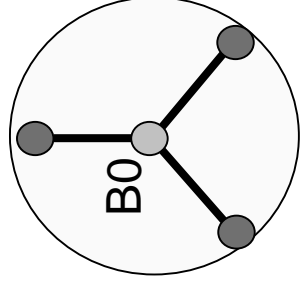
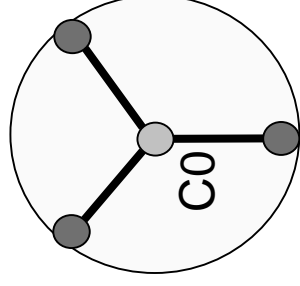
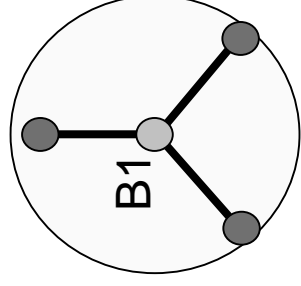
- Non-overlapping
- Proximity-based
- Size:  $k$  to  $3k-1$



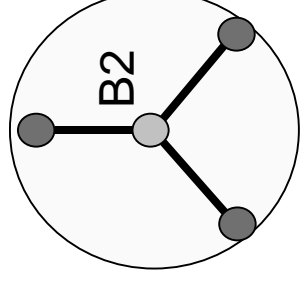
# NICE Hierarchy

## Clusters

- Non-overlapping
- Proximity-based
- Size:  $k$  to  $3k-1$

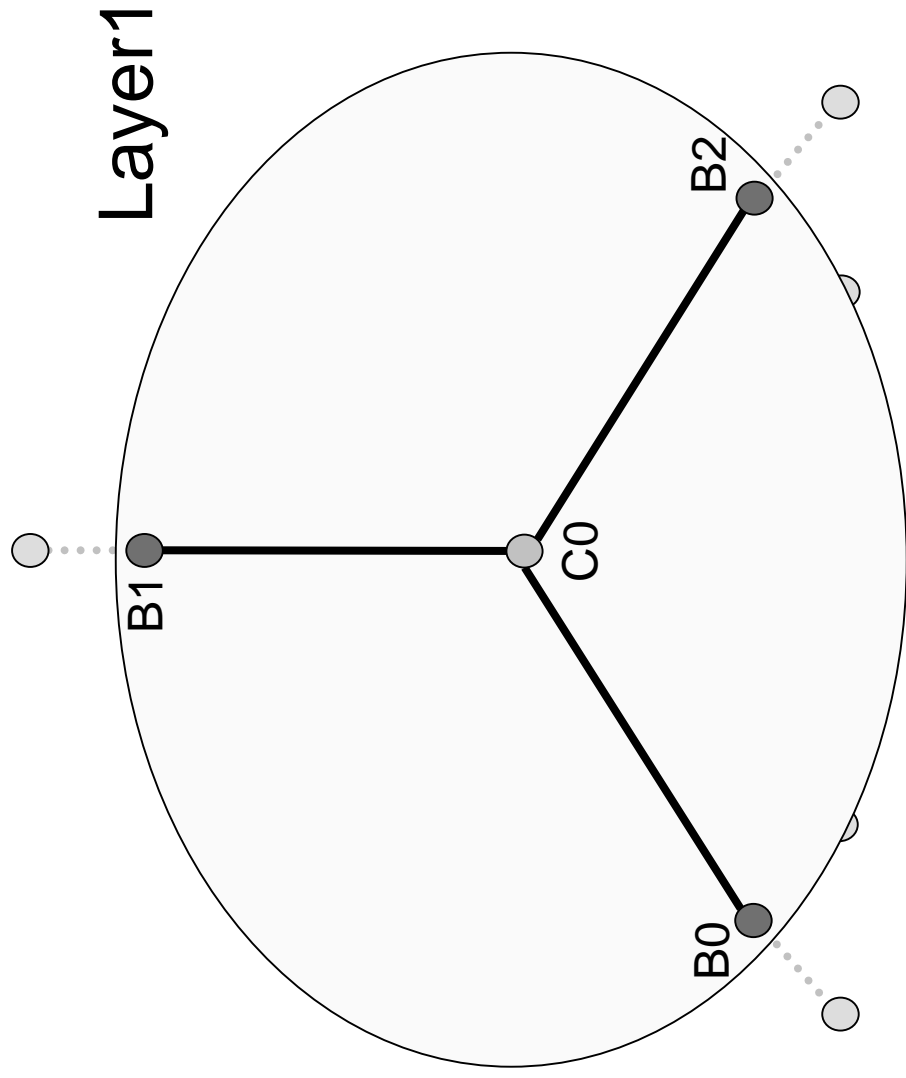


## Layer0



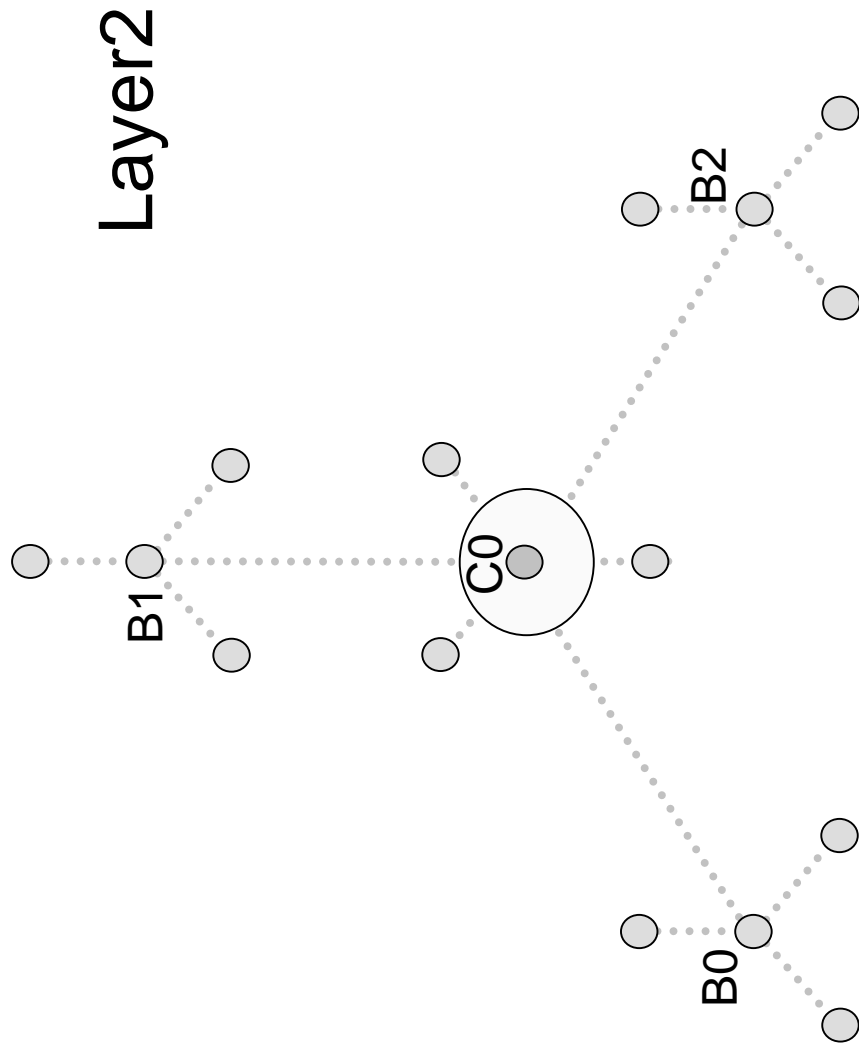
Graph-theoretic center is the cluster leader

