

Towards an Elastic Application Model for Augmenting the Computing Capabilities of Mobile Devices with Cloud Computing

Xinwen Zhang · Anugeetha Kunjithapatham · Sangoh Jeong · Simon Gibbs

Published online: 2 April 2011
© Springer Science+Business Media, LLC 2011

Abstract We propose a new *elastic application* model that enables seamless and transparent use of cloud resources to augment the capability of resource-constrained mobile devices. The salient features of this model include the partition of a single application into multiple components called *weblots*, and a dynamic adaptation of weblot execution configuration. While a weblot can be platform independent (e.g., Java or .Net bytecode or Python script) or platform dependent (native code), its execution location is transparent—it can be run on a mobile device or migrated to the cloud, i.e., run on one or more nodes offered by an IaaS provider. Thus, an elastic application can augment the capabilities of a mobile device including computation power, storage, and network bandwidth, with the light of dynamic execution configuration according to device's status including CPU load, memory, battery level, network connection quality, and user preferences. This paper presents the motivation behind developing elastic applications and their architecture including typical

elasticity patterns and cost models that are applied to determine the elasticity patterns. We implement a reference architecture and develop a set of elastic applications to validate the augmentation capabilities for smartphone devices. We demonstrate promising results of the proposed application model using data collected from one of our example elastic applications.

Keywords elastic application · cloud computing · mobile device · weblot · dynamic execution configuration

1 Introduction

Applications on smartphones traditionally are constrained by limited resources such as low CPU frequency, small memory, and a battery-powered computing environment. For example, the iPhone 3G is equipped with 412 MHz CPU, 512 MB RAM, and a battery allowing about 5 hours of talking time. The new Samsung Galaxy Android phone has 528 MHz CPU, 128 MB RAM, and battery offering about 6.5 hours of talk time. Both devices have up to 7.2 Mbps 3G data network connection. Compared to today's PC and server platforms, these devices still cannot run compute-intensive applications such as complex media processing, search, and large-scale data management and mining.

Cloud computing delivers new computing models for both service providers and individual consumers including infrastructure-as-a-service (IaaS), platform-as-a-service (PaaS), and software-as-a-service (SaaS), which enable novel IT business models such as resource-on-demand, pay-as-you-go, and utility-

X. Zhang (✉)
Huawei Research Center, Santa Clara, CA, USA
e-mail: xinwen.zhang@huawei.com

A. Kunjithapatham · S. Gibbs
Samsung Information Systems America, San Jose, CA, USA

A. Kunjithapatham
e-mail: anugeetha.k@samsung.com

S. Gibbs
e-mail: s.gibbs@samsung.com

S. Jeong
LG Electronics, 221 Yangjae-Dong, Seocho-gu,
Seoul, Korea
e-mail: sangohjeong@gmail.com

computing [6]. From the perspective of service providers, cloud computing is often viewed as a vast and scalable platform for service delivery. We suggest a new perspective, one tuned to the needs of mobile devices. We consider cloud computing as a means to extend or augment the capabilities of resource constrained devices.

There are several approaches to realize this perspective. One approach is to duplicate the runtime environment of the device in the cloud and then run the application either on the device or in the cloud. The off-device runtime environment is sometimes called a “surrogate” [18], a “clone” [9], or a cloudlet [23]. Virtual machine technology is often used to host and isolate the off-device runtime so making this approach fit well with emerging IaaS platforms such as Amazon EC2 [1]. Running a device clone in the cloud has some attractive properties such as enhanced CPU and memory resources which lead to better performance. Furthermore, applications do not need any modification—the clone and the physical device can run identical binaries. However, this approach has disadvantages too. First, the application on the clone may need to access the physical hardware on the device. For example, consider a GPS application or simply the question of how an application running in the clone interacts with the user. It is certainly possible to transfer device I/O between the device and clone environment over the network, but this may impact responsiveness and battery use. Secondly, simply replacing one processor with another fails to take full advantage of cloud compute resources. Ideally, a cloud application should be able to run in a highly parallel fashion distributed over many cloud nodes.

Thirdly, completely duplicating a device and running it on the cloud increases the complexity of device management. For example, the cloud system needs similar security protection and data privacy control as those on the device since it runs all possible applications with data from the original device.

The above considerations lead us to focus on application level augmentation instead of cloning a complete device environment. Often these applications are data-parallel with high compute-to-communication ratio. Examples include media processing, search, and data mining. Our goal is to design an architecture and related middleware to enable *elastic applications* which consist of multiple components called *weblets*, each of which can be launched on a mobile device or in the cloud. The decision of where to launch a weblet is based on application configuration and/or the status of the device such as its CPU load and battery level. Ideally the application model could also

support migration of weblets between the device and cloud platform during runtime. While offloading and delegating computing have been proposed by many researchers [8, 11, 14, 18], the novelty of our approach lies in enabling flexible and optimized elasticity by considering multiple factors including device status, cloud status, application performance measures, and user preferences (e.g., different running modes of an application including power-saving mode, high speed mode, low cost mode, offline mode, or in terms of expected application specific modes).

To enable this new application model, many challenges exist in different areas, including management of heterogeneous computing environments, data management and communication dependencies between weblets, state synchronization between weblets, and cost-effective dynamic execution configuration. The middleware should provide infrastructure for seamless and transparent execution of elastic applications and offer convenient development support. This paper first gives the concepts and typical elasticity patterns (Section 2). We then focus on the optimization of cost-effective execution configuration by considering multiple factors (Section 3), which we believe is one of the most critical and unique components of the application model. We then present a high-level description of an implemented reference framework including deployment and runtime architecture and software development kit (SDK) (Section 4). We then illustrate a set of example elastic applications developed based on our reference architecture and cloud platform (Section 5). We show some experimental results which confirm the augmentation capabilities of our approach with collected data from an elastic image processing application (Section 6). We present some related work and oversee further research themes along this novel application model at the end of this paper (Section 7 and 8).

2 Concepts & elasticity patterns

2.1 Concepts and benefits

We define elastic applications as having two properties. First, following the client/server split of traditional web applications, an elastic application is split or partitioned so that execution occurs partially on the device and partially on the cloud. Previous work has proposed many mechanisms for splitting an application into modular components for remote execution or *cyber foraging* purposes, such as [7, 8, 11, 14, 21, 25]. For elastic devices we assume application developers can determine how to organize weblets based on their functionalities and

