

Differentiated Quality of Service in Application Layer Active Networks.

Chris Roadknight and Ian W. Marshall

BT Adastral Park, Martlesham Heath, Ipswich, Suffolk, UK IP5 3RE

{christopher.roadknight, ian.w.marshall}@bt.com

Abstract.

A novel approach to quality of service control in an active service network (application layer active network) is described. The approach makes use of a distributed genetic algorithm based on the unique methods that bacteria use to transfer and share genetic material. We have used this algorithm in the design of a robust adaptive control system for the active nodes in an active service network. The system has been simulated and results show that it can offer clear differentiation of active services. The algorithm places the right software, at the right place, in the right proportions; allows different time dependencies to be satisfied and simple payment related increases in performance.

Keywords. Network Management, Genetic Algorithms, ALAN

Introduction.

To be popular with customers an active service platform must provide some clear service quality assurances. Users of an active service network supply the programs and policies required for their custom services in transport packets alongside their data. Clearly it should be possible for these users to specify the Quality of Service (QoS) using any metric that is important to them. The rate of loss of packets carrying service requests or policies, and the service response time (latency) are two obvious examples. In this paper we discuss the management of QoS in an Application Layer Active Network (ALAN) [1] that enables users to place software (application layer services) on servers embedded in the network. Despite the obvious virtual networking overheads, the resulting end to end service performance will often be significantly better than if the services executed in the user's end systems (as at present). For example, a network based conference gateway can be located so as to minimise the latency of the paths used in the conference, whereas an end system based gateway will usually be in a sub-optimal location.

For the purposes of this work we have assumed that the latency and loss associated with the network based servers is significantly greater than the latency and loss associated with the underlying network. In the case of latency this is clear - packet handling times in broadband routers are around ten microseconds, whilst the time taken to move a packet into the user space for application layer processing is a few milliseconds. In the case of loss the situation is less clear since currently servers do not drop requests, they simply time-out. However, measurement of DNS lookup [2] suggest DNS time-outs occur significantly more frequently than DNS packet losses, so we feel our assumption is reasonable.

In the next section we briefly describe our active services platform ALAN and its associated management system. We then justify our approach to QoS in an ALAN environment. We then describe a novel control algorithm, which can control QoS in the desired manner. Finally we show the results of some simulations using the novel algorithm. The results are very encouraging and illustrate for the first time that a distributed AI approach may be a productive QoS management tool in an active services network. However, further work is required before we can justify the use of our approach in a working active network.

ALAN

ALAN [1] is based on users supplying java based active code (proxylets) that runs on edge systems (dynamic proxy servers - DPS) provided by network operators. Messaging uses HTML/XML and is normally carried over HTTP. There are likely to be many DPSs at a physical network node. It is not the intention that the DPS is able to act as an active router. ALAN is primarily an active service architecture, and the discussion in this paper refers to the management of active programming of intermediate servers. Figure 1 shows a schematic of a possible DPS management architecture.

