

# An Adaptive Algorithm for Efficient Message Diffusion in Unreliable Environments

Benoît Garbinato  
Fernando Pedone  
Rodrigo Schmidt

Proceedings of the 2004 IEEE International Conference on Dependable Systems and Networks (DSN 2004)

## Our goal

We want to solve the reliable broadcast problem in large-scale distributed systems, using as few messages as possible

We want our solution to be able to adapt to changes in the system configuration

# Outline

- ④ Distributed systems
- ④ Large scale systems
- ④ Reliable broadcast
- ④ Probabilistic approach
- ④ Adaptation & Optimality

## Distributed systems

“A distributed system is one that stops you from getting any work done when a machine you’ve never even heard of crashes.”

Leslie Lamport quoted by Sape Müllender in Distributed Systems, 2nd edition. Addison-Wesley, 1993.

- ④ As slow as the slowest
  - ④ As weak as the weakest
- ⇒ In large-scale systems,  
this is even worse...

# Large scale

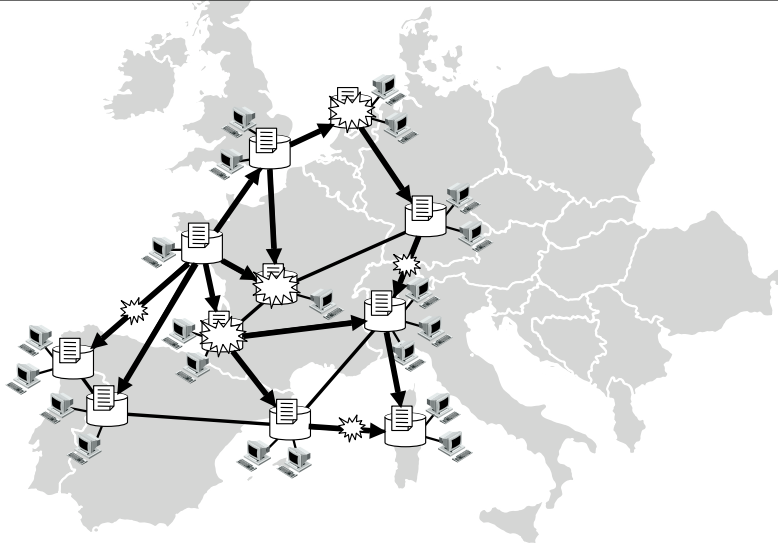
- ④ Many nodes (processes)  
and/or
- ④ Long-distances (high-latency)
  - ⇒ Changes occur frequently
  - ⇒ Difficult to act globally

# Reliability

- ④ Try to solve typical problems, e.g., consensus, reliable broadcast, atomic commitment, etc., despite the occurrence of failures
- ④ A solution depends on the actual failure model, e.g., crash-stop, crash-recovery, Byzantine, etc.

# Gossiping algorithms

Principle: only talk to some of your neighbors



Tuesday, May 23, 2006

7

## Probabilistic model

### Model definition

- ☉ The system is defined as a graph  $G = (\Pi, \Lambda)$   
 $\Pi = \{p_1, p_2, \dots, p_n\}$  is a set processes (vertices)  
 $\Lambda = \{l_1, l_2, \dots, l_m\} \subseteq \Pi \times \Pi$  is a set of links (edges)
- ☉ Processes communicate by message passing
- ☉ A failure probability  $P_i$  is associated with each process  $p_i$  and a message-loss probability  $L_j$  is associated with each link  $l_j$ . This probabilities set defines a configuration.

Tuesday, May 23, 2006

8

# Probabilistic reliability

## Problem statement

**Validity.** If a process broadcasts a message  $m$ , then it eventually delivers  $m$ .

**Agreement.** If a process delivers a message  $m$ , then every process eventually delivers  $m$  with probability  $K$ .

**Integrity.** For any message  $m$ , every process delivers  $m$  at most once, and only if  $m$  was previously broadcast by some process.

## Our goal revisited

- ④ We want to provide a solution the probabilistic reliable broadcast in large-scale systems
- ④ We want our solution to adapt to changes in the unreliability probabilities of processes & links
- ④ We want to measure the effectiveness of our adaptive strategy with respect an algorithm that minimizes the number of messages

# Optimality & adaptation

**Optimality.** A probabilistic reliable broadcast algorithm  $O$  is optimal to some configuration  $C$  w.r.t. the number of messages if there is no algorithm  $X$  such that processes executing  $X$  in  $C$  exchange fewer messages than processes executing  $O$  in  $C$ .

