# Epidemic Algorithms for replicated Database maintenance

Alan Demers et al
Xerox Palo Alto Research Center, PODC 87

Presented by: Harshit Dokania

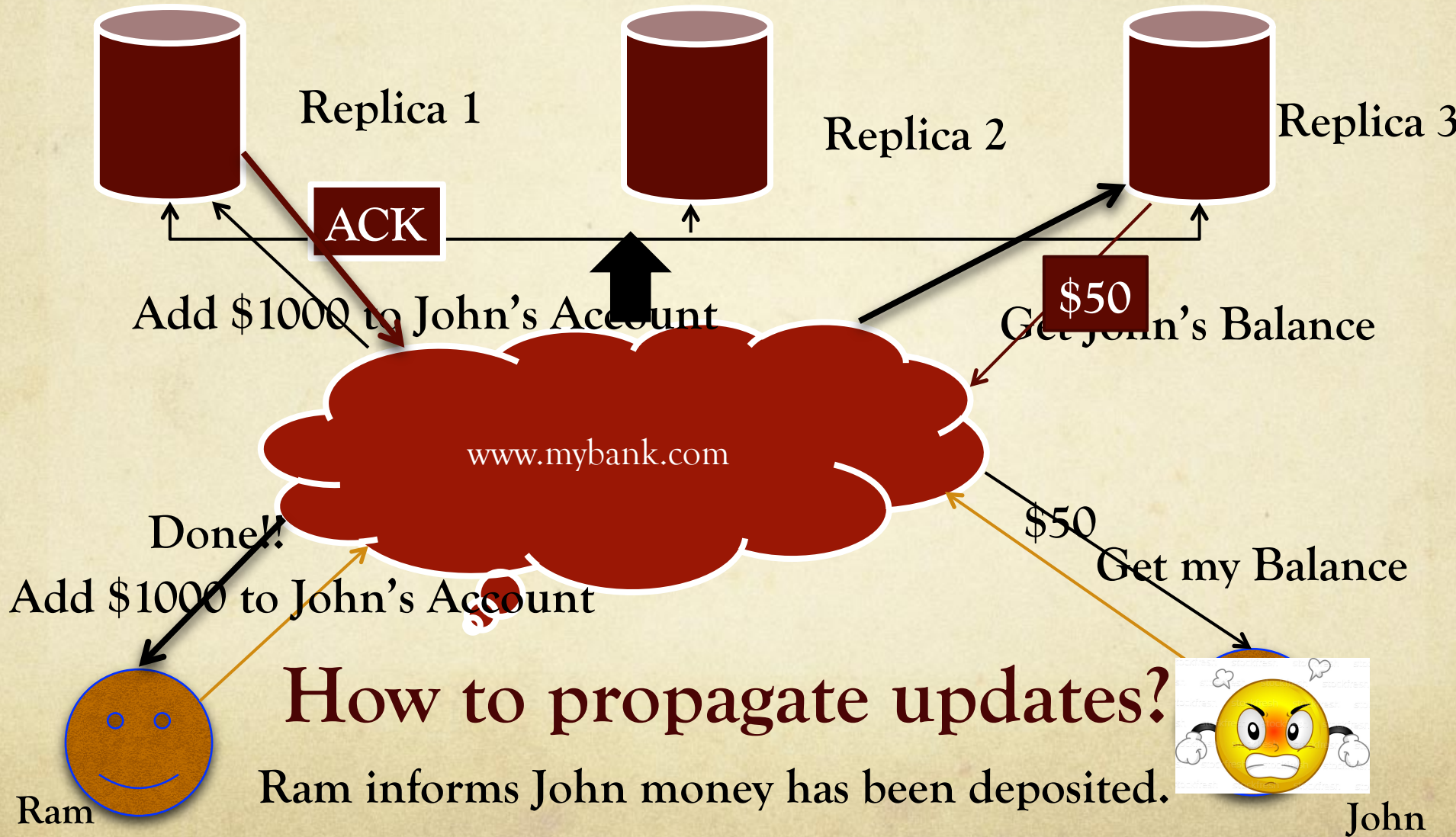# Motivation

○ **Replicated Databases**

**Availability, Fault – tolerant, Faster reads**

**Consistency**

**Eventual** Consistency : Epidemic Algorithms

**Faster reads, writes**

# Why Consistency ?

Replica 1    Replica 2    Replica 3

ACK

Add $1000 to John's Account

$50

Get John's Balance

www.mybank.com

Done!!

