



ELSEVIER

Contents lists available at ScienceDirect

## Journal of Network and Computer Applications

journal homepage: [www.elsevier.com/locate/jnca](http://www.elsevier.com/locate/jnca)

## Review

## A survey of computation offloading strategies for performance improvement of applications running on mobile devices



Minhaj Ahmad Khan

Bahauddin Zakariya University Multan, Pakistan

## ARTICLE INFO

## Article history:

Received 6 August 2014

Received in revised form

1 May 2015

Accepted 25 May 2015

Available online 3 July 2015

## Keywords:

Computation offloading

Mobile computing

Performance improvement

Mobile cloud computing

Cyberaging

## ABSTRACT

Handheld mobile devices have evolved from simple voice communication devices to general purpose devices capable of executing complex applications. Despite this evolution, the applications executing on the mobile devices suffer due to their constrained resources. The constraints such as limited battery lifetime, limited storage and processing capabilities produce an adverse impact on the performance of applications executing on the mobile devices.

Computation offloading addresses the issue of limited resources by transferring the computation workload to other systems having better resources. It may be oriented towards extending battery lifetime, enhancing storage capacity or improving the performance of an application. In this paper, we perform a survey of the computation offloading strategies correlated with performance improvement for an application. We categorize these approaches in terms of their workload distribution and offloading decisions. We also describe the evolution of the computation offloading based environment as well as a categorization of application partitioning mechanisms adopted in various contributions. Furthermore, we present a parameter-wise comparison of automated frameworks, the application domains that benefit from computation offloading and the future challenges impeding the evolution of computation offloading.

© 2015 Elsevier Ltd. All rights reserved.

## Contents

|  |    |
|--|----|
| 1. Introduction . . . . .  | 29 |
| 2. Offloading taxonomy: architectures and effectiveness . . . . .    | 29 |
| 2.1. Computation offloading architectures . . . . .                  | 29 |
| 2.2. Trade-offs for offloading decisions . . . . .                   | 30 |
| 3. Evolution of offloading and wireless technology . . . . .         | 30 |
| 4. Offloading architectures and approaches . . . . .                 | 31 |
| 4.1. Static offloading . . . . .                                     | 31 |
| 4.2. Dynamic offloading . . . . .                                    | 33 |
| 5. Application partitioning for computation offloading . . . . .     | 35 |
| 5.1. Static partitioning . . . . .                                   | 35 |
| 5.2. Dynamic partitioning . . . . .                                  | 35 |
| 6. Comparison of offloading frameworks . . . . .                     | 36 |
| 7. Application domains benefiting from offloading . . . . .          | 36 |
| 8. Current challenges for effective computation offloading . . . . . | 37 |
| 8.1. Partitioning . . . . .  | 37 |
| 8.2. Automated transparency & portability . . . . .                  | 37 |
| 8.3. Security . . . . .  | 37 |
| 8.4. Application requirements . . . . .                              | 37 |
| 9. Conclusion . . . . .  | 38 |
| References . . . . .   | 38 |

E-mail address: [mik@bzu.edu.pk](mailto:mik@bzu.edu.pk)<http://dx.doi.org/10.1016/j.jnca.2015.05.018>

1084-8045/© 2015 Elsevier Ltd. All rights reserved.

## 1. Introduction

With the advent of smartphone technologies, the mobile devices have become ubiquitous. These devices are no longer constrained to providing only communication services. Instead, these devices are capable of executing applications with diverse requirements. The processing required by these applications may range from simple mathematical computations performed by a calculator to a very complex voice recognition system.

The execution of complex applications requires the mobile devices to possess powerful resources. The scarcity of these resources has adverse effects on the ever-growing usage of the mobile devices. For instance, the statistics according to *StatCounter* show that about 30.66% of the platforms used for web browsing are the mobile systems (smartphones/tablets) (*StatCounter, 2014*). Consequently, the mobile market plays a significant role in e-commerce and sales growth. This role is however diminished by the fact that the mobile systems have limited energy and power resources. Although there have been efforts to incorporate high performance multiple core processors in smartphones, the gap between the existing and the required resources continues to grow. In this context, the computation offloading is a mechanism that enables us to bridge the gap by making intensive computations execute on large systems having sufficient resources as required by the application. This not only makes a resource constrained mobile system seem like a high-end powerful machine, but also enables to perfectly utilize the existing resources.

The computation offloading is not a novel idea as it has evolved from various paradigms incorporating distributed computing (*Dinh et al., 2013; Kumar and Lu, 2010; Sanaei et al., 2014; Fontana et al., 2013*). The performance improvement of an application is achieved by partitioning it into several subprograms each of which may be assigned to a different processor for execution. Each processor makes use of its own memory and/or shares the memory with other processors to perform computations in parallel. Subsequently, the results are returned to the processor controlling the overall execution.

A cloud computing platform is also based on the intuition of distributed computing and offers the compute services through a Service Level Agreement (SLA) on a large network usually the Internet. It differs from other computing paradigms since an assurance regarding availability of services is provided to the users. The Mobile Cloud Computing (MCC) therefore refers to provision of services through a cloud to mobile devices that are characterized with limited resources (*Dinh et al., 2013; Kumar and Lu, 2010; Sanaei et al., 2014; Fontana et al., 2013; Juntunen et al., 2012; Khan et al., 2014b; Berl et al., 2010*). The computation of a mobile application may be offloaded to another resource-rich system termed as *surrogate*. Such kind of computation offloading not only mitigates the issue of limited resources of mobile devices but also enables to harness the processing power of high-end machines that will otherwise be idle (*Barbera et al., 2013; Ou et al., 2007; Cui et al., 2013; Sanaei et al., 2012; Miettinen and Hirvisalo, 2009; Kumar et al., 2013; Satyanarayanan et al., 2009*).

In this paper, we perform a comprehensive survey of the computation offloading strategies impacting the performance of the applications executing on mobile devices. Although the computation offloading has also been aimed at saving energy required for executing an application (*Lu et al., 2013; Hong et al., 2009; Wen et al., 2012; Rudenko et al., 1998; Nurminen, 2010; Nimmagadda et al., 2009; Miettinen and Nurminen, 2010; Mayo and Ranganathan, 2004; Sinha and Kulkarni, 2011; Ge et al., 2012*), but in this paper, we mainly consider the contributions which impact the execution performance (computation speed) of applications running on mobile devices. The survey encompasses the research work for computation offloading arranged in terms of

multiple aspects including the taxonomy, strategies, evolution pattern and relevant application domains. We also present a categorization of partitioning approaches adopted in different contributions and a parameter-wise comparison of main offloading frameworks. We also discuss main issues related to computation offloading and suggest possible approaches to address these issues effectively.

The rest of the paper is organized as follows. *Section 2* describes the offloading taxonomy in terms of architectures and criteria for its effectiveness. The evolution of offloading and wireless technologies is described in *Section 3*. The offloading approaches and contributions aimed at performance improvement are surveyed in *Section 4*. A categorization of partitioning approaches used in computation offloading is given in *Section 5*. A parameter-wise comparison of the automated computation offloading frameworks is described in *Section 6*, whereas the applications benefiting from computation offloading are discussed in *Section 7*. The main issues related to an effective implementation of computation offloading are discussed in *Section 8* together with their solutions before concluding at *Section 9*.

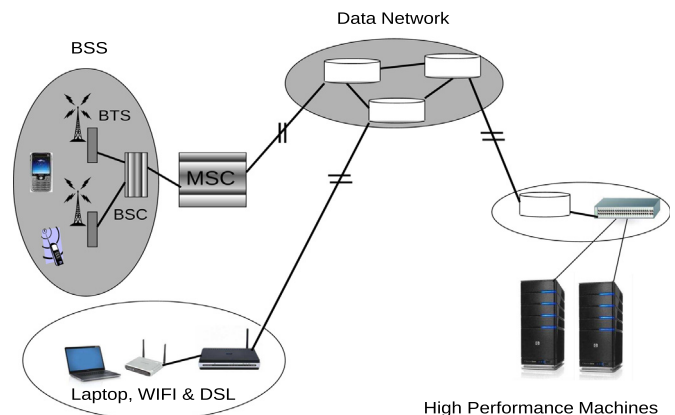
## 2. Offloading taxonomy: architectures and effectiveness

Many clients such as mobile phones or low power laptops require computation to be offloaded to powerful server machines. The decision of offloading may not always be beneficial to leverage the performance or energy requirements as a significant overhead is involved while offloading computations. This section describes succinctly the general architectures for which offloading may be required and the parameters that impact its effectiveness.

### 2.1. Computation offloading architectures

In an environment supporting computation offloading, the users with mobile devices are connected to a high performance server in different ways. The simplest form of this connection is made through Wi-Fi based networks that connect mobile devices to other machines using wireless routers as shown in *Fig. 1*. The wireless router not only connects devices to a local network but also may be connected to a DSL device thereby providing connections to remote servers through Internet.

Similarly, in a more complex form, the users with mobile devices first connect to a wireless network through devices such as Base Transceiver Station (BTS), Base Station Controller (BSC), and Mobile Switching Center (MSC) to transfer data to public data networks. The communication data is then transferred through gateways to any local network on which the high performance machines are hosted.



**Fig. 1.** Offloading architecture.

After establishing a connection with the high performance machines, the mobile devices may perform a lookup operation to search for services that may be provided by the high performance server machines. This may also be termed as the first operation initiated by the application. The application may however opt to perform the lookup operation at a later time during execution depending upon the time at which the offloading decision is made and the requirement of the application. The client machines in these environments are usually low power mobile devices, and consequently, the computation offloading strategies take into account the cost/benefit analysis in terms of the execution time and energy requirements. The server machines are mostly the high-end standalone servers, or machines connected to form a grid, cluster, cloud or a combination of these. The computers in a grid are loosely coupled, whereas those in a cluster are tightly coupled with highly efficient interconnection interfaces such as *Myrinet*. A cloud system, in contrast, uses virtualization to enable multiple operating systems so that remote users can access services offered by the cloud platform.

2.2. Trade-offs for offloading decisions

For minimization of execution time and reduction of energy, the computation offloading from a mobile device to a server machine is performed by applying a specific criterion to ensure that the offloading will be beneficial (Li et al., 2001; Xian et al., 2007; Wolski et al., 2008; Nimmagadda et al., 2010; Cuervo et al., 2010; Wang and Li, 2004b; Niu et al., 2014; Elgazzar et al., 2013). The required criteria take into account several parameters as elaborated below.

For minimizing execution time, let  $O_r$  be the overhead of runtime activities including the time for data transfer and the time for offloading code, i.e.

$$O_r = T_d + T_o, \tag{1}$$

where  $T_d$  is the time for data transfer and  $T_o$  is the time taken for offloading code (performing offloading decision, partitioning and the code transfer). Let  $T_s$  be the time to execute code on the server machine and  $T_m$  be the time to execute code on the mobile device. The computation offloading is considered effective for minimization of execution time, if we have

$$T_s + O_r < T_m. \tag{2}$$

Similarly, for energy reduction, let  $E_d$  represent the energy for data transfer and  $E_o$  represent the energy required for offloading. Let  $E_m$  represent the energy required for execution of entire application on the mobile device and  $E_r$  be the energy required for runtime activities. The computation offloading is effective for reducing requirements if

$$E_r < E_m, \tag{3}$$

where  $E_r$  is represented as

$$E_r = E_d + E_o. \tag{4}$$

3. Evolution of offloading and wireless technology

The term “offloading” has been used widely since year 1995. Its usage has evolved together with the evolution of distributed and parallel computing paradigms. Figure 2 shows the number of publications each year<sup>1</sup> citing the term offloading.

