

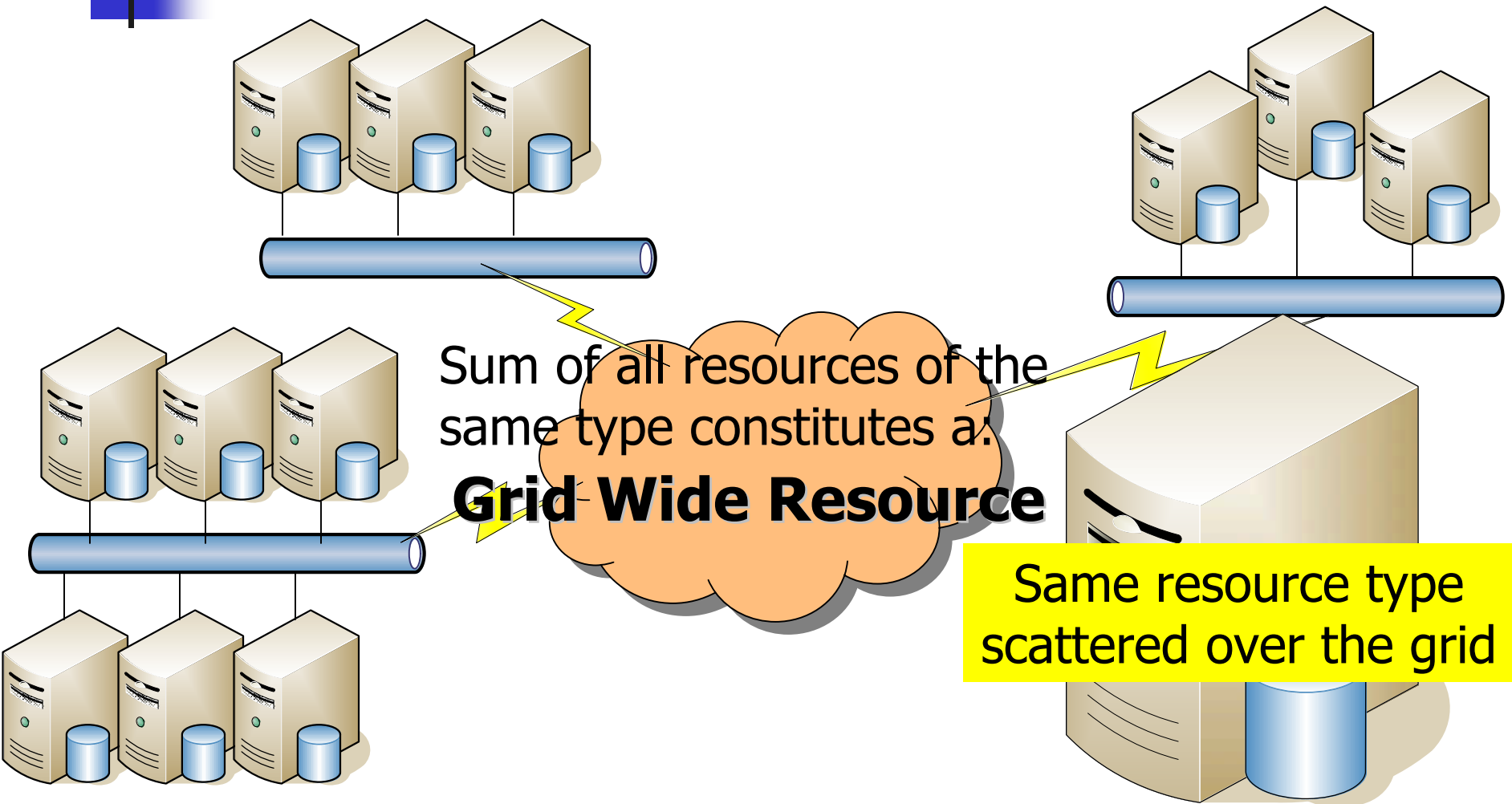


GW_iQ-P: An Efficient, Decentralized Grid-Wide Quota Enforcement Protocol

Kfir Karmon, Liran Liss and Assaf Schuster
Technion – Israel Institute of Technology

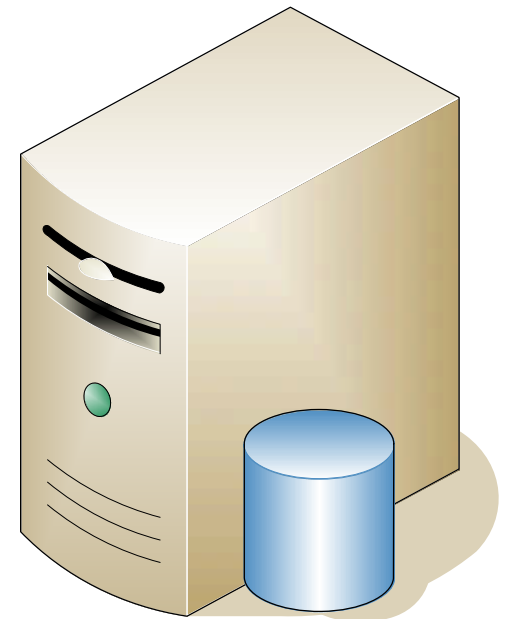
SYSTOR 2007
IBM HRL, Haifa, Israel

Background – Grid Resources



Background – Grid Resources

- Grid Wide Resources
 - CPU hours
 - Disk space
 - DB Connections
 - Outbound traffic
 - Concurrent number of CPUs
 - Allocated RAM
 - Floating Software Licenses
 - Open sockets
 - Etc...





GWiq Motivation

- A grid wide resource tends to be huge and can be exploited
- Grid Wide Quota Enforcement is vital:
 - Security: Prohibit malicious use
 - Fail Safe: Prevent resource leaks (bugs)
 - Financial: Moderate use per paid share

Centralized GWiQ Enforcement

- Central server holds the GWiQ bounds for each (user, resource) tuple
- Per request, resource usage permits are leased until the GWiQ is exhausted.





Objectives

- We strive for a **Grid Wide Quota** enforcement protocol that is:
 - **Decentralized**: No hotspots, No single point of failure.
 - **Efficient**: Overcome latency caused by grid's physical distribution.
 - **Scalable**: Can handle Mega-Grids



GWiq-P:

Grid Wide Quota enforcement Protocol



GWiq-P: Basic Concept

GWiq Enforcement



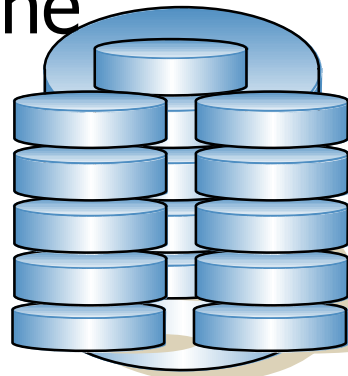
At all times the sum of all local quotas < GWiq

- Using sandboxes to enforce local quotas
- Given a attempt to access the resource:
 - If (**local-quota** \geq **request**) then
 - Grant access
 - **local-quota** = **local-quota** – **request**
 - Else
 - Freeze job execution until **local-quota** reinforced

GWiq-P: Resource Coins

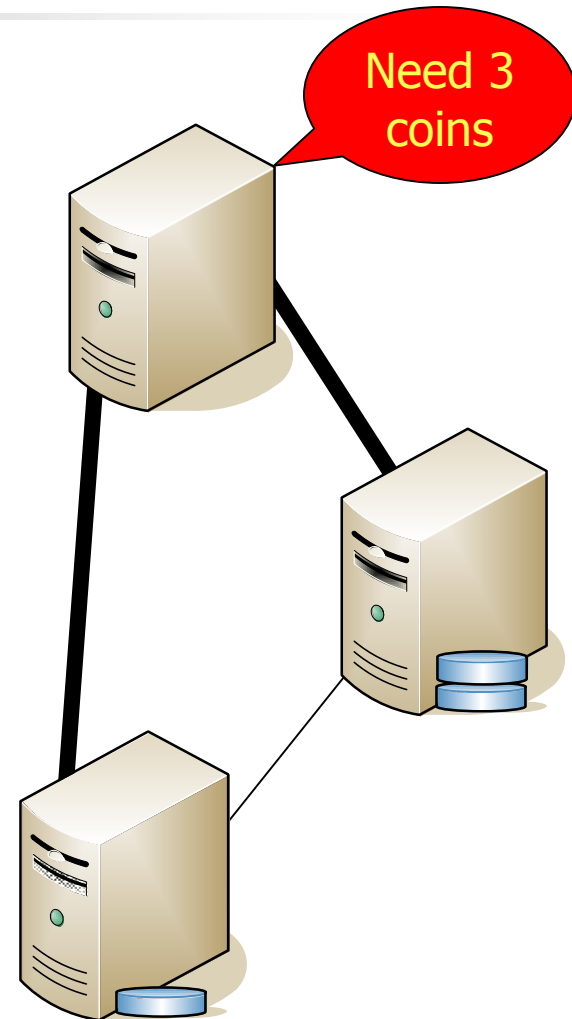
- A **resource coin** denotes the smallest consumable portion of a grid resource.
- Each (user, resource) GWiq is broken down to **coins**.
- A user's job may use the resource up to the amount that the **coins** are worth.
 - i.e. Depositing four 1MB coins grants the job (another) 4MB to use