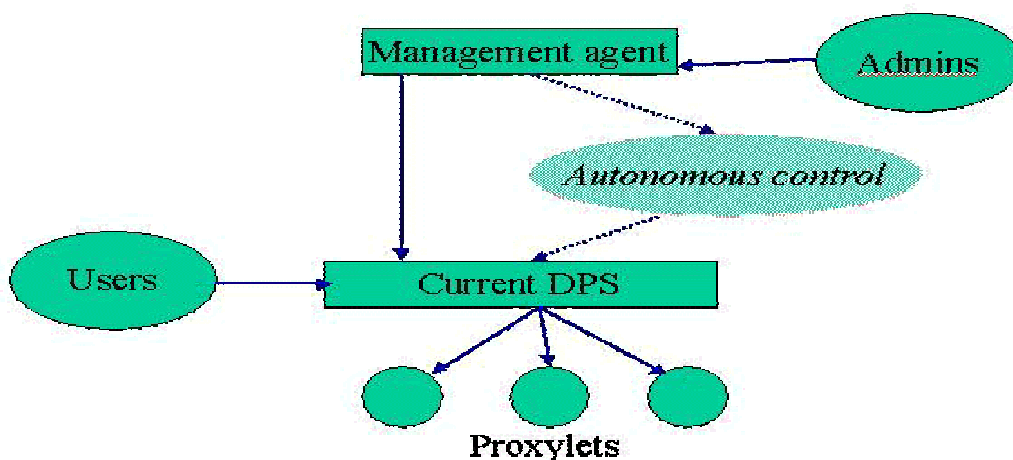


Figure 1. Schematic of proposed ALAN design

The DPS has an autonomous control system that performs management functions delegated to it via policies (scripts and pointers embedded in XML containers). Currently the control system supports a conventional management agent interface that can respond to high level instructions from system operators [3]. This interface is also open to use by users (who can use it to run programs/active services) by adding a policy pointing to the location of their program and providing an invocation trigger. Typically the management policies for the program are included in an XML metafile associated with the code using an XML container, but users can also separately add management policies associated with their programs using HTTP post commands. In addition the agent can accept policies from other agents and export policies to other agents. This autonomous control system is intended to be adaptive.

Not shown in the figure are some low level controls required to enforce sharing of resources between users, and minimise unwanted interactions between users. There is a set of kernel level routines [4] that enforce hard scheduling of the system resources used by a DPS and the associated virtual machine that supports user supplied code. In addition the DPS requires programs to offer payment tokens before they can run. In principle the tokens should be authenticated by a trusted third party. At present these low level management activities are carried out using a conventional hierarchical approach. We hope to address adaptive control of the o/s kernel supporting the DPS in future work.

Network level QoS

Currently there is great interest in enabling the Internet to handle low latency traffic more reliably than at present. Many approaches, such as intserv [5], rely on enabling

the network to support some type of connection orientation. This matches the properties of older network applications, such as telephony, well. However it imposes an unacceptable overhead on data applications that generate short packet sequences. Given that traffic forecasts indicate that by the end of the next decade telephony will be approx 5% of total network traffic, and short data sequences will be around 50% of network traffic, it does not seem likely that connection orientation will deliver optimal results.

A recent alternative has been to propose differentiated services [6], an approach that is based on using different forwarding rules for different classes of packet, and maintaining the properties of the best class by admission control at the ingress to the network. There are difficulties however.

- Admission control does not work well with short packet sequences [7]
- The proposed algorithms assume Poisson burst intervals when real traffic is in fact fractional Gaussian [8,9] and much harder to predict
- The performance benefits can only be obtained if the distribution of demand is such that only a small proportion of the traffic wishes to use the better classes [10]
- The proposed classes typically propose a low loss, low latency class that uses a disproportionate proportion of the available network resources

Despite the difficulties it is clear that differentiated services is currently the best available alternative. It therefore seems advisable to base any proposals for QoS management of active services on the diffserv approach. However, it also seems advisable to modify the approach and attempt to avoid some of the difficulties identified.

Emergent approach to differentiated active services

We propose a new approach to differentiating active services, controlled by an emergent control algorithm. Users can request low latency at the cost of high loss, moderate latency and loss, or high latency and low loss by adjusting the time to live (ttl) of the packets they send. Short ttl packets will experience high loss when the network is congested and long ttl packets will experience high delay when the network is congested. Users cannot request low loss and low delay together. This choice means that all the classes of service we support have approximately the same resource cost. As a result we do not have to consider complex admission control to ensure a favourable demand distribution, and we do not have to allocate significant resources to support a minority service. Two adaptations are possible if the performance is reduced by congestion; either the application sends less packets or the application persists until an application specific latency cut-off is reached and then terminates the session. Services such as telephony would use a low latency/high loss transport regime. This would require the application to be more loss tolerant than at present, however as mobile telephones demonstrate this is not hard to achieve. Interoperation with legacy telephones could be achieved by running loss tolerant algorithms (e.g. FEC) in the PSTN/IP gateway. We do not believe that users want an expensive low loss, low latency service. The current PSTN exemplifies this service and users are moving to VoIP as fast as they are able, despite lower quality, in order to benefit from reduced prices.

Near optimal end to end performance across the network is obtained by enabling the servers to retain options in their application layer routing table for fast path, medium path and slow path (i.e. high loss medium loss and low loss). Packets are then quickly

routed to a server whose performance matches their ttl. This avoids any need to perform flow control and force sequences of packets to follow the same route.

For this approach to work well the properties of the servers must adapt to local load conditions. Fast servers have short queues and high drop probabilities, slow servers have long queues and low drop probabilities. If most of the traffic is low latency the servers should all have short buffers and if most of the demand is low loss the servers should have long buffers. Adaptation of the buffer length can be achieved using an adaptive control mechanism [11], and penalising servers whenever a packet in their queue expires. Use of adaptive control has the additional advantage that it makes no assumptions about traffic distributions, and will work well in a situation where the traffic has significant Long Range Dependency. This then resolves the final difficulty we noted with the current network level diffserv proposals.

