

Task-based Adaptation for Ubiquitous Computing

Vahe Poladian, João Pedro Sousa, David Garlan, Bradley Schmerl, Mary Shaw
 School of Computer Science, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh PA 15213
 Phone: +1 412 268 5056 Fax: +1 412 268 3455
 poladian|jpsousa|garlan|schmerl|mary_shaw@cs.cmu.edu

Abstract— An important domain for autonomic systems is the area of ubiquitous computing: users are increasingly surrounded by technology that is heterogeneous, pervasive, and variable. In this paper we describe our work in developing self-adapting computing infrastructure that automates the configuration and re-configuration of such environments. Focusing on the engineering issues of self-adaptation in the presence of heterogeneous platforms, legacy applications, mobile users, and resource variable environments, we describe a new approach based on the following key ideas: (a) Explicit representation of user tasks allows us to determine what service qualities are required of a given configuration; (b) Decoupling task and preference specification from the lower level mechanisms that carry out those preferences provides a clean engineering separation of concerns between what is needed and how it is carried out; and (c) Efficient algorithms allow us to calculate in real time near-optimal resource allocations and reallocations for a given task.

Index Terms— Self-adaptation, ubiquitous computing, resource-aware computing, multi-fidelity applications.

I. INTRODUCTION

Self-adaptive systems are becoming increasingly important. What was once the concern of specialized systems, with high availability requirements, is now recognized as being relevant to almost all of today's complex systems [9,18]. Increasingly computing systems that people depend on cannot be taken off-line for repair – they must adapt to failures in environments that are not entirely under the control of the system implementers, and they must adjust their run-time characteristics to accommodate changing loads, resources, and goals.

One particularly important domain for self-adaptation is the area of ubiquitous computing. Today users are surrounded by technology that is heterogeneous, pervasive, and variable. It is heterogeneous because computation can take place using a wide variety of computing platforms, interfaces, networks, and services. It is pervasive through wireless and wired connectivity that pervades most of our working and living environments. It is variable because resources are subject to change:

users can move from resource-rich settings (such as workstations and high-bandwidth networks in an office) to resource-poor environments (such as a PDA in a park).

Coping with this situation requires automated mechanisms. In particular, ideally systems should be able to adapt to user

mobility, recover from service failures and degradations, and allow continuity across diverse environments. Without automated mechanisms to support this kind of adaptation, users become increasingly overloaded with distractions of managing their system configurations; alternatively, they may simply opt not to use the capabilities of their environments.

This automation raises a number of serious engineering challenges: How can one determine when reconfiguration is appropriate? Assuming reconfiguration is desirable, how does one determine a satisfactory allocation of resources, particularly if there are multiple ways to support a given computing task, or limitations on the resource pool? How can users instruct the system about the kinds of adaptation that are desired, without becoming bogged down in low-level system details? How can one add adaptation mechanisms to the everyday computing environments that users are familiar with: text editors, spreadsheets, video viewers, browsers, etc.

In this paper we describe our experience over the past five years of developing self-adapting ubiquitous computing infrastructure that automates the configuration and reconfiguration of everyday computing environments. Focusing on the engineering issues of providing self-adaptation in the presence of heterogeneous platforms, legacy applications, mobile users, and resource variable environments, Project Aura [8] has developed an approach that we believe addresses each of the questions above. The key ideas behind this work are the following: (a) Explicit representation of user tasks allows us to determine what service qualities are required of a given configuration; (b) Decoupling task and preference specification from the lower level mechanisms that carry out those preferences provides a clean engineering separation of concerns between what is needed and how it is carried out; and (c) Efficient algorithms allow us to calculate in real time near-optimal resource allocations and reallocations for a given task.

In the remainder of this paper Section II describes our work in the context of related research. Section III outlines the research challenges in making self-adaptive systems task-aware, describes the Aura architecture, and illustrates how the architecture addresses such research challenges. Section IV elaborates on the specifics of such an architecture: how the requirements and user preferences are captured for each task; the formal underpinnings of the internal representation of tasks and preferences; and the algorithm that supports automatic system configuration and self-adaptation. Section V evaluates the effectiveness of the approach and consider les-

sons learned from this work.

II. RELATED WORK

Currently, adaptive systems fall into two broad categories: fault-tolerant systems, and fidelity-aware systems. First, fault-tolerant systems react to component failure, compensating for errors using a variety of techniques such as redundancy and graceful degradation [5,12]. Such systems have been prevalent in safety-critical systems or systems for which the cost of off-line repair is prohibitive (e.g., telecom, space systems, power control systems, etc.) Here the primary goal is to prevent or delay large-scale system failure.

Second, fidelity-aware systems react to resource variation: components adapt their computing strategies so they can function optimally with the current set of resources (bandwidth, memory, CPU, power, etc.) [7,15,20,22]. Many of these systems emerged with the advent of mobile computing over wireless networks, where resource variability becomes a critical concern. While most of this research focuses on one component at a time, our work leverages on this research but tackles the problem of multi-component integration, configuration, and reconfiguration. Although somewhat related, this kind of automatic configuration is distinct from the automatic configuration being investigated in other research [19]. There, configuration is taken in the sense of *building and installing* new applications into an environment, whereas here, it is taken in the sense of *selecting and controlling* applications so that the user can go about his tasks with minimal disruption.

Resource scheduling [14], resource allocation [16,21], and admission control have been extensively addressed in research. From analytical point of view, closest to our work are Q-RAM [16], a resource reservation and admission control system maximizing the utility of a multimedia server based on preferences of simultaneously connected clients; Knapsack algorithms [23]; and winner determination in combinatorial auctions. In our work, we handle the additional problems of selecting applications among alternatives, and accounting for cost of change. Dynamic resolution of resource allocation policy conflicts involving multiple mobile users is addressed in [2] using sealed bid auctions. While our work shares utility-theoretic concepts with [2], the problem solved in our work is different. In that work, the objective is to select among a handful of policies so as to maximize an objective function of multiple users. In our work, the objective is to choose among possibly thousands of configurations so as to maximize the objective function of one user. As such, our work has no game-theoretic aspects, but faces a harder computational problem. Furthermore, our work takes into account tasks that users wish to perform.

At a coarser grain, research in distributed systems addresses global adaptation: for example, a system might reconfigure a set of clients and servers to achieve optimal load balancing. Typically, such systems use global system models, such as architectural models, to achieve these results [4,10,11]. To achieve fault-tolerance and coarse-grain adaptation (e.g. hot

component swapping,) our work builds on this, as well as on service location and discovery protocols [13,24].

III. TASK-BASED SELF-ADAPTATION

A. Task-aware Systems

A central tenet of our work is that systems are used to carry out high-level activities of users: planning a trip, buying a car; communicating with others. In today's systems those activities and goals are implicit. Users must map their tasks to computing systems by invoking specific applications (document editors, email programs, spreadsheets, etc.) on specific files, with knowledge of specific resources. In a ubiquitous computing world with shifting resources and increased heterogeneity, the cognitive load required for users to manage this manually quickly becomes untenable.

In contrast, a task-aware system makes user tasks explicit, by encoding user goals, and by providing a placeholder to represent the quality attributes of the services used to perform those tasks. So, for example, for a particular task, in the presence of limited bandwidth, the user may be willing to live with a small video screen size, while in another task reducing the frame rate would be preferable. In task-aware systems, users specify their tasks and goals, and it is the job of the *system* to automatically map them into the capabilities available in the ubiquitous environment.

Once such information is represented, a self-managing system can in principle query the task to determine both when the system is behaving within an acceptable envelop for the task, and also can choose among alternative system reconfigurations when it is not.