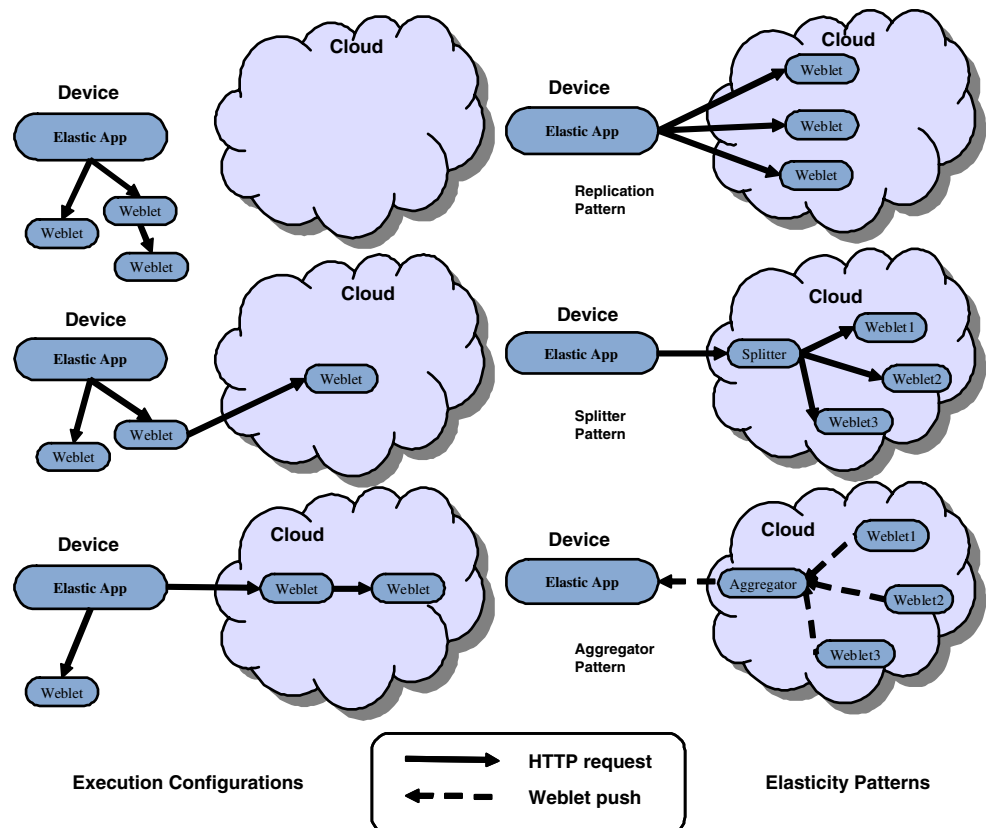
runtime behaviors such as computation demand, data dependency, and communication need, which we believe should be part of high-level design consideration of an application. Elastic middleware should provide necessary SDK and tools allowing developers to implement and test their designs. A unique requirement for elastic applications is that a weblet's functionality should not be affected by the location or environment where it is running. Essentially, the location of individual weblets should be transparent to users. One principle for partitioning applications is that each weblet should have minimum dependency on others. This is not only for robustness but decreases communication overhead between weblets during runtime.

Second, the *execution configuration* of an elastic application is not static, instead it is determined when the application is launched and potentially modified during runtime. By execution configuration, we mean the assignment of application partitions to execution units (e.g., cores or virtual machines), either on the device or in the cloud. The left hand side of Fig. 1 shows some possible execution configurations for an application using three weblets.

There are several benefits that the elastic application concept offers to mobile users and application develop-

ers deriving from coarse-grained application partitioning and dynamic configuration. First, elastic applications are not constrained by the compute capabilities of today's mobile platforms and can be configured to take advantage of multiple processing cores when available. If more compute (or storage) is needed then this can be obtained from the cloud. As devices become more powerful, compute and storage can shift back to the device. On the other hand, mobile device compute and storage need not be designed to satisfy the most demanding applications. Device resources can be modest (and less power consuming) since the more demanding applications can acquire resources from the cloud. From a performance perspective, the ability to allocate resources in the cloud and migrate functionality gives the device great flexibility. For example, performance can be increased or optimized to fit various goals (such as responsiveness, monetary cost, or power consumption). Furthermore, application components that are partitioned for migration can also be replicated. The failure then of one instance of a replicated component need not compromise the application. Also, the elastic application model offers a testbed for future technologies of mobile devices. Applications that run on the cloud today can move to the device in future products.

Fig. 1 Execution configurations and elasticity patterns



This greatly extends the lifetime of applications and reduces development costs.

2.2 Elasticity patterns

We now consider elastic applications and weblets in more detail. Our motivation for using weblets is that developers are familiar with the web application model and so can easily transition from the client/server partitioning of web applications to the more general form of partitioning found in elastic applications. Furthermore, programming methods used for web applications, for example AJAX and REST, are adapted by weblets. To see the similarities and differences of web applications and elastic applications, it is interesting to compare weblets with traditional web services. We highlight some areas for comparison in Table 1.

In designing a web application, a key issue is determining what logic will run on the server and what on the client. For early web sites, the client was mainly used for rendering and input, but now with JavaScript, AJAX, and plug-ins such as Flash and Silverlight, many tasks can be performed by the client. With elastic applications there is a similar issue, but because several weblets can be created by a single application, the topology of elastic applications is more varied. It appears these topologies fall into some common patterns, what we call *elasticity patterns*, several of these are shown on the right hand side of Fig. 1 and briefly summarized as follows.

Replication patterns: pools and shadowing Weblat replication refers to running multiple weblats with the same interface, i.e., accepting the same types of request. There are two forms of replication: *pools* and *shadowing*. Weblat pools allow an application to leverage cloud CPU cycles and augment its throughput. With this pattern, the application issues requests that are routed to weblats as they become available. Weblat pools are well suited for applications that are easily divided into similar tasks, for example processing sets

of images or scanning sets of files. Closely related to pools is shadowing in which the same request is sent to a set of replicated weblats in parallel. Shadowing can be used for fault tolerance and latency control. For example, shadowing a weblat on the device with a copy on the cloud can help the application recover from loss of network connectivity or loss of battery power. Shadowing can also enable more flexible latency control for an application, e.g., the device can use the earliest response from multiple shadowed weblats on the cloud.

Splitter pattern With the splitter pattern, a set of worker weblats perform variant implementations of a shared interface. For example, the workers may encapsulate adapters to access different social networks, or codecs to process different media formats. The application is decoupled from the various implementations by a splitter weblat that routes requests to appropriate workers. This pattern increases application extensibility since new worker weblats are added without changing the application structure. Splitting can also enhance the user experience by converging multiple services on a single device. For instance, in the case where the worker weblats access different social networks, the splitter weblat’s interface provides a unified or converged interface to a range of social networking services.

Aggregator pattern An elastic application can also aggregate computations from multiple worker weblats. In this pattern, an aggregator weblat collects information from multiple worker weblats and uses *weblat push* to relay this information to the device. For example, an application can run multiple weblats in the cloud as background threads that monitor the user’s web accounts (e.g., emails or instant messages), the aggregator weblat pushes events (such as account activity) to the device. In some cases the splitter and aggregator patterns are combined or overlaid, the splitter pushes requests to the workers while the aggregator pushes events back to the device.

Table 1 Weblats vs. web services

Weblats	Web services
HTTP (REST interface)	HTTP (REST or SOAP interface)
Single client	Many clients
Client is application root or other weblat	Clients are generally browsers or other web services
Short-lived & long-lived requests	Generally short-lived requests
Dynamic endpoints (may migrate)	Fixed endpoints
Lifetime is client dependent	Lifetime is client independent
Runs on servers or client (cloud or device)	Runs on servers
Push to client possible	Not available or non-standard

3 Cost optimization for elastic applications

3.1 Cost model

The augmented computation of an elastic application is not free but introduces costs to the mobile device and user, which depends on when and where a weblet is running and communications within weblets or between weblets and Internet. Furthermore, elastic applications can exhibit variant runtime behaviors with dynamic execution configurations, such as power consumption, monetary consumption, application performance, and even security and privacy properties. Therefore, the dynamic execution configuration of an elastic application is decided based on some cost saving objectives, which form a cost model in our framework. As Fig. 2 shows, the cost model takes inputs of sensor data from both device and cloud sides, and runs optimizing algorithms to decide execution configuration of applications. Device and cloud related data such as battery level, network conditions, device loads, cloud loads and other performance data including current latency of the application, are obtained from appropriate sensing modules. The output of the cost model is possible actions that lead to the optimal execution configuration for the application, such as allocating resources on the cloud, launching/migrating weblets on/to device and/or cloud, selecting/switching between different network interfaces, replicating and shadowing weblets on cloud, etc.