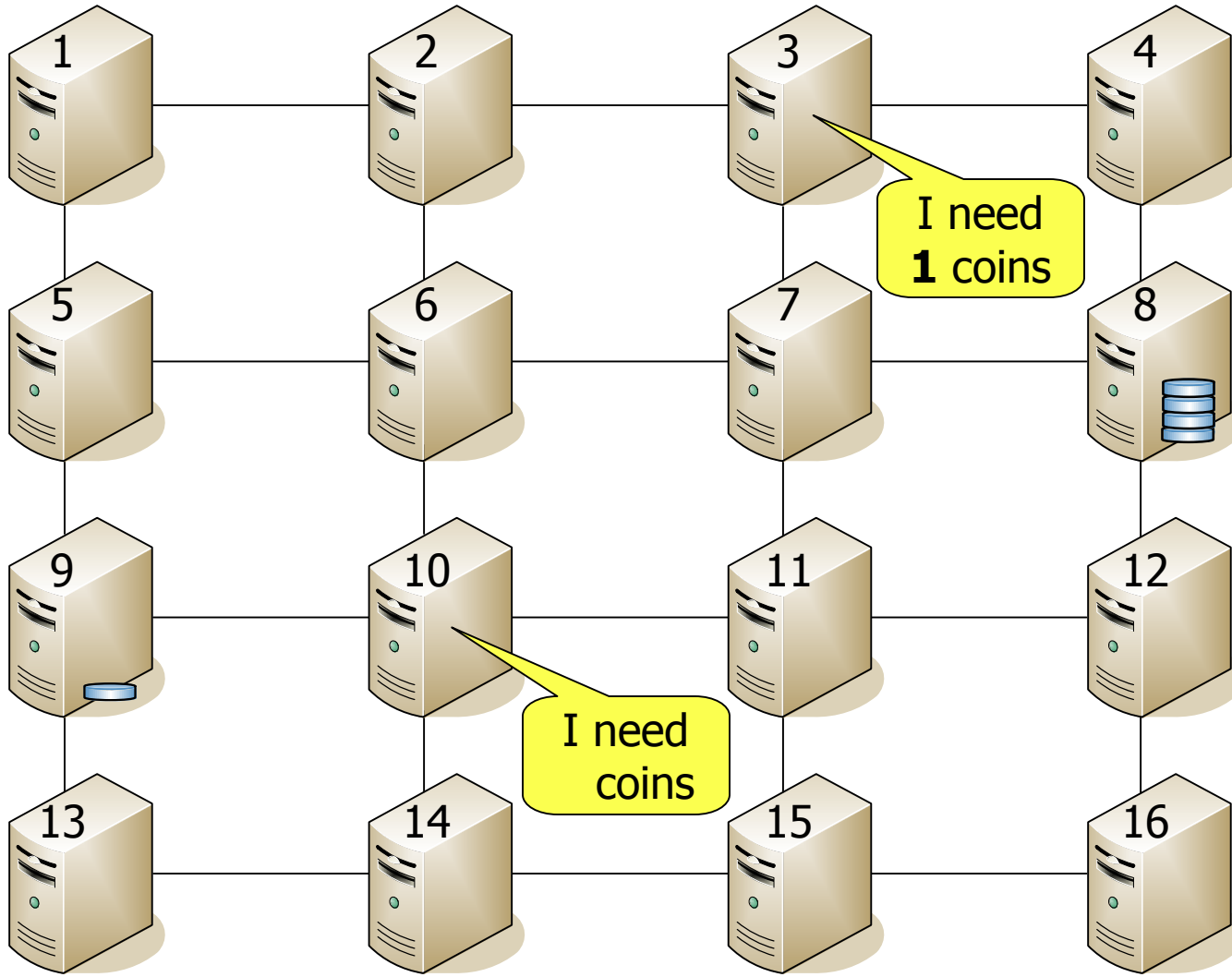
Local Quota = Hosting SBox's resident coins



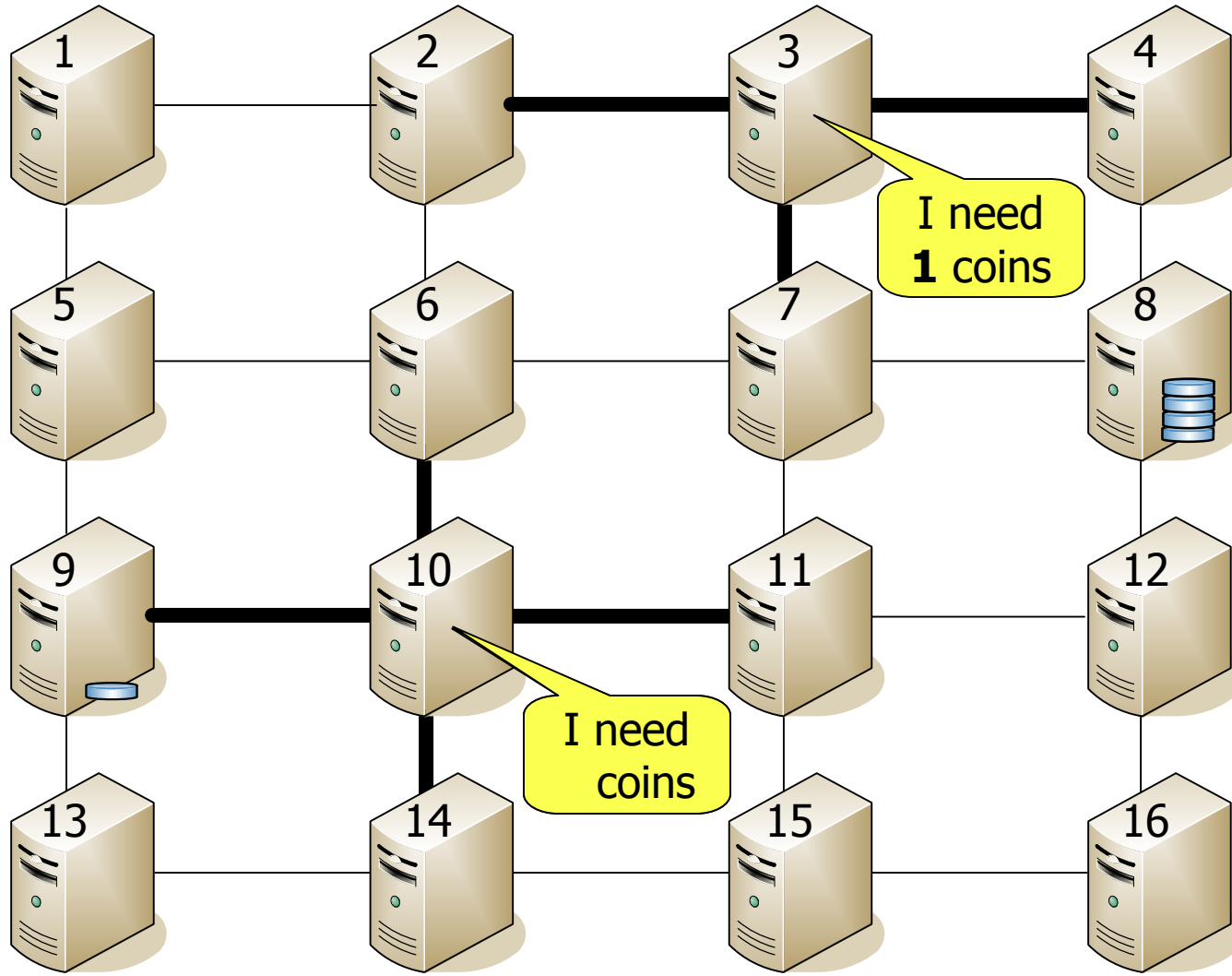
GWiq-P – Spanning Forest

- Using a BF-based alg we build a **spanning forest**.
- A **sandbox** hosting a needy job will start forming a tree around itself.
- At all times, each neighbor will join the tree to which it is closest to its root.
 - Member of one tree at a time.
- Surplus coins will be transferred to the root.