$50

Add $1000 to John's Account

Get my Balance

## How to propagate updates?
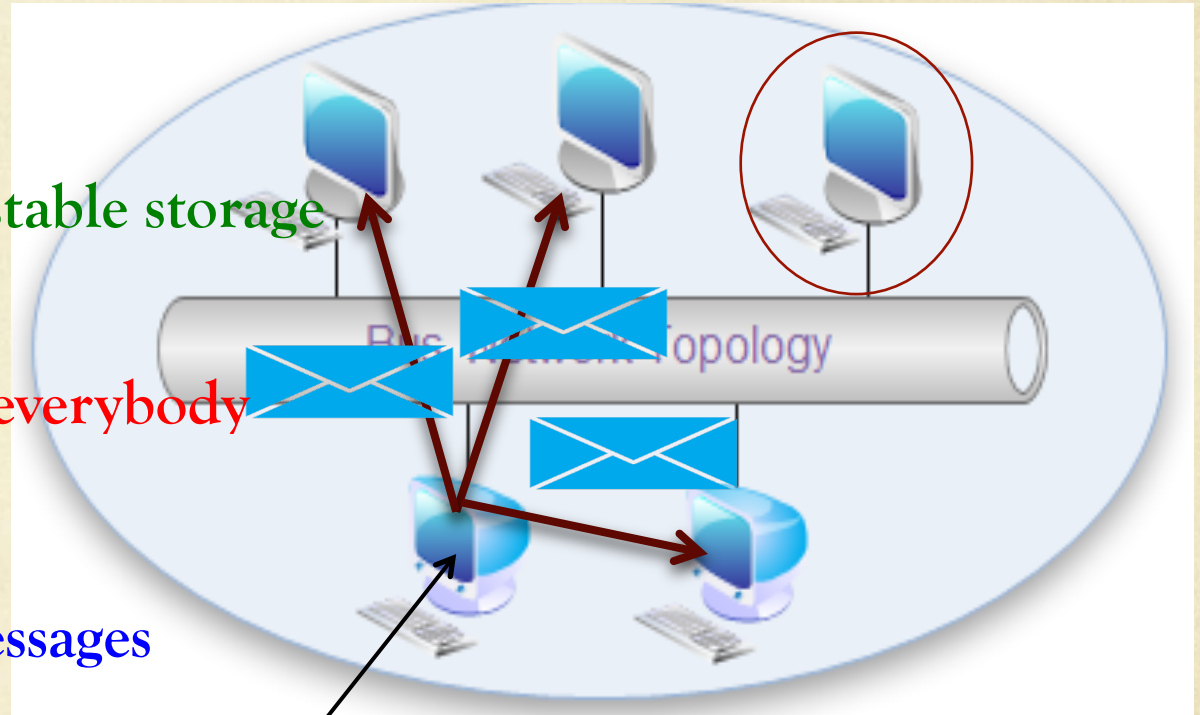
Ram informs John money has been deposited.

Ram

John

# How to Propagate Updates?

○ Direct Mail

○ Anti-entropy - Simple Epidemic

infected or susceptible sites

○ Rumor Mongering – Complex Epidemic

infected, susceptible and removed

Goal: For all s, s' ∈ S: s.valueOf = s'.valueOf

# Direct Mail

- Messages are queued, stable storage

  Reliable but...