Adaptive control

Conventional control of dynamic systems is based on monitoring state, deciding on the management actions required to optimise future state, and enforcing the management actions. In classical control the decision is based on a detailed knowledge of how the current state will evolve, and a detailed knowledge of what actions need to be applied to move between any pair of states (the equations of motion for the state space). Many control schemes in the current Internet (SNMP, OSPF) are based on this form of control. There is also a less precise version known as stochastic control, where the knowledge takes the form of probability density functions, and statistical predictions. All existing forms of traffic management are based on stochastic control, typically assuming Poisson statistics.

Adaptive control [11] is based instead on learning and adaptation. The idea is to compensate for lack of knowledge by performing experiments on the system and storing the results (learning). Commonly the experimental strategy is some form of iterative search, since this is known to be an efficient exploration algorithm. Adaptation is then based on selecting some actions that the system has learnt are useful using some selection strategy (such as a Bayesian estimator) and implementing the selected actions. Unlike in conventional control, it is often not necessary to assume the actions are reliably performed by all the target entities. This style of control has been proposed for a range of Internet applications including routing [12], security [13,14], and fault ticketing [15]. As far as we are aware the work presented here is the first application of distributed adaptive control to service configuration and management.

Holland [16] has shown that Genetic Algorithms (GAs) offer a robust approach to evolving effective adaptive control solutions. More recently many authors [17,18,19] have demonstrated the effectiveness of distributed GAs using an unbounded gene pool and based on local action (as would be required in a multi-owner internetwork). However, many authors, starting with Ackley and Littman [20], have demonstrated that to obtain optimal solutions in an environment where significant changes are likely within a generation or two, the slow learning in GAs based on mutation and inheritance needs to be supplemented by an additional rapid learning mechanism. Inman [21] pointed out that gene interchange (as observed in bacteria [22,23]) could provide the rapid learning required. This was recently demonstrated by Furuhashi [24] for a bounded, globally optimised GA. In previous work [25] we have demonstrated that a novel unbounded, distributed GA with “bacterial learning” is an effective adaptive control algorithm for the distribution of services in an active service provision system derived from the application layer active network (ALAN). In this paper we

demonstrate for the first time that our adaptive control algorithm can deliver differentiated QoS in response to user supplied metrics.

Algorithm Details.

Our proposed solution makes each DPS within the network responsible for its own behaviour. The active service network is modelled as a community of cellular automata. Each automaton is a single DPS that can run several programs (proxylets) requested by users. Each proxylet is considered to represent an instance of an active service. Each member of the DPS community is selfishly optimising its own (local) state, but this 'selfishness' has been proven as a stable model for living organisms [26]. Partitioning a system into selfishly adapting sub-systems has been shown to be a viable approach for the solving of complex and non-linear problems [27].

In this paper we discuss results from an implementation that supports up to 10 active services. The control parameters given below are examples provided to illustrate our approach. Our current implementation has up to 1000 vertices connected on a rectangular grid (representing the network of transport links between the dynamic proxy servers). Each DPS has an amount of genetic material that codes for the rule set by which it lives. There is a set of rules that control the DPS behaviour. There is also a selection of genes representing active services. These define which services each node will handle and can be regarded as pointers to the actual programs supplied by users. The service genes also encode some simple conditionals that must be satisfied for the service to run. Currently each service gene takes the form $\{x,y,z\}$ where:

- x. is a character representing the type of service requested (A-J)
- y. is an integer between 0 and 200 which is interpreted as the value in a statement of the form "Accept request for service [Val(x)] if queue length < Val(y)".

z. is an integer between 0 and 100 that is interpreted as the value in a statement of the form "Accept request for service [Val(x)] if busyness < Val(z)% "

