

An Optimized and Distributed Data Packet Forwarding Scheme in LTE/LTE-A Networks

Haneul Ko, Giwon Lee, Dongeun Suh, Sangheon Pack, and Xuemin Sherman Shen

Abstract—Long-term evolution (LTE)/LTE-advanced (LTE-A) networks have recently introduced a data packet forwarding scheme between evolved node Bs (eNBs) to reduce the signaling overhead and the delay incurred in the data path switching scheme which is a baseline handover scheme in LTE/LTE-A networks. Even with the data packet forwarding scheme, if the length of the forwarding chain is set inappropriately, the data packet forwarding scheme suffers from the degraded throughput. To attain the optimal handover performance in terms of throughput and delay and reduce the signaling overhead, we propose an optimized and distributed data packet forwarding scheme where the optimal length of the forwarding chain is obtained by a Markov decision process (MDP). Also, a low-complexity value iteration algorithm is devised to solve the optimality equation of MDP in a more practical manner. Real trace-driven evaluation results demonstrate that the proposed scheme determines the optimal length of the forwarding chain adaptively to applications' quality of service (QoS) requirements and reduces the signaling overhead and delay while achieving higher throughput in diverse environments.

Index Terms—Data packet forwarding, data path switching, LTE/LTE-A, handover, Markov decision process (MDP), trace-driven evaluation.

I. INTRODUCTION

Recent report [2] predicts that mobile data traffic will hit an annual run rate of 134.4 exabytes by 2017. More specifically, a compound annual growth rate in the data traffic from 2012 to 2017 is 66 percent. Also, the number of personal devices connected to mobile networks by 2017 will exceed 10 billion [3]. To handle the mobile data explosion problem, small cells will be densely deployed in future wireless/mobile networks and frequent handovers among small cells will be experienced [4]–[7]. In order to improve users' quality of experience (QoE) and to mitigate the signaling overhead to core network entities under frequent handovers, efficient handover procedures should be devised.

Figure 1 shows a long term evolution (LTE)/LTE-advanced (LTE-A) architecture [8], [9]. An evolved node B (eNB) is a base station providing a wireless connection between a mobile

node (MN) and LTE/LTE-A networks. A serving gateway (S-GW) acts as a mobility anchor point when the MN moves to another eNB. On the other hand, a mobility management entity (MME) provides mobility functions such as paging, tracking area list (TAL) management, and handover management. A packet data network (PDN) gateway (P-GW) provides access to PDNs by assigning an IP address to the MN and serves as a mobility anchor point for 3GPP and non-3GPP handover.

In LTE/LTE-A networks, a handover procedure is generally based on the data path switching [9]. When an MN moves from a source eNB to a target eNB and establishes a radio connection with the target eNB, the target eNB requests a data path switching to the MME and the MN can receive data packets through the optimal path (i.e., S-GW - target eNB - MN). However, frequent data path switchings cause a significant signaling overhead to core network entities. Also, since the delay for the data path switching is not negligible, the users' perceived QoE can be degraded [15]. To solve these problems, several works have been conducted [10]–[16]. For example, a data packet forwarding scheme, where data packets are forwarded through an interface between two eNBs (i.e., X2 interface in LTE/LTE-A networks), does not involve any core network entities (e.g., MME) and the handover management can be locally processed. Therefore, the delay and the signaling overhead to core network entities can be reduced. However, as the length of the data packet forwarding chain increases¹, the performance of the data packet forwarding scheme degrades due to the inevitable packet tunneling overhead. To sum up, there is a tradeoff between the data path switching and the data packet forwarding. When the data path switching is conducted, the packets are delivered through the optimal routing path. Therefore, the throughput can be maximized and the local traffic overhead can be reduced; however, the signaling overhead to core network entities increases and the delay owing to the data path switching is expected. Meanwhile, when the data packet forwarding is employed, the signaling overhead to core network entities and the delay for the data path switching can be reduced at the expense of the non-optimal routing path. Therefore, the length of the forwarding chain should be maintained to attain the optimal handover performance.

In this paper, we propose an optimized and distributed data packet forwarding scheme in LTE/LTE-A networks. Specifically, the optimal policy between the data path switching and

A preliminary version of this paper was presented at IEEE International Conference on Communications (ICC) 2014, Sydney, Australia, June 2014 [1].

This work was supported by National Research Foundation of Korea Grant funded by the Korean Government (NRF-2014R1A2A1A12066986 and NRF-2014K1A3A1A21001357).

Copyright (c) 2015 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

H. Ko, G. Lee, D. Suh, and S. Pack are with the School of Electrical Engineering, Korea University, Seoul, Korea. (e-mail:{st_basket, goodkw, fever1989, shpack}@korea.ac.kr) and X. Shen is with Department of Electrical and Computer Engineering, University of Waterloo, ON, Canada. (e-mail:xshen@bcr.uwaterloo.ca)

¹The data packet forwarding chain denotes the tunnel which is constructed through the moving route of an MN. For example, when the MN moves from eNB 2 to eNB 3 through eNB 5, the data packet forwarding chain is constructed as follows: eNB 2 - eNB 5 - eNB 3.

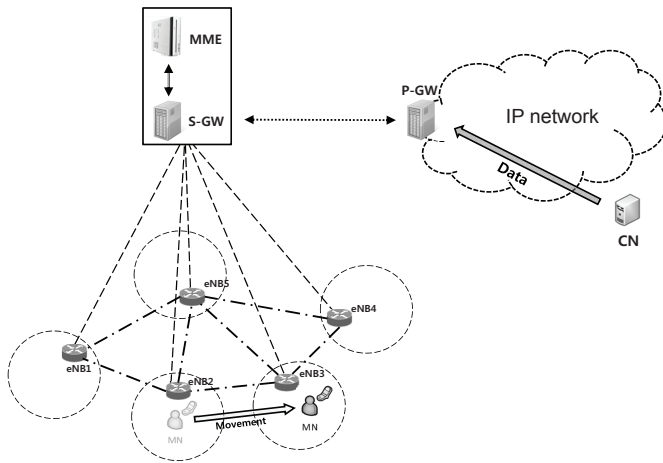


Fig. 1. LTE/LTE-A network architecture.

the data packet forwarding is formulated by a Markov decision process (MDP). In the proposed scheme, the optimal policy to maximize the reward function is obtained by a low-complexity value iteration algorithm. Real trace-driven evaluation results demonstrate that the proposed scheme determines the optimal length of the forwarding chain adaptively to applications' requirements and reduces the signaling overhead and delay while achieving higher throughput in diverse environments.

The main contribution of this paper is three-fold: 1) to the best of our knowledge, this is the first work on the optimization of the data packet forwarding scheme by means of the MDP formulation. The optimal length of the forwarding chain depending on the application's quality of service (QoS) requirement is obtained and then the handover performance in terms of the throughput and delay can be improved while reducing the signaling overhead; 2) since the conventional algorithm to solve the MDP problem has high complexity which is a major obstacle of applying MDP to practical systems, we introduce a low-complexity value iteration algorithm. The devised low-complexity value iteration algorithm allows to solve the optimality equation of MDP in a more practical manner; and 3) extensive evaluation results based on the empirical data sets are presented and analyzed to assess the performance of the proposed scheme.

The remainder of this paper is organized as follows. Related works are summarized in Section II. The comparison between the data path switching and the data packet forwarding is given in Section III and the MDP for the optimized and distributed data forwarding scheme is formulated in Section IV. Finally, real trace-driven evaluation results and concluding remarks are given in Section V and Section VI, respectively.

II. RELATED WORKS

The data path switching causes a significant signaling overhead to core network entities and long delay. To resolve these problems, several schemes were proposed in the literature. These works can be categorized into: 1) enhanced data path switching schemes [10]–[13]; 2) data packet forwarding schemes [14]–[16].