- Not everybody knows everybody

- Queue overflows

  Source site = O(n) messages

  Traffic =  # sites * Average distance b/w the sites

Anti-Entropy in Background
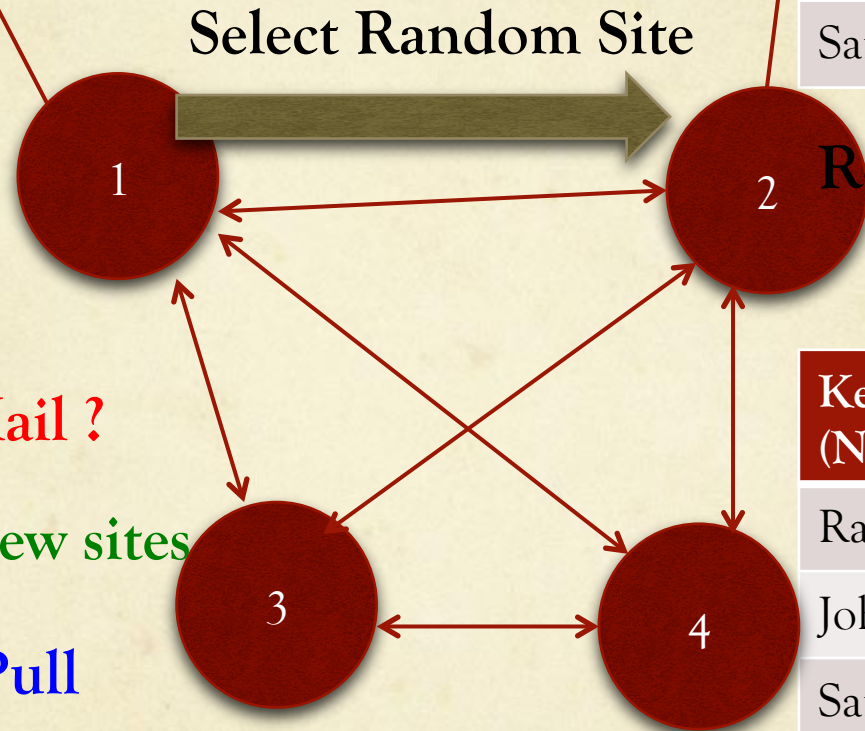
Update[v : V] = s.ValueOf ⟵——(v, GMT)

# Anti-Entropy

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 6000 | 100 |
| John | 50 | 100 |
| Sam | 1500 | 105 |

**Select Random Site**

**Resolve Conflicts**

1

2

3

4

**Slower than Direct Mail ?**

**Distribute update to few sites**

**Push, Pull, Push - Pull**

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1500 | 105 |

# Push vs Pull

○ **Pull**

$p_i$ = Probability of a site remaining susceptible after the $i^{th}$ cycle

Only if it Selects susceptible site in $i+1^{st}$ cycle

$$p_{i+1} = (p_i)^2 \approx 0$$

○ **Push**

Only if No infectious site chose to contact susceptible site

$$p_{i+1} = p_i(1-1/n)^{n(1-p_i)} = p_i e^{-1} \approx 0 \text{ (less rapidly)}$$

**But Anti-Entropy is Expensive!!!**

# Anti-Entropy is Expensive

o     Usually Databases are in "nearly" complete agreement

     <span style="color:red">Then why send entire Database across network ?</span>

o     Exchange Database only if checksum of Database disagree

     <span style="color:red">Time to update to all sites > Interval between the updates</span>

o     Recent update list that contains all new changes for time window $t'$

     <span style="color:green">$t'$ MUST > Time required to distribute the update</span>

