# A New Cache Approach Based on Graph for Web Servers

Euloge EDI, Denis TRYSTRAM, JeanMarc VINCENT

[1] *Laboratoire Informatique et Distribution,ID-IMAG*
*51, avenue Jean KUNTZMANN,*
*38330 Montbonnot Saint-Martin, FRANCE*
*{Euloge.Edi, Denis.Trystram, Jean-Marc.Vincent}@imag.fr,*
*home page:http://www-id.imag.fr*

## Abstract

*This article deals with the problem of data cache for dynamic web servers. A set of requests (which number could be large) is submitted to a server for computing. Every request is modelled by a series-parallel graph of elementary tasks. These tasks are computed and their results are submitted to be stored into a cache. Due to the limitation of the cache capacity, some strategies should be used to select the right ones. The aim of this paper is to determine, in a static session, what tasks results have to be cached in order to reduce the computing time. The problem have been formalized using graphs representation whose complexity is briefly shown to be NP-hard. Then, heuristics are proposed to select the cached data. Finally some preliminary experimental results are discussed.*

*Keywords— dynamic request, web cache, tasks graph, utility, task output.*

## 1. Introduction

The always increasing performance of the standard PCs power for solving hard problem makes it possible to address new applications which require huge computing power and storage capacity. Thus, web applications available on Internet are more and more complex. Despite networks improvements (Myrinet, gigabit Ethernet), the growing number of customers and the heaviness of applications create new problems due to limited network bandwidth and also to server workload. From the users point of view, it induces an increasing response time that is crucial to be reduced. Web caches have already been proposed as a very promizing way for improving the Internet quality of service. They reduce both the user response time and the network congestion. In spite of their low costs (compared to other techniques to gain performance on the web), their application are mainly related to static request servers [12, 9]. The previous work in [11] is an attempt to define effective cooling policies taking into account various parameters of caches. The web-cache problem is not restricted for static requests. It has been shown, in [7, 13], that the most penalizing problem for these systems is more the computing time of the response than its transfer time from the server to the clients. Putting requests data solutions into caches reduces the response generation time to the time for searching the results into the memory. Thus, it is difficult to manage the large number of data generated by dynamic Web servers. We consider a generic architecture which consists in a server according to the model {frontend, caches, computing unit} [3]. Each element of the model can be an association of standard computers (PC cluster). In this case, the requests will be parallelized and executed on a set of available processors. The originality of this work is to model the requests as a series of task graphs and to use cache techniques for reducing the user response time by avoiding redundant computations. The cached elements are the results of some elementary tasks of the parallelized requests. This approach is applied to a web server of user parametrized geographic maps.

In section 2, the modelization of the problem is presented for the static case. We obtain a Linear Program similar to a knapsack problem [14], showing that the problem is NP-hard. In the third section, an algorithm filling the cache is described whose objective is to decrease the average user response time. The last section is devoted to some experimental results of the cartographic maps web server,

## 2. Modelization

### 2.1. Request and session

An user request is sent to the server from a client interface by specifying reals parameters of execution. For example a typical request could be :

**R :** Map on **Europe15**, **disparities** between closed **cities** taking **birthrate** as statistical resource and **value 100 as reference threshold**.

We consider that this request arrives for the first time on the server. It is then analysed and transformed into a parallel program which will be executed on available processors on a PCs cluster. This program is modelled by a graph of elementary tasks whose nodes correspond to computations on the initial data and edges are constraints due to precedences induced by the data flow graph. Precedence constraints represent data accesses and thus could be expensive if two successive tasks are mapped on distinct processors [10]. The graph, associated to a request, $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is connected, directed and acyclic. Its parameters are denoted by :

