

# Message reliability and caching for publish/subscribe systems

Andreas Tanner, Gero Mühl

Technische Universität Berlin  
Einsteinufer 17, 10587 Berlin, Germany  
{tanner,gmuehl}@ivs.tu-berlin.de

March 5, 2004

Publish/subscribe

Message completeness

Temporal logic

Safety and Liveness

Axioms for p/s systems

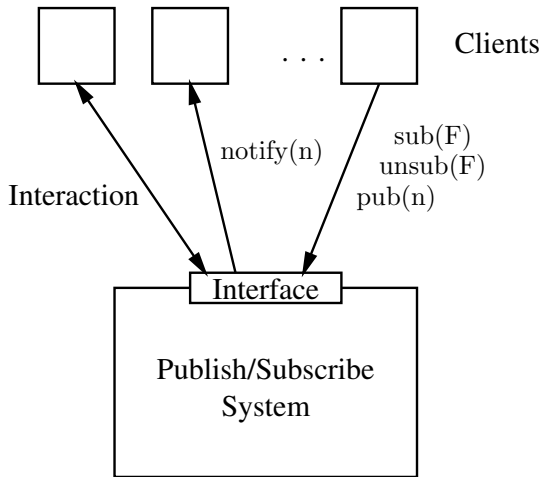
Distributed Implementaion

Rebeca

## Publish/Subscribe

- ▶ Enables loosely coupled communication using notifications.
- ▶ Two kinds of “clients”
  - ▶ Producers publish notifications
  - ▶ Consumers subscribe to notifications
- ▶ Notification service
  - ▶ Decouples producers from consumers
  - ▶ Delivers a published notification to all consumers with a matching subscription

## Interface

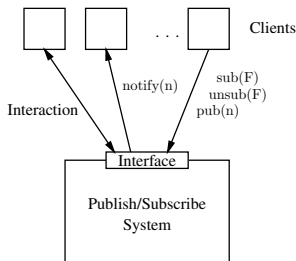


## Basic Definitions and Assumptions

- ▶ A filter  $F$  is a mapping from the set of notifications  $\mathcal{N}$  to the boolean values *true* and *false*.
- ▶ A notification  $n$  matches a filter  $F$  iff  $F(n) = \text{true}$ .
- ▶ Set of notifications matched by a filter  $F$  and by a set of filters  $\mathcal{A}$  is denoted by  $N(F)$  and  $N(\mathcal{A})$ , respectively.
- ▶ Notifications are unique and can be published only once.

## Black box view of a publish/subscribe system

- ▶ Describes the system behavior by solely looking at its interface.



- ▶ Interface Operations (set of all actions  $A$ )

$sub(Y, F)$	Client $Y$ subscribes to filter $F$
$unsub(Y, F)$	Client $Y$ unsubscribes to filter $F$
$ack(Y, F)$	System notifies client $Y$ about message completeness guarantee
$notify(Y, n)$	System notifies client $Y$ notified about $n$
$pub(X, n)$	Client $X$ publishes $n$

## Multicast

- ▶ publish/subscribe is special type of multicast
  - ▶ event-based
  - ▶ dynamically created groups
- ▶ problems classically considered for multicast:
  - ▶ message completeness
  - ▶ message order
  - ▶ etc.

## Message completeness

How can *message completeness* be defined in a publish/subscribe system?

- ▶ *first try*: once a client subscribes to a filter, he gets all matching messages published thereafter



## Message completeness

How can *message completeness* be defined in a publish/subscribe system?

- ▶ *first try*: once a client subscribes to a filter, he gets all matching messages published thereafter
- ▶ in an asynchronous system, this is not achievable if messages are not cached (subscription must first be propagated through the network)

## Message completeness

How can *message completeness* be defined in a publish/subscribe system?

- ▶ *first try*: once a client subscribes to a filter, he gets all matching messages published thereafter
- ▶ in an asynchronous system, this is not achievable if messages are not cached (subscription must first be propagated through the network)
- ▶ *better*: after a client subscribes to a filter, there is a future time at which he is guaranteed to see all messages published thereafter

## Message completeness

How can *message completeness* be defined in a publish/subscribe system?