Similarly, the research work referring to the terms “data offloading” and “computation offloading” is also increasing gradually, as shown in

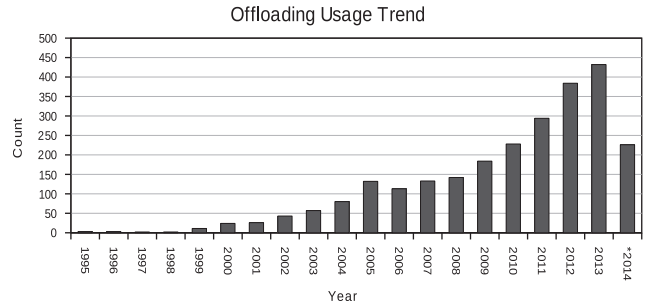


Fig. 2. Offloading usage trend.

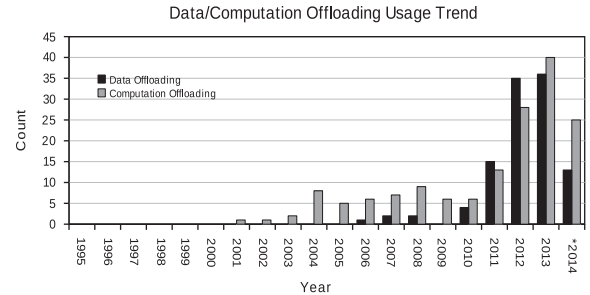


Fig. 3. Data and computation offloading usage trend.

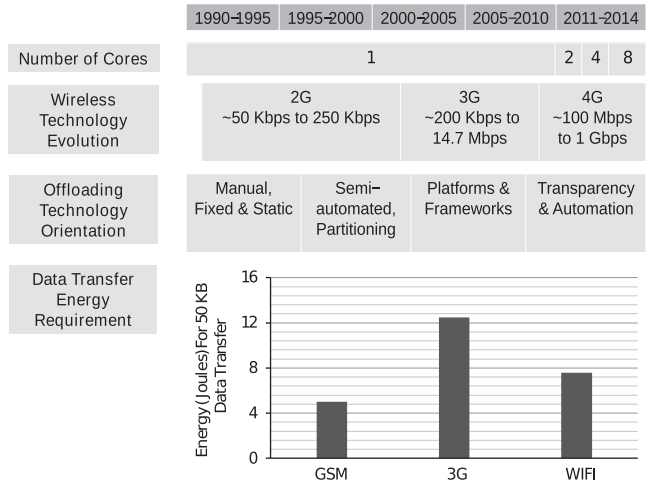


Fig. 4. Evolution of wireless technology.

Fig. 3. Most of the data offloading systems aim at storage of data to remote servers with large storage repositories. One of the objectives of the recently evolved Mobile Cloud Computing (MCC) is to provide storage facilities to the users. The synchronization of data with that existing on the cloud storage repository is also provided by MCC. Similar to data offloading, the computation offloading has also evolved to be incorporated in MCC. In general, it aims at energy minimization and performance improvement.

Figure 4 shows a quantitative and chronological evolution of several parameters related to wireless technology. The smartphones have evolved to contain multi-core based processors. Similarly, with the implementation of 3G and 4G based networks, the wireless technology is now able to offer more bandwidth than the previous generations. The orientation of offloading research has evolved from defining manual mechanisms to automated transparent offloading mechanisms. The energy requirements (Joules) as given in Balasubramanian et al. (2009) for 50 kB data transfer (download with intervals of 20 s) through GSM, 3G and

<sup>1</sup> Statistics obtained from the ACM Digital Library for duration up to July 2014.

Wi-Fi are also shown. The Wi-Fi based data transmission requires the highest amount of energy.

#### 4. Offloading architectures and approaches

We categorize computation offloading approaches into static and dynamic depending upon the time at which the decision of offloading takes place.

##### 4.1. Static offloading

As shown in Fig. 5, the static offloading approach makes use of performance prediction models or offline profiling to estimate the performance (Li et al., 2001; Xian et al., 2007; Chu et al., 2004; Ding and Li, 2003; Gurun et al., 2008; Ou et al., 2007). The application is then partitioned into client and server partitions which may subsequently be executed.

A comparison of different static offloading strategies is shown in Table 1. The comparison is performed in terms of core

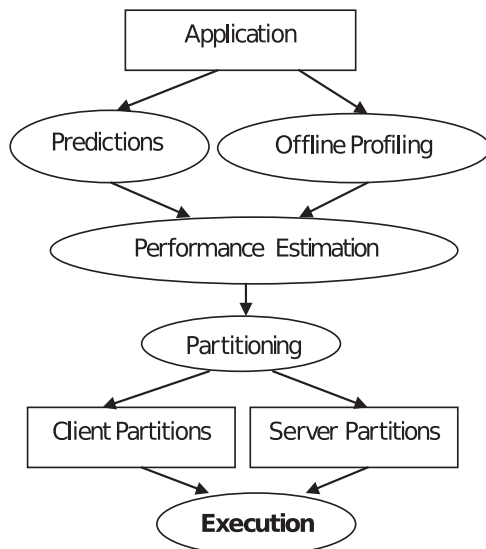


Fig. 5. Static offloading mechanism.

components (the basic component on which processing takes place), the parameters considered for offloading decision, the offloading approach and the benchmarks for which the strategy are shown to be beneficial.

The approach suggested in Li et al. (2001) first generates a cost graph for the application. The cost graph takes into account the computation time and the data to be transferred. The suggested approach then distributes the program into client and server subtasks. The data communication among the tasks being executed by hosts takes place using the primitives of *push* and *pull*. The primitives correspond to sending and receiving the modified data. The application is modeled to produce the cost graphs representing energy consumption and data communication. The sum of both these parameters is minimized by suggesting a branch-and-bound algorithm and a pruning heuristic that reduces the search space to provide a near-optimal solution. The suggested approach produces a significant improvement in execution time and energy consumption for benchmarks from Mediabench suite and *gnugo* game.

An adaptive approach presented in Xian et al. (2007) performs computation offloading by using an initial profile obtained by executing the program. If the program does not run to completion within a specified timeout, the offloading takes place and the rest of the computations are performed on some server. The minimum time required for executing the code on the mobile system is computed using the energy consumption on the local mobile processor. With the reduced energy consumption, a significant improvement in the performance is achieved for image processing benchmarks.

A framework called *Roam* which may be used for offloading of applications is suggested in Chu et al. (2004). The framework enables partitioning of an application into several components that may then be migrated to any other platform. This architecture supports heterogeneity in that the application components may be migrated to another system having a different execution environment. The approach of application offloading incorporates adaptation of three different types. The first one, *dynamic instantiation* based adaptation, partitions an application into several device dependent components. Each component has implementation for multiple platforms. The approach then takes into account the capabilities of the target system in order to select the components to be migrated. The second type, *offloading computation*, makes the

Table 1  
Comparison of static offloading strategies.

| Framework contribution   | Core component              | Parameters   | Offloading approach | Candidate applications                                  |
|--------------------------|-----------------------------|--|---------------------|---|
| Li et al. (2001)         | Cost graph                  | Computation and data transfer time                         | Static              | Mediabench & <i>gnugo</i>                               |
| Xian et al. (2007)       | Execution profile           | Energy consumption and time required for execution         | Static              | Image processing  |
| Chu et al. (2004)        | Application components      | Components categorization                                  | Static              | General applications                                    |
| Ou and Yang (2006)       | Multi-cost graph            | Computation and communication costs                        | Static              | Audio and video applications                            |
| Messer et al. (2002)     | Execution profile           | Communication cost and connectivity of nodes               | Static              | Text editor, Biomer and Voxel                           |
| Rim et al. (2006)        | Java bytecode               | Configuration based  | Static              | SciMark benchmark                                       |
| Othman and Hailes (1998) | Jobs                        | Power consumption for execution and data transfer          | Static              | General applications                                    |
| Wang and Li (2004b)      | Control flow graph          | Execution, communication, scheduling and bookkeeping costs | Static              | Image processing, speech recognition and compression    |
| Wang and Dey (2010)      | 3D rendering                | Communication and computation costs                        | Static              | Games processing  |
| Rachuri et al. (2011)    | Mobile phone sensor samples | Energy, latency and data traffic                           | Static              | Social behavior   |
| Balan et al. (2007)      | Functions based modules     | Configuration based  | Static              | Natural language, speech processing and computer vision |
| Chun et al. (2011)       | Execution profiles          | Computation cost and migration cost                        | Static              | Virus scanning and image search                         |
| Ou et al. (2007)         | Analytical model            | Surrogates coverage  | Static              | General applications                                    |
| Gurun et al. (2004)      | Performance history         | Prediction errors  | Static              | General applications                                    |
| Niu et al. (2014)        | Object relation graph       | Bandwidth, execution cost and data transfer                | Static              | Dacapo benchmark  |

applications use distributed resources by offloading components to remote servers. It is mainly required for offloading the application logic based code. The third type *transformation* makes the user interface components compatible with the target device at run-time. The decision of partitioning is however static and is made at the time of designing the application.

The application partitioning algorithm suggested in [Ou and Yang \(2006\)](#) divides the application into two main parts. The first part contains the partition that cannot be offloaded and will execute on the mobile device locally. The second part contains  $k$  partitions that can be offloaded to surrogates. The partitions are formed by modeling the computation and communication costs of the application components as a dynamic multi-cost graph. A special tightest and lightest vertex solution algorithm is then used to select a vertex in a partition. The algorithm considers the edge weights and vertex weights for partitioning. On the IBM laptop X31 and using two desktop PCs as surrogates, the application partitioning is shown to improve the performance for PI calculation, MP4 player and MP4 audio/video generation benchmarks.

A prototype platform *AIDE* suggested in [Messer et al. \(2002\)](#) makes use of three modules for profiling the application execution, partitioning and migration of code. Initially, a Java application is partitioned by providing a set of min-cut partitioning. All the partitions are then evaluated by placing one node in first partition and all others in second partition. The nodes of second partition having the highest connectivity are moved to first partition iteratively. Subsequently, the minimum cut represents partitioning with the lowest inter-partition weight with respect to the communication cost between two partitions. For a diverse set of benchmarks including the *JavaNote* (text editor), *Biomer* (molecular editor) and *Voxel* (fractal landscape), the *AIDE* platform is shown to reduce the execution time significantly.