An important part of the cost model is choosing the attributes or objectives that should be optimized. We consider the following four attributes in our current elastic application framework, while we believe new cost objectives can be integrated easily.

Power consumption Each application/weblet running on a mobile device consumes battery power by using CPU cycles, memory and radio module for communi-

cation with peer weblets on the cloud and/or external web services. The power consumption of a weblet on the device heavily depends on the I/O operations it performs [4, 27]. In addition, different communication channels, such as W-CDMA, WiFi (802.11) etc., consume different power [2, 3, 5]. Considering the above, it is evident that although launching/migrating weblets to clouds should ideally save power consumption of computation on the device, the power consumption of network interfaces may override the benefits of the migration.

Monetary cost Execution of a weblet on a cloud platform may involve a monetary cost for the application user, based on the exact resources consumed on the platform. Usually, a commercial cloud service provider measures the cost of a computing task based on the amount of CPU cycles, storage, and communication traffic (in and out) of a cloud platform [1]. The monetary cost of a weblet running on the cloud platform is determined by the size of the input data consumed by the weblet (including those from peer weblets on the device for the same application and external web services), total execution time of the weblet on the cloud platform, data size/rate for intra-cloud communication between this weblet and others within the same cloud service provider (if applicable), and any other attributes that affect these parameters, such as network status affecting data transmission rate.

Performance attributes As an elastic application potentially runs across different platforms, latency is an important design consideration. There are different aspects of latency, such as impact on the user experience when using the application's UI and network latency with different network connections and traffic status, and the application latency to finish a particular computing task. Throughput also can be an important objective for some applications. For example, an application that does image analysis to find similar pictures from a large database needs maximum throughput. To achieve this, the heavy computing tasks are launched or migrated to the cloud, although there is a tradeoff between doing this and the data communication overhead: too much communication may slow down the overall application throughput. Given this, building a good performance model is more challenging than power and monetary aspects. In general, to optimize latency, throughput and some application-specific options, CPU cycles and memory used by the weblets, along with the available network bandwidth for communication between the device and the cloud should be carefully evaluated.

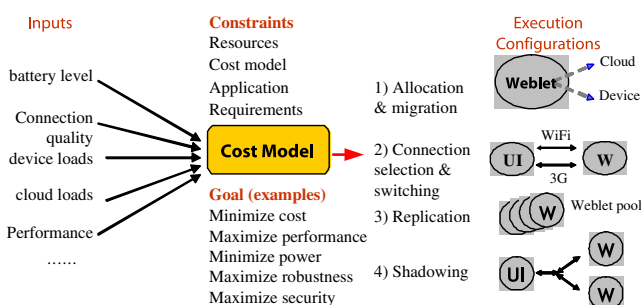


Fig. 2 Cost model of elastic applications

Security and privacy Security is increasingly concerned in web-based computing systems. A mobile device potentially contains many user secrets and privacy-sensitive data, such as: contacts, SIM information, credit card details and many other credentials that may be needed to consume web services. Naturally, a mobile user may trust her device more than the cloud platform which is controlled by a third-party service provider. As launching or migrating a weblet to the cloud may also require offloading user data to the cloud, the user security and privacy concerns are even higher with an elastic mobile device. A weblet on the device or the cloud may need to access external web services on behalf of the user. For cost modeling purposes, we need to evaluate if a weblet requires any user data and if the user has strong concerns about offloading such data to the cloud. If the user has concerns over doing this, the weblet that requires this data should be launched on the device only and never migrated. Furthermore, during runtime, if a weblet needs to acquire external user data from other web services, which usually requires user credentials (username/password, public key certificate, or any other security credentials), the weblet may have to be migrated back to the device.

3.2 Optimizing execution configuration

Once a cost model is developed for a particular application, a mechanism is needed for efficient and intelligent dynamic execution configuration, e.g., via some light-weight machine learning algorithms at the device side. In our implementation of one elastic application, we use Naive Bayesian Learning techniques to find the optimal weblet configuration (# of weblets on device and cloud), given device status (in terms of CPU, memory and network consumption), user preference (in terms of expected # of images that should be concurrently processed), and history data of the application.

As Fig. 3 shows, a vector ‘*x*’ consists of values representing device status components such as the upload bandwidth, throughput, power level, memory usage and file cache. A vector ‘*z*’ consists of values representing user’s preferred setting for cost objectives including monetary cost, power consumption, and processing speed. The configuration variable ‘*y*’ has values from 1 to *N* (max number of possible configurations), where each value maps to a specific configuration pair. Given all these data, the following expression can be applied to determine the most optimal configuration.

$$y^* = \underset{y}{\operatorname{argmax}} p(y) \prod_{i=1}^L p(x_i|y) \prod_{j=1}^M p(z_j|y) \quad (1)$$

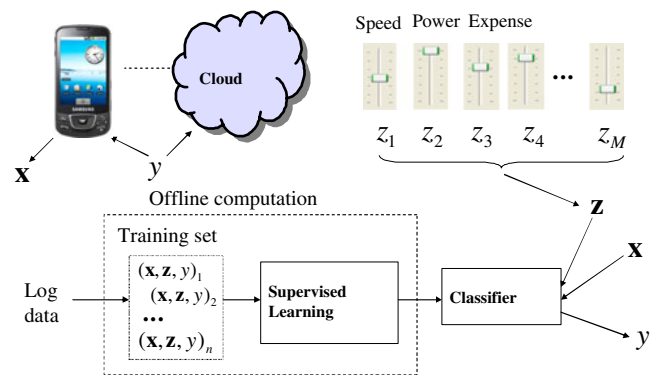


Fig. 3 Weblet scheduling through Machine Learning techniques

In the above expression, x_i is the i -th status component value that can have different number of states for each component and z_j is a j -th preference component, where $i \in \{1, 2, \dots, L\}$ and $j \in \{1, 2, \dots, M\}$, with L and M representing the number of components in the status vector and the number of components in the preference vector, respectively.

Note that it is relatively easy to determine dynamic configurations in this application since it has only one type of weblet. For a general application with multiple types of weblets, each having different runtime behaviors, the optimization can be very complex and the computation itself may override the cost savings. Considering that an elastic application can be installed and executed by many users on similar devices, a service-oriented cost optimization implementation can save computation cost for the device.

4 Reference implementation

4.1 Reference architecture

To experiment with this new application model, we have developed a reference framework including application bundle, architecture, and some example elastic applications. Our framework works with Amazon EC2 and S3. Figure 4 shows the main functional components.