The system is initialised by populating a random selection of network vertices with DPSs (active nodes), and giving each DPS a random selection of the available service genes. Requests are then entered onto the system by injecting a random sequence of characters (representing service requests), at a mean rate that varies stochastically, at each vertex in the array. If the vertex is populated by a DPS, the items join a queue. If there is no DPS the requests are forwarded to a neighbouring vertex. The precise algorithm for this varies and is an active research area, however the results shown here are based on randomly selecting a direction in the network and forwarding along that direction till a DPS is located. This is clearly sub-optimal but is easy to implement. The traffic arriving at each DPS using this model shows some Long Range Dependency (LRD), but significantly less than real WWW traffic. Increasing the degree of LRD would be straightforward. However, the necessary change involves additional memory operations that slows down the simulation and makes the results harder to interpret. In any case inclusion of significant LRD would not change the qualitative form of the main results since the algorithm is not predictive and makes no assumptions regarding the traffic pdf. Each DPS evaluates the items that arrive in its input queue on a FIFO principle. If the request at the front of the queue matches an available service gene, and the customer has included payment tokens equal to (or greater than) the cost for that service in the DPS control rules, the service will run. In the simulation the request is deleted and deemed to have been served, and the node is rewarded by a value equal to the specified cost of the service. If there is no match the request is forwarded and no reward is given. In this case the forwarding is informed by a state table maintained by the DPS using a node state algorithm. Packets with a short ttl

are forwarded to a DPS with a short queue and packets with a long ttl are forwarded to a DPS with a long queue. Each DPS is assumed to have the same processing power, and can handle the same request rate as all the others. In the simulation time is divided into epochs (to enable independent processing of several requests at each DPS before forwarding rejected requests). An epoch allows enough time for a DPS to execute 3-4 service requests, or decide to forward 30-40 (i.e. forwarding incurs a small time penalty). An epoch contains 100 time units and is estimated to represent O(100)ms. The busyness of each DPS is calculated by combining the busyness at the previous epoch with the busyness for the current epoch in a 0.8 to 0.2 ratio, and is related to the revenue provided for processing a service request. For example, if the node has processed three requests this epoch (25 points each) it would have 75 points for this epoch, if its previous cumulative busyness value was 65 then the new cumulative busyness value will be 67. This method dampens any sudden changes in behaviour. A brief schematic of this is shown in figure 2.

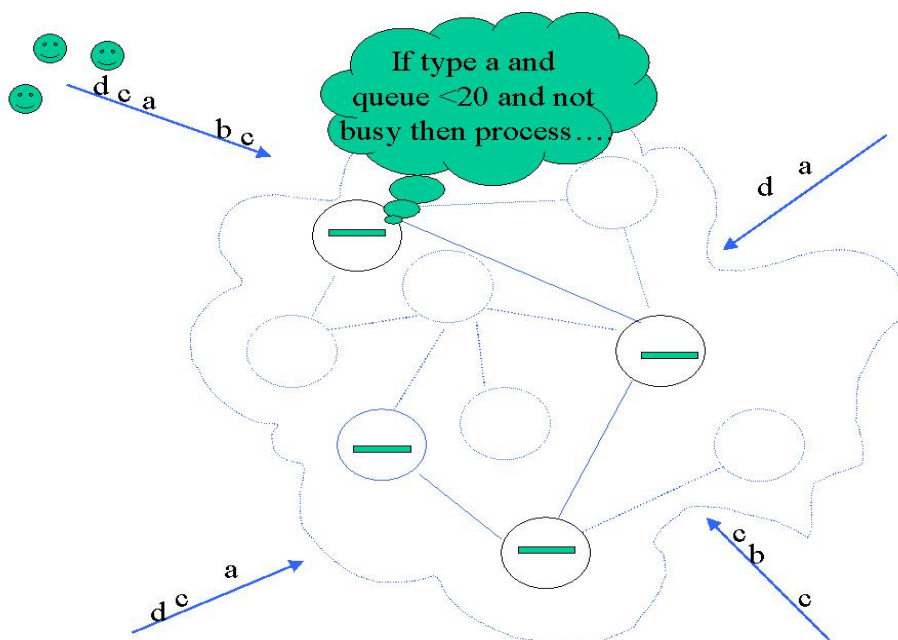


Figure 2 Future network model

The DPS also has rules for reproduction, evolution, death and plasmid migration. It is possible to envisage each DPS as a bacterium and each request for a service as food. The revenue earned when a request is handled is then analogous with the energy released when food is digested. This analogy is consistent with the metabolic diversity of bacteria, capable of using various energy sources as food and metabolising these in an aerobic or anaerobic manner.

Genetic diversity is created in at least 2 ways, mutation and plasmid migration. Mutation involves the random alteration of just one value in a single service gene, for example:

"Accept request for service A if $DPS < 80\%$ busy" could mutate to "Accept request for service C if $DPS < 80\%$ busy" or alternatively could mutate to "Accept request for service A if $DPS < 60\%$ busy".