The framework *DiET* ([Rim et al., 2006](#)) is able to make modification to Java bytecode to support offloading of methods. The mobile users request to execute an application available through service providers. The client part of the application is downloaded to the mobile device. The complex computation based methods are modified with remote procedure calls in the client part. The server reads the requests and executes the code. Moreover, the automated offloading mechanism is portable and requires no special JVM dependent instructions. For the SciMark benchmark, the suggested approach is able to produce up to 59% of speedup for the *MonteCarlo* integration method.

In [Othman and Hailes \(1998\)](#), the authors target offloading in a wireless network from a mobile device to the mobile support station (MSS). It estimates the power consumption by the CPU in case of local execution and power consumption for data/results transfer to/from the remote server together with the response time for executing on the local machine and the MSS. If it is found beneficial to use the MSS, the jobs are offloaded. Consequently, there is a significant improvement in response time for execution of different jobs offloaded to the MSS.

The strategy proposed in [Wang and Li \(2004b\)](#) implements computation offloading by partitioning the code in client and server parts. A polynomial time algorithm is suggested to achieve optimal partitioning of code for a given set of input data. For a program, a control flow graph is built where each vertex is a basic block and each edge represents dependencies. A point-to analysis is then performed to identify the memory addresses or locations during data transfer. For distribution, various constraints are used to ensure data consistency. A cost analysis that takes into account the costs required for execution, scheduling, bookkeeping and communication is used to model the problem as a minimization problem. The problem is then represented as the min-cut network flow problem and is solved using an option-clustering heuristic. On an IPAQ 3970, and a Pentium-IV based server, the suggested

offloading approach is able to reduce execution time for photo processing, graphics compression/de-compression, speech recognition and graph drawing benchmarks.

In [Wang and Dey \(2010\)](#), an approach for adapting the rendering settings for games in a mobile cloud is described. A static analysis is initially performed to select optimal settings for 3D rendering. These settings correspond to different adaptation levels where each level is associated with a total of communication and computation costs. During execution, an algorithm works to adjust the rendering settings in conformance with the existing communication and computation costs. For the game *PlaneShift* being played on a netbook, and using game server having GPU, the experimental results show an improvement in the performance in terms of the Game Mean Opinion Score (GMOS) corresponding to the gaming user experience.

A mobile phone based framework to capture the users' social behavior in a working environment is specified in [Rachuri et al. \(2011\)](#). The quantitative information such as the most sociable person in the environment and the number of interactions between two users have been useful for increasing productivity of organizations. To obtain such information, the mobile phone sensors are used to capture the behavior. The sensors sample the data at a specific rate. The samples are then processed to infer the required information. Due to the limited capability of the mobile devices, the processing is distributed among several devices. The decision of performing the computation locally or remotely is made by considering the parameters of energy, latency and data traffic. The overall task with these parameters is first divided into subtasks and a configuration for processing the task is found using the multi-criteria decision theory. With a Nokia 6120 mobile phone as a client and an Intel Xeon based server, the suggested approach is efficiently able to process the data and infer the required information.

An approach to partition the application for offloading using a language *Vivendi* is suggested in [Balan et al. \(2007\)](#). The language *Vivendi* is developed to describe the relevant specification of the application whose computation is to be offloaded. A file in the *Vivendi* language may contain the prototypes of functions that can be executed remotely. The next part of the approach incorporates *Chroma* ([Balan et al., 2003](#)) to monitor resources and predict the behavior. Subsequently, the stubs may be generated using the *Vivendi* stub generator and all function calls at corresponding points are replaced by calls to stubs. All the modules are then compiled and linked to generate an executable application. The suggested approach is able to support offloading for diverse applications including the natural language, speech and computer vision based applications.

The framework *CloneCloud* ([Chun et al., 2011](#)) facilitates the execution of a mobile application on the cloud. The *CloneCloud* initially partitions the application to make its parts execute on the mobile device and the cloud servers. A static offline analysis is performed to identify the partition. A dynamic profiler then generates profiles corresponding to different inputs. Consequently, a profile tree representing the execution traces is constructed. For each call of code, the computation cost and the migration cost in the case of local, remote or hybrid execution are computed. The optimization problem is then solved by minimizing these costs using an integer linear programming (ILP) solver. On an Android phone used as a client, and an Intel Xeon based server running mobile clones, the experimental results of clone execution show up to 20 times speedup for the applications including the virus scanning, image search and behavior profiling.

In [Ou et al. \(2007\)](#), an analytical model is presented for analyzing the performance of offloading systems. The model takes into account the distribution of surrogates and shows that in the areas well covered by surrogates, the offloading may result in

speedup in the performance. In contrast, the areas with less coverage of surrogates, the offloading does not improve the performance.

The framework *NWSLite* (Gurun et al., 2004) is used for predicting the costs of location and remote execution. Its prediction model uses a non-parametric approach. The *NWSLite* framework incorporates a large number of models each with different parameterizations. It forecasts measurements based on the performance history. The predictors are ranked with respect to the prediction errors and the best prediction model having the smallest prediction error. The *NWSLite* prediction models are executed in parallel thereby making it more efficient than the previously suggested LSQ (Noble et al., 1997) and RPF (Rudenko et al., 1999).

The authors in Niu et al. (2014) aim at improving the execution performance by using the branch-and-bound and min-cut based approaches for partitioning mobile applications. It works by performing a static analysis & profiling, followed by the generation

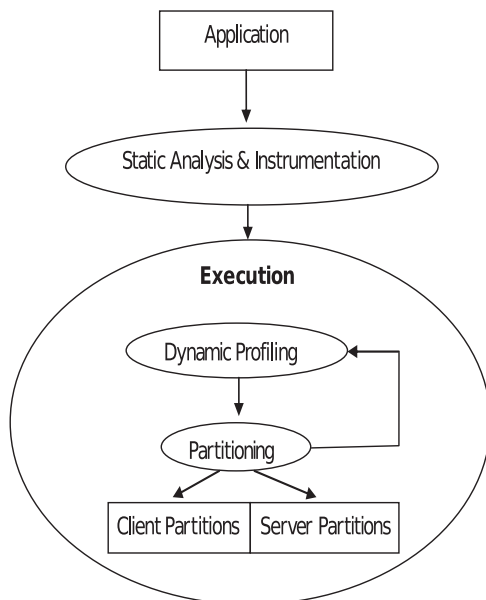


Fig. 6. Dynamic offloading mechanism.

of a weighted object relation graph (WORG), which is used to represent the objects and relations between objects. The bandwidth parameter is then used together with the WORG to partition an application into client and server parts. The branch-and-bound based algorithm produces optimal partitioning results for small applications, whereas the min-cut based approach works for large applications. Using a ThinkPad notebook for customized and the Dacapo suite benchmarks, the branch-and-bound and the min-cut based approaches produce speedups of 44.17% and 37.44%, respectively.

4.2. Dynamic offloading

As shown in Fig. 6, the dynamic offloading strategies initially perform static analysis of the code and instrumentation in order to perform dynamic/online profiling during execution (Chen et al., 2004; Chun and Maniatis, 2009; Wang and Li, 2004a; Marin, 2013; Yang et al., 2013). Based on the information obtained from dynamic profiling, the application is partitioned into client and server partitions. The execution then continues with the updated configuration.

A comparison of different dynamic offloading strategies is shown in Table 2. The comparison is performed in terms of core components, the parameters considered for offloading decision, the offloading approach and the benchmarks for which the strategy is shown to be beneficial.

In Chen et al. (2004), the authors suggest to perform compression and de-compression operations simultaneously during computation offloading. For any application requiring the data to be transferred, it reduces the penalty of data transfer. Consequently, the application performance improves if the benefit produced by the data compression (in terms of the reduced number of packets) is higher than the overall cost of data compression and de-compression. The suggested approach is shown to be effective for making decision of Java code to be compiled and executed on remote server or locally.

With the notion of augmented execution, an application may be executed on some clones of a smartphone (Chun and Maniatis, 2009). The runtime engine offloads the computation in a seamless way to another system that contains a clone of the entire system image. Consequently, the results may be integrated back to the smartphone. A special case of multiplicity based augmentation is

Table 2 Comparison of dynamic offloading strategies.

| Framework contribution          | Core component        | Parameters   | Offloading Approach | Candidate Applications                              |
|---------------------------------|-----------------------|--|---------------------|---|
| Chen et al. (2004)              | Application code      | Data transfer  | Dynamic             | Compilation and execution of Java code              |
| Chun and Maniatis (2009)        | System clones         | Fixed configuration  | Dynamic             | Data parallel applications and file system scanning |
| Wang and Li (2004a)             | Application graph     | Computation, communication, registration and scheduling costs    | Dynamic             | FFT, encode and decode benchmarks                   |
| Gu et al. (2003)                | Application graph     | Graph dependencies, network traffic, call delay and memory sizes | Dynamic             | Image and text editors                              |
| Tilevich and Smaragdakis (2002) | Application code      | Fixed configuration  | Dynamic             | Speech synthesis and MS-PowerPoint                  |
| Huerta-Canepa and Lee (2008)    | Execution profile     | Class usage and frequency  | Dynamic             | General applications                                |
| Yang et al. (2008)              | Multi-cost graph      | Communication cost and class weight                              | Dynamic             | Text recognition and translation                    |
| Nimmagadda et al. (2010)        | Real-time constraints | Network bandwidth and server speed                               | Dynamic             | Real-time surveillance                              |
| Cuervo et al. (2010)            | Application profile   | Energy consumption, bandwidth and latency                        | Dynamic             | Face recognition and games                          |
| Kremer et al. (2003)            | Application code      | Safety for remote execution                                      | Dynamic             | Face recognition                                    |
| Flinn et al. (2001)             | Application profile   | Application fidelities   | Dynamic             | General applications                                |
| Sivavakeesar et al. (2006)      | Lookup service        | Lookup latency   | Dynamic             | General applications                                |
| Wolski et al. (2008)            | Estimation model      | Network bandwidth and execution costs                            | Dynamic             | General applications                                |

presented that could work for performance improvement of data parallel applications. It requires multiple clones of the smartphone image. Similarly, a hardware based augmented execution is shown to improve the performance of scanning the file system.

In Wang and Li (2004a), the application partitioning is performed through a parametric analysis of the computation and communication costs. The problem of finding optimal partitioning is modeled as the min-cut network flow problem. The modules of the application distributed on the mobile device of the server depending upon the current value of runtime parameters. A program is first divided into modules or tasks that are executed on the server or the mobile device exclusively. A cost analysis then takes into account the computation, communication, task scheduling, and data registration costs and formulates the optimal partitioning as a single-source single-sink min-cut network flow problem. Using the mobile client HP IPAQ 3970, and a server machine having P4 processor, the results show that an effective partitioning significantly impacts the performance of several applications such as FFT, *encode* and *decode* from Mediabench and Minbench benchmarks.

