

Caching Strategies Based on Information Density Estimation in Wireless Ad Hoc Networks

Marco Fiore, *Member, IEEE*, Claudio Casetti, *Member, IEEE*, and Carla-Fabiana Chiasserini, *Senior Member, IEEE*

Abstract—We address cooperative caching in wireless networks, where the nodes may be mobile and exchange information in a peer-to-peer fashion. We consider both cases of nodes with large- and small-sized caches. For large-sized caches, we devise a strategy where nodes, independent of each other, decide whether to cache some content and for how long. In the case of small-sized caches, we aim to design a content replacement strategy that allows nodes to successfully store newly received information while maintaining the good performance of the content distribution system. Under both conditions, each node takes decisions according to its perception of what nearby users may store in their caches and with the aim of differentiating its own cache content from the other nodes'. The result is the creation of content diversity within the nodes neighborhood so that a requesting user likely finds the desired information nearby. We simulate our caching algorithms in different ad hoc network scenarios and compare them with other caching schemes, showing that our solution succeeds in creating the desired content diversity, thus leading to a resource-efficient information access.

Index Terms—Data caching, mobile ad hoc networks.

I. INTRODUCTION

PROVIDING information to users on the move is one of the most promising directions of the infotainment business, which rapidly becomes a market reality, because infotainment modules are deployed on cars and handheld devices. The ubiquity and ease of access of third- and fourth-generation (3G or 4G) networks will encourage users to constantly look for content that matches their interests. However, by exclusively relying on downloading from the infrastructure, novel applications such as mobile multimedia are likely to overload the wireless network (as recently happened to AT&T following the introduction of the iPhone [1]). It is thus conceivable that a peer-to-peer system could come in handy, if used in conjunction with cellular networks, to promote content sharing using ad hoc networking among mobile users [2]. For highly popular content, peer-to-peer distribution can, indeed, remove

bottlenecks by pushing the distribution from the core to the edge of the network.

In such an environment, however, a cache-all-you-see approach is unfeasible, because it would swamp node storage capacity with needless data that were picked up on the go. Thus, several techniques of efficiently caching information in wireless ad hoc networks have been investigated in the literature; for example, see the surveys in [3] and [4] and the related work discussed in Section II.

The solution that we propose, called Hamlet, aims at creating *content diversity* within the node neighborhood so that users likely find a copy of the different information items nearby (regardless of the content popularity level) and avoid flooding the network with query messages. Although a similar concept has been put forward in [5]–[8], the novelty in our proposal resides in the probabilistic estimate, run by each node, of the *information presence* (i.e., of the cached content) in the node proximity. The estimate is performed in a cross-layer fashion by overhearing content query and information reply messages due to the broadcast nature of the wireless channel. By leveraging such a local estimate, nodes autonomously decide what information to keep and for how long, resulting in a distributed scheme that does not require additional control messages. The Hamlet approach applies to the following cases.

- *Large-sized caches.* In this case, nodes can potentially store a large portion (i.e., up to 50%) of the available information items. Reduced memory usage is a desirable (if not required) condition, because the same memory may be shared by different services and applications that run at nodes. In such a scenario, a caching decision consists of computing for how long a given content should be stored by a node that has previously requested it, with the goal of minimizing the memory usage without affecting the overall information retrieval performance;
- *Small-sized caches.* In this case, nodes have a dedicated but limited amount of memory where to store a small percentage (i.e., up to 10%) of the data that they retrieve. The caching decision translates into a cache replacement strategy that selects the information items to be dropped among the information items just received and the information items that already fill up the dedicated memory.

We evaluate the performance of Hamlet in different mobile network scenarios, where nodes communicate through ad hoc connectivity. The results show that our solution ensures a high query resolution ratio while maintaining the traffic load very low, even for scarcely popular content, and consistently along different network connectivity and mobility scenarios.

Manuscript received June 14, 2010; revised October 12, 2010 and January 19, 2011; accepted March 18, 2011. Date of publication April 5, 2011; date of current version June 20, 2011. This work was supported in part by the Regione Piemonte through the Monitoraggio Aree Sensibili e Protette Project. The review of this paper was coordinated by Prof. J. Misić.

M. Fiore is with the Institut National de Recherche en Informatique et en Automatique (INRIA)—Centre of Innovation in Telecommunications and Integration of Services (CITI) Laboratory, Institut National des Sciences Appliquées (INSA) Lyon, 69621 Villeurbanne Cedex, France (e-mail: marco.fiore@insa-lyon.fr).

C. Casetti and C.-F. Chiasserini are with the Dipartimento di Elettronica e Telecomunicazioni, Politecnico di Torino, 10129 Torino, Italy (e-mail: claudio.casetti@polito.it; carla.chiasserini@polito.it).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TVT.2011.2136363

The remainder of this paper is organized as follows. First, we discuss the related literature in Section II and outline the system characteristics and assumptions in Section III. The Hamlet framework is then described in Section IV. The simulation scenarios considered for the performance evaluation of Hamlet are detailed in Section V, whereas the results are presented in Sections VI and VII for large- and small-sized caches, respectively. Finally, Section VIII draws some conclusions and highlights future research directions.

II. RELATED WORK

Several papers have addressed content caching and content replacement in wireless networks. In the following sections, we review the works that are most related to this paper, highlighting the differences with respect to the Hamlet framework that we propose.

A. Cooperative Caching

In [9], distributed caching strategies for ad hoc networks are presented according to which nodes may cache highly popular content that passes by or record the data path and use it to redirect future requests. Among the schemes presented in [9], the approach called HybridCache best matches the operation and system assumptions that we consider; we thus employ it as a benchmark for Hamlet in our comparative evaluation. In [10], a cooperative caching technique is presented and shown to provide better performance than HybridCache. However, the solution that was proposed is based on the formation of an overlay network composed of “mediator” nodes, and it is only fitted to static connected networks with stable links among nodes. These assumptions, along with the significant communication overhead needed to elect “mediator” nodes, make this scheme unsuitable for the mobile environments that we address. The work in [11] proposes a complete framework for information retrieval and caching in mobile ad hoc networks, and it is built on an underlying routing protocol and requires the manual setting of a networkwide “cooperation zone” parameter. Note that assuming the presence of a routing protocol can prevent the adoption of the scheme in [11] in highly mobile networks, where maintaining network connectivity is either impossible or more communication expensive than the querying/caching process. Furthermore, the need of a manual calibration of the “cooperation zone” makes the scheme hard to configure, because different environments are considered. Conversely, Hamlet is self contained and is designed to self adapt to network environments with different mobility and connectivity features.

One vehicular ad hoc network scenario is addressed in [12], where the authors propose both an information retrieval technique that aims at finding the most popular and relevant data matching a user query and a popularity-aware data replacement scheme. The latter approach ensures that the density of different content is proportional to the content’s popularity at the system steady state, thus obeying the square-root rule proposed in [13] for wired networks. We point out that the square-root rule does not consider where copies of the data are located but only how many copies are created. It is thus insufficient in

network environments whose dynamism makes the positioning of content of fundamental importance and renders steady-state conditions (as assumed in [13]) hard to be achieved.

B. Content Diversity

Similar to Hamlet, in [6], mobile nodes cache data items other than their neighbors to improve data accessibility. In particular, the solution in [6] aims at caching copies of the same content farther than a given number of hops. Such a scheme, however, requires the maintenance of a consistent state among nodes and is unsuitable for mobile network topologies. The concept of caching different content within a neighborhood is also exploited in [7], where nodes with similar interests and mobility patterns are grouped together to improve the cache hit rate, and in [8], where neighboring mobile nodes implement a cooperative cache replacement strategy. In both works, the caching management is based on instantaneous feedback from the neighboring nodes, which requires additional messages. The estimation of the content presence that we propose, instead, avoids such communication overhead.

C. Caching With Limited Storage Capability

In the presence of small-sized caches, a cache replacement technique needs to be implemented. Aside from the scheme in [8], centralized and distributed solutions to the cache placement problem, which aim at minimizing data access costs when network nodes have limited storage capacity, are presented in [14]. Although centralized solutions are not feasible in ad hoc environments, the distributed scheme in [14] makes use of cache tables, which, in mobile networks, need to be maintained similar to routing tables. Hamlet does not rely on cache tables, and thus, it does not incur the associate high communication penalty. In [15], a content replacement strategy that aims at minimizing energy consumption is proposed. To determine which content should be discarded, the solution exploits the knowledge of data access probabilities and distance from the closest provider—an information that is typically hard to obtain and is not required by Hamlet.

A content replacement scheme that addresses storage limitations is also proposed in [16]. It employs a variant of the last recently used (LRU) technique, which favors the storage of the most popular items instead of the uniform content distribution targeted by Hamlet. In addition, it exploits the cached item IDs provided by the middleware to decide on whether to reply to passing-by queries at the network layer, as well as link-layer traffic monitoring to trigger prefetching and caching. In [17], the popularity of content is taken into account, along with its update rate, so that items that are more frequently updated are more likely to be discarded. Similarly, in [18], cache replacement is driven by several factors, including access probability, update frequency, and retrieval delay. These solutions thus jointly address cache replacement and consistency, whereas in this paper, we specifically target the former issue. However, as will be pointed out, Hamlet can easily be coupled with a dedicated cache consistency scheme, e.g., [19] and [20].

D. Data Replication

Although addressing a different problem, some approaches to data replication are relevant to the data caching solution that we propose. One technique of eliminating information replicas among neighboring nodes is introduced in [21], which, unlike Hamlet, requires knowledge of the information access frequency and periodic transmission of control messages to coordinate the nodes' caching decisions. In [5], the authors propose a replication scheme that aims at having every node close to a copy of the information and analyze its convergence time. However, unlike Hamlet, the scheme implies a significant overhead and an exceedingly high convergence time, thus making it unsuitable for highly variable networks. Finally, the work in [22] adopts a cross-layer approach to data replication in mobile ad hoc networks, where network-layer information on the node movement path helps to trigger the replication before network partitioning occurs.