# NICE Hierarchy



Leaders for the higher layer  
and repeats

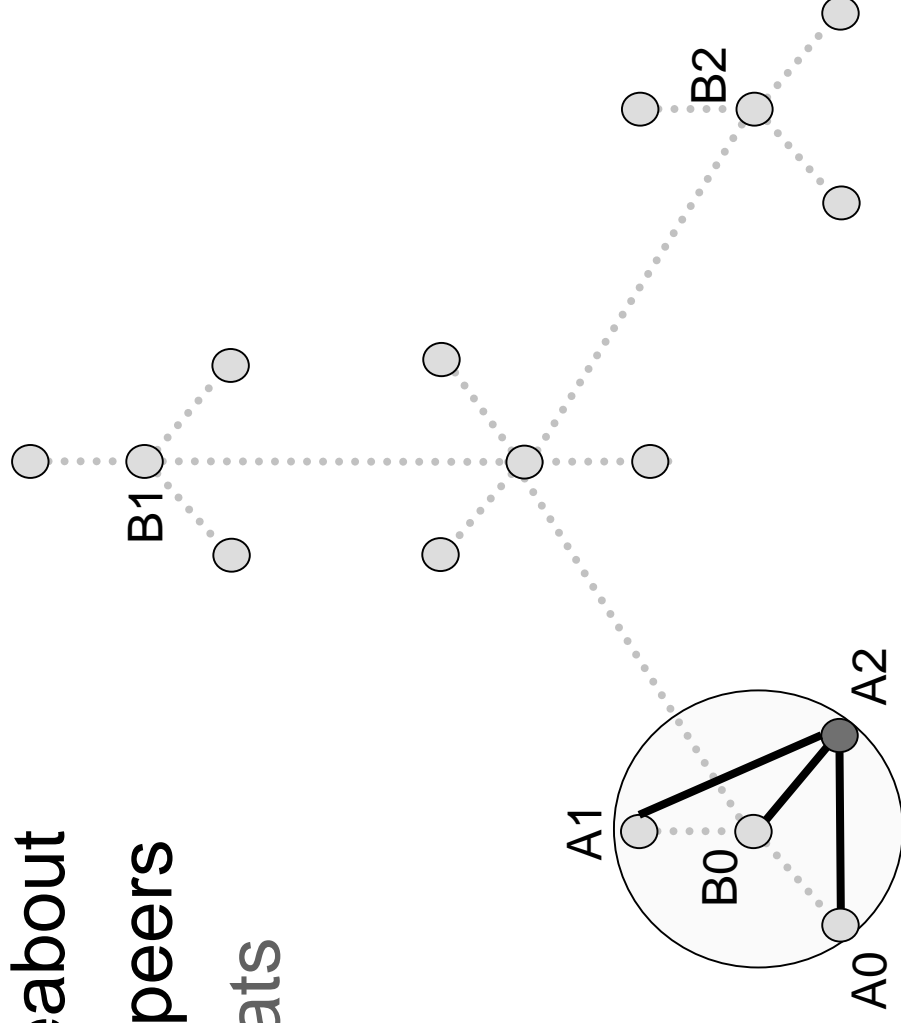
# NICEHierarchy



logNlayers

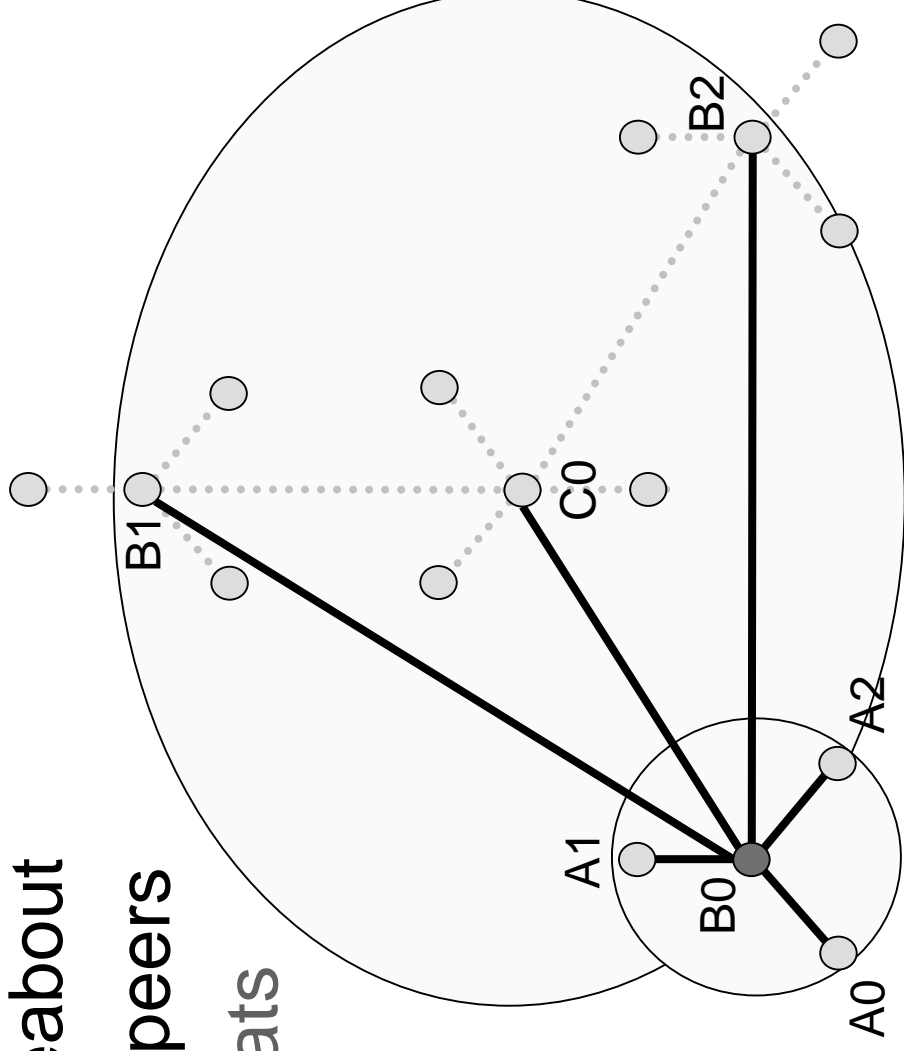
# ControlTopology

- Softstateabout
- allclusterpeers
- HeartBeats



# ControlTopology

- Softstateabout
- allclusterpeers
- HeartBeats

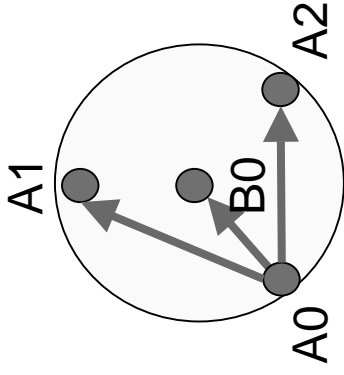
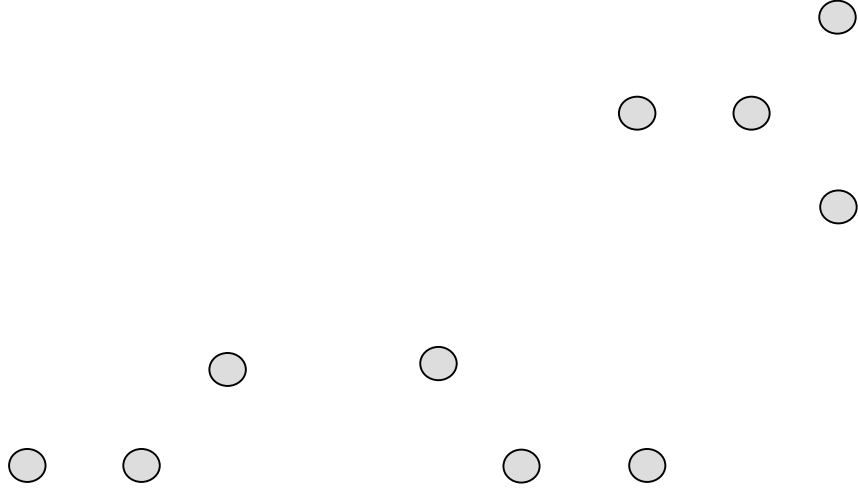


StateandControlmessageoverheads:

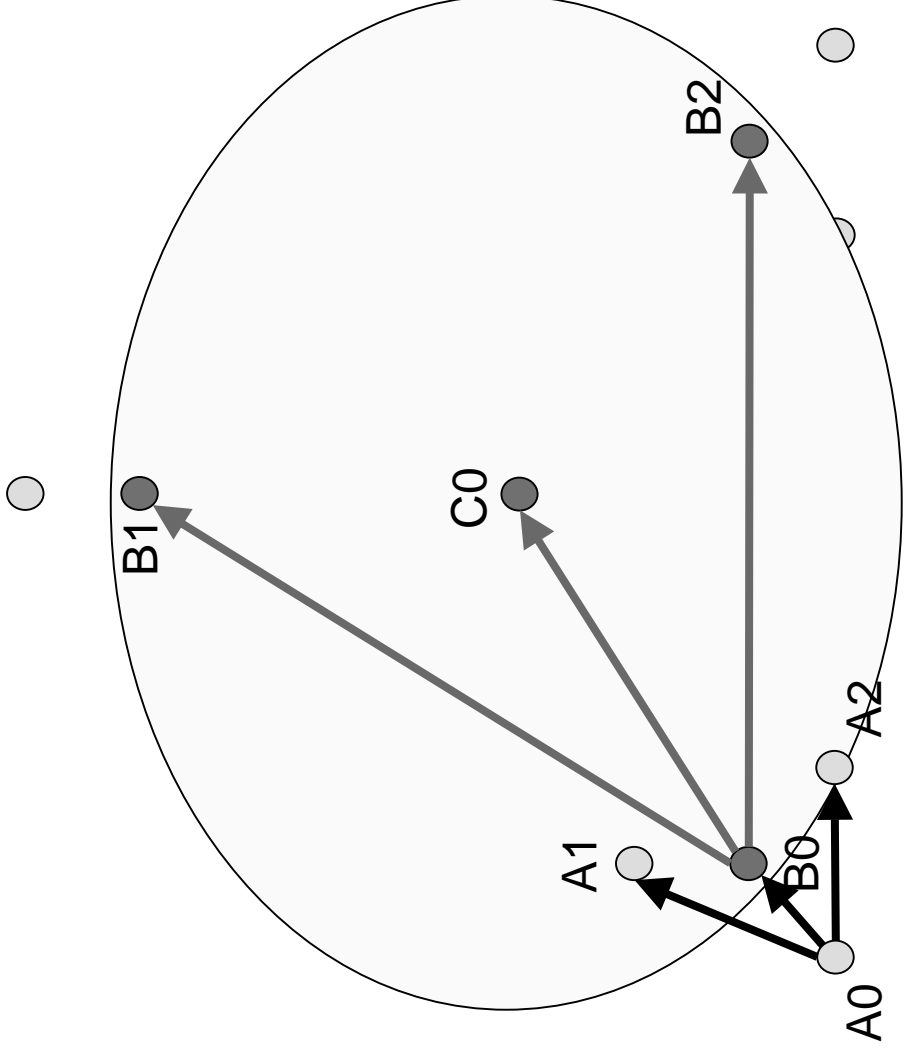
Average:Constant

Worstcase: $O(k\log N)$

# BasicDataPath

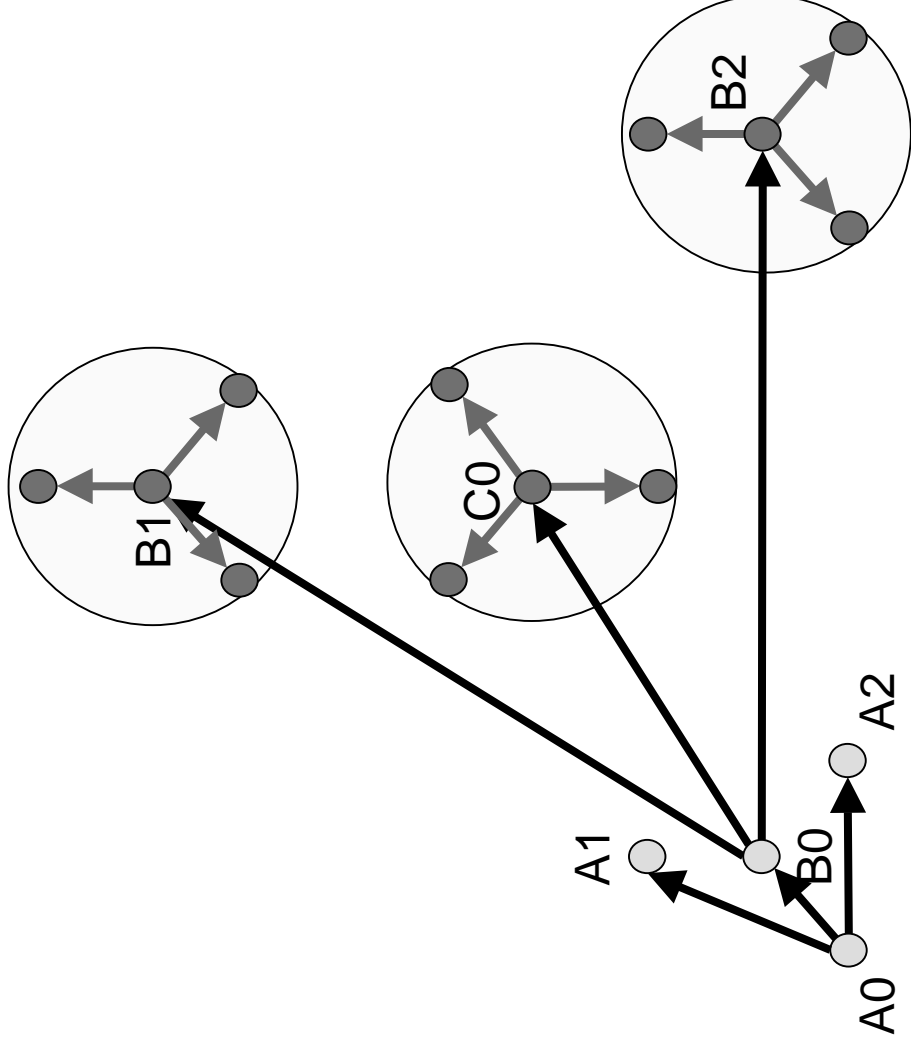


# BasicDataPath

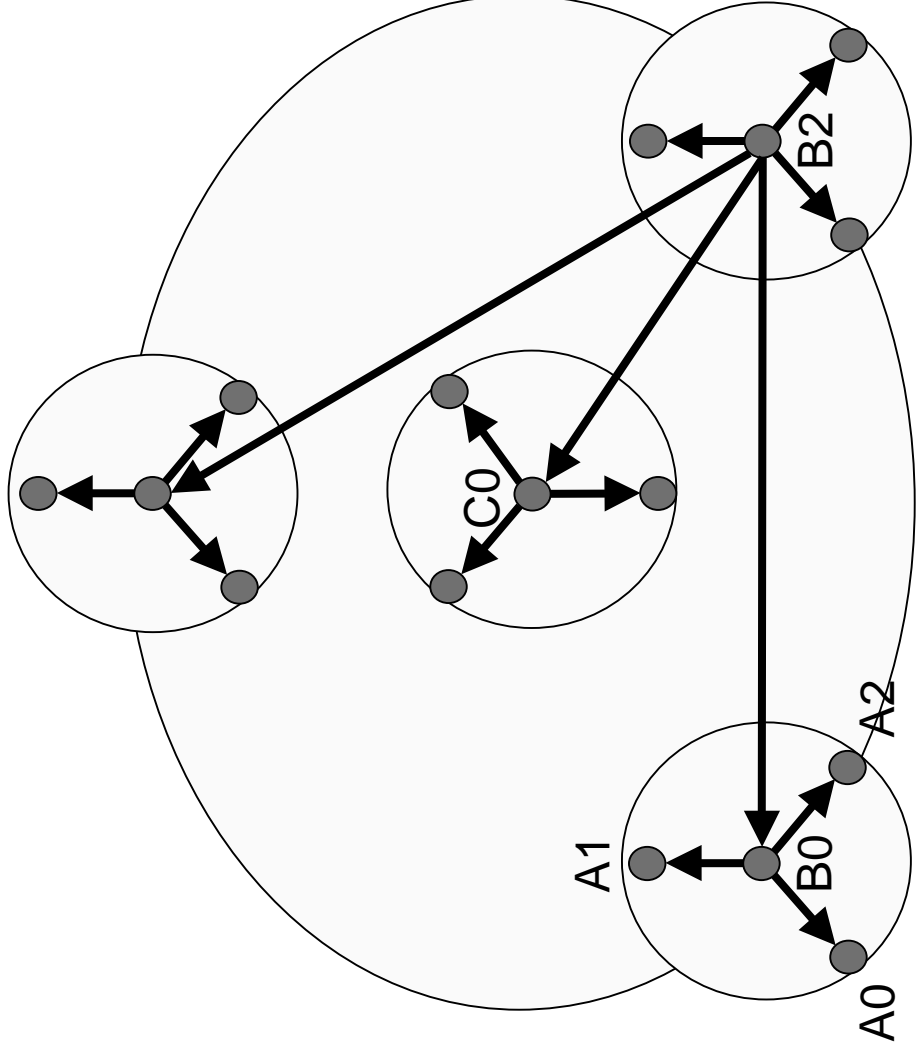




# BasicDataPath



# BasicDataPath



# NICE Invariants

- Cluster sizes between 3k - 1
- Cluster leader is the central member
  - Leaders form next higher layer

