

# A Dynamic Programming Offloading Algorithm for Mobile Cloud Computing

Haleh Shahzad

Dept. Electrical and Computer Engineering  
McMaster University  
Hamilton, Canada  
shahzad@mcmaster.ca

Ted H. Szymanski

Dept. Electrical and Computer Engineering  
McMaster University  
Hamilton, Canada  
teds@mcmaster.ca

**Abstract**—Computational offloading is an effective method to address the limited battery power of a mobile device, by executing some components of a mobile application in the cloud. In this paper, a novel offloading algorithm called ‘Dynamic Programming with Hamming Distance Termination’ (denoted DPH) is presented. Our algorithm uses randomization and a hamming distance termination criterion to find a nearly optimal offloading solution quickly. The algorithm will offload as many tasks as possible to the cloud when the network transmission bandwidth is high, thereby improving the total execution time of all tasks and minimizing the energy use of the mobile device. The algorithm can find very good solutions with low computational overhead. A novel and innovative approach to fill the dynamic programming table is used to avoid unnecessary computations, resulting in lower computation times compared to other schemes. Furthermore, the algorithm is extensible to handle larger offloading problems without a loss of computational efficiency. Performance evaluation shows that the proposed DPH algorithm can achieve near minimal energy while meeting an application’s execution time constraints, and it can find a nearly optimal offloading decision within a few iterations.

**Keywords** — *Computational Offloading; Mobile Cloud Computing; Dynamic Programming; Randomization; Hamming Distance; Energy-Efficiency;*

## I. INTRODUCTION

Computation offloading is a method where some of the computational tasks of a mobile application can be offloaded to run on remote servers in the cloud, in order to save energy [1] [2]. However, the problem of partitioning the application tasks for offloading is NP-complete in general [3]. The main goal of the offloading algorithm is to minimize the overall energy used by the mobile application, while meeting an execution time constraint.<sup>1</sup>

A task to be offloaded must be transmitted over a wireless access network, and the time-varying wireless transmission bandwidth must be considered. An adaptive offloading algorithm can determine the offloading decisions dynamically according to a changing wireless environment.

Reference [4] presented a system that enables energy-aware offloading of mobile tasks to the cloud called MAUI. Further improvements were proposed in CloneCloud [5] and

Thinkair [6]. Dynamic partitioning was presented in [7]. In all cases, the partitioning problem results in an integer programming problem which cannot be solved efficiently. A Dynamic Programming (DP) algorithm was proposed in [8], where a two dimensional DP table was used. However, this scheme did not consider an execution time constraint when computing the offloading decisions, although this time constraint is an important issue for many interactive applications [9]. Furthermore, a backtracking algorithm was needed to find the final decisions, which was time consuming.

Reference [10] presented a dynamic offloading algorithm which is based on Lyapunov optimization. The algorithm is based upon a relationship between the current solution and the optimal solution. This scheme requires a considerable amount of execution time and it needs many iterations to converge upon a solution.

Reference [11] provided a dynamic programming approach which builds a three-dimensional programming table and requires pseudo polynomial time complexity [11]. However, it doesn’t consider the energy consumed in the mobile device which is an important criteria for mobile devices.

A semidefinite relaxation approach for the offloading problem was presented in [12]. They considered a mobile cloud computing scenario consisting one nearby ‘Computing Access Point’ (CAP), and a remote cloud server(s). The algorithm first solves a linear program, and uses randomization and relaxation to generate an integer solution. The algorithm can find a near-optimal solution when using about 100 trials of relaxation, but they didn’t discuss the time complexity of their algorithm.

In this paper, an innovative dynamic programming algorithm called DPH is proposed. Dynamic programming is an optimization approach that transforms a complex problem into a sequence of simpler problems. Our DPH algorithm introduces randomization, i.e., we generate random bit strings of 0s and 1s periodically and utilize sub-strings when they improve the solution (similar to genetic optimization). We also fill the dynamic programming table in a creative way to avoid the extra computation for common sub-strings. The algorithm can find a nearly optimal solution after several iterations. It **uses a Hamming distance criterion to terminate** the search to obtain the final decision quickly. The hamming

<sup>1</sup> This work was funded in part by a grant from the Natural Sciences and Engineering Research Council of Canada.

distance termination criterion is met when a given fraction of tasks are uploaded. The final solution depends upon the wireless network transmission bandwidth and the computational power of the CAP and cloud servers.

