

# Evidence-aware Mobile Computational Offloading

Huber Flores, Pan Hui, Petteri Nurmi, Eemil Lagerspetz, Sasu Tarkoma, Jukka Manner, Vassilis Kostakos, Yong Li, and Xiang Su

**Abstract**—Computational offloading can improve user experience of mobile apps through improved responsiveness and reduced energy footprint. A fundamental challenge in offloading is to distinguish situations where offloading is beneficial from those where it is counterproductive. Currently, offloading decisions are predominantly based on profiling performed on *individual devices*. While significant gains have been shown in benchmarks, these gains rarely translate to real-world use due to the complexity of contexts and parameters that affect offloading. We contribute by proposing *crowdsensed evidence traces* as a novel mechanism for improving the performance of offloading systems. Instead of limiting to profiling individual devices, crowdsensing enables characterising execution contexts across a community of users, providing better generalisation and coverage of contexts. We demonstrate the feasibility of using crowdsensing to characterize offloading contexts through an analysis of two crowdsensing datasets. Motivated by our results, we present the design and development of EMCO toolkit and platform as a novel solution for computational offloading. Experiments carried out on a testbed deployment in Amazon EC2 Ireland demonstrate that EMCO can consistently accelerate app execution while at the same time reduce energy footprint. We also demonstrate that EMCO provides better scalability than current cloud platforms, being able to serve a larger number of clients without variations in performance. Our framework, use cases, and tools are available as open source from GitHub.

**Index Terms**—Mobile Cloud Computing, Code Offloading, Cloud Offload, Big Data, Crowdsensing, Mobile Context Modeling.



## 1 INTRODUCTION

Computational offloading is a powerful mechanism for improving responsiveness and battery efficiency, two critical factors in the acceptability and usability of mobile applications [1]. By outsourcing (parts of) computationally intensive tasks to dedicated computing infrastructure, devices can reduce burden on their own resources while benefiting from resources provided by the dedicated infrastructure. Several platforms for providing offloading support have been proposed in the literature with the main difference being the level at which offloading is enabled. Examples of proposed mechanisms include annotation [2], [3], thread-migration [4], thread-migration through data shared memory (DSM) [5], and various others [6], [7], [8].

A key challenge in computational offloading is to distinguish situations where offloading is beneficial from those where it is counterproductive. For example, offloading network intensive computing tasks is only beneficial when the network link is sufficiently fast and stable whereas offloading parts of applications that require interactivity (e.g.,

augmented reality) is often counterproductive as energy gains are offset by decreased responsiveness resulting from increased latency [9]. Making optimal decisions, however, is far from straightforward as a wide range of factors can influence the effectiveness of offloading. For example, the extent of energy savings from reduced CPU usage depend on temperature and the number of applications running on a device [10]. Similarly, network type, device model, and the capacity and current load of the cloud surrogate have a major impact on offloading performance [11], [12]. Quantifying and characterizing the effect of these factors on a *single* device is clearly infeasible due to the inherent complexity and diversity of these factors. Indeed, in our experiments we demonstrate that quantifying individual factors, such as characterizing whether a particular mobile application requires heavy processing or characterizing the effect of network type, requires a large amount of samples (data from over 100 different apps or over 6,000 RTT samples). While individual devices cannot capture such amounts easily, harnessing a *community* of devices would make the process feasible. A population of 100 devices, which is easily obtainable within a university or even a company, would be capable of collecting the required measurements within a week or less, and a community of 1000 devices, which is increasingly possible through app store deployment [13], could collect the required measurements in less than a day.

We contribute by developing *evidence-aware mobile computational offloading* (EMCO) as a novel mechanism for optimizing the performance of computational offloading frameworks. The key novelty in EMCO is the use of *crowdsensed evidence traces* to characterize the influence of different contextual parameters and other factors on

- H. Flores is with the Department of Computer Science, University of Helsinki, Finland.  
E-mail: see <http://huberflores.com>
- P. Hui is with the University of Helsinki, Finland and the Hong Kong University of Science and Technology, Hong Kong.
- P. Nurmi, E. Lagerspetz and S. Tarkoma are with the University of Helsinki, Finland.
- V. Kostakos is with the University of Melbourne, Australia.
- J. Manner is with Aalto University, Finland.
- Y. Li is with Tsinghua University, China.
- X. Su is with the University of Oulu, Finland.

Manuscript received November xx, 201x; revised Month xx, 201x.

offloading decisions [14]. Contrary to existing solutions, which either rely on static code profiling performed on individual devices [2], [3], [4], [5], [6], [11], [15], [16], [17] or on parametrized models that consider a handful of parameters such as network latency, remaining energy, and CPU speed [9], [18], the use of crowdsensing enables EMCO to quantify and characterize the effect of a wide range of parameters and how they vary over execution contexts. EMCO models the context where offloading decisions are made is through simple dimensions that are easy to scale, and determines optimal dimensions using an analytic process that characterizes the performance of offloading based on contexts captured by the community. EMCO also supports identifying reusable and profitable code invocations that can be cached locally to further accelerate performance. We evaluate EMCO through extensive benchmarks conducted on an Amazon EC2 testbed deployment. Our results demonstrate that EMCO significantly accelerates mobile app performance, resulting in over five-fold improvements in response time and over 30% reductions in energy. When pre-caching on the client is used, the magnitude of savings is even higher with up to 30-fold improvements in response rate being possible.

The contributions of the paper are summarized as follows:

- We propose *crowdsensed evidence traces* as a novel mechanism for optimizing offloading performance. By aggregating samples from a larger community of devices, it is possible to learn optimal execution contexts where to offload. We analyze two crowdsensing datasets, Carat [19] and NetRadar [20], to demonstrate that it is possible to use crowdsensing to accurately quantify the effect of different factors on offloading performance. We also demonstrate that carrying out the profiling on a single device is not feasible due to the number of samples required.
- We present the design and development of the Evidence-Aware Mobile Computational Offloading (EMCO) platform and toolkit. We demonstrate that EMCO enables apps to accelerate their response time up to five-fold, and up to thirty-fold with pre-caching. We also demonstrate that EMCO can achieve further reductions in energy footprint through optimal selection of the cloud surrogate.
- We develop LAPSI<sup>1</sup> as a novel cloud-based runtime environment that supports characterization of surrogate performance. LAPSI is based on the Dalvik virtual machine which we have directly integrated with a cloud server. This creates a lightweight compiler for executing code and provides good scaling performance over a large community of clients.
- We develop and deploy a real testbed in Amazon EC2 and show that simply by dynamically adapting the resource allocation of the surrogate, EMCO increases gains in energy from offloading by 30%. In addition, we show that EMCO considerably reduces the energy required in the communication with cloud through the use of pre-caching.

1. Finnish for child, motivated by our runtime environment's parent-child structure.

## 2 RELATED WORK

The origins of computational offloading are in cyberforaging, which aimed at distributing computations from a resource constrained (mobile) device to nearby infrastructure [21]. Since then, most of the work has focused on linking the device with remote infrastructure or a hybrid that combines remote and local connectivity. An example of the latter are Cloudlets [15] in which constrained devices act as thin clients that connect to remote infrastructure through a nearby gateway service (the cloudlet).

In most of the early works, offloading took place at application level, i.e., the execution of the entire app was shifted to another device. Current state-of-the-art solutions provide finer-grained control over the offloading process by instrumenting the apps at code level. For example, MAUI [2] proposes a strategy where the software developer annotates the code and during runtime methods are offloaded when suitable conditions are detected by the MAUI profiler. Similarly, CloneCloud [4] exploits a dynamic approach, where a code profiler extrapolates pieces of bytecode at thread level to a server. CloneCloud uses static analysis to partition code. However, this is inadequate to verify runtime properties of code, which usually cause unnecessary code offloading. COMET [5] is another framework for code offloading that offloads at thread level, but introduces a DSM mechanism that allows the device to offload unmodified multi-threads. ThinkAir [3] addresses scalability by creating virtual machines (VMs) of a complete smartphone system in the cloud. However, since the development of mobile application uses annotations, the developer must follow a brute-force approach to adapt his/her application to a specific device. Moreover, the scalability claimed by ThinkAir does not support multi-tenancy, but the system creates multiple virtual machines based on Android-x86 within the same server for parallelization. Other frameworks for computational offloading include [16], [22], [23], [24], [25]. COSMOS [6] provides offloading as a service at method level using Android-x86. The framework introduces an extra layer in an existing offloading architecture to solve the mismatch between how individual mobile devices demand computing resources and how cloud providers offer them. However, COSMOS relies on the user to manually solve this problem, which is inefficient. Moreover, end users seldom are familiar with different cloud options and their costs.

Studies of computational offloading in the wild have mostly shown that offloading *increases* computational effort rather than reduces it [26]. This is due to the large amount of mobile usage contexts, and the poor understanding of the conditions and configurations in which a device offloads to the cloud. Some recent work attempts to overcome this problem by dynamically alleviating the issues of inferring the right matching between mobile and cloud considering multiple levels of granularity [11], [27]. Similarly, CDroid [26] attempts to improve offloading in real scenarios. However, the framework focuses more on data offloading rather than computational offloading.

The EMCO framework proposed in the present paper differs from existing solutions by harnessing mobile crowdsensing for characterizing context parameters of mobile apps. By relying on samples from a community, EMCO

Dataset	Carat	NetRadar
Data source	App store deployment (Multiple locations)	Finland (Metropolitan Helsinki area)
Number of samples	$2.9 * 10^9$ (from 328,000 apps)	135,000
Start time	01/08/2013	01/01/2015
Duration	1 year	1 year

TABLE 1: Summary statistics of application usage and network connectivity datasets.

reduces the effort and time of learning when to offload. This also improves the detection of offloading opportunities as it captures a large spectrum of possibilities from the community. Lastly, unlike other work that just uses the cloud as computational infrastructure to offload, EMCO takes advantage of the capabilities of cloud infrastructure to optimize the diagnosis of offloading contexts.

### 3 FEASIBILITY OF CROWDSENSING FOR QUANTIFYING OFFLOADING CONTEXTS

The focus of our work is on improving the effectiveness of offloading through the use of crowdsensing. Before detailing our approach, we investigate two fundamental questions in the use of crowdsensing. First, we examine the accuracy and generality of crowdsensed evidence traces and demonstrate that measurements provided by a community of devices can indeed be used to characterize offloading contexts. Second, we determine the amount of measurements required to characterize execution contexts. Specifically, we demonstrate that collecting a sufficient amount of measurements is infeasible on individual devices, but can be captured even from a (reasonably) small community of devices. We carry out our analysis using two large-scale datasets of mobile crowdsensing: Carat [19] and NetRadar [20]; see Table 1 for a summary of the datasets. The former captures the processing requirements and energy use of applications, whereas the latter captures network performance related parameters across different devices. Additionally, we carry out a crowdsensing simulation to characterize cloud resource allocation performance. In summary, we show that characterization via crowdsensing captures richer context information about smartphone usage (mobile context), network (communication context) and task acceleration (execution context).

#### 3.1 Smartphone Application Usage

We first consider the feasibility of using crowdsourcing to quantify the computational performance on applications, a critical factor in deciding when to offload. We perform the analysis by considering a large-scale data set containing energy footprints of mobile apps collected through the Carat application [19]. While many offloading frameworks operate on finer level of code granularity, i.e., on a method or class level, collecting such data is infeasible as it requires a large-scale deployment of apps which have been instrumented with the appropriate mechanism. Nevertheless, we demonstrate that even when considering application level granularity, it is possible to identify the conditions in which mobile applications require intensive processing. The analyzed data contains data from approximately 328,000 apps,

collected over a period of one year. For each app, the data contains the expected % of energy drain which is calculated over the entire community of devices; see [19] for details of the Carat methodology.

We first consider the minimum number of apps that are needed for identifying conditions for resource intensiveness. We analyze the distribution of expected energy drain, as given by the Carat application, and identify convergence of the resulting energy distribution. The results of our analysis are shown in Figure 1a. From the results we can observe that measurements from around 120 different apps are required. The average energy drain of apps within the community is  $< 0.004$ , which depicts battery consumption in % per second. In our experiments we consider this value as a lower bound for the detection of resource-intensive behavior.