- ▶ *first try*: once a client subscribes to a filter, he gets all matching messages published thereafter
- ▶ in an asynchronous system, this is not achievable if messages are not cached (subscription must first be propagated through the network)
- ▶ *better*: after a client subscribes to a filter, there is a future time at which he is guaranteed to see all messages published thereafter
- ▶ *even better*: after a client subscribes to a filter, there is a future time *the client is notified about* at which he gets message completeness guarantee

## The Need for a formal treatment

- ▶ A formal specification
  - ▶ defines precisely what is expected from a correct system and
  - ▶ allows to reason about the correctness of an implementation.
- ▶ A formal treatment gives new insights that could otherwise be overlooked!
- ▶ *Linear Temporal Logic* offers a formalism suitable for characterization of the behaviour of distributed systems

## State

- ▶ *state*  $s$  is assignment  $s = s : \mathcal{V} \ni v \mapsto v_s$  of the state variables  $v \in \mathcal{V}$  to values in their domains

## State

- ▶ *state*  $s$  is assignment  $s = s : \mathcal{V} \ni v \mapsto v_s$  of the state variables  $v \in \mathcal{V}$  to values in their domains
- ▶ *interface operations*  $op$  trigger state transitions

## State

- ▶ *state*  $s$  is assignment  $s = s : \mathcal{V} \ni v \mapsto v_s$  of the state variables  $v \in \mathcal{V}$  to values in their domains
- ▶ *interface operations*  $op$  trigger state transitions
- ▶ a *trace* is a sequence of initial state  $s_0$ , followed by interface operations:

$$\sigma = s_0, op_1, op_2, \dots, op_n \quad (1)$$

## State

- ▶ *state*  $s$  is assignment  $s = s : \mathcal{V} \ni v \mapsto v_s$  of the state variables  $v \in \mathcal{V}$  to values in their domains
- ▶ *interface operations*  $op$  trigger state transitions
- ▶ a *trace* is a sequence of initial state  $s_0$ , followed by interface operations:

$$\sigma = s_0, op_1, op_2, \dots, op_n \quad (1)$$

- ▶ predicate applied to trace  $\sigma$  refers to *first state*  $s_0$  or *first operation*  $op_1$



## Temporal quantifiers

For some formula  $\phi$  and  $\sigma = s_0, op_1, op_2, \dots, op_n$ ,

- ▶  $\diamond\phi(\sigma)$  holds iff there exists  $i$  such that  $\phi$  holds for the trace  $s_i^\sigma, op_{i+1}, \dots$ ,

## Temporal quantifiers

For some formula  $\phi$  and  $\sigma = s_0, op_1, op_2, \dots, op_n$ ,

- ▶  $\diamond\phi(\sigma)$  holds iff there exists  $i$  such that  $\phi$  holds for the trace  $s_i^\sigma, op_{i+1}, \dots$ ,
- ▶  $\square\phi(\sigma)$  holds iff for all  $i$ ,  $\phi$  holds for the trace  $s_i^\sigma, op_{i+1}, \dots$ ,

## Temporal quantifiers

For some formula  $\phi$  and  $\sigma = s_0, op_1, op_2, \dots, op_n$ ,

- ▶  $\diamond\phi(\sigma)$  holds iff there exists  $i$  such that  $\phi$  holds for the trace  $s_i^\sigma, op_{i+1}, \dots$ ,
- ▶  $\square\phi(\sigma)$  holds iff for all  $i$ ,  $\phi$  holds for the trace  $s_i^\sigma, op_{i+1}, \dots$ ,
- ▶  $\circ\phi(\sigma)$  holds iff  $\phi$  holds for the trace  $s_1^\sigma, op_2, \dots$

## LTL and Concurrency

- ▶ LTL is well-suited to describe concurrent systems
- ▶ A concurrent system is replaced by a nondeterministic sequential one.
- ▶ The concurrent execution of two operations in the real system is replaced in the model by the nondeterminism of which one occurs first.
- ▶ This type of nondeterminism is conceptually different from that studied in the area of automata theory ( $\Rightarrow$  branching time).