Rath and Panwar [10] proposed a proactive data multicast scheme in which a set of candidate target eNBs is maintained based on the radio signal measurement and the data is multicasted to all eNBs before completing a handover event to reduce the service interruption time. On the other hand, Kim *et al.* [11] proposed a data multicast method based on the speed of the MN to avoid too early handover trigger for low-speed MNs and too late handover trigger for high-speed MNs. However, when the expected handover does not occur, the multicasted data are simply wasted and thus unnecessary traffic overhead can occur in these schemes [10], [11]. Guo *et al.* [12] proposed a reactive data bicasting scheme to reduce the delay during the handover by making relatively moderate modifications to the 3GPP handover procedure. In this scheme, since the data bicasting is initiated after the handover, unnecessary traffic overhead does not occur while the signalling overhead to core network entities cannot be diminished. Pacifico *et al.* [13] introduced a fast data path switching scheme to reduce the handover delay. Specifically, the target eNB can send a data path switching request message immediately after the handover command even when a handover confirmation messages is not received. However, the signalling overhead to core network entities was not addressed in this scheme.

A few data packet forwarding schemes were also proposed in [14]–[16]. Yan *et al.* [14] suggested a state-aware pointer forwarding scheme. A pointer forwarding chain between mobility anchors (MAs) is established to reduce the signaling overhead for the data path switching by considering the current mobility state of the MN in deciding whether the forwarding chain should be prolonged or refreshed. Guo *et al.* [15] proposed two local mobility management schemes (i.e., traffic forwarding with the cascading path (*TF_CP*) and traffic forwarding with the shortest path (*TF_SP*)). In *TF_CP*, only when the length of the forwarding chain is smaller than a predefined threshold, the data packet forwarding is conducted. Otherwise, the target eNB triggers the data path switching. On the other hand, in *TF_SP*, the path lengths of the data path switching and the data packet forwarding are compared and the target eNB chooses a shorter one. Lee *et al.* introduced a dynamic data packet forwarding scheme for network-based distributed mobility management (DMM) [16]. In the proposed scheme, forwarding chains are dynamically established among mobility anchors by defining a session-to-mobility ratio (SMR) that considers both the numbers of session arrivals and handovers. In other words, if the SMR is higher than a predefined threshold, the data path switching operation is conducted to reduce the forwarding (or tunneling) overhead. Otherwise, the data packet forwarding is performed to reduce the signaling overhead due to the data path switching. However, how to optimize the performance of the data packet forwarding scheme (i.e., how to obtain the optimal length of the forwarding chain) was not investigated in these previous works.

III. DATA PATH SWITCHING VS. DATA PACKET FORWARDING

A handover message flow in the data path switching scheme is described in Figure 2. First, a source eNB decides a

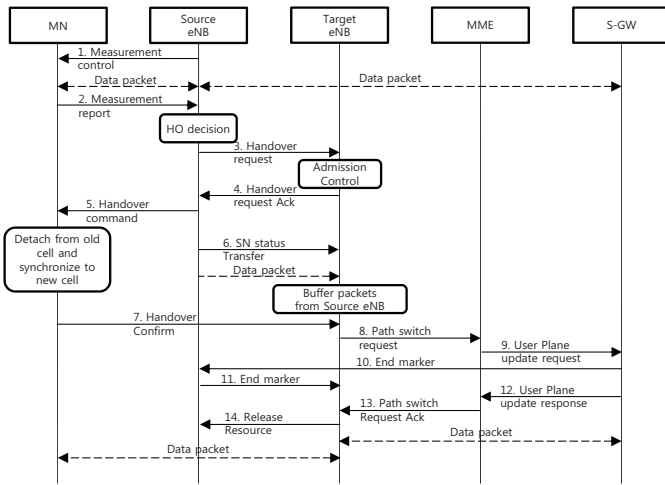


Fig. 2. Data path switching scheme.

handover through a channel measurement procedure with an MN (Steps 1-2). If the channel condition to a target eNB is better than that to the source eNB, the source eNB decides a handover to the target eNB by sending a handover request message to the target eNB (Step 3), which in turn performs admission control for the handover request message. If the handover request is accepted, the source eNB receives a handover request acknowledgement message and commands the handover to the MN (Steps 4-5). After that, the MN conducts synchronization to the target eNB. Meanwhile, the source eNB sends a serial number (SN) status transfer message and data packets to the target eNB (Step 6). When the target eNB receives a handover confirmation message from the MN (Step 7), the target eNB sends a path switch request message to the MME (Step 8). Then, the MME sends a user update request message to the S-GW (Step 9), which in turn switches the data path from the source eNB to the target eNB. On the other hand, an end marker message for releasing resource for the MN is sent to the source eNB (Step 10). Finally, the MN can receive data packets through the optimal path (i.e., S-GW - target eNB - MN). However, when the mobility of the MN is high, the significant signaling overhead for the data path switching can occur. Also, since data packets are forwarded over the detour path (i.e., S-GW - source eNB - target eNB - MN) before the optimal path is set, data packets through the optimal path can be reached earlier than the data packets through the detour path. Therefore, out of order packet delivery can occur. In addition, the delay jitter for the first packet after the data path switching may be large if the new path has long delay, which is particularly important in delay-sensitive applications such as VoIP and multimedia service [15], [17].

The above-mentioned problem can be mitigated by the data packet forwarding scheme. Figure 3 shows a handover message flow in the data packet forwarding scheme. In the data packet forwarding scheme, a target eNB does not send any data path switching request message. Instead, a local traffic forwarding chain is constructed to transmit data packets between the source eNB and the target eNB for a handover event. Therefore, data packets can be processed in a distributed

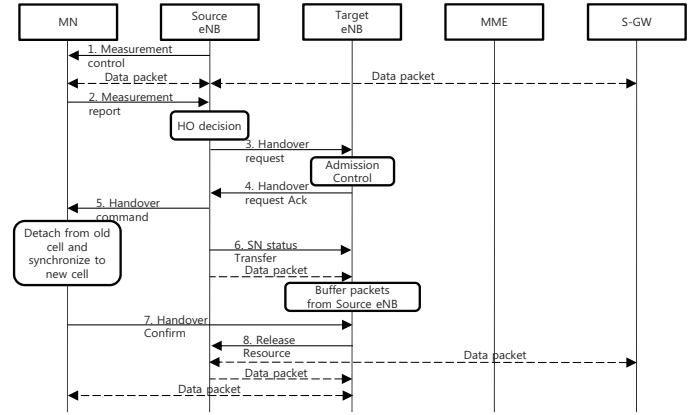


Fig. 3. Data packet forwarding scheme.

manner. However, when the session length is long in the data packet forwarding scheme, a longer data packet forwarding chain can be formed and significant forwarding/tunneling overhead over the non-optimal route is expected.

IV. OPTIMIZED AND DISTRIBUTED DATA FORWARDING SCHEME

Apparently, the data packet forwarding scheme has pros and cons compared with the data path switching scheme. For frequent handovers, the data packet forwarding scheme can reduce the signaling overhead to core network entities. However, a long data packet forwarding chain can increase the tunneling overhead and the delivery latency. Therefore, the maximum length of the data packet forwarding chain should be determined to an appropriate value. In addition, the optimal length of the data packet forwarding chain can be varied depending on the application requirements and thus an application-aware approach in determining the optimal chain length should be investigated. To this end, we formulate an MDP model² with five elements: 1) decision epochs; 2) states; 3) actions; 4) transition probabilities; and 5) rewards (costs) [18], [19]. After that, the optimality equation and a low-complexity algorithm to solve the equation are introduced. Important notations for the MDP model are summarized in Table I.