To further illustrate the practical value of our approach, we assess the fraction of apps that could benefit from offloading through two representative examples of application categories that are generally believed to benefit from offloading: photo face manipulation and games. We also separately consider puzzle and chess games, two subcategories of games with heavy computational requirements, to demonstrate that our methodology applies to smaller sets of applications. For each category, we extract all unique apps. Overall there are about 550 photo face manipulation apps, about 7,805 game apps, about 717 puzzle apps, and about 166 chess apps, i.e., each category meets the minimum requirement for estimating resource intensiveness. Each group is compared with the average energy drain of the subset, which excludes its own energy drain. From the results in Figure 1 we can observe that around (1b) 43.84%, (1c) 44.56%, (1d) 42.75% and (1e) 33% of apps implement computational operations that require higher energy drain than normal, which is a significant number of apps. This result suggests that computational offloading could be beneficial for a large fraction of apps.

As the final step of our analysis, we consider the amount of samples required from the community to accurately characterize high energy drain. According to Figure 1f, we find that about 10,000 samples are gathered to estimate high energy drain within interval of  $0.0040 \pm 0.000035$ . As Carat collects samples each time the battery of a device changes, we get on average 50 – 100 samples per device in a day. Collecting the required 10,000 samples on a single device would thus clearly be infeasible, but a community of 100 devices could achieve this in a day or two.

A limitation of our analysis is that we considered app level offloading instead of method, class or thread level migration. Intuitively, we would expect finer-granularity to further improve decision making. However, this requires instrumenting the offloading mechanisms with the ability to monitor its own local and remote execution. To enable this kind of analytics, our implementation of EMCO integrates required mechanisms to capture both local and remote execution performance; see Section 4 for details.

#### 3.2 Network Connectivity

Communication latency during application usage is another key building block for making optimal offloading decisions. Latency is typically monitored through a network profiler, which is also responsible for characterizing different

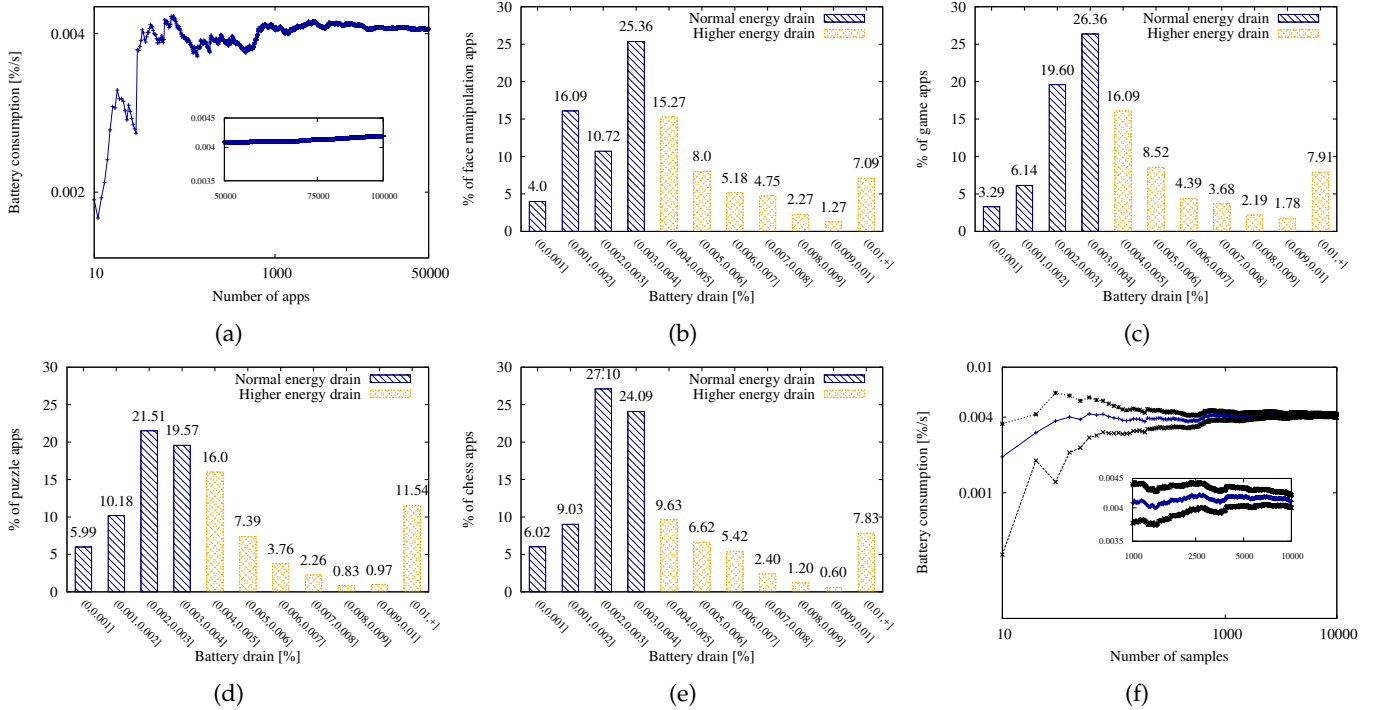


Fig. 1: Characterization of smartphone app usage via crowdsensing. (a) Number of applications for community characterization. (b, c, d, e) Smartphone apps that depict higher energy drain at different community granularity. (f) Number of samples for energy characterization.

communication interfaces, e.g., WiFi, 3G, LTE. Characterizing networking performance, however, can easily become energy-intensive, especially in suboptimal network conditions as more samples and a longer duration of monitoring are required. In this section we examine the potential of replacing the network profiler with crowdsensing. The main question we investigate is how many samples are required to characterize latency accurately at a particular time or location. We answer this question by analyzing the NetRadar dataset<sup>2</sup> [20] and characterizing the performance of 3G and LTE communications. Since latency in the cellular network also depends on the provider, we consider three different mobile operators, which are the main network providers in Finland. For anonymity, we refer to the operators using pseudonyms Operator-1, Operator-2 and Operator-3.

The topmost row of Figure 2 shows the average latency of the communication (RTT) for each provider. From the results, we can observe that the average RTT varies across the operators: 107ms (SD≈283) for Operator-1, 130ms (SD≈387) for Operator-2, and 129ms (SD≈403) for Operator-3. For LTE, the averages are 43ms (SD≈35) for Operator-1, 33ms (SD≈65) for Operator-2, and 39ms for (SD≈90) Operator-3. While there is a notable difference between 3G and LTE in terms of speed, we can observe higher stability in the communication link of LTE. In terms of quality of service, in the case of 3G, Operator-1 provides the most stable communication and Operator-3 the most unstable one. In the case of LTE, Operator-2 provides the most stable communication and Operator-1 the most unstable one, i.e., both the operator and the communication technology should be taken into

account when making offloading decisions.

Regarding the amount of data needed for characterization, we calculate the number of samples required to estimate latency. The average and median latency as a function of the number of samples for 3G and LTE are shown on the two bottommost rows of Figure 2. We can observe that the amount of data required to characterize the latency of communication depends on the stability of the communication link. For instance, for 3G, to characterize the most stable communication (Operator-1) with an interval of  $100 \pm 7.24$  at least 6,300 samples are needed. In contrast, if the same amount of samples is used to characterize the most unstable channel (Operator-3), then an interval of  $132 \pm 10.06$  is achieved. Thus, to improve accuracy and reduce error at comparable rates as Operator-1, an addition of 3,700 samples is required (10,000 in total). By generalizing the characterization time as a function of samples required, we have  $T = S / (S_t * n)$ , where  $T$  is the time that takes to complete the characterization,  $S$  is the total number of samples required for the characterization,  $S_t$  is the number of samples collected in time  $t$ , and  $n$  is the number of devices available to collect samples. In terms of sampling complexity, assume we want to characterize latency of 3G for Operator-3 and that we receive 5 samples per day from a device. A single device would then need about 2,000 days to complete the task. For a single device to gather the data in one year, it would have to collect 370 samples per day, significantly reducing battery life. Even if this was done, we would only be sampling latency at locations the device visited. In contrast, 100 devices (easily found at University Campus) can complete the characterization in about three weeks. By piggybacking the sample collection as part of

2. <https://www.netradar.org/en>

application data requests, significantly larger communities of devices could be reached, making it possible to complete the characterization in a matter of hours.

### 3.3 Augmented Resource Allocation

Computational provisioning takes place at infrastructure level through instances, which are physical or virtualized servers associated to a specific type of resource [27]. An instance follows a utility computing model, whereby consumers pay on the basis of their usage and type preferences. The cost of the instance is proportional to its type, which in turn determines the level of acceleration at which a task is executed and energy required in the communication [28]. Moreover, the type of the instance also defines its capacity to handle multiple offloading requests at once. Speeding up code execution certainly depends on how code is written. However, higher types of instances can process a task faster than lower ones as higher types can rely on larger memory span and higher parallelization in multiple processors [29].

To thoroughly characterize cloud resources, we conduct a simulation study where a large number of clients create a dynamic workload which is sent to the cloud [30]. Our simulator operates in two different operational modes: (i) concurrent and (ii) inter-arrival rate. In the former, the simulator creates  $n$  concurrent threads that offload a random computational task loaded from a pool of common algorithms found in apps, e.g., quicksort, bubblesort. Each thread represents a mobile device outsourcing a task to a remote server. This mode is utilized to benchmark the instances of the cloud. In inter-arrival rate mode, the simulator takes as parameters the number of devices (workload), the inter-arrival time between offloading requests and the time that the workload is active. This mode is utilized to produce synthetic workloads. We note that while throughput could be analyzed through server logs, such logs are (to our best knowledge) not freely available and we are not aware of any studies that would have done so. Thus, we focus instead on performing the characterization on the clients with the help of a community of devices.

To simulate crowdsensing in the wild, we configure our simulator in concurrent mode to stress the instances with a heavy workload of requests. In this experiment, each request that is created by the simulator is taken randomly from the pool. The processing required for each task is also determined randomly. The random nature of the experiment is important to verify all possible cases that can influence the processing of a task by a server. We conduct a 3-hour experiment per server to ensure coverage. We also check this characterization by relying on static load. We evaluate the influence of increasing users' load by configuring the workload from 1 to 100 in intervals of 10 users; per each server loads 1,10,20,30,40,50,60,70,80,90 and 100. Concurrent load is generated with an inter-arrival rate of 1 minute. This means that the maximum load (load=100) used for characterizing a server is  $\approx 18000$  ( $3 \times 60 \times 100$ ) requests and the minimum load (load=1) is  $\approx 180$  requests. The purpose of the 1 minute inter-arrival is to give enough cool down time to the server before stressing it again.

Figure 3 shows the results of the experiments. We can observe how the response time of the requests is distributed

through the interpercentile range of the processed load. The slope of the mean response time becomes less steep as we use more powerful instances. This suggests that the response time of a request is defined by the type of the server. Based on this property, we characterize each server into an acceleration group. We find that servers can be classified to 3 acceleration groups. This is important as servers with different costs provide the same level of acceleration when server utilization is low. However, there is also a correlation between server acceleration and utilization, the response time of an offloading request degrades based on the resource capabilities of the server as utilization increases.

Lastly, we used a minimax algorithm with a static input as workload to confirm our results. We can observe the differences between acceleration levels. A task is executed  $\approx 1.25$  times faster by a server of level 2 when compared with one of level 1. Similarly, a task is accelerated  $\approx 1.73$  times by a server of level 3 compared with level 1. The difference between levels 3 and 2 is also significant ( $\approx 1.36$  times acceleration). Our results thus demonstrate that crowdsensing is capable of characterizing both cloud throughput and acceleration across a wide range of contexts.

## 4 EMCO TOOLKIT AND ARCHITECTURE

In this section we describe the EMCO development toolkit and platform, which integrate novel mechanisms for harnessing crowdsensed evidence traces in offloading decision making. EMCO handles two types of mobile communication with the cloud: synchronous code offloading and asynchronous injection of data analytics via push notifications. The overall architecture of EMCO is shown in Figure 4. EMCO follows a model where the source code of an app resides both on the mobile and at the server, but it releases the mobile from a fixed cloud counterpart as implemented by other frameworks. EMCO encourages a scalable *provisioning as a service* approach where EMCO is deployed on multiple interconnected servers and responds on-demand to computational requests from any available server. EMCO also provides a toolkit that facilitates the integration of the mobile and cloud mechanisms in the development life-cycle. In the following we describe the different mechanisms and components of EMCO in detail.

### 4.1 Development Toolkit

To integrate offloading support into applications, EMCO provides a toolkit that consists of a GUI conversion tool that prepares the apps for client/server deployment. Our conversion tool receives as input an Android project and produces as output two projects, one for the mobile, and another for the cloud. Each project is defined with Maven-Android<sup>3</sup> so that for each project an Android Application Package (APK) file is built automatically. The tool also provides means for automatic deployment of the APK to the cloud. Our conversion tool has been designed to operate with existing frameworks and hence does not introduce any additional complexity, such as learning a new API.