NICE protocol maintains these invariants

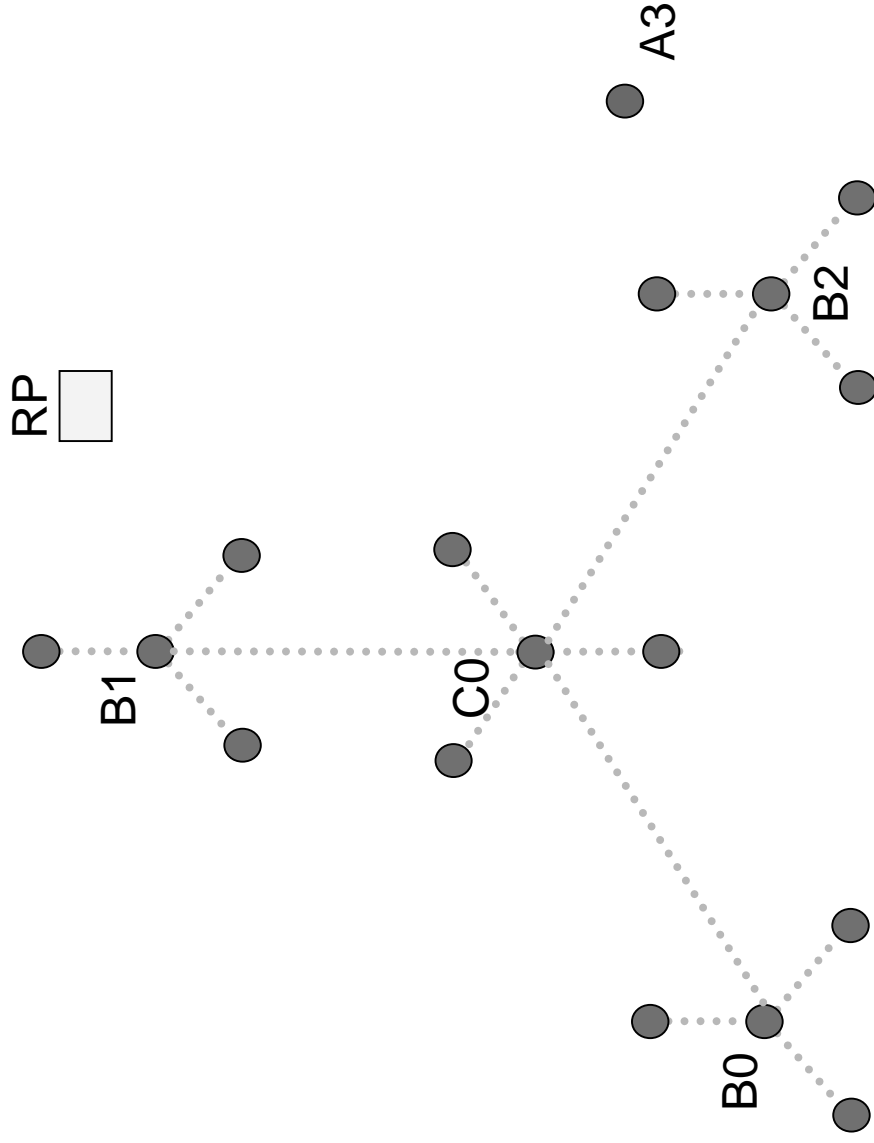


# NICE Protocol Operations

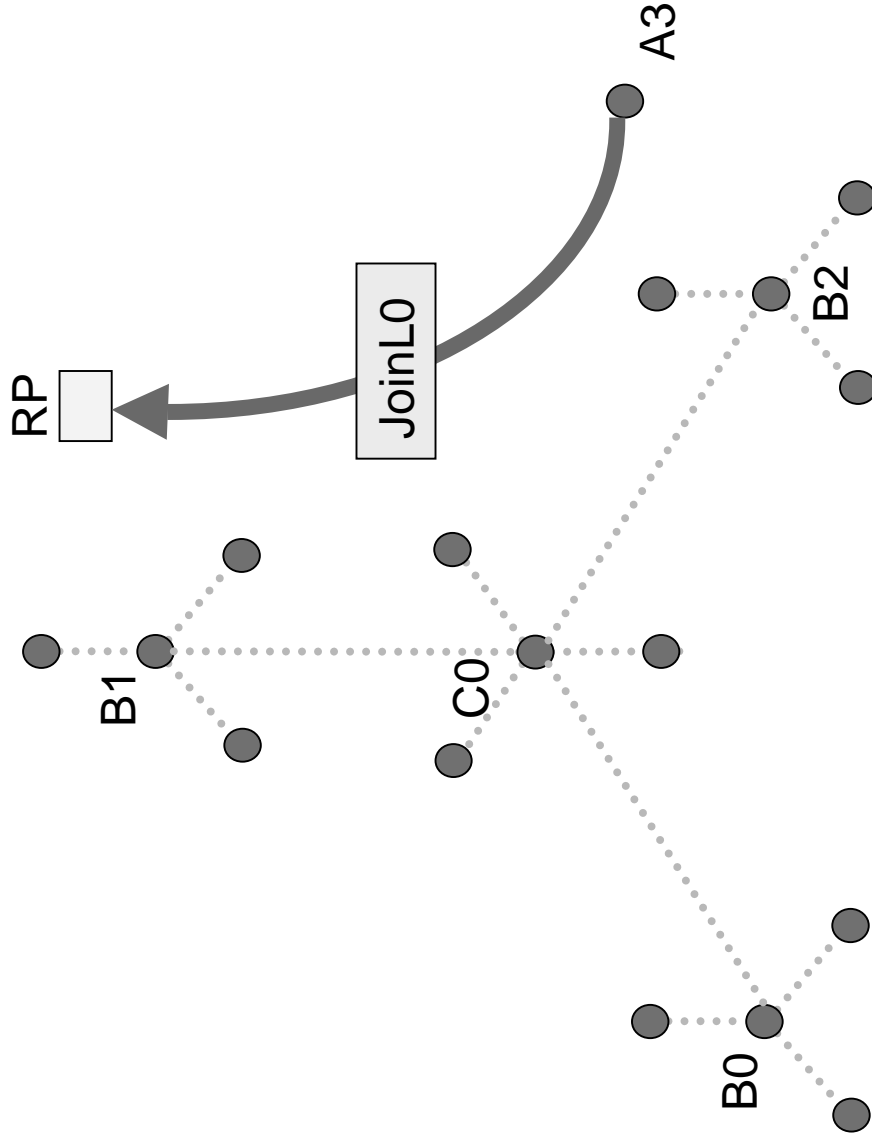
- MemberJoin
- MemberDepart
- ClusterSplit
- ClusterMerge
- ClusterRefine