## LTL Examples

- ▶  $\diamond \square A$
- ▶  $\square \diamond A$
- ▶  $A \Rightarrow \diamond B$
- ▶  $\square [A \Rightarrow \diamond B]$
- ▶  $\square [A \Rightarrow \square A]$
- ▶  $\square [\square A \Rightarrow \diamond \square B]$
- ▶  $\square [A \vee \square \neg A]$

## State Variables in p/s systems

$P_X$	set of <i>published notifications</i>
$S_Y^{\text{ack}}$	set of <i>acknowledged subscriptions</i>
$S_Y^{\text{pend}}$	set of <i>pending subscriptions</i>
$D_Y$	multiset of <i>delivered notifications</i>

Initial values of state variables:  $\emptyset$  for all of them

## Effect of interface operations on state variables

$pub(X, n)$	$P'_X = P_X \cup \{n\}$
$sub(Y, F)$	$S_Y^{pend'} = S_Y^{pend} \cup \{F\}$
$unsub(Y, F)$	$S_Y^{ack'} = S_Y^{ack} \setminus \{F\}; S_Y^{pend'} = S_Y^{pend} \setminus \{F\}$
$ack(Y.F)$	$S_Y^{pend'} = S_Y^{pend'} \setminus \{F\}; S_Y^{ack'} = S_Y^{ack} \cup \{F\}$
$notify(Y, n)$	$D'_Y = D_Y \cup \{n\}$

## Safety and Liveness

### ▶ Safety Conditions

- ▶ Something “irremediably” bad will never happen.
- ▶ Usually, phrased as an invariant of the system.
- ▶ Usually, trivially satisfied by doing *nothing*.
- ▶ General form:  $Init \Rightarrow \Box \neg A$ .
- ▶ Violation can be detected after *finite* time.
- ▶ E.g. partial correctness (Program never halts with wrong result)



## Safety and Liveness (2)

- ▶ Liveness Conditions
  - ▶ Something “good” that should happen eventually happens.
  - ▶ Usually, trivially satisfied by doing *everything*.
  - ▶ General form:  $Init \Rightarrow \Box[A \Rightarrow \Diamond B]$ .
  - ▶ Violation can be detected after *infinite* time only.
  - ▶ Example: Termination (Program eventually halts)
- ▶ Many useful system properties (e.g., total correctness) can be expressed as the intersection of safety and liveness conditions.

## Safety axioms

A *publish/subscribe system* satisfies *message complete safety* if



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow [\bigcirc \square \neg \text{notify}(Y, n, X)] \right] \quad (2)$$

## Safety axioms

A *publish/subscribe system* satisfies *message complete safety* if



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow [\bigcirc \square \neg \text{notify}(Y, n, X)] \right] \quad (2)$$



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow [\exists X. n \in P_X] \right] \quad (3)$$

## Safety axioms

A *publish/subscribe system* satisfies *message complete safety* if



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow [\bigcirc \square \neg \text{notify}(Y, n, X)] \right] \quad (2)$$



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow [\exists X. n \in P_X] \right] \quad (3)$$



$$\square \left[ \text{notify}(Y, n, X) \Rightarrow \exists F \in S_Y^{\text{pend}} \cup S_Y^{\text{ack}}. n \in N(F) \right] \quad (4)$$

## Liveness axioms

A *publish/subscribe system* satisfies *message complete liveness* if



$$\square \left[ (sub(Y, F) \wedge \neg \diamond unsub(Y, F)) \Rightarrow \diamond [ack(Y, F)] \right] \quad (5)$$

## Liveness axioms

A *publish/subscribe system* satisfies *message complete liveness* if



$$\square \left[ (sub(Y, F) \wedge \neg \diamond unsub(Y, F)) \Rightarrow \diamond [ack(Y, F)] \right] \quad (5)$$

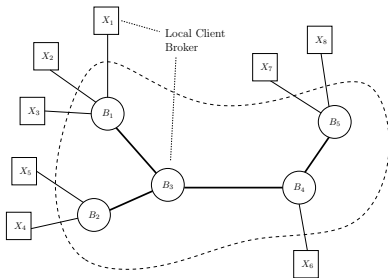


$$\square \left[ (ack(Y, F) \wedge \neg \diamond unsub(Y, F)) \Rightarrow \right. \\ \left. (\diamond pub(X, n) \wedge n \in N(F) \Rightarrow \right. \\ \left. \left. \diamond notify(Y, n, X) \right) \right] \quad (6)$$

## Correctness

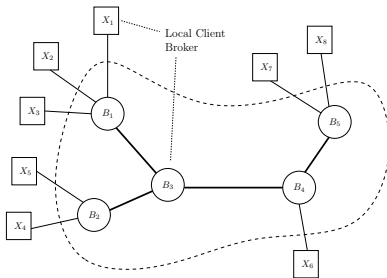
A publish/subscribe system is *correct*, if it satisfies safety and liveness.

