

CORTEX: Towards Supporting Autonomous and Cooperating Sentient Entities

P. Veríssimo, V. Cahill, A. Casimiro,
K. Cheverst, A. Friday, J. Kaiser

DI-FCUL

TR-02-1

February 2002

Departamento de Informática
Faculdade de Ciências da Universidade de Lisboa
Campo Grande, 1700 Lisboa
Portugal

Technical reports are available at <http://www.di.fc.ul.pt/tech-reports>. The files are stored in PDF, with the report number as filename. Alternatively, reports are available by post from the above address.

CORTEX¹: Towards Supporting Autonomous and Cooperating Sentient Entities

P. Veríssimo²
U. Lisboa
pjuv@di.fc.ul.pt

V. Cahill
T.C. Dublin
vinny.cahill@cs.tcd.ie

A. Casimiro
U. Lisboa
casim@di.fc.ul.pt

K. Cheverst
U. Lancaster
kc@comp.lanc.ac.uk

A. Friday
U. Lancaster
adrian@comp.lanc.ac.uk

J. Kaiser
U. Ulm
kaiser@informatik.uni-ulm.de

ABSTRACT

We are now at the point where the emergence of a new class of applications that operate independently of direct human control can be envisaged. However, this is also the crossroads between the requirements put on system support, by the advances of research on high-level models for this class of applications--- e.g. on autonomous agents and distributed AI--- and the shortcomings of current architectures and middleware models.

This paper addresses the latter problem. It does a constructive analysis of the key characteristics of the aforementioned applications, amongst which sentience, studying the interaction types to be supported, and establishing requirements on the programming model, and on distributed architecture. It then describes the approach taken by CORTEX to reach a solution, in the form of infrastructural support to construct large-scale proactive applications and thereby to validate the use of sentient objects as a viable approach to the construction of such applications.

1. INTRODUCTION

Human society, at every level, is increasingly dependent on information. Information systems such as the World-Wide Web are now massively pervasive and critical to the functioning of the global economy. For the most part, current large-scale information systems are *centralised*, in the sense that they are centrally managed and controlled, and *reactive* in that they function primarily by responding to end-user requests. We are now at the point where the emergence of a new class of large-scale *decentralised* and *proactive* applications, i.e., applications that operate independently of direct human control, can be envisaged. However, this is also the crossroads between the requirements put on system support, by the advances of research on high-level models for this class of applications--- e.g. on autonomous agents and distributed AI--- and the shortcomings of current architectures and middleware models.

Two issues render this problem a difficult one. Firstly, most of the recent research has concentrated on functional and behavioural aspects of the participants (objects, agents, etc.), but neatly there is not yet a clear answer to the non-functional requirements put on the supporting substrate. Some characteristics we can anticipate include *autonomy*, *large scale*, *geographical dispersion*, *mobility* and *evolution*. Secondly, the scenario we envisage is that the future is near when mission-critical computer systems will be ubiquitous and pervasive, and comprised of networked components that will act autonomously in responding to a myriad of inputs, to affect and control the surrounding environment. This introduces

¹ This work was partially supported by the EC, under the IST/FET programme, through project IST-2000-26031 (CORTEX - CO-operating Real-time senTient objects: architecture and EXperimental evaluation).

² Contact author: Ph: +351-21-7500103, Fax: +351-21-7500084.

additional non-functional characteristics of these applications, forming a difficult to address combination: *sentience*, *time* and *safety/security* criticality.

Amongst the conditions that are making this possible, two are particularly noteworthy - the availability of improved sensor technology supporting accurate and trustworthy visual, auditory, and location sensing [16,42]; and the emergence of paradigms for reliable, consistent and timely input of sensor data, data dissemination and data fusion, and actuation on the environment [1,2,13]. These developments will enable a new generation of applications in areas such as intelligent vehicles, mobile robotics, smart buildings, and traffic management as well as in more traditional areas such as telecommunications management, process control and C³ (command, control and communications). To accommodate growth and adaptability with respect to number of participants, integration of new services, and quality of service issues, to name but a few, new computational models are needed. These models must be more powerful than the client/server model, which does not reflect the autonomy and spontaneity of co-operating entities. Proactive applications need active components, which are able to sense their environment and spontaneously interact and co-operate with others. Moreover, the communication infrastructure supporting these applications will involve a plethora of different network types and media with widely varying attributes concerning addressing schemes, topology, bandwidth and reliability. However, what is now "the Internet" is far from being the desired solution, since albeit displaying many of these technologies, there is not structure and architecture making sense of them collectively.

CORTEX proposes to devise an architecture and a set of paradigms for the construction of applications composed of collections of what may be called *sentient objects* - mobile intelligent software components that accept input from a variety of different sensors allowing them to sense the environment in which they operate before deciding how to react. Furthermore, the latter organise themselves in autonomous, mobile and rapidly composable co-operating communities. In the future, sentient objects will pervasively be included in almost everything of our daily life. They will ubiquitously integrate all kinds of devices and interact seamlessly amongst themselves in ways that go far beyond what the client/server paradigm, which is supported by current state-of-the-art middleware [17,19,30], allows. Applications will form islands of co-operation inside a wider network universe composed from different physical networks with characteristics ranging from high speed backbones to wireless connections and deeply embedded field buses.