o     Exchange recent update list, then compare checksums

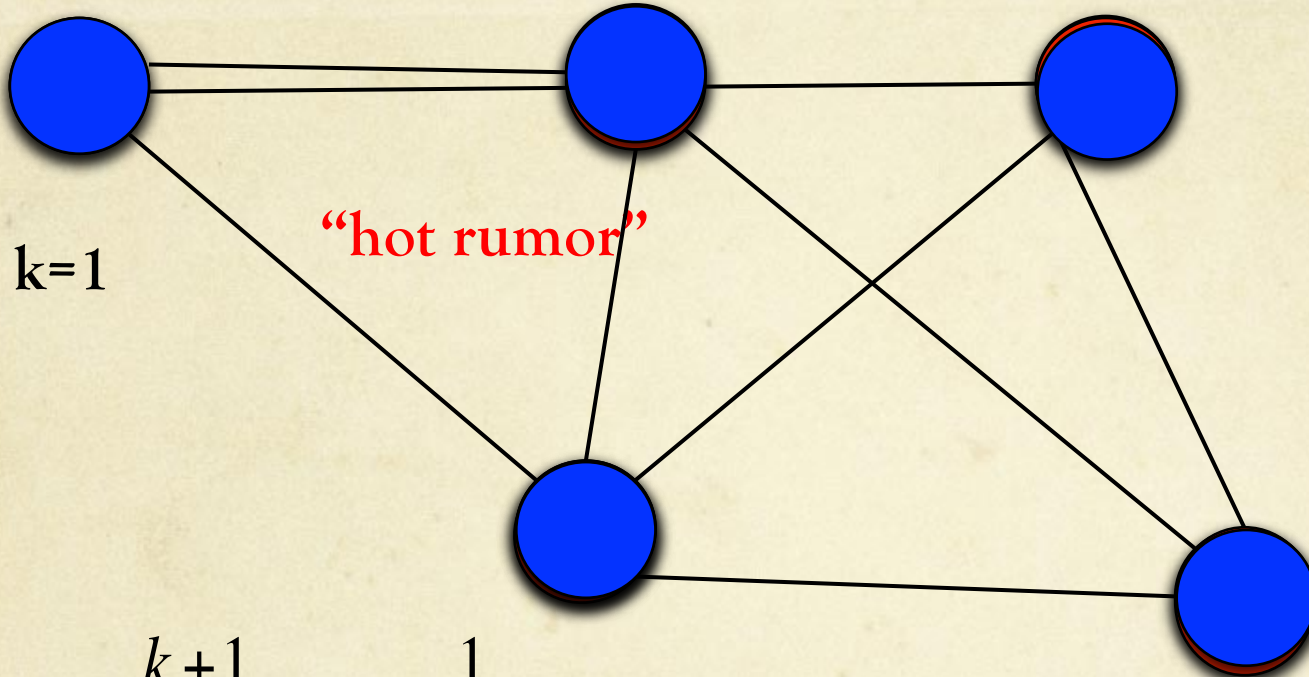     <span style="color:green">Checksums agree, Less traffic, Less Database comparison</span>

# Rumor Mongering

**Removed**

~~Infected~~

**Susceptible**

"hot rumor"

k=1

$$i(s) = \frac{k+1}{k}(1-s) + \frac{1}{k}\log s$$

$$s = e^{-(k+1)(1-s)}$$

Convergence, i = 0

s = ?    **Residue**

s

s

s

# Counter vs Probability

○ **Become removed after k unnecessary contacts**

## Response and Counters

| Counter k | Residue s | Traffic m | Converge $t_{ave}\mid t_{last}$ |
|---|---|---|---|
| 1 | 0.176 | 1.74 | 11.0 16.8 |
| 2 | 0.037 | 3.30 | 12.1 16.9 |
| 3 | 0.011 | 4.53 | 12.5 17.4 |
| 4 | 0.0036 | 5.64 | 12.7 17.5 |
| 5 | 0.0012 | 6.68 | 12.8 17.7 |

## Blind and Probabilistic

| Counter k | Residue s | Traffic m | Converge $t_{ave}\mid t_{last}$ |
|---|---|---|---|
| 1 | 0.960 | 0.04 | 19 38 |
| 2 | 0.205 | 1.59 | 17 33 |
| 3 | 0.060 | 2.82 | 15 32 |
| 4 | 0.021 | 3.91 | 14.1 32 |
| 5 | 0.008 | 4.95 | 13.8 32 |

**Convergence**

**Residue**

**Push, RM**

**Traffic**

# Push vs Pull

o  **Numerous updates**

Susceptible can find infective with high Probability

| Counter k | Residue s | Traffic m | Convergence t_ave \| t_last |
|-----------|-----------|-----------|------------------------------|
| 1 | $3.1 * 10^{-7}$ | 2.70 | 9.97   17.63 |
| 2 | $5.8 * 10^{-4}$ | 4.49 | 10.07  15.39 |
| 3 | $4.0 * 10^{-6}$ | 6.09 | 10.08  14.00 |

Pull, Response and Counters

Exchange counters  If both know the update  Increment the site with smaller counter

Push gets better than Pull with connection limit. How?

