

Multimedia Computational Offloading for 5G Mobile Edge Computing

Adhimuga Sivasakthi D (✉ dasivasakthi0607@gmail.com)

Madras Institute of Technology <https://orcid.org/0000-0001-5694-9547>

Gunasekaran Raja

Anna University Chennai

Research Article

Keywords: Computation offloading, mobile edge computing, multihopmesh communication, edge server, energy efficiency

Posted Date: February 15th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-221801/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Multimedia Computational Offloading for 5G Mobile Edge Computing

D Mr. D. Adhimuga Sivasakthi^{*1} Dr. R. Gunasekaran²,

**¹Research Scholar, Department of Computer Technology, Madras Institute Of Technology,
Anna University, Chennai, India*

*Email: dasivasakthi0607@gmail.com

² Professor & Head in Department of Computer Technology,

Anna University, Chennai , Tamilnadu, India

Abstract

10 Due to the tremendous growth of smartphone users with their massive usage of faster
11 response, delay sensitive applications lead to high traffic demands, still not yet been met by
12 current researchers. Thus, the ever-increasing computational demands are managed through
13 computational offloading, which offload intensive workloads to the small cells having
14 functionality similar to the resource-rich providers namely edge server. With the use of edge
15 server, the heavily loaded computations get executed outside of the mobile device which
16 results in minimization of mobile energy consumption and latency. Mobile edge computing is
17 an emerging paradigm of commercial infrastructure for computational offloading, which
18 enhances the power of smart mobile devices. Generally, edge provides cloud services and
19 resources to the nearest proximity of users with radio access in fifth-generation (5G)
20 networks for low latency, prompt response, and filtering. But, the frequent network failure
21 automatically degrades the performance. Hence in this paper, we integrate the functionalities
22 of Edge Server with Pico cells (ESP) for efficient traffic management and also propose a
23 distributed Multi hop Mesh Middle layer (MMM) architecture for seamless communication
24 with high resilience, minimal latency, and reduction of energy consumption. Generally, a
25 pico cell is a distributed antenna system, an alternative to a repeater used to extend wireless

26 services to 100 users. We developed an Ant Social based Vector (ASV) optimization
27 algorithm for an intelligent offload decision, which paves the way for backhaul routing
28 conflicts. Simulation results of Cloud Sim show that offloading tiny devices integrated with
29 significant additional computational capabilities would satisfy all the demands from each
30 network within a 15ms delay.

31 **Keywords:** *Computation offloading; mobile edge computing; multihopmesh communication;*
32 *edge server; energy efficiency*

33 **1. Introduction**

34 With the spectacular growth of the mobile devices, the massive demands of resource-hungry
35 and delay sensitive critical mobile services, such as video format transformation, voice
36 recognition, interactive gaming, virtual reality, augmented reality, real-time speech
37 recognition and content-based image retrieval have attracted significant attention of future 5G
38 networks. Since, these applications are delay sensitive [1] which resulted in an increasingly
39 high computing demand, which address multiuser computation and heterogeneous network
40 issues. For that, it is mandatory to have an exclusive virtual machine for each user which
41 allows resource sharing by the assumption of peak workload demands as static. At the time of
42 offloading, the demands for doing computations exceed the capacity of the mobile device.
43 Due to the limitations of the computation capacity and cache size, resource-constrained
44 applications cannot be executed efficiently on the mobile device. Most of the current work
45 focuses on datasets with short data packets (below Kbits).The researchers of future wireless
46 use cases e.g., Virtual Reality (VR), Augmented Reality (AR) and Tactile Internet (TI),
47 demands not only transmission of a large packet but requires immense computing facilities. It
48 is difficult for a mobile user with limited computation and storage resources to meet the
49 requirements of current use cases. On the other hand, battery lifetime of mobile devices is the
50 primary constraint, which indicates plethoric energy consumption [2] from local computing

51 which is not tolerable for the user to get better user experience. The file size of responses that
52 comes back from the server is smaller when compared to the file size of computation or
53 transmission. Thus, it is mandatory to keep the computational server close to the user's
54 location.

55 Notice that the distance of cloud and cloudlets are farthest from a mobile device, so
56 offloading to cloud server leads to delay in user request [3]–[5]. Alternative techniques or
57 adaptive algorithms are in need to offload intensive task to small cells. Small cell offloading
58 brings new challenges such as computation workload, file size, wireless communication
59 protocol, latency. But it subsidizes heterogeneity of computing, storage, communication
60 resources, and diversity requirement of network services. Therefore, it is more preferable in
61 situations where eNB's get overloaded in hotspots with high-density demand. For example, to
62 prevent mobile users feeling as irritation, dizzy and nauseous, VR system has to guarantee the
63 latency which should be less than the time humans start noticing the lag [6] (e.g., not more
64 than 20ms).

65 By the year 2019, IDC reported that 45% of the IoT data generated by diverse requests
66 gets processed, stored, and analyzed on edge [7]. Edge technology is first initiated by the
67 European Telecommunications Standards Institute (ETSI). Edge is an emerging technology
68 which promotes the necessity for communication and computation capabilities that were
69 adaptively managed and shared among mobile devices. The access of small size nano
70 datacenter in 4G and 5G network is termed as mobile edge. An edge device is an
71 infrastructure of a heterogeneous router, gateway, switch, access points, or a base station and
72 end-user devices. Generally, edge devices are required to have sufficient computational [8]-
73 [12] and storage resources. It aid as a complement cloud to extend resources and services. In
74 particular, response time is a primary factor for user satisfaction namely Quality of
75 Experience (QoE) [13], [14]. Through Edge with Pico cell (ESP), swift response for the real-

76 time application and expansive resource demands of the end user gets satisfied with reduced
77 data traffic and minimal delay. In order to satisfy end users, ESP can be used as data
78 delivery and sharing points, thereby computations have been performed with vast volumes of
79 saved data. Caching at the pico cells may considerably reduce access latency and save a
80 considerable amount of backhaul network traffic.

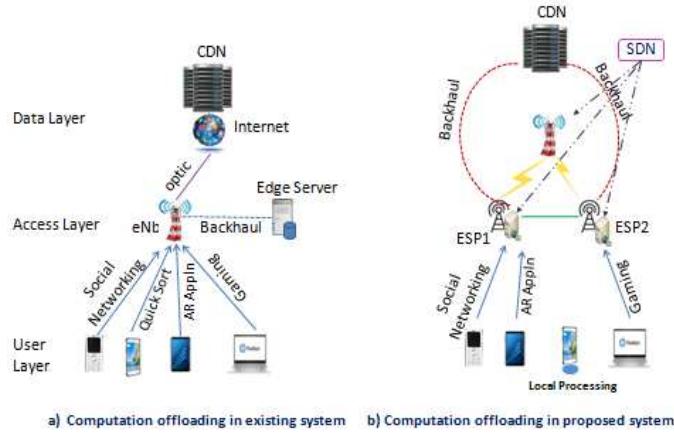
81 Hence, deploying computation and storage resources in ESP could enable a large number
82 of diverse applications with instant response. Further, reducing the volume of downward
83 traffic in the network core essentially brings energy efficiency. To offer distributed services
84 with prompt response, a novel distributed Multi hop Mesh Middle layer (MMM) architecture
85 has been proposed. As shown in Figure 1, MMM integrate the functionalities of ESP as a
86 middle layer between user and data center. In this paper, we are considering ESPs as mesh
87 topology, inter linked with multi hop communication for data filtration and seamless
88 communication. The motivation of this paper is image processing, voice recognition and
89 translation, m-gaming, cooperative spectrum sensing in cognitive radio cloud networks,
90 because of resource-intensive nature of the application and intrinsic limitations of the mobile
91 devices and communication medium.

92 For instance, we can transmit a high-resolution image simultaneously using ESP with a
93 frame rate of 120 frames per second in a network for processing such as frame detection,
94 accuracy in object recognition, virtual holograms modeling, etc. with short time span of less
95 than 20ms. Our proposed MMM enables the applications and services to employ on top of
96 multi-vendor computing platforms, also entitle the mobile operators with more number of
97 customers. We consider multi-user computation offloading problem in 5G network with a
98 scenario that 10% of 8 million people requested to view movies from Netflix at the same
99 time. To handle this scenario, we treat all requests in parallel and access an infrastructure of
100 Content Delivery Network (CDN) from ESP through backhaul as in distributed manner;

101 thereby the hosted content must not travel through various networks to reach the end users
102 mobile.

103 Our main contributions are as follows:

- 104 • We design a novel architecture called Multi hop Mesh Middle layer(MMM) to
105 support intensive computational offloading with self-healing framework.
106 • We formulate energy efficient handover and offloading in 5G with the assistance of
107 SDN.
108 • We develop an Ant Social based Vector (ASV) optimization algorithm to minimize
109 latency and overall network energy.



110 a) Computation offloading in existing system b) Computation offloading in proposed system

111 **Fig. 1** Computation offloading Scenario a) Existing and b) Proposed Framework

112 The rest of this paper is organized as follows: Section II discusses related work, Section
113 III covers system architecture, Section III.1 presents ASV in MMM Architecture, Section IV
114 Section describes experimental analysis of the proposed MMM framework and Section V
115 specifies the conclusion.

116 **2. Related Work**

117 Due to the bandwidth limitations in wireless networks there have been a few researches on
118 efficient computation in 5G with 1) static assumption of offload parameters such as dynamic
119 usage behavior, and network condition, 2) feasible elastic solutions. Generally, decision
120 making is an essential aspect which is primarily based on the prediction of computing time.

121 So, in early 2000s, the main focus is to develop algorithms for offloading decisions to
122 decide whether offloading would benefit mobile users or not [15]. Energy saving [16] is
123 another crucial issue, hence complex and energy constrained computations could be
124 offloaded to the cloud, which performs operations faster with retaining of mobile battery life
125 [17],[18] but it leads to a signal overhead for a mobile user.

126 To overcome offloading problems, Lei et al. proposed a computation offloading
127 mechanism where intensive application is executed completely at remote side, i.e., outside of
128 the mobile device is termed as computation offloading. The method proposed in [19]
129 addresses for a multi-client mobile environment in which virtual machine can be migrated
130 from a cloud datacenter to mini-datacenter (edge), and can be managed through cloud
131 surrogate for better mobile energy efficiency, but the mechanism is not guaranteed for IP
132 continuity. To overcome the IP continuity issue, Wang et al. proposed a centralized improved
133 particle swarm optimization algorithm for a single client application. The authors in [20]
134 focus on minimization of mobile energy consumption, and optimal VM and CDN placement.
135 The mechanism specified in this paper lacks composite services and other topological support
136 except for LTE. Also, Lin et al. suggested a computation off loading framework for effective
137 decision named Ternary Decision Maker (TDM) for a single client application. The scheme
138 proposed in [21] dealt with reduction of both mobile energy consumption and response time.
139 But this framework is suitable for LTE with intermittent connectivity.