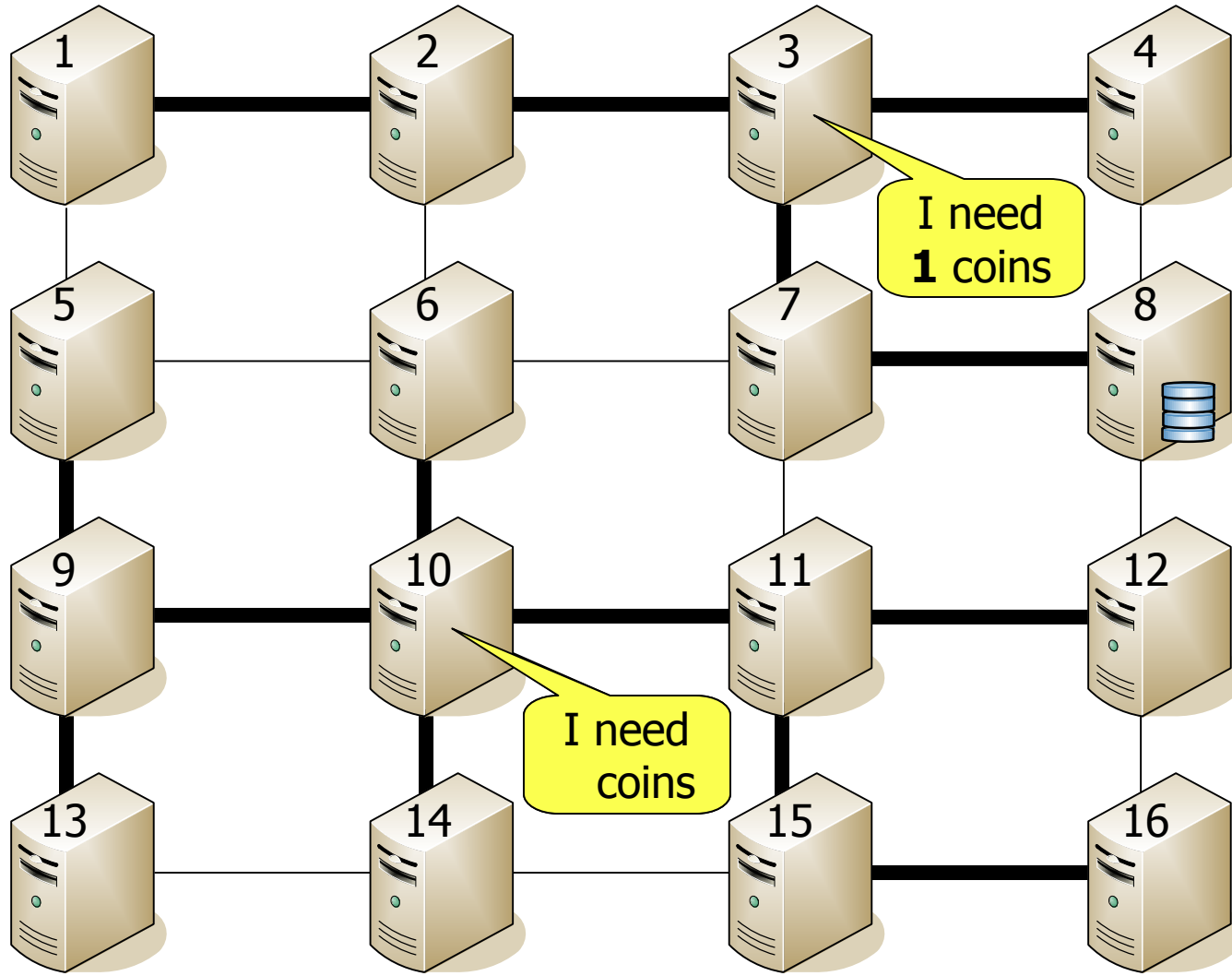




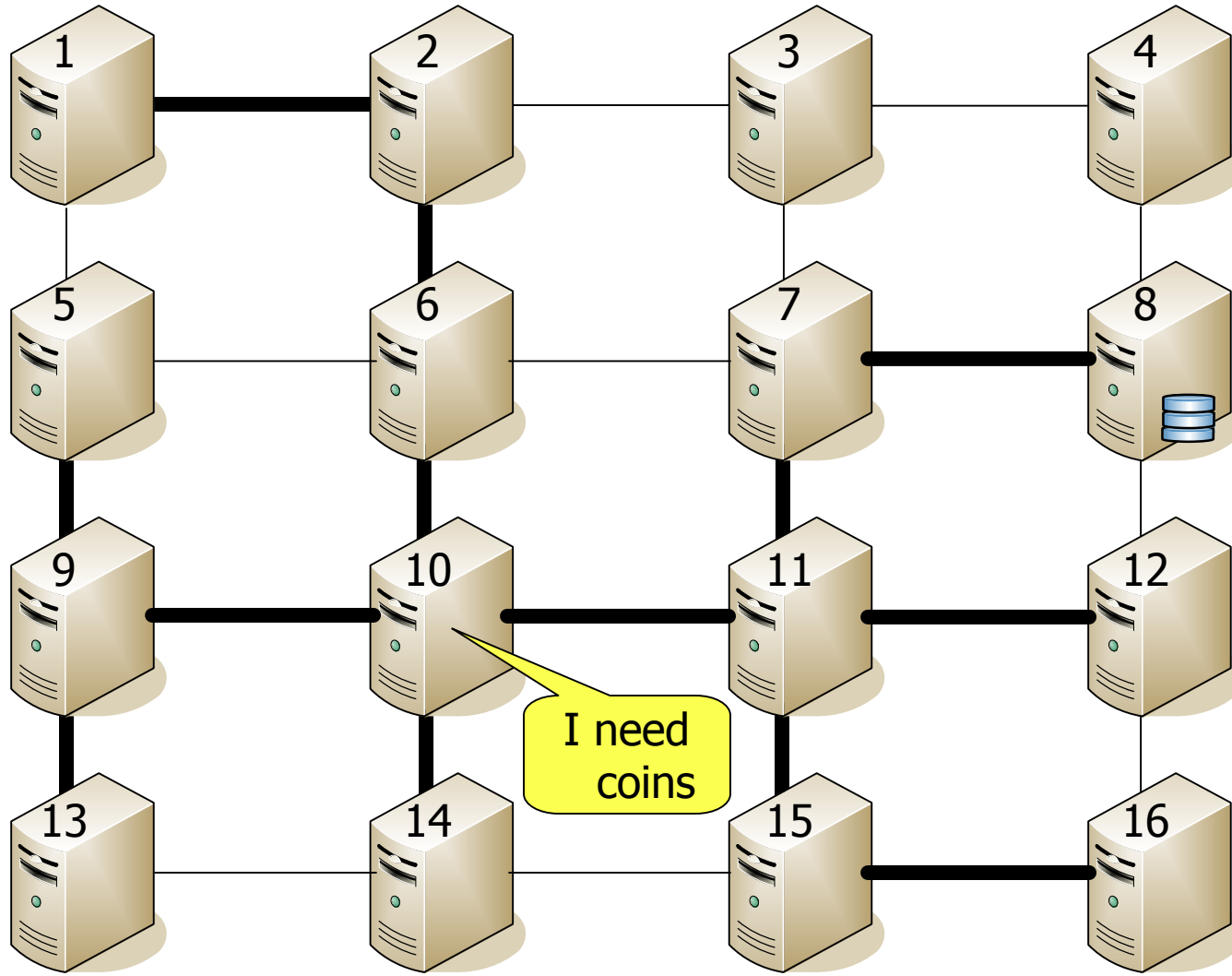
GWiq-P : In action 2/5



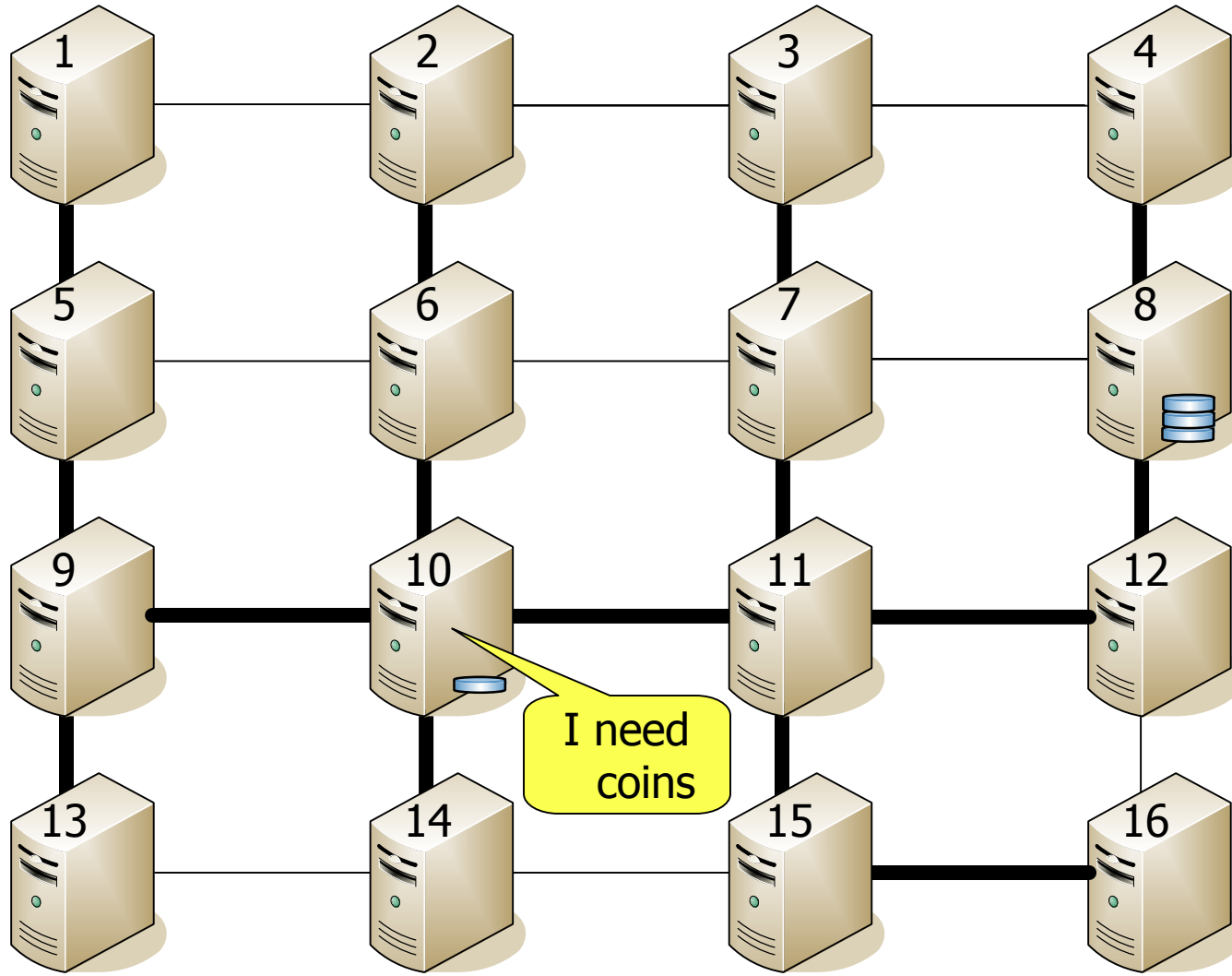
GWiq-P : In action 3/5



GWiq-P : In action 4/5



GWiq-P : In action 5/5





GWiq-P: Fault Intolerance

- Links and nodes may crash
 - Link crash → Passing coins are lost
 - Node crash → Local coins are lost
- Lost coins ⇔ Unfair GWiq reduction
- Can't be sure how many coins are on the wire
 - Did my neighbor send something?
 - Have the coins I sent reached already?
- How do we restore unknown lost coins???



FT Solution1: Transactional

- Coin transfers are done in Transactions