The conversion tool transforms all methods that fulfill offloading requirements by wrapping each task (or method)

3. <http://code.google.com/p/maven-android-plugin/>

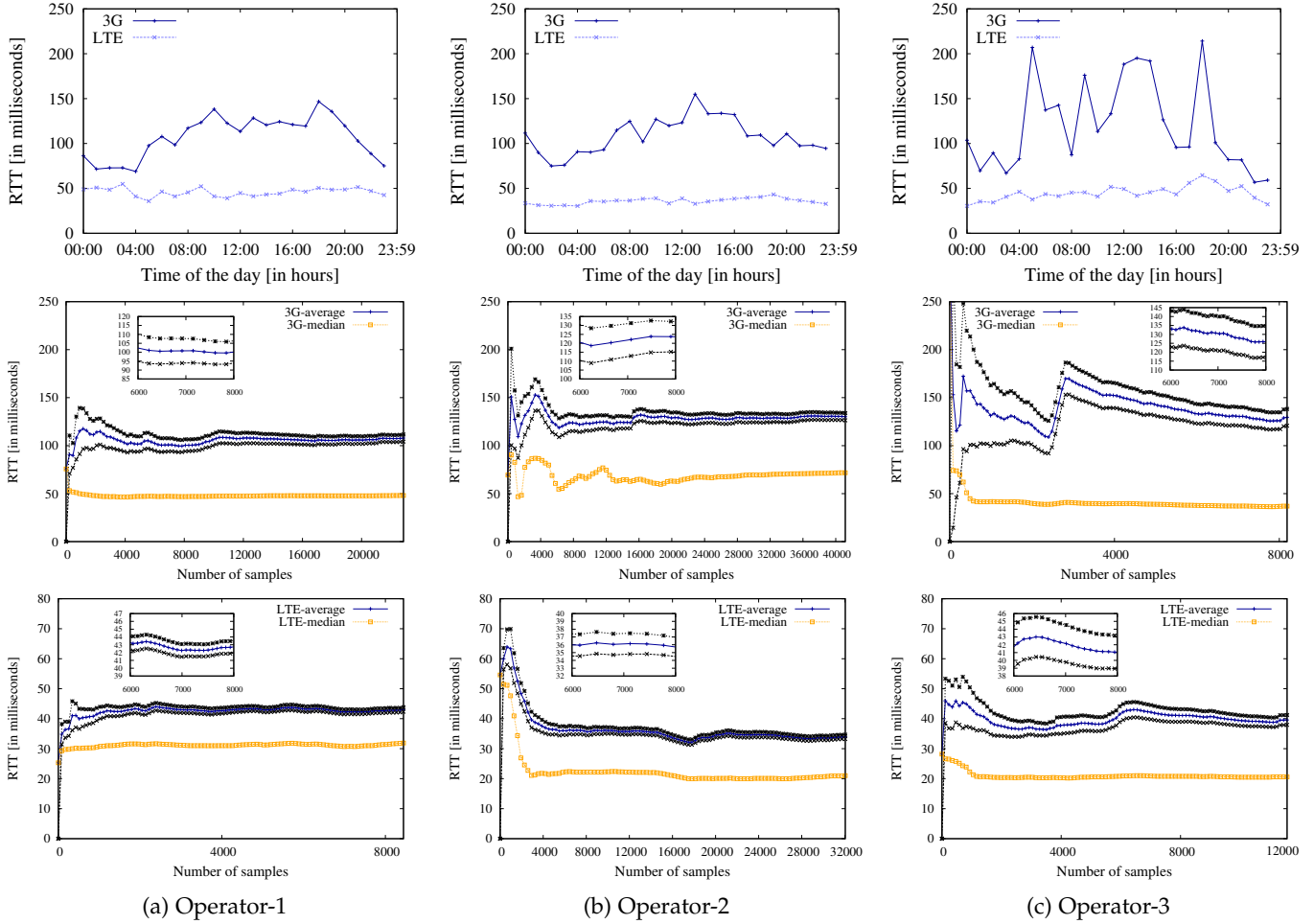


Fig. 2: Network characterization: (top and middle-top) Number of samples collected at different times of day; (middle-bottom) latency for 3G and LTE technologies (Operators 1-3); (bottom) amount of samples required to characterize latency.

into a container, which retrieves its execution details from a *descriptor* during runtime (see Sec. 4.3) [4]. The container is named with the name of the task, while the original task is renamed by adding the prefix "local". Similar instrumentation is followed in [3], [6], [8] with the main difference in our approach being the inclusion of a descriptor that governs the execution of tasks based on the community analysis. Besides equipping the app with offloading mechanisms, the conversion tool integrates the device with crowdsensing mechanisms, i.e., mechanisms for collecting context information and sharing it with the wider community. As has been shown in previous works [19], such monitoring has minimal impact on the battery lifetime of the mobile client.

## 4.2 EMCO Server

The *EMCO Server* provides computational resources for offloading and analyzes crowdsensed measurements to identify contexts where offloading is beneficial. We next detail the server's subcomponents and their key functionalities.

### Evidence Analyzer

The Evidence Analyzer analyzes evidence traces collected through crowdsensing and identifies conditions that are

beneficial for offloading. Our implementation of the Evidence Analyzer is based on a dual layer analysis pipeline. The first layer, dimension discovery, identifies factors that are most relevant for creating a *code offload descriptor* for a particular app. The second layer, classification, uses the evidence obtained by the first layer to build a classifier which integrates the decision logic used to construct the offload descriptor. The constructed classifier can then be deployed on the client to make offloading decisions directly without interfacing with the cloud.

The first layer, dimension discovery, operates on traces collected from the community. A trace consists of a list of parameters or features such as device model, app identifier, the name of the executed method, input of the method, local execution time, RTT, type of server, and remote execution time. Additionally, when the code is offloaded, each trace includes a binary variable that indicates whether the remote invocation succeeded. The invocation can fail, e.g., if the surrogate is unable to invoke the code (e.g., missing dependencies) or communicate back the result of that invocation (e.g., malformed serialization). To analyze evidence, we transform the evidence into *rate distributions* which model variations of a specific feature  $f$  over time. For each parameter, we construct separate rate distributions for

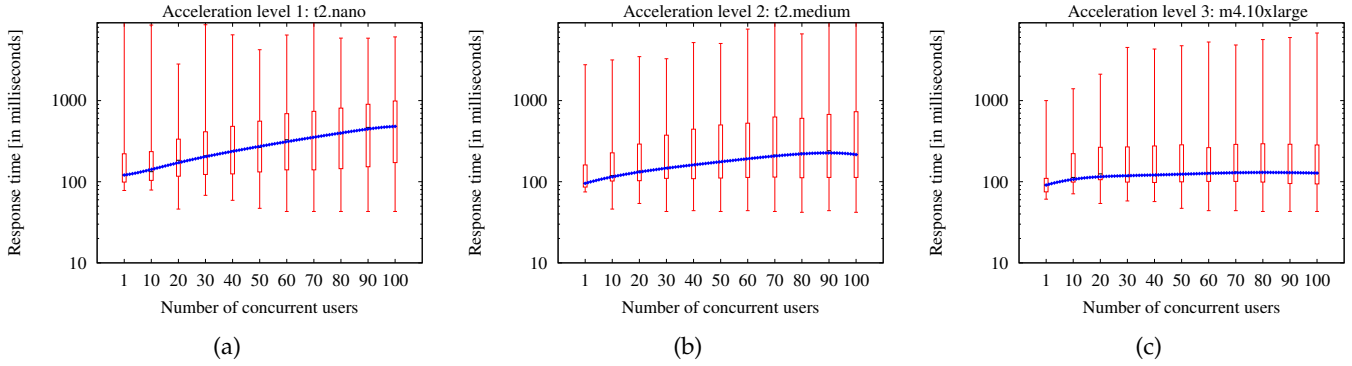


Fig. 3: Cloud-based characterization via crowdsensing gives different acceleration levels to servers based on performance degradation as more users are added.

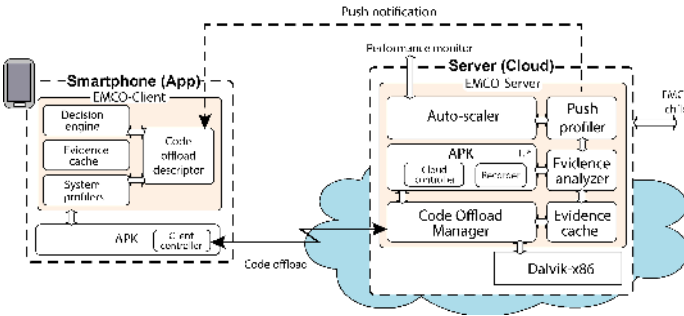


Fig. 4: Evidence-aware mobile code offloading architecture.

local and remote execution and compare them to identify when offloading is beneficial. Rate distributions are created for all obtainable parameters with the only requirement being that the feature and its value can be estimated during the execution of local and remote computations.

Formally, we construct rate distributions by converting evidence traces into a set of samples. Let  $s_t = (c, p, \hat{f})$  denote a sample taken at time  $t$  of code execution with identifier  $c$  (e.g., app id or method name) which is processed at  $p$  (local or remote) and which is characterized by the set of features  $\hat{f}$ . In our implementation we consider the features to be represented as key-value pairs, e.g., "response-time=2s", "energy-consumed=180J". We sort the samples based on time and split them into two groups, one for local execution and one for remote. For each group, successive measurements  $(s_{t_1}, fs)$  and  $(s_{t_2}, fs)$ , are converted into a rate distribution  $u = \frac{(fs_2 - fs_1)}{(t_2 - t_1)}$ , which later is associated with the rest of features to create a pair  $R = (u, \hat{f})$ . This association is important in order to keep track about the values of the other features with  $fs$ . Once the rate distributions are calculated for a given feature, pairs  $R_{local}$  and  $R_{remote}$  are compared as shown in Figure 5. The comparison consists of finding an overlap, or a *slice*, between the two distributions. Typically,  $R_{remote}$  values that are higher than  $R_{local}$  ones suggest that remote execution is counterproductive for the device in the context specified by the slice.

When slices that are counterproductive for offloading are encountered, we remove the corresponding samples from analysis and recalculate the rate distributions. We then analyze the next feature  $fs_i$  over the remaining set of

samples, i.e., we gradually refine the samples and features to identify conditions that are favorable for offloading. Once the refinement process completes, we construct a *code offload descriptor* which incorporates the most important dimensions. To create the descriptor, we compare the distance  $d$  between the local and remote distributions as this determines the level of improvement; see Figure 5. We aggregate the distances  $d$  into bins, and select the bin with the highest count for the descriptor. Note that the process thus identifies the most frequent conditions for obtaining benefits instead of selecting the ones with the highest performance increase. This is essential for avoiding overfitting to highly specific conditions which are encountered infrequently. EMCO also extracts various alternatives for surrogates based on the levels of frequency of  $d$ . Additionally, the set of rules is created at the end of the descriptor following a simple heuristic method presented in [31]. Since processing in the cloud is not constrained to service response time, we choose the heuristic solution. Naturally, any other method can be implemented and the modularity of EMCO provides an easy way to extend the functionality of the components.

Once the first layer has finished refining the evidence, the remaining samples are passed on the second layer, classification, to train a classifier which can be deployed on the mobile device as a decision engine. As classifiers we currently use C4.5 decision trees as they can be learned efficiently in the cloud from crowdsourced data and as the decisions made by the classifier can be easily explained [32]. An additional benefit from the use of decision trees is that it provides a mechanism for app developers to scrutinize which parts of their app have been offloaded and why.

The overall EMCO analysis is described in Algorithm 1 and consists of  $\approx 1500$  LOC written in Scala. To avoid overloading the surrogate with extra processing, the analysis is outsourced to a separate monitoring server, which also collects performance metrics about EMCO execution.

### Evidence Cache

To facilitate bootstrapping of offloading decisions, the server maintains a cache of previous method invocations together with results of the invocations. We take advantage of the fact that each offload request (including method name, input parameters, and other information relevant for an invocation) received by EMCO server is serializable and

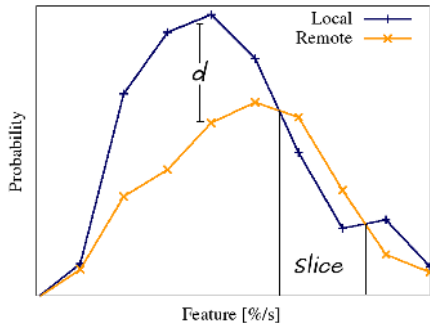


Fig. 5: Comparison of local and remote distributions.

