# Hoarding Content in M-Learning Context

Anna Trifonova and Marco Ronchetti

*Department of Information and Communication Technologies*
*University of Trento, 38050, Povo (Trento), Italy*
*{Anna.Trifonova, Marco.Ronchetti}@dit.unitn.it*

## Abstract

*With the advances in mobile technologies is now possible to support learners and teachers activities on the move. We analyzed the functionalities that should be provided by a general mobile learning platform and identified a problem that is weakly studied, namely support for offline usage of learning material, called hoarding. Hoarding can use some techniques used by caching and pre-fetching schemas, but in most cases the goal of the last two is to reduce latency time, bandwidth consumption and/or servers workload, while in hoarding the aim is to reduce the size of the hoarding set, keeping the accuracy very high. We want to study the parameters that could help hoarding algorithm to face the peculiarity in m-learning scenario and our final goal is to provide an efficient strategy, taking into account additional parameters, extracted automatically by the system.*

## 1. Introduction

E-learning is growing fast and many Universities and companies are now supporting an e-learning solution. There is no doubt that WWW is a successful educational medium. On the other hand the rush in the wireless and mobile technologies creates opportunity for new research field - 'mobile learning' that includes a wide variety of applications and new teaching and learning techniques [4]. In their tries of finding the best way to apply mobile devices in education people are experimenting with different fields. Courses modules were created throughout different projects for people with numeracy and literacy problems, for kids, for university students, for teachers, for computer science subjects, psychology or languages.

In our previous work we classified services that are specific and should be provided by a general m-learning platform [3] and later we concentrate on one of these services as a concrete problem to solve. Namely this is the hoarding of content for offline usage. Hoarding is a technique for selecting set of documents to be uploaded and used when disconnected. Related terms are caching and pre-fetching, though they are used considering online conditions and Web performance. Caching is a technique

for keeping content that has been requested by one user available on nearest server for certain amount of time so other requestors can access it faster. Pre-fetching is a technique that tries to improve the clients' experience by guessing what will be needed in the near future and caching it. Different schemes of caching and pre-fetching are proposed and the goal is to reducing network traffic, minimizing access latency, bottlenecks, servers' workload and etc. Although the goal of hoarding content for offline usage is shifted from the one of caching and pre-fetching, some techniques can be reused. However while in the online case one can balance between the accuracy of the cache and the added traffic, the situation we face the accuracy required is very high and added limitation is the memory available. The characteristics of the learning scenario expose additional information to be considered and possibility to improve the existing solutions.

## 2. The hoarding process

The hoarding process should consist of few steps that we can formalize as follows:
1. Predict the 'starting point' of current user for his/her next learning session and set its priority to maximum.
2. Create 'candidate' set of the related documents (LO).
3. Predict the most probable session path or sequence of LO the user will be following.
4. Prune the candidate set i.e. exclude the objects that will not be needed by the user, thus making it smaller.
5. Find the priority to all objects still in the hoarding set based on their importance for the next session (higher if the probability the object will be used soon is high).
6. Sort the objects, based on their priority and produce an ordered list of objects.
7. Cache, starting from the beginning of the list, putting on the device those objects with bigger priority, until available memory is filled in.

We can see that the algorithm will strongly depend on system's knowledge about the user. This knowledge includes user's learning style, natural learning habits and abilities, the level of expertise in the studying field and topic. It can be acquired by direct assessment of the user, by questionnaires and quizzes, but also by observing and

analyzing the user behavior during studying with the system, i.e. automatically discovering user's learning style, preferences, acquired knowledge and etc.

For making the things clear we can consider two separate engines. One will deal with observing the user and creating user models and the other for the hoarding. We call the first one 'User Behavior Analyzing Engine' and it should be discussed later on. The hoarding algorithm should take as input the output from the 'User Behavior Analyzing Engine' (i.e. the user models with the similarities and the differences of the particular user with the common users' behavior and the current user preferences and learning history) and additional information about the learning content itself (domain knowledge). This will be also discussed further.

Some questions appear on this stage:
- What is a 'session' in the mobile learning scenario?
- What is the best starting point for user's next session?
- What is the candidate set for pruning?
- How (based on what) to prune the hoarding set?
- How to prioritize the LO?
- How can we predict the learning path/sequence?
- What are the important parameters of the user behavior that have influence on the prediction?
- How to use those parameters for predicting and/or pruning and do they have different significance for the prediction and/or pruning process?
- How do we measure the successfulness of the automatic hoarding and how do we improve the work of the algorithms, considering these measures?