In our current framework design, a typical elastic application consists of a UI component, one or more weblets, and a manifest. Weblets are autonomous software entities that run either on the device or cloud and expose RESTful web service interfaces via HTTP. The manifest is a static XML file that contains metadata for the application. It could be used to specify any requirements and constraints for the application and the individual weblets, such as: the digital signature

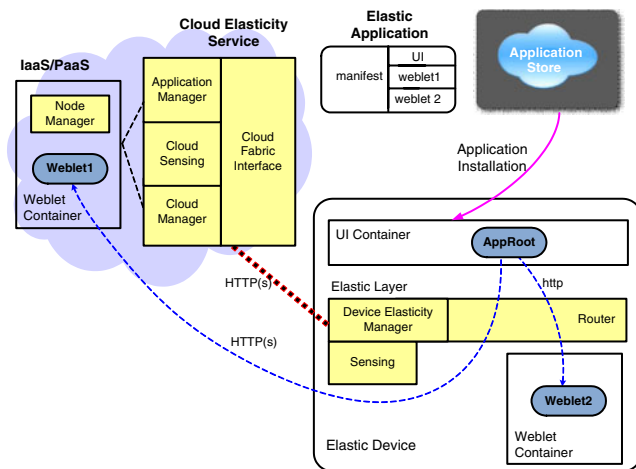


Fig. 4 Reference architecture for elastic application

needed to download/migrate the weblets, requirements for compute power, network and storage, time limits for weblet execution, maximum instances of the weblet that can be launched on the device and the cloud, if a weblet can be launched/migrated to the cloud and specifics about handling data required/generated by the application/weblets etc.

On the device side, the key component is the device elasticity manager (DEM) which is responsible for configuring applications at launch time and making configuration changes during run time. The configuration of an application includes: where the application's components (weblets) are located, whether or not components are replicated or shadowed (e.g., for reliability purposes), and the selection of paths used for communication with weblets (e.g., WiFi or 3G if such a choice exists). The router passes requests from UI components to weblets. It insulates the UI logic from weblet location. When a weblet is migrated, the router will be aware of the new location and will continue passing requests from the UI to the weblet (and passing replies back to the UI). Each device also provides sensing data from the device such as processor type, utilization, and battery state. This data is made available to the elasticity manager and is used to determine when and where a new weblet instance should be launched.

The cloud elasticity service (CES) consists of the cloud manager, application manager, and sensing information collection. The cloud manager is responsible for allocating resources from, and releasing to, underlying cloud nodes. It maintains usage information, including compute, bandwidth and storage, for the various weblets running on the cloud. The application manager provides functions to install and maintain applications

on behalf of elastic devices, and helps launch weblets on different cloud nodes. Sensing information refers to the collection of operational data on the cloud platform. These data are made available to the cloud manager to assist it in tracking usage. In addition to application performance data, sensing may collect information on cloud architecture, failures of various forms, and resource availability. As a service provider, the CES exports a web service, referred to as the cloud fabric interface (CFI) to elastic devices and applications. A node manager on each cloud node oversees resources associated with a particular node (server) within the cloud. It communicates directly with the cloud manager and application manager. Each node runs one or more weblet containers which are the weblet runtime environments hosted on an Amazon EC2 instance.

4.2 SDK development

We have implemented a preliminary SDK based on the reference architecture, which is used to develop the basic interfaces of weblets in our example applications. Using this SDK as a base, developers can build elastic applications in high-level languages such as JavaScript, Java, and C#. Currently the SDK has C# bindings; however it is easy to extend it to other languages.

A typical elastic application includes a `AppRoot` component and one or more weblets. The `AppRoot` is the part of the application that provides the user interface and issues requests to weblets. All of these are packaged into one bundle, which includes the binaries of weblets and a manifest describing the application, and most importantly, the developer-signed hash values of the individual weblets. Figure 5 shows a state diagram illustrating the lifecycle of a weblet, including the various states that a weblet can be in and the actions that cause the state transitions. A weblet is an independent functional unit of an application that is capable of compute, storage, and networking tasks. It resembles an embedded or dedicated web server and presents a web service interface (i.e., it is accessed via HTTP). In our SDK, an abstract class called `AbstractWeblet` is defined to represent the core behavior of weblets. Other specific types of weblets can be implemented as subclasses of `AbstractWeblet` and extend its methods as required. Each weblet is associated with a weblet type and identified through a unique id. Once an application has defined one or more weblet types, it can use the DEM to create instances (i.e., to create specific weblets) and issues requests to these weblets. We describe below the core interfaces supported by DEM and illustrate how to create and communicate with a weblet.

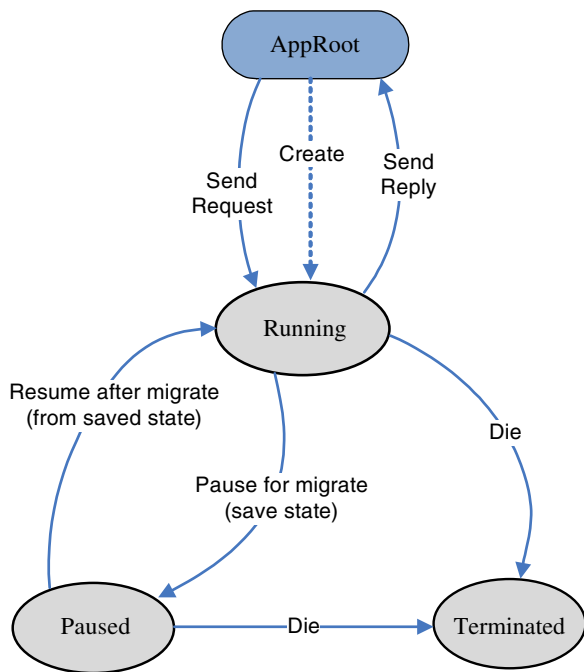


Fig. 5 Lifecycle of a weblet. A weblet is always created by the AppRoot, and can be in state of Running, Paused, or Terminated

Create a weblet In order to create a weblet, two parameters are required including: the type of the weblet and a call back function to invoke when the weblet is created. For example, the following code creates a weblet of type

```

MyWebletType:
    CreateWeblet("MyWebletType",
        OnCreateMyWeblet);
  
```

Send a request to a weblet A request can be sent to a weblet using the SendWebletRequest interface. The interface requires the following parameters: the weblet to which the request is to be sent, the actual request (query string) to be sent, and a callback function that is invoked when a reply is received. Below is sample code to achieve this:

```

SendWebletRequest(w, "RequestQuery",
    OnMyWebletReply);
  
```

Receive a response from a weblet When the DEM receives a response from a weblet it invokes the callback function indicated by the requestor along with request. With this, the requestor only needs to implement the call back function to receive the reply. For example, in continuation with the example given above, the DEM would invoke the following function to convey the reply to the application:

```

OnMyWebletReply(w, reply);
  
```

The DEM can decide to migrate a running weblet from the device to the cloud or vice-versa; weblet migration is transparent to the application. When a weblet is running on device and the DEM decides to migrate it to a cloud node, the DEM issues a Pause request to the weblet, this causes the weblet to close its request interface, release resources and save state. The DEM then sends the saved state to the cloud via the CFI. After the state has been transferred to the cloud, the weblet is resumed and restores itself from the saved state. The CFI returns the new connection information for the weblet (e.g., IP address, port, and session tokens) to the DEM so that the DEM may continue to route requests to the weblet on cloud.

5 Elastic applications

To demonstrate the elastic application model, we have developed several applications with our SDK and deployed on the reference architecture with Amazon EC2. This section explains the design and elasticity patterns of these sample applications, aiming to provide the abstract principle of designing elastic applications in general.

