

Composite Event Detection as a Generic Middleware Extension

Peter R. Pietzuch, Brian Shand, and Jean Bacon
University of Cambridge Computer Laboratory

Abstract

Event-based communication provides a flexible and robust approach to monitoring and managing large-scale distributed systems. Composite event detection extends the scope and flexibility of these systems, by allowing application components to express interest in complex patterns of events. This makes it possible to handle the large numbers of events generated in Internet-wide systems, and in network monitoring and pervasive computing applications. In this article we introduce a novel generic composite event detection framework that can be added on top of existing middleware architectures, as demonstrated in our implementation over JMS. We argue that the framework is flexible, expressive, and easy to implement. Based on finite state automata extended with a rich time model and parameterization support, it provides a decomposable core language for specifying composite events. This allows detection to be distributed automatically throughout the system, guided by distribution policies that control the quality of service. Finally, tests show that using our composite event system over JMS can reduce bandwidth consumption while maintaining low notification delay for composite events.

Builders of large-scale distributed systems need middleware services that support loosely coupled, scalable interaction between large numbers of communication partners. Event-based communication has become a new paradigm for building such systems, because of both these advantages and the simple localized application programming model that uses events as the basic communication mechanism: an event can be seen as a notification that something of interest has occurred within the system. Components act as either *event sources* that *publish* new events, or *event sinks* that *subscribe* to events by providing a specification of events of interest to them. A *publish/subscribe* (pub/sub) communication layer [1] is then responsible for disseminating events; for efficiency, it can often also filter events by topic or content, according to client specifications.

Many existing pub/sub systems [2–4] restrict subscriptions to single events only, and thus lack the ability to express interest in the occurrence of *patterns of events*. However, especially in large-scale applications, event sinks may be overwhelmed by the vast number of primitive low-level events, and would benefit from a higher-level view. Such a higher-level view is given by *composite events* (CEs) that are published when an event pattern occurs. To date, it is usually left to the event sink to implement a detector for composite events, making it unnecessarily complex and error-prone.

In this article we address the problem by proposing a middleware extension for composite event detection that works on top of a range of pub/sub systems. Our framework includes

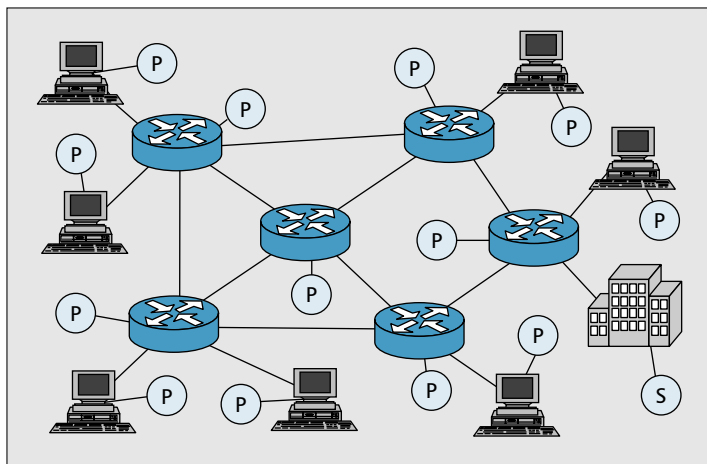
a generic language for specifying CEs and CE detectors that can detect CEs in a distributed way.

This article is an extension of earlier work [5] incorporating new material on distribution policies and automata construction. It is organized as follows. First, we motivate the necessity of CE detection in large-scale distributed systems. After related work, we discuss prerequisites of the detection framework such as the pub/sub infrastructure requirements, time model, and event model. The CE detectors and the associated core language are then presented, and we discuss distributed detection and detector placement policies. We present our implementation over JMS, show how detectors are constructed, and evaluate the performance. We finish with an introduction to higher-level specification languages and conclusions.

Motivation

Large-scale event systems need to support composite event detection in order to quickly and efficiently notify their clients of new relevant information in the network. This is particularly important for widely distributed systems, such as sensor networks, where bandwidth is limited and components are loosely coupled. In such systems, distributed CE detection can improve efficiency and robustness.

For example, in telecommunications network management, certain events raised by network devices may indicate that some problem has occurred [6]. As alarm-event notifications continue it should be possible to diagnose the specific cause. This must be carried out as quickly as possible so that faults can be repaired and customers reassured that management is aware of the problem and action is being taken. In practice, millions of events may be notified daily regarding fewer than 100 real faults. Issues of scalability, performance, and rele-



■ Figure 1. A publish/subscribe system for network systems monitoring.

vance to the end user of fault reporting are of current concern. For each possible problem an associated pattern of events can be specified. At present, diagnosis may be carried out sequentially at a network management center, comparing the events that have occurred with a number of hypotheses.

Figure 1 shows a network in which monitoring is carried out by dedicated publishers (P) connected to routers and workstations. Note that not all events indicate an actual fault that requires human intervention. For example, degradation in service at a single workstation may be transient in nature, but five events signaling reduced bandwidth at different points in the network may be symptomatic of a larger problem. A network management center (S) has to detect and react to network faults. Without CE detection, it has to subscribe to all events generated by the publishers. Many events are thus sent unnecessarily to the management center as they do not indicate faults and could be filtered out earlier on.

Composite event detectors allow the detection to occur within the network itself. Moreover, distributing the event composition allows filtering of large numbers of fault notifications to be carried out close to their sources, thus saving communication bandwidth. Since the detection happens in parallel, efficiency is also improved.

For reliability and efficiency, each CE detector should be distributed near its event sources. Otherwise, if a connection to the rest of the network failed, local notification of CEs might fail unnecessarily. Besides, sending these events off-site for detection would be a waste of bandwidth if all relevant events are known to be locally produced.

Just as a general-purpose pub/sub system supports flexible messaging, so too can a generic CE framework extend this support. Therefore, this article proposes a general-purpose middleware system for CE detection, independent of the specific underlying pub/sub infrastructure. By making CE detection interoperate closely with the underlying communication infrastructure, we obtain a system that is more efficient than an ad hoc implementation of CE detectors at the application level.

The Active Office

The Active Office is a computerized building that is aware of its inhabitants' behavior (Fig. 2). Workers wear Active Bats [7] to inform the building of their

movements at least once a minute. Other sensors monitor doors, office temperatures, electronic whiteboard usage, and lighting. A content-based pub/sub system is used so that applications can be notified of specific events, such as "location events where Peter is seen in room FE04." We used the following two application scenarios to test our CE detection framework.

Scenario 1 — The building services manager wants to know about temperature events under 15°C in an occupied room.

Scenario 2 — Jean wants the list of participants and electronic whiteboard contents of any meeting she has attended to be sent to her wireless PDA, but only if she does not log in to the workstation in her office within five minutes of the meeting.

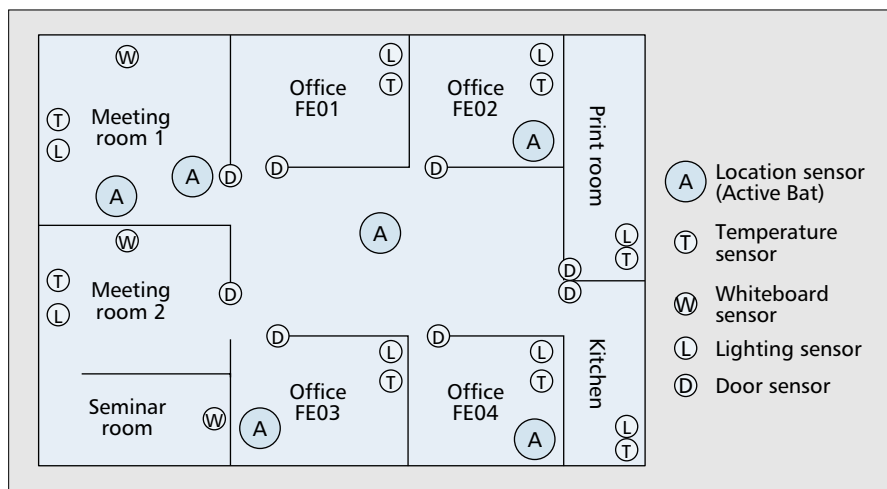
There are many advantages to using CE middleware for services in an Active Office, instead of (or perhaps as well as) offering predefined composite subscriptions on dedicated servers. Most important are the flexibility with which recipients can compose personal subscriptions, and the ease with which composite patterns can be reused and distributed close to event sources. The cost of establishing this network of CE detection broker nodes is then offset by the simplicity of configuring it for new CE subscriptions.