- Pros
 - Coins never lost on transfer
- Cons
 - Slowdown
 - Needs persistent storage
 - Node failure (temporarily?) reduces GWiQ by local quota

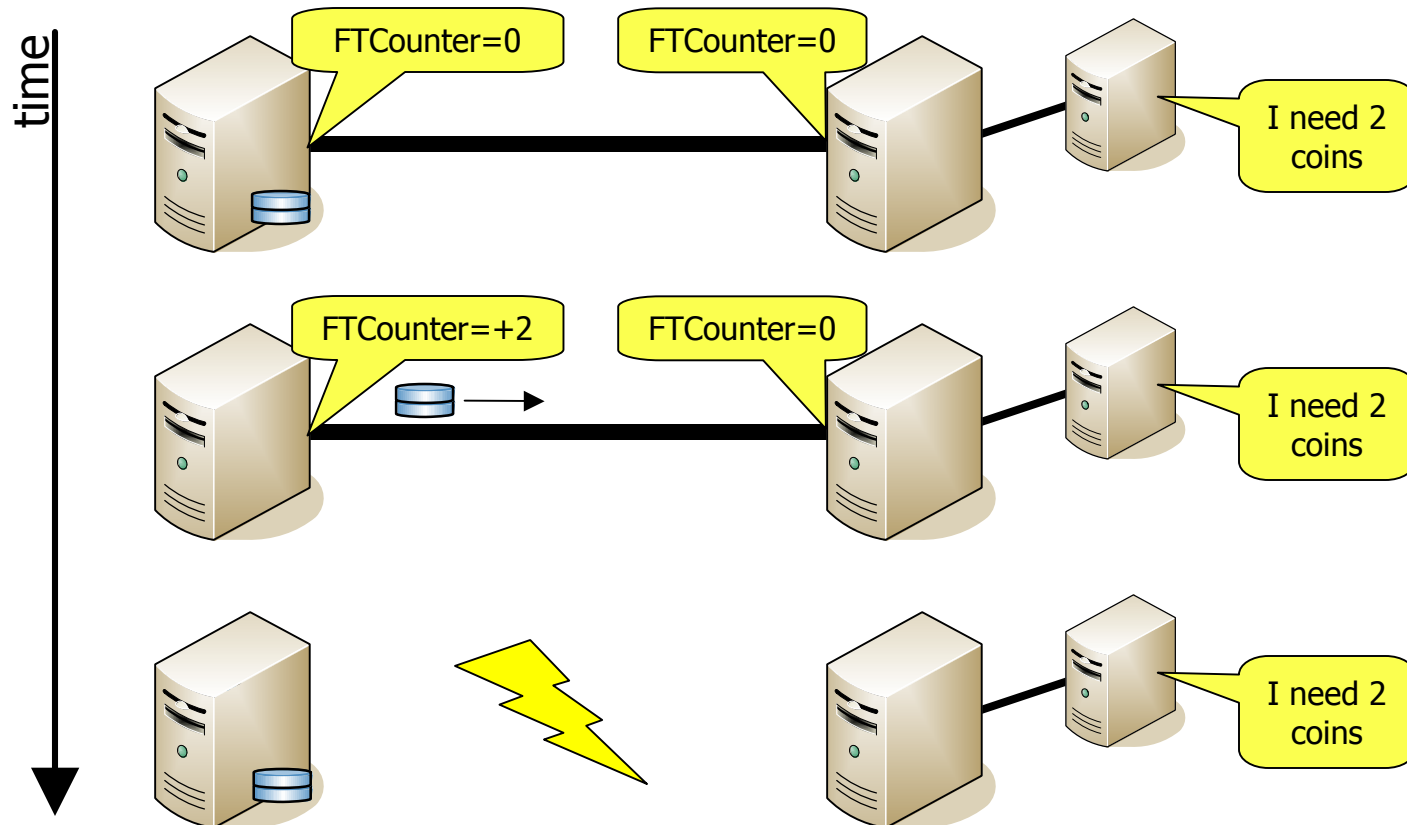


FT Solution2: Light Weight 1/6

- Machine down \Leftrightarrow All of its links are down
- Each node holds a counter for every link.
 - $\text{FTCounter} = \text{Sent coins} - \text{Received coins}$
 - $\text{FTCounter} = \text{Net number of transferred coins}$
- When a link goes down:
 - If $\text{FTCounter} < 0$
 - Create a fictive demand for $|\text{FTCounter}|$ coins
 - Else
 - $|\text{FTCounter}|$ coins are added to the node's surplus