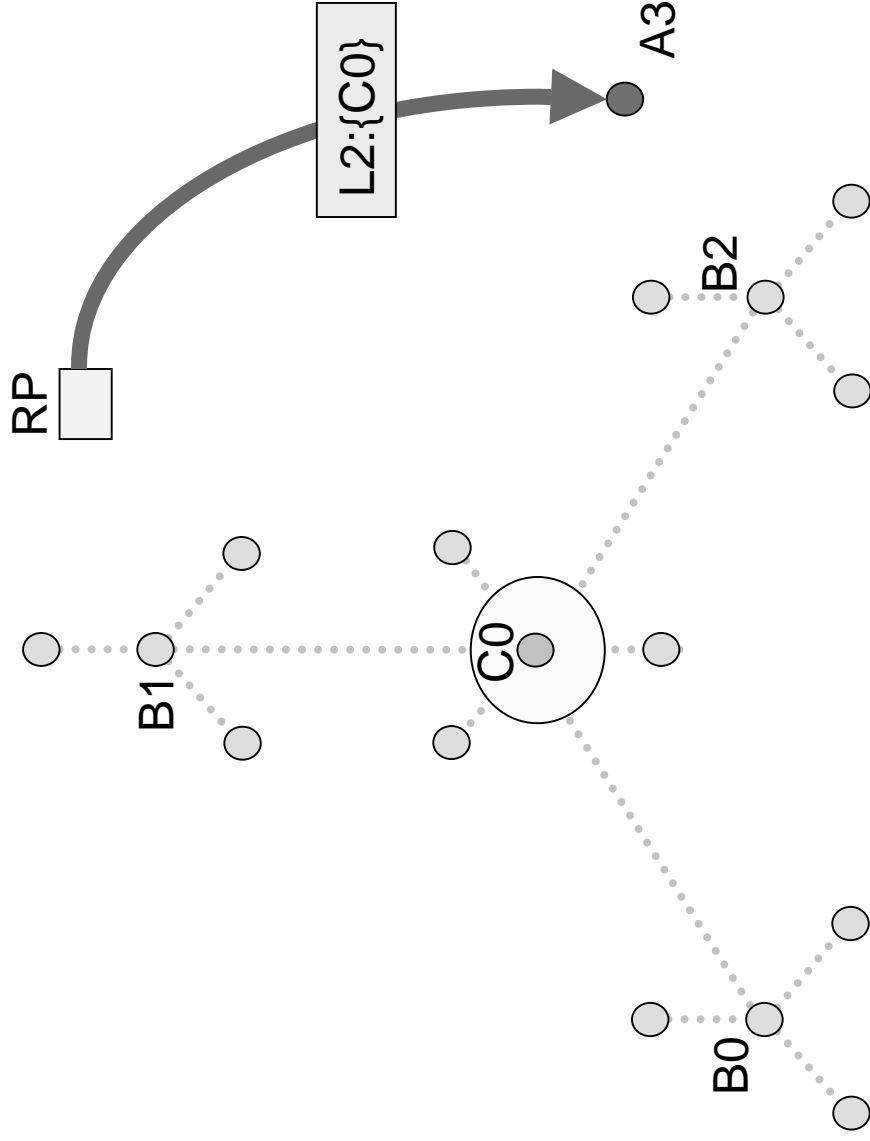
# Join Procedure



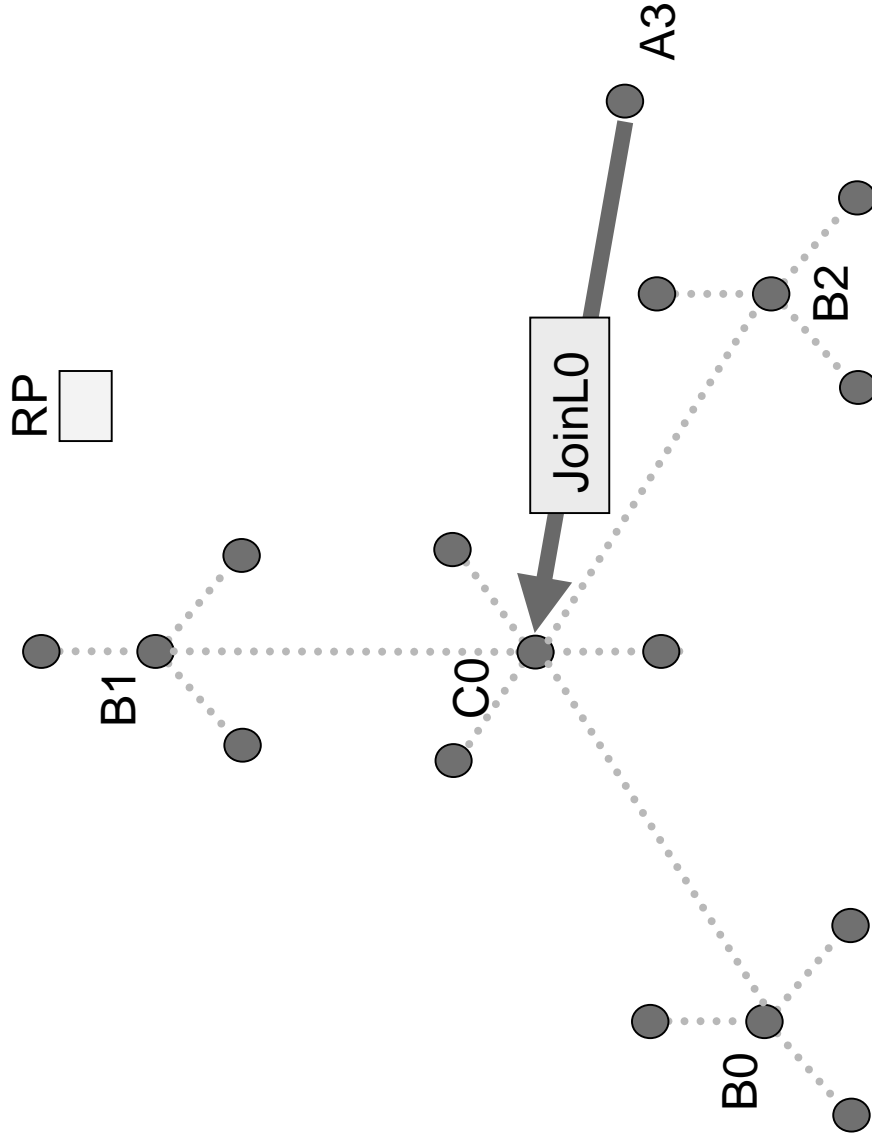
# JoinProcedure



# Join Procedure

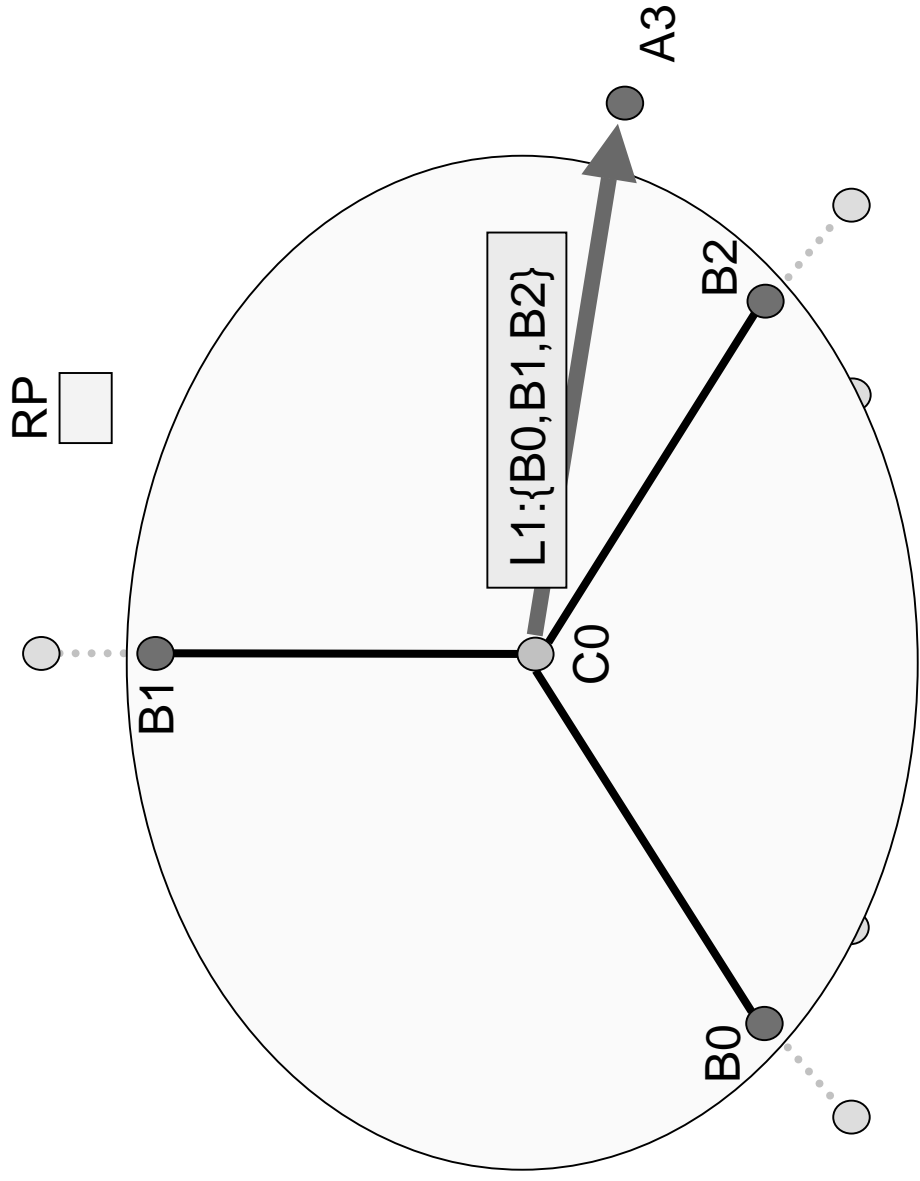


# JoinProcedure

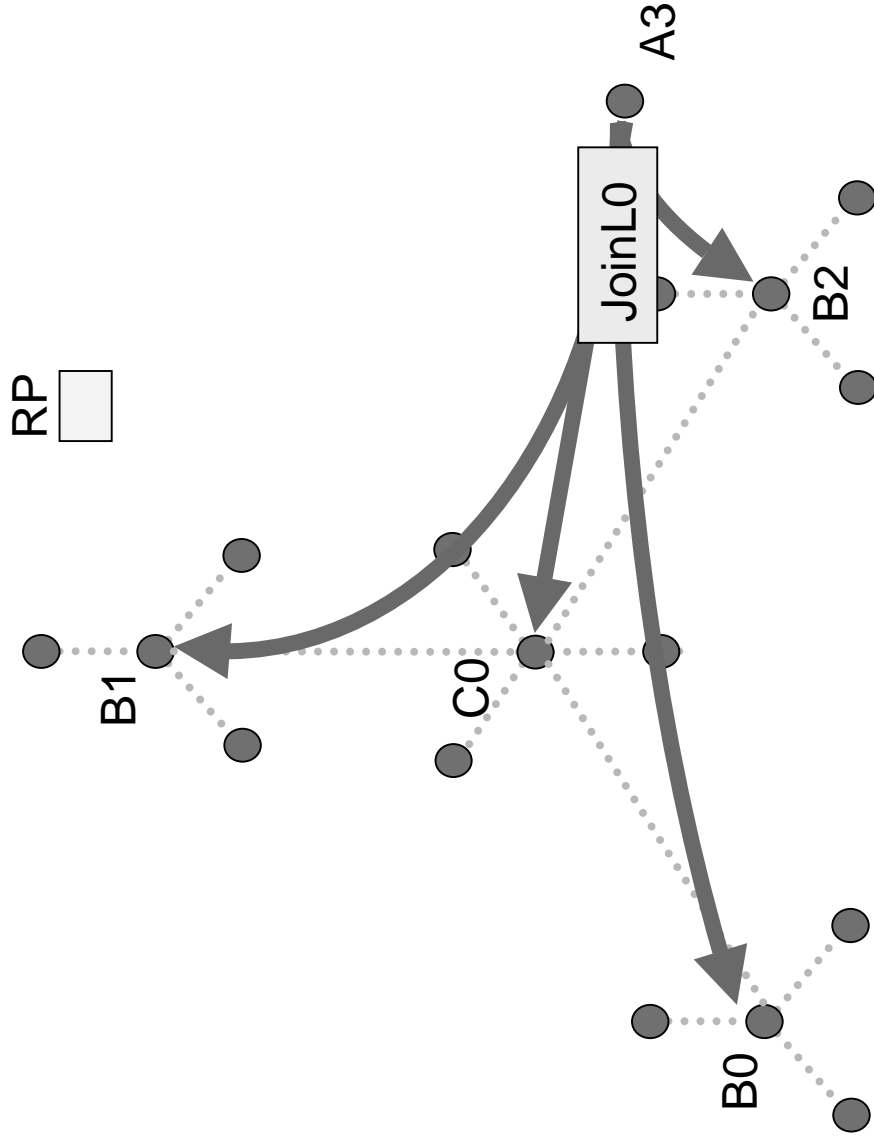




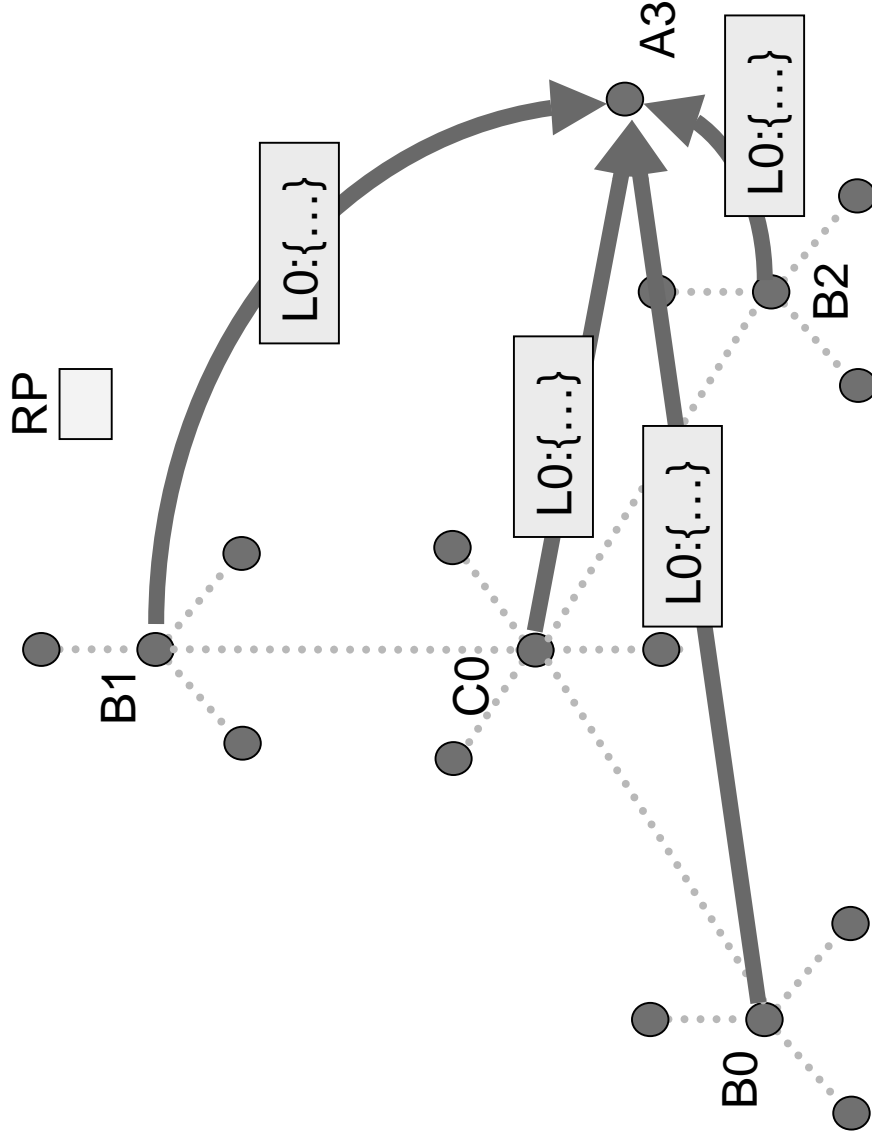