Plasmid migration involves genes from healthy individuals being shed or replicated into the environment and subsequently being absorbed into the genetic material of less healthy individuals. If plasmid migration doesn't help weak strains increase their fitness they eventually die. If a DPS acquires more than 4-6 service genes through interchange the newest genes are repressed (registered as dormant). This provides a long term memory for genes that have been successful, and enables the community to successfully adapt to cyclic variations in demand. Currently, values for queue length and cumulative busyness are used as the basis for interchange actions, and evaluation is performed every five epochs. Although the evaluation period is currently fixed there is no reason why it should not also be an adaptive variable.

If the queue length or busyness is above a threshold (both 50 in this example), a random section of the genome is copied into a 'rule pool' accessible to all DPSs. If a DPS continues to exceed the threshold for several evaluation periods, it replicates its

entire genome into an adjacent network vertex where a DPS is not present. Healthy bacteria with a plentiful food supply thus reproduce by binary fission. Offspring produced in this way are exact clones of their parent.

If the busyness is below a different threshold (10), a service gene randomly selected from the rule pool is injected into the DPS's genome. If a DPS is 'idle' for several evaluation periods, its active genes are deleted, if dormant genes exist, these are brought into the active domain, if there are no dormant genes the node is switched off. This is analogous to death by nutrient deprivation.

So if a node with the genome $\{a,40,50/c,10,5\}$ has a busyness of >50 when analysed, it will put a random rule (e.g. $c,10,5$) into the rule pool. If a node with the genome $\{b,2,30/d,30,25\}$ is later deemed to be idle it may import that rule and become $\{b,2,30/d,30,25/c,10,5\}$.

Experiments.

The basic traffic model outlined above was adjusted to enable a range of ttls to be specified. The ttls used were 4, 7, 10, 15, 20, 25, 30, 40, 50, 100 (expressed in epochs). Approximately the same number of requests were injected at each ttl. The DPS nodes were also given an extra gene coding for queue length, and penalised by 4 time units whenever packets in the queue were found to have timed out. A DPS with a short queue will handle packets with a short ttl more efficiently since the ttl will not be exceeded in the queue and the DPS will not be penalised for dropping packets. Thus if local demand is predominantly for short ttl DPS nodes with short queues will replicate faster, and a colony of short queue nodes will develop. The converse is true if long ttl requests predominate. If traffic is mixed a mixed community will develop. In figure 3 the red dots represent DPS nodes with long queues, the blue dots represent intermediate

queues and the green dots represent short queues. It is clear that the distribution of capability changes over time to reflect the distribution of demand, in the manner described above.

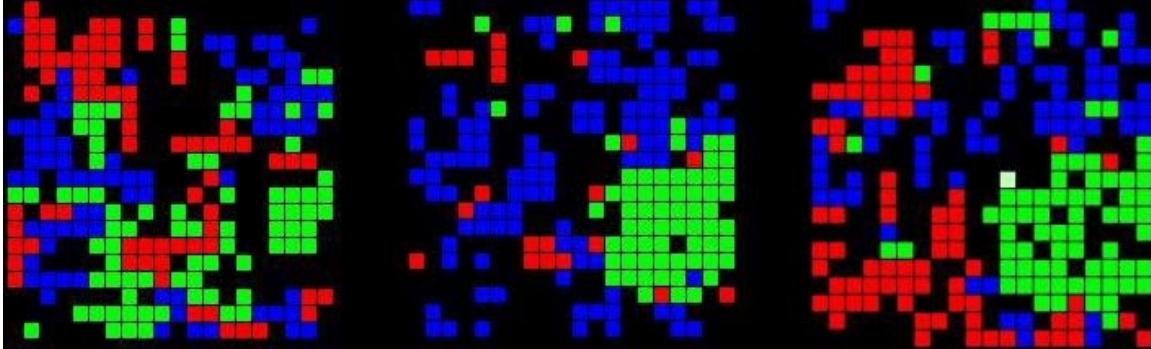


Figure 3. Distribution of DPS nodes with short medium and long queues at three different times.

Figure 4 illustrates the differentiated QoS delivered by the network of DPS nodes. The time taken to process each request is shown on the y axis and the elapsed system time is shown on the x axis. It can be seen that the service requests with shorter times to live are being handled faster than those with a longer time to live, as expected. Figure 5 shows the expected corollary. More service requests with short ttls are being dropped. This is due to them timing out, and is the essential down-side to specifying a short ttl. Although the numbers of requests at each ttl value are roughly equal, fewer short ttl requests are handled.