The rest of the paper attempts to answer some of these questions, starting from the last one (sec. 3) and discuss important things for solving some problems. In sec. 4 we look at the ways to find the student's learning sequence and discuss the lack of initial data for finding the user's starting point. Section 5 shortly argues about the difference between the general definition of 'session' in the WWW world and the one applicable in the mobile scenario. In 6 we discuss some aspects of the relations between learning objects (LO) and how their correlations can be used in the hoarding for pruning. Afterwards in sec. 7 we discuss possible ways to model the student (discussed in more details elsewhere [5]) and his/her behavior so we can 'predict' what materials will be needed during the offline period. In 8 we talk about the additional data of the learning material and the studied topic that might be useful for the hoarding. Conclusions in 9 are followed by references in 10.

## 3. Measure the goodness

An important question is to measure the goodness of the hoarding and to try to improve it every next time. Often used metric in the evaluation of caching proxies is the *hit ratio*, calculated by dividing the number of hits by the total number of uploaded predictions (cache size). It is a good measure for hoarding systems, though a better measure is the *miss ratio* - a percentage of accesses for which the cache is ineffective. Authors in [2] defined a *miss cost* as a main difference in the evaluation of caching and hoarding systems. In caching/pre-fetching systems the misses in the prediction reflect as a time penalty as the missing content should be retrieved from the web. This differs from the mobile case where with unavailable internet connection a miss in the hoard might be fatal. In order to quantify this measure it is possible to demand a user rating on every miss, using few different impact values. [2] also defines *time to first miss* measure - a simple count between the start of the disconnected operation and the first hoard miss. Note that this evaluation criterion can be used only on real-use of a system (and its hoard part). It is also strongly related to the hoarding size. Another possible measurement is the *miss-free hoard size*, defined as the minimum amount of disc space that a particular hoarding system would require to allow a complete disconnection period to take place without any misses.
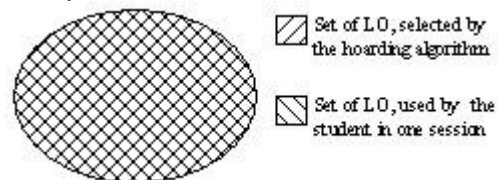


*Figure 1: The ideal hoarding set*

The goal of the hoarding algorithm is to maximize the 'hit rate' and at the same time to minimize the 'miss rate'. In other words the ideal case is to achieve hit_rate=100% and miss_rate=0%, which would mean than the hoarding set contains *all* and *only* the items that the user needs during her/his studying session as shown on figure 1.
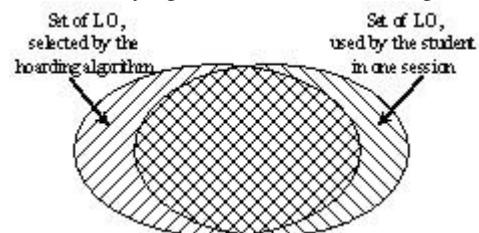


*Figure 2: The expected picture*

Though the ideal picture is to select *all* and *only* those items that the user will use it is obvious that in real system such an ideal situation is almost impossible to reach. Most probably we will have some (desirably big) overlapping between the cached by the hoarding algorithm LO and those LO really requested by the learner.

If *miss cost* measure mentioned before is used it might be better to try to minimize overall cost of all misses.

As mentioned before the hoarding module should be able to analyze how successfully the previous hoarding was done for improving further prediction. For this we need to check which parameters or combinations of parameters of the user model and/or domain knowledge have bigger impact on the goodness of the algorithm.

By analyzing the goodness of the prediction of the hoarding algorithm we can try to tune its work. For example if a user indicates a LO miss as fatal the algorithm should check why this LO was not cached, e.g. if this entry was pruned or was given a small priority, and later the 'rules' for pruning and/or prioritizing should be reconsidered accordingly.

## 4. Learning sequences and cold start problem

Though it should be possible to extract specific knowledge about the user behavior and to try to predict students' future steps, on the first user access to the system it (the system) is totally unaware of the properties and preferences of this specific user.