However, a number of important research questions arise, and the way we answer them strongly influences the way we look at and build task-aware systems:

- How do we represent a task? What encoding schemes can best be used to capture the user's requirements for system quality?
- How should we characterize the knowledge for mapping a user task to a system's configuration? As a user moves from task to task, different configurations will be appropriate, even for the same set of applications.
- Should we trigger an adaptation as soon as an opportunity for improvement is detected, or should we factor in how distracting the change will be to the user against how serious the fault is?
- Is the binary notion of fault enough, or do we need to come up with a measure of fault "hardness" – a continuum between "all is well," and "the system is down?"
- What is the length of time that the user is expected to carry out the current task? What are likely other tasks that the user will work on next?

Over the past five years we have been experimenting with various answers to these questions. Centered on a large ubiquitous computing research project, Project Aura [8], we have evolved a system that, in brief, addresses these questions as follows:

- We represent *a task as a set of services*, together with a set of quality attribute preferences expressed as multi-dimensional utility functions, possibly conditioned by context conditions.
- We define *a vocabulary for expressing requirements*, which delimits the space of requirements that the automatic reconfiguration can cover. The set of requirements for a particular task expresses which services are needed from the system, as well as the fidelity constraints that make the system adequate or inadequate for the task at hand. The required services are dynamically mapped to the available components and the fidelity constraints are mapped into resource-adaptation policies.
- We incorporate the notion of *cost of reconfiguration* into the evaluation of alternative reconfigurations. This cost captures user's intolerance for configuration changes by the infrastructure. A high cost of reconfiguration will make the system highly stable, but frequently less optimal; a low cost of configuration will permit the system to change frequently, but may introduce more user distraction from reconfigurations.
- We invert the notion of fault by adopting an econometric-based notion of *task feasibility*: ranging from 0 (the task is not feasible under the current system conditions) to 1 (system is totally appropriate for the current task). This enables an objective evaluation of configuration alternatives, regardless of the sources of change (both changes to the task, and also to the availability of resources and components).

We now describe the system architecture that permits such task-based self-adaptation, and elaborate on these decisions.

B. The aura layers

The starting point for understanding Aura is a layered view of its infrastructure together with an explanation of the roles of each layer with respect to task suspend-resume and dynamic adaptation. Table 1 summarizes the relevant terminology.

TABLE 1. TERMINOLOGY.

<i>task</i>	An everyday activity such as preparing a presentation or writing a report. Carrying out a task may require obtaining several <i>services</i> from an <i>environment</i> , as well as accessing several <i>materials</i> .
<i>environment</i>	The set of <i>suppliers</i> , <i>materials</i> and <i>resources</i> accessible to a user at a particular location.
<i>service</i>	Either (a) a service type, such as printing, or (b) the occurrence of a service proper, such as printing a given document. For simplicity, we will let these meanings be inferred from context.
<i>supplier</i>	An application or device offering <i>services</i> – e.g. a printer.
<i>material</i>	An information asset such as a file or data stream.
<i>capabilities</i>	The set of <i>services</i> offered by a <i>supplier</i> , or by a whole <i>environment</i> .
<i>resources</i>	Are consumed by <i>suppliers</i> while providing <i>services</i> . Examples are: CPU cycles, memory, battery, bandwidth, etc.
<i>context</i>	Set of human-perceived attributes such as physical location, physical activity (sitting, walking...), or social activity (alone, giving a talk...).
<i>user-level state of a task</i>	User-observable set of properties in the <i>environment</i> that characterize the support for the task. Specifically, the set of

services supporting the task, the user-level settings (preferences, options) associated with each of those services, the *materials* being worked on, user-interaction parameters (window size, cursors...), and the user's preferences with respect to quality of service tradeoffs.

<i>layer</i>	<i>mission</i>	<i>roles</i>
Task Management	<i>what</i> does the user need	<ul style="list-style-type: none"> • monitor the user's task, context and preferences • map the user's task to needs for services in the environment • complex tasks: decomposition, plans, context dependencies
Environment Management	<i>how</i> to best configure the environment	<ul style="list-style-type: none"> • monitor environment capabilities and resources • map service needs, and user-level state of tasks to available suppliers • ongoing optimization of the utility of the environment relative to the user's task
Env.	support the user's task	<ul style="list-style-type: none"> • monitor relevant resources • fine grain management of QoS/resource tradeoffs

TABLE 2. SUMMARY OF THE SOFTWARE LAYERS IN THE INFRASTRUCTURE.

The infrastructure exploits knowledge about a user's tasks to automatically configure and reconfigure the environment on behalf of the user. First, the infrastructure needs to know *what* to configure for; that is, what a user needs from the environment in order to carry out his or her tasks. Second, the infrastructure needs to know *how* to best configure the environment: it needs mechanisms to optimally match the user's needs to the capabilities and resources in the environment.

In our architecture, each of these two subproblems is addressed by a distinct software layer: (1) the **Task Management** layer determines *what* the user needs from the environment at a specific time and location; and (2) the **Environment Management** layer determines *how* to best configure the environment to support the user's needs.

Table 2 summarizes the roles of the software layers in the infrastructure. The top layer, *Task Management* (TM), captures knowledge about user tasks and associated intent. Such knowledge is used to coordinate the configuration of the environment upon changes in the user's task or context. For instance, when the user attempts to carry out a task in a new environment, TM coordinates access to all the information related to the user's task, and negotiates task support with Environment Management (EM). Task Management also monitors explicit indications from the user and events in the physical context surrounding the user. Upon getting indication that the user intends to suspend the current task or resume some other task, TM coordinates saving the user-level state of the suspended task and reinstatiates the resumed task, as appropriate. Task Management may also capture complex representations of user tasks (out of scope of this paper) including task decomposition (e.g., task A is composed of subtasks B and C), plans (e.g., C should be carried out after B), and context dependencies (e.g., the user can do B while sitting or walking, but not while driving).

The *EM layer* maintains abstract models of the environment. These models provide a level of indirection between the user’s needs, expressed in environment-independent terms, and the concrete capabilities of each environment.

This indirection is used to address both heterogeneity and dynamic change in the environments. With respect to heterogeneity, when the user needs a service, such as speech recognition, EM will find and configure a “supplier” for that service among those available in the environment. With respect to dynamic change, the existence of explicit models of the capabilities in the environment enables automatic reasoning when those capabilities change dynamically. The Environment Management adjusts such a mapping automatically in response to changes in the user’s needs (adaptation initiated by TM), and changes in the environment’s capabilities and resources (adaptation initiated by EM). In both cases adaptation is guided by the maximization of a *utility function* representing the user’s preferences (see Section 0).

The *Environment layer* comprises the applications and devices that can be configured to support a user’s task. Configuration issues aside, these suppliers interact with the user exactly as they would without the presence of the infrastructure. The infrastructure steps in only to automatically configure those suppliers on behalf of the user. The specific capabilities of each supplier are manipulated by EM, which acts as a translator for the environment-independent descriptions of user needs issued by TM.

By factoring models of user preferences and context out of individual applications, the infrastructure enables applications to apply the adaptation policies appropriate for each task. That knowledge is very hard to obtain at the application level, but once it is determined at the user level – by Task Management – it can easily be communicated to the applications supporting the user’s task.

Each layer reacts to changes in user tasks and in the environment at a different granularity and time-scale. Task Management acts at a human perceived time-scale (minutes), evaluating the adequacy of sets of services to support the user’s task. Environment Management acts at a time-scale of seconds, evaluating the adequacy of the mapping between the requested services and specific suppliers. Adaptive applications (fidelity-aware and context-aware) choose appropriate computation tactics at a time-scale of milliseconds. A detailed description of the architecture, including the formal specification of the interactions between the components in the layers defined above, is available in [25].

C. Examples of Self-Adaptation

To clarify how this design works, we illustrate how the infrastructure outlined in Section B handles a variety of examples of self-adaptation, ranging from traditional repair in reaction to faults, to reactions to positive changes in the environment, to reactions to changes in the user’s task.