5.1 Elastic image processing

The simplest is an image processing application in which various filtering operations are applied to set of images. Following the replication pattern, a weblet pool is created on the cloud; images are then processed in parallel by pool members. The application can adjust the size of the pool, so it is possible to compare throughputs for different execution configurations. The application supports 3 workloads (number of images to process). Load 1 has 1 image, load 2 has 4 images, and load 3 has 16 images. The images are 24-bit color images of size 240 × 360. Figure 6 shows a snapshot when the application is running on a Samsung Galaxy smart phone with Android 1.6.

The application implements many of the elastic application concepts described in the previous sections and enables the user to do the following through the UI: (1) Specify if the device is online (i.e., if it should make use of the cloud to run weblets); (2) specify the number of weblets to run on the cloud; (3) specify the kind of filtering to be done on the images; and (4) specify the load or the number of images to process concurrently.

This image processing application consists of only one type of weblet (ImageWeblet), as shown in Fig. 7. Its functionality is to perform image filtering as specified by the user. The weblet is replicated on the

Fig. 6 Snapshot of elastic image processing application



device and the cloud, as and when required. The total number of weblet instances spawned depends on the load specified and the user specified value for number of weblets in the cloud.

5.2 Elastic augmented reality: object identification and replacement

A second example is a form of augmented reality (AR) in which real-world objects are detected and enhanced. This application runs tracking and rendering on the device and uses the splitter pattern with a set of matcher weblets on the cloud. Each matcher searches for different objects within video frames. The splitter collects information on identified objects and relays this to the device for rendering. By running the matchers in the cloud, many more objects can be detected (per unit time) than when the application runs fully on the device.

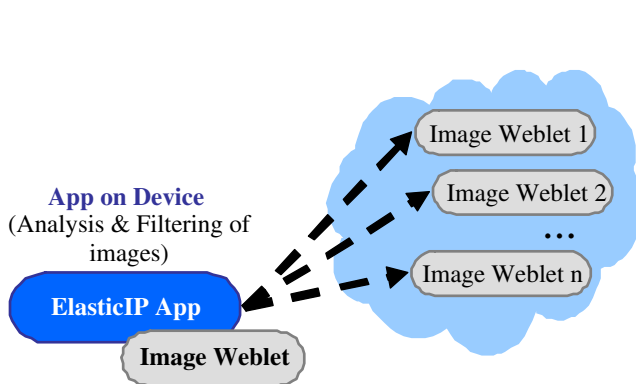


Fig. 7 Weblet topology of elastic image processing



Fig. 8 Snapshot of elastic AR: object identification and replacement

A snapshot of this application is shown in Fig. 8. On a high level, the application works as follows. First, the image/video frame from the mobile device camera is captured, and the feature points for the image in the frame are extracted. With these feature points, planar objects are then identified, by matching the extracted feature points to those of known planar objects in a database. The recognized planar object is then used to select a replacement image in the user’s choice of language. Finally, the replacement image is provided to the device and is overlaid on top of the current image. For this scenario, the users are able to specify the language of choice for the user and the number of weblets to run on the cloud through the UI. This augmented reality application would consist of three types of weblets, as illustrated in Fig. 9: The first type of weblet (*Tracker*) performs feature extraction on the live feed; the second type of weblet (*Matcher*) matches extracted features with images in a database. As this is a compute-intensive task, multiple instances of this weblet type are spawned on the cloud; and

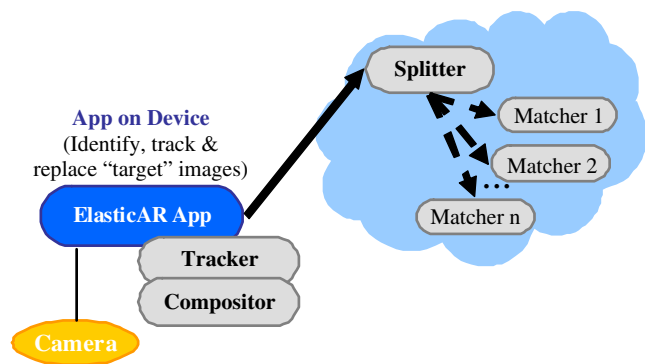


Fig. 9 Weblet topology of elastic AR: object identification and replacement

the third type of weblet (*Compositor*) performs the image replacements on the device.

5.3 Elastic augmented reality: augmented video

Our third sample application is another elastic AR application which enables an augmented video scenario, where in a user visiting a building (e.g., a technology expo or a museum) is able to seamlessly access information about points of interest (POIs)/demos in the area. When the user holds up his phone camera, information about points of interest in the area (title, description of demos and the number of people in the demo area) is overlaid on the camera view, represented using different types of icons and labels.

Figure 10 shows the snapshot of the application when running on Samsung Galaxy. When running, the application obtains the accelerometer, compass and GPS readings from the phone thus the current position and direction of the user are identified. The application then obtains a list of POIs/demos in the target area by submitting its location information to a web service interface running on cloud. The web service also collects the number of people near each POI, which would be done periodically by monitoring the location/coordinates of people in the target area through their devices; After these, the POIs near the user’s position and within the camera’s field of view are determined, and the POIs and information about the number of people in their proximity are overlaid on the user’s screen.

Figure 11 shows the weblet topology of this application, which consists of four types of weblets. The first type of weblet (*UserTracker*) identifies the po-



Fig. 10 Snapshot of elastic AR: augmented video

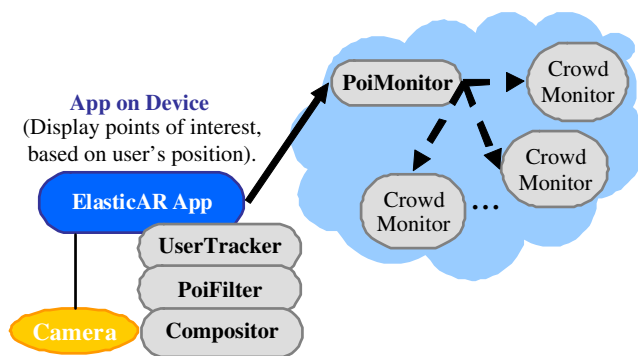


Fig. 11 Weblet topology of elastic AR: augmented video

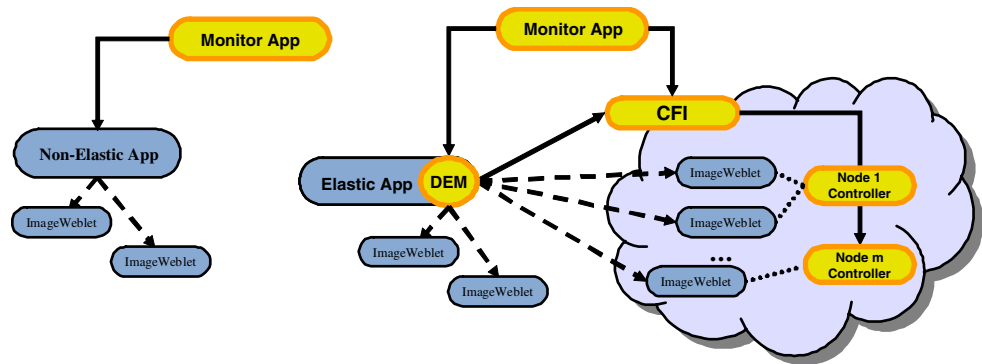
sition and direction of the user. The second type of weblet (*PoiMonitor* and *CrowdMonitor*) collects device information and calculates the number of people near the POIs. Multiple instances of this weblet type are spawned on the cloud. The third type of weblet (*PoiFilter*) determines the POIs in the user’s vicinity and field of camera view. The fourth type of weblet (*Compositor*) overlays the information about the POIs on the device’s screen.