A comparison between our simulation results and reference [12] shows that our proposed algorithm can find a nearly optimal solution while satisfying a computational time constraint, with a lower computation time. Our proposed algorithm is efficient and can handle larger problems with much lower complexity compared to other dynamic programming algorithms.

The remainder of the paper is organized as follow: Section II provides the system model and problem formulation. Section III presents the proposed algorithm which is based on the dynamic programming table. Evaluation and simulation results are shown in section IV and the paper concludes in section V.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

Consider an application consisting of some unoffloadable (i.e., local) tasks and  $N$  offloadable tasks. Normally, local tasks include those that directly handle user interaction, access local I/O devices or access specific information on the mobile device. Therefore, local tasks must be locally processed by the mobile user. We can merge all the local tasks into one task [13]. In [14], an example of a face recognition problem which is composed of eight offloadable tasks and one local task is presented.

### A. Network Model

We consider a handheld mobile device with  $N$  independent tasks that can be executed locally or transferred to cloud for execution as shown in Fig. 1. We assume that a WiFi wireless network is available for the mobile device, but the network transmission bandwidth can change dynamically. Typically, wireless interference and network congestion will dynamically change the network transmission bandwidth. The mobile device needs to decide whether each task should be processed locally or offloaded, according to the current wireless network transmission bandwidth.

The time taken to transfer a task between a mobile device and the cloud through a wireless link is an important issue since a total execution time constraint for all tasks exists. Therefore, the dynamic programming algorithm must consider the current wireless network bandwidth when computing a decision.

### B. Problem Formulation

For task  $i$ , let  $M_i \in \{0, 1\}$  be an execution indicator variable. Let  $M_i = 1$  if task  $i$  is executed at the mobile device and 0 otherwise. If it is executed locally, the energy consumption is  $El_i$ .  $Er_i$  is the energy consumption of the mobile device when the task  $i$  is executed on the cloud, and  $Et_i$  is the energy-cost of transmitting task  $i$  to the cloud server.

Variable  $Tl_i$  is the local execution time to process task  $i$ , and  $Tr_i$  is the remote execution time when task  $i$  is executed in

remote cloud server. Variable  $Tt_i$  is the transmission time to transfer task  $i$  to the cloud server.

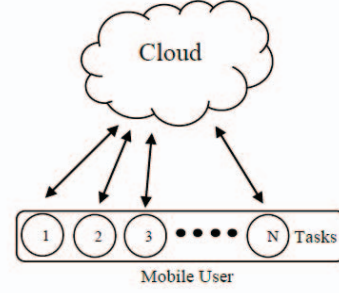


Fig.1 Network Model

It is clear that the transmission energy used to upload each task will depend on the network transmission bandwidth. Therefore, changes in wireless network bandwidth will affect the offloading decision. For example, if we assume that transmission time of each task is equal to the size of each task divided by the network transmission rate, then any variation in the transmission rate will affect the final decision of whether to offload this task or not.

The energy consumption function and its corresponding execution time are defined in (1) and (2):

$$E = \sum_{i \in N} (M_i El_i + (1 - M_i) Er_i + (1 - M_i) Et_i) \quad (1)$$

$$T = \sum_{i \in N} (M_i Tl_i + (1 - M_i) Tr_i + (1 - M_i) Tt_i) \quad (2)$$

The execution time  $T$  of all tasks must satisfy the following condition, where  $T_{\text{constraint}}$  is the execution time requirement.

$$T \leq T_{\text{constraint}} \quad (3)$$

For simplicity, we use  $M = [M_1, M_2, \dots, M_N]$  to denote a vector of binary offloading decisions. The problem that we want to solve is as follow:

$$\begin{aligned} & \min E \\ & \text{Subject to: } T \leq T_{\text{constraint}} \end{aligned} \quad (4)$$

The number of combinations of binary values  $M_i$  to search for the optimal solution grows exponentially with the number of tasks. Our goal is to determine which tasks should be offloaded to the cloud server to minimize energy while meeting a mobile application's execution time constraint.

## III. PROPOSED ALGORITHM BASED ON THE DYNAMIC PROGRAMING TABLE

### A. Innovative Way of Filling the Table

Our proposed algorithm is called ‘Dynamic Programming with Hamming Distance Termination’ (DPH). In this scheme, we use an  $N \times N$  table to store the bit-streams that show which tasks should be offloaded (where  $N$  is the number of tasks). For the first step, a random bit stream is generated that determines a first solution. This stream is assigned to the table such that 1s are assigned to the next horizontal cell, and the 0s are assigned to the next vertical cell. If the first bit of the stream is 1, the starting cell is (1, 2) and if the first bit of the stream is 0, the starting cell is (2, 1). This approach will avoid extra computations for common bit strings.