A. Decision Epochs

Figure 4 shows the timing diagram for the MDP model. A sequence $T = \{1, 2, 3, \dots, E\}$ represents the time epochs when successive decisions are made. A random variable E is the termination time of the session, which follows a geometric distribution with mean $1/\lambda_s$ [20]. Random variables X_t and Y_t denote the state and the action chosen at the decision epoch $t \in T$, respectively. The duration of each decision epoch is τ . In the proposed scheme, an action (i.e., the data packet forwarding or the data path switching) is carried out only when

²The MDP model is a mathematical framework to model a decision making in situations where outcomes are partially random and partially under the control of the decision maker. Thus, the MDP model is suitable for determining the optimal chain length in the data packet forwarding scheme.

TABLE I
SUMMARY OF NOTATIONS.

Notation	Meaning
E	Termination time of the session
X_t	State at the decision epoch t
Y_t	Action chosen at the decision epoch t
τ	Duration of each decision epoch
S	State space
H	Handover trigger phase state
F	Forwarding list state
C	Connection information state
A	Action set
P_{ij}	Probability that an MN moves from eNB i to eNB j
T_{PF}	Unit traffic overhead to transmit a data packet through the data packet forwarding chain
T_{S1}	Unit traffic overhead to conduct the data path switching
D_T	Throughput sensitivity of the session
D_D	Delay sensitivity of the session

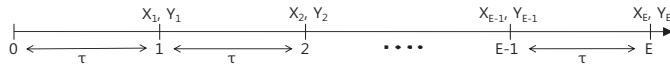


Fig. 4. Timing diagram.

the MN performs a handover to the target eNB, and therefore no actions occur at the decision epochs without handover events.

B. States

The state space S is defined as

$$S = H \times F \times C \quad (1)$$

where H denotes the handover trigger phase, F denotes the vector set of possible forwarding lists, and C means the vector set that describes the connection information of each eNB and the location of an MN.

First, H can be defined as

$$H = \{0, 1, 2\} \quad (2)$$

where $h (\in H) = 0$ and $h = 1$ represent that the MN is in the non-handover trigger and the handover trigger phases, respectively. On the other hand, $h = 2$ refers to the situation when the target eNB receives a handover confirmation message (i.e., the situation right after handover to the target eNB).

A forwarding list consists of eNBs that forward data packets to the MN and thus the forwarding list can be changed either by data path switching or handover events. The vector set of possible forwarding lists, F , is described as

$$F = \{F^1, F^2, F^3, \dots, F^{max}\} \quad (3)$$

where max is the total number of possible forwarding lists. When the total number of eNBs is N , a forwarding list with the length of M can be constructed by selecting $M (< N)$ eNBs among N eNBs and enumerating the selected M eNBs.

The numbers of the selections and enumerations are given by $\binom{N}{M}$ and $M!$, respectively. Then, max can be obtained by

$$max = \begin{cases} 1 + \sum_{M=1}^{N-1} \binom{N}{M} M!, & \text{if } N \neq 1 \\ 1, & \text{if } N = 1 \end{cases} \quad (4)$$

Let F^k be the vector representing the k th possible forwarding list. Then, F^k is defined as

$$F^k = [f_1^k, f_2^k, f_3^k, \dots, f_m^k] \quad (5)$$

where $1 \leq k \leq max$, $1 \leq m \leq N - 1$, and f_r^k ($1 \leq r \leq m$) is the identification of the r th eNB in the k th forwarding list. m is the dimension of the vector F^k , which represents the number of tunnels in the data packet forwarding chain F^k . Note that $F^1 = []$ means that the MN is located at the eNB where a session is initiated. In other words, $F^1 = []$ is an empty forwarding list.

On the other hand, C is described by

$$C = \{C^1, C^2, C^3, \dots, C^N\} \quad (6)$$

where C^i denotes the vector representing the connection between eNB i and other eNBs. Then, C^i is given by

$$C^i = [c_1^i, c_2^i, c_3^i, \dots, c_N^i] \quad (7)$$

where c_j^i represents whether eNB j is connected to eNB i , and c_j^i is defined as

$$c_j^i = \begin{cases} 1, & \text{if eNB } j \text{ is connected to eNB } i \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

For example, if N is 5 and eNB 1 is connected to eNBs 2 and 3, C^1 is given by $[0, 1, 1, 0, 0]$.

C. Actions

When the target eNB receives a handover confirmation message (i.e., $h = 2$), the target eNB decides whether to send a data path switch request message or not. Therefore, when $h = 2$, the target eNB takes an action (i.e., the data path switching or the data packet forwarding) based on the current state information. The action set can be described by $A = \{PF, PS\}$ where PF and PS represent the data packet forwarding and the data path switching, respectively.

D. Transition Probabilities

A state transition of C can be done only when $h = 1$. On the other hand, a state transition of F can occur when $h = 1$ or $h = 2$, and the state transition of F is affected by the chosen action a . That is, the transition probabilities of C and F are affected by H . For a chosen action a , the transition probability from the current state $s = [h, F^k, C^i]$ to the next state, $s' = [h', F^{k'}, C^j]$, can be described by

$$P[s'|s, a] = P[h'|h] \times P[F^{k'}|F^k, a, h] \times P[C^j|C^i, h]. \quad (9)$$

The transition probability of H can be derived as follows. We assume that the residence time in an eNB follows an exponential distribution with mean $1/\lambda_h$. Then, the transition

probability from $h = 0$ to $h = 1$ is given by $\lambda_h \tau$ [15]. Therefore, when $h = 0$, the transition probability from h to h' is defined as

$$P[h'|h = 0] = \begin{cases} 1 - \lambda_h \tau, & \text{if } h' = 0 \\ \lambda_h \tau, & \text{if } h' = 1 \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

Since $h = 1$ means that the MN is in the handover trigger phase, h' should be always 2, i.e., the target eNB receives a handover confirmation message. On the other hand, when $h = 2$, h' should be always 0 since the handover procedure is terminated. This derivation is reasonable because the decision epoch is sufficiently short and thus no consecutive handover events can occur. In short, the transition probabilities for $h = 1$ and $h = 2$ are respectively summarized as

$$P[h'|h = 1] = \begin{cases} 1, & \text{if } h' = 2 \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

and

$$P[h'|h = 2] = \begin{cases} 1, & \text{if } h' = 0 \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

When $h = 0$, F cannot be changed. Therefore, $P[F^{k'}|F^k, a, h = 0]$ is simply given by

$$P[F^{k'}|F^k, a, h = 0] = \begin{cases} 1, & \text{if } F^{k'} = F^k \\ 0, & \text{otherwise.} \end{cases} \quad (13)$$

When $h = 1$ and an MN moves from eNB i to eNB j , the transition probability is determined by F^k . If F^k has an element j (i.e., the current forwarding chain includes eNB j), $F^{k'}$ includes all elements of F^k ahead of j while the element j and the elements behind j are deleted. For example, if $F^k = [2, 5, 1]$ and the MN moves to eNB 5, $F^{k'}$ is given by [2]. On the other hand, if F^k does not have the element j , $F^{k'}$ can be obtained by adding the element i to F^k . For example, if $F^k = [2, 5, 1]$ and the MN moves from eNB 3 to eNB 6 (i.e., $i = 3$ and $j = 6$), $F^{k'}$ is given by [2, 5, 1, 3]. Consequently, the transition probability of F when $h = 1$ can be formulated as (14) and (15) at the top of the next page, where m and m' are the dimensions of F^k and $F^{k'}$, respectively. On the other hand, when $h = 2$, $F^{k'}$ is decided by the chosen action. If the chosen action is PS , $F^{k'}$ is always F^1 , i.e., the forwarding list becomes empty. Meanwhile, if the chosen action is PF , F is not changed. Therefore, $P[F^{k'}|F^k, a, h = 2]$ is given by

$$P[F^{k'}|F^k, a = PS, h = 2] = \begin{cases} 1, & \text{if } F^{k'} = F^1 \\ 0, & \text{otherwise} \end{cases} \quad (16)$$

and

$$P[F^{k'}|F^k, a = PF, h = 2] = \begin{cases} 1, & \text{if } F^{k'} = F^k \\ 0, & \text{otherwise.} \end{cases} \quad (17)$$