In the long term, society will substantially rely on this technology. To reduce vulnerability and provide a robust failure resilient environment, middleware is required that understands the metaphors of the high-level models, yet keeps the underlying system (still a fragile computer and network system) under correct operational envelopes. The paper addresses this problem. It does a constructive analysis of the key characteristics of the aforementioned applications, amongst which sentience, studying the interaction types to be supported, and establishing requirements on the programming model, and on the distributed architecture. It then describes the approach taken by *CORTEX* to reach a solution, in the form of infrastructural support to construct large-scale proactive applications and thereby to validate the use of sentient objects as a viable approach to the construction of such applications.

2. FUNDAMENTAL CHALLENGES IN SUPPORTING SENTIENT APPLICATIONS

A key enabling technology to realise the vision of ubiquitous computing and proactive applications, is an intelligent middleware supporting appropriate computational models for the envisaged generation of applications. Such middleware must support growth and adaptability to new technologies, and has to provide the hooks for these applications to enforce non-functional quality attributes like reliability and timeliness. In particular, the middleware has to cope with applications that have some or all of the following characteristics:

- Sentience – the ability to perceive the state of the surrounding environment, through the fusion and interpretation of information from possibly diverse sensors;

- Autonomy – components of these applications will be capable of acting in a decentralised fashion, based solely on the acquisition of information from the environment and on their own knowledge;
- Large scale - typical applications may be composed of billions of interacting hardware and software components;
- Time criticality - these applications will typically interact with the physical environment, and will have to cope with its pace, regardless of adverse conditions due to scale and technology shortcomings;
- Safety criticality – typical applications will interact with human users, whose well-being will frequently rely on them;
- Geographical dispersion - unlike current embedded systems, typical applications will integrate components that are scattered over buildings, cities, countries, and continents;
- Mobility – furthermore, they must possess the ability to move between hosts possibly of different networks, while remaining in continuous operation
- Evolution – these applications will have to cope with changing conditions during their lifetimes. Not only must the applications be designed to evolve, but their underlying support must also be adaptable.

Traditional approaches to the design of time and safety critical distributed applications cannot handle the complexity inherent in the scale and geographic dispersion of these new applications. On the other hand, new promising approaches, such as *autonomous decentralised systems* - a subject of active research during the past few years [18,29], are beginning to emerge. The suitability of autonomous decentralised systems is being tested in current attempts to develop applications in areas such as air traffic control, with the free-flight approach, and in the Telecommunications Intelligent Network Architecture (TINA) effort [33].

However, whereas basic technologies exist that make autonomous decentralised systems a possibility, appropriate architectures and paradigms for the construction of the relevant applications are required. Consider applications composed of collections of *sentient objects*: they must be able to discover and interact with each other and with the physical world in ways that demand predictable and sometimes guaranteed quality of service (QoS), encompassing both timeliness and reliability guarantees. Achieving predictability is made difficult by the characteristics of the changing environment in which these objects operate, including an unstable and mobile object population, unpredictable network load, varying connectivity, and the presence of failed system components. Thus, the construction of applications from sentient objects must take account of the fundamental trade-off between the existence of a dynamic environment and the need for predictable operation. To date, no comprehensive technology appropriate to the design and implementation of such applications exists.

Sentient objects will exist at very different levels of abstraction. At the lowest level such objects might represent simple sensors or actuators capable of generating or consuming events. At a slightly higher level of abstraction, a tightly-coupled embedded system that integrates many such simple objects connected via a field bus might represent a single sentient object that is itself capable of generating and/or consuming events as a component of a larger system. Dynamically varying collections of such objects may need to co-operate to form higher-level components or applications that operate over both local and wide area networks including wireless networks.

2.1. Communication, Co-ordination and Control

Independently of the level of abstraction at which we are working, three fundamental problems have to be addressed in order to support applications based on sentient objects: dissemination of information to create common knowledge, mutual awareness and a basis for local decisions; achieving co-ordination amongst peer objects in order to carry out actions in a consistent way; acting upon the environment, changing its state as a result of proactive or reactive decisions.

Consider a traffic scenario where vehicles communicate to provide a look-ahead warning service for vehicles coming from behind. If a vehicle detects an obstacle it sends an alert message which, in turn, the receiving participants can exploit to set new cruising parameters. To guarantee reliable distribution of the alert message and overcome range or reception problems, vehicles receiving this message should themselves further disseminate the message. In such a scenario, we face the following problems:

- The scope of information dissemination is dynamically determined by spatial parameters, i.e. those vehicles directly affected by the obstacle on the road.
- Communication is anonymous, hence group membership is implicit and reliable assessment of who received the message is difficult.
- The information is only valid in a restricted area.
- Many vehicles try to send similar messages, but the system should prevent the communication medium from being overloaded.
- Vehicles, which receive the message, must decide whether it is necessary to continue propagation or to stop.

Following the sentient object model, vehicles can use additional sensor information, like visual input or location information derived from GPS and/or short-range radar distance measurements, to derive decisions. The challenge is to elaborate how the different system levels that are involved will interact. The alert message would raise awareness between vehicles. The next step would be to start co-operation between the vehicles to allow a well-adjusted behaviour of the vehicles or at least to prevent crashes. This only affects a well-defined number of vehicles, which are close to each other. In the simplest case it involves the interaction between one vehicle and its successor. While awareness may be realised as a best effort facility, co-operation needs a guaranteed quality of service between communicating entities. Co-operative actions require common knowledge about the system history, and a consistent view about its future evolution, which also has to be achieved in a timely and reliable manner. This may imply the causal ordering of events, and their timely distribution in a reliable and ordered way. Some co-operation activities, with stronger consistency requirements may also require the availability of consistent membership information, e.g., views of co-operating groups of objects.