Related Work

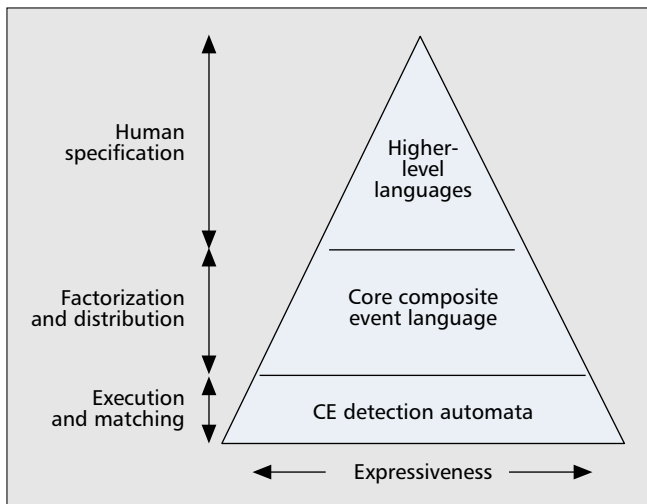
Historically, CE detection first arose in the context of triggers in active databases. Early languages for specifying composite events follow the Event-Condition-Action (ECA) model and resemble database query algebras with an expressive yet complex syntax. In general, the detection process is not distributed.

In the Ode object database [8], CEs are specified with a regular-expression-like language and detected using finite state automata (FSA). Equivalence between the CE language and regular expressions is shown. Since a CE has a single timestamp of the last event that led to its detection, a total event order is created that makes it difficult to deal with clock synchronization issues. Pure FSAs do not support parameterized events.

CE detectors based on Petri Nets are used in the SAMOS database [9]. Colored Petri Nets can represent concurrent



■ Figure 2. An active office environment.



■ Figure 3. Components of the composite event detection framework.

behavior and manage complex data such as event parameters during detection. However, even for simple expressions, they quickly become complicated. SAMOS does not support distribution and has a simple time model that is not suitable for distributed systems.

The motivation for Snoop [10] was to design an expressive CE specification language with powerful temporal support. A CE detector is a tree that reflects the structure of the event expression. Its nodes implement language operators and conform to a particular *consumption policy*. A consumption policy influences the semantics of an operator by resolving which events are consumed from the event history in case of ambiguity. For example, under a recent policy only the most recently occurring event is considered; others are ignored. Detection propagates up the tree with the leaves of the tree being primitive event detectors. A disadvantage is that the nodes are essentially Turing-complete, making it difficult to formalize their semantics and reason about their behavior. The use of consumption policies can be nonintuitive and operator-dependent.

In [11] Schwiderski presents a distributed CE architecture based on the 2g-precedence model for monitoring distributed systems. This model makes strong assumptions about the clock granularity in the system and thus does not scale to large, loosely coupled distributed systems. The language and detection algorithm used are similar to Snoop and suffer from the same shortcomings. It addresses the issue of events being delayed during transport by *evaluation policies*: *asynchronous evaluation* enables a detector to consume an event as soon as it arrives, sometimes leading to incorrect detection, whereas *synchronous evaluation* forces a detector to delay evaluation until all earlier events have arrived and assumes a heartbeat infrastructure. Although detection is distributed, no decision on the efficient placement of detectors in the network is made.

The GEM system [12] has a rule-based event monitoring language. It follows a tree-based detection approach and assumes a total time order. Communication latency is handled by annotating rules with tolerable delays. Such an approach is not feasible in an environment with unpredictable delays.

Research efforts in pervasive computing have led to CE languages that are intuitive to use in environments such as the Active Office. The work by Hayton [13] on CEs in the Cambridge Event Architecture (CEA) [14] is similar to ours in the sense that it defines a language that nonprogrammers can use to specify occurrences of interest. Hayton uses push-down FSAs to handle parameterized events. However, the language itself can become nonintuitive as the semantics of some oper-

ators is not obvious. Even though detectors can use ces as their input, distributed detection is not dealt with explicitly. As in previous work, scalar timestamps are used. Distributed pub/sub architectures such as Hermes [4], Gryphon [3, 15], and Siena [2] only provide parameterized primitive events and leave the task of CE detection to the application programmer. Siena supports restricted *event patterns*, but does not define a complete pattern language.

In our CE detection framework, we adopt the interval timestamp model introduced in [16]. The partial order of timestamps in a distributed system is made explicit by having timestamps associated with an uncertainty interval. A Common Object Request Broker Architecture (CORBA)-based detection architecture is presented in [16] that implements this time model. The notion of *event stability* is defined in order to handle communication delays. We extend this to cope with delays in wide area systems.

Design and Architecture

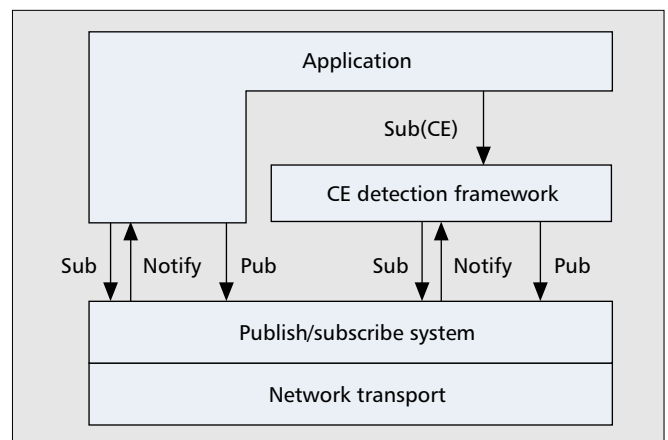
The CE detectors in our framework recognize concurrent patterns of simpler events, generating a CE whenever a match is found. The component layers of our detection architecture are illustrated in Fig. 3. Distributed CE detectors are compiled from expressions in our *core CE language*. Patterns can be specified using higher-level languages, which are first translated into the core CE language before compilation and execution.

The CE framework relies on and interacts with the underlying event system in order to detect complex patterns of events. This section outlines the prerequisites for this interaction: an interface to a pub/sub infrastructure, and formal models of events and time. Given these prerequisites, the full expressive power of our CE languages can be used.

Publish/Subscribe Infrastructure Support

One of our design goals was to keep the CE detection framework strictly separated from the pub/sub infrastructure used. The interface to the event system (Fig. 4) makes only minimal assumptions about the functionality supported, allowing our framework to be deployed on a large variety of pub/sub systems. Our current testbed uses the Java Message Service (JMS) [17], but other pub/sub systems could equally be used: earlier work was based on Hermes [4], a distributed event-based middleware architecture, and CORBA events would also be suitable.

In addition to the time and event model described below,



■ Figure 4. Interface between the CE detection framework and the pub/sub system.

```

public interface DistCEDServiceInf {
    public void registerCEType(CEType type,
        CEPublisherInf publisher);
    public void unregisterCEType(CEType type,
        CEPublisherInf publisher);
    public CEInf createCE(CEType ceType);
    public CEType createCEType(String typeName);
    public void publish(CEInf ce,
        CEPublisherInf publisher);
    public void subscribe(CEType type,
        CESubscriberInf subscriber, CEQoSInf qos,
        CESubscriberCallbackInf callback);
    public void unsubscribe(CEType type,
        CESubscriberInf subscriber);
}

```

■ Figure 5. The Java interface to the CE detection service.

the underlying pub/sub system needs to support publication of primitive events by event sources, subscription to these events by event sinks, and relaying of events from sources to sinks. Many systems also filter events en route for efficiency; our CE framework uses this if available, but no particular publication or subscription model is assumed. Our event model uses the abstraction of a *describable event set* as an atom for CE detection. If the pub/sub system supports content-based filtering, a describable event set will be defined by a parameterized filtering expression. In a topic-based system, it will conform to a certain event type only.

In particular, the pub/sub system does not need to be aware of CE types. As illustrated in Fig. 4, application event sources submit CE subscriptions to the CE detection layer. Any CEs then detected by a CE detector are published to the pub/sub system disguised as primitive events. It is then the responsibility of the pub/sub system to disseminate these encapsulated CE occurrences to all interested event sources. The same mechanism is used for the communication between distributed event detectors.

Composite Event Detection Framework

The Java interface to the CE detection service, presented to applications, is shown below in part. Applications may use this for all event services, or contact the underlying pub/sub infrastructure directly for primitive event subscriptions.

A summary of the application programming interface (API) is given in Fig. 5. Before an event type can be published it must be registered with the CE detection service so that an appropriate type/topic is created in the underlying pub/sub system. After that, a new event instance can be created using the `createCE` method. The `publish` method will pass the publication down to the pub/sub system. A call to `subscribe` subscribes to primitive or composite events. A CE subscription may trigger the instantiation of new CE detectors.

Time Model — Each event in our framework has an associated timestamp, denoting when it occurred. In a large-scale system, it may often be impossible to decide which of two events occurred first. Therefore, we assume that there is a partial order relation on a timestamp, $<$, showing which events definitively occurred before others. This is extended to a total order, \prec , using a tie-breaker convention (Appendix A), allowing events to be treated as a well ordered sequence of symbols for detection.

This may be illustrated using a two-part interval timestamp for events rather than a single conventional timestamp. These interval timestamps are used implicitly throughout the rest of this work. They can represent the clock uncertainty of a distributed time service such as NTP, and also the time interval associated with a CE. The intervals can factor in the estimat-

ed receiver-specific delay on receiving UTC, including radio transmission lag or network delays.