## Distributed Implementation



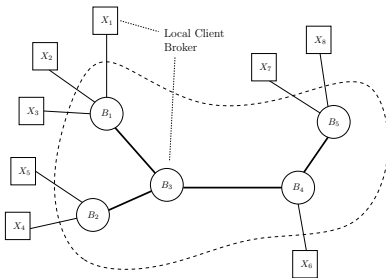


## Distributed Implementation



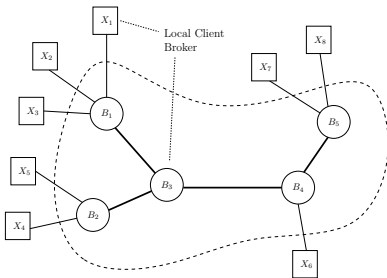
- Set of brokers  $B_1, \dots, B_n$  serving as access points.

## Distributed Implementation



- ▶ Set of brokers  $B_1, \dots, B_n$  serving as access points.
- ▶ Brokers are concurrent processes cooperating by message passing.

## Distributed Implementation



- ▶ Set of brokers  $B_1, \dots, B_n$  serving as access points.
- ▶ Brokers are concurrent processes cooperating by message passing.
- ▶ Each broker  $B$  manages a mutually exclusive set of local clients  $L_B$  and only communicates directly with its neighbor brokers  $N_B$ .

## Assumptions

- ▶ Broker topology assumed to be acyclic.
- ▶ Channels are reliable (no corrupted, duplicated, lost, or spurious messages).
- ▶ Message latency is bounded.
- ▶ Communication with clients conceptually treated as message passing but assumed to be instantaneous.

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:
  - ▶ *forward*( $n$ ) (used to disseminate notifications)



## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:
  - ▶  $forward(n)$  (used to disseminate notifications)
  - ▶  $admin(\mathcal{S}, \mathcal{U})$  (used to update routing tables)

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:
  - ▶  $forward(n)$  (used to disseminate notifications)
  - ▶  $admin(\mathcal{S}, \mathcal{U})$  (used to update routing tables)
  - ▶  $admin\_ack(\mathcal{S}, \mathcal{U})$  (used for acknowledgement between brokers)

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:
  - ▶  $forward(n)$  (used to disseminate notifications)
  - ▶  $admin(\mathcal{S}, \mathcal{U})$  (used to update routing tables)
  - ▶  $admin\_ack(\mathcal{S}, \mathcal{U})$  (used for acknowledgement between brokers)
- ▶ Processing of  $forward$  and  $pub$  messages hardwired.

## Content-Based Routing Framework

- ▶ Each broker  $B$  manages a routing table  $T_B$  comprising a set of routing entries.
- ▶ A routing entry is a pair  $(F, D)$  of a filter  $F$  and a destination  $D \in L_B \cup N_B$ .
- ▶ Brokers exchange three kinds of messages:
  - ▶  $forward(n)$  (used to disseminate notifications)
  - ▶  $admin(\mathcal{S}, \mathcal{U})$  (used to update routing tables)
  - ▶  $admin\_ack(\mathcal{S}, \mathcal{U})$  (used for acknowledgement between brokers)
- ▶ Processing of  $forward$  and  $pub$  messages hardwired.
- ▶ Processing of  $admin$  messages is customized by implementing an instance of the `administer` procedure.

## Notification Forwarding

- ▶ If a broker receives a *forward*( $n$ ) / *pub*( $n$ ) message from a neighbor / local client, it sends a *forward*( $n$ ) / *notify*( $n$ ) message to all neighbors / local clients  $D$  for which there is a routing entry  $(F, D)$  with  $n \in N(F)$ .