An architecture of an inference engine is proposed in Gu et al. (2003) for deciding the time of offloading and the application partition to be offloaded. The inference engine employs a fuzzy model and is implemented in the AIDE framework (Messer et al., 2002). Each class of a Java application is represented as a node in a weighted graph. Each class is annotated with a flag describing whether or not the class may be offloaded to a server. The inference engine uses a min-cut based algorithm to find all 2-way cuts of the weighted graph. The nodes in the graph that may not be migrated to the surrogate are merged in the partition which will be executed on the mobile device. The other nodes are merged taking into account the dependencies and the metrics of network traffic, function call delay and memory size. The experiments performed for the evaluation of an image editor, a text editor and a molecular editor show that the suggested approach minimizes the traffic requirements while working with a very small offloading overhead.

An automated approach of partitioning a Java application for remote execution is presented in Tilevich and Smaragdakis (2002). A platform called *J-orchestra* is developed to perform replacement of the object code i.e. bytecode of method calls with the remote invocation. It divides an application into a client-server based model whose most of the I/O operations are performed on the client machine and the rest of the execution takes place on the server machine. With an IPAQ PDA, the *J-orchestra* has been shown to automatically distribute applications such as speech synthesis and MS PowerPoint.

The approach presented in Huerta-Canepa and Lee (2008) provides an adaptable offloading mechanism based on the application's execution behavior. A history of the execution pattern is maintained and is later used for making offloading decision. The static offloading policy offloads the most used classes, whereas the dynamic offloading moves only the invoked classes. The decision of offloading, i.e. *static*, *dynamic*, *no action*, or *profile* is made for each resource. Subsequently, the most common decision is opted for implementation. On PDAs, the offloading approach makes the application execute faster than local execution and is beneficial for applications with large execution times.

An offloading service for mobile handsets which may be used during mobility is presented in Yang et al. (2008). Initially, the resource information is collected and is followed by partitioning of application execution on the local system and the surrogate. The discovery of a suitable surrogate is made using the instantiation of classes for remote execution. The instrumented classes are then offloaded to the surrogates. The application partitioning uses a multi-cost graph, each of whose vertices is a class. The problem of

graph partitioning is then solved by using a  $k+1$  partitioning algorithm. The proposed algorithm takes into account the weight of one class together with the weights of one-hop weights while minimizing the communication cost. On an HP IPAQ PDA, the suggested approach is applied to the *autoTranslator* software to recognize text in German language and translate it to English. The approach performs 3–5 times better than the randomly selected and the highest transfer rate based algorithms.

In Nimmagadda et al. (2010), an approach for object recognition and tracing is presented, which may be used in the real-time surveillance systems. The approach performs computation offloading on the basis of real-time constraints. These constraints use various ranges of network bandwidth and server speed to make the offloading decision of executing code locally on a robot or remotely on a server.

The MAUI framework Cuervo et al. (2010) supports fine-grained offloading of code in an automated way. To accomplish the portability of applications, two versions are created corresponding to execution on the mobile phone and the server. The MAUI architecture contains decision engine, proxy and profiler on both the client and the server. The server part also contains the coordinator component to create an instance of the partitioned application. Initially, the methods to be offloaded are annotated by the programmer. These methods are identified by MAUI through Reflection API. Subsequently, the state of the application required for transfer or return to/from the server is identified. The MAUI profile provides feedback regarding energy consumption, bandwidth and latency, etc. to the MAUI solver that in turn decides whether or not the code should be offloaded to the server. The solver models it as an optimization problem for minimizing the energy consumption subject to various latency constraints. Using MAUI, the code offloading for face recognition, video game and chess game is shown to improve the execution time.

The application partitioning by performing code analysis is suggested in Kremer et al. (2003). The subtasks that are safe for remote execution are first identified. Subsequently, an analysis is performed to estimate the actual gains after offloading. Finally, two versions corresponding to execution on local and remote machine are generated. The suggested approach is implemented in the SUIF2 compiler (Aigner et al., 2000), and is able to achieve almost 13 times and 15 times speedup in the performance of face recognition code on Skiff and IPAQ mobile appliances.

In Flinn et al. (2001), the architecture of a framework *Spectra* is presented. The *Spectra* framework does not require the application to describe the resources to be used, instead, it can predict the application behavior for future execution. It is implemented as part of the Aura framework (Sousa, 2002) and uses the application fidelities as parameters to decide to perform execution on local and remote machines exclusively or hybridly. The CPU availability, network bandwidth, battery energy and data access costs are estimated by monitors to predict the application behavior. The *Spectra* framework then selects the best location and fidelity for application execution while taking as an input the application description and the application behavior parameters. Using a Pocket PC with an SA-1100 processor as a client and an IBM T20 Laptop as a server, the *Spectra* framework is shown to select the best option for local, remote or hybrid execution.

In Sivavakeesar et al. (2006), two strategies of service discovery for offloading applications are presented. These strategies are based on flooding and unicasting. Every device is represented by a node and is associated to a lookup server that is used to store service description. When a service is required by a node, a service lookup is performed. The scope of the search (in terms of the area) for the server machine is increased gradually if no response is received from the lookup server. With flooding, the lookup message is broadcast, in contrast to unicast, which is useful for

large environments. The experimental results show that the service discovery based approach for cyberaging applications is able to reduce the latency of the service lookup operation.

An approach for deciding offloading between the local and the remote system by making use of the bandwidth parameter is provided in [Wolski et al. \(2008\)](#). The problem of estimating the local and remote execution costs is modeled as a statistical decision problem. The remote execution cost is computed as a function of the bandwidth available for transfer of data between the local and the remote systems. The Bayesian approach is then used to solve the problem and make the prediction regarding the offloading decision.

## 5. Application partitioning for computation offloading

Together with the evolution of wireless technology, the research in the field of computation offloading has also evolved vigorously. As discussed earlier, an effective computation offloading technique may significantly impact the performance. The computation offloading incorporates various steps and analyses to ensure performance gain. One of the major steps used in computation offloading is application partitioning which distributes code for local and remote execution. The application partitioning may be categorized into static (application specific, framework based and offline profile based) and dynamic as shown in [Table 3](#), and elaborated in this section.

### 5.1. Static partitioning

For offloading computation to a remote machine, a static partitioning approach is adopted when application's code modules are fixed to be executed on local or remote machines. The static partitioning may be implemented through an application specific, a framework based or an offline profile based strategy.

For a few partitioning strategies ([Jo et al., 2014](#); [Liu et al., 2014](#); [Kovachev et al., 2014](#); [Park et al., 2014](#)), the parts of an application (such as AES encryption, image processing, multimedia services and Javascript code) are pre-defined to be executed on local or remote machines. These strategies set the portions of code depending upon the application. Similarly, for offloading strategies suggested in [Lomotey \(2013\)](#) and [Toma and Chen \(2013a\)](#), the partitioning works for [Matthews et al. \(2011\)](#) related applications and frame-based tasks, respectively, whereas in [Eom et al. \(2012\)](#), [Kovachev et al. \(2012\)](#), various performance parameters are used to fix application based partitioning. The approach given in [Imai and Varela \(2011\)](#) uses a mathematical model for improving face detection. For GPS services, the application specific partitioning uses signal processing stages and navigation methods ([Ramos et al., 2011](#)), whereas for mobile games, fixed partitioning is adopted ([Liu et al., 2010](#)). For a surveillance system, a hierarchical partitioning approach is given in [Tsiatsis et al. \(2005\)](#).

For some framework based strategies ([Gangil et al., 2013](#); [Xia et al., 2014](#); [Saarinen et al., 2012](#); [Zhang et al., 2011](#); [Silva et al., 2008](#); [Ni et al., 2005](#); [Saarinen et al., 2012](#)), the fixed partitioning mechanism is usually driven by programmers. In [Liu et al. \(2005\)](#), an operating system to support distributed execution of java bytecode through static partitioning is described. The partitioning requires programmer annotations to decide the portions of code to be distributed. Similarly, the frameworks with fixed partitioning for collaborative or coalition based execution ([Nogueira and Pinho, 2005](#); [Kurkovsky and Bhagyavati, 2004](#)) are also proposed. Different API functions to support offloading are suggested in [Mtibaa et al. \(2013\)](#). The framework proposed in [Verbelen et al. \(2012\)](#) requires the developer to annotate classes which must be offloaded. The offloading approach in [Gordon et al. \(2012\)](#) partitions the application into user interface and computation based components through the proposed framework.

The offline profile based static partitioning uses a set of parameters and evaluates them before actually executing the application. The application partitioning approach given in [Niu et al. \(2014\)](#) uses branch-and-bound and min-cut based algorithms together with the bandwidth parameter. Similarly, the genetic and machine-learning based approaches are suggested in [Balakrishnan and Tham \(2013\)](#), [Folino and Pisani \(2013\)](#), and [Eom et al. \(2013\)](#) which take into account the resource status, network parameters and data to be transferred. In [Elgazzar et al. \(2013\)](#), the operations of a web service are profiled to generate a resource consumption profile which is subsequently used for performing computation offloading. For executing Javascript code, a profiler and a points-to analysis are suggested for helping developers to decide the portions of code to be offloaded ([Wang et al., 2013](#)). The approach in [Ahnn and Potkonjak \(2013\)](#) maps application partitioning as a minimization problem while taking into account performance estimate and communication cost. Similarly, a dynamic programming based algorithm ([Toma and Chen, 2013b](#)) uses the estimated execution time for offloading tasks which satisfy a specific set of constraints. Other approaches adopted in [Gurun et al. \(2008\)](#), [Han et al. \(2008\)](#), [Ou et al. \(2007\)](#), and [Li et al. \(2002\)](#) also make use of similar parameters and conditions for partitioning applications for computation offloading.

### 5.2. Dynamic partitioning

Many offloading strategies are able to adapt partitioning of code dynamically by taking into account several parameters ([Giurgiu et al., 2012](#); [Abebe and Ryan, 2012](#); [Gao et al., 2012](#)). These parameters are evaluated using profiling and performance prediction based mechanisms which manifest the possible behavior of an application. To profile execution of an application, the code is first instrumented and then analyzed for performance prediction.