Figure 6 illustrates three interval timestamps t_1 , t_2 , and t_3 . Here, $t_1 < t_2$, $t_1 < t_3$, but $t_2 \not\prec t_3$ and $t_3 \not\prec t_2$. On the other hand, using the total order, $t_1 \prec t_2 \prec t_3$; these operators are formally defined in Appendix A.

Event Model — Events provide notification of observations in a distributed system. *Primitive events* represent observations from outside the event system, while *composite events* represent patterns of events. The constituents of a CE may be primitive events, or simpler CEs. Despite this distinction, all events are treated homogeneously; we assume only that events have timestamps, and can be consistently ordered for each subscriber (e.g., by interval timestamp, source IP address, and local event generation count).

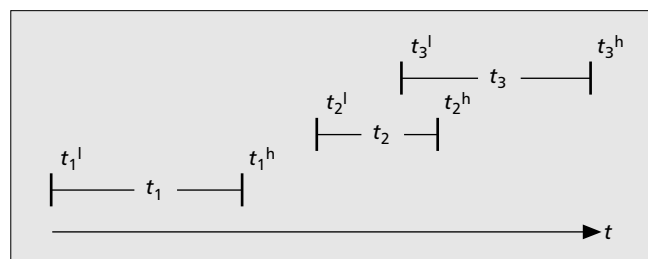
In an Active Office [7], primitive events might be “The door opens” and “Peter is seen in the room.” Similarly, “The door opens, then Peter is seen in the room” could be a CE. Thus, event sequences of interleaved primitive and composite events can be used to formalize the detection of CEs (Appendix A). Furthermore, our event ordering still supports distributed detection, since each CE detector’s subscription is used to sequence only the events needed for its CE pattern.

Composite Event Detection

The CE detectors in our framework are simple automata, with a regular structure. Unlike conventional FSAs, these automata provide support for a rich time model and parameterization, as well as the ability to detect concurrent event patterns. A novel language is used to express these patterns; this *core CE language* can then be compiled into automata for matching.

Distribution support is important for communication efficiency; we discuss how each pattern may be factorized into subexpressions. These subexpressions can then be matched independently on distributed nodes; these mobile detectors were introduced in an earlier article [18]. Patterns may also be more intuitively defined using higher-level specification languages, described later. However, this is only a matter of convenience, not expressiveness; any patterns described in a higher-level language are first translated into the core CE language before being compiled into automata. Figure 3 illustrates the relationship between these different aspects of the CE framework.

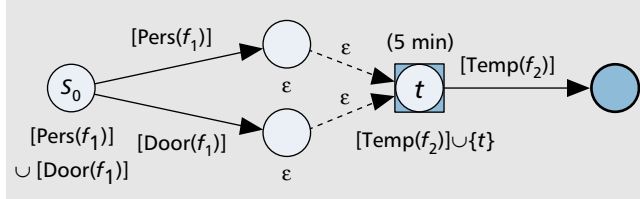
For example, in scenario 1, the building services manager subscribes to temperature events under 15°C in an occupied room. (For simplicity, we consider the movement, door, and temperature events of only a single room, although multiple rooms could be represented in a single CE expression using parameterization. We also prefilter the `PersonEvents`, limiting repeat notifications to 1/min.) We consider a room to be occupied if it has exhibited movement or door events within the last 5min. In our core CE language, we might represent the primitive events (as exposed by JMS) using:



■ Figure 6. Illustration of interval timestamps for events.

- [PersonEvent(location='Office FE02')]
- [DoorEvent(location='Office FE02')]
- [TempEvent(location='Office FE02' AND temp<15)]

Scenario 1 could be written (using [Pers(f_1)], [Door(f_1)], [Temp(f_2)] for the expressions above, for brevity) as: ([Pers(f_1)] | [Door(f_1)], [Temp(f_2)])_{t=5min}. This would be compiled into the following automaton, for use as a detector:

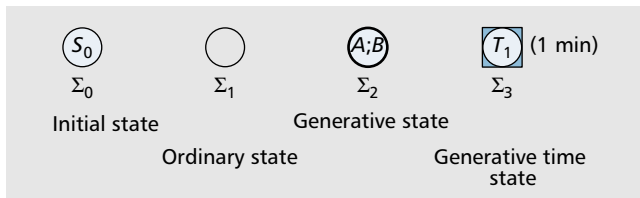


We based our core language and automata on regular expressions and FSAs for a number of reasons. First, their expressive power is well understood, but they require only limited predictable resource usage, and are thus a safer tool for distributed detection than a more general language. Still, they are powerful, and frequently used for pattern detection and matching. Furthermore, regular expressions may also easily be factorized into subexpressions for distributed detection of independent expressions. Finally, we felt it would be more sensible to extend the commonly accepted regular expression operators with necessary additions rather than arbitrarily define new operators, with the concomitant risks of redundancy or incompleteness. Our core CE language and automata therefore only minimally extend regular expressions and FSAs, to allow temporal relationships, input filtering, and parallel detection to be expressed.

Composite Event Detectors

The automata that detect CEs contain a finite number of states and state transitions, but each state also maintains the timing information of the previous symbol detected. In a given state, the automaton decides when to make the transition to another state by considering new input symbols only from a per state describable subset of the global input sequence I .

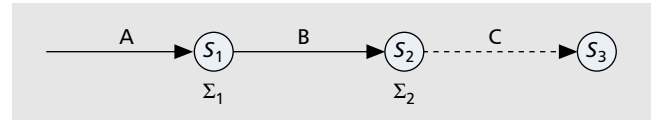
Structure of Automata — Our automata have two types of state: *ordinary* and *generative*. A generative state causes a new event to be created, either a composite of the events matched so far (with a specified type), or an instantaneous time event in the future (with a freshly allocated local identity). The timestamp of the composite event will start at the earliest start time of the constituent events, and end at the latest end time. A time event may be used later in the automaton to progress or fail after a given timeout. Each state has an *input domain of describable events*, the family of events it can match. When in a given state, the automaton processes only those new events that lie within the state's domain. The diagram below shows four states: an initial (ordinary) state, an ordinary state, a generative state for a composite event of type $A;B$, and a generative state for a time event. The input domains are $\Sigma_0 \dots \Sigma_3$.



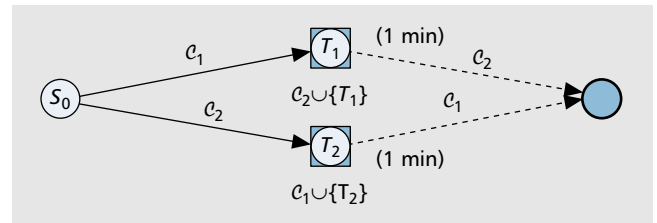
Each state can have any number of outward transitions. There are two types of transition, *strong* and *weak*, which can match events that strongly or weakly follow the previously

detected event. These correspond to the partial and total event orderings $<$ and \prec , respectively. Each transition has a describable family of events attached, any of which will cause it to be taken.

New events in the input domain of a state but not in any transitions will cause the match to fail. These new events must strongly follow the previous event if all outgoing transitions are strong, or weakly follow otherwise. If there are two or more matching transitions, they will be followed nondeterministically. When a state with no outgoing transitions is reached, an event is generated if it is generative; then the machine (or the current nondeterministic branch) immediately terminates. The diagram below illustrates both strong and weak state transitions. If a, b, c are the events that matched A, B and C , $a < b$ and $b \prec c$. Furthermore, $b \in B$ is the first event in the input stream $IS1$ for which $a < b$; similar constraints apply to c .



Limitations of Automata — The extended automata address many of the disadvantages of standard FSAs. First, temporal support is provided by explicit event timestamps and special timer events. Concurrent events are also supported; the following automaton generates a new event when composite events $c_1, c_2 \in \mathcal{D}$ occur in parallel within 1 min of each other. c_1 might represent “Peter is seen in the building but not in his office” and c_2 “Peter’s phone rings.” The resulting event could be used to divert the call to wherever Peter was last seen.



Conventional FSAs have other limitations too. Most important, they cannot handle event interrelationships such as event parameterization. For example, to detect how long each door in a building is left open, a mechanism is needed to express free parameters that apply the same expression to all rooms: detect opening (x) followed by closing (x). Our framework can resolve this issue by filtering the CE attributes of all opening and closing event pairs as soon as they are detected, reporting only matching pairs. This is still efficient, since the unnecessary composites are discarded as soon as they are detected, and every possible pairing would have been considered.

Finally, when nondeterministic FSAs are made deterministic, the number of states can grow exponentially. Although our automata potentially exhibit this behavior, it does not happen in practice: since distribution takes place at the level of CE expressions, not automata, resolution to deterministic automata is not required; instead, a list of active states is held. Furthermore, in typical composite expressions this list is usually short, since distributed detection makes parallel detection of independent subexpressions the norm.

Formal Definition of Automata — Each automaton consists of a set of states S , state domains $a_S: S \rightarrow \mathcal{D}$, and strong and weak transition domains $a_{TS}, a_{TW}: S \times S \rightarrow \mathcal{D}$. There is also a start

state $S_0 \in S$. Finally, $G \subseteq S \times (T \uplus \mathcal{D})$ defines the generative states (an extension of accepting states) and their actions.