III. SYSTEM OUTLINE AND ASSUMPTIONS

Hamlet is a fully distributed caching strategy for wireless ad hoc networks whose nodes exchange information items in a peer-to-peer fashion. In particular, we address a mobile ad hoc network whose nodes may be resource-constrained devices, pedestrian users, or vehicles on city roads. Each node runs an application to request and, possibly, cache desired information items. Nodes in the network retrieve information items from other users that temporarily cache (part of) the requested items or from one or more gateway nodes, which can store content or quickly fetch it from the Internet.

We assume a content distribution system where the following assumptions hold: 1) A number I of *information items* is available to the users, with each item divided into a number C of *chunks*; 2) user nodes can overhear queries for content and relative responses within their radio proximity by exploiting the broadcast nature of the wireless medium; and 3) user nodes can estimate their distance in hops from the query source and the responding node due to a hop-count field in the messages. Although Hamlet can work with *any* system that satisfies the aforementioned three generic assumptions, for concreteness, we detail the features of the specific content retrieval system that we will consider in the remainder of this paper.

The reference system that we assume allows user applications to request an information item i ($1 \leq i \leq I$) that is not in their cache. Upon a request generation, the node broadcasts a *query message* for the C chunks of the information item. Queries for still missing chunks are periodically issued until either the information item is fully retrieved or a timeout expires.

If a node receives a fresh query that contains a request for information i 's chunks and it caches a copy of one or more of the requested chunks, it sends them back to the requesting node through *information messages*. If the node does not cache (all of) the requested chunks, it can rebroadcast a query for the missing chunks, thus acting as a forwarder. The exact algorithm that is followed by a node upon the reception of a query message is detailed in the flowchart in Fig. 1(a).

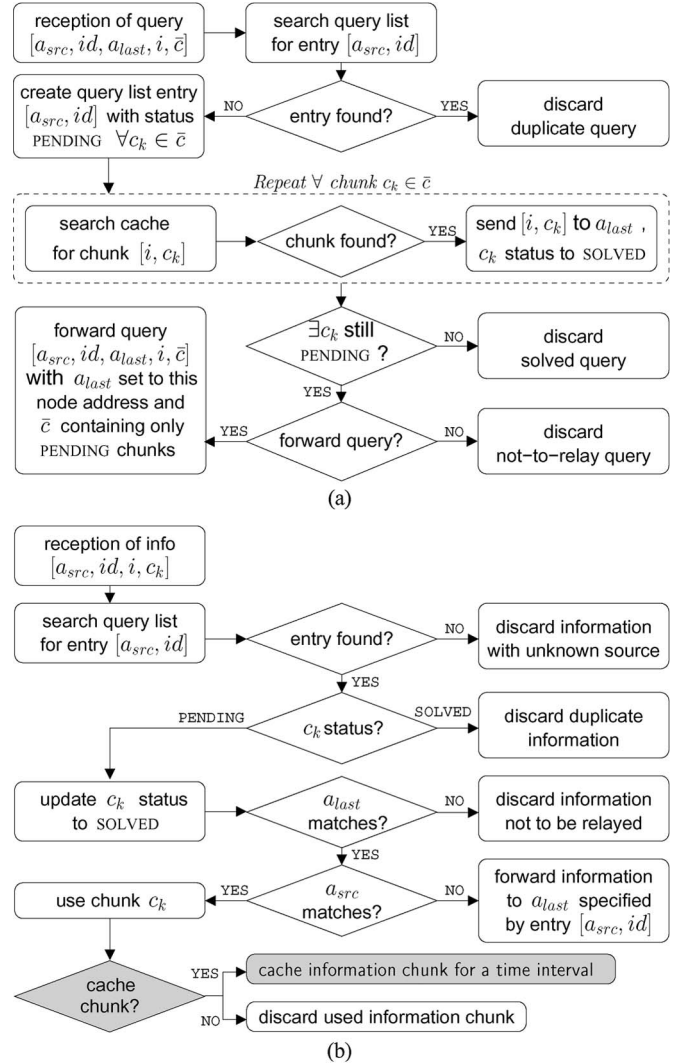


Fig. 1. Flowcharts of the processing of (a) query and (b) information messages at user nodes. We denote the address of the node that generated the query as a_{src} , the query identifier as id , the address of the last node that forwarded the query message as a_{last} , and the set of queried chunks as \bar{c} . The functional blocks that are the focus of this paper are highlighted in (b).

Once created, an information message is sent back to the query source. To avoid a proliferation of information copies along the path, the only node that is entitled to cache a new copy of the information is the node that issued the query. Information messages are transmitted back to the source of the request in a unicast fashion, along the same path from which the request came. To this end, backtracking information is carried and updated in query messages. Nodes along the way either act as relays for transit messages (if they belong to the backtracking node sequence) or simply overhear their transmission without relaying them. Fig. 1(b) depicts the flowchart of the operations at a node that receives a message that contains an information chunk.

A node that receives the requested information has the option to cache the received content and thus become a provider for that content to the other nodes. Determining a strategy of taking such caching decisions is the main objective of this paper, and as such, the corresponding decision blocks are highlighted in Fig. 1(b).

We point out that Hamlet exploits the observation of query and information messages that are sent on the wireless channel as part of the operations of the content-sharing application, e.g., the previously outlined approach. As a consequence, Hamlet does not introduce any signaling overhead.

Furthermore, several optimizations can be introduced to improve the aforementioned basic scheme for the discovery of content. Although our focus is not on query propagation, it is important to take the query process into account, because it directly determines the network load associated with the content retrieval operation. While deriving the results, we consider the following two approaches to query propagation.

- 1) *Mitigated flooding*. This approach limits the propagation range of a request by forcing a time to live (*TTL*) for the query messages. In addition, it avoids the forwarding of already-solved requests by making the nodes wait for a *query lag time* before rebroadcasting a query;
- 2) *Eureka* [23]. This approach extends mitigated flooding by steering queries toward areas of the network where the required information is estimated to be denser.

Note that this paper focuses on cooperative caching and we do not tackle information consistency; thus, we do not take into account different versions of the content in the system model. We note, however, that the previous version of this paper [24] jointly evaluated Hamlet with a basic scheme for weak cache consistency based on an epidemic diffusion of an updated cache content and we showed that weak consistency can be reached, even with such a simple approach, with latencies on the order of minutes for large networks. If prompt solutions are sought, Hamlet lends itself to be easily integrated with one of the existing consistency solutions found in the literature (e.g., [19], [20], [25], and [26]). In particular, these works propose push, pull, or hybrid approaches to achieve different levels of cache consistency. In the case of Hamlet, a *push* technique can be implemented through the addition of invalidation messages broadcast by gateway nodes, whereas information providers can *pull* an updated content (or verify its freshness) before sending the information to querying nodes. In either case, no major modification of the Hamlet caching scheme is required: the only tweaking can consist of resetting the estimation of the information presence upon the notification/detection of an updated version to ease the diffusion of the new information.

IV. HAMLET FRAMEWORK

The Hamlet framework allows wireless users to take *caching decisions* on content that they have retrieved from the network. The process that we devise allows users to take such decisions by leveraging a node's local observation, i.e., the node's ability to overhear queries and information messages on the wireless channel. In particular, for each information item, a node records the distance (in hops) of the node that issues the query, i.e., where a copy of the content is likely to be stored, and the distance of the node that provides the information. Based on such observations, the node computes an index of the *information presence* in its proximity for each of the I items. Then, as the node retrieves content that it requested, it uses the presence index of such an information item to determine

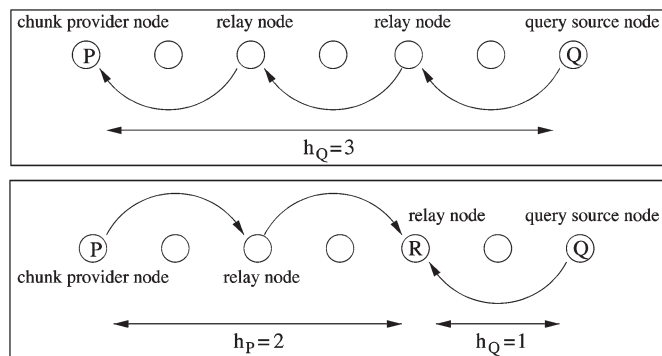


Fig. 2. Q and P denote, respectively, a node that issues a query and a node that provides the requested content. Node R in the lower plot is a relay node, overhearing the exchanged messages. The upper and lower plots, respectively, represent the case 1 h_Q value for the provider node P and the case 2 h_Q and h_P values for the relay node R with respect to the query source Q and the provider P .

whether a copy of the content should be cached, for how long, and possibly which content it should replace. By doing so, a node takes caching decisions that favor high content diversity in its surroundings, inherently easing the retrieval of data in the network. Note that our technique works on a per-item basis, and its results apply to all chunks that belong to the same content.

In the following sections, we first detail how a node estimates the presence of information chunks in its proximity. Next, we separately describe the role of the information presence index in caching decisions for nodes with large- and small-sized caches. In the former case, the information presence index determines the cache content drop time, whereas in the latter case, it drives the cache content replacement.

A. Information Presence Estimation

We define the *reach range* of a generic node n as its distance from the farthest node that can receive a query generated by node n itself. As an example, in an ideal setting in which all nodes have the same radio range, the reach range is given by the product of the *TTL* and the node radio range. Next, we denote by f the frequency at which every node estimates the presence of each information item within its *reach range*, and we define as $1/f$ the duration of each estimation step (also called time step hereafter).

A node n uses the information that was captured within its reach range during time step j to compute the following two quantities: 1) a provider counter by using application-layer data and 2) a transit counter by using data that were collected through channel overhearing in a cross-layer fashion. These counters are defined as follows.

- *Provider counter* $d_{ic}(n, j)$. This quantity accounts for the presence of new copies of information i 's chunk c , delivered by n to querying nodes within its reach range, during step j . Node n updates this quantity every time it acts as a provider node (e.g., node P in the upper plot of Fig. 2).
- *Transit counter* $r_{ic}(n, j)$. This quantity accounts for the presence of new copies of information i 's chunk c , transferred between two nodes within n 's reach range and

received (or overheard) by n , during step j . Node n thus updates this quantity if it receives (or overhears) an information message,¹ e.g., node R in the lower plot of Fig. 2; thus, the transit counter is the only data structure that needs cross-layer access, i.e., the number of information copies whose transit the node has overheard at lower layers (and subsequently inspected).