# JoinProcedure



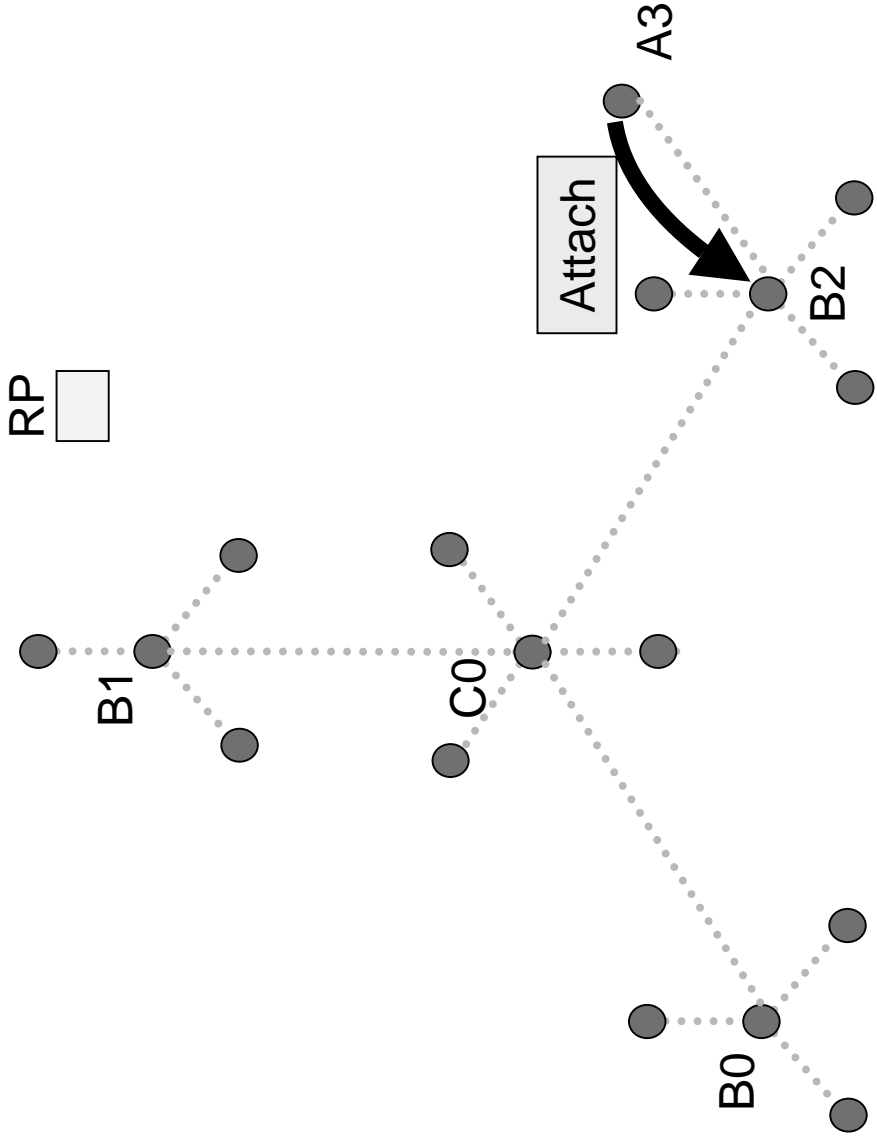
# JoinProcedure



# Join Procedure



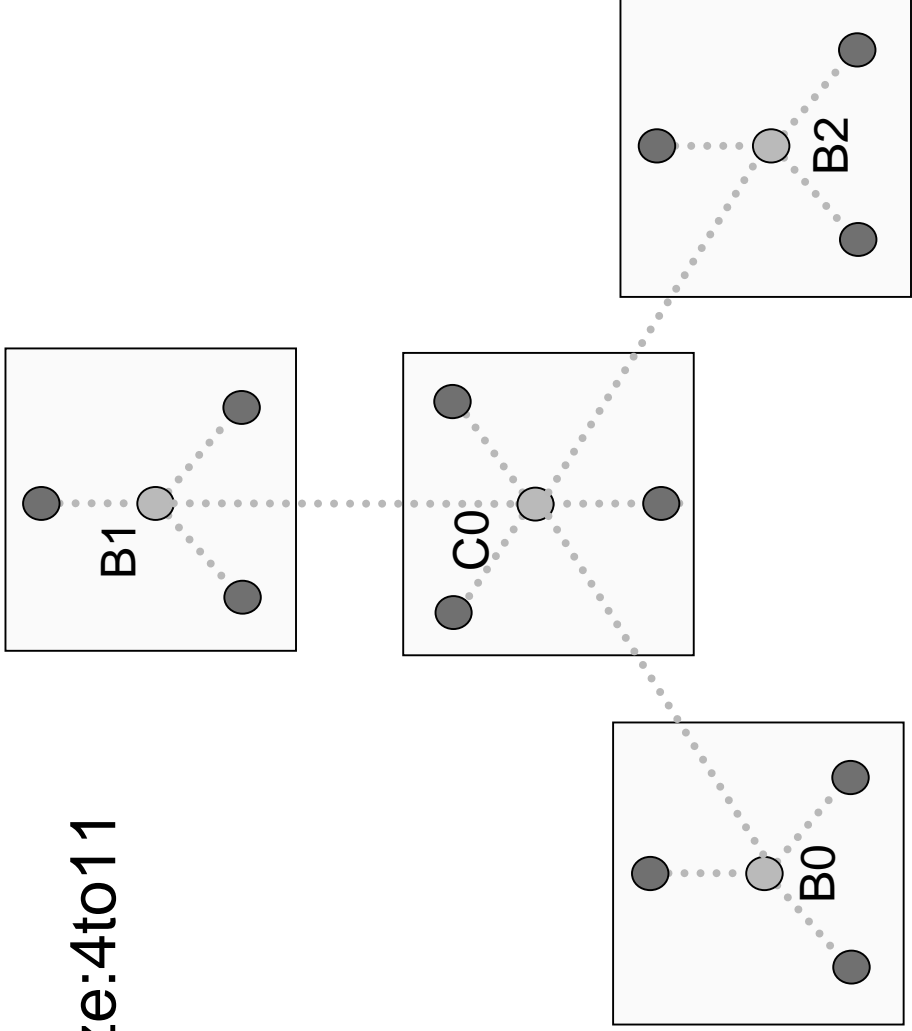
# Join Procedure



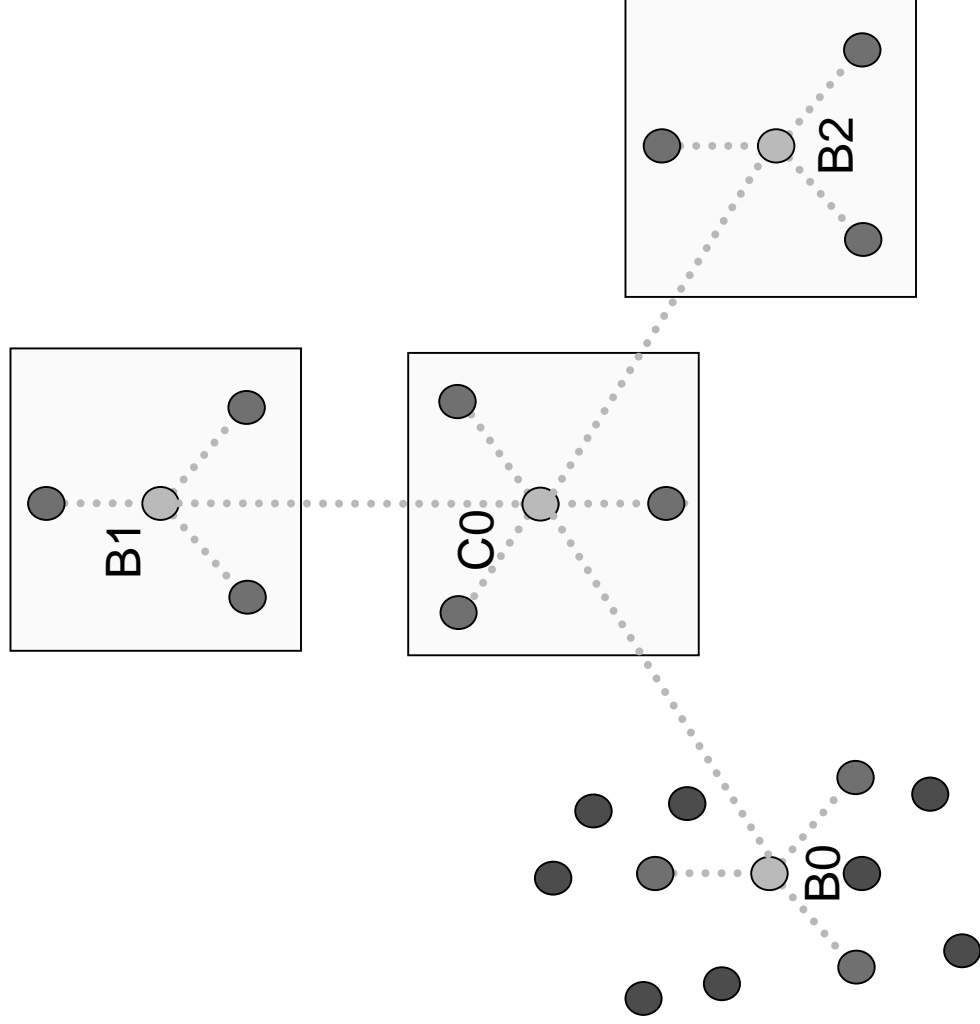
- Overhead:  $O(\log N)$  RTTs and  $O(\log N)$  messages
- Optimizations possible