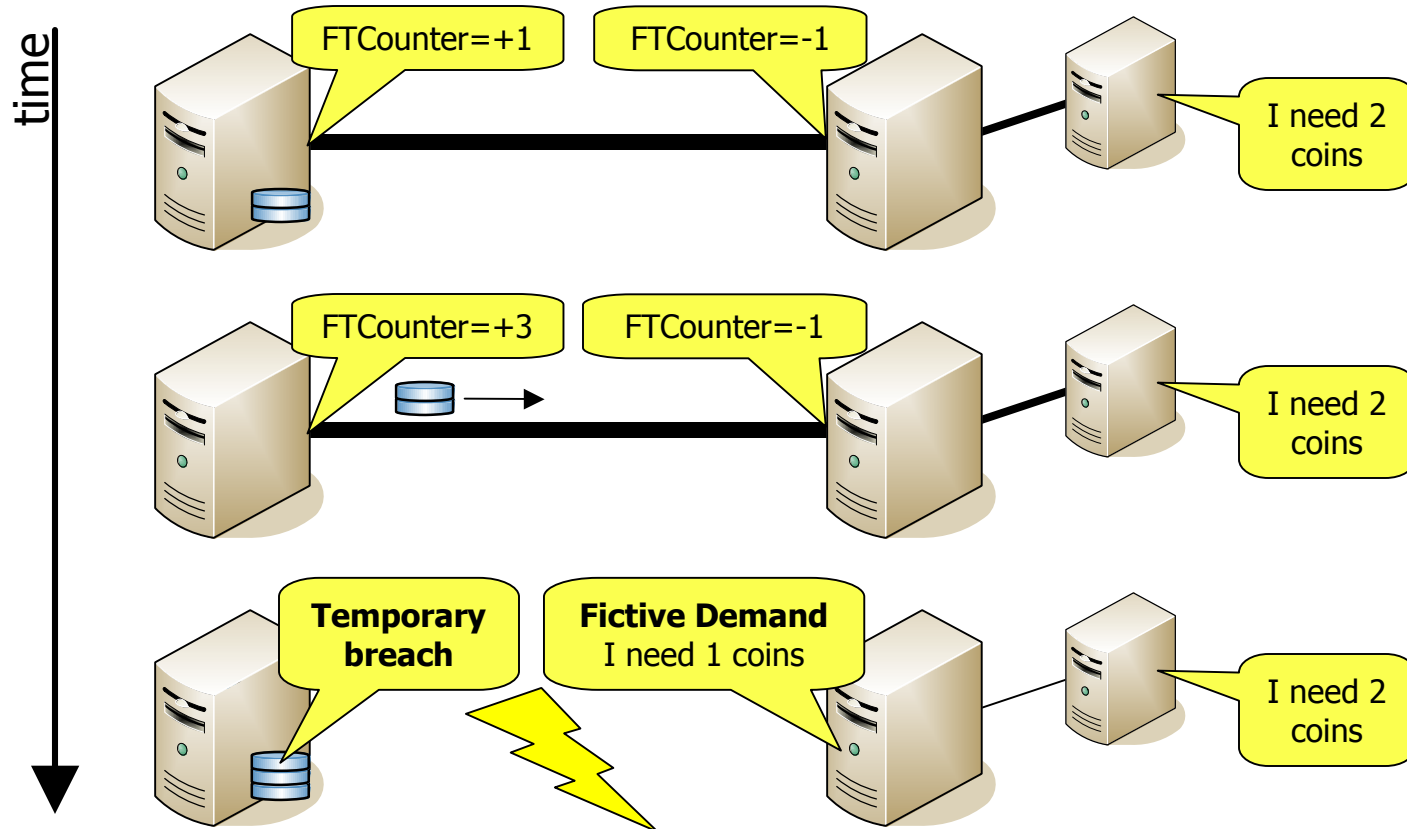
FT Solution2: Light Weight 2/6

- Example 1:



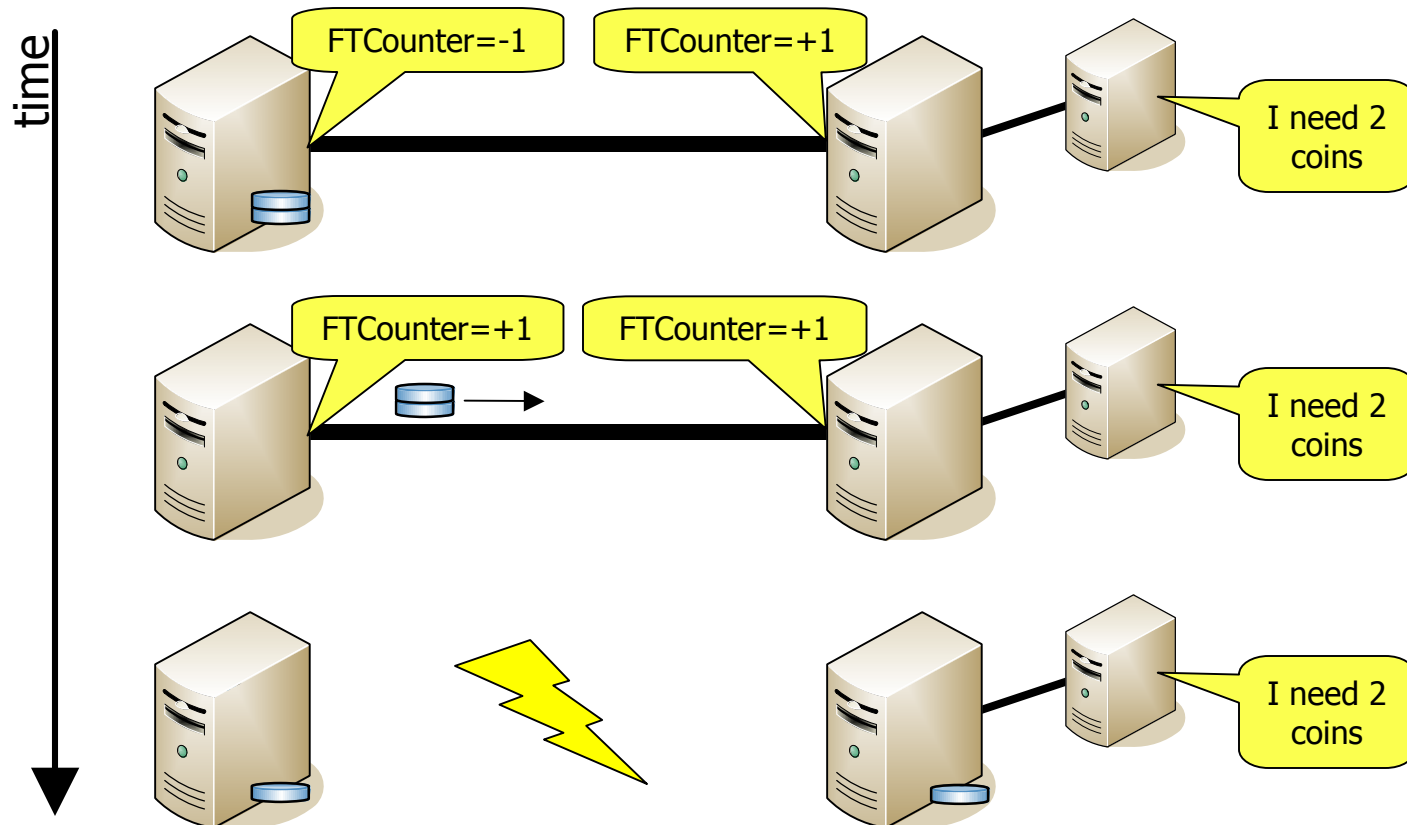
FT Solution2: Light Weight 3/6

■ Example 2:



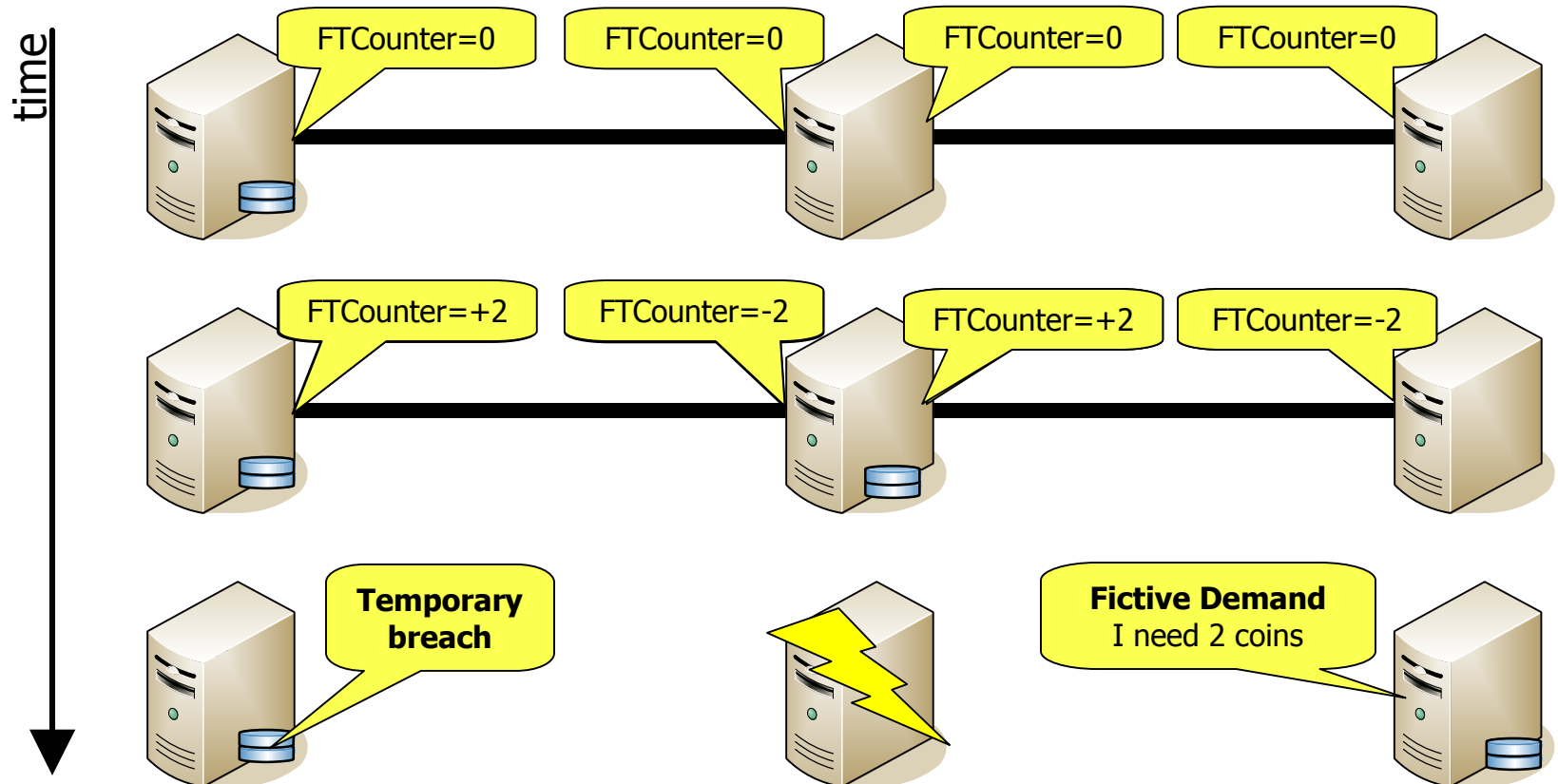
FT Solution2: Light Weight 4/6

■ Example 3:



FT Solution2: Light Weight 5/6

■ Example 4:





FT Solution2: Light Weight 6/6

- Pro

- No latency
- No need for persistent storage
- GWiQ never unfairly reduced

- Con

- May introduce temporary GWiQ breaches



FT Solution3: Hybrid 1/2

- Use the **Light-Weight** solution regularly
- Only con to deal with: GWiQ breaches
 - Caused by loaded FTCounters
- Issue FTCounter balancing **Transactions**
 - If $|\text{FTCounter}_{i,j}| > \text{Threshold}$
 - Disallow link usage. Start transaction.
 - $\text{FTCounter}_{i,j} = \text{FTCounter}_{j,i} = 0$
 - End transaction. Resume link usage.
 - Or, issue periodically.



FT Solution3: Hybrid 2/2

■ Pro

Due to
Transactions

- Reduce breaches' size

- No slowdowns (mostly)

Due to
Light Weight

- No persistence storage (mostly)

- Node crash fully remedied immediately also

■ Con

- See other side of the 'Pro'-coin ;)

- Play with the **tradeoff** using parameterization

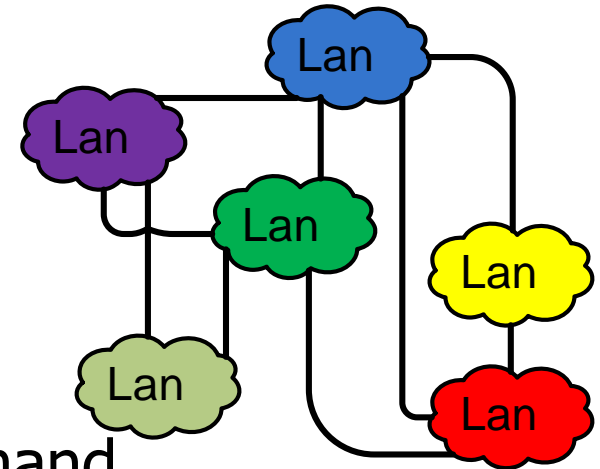


GWiq-P: Properties

- Small trees form around demand
 - Requests are remedied **locally**
- Coins are drawn towards 'hot' areas
 - Auto-Adaptable
- Fully distributed
 - No hot-spots & single points of failure
 - Low latency
 - No Congestion
 - Infinitely Scalable
- Fault Tolerant

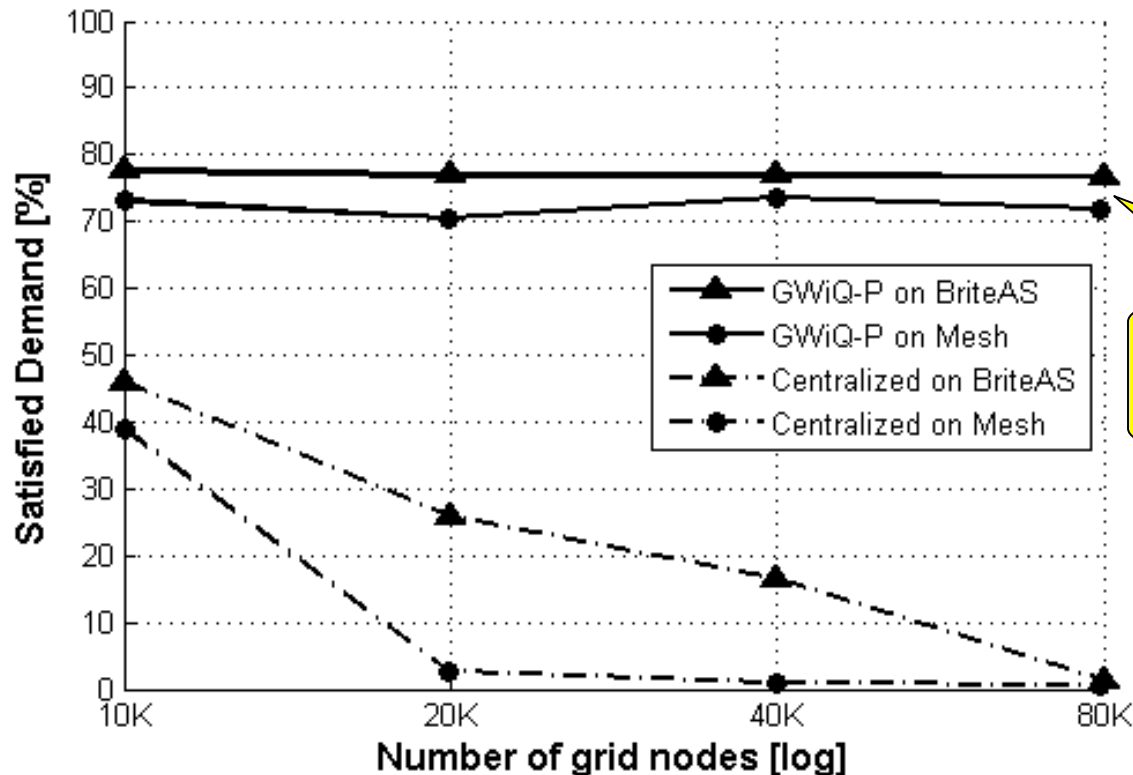
Simulations' (default) Properties

- **Topology** = BriteAS
 - Fast LANs, slower intercon/
- **NetSize** = 10K
- **Q/D=1**
 - Q for GWiQ; D for Overall demand
- **Change Rate** = $1\% * D / E[\text{EdgeDelay}]$
- **Demanders** = $1\% * \text{NetSize}$
- **Fail Rate** = $1\% * \text{Edges} / E[\text{EdgeDelay}]$
 - Applicable for FT scenarios