The provider and transit counters are updated through the hop count information that is included in the query and information message header. The exact procedure is as detailed follows.

- 1) If node n generates a reply information message that contains chunks of information item i , as an answer to a query for some chunk c that it owns, then a new copy of such chunks is possibly cached at the node that generated the query. Node n must therefore account for the presence of such a new copy at a distance h_Q , which is equal to the number of hops that were covered by the query (see the upper plot of Fig. 2). The provider counter is updated by a quantity that is inversely proportional to the distance h_Q as follows:

$$d_{ic}(n, j) = d_{ic}(n, j) + \frac{1}{h_Q}. \quad (1)$$

Note that the larger the h_Q , i.e., the farthest the new chunk copy, the lesser the added contribution.

- 2) If node n receives or overhears a new transit information i message, which contains a chunk c whose query status was pending, it must then account for the presence of the following copies: 1) a new copy of the chunk that will be cached by a node at a distance of h_Q hops and 2) an existing copy that is cached at a distance of h_P hops (see Fig. 2, lower plot). Thus, following the approach in (1), the transit counter is updated as follows:

$$r_{ic}(n, j) = r_{ic}(n, j) + \frac{1}{h_P} + \frac{1}{h_Q}. \quad (2)$$

Again, the larger h_P is, i.e., the farthest the information provider is, the lesser the contribution of the existing copy becomes.

- 3) The last case accounts for the reception or overhearing of an information message whose contribution must not (or cannot) be related to a corresponding query. This condition may happen for the following two reasons: 1) The corresponding query was already solved (hence, the message is considered duplicated information), or 2) node n had not received the corresponding query (e.g., node n just moved within the radio range of nodes in the return path of the information message and missed the query). In either case, only the contribution due to the presence of the copy at the provider is considered, hence

$$r_{ic}(n, j) = r_{ic}(n, j) + \frac{1}{h_P}. \quad (3)$$

Based on the aforementioned quantities, node n can compute a *presence index* of chunk c of information i , as observed during step j within node n 's reach range. We refer to such a value as $p_{ic}(n, j)$ and define it as

$$p_{ic}(n, j) = \min \{1, d_{ic}(n, j) + r_{ic}(n, j)\}. \quad (4)$$

According to (4), $p_{ic}(n, j)$ comprises the range [0, 1]. A zero-value means that the presence of chunk c of information i was not sensed by n during time step j . Instead, if the chunk is cached one hop away from n , $p_{ic}(n, j)$ is equal to one; this case is the “best,” where the chunk would directly be available to n if needed. Intermediate values between 0 and 1 are recorded when n observes chunks that are cached more than one hop away. Note that multiple contributions of the last kind can sum up to a maximum information presence $p_{ic}(n, j) = 1$, because we rate the dense presence of chunks a few hops away as valuable as a single chunk at a one-hop distance. The information presence index thus computed plays a crucial role in taking caching decisions in both large- and small-sized caches, as described in the next sections.

B. Large-Sized Caches: Computation of the Content Drop Time

We first consider the case in which nodes have a large-sized cache, enough to potentially store a large portion (i.e., 50% or more) of the content that they request. As aforementioned, consuming all the storage resources of a node to cache the retrieved data is not desirable, because the same memory may be shared by different services and applications that run at nodes. Thus, here, we exploit the aforementioned information presence estimate to determine a *cache drop time* after which the retrieved information items are removed from the memory: the goal is to reduce the cache utilization without affecting the performance of the content distribution system.

We denote by $\chi_i(n, j)$ the cache drop time that node n computes at the end of time step j for information item i . Such a drop time applies to all chunks, belonging to information item i , that will be received during time step $(j + 1)$. To compute $\chi_i(n, j)$, node n estimates an overall probability of information presence, by composing the presence indices $p_{ic}(n, j)$ of all chunks of information i , as follows.

Because $p_{ic}(n, j)$ are samples of the chunk presence, node n first needs to quantify the amount of time for which these samples are meaningful. If all nodes run Hamlet, the best guess that node n can take to determine the presence index is to use its local estimate of the cache drop time $\chi_i(n, j - 1)$, assuming that it is not very different from its neighbors'.

Consistent with this reasoning, the contribution of a presence index computed at step k should only be considered for a time $\chi_i(n, k - 1)$. However, discarding contributions *exactly* after a time $\chi_i(n, k - 1)$ leads to an ON/OFF behavior and yields discontinuities in the caching decision process. Moreover, a hard contribution removal threshold is inconsistent with the uncertainty in the knowledge of the neighbors' caching times; the value $\chi_i(n, k - 1)$ used by node n may differ from the

¹If the chunk is segmented into more than one Internet Protocol (IP) packet, overhearing the first segment provides access to the chunk header and the needed information.

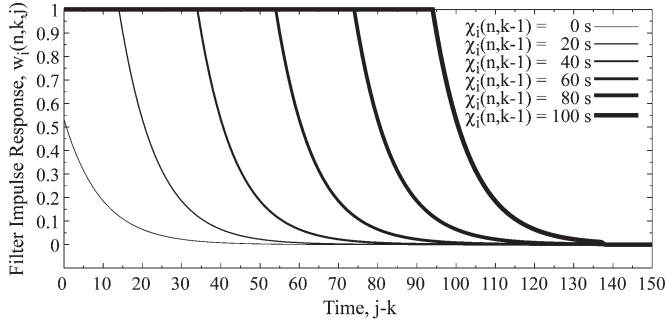


Fig. 3. Filter impulse responses $w_i(n, k, j)$ for different values of $\chi_i(n, k - 1)$, where $k = 1$, $\alpha = 0.9$, and $W = 0.5$. For $k = 1$, the time axis marks the time steps since the chunk has been cached.

cache drop time computed by the nodes within n 's reach range, particularly if they are several hops away.

To account for these factors, we smooth the contributions through an ad hoc filter. At time step j (whose duration is $1/f$), node n weighs each past index $p_{ic}(n, k)$, $k < j$ by a *smoothing factor* $w_i(n, k, j)$, which is defined as

$$w_i(n, k, j) = \begin{cases} 1, & \text{if } j - k \leq \Delta(n, k) \\ \alpha^{j-k-\Delta(n, k)}, & \text{otherwise} \end{cases}$$

$$\text{with } \Delta(n, k) = \lfloor f\chi_i(n, k - 1) - \log_\alpha W \rfloor. \quad (5)$$

This case means that, at time step j , contribution $p_{ic}(n, k)$ maintains its value unchanged if no more than $\Delta(n, k)$ time steps have passed, because the presence index was recorded. Otherwise, the filter forces an exponential decrease of the index value. In (5), the exponential decrease factor is denoted by α , whereas $\Delta(n, k)$ is such that, after a time $\chi_i(n, k - 1)$, the smoothed presence index retains a fraction W of its original value.² In addition, $w_i(n, k, j)$ depends only on the information identifier i , i.e., it has the same value for all chunks c that belong to information i , because the caching time is determined on a per-item basis. For clarity, examples of filter impulse responses for different values of $\chi_i(n, k - 1)$ are shown in Fig. 3.

Next, let us consider time step j and the tagged node n . Node n estimates how complete a single information item i is in its surroundings, due to the contributions measured during a given step k , by summing up presence indices $p_{ic}(n, k)$ that refer to all chunks c of i and smoothing them by the factor $w_i(n, k, j)$ to account for the likelihood that some content has been dropped. We have

$$\phi_i(n, k, j) = w_i(n, k, j) \frac{1}{C} \sum_{c=1}^C p_{ic}(n, k). \quad (6)$$

In other words, (6) reflects the degree of completeness of information i that node n can expect at time step j , only considering contributions that were collected during time step k . We stress that $0 \leq \phi_i(n, k, j) \leq 1$.

The overall presence index for information item i , as shown by node n at time step j , can be computed by summing up the

²Note that multiplying the time interval $\chi_i(n, k - 1)$ by f allows us to express $\Delta(n, k)$ in a number of time steps.

contributions $\phi_i(n, k, j)$ over all steps k as

$$p_i(n, j) = \min \left\{ 1, \sum_{k=j-\tau}^j \phi_i(n, k, j) \right\}. \quad (7)$$

Note that, in (7), τ is the filter memory and contributions that are older than τ time steps are ignored. Thus, to compute (7), a node has to store at most $2\tau + 1$ values for each information item, i.e., the sum in (6) and $\chi_i(n, k)$ for $\tau + 1$ and τ different time steps, respectively. As specified in the performance evaluation, τ can be set to a very small value so that the presence index computation requires minimal memory usage.

Finally, by denoting the maximum cache drop time by M_C , the caching time for the chunks that belong to information i is obtained as

$$\chi_i(n, j) = (1 - p_i(n, j)) M_C. \quad (8)$$

According to (8), in the extreme situation where the entire information i is estimated to be cached within node n 's reach range, i.e., $p_i(n, j) = 1$, the retrieved content will not be stored by the node; on the contrary, when node n observes a complete lack of content i within its reach range, i.e., $p_i(n, j) = 0$, the caching time will be equal to M_C .

C. Small-Sized Caches: Content Replacement

When equipped with a small-sized cache, nodes cannot store all content that they request but are forced to choose which items to keep and which items to discard every time newly retrieved data fill up their memory. In this case, computing cache drop times is clearly not a solution, because the lingering of items in cache is primarily determined by the rate of reception of new content. Therefore, in the presence of limited dedicated storage resources, we exploit the information presence estimate to define a content replacement policy that favors a balanced distribution of data over the network so that all content is as "close" as possible to a requesting node.