In [Chuang et al. \(2013\)](#), a programming model with an event-driven approach for providing elastic execution of applications is

**Table 3**  
Comparison of partitioning approaches adopted for computation offloading.

| Partitioning category                     | References   |
|---|--|
| Application specific static partitioning  | <a href="#">Jo et al. (2014)</a> , <a href="#">Liu et al. (2014)</a> , <a href="#">Kovachev et al. (2014)</a> , <a href="#">Park et al. (2014)</a> , <a href="#">Lomotey (2013)</a> , <a href="#">Toma and Chen (2013a)</a> , <a href="#">Eom et al. (2012)</a> , <a href="#">Kovachev et al. (2012)</a> , <a href="#">Imai and Varela (2011)</a> , <a href="#">Ramos et al. (2011)</a> , <a href="#">Liu et al. (2010)</a> , <a href="#">Zhang et al. (2009)</a> , and <a href="#">Tsiatsis et al. (2005)</a>                         |
| Framework/API based static partitioning   | <a href="#">Xia et al. (2014)</a> , <a href="#">Kurkovsky and Bhagyavati (2004)</a> , <a href="#">Gangil et al. (2013)</a> , <a href="#">Mtibaa et al. (2013)</a> , <a href="#">Verbelen et al. (2012)</a> , <a href="#">Gordon et al. (2012)</a> , <a href="#">Saarinen et al. (2012)</a> , <a href="#">Zhang et al. (2011)</a> , <a href="#">Saarinen et al. (2012)</a> , <a href="#">Silva et al. (2008)</a> , <a href="#">Ni et al. (2005)</a> , <a href="#">Liu et al. (2005)</a> , and <a href="#">Nogueira and Pinho (2005)</a> |
| Offline profile based static partitioning | <a href="#">Niu et al. (2014)</a> , <a href="#">Balakrishnan and Tham (2013)</a> , <a href="#">Elgazzar et al. (2013)</a> , <a href="#">Eom et al. (2013)</a> , <a href="#">Wang et al. (2013)</a> , <a href="#">Ahnn and Potkonjak (2013)</a> , <a href="#">Folino and Pisani (2013)</a> , <a href="#">Toma and Chen (2013b)</a> , <a href="#">Gurun et al. (2008)</a> , <a href="#">Han et al. (2008)</a> , <a href="#">Ou et al. (2007)</a> , and <a href="#">Li et al. (2002)</a>  |
| Dynamic partitioning                      | <a href="#">Chuang et al. (2013)</a> , <a href="#">Marin (2013)</a> , <a href="#">Yang et al. (2013)</a> , <a href="#">Giurgiu et al. (2012)</a> , <a href="#">Abebe and Ryan (2012)</a> , <a href="#">Zhang et al. (2012)</a> , <a href="#">Gao et al. (2012)</a> , <a href="#">Han et al. (2006a,b)</a> , <a href="#">Cai et al. (2013)</a> , <a href="#">Chuang et al. (2013)</a> , <a href="#">Trifunovic et al. (2014)</a> , and <a href="#">Shiraz et al. (2014)</a>   |



**Table 4**  
Comparison of the automation, optimization problem solving, replication granularity, fine-grained and native method call support based characteristics of the offloading frameworks.

| Framework                                    | Automation | Optimization problem solving | Replication granularity  | Fine-grained | Native method call |
|--|------------|------------------------------|--------------------------|--------------|--------------------|
| CloneCloud (Chun et al., 2011)               | High       | Highly asynchronous          | Partial threads          | Yes          | Yes                |
| MAUI (Cuervo et al., 2010)                   | Low        | Low asynchronous             | Low-level (fine-grained) | Yes          | No                 |
| SociableSense (Rachuri et al., 2011)         | Low        | Asynchronous                 | Module-level             | No           | No                 |
| Spectra (Flinn et al., 2001)                 | High       | Asynchronous                 | Task-level               | No           | No                 |
| Framework in Yang et al. (2008)              | Medium     | Asynchronous                 | Components               | No           | Yes                |
| Roam (Chu et al., 2004)                      | High       | Asynchronous                 | Component/Roamlet        | No           | No                 |
| AIDE (Messer et al., 2002)                   | Medium     | Highly asynchronous          | Class                    | No           | Yes                |
| DiET (Rim et al., 2006)                      | Medium     | Asynchronous                 | Class methods            | No           | No                 |
| J-Orchestra (Tilevich and Smaragdakis, 2002) | High       | Highly asynchronous          | Class methods            | No           | Yes                |

suggested. Its dynamic migration mechanism distributes the execution among multiple nodes depending upon the workload requirements. A framework for dynamically adapting execution on a collection of smartphones is suggested in Marin (2013). Similarly, the authors in Yang et al. (2013) propose dynamic partitioning using genetic algorithm for mobile data streams. The approach proposed in Zhang et al. (2012) initially detects movable classes and then offloads by profiling classes during execution. In Han et al. (2006b), the partitioning is mapped to min-cut problem, whereas a few components are replicated for minimizing component migration at runtime. Other offloading frameworks and mechanisms (Cai et al., 2013; Chuang et al., 2013; Trifunovic et al., 2014; Shiraz et al., 2014) use online profiles while considering various parameters for performing code partitioning dynamically.

## 6. Comparison of offloading frameworks

Table 4 describes a comparison of the automated offloading frameworks in terms of the parameters of automation, optimization problem solving, replication granularity, fine-grained offloading and native method call support. For automation, the frameworks CloneCloud, Spectra, Roam and J-Orchestra provide offloading in a highly automated manner. This requires less interaction of the programmer as compared to those having low automated offloading support. Similarly, the frameworks CloneCloud, AIDE, and J-Orchestra solve the optimization problem in a highly asynchronous manner with regard to execution of the application. The replication granularity refers to the main component that is replicated or transferred for remote execution. The fine-grained component support is provided in the CloneCloud and MAUI frameworks. Moreover, a few frameworks including the CloneCloud, framework in Yang et al. (2008), AIDE and J-Orchestra also support native method calls.

A comparison of the working mechanism in terms of the analysis performed, dynamic profiling, late binding and trusted execution of the automated frameworks is given in Table 5 (Wen et al., 2012). All the frameworks make use of a static analysis which is performed before execution of the application. The frameworks CloneCloud, MAUI, Roam and AIDE incorporate dynamic profiling to obtain information during execution of the application and perform adaptation accordingly. The late binding for offloading refers to the offloading implemented at a later time during execution of the application. It is performed by the CloneCloud, MAUI, SociableSense (Yang et al., 2008), Roam and AIDE frameworks. Currently, none of these frameworks ensures a trusted execution to provide secure, reliable and authenticated access for offloaded applications.

Table 6 (Rudenko et al., 1998) provides a comparison of the offloading frameworks in terms of their applications, trade-off parameters, optimization and dynamic adaptation strategies. The

**Table 5**  
Comparison of the static analysis, dynamic profiling, late binding and trusted execution based characteristics of the offloading frameworks.

| Framework                                    | Static analysis | Dynamic profiling | Late binding (offloading) | Trusted execution |
|--|-----------------|-------------------|---------------------------|-------------------|
| CloneCloud (Chun et al., 2011)               | Yes             | Yes               | Yes                       | No                |
| MAUI (Cuervo et al., 2010)                   | Yes             | Yes               | Yes                       | No                |
| SociableSense (Rachuri et al., 2011)         | Yes             | No                | Yes                       | No                |
| Spectra (Flinn et al., 2001)                 | Yes             | No                | No                        | No                |
| Framework in (Yang et al., 2008)             | Yes             | No                | Yes                       | No                |
| Roam (Chu et al., 2004)                      | Yes             | Yes               | Yes                       | No                |
| AIDE (Messer et al., 2002)                   | Yes             | Yes               | Yes                       | No                |
| DiET (Rim et al., 2006)                      | Yes             | No                | No                        | No                |
| J-Orchestra (Tilevich and Smaragdakis, 2002) | Yes             | No                | No                        | No                |

CloneCloud, MAUI, DiET and J-Orchestra are useful for general scientific applications, whereas the frameworks Roam and AIDE are shown to be effective for image and graphics processing. Similarly, the framework in Yang et al. (2008) and Spectra are shown to work on voice and character recognition based applications. The SociableSense is specific for applications requiring processing on social interaction in an organization. The trade-off parameters are the elements considered while optimizing the offloading decision. In general, most of the frameworks use the execution time, energy consumption and communication overhead as the main trade-off parameters. While optimizing the decision problem, different heuristics based on the min-cut,  $k+1$  partitioning, and integer linear programming (ILP) are used in most of the offloading frameworks. The frameworks also require dynamic adaptation for offloading decisions during execution of the application. The CloneCloud, MAUI and AIDE frameworks use execution pattern for runtime adaptation. Similarly, the framework in Yang et al. (2008) performs adaptation using the speedup obtained through offloading. The Roam framework uses the target device platform based runtime adaptation, whereas the DiET framework requires user configuration for runtime adaptation.

## 7. Application domains benefiting from offloading

The computation offloading has proved to be beneficial for a large number of applications lying in several domains. A domain-wise categorization of research work is shown in Table 7. A large part of the research work has targeted the applications lying in the domains of mathematics and graphics/image processing. Likewise, the games and multimedia based applications are also targeted and their number continues to grow together with the evolution

**Table 6**

Comparison of the applications, trade-off parameters, optimization and dynamic adaptation mechanisms of the offloading frameworks.

| Framework                                    | Applications                                 | Trade-off parameters                                     | Optimization strategy                         | Dynamic adaptation strategy                |
|--|--|--|---|--|
| CloneCloud (Chun et al., 2011)               | Scientific                                   | Execution speed, energy and data transfer                | Integer Linear Programming (ILP)              | Profile tree based                         |
| MAUI (Cuervo et al., 2010)                   | Scientific                                   | Energy & execution speed with data transfer              | 0-1 ILP                                       | Call graph based                           |
| SociableSense (Rachuri et al., 2011)         | Social Interaction                           | Accuracy, energy, latency and data traffic               | Multi-criteria decision theory                | Learning based                             |
| Spectra (Flinn et al., 2001)                 | Voice recognition                            | Latency, battery life and fidelity                       | Fidelity solver                               | None                                       |
| Framework in (Yang et al., 2008)             | Language translation & character recognition | Response time, communication, CPU and memory             | (k + 1) partitioning algorithm                | Speedup based                              |
| Roam (Chu et al., 2004)                      | Games & graphics                             | Capabilities of target devices and user interface design | Component-based partitioning                  | Target device capabilities based mechanism |
| AIDE (Messer et al., 2002)                   | Image and text processing                    | Processor load, memory and communication                 | Min-cut based heuristic                       | Execution graph based                      |
| DIET (Rim et al., 2006)                      | Mathematical applications                    | User directives based                                    | User configuration based                      | User configuration based                   |
| J-Orchestra (Tilevich and Smaragdakis, 2002) | General applications                         | Input/output, disk processing and GUI                    | User directives based parameters of I/O usage | None                                       |

**Table 7**

Domain-wise categorization of the research work related to computation offloading.

|  |  |
|--|--|
| Multimedia                                 | Li et al. (2001), Wang and Li (2004a), and Ou and Yang (2006)  |
| Games                                      | Li et al. (2001), Chu et al. (2004), Cuervo et al. (2010), and Wang and Dey (2010)   |
| Graphics and image processing              | Xian et al. (2007), Chen et al. (2004), Gurun et al. (2004), Messer et al. (2002), Gu et al. (2003), Yang et al. (2008), Wang and Li (2004b), Wang and Dey (2010), Balan et al. (2007), and Chun et al. (2011)                       |
| Mathematical computations                  | Wang and Li (2004a), Chen et al. (2004), Chun and Maniatis (2009), Ou and Yang (2006), Rim et al. (2006), and Wang and Li (2004b)  |
| Artificial Intelligence based applications | Chen et al. (2004), Tilevich and Smaragdakis (2002), Yang et al. (2008), Nimmagadda et al. (2010), Cuervo et al. (2010), Kremer et al. (2003), Wang and Li (2004b), Flinn et al. (2001), Balan et al. (2007), and Chun et al. (2011) |
| Health & social applications               | Matthews et al. (2011), Kundu et al. (2007), and Rachuri et al. (2011)   |
| Database, file system or GPS processing    | Chen et al. (2004), Chun and Maniatis (2009), and Mtibaa et al. (2013)   |

of wireless technology. The applications related to Artificial Intelligence and social behavior are also being offloaded as they involve complex learning based computations. The applications with database processing, file system and GPS processing have also been implemented through offloading to improve their performance.