To set the stage, suppose that Fred is engaged in a conversation that requires real-time speech-to-speech translation. For that task, assume the Aura infrastructure has assembled three

services: speech recognition, language translation, and speech synthesis. Initially both speech recognition and synthesis are running on Fred’s handheld. To save resources on Fred’s handheld, and since language translation is computationally intensive, but has very low demand on data-flow (the text representation of each utterance), the translation service is configured to run on a remote server.

Fault tolerance. Suppose now that there is loss of connectivity to the remote server, or equivalently, that there is a software crash that renders it unavailable. Live monitoring at the EM level detects that the supplier for language translation is lost. The EM looks for an alternative supplier for that service, e.g., translation software on Fred’s handheld, activates it, and automatically reconfigures the service assembly.

Resource/fidelity-awareness. Computational resources in Fred’s handheld are allocated by the EM among the services supporting Fred’s task. For computing optimal resource allocation, the EM uses each supplier’s spec sheet (relating fidelity levels with resource consumption), live monitoring of the available resources, and the user’s preferences with respect to fidelity levels. Suppose that during the social part of the conversation, Fred is fine with a less accurate translation, but response times should be snappy. The speech recognizer, as the main driver of the overall response time, gets proportionally more resources and uses faster, if less accurate, recognition algorithms. When the translation service is activated on Fred’s handheld in response to the fault mentioned above, resources become scarcer for the three services. However, having the knowledge about Fred’s preferences passed upon service activation, each supplier can react appropriately by shifting to computation strategies that save response times at the expense of accuracy [1].

Soft fault (negative delta). Each supplier issues periodic reports on the Quality of Service (QoS) actually being provided – in this example, response time and estimated accuracy of recognition/translation. Suppose that at some point during the conversation, Fred brings up his calendar to check his availability for a meeting. The suppliers for the speech-to-speech translation task, already stretched for resources, reduce their QoS below what Fred’s preferences state as acceptable. The EM detects this soft fault, and replaces the speech recognizer by a lightweight component, that although unable to provide as high a QoS as the full-fledged version when resources are plentiful, performs better under sub-optimal resource availability.¹

Soft fault (positive delta). Suppose that at some point, the language translation supplier running on the remote server becomes available again. The EM detects the availability of a new candidate to supply a service required by Fred’s task, and compares the estimated utility of the candidate solution against the current one. If there is a clear benefit, the EM automatically reconfigures the service assembly. In calculating the benefit, the EM factors in a cost of change, which is

¹ Additionally, the EM uses these periodic QoS reports to monitor the availability of the suppliers, in a heartbeat fashion.

also specified in the user’s preferences associated with each service. This mechanism introduces hysteresis in the reconfiguration behavior, thus avoiding oscillation between closely competing solutions.

Task QoS requirements change. Suppose that at some point Fred’s conversation enters a technical core for which translation accuracy becomes more important than fast response times. The TM provides the mechanisms, if not to recognize the change automatically based on Fred’s social context, at least to allow Fred to quickly indicate his new preferences; for instance, by choosing among a set of preference templates. The new preferences are distributed by the TM to the EM and all the suppliers supporting Fred’s task. Given a new set of constraints, the EM evaluates the current solution against other candidates, reconfigures, if necessary, and determines the new optimal resource allocation. The suppliers that remain in the configuration, upon receiving the new preferences, change their computation strategies dynamically; e.g., by changing to algorithms that offer better accuracy at the expense of response time.

Task suspend/resume. Suppose that after the conversation, Fred wants to resume writing one of his research papers. Again, the TM provides the mechanisms to detect, or for Fred to quickly indicate, his change of task. Once the TM is aware that the conversation is over it coordinates with the suppliers for capturing the user-level state of the current task, if any, and with the EM to deactivate (and release the resources for) the current suppliers. The TM then analyses the description it saved the last time Fred worked on writing the paper, recognizes which services Fred was using and requests those from the EM. After the EM identifies the optimal supplier assignment, the TM interacts with those suppliers to automatically recover the user-level state where Fred left off. See [25] for a formal specification of such interactions.

Task service requirements change. Suppose that while writing his paper, Fred recognizes that it would be helpful to refer to a presentation he gave recently to his research group. The TM enables Fred to explicitly aggregate viewing the presentation to the ongoing task. As soon as a new service is recognized as part of the task, the TM requests an incremental update to the EM, which computes the optimal supplier and resource assignment for the new task definition, and automatically performs the required reconfigurations. Similarly, if Fred decides some service is no longer necessary for his task, he can let the TM know, and the corresponding (incremental) deactivations are propagated to the EM and suppliers. By keeping the TM up-to-date with respect to the requirements of his tasks, Fred benefits from both the automatic incremental reconfiguration of the environment, and from the ability to suspend/resume exactly the set of services that he considers relevant for each task.

IV. SUPPORTING TASK-BASED ADAPTATION

A. Defining Task Requirements

The user expresses the requirements for a task by specifying

the services needed and the associated preferences. A shared vocabulary of services and service-specific quality dimensions must exist between the user and the system. Developing such a vocabulary is a subject of related research and out of the scope of this paper (see for instance [6]), but we give insights to the essential characteristics of such a vocabulary [25,26].

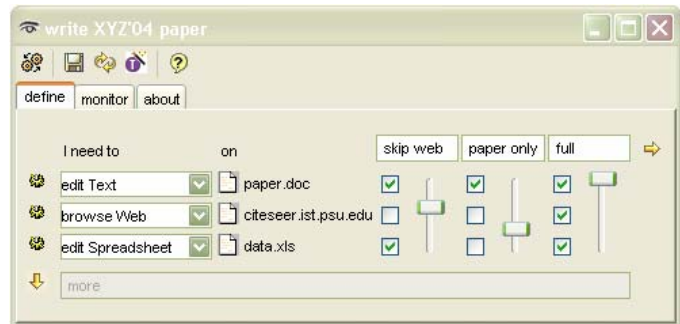


Fig. 1. Fred’s task definition for writing XYZ’04 paper

For instance, to address user mobility across different machines, we use terms that are generic enough to be meaningful on different platforms. For example, a task may capture the fact that the user needs to *edit text*, as opposed to capturing the fact that he needs to use Microsoft Word.

To make these ideas concrete, suppose that Fred is about to start writing a new paper. Fred starts by pressing the down arrow at the bottom of an empty task definition window and selecting *edit text* (Fig. 1). The text editor activated by the infrastructure brings up a (default) blank document and Fred starts working. As Fred browses the web, he decides to associate an especially relevant page with the task, so that it is brought up automatically every time the task is resumed. For that, Fred simply drags the page shortcut out of the browser and into the *more* field of the task window (the default *browse web* appears automatically). Later, Fred decides to start analyzing the performance data on a spreadsheet. Again, Fred simply drags the file produced by the data gathering tool, from the file system explorer into the *more* field and selects *edit spreadsheet* for it.

Note that the infrastructure imposes no constraints on the user’s work. This comes from recognizing that many user activities are spontaneous and short lived, and need not be classified as pertaining to a particular task. However, once the user recognizes an enduring association with a task, the infrastructure makes it easy to update the task definition on the fly.

In addition to specifying the services required by each task, the user may specify preferences with respect to how the environment should be configured. User preferences (and their formal representation, *utility functions*) used in our work have three parts: first, *configuration preferences* capture preferences with respect to the set of services to support a task. Second, *supplier preferences* capture which specific suppliers are preferred to provide the required services; and third, *QoS preferences* capture the acceptable Quality of Service (QoS) levels and preferred tradeoffs.