The rationale of our content replacement strategy is very similar to the approach employed for the cache drop time computation. Again, we start by identifying the amount of time for which the index $p_{ic}(n, j)$ must be considered valid and define a new smoothing factor $\hat{w}_i(n, k, j)$ to that end as

$$\hat{w}_i(n, k, j) = \begin{cases} 1, & \text{if } j - k \leq \Gamma(n, k) \\ 0, & \text{otherwise} \end{cases}$$

$$\text{with } \Gamma(n, k) = \lfloor f\hat{\chi}_i(n, k - 1) \rfloor. \quad (9)$$

Note that the factor $\hat{w}_i(n, k, j)$, weighting the contribution of samples $p_{ic}(n, k) \forall c$ and at time step j , is necessarily different from the factor employed in the previous section. In particular, we no longer filter the weight index, because the ON/OFF behavior that we previously wanted to avoid now becomes a desirable effect. In fact, it allows us to mimic the sudden removal of chunks from the memory upon the arrival of new content, which is typical of a small-sized cache. In addition, the estimated caching time $\hat{\chi}_i(n, k)$ has a different formulation,

whose details will be discussed at the end of this section, for clarity.

We can thereafter define the completeness of item i , estimated by node n from samples observed at time step k , as

$$\hat{\phi}_i(n, k, j) = \hat{w}_i(n, k, j) \frac{1}{C} \sum_{c=1}^C p_{ic}(n, k) \quad (10)$$

and the overall presence index as

$$\hat{p}_i(n, j) = \sum_{k=j-\tau}^j \hat{\phi}_i(n, k, j). \quad (11)$$

With respect to the equivalent formulation for the case of large-sized caches in (7), the index in (11) is not bounded by 1. Indeed, $\hat{p}_i(n, j)$ is an estimate of the *total amount* of information i in the reach range of node n at time step j .

We thus leverage $\hat{p}_i(n, j)$ as the metric for a content replacement strategy. Upon the reception, at time step j , of new data to be cached and exceeding the free storage memory, node n discards chunks³ of information i associated with the highest $\hat{p}_i(n, j)$ until the remaining data fit the storage constraints. In other words, each node tries to keep content that is estimated to be rarer in its surroundings while dropping data that are evaluated to be already commonly available in the area.

How we can compute the estimated caching time $\hat{\chi}_i(n, k)$ remains to be defined. To this end, the best guess that a node can take is to assume that users in its neighborhood attribute similar caching priorities to the information items, i.e., the ordering of $\hat{p}_i(n, k)$, $\forall i$, that they compute is in agreement with the ordering of the node. Under this assumption, a node n can estimate the amount of time for which a neighboring user will cache a newly received chunk of item i to be inversely proportional to the $\hat{p}_i(n, k)$ that it has locally computed. As aforementioned, we define a maximum cache permanence time M_C , after which, a chunk is discarded to avoid stored information from becoming stale and node movement from leading to inconsistencies with respect to previous information presence ratings. The estimated caching time is then computed as

$$\hat{\chi}_i(n, k) = \left(1 - \frac{\hat{p}_i(n, k)}{\max_i \{\hat{p}_i(n, k)\}} \right) M_C. \quad (12)$$

To provide an interpretation for (12), consider the case of chunks of the most common information item in the area. Because, as aforementioned, a node discards a chunk of information i associated with the highest $\hat{p}_i(n, j)$, the estimated caching time for such a chunk is set to 0 in (12), and caching times of much less popular chunks are, instead, estimated to be much longer (up to M_C).

As the estimated presence decreases, the chance that chunks find space in the cache of requesting nodes grows, reaching the maximum estimated caching time M_C if the information is completely absent from the area, i.e., when $\hat{p}_i(n, j) = 0$.

Note that (12) does not depend on the content query rate as a precise design choice. Indeed, it is not likely that a node

³Assuming contents with few chunks, we do not distinguish among chunks from the same item, which are thus randomly selected for removal.

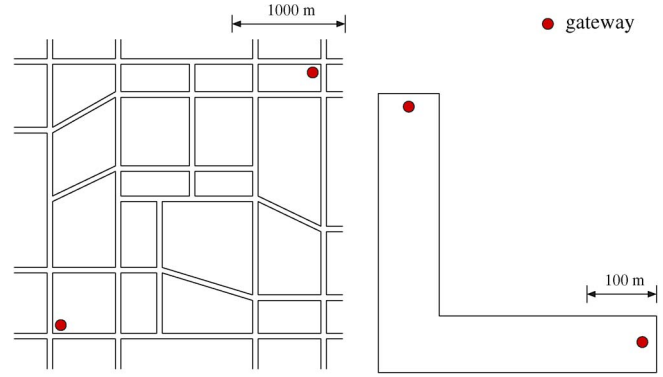


Fig. 4. Simulation scenarios: City (left) and Mall (right).

knows the frequency of requests that were generated for each information item by the different users in the network and assuming such knowledge would noticeably limit the feasibility of the solution. The performance evaluation that we conducted shows that the lack of query rate awareness does not prevent Hamlet from achieving its goals.

V. SIMULATION SCENARIOS AND METRICS

We tested the performance of Hamlet through *ns2* simulations under the following three different wireless scenarios: 1) a network of vehicles that travel in a city section (referred to as *City*); 2) a network of portable devices carried by customers who walk in a mall (*Mall*); and 3) a network of densely and randomly deployed nodes with memory limitations (*memory-constrained nodes*). The three scenarios are characterized by different levels of node mobility and network connectivity.

In the *City* scenario, as depicted in Fig. 4, vehicle movement is modeled by the intelligent driver model with intersection management (IDM-IM), which takes into account car-to-car interactions and stop signs or traffic lights [27]. We simulated a rather sparse traffic, with an average vehicle density of 15 veh/km over a neighborhood of 6.25 km². The mobility model settings, forcing vehicles to stop and queue at intersections, led to an average vehicle speed of about 7 m/s (i.e., 25 km/h). We set the radio range to 100 m in the vehicular scenario, and by analyzing the network topology during the simulations, we observed an average link duration of 24.7 s and a mean of 45 disconnected node clusters concurrently present over the road topology. The *City* scenario is thus characterized by scattered connectivity and high node mobility.

The *Mall* scenario is represented in Fig. 4 as a large L-shaped open space of 400 m of length on the long side, where pedestrian users can freely walk. In this scenario, we record an average of 128 users who walk at an average speed of 0.5 m/s according to the random-direction mobility model with reflections [28]. The node radio range is set to 25 m, leading to an average link duration equal to 43 s, with a mean of ten disconnected clusters of users present at the same time in the network. The connectivity level in the *Mall* is thus significantly higher than in the *City*, whereas node mobility is much lower.

The memory-constrained scenario is similar to the scenario employed for the performance evaluation of the cache

replacement schemes in [9] and [14]. It is composed of 300 wireless nodes deployed over a square area of a side equal to 200 m. Nodes can be static, positioned according to a uniform random distribution, or mobile, wandering according to a random-direction mobility model with reflections. The node speed is uniformly distributed in the range $[0.5v_m, 1.5v_m]$, where v_m is the average node speed—a varying parameter in our simulations. The node radio range is set to 20 m, resulting, for static nodes, in a fully connected network.

In all the scenarios, we deploy two fixed gateway nodes at opposite ends of the topology. Each gateway permanently stores 1/2 of the information items, whereas the other half is provided by the other gateway. Initially, nodes have an empty cache; they randomly request any among the I items that are not in their cache according to a Poisson process with parameter $\lambda_i = \Lambda q_i$ ($1 \leq i \leq I$). Λ is the query generation rate per node, whereas q_i represents the content popularity level (i.e., the probability that, among all possible content, a node requests item i). The *TTL* value for query messages is set to ten and five hops for the case of large- and small-sized caches, respectively, and the query lag time is 50 ms. Note that the impact of all the aforementioned query propagation parameters on the information-sharing behavior has been studied in [23]; here, we only consider what has been identified as a good parameter setting.

With regard to the Hamlet parameters, the estimation frequency is such that $1/f = 0.2M_C$; indeed, through extensive simulations, we observed that the impact of f is negligible, as long as $1/f$ is not greater than 20% of the maximum caching time. As we fix $\tau = fM_C$, this setting of f leads to a value of τ as small as 5. Then, we have $\alpha = 0.9$ and $W = 0.5$; indeed, we have verified that this combination yields a smoother behavior of the presence index $p_i(n, j)$. The values of the remaining parameters are separately specified for large- and small-sized caches.

The information-sharing application lies on top of a User Datagram Protocol (UDP)-like transport protocol, whereas, at the media access control (MAC) layer, the IEEE 802.11 standard in the promiscuous mode is employed. No routing algorithm is implemented: queries use a MAC-layer broadcast transmission, and information messages find their way back to the requesting node following a unicast path. Information messages exploit the request to send/clear to send (RTS/CTS) mechanism and MAC-level retransmissions, whereas query messages of broadcast nature do not use RTS/CTS and are never retransmitted. The channel operates at 11 Mb/s, and signal propagation is reproduced by a two-ray ground model. Simulations were run for 10 000 s.

In the aforementioned scenarios, our performance evaluation hinges upon the following quite-comprehensive set of metrics that are aimed at highlighting the benefits of using Hamlet in a distributed scenario:

- 1) the ratio between solved and generated queries, called solved-queries ratio;
- 2) the communication overhead;
- 3) the time needed to solve a query;
- 4) the cache occupancy.

We have further recorded the spatiotemporal distribution of information and the statistics of information survival, because they help in quantifying the effectiveness of Hamlet in preserving access to volatile information. As aforementioned, we did not explore the problem of cache consistency, because such an issue is orthogonal to this paper.

VI. EVALUATION WITH LARGE-SIZED CACHES

Here, we evaluate the performance of Hamlet in a network of nodes with large storage capabilities, i.e., with caches that can store up to 50% of all information items. Because such characteristics are most likely found in vehicular communication devices, tablets, or smartphones, the network environments under study are the City and Mall scenarios. As discussed in Section IV, in this case, the Hamlet framework is employed to compute the caching time for information chunks retrieved by nodes, with the goal of improving the content distribution in the network while keeping the resource consumption low.