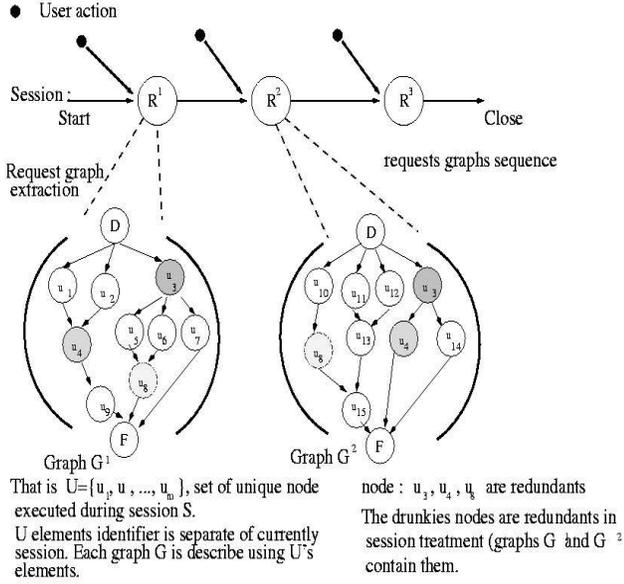$\mathcal{V} = \{v_1, v_2, \cdots\}$ *the finite set of nodes ;*

Graph $G^1$

Graph $G^2$

That is U={$u_1$, $u_2$, ..., $u_m$}, set of unique node executed during session S.
U elements identifier is separate of currently session. Each graph G is describe using U's elements.

node : $u_3$, $u_4$, $u_8$ are redundants

The drunkies nodes are redundants in session treatment (graphs G $^1$ and G $^2$ contain them.

Fig. 1. **Abstraction of user session graphs.**

$\mathcal{E}$ the set of edges ;

each node, a task of the program, corresponds to a procedure call with its parameters. Quantitative parameters for the task execution are :

  $t(v_j)$ : the processing time ;

  $s(v_j)$ : the size of the result $d(v_j)$ of the procedure.

A user session is defined by a series of requests sent by the same user during one connection to the server (figure 1). In a session, a new request is sent to the server after reception of the result of the previous one. In a user session, because successive computations operate on the same set of data, some routines could be executed several time. For example, a user may ask for the previous request **R** and after having received the result decides to request the same computation with a threshold of 80 instead of 100. Of course, the computation is very similar to the previous one storing intermediate results from **R** usually improves highly the new response time.

Let $\mathcal{U}$ be the set of all computation tasks of a session, $\mathcal{U} = \{u_1, u_2, \cdots, u_m\}$ is a set of couples (procedure-name, list-of-parameters). Usually $\mathcal{U}$ is very large (typically session is composed by seventy requests having four routines). A result $d(u_i)$, a processing time $t(u_i)$ and a size $s(u_i)$ are associated to each task $u_i$. It is supposed that these quantities are independent from each others. Then a session can be modelled by a sequence of graphs, $\mathcal{S} = \{G^1, \cdots, G^N\}$, whose nodes are tasks in $\mathcal{U}$. A task that appear several times in $\mathcal{S}$ is called redundant. For each task $u_j \in \mathcal{U}$ the occurrence number of task $u_j$ in the session $\mathcal{S}$ is denoted by $n(u_j)$ ($n(u) = 1$ if task $u$ appears only once).

A cache $\mathcal{C}$ is a memory device that could store data on demand. It is characterized by its capacity $K$ and its policy for data management. For tasks whose occurrence number is greater than one, it could be of interest to store its result in $\mathcal{C}$. Of course the data management policy should take into account the processing time of the task, the size of the result, or whatever parameter.

## 2.2. Optimization criteria

A session $(R^1, R^2, \cdots, R^N)$ of requests is associated to the graph sequence $\mathcal{S} = \{G^1, \cdots, G^N\}$, the set of tasks $\mathcal{U} = \{u_1, u_2, \ldots, u_m\}$ and also $n = \{n(u_1), n(2), \ldots, n(u_m)\}$ , the occurrence vector of elements of $\mathcal{U}$ in $\mathcal{S}$.

**Decision variable,** $\alpha$ is a binary vector of cardinality $m$ defining which results of elements of $\mathcal{U}$ are put in the cache.

$$\alpha(u) = \left\{ \begin{array}{l} 1 \text{ if the result of task } u \text{ is cached ;} \\ 0 \text{ if not.} \end{array} \right.$$

Two main optimization criteria are of interest :