# ClusterSplit

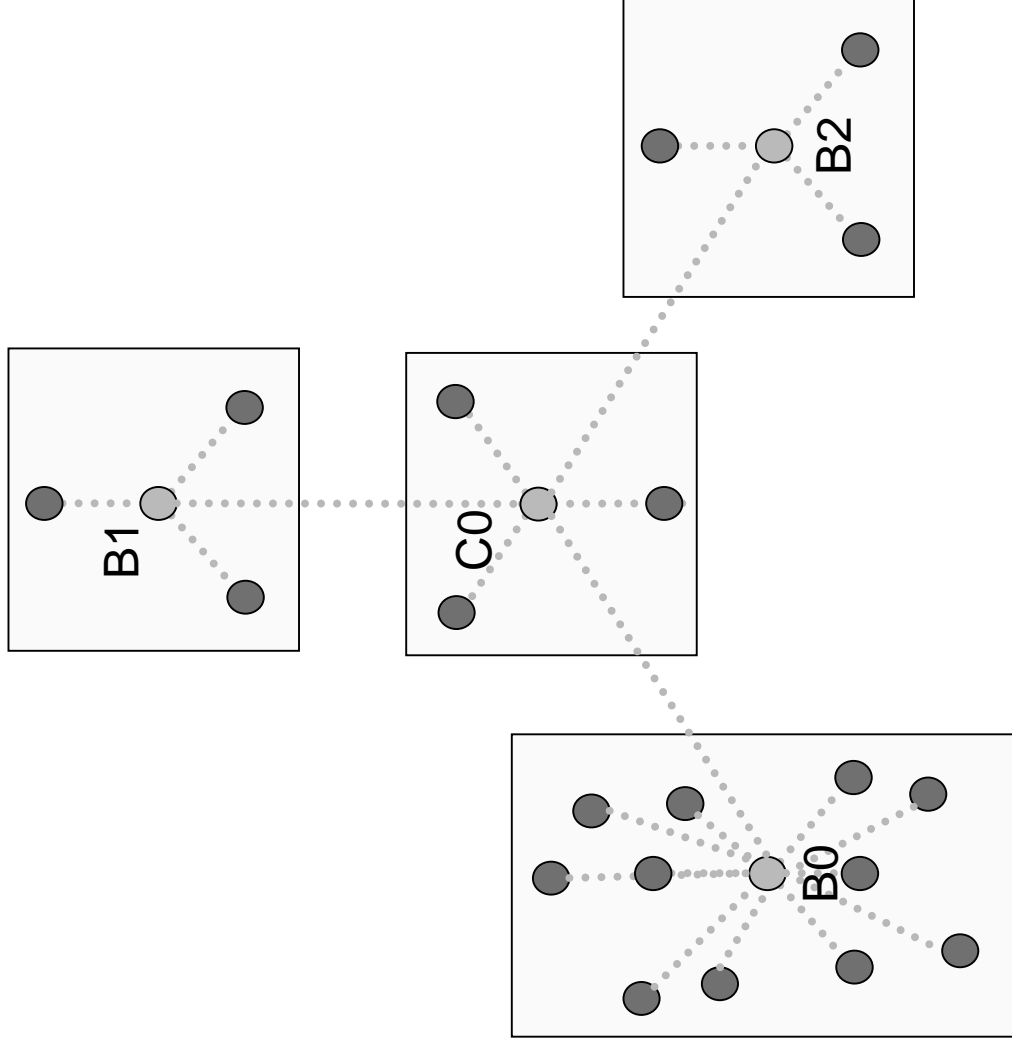
Clustersize:4to11



# ClusterSplit

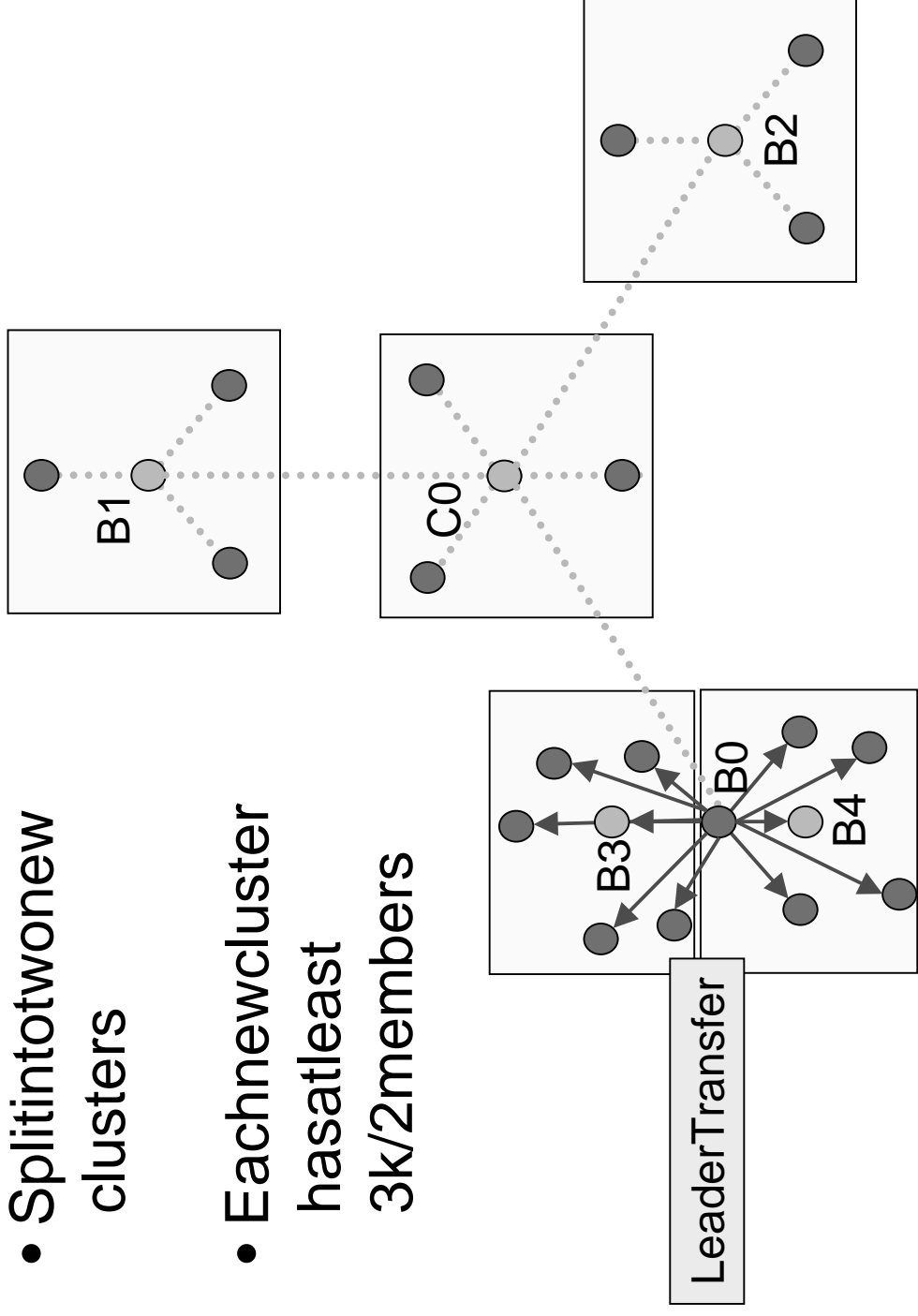


# ClusterSplit



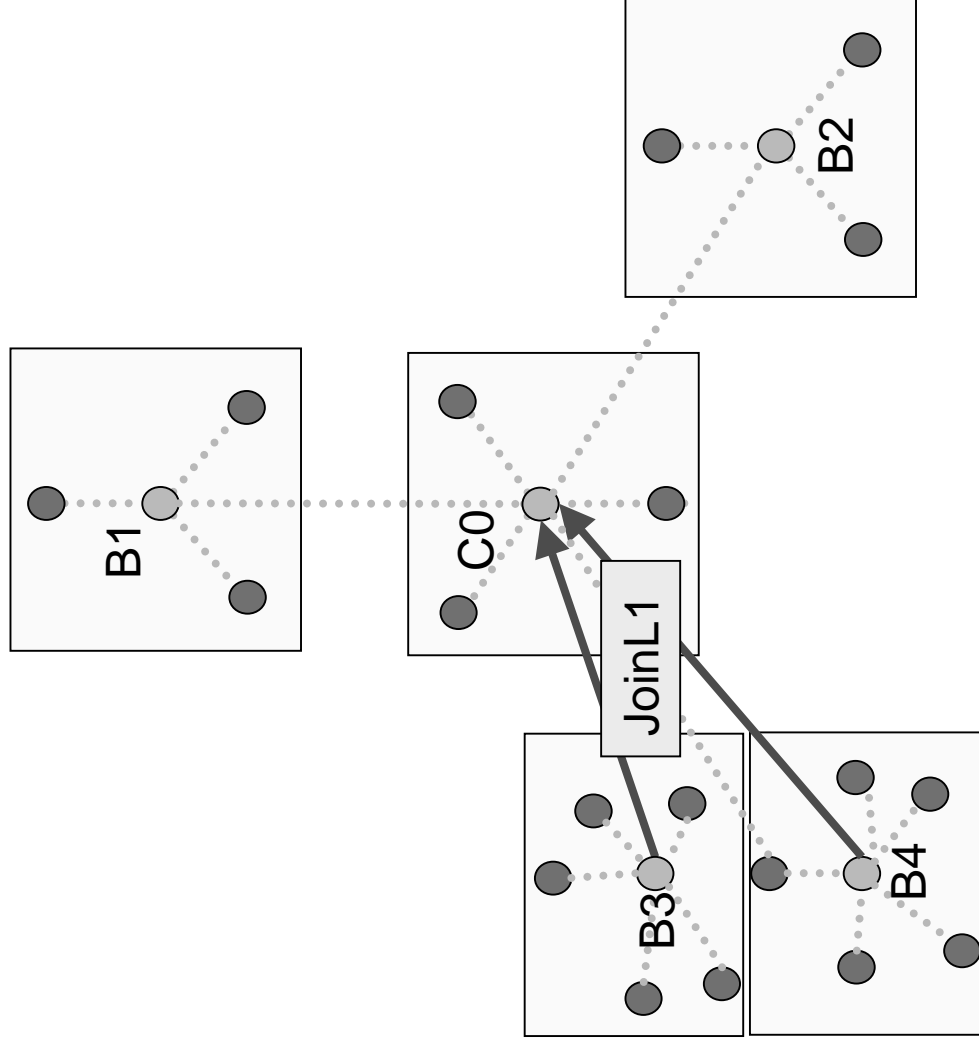
# ClusterSplit

- Split into two new clusters
- Each new cluster has at least 3k/2 members

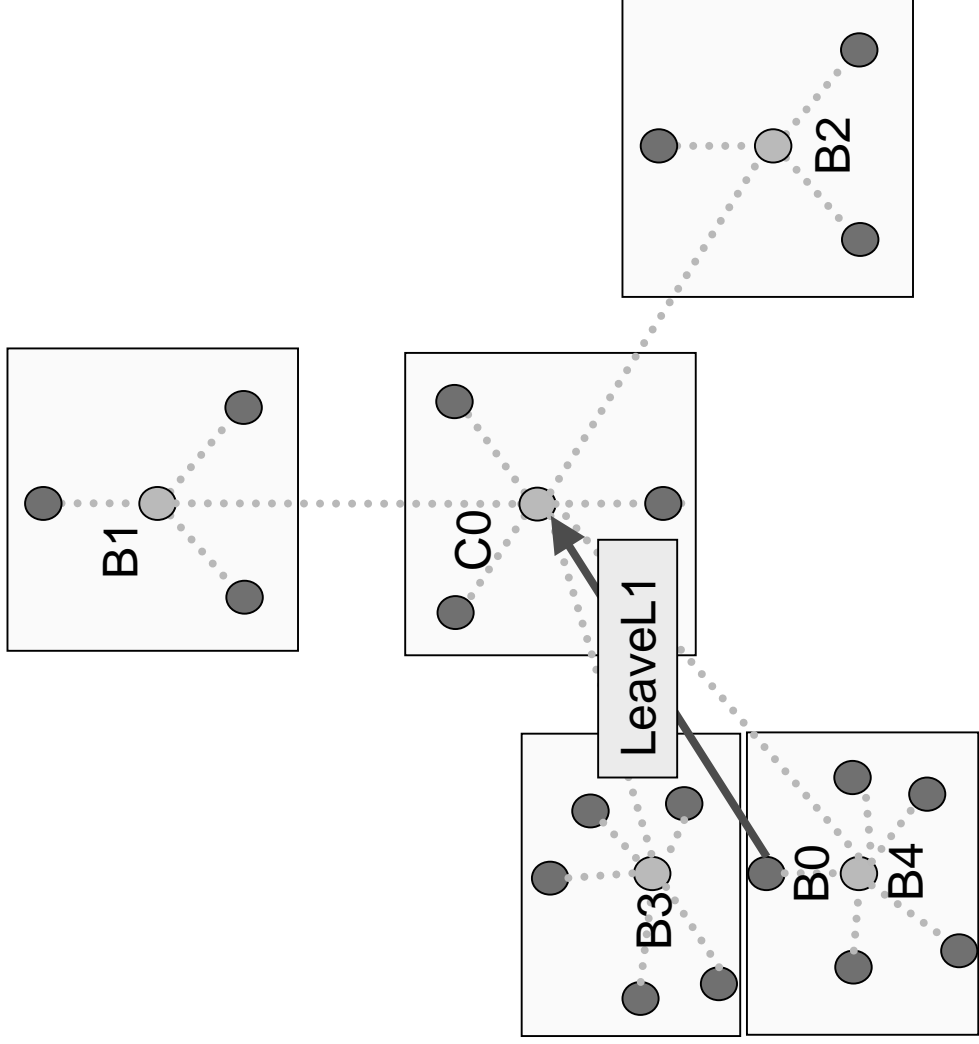




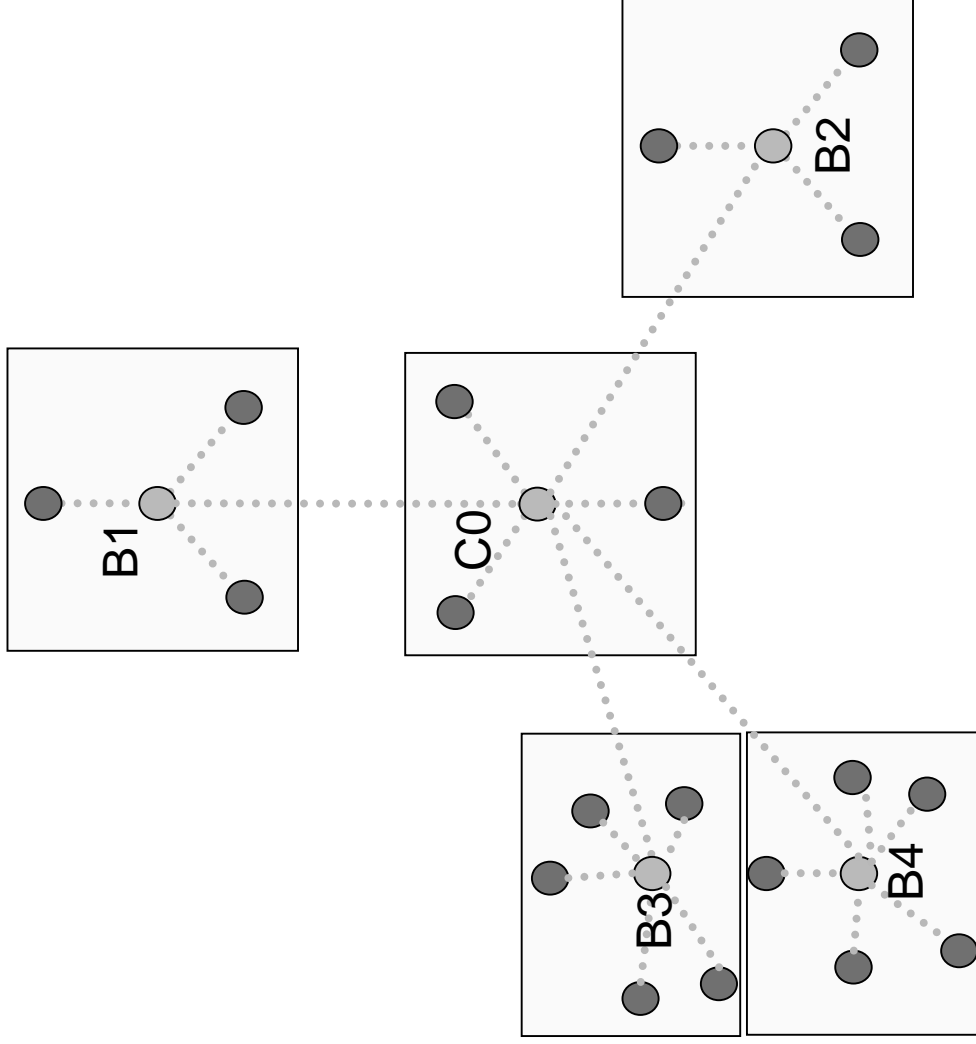
# ClusterSplit



# ClusterSplit



# ClusterSplit

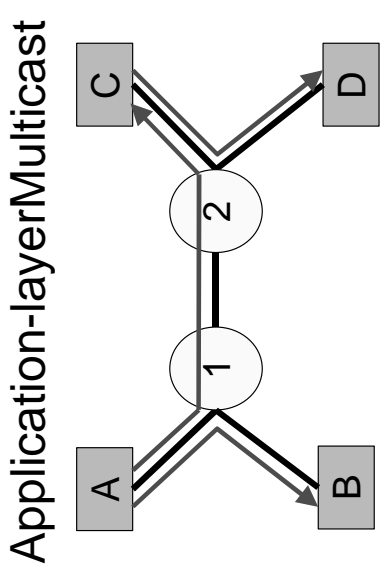


# Results

- Simulations
  - 10,000nodeTransit -Stubgraphs
  - Groupsizes upto 2048
  - Comparisonswith Narada [CMU]
- Wide-areaExperiments
  - Membersat8sites
  - Groupsizes upto 96
- Dynamicjoinsand(ungraceful)leaves
- Constantratedatasource

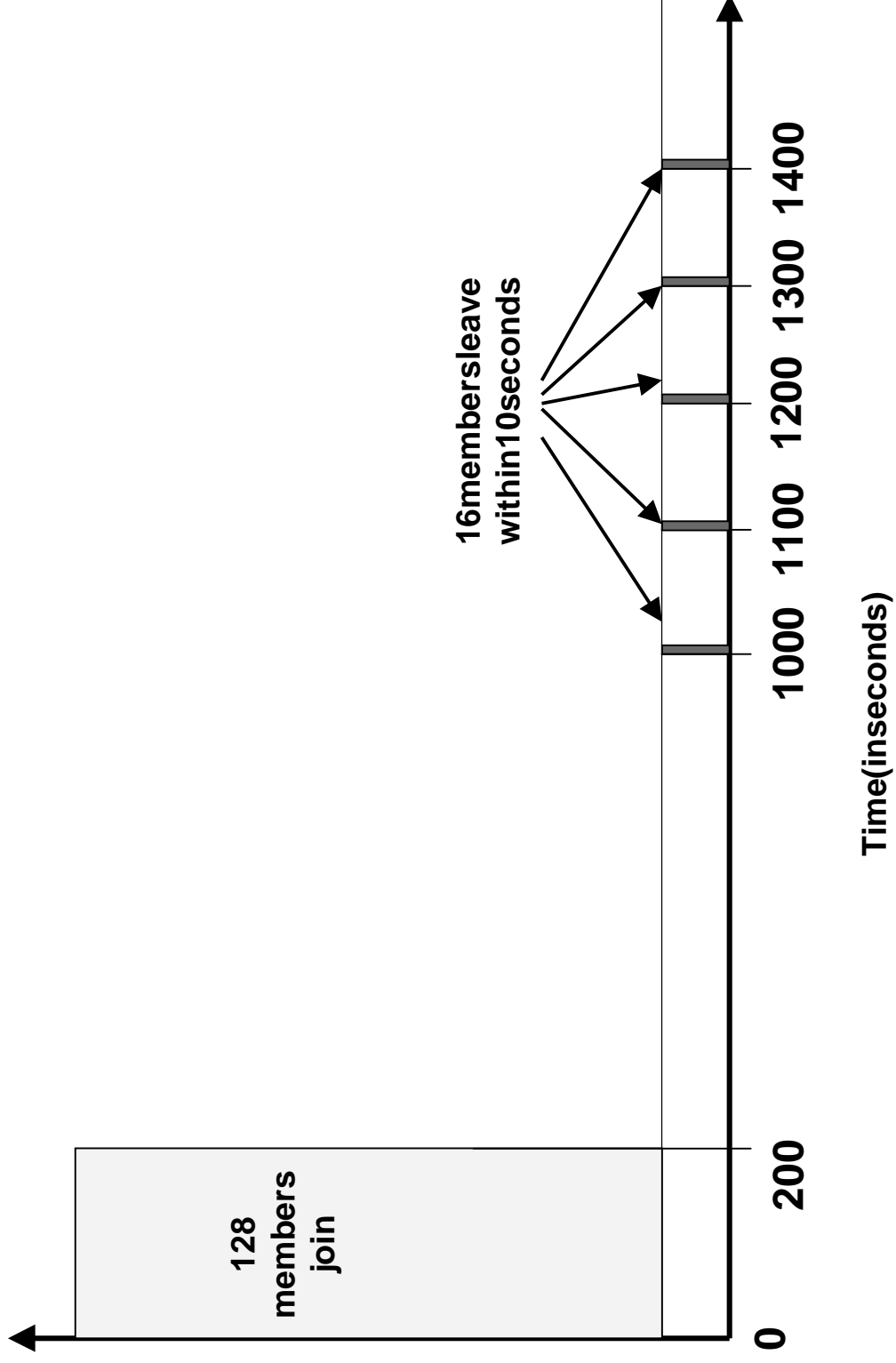
# Evaluation Metrics

- **TreeQuality:Stress**
  - Number of copies of the same data packet on a link/router
  - Example: Stress on link[A -1]=2
- **TreeQuality:Stretch**
  - Ratio of the overlay latency to the direct unicast latency
  - Example: Stretch for receiver D=5/3



- **State at end -hosts**
  - Control overheads
- **Robustness**
  - Host failures