Thus, the communication support must be able to deliver a wide range of qualities of service, in terms of both data exchange and membership services. State-of-the-art group communication protocols or generative anonymous communication [8] based on publisher/subscriber models definitely do not tackle these problems [28,32,35] and it is an open question whether these problems can be solved in the basic communication system alone.

2.2. Heterogeneity, Hierarchy and Scope

Again considering the example application outlined above, a hierarchy of communication networks will be present inside a vehicle to eventually convert the decision into deceleration or warning signals. Thus, the cruising parameters resulting from the higher level co-ordination among vehicles have to be set and controlled by networks of intelligent sensors and actuators, that we generically call *controller area networks (CANs)*. In more general terminology, islands of control must co-operate via gateways in a timely and reliable manner, through the global wide area network (*WAN*). This motivates a crucial aspect of the *CORTEX* architecture, what we call a *WAN-of-CANs* structure.

The application example unveils yet another important issue, which is related to limit the range and control the quality of information propagation in the global system. Let us assume the notion of a *zone*. A zone firstly identifies a natural border for the propagation of broadcast messages. This may be the range of a wireless transmitter or a single network section. To send a message beyond this border usually requires extra effort, like relay stations or gateways. In the obstacle-warning example, the zone is defined as a specific distance from the obstacle and this information is used to confine propagation.

An important issue for co-operation is to know what QoS can be sustained by the zones in which participants reside. Typically, a single CAN represents a zone with a very high level of predictability compared to a zone in a wireless network [22,27,36,39,43]. If participants reside in different zones,

communication has to adapt to the weakest guarantees including the losses over the gateways. In a mobile environment where migration from one zone to another zone is likely to happen, it is a great challenge to devise communication mechanisms that dynamically adapt to these changing QoS attributes while maintaining a certain level of guarantee. Paradigms like zoning and topology awareness are relevant in our context, since they allow the heterogeneity of the underlying support to be accommodated, while not necessarily making it visible to the layers above [31,38].

2.3. Predictability and Adaptability

Underlying all of these considerations is the fundamental challenge of coping with the uncertainty of synchrony. In principle, this can be achieved by adaptation. However, while there is an increasing body of research on QoS adaptation [4,5,6], most work has focused on protocol or application-level heuristics, and does not provide any guarantees on how well the system adapts. The applications we intend to support require *predictability* about timeliness. This means that even if the timeliness of the system is degrading, it should do so in a predictable way. In consequence, the coverage of timeliness assumptions should remain stable throughout the application's lifetime [34].

Some more demanding applications will require *guarantees* about timeliness objectives, that is, not only the coverage but also the assumed timeliness bounds should hold. Since timing faults are difficult to prevent in the kinds of complex and large-scale systems that we are considering, this presents us with a fundamental challenge of *avoiding contamination*, i.e., incorrect logical behaviour, when timing faults do occur. This has been shown to be a significant problem even in systems where synchrony expectations are minimal [9]. Ensuring timely system operation despite timing faults, on the other hand, requires timing fault tolerance mechanisms.

The model must recognise that synchronism is not a homogeneous property of a system, neither in time nor in space, and be able to support incremental levels of fault-tolerant real-time behaviour in subsystems of the architecture, from soft to hard real-time. These are challenging problems in large-scale systems with uncertain synchrony, especially where wireless communication is employed. We intend to use and build on previous results on partial synchrony systems, such as the timed asynchronous and quasi-synchronous models [11,37,40].

2.4. Scalability

Scalability represents a crucial transparency property concerning the ability to accommodate growth in a large-scale distributed system. Thus, connecting more participants to the system dynamically, including adding entire additional networks, or providing new services, should not be prevented by factors originating in the system design. The notion of *anonymous event-based computing* is central to addressing the needs of scalable systems in *CORTEX*.

Nevertheless, supporting non-functional attributes like timeliness and reliability guarantees adds new and challenging dimensions to scalability. Consider a simple case. Under purely functional aspects, adding a new subscriber in an anonymous communication system is completely transparent for the publisher of information. However, if temporal guarantees are to be met, the QoS of the underlying network becomes decisive (e.g., timeliness, reliability). For example, a point-to point network will require additional messages to be sent for every new subscriber affecting overall transmission time. Additionally, a distant new subscriber will add significantly to the global delivery delay of a reliable multicast.

CORTEX aims at providing appropriate abstractions to express awareness about the uncertainty and variations of physical message transmission. The recursive WAN-OF-CAN concept and, on a higher level of abstraction, the notion of zones contribute to this goal. Furthermore, we allow applications to exploit this information proactively, e.g., by trading precision against timeliness of information. We

address this problem in the context of adaptation in a partial synchrony model, whilst ensuring stability of assumption coverage.

2.5. Fault Tolerance and Security

Real-life systems deriving from the *CORTEX* approach will require measures enhancing their dependability, both from the fault tolerance and the security aspects. The confluence of ubiquitous wireless and anonymous networking, and of powerful embedded computer systems, provides an interesting spectrum of computer devices and information repositories that poses challenging security and reliability problems.

The problems concerning co-ordination, control, predictability, adaptability, and scalability form the body of research we intend to address in *CORTEX*, and constitute great challenges on their own. Far from this being a disclaimer, we believe that, having solved these challenges, we will have paved the way for the incorporation of known and emerging paradigms in modular and distributed fault tolerance for both large-scale systems and small-scale real-time systems, and in cryptographic multiparty communication and processing.