The right-hand side of Fig. 1 defines Fred’s configuration preferences for the task: that is, alternative operation modes and their order of preference. The (default) *full* configuration includes all the activities defined for the task. In addition, Fred also specifies the *skip web* degraded-mode configuration for when the circumstances are such that either a browser or connection are not available, or that the quality of service is so poor (for instance, due to low bandwidth) that Fred would rather focus on the other activities. Fred also permits the *paper only* configuration for last resort circumstances, for instance when having only a handheld with extremely limited resources. Note that Fred can define as many or as few operating modes as he feels appropriate, and indicate his relative preference for each by sliding the corresponding bar.

Suppose that for typing the notes (*edit text* service), Fred prefers *MSWord* over *Emacs*, and is unwilling to use the *vi* editor at all. This is an example of *supplier preferences*. Note that representing supplier preferences by discriminating the supplier type is a compact representation for the preferences with respect to the availability of desired features, such as spell checking or richness of editing capabilities, as well as to the user’s familiarity with the way those features are offered. For the sake of space, the user interface for specifying supplier preferences is not shown, but it is similar to the tabular form shown in Fig. 2.

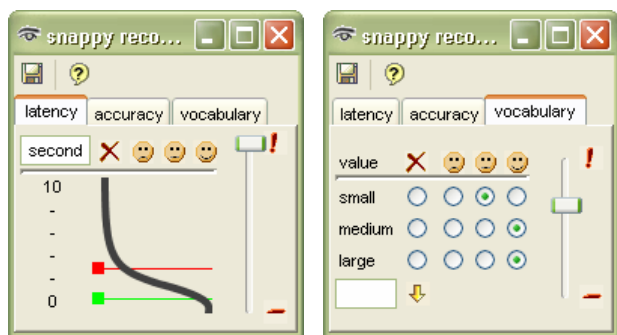


Fig. 2. QoS preferences for the *speech recognition* service

Suppose now that Fred will be browsing the web over a wireless network link. Suppose that the bandwidth suddenly drops: should the browser preserve the full quality of web pages at the expense of download time, or reduce the quality, for instance by skipping images? The answer depends on Fred’s *QoS preferences* for the current task. For browsing citations, Fred probably will be fine with dropping images and banners, with benefits in response times. However, for browsing a museum’s site on painting or online mapping, Fred may prefer full page quality to be preserved at the expense of download times.

As another example of QoS tradeoff, Fig. 2 shows an example of QoS preferences for the speech recognition service. The service has three dimensions: *latency*, *accuracy* and *vocabulary* size. The first two are numeric: the latency of recognizing each utterance is expressed in seconds, and accuracy reflects the percentage of words that are recognized accurately (not shown, for the sake of space). The user manipulates the

good and bad thresholds by dragging the green (lighter) and red (darker) handles, respectively.² Note that the utility space is represented simply using four intervals: from the lowest where the user prefers the configuration not to be considered, represented by a cross (X), to the highest corresponding to satiation, represented by a happy face (😊). The slide bar associated to each dimension captures how important, that is how much the user cares, about variations along that dimension.

We don’t expect every user to interact with the system at this kind of detail for every task. Rather, the infrastructure provides a set of templates for each service type, corresponding to common situations. For instance, the speech recognition service includes the *snappy recognition* template shown in Fig. 2, as well as the *accurate recognition* template, where the latency thresholds are relaxed, and accuracy and vocabulary more strict. The user can choose which preference template to apply to each service when defining a task (Fig. 1) or, by selecting customized tuning, manipulate preferences directly.

1) Representing User Preferences

To make preferences easier to both elicit and process, we make two simplifying assumptions. First, preferences are modeled independently of each other. In other words, the utility function for each aspect captures the user’s preferences for that aspect independently of others. Second, preferences fall into two categories: those characterized by enumeration, and those characterized by numeric values. Supplier preferences are characterized by enumeration (e.g. *MSWord*, *Emacs*, or *other*), and so are QoS dimensions such as audio fidelity (e.g. *high*, *medium* and *low*). For these, the utility function takes the form of a discrete mapping to the *utility space* (see below).

For preferences characterized by numeric values, we distinguish two intervals: one where the user considers the quantity to be good enough for his task, the other where the user considers the quantity to be insufficient. *Sigmoid* functions, which look like smooth step functions, characterize such intervals and provide a smooth interpolation between the limits of those intervals (see Fig. 2). Sigmoids are easily encoded by just two points: the values corresponding to the knees of the curve that

```
<utility combine="product">
  <QoSdimension name="latency" type="float">
    <function type="sigmoid" weight="1">
      <thresholds good="1" bad="3" unit="second"/>
    </function>
  </QoSdimension>
  <QoSdimension name="accuracy" type="float">
    <function type="sigmoid" weight="0.7">
      <thresholds good="90" bad="40" unit="percent"/>
    </function>
  </QoSdimension>
  <QoSdimension name="vocabulary" type="enum">
    <function type="table" weight="0.7">
      <entry x="small" f_x="0.8"/>
      <entry x="medium" f_x="1"/>
      <entry x="large" f_x="1"/>
    </function>
  </QoSdimension>
</utility>
```

Fig. 3. Internal representation of the QoS preferences in T 2

² The upper limit of the scale adjusts automatically between the values 10, 50, 100, 500, and 1000, further changes being enabled by a change in unit.

define the limits *good* of the good-enough interval, and *bad* of the inadequate interval. The case of “more-is-better” qualities (e.g., *accuracy*) are as easily captured as “less-is-better” qualities (e.g., *latency*) by flipping the order of the *good* and *bad* values. In case studies evaluated so far, we have found this level of expressiveness to be sufficient.

Fig. 3 shows the internal representation of the preferences captured in Fig. 2. Note that the infrastructure creates user interfaces like the one in Fig. 2 dynamically, based on the internal representation, which in turn is updated by manipulating the representations in the interface.

B. Formal Underpinnings

This section describes how user preferences, as defined in the previous section, guide the automatic configuration and reconfiguration of the environment. Our approach is based on finding the best match between the user’s needs and preferences for a specific task, and the environment’s capabilities. This framework is used both to find the optimal initial configuration, and to address the ongoing optimization of the support for the user’s tasks.

In practice, finding such a match corresponds to a constrained maximization problem. The function to be maximized is a *utility function* that denotes the user preferences, and the *constraints* are the environment’s capabilities and available resources. The result of the maximization is an abstract measure of the feasibility of carrying out the task, given the current conditions in the environment.

Utility space. Utility functions map the *capability space* (see below) onto the utility space. The latter is represented by the real number interval $[0, 1]$. The user will be happier, the higher the values in the utility space. The value 0 corresponds to the environment being unacceptable for the task; and 1 corresponds to user satiation, in the sense that increasing the capabilities of the environment will not improve the user’s perception of feasibility of the specific task.

Capability space. The capability space C_s corresponding to service s is the Cartesian product of the individual quality dimensions d of the service:

$$C_s \triangleq \otimes_{d \in \text{QoS dim}(s)} \text{dom}(d)$$

For example, possible quality dimensions for the *play video* service are *frame update rate*, the *frame size*, and *audio quality*. Thus, the capability space of video playing is three-dimensional. Cartesian product is used to combine the capability space of two services. For distinct services s and t , their combined capability space is formally expressed as:

$$C_{s \cup t} \triangleq C_s \otimes C_t$$

For example, a *web browsing* service has two quality dimensions: *latency* and *page richness*, and *video playing* has 3 dimensions of quality. Thus joint capability space of *video playing* and *web browsing* has 5 quality dimensions.

Typically, an application supports only a subset of the capability space corresponding to its various fidelities of output. In practice, approximating this subset using a discrete enumeration of points provides a reasonable solution, even if the

corresponding capability space is conceptually continuous. For example, while it makes sense to discuss a video stream encoding of decimal frames per second, typically video streams are encoded at integer rates. Despite discrete approximation, our approach does allow the handling of a rich capability space. For example, the capability space of a specific video player application can have 90 points, which is made possible by combining 5 frame rates, 6 frame sizes, and 3 audio qualities. Such a capability space can be made possible by encoding the same video in multiple frame rates, frame size, and audio quality, and possibly leveraging application-specific features such as video smoothing.