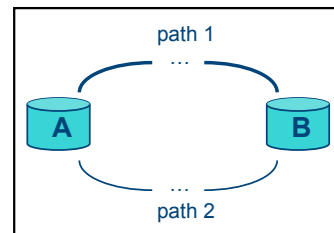
**Adaptation.** A probabilistic reliable broadcast algorithm  $A$  is adaptive to some configuration  $C$  if and only if the number of messages exchanged by processes executing  $A$  in  $C$  in response to a reliable broadcast is **eventually** equal to the number of messages exchanged by processes executing  $O$  in  $C$

## Optimal algorithm (1)

① Use **adaptive** gossiping, not **random** gossiping (traditional)

② **Example:**

- Failure probability of path 1 =  $L$
- Failure probability of path 2 =  $\alpha L$ , with  $\alpha > 1$  (path 1 more reliable than 2)



③ Probability  $K$  of reaching B from A

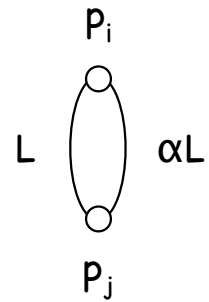
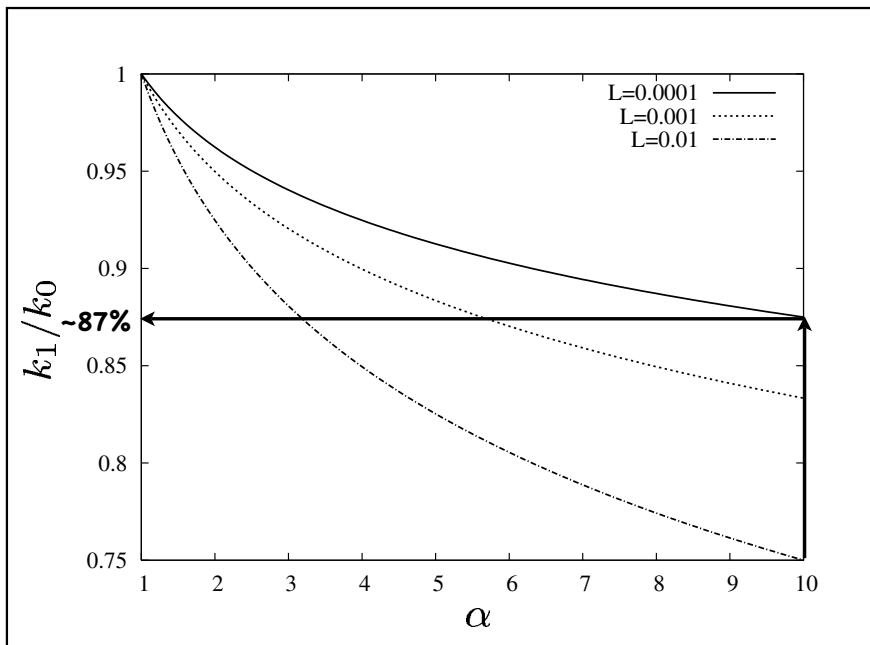
Random