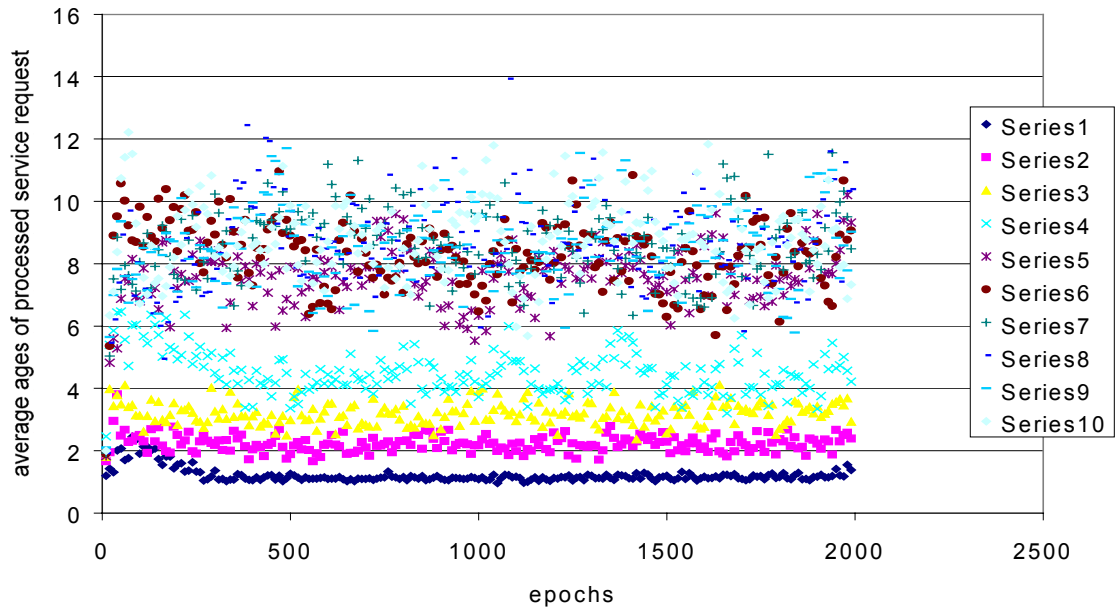


Figure 4. Different latencies for requests with differing times to live.

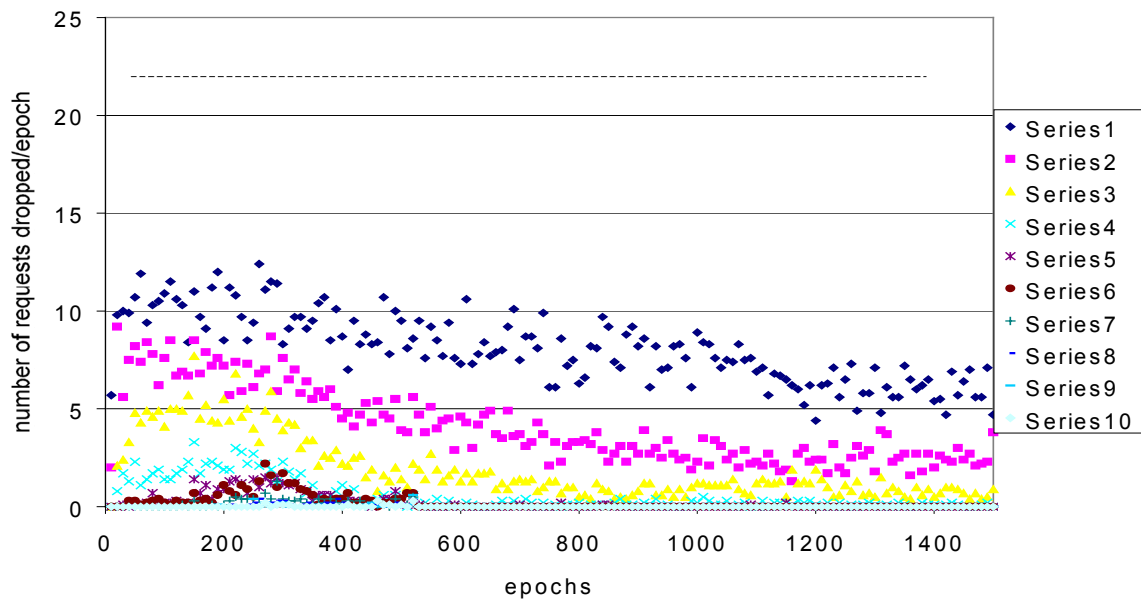


Figure 5. Different dropping rates for requests with differing times to live.