Since C is changed only when the MN moves out the current eNB, C is not affected when $h = 0$ or $h = 2$. Therefore, $P[C^j|C^i, h = 0]$ and $P[C^j|C^i, h = 2]$ are respectively given by

$$P[C^j|C^i, h = 0] = \begin{cases} 1, & \text{if } C^j = C^i \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

and

$$P[C^j|C^i, h = 2] = \begin{cases} 1, & \text{if } C^j = C^i \\ 0, & \text{otherwise.} \end{cases} \quad (19)$$

On the other hand, when $h = 1$, the transition probability of C can be derived as follows. When P_{ij} denotes the probability that the MN moves from eNB i to eNB j^3 , the transition probability $P[C^j|C^i, h = 1]$ is obtained from

$$P[C^j|C^i, h = 1] = \begin{cases} P_{ij}, & \text{if } c_j^i = 1 \\ 0, & \text{otherwise.} \end{cases} \quad (20)$$

E. Reward and Cost Functions

To define the reward and cost functions, we consider both the network perspective and the user perspective. In terms of the network perspective, the local traffic overhead and the core traffic overhead can be taken into account to define the reward and cost functions. On the other hand, in terms of the user perspective, the throughput and the delay jitter are considered to represent the reward and cost functions. To sum up, the total reward function, $r(s, a)$, can be defined as

$$r(s, a) = \omega_1 r_n(s, a) + (1 - \omega_1) r_u(s, a) \quad (21)$$

where $r_n(s, a)$ and $r_u(s, a)$ are the reward functions with respect to the network and user perspectives, respectively. ω_1 ($0 \leq \omega_1 \leq 1$) is a weighted factor to adjust the importance of the network and user perspectives.

First, $r_n(s, a)$ can be expressed as

$$r_n(s, a) = \omega_2 f_{LT}(s, a) - (1 - \omega_2) g_{CT}(s, a) \quad (22)$$

where $f_{LT}(s, a)$ and $g_{CT}(s, a)$ are the reward function due to the reduction of the local traffic overhead and the cost function due to the core traffic overhead, respectively. ω_2 ($0 \leq \omega_2 \leq 1$) is a weighted factor to decide the importance between the local traffic overhead and the core traffic overhead. Since the local traffic overhead is proportional to the number of tunnels, m , $f_{LT}(s, a)$ can be defined as

$$f_{LT}(s, a) = \begin{cases} m \times T_{PF}, & \text{if } a = PS, h = 2 \\ 0, & \text{otherwise} \end{cases} \quad (23)$$

where T_{PF} is the unit traffic overhead to transmit a data packet through the data packet forwarding chain. On the other hand, three additional control messages are needed to conduct a data path switching as shown in Figure 2. Hence, $g_{CT}(s, a)$ can be defined as

$$g_{CT}(s, a) = \begin{cases} 3 \times T_{S1}, & \text{if } a = PS, h = 2 \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

where T_{S1} is the unit traffic overhead to conduct the data path switching.

The reward function with respect to the user perspective can be written as

$$r_u(s, a) = \omega_3 f_T(s, a) - (1 - \omega_3) g_D(s, a) \quad (25)$$

where $f_T(s, a)$ and $g_D(s, a)$ are the reward function on the user throughput and the cost function on the delay jitter,

³ P_{ij} can be obtained from the empirical data.

$$\begin{aligned} P[f_1^{k'} = f_1^k, f_2^{k'} = f_2^k, \dots, f_{m+1}^{k'} = i | F^k, a, h = 1] &= 1, \text{ if } c_j^i = 1 \text{ and } j \notin F^k \\ P[F^{k'} | F^k, a, h = 1] &= 0, \text{ otherwise} \end{aligned} \quad (14)$$

and

$$\begin{aligned} P[f_1^{k'} = f_1^k, f_2^{k'} = f_2^k, \dots, f_{m'}^{k'} = f_{m'}^k | F^k, a, h = 1] &= 1, \text{ if } c_j^i = 1 \text{ and } f_{m'+1}^k = j \\ P[F^{k'} | F^k, a, h = 1] &= 0, \text{ otherwise.} \end{aligned} \quad (15)$$

respectively. ω_3 ($0 \leq \omega_3 \leq 1$) is a parameter to weight the user throughput and the delay jitter. Since the degradation of throughput is proportional to the number of tunnels, $f_T(s, a)$ can be defined as

$$f_T(s, a) = \begin{cases} m \times D_T, & \text{if } a = PS, h = 2 \\ 0, & \text{otherwise} \end{cases} \quad (26)$$

where D_T denotes the throughput sensitivity of the session. D_T can be set depending on the application type. For example, for throughput-sensitive applications (e.g., HTTP and FTP), a large value of D_T can be set. Meanwhile, $g_D(s, a)$ is defined as

$$g_D(s, a) = \begin{cases} D_D, & \text{if } a = PS, h = 2 \\ 0, & \text{otherwise} \end{cases} \quad (27)$$

where D_D is the delay sensitivity of the session. For delay-sensitive applications (e.g., VoIP and multimedia streaming), a large value of D_D should be used.

F. Optimality Equation

Let $v(s)$ be the maximum expected total reward when the initial state is s . Then, $v(s)$ can be described as [18]

$$v(s) = \max_{\pi \in \Pi} v^\pi(s) \quad (28)$$

where $v^\pi(s)$ is the expected total reward between the first decision epoch and the last decision epoch when the policy π with an initial state s is given. $v^\pi(s)$ can be obtained from [18]

$$v^\pi(s) = E_s^\pi \left[E_E \left\{ \sum_{t=1}^E r(X_t, Y_t) \right\} \right] \quad (29)$$

where E_s^π represents the expectation with the policy π and initial state s . E_E denotes the expectation with the respect to the random variable E . As mentioned before, since the termination time of a session, E , follows a geometric distribution with mean $1/\lambda_s$. Therefore, $v^\pi(s)$ can be rewritten as [20]

$$v^\pi(s) = E_s^\pi \left\{ \sum_{t=1}^{\infty} (1 - \lambda_s)^{t-1} r(X_t, Y_t) \right\} \quad (30)$$

where $(1 - \lambda_s)$ is a discount factor in the MDP model and $0 \leq (1 - \lambda_s) < 1$.

Finally, the optimality equation is given by [18]

$$v(s) = \max_{a \in A} \left\{ r(s, a) + \sum_{s' \in S} (1 - \lambda_s) P[s' | s, a] v(s') \right\}. \quad (31)$$

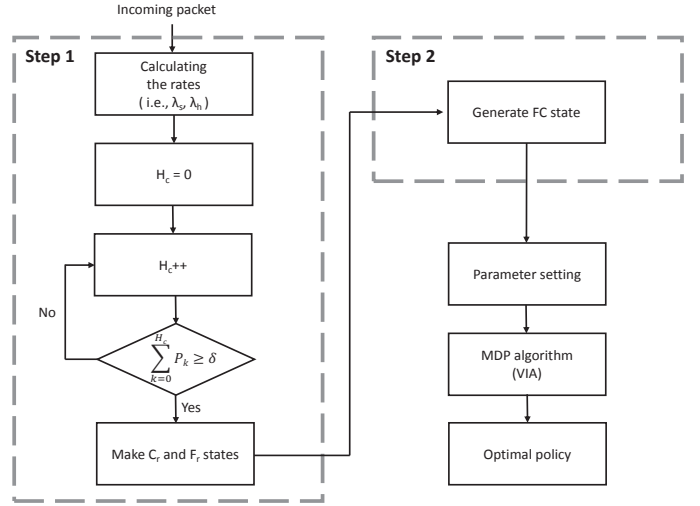


Fig. 5. Low-complexity value iteration algorithm.