The current state of an automaton is $C \subseteq S \times T \times P(E)$ where T is the set of possible timestamps. In other words, the current state consists of a number of triples, each representing a state, a timestamp, and a list of detected events. From the perspective of the automaton, the list of events is opaque, except that extra detected events may be added to it when a transition is made.

Core Composite Event Language

A CE language allows expression of CE patterns. In this section we introduce our core CE language, which can easily be compiled into automata but is still human readable, and outline its grammar. This language also defines the level at which subexpressions are chosen for distributed detection. Appendix B contains the transformation from expressions into automata, and gives precise operator semantics. The operators of the core CE language extend those found in regular languages, namely concatenation, alternation, and iteration, with operators for timing control, parallelization, and weak/strong event sequencing. In contrast to other CE languages, we avoid redundant operators to simplify analysis.

Atoms — $[A, B, C, \dots \subseteq \Sigma_0]$. Atoms detect individual events in the input stream. Here, only events in $A \cup B \cup C \cup \dots$ will be successfully matched. Other events in Σ_0 will cause failed detection, and events outside Σ_0 will be ignored. We abbreviate negation using $[\neg E \subseteq \Sigma]$ for $[\Sigma \setminus E \subseteq \Sigma]$, and also write $[E]$ instead of $[E \subseteq E]$. (Negation ensures any other events in Σ will stop detection, such as timeouts or stopper events.)

Concatenation — $\mathcal{C}_1\mathcal{C}_2$. Detects expression \mathcal{C}_1 *weakly* followed by \mathcal{C}_2 .

Sequence — $\mathcal{C}_1; \mathcal{C}_2$. This detects expression \mathcal{C}_1 *strongly* followed by \mathcal{C}_2 . Thus, \mathcal{C}_1 and \mathcal{C}_2 must not overlap in a sequence, but may in a concatenation.

Iteration — \mathcal{C}_1^* . Detects any number of occurrences of expression \mathcal{C}_1 . If \mathcal{C}_1 detects a symbol that causes it to fail, \mathcal{C}_1^* will fail too. (So $[A][A \subseteq \{A, B\}][C]$ would match input AAC but not $AABC$.)

Alternation — $\mathcal{C}_1 | \mathcal{C}_2$. This expression will match if either \mathcal{C}_1 or \mathcal{C}_2 is matched.

Timing — $(\mathcal{C}_1, \mathcal{C}_2)_{T_1 = \text{timespec}}$. The timing operator detects event combinations within, or not within, a given interval. The second expression \mathcal{C}_2 can then use T_1 in its event specification.

Parallelization — $\mathcal{C}_1 || \mathcal{C}_2$. Parallelization detects two CEs in parallel, and succeeds only if both are detected. Unlike alternation, any order is allowed, and the events may overlap in time.

The following examples illustrate the use of the core CE language to describe CEs. Let B be the events corresponding to “Brian enters the room,” let P be “Peter enters the room,” and let A be “anyone enters the room.”

1. Brian enters the room followed by Peter: $[B]; [P]$
2. Brian enters the room before Peter: $[B \subseteq \{B, P\}]$
3. Brian enters and Peter follows within an hour: $([B], [P \subseteq \{P, T_1\}])_{T_1 = 1h}$
4. Someone else enters the room when Brian is away: $[B] [\neg B \subseteq A] [B]$

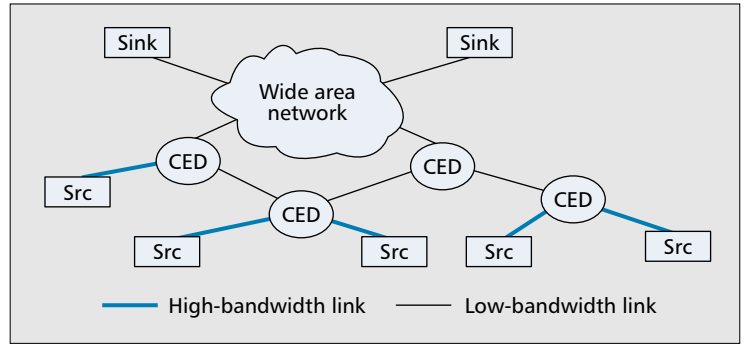


Figure 7. Illustration of distributed composite event detection.

Distributed Detection

In a large-scale distributed application, events are published at geographically dispersed sites. A centralized composite event detector would have to subscribe to all primitive events that are part of a CE expression in order to detect occurrences of composite events. This could become a bottleneck and a single point of failure.

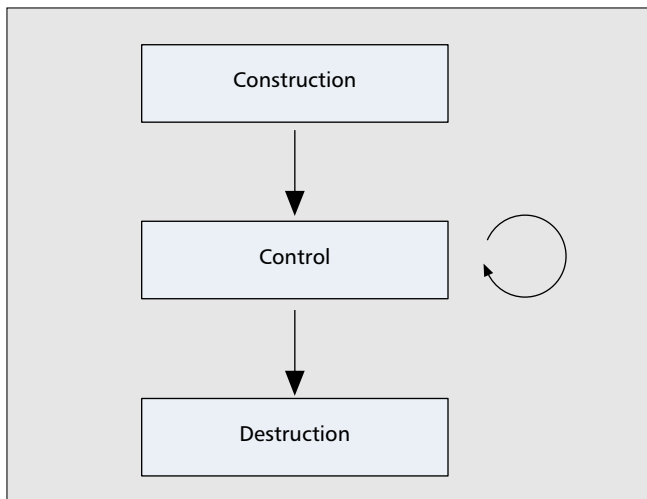
Instead, our framework provides a mechanism for distributing CE detectors. Detectors can be installed at various locations in the network and cooperate with each other. This cooperation is achieved by decomposing CE expressions stated in the core language into subexpressions that are then detected by detectors running at different nodes. Figure 7 shows an example of a network of cooperating CE detectors. The detectors are located close to event sources that publish events at a high rate, thus requiring high-bandwidth links. After CE detection, bandwidth consumption is reduced since CEs occur less frequently. Composite events are then sent to remote event sinks over a low-bandwidth wide area network. No CE detector is overwhelmed by the rate of primitive events, as it subscribes to at most two event sources.

The main difficulty when distributing detectors is to decide on their optimal placement within the system. This is complicated by the fact that the reasons for distributing detectors are potentially conflicting. For example, to minimize bandwidth usage, existing detectors should be reused for subexpressions as much as possible — even between applications, if appropriate. However, if minimum latency is required, detectors should be replicated at various regions in the network, which leads to higher bandwidth consumption. As a result, an optimal solution must be a trade-off that takes the static and dynamic characteristics of the system and the requirements of the application into account.

In our framework, mobile CE detectors detect CEs in a distributed fashion. A distribution policy ensures that detectors are installed at sensible locations, and specifies a policy for their movement and behavior during their lifetime. Network delays that can lead to incorrect detection are addressed by a detection policy.

Mobile Composite Event Detectors

We introduce the concept of a *mobile CE detector* to add distributed detection to our framework. A mobile CE detector is an agent-like entity encapsulating an automaton that detects an expression from our core language. Its actions are governed by distribution and detection policies, discussed below. It subscribes to event sources to receive event input streams and publishes the CEs detected by the automaton. The detector is capable of moving from one location to another in the network. This assumes the existence of a logical overlay network of nodes that supports the migration of components. Consequently, our work is built on top of a network of *event brokers* (e.g., corresponding to Hermes brokers or JMS



■ Figure 8. Lifetime of a mobile composite event detector.

nodes), where each event broker can host one or more mobile detectors.

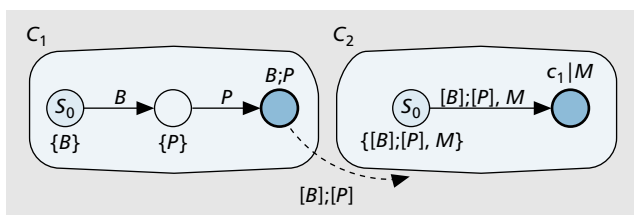
When a client submits a new CE subscription, a new mobile CE detector is created at an existing broker that is then responsible for the detection of this expression. The lifetime of a mobile CE detector is summarized in Fig. 8. At *construction* time, it sets up the detection of the new CE expression. Once CEs for the new expression are being detected, it enters a *control phase*, during which it optimizes the detection process by adapting to dynamic changes in the environment, and making sure it maintains compliance with its policy. Finally, it enters the *destruction phase* when it is no longer needed because, for example, other detectors have taken over or all event clients have unsubscribed.

During its lifetime, a mobile CE detector can carry out several actions:

- It can *create new automata* for the detection of new CE expressions or any of its subexpressions.
- For distributed detection, it can decompose the CE expression along its abstract syntax tree and *delegate detection* of subexpressions to other (already existing) detectors.
- It can *migrate* to another node in the network that is, for example, closer to the event sources to which it has subscribed.
- Finally, the detector can *destroy* itself once it is no longer required.