The problem, known as 'cold start', can be faced by assessing the learner knowledge through a quiz and/or questionnaire and making some assumptions.
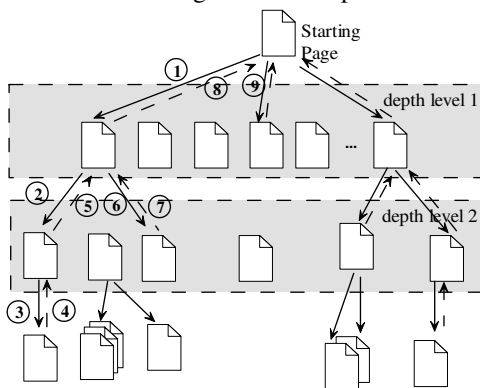


*Figure 3: Browsing path over a web-based material*

Basically the user browsing path over a web-based material (and in particular on any web-based learning source) can be viewed as a hierarchy structure (tree or directed graph). The user follows the links (the edges on the picture) from one page to another, or can go back to a previously viewed page (see fig. 3). Thus based on the knowledge about the learning material structure, the system can be aware of the most possible starting point of the student and suppose that the user will be following the links in the pages in consecutive order.

Also based on the observations on all previous users the system can estimate the average depth, in which the students browse during their first session.

In the context of caching the content on the first user access the system should upload as much as possible data trying to satisfy all user's requests.
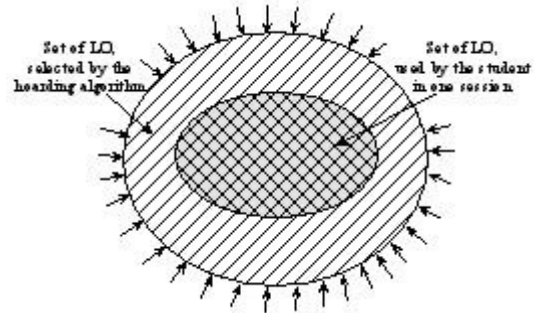


*Figure 4: The hoarding starting step*

Later the system can try to detect user's expertise level on a topic (by questionnaire for example) and narrow the hoarding set using some domain knowledge (e.g. if $LO_i$ should be proposed to advanced users, while current user is a beginner, the algorithm should exclude $LO_i$ from the hoarding set). When no other rules can be applied to decrease the size of the hoarding set the LO left might be randomly uploaded until the memory limit is reached.

## 5. Defining Session Length

In the Internet world a session is defined as "a continuous period of time during which a user's browser is viewing Web pages or a Web application within the same server or domain" (source - MSDN Library). It is series of transactions or clicks on web pages links made by a single user. There are different criteria to decide if a session is over or not. Two of the most commonly used are session length and the inactivity period of the user. For the first method the time limit for the session length is set to certain value and the activities later than this limit (counting from the start of the session) are considered a new session. In the second method if user activity stops for certain period of time on the resumption of the activity by the same user new session is considered started.

On the other hand for hoarding in mobile system the importance falls on the time between two possibilities to synchronize with the main server. In this sense we find more useful in this context to define a session as *the time between two synchronizations* of the mobile device with the main online system. The default session length might be one day, as commonly synchronization is done once per day, but during the system usage other session length might be observed and explicitly set for every user.

## 6. Links and correlation between LO

As mentioned earlier one of the steps of the hoarding algorithm is to construct 'candidate' set of objects, related (linked) to the starting point or to other objects that were predicted to be viewed. When using a web-based material the user clicks on the links of one page to go to another one and can either continue to browse further or can go

back to a previously viewed page. The links between the pages give us the structure of the web site (a learning material in particular), thus we can extract the relation between the LO, for example by parsing the pages.
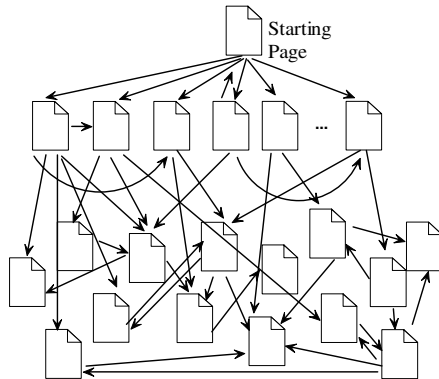


*Figure 5: Web-based material structure*

The links might be either bi-directional or not. We can build a LO correlation table in the following way:

```
for (every LO) {
 create a row;
 for (i=1, number_of_LO, i++) {
  if current_LO contains link to LO_i
    set cell_i = 1;
  else set cell_i=0;
 }
}
```