3. PROGRAMMING PARADIGM

Mobile sentient objects have autonomous behaviour resulting from interactions with the physical environment, i.e. driven by sensor inputs, as well as from the internal state of the objects. Moreover, they must be able to discover and interact with each other in ways that may lead to unpredictable interaction patterns depending, for example, on their geographical proximity.

Fundamentally, the *CORTEX* programming model describes the facilities that will be provided to application developers responsible for the construction of *proactive applications* that employ mobile *sentient objects*. At the heart of the *CORTEX* programming model is an anonymous event-based communication model, which we discuss later. Using a non-blocking event-based model, we are able to achieve autonomous sentient behaviour that is independent of the problems associated with traditional blocking communication paradigms (such as RPC [3,15,20,25]). The programming model includes mechanisms for the specification of constraints on the propagation and delivery of events, and the means to express incremental real-time and reliability guarantees, in the form of QoS properties. QoS is taken as a metric of predictability in terms of timeliness and reliability. The model will necessarily make the heterogeneity of the underlying physical system visible at some level of abstraction, for example, as a hierarchy of zones capable of delivering specified levels of QoS.

From the programmer's perspective, the system model is therefore composed of *the environment* and *a set of sentient objects* that interact with it. *CORTEX* adopts the *active environment* metaphor. The environment is thus that part of the system which, whilst being active, is limited to disseminating information about its current state and/or events that take place in the actual physical environment, to objects of the system. The environment may also be acted upon or modified by these same objects. In contrast to the environment, sentient objects are the active, mobile and autonomous entities in the system and are capable of taking decisions, and influencing both the environment and other objects. The programming model supports several different aspects of the behaviour of sentient objects including:

- *acquiring* information from the environment and other objects (the sentience aspect);
- *reacting* to possibly unexpected situations (the autonomy aspect);
- and *modifying* the state of the environment (the control aspect).

Unlike traditional distributed applications in which communication between the objects that comprise an application usually requires that the object that sends a message knows the identity of the object or objects to which the message is to be sent, sentient objects will often need to send messages to a set of

other objects whose identities are not known to the sender and which can only be determined at the time that the message is actually sent. For example, an object may need to send a message to all the objects that are nearby at the time that the message is sent.

While the basic concept of an event-based communication paradigm is simple and indeed hopefully intuitive, there are a number of difficult issues to be tackled if the paradigm is to be employed successfully in large scale proactive applications in which we can expect very large volumes of event announcements to be generated at a very high rate. Of particular importance are techniques to ensure that event announcements are only propagated to objects that have already expressed interest in events of the corresponding type and then only if the specific events are relevant in the object's current location and at the current time.

Filters provide a basic mechanism to allow objects to express interest, or lack thereof, in events of a certain type or containing certain (combinations of) parameter values. Essentially an object subscribing interest in events of a particular type should be able to provide a filter describing which occurrences of events of that type it wants to be notified of. Filters alone are, however, not sufficient. With filters an object may still receive notifications of occurrences of events in a part of the system with which it is not currently concerned.

Zones introduce a means of scoping or limiting the propagation of event notifications in the system. Objects can be organised into zones where a zone can be seen simply as a collection of objects and event notifications are only propagated within the zone of the object raising the event. Objects are organised into zones at the discretion of the application programmer based on functionality, geographical location or physical location on the network.

While filters and zones allow an object to specify, at some level, which event notifications it is interested in, they do not address non-functional requirements related to delivery of such notifications. This will be achieved in the *CORTEX* programming model through the introduction of a generic means of expressing QoS properties encompassing timeliness and reliability, including consistency and ordering of event notifications. Such a mechanism allows application programmers to specify QoS requirements in terms that are meaningful for particular application areas. The application-level QoS parameters can then be mapped onto the system-level QoS parameters characterising the service levels that can be supported by the underlying physical infrastructure at a given point in time.

4. INTERACTION MODEL

Interaction comprises the aspects of communication and co-ordination. An event-based programming model naturally leads to the spontaneous generation of messages rather than a request/response style of communication. This fact suggests the use of an anonymous generative paradigm of communication [8,21,32,35], based on typed communication channels that connect producers and consumers of events according to a publisher/subscriber model. In addition, it is well suited to autonomy, since it does not force an explicit transfer of control, nor synchronization between producers of events and their consumers. The interaction model of *CORTEX* is thus centered on an anonymous generative communication abstraction reflecting the needs of object autonomy and system robustness and evolution. We extend and modify existing approaches in two major directions:

- we consider the fact that sentient objects not only communicate via the network but also, indirectly, through the environment, when they act on it. Thus the environment constitutes an interaction and communication channel and is in the control and awareness loop of the objects.
- we address new adaptable ways of guaranteeing temporal properties of interactions, in the presence of uncertain timeliness of the environment. It is intended to exploit context awareness of sentient objects to reach this goal.

