# Autonomous Information Access to Assure Service Level Agreement in Faded Information Field[*]

**Zhensu Lu[1] and Xiaodong Lu[2]**

[1]Sanya University, Sanya, Hainan, China
[2]Electronic Navigation Research Institute, Tokyo, Japan
luzs@lzu.edu.cn,      luxd@enri.go.jp

**Abstract -** The user differentiation is required to attain heterogeneous service levels in such terms as accessible information volume, response time, or system availability. To fulfill requirement of user differentiation to attain heterogeneous service levels, based on the Faded Information Field (FIF) system architecture, autonomous real-time navigation technology is proposed by navigating users' mobile agents to low-congestion nodes to assure the maximum response time. As a result, the requests are able to avoid local congestion zones without a centralized dispatcher or workload manager, leading to globally accomplishing the heterogeneous Service Level Agreements of each user even in case of subsystem failure and rapidly changing environment. The effectiveness of the proposed technology has been proved through simulation, and the results show that the proposed technology increases an average of 60% the service satisfaction ratio.

**Index Terms** - FIF, SLA, push/pull mobile agent (Push/Pull-MA), Real-Time Navigation

## 1. Introduction

With the advances in communication technologies and the decreasing costs of computers, information provision and access by providers and users have certain constraints. For example, to improve high response time for users may violate service providers' timeliness for information update. Therefore, Service Level Agreement (SLA) has been widely deployed between information service providers and infrastructure providers [1]. It allows the provider to measure its own performance and improve itself over time by identifying and defining the customer's needs; resulting in distinguished business from the other competitors. As a customer, one may select an appropriate service provider by comparing SLA of different service providers.

A prominent research on Service Level Management (SLM) is being carried by IBM with their Autonomic Computing [2] and Grid [3] initiatives. However, in these systems the workload management is based on a centralized component, subject to fault which would drive to the fatal consequence of bringing the whole system down. Sun Microsystems is carrying research on SLA for Data Centers [4][5], and proposes an architecture based on a separated management network. Their work is focused on data monitoring and metrics, but no discussion on the management network reliability is done.

The main objective of conventional systems was to assure heterogeneous needs for average response time. In this paper, the problem of assuring the maximum response time is addressed, and a decentralized resource management technology to assure pull mobile agent real-time navigation based on Faded Information Field (FIF) system architecture is proposed. From the cooperation between the users' mobile agents and the nodes in the system, the navigation technology eludes congestion zones to assure real-time property under unbalanced load situations. The technology is founded at the system structure level, without a centralized component, that assures the availability and real-time of the system under the assumption of heterogeneous service levels.

The structure of the paper is organized as follows. In the next section, the system architecture of FIF and node structure with respect to proposed solution is presented. In section 3, information access control and the activities of pull mobile agent in real-time navigation are described. The simulation results in section 4 show the improvement and effectiveness of proposed technology. Finally, section 5 draws the conclusion.

## 2. System Architecture

### A. Faded Information Field

The main goal of the FIF is to guarantee the assurance of autonomous information service provision and utilization [6]. Service providers (SP) trace the demand trend for information, and allocate to accepting storing nodes the most accessed segment of their information services. The storing nodes then, in a recursive pruning process, further allocate the information services to adjacent nodes. As a result, the multi-level distributed information services area is created. Users with different requirements for information can be satisfied at different levels in the FIF. Consequently, the cost of service utilization (access time) and cost of service provision (update) are balanced by allocating closer to the majority of the users the most accessed part of the information services.

In the FIF system, information contents are uniquely defined by Content Code (CC). The information contents are further specified by its Characteristic Codes (CHs). Node, Push-MA and Pull-MA are three autonomous subsystems, mainly responsible for information storing, allocation and utilization, respectively.

### B. Node Structure

The users are classified in categories according to their

---

requested service level. For example, service provider might want to limit access to privileged information or prioritize some users' requests. The concept of Pull-MA is extended to allow multiple categories of users, such as Gold-MA, Silver-MA and so on. For each category, there might exist different Service Level Objectives (SLO), that each mobile agent will manage by itself in order to satisfy the users' requirements for response time and information volume consumption. To support the technology presented in this paper, two more fields are added to the message format. The field SL describes the service level that the user expects from the system while field DL specifies deadline to be satisfied. This field is introduced to support satisfying the real-time requirement. The initial deadline depends on the service level. This field is being updated by the node as the Pull-MA passes through. Negotiation between the Pull-MA and the node is done to determine if the request can be satisfied in terms of remaining deadline time and processing time of the node.
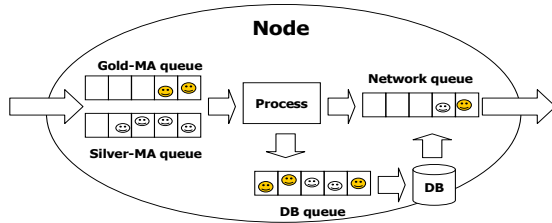