An application profile specifies a discrete enumeration of the capability points supported by an application and corresponding resource demand for each point. Note that specific *mechanisms* for obtaining and expressing application profiles exist. As demonstrated in [20], resource demand prediction based on *historical* data from experimental profiling is both feasible and accurate. Further, *metadata* and *reflection* can be used to express application profiles [3].

Application profiles describe the relationship between the capability points supported by applications, and the corresponding resource requirements. Formally, the quality resource mapping of supplier p is a partial function from the capability space of service s to the resource space: $\text{QoSProf}_p : C_s \mapsto R$. The range of the function is the subset of the capability space that is supported by the supplier.

Resource Space. The resource space R is the Cartesian product of the individual resource dimensions r of the entire environment E :

$$R \triangleq \otimes_{r \in \text{RES dim}(E)} \text{dom}(r)$$

Examples of resource dimensions are: CPU cycles, network bandwidth, memory, and battery. The actual number of resource dimensions is dependent on the environment.

Utility Functions. There is one utility function for each alternative configuration for a given task. The feasibility of the task corresponds to the best utility among the alternatives, weighted by the user’s preference for each alternative. The utility function for each configuration has two components, reflecting QoS and supplier preferences, respectively.

QoS preferences specify the utility function associated with each QoS dimension. The names of the QoS dimensions are part of the vocabulary shared between the user and the system. The utility of service s as a function of the quality of service is given by:

$$U_{\text{QoS}}(s) \triangleq \prod_{d \in \text{QoS dim}(s)} F_d^{c_d}$$

where for each QoS dimension d of service s , $F_d : \text{dom}(d) \rightarrow (0,1]$ is a function that takes a value in the domain of d , and the weight $c_d \in [0,1]$ reflects how much the user cares about QoS dimension d . As an example, *video playing* has a QoS dimension of *frame update rate*. The function $F_{\text{frameRate}}$ gives utility for various frame rates, and $c_{\text{frameRate}}$ specifies the weight of frame rate.

To evaluate the assignment of specific suppliers, we employ a supplier preference function, which is a discreet function that assigns a score to a supplier, based on its type. Also, we account for the *cost of switching* from one supplier to another at run time.

Precisely, the utility of the supplier assignment for a set a of requested services is:

$$U_{Supp}(a) \hat{=} \prod_{s \in a} h_s^{x_s} \cdot F_s^{c_s}$$

where for each service s in the set a , $F_s : Supp(s) \rightarrow (0,1]$ is a function that appraises the choice for the supplier for s ; and the weight $c_s \in [0,1]$ reflects how much the user cares about the supplier assignment for that service.

The term $h_s^{x_s}$ above expresses a change penalty as follows: h_s indicates the user's tolerance for a change in supplier assignment: a value close to 1 means that the user is fine with a change, the closer the value is to zero, the less happy the user will be. The exponent x_s indicates whether the change penalty should be considered ($x_s=1$ if the supplier for s is being exchanged by virtue of dynamic change in the environment) or not ($x_s=0$ if the supplier is being newly added or replaced at the user's request).

The overall utility is the product of the QoS preference and supplier preference. The overall utility over a set a of suppliers is:

$$U_{overall}(a) = \prod_{s \in a} h_s^{x_s} \cdot F_s^{c_s} \left(\prod_{d \in QoS \dim(s)} F_d^{c_d} \right)$$

1) The Optimization Problem

The optimization problem is to find a supplier assignment a , and for each supplier p in this assignment, a capability point such that the utility is maximized:

$$\arg \max_{\substack{p_s \in Supp(s) \\ f_p \in dom(d)}} \prod_{s \in a} h_s^{x_s} \cdot F_s^{c_s}(p_s) \cdot \left(\prod_{d \in QoS \dim(s)} F_d^{c_d}(f_{p,d}) \right)$$

The maximization is over a set of constraints, which we express below. The capability constraint stating that the chosen point $f_{p,d}$ is in the capability space for supplier p is as follows:

$$\forall_{p \in Supp(s)} f_p = \otimes_{d \in QoS \dim(s)} f_{p,d} \in C_p$$

And to ensure that the resource constraints are met:

$$\sum_{p \in Supp(s)} QoSprof_p(f_p) \leq |R|$$

where summation is in the vector space \mathbf{R} of resources, and the inequality is observed in each dimension of that space. In non-mathematical terms, this constraint expresses the fact that the aggregate resource demand by all the suppliers can not exceed the resource supply.

C. Algorithm and Analysis

In this section we solve the optimization problem. The optimization algorithm must be efficient to be usable at runtime. Two metrics we are interested in are the latency of computing an answer to a given instance of the problem, and in the com-

```

HashMap SuppPrefs; // supplier preferences
HashMap QoSPrefs; // qos preferences
HashMap SuppReg; // registered suppliers

Config BestConfig(Set reqstdSvcs){
// 1. QUERY
Map suppListsBySvc;
for each svc in reqstdSvcs do{
List suppList = null;
Pref suppPref = SuppPrefs.get(svc);
// query for supp based on svc type
suppList = query(SuppReg, svc, suppPref);
suppListsBySvc.put(svc.type, suppList);
}
// 2. GENERATE configs, compute supp pref
List configs = GenConfigs(suppListsBySvc);
configs = sort(configs);
// 3. EXPLORE the QoS space
int indexBestConfig;
float overallUtilBest = 0.0;
for each i from configs.size-1 to 0 do {
Config cCur = configs.get(i);
if (overallUtilBest > cCur.suppPrefUtil)
break;
cCur = searchQoS(cCur, QoSPrefs);
if (cCur.overallUtil > overallUtilBest){
indexBestConfig = i;
overallUtilBest = cCur.overallUtil;
}
}
return configs.get(indBestConfig);
}

GenConfigs(Map suppListsBySvc){
List configs = new List(MAX_INT);
int depth = 0;
Config partialConfig = null;
GenConfigsRecur(depth, configs,
suppListsBySvc, partialConfig);
}

GenConfigsRecur(int d, List configs,
Map suppListsBySvc, Config partialConfig){
if (d == suppListsBySvc.size()){
configs.add(new Config(partialConfig));
return;
}
List suppList = suppListsBySvc.getByInd(d);
for each supp in suppList do{
partialConfig.add(d, suppList);
GenConfigsRec(d+1, configs,
suppListsBySvc, partialConfig);
partialConfig.remove(d);
}
}

```

Fig. 4. Pseudocode of the algorithm

putational overhead of the algorithm.

1) The Algorithm

The algorithm works in three phases: (1) query, (2) generate, and (3) explore. In the first phase, it queries for relevant suppliers for each service in the task. In the second phase, it combines suppliers into configurations and ranks them according to the supplier preference only. In the third phase, it explores the quality space of the configurations. The pseudocode for the algorithm is shown in Fig. 4.

The double product term of the utility formula in B.1) allows for a clever exploration strategy. The outer product is the supplier preference score. It can be computed at the time the supplier assignment is known (in phase 2), and can be used as an upper bound for overall utility during the explore phase. Since overall utility is the product of supplier preference and QoS preference, and the latter is bounded by one, then maximum overall utility is bounded by supplier preference. The break in the loop in BestConfig takes advantage of that fact.

Consider a simple example. Assume that two services are

requested. For each service, there are two possible suppliers: a_1 and a_2 for the first service, b_1 and b_2 for the second, yielding 4 possible configurations as shown in **Table 3**. The search space can be divided into 4 quadrants, each representing the capability space of a specific configuration. We are searching for a point with the highest utility.