In the table below we can see that $LO_1$ contains link to $LO_2$ and to $LO_n$, but not to $LO_3$. The link is bi-directional for $LO_2$ and $LO_3$. In this way we can easily observe the set of LO that the user will be possibly requesting if he/she decides to browse deeper in the site, i.e. to go one level of depth further. Those are the objects directly linked to a particular object.

|          | $LO_1$ | $LO_2$ | $LO_3$ | ... | $LO_n$ |
|----------|--------|--------|--------|-----|--------|
| $LO_1$   | x      | 1      | 0      |     | 1      |
| $LO_2$   | 0      | x      | 1      |     | 1      |
| $LO_3$   | 1      | 1      | x      |     | 0      |
| ...      |        |        |        | x   |        |
| $LO_n$   | 1      | 0      | 1      |     | x      |

From the table above we can construct the 'candidate' set of LO for every next level of hoarding. Later this candidate set will be pruned (its size can be decreased by dropping some objects that are not likely to be requested).

On the other hand we can analyze the correlation between the objects, based on their concomitant usage in other user sessions. For example association rules can be discovered over all users' sessions containing an upper level LO. We can take into account only 'very strong' connections, i.e. associations discovered with confidence near to 1 and big support value. Note that it is expected that not a lot of such associations will be found, as the common scenario is to have big variety of LO and also big diversity of students' knowledge, interests and learning preferences. The rules extracted in this way will be of the type $LO_i \Rightarrow LO_j$ : conf=0.99 sup>0.5 which we can read as "Almost every time when the $LO_i$ was viewed also $LO_j$

was viewed in the same session (where $LO_i$ can be an example problem and $LO_j$ the solution given by the lecturer)".

Example:

|            | $LO_1$ | $LO_2$ | $LO_3$ | $LO_4$ | $LO_5$ | $LO_6$ |
|------------|--------|--------|--------|--------|--------|--------|
| $Session_1$ | 0      | 0      | 0      | 1      | 1      | 1      |
| $Session_2$ | 1      | 1      | 1      | 0      | 0      | 0      |
| $Session_3$ | 0      | 0      | 0      | 1      | 1      | 1      |
| $Session_4$ | 0      | 0      | 1      | 0      | 1      | 1      |
| $Session_5$ | 1      | 1      | 1      | 0      | 0      | 0      |
| $Session_6$ | 1      | 0      | 1      | 0      | 0      | 0      |
| $Session_7$ | 1      | 0      | 0      | 0      | 1      | 1      |

\* Where 1 means that $LO_i$ was viewed during $Session_j$ not taking care of the sequencing.

Association rules [1] algorithm discovers with confidence=1 the following relations:

$LO_1 \Rightarrow LO_6$ ; $LO_2 \Rightarrow LO_1$ ; $LO_2 \Rightarrow LO_3$ ;

$LO_4 \Rightarrow LO_5$ ; $LO_5 \Rightarrow LO_6$ ; $LO_6 \Rightarrow LO_5$ .

Association rules can be discovered also in more limited number of sessions (not all at a time), for example search for correlated objects only in the sessions of users that ware classified in the same group as the user for which the current hoarding set is being prepared. Considering the example above if we apply a clustering [1] algorithm (for example k-means), the algorithm produces 2 clusters from the above shown data. Applying association rules only to the sessions in the same cluster we get some additional associations. The clusters and discovered associations are as follows:

| Cluster    | Instances   | Associations          |
|------------|-------------|-----------------------|
| $cluster_0$ | $Session_1$ | $LO_1 \Rightarrow LO_5$ |
|            | $Session_3$ | $LO_3 \Rightarrow LO_5$ |
|            | $Session_4$ | $LO_3 \Rightarrow LO_6$ |
|            | $Session_7$ | $LO_4 \Rightarrow LO_6$ |
| $cluster_1$ | $Session_2$ | $LO_1 \Rightarrow LO_3$ |
|            | $Session_5$ | $LO_3 \Rightarrow LO_1$ |
|            | $Session_6$ |                       |