$$K = 1 - L^{k_0/2} (\alpha L)^{k_0/2} = 1 - (\sqrt{\alpha} L)^{k_0}$$

Adaptive

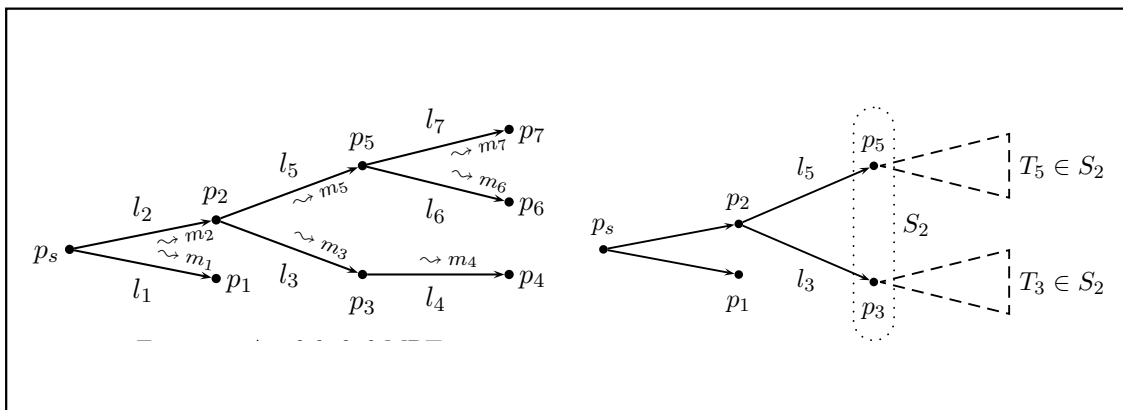
$$K = 1 - L^{k_1}$$

# Optimal algorithm (2)



# Maximum reliability tree

The maximum reliability tree is a spanning tree containing the most reliable path in  $G$  connecting all processes in  $\Pi$



# The reach function

Given a tree  $T_i$  and a vector  $\vec{m}_i$ , the reach function computes the probability that all processes in  $T_i$  are reached by at least one message

$$reach(T_i, \vec{m}_i) = \begin{cases} 1 & \text{if } T_i = \perp \\ \prod_{T_j \in S_i} (1 - [1 - (1 - P_i) \times (1 - L_j) \times (1 - P_j)]^{m_i[j]}) \times reach(T_j, \vec{m}_j) & \text{otherwise} \end{cases}$$

$$reach(T_i, \vec{m}_i) = \prod_{j=1}^{n-1} 1 - [1 - (1 - P_{pred(j)}) \times (1 - L_j) \times (1 - P_j)]^{m_i[j]}$$

# Optimization problem

With  $\lambda_j = 1 - (1 - P_{pred(j)}) \times (1 - L_j) \times (1 - P_j)$ , we end up with the following optimization problem:

$$\begin{aligned} \text{minimize} \quad & c(\vec{m}) = \sum_{j=1}^{|\vec{m}|} m[j] \\ \text{subject to} \quad & r(\vec{m}) = \prod_{j=1}^{|\vec{m}|} 1 - \lambda_j^{m[j]} \geq K \end{aligned}$$



# The optimize function

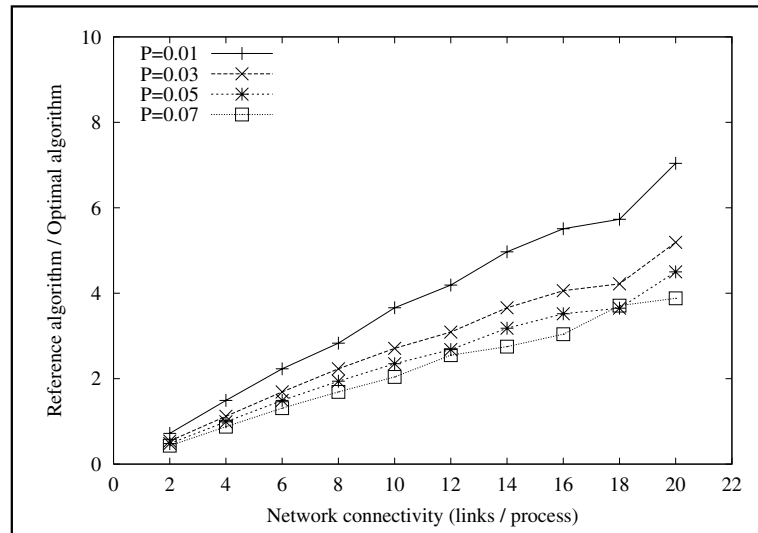
We use a greedy algorithm, since our optimization problem is itself greedy

```
function optimize(mrt, K)  
   $\vec{m} \leftarrow (1, 1, 1, \dots, 1)$   
  while  $r(\vec{m}) < K$  do  
    let  $\vec{u}_j$  be such that  $\frac{r(\vec{m} + \vec{u}_j)}{r(\vec{m})}$  is maximum  
     $\vec{m} \leftarrow \vec{m} + \vec{u}_j$   
  return  $\vec{m}$ 
```