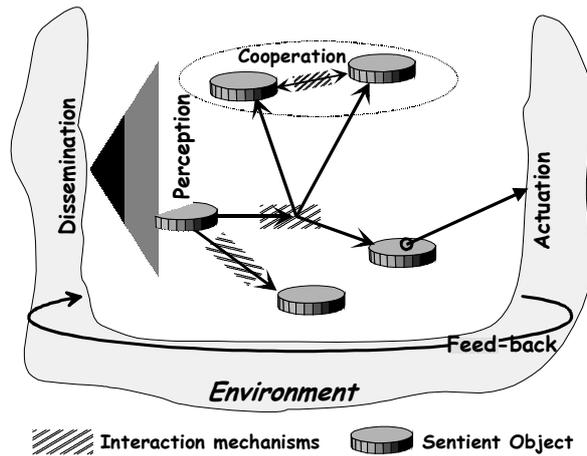


Figure 1: Events and object interactions in CORTEX

CORTEX must support several kinds of interactions (see Figure 1):

- **Environment-to-object interactions** take the form of unsolicited dissemination of the state of the former, and/or notification about events taking place therein. The transformation of events to state within the realm of the active environment components is not precluded, as a way to preserve the memory of past events.
- **Object-to-object interactions** serve two purposes. The first is related with complementing the assessment of each individual object about the state of the surrounding space, which includes environment components and the objects "within reach", that is, capable of influencing its next decisions. The second is related to collaboration, in which the object tries to influence other objects into contributing to a common goal, or into reacting to an unexpected situation.
- **Object-to-environment interactions** comprise the deliberate attempt at forcing a change in the state of the environment. This may come as a consequence of the pursuance of the object's own objectives, or of the reaction to unexpected situations created by the environment or other objects.

Given the highly interactive nature of the envisaged applications, and the fact that actions will be dictated to a great extent by assessment of the state of the environment, *CORTEX* falls under the typical constraints placed on distributed real-time systems [23,24]: some actions must be expressed in terms of timeliness properties, and/or require synchronisation between objects; most of the information provided by the environment consists of the state of real-time or time-value entities, whose value is a function of time, and thus makes any implied actions time-sensitive, even if just concerned with a value. Thus, the communication abstractions must support *predictable timing behaviour*.

CORTEX exploits *context and environmental awareness*, that is, use of internal and external context information to facilitate object interaction in changing situations. Objects, for example whilst moving, may be confronted with unpredictable communication needs or unanticipated interaction patterns with other objects and with the environment. Context awareness in terms of "which network am I in?", "how many hops away are my partners", "what delay am I to expect for this message", as well as the detection of timing failures or the assessment of membership all constitute examples of context awareness. Adaptation can then be attempted, both by the objects themselves, and by the underlying protocols, on behalf of the former. In a similar way, awareness of environmental information may be necessary for certain types of actions that require tight collaboration and synchronization among objects, or a high level of precision in interactions with the environment. As an example, we may consider objects equipped with special sensors that disseminate their perception of the environment to other objects that can now share this information without being capable of sensing this particular aspect themselves.

The main issue introduced by co-operation in the interaction model is the predictability of the *co-ordination mechanisms* necessary to carry out joint actions. Since objects operate in a real world environment, co-ordination has to be achieved under temporal constraints. As a minimum, this requires

timeliness of communication, including those primitives achieving consensus, ordering, and so forth. Additionally, whenever it is necessary to preserve causality relations [26], a temporal model of causality has to be pursued, in order to prevent anomalous behaviour.

In order to address these issues, *CORTEX* combines the group communication and the anonymous communication paradigms in a flexible way, allowing non-functional properties such as the required degree of synchrony and the reliability of communication to be specified on a per group basis.

Achieving predictable timing behaviour is a hard task in large scale, heterogeneous systems that cannot be made strictly synchronous at reasonable costs in transmission delay and bandwidth. However, and despite some of the adverse conditions just described, applications have to exhibit a certain degree of *predictability*. Approaches to this problem under uncertain operating conditions have been addressed in the mission-critical systems arena. Systems would normally have pre-defined operational envelopes, to which they would switch, in a best effort to achieve their goal [44,45]. In more general terms, this is also the track followed by the QoS adaptation body of research. This adaptation is generally done in an ad-hoc manner, and may sometimes not bring the system to an optimal tuning.

In contrast, we address the above-mentioned hard problems in the time domain under the light of *partial synchrony* models, which can withstand varying timeliness or synchrony conditions, and the occurrence of timing failures. Our approach follows recent work on timing and QoS failure detection oracles [41] under partial synchrony models that reason in terms of the $\langle \textit{assumption, coverage} \rangle$ binomial. This may help provide a precise definition of predictability, in terms of an assurance to which a probability is attached, and thus provide conditions for objects to make justifiable tradeoffs between maintaining their original goals with a reduced probability of success, or relaxing their goals whilst maintaining the initial probability.

5. SYSTEM ARCHITECTURE

The architecture of *CORTEX* must recognize two facts: much of the real infrastructure may actually be not know at system deployment time, requiring the capacity for discovery of topology, services and so forth; components are of an extremely heterogeneous nature, both in technological and exploitation terms (wired vs. wireless, public vs. private).

CORTEX features an *abstract network architecture* that reflects the hierarchical structure of large-scale heterogeneous networks, while defining the necessary mappings from this abstract description to real networks, including sensor/actuator busses and wireless links. In the architecture, non-functional properties are translated to QoS requirements, specified at the level of the interaction model abstractions. We define *gateways* as crucial architecture components, which serve as brokers for both the functional and non-functional (e.g. QoS) properties of the subsystems they hide.