---

#### Algorithm 1 Evidence analysis

---

**Require:** : path to the *traces*

- 1: **for** each *app* in *traces* **do**
  - 2:     Create samples.
  - 3:     Specify features to compare  $\langle fs \rangle$
  - 4:     **for** each *feature* in *fs* **do**
  - 5:         Compute rate distribution using samples.
  - 6:          $\langle R_{local}, R_{remote} \rangle$
  - 7:         Compare  $R_{local}$  with  $R_{remote}$
  - 8:         Remove slices where  $R_{remote}$  is higher
  - 9:         Update samples
  - 10:     **end for**
  - 11:     Calculate  $d$  using filtered samples.
  - 12:     Create *code offload descriptor*
  - 13:     EMCO policy creates the descriptor based on higher frequency of  $d$
  - 14:     Push descriptor to subscribed smartphone app
  - 15: **end for**
- 

calculate a digest key from the output of the serialization (SHA-1 checksum). We then index the resulting invocation results with the checksums. Whenever a new request is received, the server compares the checksum of the new request against the previously seen requests and returns the results from the cache whenever a match is found.

#### Code Offload Manager

The *Code Offload Manager* is responsible for handling code offload requests arriving at the server. The component extracts from each incoming request an *object package* which contains details about the app and details of invocation (method, type of parameters, value of parameters, and so forth). As outlined above, we first calculate a checksum of the request and compare it against the evidence cache. When a match is found, we return the result associated with the match. Otherwise the offload manager identifies the app and its associated (Dalvik) process in the server. The connection is forwarded to the process so that the *Cloud Controller* associated with the APK file instance can invoke the method via Java reflection. The result obtained is packed and sent back to the mobile device. In parallel, a trace is created by recording the execution of the method, e.g., execution time, type of server, CPU load, etc., and stored along with the data of the request and its result. Finally, the trace is passed to the *Evidence Analyzer* for analysis.

#### Push Profiler

The *Push Profiler* is the messaging component of the cloud server. Our implementation is based on GCM (Google Cloud Messaging) service<sup>4</sup> which is a notification service provided by Google to send messages to Android devices. We rely on an existing technique as messaging mechanisms are widely used by cloud vendors to trigger events on the mobiles to synchronize data with their cloud counterparts, e.g., Gmail. These mechanisms are well integrated with current handsets and they have been carefully optimized to have low energy consumption [33].

The EMCO server is registered to GCM service using a key obtained from the Google API. Multiple servers associated with one deployment of EMCO can use the same key. A mobile app subscribes to EMCO to receive notifications, obtaining a sender ID as a result. The ID is valid until the app unsubscribes (or GCM service is refreshed). EMCO then sends a message to the smartphone client which consists of a push request to the GCM service. The request consists of the *API key*, the *sender ID* and the message payload. Requests are queued for delivery (with maximum of five attempts) or stored if the mobile is offline. Once the message reaches the device, the Android system executes a *broadcast intent* for passing the raw data to the *code offload descriptor*.

EMCO relies on push requests to construct code offload descriptors for smartphone apps. As each message sent through GCM is limited to 1024 bytes, multiple messages are needed to form a descriptor. Messages are sent to the mobile based on different events, for instance, dynamic cloud allocation, periodic trace analysis, location detection of the mobile, etc. Since GCM is a popular service, it does not guarantee the delivery of a message and is unreliable to be used in real-time applications. To counter this problem, EMCO provides a generic interface, which can be used to integrate easily other push technologies, e.g., MQTT (Message Queue Telemetry Transport). We have also developed our own push server<sup>5</sup> based on XMPP, which can overcome the issues of proprietary push technologies and performance reliability of the notifications.

#### Auto Scaler

The *Auto Scaler* enables EMCO to support multi-tenancy. It uses the *AMI EC2 tools* wrapped into Java, to manage dynamic resource allocation of servers. The mechanism allows the framework to scale dynamically in a public or private cloud infrastructure. Note that while public cloud vendors provide mechanisms for scaling, e.g., Amazon Autoscale and Amazon Lambda, these mechanisms are designed for scaling SOA applications on demand. Resource requirements in computational offloading are driven by users' perception of app performance, which is governed by throughput of the cloud infrastructure. This differs from traditional SOA applications that focus on energy savings and reduced computational time without concern for throughput variations. EMCO addresses this gap by utilizing a hierarchical control and supervision scheme as shown in Figure 6. In our scheme, EMCO server acts as a parent, and a separate monitor server collects performance metrics from

4. [developer.android.com/google/gcm/](https://developer.android.com/google/gcm/)

5. <https://github.com/huberflores/XMPPNotificationServer>



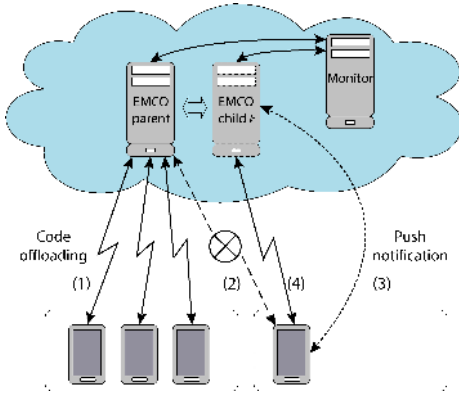


Fig. 6: Multi-tenancy and horizontal scaling in EMCO.

the available EMCO servers (based on collectD<sup>6</sup>). When the parent reaches its service utilization limits, it creates another child. In this case, the list of subscribed devices is split between the parent and child, and a push notification message is sent to each device to update the information of its assigned surrogate. In this manner, code offloading requests are balanced among the servers. Respectively, a child can become a second level parent when its service utilization limits are reached. When a child is destroyed, e.g., when it has served all requests and there are no incoming requests, the subscribed mobile devices are passed to the parent.

#### Surrogate Runtime Environment

EMCO Server is deployed on top of Lightweight App Performance Surrogate Image (LAPSI), a customized Dalvik-x86<sup>7</sup> virtual machine for Android. LAPSI has been built by downloading and compiling the source code of Android Open Source Project (AOSP) over the instance to target a x86 server architecture, and removing the Applications and Application Framework layers from the Android software architecture. Contrary to other frameworks [3], [4], EMCO does not rely on the virtualization of the entire mobile platform to execute the code remotely (such as with Android-x86). Full virtualization has been shown to consume substantial CPU resources and to slow down performance of code execution [34]. Moreover, full virtualization limits the number of active users offloading to a single server. In our case we also found that maintaining a complete execution environment requires substantial storage on the server (>100 GB) to install the software which includes AOSP code and Android SDK, among others. Moreover, once the OS is running in the server, all default features are activated (e.g., Zygote and GUI Manager) which are unnecessary for operating as surrogate and which result in higher CPU utilization of the surrogate. To avoid these overheads, LAPSI uses a Dalvik compiler extracted from the mobile platform which is deployed directly onto the underlying physical resources. This reduces the storage size of our surrogate to 60 GB. Moreover, our surrogate does not activate any default processes from the OS. Besides reducing overall system load, a major advantage of our approach

is that offload requests getting stuck or lost can be killed individually without restarting the entire offloading system.

### 4.3 EMCO Smartphone Client

The *EMCO Client* manages the offloading process on mobile devices. We next detail the subcomponents of the EMCO client and their key functionalities.

#### Context Profilers

Profilers are local monitors that gather information about parameters that can potentially influence the outcome of the offloading process. Examples of profilers include system profilers (e.g., network and CPU monitors usage), device monitors (e.g., battery state), and sensor data profilers (e.g., location or remaining battery). We rely on built-in Android APIs to get the status of each resource. In addition, EMCO also includes a code profiler that quantifies the execution of a task (i.e., a method or a process). A task is quantified in terms of response time and energy. The latter is estimated using PowerTutor [35], which has been shown to have low estimation error [19]. We also collect details about the task, e.g., cast type and parameters which allows more fine-grained analytics of offloading performance.

#### Code Offload Descriptor

The code offload descriptor controls the management (creation, updating, and deletion) of the data analytics injected from the cloud for a particular app running on the device. The descriptor consists of a characterization of the computational offloading process in JSON format which is stored as part of the app space; see Snippet 1 for an example. The descriptor effectively characterizes the logic that is used to determine which parts of the app are suitable of offloading, which acceleration levels are suitable for executing these parts, and which contexts are beneficial for offloading.

Formally, let  $D$  represent the dimensions contained in the descriptor,  $D = \{d_i \in D : 1 \leq i \leq N\}$ . Each  $d_i \in D$  consists of sets  $S$  of attributes  $s^{d_i}$ , where  $S^{d_i} = \{s_j^{d_i} \in S^{d_i} : 0 \leq j \leq N\}$ , and their values  $a^{d_i}$ , where  $A^{d_i} = \{a_k^{d_i} \in A : 1 \leq k \leq N\}$ . Our current implementation assumes attributes are represented as key-value pairs to ensure their storage and processing is lightweight. The descriptor can also contain other types of parameters, such as an indicator of whether the results of a specific method request are reusable and thus a candidate for caching, or parameters that describe user preferences regarding which surrogate to use.

#### Decision Engine

The decision engine predicts whether it will be beneficial to offload code at any particular moment. The engine relies on the offloading dimensions defined in the code offload descriptor and a C4.5 decision tree classifier learned on the EMCO server. We emphasize that EMCO is agnostic to the underlying decision engine and other techniques, including fuzzy logic, probabilistic models, logic-based reasoners and so forth, can equally well be integrated as with EMCO. The motivation for using C4.5 is that it can be effectively learned on the cloud, the output can be serialized and communicated back to the client with minimal data costs,

6. <http://collectd.org/>

7. Released as public in Ireland region of Amazon EC2 as *ami-b2487bd4*

**Snippet 1** Computational offloading descriptor in JSON

```

{
  "mobileApplication": "Chess",
  "deviceID": "i9300",
  "what-to-offload": [
    {
      "candidates-methods": "miniMax, betaPruning", %META-DATA
    }
  ],
  "when-to-offload": [
    {
      "bandwidth-threshold": "50", %META-DATA
      "interface": "LTE", %META-DATA
    }
  ],
  "where-to-offload": [
    {
      "option1-surrogates": "{(miniMax, m1.large),
        (betaPruning, m2.4xlarge)}", %META-DATA
      "option2-surrogates": "{(miniMax, m1.micro)}", %META-DATA
    }
  ],
  "location": [
    {
      "GPS-coordinates": "60.20453, 24.96185", %META-DATA
      "Preferable SSID": "HUPnet", %META-DATA
    }
  ]
}

```

the decision making mechanism is lightweight to run on the client, and the decision paths can be easily scrutinized and communicated to developers whenever needed [32].

*Evidence Cache*

Similarly to the server, the EMCO client utilizes a cache to maintain results of methods that are executed in the cloud. We store the results of method invocations into a hash table as key-value pairs, where the key is the name of the app appended with the name of the method, and the value contains the object sent from the cloud, the type in which the object should be cast during runtime and a time to live (TTL) value. The TTL value prevents to store permanently a result in the mobile. Thus, the result is assigned with a temporal lifespan. In this context, if the result is not used during that lifespan, then it is automatically removed from the hash table. The TTL value is updated every time the cache result is utilized in order to increase its lifespan.

**4.4 Worked Examples**

To illustrate EMCO in action, we consider two example scenarios, one related to a chess app that allows playing against an AI agent and the other related to a backtracking-based 3D model loading app; see Sec. 5.1 for description of the apps. We consider context descriptors associated with these apps shown in Snippets 2 and 3, respectively. The values of the context variables in the examples correspond to the average values of the parameters in the crowdsensing experiments described in Section 3 and are representative of the real-world circumstances where offloading takes place.