# Optimal algorithm

```
1: To execute broadcast(m) do  
2:    $mrt_k \leftarrow mrt_k(G, C)$   
3:   propagate(m,  $mrt_k$ ,  $p_k$ )  
4:   deliver(m)  
  
5: when receive (m,  $mrt_j$ ) for the first time  
6:   propagate(m,  $mrt_j$ ,  $p_k$ )  
7:   deliver(m)  
  
8: function propagate(m,  $mrt_j$ ,  $p_k$ )  
9:    $\vec{m}_j \leftarrow \boxed{\text{optimize}(mrt_j, K)}$   
10:  for all subtree  $T_i \in S_{j,k}$  do  
11:    repeat  $\vec{m}_j[i]$  times  
12:      send (m,  $mrt_j$ ) to  $p_i$ 
```

# Performance of the optimal algorithm

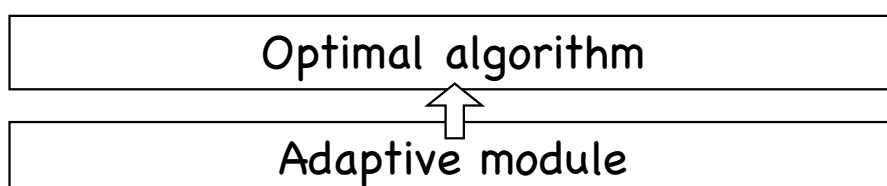


Tuesday, May 23, 2006

19

## From optimal to adaptive

- Our algorithm is proven to be optimal when it has exact knowledge of the system configuration
- By replacing such (unrealistic) exact knowledge with a module that approximates the changing configurations, we get an adaptive algorithm
- As soon as the approximation module converged towards the actual system configuration, we are again optimal



Tuesday, May 23, 2006

20

Adaptive algorithm

```

1: INITIALIZATION:
2:   for all  $p_i \in \Pi$  do
3:     initializeReliability( $C_k[p_i]$ )
4:      $C_k[p_i].d \leftarrow \infty$ 
5:      $C_k[p_i].seq \leftarrow 0$ 
6:      $C_k[p_i].suspected \leftarrow 0$ 
7:      $\Delta_k[p_i] \leftarrow \delta$ 
8:      $C_k[p_k].d \leftarrow 0$ 
9:    $\Lambda_k \leftarrow \{l_{k,i} \mid p_i \in \text{neighbors}(p_k)\}$ 
10:  for all  $l_i \in \Lambda_k$  do
11:    initializeReliability( $C_k[l_i]$ )
12:     $C_k[l_i].d \leftarrow 0$ 
13: TO UPDATE ( $\Lambda_k, C_k$ ):
14:  every  $\Delta_k[p_k]$  do :
15:     $C_k[p_k].seq \leftarrow C_k[p_k].seq + 1$ 
16:    for all  $p_i \in \text{neighbors}(p_k)$  do
17:      send ( $\Lambda_k, C_k$ ) to  $p_i$ 
18:  when received ( $\Lambda_j, C_j$ ) from  $p_j$  do {Event 1}
19:    adjust  $\leftarrow C_k[p_j].suspected -$ 
20:      ( $C_j[p_j].seq - C_k[p_j].seq$ )
21:     $C_k[p_j].suspected \leftarrow 0$ 
22:    if adjust > 0 then
23:      increaseReliability( $C_k[l_{k,j}], \text{adjust}$ )
24:      if adjust > 1 then  $\Delta_k[p_j] \leftarrow \Delta_k[p_j] + \delta$ 
25:    if adjust < 0 then
26:      decreaseReliability( $C_k[l_{k,j}], |\text{adjust}|$ )
27:    for all  $p_i \in \Pi$  do
28:      selectBestEstimate( $C_k[p_i], C_j[p_i]$ )
29:    for all  $l_i \in (\Lambda_k \cap \Lambda_j)$  do
30:      selectBestEstimate( $C_k[l_i], C_j[l_i]$ )
31:    for all  $l_i \in \Lambda_j - (\Lambda_k \cap \Lambda_j)$  do
32:       $C_k[l_i] \leftarrow C_j[l_i]$ 
33:       $C_k[l_i].d \leftarrow C_k[l_i].d + 1$ 
34:       $\Lambda_k \leftarrow \Lambda_k \cup \Lambda_j$ 
35: 34: when not[updated  $C_k[p_j], p_j \neq p_k$ , in the last  $\Delta_k[p_j]$ ] do {Event 2}
36:     $C_k[p_j].d \leftarrow C_k[p_j].d + 1$ 
37:    if  $p_j \in \text{neighbors}(p_k)$  then
38:       $C_k[p_j].suspected \leftarrow C_k[p_j].suspected + 1$ 
39:      decreaseReliability( $C_k[p_j], 1$ )
40:      decreaseReliability( $C_k[l_{k,j}], 1$ )
41: 40: every  $\Delta_{tick}$  do {Event 3}
42:    increaseReliability( $C_k[p_k], 1$ )
43: 42: when recovering from a crash lasting  $n \times \Delta_{tick}$  do {Event 4}
44:    decreaseReliability( $C_k[p_k], n$ )