140 To address the issues on the preservation of battery power through non-intermittent
141 connectivity, Ragona et al. proposed a centralized framework (ECO-CLOUD) using mobile
142 cloud computing through Near Field Communication (NFC) interfaces for a multi-client
143 application [22].A distributed multisite computation offloading using genetic algorithm has
144 been proposed by Zeinab et al. for a multi-user environment. The authors in [23] targets to
145 optimize the overall energy consumption and response time by performing efficient

146 offloading. Still the framework does not support differences in access delays, proximities,
147 computational capability or monetary cost. Besides, Roy et al. proposed a multi usermobile
148 cloud computing algorithm for the optimum selection of cloudlet. The authors in [24]
149 emphasis on the reduction of both battery power and latency time. This paper discusses about
150 the rare availability of cloudlet, with complex offloading decision and time consuming
151 selection scheme.

152 Mao et al. proposed a centralized Energy Efficient Computation Offloading (EECO)
153 algorithm to solve the time-consuming issue which jointly optimizes partial offloading and
154 radio resource allocation by incorporating the multi-access characteristics. The authors in
155 [25] consider the utilization of energy for both computation and transmission. However, the
156 theoretical analysis considers only priority-based channel allocation. Also, Yang et al.
157 proposed a distributed joint optimization algorithm for both the problems of multiclass
158 multichannel dimensional knapsack and performance function matrix. The authors in [26]
159 provides mechanisms to manage competition among multi-users, effective computation
160 partitioning. But it fails to addresses the maximization of average application throughput,
161 decision making and consumption of mobile energy.

162 Jason et al. proposed architecture for centralized framework called AR edge computing
163 for a multi-client application to perform partial remote live offloading. The authors in [27]
164 addresses the issues of reduction in end-end latency and cost. However, the mathematical
165 analysis is makes use of homogeneous network. To achieve efficient offloading, Demiriset al.
166 presented a framework of demand based computational offloading to reduce monetary cost
167 and to maximize energy savings. The authors in [28] dealt with balancing the competing
168 demands and managing infrequent access of an individual user. Also, Ganz et al., presented a
169 survey of computations on the heterogeneity of data for a single user, which need not be
170 applied to only one particular information processing technique.

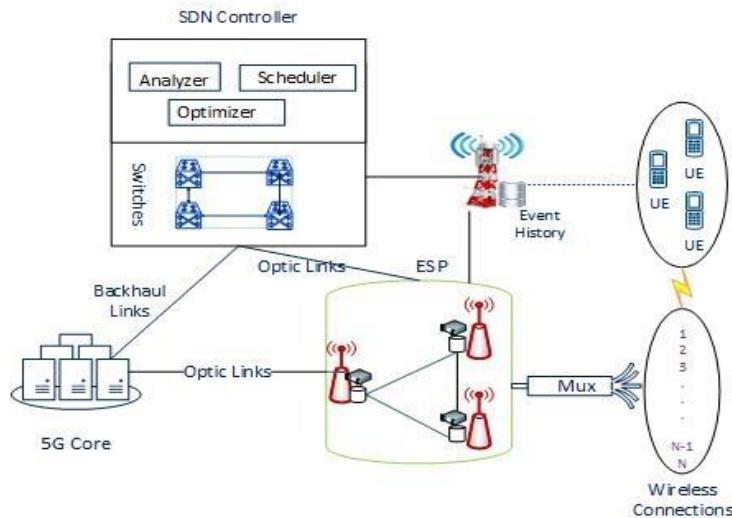
171 Besides, the technique in [29] is adaptive to changes in the input data and handling of
172 multimodal data in eNB without the need for domain knowledge. Cao et al. discussed game
173 theoretic offloading mechanisms in eNodeB (eNB)with the assistance of edge. The
174 expectation of multichannel contention environment [30] is to balance the transmission and
175 computation cost by which it can achieve least execution cost. But the mathematical model
176 works better for homogeneity in terms of their computation ability and their computation
177 tasks. In our paper, we propose a mesh multi hop self-healing ESP for 5G networks.

178 **3. MMM Architecture for Computation Offloading**

179 To improve the coverage of multiuser, pico cells are densely deployed in hotspots areas
180 around 20 meters' distance such as shopping malls, retail spaces, stock exchanges, train
181 stations, etc. which can support for a maximum of 100 users at a time. Due to fast,
182 inexpensive and vast coverage with reduced outdoor interference, we use picocells for doing
183 computations at the time of handover. So that, it allows operators to provide excellent
184 services at dynamic environment, improving customer perception and satisfaction. Also, the
185 higher signal quality of picocell leads to better throughput because it allows the network to
186 use a more spectrally efficient transmission scheme where more bits can be transmitted at the
187 same time.

188 Thus, by integrating the functionalities of Edge Server in Picocells (ESP) as a multihopin
189 mesh network provides high resilience. i.e., any sort of fault or failure in a node will not be
190 affected or make the entire network to be down. Moreover, connecting the operators of a core
191 network (CN) is done via backhaul links, which could reduce CAPEX and OPEX. In our
192 paper, we aim to achieve not only faster but an effective resource and backhaul sharing of
193 computation for multi-users with tolerable delay. Besides, we develop a distributed Ant
194 Social based Vector optimization algorithm (ASV) for seamless communication and
195 computation with reduction of signal overhead.

196 To enhance backhaul routing, we use learning technique for the design of efficient user
 197 pattern. The user pattern comprised of detailed information: 1) user behavior such as
 198 distance, age, education, culture and type of information sharing, 2) event history such as
 199 frequent site access, time and usage and 3) social affiliation. Based on this information, eNB
 200 allocate ESPs and channels for user. For efficient computation offloading.



201
 202 **Fig.2** MMM Architecture for Computation Offloading
 203 We handover computations of eNB to ESP with the assistance of Software Defined
 204 Networks (SDN).With the use of SDN, it is possible to segregate both data and control
 205 planes. Thereby, it brings mobile devices with high resilience, also overcome eNBs frequent
 206 feedback and signal overhead.

207 In our paper, to handle the mobile device with faster and effective computation, we
 208 consider eNB as a master and ESPs serves as workers. Each picocell is embedded with high-
 209 level edge cloud components namely transcoder, buffer, cache, and local database (Trained
 210 data). As shown in Figure 2, whenever mobile users are requesting for a social media, we
 211 handover the real-time demand of eNB to ESPs and so eNB can proceed with any other
 212 parallel operations for any user. Due to the non-intervention approach of eNB, energy
 213 overhead has been reduced.

214 We manage the entire communication and computation between eNB and ESP with the
215 assistance of an SDN controller, which controls both eNB and ESP for seamless
216 communication. Also, SDN monitors the inter-process communication between eNB and
217 ESP. There is direct access between ESP to CDN till the corresponding response delivered to
218 the user. The same content can be delivered (broadcast) through multiplexer via a frequency
219 reusable wireless network. Thus, we introduce an architecture named Multihop Mesh Middle
220 layer (MMM) to improve the total execution time, by making the time of computation in ESP
221 as smaller than local. Initially, the classification of traffic can be carried out by an analyzer
222 which segregates the user demands with an intensive and non-intensive task based on
223 computation time or energy required. Then, scheduler can predict the user's complete usage
224 behavior such as a type of network, type of frequent accesses of the site, timing, interest,
225 mode of viewing (stream/online), number of accesses, and network usage using event history,
226 which helps for prefetching of content from CDN. Internal cache of ESP is used to hold the
227 snapshot of next task access while the current is in execution which helps to reduce latency
228 and computation time.

229 Finally, optimizer supports optimal routing and backhaul conflicts through a distributed
230 ASV algorithm. Various devices such as Apple, Samsung, Nokia etc., are escorting
231 with different forms of connectivity, such as 4G, Wi-Fi or next-generation radio technologies.
232 Besides, the messages of social media are small, encrypted and come in different forms of
233 protocols such as http, xml, REST and JSON.

234 To manage various protocols, low latency aggregation points are required to distribute
235 messages. In particular, the location of mobile aggregation sites can be either an indoor (e.g.
236 hospital, airport, large corporate HQ), or an outdoor for a particular public coverage scenario
237 (e.g. stadium, shopping mall, railway station) which will be controlled by SDN which
238 provide radio coverage to the premises.

239 Thus, the integration of ESP in small cells empowers the direct delivery of locally-
 240 relevant, fast services from NB. Besides, our MMM provides the capability to resolve multi-
 241 user computation offloading challenges.

242 To minimize the computational offloading problem in MMM Network formulated is
 243 given below

244 1. Service – oriented optimization problem

245 2. CN-Oriented Optimization Problem

246 **1. Service – oriented optimization problem**

247 In the service-oriented optimization problem to maximize the sum of all tasks' utility, and the
 248 problem can be formulated as

249
$$\max_{\alpha_{ij}} \sum_{i=1}^M \sum_{task}^i Y$$

250 s.t.

$$\begin{aligned}
 C1: & \Delta_i \leq d_i - a_i, \forall i \\
 C2: & \sum_{j=1}^N \rho_{ij}(T_{ij} + G_{ij}) \leq T_i, \forall i \\
 C3: & Y_{task}^i \geq 0, \forall i, \\
 C4: & Y_{CN}^i \geq 0, \forall j, \\
 C5: & \sum_i \rho_{ij} \leq 1 \forall j \\
 C6: & \sum_i \rho_{ij} \leq 1 \forall i \\
 C7: & \rho_{ij} = \{0,1\}, \forall i, j.
 \end{aligned} \tag{1}$$

252 here, C1 represents the actual waiting time constraint where Δ_i is the actual waiting time.
 253 Besides, C2 gives the service delay constraint. C3 ensures the minimum utility of each task
 254 and C4 restricts the minimum utility of every CN. C5 indicates that each CN can compute at
 255 most one task one time, and C6 shows that each task can be allocated to at most one CN.
 256 Finally, C7 represents the binary constraint.

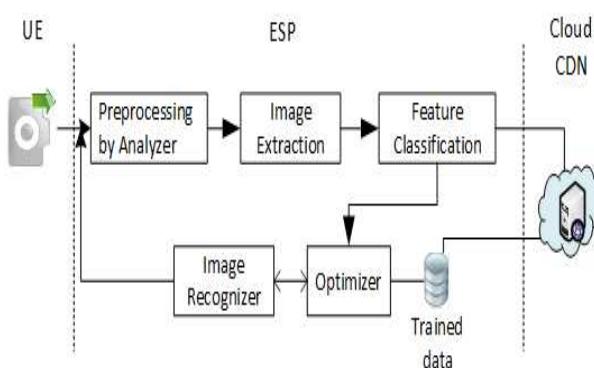
257 **2. CN-Oriented Optimization Problem:**

258 Contrary to the service-oriented problem, the CN-oriented problem focuses on maximizing
 259 the utility of CNs, and the objective function can be represented as

260
$$\max_{\rho_{ij}} \sum_{j=1}^N Y_{CN}^j \quad (2)$$

261 Subjects to constraints C1-C7. Obviously, the service-oriented and CN-oriented optimization
 262 problems presented above involve 0-1 programming, and the general solution for this kind of
 263 problem is branch and bound. The branch and bound is a centralized optimization method,
 264 which indicates that global information of all tasks and CNs is needed during the decision
 265 process. However, the information exchange will overload with the increase of network size.
 266 One effective way to solve this problem is the matching game, which performs in a
 267 distributed manner

268 For example, Image-Editing applications in Figure3, most of the stages are typically carried
 269 out in the ESP. Initially, the captured image from the mobile device can be uploaded over
 270 wireless networks. Then, image-edition procedure is completely carried out by ESP, and then
 271 the response is returned to the user from the ESP. Responses are in some user understandable
 272 media form (e.g., augmented reality). Such MMM architecture enhances the performance of
 273 overall network latency and computing time.