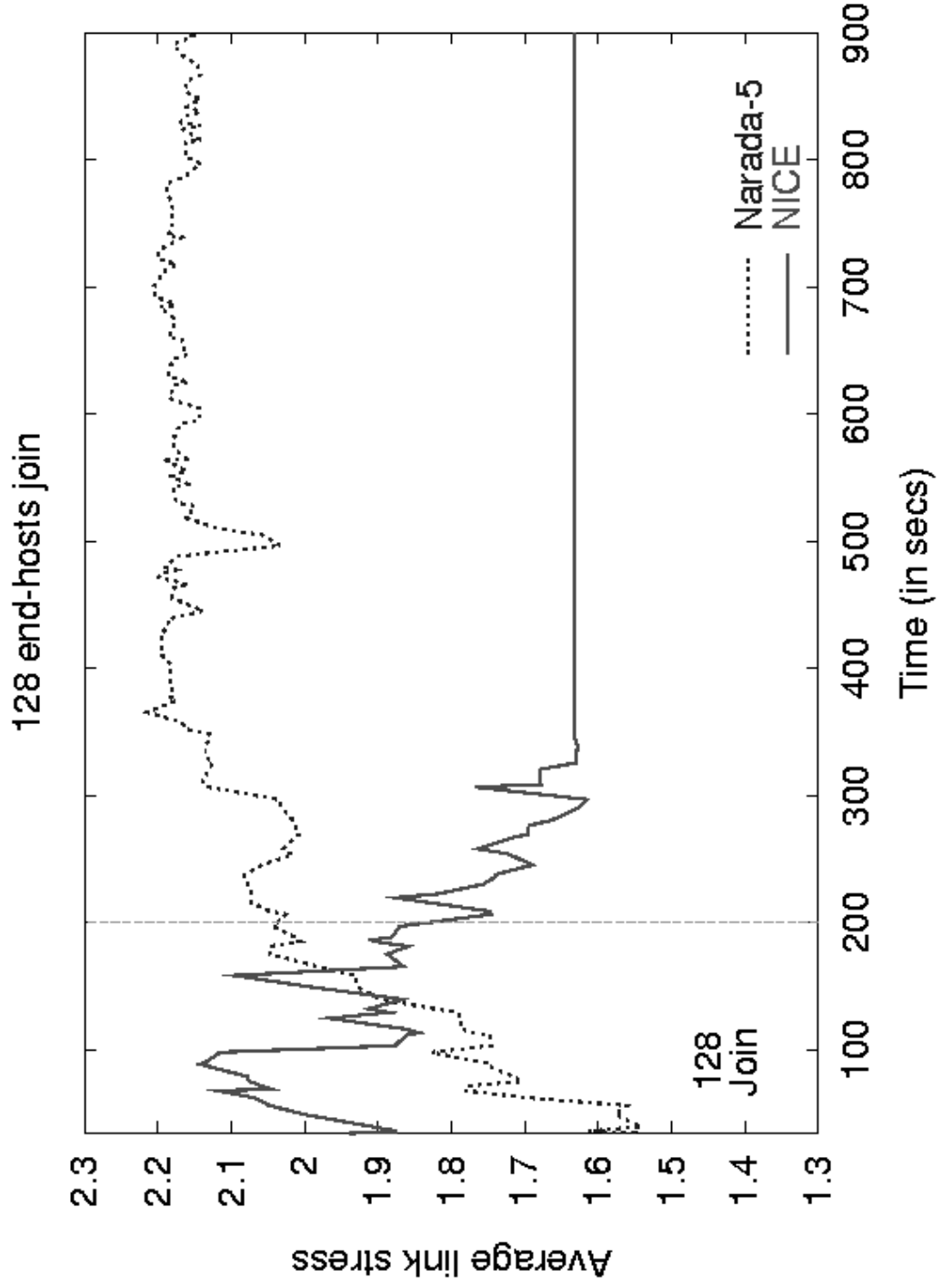
# Example Scenario



# TreeQuality: Stress

First200seconds...

Resource  
usage at  
links

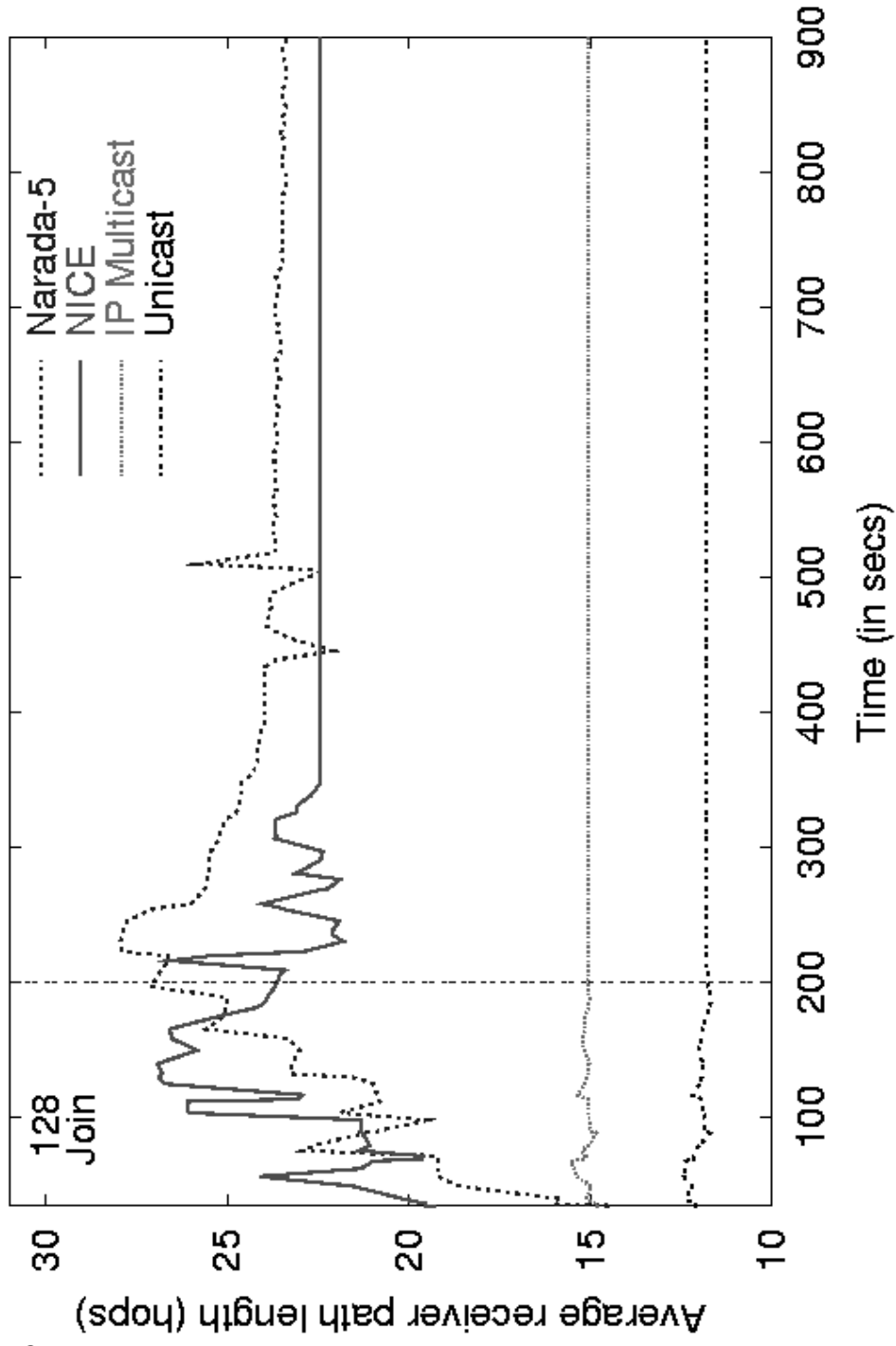


# TreeQuality: Stretch

First200seconds...