Mobile CE detectors can communicate with each other using the pub/sub infrastructure. This is usually implemented as a special event type or topic to which all detectors subscribe. A system administrator can influence the strategy of all mobile CE detectors by using this topic to change policy decisions of detectors.

Consider the Active Office application scenario introduced earlier. Let B be the event type corresponding to “Brian enters the room,” P be “Peter enters the room,” and M be “a meeting takes place in the room.” A user is interested in occurrences of “Brian enters the room followed by Peter.” The corresponding mobile CE detector \mathcal{C}_1 for the expression $[B]; [P]$ is shown below.



When another user subscribes to occurrences of the CE “Brian enters the room followed by Peter or a meeting takes place,” this new expression $([B]; [P])|M$ can be rewritten as $\mathcal{C}_1|M$. Therefore, the new detector \mathcal{C}_2 can reuse the existing detector \mathcal{C}_1 by subscribing to $[B]; [P]$. The communication between the two detectors happens exclusively through the underlying pub/sub system.

Distribution Policies

The behavior of a mobile CE detector with respect to its actions is governed by a *distribution policy*, a set of heuristics to be followed by the detector. We identify three independent dimensions that help to limit the space for defining distribution policies.

Decomposition — The degree of decomposition of the composite event expression must be stated in the policy (with optional hints from the application). In order to reuse existing detectors in the system, an expression may have to be decomposed into subexpressions. Decomposition may increase the reliability of detection if multiple detectors are detecting overlapping expressions. For load balancing reasons, a complex expression may be decomposed into manageable subexpressions. The degree of decomposition ranges from no decomposition to full decomposition, where every possible subexpression is factored out. Some policies allow decomposition only when there already exist detectors that can be reused for a subexpression.

Reuse — This dimension specifies to what extent already existing detectors are reused for a new composite event expression or any of its subexpressions. Not reusing existing detectors can result in more reliability, whereas maximum reuse will save bandwidth and computational effort. In situations in which detection latency is important, only local detectors that are in close proximity should be reused.

Locality — The location of new mobile CE detectors must be determined. For certain scenarios bandwidth usage can be reduced by moving detectors as close to primitive event sources as possible. Primitive events that constitute a CE may be of interest only to the CE detector and should thus not be widely disseminated throughout the entire system unnecessarily. This is called *publisher locality*. The opposite approach is to put new CE detectors close to application components that subscribe to them to improve reliability and detection latency. This leads to a policy with *subscriber locality*.

In practice, only certain combinations of these three dimensions will result in useful distribution policies. Table 1 summarizes five example policies, each of which attempts to optimize a different metric in the composite event framework.

Minimum Latency Policy — The detection latency is minimized by placing new detectors as close to subscribers as possible. Composite event expressions should not be decomposed into subexpressions as this would increase the detection latency. Similarly, an existing detector should only be reused if it is close to the subscriber and detects exactly the required CEs.

Minimum Bandwidth Policy — Bandwidth consumption is minimized by placing the detectors close to the primitive event publishers, leveraging the filtering aspect of CE detectors. In addition, existing detectors should be used as much as possible so that no new traffic is generated. The reuse of subexpressions may lead to decomposition.

Minimum Impact Policy — This policy minimizes the impact new detectors have on the entire system. This involves minimizing bandwidth, as before, but also means that computational load should be spread out evenly among detectors. Therefore, new detectors do not have locality, but existing detectors should be maximally reused.

Minimum Load Policy — The fourth policy minimizes the load on composite event detectors by decomposing an expression into the smallest possible subexpressions and distributing them evenly among detectors in the system. It attempts to reuse already existing detectors.

Maximum Reliability — The last policy makes CE detection more resistant to node failure by instantiating redundant detectors for extra reliability. Old detectors are reused only when at least two already exist, new detectors are created otherwise. (This “at least 2” partial reuse policy lies between no reuse and full reuse in Fig. 9.) To limit extra points of failure, detectors are decomposed for reuse only, and no locality restrictions are imposed on new detectors.

Note that a distribution policy is associated with a particular CE expression, so that every mobile CE detector can have its own policy. This enables event subscribers to specify a desired distribution policy at subscription time depending on application requirements.

The effectiveness of distribution policies can be enhanced when mobile CE detectors are able to obtain network- and system-specific parameters such as the current load of a broker node or the communication latency to a particular publisher. A mobile CE detector may use this information to optimize detection in compliance with its distribution policy.

Detection Policies

In a distributed system, events from different event sources travel along separate network routes to a mobile CE detector. Even if we assume that the network itself does not reorder events, out-of-order arrival of events at the detector can occur because of the different associated network delays. Whenever a new event arrives, it has to be inserted at the correct position in the totally ordered event input stream before the stream is fed into the automaton.

The problem is to decide when the next event in the event input stream can be safely consumed by the automaton without risking that an event with an older timestamp is still being delayed by the network. Premature consumption could lead to incorrect detection or nondetection of a CE. Thus, each CE subscription is annotated with a detection policy that specifies when a detector can consume an event from an event input stream.

Best Effort Detection — A best effort detection policy states that events are consumed from event input streams without delay. Whenever an event is available, it causes a state transition (or failure) in the automaton. Although this policy may lead to incorrect detection, it can be applied by applications that are sensitive to detection delay and willing to ignore false positives.

Guaranteed Detection — Under a guaranteed detection policy, an event is consumed from an event input stream only once it has become stable¹ [16]. The consumption of only stable

¹ An event is stable if there is no other event with an earlier timestamp in the system that should be part of this event input stream and thus consumed instead.

Policy name	Decomposition	Reuse	Locality
Minimum latency	None	With locality only	Subscribers
Minimum bandwidth	For reuse only	Maximum	Publishers
Minimum impact	For reuse only	Maximum	None
Minimum load	Maximum	Maximum	None
Maximum reliability	For reuse only	At least 2	None

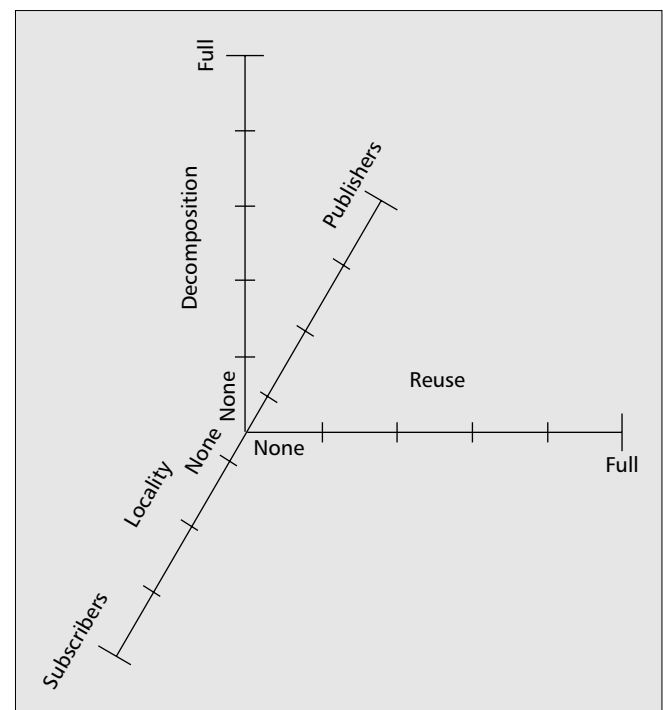
■ Table 1. Example of three distribution policies.

events ensures that no spurious CEs are detected. A detector knows that an event is stable after another event with a later timestamp from the same event source has been inserted in the event input stream. An event source that does not publish events at a high enough frequency can publish dummy *heartbeat events* that are used to “flush the network.”

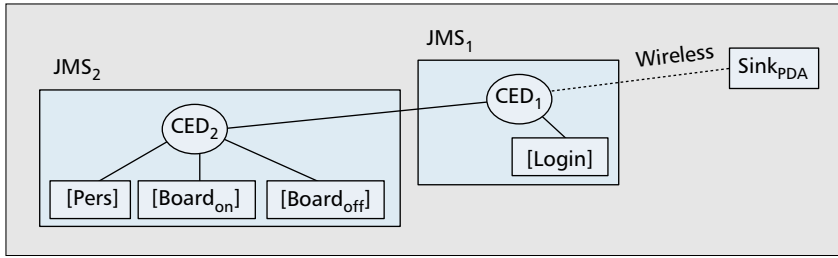
In an asynchronous distributed system, a guaranteed detection policy potentially introduces an unbounded delay at the detector. For instance, an event source might fail or decide not to cooperate by not sending heartbeat events. To avoid this problem, we are currently investigating a *probabilistic stability metric*. As opposed to a simple binary stability measure, a detector attempts to model the probability that a particular event in an event input stream is stable, and the event is only consumed if its stability metric is above a certain threshold.

Implementation using JMS

This section describes the implementation and performance results of our CE framework over JMS using JORAM [19], an open-source implementation of the JMS API. Our prototype implementation is available on the Web [20]. Application programs can publish and subscribe to composite events using the `DistCEDServiceInf` interface, presented in Fig. 5, provided by the event brokers in the system. In the pub/sub messaging model supported by JMS, a publisher registers a topic with a particular JMS provider, such as a JORAM or J2EE server.