274

275 **Fig. 3** An example of mobile image editing with ESP

276 Steps performed by ESP for image editing:

277 1. User uploaded the updated image captured from a mobile device
278 2. ESP performs preprocessing and then extracts the features from the image
279 3. Trained data in CDN will be fetched, and stored in the cache
280 4. Optimizer checks for the classified features with the trained data
281 5. If there is a match occurs, then ESP returns the edited image to user for approval and
282 then updates CDN with the same. Otherwise sends the error message to user.

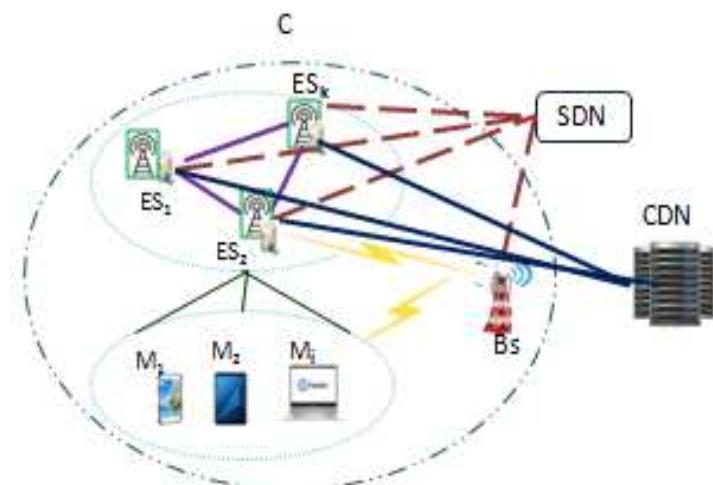
283 ***3.1. ASV in MMM Offloading***

284 Generally, heuristic algorithms can be applied to a wide range of any optimization problems
285 to yield near-optimal solutions with minimal execution time. Thus, we develop an ASV
286 algorithm for optimal path selection. To do so, our proposed algorithm inspired the behavior
287 of ants for seeking food between source locations to the repository. Generally, a pheromone
288 trial is a liquid, lays down by ants during they travel. It acts as trial for further ants to
289 proceed. By the trail values, it is easy to identify the shortest and optimal path. Since the
290 value of trail gets evaporated over time, the only shorter path will attract all ants to deposit
291 more pheromone values in that path.

292 These positive feedbacks of ants will be updated in ta bu (target) list for optimal selection.
293 Finally, the shortest path with high pheromone value will be considered as optimal. When
294 this artificial intelligence concept is applied to NP-hard problems, the term "trail value" or
295 pheromone refers to the latency and its transmission cost from the user location to destination
296 as CDN. The main idea is spreading the utility gained from finishing the computation
297 throughout processor time it requires based on distance or latency.

298 For instance, consider a scenario of a single user requesting for social networking. User is
299 a source node, and CDN is a destination. With each ESP has a list which will be updated by
300 ants with their respective latency values to the destination. Ants will travel from source to
301 destination by laying pheromone in their path with the constraint as ants won't revisit any

302 node twice. Now, associatively compare the collected value in the list of each ESP with the
 303 latency. Finally, the path with less latency will be chosen as optimal for processing the
 304 offload. Our proposed ASV algorithm is also applicable for multiuser diverse application
 305 needs. Also, ESPs are in the form of mesh network, which yields optimal routing. Generally,
 306 by offloading to shorter distance ESPs we can achieve less latency, power, and transmission
 307 or monetary cost.



308

309 **Fig. 4** System Model for computation

310 Hence, it avoids migration of ESP from one location to another based on user's
 311 proximity. Therefore, the ESP of a minimum latency path has been selected for offloading.
 312 Any ESP can execute and share the partial processing with its neighbor ESPs of another
 313 region more easily without transmission delay. Thus, high resilience with multi sharing is
 314 applicable readily and at least cost.

315 Assume a sample of an application as in Figure 4 has few numbers of clusters denoted as
 316 C_i , for $1 \leq i, k \leq n$, and each cluster is modeled as fourtuples $C_i = \{M_i, B_s, E_{S_k}, \gamma_i\}$, where M_i ,
 317 B_s denotes the number of users and base station respectively. Whereas E_{S_k} represents a
 318 number of ESP sin a cluster, γ_i signifies the annotation of the component, indicating whether
 319 components can be executed as locally or on ESP. A mobile device of k user is modeled as
 320 eight tuples with $M_k = \{Z_{in}, Z_{out}, N, \alpha, P_x, P_{idle}, P_y, P_z\}$, where Z_{in}, Z_{out} denoted as size of

321 computation for input and output data, N represents maximum capacity of workload
 322 (maximum workload assigned which measured in Million Instructions), α denotes the
 323 demand of the user, P_x, P_y denotes the power consumption at local computation and
 324 communication by the mobile device, P_{idle} , P_z represents the power consumption at idle time
 325 and display of the device.

326 *3.1.1. Local Computing*

327 If $\gamma_i = 0$, then the component would be executed locally in a mobile device and thus the
 328 energy consumption E and latency L can be calculated using

$$329 \quad E^{loc}_i = Z_{in} (T^{loc}_i P_x) + P_z \quad (1)$$

$$330 \quad L^{loc}_i = T^{loc}_i + \xi \quad (2)$$

331 where T^{loc}_i denotes time taken for local computing, ξ is a constant related to hardware
 332 architecture, and we have

$$333 \quad P_x = \alpha / N.$$

334 *3.1.2. ESP Computing*

335 If computation energy of edge server for a component is less than E^{loc}_i , then we have $\gamma_i = 1$,
 336 and so the component could be offloaded to ESP for computation. When the delay sensitive
 337 component is offloaded to ESP, then the overhead data rate due to the downlink data
 338 transmission T^{ov} of M users via wireless channel can be represented as

$$339 \quad T^{ov}_i = \log_2 (1 + t_n g_{n,s} / P_0) \quad (3)$$

340 where, $t_n, g_{n,s}$ denotes transmission power of user and channel gain between user n and base
 341 station s , P_0 is the power of spectrum density.

342 Let the mobile device with the computational ability of user as C_u and base station as C_b ,
 343 which is quantified by the total number of CPU cycles per second. Then the execution latency
 344 L^{ESP}_i for the task can be computed as

$$345 \quad L^{ESP}_i = T^{trans}_i * C_{ui} / C_{bi} \quad (4)$$

346 where C_{ui} , C_{bi} denotes the computation capability of user and base station via a wireless
347 channel, and T^{trans} is the transmission delay, so we have $T^{trans_i} = Z_{in}/T^{ov_i}$.

348 Let n_i is dynamic power consumption of a CPU cycle, and we have $n_i = Z_{in} * C_u$, then the
349 energy consumption at over head E^{ov} can be represented as

350
$$E^{ov_i} = t_n T^{trans_i} + \delta \quad (5)$$

351 where δ is constant energy due to a mobile device which could hold the channel for a while
352 even after transmission of data. Thus, the computational energy of ESP with latency will be
353 expressed as

354
$$E^{ESP_i} = n_i (P_y + P_{idle} + P_z) + T^{ov_i} (P_y + P_{idle} + P_z) E^{ov_i} / B_i \quad (6)$$

355 where B_i is the spectral bandwidth, T^{ov_i} , E^{ov_i} is the over head time and energy for obtaining
356 and releasing ESP. Computations can be done either locally or in ESP.
357 In this paper, our proposed ASV is formulated as a mixed integer, and non-convex problem
358 along with the objective of maximizing energy savings and latency is expressed as follows

359
$$F(x) = \max \sum_{i=1}^n E_i \cdot \gamma_i \quad (7)$$

360 Such that $\sum_{i=1}^n \gamma_i (T^{trans_i} + T^{ov_i}) + T^{loc_i} \gamma_i \leq \alpha \quad (8)$

361 Along with $\sum_{i=1}^n \gamma_i L_{ESP_i} \leq \text{lag} \quad (9)$

362 where α is the total local execution time of the mobile application, lag represents the
363 maximum tolerable delay for a sensitive application, E_i is the difference between a
364 component local computation energy and edge computation energy. Thus, we have the
365 choices of offloading decisions, and the computational complexity of M mobile users and T
366 computation tasks for ESP is $O(n^M)$.

367 In our proposed, computations can be done by ESP, due to which long packet
368 transmission and long thread of feedback would be reduced. Which in turn reduce the issues
369 of battery drain, monetary cost, and latency, particularly low-speed connection. Our objective
370 is minimizing overall power consumption with tolerable delay. In particular, delay is the

371 integration factor of communication, computation and network delay. Due to the integrated
372 functionalities of ESP with a mesh network, 20% of above-mentioned delays will be reduced.
373 In Demand Vector Classification (DVC) algorithm, based on the user's demand, traffic can
374 be segregated as real time and non-real time. DVC is an effective technique which provides
375 optimal and straightforward solutions under all classification techniques.

Algorithm 1: DVC algorithm

Input : cluster of user requests

Output : Decision Vector γ

```
1      Start demand_classification
2      Initialize threshold, n
3      for each request n do
4          compute Eloc , EESP
5          if Eloc> EESP&& EESP≤ threshold
6               $\gamma_i = 1$ 
7          Else
8               $\gamma_i = 0$ 
9          end if
10         end for
11      Stop demand_classification
```

376

377 Thus, analyzer will concurrently evaluate the energy consumed by the mobile device for
378 the current task. Also, the max value in the energy level is assumed to be the threshold and
379 idle value when mobile is inactive as 0.5. In order to take an effective decision, equation 1
380 and 6 evaluates the energy of the current task to be executed in device. Generally, handover
381 will be associatively managed and processed. If the evaluated value is higher, then eNB will
382 offload the task to ESPs, which is nearer to the user's location with the current value as $\gamma_i = 1$;
383 otherwise, eNB allows to perform computations locally with the value as $\gamma_i = 0$.

384 The proposed algorithm has three phases

- 385 • Initiation phase
- 386 • Discovery phase and
- 387 • Handover phase

388 In the initiation phase, classification of real-time and non-real time computations can be
389 taken place based on user demands. Discovery phase is also known as decision phase, which
390 decide whether to offload computation to ESP or allow device to execute locally. Once the
391 demand is identified, then discovery phase has to determine better ESP {n=n1...nx} with fair
392 location based on the event history of a profiler which is represented as z. Since ESPs are in
393 mesh type, it enables multi sharing and resilience with the entire satisfaction of users with n
394 nodes and m ants.