The basic infrastructure is composed of a global wide area network (WAN) that comprises substructures subsumed by the abstraction of a Controller Area Network (CAN). The WAN comprises all that makes the globally available, mostly (but not only) wired, mostly public, wide range network infrastructure. A CAN module (CAN is taken to generically mean a small, control or fieldbus type network, of which the CAN [7] standard is only one representative) represents a confined environment in which a certain quality of communication in terms of bandwidth, transmission delays, and reliability can be enforced. The *WAN-of-CAN* structure allows a hierarchical composition of heterogeneous environments with respect to timeliness: at the lowest level we may find networks with highly predictable communication, controlling physical devices such as sensors and actuators.

This WAN-of-CAN structure is assembled by means of *gateways*. A gateway is a crucial architectural construct that provides the propagation of QoS constraints on event flows, and on the events proper, while ensuring timeliness confinement between parts of the architecture, namely in what concerns CAN modules. From an architectural point of view, gateways can be seen as artefacts that provide a representation of a certain environment to the outside world. Therefore, they must provide means to

specify how this representation will be established and how the events will flow from, and to the outside environment.

CORTEX, such as many other complex systems, offers a set of basic support services through a middleware layer. Today, the focus of middleware is on interoperability by providing functionally compatible interfaces. A number of mobile-specific distributed systems services have been developed in recent years that aim to operate in challenging mobile environments [10,14]. However, to date such services are not designed to offer sufficient levels of dependability or support the highly asynchronous interaction model required by the *CORTEX* computational paradigm. In more detail, when targeting mission and safety critical applications, e.g. traffic management systems, predictability under widely varying load and fault conditions becomes an additional decisive requirement. It is the conflict between technical conditions and the application requirements that makes predictability one of the greatest challenges for the middleware of the future.

The ability to enforce and check timeliness of actions with given coverage assumptions is a must to achieve dependable execution in face of uncertain timeliness. This requires the availability of a number of basic services, such as timing failure detection and clock synchronization. Similarly, discovery services are mandatory, not only in terms of topology, but also in terms of services offered by the infrastructure. Together, these are the distinctive services supplied by the *CORTEX* middleware.

6. APPLICATION SCENARIOS

In this section, we illustrate with a few scenarios the relevance of the *CORTEX* architecture.

Supporting field workers in the electricity industry:

Field workers in the electricity supply industry work in a highly distributed safety critical and real-time environment. Current best working practice is based on a centralised co-ordination body (the control centre). However, such centralisation inevitably proves to be a bottleneck and potential point of failure during periods of high activity, such as during lightning storms. Furthermore, the highly mobile field engineers are often unable to establish contact with the control centre in a timely fashion, although they may well be able to establish ad-hoc dialogues with neighbouring colleagues. By enhancing collaboration between colleagues and replicating the view of the current network state to all field engineers, there is the potential to distribute operational control and co-ordination [12].

One example of the application of *CORTEX* would be to put intelligence and monitoring capabilities into the power distribution network infrastructure itself (e.g. at substation switches). This would allow predictable action by providing different levels of dependability in isolated parts of the network. A switch might, for example, take autonomous action to ensure fail-safe behaviour under certain conditions. Alternatively, these ‘sentient switches’ may take a proactive role in collaborating with field engineers directly. Such collaboration will, we believe, facilitate the establishment of pockets or zones of co-ordination, enabling useful work to be performed, despite the inability to achieve direct communication with a centralised control centre.

Mountain rescue:

Mountain rescue workers are constantly faced with search and rescue operations in which they are called upon to locate stranded, and possibly injured, people in extremely hostile conditions. Such environments also place stringent constraints on mobile computational devices, such as weight, battery life and communications availability. To affect a successful search of a mountain-side requires a co-ordinated effort by a team of rescuers who are themselves vulnerable to hostile weather conditions, can become separated and even injured. In such scenarios, tracking the location of search team members is required in order to provide both the control centre and rescuers with an awareness of the location of those team members involved in the rescue.

Sharing of location information poses a technical challenge; the availability of the communications infrastructure and partitioning of the search team is greatly affected by the topology of the mountain terrain and unpredictability of the prevailing weather conditions.

However, clearly ad-hoc networking can enable collaborations between neighbouring team members and promote the sharing of information, such as location and medical telemetry, to enhance the effectiveness of a typical search and rescue operation. Moreover, information gathered in the field can be related back to remote experts at the base or local accident and emergency departments.

The dynamic nature of such collaborations is poorly supported by existing distributed systems given the general bias towards a reliable and fixed communications infrastructure. The *CORTEX* paradigm fits well to this application domain by supporting the notion of zones that, in this scenario, could represent zones of network availability.

Next generation cars:

As a more futuristic example, consider the reaction of a queue of cars to an accident on a typically busy motorway carriageway. A driver in the queue will be forced to brake suddenly, but when the next driver reacts, she will brake hard enough to stop her car within the remaining braking distance. Each car in the queue reacts similarly. The usual outcome of this behaviour is that a number of drivers will not have sufficient braking distance left and a multiple car collision will occur.

Using the *CORTEX* paradigm, a more co-ordinated approach could be achieved in which vehicles publish events (such as the fact that they are performing an emergency brake or that the car will be stationary in x ms) to other interested parties (e.g. cars following within a certain distance). Sentient objects (located within other cars in the queue) can ask to receive the braking event, and when notified can take appropriate braking action (publishing their own braking events). In this way, the entire queue of vehicles can be brought to a halt in a progressive and controlled manner.

7. CONCLUSIONS

In this paper we addressed the problem of how to handle the requirements of an emerging class of applications that operate independently of direct human control, using the necessary support in terms of system architecture and middleware models. We discussed the fundamental challenges in supporting this kind of applications, which derive from several key characteristics of the latter, including sentience, autonomy, large scale, geographical dispersion, mobility and evolution.