For the context associated with the chess app (Snippet 2), WiFi and LTE have similar latencies. Hence, EMCO's decision of whether to offload (and if so, where to offload) largely depends on the computational capability of the surrogate and the mobile device. In Snippet 3, on the other hand, WiFi has high latency (147ms) which would result in LTE being used as the communication technology (latency 39ms). In both scenarios the surrogate used for offloading

**Snippet 2** Chess experimental context.

```

label=Ctx1, device=i9300, CPU=20%
LTE=40ms, 3G=139ms, WiFi=27ms
method=Minimax, parameters="chessboard-15",
energy=533J,
surrogates=["m1.small" (response time (load=1)),
  "m1.medium" (load=1), "m1.large" (load=1),
  "m3.medium" (load=3), "m3.large" (load=2),
  "m3.xlarge" (load=6), "m3.2xlarge" (load=10),
  "m2.4xlarge" (load=8)
]

```

**Snippet 3** Backtracking experimental context.

```

label=Ctx0, device=i9250, CPU=43%
LTE=39ms, 3G=138ms, WiFi=147ms
method=ObjLoader, parameters="model.3 ds",
energy=212J,
surrogates=["m1.small" (load=1),
  "m1.medium" (load=1), "m1.large" (load=11),
  "m3.medium" (load=30), "m3.large" (load=10),
  "m3.xlarge" (load=22), "m3.2xlarge" (load=70),
  "m2.4xlarge" (load=50)
]

```

can be chosen based on estimated response time to minimize latency. Note that since both scenarios represent computationally intensive tasks, the effect of network bandwidth is not apparent. For data-intensive, real-time tasks, such as face manipulation on a live camera stream, the combined effect of bandwidth, latency, energy, and remote surrogate response time need to be taken into account. We have included the type of the device and its current CPU usage as part of the examples as both the savings from offloading and the energy footprint of the method depend on device characteristics and CPU load.

**Snippet 4** Two alternative Chess contexts.

```

label=Ctx2, device=i9300, CPU=80%
LTE=73ms, 3G=365ms, WiFi=80ms
method=Minimax, parameters="chessboard-30",
energy=97J

```

```

label=Ctx3, device=i9300, CPU=40%
LTE=65ms, 3G=411ms, WiFi=287ms
method=Minimax, parameters="chessboard-19",
energy=528J

```

Another benefit of EMCO is that it can react to dynamic changes in the offloading context. For example, the offloading decision can be adapted based on the location where the app is being used. To illustrate this, we consider two additional contexts shown in Snippet 4. The latencies in these contexts correspond to average latencies within two distinct cities (Ctx2 and Ctx3) and have been derived from our study in Section 3.2. We also consider different stages of the game, chessboard-19 and chessboard-30, where 19 and 30 correspond to the number of moves made. The former corresponds to an intermediary state which has many possible moves whereas the latter corresponds to a late stage of the game with fewer possible moves and hence lower computational requirements; see Figure 7. We next

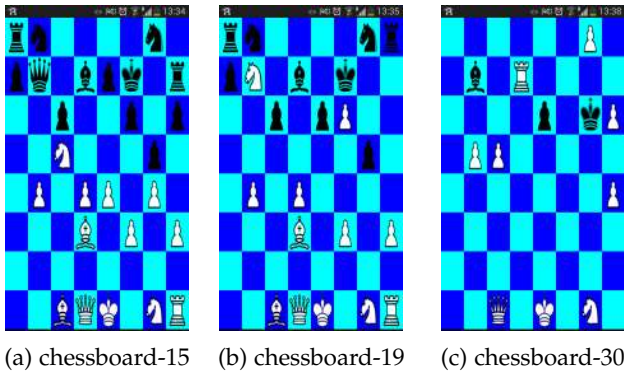


Fig. 7: Three execution states for the chess app: (a), (b) intermediate stage (15th and 19th moves) that are computationally intensive, and (c) a later stage of the game (30th move) which is computationally lightweight to solve.

analyze how EMCO behaves when the configuration of the chess application changes from Ctx1, to Ctx2 and Ctx3.

Assume first that the context of the chess application changes to Ctx2. This triggers an automated update of the values of the network interfaces using information from the community of devices running EMCO. We assume that the status of the cloud infrastructure remains the same. In Ctx2, LTE performs better than WiFi whereas in Ctx1 WiFi performed better. Consequently, EMCO would trigger a change of network interface to LTE. The example in Ctx2 corresponds to a later stage of the game (See Figure 7b) and thus the computational requirements are likely to be lower than in Ctx1. While in theory the request could be served on the local device, the high CPU load on the client in the example context would trigger EMCO to offload the request. As the computational requirements are likely to be low, any available surrogate can be chosen without significant influence on performance. Alternatively, the diagnosis generated by EMCO can be useful for optimizing resource allocation [36] to select the server with the cheapest cost that provides the same quality of service. Consider next a change from Ctx2 to Ctx3, the second in Snippet 4. The context corresponds to an intermediate state of the game and hence computing requirements are higher than in Ctx2. This would result in EMCO optimizing surrogate selection, e.g., to use a m2.4xlarge surrogate to optimize response time. In summary, these examples highlight EMCO's capability to make fine-grained adjustments depending on the currently prevailing offloading context and determine not only when to offload, but to which surrogate to offload.

## 5 EVALUATION

EMCO has been designed to improve the user experience of mobile apps through better offload decision making. This results in reduced energy footprint and improved responsiveness. This section evaluates EMCO, demonstrating that it achieves its design goals and improves on current state-of-the-art solutions. Our main experiments have been carried out through a real-world testbed deployment of EMCO. We also consider simulation-based experiments carried out using crowdsensing and offline profiling of applications.

### 5.1 Experimental Setup and Testbed

Our main experiment considers a typical scenario where a mobile app is associated with cloud infrastructure to improve energy and performance on the users' devices. We consider two simple examples of such apps: a backtracking-based 3D model loader app and a chess app with an AI agent. These apps were designed to incorporate common computational operations that are used in apps with heavy computational requirements and battery footprint. Our implementations are freely available for research purposes.

*Backtracking*<sup>8</sup> overlays a 3D model on the surface of the mobile camera. A user can interact with a model once it is displayed, e.g., moving and turning. Our implementation is equipped with ten models (in 3ds format) with the size of each model being between 500KB and 2000KB. Loading and displaying the mesh of each model requires a significant amount of processing and hence each model is loaded individually as required by the user. The heavy processing requirements of model loading make the app a good candidate for offloading. Loading the model in the cloud enables splitting it into parts so that the only parts the user interacts with are sent to the mobile device. Examples of similar apps include Augment <http://www.augment.com/augmented-reality-apps/> and Glovius 3D CAD viewer <https://www.glovius.com/>.

*Chess*<sup>9</sup> allows user to play (chess) against an AI agent on multiple levels of difficulty. The agent uses a minimax algorithm to determine the best move to make. The portion of code that corresponds to the minimax logic is the most intensive task in the application and hence a good candidate for offloading. Note, that the processing requirements of the minimax vary considerably depending on the state of the chess board with early and late moves typically being less intensive than intermediate ones; see Figure 7.

To conduct our experiments, we have deployed EMCO in Amazon (Ireland region / eu-west-1) using latest general purpose instances<sup>10</sup> (m3.medium, m3.large, m3.xlarge, m3.2xlarge) and a memory-optimized instance (m2.4xlarge). As client devices we consider two smartphone models: i9300 (Samsung Galaxy S3) and i9250 (Google Nexus). We follow the principle that an offloading request should both improve (or at least preserve) app responsiveness while reducing energy consumption of the client [29]. The appropriate level of responsiveness depends on the type of application. For example, real-time games or video applications would expect to have a speed of 30-60 frames per second, while interactive non-realtime applications can have response times up to one second. The energy footprint of the client is measured using a Monsoon Power Monitor<sup>11</sup>.

### 5.2 Characterizing Offloading Context

As the first step of our evaluation we consider a benchmark where we collect evidence for characterizing the execution

8. <https://github.com/huberflores/MeshOffloading>

9. <https://github.com/huberflores/CodeOffloadingChess>

10. We have also considered general purpose instances of previous generation: m1.small, m1.medium, and m1.large. These instances consistently resulted in worse performance than local execution on the mobile client and hence we have excluded them from our discussion.

11. <http://www.monsoon.com/LabEquipment/PowerMonitor/>

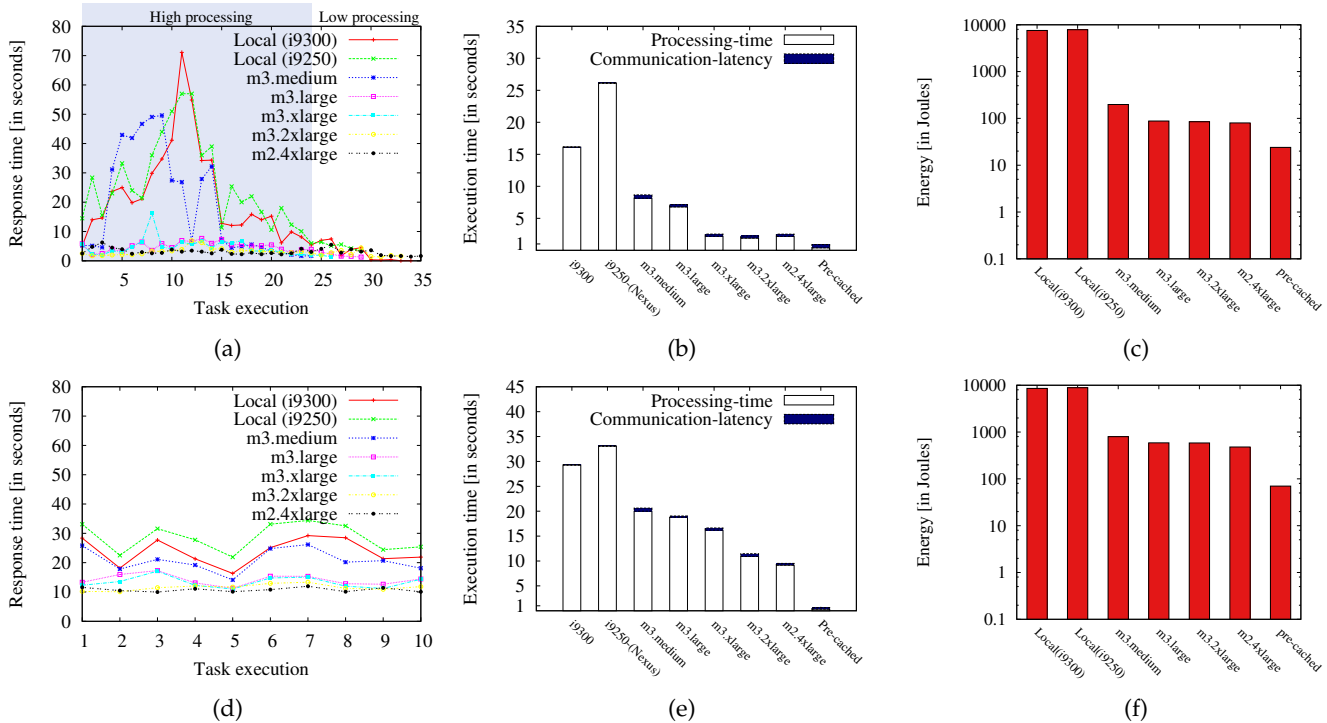


Fig. 8: App performance: (a), (b), and (c) from the chess application and (d), (e), and (f) from the backtracking app. The first column of each row corresponds to response time with a specific execution environment (surrogate or local) and app state (the number of moves since game start for chess and the index of the model being loaded for backtracking). The second column is the average execution time of the task, and the third column show the energy consumption of the client.

of our test apps. We execute both apps 48 times on each of the five surrogates (i.e., 240 remote execution trials) and 100 times on each mobile device (i.e., 200 local execution trials). For the chess app, the number of offloading requests depends on the length of the chess match. In total we observed 7200 requests, i.e., approximately 30 moves per game. For the backtracking app we loaded all 10 models during each trial, resulting in 2400 offloading requests (48 trials  $\times$  5 surrogates  $\times$  10 models). The surrogate used for offloading requests was selected in a round-robin fashion from the instances which were available.

The results of our benchmarks are shown in Figure 8. For both apps we observe that offloading generally improves both response time and energy footprint. However, we can also observe that the magnitude of the benefits depends on the execution context. This is particularly true for the chess app where we can see a large variation in response time depending on the complexity of the game state (Figure 7). For later stages of the game the benefits of offloading start to be marginal compared to local execution - unless the request is found in the pre-cache which results in latency becoming the dominant factor for offloading decisions.

### 5.3 Response Time and Energy Footprint

We next evaluate the gains in response time and energy obtained by using EMCO. As part of the experiment we also assess the performance gains provided by caching results of offloading requests and reusing them. We perform the evaluation by executing the apps locally on the device and offloading them into all available surrogates and measuring the associated response time and energy consumption.

Figure 8c and Figure 8f show the energy consumed in each offloading case for each application, respectively. From the results we can observe that the effort to offload is different depending on the time of active communication. From the figure we can also observe caching to further reduce energy consumption and response time.