■ Figure 9. Design space for distribution policies.



■ Figure 10. Implementation of scenario 2 using the CE framework.

Whenever a message is published on the topic, topic subscribers are notified by the JMS provider via a callback mechanism. Content-based filtering on the fields in the message header is supported.

Although a JMS provider can be a distributed service, most current implementations are centralized, although they may provide redundancy through replication and clustering. Clients may need to connect to several providers, such as local and remote message servers. Therefore, the binding of our CE framework to JMS does not assume that all events (primitive or composite) have a single JMS provider. Instead, our implementation uses a JNDI directory to look up the JMS server for a particular topic. For CEs, we use this to ensure that we establish only a single CE detector for a given CE type, since all such detectors will produce the same events. Since the directory may itself be distributed, this does not imply unnecessary centralization.

To support automatic distribution of CE detection, all event brokers subscribe to a common *administration topic* (DistCEDAdminTopic) that is hosted by an *admin JMS server*. When a new CE expression needs to be detected, the event brokers collectively decide how and where to instantiate the mobile CE detector (with the expression's automata). The locations of newly created detectors are registered with the JNDI directory. In the following experiments, the distribution policy is a simple choice function derived from the hash of the CE type name, although the more complex policies outlined earlier would also be possible.

Evaluation and Results

To test the CE framework implementation on JMS, we simulated Active Office scenario 2 described earlier. The movement of people was treated as a Markov process, with a probability matrix describing the likely movements in each time interval. We used the office layout shown in Fig. 2 with nine rooms and 15 occupants. Eight of the occupants were classed as residents, predisposed to use the offices, while the remainder were visitors, preferring the meeting rooms. The event sinks in the scenario were PDAs connect-

ed by an (expensive) wireless link with limited bandwidth. The goals of the experiment were to minimize the usage of that link and achieve low notification delay for CES.

The CE subscription \mathcal{C} presented to our CE framework as the subscription submitted by the event sinks was as follows:

$$\mathcal{C} \equiv ([\mathcal{C}_1(f_1)], [T_1] \subseteq \{T_1, \text{Login}(f_2)\})_{T_1=5\text{min}} \quad (1)$$

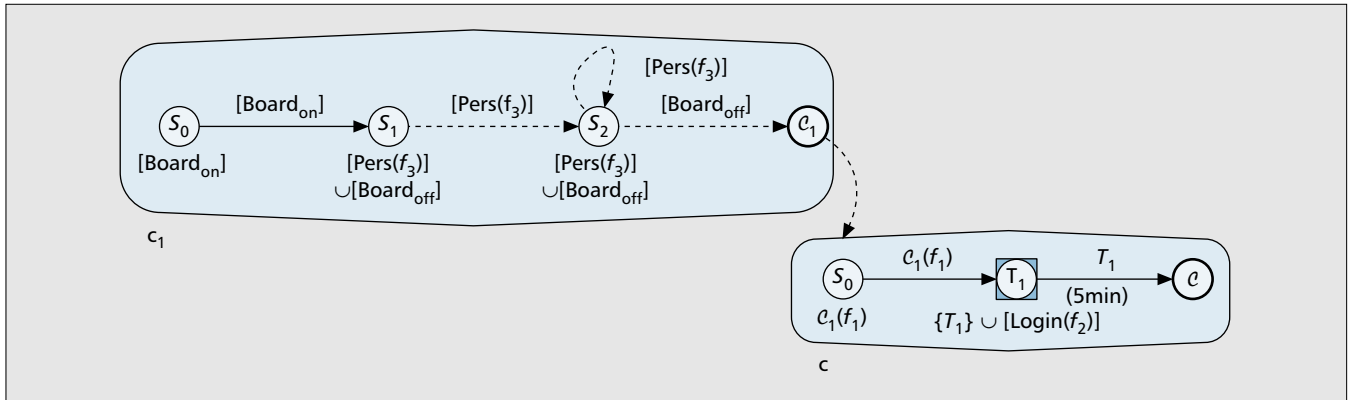
$$\mathcal{C}_1 \equiv [\text{Board}_{\text{on}}] [\text{Pers}(f_3)] [\text{Pers}(f_3)]^* [\text{Board}_{\text{off}}] \subseteq \{\text{Pers}(f_3), \text{Board}_{\text{off}}\} \quad (2)$$

where f_{1-3} are JMS filter expressions, omitted for brevity. Figure 10 shows how the detection of \mathcal{C} was distributed over two event brokers by the CE framework. The detector CED_2 was responsible for the subexpression \mathcal{C}_1 . All primitive events to which it subscribed and the resulting CEs were located on the server JMS_2 . CED_1 then detected the complete expression \mathcal{C} and output its CEs to a different server JMS_1 .

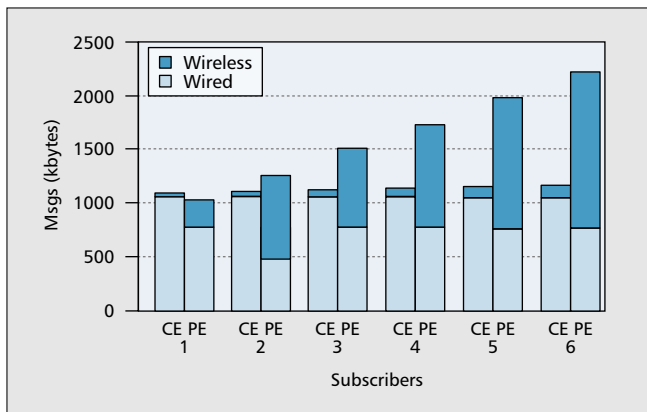
The automata used to detect expressions \mathcal{C} and \mathcal{C}_1 are shown in Fig. 11; these were obtained by applying the transformations of Appendix B to the expressions, and optimizing the results by removing empty or redundant transitions. For example, if the detector for \mathcal{C}_1 in Fig. 11 received a $[\text{Board}_{\text{on}}]$ event, it would accept the event and advance to state S_1 , and also instantiate a new detector in state S_0 . If $[\text{Board}_{\text{off}}]$ were received next, the detector in state S_1 would process this but fail to match it, since $[\text{Board}_{\text{off}}]$ is in the input alphabet of S_1 but lies on no outgoing transitions, so the failed partial match would be discarded. On the other hand, the detector in state S_0 would ignore $[\text{Board}_{\text{off}}]$ events since these lie outside its input alphabet.

Conversely, the sequence of events $[\text{Board}_{\text{on}}] [\text{Pers}(\text{name}=\text{'Peter'}, \text{location}=\text{'Meeting Room 1'})] [\text{Board}_{\text{off}}]$ would be matched by \mathcal{C}_1 , generating a new event which the pub/sub system would forward to \mathcal{C} .

We compared our CE framework (CE) against a JMS-only solution (PE), in which the wireless PDAs subscribed to all the primitive events and performed the CE detection themselves in an ad hoc manner. Figure 12 shows the total data transferred over the wireless and wired networks with a changing number of subscribers. As expected, there is a small overhead when using our CE framework for a single subscriber. However, as the number of subscribers increases, less data needs to be sent over the wireless network because CE detectors can be reused. For six subscribers,



■ Figure 11. Automata for CE expressions \mathcal{C} and \mathcal{C}_1 .



■ Figure 12. Amount of data sent.

our CE framework generates 53 percent of the total traffic generated by the primitive ad hoc solution, and only 8 percent of the wireless traffic. Note that the traffic over the wired network stays roughly constant as it is mainly caused by primitive event sources sending messages to the JMS servers.

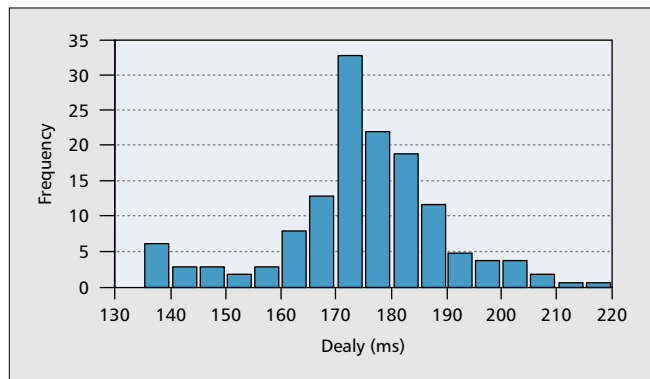
The additional notification delay introduced by our CE framework is small. The plot in Fig. 13 shows the distribution of delay it takes for a subscriber to be notified of an occurrence of \mathcal{C} after the CE logically happened in the system (i.e., its last primitive event was published). The notification delay stays below 220 ms and is fairly constant during the course of the experiment.

Further Work: Higher-Level Specification Languages

When designing a language for the specification of composite events for ubiquitous applications, two conflicting requirements arise. Primarily, the language should facilitate the implementation of efficient detectors and be decomposable for distributed detection (i.e., the language should be optimized to be *machine-processable*). On the other hand, the syntax and semantics of the CE language should be clean and intuitive so that it is *human-processable*. Therefore, we introduce the idea of *higher-level specification languages* for humans to express composite events in a natural and domain-dependent way. These languages are then compiled down into our automata. Whereas our core CE language is optimized for machine detection, the higher-level languages focus on CE specification by end users or programmers. The following are three examples of such languages.