We first compare Hamlet's performance to the results obtained with a deterministic caching strategy, called *DetCache*, which simply drops cached chunks after a fixed amount of time. Then, we demonstrate the effectiveness of Hamlet in the specific task of information survival. In all tests, we assume $I = 10$ items, each comprising $C = 30$ chunks. All items have identical popularity, i.e., all items are requested with the same rate $\lambda = \Lambda/I$ by all network nodes. The choice of equal request rates derives from the observation that, in the presence of nodes with a large-sized memory, caching an information item does not imply discarding another information item; thus, the caching dynamics of the different items are independent of each other and only depend on the absolute value of the query rate. It follows that considering a larger set of items would not change the results but only lead to more time-consuming simulations.

Each query includes 20 B plus 1 B for each chunk request, whereas information messages include a 20-B header and carry a 1024-B information chunk. The maximum caching time M_C is set to 100 s, unless otherwise specified. Queries for chunks that are still missing are periodically issued every 5 s until either the information is fully retrieved or a timeout that is set to 25 s expires.

A. Benchmarking Hamlet

We set the deterministic caching time in *DetCache* to 40 s, and we couple *DetCache* and Hamlet with both the mitigated flooding and Eureka techniques for query propagation. We are interested in the following two fundamental metrics: 1) the ratio of queries that were successfully solved by the system and 2) the amount of query traffic that was generated. The latter metric, in particular, provides an indication of the system effectiveness in preserving locally rich information content: if queries hit upon the sought information in one or two hops, then the query traffic is obviously low. However, whether such a wealth of information is the result of a resource-inefficient cache-all-you-see strategy or a sensible cooperative strategy, e.g., the approach fostered by Hamlet, remains to be seen. Thus, additional metrics that are related to cache occupancy

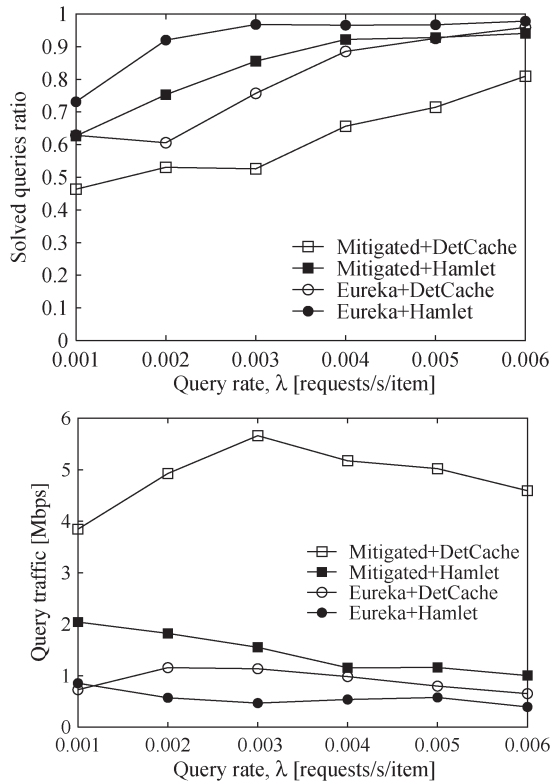


Fig. 5. City: solved-queries ratio (top) and query traffic (bottom) obtained with different schemes versus content request rate.

TABLE I
AVERAGE OCCUPANCY OF THE NODE CACHES, EXPRESSED AS A PERCENTAGE OF THE CHUNKS TOTAL NUMBER FOR $\lambda = 0.003$

Cache usage	Mitigated + DetCache	Mitigated + Hamlet	Eureka + DetCache	Eureka + Hamlet
City	12.3	8.08	13.8	15.1
Mall	15.6	9.9	14.3	12.2

and information cache drop time must be coupled with the aforementioned metrics.

Fig. 5 shows the solved-queries ratio (top plot) and the amount of query traffic (bottom plot) as λ varies in the City scenario. When DetCache is used, the higher the query rate, the larger the number of nodes that cache an information item. This case implies that content can be retrieved with higher probability and also that it is likely to be found in the proximity of the requesting node, thus reducing the query traffic per issued request. Note that, due to its efficient query propagation mechanism, Eureka reduces the propagation of useless queries (and, hence, collisions), yielding a higher solved-queries ratio than mitigated flooding. However, it is evident that deterministic caching does not pay off as much as cooperative caching does in Hamlet. Table I shows that the average occupancy of node caches in Hamlet is comparable to the values observed with DetCache. Thus, it is the quality, not the quantity, of the information cached by Hamlet that allows it to top a sophisticated propagation scheme such as Eureka as far as the solved-queries ratio is concerned.

The positive effect of the caching decisions can also be observed in Fig. 5 in terms of the reduced overhead and latency

TABLE II
AVERAGE QUERY SOLVING TIME (IN SECONDS), WITH $\lambda = 0.003$

Cache Time	Mitigated + DetCache	Mitigated + Hamlet	Eureka + DetCache	Eureka + Hamlet
City	299.442	140.419	157.228	57.617
Mall	24.242	13.795	33.420	24.816

in solving queries. Indeed, Hamlet reduces the overhead by shortening the distance between requesting nodes and desired information content. Similarly, Table II shows how sensible caching choices can significantly reduce the time required to solve queries, again due to the homogeneous availability of information that they generate in the network.

Further proof of such virtuous behavior by Hamlet is provided in Fig. 6, where mitigated flooding is used for query propagation. The figure depicts the time evolution of content presence over the road topology for one information item; in particular, the z -axis of each plot shows the fraction of different chunks that comprise an information item that are present in a squared area of 600 m². On the one hand, it can be observed that mitigated flooding with DetCache creates a sharp separation between the area where the content source resides, characterized by high item availability, and the region where, due to vehicular traffic dynamics, information-carrying nodes rarely venture. On the other hand, Hamlet favors the diffusion of content over the entire scenario so that nodes in areas away from the information source can also be served.

Fig. 7 refers to the Mall scenario. The poor performance of Eureka in this case is due to the lack of information items over large areas of the Mall scenario, resulting in queries not being forwarded and, thus, remaining unsolved [23]. Interestingly, Hamlet greatly reduces the query traffic for any λ , although providing a much higher solved-queries ratio. With regard to the caching occupancy, because Hamlet leads to results that are comparable with the results obtained with DetCache (see Table I, Mall scenario), it can be asserted that the performance gain achieved through Hamlet is due to the more uniform content distribution across node caches. Finally, Table II confirms that such an improved availability of information shortens the waiting time to receive requested items.

When comparing results obtained from the Mall and City scenarios, we note that the solved-queries ratio is consistently lower. We recall that vehicular mobility in the City environment is characterized by scattered connectivity but high node speed, whereas the Mall environment provides a better network connectivity level but reduced node mobility. The low node mobility in the Mall keeps items away from the sources of unpopular items for long periods of time. Thus, the probability of solving requests for such rare content is low, unless an efficient caching scheme allows nodes to preserve at least a few copies of every item in every neighborhood, as Hamlet does. It is also worth pointing out that, with respect to the City environment, the Mall includes a smaller number of nodes; thus, fewer queries are issued, and a much smaller amount of query traffic is generated.

Finally, we may wonder how well Hamlet performs with respect to DetCache when the cache time employed by the latter approach is set to a value other than 40 s. Through extensive

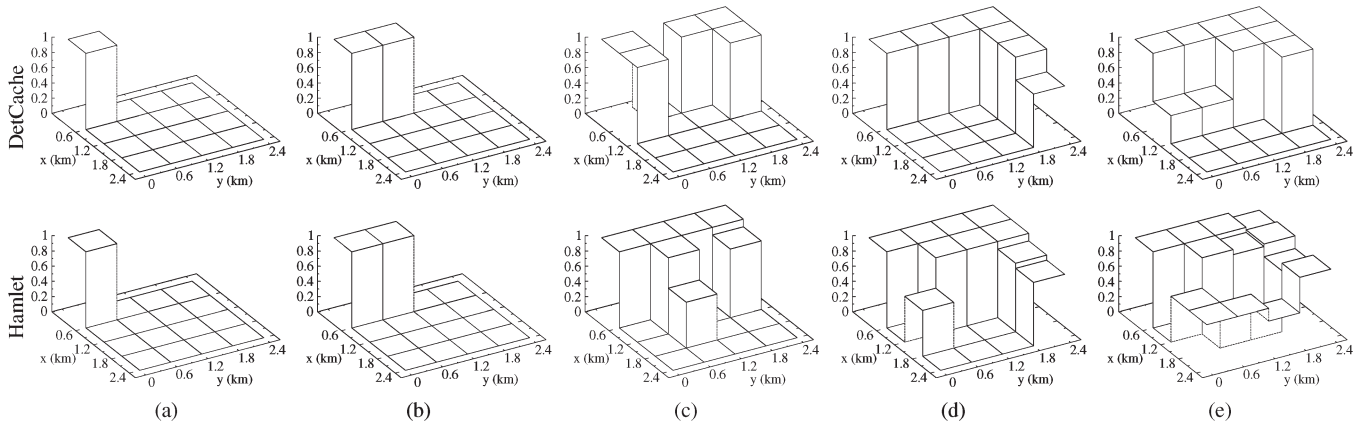


Fig. 6. City: Space-Time evolution of one information item during the first 400 s of simulation, with mitigated flooding and DetCache (top) and Hamlet (bottom). The z-axis shows the content completeness in each spatial slot of 600 m², with a value of 1 meaning that all of the items' chunks can be found in the slot. (a) t = 0 s. (b) t = 100 s. (c) t = 200 s. (d) t = 300 s. (e) t = 400 s.