A 2D  $8 \times 8$  table is shown in table I. To clarify, assume that  $N = 8$  and the first random stream is 11100110 (black numbers) or 00110110 (red numbers), (2 examples are given). Assume that the second random bit stream in each case is 11000111. The starting cell of the second stream is (1, 2) since the first bit is 1. By following the aforementioned rules to fill the table, the resulting green stream is shown in table I.

Table I. How to fill the table for 2 examples.

	1	1	1				
0			0				
0	1	1	0	1	1		
		0	1	1	0		
			0				

	1	1	1				
0		0	0				
0	1	1/0	0	1	1		
		0	1	1	0/1		
				0			

Whenever a bit stream is generated randomly, we calculate the consumed energy and time of each cell (i.e., each task) in the table, and also at the same time calculate the total energy and execution time of this bit stream. However, if a random bit stream is generated which has some common cells with an existing string in the table, we only calculate the total energy of new string until the first common cell and then compare this new total energy with the existing total energy at this cell.

If the new total energy at this specific cell is less than the previous one, we keep the new sub-string and delete the old sub-string, and replace the total-energy and cell-energy of this cell with new amounts. We then update the energy and execution time of the remaining cells for the existing bit stream, based on the new values at this common cell.

Otherwise, if the total energy of the existing bit stream is less than that of the new bit stream at the common cell, we will perform the same procedure while keeping the existing stream.

Every time a new stream is generated, we keep tracking the arrangement of the stream in the table. We terminate and accept a solution which has Hamming distance larger than a given threshold from an all 1 stream. The all 1 stream denotes the case where all components are executed locally. For example, the algorithm can terminate after  $K=20$  iterations, or

when 70% of the tasks have been offloaded. This heuristic termination criterion yields good results.

### B. Algorithm of the Proposed Scheme

The algorithm of our proposed scheme is shown in table II:

Table II. Algorithm of Proposed Method (DPH)

1. Initialize Energy and Time matrixes and set the time constraint ( $T_{constraint}$ ) and Transmission Rate
2. **for** iteration = 1 to iteration\_num
3. generate a random bit stream
4. check the first bit to specify the starting cell in the table
5. **for** i = 1 to N-1
6. Put each bit of the bit stream in the correct position in table
7. Calculate the self-Energy and time of each cell and the total energy and time.
8. **if** this specific cell in table is visited before  
compare the new Total Energy of this cell with the previous one
9. **if** the new Total Energy of the cell is less than the previous one
10. Replace the total energy and time of this cell with the new calculated amounts.
11. Update the remaining amounts in the remaining cells of the previous bit stream based on the new amount of this common cell.
12. Calculate the energy and time of the remaining bits of the new bit stream.
13. Track the position of all bits in the table in a matrix
14. **else**
15. Keep the previous total energy and time in the cell.
16. Calculate the Energy and time of the remaining cells of the new stream based on the existing amount of this cell.
17. Track the position of all bits in the table in a matrix
18. **End if**
19. **End if**
20. **End for**
- 21.
22. **if** Number of bits in table =  $N$  &  $E_{total} < E_{min}$  &  $T_{total} < T_{constraint}$  & hamming distance criterion is met
23. **return**  $E_{total}$ ,  $T_{total}$
24. **end if**
25. **End for**

## IV. EVALUATION

The proposed DPH algorithm is programmed using MATLAB. We adopt the mobile device characteristics from [12], which is based on Nokia N900, and set the number of tasks as  $N = 10$ .

The local computation time is  $4.75 \times 10^{-7}$  s/bit and the local processing energy consumption is  $3.25 \times 10^{-7}$  J/bit, when the x264 CBR encode application (with 1900 cycles/byte) is considered as task (i) in our simulations, as given in [12].

The energy consumption for transmitting and receiving data of the mobile user are both  $1.42 \times 10^{-7}$  J/bit. The input and output data sizes of each task are assumed to be uniformly distributed from 10 to 30MB and from 1 to 3MB respectively, as in [12]. When tasks are offloaded to the cloud, the peak transmission rate is 15 Mbps and transmission time of each task is equal to the size of each task / transmission rate. The

CPU rate of the cloud is  $10 \times 10^9$  cycle/s. Based on [16], the power consumption of this mobile device in idle mode (i.e. when tasks are executed in cloud) is 30 mw.