The Pretty Language — The “pretty” language has a verbose syntax similar to many current rule-based specification languages. It does not have a minimal set of operators. CE specifications in the pretty language, such as “Event A” followed by “Event B” within “1h” resemble English language statements, making it easier for nonprogrammers to express CEs.

Programming Language Binding — A binding of CEs to a programming language such as C++ or Java attempts to hide CE specification by integrating it into the programming language, making its usage easier for programmers. This can be achieved with a sequence of method calls on event objects that build a CE expression: `eventA.after(eventB.repeated(3))`. At runtime, these method calls are translated into a core CE language expression.



■ Figure 13. Delay distribution of \mathcal{C} .

Graphical Composition — In the Active Office, users may interact with the system at runtime by specifying its behavior with rules based on ces such as “Turn off the office light after 7 p.m.” A graphical composition tool could be used based on a simple model familiar to users. For instance, CE streams could be visualized as water flows in pipes, allowing different types of piping to be composed to build CEs.

Conclusions

In a world with many mobile entities and complex Internet-based applications, events will become the dominant communication paradigm. CE detection in these large-scale systems provides a means of managing the complexity of a vast number of events. We consider our work a first step in facing this challenge, providing novel scalable middleware services such as generic CE detection.

In this article we have presented a general CE detection framework as an extension of an existing pub/sub middleware. The framework assumes a realistic interval-based time model, and its event model makes few assumptions about the pub/sub communication infrastructure employed. Our CE detectors are an easily implementable extension of conventional FSAs. They can handle timestamps, concurrent events, and come with a core CE language that is expressive and decomposable. Higher-level specification languages can provide more domain-specific ways to specify composite events. The abstraction of mobile CE detectors allows distributed CE detection, making the framework more scalable and robust. We introduce the concept of distribution and detection policies that control the distributed behavior of detectors. Finally, the implementation of our CE framework over JMS demonstrates that it can improve performance in a real pub/sub application, compared to client-side JMS subscriptions.

References

- [1] P. T. Eugster *et al.*, “The Many Faces of Publish/Subscribe,” *ACM Comp. Surveys*, vol. 35, no. 2, 2003, pp. 114–31.
- [2] A. Carzaniga, D. S. Rosenblum, and A. L. Wolf, “Design and Evaluation of a Wide-Area Event Notification Service,” *ACM Trans. Comp. Sys.*, vol. 19, no. 3, Aug. 2001, pp. 332–83.
- [3] IBM T J Watson Res. Ctr., *Gryphon: Publish/Subscribe Over Public Networks*, White paper, 2002.
- [4] P. R. Pietzuch and J. M. Bacon, “Hermes: A Distributed Event-Based Middleware Architecture,” *Proc. 1st Int'l. Wksp. Distrib. Event-Based Sys.*, Vienna, Austria, July 2002, pp. 611–18.
- [5] P. R. Pietzuch, B. Shand, and J. Bacon, “A Framework for Event Composition in Distributed Systems,” *Proc. 4th Int'l. Conf. Middleware, LNCS*, vol. 2672, Rio de Janeiro, Brazil, June 2003, pp. 62–82.
- [6] M. Steinder and A. Sethi., “The Present and Future of Event Correlation: A Need for End-to-end Service Fault Localization,” *Proc. IIS SCI: World Multi-Conf. Systemics Cybernetics Informatics*, Orlando, FL, 2001.
- [7] M. Addlessee *et al.*, “Implementing a Sentient Computing System,” *IEEE Comp.*, vol. 34, no. 8, Aug. 2001, pp. 50–56.
- [8] N. H. Gehani, H. V. Jagadish, and O. Shmueli, “Event Specification in an Active Object-Oriented Database,” *Proc. ACM SIGMOD Int'l. Conf. Mgmt. of Data*, 1992, pp. 81–90.

- [9] S. Gatzia and K. R. Dittich, "Detecting Composite Events in Active Database Systems Using Petri Nets," *Proc. 4th RIDE-AIDS*, 1994, pp. 2–9.
- [10] S. Chakravarthy and D. Mishra, "Snoop — An Expressive Event Specification Language For Active Databases," Tech. rep. UF-CIS-TR-93-007, Dept. of Comp. and Info. Sci., Univ. of FL, 1993.
- [11] S. Schwiderski, "Monitoring the Behavior of Distributed Systems," Ph.D. thesis, Comp. Lab., Univ. of Cambridge, 1996.
- [12] M. Mansouri-Samani and M. Sloman, "GEM: A Generalised Event Monitoring Language for Distributed Systems," *IEE/IOP/BCS Distrib. Sys. Eng. J.*, vol. 4, no. 2, June 1997, pp. 96–108.
- [13] R. Hayton, "OASIS — An Open Architecture for Secure Interworking Services," Ph.D. thesis, Comp. Lab., Univ. of Cambridge, 1996.
- [14] J. Bacon *et al.*, "Generic Support for Distributed Applications," *IEEE Comp.*, Mar. 2000, pp. 68–77.
- [15] G. Banavar *et al.*, "Information Flow Based Event Distribution Middleware," *Middleware Wksp. ICDCS '99*, 1999, pp. 114–21.
- [16] C. Liebig, M. Cilia, and A. Buchmann, "Event Composition in Time-dependent Distributed Systems," *Proc. 4th Int'l. Conf. Cooperative Info. Sys.*, 1999, pp. 70–78.
- [17] Sun, Java Message Service, <http://java.sun.com/products/jms/>, 2001.
- [18] P. R. Pietzuch and B. Shand, "A Framework for Object-Based Event Composition in Distributed Systems," Presented at PhDOS Wksp. (ECOOOP '02), <http://www.cl.cam.ac.uk/Research/SRG/opera/pub/phd00s02-ced.pdf>, Malaga, Spain, June 2002.

- [19] ObjectWeb Open Source Middleware, JORAM Java Open Reliable Asynchronous Messaging, 3.2.0 rel., <http://www.objectweb.org/joram>, Oct. 2002.
- [20] The Opera Group, DistCED Prototype Implementation, <http://www.cl.cam.ac.uk/Research/SGR/opera/projects/DistCED>, May 2003.

Biographies

PETER PIETZUCH (Peter.Pietzuch@cl.cam.ac.uk) is a final year Ph.D. student at the University of Cambridge Computer Laboratory. Before starting his Ph.D., he received a B.A. degree in computer science from Cambridge University. His primary research interest is middleware support for large-scale distributed systems, focusing on event-based middleware and peer-to-peer application-level routing.

BRIAN SHAND (Brian.Shand@cl.cam.ac.uk) is a Ph.D. student at the University of Cambridge Computer Laboratory. His research currently focuses on trust-based middleware for distributed computational services. Before this, he investigated distributed objects for image processing applications, for an M.Sc. at the University of Cape Town.

JEAN BACON [SM] (Jean.Bacon@cl.cam.ac.uk) is a reader in distributed systems at the University of Cambridge Computer Laboratory. She has acted as Editor in Chief of *IEEE Distributed Systems Online* since its start in 2000. She is a member of the Board of Governors of the IEEE Computer Society.

Appendix A: Formalizing the Time and Event Models

Definition of Interval Timestamps

This appendix formalizes our notion of an interval timestamp; the next presents our model of the event subscriptions available to the CE service, in terms of describable events and event input sequences.

Conventional timestamps are often inappropriate for distributed event systems. In a distributed system, node clocks may have unknown jitter within a known synchronization distance. As a result, if two nodes detect events A and B , respectively, it may be impossible to decide which occurred first. An interval timestamp, consisting of a start time and an end time, can make this ambiguity explicit, yet remains consistent with the physical time order of the events [16].

Let $t = [t^l; t^h]$ be an interval timestamp with start and end times t^l and t^h ($t^l \leq t^h$). We define the order relations $<$ and \prec and union operator \cup as

$$t_1 < t_2 \triangleq t_1^l < t_2^l \quad (3)$$

$$t_1 \prec t_2 \triangleq (t_1^l < t_2^l) \vee (t_1^h = t_2^h \wedge t_1^l < t_2^l) \quad (4)$$

$$t_1 \cup t_2 \triangleq [\min(t_1^l, t_2^l); \max(t_1^h, t_2^h)] \quad (5)$$

Formalizing Describable Events and Input Sequences

Users of event systems subscribe for notification of relevant events. Our CE detectors use the same subscription mechanism to describe which events they need to receive. In a sense, therefore, subscriptions (and the associated filter expressions) represent the atomic input streams available to CE detectors.