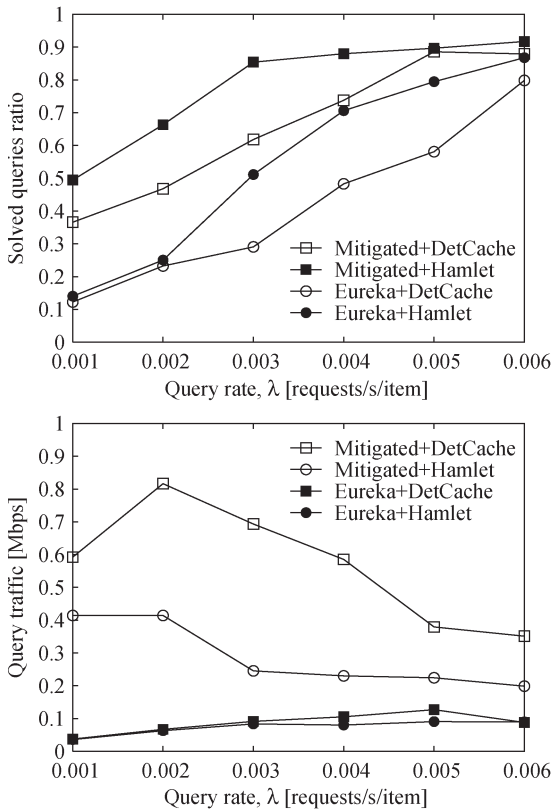


Fig. 7. Mall: Solved-queries ratio (top) and query traffic (bottom) with different schemes versus content request rate.

simulations, we found that DetCache achieves the best tradeoff between solved-queries ratio and cache occupancy when its caching time is set to the average value of cache drop time obtained using Hamlet [24]. This case is a further proof of how Hamlet can autonomously determine an optimal cache drop time at each node.

B. Information Survival

Next, we look at how long information survives in self-organized networks without centralized control. We say that, at a given time instant, a particular item has *survived* if each of its chunks is cached by at least one node in the network.

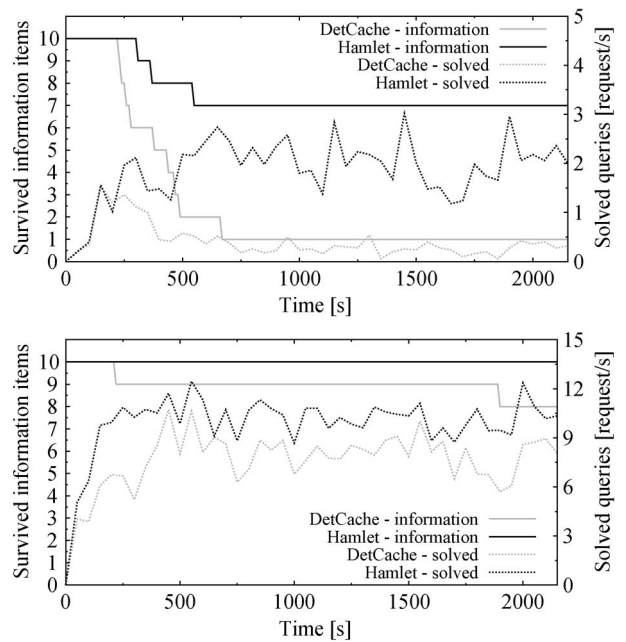


Fig. 8. Information survival in the Mall (top) and City (bottom) scenarios. The temporal behavior of the survived information and solved queries when the gateway nodes are switched off at $t = 200$ s.

Fig. 8 shows the time evolution of the number of survived items when the gateway nodes are turned off after $t = 200$ s in the Mall and the City. The same plots also present the overall solved query rate as time elapses, as observed at 50-s time discretization steps. We consider that DetCache and Hamlet are coupled with mitigated flooding, and we set λ to 0.003 request/s.

In the Mall, the adaptive caching time introduced by Hamlet helps most of the items survive in the system. In the City, the high node mobility supports the circulation of information and, thus, its survival after the gateways shut down. However, even through mobility, DetCache allows some of the items to vanish, whereas Hamlet succeeds in keeping all of items around.

In the Mall, we further explore the reaction of DetCache and Hamlet to different gateway switch-off times. In Fig. 9, for each gateway switch-off time, we show the evolution of the number of survived items at some landmark time instants

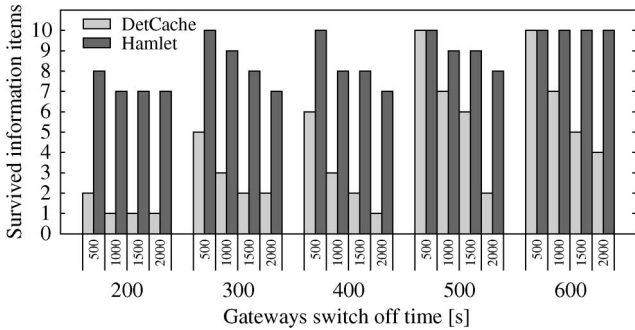


Fig. 9. Mall: Information survival for different gateway switch-off times. The smaller numbers on the x -axis indicate the landmark time (in seconds) to which the number of survived items refers.

(computed from the start of the simulation). Clearly, the later the gateways are shut down, the higher the probability of information survival, because the information has more time to spread through the network. We observe that Hamlet can maintain information presence equal to 100% if the information is given enough time to spread, i.e., gateways are disabled after 600 s or more, whereas DetCache loses half the items within the first 2000 s of simulation. We could wonder whether caching times give the edge to either Hamlet or DetCache. However, the average caching time in Hamlet ranges from 37 s to 45 s, depending on the gateway switch-off times and on the specific information item considered. These values are very close to the DetCache caching time of 40 s, showing that Hamlet improves information survival by better distributing content in the network and not by simply caching them for longer periods of time.

VII. EVALUATION WITH SMALL-SIZED CACHES

We now evaluate the performance of Hamlet in a network where a node cache can accommodate only a small portion of the data that can be retrieved in the network. As an example, consider a network of low-cost robots that are equipped with sensor devices, where maps that represent the spatial and temporal behavior of different phenomena may be needed by the nodes and have to be cached in the network. We thus consider the memory-constrained scenario introduced in Section V and employ the Hamlet framework to define a cache replacement strategy, as detailed in Section IV.

In such a scenario, the caching dynamics of the different information items become strongly intertwined. Indeed, caching an item often implies discarding different previously stored content, and as a consequence, the availability of one item in the proximity of a node may imply the absence of another item in the same area. Thus, in our evaluation, it is important to consider a large number of items, as well as to differentiate among these items in terms of popularity. We consider an overall per-node query rate $\Lambda = 0.1$ and sets of several hundreds of items. We assume that popularity levels q_i are distributed according to the Zipf law, which has been shown to fit popularity curves of content in different kinds of networks [29]. When not stated otherwise, the Zipf distribution exponent is set to 0.5. Such a value was selected, because it is close to the values observed in the real world [29], and the skewness that it introduces in the

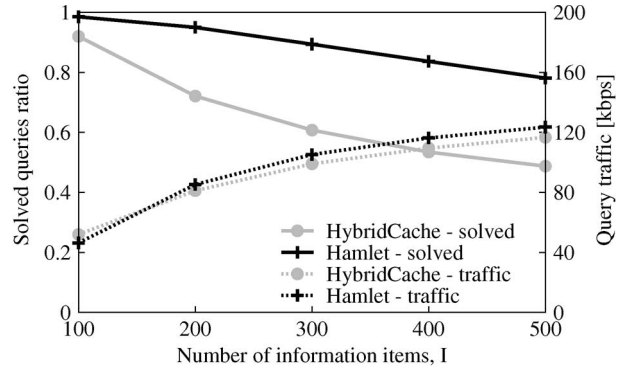


Fig. 10. Static memory-constrained nodes: Solved-queries ratio and query traffic as the information set size varies, with HybridCache and Hamlet.

popularity distribution is already sufficient to make differences emerge between the caching schemes that we study. In any case, we provide an analysis of the impact of the Zipf exponent at the end of this section.

We assume that nodes can cache at most ten items, which correspond to a percentage between 2% and 10% of the entire information set, depending on the considered value of I . In addition, we set $C = 1$ to account for the smaller size of information items typically exchanged by memory-constrained nodes and M_C to 300 s, because the increased network connectivity prolongs the reliability of information presence estimation.

Here, we compare Hamlet with the well-known HybridCache cache replacement technique [9]. In HybridCache, a node that requests an item always caches the received data. Instead, a node on the data path caches the information if its size is small; otherwise, it caches the data path, provided that the content copy is not very far away. When the maximum cache size is capped, content in excess is dropped according to a metric based on the number of requests observed for the different items. Because HybridCache does not exploit information presence estimation, it is less demanding than Hamlet in terms of computation and memory capabilities.

We couple both schemes with mitigated flooding. While deriving the results, we noted that caching the data paths leads to poor performance due to the high cache replacement frequency in the simulated scenarios. Therefore, we set the HybridCache parameters so that the following two conditions are satisfied: 1) The size of the data never results in data path caching but always in information caching, and 2) mitigated flooding is always employed for query forwarding. In addition, to reduce the number of query transmissions in the network, queries for missing chunks are not reissued, and both Hamlet and HybridCache are coupled with the preferred group broadcasting (PGB) technique [30].

A. Benchmarking Hamlet

Let us first focus on the memory-constrained scenario outlined in Section V with static nodes. Fig. 10 presents the solved-queries ratio and the overall query traffic versus the information set size. We observe that Hamlet reacts better to the growth of the number of items than HybridCache, without incurring any

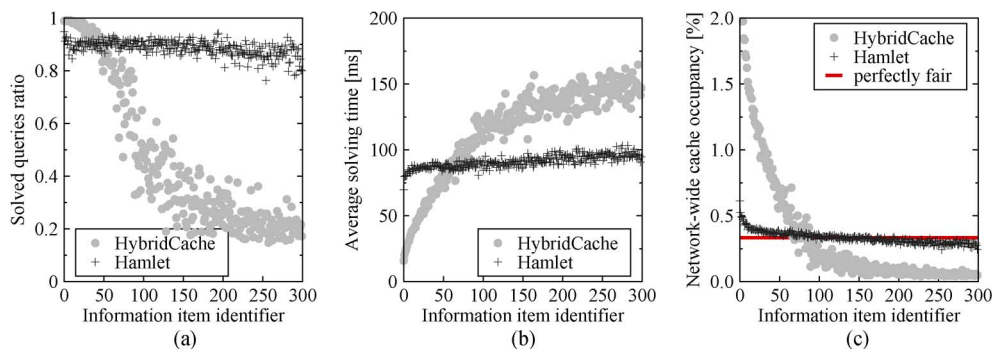


Fig. 11. Static memory-constrained nodes. (a) Query-solving ratio, (b) time, and (c) average networkwide cache occupancy for each item when using HybridCache and Hamlet, with $I = 300$. In (c), the red horizontal line represents perfect fairness in cache occupancy among different items.