End-to-end  
latency to  
receivers

128 end-hosts join

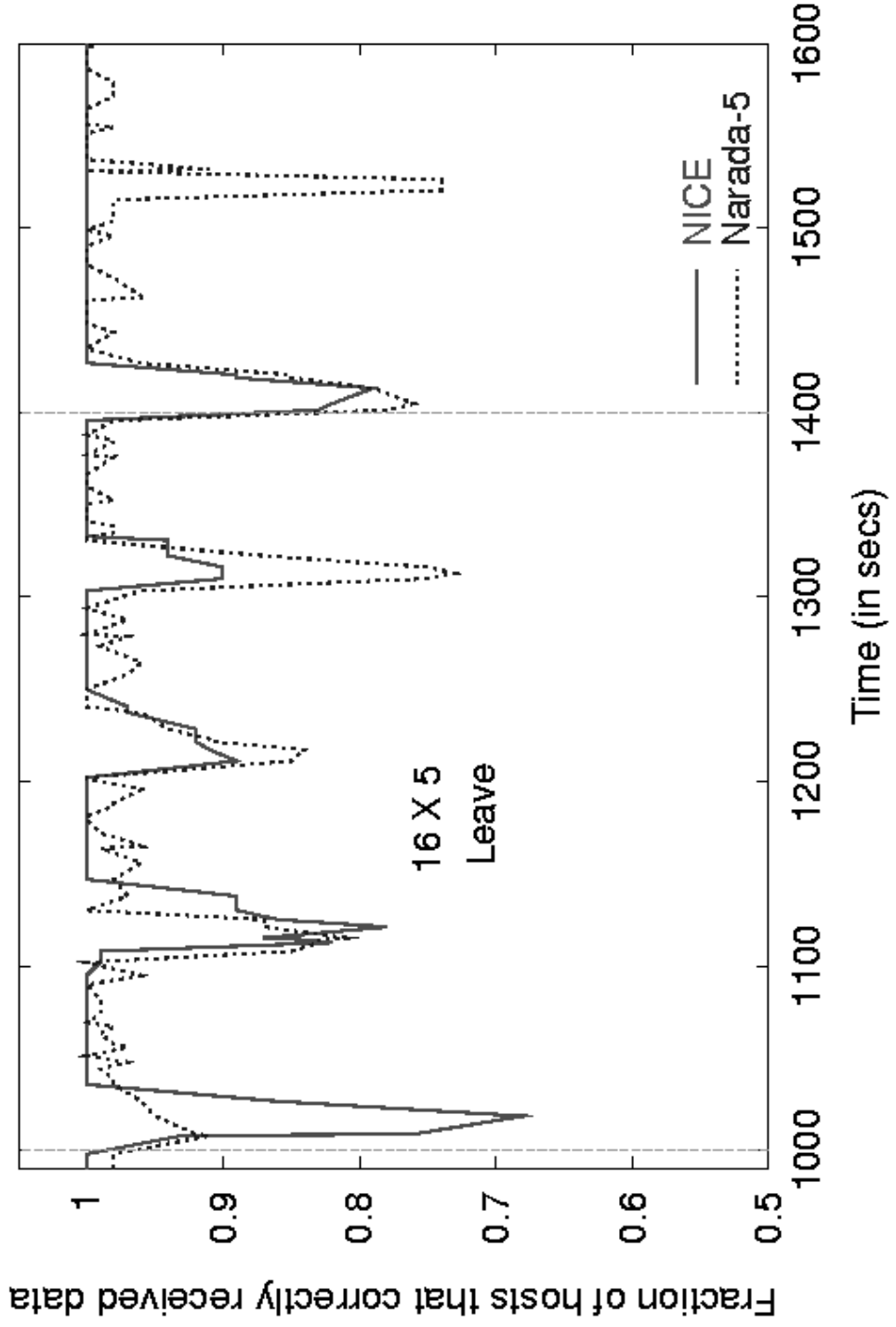




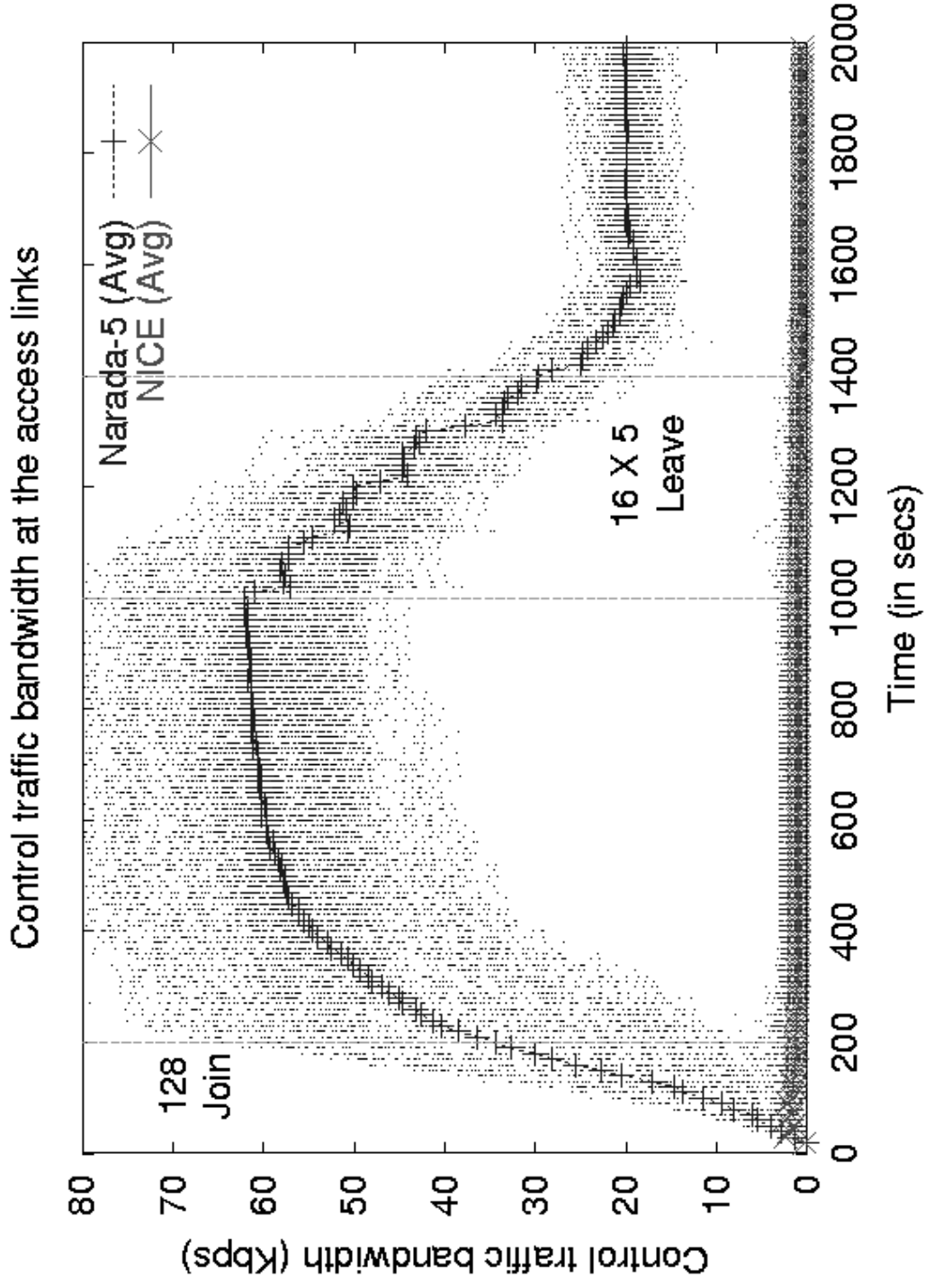
# FailureRecovery

After 1000 secs

Periodic leaves in sets of 16



# ControlOverheads



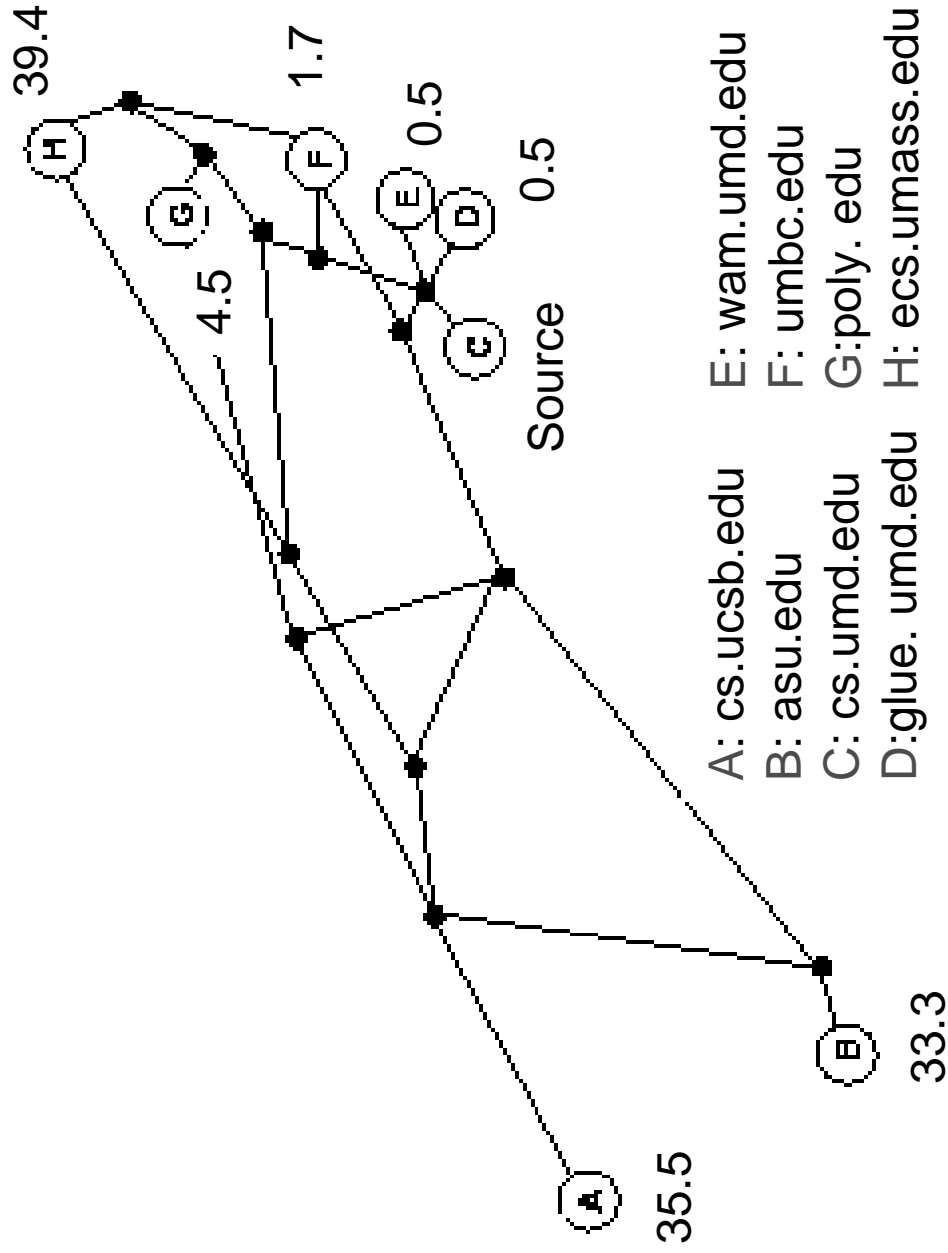
# Control Overheads

Group Size	Narada-30	NICE
32	<b>9.23</b>	<b>1.03</b>
128	<b>65.62</b>	<b>1.19</b>
512	<b>199.96</b>	<b>1.93</b>
1024	-	<b>2.81</b>
1560	-	<b>3.28</b>
2048	-	<b>5.18</b>

Bandwidth overheads averaged over all network routers

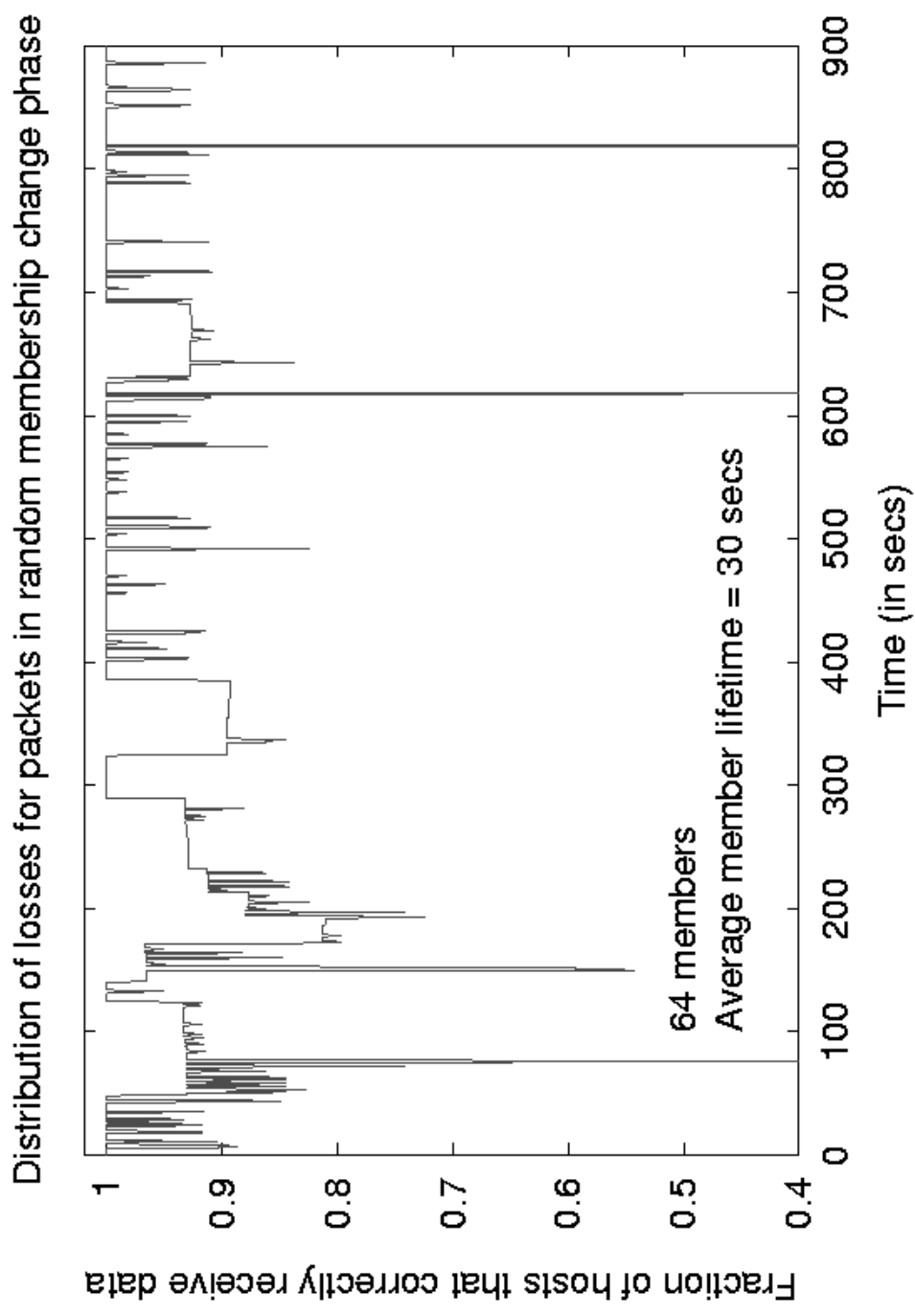


# Wide-area Testbed



# FailureRecovery

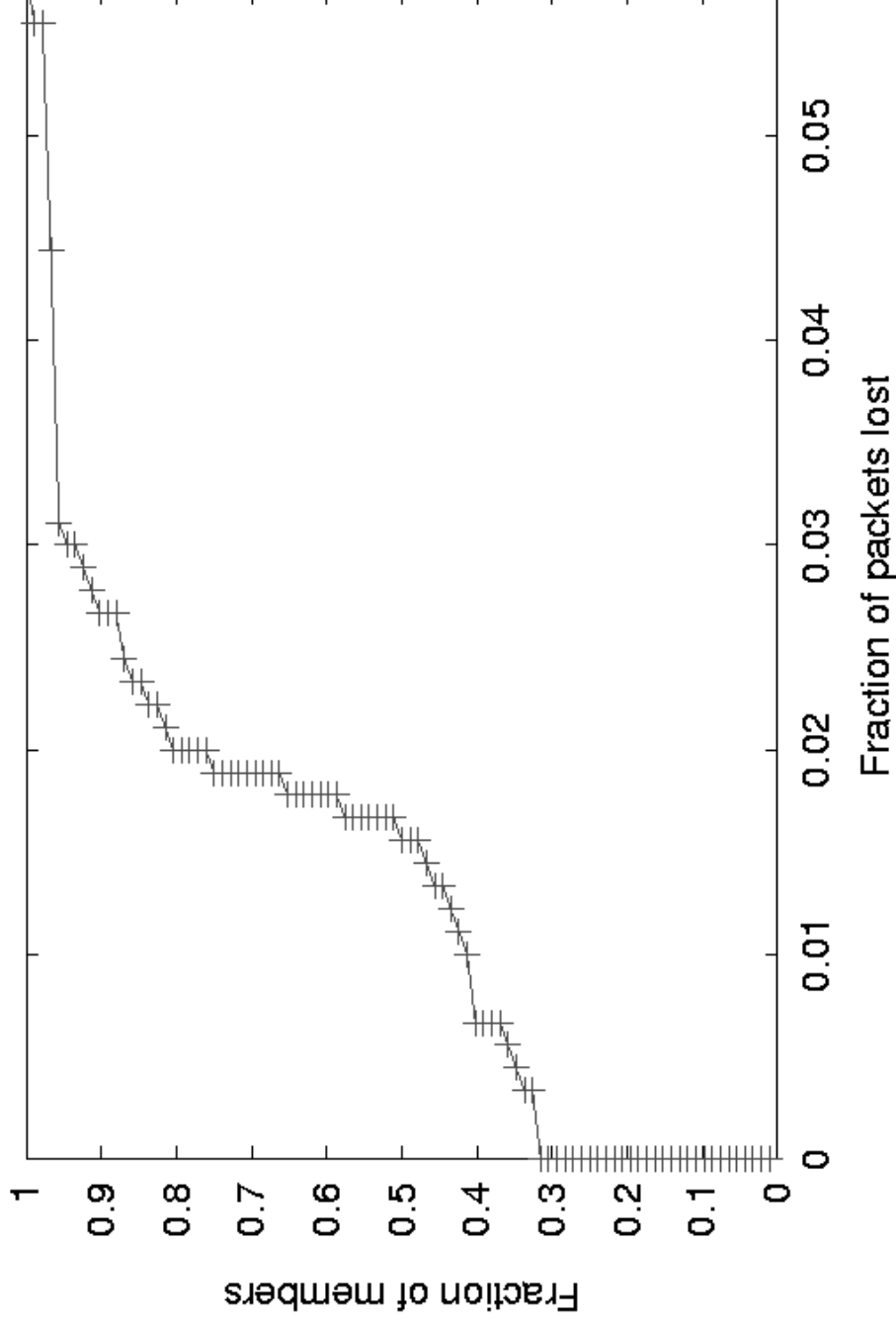
Includes the effect of network losses



# Failure Recovery

Includes the effect of network losses

Cumulative distribution of losses at members in random membership change phase



# RelatedWork

- Mesh-first
  - Narada, Gossamer
- Tree-first
  - Yoid, HMTP
- Implicit
  - Scribe, Bayeux, CAN -multicast, Delaunay-Triangulation

“AComparativeStudyofApplicationLayerMulticastProtocols ” ,

S. Banerjee andB. Bhattacharjee

- Availableat:<http://www.cs.umd.edu/~suman/publications.html>



# Current Work

- Detailed analysis of tree quality
  - Stress and stretch
- Implementing applications
  - Video delivery



# Conclusions

- NICE scale to large member groups
  - Low control overhead
  - Does not sacrifice tree quality or robustness
- Scalability using hierarchy

<http://www.cs.umd.edu/projects/nice>