395 The optimizer uses Dynamic Offloading Decision (DOD) algorithm to determine an
396 optimal path from the possible nodes, by which user will get service either from ESP or
397 locally. Also, a scheduler will get updated with the current status of ESP computations, which
398 will be monitored by eNB. The path selection is purely based on proximity (distance) of
399 users. Also, the forwarding nodes resolve their decision and routing conflicts with tolerant
400 delay. Through proposed ASV based path selection algorithm, long communication delay due
401 to a poor communication link will be effectively managed. Hence predicted pheromone (trial)
402 could be updated in the optimal list.

Algorithm 2: DOD algorithm

Input: cluster of user requests, γ_i

Output: path ∂

```
1 Start path_selection
2 Initialize nodes as n, computations as m
3 for i= 1 to n do
4   if  $\gamma_i = 1$ 
5     Get z and predict user behavior
6     for j = 1: m do
7       if  $E_{ij} > 0$ 
8         Compute  $T^{trans}$ ,  $T^{ov}$ ,  $L^{ov}$ ,  $T^{loc}$ 
9         if  $\min \{T^{trans}_{ij} + T^{ov}_{ij}\} + T^{loc}_{ij} \leq \alpha$ 
10         $\partial = \partial \cup \{i\}$ 
11        Execute ASV algorithm
12      else
13         $\gamma_i = 0$ 
14        Compute  $L^{loc}$ 
15      end if
16       $\partial = \partial \cup \{j\}$ 
```

```

17    end if
18    end for
19    end if
20    end for
21    updated path  $\partial$ 
22    Stop path_selection

```

403

404 Generally, each eNB maintains the details of event history of each cluster and their
 405 corresponding nodes. Through which the social relations of a user (user behavior) can be
 406 predicted, and so opportunistic network with close proximity for computation can be
 407 determined. Besides, an effective handover and communication could be taken place between
 408 eNB to ESP. Our proposed is a multihop mesh type; single node failure will not make the
 409 entire network to be down (resilience). In addition, the handover phase will receive requests
 410 of user, that will be segregated and assigned (hand over) to ESP so that communication
 411 between users and computations can be done by ESP directly, which in turn reduce delays.
 412 The non-intervention process of eNB is termed as handover, and all the activities between
 413 eNB and ESP can be controlled through SDN.

Algorithm 3: ASV algorithm

Input: cluster of user requests, path ∂

Output: Optimal path S

```

1      Start optimal_path_detection
2      Initialize t computations, n nodes and m ants for each
3      partial_solution= $\partial$ 
4      while t  $\neq$  0 do
5          for i= 1to n
6              for j=1 to m
7                  Find available vertices v
8                  Initialize jth ant in the ith computations
9                  if available (v) > 1
10                 compute LESP
11                 partial_solution= min (LESP)
12                 else
13                     Ant die
14                 Return
15                 endif
16             end for
17             sort partial_solution in ascending order
18             move to the next available node

```

```
19     if partial_solution ≤ lag  
20     S = partial_solution  
21     Add the solution and report to the Scheduler  
22     for each ant that completed a solution do  
23         update S for each edge that the ant versed  
24     end for  
25     endif  
26     return S  
27     end for  
28     Stop optimal_path_detection
```

414

415 The ASV algorithm, initially make the optimal path S as an empty set. If the path has at
416 least one node then traverse for t computations, till an optimal solution S is found
417 successfully. If there is no next node, then the path of the current node is optimal. If there is
418 more than one next node to process, then compute latency of each node. Finally, compare the
419 latency values of each node in the path with other nodes. If the latency value is less than and
420 equal to the specified lag time (i.e., 20ms), then mark that path as optimal S and assign the
421 minimum latency value as the best trail value. The trial value is updated automatically like
422 the ant tabu table repository process. If a node performs partial computations and it fails,
423 remaining computations will be carried out by other ESPs with the assistance of mesh
424 network to perform successful task completion.

425 **3.2. Open Research Challenges**

426 In our paper, we consider primary challenges of generalized architecture and optimal
427 offloading. Our results show that it would be an alternative, reliable and an effective
428 offloading for both static and dynamic environments than other existing works.

429 **3.2.1. Architecture and Deployment**

430 The specification of how to perform computation offloading and the deployment location of
431 edge server within the network are not mentioned. So, investigating the optimal site selection
432 is a problem in edge server. Most of the research articles theoretically addresses edge server
433 concept, which can then be analyzed through the modern learning techniques along with

434 novel algorithms. However, the critical issues are to be validated by numerical analysis or by
435 simulations, and also essential to validate key principles and findings through simulation
436 under more complex realistic situations and scenarios. Besides, there is no clear picture about
437 the type of simulators and the place where it will be best suited. At the same time, massive
438 trials and further experiments in emulated networks or real networks are mandatory to move
439 edge server concept closer to reality.

440 Thus, choosing the optimal placement by deploying picocells with the integration of edge
441 (ESP) requires computational resource provisioning along with the deployment budget. Also,
442 determining the ESP density to cope up with their demands is also closely related to the
443 infrastructure deployment cost. Caching at the ESP named Computation Caching (CC),
444 supports sharing of multi-user to enhance the user experience and also handles the
445 management of load efficiently by providing results to end users directly without fetching
446 their required tasks beforehand. Unlike content caching, CC presents several new challenges.

447 First is a diverse type of computing tasks, which depends on the computing environment.
448 Though few contents of frequently used data are cacheable for reuse by other devices, data
449 belongs to personal computing is not cacheable and so it must often be executed in real time.
450 Second, it is not practical to build user patterns locally at each server; instead, learning the
451 intelligent techniques and methods over extensive sets can provide a broader view on the user
452 patterns.

453 *3.2.2. Offloading Decision*

454 In green networking, the energy consumption at the ESP has to be taken into account during
455 the decision. Most of the research articles dealing with the offloading decision with the
456 assumption of static scenarios, i.e., the location of UEs are fixed before and during the
457 offloading. Besides, the consumption of energy due to transmission will drastically change
458 the offloaded data. Also, channel quality drops due to low movement or fading. Hence, it is

459 necessary to propose new advanced methods such as data analytics, machine learning
460 techniques, etc. for efficient offloading decisions. For instance, collaborating various
461 prediction techniques on the UEs mobility and channel quality during the offloading provides
462 better estimation of offloading for varying conditions. A major challenge in offload decision
463 belongs to the consideration of backhaul alleviation between the ESP and ability to reflect
464 varying load and parameters during the offloading decision. This paradigm motivates for
465 rethinking and reshaping research effort from single downlink to the mixed downlink and
466 uplink in the future.

467 *3.2.3. Mobility Management*

468 Mobility is a critical issue to ensure service continuity for dynamic mobile users. Moreover,
469 efficient management in the utilization of communication and computation services are also a
470 challenging issue from the service orchestration perspective. In order to manage the recent
471 trends of learning techniques optimally, analytical mechanisms can be used by ESP to
472 estimate, predict and manage the movement of users through personal preference
473 information, which in turn improves the user experience. Besides, to achieve better user
474 computation experience, proposed offloading techniques consider mobility-aware scheduling
475 policies at the ESP. This approach introduces a set of interesting research problems including
476 mobility-aware online prefetching of user computation data, server scheduling, and fault-
477 tolerance computation.

478 *3.2.4. Security and Privacy*

479 The distributed deployment feature leads to frequent site attacks hence requires stringent
480 security policies/implementation of trust management systems. In order to gain access to the
481 platform, derive the information needed regarding user proximity and network analytics.
482 Besides, service providers would like to attain user information to update their services; then
483 there is a great challenge to the development of privacy protection mechanisms.

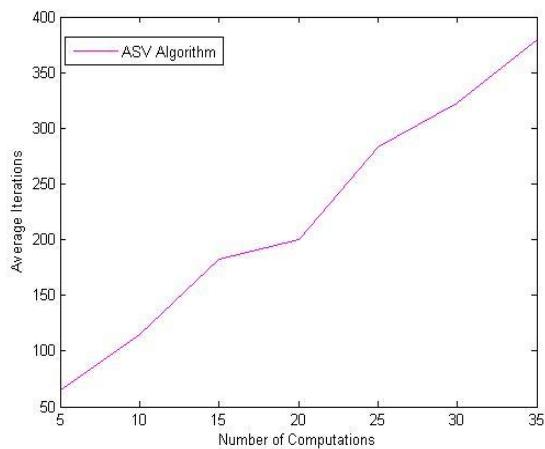
484 **4. Experimental Results**

485 As a proof of concept, we evaluate our framework in Cloud Sim with 3 ESPs and
 486 components such as a cache for prefetch of trained data and high-end Xen hypervisor server
 487 for computation. The proposed ESP mesh can support around 250 to 300 nodes. For instance,
 488 we are taking image editing process with 250 nodes, in which each server of ESP with 1Gb
 489 RAM, 1.6 TB storage capacity, an average bandwidth of 20MHz and an average image size
 490 of 1000 KB are taken. Our evaluation metrics from Table I are the real time parameters based
 491 on load per cluster. The objective is to show how the amount of offloaded tasks from the
 492 mobile devices will impact the total power consumption and the incurred delay in the system.
 493 In this paper, ESPs are in the form of a mesh topology, to reduce power consumption and
 494 increase coverage. Thus, the computations for a multi-user can be opportunistically offloaded
 495 to ESPs through wifi. Besides, the tasks during handover can be controlled through SDN.

496 **Table 1:** Simulation Parameters

Description	Value
Number of pico cells	30
Bandwidth of pico cells	20 MHz
Distance between pico to User	50m
Distance between pico to ESP	1 Km
Computation capacity of ESP	30MFLOPS
Storage capacity of ESP	100MB
Coverage diameter per ESP	100m
Number of devices per ESP	[10-30]
Number of computations per device	4
Transmission power of user device	0.5w
Computation capacity per device	[1-100] MIPS
Storage capacity per device	[1-100] Mb
Inter edge round time trip	10 ms
Simulation Iterations	20
Latency	(15, 20, 70) ms
Interval Time	100 ms
Execution time	30 ms
Resource requests	[2, 5, 7] MFLOPS
Packet size	[280,410,700] KB

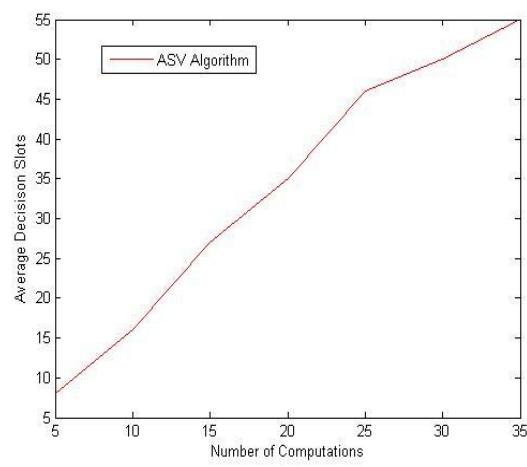
498 The effectiveness of the proposed scheme is evaluated through the experiments with
 499 varying numbers of requests. Simulation in CloudSim shows that ESP is serving most of the
 500 requests. In order to verify the effectiveness, we consider a scenario of a 5G heterogeneous
 501 ultra-dense network, where a set of 3 pico cells with a distance of 50m to users are distributed
 502 with a bandwidth of 20MHz. Each small cell is surrounded by three users with the
 503 transmission power of 0.5W, capacity of both computation and storage are 100MIPS and
 504 100Mbps respectively. Also each user could be accomplished with at most four types of
 505 computations for 100 iterations.