## 8. Current challenges for effective computation offloading

Despite the long term evolution of the offloading techniques, several issues are yet to be resolved. The most challenging issues including partitioning, automated transparency & portability, security, and application requirements are discussed below together with their possible solutions.

### 8.1. Partitioning

The computation offloading requires the application code to be partitioned into client and server parts for local and remote execution, respectively. The partitioning takes into account several parameters including costs of data transfer and computation time, however the optimal partitioning is an NP-complete problem. Consequently, different heuristics with fixed constraints are employed in the partitioning strategies.

For an effective offloading implementation, the partitioning problem needs to be solved in a quasi-automated manner requiring directives from the programmer as well as automated distribution of modules. In this regard, the scheduling techniques for heterogeneous systems (Sih and Lee, 1993; Khan, 2012; Topcuouglu et al., 2002) may be incorporated to minimize the total execution time.

### 8.2. Automated transparency & portability

The frameworks implemented for computation offloading yet lack the automated transparency so that the surrounding environment is detected and the computation offloading takes place in a seamless manner (Sanaei et al., 2012; Sanaei et al., 2014; Cui et al., 2013; Chuang et al., 2013; Gordon et al., 2012). This is a complex task as it requires an implementation of a standard protocol that will perform lookup services and other functionalities depending upon the environment while taking into account its constraints. An implementation of the standard protocol for a diverse collection of devices and environments will render it portability as well.

### 8.3. Security

With computations being offloaded to remote machines/servers, the security of data and environment for the remote systems needs to be ensured (Winkler, 2011; Sanaei et al., 2014; Khan et al., 2014a,b; Kumar et al., 2013; Zhang et al., 2009). This requires restraining the types of operations that may be offloaded for remote execution. A limited set of permissible operations may be provided by implementing a virtual machine and making the remote component execute in the environment provided by the virtual machine (Goldberg et al., 1996). Moreover, different authorization and authentication mechanisms may be incorporated in order to ensure security of data on the cloud (Khalid et al., 2013; Antonopoulos and Gillam, 2010).

### 8.4. Application requirements

The applications being executed on mobile devices are not only growing in size but also in terms of complex operations. The widely used multimedia applications including the VoIP, online

streaming, and video/audio chat require the mobile devices to improve the energy requirements, graphics rendering and the execution time. Moreover, these applications require real-time processing. Consequently, it is not possible to offload all the modules remotely. In this regard, the caching techniques and implementation of a specialized hardware such as a Digital Signal Processor (DSP) (Lin et al., 2004) or a System-on-Chip (SoC) (Holzenspies et al., 2007) may be beneficial for an effective offloading.

## 9. Conclusion

This paper presents a comprehensive survey of the research work conducted on computation offloading which aims at performance improvement of applications executing on the resource constrained mobile devices. The limited resources of mobile devices require the intensive computations to be offloaded in order to mitigate the issues of slow execution and low energy. Some of the offloading strategies work in a fixed static manner while others are able to perform offloading in accordance with the dynamic behavior of the application. We perform a comparative analysis of these strategies as well as the automated frameworks implemented to support computation offloading.

We also survey the evolution of mobile technologies and also compare different partitioning mechanisms used for distributing code between local and remote machines. The research work is also categorized in terms of the application domains for which the computation offloading is shown to be effective. Moreover, the main issues related to computation offloading: partitioning, automated transparency & portability, security, and application requirements are discussed, and their possible solutions are also proposed.

## References

- Abebe E, Ryan C. Adaptive application offloading using distributed abstract class graphs in mobile environments. *J Syst Softw* 2012;85(12):2755–69.
- Ahn JH, Potkonjak M. mHealthMon: toward energy-efficient and distributed mobile health monitoring using parallel offloading. *J Med Syst* 2013;37(5):1–11.
- Aigner G, Diwan A, Heine D, Lam M, Moore D, Murphy B, et al., The suif2 compiler infrastructure; 2000.
- Antonopoulos N, Gillam L. *Cloud computing: principles, systems and applications*. 1st ed. London, UK: Springer Publishing Company, Incorporated; 2010.
- Balakrishnan P, Tham C-K. Energy-efficient mapping and scheduling of task interaction graphs for code offloading in mobile cloud computing. In: Proceedings of the 2013 IEEE/ACM sixth international conference on utility and cloud computing, UCC '13. Washington, DC, USA: IEEE Computer Society; 2013. p. 34–41.
- Balan RK, Satyanarayanan M, Park SY, Okoshi T. Tactics-based remote execution for mobile computing. In: Proceedings of the first international conference on mobile systems, applications and services, MobiSys '03. New York, NY, USA: ACM; 2003. p. 273–86.
- Balan RK, Gergle D, Satyanarayanan M, Herbsleb J. Simplifying cyber foraging for mobile devices. In: Proceedings of the fifth international conference on mobile systems, applications and services, MobiSys '07. New York, NY, USA: ACM; 2007. p. 272–85.
- Balasubramanian N, Balasubramanian A, Venkataramani A. Energy consumption in mobile phones: a measurement study and implications for network applications. In: Proceedings of the ninth ACM SIGCOMM conference on Internet measurement conference, IMC '09. New York, NY, USA: ACM; 2009. p. 280–93.
- Barbera M, Kosta S, Mei A, Stefa J. To offload or not to offload? the bandwidth and energy costs of mobile cloud computing. In: INFOCOM, 2013 Proceedings. IEEE, Turin, Italy; 2013. p. 1285–93. <http://dx.doi.org/10.1109/INFOCOM.2013.6566921>.
- Berl A, Gelenbe E, Di Girolamo M, Giuliani G, De Meer H, Dang MQ, et al. Energy-efficient cloud computing. *Comput J* 2010;53(7):1045–51. <http://dx.doi.org/10.1093/comjnl/bxp080>.
- Cai H, Zhang W, Zhang Y, Huang G. Sm@rt offloader: supporting adaptive computation offloading for android applications. In: Proceedings demo & #38; poster track of ACM/IFIP/USENIX international middleware conference, MiddlewareDPT '13. New York, NY, USA: ACM; 2013. p. 3:1–3:2.
- Chen G, Kang B-T, Kandemir M, Vijaykrishnan N, Irwin M, Chandramouli R. Studying energy trade offs in offloading computation/compilation in java-enabled mobile devices. *IEEE Trans Parallel Distrib Syst* 2004;15(9):795–809.
- Chu H-h, Song H, Wong C, Kurakake S, Katagiri M. ROAM, a seamless application framework. *J Syst Softw* 2004;69(3):209–26.
- Chuang W-C, Sang B, Yoo S, Gu R, Kulkarni M, Killian C. Eventwave: programming model and runtime support for tightly-coupled elastic cloud applications. In: Proceedings of the fourth annual symposium on cloud computing, SOCC '13. New York, NY, USA: ACM; 2013. p. 21:1–21:16.
- Chun B-G, Maniatis P. Augmented smartphone applications through clone cloud execution. In: Proceedings of the 12th conference on hot topics in operating systems, HotOS'09. Berkeley, CA, USA: USENIX Association; 2009. p. 8.
- Chun B-G, Ihm S, Maniatis P, Naik M, Patti A. CloneCloud: elastic execution between mobile device and cloud. In: Proceedings of the sixth conference on computer systems, EuroSys '11. New York, NY, USA: ACM; 2011. p. 301–14.
- Cuervo E, Balasubramanian A, Cho D-k, Wolman A, Saroiu S, Chandra R, et al. MAUI: making smartphones last longer with code offload. In: Proceedings of the eighth international conference on mobile systems, applications, and services, MobiSys '10. New York, NY, USA: ACM; 2010. p. 49–62. (<http://dx.doi.org/10.1145/1814433.1814441>).
- Cui Y, Ma X, Wang H, Stojmenovic I, Liu J. A survey of energy efficient wireless transmission and modeling in mobile cloud computing. *Mob Netw Appl* 2013;18(1):148–55. <http://dx.doi.org/10.1007/s11036-012-0370-6>.
- Ding Y, Li Z. A compiler analysis of interprocedural data communication. In: Proceedings of the 2003 ACM/IEEE conference on supercomputing, SC '03. New York, NY, USA: ACM; 2003. p. 11.
- Dinh HT, Lee C, Niyato D, Wang P. A survey of mobile cloud computing: architecture, applications, and approaches. *Wireless Commun Mob Comput* 2013;13(18):1587–611. <http://dx.doi.org/10.1002/wcm.1203>.
- Elgazzar K, Martin P, Hassanein HS. Empowering mobile service provisioning through cloud assistance. In: Proceedings of the 2013 IEEE/ACM sixth international conference on utility and cloud computing, UCC '13. Washington, DC, USA: IEEE Computer Society; 2013. p. 9–16.
- Eom H, Stjute P, Figueiredo R, Tickoo O, Ilikkal R, Iyer R. Snarf: a social networking-inspired accelerator remoting framework. In: Proceedings of the first edition of the MCC workshop on mobile cloud computing, MCC '12. New York, NY, USA: ACM; 2012. p. 29–34.
- Eom H, Juste PS, Figueiredo R, Tickoo O, Ilikkal R, Iyer R. Machine learning-based runtime scheduler for mobile offloading framework. In: Proceedings of the 2013 IEEE/ACM sixth international conference on utility and cloud computing, UCC '13. Washington, DC, USA: IEEE Computer Society; 2013. p. 17–25.
- Flinn J, Narayanan D, Satyanarayanan M. Self-tuned remote execution for pervasive computing. In: Proceedings of the eighth workshop on hot topics in operating systems, 2001; 2001. p. 61–6.
- Folino G, Pisani FS. A framework for modeling automatic offloading of mobile applications using genetic programming. In: Proceedings of the 16th European conference on applications of evolutionary computation, EvoApplications'13. Berlin, Heidelberg: Springer-Verlag; 2013. p. 62–71.
- Fontana R, Mparmpopoulou K, Grosso P, Van der Veldt K. Monitoring greenclouds: evaluating the trade-off between performance and energy consumption in das-4; 2013.
- Gangil A, Dhurandher SK, Obaidat MS, Singh V, Bhatia S. Moclo: a cloud framework for mobile devices. In: Proceedings of the 2013 IEEE international conference on green computing and communications and IEEE Internet of things and IEEE cyber, physical and social computing, GREENCOM-ITHINGS-CPSCOM '13. Washington, DC, USA: IEEE Computer Society; 2013. p. 632–7.
- Gao B, He L, Liu L, Li K, Jarvis SA. From mobiles to clouds: developing energy-aware offloading strategies for workflows. In: Proceedings of the 2012 ACM/IEEE 13th international conference on grid computing, GRID '12. Washington, DC, USA: IEEE Computer Society; 2012. p. 139–46.
- Ge Y, Zhang Y, Qiu Q, Lu Y-H. A game theoretic resource allocation for overall energy minimization in mobile cloud computing system. In: Proceedings of the 2012 ACM/IEEE international symposium on low power electronics and design, ISLPED '12. New York, NY, USA: ACM; 2012. p. 279–84.
- Giurgiu I, Riva O, Alonso G. Dynamic software deployment from clouds to mobile devices. In: Proceedings of the 13th international middleware conference, Middleware '12. New York, NY, USA: Springer-Verlag New York, Inc.; 2012. p. 394–414.
- Goldberg I, Wagner D, Thomas R, Brewer EA. A secure environment for untrusted helper applications confining the wily hacker. In: Proceedings of the sixth conference on USENIX security symposium, focusing on applications of cryptography, SSYM '96, vol. 6. Berkeley, CA, USA: USENIX Association; 1996. p. 1.
- Gordon MS, Jamshidi DA, Mahlke S, Mao ZM, Chen X. Comet: code offload by migrating execution transparently. In: Proceedings of the 10th USENIX conference on operating systems design and implementation, OSDI'12. Berkeley, CA, USA: USENIX Association; 2012. p. 93–106.
- Gu X, Nahrstedt K, Messer A, Greenberg I, Milojicic D. Adaptive offloading inference for delivering applications in environments. In: Proceedings of the first IEEE international conference on pervasive computing and communications, PERCOM '03. Washington, DC, USA: IEEE Computer Society; 2003. p. 107.
- Gurun S, Krintz C, Wolski R. NWSLite: a light-weight prediction utility for mobile devices. In: Proceedings of the second international conference on mobile systems, applications, and services, MobiSys '04. New York, NY, USA: ACM; 2004. p. 2–11.
- Gurun S, Wolski R, Krintz C, Nurmi D. On the efficacy of computation offloading decision-making strategies. *Int J High Perform Comput Appl* 2008;22(4):460–79.