Let $E = \{e_1, e_2, \dots\}$ be the space of possible events in the system. Each event e has timestamp $T(e)$ and a unique identifier $u(e)$ ordered by $<$. We write $e_1^{t_1}$ to show $T(e_1) = t_1$. Events are then ordered consistently with their timestamps:

$$\forall e_1, e_2 \in E, \quad e_1 < e_2 \triangleq T(e_1) < T(e_2) \quad (6)$$

$$e_1 \prec e_2 \triangleq (T(e_1) < T(e_2)) \vee (T(e_1) = T(e_2) \wedge u(e_1) < u(e_2)) \quad (7)$$

The space of events may be further categorized. The special *empty event* $\epsilon \in E$ is always detected. *Time events* $E_T \subseteq E$ are made to occur at a given future instant or after a certain interval when timers expire. With instantaneous timestamps ($t^l = t^h$), they help detect composite events with time restrictions. If not supported by the pub/sub infrastructure, CE detectors can generate them as needed.

Event systems often allow us to differentiate types of event, by subscribing to subspaces of the event space E (e.g., "events where a door opens," or "events where FE04's door opens"). These sets of events are denoted by upper case letters: E, A, B . (In pub/sub systems, these are often called event types.) Individual event instances, on the other hand, take lower case letters: e_1, e_2, a, b .

Subscriptions also need certain properties to be useful for CE detection. For example, if subscriptions A and B are valid, it should be possible to detect events matching both or either of subscriptions $A \cap B$ or $A \cup B$. (If this is not supported by the underlying event framework, it can be simulated by detectors if the event input streams are well ordered together under the total order \prec .)

There should be a maximal subscription E_D of all events that can be matched. There is also a subscription to detect any predefined matchable event alone (and the special empty event ϵ is always matched); see Eq. 8. Finally, CE detectors should also be able to detect events matching one subscription but not another.

Each subscription can be associated with the set of events that would match it. The theoretical collection of all these subscription sets is the family of *describable event sets* $\mathcal{D} \subseteq \mathcal{P}(E)$. This is a special collection of subsets of E : those that can be detected within the CE framework. \mathcal{D} is closed under finite union, finite intersection, and element complementing relative to E_D :

$$\text{Let } E_D = \bigcup_{E \in \mathcal{D}} E. \text{ Then } E_D \in \mathcal{D}, \epsilon \in E_D, \text{ and} \quad (8)$$

$$\forall e \in E_D, \{e, \epsilon\} \in \mathcal{D}$$

The automata that detect CEs need to be able to treat incoming events as a well ordered stream in order to match sequential patterns of events. By totally ordering events with \prec

this can be achieved, resulting in the *global event input sequence* $I = (e_1, e_2, e_3, \dots)$, where $e_n < e_{n+1} \forall n \in N$.

However, not all events are relevant to all patterns or at all stages of a particular pattern. Describable event sets provide partial views of the input events, selecting subsequences of I .

Thus, CE detectors can restrict their view of the input sequence to only the relevant symbols. For example, if $E \in \mathcal{D}$ $I_E = (e_{E1}, e_{E2}, e_{E3}, \dots)$ denotes the subsequence consisting of elements of E .

Appendix B: Generating Composite Event Detectors

This appendix details how expressions in our core CE language are transformed into CE detection automata, as outlined earlier. The grammatical components of our core language are listed below, with corresponding automata.

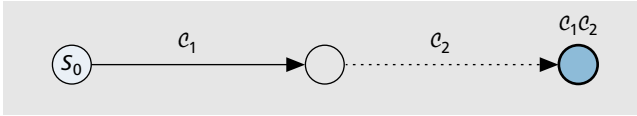
Atoms: $[A, B, \dots \subseteq S_0]$. Atoms detect individual events in the input stream. The resulting automaton considers an input stream of all events that are elements of S_0 . Any input event in $A \cup B \cup \dots$ will be successfully detected; an event in S_0 but not in $A \cup B \cup \dots$ is a failure that stops expression matching (Fig. 14a).

Negation: $[\neg E \subseteq \Sigma] \triangleq [\Sigma \setminus E \subseteq \Sigma]$.

Trivial Input: $[E] \triangleq [E \subseteq E]$.

Concatenation: $\mathcal{C}_1\mathcal{C}_2$. Detects expression \mathcal{C}_1 weakly followed by \mathcal{C}_2 . In the diagram, the shaded boxes are automata matching \mathcal{C}_1 and \mathcal{C}_2 . An empty transition is then added for each generative state of \mathcal{C}_1 or \mathcal{C}_2 , and those states become ordinary (Fig. 14b). If \mathcal{C}_1 's or \mathcal{C}_2 's detection were distributed, each sub-machine could be replaced by a single transition.

Removing empty transitions² gives



² Outgoing transitions from the second submachine's start state inherit the strength or weakness of the empty transition, but keep their original labelings.

Sequence: $\mathcal{C}_1; \mathcal{C}_2$. This detects \mathcal{C}_1 strongly followed by \mathcal{C}_2 . Thus, \mathcal{C}_1 and \mathcal{C}_2 must not overlap in a sequence, but they may in a concatenation (Fig. 14c).

Iteration: \mathcal{C}_1^* . Detects any number of occurrences of expression \mathcal{C}_1 . If \mathcal{C}_1 detects a symbol that causes it to fail, the composite machine \mathcal{C}_1^* stops detecting iterations, even when \mathcal{C}_1 is distributed to another node (Fig. 14d).

Alternation: $\mathcal{C}_1 | \mathcal{C}_2$. This expression will match if either \mathcal{C}_1 or \mathcal{C}_2 is matched by the input stream. This may result in non-determinism for the number of input symbols that are matched by both \mathcal{C}_1 and \mathcal{C}_2 (Fig. 14e).

Timing: $(\mathcal{C}_1, \mathcal{C}_2)T_1 = \text{timespec}$. The timing operator can be used to detect event combinations within, or not within, a given interval. In the above expression, event T_1 will be generated at a certain time after \mathcal{C}_1 is detected, either a relative time, such as a minute later, or an absolute time. The second expression \mathcal{C}_2 may then use T_1 as an event specification in detecting CEs. Furthermore, \mathcal{C}_2 is extended so that all states include T_1 in their input domain. For distribution or reuse, the modified \mathcal{C}_2 detector is treated as distinct from the original, and it should be on the same node as the $(\mathcal{C}_1, \mathcal{C}_2)T_1$ detector (Fig. 14f).

Parallelization: $\mathcal{C}_1 || \mathcal{C}_2$. This can be used to detect composite expressions \mathcal{C}_1 and \mathcal{C}_2 in parallel. The diagram assumes that separate detectors for \mathcal{C}_1 and \mathcal{C}_2 already exist. They must be separate to maintain the two independent timestamps needed for proper order restrictions on the two input sequences (Fig. 14g).

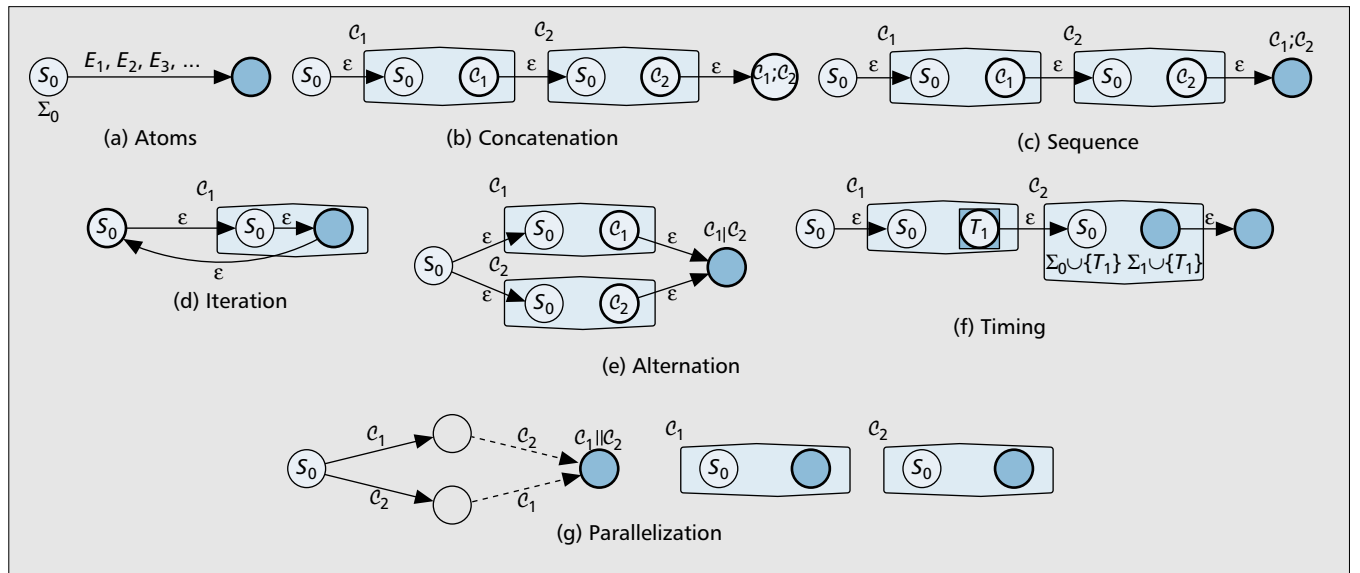


Figure 14. Composite event detectors.