Simulations 1/6

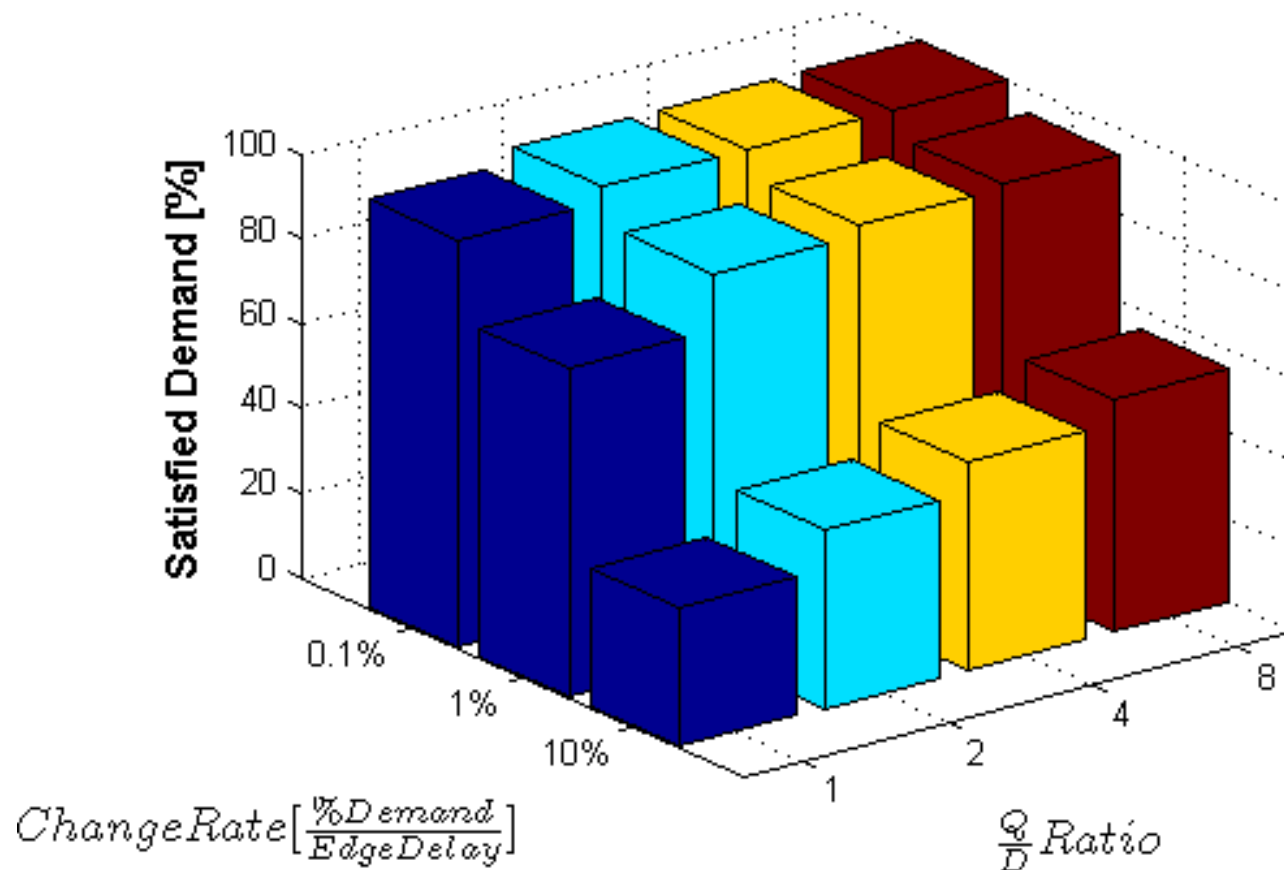
- Topology: BriteAS
- Change rate: $1\% \cdot D/E[\text{EdgeDelay}]$
- $Q/D=1$
- Demanders: 1%



GWiQ-P scalability
due to **locality**

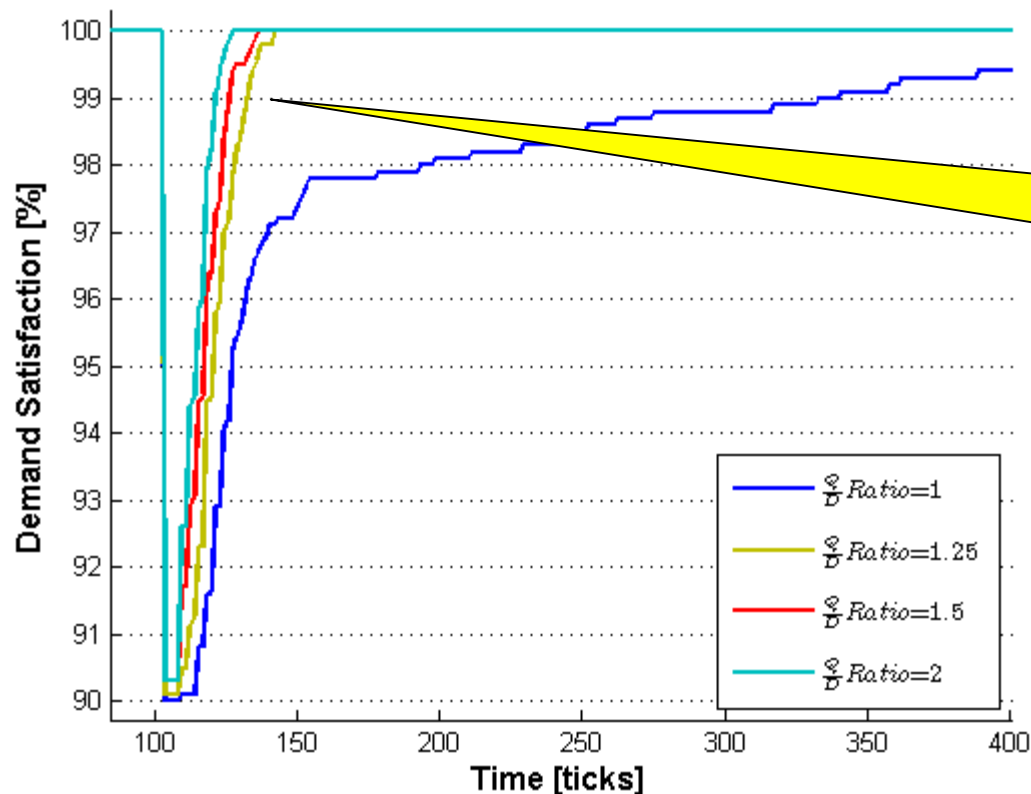
Simulations 2/6

- Topology: 10K BriteAS
- Demanders: 1%



Simulations 3/6

- Topology: 10K BriteAS
- $E[\text{EdgeDelay}] \sim 3 \text{ ms}$
- Change rate: One Time $1\% \cdot D$
- Demanders: 1%