**The server overload,** $W$, is the total computation time of redundant tasks during one session. If a task result is cached, the task is computed only once. If not, the task is computed each time. So,

$$W(\alpha) = \sum_{u \in \mathcal{U}} (n(u) - 1)(1 - \alpha(u))t(u).$$

**User response time,** $\mathcal{T}(\alpha)$, the critical path length of a session, is the sum of the critical paths of all the subgraphs representing the requests in the session. Denote by $CC^i(\alpha)$ the critical path (set of nodes) of the $i^{th}$ request graph given cache content. We have :

$$\mathcal{T}(\alpha) = \sum_{u \in \mathcal{U}} \left( \sum_i \mathbb{1}_{CC^i(\alpha)}(u) - \mathbb{1}_{non-first}(u) \right) (1 - \alpha(u))t(u),$$

with

$$\mathbb{1}_{non-first}(u) = \left\{ \begin{array}{l} 1 \text{ if } u \text{ is on the critical path} \\ \quad \text{but its first occurrence} \\ \quad \text{is not on the critical path;} \\ 0 \text{ if not.} \end{array} \right.$$

**Execution hypothesis :**
1. graphs representing requests are series-parallel (SP-graph)[5] ;
2. widths of generated graphs are lower than the number of computation nodes of the cluster. Processing time of a request is equal to the critical path length of its generated graph (and thus can be computed in polynomial time);

3. the total sequence of the requests is known before the execution ;

4. a location in the cache is written only once during the session : at the end of the first execution of the corresponding task.

The first hypothesis is related to the graph structure. In our target application, the graph structure appears via recursive calls and then are SP-graphs. The second hypothesis is based on the capability of the PC-cluster. During the parallelization, the maximum number of parallel tasks is controlled to be less than processor number. The third hypothesis consider that we are in a very optimistic case. Where all knowledge is available at the beginning. The last hypothesis is the restriction of a dynamic policy to a static case. The requests are known in advance, so the tasks to be excuted are known too. We can then put them in cache after computed them for the first time. We get then an upper bound for the cache performance.

**Constraints :** The set of constraints $(C)$ is defined by :

$$(C) = \begin{cases} \sum_{u \in \mathcal{U}} \alpha(u)s(u) \leq K & \text{(capacity)} \\ \alpha(u_j) \in \{0, 1\} & \text{(integrity)} \end{cases}$$

**Objective : O1** minimize the server overload $W(\alpha)$ under constraints $(C)$.

It is easy to remark that problem O1 be identified as a variant of well-known knapsack problem [14] ; that is known to be a $\mathcal{NP}$-hard problem.

**Objective : O2** minimize the user mean response time for a session $\mathcal{T}(\alpha)$ under constraint $(C)$.

**Proposition 1** *Problem O2 is NP-hard.*

Problem O2 can be easily reduced from problem O1. Consider a particular instance of problem O2 that consists of a sequence of $N$ chains. As

$$T(\alpha) = W(\alpha) + \sum_{u \in \mathcal{U}} t(u),$$

minimizing $T(\alpha)$ needs to minimize $W(\alpha)$ which is solving problem O1 which is NP-hard.

As the processor number is sufficiently large, hypothesis (2), the critical path of the session is the hard constraint for user response time.

## 3. Minimizing the user response time

As seen previously, tasks on critical paths really constrained the execution time of the session. Computing the critical path of a graph is usually achieved in $\mathcal{O}(n^2 log(n))$ times, where $n$ is the number of nodes. In our case, of SP-graphs, it is straitforward to see that it can be reduced to $\mathcal{O}(n)$.

### 3.1. Selection algorithm

A greedy algorithm is designed according to a task selection policy. The basis of the heuristic that we propose is to choose which task results to put in the cache on the critical path of the session. So if there is only one task result to store, it should guaranty that the critical path length is reduced. This algorithm computes the critical path of the session after each entry in the cache, because the critical path does not necessarily remain the same.

---
**Algorithm 1** Greedy heuristic algorithm.
---
**for all** $u \in \mathcal{U}$ **do**
    $\alpha(u) = 0$; {Initialization of cache index}
    $AS = K$; {Initialization of cache size : available space for caching}
**end for**
**repeat**
    compute set AT of tasks $u \in \mathcal{U}$ such that :
        $\alpha(u) = 0$;
        $u$ is on the critical path given cache state $\alpha$;
        $u$ is redundant $(n(u) > 1)$;
        $s(u) \leq AS$;
    select $u \in AT$ according to a given criterion;
    $\alpha(u) = 1$;
    $AS \leftarrow AS - s(u)$;
**until** $AT = \emptyset$;
---

The algorithm tends to reduce the critical path length of the graph at each step of the loop. It results an execution time $T(\alpha)$ that is also the time to compute the session given $\alpha$.