506

507

Fig. 5 Number of computations vs Average Iterations

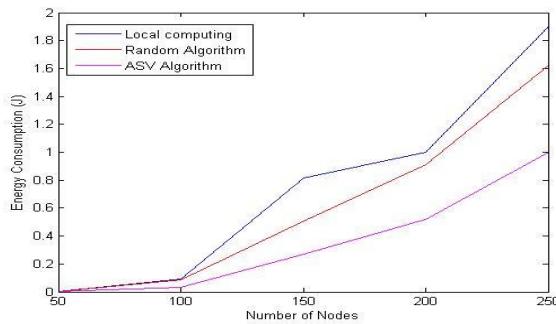


508

509

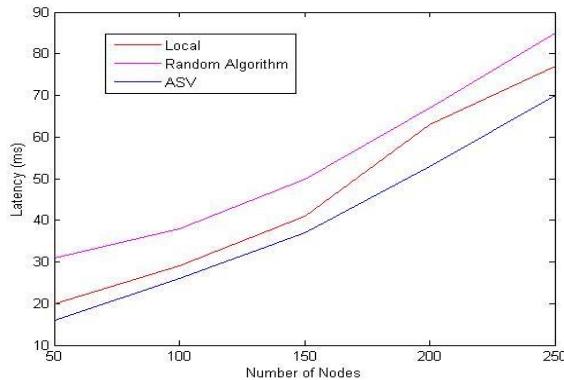
Fig. 6 Average Decision slots for convergence with different number of computations

510 Regarding the kind of computations, the size of an input data is set to be 250KB to 800
 511 KB. Hence, we conclude that lengthier size of the input packet, then longer the computations,
 512 with high energy cost. In particular, our proposed DOD offloading scheme exhibits parallel
 513 associative computations with reduced latency and energy. Multi sharing of the core cloud is
 514 showing an increase with increasing number of requests to support the ubiquitous coverage
 515 for the services. Figure 5 and Figure 6 shows the consumption of energy and average
 516 decisions of offload sing ASV with respect to number of computations. In our scenario, we
 517 are taken 3 ESPs and three users per ESP with simultaneous access of 4 tasks. So, totally 36
 518 computations can be done in ESP simultaneously with varied iterations. It is clear that the
 519 computations and its iterations are inversely proportional to each other. ESPs can perform
 520 computations simultaneously; thus energy consumption is less compared to mobile device
 521 execution.



522

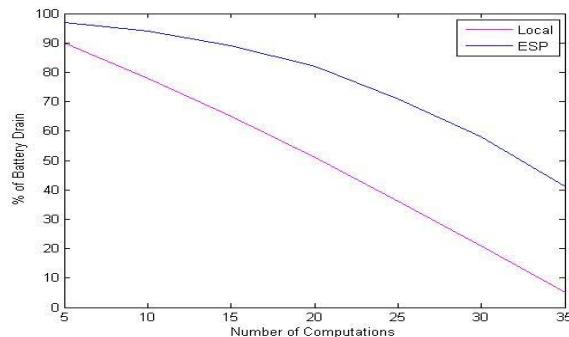
523 **Fig. 7** Comparison of energy consumption using ASV with existing mechanisms



524

525 **Fig. 8** Comparison of latency constraints in proposed ASV with existing mechanisms

526 Figure 7 and Figure 8 depicts energy consumption and latency of intensive computations
 527 that are carried out by ASV is less compared to the computations done by local and random
 528 algorithm. It is clear that the computation for up to 100 nodes of varied type's results in
 529 reduced latency by 2%, than with local and random computing; if it exceeds 120 nodes then
 530 local and random computing cannot effectively produce results, there by affects the Quality
 531 of Experience (QoE). But our proposed ASV provides a faster response with improved QoE.
 532 We can see from Figure 9 that local computations will drain the user's battery much faster
 533 than ESP. Since computations can be done in ESP in an effective manner, the battery can be
 534 prolonged. The user will get satisfied by the faster intensive computations with retaining
 535 power in their device.



536
 537 **Fig.9** Battery drain in ASV based on the number of computations
 538 Besides, increasing the percentage of requests by the mobile devices in turn will reduce
 539 total power consumption in contrast to the scenarios that need the 5G core to fulfill.
 540 However, application and content providers are challenged by the latency of the network
 541 when connecting to the CDN, which has to be resolved by today's emerging intelligent
 542 mechanism. Thus, in order to obtain a tolerable latency, combination of proximity and
 543 intelligence has been embedded in our proposed ESP.

544 **5. Conclusions**

545 Thus, the proposed MMM architecture for computation offloading provides more
 546 computational resources for faster execution. Moreover, it corroborates the need of ESP for

547 efficient computation offloading, where the distance of information-propagations are tens to
548 hundreds of meters with higher computation capability. Our solution for effective offloading
549 with ESP is unique, is not only bringing computation resources closer to mobile user's, but
550 also decouples the problem of identifying user needs through the use of prediction and
551 learning technique. Hence, our proposed MMM architecture efficiently handles both the
552 offloading and handover. In addition, it efficiently sustains with the optimal routing through
553 intelligent ASV algorithm. In particular, ESP is an emerging mechanism that can be best
554 suited in emergency situations, battlefield surveillance, retail places, stadium, shopping malls,
555 high-speed mobile video applications access to transport etc. Thus, our proposed offloading
556 scheme is a self-healing network that could be intelligently managed by licensed reusable
557 channels. As a future work, we plan to create a test bed for offloading scheme, and then to
558 expand this concept to other areas such as trust management among others.

559 **Acknowledgment**

560 D. AdhimugaSivasakthi and Gunasekaran Raja gratefully acknowledge for the support of
561 NGNLab, Department of Computer Technology, Anna University, MIT Campus, Chennai.

562 **Compliance with Ethical Standards**

563 **Funding:** No funding has been received.

564 **Disclosure of potential conflicts of interest:** Authors declare that they have no conflict of
565 interest.

566 **Ethical approval:** This article does not contain any studies with human participants
567 performed by any of the authors.

568 **Consent to participate:** Not applicable

569 **Consent for publication:** Not applicable

570 **Code availability:** Not applicable

571 ***Declaration of interests:*** The authors declare that they have no known competing financial
572 interests or personal relationships that could have appeared to influence the work reported
573 in this paper

574 **References**

- 575 [1] Golkarifard, M., Yang, J., Movaghar, A., & Hui, P. (2017). A Hitchhiker's guide to
576 computation offloading: Opinions from practitioners. *IEEE Communications
577 Magazine*, 55(7), 193-199.
- 578 [2] Khoda, M. E., Razzaque, M. A., Almogren, A., Hassan, M. M., Alamri, A., & Alelaiwi,
579 A. (2016). Efficient computation offloading decision in mobile cloud computing over
580 5G network. *Mobile Networks and Applications*, 21(5), 777-792.
- 581 [3] Mach, P., & Becvar, Z. (2017). Mobile edge computing: A survey on architecture and
582 computation offloading. *IEEE Communications Surveys & Tutorials*, 19(3), 1628-1656.
- 583 [4] Wang, S., Zhang, X., Zhang, Y., Wang, L., Yang, J., & Wang, W. (2017). A survey on
584 mobile edge networks: Convergence of computing, caching and communications. *Ieee
585 Access*, 5, 6757-6779.
- 586 [5] Mao, Y., You, C., Zhang, J., Huang, K., & Letaief, K. B. (2017). A survey on mobile
587 edge computing: The communication perspective. *IEEE Communications Surveys &
588 Tutorials*, 19(4), 2322-2358.
- 589 [6] <http://blogs.valvesoftware.com/abrash/latency-the-sine-qua-non-of-ar-and-vr/>.
- 590 [7] <https://www.bayshorenetworks.com/blog/breaking-down-idc-top-10-iot-predictions-for->
591 2017
- 592 [8] Mao, Y., Zhang, J., & Letaief, K. B. (2016). Dynamic computation offloading for
593 mobile-edge computing with energy harvesting devices. *IEEE Journal on Selected
594 Areas in Communications*, 34(12), 3590-3605.

- 595 [9] Bhattacharya, A., & De, P. (2017). A survey of adaptation techniques in computation
596 offloading. *Journal of Network and Computer Applications*, 78, 97-115.
- 597 [10] Chen, X., Jiao, L., Li, W., & Fu, X. (2015). Efficient multi-user computation offloading
598 for mobile-edge cloud computing. *IEEE/ACM Transactions on Networking*, 24(5),
599 2795-2808.
- 600 [11] Kottursamy, K., Raja, G., Padmanabhan, J., & Srinivasan, V. (2017). An improved
601 database synchronization mechanism for mobile data using software-defined
602 networking control. *Computers & Electrical Engineering*, 57, 93-103.
- 603 [12] Kottursamy, K., Raja, G., & Saranya, K. (2016). A data activity-based server-side cache
604 replacement for mobile devices. In *Artificial Intelligence and Evolutionary
605 Computations in Engineering Systems* (pp. 579-589). Springer, New Delhi.
- 606 [13] Elgazzar, K., Martin, P., & Hassanein, H. S. (2014). Cloud-assisted computation
607 offloading to support mobile services. *IEEE Transactions on Cloud Computing*, 4(3),
608 279-292.
- 609 [14] Liu, W., Gong, W., Du, W., & Zou, C. (2017, February). Computation offloading
610 strategy for multi user mobile data streaming applications. In *2017 19th International
611 Conference on Advanced Communication Technology (ICACT)* (pp. 111-120). IEEE.
- 612 [15] Saguna and Intel, "Using mobile edge computing to improve mobile network
613 performance and profitability," White paper, 2016
- 614 [16] Hu, Y. C., Patel, M., Sabella, D., Sprecher, N., & Young, V. (2015). Mobile edge
615 computing—A key technology towards 5G. *ETSI white paper*, 11(11), 1-16.
- 616 [17] Bharathi U., Mercy Faustina J., Elavarasi B & Gunasekaran R. (2017)Energy Efficient
617 Data Offloading Mechanism,' in proc. IEEE International Conference on
618 Telecommunication, Power Analysis and Computing Techniques(ICTPACT), Apr-6-8,