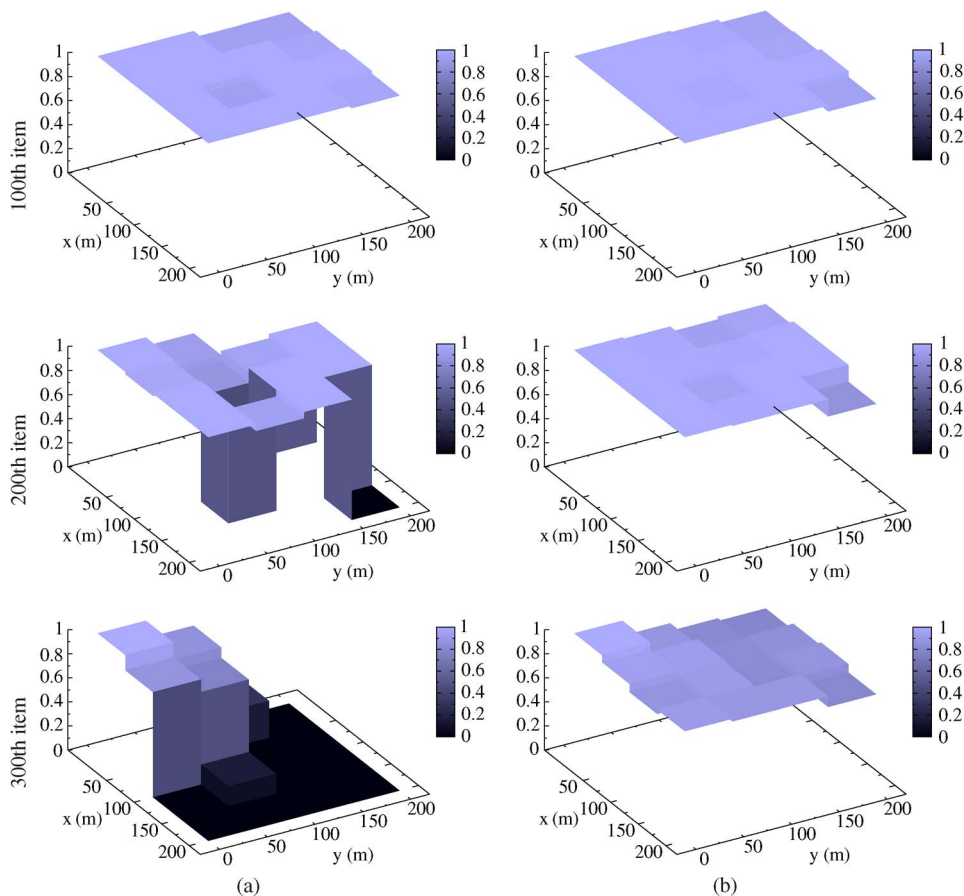


Fig. 12. Static memory-constrained nodes: Spatial distribution of the 100th, 200th, and 300th items, averaged over time, for Zipf distribution exponents under HybridCache and Hamlet, with $I = 300$. The z -axis in the plots shows the mean content completeness in each spatial slot, with a value of 1, meaning that the entire content can be found in the same spatial slot. (a) HybridCache. (b) Hamlet.

penalty in terms of network load, as shown by the similar query traffic generated by the two schemes.

Observing the performance of Hamlet and HybridCache on a per-item basis allows a deeper understanding of the results. In Fig. 11(a), we show the solving ratio of the queries for each item when $I = 300$. Along the x -axis, items are ordered in decreasing order of popularity, with item 1 representing the most sought-after information and item 300 the least requested information. Unlike Hamlet, HybridCache yields extremely skewed query solving ratios for the different content; a similar observation also applies to the time needed to solve queries, as shown in Fig. 11(b). The explanation for such behavior lies

in the distribution of information in the network. Fig. 11(c) depicts the average percentage of memory used to cache a given item, aggregated over all network nodes. As expected from the previous results, HybridCache fosters the storage of popular content, whereas it disregards content that is less requested, even if it represents two thirds of the whole information set. Instead, Hamlet achieves, in a completely distributed manner, a balanced networkwide utilization of node caches. Indeed, the results of Hamlet are very close to the most even cache occupancy that we can have, represented by the horizontal red line in the plot and corresponding to the case where the total network storage capacity is equally shared among the I items.

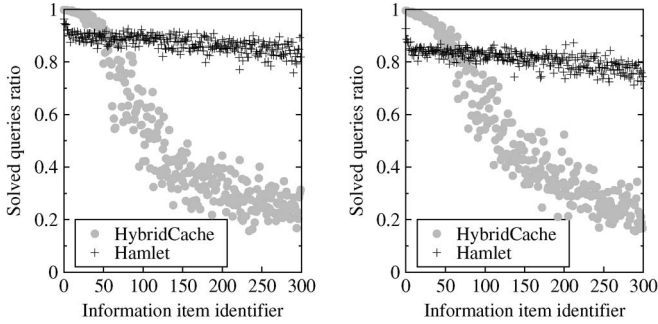


Fig. 13. Memory-constrained mobile nodes: Query-solving ratio for each information item when using HybridCache and Hamlet, with $I = 300$. The plots refer to v_m that is equal to 1 m/s (left) and 15 m/s (right).

Furthermore, it is not only the sheer quantity of data that makes a difference but its spatial distribution also plays a major role. If several nodes cache a rare item but they are all very close to each other, queries that were generated in other areas of the network take more hops to be satisfied. This case happens with HybridCache, as proven by the spatial distribution of the 100th, 200th, and 300th items, as shown in Fig. 12(a). Conversely, the spatial distribution achieved by Hamlet, as shown in Fig. 12(b), is more uniform, leading to a faster more likely resolution of queries.

We now compare the performance of HybridCache and Hamlet in the scenario with memory-constrained mobile nodes. We test the two schemes when $I = 300$ and for an average node speed v_m equal to 1 and 15 m/s.

The solved-queries ratio recorded with HybridCache and Hamlet on a per-item basis are shown in Fig. 13. Comparing the left and right plots, we note that the node mobility, even at high speed, does not seem to significantly affect the results due to the high network connectivity level. The spatial redistribution of content induced by node movements negatively affects the accuracy of Hamlet's estimation process, which explains the slight reduction in the solved query ratio at 15 m/s. That same movement favors HybridCache, at least at low speed, because it allows unpopular information to reach areas that are far from the gateway. However, the difference between the two schemes is evident, with Hamlet solving an average of 20% requests more than HybridCache, when nodes move at 15 m/s.

Note that, for the query resolution delay and the average cache utilization at the network nodes, we obtained qualitatively similar results as in the static case, with Hamlet achieving more homogeneous solving times and fairer distribution of content in the network than HybridCache.

B. Impact of the Zipf Distribution Skewness

Finally, we study the impact of the Zipf distribution exponent on the performance of the cache replacement strategies. We recall that an exponent that is equal to zero implies perfect homogeneity, i.e., Zipf distribution that degenerates into a uniform distribution, whereas the difference in popularity among content becomes much more unbalanced as the exponent grows. We focus on a network where ten items are available and each node can cache at most one complete item. The choice of this setting is mandated by the fact that, in the presence of

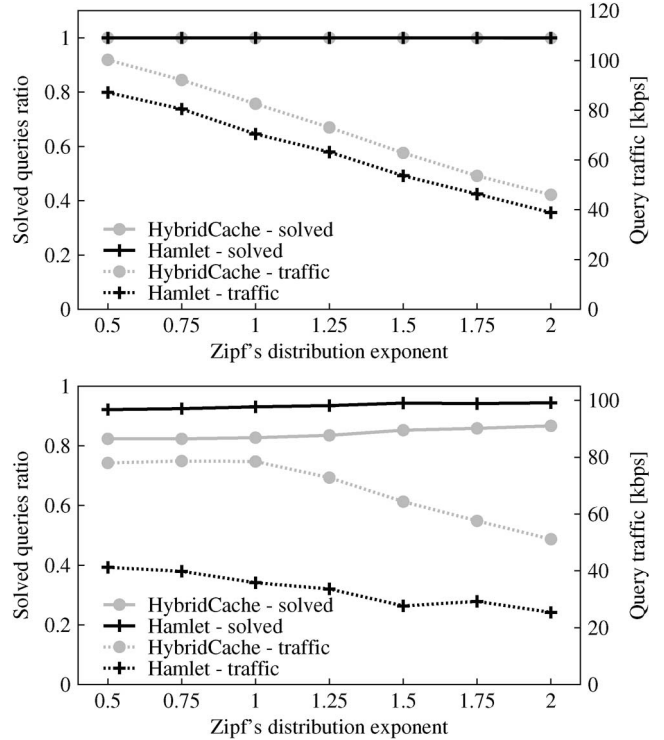


Fig. 14. Memory-constrained static (top) and mobile (bottom) nodes: Solved-queries ratio and query traffic as the Zipf distribution exponent varies when using HybridCache and Hamlet, with $I = 10$.

hundreds of different items, unbalanced popularity distributions (i.e., exponents higher than 0.5) lead to very low λ_i for the 100 or so least popular items, thus making requests for such content extremely rare.

Fig. 14 depicts the evolution of the solved-queries ratio and the query traffic as the Zipf exponent ranges vary. By comparing the two plots, we note that the presence of mobility ($v_m = 1$ m/s) leads to a higher number of unsolved requests and in a larger amount of traffic generated within the network under HybridCache, because queries propagate far from the source without finding the desired item. However, what is most interesting is how the network load tends to decrease as the Zipf exponent grows, both in the absence and presence of node movements. On the one hand, higher values of the exponent lead to more unbalanced query rates, with very few items that are extremely popular and a long tail of seldom-accessed data. Being requested so often, popular items become commonly found in nodes caches, and the relative queries are solved faster, generating small traffic. On the other, when the Zipf exponent is small, the distribution of queries is more balanced, with information more evenly distributed in the network. This case implies that items can usually be found but are hardly cached very close to the requesting node. Thus, the different items are all requested at a fairly high rate but are not immediately found, generating larger query traffic.