- Han S, Zhang S, Zhang Y. Energy saving of mobile devices based on component migration and replication in pervasive computing. In: Proceedings of the third international conference on ubiquitous intelligence and computing, UIC'06. Berlin, Heidelberg: Springer-Verlag; 2006. p. 637–47.
- Han S, Zhang S, Zhang Y. A generic software partitioning algorithm for pervasive computing. In: Proceedings of the first international conference on wireless algorithms, systems, and applications, WASA'06. Berlin, Heidelberg: Springer-Verlag; 2006. p. 57–68.
- Han S, Zhang S, Cao J, Wen Y, Zhang Y. A resource aware software partitioning algorithm based on mobility constraints in pervasive grid environments. *Future Gener Comput Syst* 2008;24(6):512–29.
- Holzspies P, Smit GJM, Kuper J. Mapping streaming applications on a reconfigurable MPSoC platform at run-time. In: International symposium on system-on-chip, 2007; 2007. p. 1–4.
- Hong Y-J, Kumar K, Lu Y-H. Energy efficient content-based image retrieval for mobile systems. In: IEEE international symposium on circuits and systems, ISCAS 2009; 2009. p. 1673–76.
- Huerta-Canepa G, Lee D. An adaptable application offloading scheme based on application behavior. In: 22nd international conference on advanced information networking and applications—workshops, 2008. AINAW 2008; 2008. p. 387–92.
- Imai S, Varela CA. Light-weight adaptive task offloading from smartphones to nearby computational resources. In: Proceedings of the 2011 ACM symposium on research in applied computation, RACS '11. New York, NY, USA: ACM; 2011. p. 146–52.
- Jo H, Hong S-T, Chang J-W, Choi DH. Offloading data encryption to gpu in database systems. *J Supercomput* 2014;69(1):375–94.
- Juntunen A, Kempainen M, Luukkainen S. Mobile computation offloading—factors affecting technology evolution. In: 2012 International conference on mobile business (ICMB 2012), Delft, The Netherlands; 21–22 June 2012, p. 12.
- Khalid U, Ghafoor A, Irum M, Shibli MA. Cloud based secure and privacy enhanced authentication & authorization protocol. *Proc Comput Sci* 2013;22:680–688; 17th international conference in knowledge based and intelligent information and engineering systems—KES2013. (<http://www.sciencedirect.com/science/article/pii/S1877050913009423>).
- Khan AU, Othman M, Ali M, Khan AN, Madani SA. Pirax: framework for application piracy control in mobile cloud environment. *J Supercomput* 2014a;68(2):753–76.
- Khan AuR, Othman M, Madani SA, Khan SU. A survey of mobile cloud computing application models. *IEEE Commun Surv Tutor* 2014b;16(1):393–413.
- Khan MA. Scheduling for heterogeneous systems using constrained critical paths. *Parallel Comput* 2012;38(4–5):175–93 <http://www.sciencedirect.com/science/article/pii/S0167819112000105>.
- Kovachev D, Yu T, Klamra R. Adaptive computation offloading from mobile devices into the cloud. In: Proceedings of the 2012 IEEE 10th international symposium on parallel and distributed processing with applications, ISPA '12. Washington, DC, USA: IEEE Computer Society; 2012. p. 784–91.
- Kovachev D, Cao Y, Klamra R. Building mobile multimedia services: a hybrid cloud computing approach. *Multimed Tools Appl* 2014;70(2):977–1005.
- Kremer U, Hicks J, Reh J. A compilation framework for power and energy management on mobile computers. In: Proceedings of the 14th international conference on languages and compilers for parallel computing, LCPC'01. Berlin, Heidelberg: Springer-Verlag; 2003. p. 115–31.
- Kumar K, Lu YH. Cloud computing for mobile users: can offloading computation save energy?. *Computer* 2010;43(4):51–6.
- Kumar K, Liu J, Lu YH, Bhargava B. A survey of computation offloading for mobile systems. *Mob Netw Appl* 2013;18(1):129–40.
- Kundu S, Mukherjee J, Majumdar AK, Majumdar B, SekharRay S. Algorithms and heuristics for efficient medical information display in pda. *Comput Biol Med* 2007;37(9):1272–82.
- Kurkovsky S, Bhagyavati MS, Ray A. A collaborative problem-solving framework for mobile devices. In: Proceedings of the 42nd annual southeast regional conference, ACM-SE 42. New York, NY, USA: ACM; 2004. p. 5–10.
- Li Z, Wang C, Xu R. Computation offloading to save energy on handheld devices: a partition scheme. In: Proceedings of the 2001 international conference on compilers, architecture, and synthesis for embedded systems, CASES '01. New York, NY, USA: ACM; 2001. p. 238–46.
- Li Z, Wang C, Xu R. Task allocation for distributed multimedia processing on wirelessly networked handheld devices. In: Proceedings of the 16th international parallel and distributed processing symposium, IPDPS '02. Washington, DC, USA: IEEE Computer Society; 2002. p. 312.
- Lin L-C, Lin T-J, Lee C-C, Chao C-M, Chen S-K, Liu C-H, et al. A novel programmable digital signal processor for multimedia applications. In: The 2004 IEEE Asia-Pacific conference on circuits and systems, 2004. Proceedings, vol. 1; 2004. p. 121–4. (<http://dx.doi.org/10.1109/APCCAS.2004.1412707>).
- Liu H, Roeder T, Walsh K, Barr T, Sire EG. Design and implementation of a single system image operating system for ad hoc networks. In: Proceedings of the third international conference on mobile systems, applications, and services, MobiSys '05. New York, NY, USA: ACM; 2005. p. 149–62.
- Liu J, Kumar K, Lu Y-H. Tradeoff between energy savings and privacy protection in computation offloading. In: Proceedings of the 16th ACM/IEEE International Symposium on Low Power Electronics and Design, ISLPED '10. New York, NY, USA: ACM; 2010. p. 213–8.
- Liu W, Chen J-J, Toma A, Kuo T-W, Deng Q. Computation offloading by using timing unreliable components in real-time systems. In: Proceedings of The 51st annual design automation conference on design automation conference, DAC '14. New York, NY, USA: ACM; 2014. p. 39:1–39:6.
- Lomotey RK, Deters R. Reliable services composition for mobile consumption in mhealth. In: Proceedings of the fifth international conference on management of emergent digital EcoSystems, MEDES '13. New York, NY, USA: ACM. 2013. p. 182–6.
- Lu Y, Zhou B, Tung L-C, Gerla M, Ramesh A, Nagaraja L. Energy-efficient content retrieval in mobile cloud. In: Proceedings of the second ACM SIGCOMM workshop on mobile cloud computing, MCC '13. New York, NY, USA: ACM; 2013. p. 21–6.
- Marin R-C. Hybrid contextual cloud in ubiquitous platforms comprising of smart-phones. *Int J Intell Syst Technol Appl* 2013;12(1):4–17.
- Matthews J, Chang M, Feng Z, Srinivas R, Gerla M. Powersense: power aware dengue diagnosis on mobile phones. In: Proceedings of the first ACM workshop on mobile systems, applications, and services for healthcare, mHealthSys '11. New York, NY, USA: ACM; 2011. p. 6:1–6:6.
- Mayo RN, Ranganathan P. Energy consumption in mobile devices: why future systems need requirements-aware energy scale-down. In: Proceedings of the third international conference on power-aware computer systems, PACS'03. Berlin, Heidelberg: Springer-Verlag; 2004. p. 26–40.
- Messer A, Greenberg I, Bernadat P, Milojicic D, Chen D, Giulii TJ, et al. Towards a distributed platform for resource-constrained devices. In: Proceedings of the 22nd international conference on distributed computing systems (ICDCS'02), ICDCS '02. Washington, DC, USA: IEEE Computer Society; 2002. p. 43–51.
- Miettinen AP, Hirvisalo V. Energy-efficient parallel software for mobile hand-held devices. In: Proceedings of the first USENIX conference on hot topics in parallelism, HotPar'09. Berkeley, CA, USA: USENIX Association; 2009. p. 12.
- Miettinen AP, Nurminen JK. Energy efficiency of mobile clients in cloud computing. In: Proceedings of the second USENIX conference on hot topics in cloud computing, HotCloud'10. Berkeley, CA, USA: USENIX Association; 2010. p. 4.
- Mtibaa A, Fahim A, Harras KA, Ammar MH. Towards resource sharing in mobile device clouds: power balancing across mobile devices. *SIGCOMM Comput Commun Rev* 2013;43(4):51–6.
- Ni Y, Kremer U, Stere A, Iftode L. Programming ad-hoc networks of mobile and resource-constrained devices. *SIGPLAN Not* 2005;40(6):249–60.
- Nimmagadda Y, Kumar K, Lu Y-H. Energy-efficient image compression in mobile devices for wireless transmission. In: Proceedings of the 2009 IEEE international conference on multimedia and expo, ICME'09. Piscataway, NJ, USA: IEEE Press; 2009. p. 1278–81.
- Nimmagadda Y, Kumar K, Lu Y-H, Lee CSG. Real-time moving object recognition and tracking using computation offloading. In: 2010 IEEE/RSJ international conference on intelligent robots and systems (IROS); 2010. p. 2449–55.
- Niu J, Song W, Atiquzzaman M. Bandwidth-adaptive partitioning for distributed execution optimization of mobile applications. *J Netw Comput Appl* 2014;37:334–47. <http://dx.doi.org/10.1016/j.jnca.2013.03.007>.
- Noble BD, Satyanarayanan M, Narayanan D, Tilton JE, Flinn J, Walker KR. Agile application-aware adaptation for mobility. *SIGOPS Oper Syst Rev* 1997;31(5):276–87.
- Nogueira L, Pinho LM. Dynamic qos-aware coalition formation. In: Proceedings of the 19th IEEE international parallel and distributed processing symposium (IPDPS'05)—workshop 2, vol. 03, IPDPS '05. Washington, DC, USA: IEEE Computer Society; 2005. p. 135:1.
- Nurminen J. Parallel connections and their effect on the battery consumption of a mobile phone. In: Consumer communications and networking conference (CCNC), 2010 seventh IEEE; 2010. p. 1–5.
- Othman M, Hailles S. Power conservation strategy for mobile computers using load sharing. *SIGMOBILE Mob. Comput Commun Rev* 1998;2(1):44–51. <http://dx.doi.org/10.1145/584007.584011>.
- Ou S, Yang K, Liotta A. An adaptive multi-constraint partitioning algorithm for offloading in pervasive systems. In: Proceedings of the fourth annual IEEE international conference on pervasive computing and communications, PERCOM '06. Washington, DC, USA: IEEE Computer Society; 2006. p. 116–25.
- Ou S, Yang K, Zhang J. An effective offloading middleware for pervasive services on mobile devices. *Pervasive Mob Comput* 2007;3(4):362–85.
- Ou S, Yang K, Liotta A, Hu L. Performance analysis of offloading systems in mobile wireless environments. In: IEEE international conference on communications, 2007. ICC '07; 2007. p. 1821–26. (<http://dx.doi.org/10.1109/ICC.2007.304>).
- Ou S, Yang K, Liotta A, Hu L. Performance analysis of offloading systems in mobile wireless environments. In: IEEE international conference on communications, 2007. ICC '07; 2007. p. 1821–6.
- Park S, Chen Q, Han H, Yeom HY. Design and evaluation of mobile offloading system for web-centric devices. *J Netw Comput Appl* 2014;40:105–15. <http://dx.doi.org/10.1016/j.jnca.2013.08.006>.
- Rachuri KK, Mascolo C, Musolesi M, Rentfrow PJ. Sociablesense: Exploring the trade-offs of adaptive sampling and computation offloading for social sensing. In: Proceedings of the 17th annual international conference on mobile computing and networking, MobiCom '11. New York, NY, USA: ACM; 2011. p. 73–84.
- Rim H, Kim S, Kim Y, Han H. Transparent method offloading for slim execution. In: 2006 1st international symposium on Wireless Pervasive Computing; 2006. p. 1–6.
- Ramos HS, Zhang T, Liu J, Priyantha NB, Kansal A. Leap: a low energy assisted gps for trajectory-based services. In: Proceedings of the 13th international conference on ubiquitous computing, UbiComp '11. New York, NY, USA: ACM; 2011. p. 335–44.