The first step to address the problem consists in defining a programming model that provides the necessary means for application developers to construct *proactive applications* that employ mobile *sentient objects*. We described the fundamental aspects of the *CORTEX* programming model, proposing an anonymous event-based communication model and the possibility of specifying constraints and guarantees in the form of QoS properties. We also introduced a few concepts important for the programming model, like *filters* and *zones*.

The interaction among the objects in the system, which comprise the aspects of communication and co-ordination, are dealt within the proposed interaction model. We also propose the fundamentals of a system architecture, which defines the components that are necessary to implement the communication abstractions identified in the interaction model. A key aspect of this architecture is that it must be capable to handle the extremely heterogeneous nature of the components it is composed of.

The *CORTEX* project already took the first steps towards the definition of the fundamental paradigms and solutions needed to address this class of sentient and proactive applications. While we have not yet closed the doors to new ideas and improvements, we are on our way to construct several proof-of-concept prototypes to illustrate the feasibility of our ideas. We have been using several application scenarios to drive our work, of which the ones presented in this paper are just a fraction. We expect to publish further results of our ongoing work in a near future.

BIBLIOGRAPHY

- [1] N. Ackroyd and R. Lorimer. Global Navigation: a GPS User's Guide. *Lloyd's of London*, 1994, 2nd Ed.
- [2] R. Azuma. Tracking Requirements for Augmented Reality. *Communications of the ACM*, 36(7), pp.50-51, Jul 1993.
- [3] Bakre, A., and B.R. Badrinath. M-RPC: A Remote Procedure Call Service for Mobile Clients, Technical Report WINLAB TR-98, Department of Computer Science, Rutgers University, U.S. June 1995.
- [4] Scott Brandt, Gary Nutt, Toby Berk and James Mankovich. A Dynamic Quality of Service Middleware Agent for Mediating Application Resource Usage. *Proceedings of the 19th IEEE Real-Time Systems Symposium*. pp.307-317. Madrid, Spain. Dec 1998.
- [5] I. Busse, B. Deffner and H. Schulzrinne. Dynamic QoS Control of Multimedia Applications based on RTP. *Computer Communications*,. 19(1), Jan 1996.
- [6] A. Campbell and G. Coulson. A QoS adaptive transport system: Design, implementation and experience. *Proceedings of the Fourth ACM Multimedia Conference (MULTIMEDIA'96)*. pp.117-128. New York, NY, USA. Nov 1996.
- [7] CAN - Controller Area Network for high-speed communication. Int'l Std.11898- Road vehicles - Interchange of digital information . ISO, 1993.
- [8] N. Carriero, D. Gelernter. Linda in Context, *Communications of the ACM*, 32, 4, April 1989, pp 444-458.
- [9] T. Chandra and S. Toueg. Unreliable Failure Detectors for Reliable Distributed Systems. *Journal of the ACM*, 43(2):225-267, March 1996.
- [10] Keith Cheverst. Development of a Group Service to Support Collaborative Mobile Groupware, Ph.D. Thesis, Computing Department, Lancaster University, Bailrigg, Lancaster, LA1 4YR, U.K., April 1999.
- [11] F. Cristian and C. Fetzer. The Timed Asynchronous System Model. *Proceedings of the 28th Annual International Symposium on Fault-Tolerant Computing*. pp.140-149. Munich, Germany. June 1998.
- [12] Davies, N., Friday, A., Blair, G.S., and Cheverst, K. Distributed Systems Support for Adaptive Mobile Applications. *ACM Mobile Networks and Applications*, special issue *Mobile Computing – System Services*, 4(5), 1996.
- [13] G.W. Fitzmaurice. Situated Information Spaces and Spatially Aware Palmtop Computers. *Communications of the ACM*, 36(7), pp.38-49, Jul 1993.
- [14] Adrian Friday. Infrastructure Support for Adaptive Mobile Applications, Ph.D. Thesis, Computing Department, Lancaster University, Bailrigg, Lancaster, LA1 4YR, U.K., September 1996.
- [15] Mads Haahr, Raymond Cunningham and Vinny Cahill. Supporting CORBA Applications in a Mobile Environment, in *Proceedings of the 5th ACM/IEEE International Conference on Mobile Computing and Networking (MobiCom)*, 1999.
- [16] A. Harter and A. Hopper. A Distributed Location System for the Active Office. *IEEE Network*, 8(1), 1994.
- [17] M. Horstmann and M. Kirtland. DCOM Architecture. <http://www.microsoft.com/jini/specs/>
- [18] H. Ihara and K. Mori. Autonomous Decentralized Computer Control Systems. *IEEE Computer*, 17(8):57-66, August 1984.
- [19] JINI Technology 1.1. Specification, Sun Microsystems, <http://www.sun.com/jini/specs/>
- [20] Joseph, A., A. deLepinasse, J. Tauber, D. Gifford, and M.F. Kaashoek. Rover: A Toolkit for Mobile Information Access, *Proc. 15th ACM Symposium on Operating System Principles (SOSP)*, Copper Mountain Resort, Colorado, U.S., 3-6 December 1995. ACM Press, Vol. 29, Pages 156-171.
- [21] J. Kaiser, M. Mock. Implementing the Real-Time Publisher/Subscriber Model on the Controller Area Network (CAN), *Proc. of the 2nd Int. Symp. on Object-oriented Real-time distributed Computing (ISORC99)*, Saint-Malo, France, May 1999.
- [22] H. Kopetz and G. Grünsteidl. TTP - A Time-Triggered Protocol for Fault-Tolerant Real-Time Systems, Research Report 12/92, Inst. f. Techn. Informatik, Tech. Univ. of Vienna, 1992.
- [23] P. Verissimo and Luís Rodrigues. Distributed Systems for System Architects, Kluwer Academic Publishers, 2001.