VIII. CONCLUSION

We have introduced Hamlet, which is a caching strategy for ad hoc networks whose nodes exchange information items in a peer-to-peer fashion. Hamlet is a fully distributed scheme

where each node, upon receiving a requested information, determines the cache drop time of the information or which content to replace to make room for the newly arrived information. These decisions are made depending on the perceived “presence” of the content in the node’s proximity, whose estimation does not cause any additional overhead to the information sharing system. We showed that, due to Hamlet’s caching of information that is not held by nearby nodes, the solving probability of information queries is increased, the overhead traffic is reduced with respect to benchmark caching strategies, and this result is consistent in vehicular, pedestrian, and memory-constrained scenarios. Conceivably, this paper can be extended in the future by addressing content replication and consistency. The procedure for information presence estimation that was developed in Hamlet can be used to select which content should be replicated and at which node (even if such a node did not request the content in the first place). In addition, Hamlet can be coupled with solutions that can maintain consistency among copies of the same information item cached at different network nodes, as well as with the versions stored at gateway nodes.

REFERENCES

[1] J. Wortham (2009, Sep.). Customers Angered as iPhones Overload AT&T. *The New York Times*. [Online]. Available: <http://www.nytimes.com/2009/09/03/technology/companies/03att.html>

[2] A. Lindgren and P. Hui, “The quest for a killer app for opportunistic and delay-tolerant networks,” in *Proc. ACM CHANTS*, 2009, pp. 59–66.

[3] P. Padmanabhan, L. Gruenwald, A. Vallur, and M. Atiquzzaman, “A survey of data replication techniques for mobile ad hoc network databases,” *VLDB J.*, vol. 17, no. 5, pp. 1143–1164, Aug. 2008.

[4] A. Derhab and N. Badache, “Data replication protocols for mobile ad hoc networks: A survey and taxonomy,” *IEEE Commun. Surveys Tuts.*, vol. 11, no. 2, pp. 33–51, Second Quarter, 2009.

[5] B.-J. Ko and D. Rubenstein, “Distributed self-stabilizing placement of replicated resources in emerging networks,” *IEEE/ACM Trans. Netw.*, vol. 13, no. 3, pp. 476–487, Jun. 2005.

[6] G. Cao, L. Yin, and C. R. Das, “Cooperative cache-based data access in ad hoc networks,” *Computer*, vol. 37, no. 2, pp. 32–39, Feb. 2004.

[7] C.-Y. Chow, H. V. Leong, and A. T. S. Chan, “GroCoca: Group-based peer-to-peer cooperative caching in mobile environment,” *IEEE J. Sel. Areas Commun.*, vol. 25, no. 1, pp. 179–191, Jan. 2007.

[8] T. Hara, “Cooperative caching by mobile clients in push-based information systems,” in *Proc. CIKM*, 2002, pp. 186–193.

[9] L. Yin and G. Cao, “Supporting cooperative caching in ad hoc networks,” *IEEE Trans. Mobile Comput.*, vol. 5, no. 1, pp. 77–89, Jan. 2006.

[10] N. Dimokas, D. Katsaros, and Y. Manolopoulos, “Cooperative caching in wireless multimedia sensor networks,” *ACM Mobile Netw. Appl.*, vol. 13, no. 3/4, pp. 337–356, Aug. 2008.

[11] Y. Du, S. K. S. Gupta, and G. Varsamopoulos, “Improving on-demand data access efficiency in MANETs with cooperative caching,” *Ad Hoc Netw.*, vol. 7, no. 3, pp. 579–598, May 2009.

[12] Y. Zhang, J. Zhao, and G. Cao, “Roadcast: A popularity-aware content sharing scheme in VANETs,” in *Proc. IEEE Int. Conf. Distrib. Comput. Syst.*, Los Alamitos, CA, 2009, pp. 223–230.

[13] E. Cohen and S. Shenker, “Replication strategies in unstructured peer-to-peer networks,” in *Proc. ACM SIGCOMM*, Aug. 2002, pp. 177–190.

[14] B. Tang, H. Gupta, and S. Das, “Benefit-based data caching in ad hoc networks,” *IEEE Trans. Mobile Comput.*, vol. 7, no. 3, pp. 289–304, Mar. 2008.

[15] W. Li, E. Chan, and D. Chen, “Energy-efficient cache replacement policies for cooperative caching in mobile ad hoc network,” in *Proc. IEEE WCNC*, Kowloon, Hong Kong, Mar. 2007, pp. 3347–3352.

[16] M. K. Denko and J. Tian, “Cross-layer design for cooperative caching in mobile ad hoc networks,” in *Proc. IEEE CCNC*, Las Vegas, NV, Jan. 2008, pp. 375–380.

[17] H. Chen, Y. Xiao, and X. Shen, “Update-based cache replacement policies in wireless data access,” in *Proc. BroadNets*, Boston, MA, Oct. 2005, pp. 797–804.

[18] J. Xu, Q. Hu, W.-C. Lee, and D. L. Lee, “Performance evaluation of an optimal cache replacement policy for wireless data dissemination,” *IEEE Trans. Knowl. Data Eng.*, vol. 16, no. 1, pp. 125–139, Jan. 2004.

[19] J. Cao, Y. Zhang, G. Cao, and L. Xie, “Data consistency for cooperative caching in mobile environments,” *Computer*, vol. 40, no. 4, pp. 60–66, Apr. 2007.

[20] N. Dimokas, D. Katsaros, and Y. Manolopoulos, “Cache consistency in wireless multimedia sensor networks,” *Ad Hoc Netw.*, vol. 8, no. 2, pp. 214–240, Mar. 2010.

[21] T. Hara, “Effective replica allocation in ad hoc networks for improving data accessibility,” in *Proc. IEEE INFOCOM*, Anchorage, AK, Apr. 2001, pp. 1568–1576.

[22] K. Chen, S. H. Shah, and K. Nahrstedt, “Cross-layer design for data accessibility in mobile ad hoc networks,” *Wireless Pers. Commun.*, vol. 21, no. 1, pp. 49–76, Apr. 2002.

[23] M. Fiore, C. Casetti, and C.-F. Chiasserini, “Information density estimation for content retrieval in MANETs,” *IEEE Trans. Mobile Comput.*, vol. 8, no. 3, pp. 289–303, Mar. 2009.

[24] M. Fiore, F. Mininni, C. Casetti, and C.-F. Chiasserini, “To cache or not to cache?” in *Proc. IEEE INFOCOM*, Rio de Janeiro, Brazil, Apr. 2009, pp. 235–243.

[25] T. Hara, “Replica allocation methods in ad hoc networks with data update,” *ACM Mobile Netw. Appl.*, vol. 8, no. 4, pp. 343–354, Aug. 2003.

[26] W. Li, E. Chan, D. Chen, and S. Lu, “Maintaining probabilistic consistency for frequently offline devices in mobile ad hoc networks,” in *Proc. IEEE ICDCS*, Montreal, QC, Canada, Jun. 2009, pp. 215–222.

[27] M. Fiore, J. Haerri, F. Filali, and C. Bonnet, “Vehicular mobility simulation for VANETs,” in *Proc. IEEE ANSS*, Norfolk, VA, Mar. 2007, pp. 301–309.

[28] E. M. Royer, P. M. Melliar-Smith, and L. E. Moser, “An analysis of the optimum node density for ad hoc mobile networks,” in *Proc. IEEE ICC*, Helsinki, Finland, Jun. 2001, pp. 857–861.

[29] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, “Web caching and Zipf-like distributions: Evidence and implications,” in *Proc. IEEE INFOCOM*, New York, Mar. 1999, pp. 126–134.

[30] V. Naumov, R. Baumann, and T. Gross, “An evaluation of intervehicle ad hoc networks based on realistic vehicular traces,” in *Proc. ACM MobiHoc*, Florence, Italy, May 2006.



Marco Fiore (M’05) received the M.Sc. degrees from the University of Illinois at Chicago and the Politecnico di Torino, Torino, Italy, in 2003 and 2004, respectively, and the Ph.D. degree from the Politecnico di Torino in 2008.

He was a Visiting Researcher with Rice University, Houston, TX, in 2006 and 2007 and with the Universitat Politècnica de Catalunya, Barcelona, Spain, in 2008. He is currently an Assistant Professor with Institut National des Sciences Appliquées (INSA) Lyon, Villeurbanne, France, where he is also an Institut National de Recherche en Informatique et en Automatique (INRIA) Researcher within the Smart Wireless Networking Team, which is hosted by the Centre of Innovation in Telecommunications and Integration of Services (CITI) Laboratory. His research interests include mobile networking, particularly vehicular networks, mobility modeling and analysis, opportunistic communication, and mobile network security.



Claudio Casetti (M’05) received the M.Sc. degree in electrical engineering and the Ph.D. degree in electronic engineering from the Politecnico di Torino, Torino, Italy, in 1992 and 1997, respectively.

He is currently an Assistant Professor with the Dipartimento di Elettronica e Telecomunicazioni, Politecnico di Torino. He is a coauthor of more than 120 journal papers and conference proceedings in networking and is the holder of three patents. His research interests include ad hoc wireless networks and vehicular networks.



Carla-Fabiana Chiasserini (M'98–SM'09) received the M.Sc. (*summa cum laude*) degree in electrical engineering from the University of Florence, Florence, Italy, in 1996 and the Ph.D. degree from the Politecnico di Torino, Torino, Italy, in 1999.

Since 1999, she has been with the Dipartimento di Elettronica e Telecomunicazioni, Politecnico di Torino, where she is currently an Associate Professor. From 1998 to 2003, she was a Visiting Researcher with the University of California at San Diego, La Jolla. Her research interests include the architecture, protocols, and performance analyses of wireless networks. She has published more than 180 papers in prestigious journals and leading international conference proceedings and is the coholder of three patents.

Dr. Chiasserini is the Associated Editor for several journals, including the *IEEE Wireless Communications Magazine* and the *IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS*.