Two sites contact the same recipient  One gets rejected  Still gets the update
Two sites contact the same infected site  One gets rejected  Only 1 site updated

# Direct Mail vs Rumor

o     **Both has no guarantee**

      - **Anti-entropy is the key**

o     **But what if Anti- entropy finds conflicts ?**

      - **Let it be...**

      - **Clearinghouse uses Direct Mail for redistribution**

**Worst case : DM Manages to deliver an update half the sites**
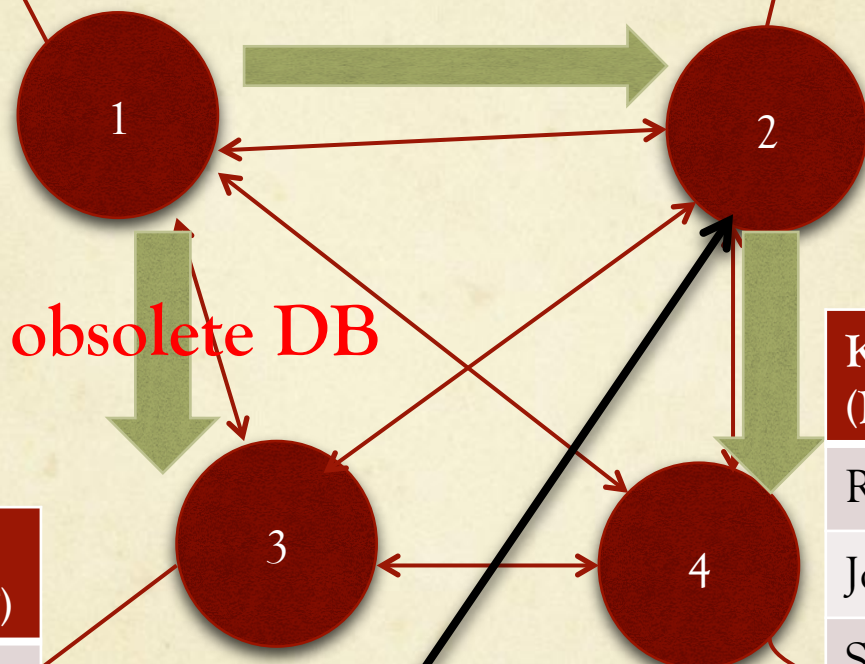
      **Later AE-DM generates $O(n^2)$ messages**

      **or**

      **AE-RM generates $O(n)$ messages**

# How do we delete?

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

**Resurrection with obsolete DB**

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

1

2

3

4

**Delete my account**

John

# Death Certificates

Delete John, 120

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| Sam | 1200 | 100 |
| Sam | 1200 | 100 |

| Key (Name) | Value (Bal) | Tim (GM |
|---|---|---|
| Ram | 5000 | 110 |
| John | 1050 | 110 |
| Sam | 1200 | 100 |

Delete John, 120

**But when do we delete Death Certificates**

**Looks Easy !!!**

| Key (Name) | Value (Bal) | Time (GMT) |
|---|---|---|
| Ram | 5000 | 110 |
| Sam | 1200 | 100 |

**Delete my account**

John

# Delete Death Certificates

o   Delete them after 30 days

    Still has risk of <span style="color:red">resurrection</span>

o   Sarin and Lynch propose Chandy-Lamport snapshot algorithm

    Records <span style="color:blue">Application</span> messages when sees <span style="color:blue">duplicate Marker</span> messages

    <span style="color:green">Snapshot</span> ensures Death certificate received at all sites

    But what when site <span style="color:red">permanently fails?</span>

# Dormant Death Certificates

o    Delete death certificates at most sites

    -But retain Dormant death certificates at some retention sites "r"

o    Obsolete update encounters Dormant certificate
    Activate Death certificate

    -Original timestamp, Activation timestamp

o    Much larger threshold for Dormant Death Certificates

o    But lost due to permanent server failures

$$p_{fail} = 2^{-r}$$

# Spatial Distribution

o      **Saves significant amount of traffic**

o      **probability of connecting site at distance $d = d^{-a}$**

         **when a = 2 ; <span style="color:green">convergence</span> is polynomial in <span style="color:green">log n</span>**

         **Generalized for CIN with distribution <span style="color:green">$d^{-2D}$</span>**

           <span style="color:red">**dependent on dimension of mesh 'D'**</span>

o      **$Q_s(d)$ = sites at distance d or less from s**

         <span style="color:blue">**arbitrary network**</span>    <span style="color:green">**Best Scaling**</span>