6 Experimental validation

We validate the elasticity of our framework by using the aforementioned image processing application as benchmark. This application consists of only one type of weblet called *ImageWeblet*. Its functionality is to perform image filtering with an algorithm specified by the user. The weblet is replicated on the device and the cloud, as and when required. The total number of weblet instances spawned depends on application load and the number of weblets in the cloud, both specified by the user. The application UI enables the user to do the following configurations during runtime: online (can launch weblet at cloud) or offline (all weblets are running on the device) mode of the application, number of weblets to run on the cloud (if in online mode), the filtering algorithm to be used, and the number of images (workload) to process at the same time. The images used in by the experiment are 24-bit color with size 240×360 .

The goal of our validation is to compare the performance of an elastic device (ED) and a non-elastic device (NED) running the same image processing application. Figure 12 shows an overview of the demo configuration and system setup. For the elastic device, the application uses an in-house cloud comprising of 8

Fig. 12 Experiment configuration for elastic image processing application



Linux boxes. A non-elastic version of the application is also running independently in order to compare it with the elastic version. Essentially, the non-elastic version uses only the device to run weblets, whereas the elastic version uses both the device and the in-house cloud. The setup also includes PCs to host the CFI and a performance monitor application. The CFI is implemented with PHP scripts on a Linux server with Apache and MySQL.

The performance monitor collects several measurements, including the available upload/download bandwidth (KB/s), application workload (number of images to be processed) and throughput (the number of image tiles processed/s), average CPU usage (%), and available memory (MB), from the test device and from the cloud. In addition, it also maintains information about the total number of weblets started for the application and the individual number of weblets running on the device and the cloud.

Each configuration has a unique composition of device weblets and cloud weblets. We set the maximum number of weblets as 16 and consequently, more than 100 different configurations are possible. The configuration specifying 1 device weblet & 0 cloud weblets is considered the default configuration for the non-elastic device. Among all possible configurations, we chose the 74 configurations where the number of device weblets is less than or equal to 4 (due to limitations with CPU utilization) for analysis. For each configuration, the data was collected 20 times and the average values were considered for final comparisons.

Figure 13 shows the performance of the elastic device over 74 configurations. In comparison with the throughput of about 6 tiles/s for the default/non-elastic device configuration (1 device weblet, 0 cloud weblets), the throughputs of all other configurations are better. We can observe that the throughput for the configuration with 0 device weblet and 16 cloud weblets

has the highest throughput among the 74 configurations tested. The configuration with 16 device weblets has the best performance, as there are a total of 16 images in load 3. A surprising observation is that the configuration with 8 weblets performed better than configurations with 9–15 weblets (a result of internal application logic). This indicates that an intelligent weblet scheduling is essential to identify the most efficient weblet configuration.

CPU usage is more predictable overall, in that more device weblets lead to more CPU usage. However, the trend is interesting when comparing the number of device weblets. For configurations with up to 2 device weblets, running more cloud weblets leads to more CPU usage. For configurations with 3 and 4 device weblets, a general trend is that running more cloud weblets reduces the CPU usage. By combining CPU usage data in Fig. 14 with the throughput data in Fig. 13, we are able to identify the configurations that lead to low CPU usage and high throughput: for instance, configurations (0,2), (0,3) and (0,4) have lower CPU usage (than that of a non-elastic device) and higher

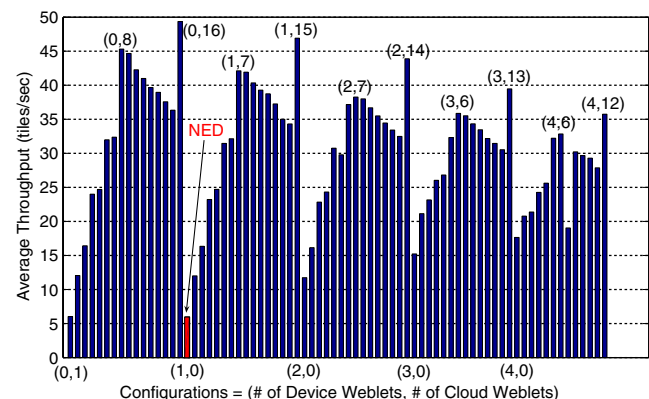


Fig. 13 Throughputs vs. configurations

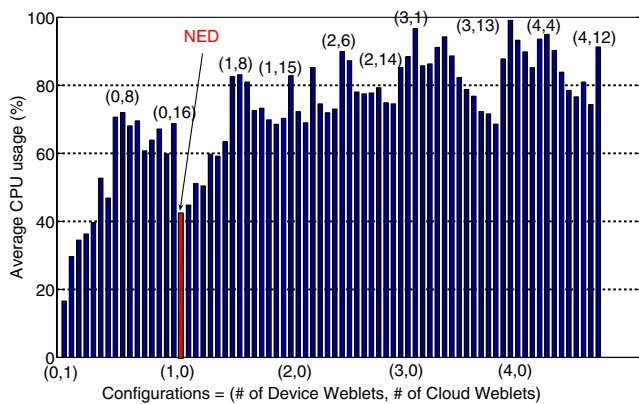


Fig. 14 CPU usage vs. configurations

throughput. This results in more available CPU cycles for other applications and improves multi-tasking capabilities.

Figure 15 shows interesting but not easily comprehensible results regarding the available memory versus configurations. Certain configuration such as (0,1), (0,2), (0,3), (1,0), (1,1), (1,2), (2,0), (2,1), (2,2), (3,0), (3,1), and (4,0) have much available memory. Most of the configurations have only up to 4 total weblets and using only the device weblets consumed only a little memory up to 4 device weblets. Meanwhile, other configurations up to (3,13) have similar available memory, but there is significant variation between (4,1) and (4,12). It is not clear why the system behaves that way, but it could be related to cache operations and memory paginations. This will need further investigations. Combining this result with Fig. 13 can lead to a memory-constrained optimal configuration. Of course, it would also be conceivable to find a good configuration constrained on both memory and CPU.

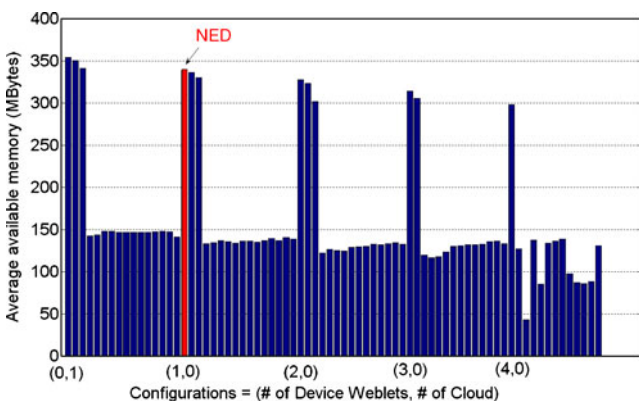


Fig. 15 Available memory vs. configurations

7 Related work

7.1 Execution offloading

The elastic model builds on previous work in the areas of remote execution and application offloading. Cyber foraging [7, 8, 18, 25] is a common approach explored by many to augment the capability of resource-constrained mobile devices. The basic idea is to dynamically discover and make use of nearby resources, aka surrogates, to offload the execution of an application or parts of an application running on a mobile device. Compared to these approaches, elastic application model has more flexible deployment patterns to parallelize tasks on multiple remote cores.