- 619 [18] Kumar, K., Liu, J., Lu, Y. H., & Bhargava, B. (2013). A survey of computation
620 offloading for mobile systems. *Mobile networks and Applications*, 18(1), 129-140.
- 621 [19] Jiao, L., Friedman, R., Fu, X., Secci, S., Smoreda, Z., & Tschofenig, H. (2013, July).
622 Cloud-based computation offloading for mobile devices: State of the art, challenges and
623 opportunities. In *2013 Future Network & Mobile Summit* (pp. 1-11). IEEE.
- 624 [20] Wang, S., Zhou, A., Hsu, C. H., Xiao, X., & Yang, F. (2015). Provision of data-
625 intensive services through energy-and QoS-aware virtual machine placement in national
626 cloud data centers. *IEEE Transactions on Emerging Topics in Computing*, 4(2), 290-
627 300.
- 628 [21] Lin, Y. D., Chu, E. T. H., Lai, Y. C., & Huang, T. J. (2013). Time-and-energy-aware
629 computation offloading in handheld devices to coprocessors and clouds. *IEEE Systems
630 Journal*, 9(2), 393-405.
- 631 [22] Ragona, C., Granelli, F., Fiandrino, C., Kliazovich, D., & Bouvry, P. (2015, December).
632 Energy-efficient computation offloading for wearable devices and smartphones in
633 mobile cloud computing. In *2015 IEEE Global Communications Conference
634 (GLOBECOM)* (pp. 1-6). IEEE.
- 635 [23] Goudarzi, M., Movahedi, Z., & Nazari, M. (2016, July). Efficient multisite computation
636 offloading for mobile cloud computing. In *2016 Intl IEEE Conferences on Ubiquitous
637 Intelligence & Computing, Advanced and Trusted Computing, Scalable Computing and
638 Communications, Cloud and Big Data Computing, Internet of People, and Smart World
639 Congress (UIC/ATC/ScalCom/CBDCom/IoP/SmartWorld)* (pp. 1131-1138). IEEE.
- 640 [24] Mukherjee, A., De, D., & Roy, D. G. (2016). A power and latency aware cloudlet
641 selection strategy for multi-cloudlet environment. *IEEE Transactions on Cloud
642 Computing*, 7(1), 141-154.

- 643 [25] Zhang, K., Mao, Y., Leng, S., Zhao, Q., Li, L., Peng, X., ... & Zhang, Y. (2016).
644 Energy-efficient offloading for mobile edge computing in 5G heterogeneous
645 networks. *IEEE access*, 4, 5896-5907.
- 646 [26] Yang, L., Cao, J., Wang, Z., & Wu, W. (2017, August). Network aware multi-user
647 computation partitioning in mobile edge clouds. In *2017 46th International Conference*
648 *on Parallel Processing (ICPP)* (pp. 302-311). IEEE.
- 649 [27] Schneider, M., Rambach, J., & Stricker, D. (2017, March). Augmented reality based on
650 edge computing using the example of remote live support. In *2017 IEEE International*
651 *Conference on Industrial Technology (ICIT)* (pp. 1277-1282). IEEE.
- 652 [28] Soh, H., & Demiris, Y. (2014). Incrementally learning objects by touch: Online
653 discriminative and generative models for tactile-based recognition. *IEEE transactions*
654 *on haptics*, 7(4), 512-525.
- 655 [29] Ganz, F., Puschmann, D., Barnaghi, P., & Carrez, F. (2015). A practical evaluation of
656 information processing and abstraction techniques for the internet of things. *IEEE*
657 *Internet of Things journal*, 2(4), 340-354.
- 658 [30] Cao, H., & Cai, J. (2017). Distributed multiuser computation offloading for cloudlet-
659 based mobile cloud computing: A game-theoretic machine learning approach. *IEEE*
660 *Transactions on Vehicular Technology*, 67(1), 752-764.