Fig. 1   Node structure

The node structure is presented in Fig.1. There can be seen the main modules, and each one possess its own queues. The main modules are the mobile agent execution environment, the database (DB) that contains the information, and the network device. Each node monitors the Pull-MA activity. More precisely, monitors for how long the agents spend time inside the node, and keeps record of the maximum time in node within a monitoring period $\Delta T_m$. It is measured for both cases of DB access and forward case. The Maximum DB Access Time (MAT) and Maximum Forward Time (MFT) are calculated for all user categories.

At the first time of being queued, the Pull-MAs are stamped with the current time $T_{in}$. When leaving the node at the last stage of the process, the total time in the node is calculated by subtracting the current time in the node $T_{out}$ to the time stamped in the Pull-MA, also taking into account the transmission time $T_{trans}$ That is, $\Delta T_a = T_{out} - T_{in} + T_{trans}$ stands for the time spend in node of a specified agent $a$. Then, the MAT is defined as the maximum value $\Delta T_a$ for all agents that accessed the DB within the monitoring period $\Delta T_m$. MFT is defined equally for the agents that did not access the DB and were forwarded instead.

## 3.  Autonomous Real-time Navigation

In FIF system, whenever edge nodes are accessed by client through Pull-MA, a situation may arise where local peak of requests congests partially the network, while other resources are spare. Fair distribution of requests may be managed by centralized dispatcher, but it may become single point of failure. Moreover, it limits the throughput of requests which is quite inappropriate for the high assurance system. Looking at the life cycle of Pull-MA, different activities from its initiation till termination are described as follows.

### A.   Pull-MA Execution

Pull-MA incorporates the request of user with respective service level requirements. Priority scheduling at the initial agent processor ensures the different access time requirements. The agents will be scheduled in separated queues depending on their level or category. The scheduler will elect an agent from the Golden-MA queue with higher probability than from the Silver-MA queue. The ratio of electing from the Golden-MA queue can change dynamically to assure timeliness of Golden-MA.

The Pull-MA requests are queued for execution depending on the service level. There are as many queues as user categories, and the requests are scheduled by priority. Each queue is a plain FIFO independent of the other. The more priority, more probability to be scheduled that queue. There can be identified two choices broadly. Forward the agent or satisfy the request, granting access to the database. In case of forwarding, the next node is determined based on the navigation algorithm explained in section 3.C. After that, the Pull-MA is queued in the network in order to forward to other node, or send the requested information back to the user.

### B.   Node Selection

Choosing the candidate nodes to migrate is one of the steps in the Pull-MA processing. There are two elements to be considered. One is the information detail required by the Pull-MA related to the information available in the node. The other is the MAT of the nodes in the locality relative to the agent deadline. At first, the information detail (CH) is considered. If the CHs required by the agent are greater than the CHs available, then the request can only be satisfied at upper nodes. The candidates are all the neighbours at the upper layer. The *CheckMAT(candidates)* operation selects the candidate nodes based on the following conditions:

$$DL - MFT > next.MAT$$
$$node.MAT > next.MAT$$

The first condition assess if the deadline (DL) will be accomplished in the next node, taking into account the transmission time. The second condition is meant for driving the requests toward less congested nodes. The situation may arise in which the agent is forwarded to a node with more information than required, but that node cannot satisfy the deadline. In such case, the possible candidates can include nodes at both lower and upper layers.

Next, the Pull-MAs granted DB access is scheduled at the DB queue. In that queue a time-priority, category-independent scheduling is used; Specifically, Earlier Deadline First (EDF) scheduling is performed. The EDF is a dynamic queue discipline widely used for real-time scheduling. The requests are reordered by remaining deadline time.

## C. Pull-MA Navigation

The Pull-MA requests can obtain information about MAT of neighbour nodes. However, if all requests decide to go to the lowest MAT node, and considering that the nodes in the same layer will forward the requests to likely that same free node, that will result on a burst of congestion. The requests that were looking for a low MAT value found themselves in a suddenly congested node, and unable to satisfy the deadline in the worst case. A way to solve this problem is to do random navigation. The Pull-MA will select the candidates to migrate to those nodes that the MAT is lower than the deadline. Then, just choose one randomly. While a uniform distribution for the random node selection can be done, the fact is that it assumes that 1) the users access the system also with uniform distribution, and 2) the connectivity between nodes is also fairly distributed; two assumptions that cannot be guaranteed in real situations. An intermediate solution is proposed, between best MAT and uniform random distribution navigation. Based on the MAT values of the candidates, the probability of being selected is calculated as follows: At first, the candidates are ordered by increasing MAT. Probability $p$ for node at position $i$ is defined as:

$$P_i = \frac{MAT_{N-i}}{\sum_{n=0}^{N-1} MAT_n}$$

Where $N$ is the number of candidates. As can be seen, the probabilities are assigned in inverse order respect MAT ordering. The result is that the node with less MAT will have more probability to be selected. The result of these schemes is that the Silver-MA requests, with relaxed time requirements, are usually satisfied as soon as they reach the required CH level. On the other hand, the Gold-MA requests might encounter the lower layers congested by silver requests. Following the algorithms, they will navigate to upper layers reaching non-congested nodes, which can serve within the real-time requirements of gold requests. Assuming a request for information volume $v$, which navigates through $n$ nodes, and network transmission time $T_{trans}$. The total service time can be approximated as:

$$T = \sum_{i=1}^{n} (T_{proc} + T_{net} + T_{trans}) + T_{db}$$

## D. MAT Propagation

Periodically, each node sends a multicast message to the neighbouring nodes with its MAT value. Following this procedure, each node knows the current status of the neighbours. This information is vital for the Pull-MA navigation, as the time that will be required to process the request must not be greater than the request deadline. The messaging is asynchronous, so the timing for receiving and even sending information is different for each node. A key decision is the propagation period ($T_p$). If $T_p$ is too small, there might be a significant impact on the performance of the node. If it is too big, the information being reported might be too old, leading to failure in adapting to changes of the workload. The value of $T_p$ is to be determined experimentally, as it depends on the expected workload and servers' performance.

The MAT propagation is done within the locality of inter-connected nodes, so the number of messages being exchanged is small. This allows that $T_p$ can be small enough to report the current status of the node without degrading the node performance. Moreover, the system can be scaled to any number of nodes without compromising the accuracy.

## 4. Performance Evaluation

The objective of this simulation is to evaluate the improvement of satisfaction ratio for gold users under a local congestion of the system and in coexistence with a large number of silver users. For comparison, the conventional non-forwarding approach is adopted. The simulation results show the quantitative improvement of proposed navigation technology. The total satisfaction ratio is evaluated for navigation and non-forwarding cases.

## A. System Model

TABLE I   Parameters of simulation

| Parameter | Value |
|---|---|
| Agent forwarding | 1 ms |
| Agent information provision | 50 ms |
| MA Processing ratio δ | 100 MAs/sec |
| Gold/Silver MA scheduling | 75%, 25% |
| Gold-MA deadline | 70 ms |
| Silver-MA deadline | 2 sec |
| MAT monitoring / forwarding | every 1 sec |

The proposed technology is evaluated by considering a simple information environment, with 4 levels of information, and 17 nodes, that are assumed to be uniform in processing. The simulations consider the following common parameters, shown in TABLE I. As can be seen, the deadline for golden requests is very tight. The silver deadline is comparatively large, but is a delay that can be accepted by most of the non-demanding users. The system is evaluated under an increasing number of requests, and also under different information volume discrimination patterns.

## B. User Model

The preference of the users for the information levels is shown in TABLE II. The Silver/Gold information difference evaluated between 0 and 3, so there are 4 patterns of silver users. Diff 0 means that the silver can access all the information, just like golden-ma. Diff 3 is the case that silver-ma can only access CH1 information.

TABLE II   The users' preference

| User level | Ratio | CH1 | CH2 | CH3 | CH4 |
|---|---|---|---|---|---|
| Silver diff 0 | 75% | 47% | 30% | 17% | 6% |
| Silver diff 1 | 75% | 50% | 31% | 19% | - |
| Silver diff 2 | 75% | 62% | 38% | - | - |
| Silver diff 3 | 75% | 100% | - | - | - |
| Gold | 25% | 47% | 30% | 17% | 6% |

The user distribution, i.e., the arrival rate at the edge nodes, follows a binomial distribution. Following the binomial distribution, there will be some nodes in a locality that are heavily loaded, but others with many remaining resources. Proposed technology works in such a transient situation, forwarding the requests from a congested locality towards areas of low congestion.
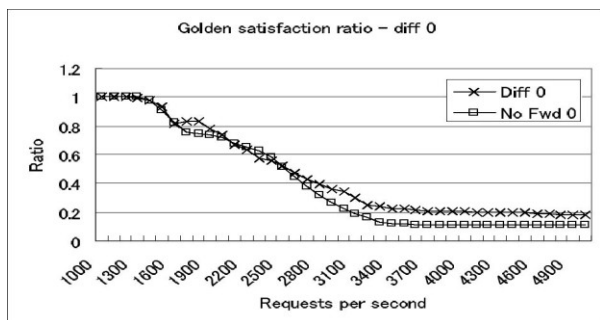
*C. Satisfaction Ratio*