Existing offloading methods bind mobile devices to any available server in the deployment and it remains static server during application usage. In contrast, EMCO implements a mechanism that allows requests to be treated independently, which means that can be handled by multiple servers based on its processing requirements. To compare EMCO against current state-of-the-art, we have developed a generic code offloading framework which (i) makes offloading decisions based on network latency and estimated energy, but (ii) has no support for diagnosing cloud performance and allocates requests to a random server instead (in the experiments these correspond to m3.medium for the chess app and m3.2xlarge for the backtracking app). Our baseline does not directly follow any single system, but incorporates mechanisms commonly found in current offloading frameworks<sup>12</sup> [2], [3], [6], [8].

The response times of EMCO and the baseline are shown in Figure 9a. For the chess app (top figure), we can observe that EMCO is capable of adjusting the offloading environment according to characteristics of the computing task and current context, which results in a stable and fast response time throughout. The performance of the baseline suffers during complex intermediate states as the system cannot

12. <https://github.com/huberflores/CodeOffloadingAnnotations>

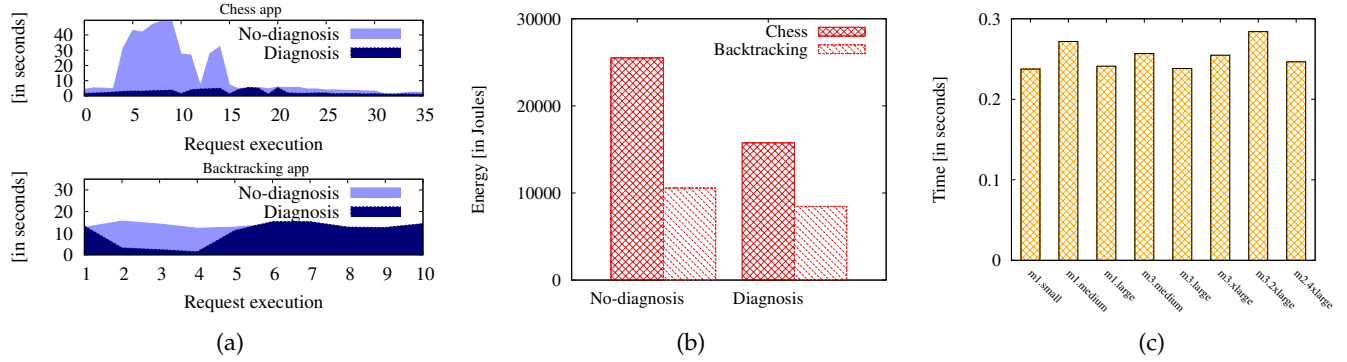


Fig. 9: (a): Request response time of both scenarios with and without diagnosis, (b): energy spent by the client, and (c): response time of pre-cached results in different instances types.

adapt to the changing computational requirements. Furthermore, the baseline is incapable of adjusting its performance to different load conditions on the surrogate which further causes fluctuations in performance. In particular, the performance of the lower-end general purpose instances, m3.medium and m3.xlarge, degrades in heavy workload conditions (Figure 10a) and thus ideally the app must offload to m3.large, m3.2xlarge or m2.4xlarge. Since our goal is to minimize response time and energy, EMCO offloads in these situations to m2.4xlarge. By using the diagnosed configuration, the chess app accelerates its responsiveness five-fold compared to local execution. Moreover, dynamic surrogate selection reduces 30% and 10% of the energy required in the communication with the cloud when compared with random surrogate selection in the chess and backtracking app scenarios, respectively.

For the backtracking app (bottom of Figure 9a) offloading should ideally use as surrogate m3.medium, m3.large, m3.xlarge, m3.2xlarge or m2.4xlarge. However, under heavy workload conditions, the app should be offloaded to m3.large as it is the lowest general purpose instance which is capable of providing major gains in throughput in such circumstances. Since the generic framework is allocated with a powerful machine randomly, it can provide similar performance especially when the number of surrogates to choose from is small. However, the likelihood of achieving that same level of performance every time depends on the available cloud resources. Moreover, despite the fact that both frameworks provide similar response times rates to users, only EMCO is able to achieve further acceleration through the identification of reusable offloading results from its crowdsensing process. Thus, by using pre-cached functionality, EMCO accelerates about 30 times the response time of applications installed in devices such as i9250.

Figure 9b compares the energy consumption of the client under EMCO and the baseline system. We observe accumulative gains in energy obtained by EMCO in a continuous offloading environment. This also can be seen as the difference between plain instrumentation and instrumentation with data analytic support (EMCO). Based on these results, we can conclude that EMCO is more effective as it saves more energy in comparison with existing offloading frameworks. EMCO's decision engine also appears to be lighter in comparison to the mechanisms used by MAUI and related

frameworks. Specifically, when we compare the MAUI linear programming (LP) mechanism with our decision engine using the same amount of parameters, the results show that EMCO consumes  $\approx 12\%$  less energy. The processing steps required to take a decision using LP depends on the number of parameters introduced in the engine.

The utilization of a powerful server is naturally associated with a cost. To mitigate the cost of running increasingly powerful servers, after pre-cache results are stored, they can be shared with lower capability servers, which can then respond to duplicated offloading requests. For instance, pre-cached results processed in m3.2xlarge can be transferred to m1.small, such that m1.small can then serve the result without any processing, yielding rates comparable to a m3.2xlarge. We conduct multiple experiments to determine the response time of a reusable result using different instances types. Figure 9c shows the results. We can observe that the response time of a reusable result is stable on different instances types on Amazon, average  $\approx 0.253$ , standard deviation (SD)  $\approx 0.016$ . We also notice that a cached result served from any surrogate outperforms a fresh result computed by the most powerful surrogates shown in Figures 8b and 8e.

## 5.4 Generalizability

We next demonstrate that the benefits of EMCO translate to a wider range of apps and that EMCO is capable of optimizing their performance in a wide set of execution contexts. We accomplish this through an evaluation of 10 apps in simulated offloading contexts. In our evaluation, each app creates virtual offloading requests whose context parameters are determined using the crowdsensed datasets considered in our work. By simulating the context we can ensure the parameters are representative of real world situations encountered by users. Simulation also helps to provide reproducibility for our analysis.

For the evaluation, we chose 10 open source apps from the F-droid<sup>13</sup> repository; see Table 2. We chose the most popular applications from the repository that had computationally intensive parts that benefit from offloading. We used a combination of community analysis and static profiling to determine which apps contained computationally

13. <https://f-droid.org/>

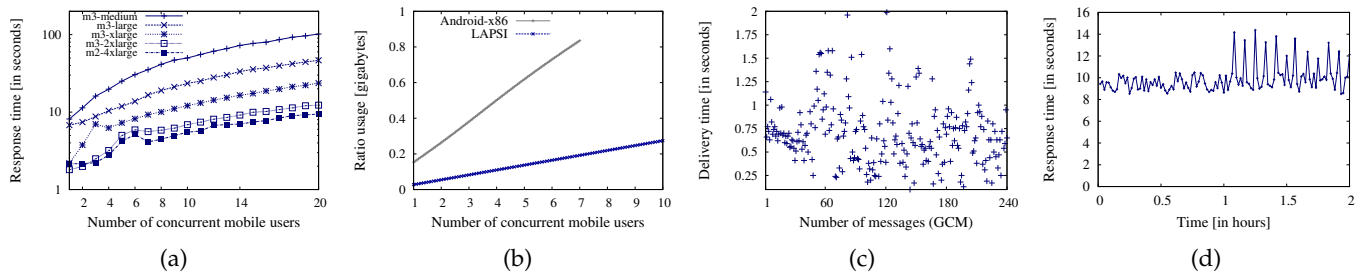


Fig. 10: (a) EMCO’s ability for handling concurrent users; (b) Memory usage comparison between Android-x86 and LAPSI; (c) Response time of app reconfiguration via push notification; (d) CPU load of a EMCO server when handling up to 20 concurrent requests, first hour without execution monitoring second hour with monitoring.

intensive tasks. Specifically, Carat community analysis [19] was first used to assess the overall performance of the app. Only apps with a rate of 0.004 or higher were considered, which translates to roughly 7 hour battery lifetime for the device while the app is running continuously. Next, we performed static profiling over the app’s code and validated its resource intensiveness using the *AT&T Video Optimizer*<sup>14</sup>. For each app, we chose all so-called outlier methods, which are methods with a high processing step count and whose execution time exceeds the average processing requirement of the rest of the app significantly. Note that, as we use execution time and processing counts for choosing which methods to simulate, this also helps to ensure the methods can indeed be accelerated. As an example, app1 contains one outlier method that requires  $\approx 2$  seconds to complete, while all remaining methods can be executed in 800 milliseconds or less. The final set of apps were from the following categories: puzzles (app1, app3, app6, app7, app8), games (app2, app4, app5, app10) and tools (app9); see Table 2 for the package names of the apps.

For each computationally intensive method we create virtual offloading requests where network latency was determined based on the NetRadar dataset (see Section 3). We assume that the cloud deployment consists of three instances (t2.nano, t2.medium, and m4.10xlarge). To determine the response time of the application, we calculated an expected runtime for the app and combined it with latency values of the simulated request. The expected runtime was determined by (i) calculating a bound for the execution of the method through static profiling, and (ii) scaling the estimated runtime with an acceleration level corresponding to the currently selected cloud instance derived through benchmarking of each cloud instance; see [30] for validation of this method. For each method we created 1,000 virtual requests, and calculated response time using both the diagnosis approach of EMCO and a no-diagnosis variant used in existing systems. The results in [30] show that, while the acceleration levels of cloud instances contain variation, the magnitude of performance improvements (or decrease) is stable between different types of cloud instances.

The results of our evaluation are shown in Table 2. In the evaluation we consider the upper bound on execution time obtained through static profiling as the response time requirement of each app, and calculate the fraction of times

this can be achieved. EMCO is capable of satisfying the requirements in 78.88% of cases, whereas existing techniques succeed only in 47.28% of the requests. From the table we can also observe that the extent of delay in failed cases is significantly smaller for EMCO (454ms) than for the no-diagnosis baseline (843ms). In other words, EMCO not only performs significantly better in terms of meeting the response time requirements, but significantly reduces delays in cases where the requirements cannot be met.

To put the results into context, we performed a cost-benefit analysis by deriving utilities for EMCO and the baseline solution. First we calculate the overall cost of offloading by associating a monetary cost for serving each simulated request. The cost was determined so that higher cost reflects access to better resources. As each request was offloaded, we normalized the total cost of all requests to the range  $[0, 1]$  so that a value of 0 corresponds to always using the slowest instance, and a cost of 1 to always using the fastest instance. For normalizing the costs, we assign weights 0.5, 1.0, and 150.0 to the three instance types, which reflects the current pricing structure of Amazon Ireland and serves as representative example of real-world costs. To determine an utility for each request, we assume meeting the response time requirement provides a reward  $R$ , whereas failing to meet the requirement results in a penalty  $P = \alpha tR$ . Here  $\alpha$  is a scaling factor that determines the importance of succeeding relative to failing to meet the user specified response rate requirements. The variable  $t$  is the extent at which the request exceeds the response rate requirement. Using the values in Table 2, it can be shown that EMCO provides better utility-cost ratio than existing approaches whenever  $\alpha > 0.0808$ , i.e., whenever the penalty for failing the response is at least 8% of the value of succeeding. If the penalty of exceeding the response time is very low, offloading to the slowest instance is naturally the best policy. To this end, the higher cost of EMCO compared to the no-diagnosis case is a direct consequence of having higher success rate, and hence using more resources. In practice we would expect the penalty of failure exceeding the reward of meeting the response time requirement (i.e.  $\alpha > 1$ ), implying that EMCO performs better in most cases of practical interest.