Figures

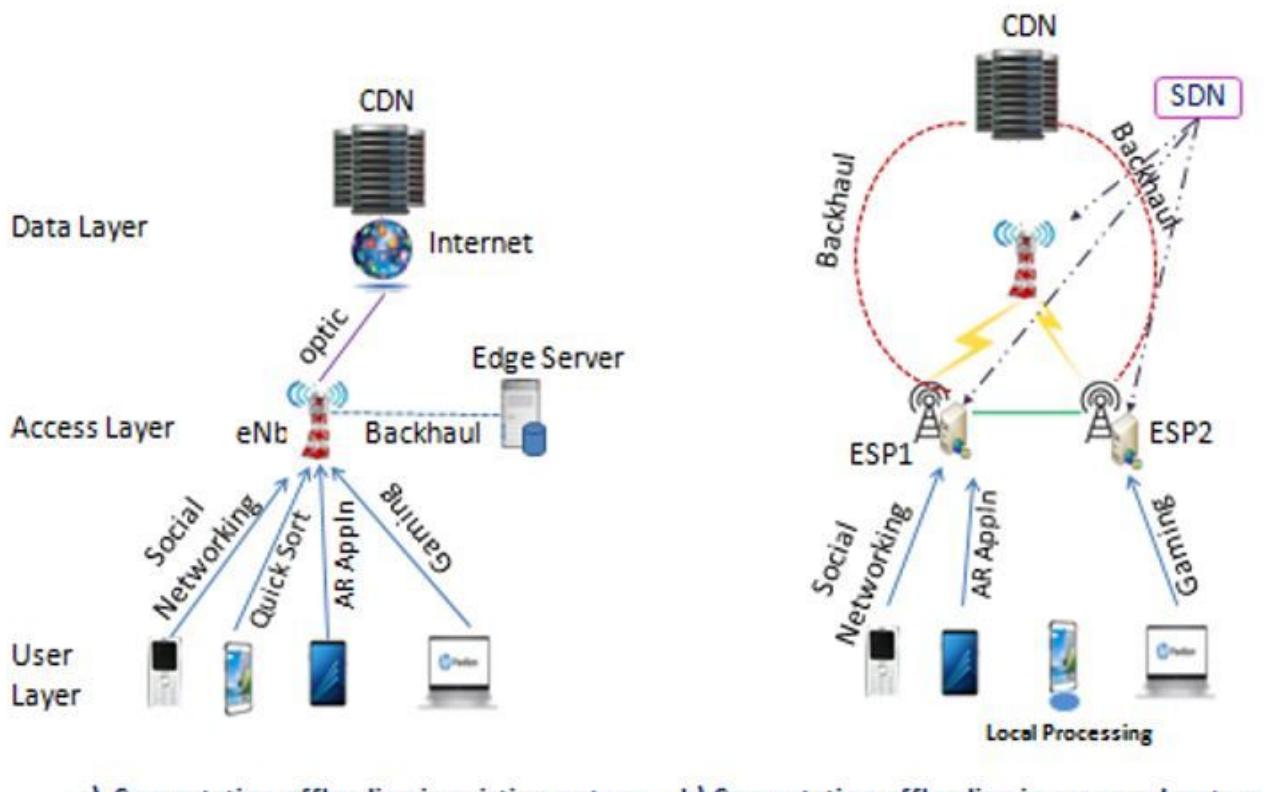


Figure 1

Computation offloading Scenario a) Existing and b) Proposed Framework

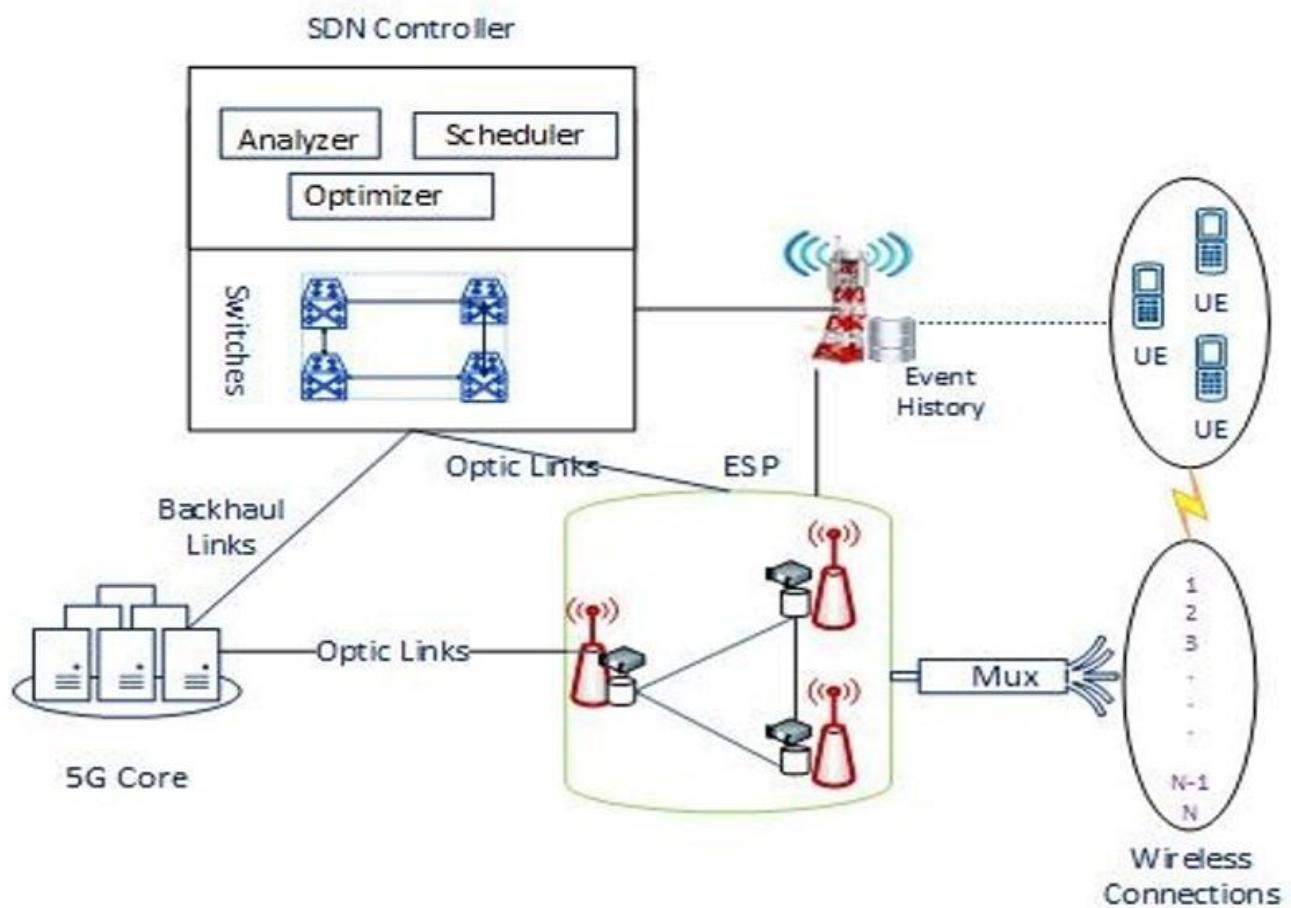


Figure 2

MMM Architecture for Computation Offloading

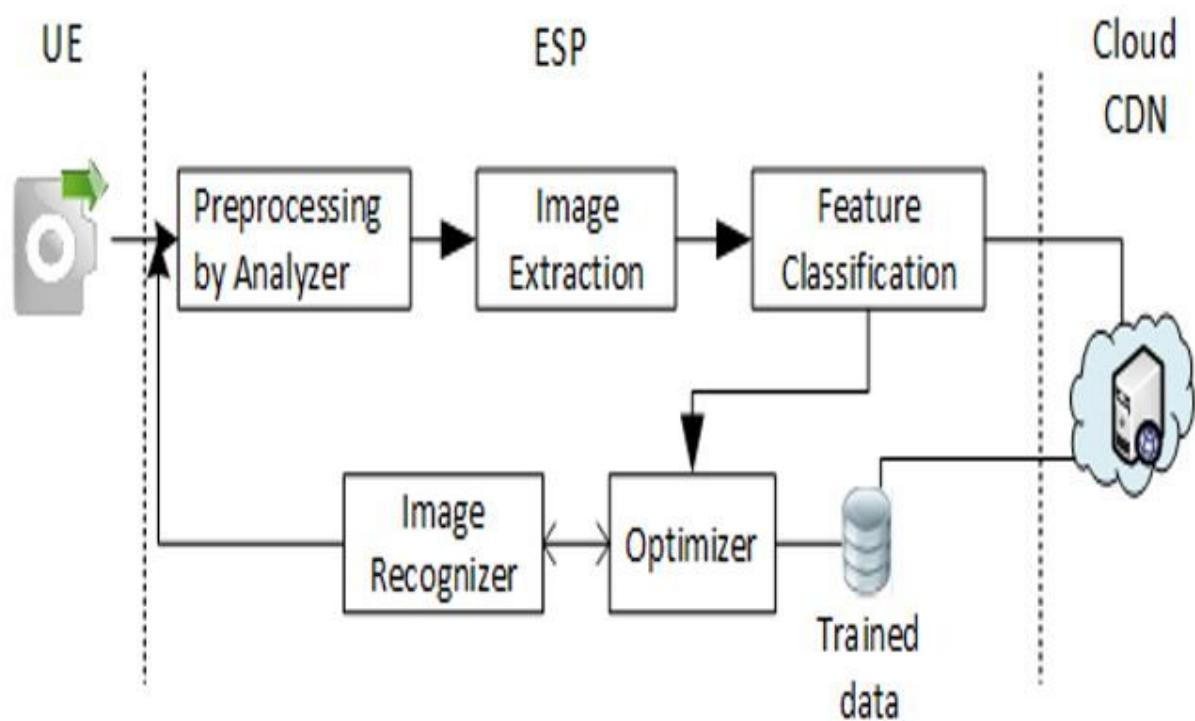


Figure 3

An example of mobile image editing with ESP

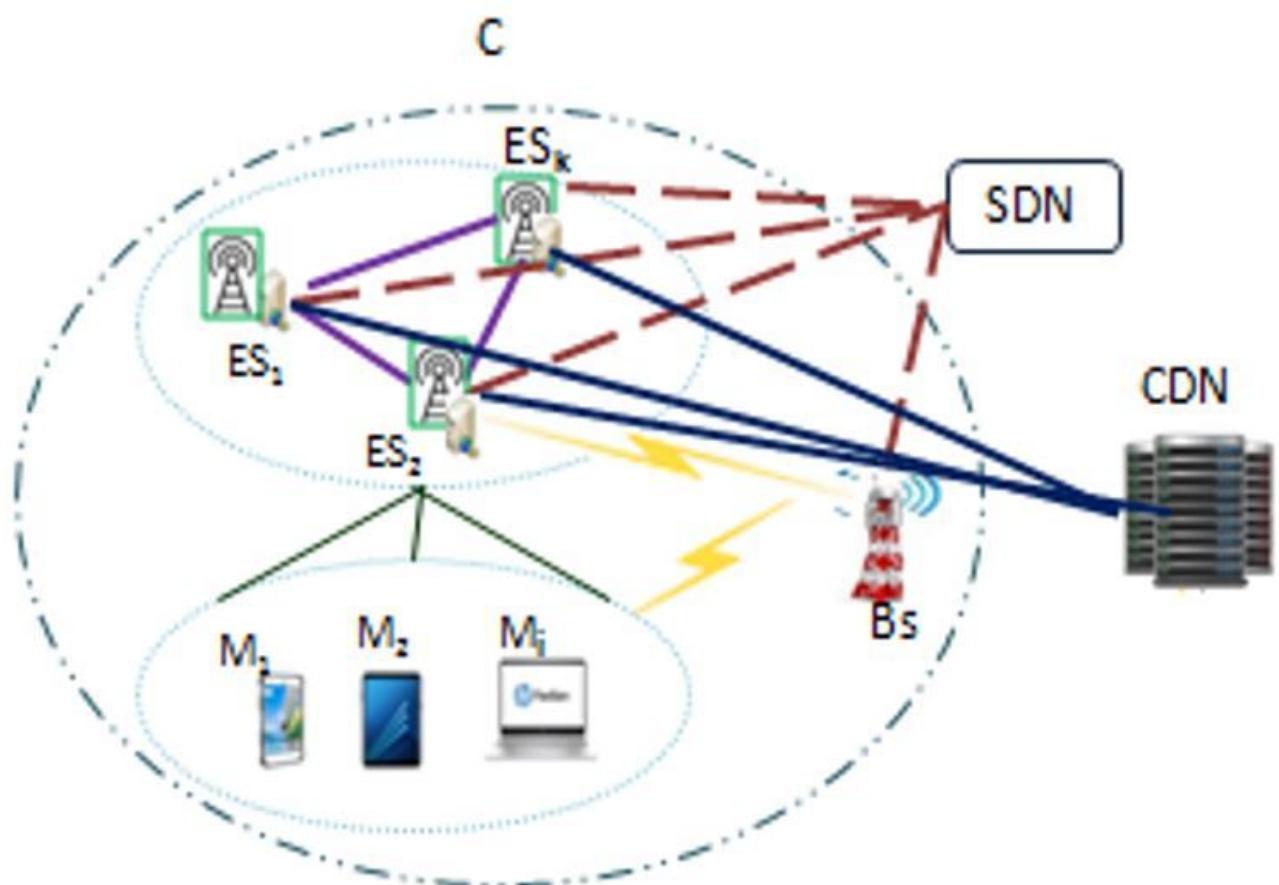


Figure 4

System Model for computation

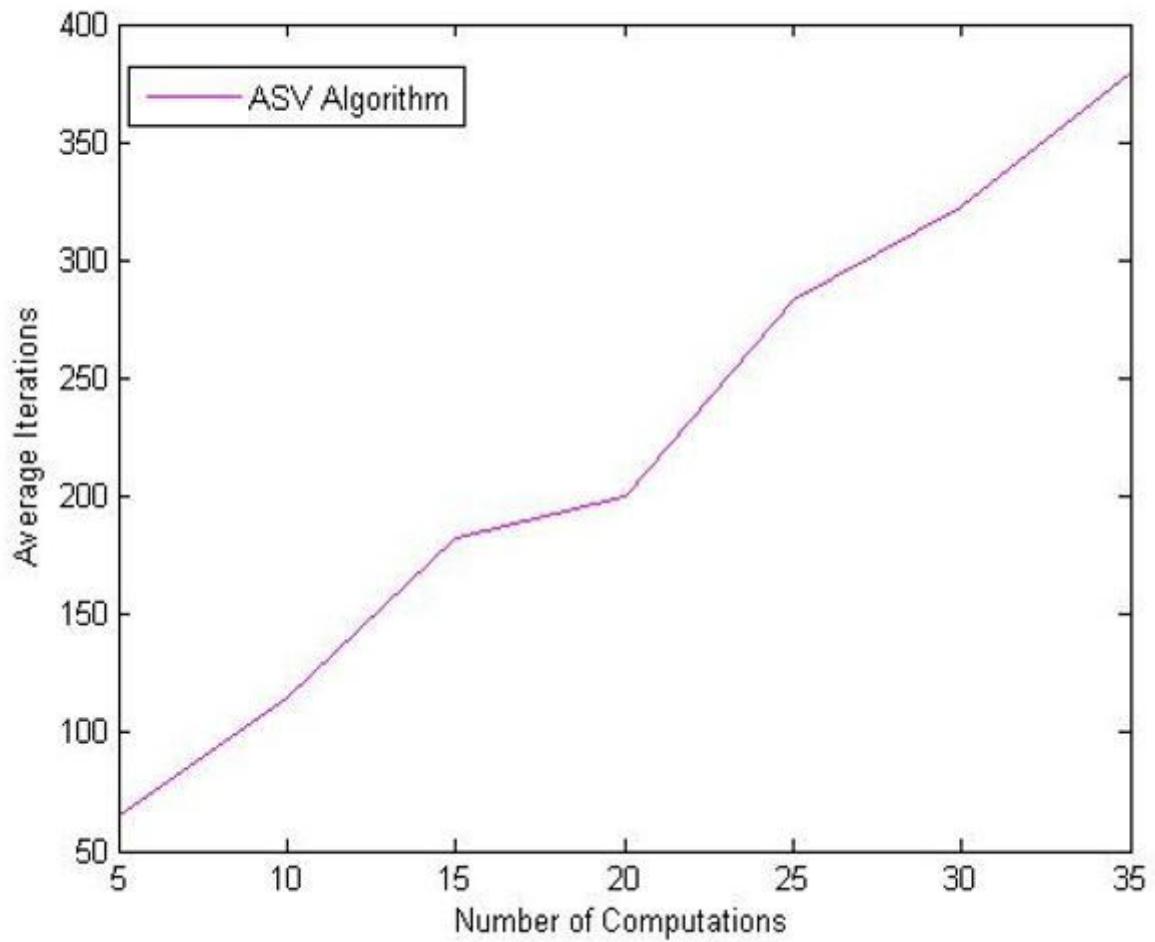


Figure 5

Number of computations vs Average Iterations