As noted, the maximum utility that can be achieved within each quadrant is bounded by the supplier preference portion of utility. These observations help provide a stop condition for the search: once a point is found that has overall utility of Λ , there is no need to explore configurations with supplier preference portion of utility of less than Λ .

TABLE 3. THE STRUCTURE OF THE SEARCH SPACE

a_1, b_1	a_1, b_2
a_2, b_1	a_2, b_2

In Table 3, the shading of each quadrant reflects the hypothetical values of supplier preference portion of utility for each configuration: the darker the shade, the higher the value. Assume these values are: .8, .6, .4, and .2. Each of these values is an upper bound for maximum *overall utility* possible from the respective quadrant. We explore inside the quadrants, starting from the darkest. If the maximum utility for the quadrant a_1, b_1 is higher than 0.6, then at this time we know the best point in the entire space is found, and can stop the search. If not, we continue the search in quadrant a_2, b_1 , and so on.

Exploring the quality space of a configuration is a variant of a 0-1 Knapsack problem, called multiple dimensional, multiple choice 0-1 Knapsack. Multiple dimensions refer to the multiple constraints that are present in the problem. Multiple choice refers to choosing one among a set of similar items. In our problem, resources map to knapsack dimensions and the capability space of one service maps to one set of similar items. This is a well-studied problem in the optimizations research, and is at the core of such optimization problems as winner determination in combinatorial algorithms. [16,23] show the problem to be NP-complete, and give approximation algorithms. [23] gives an exact solution that is demonstrably fast on inputs drawn from certain probability distributions.

One of the approximating algorithms to the problem uses utility to resource ratio as a metric for ranking the capability points, it then applies greedy branch-and-bound and LP-relaxation to find a near-optimal answer. In the multiple resource case, quadratic weighted-average is used to compute a single resource currency from multiple resources, and the solution to the single resource case is reused iteratively [16].

In our solution, SearchQoS invokes a third-party library called Q-RAM, the package described in [16].

2) Analysis

To analyze the running time of the algorithm, let:

\mathbf{n} be the number of requested services

\mathbf{P} be the total number of available suppliers

\mathbf{p} be the number of suppliers for a given service type

\mathbf{q} be the size of the capability space of a supplier.

\mathbf{P} and \mathbf{p} describe the *richness* of the environment, and can potentially increase as more applications, hardware, and de-

VICES are made available. \mathbf{q} describes the capability richness of a supplier. We can assume that the size of the user task is limited to a small number of applications. Thus \mathbf{n} is bounded.

Next we analyze the running time of the three phases.

The query phase retrieves items from a hashtable. Retrieving one item is logarithmic in the size of the hashtable. \mathbf{n} retrievals from a hashtable of size \mathbf{P}/\mathbf{p} take $O(\mathbf{n} \cdot \log(\mathbf{P}/\mathbf{p}))$.

The generate phase is a recursion of depth \mathbf{n} , with a loop of size \mathbf{p} at each level. Thus, it takes $O(\mathbf{p}^{\mathbf{n}})$.

The explore phase in the worse case takes $O(\mathbf{p}^{\mathbf{n}}) \cdot O(\text{searchQoS})$. The size of the QoS space of a configuration of \mathbf{n} suppliers each of which has a capability space of size \mathbf{q} is $O(\mathbf{q}^{\mathbf{n}})$. Approximation algorithm we use can search that space in time $O(\mathbf{n} \cdot \mathbf{q} \cdot \log \mathbf{q})$ [16,23]. Thus the explore phase takes $O(\mathbf{p}^{\mathbf{n}}) \cdot O(\mathbf{n} \cdot \mathbf{q} \cdot \log \mathbf{q})$ in the worst case, and dominates all other terms. The first term, $O(\mathbf{p}^{\mathbf{n}})$, presents a possible scalability bottleneck.

Let us demonstrate how the exploration strategy described earlier helps tackle that bottleneck. Recall the break condition in the explore phase, illustrated in the example introduced in IV.C.1). The number of configurations that are explored will depend on the distribution of the supplier preference values, and Λ , the highest achievable utility value. Let's assume an average number of suppliers per service $\mathbf{p} = 10$, and a specific distribution of supplier preference values that is uniform, i.e. the most preferred supplier scores 0.9^0 , the next one scores: 0.9^1 , etc. In Table 4, we show the number of configurations generated, and the number of configurations that are actually explored depending on the value of maximum achievable utility, Λ , and number of services in the task, \mathbf{n} .

TABLE 4. NUMBER OF CONFIGURATIONS GENERATED AND EXPLORED FOR VARIOUS VALUES OF \mathbf{n} , AND Λ , MAXIMUM UTILITY ACHIEVED

	$\mathbf{n}=1$	2	3	4	...	8
Generated	10	10^2	10^3	10^4	...	10^8
$\Lambda = .9$	2	3	4	5	...	8
$\Lambda = .81$	3	6	10	15	...	36
$\Lambda = .73$	4	10	20	35	...	120
$\Lambda = .66$	5	15	35	70	...	330

The first row shows the number of services. The second row shows the number of configurations generated, which is $\mathbf{p}^{\mathbf{n}}$, in this case, $10^{\mathbf{n}}$. In each subsequent row, we show the number of configurations that are sufficient to explore, if the maximum utility shown in the first column in that row is actually achieved by some configuration. For instance, for a task with 4 requested services, even if the maximum utility achievable is as modest as $\Lambda = 0.6$, then the number of supplier configurations explored is 126, which is two orders of magnitude smaller than the 10^4 , the total number of configurations.

3) Reconfiguration

The algorithm also handles reconfiguration. When there is a running configuration, the utility from the best computed configuration is compared with the *observed* utility of the running

configuration, and a switch is made if the latter is lower than the former. Because of the cost of change term in the supplier preferences formula, some of the suppliers in the running configuration may get advantage.

V. EVALUATION

A. Case Study

In this section we report on a case study of automatically configuring an environment for the task of reviewing a documentary video clip. The task requires three services: *video playing*, *text editing*, and *browsing*. The user watches the clip, takes notes, while browsing the net for information. The quality dimensions are as follows: for video playing: frame rate, frame size, and audio quality; for browsing: latency of loading pages, and richness of the pages (pages have graphics that are not required for the task can be helpful); for text editing: none.

We performed the case study in two steps. In the first step we collected application profile data, specified preferences, and identified resource limits. In the second step, we ran a prototype implementation of the algorithm.

1) Input Data Collection

As an experimental platform, we chose an IBM Thinkpad 30 laptop, equipped with 256 MB of memory, 1.6 Ghz CPU, WaveLAN card, and Windows XP Professional. In power saving mode, the CPU can run at a percentage of the maximum speed, effectively creating a tight CPU constraint.

The model requires three inputs: (1) user preferences, (2) application profiles, and (3) resource availability. For the purposes of this experiment, we used synthetic preferences intended to be representative of the tasks. We identified several applications that supported various facets of the task. Those applications were installed on the laptop. To obtain application profiles, we measured resource usage corresponding to a small set of capability points. We performed this profiling *offline*, with each supplier running separately. Resource availability is as follows: 400 MHz of processing power, when the CPU is running at $\frac{1}{4}$ of the baseline speed; 64 MB of free memory after excluding the memory taken by the operating system and other essential critical systems; and 512 Kbps of bandwidth, provided by an 802.11 wireless access point backed by a DSL line.

The applications used in the experiment were:

- **Video players:** RealOne and Windows Media,
- **Text editors:** TextPad, WordPad, Notepad, Microsoft Word, and GNU Emacs,
- **Browsers:** Internet Explorer, Netscape, and Opera,

These suppliers allow a total of $30 = 2 * 5 * 3$ configurations.

We measured CPU and physical memory load using Windows Performance Monitor. We used *percent processor time*, *working set* counters of the *Process* performance object to measure CPU and memory utilization respectively. We took the sampling *average* over a period of time. The performance monitoring API does not provide per process network statistics, so the mechanism for measuring bandwidth demand was different in each case, as explained below.