## 5.5 Scalability

We next consider the scaling performance of the server under heavy load of incoming requests. We carry out our

14. <https://developer.att.com/video-optimizer>

Id/App package/No.Methods	Resp. Time Req. [s]	Fulfilled [%]		Cost [0-1]		Exceeded by [s]	
		No-diag	Diag	No-diag	Diag	No-diag	Diag
app1 / org.ligi.gobandroid / 1	2	56.8%	90.5%	0.35	0.49	0.67	0.5
app2 / emre.android.atomix / 3	1	19.5%	47.7%	0.33	0.88	0.84	0.38
app3 / com.uberspot.a2048 / 2	2	58.9%	89.4%	0.34	0.50	0.92	0.49
app4 / org.bobstuff.bobball / 1	2.5	69%	92.9%	0.35	0.32	0.73	0.33
app5 / org.scoutant.blokish / 4	1.5	25.4%	56.1%	0.34	0.87	0.99	0.43
app6 / name.boyle.chris.sgt puzzles / 1	2	50.8%	88%	0.33	0.60	0.86	0.56
app7 / io.github.kobuge.games.minilens / 3	2	61.6%	89.7%	0.34	0.50	0.88	0.51
app8 / com.petero.droidfish / 2	2.5	69.1%	94.5%	0.35	0.38	0.98	0.41
app9 / fr.ubordeaux.math.paridroid / 5	2	44.2%	91.5%	0.36	0.64	0.78	0.61
app10 / net.minetest.minetest / 1	1	17.5%	48.5%	0.36	1.0	0.78	0.32
<b>Average for all apps</b>	<b>1.85</b>	<b>47.28%</b>	<b>78.88%</b>	<b>0.34</b>	<b>0.62</b>	<b>0.843</b>	<b>0.454</b>

TABLE 2: Crowd-sensed context evaluation, results of 1000 requests per method and ratio of request times below responsiveness requirements (% fulfilled).

evaluation by straining the cloud infrastructure with synthetic requests created by our trace simulator (described in Sec. 3.3) and measuring throughput as a function of load. In our benchmarks, we only consider server instances that provide better response time than local execution (i.e., m1.small, m1.medium, and m1.large were excluded). As part of the evaluation we also compare the scaling performance of our LAPSI surrogate against a deployment using Android-x86 surrogates.

Figure 10a shows the response rate of the different cloud instances as the number of concurrent users increases. We observe that the capability to handle concurrent requests mirrors the available resources with the response rate slowing linearly as the number of concurrent requests increases. As expected, we can also observe that the magnitude of the slope depends on the available resources of the server with the most powerful servers having the smallest coefficients. Three surrogates (m3.xlarge, m3.2xlarge, m2.4xlarge) can support over 20 concurrent users before the response rate drops below local execution, whereas two of the surrogates (m3.medium, m3.large) can support 10 concurrent users before dropping below local execution. The response times include also latency induced by EMCO, which typically is around 100 milliseconds. Figure 10b shows the memory use of LAPSI surrogates compared to Android-x86 surrogates. The memory use of the Android-x86 surrogates rapidly increases, with 7 concurrent users requiring over a gigabyte of writable memory. The total memory footprint (including shared memory) of these 7 users fills the 8 GB memory of a m3.large instance, so that no more Android-x86 surrogates can run on the instance. In contrast, with our LAPSI deployment the increase in memory use is modest. To put our result into context, an Android-x86 deployment would need to boot up a new instance to serve each 7 users wishing to use the same application at the same time, whereas our deployment can serve over 28 users simultaneously with the same amount of resources. This makes LAPSI better suited for deployments targeting communities of users. In EMCO, horizontal scaling is achieved using a notification server as load balancer. The notifications allow the mobile to reconfigure the surrogate in which it is currently subscribed. As part of the experiment we also benchmark our push notification mechanism by measuring its delivery rate. We fix the message size of the notifications to 254 bytes, which is the minimum amount of data required to reconfigure

the surrogate. We first send 15 consecutive messages (once per second), after which we have a half an hour period of inactivity. We repeat these phases over a 8 hour period, resulting in a total of 240 messages. The sleep periods are used to mitigate the possibility of the cloud provider associating our experiments with denial of service attack and the duration of the experiments was designed to collect samples at different times of day.

Figure 10c shows the response time results for the notification messages. The average time is 0.75 seconds (median 0.66, SD 0.69). Note that the response time depends on the underlying messaging protocol and the availability of the server at the cloud provider [33]. The overall delay in surrogate migration depends also on the events taking place on the server side. The slowest case is when a new child node is created (i.e., when server utilization reaches maximum limit) which results in approximately one minute reconfiguration time. This consists of allocating a new server by the cloud provider (approximately 40 seconds on Amazon) and reconfiguring clients through push notifications. When a child node is removed (i.e., when server utilization falls below a lower limit), reconfiguration is faster as the subscribed devices are simply moved to the parent node and the only delay results from push notifications.

Finally, we demonstrate the benefits of using a separate monitoring server as part of the EMCO architecture. Figure 10d shows the performance of a m2.4xlarge surrogate during a 2h period with 20 concurrent users. During the first hour, no performance metrics are gathered, whereas during the second hour the surrogate was configured to collect performance measures. From the figure we can clearly see that monitoring cloud execution burdens the surrogate, degrading its performance, whereas the performance of the surrogate is stable when no metrics are gathered. Accordingly, our experiment demonstrates that separating the monitoring of cloud execution to a separate server enables faster and more stable response time.

## 5.6 Optimality of Decision Making

As the final step of evaluation we consider the optimality of the EMCO decision engine. To carry out the evaluation, we consider our benchmark evaluation in Section 5.2 and establish a ground truth by comparing the computational effort of the device between remote and local execution. We use crowdsourced data to characterize execution contexts to allow EMCO to make informed decisions.

We conduct our evaluation using five-fold cross-validation, i.e., all measurements in four folds are used to train our decision engine (i.e., C4.5 decision tree classifier) and its predictions are compared to ground truth for the remaining fold. In 99.82% of the situations, EMCO was correctly able to choose the option with the smallest energy cost between local and remote execution. Generally the performance depends on the extent of intersecting areas in the rate distributions, which represent conditions where the benefit of offloading is uncertain or counterproductive. In our experiments, the amount of training data was sufficient for learning a model that has few intersections. During early stages of the app, the performance is likely to be lower (worst-case performance 77%). However, as shown in Section 3, a community of moderate size typically produces the necessary amount of measurements for minimizing intersections arising from uncertainty within a day or two.

## 6 DISCUSSION

A key benefit of EMCO is that it facilitates bootstrapping of offloading decision making. While the amount of data that is required for a characterizing the execution of a particular app depends on the complexity of the app, its usage contexts, and execution environment, in most cases EMCO can provide substantial improvements in the bootstrapping process. Indeed, the more variability there is in the execution, the more substantial the benefits EMCO can bring through community analysis. We note that, while not a core focus for the present paper, another benefit of EMCO is that the number of intersecting segments between local and remote distributions provides an estimate of the uncertainty in the app execution, which in turn can be used to estimate the point where sufficient amount of data has been collected. Estimating this on individual devices is difficult as one device typically only experiences a small set of execution contexts and the performance of the system is likely to degrade when a previously unseen context occurs.

As with any other history-based approach, the performance of EMCO gradually improves as more evidence is gathered from the community of devices. When no data is available to bootstrap the system, the decision inference considers only contextual parameters collected in real-time from the context profiles similarly to other proposed frameworks. However, even in such case, major performance improvements can be obtained compared to other frameworks. For example, our dynamic surrogate migration mechanism is capable of providing access to surrogates in a service-oriented style. We have also experimentally demonstrated that, once sufficient amounts of evidence have been collected, our framework outperforms existing solutions by providing better support for making offloading decisions. As most of the evidence analysis takes place in the cloud, the effort required by the mobile device to execute EMCO is minimal compared with other frameworks. For example, both linear programming [2] and Hidden markov models [37] require processing raw data from system profilers in order to determine when to offload. In contrast, EMCO avoids raw data processing as the constructed models are sent from the cloud to mobile devices. While transferring data to the cloud requires some energy, these costs typically

are minimal [19] and their effect on battery can be mitigated by piggybacking data as part of application requests.

Recently, the *Android runtime* (ART)<sup>15</sup> has been introduced as a novel runtime compiler that is intended as replacement for Dalvik. ART has been designed to operate using Dex bytecode, which is the same format as used by Dalvik and hence LAPSI and EMCO can be easily ported to ART. Compared to our LAPSI implementation, the main difference with ART is that it implements *Ahead-of-Time* (AOT) compilation instead of *Just-in-Time* (JIT). This means that an application that uses ART is compiled once during the first execution, and the static portions of the app will not be compiled again for subsequent executions. This is closely related to the pre-caching functionality used as part of EMCO. However, our solution is more general as it has been designed to support also collaborative reuse of method invocations between devices.

Our results have shown the potential of migrating cloud surrogates dynamically based on concurrent load, execution context, and other factors. This differs from existing systems which either provide no means of effectively matching the client with the cloud infrastructure or make the decision the responsibility of the app designer [6]. In the former case, the effectiveness depends on the optimality of the initial allocation, whereas in the latter case the effectiveness is reliant on the designers knowledge of the intricacies of offloading contexts. Our use of crowdsensing to monitor cloud execution provides more flexibility in the process, allowing us to dynamically determine the optimal acceleration level and to adapt to varying loads on the cloud size. EMCO also lends itself to more elaborate cloud allocation, e.g., by associating a cost (or utility) for energy, surrogate type, and responsiveness, we can minimize the cost of surrogate allocation for each client.

EMCO can be applied in many variable context scenarios. For example, an interstate bus line or a high-speed train with onboard WiFi is susceptible to constantly changing network conditions, which can result in poor app responsiveness. Routing offloading requests of connected clients with EMCO mitigates response time variations and provides a superior experience for latency-critical applications. As another example, EMCO can diagnose situations for opportunistic offloading. For example we can use the community to identify application combinations that increase CPU usage while running in the background. From our study presented in section 3.1, we can observe 60% – 70% CPU use with 13.65% of apps (TripAdvisor, Whatsapp, Facebook, Waze, among others). By detecting currently active apps, we can utilize offloading for less intensive tasks to reduce energy drain caused by heavy CPU usage. However, since resource allocation has a higher impact in computational offloading compared to other opportunities, we focus our diagnosis on resource allocation. A wider range of context values for opportunistic offloading can be handled, e.g., using a decision tree [32] and ranking the context values in order of their impact on offloading.

EMCO also implements an energy profiler that estimates the energy required from a portion of code. However, estimation of energy consumption with any kind of

15. <https://source.android.com/devices/tech/dalvik/>



software tool, e.g., PowerTutor, suffers from an estimation error. EMCO overcomes this problem by supplementing local estimates made on the device with community-based estimates that aggregate over a large community of devices. By collecting many samples from the same task, EMCO is able to reduce the error to determine when a computational task requires high energy consumption for the device or not.

## 7 SUMMARY

The present paper contributed by proposing crowdsourced evidence traces as a novel mechanism to optimize offloading decision making. As our first contribution we performed an analysis of two crowdsensed datasets to demonstrate that crowdsensing is a feasible mechanism for characterizing parameters relevant for offloading decision making. As our second contribution, we proposed a toolkit and platform, EMCO, that analyzes these traces to identify reusable and generic contexts where offloading is beneficial. Further performance gains are achieved through caching and a novel cloud surrogate LAPSI that we have developed and integrated as part of EMCO. Experiments carried out on a testbed deployment in Amazon EC2 Ireland demonstrated that EMCO can consistently accelerate app execution while at the same time reducing application footprint. In addition, we evaluate offloading optimization through crowd-sensed context diversity over multiple off-the-shelf applications. We also demonstrated that EMCO provides better scalability than current cloud platforms, being able to serve a larger number of clients without variations in performance. We provide our framework, use case and tools as open source in GitHub along with the LAPSI located in Amazon EC2.

## ACKNOWLEDGMENTS

The authors thank the anonymous reviewers for their insightful comments. This research has been supported, in part, by projects 26211515 and 16214817 from the Research Grants Council of Hong Kong.