The selection criterion are very important for the cache efficiency. We propose below a heuristic that allows the reduction of the response time. Calling $T(\alpha_{opt})$ the optimal execution time related to the size of the cache, we have the inequality

$$T(\alpha_{opt}) \leq T(\alpha) \leq T(0)$$

The selection criterion fixes a local priority for the storage of tasks results in the cache. So we have to build this selection function such as it achieves a "good" performance.

### 3.2. Tasks utility

As for the knapsack problem, a utility value $g(u)$ (usually called a profit) is associated to each task $u$ at time when the algorithm selects a task result to put it into the cache. Several indicators could be used as utility :

**Size** The size of task result is a static utility parameter, according to this utility, tasks on the current critical path

are selected in an increasing order of their result size.

$$g_1(u) = s(u).$$

This idea of using such a heuristic is comes from underlying hit-rate model. What is important is the number of cache hit. Storing small size results increases the global number of tasks results in the cache.

**Execution time** Here, the utility of a task is the economy computation time of the computed on the current critical path.

$$g_2(u) = \left( \sum_i \mathbb{1}_{CC^i(\alpha)}(u) - \mathbb{1}_{non-first}(u) \right) t(u).$$

Tasks on the current critical path are selected in a decreasing order of their contribution to the critical path length. In relation with objective O2, the tasks with a huge computation time can intuitively reduce the session computation time. This strategy should perform well, but could be limited by the cache capacity.

**Ratio time_size** The utility of a tasks is the ratio between the computation economy when caching it to the size of its result $s(u)$. Tasks on the current critical path are selected in a decreasing order of their time_size ratio.

$$g_3(u) = \frac{t(u)}{s(u)} \left( \sum_i \mathbb{1}_{CC^i(\alpha)}(u) - \mathbb{1}_{non-first}(u) \right).$$

In relation with heuristics for the knapsack, this strategy selects at each stage, the task result which achieves the best computation time per unit of space in the cache.

## 4. A parallel web server of dynamic maps

### 4.1. Server architecture and requests structure

We apply our heuristics to a geographical map Web server presented in figure(2). Europe is represented by the countries association. For example : Europe_15 correspond to the set of the $15^{th}$ first countries of Europe community. Each country is subdivided in elementary units representing the cities (NUTS4) of the given nations. The basics units are aggregated to create several cutting level (regions, departments, ...) and for every unit an inventory of various statistical resources (population, birthrate, ...) is available. This cutting can be modified by the user and deviations calculated compared to precise criteria. The server has a web interface on the web from which all requests are submitted. The goal of this server is to provide to users an effective and interactive tool for mapping demographic data sets on a defined geographical area [15]. Since input files are huge (example : 791 Mo for an inventory on 116000 units)[1], cache strategies

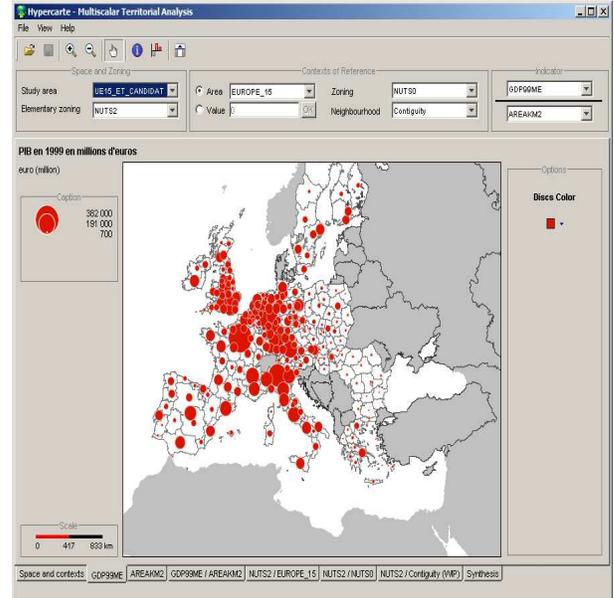[1] This size is given for European space subdivided in at a NUTS4 level.



Fig. 2. **Cartography server. Clients formulate requests via this interface.**

should be efficient. Future works will consider data set on all 25 states Europe with a statistical inventory on communal level (NUTS5) that is about one million of units. Five types of requests are possible from this interface.