G. Low-complexity value iteration algorithm

Generally, a value iteration algorithm (VIA) can be used to solve the optimality equation and to obtain the optimal policy. Each iteration of VIA is performed in $O(|A||S|^2)$ where A is the action set and S is the state space [20], [21]. The state space in our MDP model consists of H , C , and F . The size of C is proportional to the number of eNBs (i.e., N) as shown in (6). Also, the size of F is decided by N as shown in (4). Since there are a lot of eNBs in LTE networks⁴, if all eNBs in LTE networks are considered in VIA, high complexity of VIA is inevitable. This long iteration time is a major obstacle of applying MDP to the practical systems [23]. Due to this reason, we introduce a low-complexity value iteration algorithm as shown in Figure 5. The low-complexity value iteration algorithm reduces C and F state spaces which are main contributors in increasing the complexity of VIA in our MDP model. Reducing the state space is conducted by two steps: 1) making the reduced vector set on the connection information of each eNB and the location of an MN (denoted by C_r) and the reduced vector set of possible forwarding lists (denoted by F_r), which is represented by **Step 1** in Figure 5; 2) generating FC state, which is represented by **Step 2** in Figure 5.

Since the number of eNBs that an MN moves across with an ongoing session is limited, eNBs far from the current eNB can be excluded in the state space to reduce the complexity.

⁴As reported in [22], the number of eNBs managed by SKT, which is the largest cellular network operator in South Korea, is about 10^5 .

In other words, only sufficiently close eNBs to the current eNB are considered to make C_r and F_r states. Let H_c be the expected maximum number of handovers that the MN moves across with an ongoing session. To derive H_c , P_k is defined as the probability that the MN moves across k eNBs during the session duration time. If the residence time and the session duration time follow exponential distributions and their means are estimated by the central system in a statistical manner [24], P_k can be easily obtained by [25]. If the summation of P_k from $k = 0$ to $k = H_c$ is equal to or larger than a pre-defined threshold δ (i.e., $\sum_{k=0}^{H_c} P_k \geq \delta$) and δ is sufficiently large, it is likely that the MN conducts handovers no more than H_c for most sessions. Therefore, H_c can be defined as

$$H_c = \arg \min_x \left(\sum_{k=0}^x P_k \geq \delta \right) \quad (32)$$

where $\arg \min(\cdot)$ returns the minimum argument satisfying the given inequality.

When H_c is given, only eNBs that the MN can move with handovers less than or equal to H_c can be included in the state space. Therefore, the number of handover candidate eNBs, N_c , can be defined by

$$N_c = n \left(\bigcup_{k=0}^{H_c} K_s^k \right) \quad (33)$$

where $n(\mathbb{S})$ returns the number of elements of the set \mathbb{S} and K_s^k denotes the set of eNBs that an MN can move with k handovers from the initial eNB s .

With N_c , C_r can be expressed as

$$C_r = \{C_r^1, C_r^2, C_r^3, \dots, C_r^{N_c}\}. \quad (34)$$

Also, when the number of eNBs is limited by N_c , the total number of forwarding lists, max_c , can be defined as

$$max_c = \begin{cases} 1 + \sum_{M=1}^{N_c-1} \binom{N_c}{M} M!, & \text{if } N_c \neq 1 \\ 1, & \text{if } N_c = 1. \end{cases} \quad (35)$$

Then, F_r , is described as

$$F_r = \{F_r^1, F_r^2, F_r^3, \dots, F_r^{max_c}\}. \quad (36)$$

Note that F_r is constrained by C_r . When the current state of C_r , C_r^i , is decided, the elements of F_r including i cannot be the current state. For example, when the current state is C_r^1 , i.e., the MN is in eNB 1, the elements of F_r such as $F_r^i = [1, 2]$ and $F_r^j = [1]$, where i and j are arbitrary integer numbers, cannot be the current states. By considering this correlation, C_r and F_r can be integrated to one state space FC to further reduce state space (Step 2 in Figure 5). FC denotes the vector set of integrating F_r and C_r , which can be described as

$$FC = \{FC^1, FC^2, FC^3, \dots, FC^{N_{fc}}\} \quad (37)$$

where FC^k denotes the vector set representing the k th integrated state and N_{fc} is the total number of elements of FC .

Finally, the parameters (i.e., ω_1 , ω_2 , ω_3 , D_D , and D_T) are set and VIA in Algorithm 1 where $\|v\| = max\{v(s)\}$

for $s \in S$ is conducted to obtain the near optimal policy ξ . Note that the policy that satisfies the optimality equation in (31) becomes the near optimal policy, ξ , i.e., $\xi(s) = \arg \max_{a \in A} \left\{ r(s, a) + \sum_{s' \in S} (1 - \lambda_s) P[s'|s, a] v(s') \right\}$. Therefore, we can obtain the near optimal policy and the expected total reward from Algorithm 1.

Algorithm 1 Value iteration algorithm.

- 1: Set $v^0(s) = \mathbf{0}$ for each state s . Specify $\epsilon > 0$, and set $k = 0$.
 - 2: For each state s , compute $v^{k+1}(s)$ by
$$v^{k+1}(s) = \max_{a \in A} \left\{ r(s, a) + \sum_{s' \in S} (1 - \lambda_s) P[s'|s, a] v^k(s') \right\}$$
 - 3: If $\|v^{k+1}(s) - v^k(s)\| < \epsilon \lambda_s / (2(1 - \lambda_s))$, go to step 4. Otherwise, increase k by 1 and return to step 2.
 - 4: For each state $s \in S$, compute the stationary optimal policy
$$\xi(s) = \arg \max_{a \in A} \left\{ r(s, a) + \sum_{s' \in S} (1 - \lambda_s) P[s'|s, a] v^{k+1}(s') \right\}$$
 and stop.
-

Note that when the state space is reduced, handover events exceeding H_c cannot be considered and thus the expected total reward of the low-complexity VIA may be different from that of the original VIA. However, such disparity can be sufficiently mitigated by adjusting the threshold δ , which will be elaborated in Section V-A.

V. REAL TRACE-DRIVEN EVALUATION RESULTS

For performance evaluation, we compare the proposed scheme, S_{MDP} , with the following four schemes: 1) S_{PS} where the target eNB always conducts the data path switching, 2) S_{PF} where the data packets are always forwarded through the forwarding chain, 3) $S_{K=2}$ where the target eNB conducts the data path switching only when the forwarding chain length exceeds 2, and 4) $S_{K=3}$ where the target eNB conducts the data path switching only when the forwarding chain length exceeds 3 [15]. Default values of weighted factors (i.e., ω_1 , ω_2 , and ω_3) are set to 0.5. We consider both throughput-sensitive and delay-sensitive applications. For throughput-sensitive applications (i.e., data session), D_T and D_D are set to 3 and 1, respectively. On the other hand, for delay-sensitive applications (i.e., voice session), D_T and D_D are set to 1 and 3, respectively⁵. It is assumed that the average number of hops between the eNB and the S-GW is 10 which is based on 3GPP specification [15]. Based on the Ethernet specification, the data packet size and control packet size are set to 1518 bytes and 64 bytes, respectively. By considering the packet size and the average number of hops, T_{PF} is set to 2.37 (i.e., $1518 / (64 \times 10)$). Also, T_{S1} is normalized as 1.

To derive the handover-to-session ratio $\rho = \lambda_h / \lambda_s$, we use the dataset in [26] that contains 142 days of mobile phone records (i.e., session and cell transition information). Specifically, the handover-to-voice session ratio, ρ_v , is computed

⁵For throughput-sensitive applications and delay-sensitive applications, a large value of D_T and D_D should be used to indicate application characteristic, respectively.