CloneCloud takes the approach of cloning the entire user’s mobile device environment on a remote server. Applications can then be quickly restarted on or migrated to the remote machine when the user’s machine is running low on resources [9]. Similar virtual machine-based approach is used by cloudlet [23]. As mentioned in Section 1, our elastic application model offloads computing tasks in more fine-grained level such that it leverages the parallel computing advantage of cloud resources.

Some research work extend existing programming language and application runtime middleware to transform applications into distributed systems [11, 14, 21]. Adaptive Offloading [11] leverages Java’s object oriented design to identify possible partitions for a Java application and modified the JVM to support such partitioning. Coign [14] makes use of the location transparency supported by COM and converts an application built from COM components into a distributable application. R-OSGi [21] extends the centralized module management functionality supported by the OSGi specification to enable an OSGi application to be transparently distributed across multiple machines. The main limitation with these approaches is that they are tied to one particular language or specification and hence not suitable for a wide range of applications. Compared to these approaches, our proposed elastic application model is programming language independent, and can be extended to many existing application middleware.

Virtual machine migration [19, 26] and VM-based cloudlet [23] are complementary approaches to enable users to seamlessly access their applications and data across multiple and heterogeneous devices in general. It also enables users to instantly continue/restore an application on a different device, when their current machine is running low on resources.

7.2 Cost-aware configuration

There have been literatures dealing with configuration methods based on resource estimation. In [17, 24], a resource-aware configuration method decides the configuration based on user's quality of service requirement, resource and service availability, and application fidelity as a function of resources. It tries to maximize a product-based utility function so that the aggregate resource demand cannot exceed the resource supply. A machine learning approach is introduced in [32] to capture the complex nonlinear relationship between resource properties and computing power. It provides resource selection for a job in Grid scheduling to have the maximum utility of CPU. However, these works are not specifically targeted for remote execution on mobile devices. The tactics-based remote execution in [8] aims to select the best tactic, the useful knowledge about an application relevant to remote execution, using resource prediction and resource monitoring. It tries to maximize the latency-fidelity metric in the tactic selection. In [10], a similar work is proposed, using a product-based decision criterion for remote execution. It considers only three metrics of execution time, energy usage, and fidelity.

Narayanan et al. [20] use historical application logging data to predicate the fidelity of an application, which decides its resource consumption. However, in this work, only aspects of device hardware and application inputs are considered. In our cost model, we consider more comprehensive factors not only on device side, but also on cloud side. Uniquely, we incorporate user preferences in terms of cost objectives. Gurun et al. [13] extend the network weather service (NWS) toolkit in grid computing to predict offloading, which can be leveraged as an implementation mechanism for our cost model.

In [12], a Fuzzy Control model-based offloading inference engine is introduced to solve when to trigger adaptive offloading and how to partition an application. However, the decision criterion is based only upon the memory, not considering multiple factors. Our approach provides an optimized elasticity by considering multiple factors as costs, including device status, Cloud status/usage, application performance measures, and user preference to the cost factors.

Comparing with these approaches our work is based on the assumption that cloud has huge resources, thus releasing the resource estimation requirement in the cloud from the decision of weblet scheduling. There are some literatures regarding cost analysis on the cloud side-only. The tradeoffs between the cloud computing

and desktop grids are provided in [15]. The total cost of ownership and utilization cost is introduced in [16] to evaluate the economic efficiency of the Cloud. A workload balancing approach [30] is proposed between public Cloud and private Cloud for cost-saving. In [28], the monetary cost of leasing CPU time from commercial Clouds is compared with that of purchasing and using a server cluster of equivalent capability.

8 Conclusions and future work

We propose an elastic application programming model aiming to remove the constraints of specific mobile platforms by providing a distributed framework that extends the device into the cloud. The salient feature of this model is that it offers a range of elasticity patterns between resource constrained devices and Internet-based clouds. Each pattern in turn can be realized by several execution configurations. A comprehensive cost model is used to dynamically adjust execution configurations thus optimizing application performance in terms of a set of objectives. We present the high level design of elasticity framework and primitive experimental results with an example application.

8.1 Future work

There are aspects of this work that need further research efforts. We highlight some of them at the end of this paper.

Data and state synchronization As aforementioned in the elasticity patterns, weblets of a single application may share application data and state. For example, different weblets may require the same data from the device for their input, or they may update the same data during runtime. Since weblets run in different locations, it is desirable to replicate data to increase performance, but then data integrity and synchronization become issues. Alternatively, data synchronization can be explicitly performance by applications, or implicitly by framework architecture and transparent to applications. In the first case, an elastic application handles its own data management including storage and synchronization between device and cloud nodes. The advantage is flexibility: a user or application developer can select the data storage mechanism on the cloud. However, this leaves data access handling to developers, and the user may need to manually initiate synchronization during runtime. In the architecture-based approach, application data are duplicated and

synchronized by the elasticity architecture, such that the applications are not aware of data location. APIs can be defined to access (read and write) data via middleware, which hide the details of data management including synchronization and backup. This releases the burden of data management from application developers, while heterogeneous data storage mechanisms at device and cloud side give challenges to middleware design.

Communication between weblets In our reference architecture, weblet requests are initiated on the device side and can propagate to the cloud to be passed from one cloud weblet to another (as in the *splitter pattern*). To support more flexible elasticity patterns, a mechanism is needed to allow a cloud-residing weblet to invoke requests of device weblets. This problem becomes challenging when the device is mobile, it may switch between different network channels, e.g., between WiFi or 3G network, or even between different wireless network providers, and it may be running behind a firewall or using NAT. Communication beyond organization boundaries is another challenging issue to be solved. Increasingly smartphones run enterprise applications and connect to intranet servers. Typically VPN software is installed on these devices. How to enable secure and flexible communications between weblets on enterprise-owned mobile devices and cloud servers needs further research efforts.

Media channel between weblets Although HTTP is light-weight and flexible, and is used in our reference architecture and example applications, it is not a good option for media streaming or distributed visual processing between device and cloud. Tasks requiring significant processing or data storage, such as visual object recognition or rendering complex 3D models, can be performed in the cloud rather than on the device. However for cloud computing to be used in highly interactive and visually rich applications, there is a need for a high-speed and low-latency transfer method for structured visual data between the device and cloud. For example, consider a 3D game where some scene elements are rendered on the cloud and then sent, along with camera and depth information, to the device for mixing with locally rendered elements. Early work in this direction includes frame-oriented 2D graphics protocols (e.g., RDP, RFB, VNC), protocols for remote rendering of 3D graphics (e.g., X11 extensions for OpenGL) and protocols for encoding segmented video (MPEG4). Generally these protocols involve a generic decoder, i.e., no application-specific logic is required for decode and display. For situations

where application logic is split between the device and the cloud, and visual processing takes place on both sides, new protocols are needed to exchange partially rendered and partially processed data.