In addition to varying the latency and loss associated with service requests users may also wish to vary the price they are willing to pay. In the basic algorithm it was assumed that the network provider allocated a reward to each DPS for processing a service request. We investigated the impact of allowing the DPS to collect a greater

reward. In the modified model the DPS is rewarded by the amount of tokens the user includes with the request. The traffic input was adjusted so that requests for different services carried different amounts of payment tokens. Initially the DPS nodes were rewarded equally (25 'tokens') for each of three services A, B and C. After 500 epochs the rate of reward is changed so that each DPS is rewarded 4 times as much for processing service C (40 tokens) as it is for processing service A (10 tokens), with B staying at 25. This is equivalent to offering users a choice of three prices for a single service. Fig 6 shows the latency of service requests for the 3 different service types.

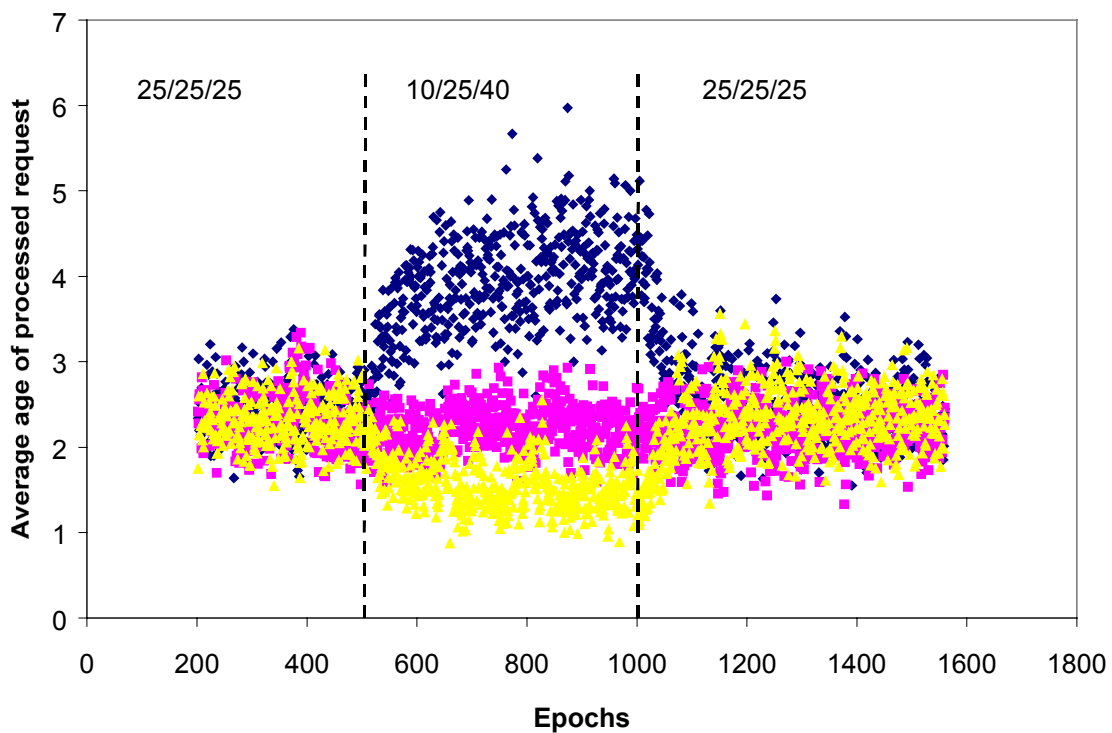


Figure 6. Effects of different charging levels on age related QoS

It is apparent that within 100 epochs the average latency for providing service C is reduced while the latency for A is increased. Fig 7 shows that requests for service A are also dropped (due to timing out) more than requests for service B and C. Before the change in reward the numbers of DPSs handling each service were similar. After the reward rate change the plasmids for handling services C and B have spread much more

widely around the network at the expense of the plasmid for the relatively unrewarding service A. After 1000 epochs the rate of requests for all three services was returned to the original state. It can be seen, in both figures, that equality in quality of service, both in terms of loss rate and latency, quickly returned.