TABLE II
 N_c AND P_h

Session type	δ	H_c	P_h
Voice session	0.7	0	0.7491
	0.9	1	0.9659
	0.95	2	0.9969
Data session	0.7	2	0.7547
	0.9	5	0.9919
	0.95	6	0.9981

as 0.296 whereas the handover-to-data session ratio, ρ_d , is obtained as 1.71. Also, the network topology (e.g., the number of eNBs and the connectivity among eNBs) is configured based on [26]. In the dataset [26], the number of eNBs is 3510. Also, P_{ij} is set by counting the number of handovers between eNB i and eNB j . τ is set to 1 sec. The default value of δ is 0.95 and ε in the value iteration algorithm is set to 0.001. On the other hand, since the results for effect of T_{PF} and T_{S1} are similar as the results of D_T and D_D , respectively, the results for effect of T_{PF} and T_{S1} are not depicted.

A. Effect of δ

The low-complexity VIA is based on the observation that the number of handovers is bound to H_c in most cases. If handovers more than H_c occur, the accuracy of the low-complexity VIA can be affected. Therefore, we need to show that the accuracy of the low-complexity VIA is sufficiently high. To this end, we define the probability P_h that the number of handovers at each decision epoch is less than or equal to H_c . Intuitively, higher P_h represents that the low-complexity VIA has higher accuracy. P_h can be calculated as follows. Since the average session duration is $1/\lambda_s$, the average number of decision epochs during the session duration is given by $\lfloor 1/(\lambda_s\tau) \rfloor$ where $\lfloor \cdot \rfloor$ is the floor function. On the other hand, the occurrence probability of handover at each decision epoch is $\lambda_h\tau$ by (10). Consequently, P_h can be computed as

$$P_h = \sum_{i=0}^{H_c} \binom{\lfloor 1/(\lambda_s\tau) \rfloor}{i} (\lambda_h\tau)^i (1 - \lambda_h\tau)^{\lfloor 1/(\lambda_s\tau) \rfloor - i}. \quad (38)$$

Table II shows P_h under different values of δ and H_c . As shown in Table II, when δ is set to a sufficiently high value (e.g., 0.9 or 0.95), P_h is adequately large. In other words, handovers more than H_c rarely occur when δ is larger than 0.9. Meanwhile, when δ is 0.7, P_h is relatively small, which indicates that the low-complexity VIA has unfavorable accuracy. Therefore, the default value of δ is set to 0.95 in the following results.

B. Effect of $(1 - \lambda_s)$

Figure 6(a) shows the expected total reward as a function of the discount factor, $(1 - \lambda_s)$, in delay-sensitive applications. Note that a larger discount factor refers to longer simulation time [20]. Therefore, the expected total rewards of S_{PS} and $S_{K=2}$ decrease as the discount factor increases due to the increased number of handovers and data path switching

operations⁶. Note that when the delay-sensitive application is considered, the negative effect by the data path switching (e.g., delay jitter) is more significant. On the other hand, S_{MDP} and S_{PF} are rarely affected by the increased number of handovers since data packets are forwarded over a locally established forwarding chain. In addition, S_{MDP} follows almost the same actions as S_{PF} since the data packet forwarding is beneficial in increasing the total reward. On the other hand, since H_c is set to 2 for the voice session (i.e., delay-sensitive application), $S_{K=3}$ does not conduct any PS actions and thus $S_{K=3}$ has the same expected total reward as S_{PF} in delay-sensitive applications.

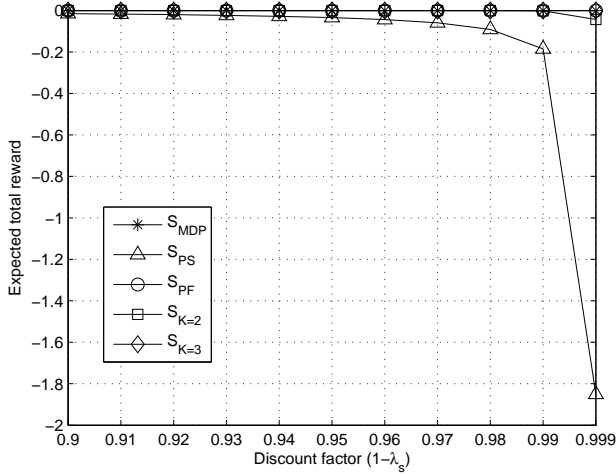
Figure 6(b) shows the expected total reward as a function of the discount factor, $(1 - \lambda_s)$, for throughput-sensitive applications. It can be shown that the expected total rewards of all schemes except S_{PF} increase as the discount factor increases. This can be explained as follows. When the simulation time is long, more PS actions can be conducted. Since more PS actions can further increase the throughput by optimizing the data packet delivery path, longer simulation time provides more expected total rewards. Meanwhile, S_{PF} does not conduct any PS actions and therefore S_{PF} is independent of the increased simulation time and does not have any reward or cost (i.e., the expected total reward of S_{PF} is always 0). Also, it can be seen that the expected total reward of S_{MDP} is comparable to that of S_{PS} when $(1 - \lambda_s)$ is 0.9 ~ 0.95. On the other hand, when $(1 - \lambda_s)$ is 0.96 ~ 0.99, the expected total reward of S_{MDP} is almost the same as that of $S_{K=2}$. This can be explained as follows. When the average data session duration is short (i.e., $(1 - \lambda_s)$ is 0.9 ~ 0.95), only few handover events occur during the session time and thus the target eNB always conducts the data path switching to obtain higher throughput in throughput-sensitive applications. In other words, since handover events rarely occur, the signaling overhead to core entities due to the data path switching is not quite high even when the target eNB always conducts the data path switching. Meanwhile, when the average data session duration is long (i.e., $(1 - \lambda_s)$ is 0.96 ~ 0.99), more handover events occur. Therefore, the signaling overhead to core entities due to the data path switching can be significant if the target eNB always conducts the data path switching. Thus, S_{MDP} conducts the data path switching when the forwarding chain length exceeds 2. Interestingly, when more frequent handover events are considered, i.e., when $(1 - \lambda_s)$ is 0.99 ~ 0.999, it can be found that S_{MDP} outperforms all other schemes apparently. This is because S_{MDP} always conducts appropriate and adaptive actions regardless of the number of handover events whereas other schemes follow the static actions.

C. Effect of ω_2

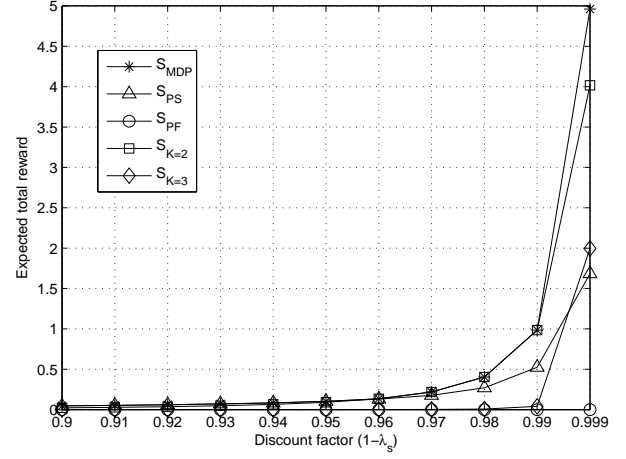
From Figure 7, it can be seen that S_{MDP} operates adaptively even when the weighted factor, ω_2 , is changed⁷. When ω_2 is small (i.e., 0 and 0.2), the expected total reward of S_{MDP}

⁶Due to scale of y-axis, the decrement of the expected total reward of $S_{K=2}$ is not shown well in Figure 6(a).

⁷Since there are no special tendencies as ω_1 and ω_3 are changed, we do not include these results.



(a) Delay-sensitive applications



(b) Throughput-sensitive applications

Fig. 6. Expected total reward vs. $(1 - \lambda_s)$ ($\lambda_h = 0.00296$ in (a) / $\lambda_h = 0.0171$ in (b) [26]).