## Notification Forwarding

- ▶ If a broker receives a *forward*( $n$ ) / *pub*( $n$ ) message from a neighbor / local client, it sends a *forward*( $n$ ) / *notify*( $n$ ) message to all neighbors / local clients  $D$  for which there is a routing entry  $(F, D)$  with  $n \in N(F)$ .
- ▶ But a notification is never passed back to the neighbor it was received from.

## Notification Forwarding

- ▶ If a broker receives a *forward*( $n$ ) / *pub*( $n$ ) message from a neighbor / local client, it sends a *forward*( $n$ ) / *notify*( $n$ ) message to all neighbors / local clients  $D$  for which there is a routing entry  $(F, D)$  with  $n \in N(F)$ .
- ▶ But a notification is never passed back to the neighbor it was received from.
- ▶ As the topology is acyclic, duplicate notifications are avoided.

## Notification Forwarding

- ▶ If a broker receives a *forward*( $n$ ) / *pub*( $n$ ) message from a neighbor / local client, it sends a *forward*( $n$ ) / *notify*( $n$ ) message to all neighbors / local clients  $D$  for which there is a routing entry  $(F, D)$  with  $n \in N(F)$ .
- ▶ But a notification is never passed back to the neighbor it was received from.
- ▶ As the topology is acyclic, duplicate notifications are avoided.
- ▶ As every *notify*( $Y, n$ ) has a preceding *pub*( $X, n$ ), no spurious notifications are delivered.



## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on administer procedure.

## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on administer procedure.
- ▶ Paper gives criterias (*valid routing algorithms*) that when met, lead to correct system

## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on `administer` procedure.
- ▶ Paper gives criterias (*valid routing algorithms*) that when met, lead to correct system
- ▶ Various routing algorithms can be implemented via `administer`:

## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on `administer` procedure.
- ▶ Paper gives criterias (*valid routing algorithms*) that when met, lead to correct system
- ▶ Various routing algorithms can be implemented via `administer`:
  - ▶ Flooding

## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on `administer` procedure.
- ▶ Paper gives criterias (*valid routing algorithms*) that when met, lead to correct system
- ▶ Various routing algorithms can be implemented via `administer`:
  - ▶ Flooding
  - ▶ Simple routing

## Handling of *sub* and *unsub* messages

- ▶ Handling of subscription changes is parametrized on `administer` procedure.
- ▶ Paper gives criterias (*valid routing algorithms*) that when met, lead to correct system
- ▶ Various routing algorithms can be implemented via `administer`:
  - ▶ Flooding
  - ▶ Simple routing
  - ▶ Routing with filter merging

## Rebeca

- ▶ started by Ludger Fiege and Gero Mühl, TU Darmstadt as Java-based implementation of p/s system
- ▶ implemented along the lines of formal framework
- ▶ based on microkernel architecture with routing component as kernel
- ▶ ported to C# by Andreas Ulbrich
- ▶ uses events and delegates, seamlessly integrating into .NET framework
- ▶ extended (acknowledgements, caching) in student project winter 2003/2004 at TU Berlin

## Rebeca

- ▶ started by Ludger Fiege and Gero Mühl, TU Darmstadt as Java-based implementation of p/s system
- ▶ implemented along the lines of formal framework
- ▶ based on microkernel architecture with routing component as kernel
- ▶ ported to C# by Andreas Ulbrich
- ▶ uses events and delegates, seamlessly integrating into .NET framework
- ▶ extended (acknowledgements, caching) in student project winter 2003/2004 at TU Berlin
- ▶ ... Rebeca promises to become the first ever publish/subscribe system whose correctness can be formally proven!



## Future Work

- ▶ Routing in cyclic topologies
- ▶ Funnel functions
- ▶ Fault tolerance
- ▶ Adaptivity
- ▶ Integration of routing and composite event detection
- ▶ Streaming operators
- ▶ Quality of Service

**Thank You.**

Questions?

Andreas Tanner

FG Intelligent Networks and Management of Distributed Systems

Technische Universität Berlin

[tanner@ivs.tu-berlin.de](mailto:tanner@ivs.tu-berlin.de)