The above associations (like $LO_1 \Rightarrow LO_5$) show that if object $LO_1$ is to be selected for hoarding there is big probability that also $LO_5$ will be needed during the same session. Moreover associations of the type $LO_5=1$ & $LO_6=1 \Rightarrow LO_2=0$ can also be found, showing that if the user will be viewing $LO_5$ and $LO_6$ most probable $LO_2$ will not be viewed, thus can be excluded from the hoarding set or at least its priority can be set to lower level.
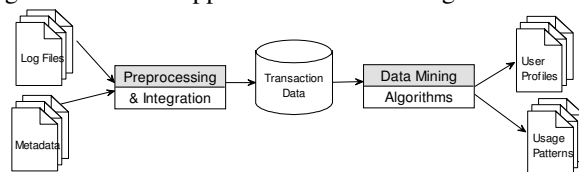
For the example above we considered only associations with confidence=1 and any support greater than 0, but in a real situations the best values for these parameters should be discovered experimentally.

The confidence value of the discovered associations LO can help in placing the items of the in an ordered list.

Also other data mining [1] and/or machine learning algorithms should be considered and tested to see their appropriateness for the hoarding process and how they can be combined best.

## 7. User modeling

There are different ways to model user behavior depending on the application and its needs. In the context of hoarding we recognize two groups of characteristics that the algorithm will use differently. Schematically call the first 'user behavior' and the second 'user knowledge'. Depending on the mobile learning system it is possible that not all the parameters can be discovered or they might be discovered through different techniques. The data about the user might be obtained by (any combination of) questionnaires, tests and quizzes or automatically by tracking the user and analyzing the log files. The process for retrieving automatically the user patterns consists of few steps, shown in the figure below. The first step is the preparation of the data for analyzing. For this step the log files should be preprocessed and integrated into a database. Afterwards different knowledge extraction algorithms can be applied to find interesting relations.



**Figure 6:** *Architecture for deriving user profiles*

*User behavior* can be described in terms of browsing styles (e.g. consecutive, random, interest driven); preferred type of educational media (e.g. prefers video to combination of text and pictures); etc. Based on it we can group the learners and analyze the similarities and differences between the groups and between the members of the same group. This should help to predict what will be needed, i.e. this data will be used to fill-in the hoarding set.

On the other hand the *user knowledge* profile should consist of things the system knows about what the user already knows, like the system awareness of the user's competence in certain subject (i.e. beginner, intermediate, advanced) or list of topics already covered previously. In contrast of *user behavior* the *user knowledge* profile will be user for pruning entries from the hoarding set, i.e. for excluding objects to decrease the size of the hoard.

We can distinguish static data about the user and dynamically changing data. The static data is for example the user gender, mother tongue, nationality, etc. Dynamic data is our current knowledge about the changeable over time user parameters and should be reviewed in certain periods of time. For example the user browsing pattern might change drastically few days before an exam date, thus the hoarding system should be able to quickly recognize such changes and react accordingly.

## 8. Metadata | Domain Knowledge

The metadata represents specific domain knowledge or knowledge about the specific learning material. Metadata is in general provided by the educator or the learning material creator. It will generally vary from one application to another, but can be used by the hoarding algorithm to improve it work.

One direction is to help in solving the 'cold start' problem by providing specific knowledge about the learning material structure, like 'initial point', provisioned common learning path, or connections between individual LO. The relationships between LO can also influence different weights of the parameters that are forming the priorities of the LO while hoarding.

## 9. Conclusions

In this paper we have described the hoarding problem for a mobile user without connection. The problem is how to support work on a mobile device when it is impossible to load on the mobile device all the data that comprise the full knowledge.

We have outlined a general algorithm, and we have posed a number of questions that need to be answered in order to solve the problem. We have also attempted to give some first answers to these questions. Our work is still in progress.

## 10. References

[1] Hand at al., "Principles of Data Mining", *Massachusetts Institute of Technologies, 2001, ISBN 0-262-08290-X*

[2] Kuenning, Popek, "Automated Hoarding for Mobile Computers", *Proc. of 16th ACM Symposium on Operating Systems Principles*, St. Malo, France, Oct. 1997.

[3] Trifonova, Ronchetti, "A General Architecture to Support Mobility in Learning", *Proc. of the 4th IEEE ICALT 2004*, Aug. 30 – Sept. 1, 2004, Joensuu, Finland.

[4] Trifonova, Ronchetti, "Where is Mobile Learning Going?", *Proc. of E-Learn 2003*, Phoenix, AZ, USA, Nov. 7-11, 2003.

[5] Trifonova, "User Behavior Observations for Offline Delivering of Materials in M-learning", *Tech. Report DIT-04-104*