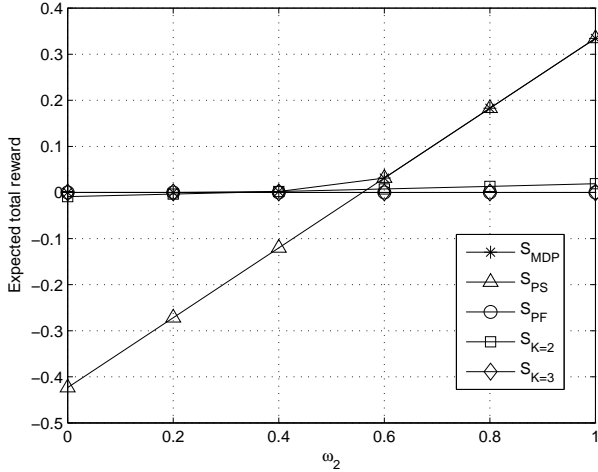


Fig. 7. Expected total reward vs. ω_2 ($D_T = D_D = 1$).

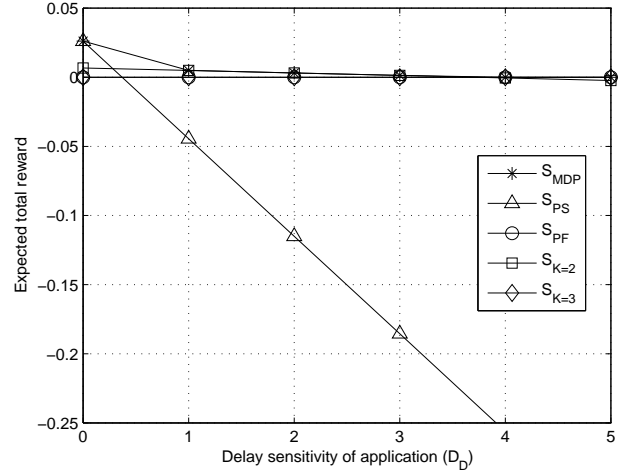


Fig. 8. Expected total reward vs. D_D ($D_T = 1$).

is the same as that of S_{PF} . This is because small ω_2 represents that the core traffic overhead affects more significantly to the expected total reward than the local traffic overhead does. In such a situation, the reward $f_{LT}(s, a)$ is relatively smaller than the cost $f_{CT}(s, a)$ if the PS action is chosen. Therefore, S_{MDP} always chooses the PF action for small ω_2 . On the other hand, when ω_2 is 0.4, the data packet forwarding and the data path switching provide more total rewards when the forwarding chain lengths are 1 and 2, respectively. Therefore, S_{MDP} chooses the PF and PS actions when the forwarding chain lengths are 1 and 2, respectively. Consequently, the expected total reward of S_{MDP} is almost the same as that of $S_{K=2}$. Meanwhile, when ω_2 has a large value (e.g., 0.6, 0.8, and 1), the PS action always gives more rewards and less costs, and thus S_{MDP} conducts the data path switching in most cases.

D. Effect of D_D

The effect of D_D is shown in Figure 8. As D_D increases (i.e., the application becomes more sensitive to the delay), the cost for the chosen PS action also increases. Therefore, the expected total rewards of S_{MDP} , S_{PS} and $S_{K=2}$ decrease as D_D increases. However, since S_{MDP} does not choose the PS action when D_D is large, the expected total reward of S_{MDP} is larger than that of S_{PF} . On the other hand, since S_{PS} and $S_{K=2}$ do not take the delay sensitivity of the application into consideration, the expected total rewards of S_{PS} and $S_{K=2}$ decrease continuously regardless of the delay sensitivity of the application. In addition, since $S_{K=2}$ chooses less PS actions than S_{PS} , the reduction ratio of the expected total reward of $S_{K=2}$ is smaller than that of S_{PS} .

E. Effect of D_T

From Figure 9, it can be shown that S_{MDP} operates adaptively as D_T is changed. Since increasing D_T means that

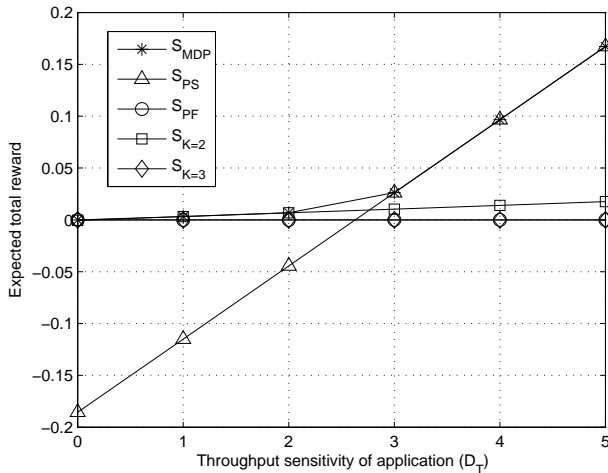


Fig. 9. Expected total reward vs. D_T ($D_D = 2$).

the application is more sensitive to the throughput, the reward increases when the PS action is chosen thanks to the optimal data forwarding path. Therefore, S_{MDP} chooses the PS action when D_T is large. On the other hand, $D_T = 0$ means that the throughput is not considered at all and thus S_{MDP} chooses the PF action when $D_T = 0$. When the D_T is set to 1 or 2 (i.e., the sensitivities to the delay and throughput are comparable), S_{MDP} has the same policy as $S_{K=2}$. This is because the PF action and the PS action give more rewards (or less costs) when the forwarding chain lengths are 1 and 2, respectively. Meanwhile, when D_T is large (i.e., 3, 4, and 5), S_{MDP} always chooses the PS action for higher throughput.

VI. CONCLUSION

In this paper, we proposed an optimized and distributed data forwarding scheme. To find out the optimal policy between the data path switching and data packet forwarding, we formulated MDP in which the reward and cost functions in throughput/delay-sensitive applications are defined with the respect to the network/user's QoE. Also, a low-complexity value iteration algorithm is devised to solve the optimality and therefore we believe the proposed scheme can be applied to practical systems. Real trace-driven evaluation results demonstrate that the proposed scheme can increase the expected total reward compared with other schemes and can achieve the adaptive performance optimization. In addition, we showed that the proposed scheme can adaptively set the optimal forwarding chain depending on the application types (i.e., throughput/delay-sensitive applications). In our future work, we will extend the proposed scheme to reflect the local and core network traffic loads and evaluate the performance of the extended scheme in terms of mobile traffic offloading. Also, we will investigate how to extend the proposed forwarding scheme under high vehicular mobility.

REFERENCES

[1] H. Ko, G. Lee, and S. Pack, "Optimized and Distributed Data Packet Forwarding in LTE/LTE-A Networks," in *Proc. IEEE International Conference on Communications (ICC) 2014*, pp. 2502-2507, Jun. 2014.

[2] Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2012-2017, [Online]. Available: http://www.cisco.com/en/US/solutions/collateral/ns341/ns525/ns537/ns705/ns827/white_paper_c11-520862.pdf.

[3] Ericsson, "Traffic and Market Report," Jun. 2012.

[4] G. Iosifidis, L. Gao, J. Huang, and L. Tassiulas, "A Double-Auction Mechanism for Mobile Data-Offloading Markets," *IEEE/ACM Transactions on Networking*, to appear.

[5] H. Fu, P. Lin, and Y. Lin, "Reducing Signaling Overhead for Femtocell/Macrocell Networks," *IEEE Transactions on Mobile Computing*, vol. 12, no. 8, pp. 1587-1597, Aug. 2013.

[6] K. Zheng, Y. Wang, C. Lin, X. Shen, and J. Wang, "A Graph-based Interference Coordination Scheme in OFDMA Femtocell Networks," *IET Communications*, vol. 5, no. 17, pp. 2533-2541, Nov. 2011.

[7] Y. Liu, L.X. Cai, X. Shen, and H. Luo, "Deploying Cognitive Cellular Networks under Dynamic Resource Management," *IEEE Wireless Communications*, vol. 20, no. 2, pp. 82-88, Apr. 2013.