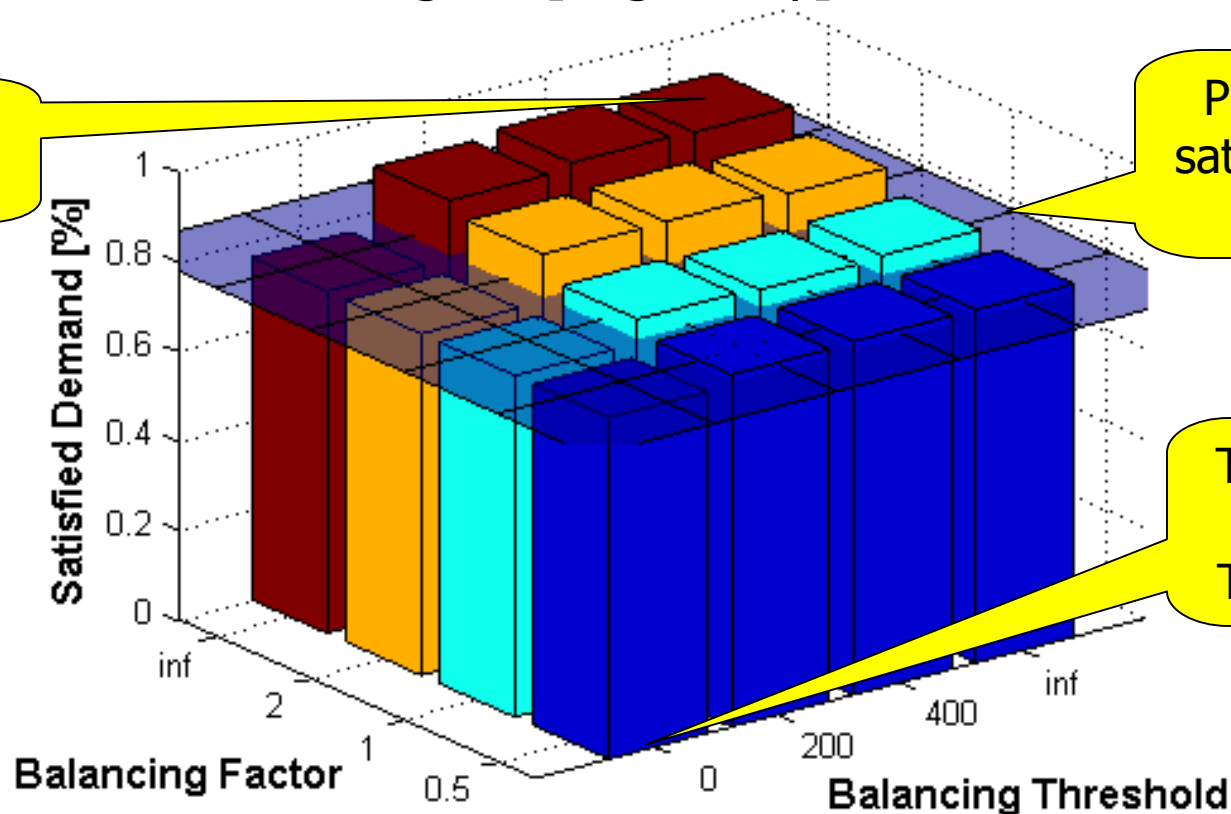


Return to
99% sat after
 $\approx 30\text{ms}$

Simulations 4/6

- Topology: 10K BriteAS
- Change rate: $1\% * D/E[\text{EdgeDelay}]$
- Fail rate: $1\% * \text{Edges}/E[\text{EdgeDelay}]$
- $Q/D=1$
- Demanders: 1%

Excess coin
exploitation

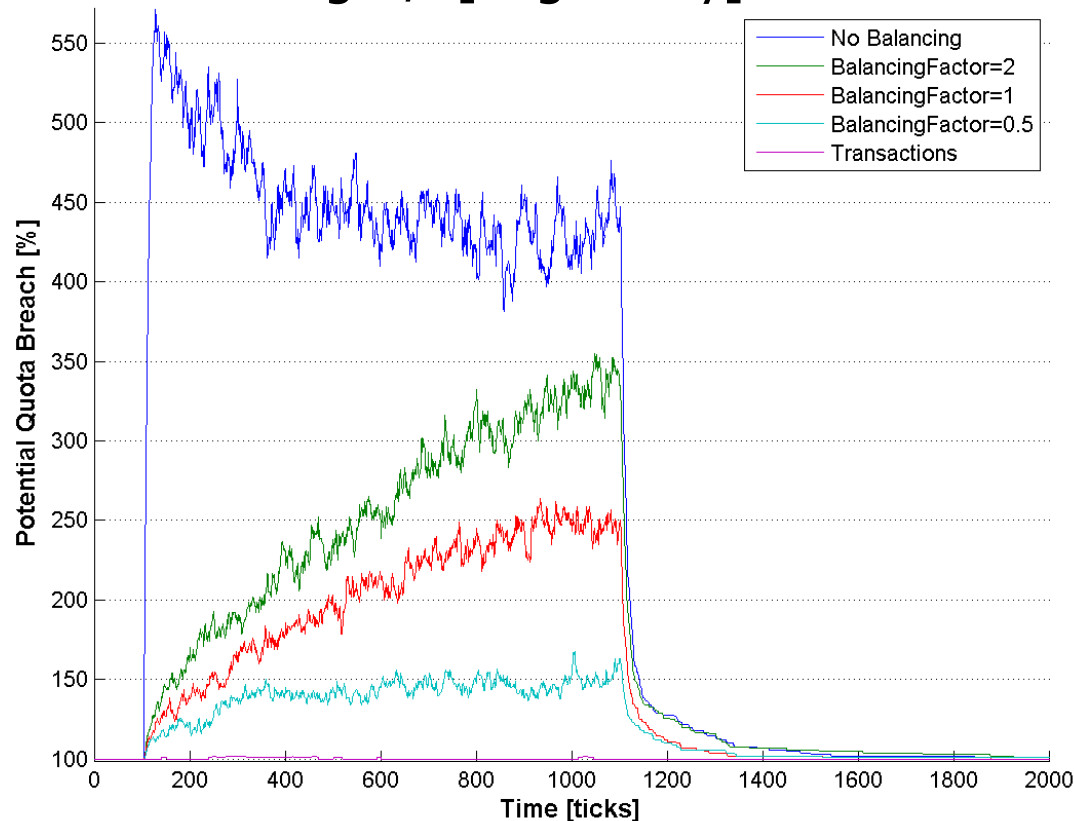


Plane depicts
sat in 'no-faults'
scenario

Threshold=0
 \Leftrightarrow
Transactions

Simulations 5/6

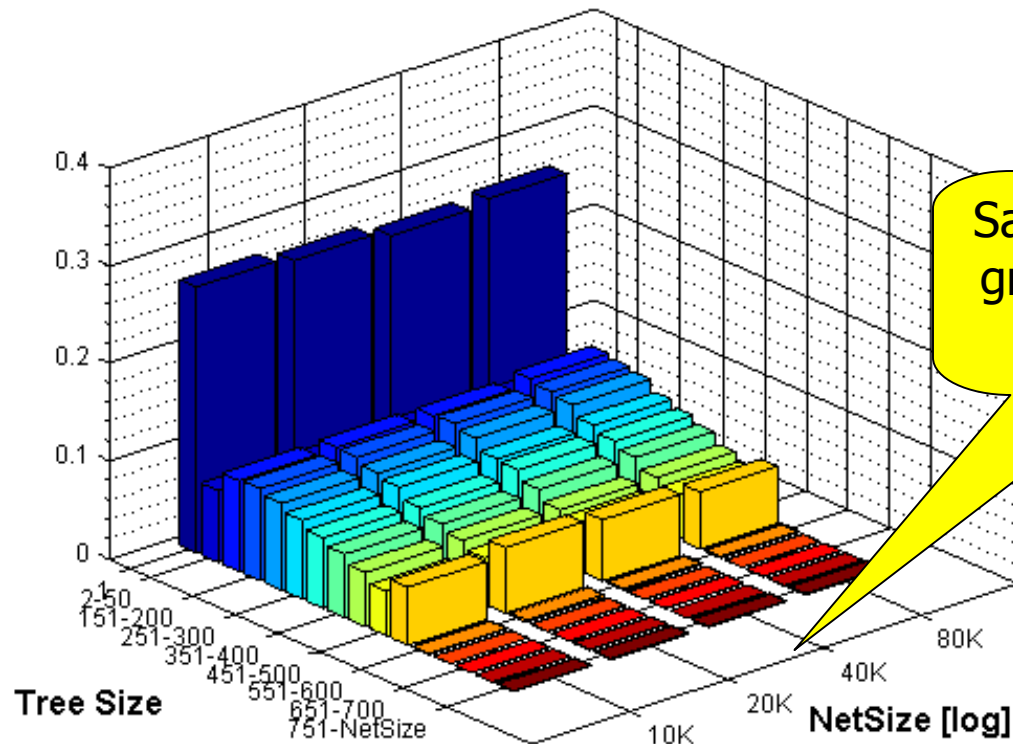
- Topology: 10K BriteAS
- Change rate: $1\% * D/E[\text{EdgeDelay}]$
- Fail rate: $1\% * \text{Edges}/E[\text{EdgeDelay}]$
- $Q/D=0.5$
- Demanders: 1%



Simulations 6/6

- Topology: BriteAS
- Change rate: $1\% * D/E[\text{EdgeDelay}]$
- Fail rate: $1\% * \text{Edges}/E[\text{EdgeDelay}]$
- $Q/D=1$
- Demanders: 1%

Nodes' Tree-Sizes Distribution



Same distribution,
growing network.
Locality in FT



Conclusion

- We displayed **GWiq-P**, a **Grid Wide Quota** enforcement **P**rotocol
- **GWiq-P** is infinitely **scalable**
- **GWiq-P** is fully **distributed**
- **GWiq-P** is **local** hence very efficient
- **GWiq-P** is **fault tolerant**



Thank You!

Q&A

Contact: Kfir Karmon
karmon@cs.technion.ac.il