For a representative clip to watch, we obtained a two minute trailer of a movie in Windows native .wmv and Real Networks native .rpm in several different bit-rates. Where cross-player compatibility is supported, we obtained additional capability points. For example, RealOne plays .wmv format. Also, players provide quality knobs, allowing improved quality in exchange for higher CPU utilization. For example, Windows Media player supports video smoothing that provides higher frame rate than the rate encoded in the stream. For each player, 32 points quality points were sampled. To measure bandwidth demand, we consider the bit-rate of the stream, and cross-check with the application-reported value. We observed that the CPU consumption of different players can be widely different for the same quality point.

We measured CPU and memory used while typing and formatting text for 2 minutes with each text editor. We observed that the memory consumption of different text editors can be widely different.

All browsers surveyed support a text-only mode, providing two points in the page richness dimension. To obtain different levels of latency, we used a bandwidth-limiting http proxy, and pointed the browser to the proxy. We measured latency by allowing the following bandwidth limits: 28, 33, 56, 128, 256, 512 Kbps. Our script included a sequence of approximately 15 pages with a mix of both text graphics on the internet. Each test started with a clean browser cache. 16 quality points were sampled. We observed that the browsers have very similar resource consumption patterns.

Although we realize that the methods for obtaining resource consumption measures are not precise, we believe that they yield good enough approximations for this feasibility analysis.

Note that the capability space of a configuration of suppliers has approximately 500 points ($32 * 16 * 1$), based on the samples taken. 30 configurations together provide a capability space of approximately 15,000 points.

2) Prototype Evaluation

The algorithm is guaranteed to find an optimal assignment of suppliers. Furthermore, it will obtain the optimal set of quality points for the suppliers, as long as Q-RAM finds the optimal point inside each quadrant. Whenever Q-RAM returns a near-optimal answer, our algorithm will return a near-optimal set of quality points.

Additionally, we evaluated a prototype implementation of the algorithm according to two metrics: (1) latency, and (2) system overhead. Latency measures the time it takes to compute an answer, from the time that a client program requests it. Overhead measures percent CPU and memory utilization of the algorithm. To adapt the configuration in response to environment changes, it is necessary to run the algorithm periodically. Thus, overhead of periodic invocation is a useful metric.

The latency of computing the best configuration averaged over 10 trials was 531 ms. In the query and generate phases, the algorithm spends less than 10 ms each. In the explore phase, it spends just under 500 ms (approximately 10 ms was due to parsing the request, and formatting the answer). The bulk of the time in the explore phase was due to external proc-

ess invocation and file input-output (Q-RAM package is an external executable). Thus, the latency can be significantly reduced by linking into Q-RAM in-process.

We invoke the algorithm a total of 50 times in 5 second intervals over a period of 250 seconds, and measure *average CPU utilization*. Average CPU utilization is 3.8%. This overhead is fairly low, and can be further lowered by running the algorithm less frequently, e.g. once per 10 or 25 seconds.

Memory usage of the process running the algorithm is approximately 8.8 MB. While this is a significant overhead, most of it is due to the Java virtual machine

B. Lessons Learned

The form of self-adaptation that we address in this work is targeted at support for everyday computing in ubiquitous computing environments. To accomplish this it is essential. However, using existing applications in existing environments is a challenge since these applications are not in general designed for self-configuring capabilities, such as those that we are attempting to provide.

Our approach to handling legacy applications has been to write wrappers around existing application to present those applications as suppliers for particular services in Aura. These wrappers mediate communication with Aura so that the application state can be set and retrieved, and so that resource usage can be monitored. However, our experience in writing over a dozen wrappers for applications in Windows and Linux is that obtaining the information we need is often challenging, even if the applications are developed by the same vendors.

To summarize the requirements on suppliers in the Aura environment we divide them into two categories:

1. **Mobility:** To allow applications to be configured for use with user tasks, we require a way to be able to get and set the state of those applications. This allows not only mobility, but coarse-grained adaptation where we need to change a supplier.
2. **Adaptation:** To facilitate finer-grained adaptation, we need to access the resource usage statistics of applications, and to be able to restrict their resource usage.

1) Mobility

As described in Section III.B, the task layer requires that the user-level state of application be retrieved and set. This is so that task state can be transferred between environments, and also so that tasks can be suspended and resumed by a user. This requires some consensus about the meaning of the state of a particular generic service, such as text editing. However, the Aura infrastructure also allows applications to augment the state with additional properties. If these properties do not have meaning for other suppliers of the same service, we require that those suppliers preserve those properties.

In our experience with developing suppliers, we have had mixed success with getting state information from applications. While more recent applications allow reflective access to get and set this information through programmatic interfaces such as .NET, it is not as easy with older applications. Even when applications provide a programmatic interface, it is

possible that they do not expose the required information. For example, to restore the state of web browsing, it is desirable to set and get the history of the browser so that backward and forward browsing state can be maintained. Internet Explorer, while providing a COM interface for setting the current web page, does not provide these additional APIs.

The fact that we require this information from applications may at first seem like a lot to ask of application developers. However, many applications are increasingly allowing access to this information. In our experience, it has always been possible to get and set some form of the state; the difficult issue is in the varying mechanisms that we have had to use, and the scope of the information that we have access to.

Furthermore, our experience with writing suppliers has demonstrated the need for applications to be able to “sandbox” the set of materials from one task for another. For example, suppose that one user task requires two spreadsheets to be edited, while another task requires on spreadsheet to be edited. If these tasks are resumed simultaneously, the originating task of each spreadsheet needs to be remembered. Microsoft Excel supports this, because there can be multiple instances of Excel running, each with their own set of files. On the other hand, Microsoft Powerpoint makes this difficult, because only one application instance is allowed, therefore it is not possible to keep track of the origin of files.

2) Adaptation

To allow for adaptation to changing resources, in Section IV.B we described a formalism that can provide optimal configurations or reconfigurations based on the available supply of resources, and the expected resource usage of the applications. For the mathematical formalism to work in practice, Aura infrastructure requires information about resource usage and quality of service of applications, resource supply in the environment, as well as certain level of cooperation from applications about expected resource usage.

In practice, the following set of requirements need to be satisfied in order for the mathematical model of Aura to work in practice: (1) ability to monitor application-provided quality of service, (2) ability to monitor and report application resource usage, (3) ability to monitor available resource supply, and (4) ability to enforce resource usage limits on applications.

Let us discuss how each of these requirements can be translated into design guidelines for application and system developers.

Requirement 1 can be satisfied directly by application developers by exposing rich APIs that report the quality of service. For example, many of the commercial and open-source video/media players report the richness of the stream in bit rates, the frame update rate of the video, the size of the frame, the color depth, etc.

Requirements 2 and 3 can be satisfied by a shared service that is either provided by the operating system, or third-party middleware. While modern-operating systems generally provide reasonable performance and resource monitoring hooks, there is room for improvement. However, some resources can be more difficult to account for on a per-process basis (e.g.,

battery). Notice that there are research operating systems, e.g., Nemesis [17] that specifically provide accurate resource accounting and enforce resource usage limitations.

With respect to requirement 4, we believe that a two-fold approach is needed. First, applications can provide various adaptation strategies (e.g., more CPU-intensive video stream decoding or less CPU-intensive decoding); we believe that applications should provide the ability to comply with resource usage limitations. Some of the video players on the market provide such ability directly, e.g., with respect to network bandwidth. However, it is also desirable to have an operating system-provided mechanism for ensuring that resource limitations are enforced if an application proves to be uncooperative. For example, some video players aggressively prefetch and saturate use all available bandwidth, despite being told to use low bit rate stream. In such cases, a mechanism external to application (e.g., bandwidth throttling) can enforce resource limitations imposed on applications.³

VI. CONCLUSION AND FUTURE WORK

In this paper we have described an approach to self-configuring capabilities for everyday computing environments. Motivated by the challenges of supporting heterogeneity, resource variability, mobility, ubiquity, and task-specific user requirements, we have developed a self-adaptation infrastructure that has three distinctive features. It allows explicit representation of user tasks, including preferences and service qualities. It provides an environment management capability to translate user-oriented task and preference specifications into resource allocations that match the intended environment. Finally, it provides a formal basis for understanding the resource allocation and derived algorithms that support optimal allocation at run time.