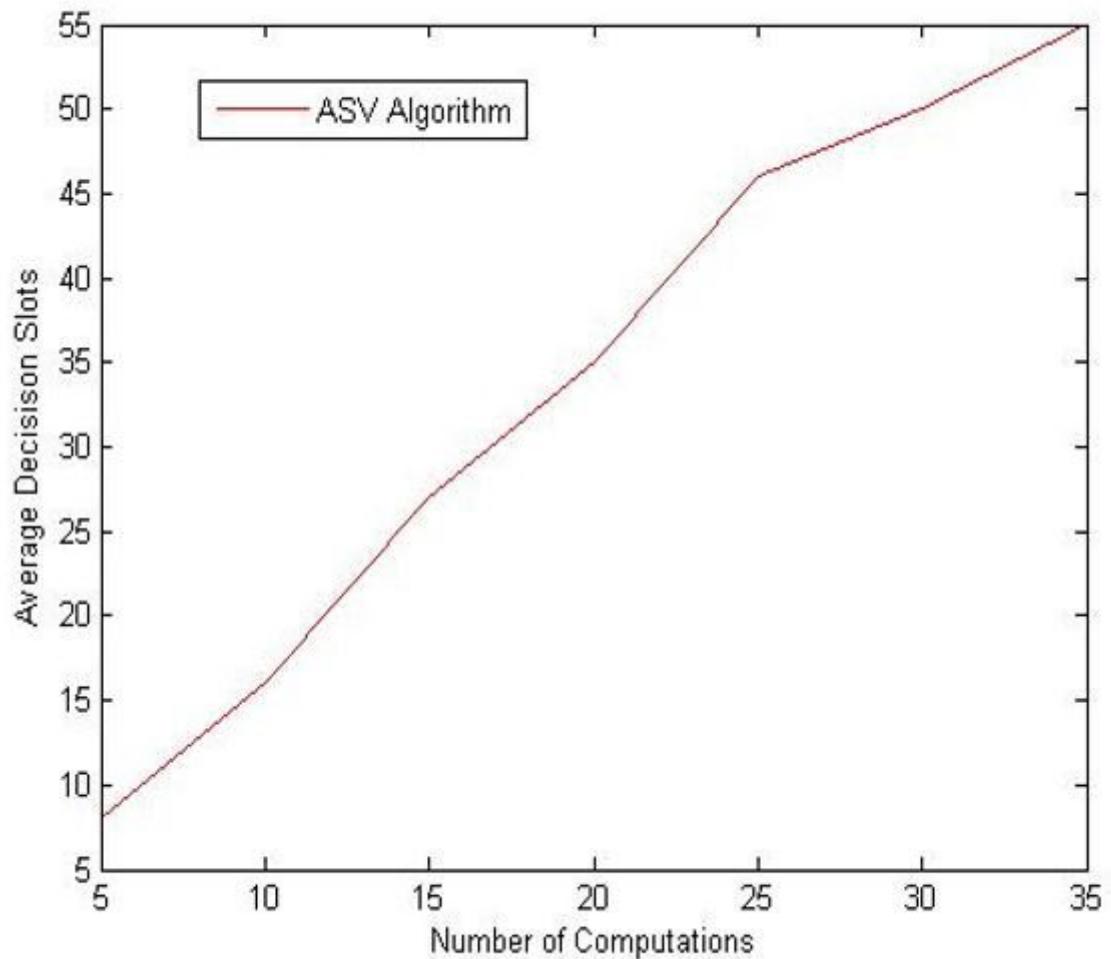


Figure 6

Average Decision slots for convergence with different number of computations

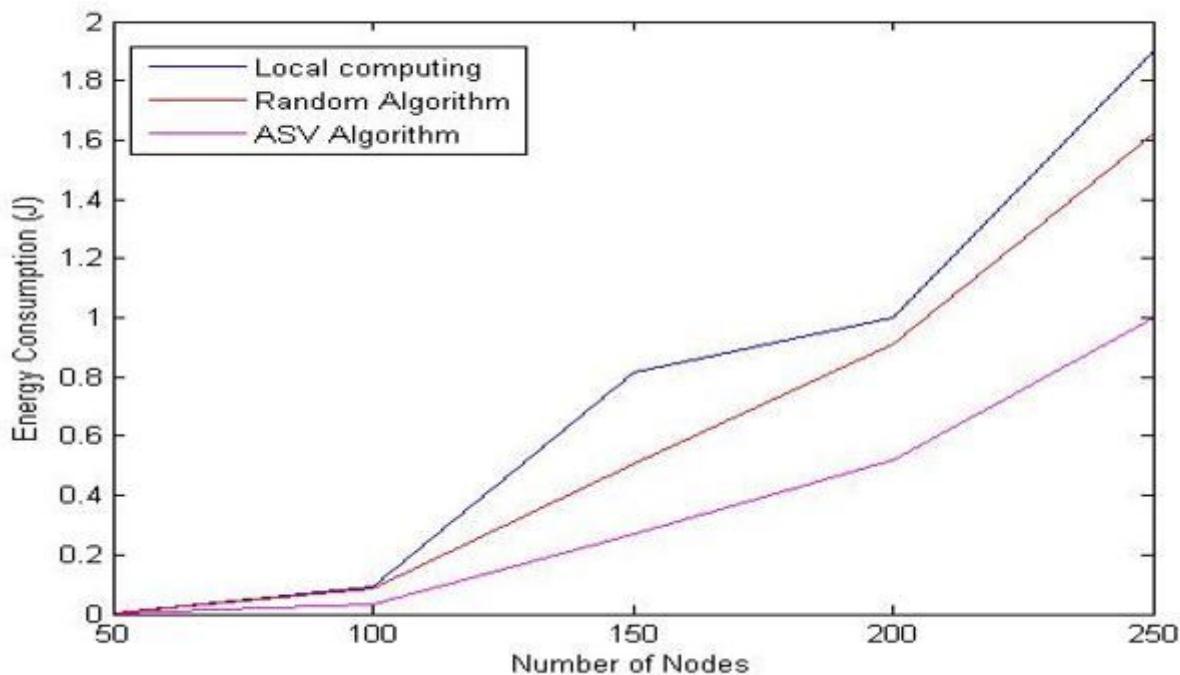


Figure 7

Comparison of energy consumption using ASV with existing mechanisms

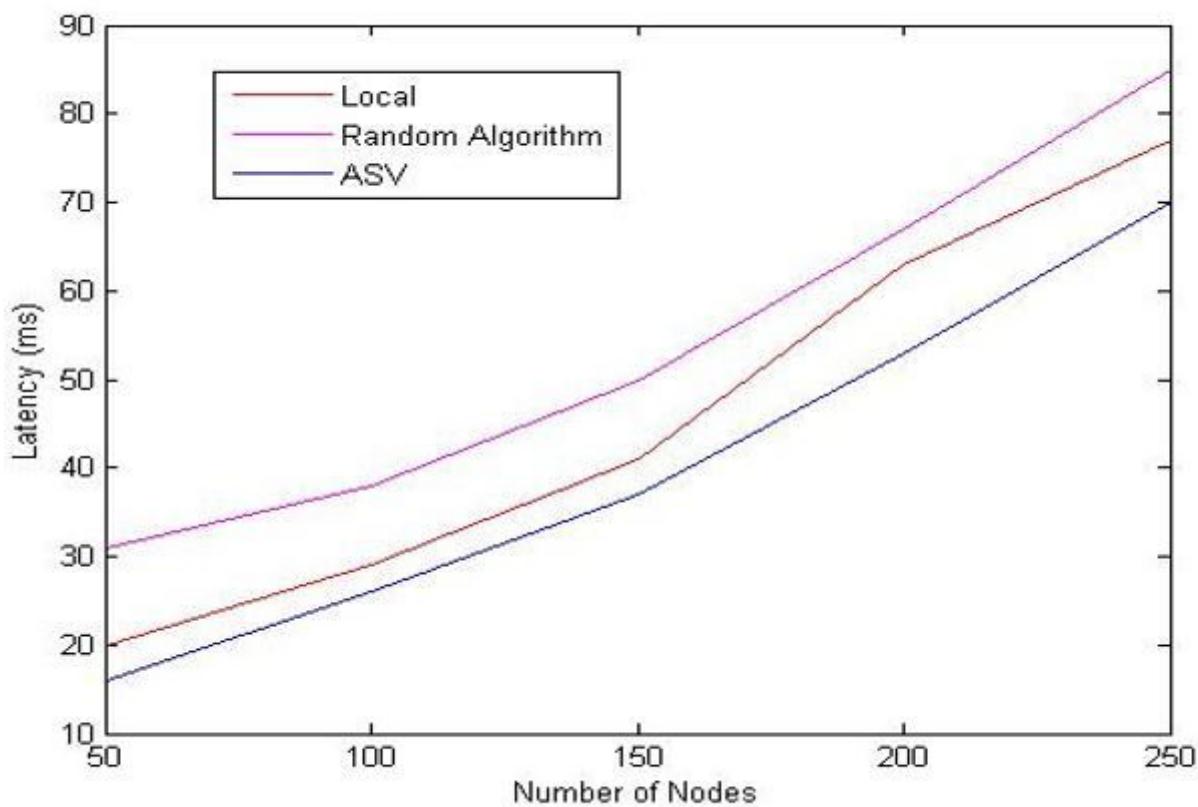


Figure 8

Comparison of latency constraints in proposed ASV with existing mechanisms

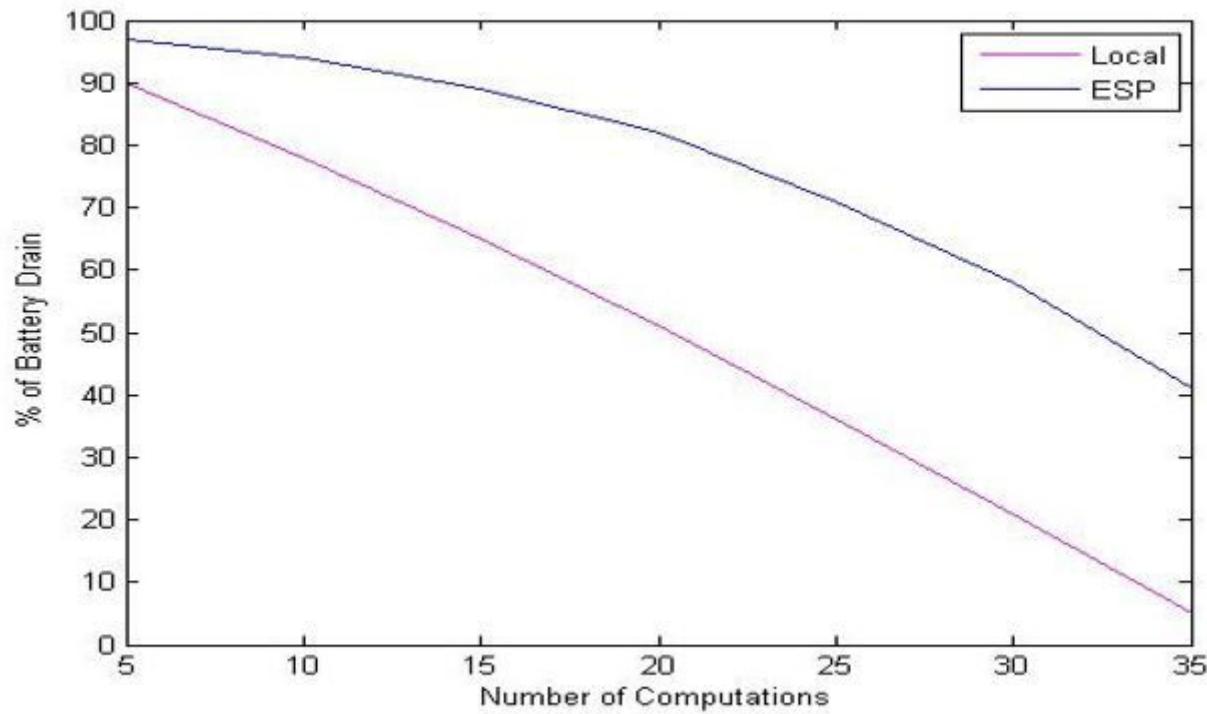


Figure 9

Battery drain in ASV based on the number of computations