# Using Spatial Distribution

o   **Anti-Entropy**

   Uniform distribution ≈ 80 conversations on critical links

   Expected traffic per link per cycle < 6 conversations

   $Q_s(d)$ adapts well to "local dimension"

o   **Rumor Mongering**

   Making it less sensitive for sudden increases in $Q_s(d)$

   Each site s builds a list of sites sorted by distance

   For a = 2 again complexity is O($d^{-2d}$)

# Simulation Results: AE

### Push-Pull, No Connection limit

| Distribution | $t_{last}$ | $t_{ave}$ | Compare Traffic Avg | Critical |
|---|---|---|---|---|
| Uniform | 7.81 | 5.27 | 5.87 | 75.74 |
| a = 1.2 | 10.04 | 6.29 | 2.00 | 11.19 |
| a = 1.4 | 10.31 | 6.39 | 1.93 | 8.77 |
| a = 1.6 | 10.94 | 6.70 | 1.71 | 5.72 |
| a = 1.8 | 11.97 | 7.21 | 1.52 | 3.74 |
| a = 2.0 | 13.32 | 7.76 | 1.36 | 2.38 |

### Push-Pull, Connection limit 1

| Distribution | $t_{last}$ | $t_{ave}$ | Compare Traffic Avg | Critical |
|---|---|---|---|---|
| Uniform | 11.00 | 6.97 | 3.71 | 47.54 |
| a = 1.2 | 16.89 | 9.92 | 1.14 | 6.39 |
| a = 1.4 | 17.34 | 10.15 | 1.08 | 4.68 |
| a = 1.6 | 19.06 | 11.06 | 0.94 | 2.90 |
| a = 1.8 | 21.46 | 12.37 | 0.82 | 1.68 |
| a = 2.0 | 24.64 | 14.14 | 0.72 | 0.94 |

Increase in convergence < factor of 2

Decrease in Average traffic ≈ factor of 4    Frequent anti-entropy

Massive decrease in compare traffic for transatlantic links > factor of 30

Connection limit reduces the compare traffic/cycle but increases the number of cycles

# Simulation Results: RM

| Distribution | k | $t_{last}$ | $t_{ave}$ | Compare traffic Avg | Bushey | Update traffic Avg | Bushey |
|---|---|---|---|---|---|---|---|
| Uniform | 4 | 7.83 | 5.32 | 8.87 | 114.0 | 5.84 | 75.87 |
| a = 1.2 | 5 | 10.14 | 6.33 | 3.20 | 18.0 | 2.60 | 17.25 |
| a = 1.4 | 6 | 10.27 | 6.31 | 2.86 | 13.0 | 2.49 | 14.05 |
| a = 1.6 | 7 | 11.24 | 6.90 | 2.94 | 9.80 | 2.27 | 10.54 |
| a = 1.8 | 8 | 12.04 | 7.24 | 2.40 | 5.91 | 2.08 | 7.69 |
| a = 2.0 | 9 | 13.00 | 7.74 | 1.00 | 3.44 | 1.90 | 5.94 |

Feedback, Counter, Push-Pull, No connection limit

With increase in "a", k increases gradually  Convergence time increases

Decrease in Average traffic ≈ factor of 3

Massive decrease in compare traffic for transatlantic links > factor of 30

# Discussion

o    Anti- Entropy is <span style="color:green">robust</span> than Rumor Mongering

    Rumors may become <span style="color:red">inactive</span> leaving sites <span style="color:red">susceptible</span>

o    Push $\xi$ Pull much sensitive to spatial distribution than Push-Pull (RM)

       k = <span style="color:red">36</span> for a = 1.2 (Push, Feedback, No connection limit, Counter)

o    Anti - Entropy with distribution of $d^{-2}$ was implemented at CIN

       Massive <span style="color:green">improvement</span> in network load $\xi$ consistency

# Cons/Questions

○ **Storage**, Death Certificates

○ **Irregular** Network Topologies

- **Dynamic Hierarchy**