[8] S. Sesia, I. Toufik, and M. Baker, *LTE - The UMTS Long Term Evolution: From Theory to Practice*. John Wiley and Sons, 2009.

[9] 3GPP TS 36.300, "Evolved Universal Terrestrial Radio Access (E-UTRA) and Evolved Universal Terrestrial Radio Access Network (E-UTRAN); Overall Description; Stage 2," version 12.0, Dec. 2013.

[10] A. Rath and S. Panwar, "Fast Handover in Cellular Networks with Femtocells," in *Proc. IEEE International Conference on Communications (ICC) 2012*, pp. 2752-2757, Jun. 2012.

[11] D. Kim, H. Lee, N. Kim, and H. Yoon, "A Velocity-based Bicast Handover Scheme for 4G Mobile System," in *Proc. International Wireless Communications and Mobile Computing Conference (IWCMC) 2008*, pp. 147-152, Aug. 2008.

[12] T. Guo, A. Quddus, and R. Tafazolli, "Seamless Handover for LTE Macro-Femto Networks Based on Reative Data Bicast," *IEEE Communication Letters*, vol. 16, no. 11, pp. 1788-1791, Nov. 2012.

[13] D. Pacifico, M. Pacifico, and C. Fischione, "Improving TCP Performance during the Intra LTE Handover," in *Proc. IEEE Global Telecommunications Conference (GLOBECOM) 2009*, pp. 1-8, Nov/Dec. 2009.

[14] Z. Yan and J. Lee, "State-Aware Pointer Forwarding Scheme with Fast Handover Support in a PMIPv6 Domain," *IEEE Systems Journal*, vol. 7, no. 1, pp. 92-101, Mar. 2013.

[15] T. Guo, N. Wang, and R. Tafazolli, "Local Mobility Management for Networked Femtocells Based on X2 Traffic forwarding," *IEEE Transactions on Vehicular Technology*, vol. 62, no. 1, pp. 326-340, Jan. 2013.

[16] J. Lee, Z. Yan, J. Bonnin, and X. Lagrange, "Dynamic Tunneling for Network-based Distributed Mobility Management Coexisting with PMIPv6," in *Proc. IEEE International Symposium on Personal, Indoor, and Mobile Radio Communication (PIMRC) 2013*, pp. 2995-3000, Sep. 2013.

[17] X. Xu, Y. Cui, J. Liu, W. Wang, Y. Tan, G. Liu, and C. Zhu, "Seamless Handover Scheme for Real-Time Multimedia Services in PMIPv6," in *Proc. IEEE International Conference on Communications and Networking in China (CHINACOM) 2013*, pp. 834-839, Aug. 2013.

[18] M. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley and Sons, 2009.

[19] E.A. Feinberg and A. Shwartz, *Handbook of Markov Decision Processes: Methods and Applications*. Kluwer Academic Publishers, 2002.

[20] E. Stevens, Y. Lin, and V. Wong, "An MDP-based Vertical Hand-off Decision Algorithm for Heterogeneous Wireless Networks," *IEEE Transactions on Vehicular Technology*, vol. 57, no. 2, pp. 1243-1254, Mar. 2008.

[21] M. Littman, T. Dean, and L. Kaelbling, "On the Complexity of Solving Markov Decision Problems," in *Proc. International Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 394-402, Aug. 1995.

[22] OpenCellID, [Online]. Available: <http://www.opencellid.org/>

[23] T. Sun, Q. Zhao, and P. Luh, "Incremental Value Iteration for Time-Aggregated Markov-Decision Process," *IEEE Transactions on Automatic Control*, vol. 52, no. 11, pp. 2177-2182, Nov. 2007.

[24] S. Pack and Y. Choi, "Fast Handoff Scheme Based on Mobility Prediction in Public Wireless LAN Systems," *IEE Proceedings Communications*, vol. 151, no. 5, pp. 489-495, Oct. 2004.

[25] Y. Lin, "Reducing Location Update Cost in a PCS Network," *IEEE/ACM Transactions on Networking*, vol. 5, no. 1, pp. 25-33, Feb. 1997.

[26] M. Ficek and L. Kencl, "Inter-Call Mobility Model: A Spatio-Temporal Refinement of Call Data Records Using a Gaussian Mixture Model," in *Proc. IEEE International Conference on Computer Communications (INFOCOM) 2012*, pp. 469-477, Mar. 2012.



Haneul Ko received the B.S. degree from Korea University, Seoul, Korea, in 2011. He is currently an M.S. and Ph.D. integrated course student in School of Electrical Engineering, Korea University, Seoul, Korea. His research interests include 5G network architecture, mobility management, and Future Internet.



Xuemin (Sherman) Shen received the B.Sc.(1982) degree from Dalian Maritime University (China) and the M.Sc. (1987) and Ph.D. degrees (1990) from Rutgers University, New Jersey (USA), all in electrical engineering. He is a Professor and University Research Chair, Department of Electrical and Computer Engineering, University of Waterloo, Canada. Dr. Shens research focuses on resource management in interconnected wireless/wired networks, wireless network security, wireless body area networks, vehicular ad hoc and sensor networks. Dr.

Shen serves/served as the Editor-in-Chief for IEEE Network, Peer-to-Peer Networking and Application, and IET Communications; a Founding Area Editor for IEEE Transactions on Wireless Communications; an Associate Editor for IEEE Transactions on Vehicular Technology, Computer Networks, and ACM/Wireless Networks, etc.; and the Guest Editor for IEEE JSAC, IEEE Wireless Communications, IEEE Communications Magazine, and ACM Mobile Networks and Applications, etc. Dr. Shen is a registered Professional Engineer of Ontario, Canada, an IEEE Fellow, an Engineering Institute of Canada Fellow, a Canadian Academy of Engineering Fellow, and a Distinguished Lecturer of IEEE Vehicular Technology Society and Communications Society.



Giwon Lee received the B.S. and M.S. degrees from Korea University, Seoul, Korea, in 2009 and 2011, respectively. He is currently an Ph.D. course student in School of Electrical Engineering, Korea University, Seoul, Korea. His research interests include mobile cloud computing, software defined networking, vehicular networks, and Future Internet.



Dongeun Suh received the B.S. degrees from Korea University, Seoul, Korea, in 2012. He is currently an Ph.D. course student in School of Electrical Engineering, Korea University, Seoul, Korea. His research interests include adaptive multimedia streaming, mobile data offloading, and software-defined networking.



Sangheon Pack received the B.S. and Ph.D. degrees from Seoul National University, Seoul, Korea, in 2000 and 2005, respectively, both in computer engineering. In 2007, he joined the faculty of Korea University, Seoul, Korea, where he is currently an Associate Professor in the School of Electrical Engineering. From 2005 to 2006, he was a Postdoctoral Fellow with the Broadband Communications Research Group, University of Waterloo, Waterloo, ON, Canada. He was the recipient of KICS (Korean Institute of Communications and Information Sciences)

Haedong Young Scholar Award 2013, IEEE ComSoc APB Outstanding Young Researcher Award in 2009, LG Yonam Foundation Overseas Research Professor Program in 2012, and Student Travel Grant Award at the IFIP Personal Wireless Conference (PWC) 2003. From 2002 to 2005, he was a recipient of the Korea Foundation for Advanced Studies Computer Science and Information Technology Scholarship. He was a publication co-chair of IEEE INFOCOM 2014, a co-chair of IEEE VTC 2010-Fall transportation track, a co-chair of IEEE WCSP 2013 wireless networking symposium, a TPC vice-chair of ICOIN 2013, and a publicity co-chair of IEEE SECON 2012. He is an editor of Journal of Communications Networks (JCN) and a senior member of the IEEE. His research interests include Future Internet, SDN/ICN/DTN, mobility management, mobile cloud networking, multimedia networking, and vehicular networks.