## REFERENCES

- [1] S. Ickin, K. Wac, M. Fiedler, L. Janowski, J. Hong, and A. K. Dey, "Factors influencing quality of experience of commonly used mobile applications," *IEEE Communications Magazine*, vol. 50, no. 4, pp. 48–56, 2012.
- [2] E. Cuervo, A. Balasubramanian, D.-k. Cho, A. Wolman, S. Saroiu, R. Chandra, and P. Bahl, "Maui: making smartphones last longer with code offload," in *Proceedings of the ACM International Conference on Mobile systems, applications, and services (MobiSys 2010)*, (San Francisco, CA, USA), June 15–18, 2010.
- [3] S. Kosta, A. Acinas, P. Hui, R. Mortier, and X. Zhang, "Thinkair: Dynamic resource allocation and parallel execution in the cloud for mobile code offloading," in *Proceedings of the Annual IEEE International Conference on Computer Communications (INFOCOM 2012)*, (Orlando, Florida, USA), March 25–30, 2012.
- [4] B.-G. Chun, S. Ihm, P. Maniatis, M. Naik, and A. Patti, "Clonecloud: elastic execution between mobile device and cloud," in *Proceedings of the Annual Conference on Computer systems (EuroSys 2011)*, (Salzburg, Austria), April 10–13, 2011.
- [5] M. S. Gordon, D. A. Jamshidi, S. Mahlke, Z. M. Mao, and X. Chen, "Comet: code offload by migrating execution transparently," in *Proceedings of the USENIX Conference on Operating Systems Design and Implementation (OSDI 2012)*, (Hollywood, CA, USA), October 08–10, 2012.
- [6] C. Shi *et al.*, "Cosmos: Computation offloading as a service for mobile devices," *Proceedings of the ACM International Symposium of Mobile Ad Hoc Networking and Computing (MobiHoc 2014)*, August 11–14, 2014.
- [7] A. Saarinen, M. Siekkinen, Y. Xiao, J. K. Nurminen, M. Kempainen, and P. Hui, "Smartdiet: offloading popular apps to save energy," in *Proceedings of the ACM Conference on Applications, technologies, architectures, and protocols for computer communication (SIGCOMM 2012)*, (Helsinki, Finland), August 13–17, 2012.
- [8] H. Flores, R. Sharma, D. Ferreira, V. Kostakov, J. Manner, S. Tarkoma, P. Hui, and Y. Li, "Social-aware hybrid mobile offloading," *Pervasive and Mobile Computing*, 2016.
- [9] M. Conti, M. Kumar *et al.*, "Opportunities in opportunistic computing," *IEEE Computer*, vol. 43, no. 1, pp. 42–50, 2010.
- [10] E. Peltonen, E. Lagerspetz, P. Nurmi, and S. Tarkoma, "Energy modeling of system settings: A crowdsourced approach," in *Proceedings of International Conference on Pervasive Computing and Communications (PerCom 2015)*, (St. Louis, MO, USA), March 23–27, 2015.
- [11] H. Flores and S. Srirama, "Mobile code offloading: should it be a local decision or global inference?" in *Proceedings of the ACM International Conference on Mobile systems, applications, and services (MobiSys 2013)*, (Taipei, Taiwan), June 25–28, 2013.
- [12] Y. Kwon, S. Lee, H. Yi, D. Kwon, S. Yang, B.-g. Chun, L. Huang, P. Maniatis, M. Naik, and Y. Paek, "Mantis: Efficient predictions of execution time, energy usage, memory usage and network usage on smart mobile devices," *IEEE Transactions on Mobile Computing*, vol. 14, no. 10, pp. 2059–2072, 2015.
- [13] S. Sigg, E. Lagerspetz, E. Peltonen, P. Nurmi, and S. Tarkoma. (2016) Sovereignty of the Apps: There's more to Relevance than Downloads. <https://arxiv.org/abs/1611.10161>.
- [14] H. Ma *et al.*, "Opportunities in mobile crowd sensing," *IEEE Communications Magazine*, vol. 52, no. 8, pp. 29–35, 2014.
- [15] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, "The case for vm-based cloudlets in mobile computing," *IEEE Pervasive Computing*, vol. 8, no. 4, pp. 14–23, 2009.
- [16] M.-R. Ra *et al.*, "Odessa: enabling interactive perception applications on mobile devices," in *Proceedings of the ACM International Conference on Mobile systems, applications, and services (MobiSys 2011)*, (Washington, DC, USA), June 28 – July 1, 2011.
- [17] X. Gu, K. Nahrstedt, A. Messer, I. Greenberg, and D. Milojicic, "Adaptive offloading for pervasive computing," *IEEE Pervasive Computing*, vol. 3, no. 3, pp. 66–73, 2004.
- [18] K. Kumar and Y.-H. Lu, "Cloud computing for mobile users: Can offloading computation save energy?" *Computer*, vol. 43, no. 4, pp. 51–56, 2010.
- [19] A. J. Oliner, A. P. Iyer, I. Stoica, E. Lagerspetz, and S. Tarkoma, "Carat: Collaborative energy diagnosis for mobile devices," in *Proceedings of the ACM Conference on Embedded Networked Sensor Systems (SenSys 2013)*, (Rome, Italy), November 11–14, 2013.
- [20] S. Sonntag, J. Manner, and L. Schulte, "Netradar-measuring the wireless world," in *Proceedings of IEEE International Symposium on Modeling & Optimization in Mobile, Ad Hoc & Wireless Networks (WiOpt 2013)*, (Tsukuba Science City, Japan), May 13–17, 2013.
- [21] R. Balan, J. Flinn, M. Satyanarayanan, S. Sinnamohideen, and H.-I. Yang, "The case for cyber foraging," in *Proceedings of ACM SIGOPS European Workshop 2002*, (Saint-Emilion, France), July 01, 2002.
- [22] N. Fernando, S. W. Loke, and W. Rahayu, "Mobile cloud computing: A survey," *Future Generation Computer Systems*, vol. 29, no. 1, pp. 84–106, 2013.
- [23] P. Bahl, R. Y. Han, L. E. Li, and M. Satyanarayanan, "Advancing the state of mobile cloud computing," in *Proceedings of ACM MobiSys Workshop 2012*, (Low Wood Bay, Lake District, United Kingdom), June 25–29, 2012.
- [24] M. A. Hassan *et al.*, "Pomac: Properly offloading mobile applications to clouds," *Energy (J)*, vol. 25, p. 50, 2014.
- [25] Y. Kwon, H. Yi, D. Kwon, S. Yang, Y. Cho, and Y. Paek, "Precise execution offloading for applications with dynamic behavior in mobile cloud computing," *Pervasive and Mobile Computing*, 2015.
- [26] M. V. Barbera, S. Kosta, A. Mei, V. C. Perta, and J. Stefa, "Mobile offloading in the wild: Findings and lessons learned through a real-life experiment with a new cloud-aware system," in *Proceedings of the Annual IEEE International Conference on Computer Communications (INFOCOM 2014)*, (Toronto, Canada), April 27 – May 2, 2014.
- [27] H. Flores and S. Srirama, "Adaptive code offloading for mobile cloud applications: Exploiting fuzzy sets and evidence-based

learning," in *Proceedings of ACM MobiSys Workshop 2013*, (Taipei, Taiwan), June 25–28, 2013.

- [28] F. A. Silva *et al.*, "Benchmark applications used in mobile cloud computing research: a systematic mapping study," *The Journal of Supercomputing*, vol. 72, no. 4, pp. 1431–1452, 2016.
- [29] H. Flores, P. Hui, S. Tarkoma, Y. Li, S. Srirama, and R. Buyya, "Mobile code offloading: From concept to practice and beyond," *IEEE Communications Magazine*, vol. 53, no. 4, pp. 80–88, 2015.
- [30] H. Flores, X. Su, V. Kostakos, J. Riekkii, E. Lagerspetz, S. Tarkoma, P. Hui, Y. Li, and J. Manner, "Modeling mobile code acceleration in the cloud," in *Proceedings of the Annual IEEE International Conference on Distributed Computing Systems (ICDCS 2017)*, (Atlanta, GA, USA), June 5–8, 2017.
- [31] K. Nozaki, H. Ishibuchi, and H. Tanaka, "A simple but powerful heuristic method for generating fuzzy rules from numerical data," *Fuzzy sets and systems*, vol. 86, no. 3, pp. 251–270, 1997.
- [32] E. Peltonen, E. Lagerspetz, P. Nurmi, and S. Tarkoma, "Constella: Crowdsourced system setting recommendations for mobile devices," *Pervasive and Mobile Computing*, vol. 26, pp. 71–90, 2016.
- [33] H. Flores and S. N. Srirama, "Mobile cloud messaging supported by xmpp primitives," in *Proceedings of ACM MobiSys Workshop 2013*, (Taipei, Taiwan), June 25–28, 2013.
- [34] M. Toyama, S. Kurumatani, J. Heo, K. Terada, and E. Y. Chen, "Android as a server platform," in *Proceedings of IEEE Consumer Communications and Networking Conference (CCNC 2011)*, (Las Vegas, Nevada, USA), January 9–12 2011.
- [35] M. Gordon *et al.*, "Power tutor, a power monitor for android-based mobile platforms," 2009.
- [36] M. Mazzucco, M. Vasar, and M. Dumas, "Squeezing out the cloud via profit-maximizing resource allocation policies," in *IEEE International Symposium on Modeling, Analysis & Simulation of Computer and Telecommunication Systems (MASCOTS 2012)*, (Arlington, Virginia, USA), August 7–9, 2012.
- [37] K. K. Rachuri, C. Mascolo, M. Musolesi, and P. J. Rentschler, "Sociablesense: exploring the trade-offs of adaptive sampling and computation offloading for social sensing," in *Proceedings of the Annual International Conference on Mobile computing and Networking (MobiCom 2011)*, (Las Vegas, Nevada, USA), September 19–23 2011.



**Huber Flores** is a postdoctoral research scientist in the Department of Computer Science at the University of Helsinki, Finland. He obtained his PhD in computer science from the Faculty of Mathematics and Computer Science, University of Tartu, Estonia. He is also member of ACM (SIGMOBILE) and IEEE. His major research interests include mobile offloading, mobile middleware architectures, and mobile cloud computing. Email: huber.flores@helsinki.fi



**Pan Hui** received his PhD from the Computer Laboratory at University of Cambridge. He is the Nokia Chair in Data Science and Professor of Computer Science at the University of Helsinki, the director of HKUST-DT System and Media Lab at the Hong Kong University of Science and Technology, and an adjunct professor of social computing and networking at Aalto University. He worked for Intel Research Cambridge and Thomson Research Paris and was a Distinguished Scientist at Deutsche Telekom Laboratories. He has published more than 200 research papers and with over 13,000 citations. He is an associate editor for IEEE Transactions on Mobile Computing and IEEE Transactions on Cloud Computing and an ACM Distinguished Scientist. Email: panhui@cse.ust.hk



**Petteri Nurmi** is a Senior Researcher at the University of Helsinki and Science Advisor at Moprim. He received his PhD in Computer Science from the University of Helsinki in 2009 and he is an adjunct professor of Computer Science at University of Helsinki since 2012. His research focuses on sensing, covering system design, algorithms, theoretical models, and new application areas. He has published over 80 scientific publications. Email: petteri.nurmi@helsinki.fi



**Emil Lagerspetz** is a postdoctoral scholar at the University of Helsinki. He completed his PhD in 2014 at the University of Helsinki. His research pertains to the analysis of large datasets, mobile computing, and energy efficiency. Dr. Lagerspetz has received the University of Helsinki Doctoral Dissertation Award in 2015, Jorma Ollila Grant in 2015, and the Marc Weiser Best Paper Award at IEEE PerCom in 2015. Email: emil.lagerspetz@helsinki.fi



**Sasu Tarkoma** is a Professor of Computer Science at University of Helsinki since 2009. He received his Ph.D. degree in 2006 in Computer Science from the University of Helsinki. He has worked in the IT industry as a consultant and chief system architect as well as principal researcher and laboratory expert at Nokia Research Center. He has over 140 scientific publications, 4 books and 4 US Patents. Email: sasutarkoma@helsinki.fi



**Jukka Manner** received his PhD computer science from the University of Helsinki in 2004. He is a full professor (tenured) of networking technology at Aalto University, Department of Communications and Networking (Comnet). His research and teaching focuses on networking, software and distributed systems, with a strong focus on wireless and mobile networks, transport protocols, energy efficient ICT and cyber security. Email: jukka.manner@aalto.fi



**Vassilis Kostakos** is Professor of Human Computer Interaction at the School of Computing and Information Systems at the University of Melbourne. He has held appointments at the University of Oulu, University of Madeira and Carnegie Mellon University. He has been a Fellow of the Academy of Finland Distinguished Professor Programme, and a Marie Curie Fellow. He conducts research on ubiquitous and pervasive computing, human-computer interaction, and social and dynamic networks. Email: vassilis.kostakos@unimelb.edu.au



**Yong Li** is currently a Faculty Member of the Department of Electronic Engineering, Tsinghua University. He received his PhD in electronic engineering from Tsinghua University, Beijing, China in 2012. He has previously worked as Visiting Research Associate with Telekom Innovation Laboratories, The Hong Kong University of Science and Technology. His research interests are in the areas of networking and communications, including mobile opportunistic networks, device-to-device communication, software-defined networks, network virtualization, and future Internet. Email: liyong07@tsinghua.edu.cn



**Xiang Su** is a postdoctoral researcher in computer science in the Center of Ubiquitous Computing at the University of Oulu. Su received a PhD in technology from the University of Oulu. His research interests include semantic technologies, the Internet of Things (IoT), knowledge representations, and context modeling and reasoning. Email: xiang.su@oulu.fi