### 4.2. Parallel programming environnement

ATHAPASCAN [6] has been used as parallel programming environment. In ATHAPASCAN environnement, the application designer develops his program by recursive procedure calls named tasks. Because procedure calls could be forked, the execution could be done in parallel. The forking process determines the task granularity. Given a scheduling strategy, the runtime environment schedules automatically those tasks on a parallel machine. The data flow graph is an instance of the precedence relation graph between all elementary tasks of the program. So, by construction, we get a series-parallel task graph.

### 4.3. Cache implementation

As cache environment, we use CaLi [8], a C++ framework elaborated to allow a fast implementation of distributed or local cache systems. The library follows the ideas of the *generic programming* paradigm, providing users with a high efficiency while keeping a great flexibility. CaLi is independent of the running application and management of the cache is specified separately by the user (type of items to cached, insert and remove policy). CaLi could be used with distributed caches. In this case it uses MPI to achieve strong coherence between the distributed caches. Our environnement

TABLE I
**Tasks results are put into cache in increasing order of their size.**

| Csize (b) | timeWoC(s) | timeWC (s) | teconomy | HitTx | Citem | CacheC |
|---|---|---|---|---|---|---|
| 1024 | 28.9 | 28.891 | 0.05 | 0.009 | 9 | 961 |
| 4096 | 28.9 | 28.890 | 0.06 | 0.03 | 31 | 3918 |
| 16384 | 28.9 | 28.82 | 0.27 | 0.092 | 94 | 15875 |
| 65536 | 28.9 | 28.49 | 1.45 | 0.22 | 225 | 63674 |
| 262144 | 28.9 | 26.81 | 7.26 | 0.33 | 343 | 253890 |
| 1048576 | 28.9 | 19.8 | 31.5 | 0.42 | 441 | 995034 |
| 4194304 | 28.9 | 15.49 | 46.4 | 0.451 | 464 | 1521173 |
| 16777216 | 28.9 | 15.49 | 46.4 | 0.451 | 464 | 1521173 |

TABLE II
**Tasks results are put into cache in decreasing order of their cumulative execution time on the critical path.**

| Csize (b) | timeWoC(s) | timeWC (s) | teconomy | HitTx | Citem | CacheC |
|---|---|---|---|---|---|---|
| 1024 | 27.26 | 27.23 | 0.11 | 0.19 | 2 | 970 |
| 4096 | 27.26 | 27.21 | 0.206 | 0.19 | 2 | 3931 |
| 16384 | 27.26 | 27 | 0.955 | 0.47 | 4 | 15996 |
| 65536 | 27.26 | 26.75 | 1.88 | 0.43 | 4 | 63957 |
| 262144 | 27.26 | 25.04 | 8.17 | 1.4683 | 14 | 256028 |
| 1048576 | 27.26 | 18.42 | 32.36 | 10.2 | 103 | 1003507 |
| 4194304 | 27.26 | 14.67 | 46.42 | 45.1 | 464 | 1521173 |
| 16777216 | 27.26 | 14.67 | 46.42 | 45.1 | 464 | 1521173 |

TABLE III
**Tasks results are put into cache in decreasing order of their ratio time-size.**

| Csize (b) | timeWoC(s) | timeWC (s) | teconomy | HitTx | Citem | CacheC |
|---|---|---|---|---|---|---|
| 1024 | 28.86 | 28.82 | 0.17 | 0.22 | 2. | 964 |
| 4096 | 28.86 | 28.73 | 0.496 | 0.66 | 6 | 3987 |
| 16384 | 28.86 | 28.43 | 1.52 | 1.0 | 10 | 15967 |
| 65536 | 28.86 | 27.41 | 5.08 | 2.13 | 21 | 63980 |
| 262144 | 28.86 | 24.92 | 13.7 | 6.49 | 65 | 256018 |
| 1048576 | 28.86 | 18.35 | 36.42 | 24 | 244 | 1003516 |
| 4194304 | 28.86 | 15.50 | 46.3 | 45.1 | 464 | 1521173 |
| 16777216 | 28.86 | 15.47 | 46.4 | 45.1 | 464 | 1521173 |

is designed for parallelized and distributed cache strategies. In this paper, only centralized cache have been considered.