#### A. Simulation Results

In this section, we provide the simulation results which are done in Matlab. Reference [12] assumed that the total cost to minimize is expressed in Joules, and is defined as the weighted sum of (i) the total energy consumption, (ii) the costs to process the tasks in the cloud, and (iii) the corresponding transmission and processing delays. The weight factor of delays in the weighted sum was defined as  $\rho = 1$  J/s, and  $\beta$  (J/bit) is the weight factor in the processing cost for a task to be offloaded to the cloud (i.e. the weighted processing cost for task  $i$  equals to  $\beta C_{ci}$ ).

Ref. [12] reports experimental results for 2 cases: (i) the traditional case LC (Local-Cloud) where tasks can be offloaded to the cloud, and (ii) the LAC (Local-Access-Cloud) case which allows offloading to the nearby CAP or to the remote Cloud. Figure 2 in [12] shows the total cost versus the parameter  $\alpha$  where  $\alpha$  is the weight of the CAP processing cost. When  $\alpha$  is large, there is no CAP and the performance of LC and LAC converge together. To be consistent, we set  $\beta = 5 \times 10^{-7}$  J/bit and we set the value of cost  $C_{ci}$  to be equal to the input data size  $D_{in}(i)$  (times J/bit).

Fig. 2 shows our comparison of the total cost between the LC and LAC algorithms in [12], and the DPH algorithm for the LC case. The results of the DPH algorithm are ‘nearly optimal’ and are only slightly larger (about 5%) than the results in [12] (which uses linear programming followed by 100 randomization/relaxation trials).

Fig. 3 presents the total cost for 3 cases, (i) all-local execution, (ii) all-offloaded, and (iii) our DPH scheme. (We set  $T_{\text{constraint}} = 700$  and the wireless transmission rate changes between 3 to 10 Mbps). The energy consumption and execution time of the DPH is computed using (1) and (2), respectively. Since the DPH algorithm is a dynamic algorithm, it finds an offloading decision according to the wireless transmission rate. Fig. 3 minimizes the total cost, using the same weights  $\beta$  and  $\rho$  as [12].

Fig. 3 shows that the total cost for DPH is slightly better than the all-offloaded case, for this choice of parameters. For low transmission rates, all-local execution consumes less energy than DPH and all-offloaded, since the cost of transmitting tasks to the cloud is relatively high. However, all-local execution exceeds the execution time constraint as shown in Fig. 4. Hence, the best solution is DPH.

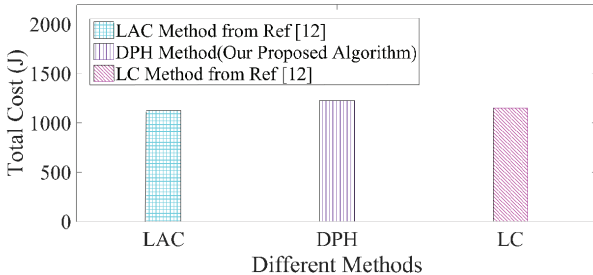


Fig. 2. Comparison of the total Cost for different methods.

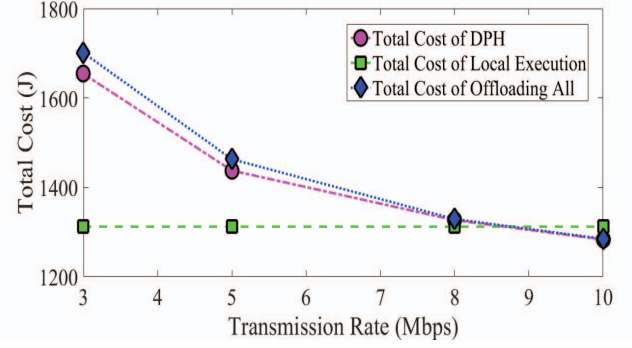


Fig. 3. Comparison of the total cost of (i) all-local execution, (ii) all-offloaded and (iii) our DPH algorithm, when  $T_{\text{constraint}} = 700$ .