```

```

1: Initialization
2:    $U \leftarrow 100$  {precision of probabilistic intervals}
3: function initializeReliability(estimate)
4:   with estimate do
5:     for all  $u = 1..U$  do
6:        $P_{F|B}[u] \leftarrow \frac{2u-1}{2U}$  {probabilistic intervals}
7:        $P_B[u] \leftarrow \frac{1}{U}$  {with equal initial beliefs}
8: function decreaseReliability(estimate, factor)
9:   with estimate repeat factor times
10:  for all  $u = 1..U$  do
11:     $P_B[u] \leftarrow \frac{P_B[u] \times P_{F|B}[u]}{\sum_{v=1}^U P_B[v] \times P_{F|B}[v]}$ 
12: function increaseReliability(estimate, factor)
13:   with estimate repeat factor times
14:  for all  $u = 1..U$  do
15:     $P_B[u] \leftarrow \frac{P_B[u] \times (1 - P_{F|B}[u])}{\sum_{v=1}^U P_B[v] \times (1 - P_{F|B}[v])}$ 

```

Tuesday, May 23, 2006

21

## Bayesian inference (1)

$$P_{B|F}[u] = \frac{P_{F|B}[u] \times P_B[u]}{\sum_{v=1}^U P_{F|B}[v] \times P_B[v]}$$

Tuesday, May 23, 2006

22

# Bayesian inference (2)

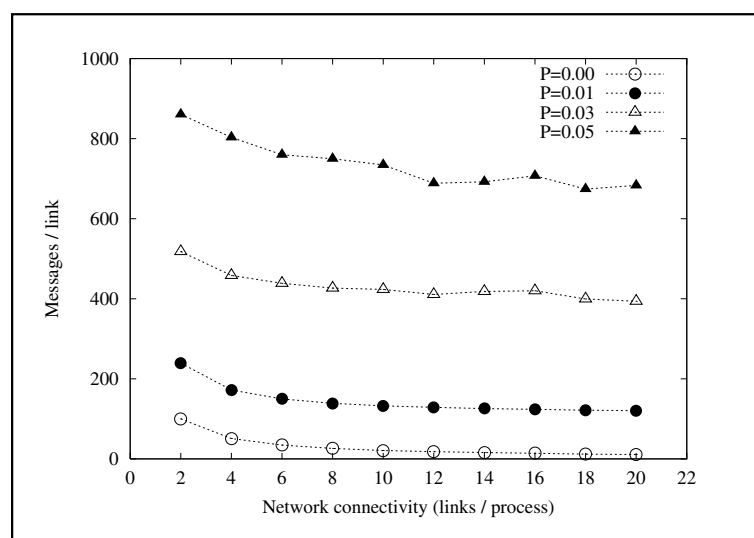
$u$	$C_k[p_i].P_{F B}[u]$	$C_k[p_i].P_B[u]$
1	[0.0 , 0.2)	0.2
2	[0.2 , 0.4)	0.2
3	[0.4 , 0.6)	0.2
4	[0.6 , 0.8)	0.2
5	[0.8 , 1.0]	0.2

(a) Initial configuration

$u$	$C_k[p_i].P_{F B}[u]$	$C_k[p_i].P_B[u]$
1	[0.0 , 0.2)	0.04
2	[0.2 , 0.4)	0.12
3	[0.4 , 0.6)	0.20
4	[0.6 , 0.8)	0.28
5	[0.8 , 1.0]	0.36

(b) After a failure suspicion

# Adaptation convergence



# Question?