## 4.4. Experiments

The goal is to integrate this cache management into an actual web portal. Experiments have been done on real case to evaluate the efficiency of these three policies on the running of the cartography server. We run the server on a symmetric biprocessor(AMD Athlon(tm), 1194 MHz, 256 KB cache size). For the experiment we process different sessions. A session is a sequence of randomly chosen requests among the set of all possible requests. A session is automatically transformed in a task graph. Session containing 17 requests have been processed and each of them generates over than 1000 elementary tasks. A single task is computed by one of the smp node. The size of the tasks results are between 218 bytes and 31 Kbytes. Only tasks with significant computation time (larger than $10^{-2}$s) comparing to data manipulation time by the cache are submitted to the cache. For cache size between 1024 bytes and 16 Mbytes, one processes 50 random sessions for each task utility policy. For each session one notes the value of the different statistical performance indicator.

In the following tables, *Csize* is the cache size, *timeWoC* is the execution time without caching any tasks, *HitTx* is the hit rate, *timeWC* is the execution time with HitTx, *CacheC* is the total size of tasks result in the cache, and *Citem* the number of task result in the cache

**Size as utility** ($g_1$) (in table(I) ). This strategy based on the tasks result size allows a good hit rate. The computation time economy is not so good because the task results with a small size are not necessarily the most critical for the execution time. Wasted space in the cache is hight because only tasks results with big size are out of the cache and can't be packed in the residual space of the cache. Compared to the caching policy define in [18, 17, 16], this shows that it is not the number of hits which is important for this kind of application.

**The utility of a task result is critical execution time** ($g_2$) (in table(II) ). The number of tasks results cached is small, so the hit rate do not increase quickly. But those which are cached, affect the critical path length. As

our intuition, the computation time economy is good in relation with the content of the cache.

**The utility of a task result is the ratio time_size** ($g_3$) (in table(III) ). This strategy takes into account both the size and the computation economy. So it achieve both best hit rate and computation time economy. In our situation the wasted cache space is small.

From curves in figure 3, we operate the cache by choosing an adapted cache size. The best strategy is for the ratio time_size utility function. It achieves the best use of the available space in the cache. Curves of time and size strategies are very closed. The computation economy strategy is efficient but does not take into account the size of the cache. All these strategies are equivalent when the cache is large enough. Execution time is then the time to compute at least one time the different tasks. In our situation, a cache size of about 6Mb will achieve a good performance rate. The contrary conclusion about the consideration of the hit rate and the computation time economy comes from the fact that in this model, one put in cache the elementary tasks results witch can be a small part of one request computation. That will not be the same if we put in cache the request full solution because a hit will incurs a request solution. In the last case, the impact of the hit rate on the computation economy will be important. The computation economy strategy is efficient but it does not take into account the size of the cache. All these strategies are equivalent when the cache is large enough. The execution time is then the time to compute at least one time the different tasks.
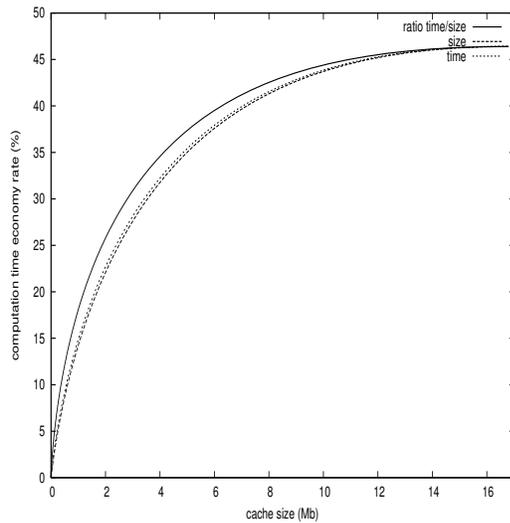
Fig. 3. **Comparing the three strategies.**

## 5. Concluding remarks

In this paper, an algorithm to improve the user response time for caching dynamics requests have been proposed. This algorithm is based on a task weighting heuristic for known session (which means that the serie of requests to be processed is known in advance). We propose here a new way of optimizing the scheduling by avoiding redundant computation. These algorithms for static session have been tested for several strategies. The ratio time_size strategy appears to be significantly better than time based or size based strategies. Complementary studies should be done on various scenarios to inforce this result. The software framework developped in this paper could easily be adapted for this purpose. The task graph model could be completed by annotations of precedence relations. This could be used to make the computation of the makespan more accurate and so data transfers inside the cluster could be optimized. Simultaneously, cache strategy could take in account the cluster node which process the elementary tasks of the graph. We are currently working on this improvement. The continuation of this work concerns the application of this technique for dynamic sessions (on-line caching). It consists in extending the utility heuristic for dynamic sessions, evaluate it on-line for a real case. The cartography server has been designed in order to achieve such a strategy.