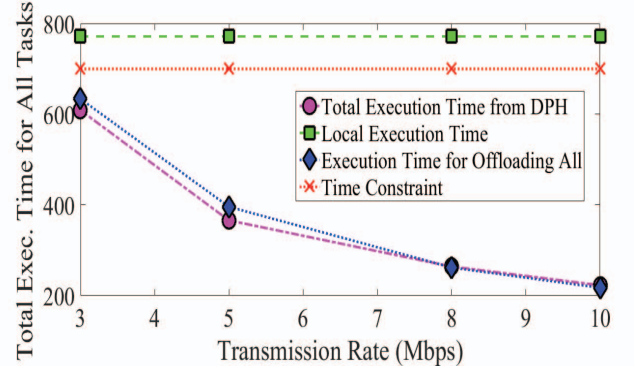


Fig. 4. Comparison of the total execution time of (i) local execution, (ii) offloading all and (iii) our DPH algorithm when  $T_{\text{constraint}} = 700$ .

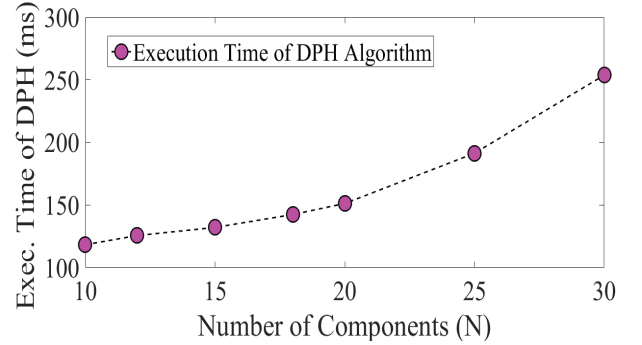


Fig. 5. Execution Time of DPH for different number of components

Fig. 4 plots the total execution time versus the wireless transmission rate. As the rate increases, the benefits of offloading grow, and the total execution time in the mobile device decreases. The DPH algorithm yields the best results compared to the cases (i) all-local and (ii) all-offloaded. At high rates, the case offload-all and DPH yield similar results.

Fig. 5 shows the execution time of DPH algorithm when number of tasks changes from 10 to 30. The execution time grows nearly linearly with the number of tasks.

Fig. 6 compares the performance under 3 cases: (i) all-local execution, (ii) cloud offloading, and (iii) CAP offloading. The cloud case assumes 2 relatively low transmission rates, BW = 5 or 20 Mbps. The CAP parameters



are (i) CPU rate =  $5 \times 10^9$  cycle/s, (ii) a high transmission rate of 72 Mbps as in [12]. The case (ii) cloud offloading with a BW = 20 Mbps outperforms all others, due to the high computational capacity in the cloud.

Fig. 7 repeats the tests in Fig. 6, when the task sizes are half that in Fig. 6. The IO sizes of each task are uniformly distributed from 5 to 15 MB and from 1 to 3 MB respectively. When the CAP has a high BW of 72 Mbps and the task sizes are small, the performance of the CAP improves and the case (iii) CAP offloading can outperform case (ii).

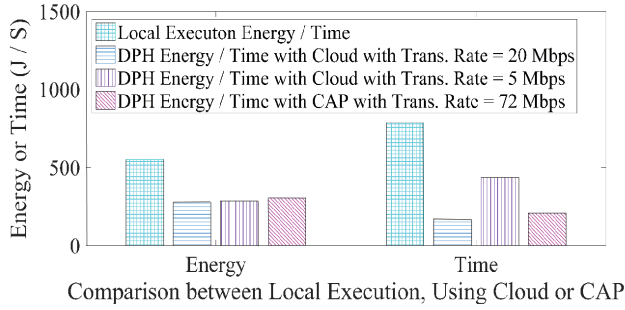


Fig. 6. Comparison of the Energy and Time between local execution, using cloud or CAP.

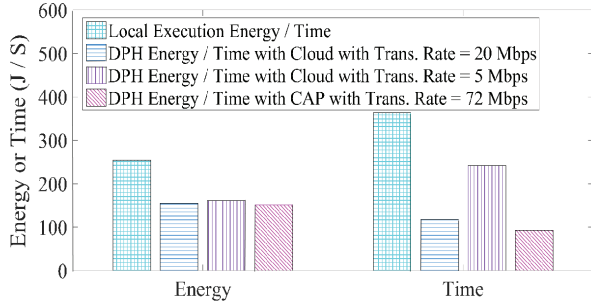


Fig. 7. Comparison of the Energy and Time between local execution, using cloud or CAP for a small set of data sizes.