- [24] H. Kopetz. Real-Time Systems, Design Principles for Distributed Embedded Applications, Kluwer Academic Publishers, 1997.
- [25] Kümmel, S., A. Schill, and G. Volkmann. RPC over Advanced Network Technologies: Evaluation and Experiences. *Proc. of the 3rd International Workshop on Services in Distributed Networked Environments (SDNE)*, Macau, China, 3-4 June 1996. IEEE Computer Society Press.
- [26] L. Lamport. Time, Clocks and the Ordering of Events in a Distributed System, Z. Yand and T. Marsland (Eds.), *Global States and Time in Distributed systems*, IEEE Computer Society Press, 1994.
- [27] M.A. Livani, J. Kaiser, W.J. Jia. Scheduling Hard and Soft Real-Time Communication in the Controller Area Network (CAN), *23rd IFAC/IFIP Workshop on Real Time Programming*, Shantou, China, June 1998.
- [28] S. Maffei. iBus - The Java Intranet Software Bus, Olsen&Associates, www.olsen.ch, 1997.
- [29] K. Mori. Autonomous decentralized Systems: Concepts, Data Field Architectures, and Future Trends, *Int. Conference on Autonomous Decentralized Systems (ISADS93)*, 1993.
- [30] Object Management Group. The Common Object Request Broker: Architecture and Specification. *OMG Document 96-03-04*, July 1995.
- [31] Karl O'Connell, Tom Dinneen, Stephen Collins, Brendan Tangney, Neville Harris and Vinny Cahill. Techniques for Handling Scale and Distribution in Virtual Worlds. *Proceedings of the 7th ACM SIGOPS European Workshop*. pp.17-24. Connemara, Ireland, Sep 1996.
- [32] B. Oki, M. Pfluegl, A. Seigel, D. Skeen. The information Bus®- An Architecture for Extensible Distributed Systems, *14th ACM Symposium on Operating System Principles*, Asheville, NC, Dec 1993, pp.58-68.
- [33] Overall Concepts and Principles of TINA. TINA Baseline, TB_MDC.018_1.0_94, February 1995.
- [34] D. Powell. Failure Mode Assumptions and Assumption Coverage. *Digest of Papers, The 22nd International Symposium on Fault-Tolerant Computing Systems*. pp.386-395. Boston, USA. July 1992.
- [35] R. Rajkumar, M. Gagliardi, L. Sha. The Real-Time Publisher/Subscribe Inter-Process Communication Model for Distributed Real-Time Systems: Design and Implementation, *IEEE Real-time Technology and Applications Symposium*, June 1995.
- [36] J. Rufino, P. Verissimo, C. Almeida, L. Rodrigues. Fault-Tolerant Broadcasts in CAN, *Digest of Papers, The 28th International Symposium on Fault-Tolerant Computing Systems*, Munich, Germany, June 1998.
- [37] B. Sabata, S. Chatterjee, M. Davis, J. Sydir, T. Lawrence. Taxonomy of QoS Specification, *Proc. of the IEEE Third International Workshop on Object-Oriented Real-Time Dependable Systems (WORDS)*, Newport beach, CA, February 1997
- [38] Gradimir Starovic, Vinny Cahill and Brendan Tangney. An Event Based Object Model for Distributed Programming. *OOIS (Object-Oriented Information Systems) '95*. pp.72-86. Dec 1995.
- [39] K. Tindell, A. Burns. Guaranteed Message latencies for Distributed Safety-Critical Hard Real Time Control Networks, Technical Report YCS229, Dept. of Comp. Science, University of York, May 1994.
- [40] P. Verissimo and C. Almeida. Quasi-synchronism: a step away from the traditional fault-tolerant real-time system models, *Bulletin of the Technical Committee on Operating Systems and Application Environments (TCOS)*, 7(4):35-39, Winter 1995.
- [41] P. Verissimo, A. Casimiro and C. Fetzer. The Timely Computing Base: Timely actions in the presence of uncertain timeliness. In *Proceedings of the International Conference on Dependable Systems and Networks*, pages 533-542, New York City, USA, June 2000. IEEE Computer Society Press.
- [42] A. Ward, A. Jones and A. Hopper. A New Location Technique for the Active Office. *IEEE Personal Communications*, 4(5), 1997.
- [43] K.M. Zuberi and K. G. Shin. Non-Preemptive Scheduling of messages on Controller Area Network for Real-Time Control Applications, Technical Report, University of Michigan, 1995.
- [44] E. Jensen and J. Northcutt. Alpha: A non-proprietary os for large, complex, distributed real-time. In *Procs. of the IEEE Workshop on Experimental Distributed Systems*, pages 35-41, 1990, Alabama, USA.
- [45] Verissimo, P. and Barrett, P. and Bond, P. and Hilborne, A. and Rodrigues, L. and Seaton, D. The Extra Performance Architecture (XPA). In *Delta-4 - A Generic Architecture for Dependable Distributed Computing*, D. Powell ed., pages 211-266, Springer Verlag, ESPRIT Research Reports Series, 1991.