## References

[1] M.D. HAMILTON, P. McKEE, I. MITRANI : Optimal caching policies for Web objects. In Proceedings of the 9th International Conference High-Performance Computing and Networking (HPCN Europe 2001), Amsterdam, The Netherlands, 25-27 June 2001 Hertzberger, L.O., Hoekstra, A.G. and Williams, R. (eds.) Lecture Notes in Computer Science,vol 2110, pp. 94-103

[2] S. PODLIPNIG, L. BÖSZÖRMENYI : A survey of web cache replacement strategies. ACM computing survey 2003, vol 35, n.4, pp. 374–398

[3] E. EDI, D. TRYSTRAM, JM. VINCENT : Amélioration de performance de serveur Web de requêtes dynamiques. Colloque Africain sur la Recherche en Informatique, CARI'02 (2002), pp. 121–126

[4] A. JAKOBY, M. LISKIEWICZ, R. REISCHUK : Space Efficient Algorithms for Directed Series-parallels Graphs. In Electronic Colloquium on Computational Complexity 2002, Report vol 21

[5] C.W. HO, S.Y. HSIEH, G.H. CHEN : Parallel Decomposition of generalized Series-parallels Graphs. In Journal of Information Science and Engineering 1999, vol 15, pp. 407–417

[6] M. DOREILLE : Thèse : Athapascan-1 : Vers un Modèle de Programmation Parallèle Adapté au Calcul Scientifique. Institut National Polytechnique de Grenoble 1999

[7] H. ZHU, B. SMITH, T. YANG : Hierarchical Resource Management for Web Server Clusters with Dynamic Content. Proceedings of the 1999 ACM SIGMETRICS international conference on Measurement and modeling of computer systems, Atlanta, Georgia, United States, pp. 198–199

[8] http://icis.pcz.pl/∼zola/CaLi CaLi – Generic computational buffers library2004.

[9] E.P. MARKATOS : Main memory caching of Web document. Fifth International World Wide Web Conference,May 6-10, 1996, Paris, France

[10] Y.K. KWOK, I. AHMAD : Dynamic Critical-Path Scheduling: An Effective Technique for Allocating Task Graphs to Multiprocessors. IEEE Transactions on Parallel and Distributed Systems 1996, vol 7, pp. 506–521

[11] L. CHERKASOVA, G. CIARDO : Role of Aging, Frequency and Size in a Web Cache Replacement Policies. Proceedings of the 9th International Conference on High-Performance Computing and Networking 2001, LNCS 2110, pp. 114–123

[12] M. COLAJANI, P.S. YU, M. DIAS : Analysis of Task Assignment Policies in Scalable Distributed Web-Server Systems. IEEE Transaction Parallel Distributed System 1998, vol 9, n.6, pp.585–600

[13] B. SMITH, A. ACHARYA, T. YANG : Exploiting result equivalence in caching dynamic Web content. Proceedings of USENIX Symposium on Internet Technologies and Systems, Boulder, Colorado, USA, October 1999.

[14] H. KELLERER, U. PFERSCHY, D. PISINGER : Knapsack problems. Springer-Verlag 2004

[15] H. Mathian, C. Grasland, J-M. Vincent : Multiscalar analysis and map generalization of discrete social phenomena Problems and Political Consequences Statistical Journal of the U.N. Economics Commission for Europe 2000, vol 17, n.2,pp. 157–188

[16] G. PIERRE : Conception d'un Système de Cache Adapté aux Spécificités des Utilisateurs Proceedings of the 1998 NoTeRe colloquium, octobre 1998

[17] S. DAVID, T. DAVID Probabilistic Methods for Web Caching Performance Evaluation, Elsevier Science Publishers B. V. 2001, vol 46, n.2-3,

[18] C. PERRIN : Gestion de cache web sur un serveur en grappe INRIA, projets APACHE et SIRAC 2000