## V. CONCLUSION

A mobile device must decide which computational tasks of a mobile application should be offloaded in order to minimize energy consumption while satisfying an execution time constraint. We have proposed an efficient heuristic algorithm called DPH to solve this optimization problem, which uses dynamic programming combined with randomization. It also uses a hamming distance as a termination criterion.

Simulation results show that the proposed DPH algorithm can find nearly optimal solutions and it can be easily handle larger problems without losing computational efficiency. The DPH algorithm can be used dynamically, to adapt to the

changes in the network transmission rate. The algorithm will tend to offload as many tasks as possible when the network performance is good, resulting in a rapid convergence to a near optimal solution with a very fast execution time.

## REFERENCES

- [1] Z. Li, C. Wang, and R. Xu, "Computation offloading to save energy on handheld devices: a partition scheme," in Proc. International Conf. Compilers, Architecture, Synthesis Embedded Syst., pp. 238–246, 2001.
- [2] P. Rong and M. Pedram, "Extending the lifetime of a network of batterypowered mobile devices by remote processing: a Markovian decisionbased approach", Design Automation Conf., pp. 906–911, 2003.
- [3] X. Gu, K. Nahrstedt, A. Messer, I. Greenberg, and D. Milojicic, "Adaptive offloading for pervasive computing," IEEE Pervasive Comput., vol. 3, no. 3, pp. 66–73, 2004.
- [4] E. Cuervo, A. Balasubramanian, D.-k. Cho, A. Wolman, S. Saroiu, R. Chandra, and P. Bahl, "MAUI: Making smartphones last longer with code offload," in Proc. ACM International Conference on Mobile Systems, Applications, and Services (MobiSys), pp. 49–62, 2010.
- [5] B.-G. Chun, S. Ihm, P. Maniatis, M. Naik, and A. Patti, "Cloudfunder: Elastic execution between mobile device and cloud," in Proc. ACM Conference on Computer Systems (EuroSys), pp. 301–314, 2011.
- [6] S. Kosta, A. Aucinas, P. Hui, R. Mortier, and X. Zhang, "Thinkair: Dynamic resource allocation and parallel execution in the cloud for mobile code offloading," in Proc. IEEE International Conference on Computer Communications (INFOCOM), pp. 945–953, 2012.
- [7] J. Niu, W. Song, L. Shu, and M. Atiquzzaman, "Bandwidth-adaptive application partitioning for execution time and energy optimization," in Proc. IEEE International Conference on Communications (ICC), pp. 3660–3665, 2013.
- [8] Y. Liu, M. J. Lee, "An Effective Dynamic Programming Offloading Algorithm in Mobile Cloud Computing System", IEEE WCNC'14, pp.1868 – 1873, 2014.
- [9] H. Lagar-Cavilla, N. Tolia, E. De Lara, M. Satyanarayanan, and D. O'Hallaron, "Interactive resource-intensive applications made easy," ACM/IFIP/USENIX International Conf. Middleware, pp. 143–163, 2007.
- [10] D. Huang, P. Wang, D. Niyato, "A Dynamic Offloading Algorithm for Mobile Computing", IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, VOL. 11, NO. 6, 2012.
- [11] A. Toma, J. Chen, "Computation Offloading for Frame-Based Real-Time Tasks with Resource Reservation Servers", IEEE Euromicro Conference on Real-Time Systems, pp. 103-112, 2013.
- [12] M. Chen, B. Liang, M. Dong, "A Semidefinite Relaxation Approach to Mobile Cloud Offloading with Computing Access Point", IEEE Signal Processing Advances in Wireless Communications, pp. 186-190, 2015.
- [13] X. Gu, K. Nahrstedt, A. Messer, I. Greenberg, and D. Milojicic, "Adaptive offloading for pervasive computing," IEEE Pervasive Comput., vol. 3, no. 3, pp. 66–73, 2004.
- [14] Available: <http://darnok.org/programming/face-recognition/>.
- [15] A. Kammerdiner, P. A. Krokhmal, P. M. Pardalos, "On the Hamming distance in combinatorial optimization problems on hypergraph matchings", Springer Optimization Letters, Vo. 4, pp. 609-617, 2010.
- [16] I. Jantunen, J. Jantunen, H. Kaaja, S. Boldyrev, L. Wang, and J. Hämäläinen, "System Architecture for High-speed Close-proximity Low-power RF Memory Tags and Wireless Internet Access", International Journal on Advances in Telecommunications, vol. 4, 2011.