- Rudenko A, Reiher P, Popek GJ, Kuenning GH. Saving portable computer battery power through remote process execution. *SIGMOBILE Mob. Comput Commun Rev* 1998;2(1):19–26. <http://dx.doi.org/10.1145/584007.584008>.
- Rudenko A, Reiher P, Popek GJ, Kuenning GH. The remote processing framework for portable computer power saving. In: Proceedings of the 1999 ACM symposium on applied computing, SAC '99. New York, NY, USA: ACM; 1999. p. 365–72.
- Saariainen A, Siekkinen M, Xiao Y, Nurminen JK, Kempainen M, Hui P. Smartdirt: offloading popular apps to save energy. *SIGCOMM Comput Commun Rev* 2012;42(4):297–8.
- Saariainen A, Siekkinen M, Xiao Y, Nurminen JK, Kempainen M, Hui P. Can offloading save energy for popular apps?. In: Proceedings of the seventh ACM international workshop on mobility in the evolving Internet architecture, MobiArch '12. New York, NY, USA: ACM; 2012. p. 3–10.
- Sanaei Z, Abolfazli Z, Gani A, Khokhar RH. 2012. Tripod of requirements in horizontal heterogeneous mobile cloud computing. *CoRR abs/1205.3247*.
- Satyanarayanan M, Bahl P, Caceres R, Davies N. The case for vm-based cloudlets in mobile computing. *IEEE Pervasive Comput* 2009;8(4):14–23.
- Shiraz M, Ahmed E, Gani A, Han Q. Investigation on runtime partitioning of elastic mobile applications for mobile cloud computing. *J Supercomput* 2014;67(1):84–103.
- Sih GC, Lee EA. A compile-time scheduling heuristic for interconnection-constrained heterogeneous processor architectures. *IEEE Trans Parallel Distrib Syst* 1993;4(2):175–87.
- Silva JAN, Veiga L, Ferreira P. Spade: scheduler for parallel and distributed execution from mobile devices. In: Proceedings of the sixth international workshop on middleware for pervasive and ad-hoc computing, MPAC '08. New York, NY, USA: ACM; 2008. p. 25–30.
- Sinha K, Kulkarni M. Techniques for fine-grained, multi-site computation offloading. In: Proceedings of the 2011 11th IEEE/ACM international symposium on cluster, cloud and grid computing, CCGRID '11. Washington, DC, USA: IEEE Computer Society; 2011. p. 184–94. (<http://dx.doi.org/10.1109/CCGrid.2011.69>).
- Sivavakeesar S, Gonzalez O, Pavlou G. Service discovery strategies in ubiquitous communication environments. *IEEE Commun Mag* 2006;44(9):106–13. <http://dx.doi.org/10.1109/MCOM.2006.1705986>.
- Sousa JaP, Garland D. Aura: an architectural framework for user mobility in ubiquitous computing environments. In: Proceedings of the IFIP 17th world computer congress—TC2 Stream/3rd IEEE/IFIP conference on software architecture: system design, development and maintenance, WICSA 3. Denter, The Netherlands: Kluwer B.V.; 2002. p. 29–43. (<http://dl.acm.org/citation.cfm?id=646546.693943>).
- StatCounter, Comparison from 2011 to 2014. Stat Counter—Global Stats; 2014. (<http://gs.statcounter.com/#desktop+mobile+tablet-comparison-ww-yearly-2011-2014>).
- Tilevich E, Smaragdakis Y. J-Orchestra: automatic java application partitioning. In: Magnusson B, editor. *ECOOP. Lecture notes in computer science, vol. 2374*. London, UK: Springer; 2002. p. 178–204.
- Toma A, Chen J-J. Computation offloading for frame-based real-time tasks with resource reservation servers. In: Proceedings of the 2013 25th Euromicro conference on real-time systems, ECRTS '13. Washington, DC, USA: IEEE Computer Society; 2013. p. 103–12.
- Toma A, Chen J-J. Computation offloading for real-time systems. In: Proceedings of the 28th annual ACM symposium on applied computing, SAC '13. New York, NY, USA: ACM; 2013. p. 1650–1.
- Topcuoglu H, Hariri S, Wu M-y. Performance-effective and low-complexity task scheduling for heterogeneous computing. *IEEE Trans Parallel Distrib Syst* 2002;13(3):260–74.
- Trifunovic S, Picu A, Hossmann T, Hummel KA. Adaptive role switching for fair and efficient battery usage in device-to-device communication. *SIGMOBILE Mob Comput Commun Rev* 2014;18(1):25–36.
- Tsiatsis V, Kumar R, Srivastava MB. Computation hierarchy for in-network processing. *Mob Netw Appl* 2005;10(4):505–18.
- Verbelen T, Simoens P, DeTurck F, Dhoedt B. Aiolos: middleware for improving mobile application performance through cyber foraging. *J. Syst. Softw* 2012;85(11):2629–39.
- Wang S, Dey S. Rendering adaptation to address communication and computation constraints in cloud mobile gaming. In: Global telecommunications conference (GLOBECOM 2010), 2010. IEEE, Miami, FL, USA; 2010. p. 1–6.
- Wang C, Li Z. Parametric analysis for adaptive computation offloading. *SIGPLAN Not* 2004a;39(6):119–30.
- Wang C, Li Z. A computation offloading scheme on handheld devices. *J Parallel Distrib Comput* 2004b;64(6):740–6. <http://dx.doi.org/10.1016/j.jpdc.2003.10.005>.
- Wang X, Liu X, Huang G, Liu Y. Appmobicloud: Improving mobile web applications by mobile-cloud convergence. In: Proceedings of the fifth Asia-Pacific symposium on internetware, Internetware '13. New York, NY, USA: ACM; 2013. p. 14:1–14:10.
- Wen Y, Zhang W, Luo H. Energy-optimal mobile application execution: Taming resource-poor mobile devices with cloud clones. In: INFOCOM, 2012 Proceedings IEEE; 2012. p. 2716–20.
- Winkler JR. *Securing the cloud: cloud computer security techniques and tactics*. Waltham, MA, USA: Syngress Publishing; 2011.
- Wolski R, Gurun S, Krintz C, Nurmi D. Using bandwidth data to make computation offloading decisions. In: IEEE international symposium on parallel and distributed processing, 2008. IPDPS 2008; 2008. p. 1–8. (<http://dx.doi.org/10.1109/IPDPS.2008.4536215>).
- Xia F, Ding F, Li J, Kong X, Yang LT, Ma J. Phone2cloud: exploiting computation offloading for energy saving on smartphones in mobile cloud computing. *Inf Syst Front* 2014;16(1):95–111.
- Xian C, Lu Y-H, Li Z. Adaptive computation offloading for energy conservation on battery-powered systems. In: 2007 international conference on parallel and distributed systems, vol. 2; 2007. p. 1–8.
- Yang K, Ou S, Chen H-H. On effective offloading services for resource-constrained mobile devices running heavier mobile internet applications. *IEEE Commun Mag* 2008;46(1):56–63.
- Yang K, Ou S, Chen H-H. On effective offloading services for resource-constrained mobile devices running heavier mobile internet applications. *IEEE Commun Mag* 2008;46(1):56–63.
- Yang L, Cao J, Yuan Y, Li T, Han A, Chan A. A framework for partitioning and execution of data stream applications in mobile cloud computing. *SIGMETRICS Perform Eval Rev* 2013;40(4):23–32.
- Zhang X, Schiffman J, Gibbs S, Kunjithapatham A, Jeong S. Securing elastic applications on mobile devices for cloud computing. In: Proceedings of the 2009 ACM workshop on cloud computing security, CCSW '09. New York, NY, USA: ACM; 2009. p. 127–34.
- Zhang X, Kunjithapatham A, Jeong S, Gibbs S. Towards an elastic application model for augmenting the computing capabilities of mobile devices with cloud computing. *Mob Netw Appl* 2011;16(3):270–84.
- Zhang Y, Huang G, Liu X, Zhang W, Mei H, Yang S. Refactoring android java code for on-demand computation offloading. *SIGPLAN Not* 2012;47(10):233–48.