While providing a good starting point, this work also suggests a number of important future directions. First is the extension of task specification so it can express richer notions of task, such as work flow, cognitive models, and goal driven task realization. Second is the extension of resource allocation algorithms to take advantage of future predictions. This entails much richer notions of utility, such as those prescribed by Options Theory. Finally, there are many directions that one could pursue in the area of user interface design to make it even easier for users to create and reuse task descriptions.

ACKNOWLEDGEMENTS

This research is supported by the National Science Foundation under Grant ITR-0086003, by the Sloan Software Industry Center at Carnegie Mellon, and by the High Dependability Computing Program from NASA Ames cooperative agree-

ment NCC-2-1298.

REFERENCES

1. R.K. Balan, J.P. Sousa, M. Satyanarayanan. Meeting the Software Engineering Challenges of Adaptive Mobile Applications. *Tech. Report, CMU-CS-03-11*, Carnegie Mellon U., Pittsburgh, PA, 2003.
2. L. Capra, W. Emmerich and C. Mascolo. A Micro-Economic Approach to Conflict Resolution in Mobile Computing. *In Proc Foundations of Software Engineering (ACM SIGSOFT/FSE)*, 2002.
3. L. Capra, W. Emmerich and C. Mascolo. Reflective Middleware Solutions for Context-Aware Applications. *Proc Int'l Conference on Metalevel Architectures and Separation of Crosscutting Concerns (REFLECTION)*, 2001. LNCS (Springer-Verlag).
4. Cheng, S.W. et al. Software Architecture-based Adaptation for Pervasive Systems. *Proc of the International Conf. on Architecture of Computing Systems: Trends in Network and Pervasive Computing*, April 2002. Springer LNCS Vol. 2299, 2002.
5. F. Cristian. Understanding Fault-Tolerant Distributed Systems. *Communications of the ACM*, 34(2):56-78, 1991.
6. The DAML Services Coalition (multiple authors), "DAML-S: Web Service Description for the Semantic Web", *Proc Int'l Semantic Web Conference (ISWC)*, 2002.
7. J. Flinn, E. de Lara, et al. Reducing the Energy Usage of Office Applications. *Proc. IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, 2001.
8. D. Garlan, D. Siewiorek, A. Smailagic, and P. Steenkiste. Project Aura: Towards Distraction-Free Pervasive Computing. *IEEE Pervasive Computing*, Volume 21, Number 2, April-June, 2002.
9. D. Garlan, J. Kramer, J., and A. Wolf (eds). Proceedings of the ACM SIGSOFT Workshop on Self-Healing Systems (WOSS'02), Charleston, SC, November 18-19, 2002.
10. D. Garlan, S.W. Cheng, B. Schmerl. Increasing System Dependability through Architecture-Based Self-repair. *Architecting Dependable Systems*, R. Lemos, C. Gacek, A. Romanovsky (Eds), Springer-Verlag, 2003.
11. Georgiadis, I., Magee, J., Kramer, J. Self-Organising Software Architectures for Distributed Systems. *Proc. ACM SIGSOFT Wksp on Self-Healing Sys. (WOSS'02)*. Nov. 2002.
12. M. Hiltunen, R. Schlichting. Adaptive Distributed and Fault-Tolerant Systems, *International Journal of Computer Systems Science and Engineering*, 11(5):125-133, September, 1996.
13. Jini. www.jini.org. Accessed: Sep. 2003.
14. M. Jones, D. Rosu, M. Rosu. CPU Reservations and Time Constraints: Efficient, Predictable Scheduling of Independent Activities. *Proc ACM Symp Operating Systems Principles (SOSP)*, 1997.
15. E. de Lara, D. S. Wallach, W. Zwaenepoel. Puppeteer: Component-based Adaptation for Mobile Computing. *Proc. 3rd USENIX Symposium on Internet Technologies and Systems (USITS)*, 2001.
16. C. Lee, et al. A Scalable Solution to the Multi-Resource QoS Problem. *Proc. IEEE Real-Time Systems Symposium (RTSS)*, 1999.
17. I. M. Leslie, D. McAuley, R. Black, T. Roscoe, P. Barham, D. Evers, R. Fairbairns, and E. Hyden. The Design and Implementation of an Operating System to Support Distributed Multimedia Applications. *IEEE Journal on Selected Areas In Communications*, 14(7):1280-1297, Sept 1996.
18. J. Kephart, M. Parashar (eds) *Proc. International Conference on Autonomous Computing*, May 17-18, 2004. New York, NY. IEEE Press.
19. F. Kon, et al. Dynamic Resource Management and Automatic Configuration of Distributed Component Systems. *Proc. USENIX Conference on OO Technologies and Systems (COOTS)*, 2001.
20. D. Narayanan, J. Flinn, M. Satyanarayanan. Using History to Improve Mobile Application Adaptation. *Proc. 3rd IEEE Workshop on Mobile Computing Systems and Applications (WMCSA)*, 2000.
21. R. Neugebauer and D. McAuley. Congestion Prices as Feedback Signals: An Approach to QoS Management. *Proc. ACM SIGOPS European Workshop*, 2000.

³ Notice that we are not advocating here the need for applications to interfere with the low level scheduling of resources by the operating system. We simply advocate that on the level timescale of seconds the resource usage by applications should be consistent with the expected quality of service delivered.

22. B. Noble, et al. Agile Application-Aware Adaptation for Mobility. *Proc. of the 16th ACM Symp. on Operating Systems Principles (SOSP'97)* October 1997. *Operating Systems Review* 31(5), ACM Press, 276-287.
23. D. Pisinger. An exact algorithm for large multiple knapsack problems. *European Journal of Operational Research*, 114, (1999).
24. Service Location Protocol. www.srvloc.org. Accessed: Sep. 2003.
25. J.P. Sousa, D. Garlan. The Aura Software Architecture: an Infrastructure for Ubiquitous Computing. *Tech. Report CMU-CS-03-183*, Carnegie Mellon U., Pittsburgh, PA, 2003.
26. J.P. Sousa, D. Garlan. Beyond Desktop Management: Scaling Task Management in Space and Time. *Tech. Report CMU-CS-04-160*, Carnegie Mellon U., Pittsburgh, PA, 2004.

Vahe Poladian is a doctoral candidate at Carnegie Mellon University. His research interests are in applications of microeconomics, utility theory, and decision theory in configuring software on mobile computing platforms. He received his BS in Computer Science and Mathematics from Macalester College.

João Pedro Sousa is a doctoral candidate at Carnegie Mellon University. His research interests are in the area of increasing the benefit to users of Ubicomp environments. He received his Masters of Software Engineering at Carnegie Mellon University.

David Garlan is a professor of computer science at Carnegie Mellon University. His research interests include software architectures, formal methods, self-healing systems, and task-based computing. He received his PhD in computer science from CMU.

Bradley Schmerl is a systems scientist at CMU. His research interests include dynamic adaptation, software architectures, and software engineering environments. He received his PhD in computer science from Flinders University in South Australia.

Mary Shaw is the Alan J. Perlis professor of computer science at Carnegie Mellon University. Her research interests include value-based software engineering, everyday software, software engineering research paradigms, and software architecture. She received her PhD in computer science from Carnegie Mellon University.