Weblet migration Code and computation migration is a traditional problem in many systems [8, 29]. To enhance the mobile user experience, our model supports migration without the need to migrate code. During installation of elastic applications, the code used by weblets is installed on both the device and the cloud. When a weblet is required to migrate from device to cloud, a new weblet instance is allocated on a cloud node, and the runtime state is copied from the device weblet to this new weblet. We believe this state migration is more efficient than migrating a weblet's memory image and state information. To support this type of migration, a weblet is not migrated when in an arbitrary state. Instead the weblet closes any pending requests and then saves state information in preparation for transfer. After migration the weblet loads the saved state and resumes its operation. With this approach, the specification and representation of a weblet state are critical. Basically, the state information should include the current task status, its working data, and handles to any other weblets with which it communicates. The state should also ensure that the physical location of the new weblet does not affect existing communication channels between other weblets and external parties. For this purpose, a routing-like mechanism should be provided by the architecture and supported by the middleware. A weblet can then have some well-known name for use by the application, while the binding between the name and a physical weblet entry point (e.g., a URL) is dynamic.

Trust and security The elastic application model and middleware should provide a mechanism to authenticate weblets belonging to a single application. Authentication is the prerequisite to building secure communication between weblets. Also, session management is essential, especially weblet behaviors at cloud side should be accounted, e.g., to give the mobile user the resource usage and cost of the application. In our reference architecture, we have designed a lightweight protocol to distribute shared secrets and session keys between weblets for mutual authentication purposes [31]. Beyond this, there are some challenging problems for elastic applications. First of all, a mobile user needs trust to launch weblets on a public cloud, especially when the computation and network traffic incur monetary bills to the user. This demands that the computing environments in the cloud should be verifiable by a user

or a trusted party, e.g., to ensure there is no hidden or even malicious code running beside weblets. Similarly, the quality of service from cloud providers should be verifiable. Furthermore, a mobile user should be assured that the weblets running in the cloud are the ones that she has installed and their integrity can be verified via trusted mechanisms. We believe that extending the trusted computing base (TCB) of the mobile device to some necessary but minimum cloud service is necessary to satisfy these security requirements [22].

References

1. Amazon EC2. <http://aws.amazon.com/ec2/>. Accessed 21 Mar 2011
2. RFMD product data sheet. <http://www.rfmd.com/products>. Accessed 21 Mar 2011
3. Wifi power consumption analysis. <http://nesl.ee.ucla.edu/fw/documents/reports/2007/poweranalysis.pdf>. Accessed 21 Mar 2011
4. Samsung flash memory databook. <http://www.samsung.com/global/business/semiconductor/packageInfoList.do>. Accessed 21 Mar 2011
5. Analog devices data sheet. http://www.analog.com/en/content/psearch_landing/fca.html. Accessed 21 Mar 2011
6. Armbrust M, Fox A, Griffith R, Joseph AD, Katz RH, Konwinski A, Lee G, Patterson DA, Rabkin A, Stoica I, Zaharia M (2009) Above the clouds: a berkeley view of cloud computing. Technical Report UCB/ECS-2009-28. EECS Department, University of California, Berkeley
7. Balan R, Flinn J, Satyanarayanan M, Sinnamohideen S, Yang H (2002) The case for cyber foraging. In: Proc. of the 10th ACM SIGOPS European workshop
8. Balan RK, Satyanarayanan M, Park S, Okoshi T (2003) Tactics-based remote execution for mobile computing. In: Proc. of the 1st international conference on mobile systems, applications, and services, pp 273–286
9. Chun B-G, Maniatis P (2009) Augmented smartphone applications through clone cloud execution. In: USENIX HotOS XII
10. Flinn J, Park S, Satyanarayanan M (2002) Balancing performance, energy, and quality in pervasive computing. In: Proc. of the international conference on distributed computing systems
11. Gu X, Messer A, Greenberg I, Milojicic D, Nahrstedt K (2004) Adaptive offloading for pervasive computing. *IEEE Pervasive Computing* 3(3):66–73
12. Gu X, Nahrstedt K, Messer A, Greenberg I, Milojicic D (2003) Adaptive offloading inference for delivering applications in pervasive computing environments. In: Proc. of IEEE international conference on pervasive computing and communications
13. Gurun S, Krintz C, Wolski R (2004) Nwslite: a light-weight prediction utility for mobile devices. In: Proc. of international conference on mobile systems, applications, and services
14. Hunt GC, Scott ML, Hunt GC, Scott ML (1999) The coign automatic distributed partitioning system. In: Proc. of the 3rd symposium on operating systems design and implementation, pp 187–200
15. Kondo D, Javadi B, Malecot P, Cappello F, Anderson DP (2009) Cost-benefit analysis of cloud computing versus desktop grids. in: proc. of the ieee international symposium on parallel & distributed processing
16. Li X, Li Y, Liu T, Qiu J, Wang F (2009) The method and tool of cost analysis for cloud computing. In: Proc. of IEEE international conference on cloud computing
17. Poladian V, Sousa JP, Garlan D, Shaw M (2004) Dynamic configuration of resource-aware services. In: Proc. of international conference on software engineering
18. Porras ORJ, Kristensen MD (2008) Dynamic resource management and cyber foraging. Chapter *Middleware for Network Eccentric and Mobile Applications*. Springer Press
19. Kozuch M, Satyanarayanan M (2002) Internet suspend/resume. In: Proc. of the 4th IEEE workshop on mobile computing systems and applications
20. Narayanan D, Flinn J, Satyanarayanan M (2000) Using history to improve mobile application adaptation. In: Proc. of the 3rd IEEE workshop on mobile computing systems and applications
21. Rellermeyer JS, Alonso G, Roscoe T (2007) R-osgi: distributed applications through software modularization. In: Proc. of the ACM/IFIP/USENIX international conference on middleware
22. Santos N, Gummadi KP, Rodrigues R (2009) Towards trusted cloud computing. In: Proc. of the workshop on hot topics in cloud computing
23. Satyanarayanan M, Bahl P, Caceres R, Davies N (2009) The case for vm-based cloudlets in mobile computing. *IEEE Pervasive Computing* 8(4):14–23
24. Sousa JP, Balan RK, Poladian V, Garlan D, Satyanarayanan M (2008) User guidance of resource-adaptive systems. In: Proc. of international conference on software and data technologies
25. Sousa J, Garlan D (2002) Aura: an architectural framework for user mobility in ubiquitous computing environments. In: Proc. of the 3rd Working IEEE/IFIP conference on software architecture
26. Travostino F (2006) Seamless live migration of virtual machines over the man/wan. In: Proc. of the ACM/IEEE conference on Supercomputing
27. Vijaykrishnan N, Kandemir M, Irwin M, Kim H, Ye W (2000) Energy-driven integrated hardware-software optimizations using simplepower. In: Proc. of the int. symposium on computer architecture
28. Walker E (2009) The real cost of a CPU hour. *IEEE Comput* 42(4):35–41
29. Xian C, Lu YH, Li Z (2007) Adaptive computation offloading for energy conservation on battery-powered systems. In: ICPADS
30. Zhang H, Jiang G, Yoshihira K, Chen H, Saxena A (2009) Intelligent workload factoring for a hybrid Cloud computing model. In: Proc. of the congress on services
31. Zhang X, Schiffman J, Gibbs S, Kunjithapatham A, Jeong S (2009) Securing elastic applications on mobile devices for cloud computing. In: Proc. of ACM cloud computing security workshop
32. Zhao G, Shen Z, Miao C, Wan C (2009) ELM-based intelligent resource selection for Grid scheduling. In: Proc. of IEEE international conference on machine learning and applications