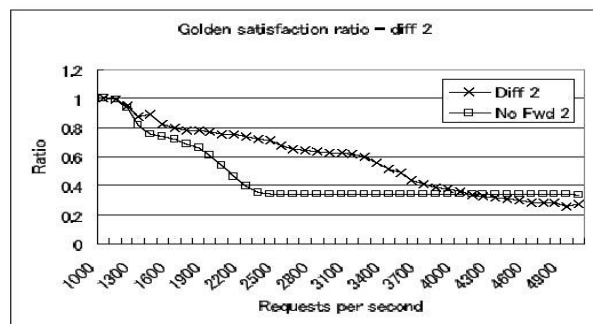Fig. 2   Golden satisfaction ratio-diff0



Fig. 3   Golden satisfaction ratio-diff2

The satisfaction ratio is defined as the percentage of agents that could send back to the user the required information within the deadline. We evaluated the forwarding technology compared with the greedy placement (requests are satisfied as soon as reaching the corresponding layer), for each of the users preference model shown in TABLE II. By comparing the results for each silver difference level, it can be observed that for no differentiation or one level of differentiation, there is only slight improvement in the satisfaction ratio. For difference 0 (Fig.2), it was observed that from around 3000 requests per second, the gold satisfaction ratio improves 10%. The satisfaction ratio of difference 1 is similar with difference 0, because the only change is that for difference 1, the requests cannot access the SP node, the top-most layer.

The major benefit of the forward technology can be seen for difference levels 2 (Fig.3) and 3. For difference 2, the satisfaction ratio is above 70% until 2500 requests per second (35% of improvement regarding no forwarding system). In the case of difference 3, the system is able to maintain more than 90% satisfaction ratio until 2500 requests per second (at most 35% of improvement), and more than 70% satisfaction until 3500 requests per second. It is noticeable that the satisfaction

ratio becomes worse than without the forwarding technology under heavy load. It can be thought that the cause is the overall overhead in the system due to the forwarding mechanism.

*D. Navigation Method Comparison*

Basically, the greedy node selection means that the Pull-MA will choose the next node with less MAT value. Uniform random method means that the node will dispatch the Pull-MAs to the upper nodes in the same probability without considering their MAT. Finally, the MAT random is proposed node selection technology, where the probability of choosing each node is proportional to the upper nodes' MAT value. No forward method means that the requests are processed as soon as the required information volume layer is reached. Forward means that the requests will be sent to other nodes looking for better response times. The evaluation measures the satisfaction ratio within the same parameters of the previous simulations. It can be seen in TABLE III that the proposed forward navigation can get better satisfaction ratio than conventional processes.

TABLE III   Navigation methods comparison

| Node selection | No forward | Forward |
|---|---|---|
| Greedy | 17.31% | 32.82% |
| Uniform random | 38.41% | 61.65% |
| MAT random | 42.53% | 67.80% |

## 5.  Conclusion

In this paper, the mechanisms that allow the access differentiation depending on the user's category are described. Furthermore, in a temporal local congestion, the technology to navigate towards low-congestion nodes in order to assure the maximum response time under dynamically changing situations is proposed. In the FIF architecture, service provision, utilization and maintenance are performed autonomously by agent entities that locally coordinate to achieve their own objectives. The evaluation results have revealed the effectiveness of the proposed technology to assure the service level in terms of satisfaction ratio, maximum access time and available information depending on the agreement between the service provider and the users.

## 6. References

[1]    A. Dan, H. Ludwig and G. Pacifici. "Web Service Differentiation with Service Level Agreements", IBM Corporation, white paper, (2003).
[2]    J.O. Kephard and D.M. Chess. "The vision of Autonomic Computing", *IEEE Computer*, vol.36, no.1, pp.41-50, (2003).
[3]    I. Foster, C. Kesselman and S. Tuecke. "The Anatomy of the Grid: Enabling Scalable Virtual Organizations", *International J. Supercomputer Applications*, vol.15, no.3, (2001).
[4]    E. Wustenhoff. "Service Level Agreement in the Data Center", SUN Microsystems, http://www.sun.com/, (2002).
[5]    E. Wustenhoff. "Service Level Management in the Data Center", SUN Microsystems, http://www.sun.com/, (2002).
[6]    X.D. Lu, K. Moriyama, I. Luque, M. Kanda, Y. Jiang, R. Takanuki and K. Mori. "Timeliness and Reliability Oriented Autonomous Network-Based Information Services Integration in Multi-Agent Systems", *IEICE Trans. Inf. & Syst.*, vol.E88-D, no.9, pp.2089-2097, (2005).