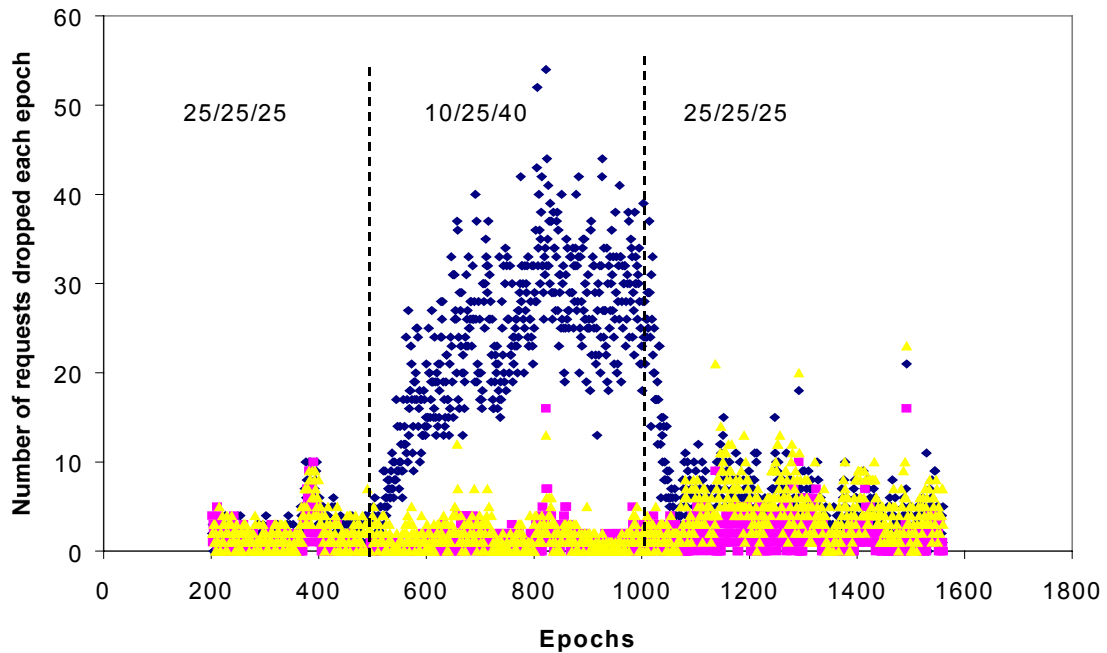


Figure 7. Effects of different charging levels on dropping of requests

These last results indicate that the control method could potentially be used for a range of user specified parameters. We see no reason why other parameters of interest could not be added to the model, and are very encouraged by the initial results. In particular we note that the latencies and loss rates are comparable to those obtained in many conventional approaches to differentiated services, but many of the difficulties concerning admission control have been avoided.

Conclusions.

Our initial results show that the long-term self-stabilising, adaptive nature of bacterial communities are well suited to the task of creating a stable community of

autonomous active service nodes that can offer consistent end to end QoS across a network. The methods used for adaptation and evolution enable probabilistic guarantees for metrics such as loss rate and latency similar to what can be achieved using more conventional approaches to differentiated services.

References.

- [1] Fry M and Ghosh A, "Application Layer Active Networking" *Computer Networks*, 31, 7, pp. 655-667, 1999.
- [2] Giaffreda R, Walker N and Ruiz RY, "Performance of the DNS name resolution infrastructure", Proc. IEE colloquium on "*Control of Next Generation Networks*", London, Oct 1999.
- [3] Marshall IW, Hardwicke J, Gharib H, Fisher M and Mckee P, "Active management of multiservice networks", Proc. IEEE NOMS2000 pp981-3
- [4] Waddington DG and Hutchison D, "Resource Partitioning in General Purpose Operating Systems, Experimental Results in Windows NT", *Operating Systems Review*, 33, 4, 52-74, Oct 1999.
- [5] IETF Intserv charter: <http://www.ietf.org/html.charters/intserv-charter.html>
- [6] IETF Diffserv charter: <http://www.ietf.org/html.charters/diffserv-charter.html>
- [7] Jacobson V, "Support for differentiated services in router designs", Proc. Network modelling in the 21st Century, Royal Society Dec 1999.
- [8] Paxson V and Floyd S, "Wide Area Traffic: The Failure of Poisson Modelling", *IEEE/ACM Transactions on Networking*, 3, 3, pp 226-244, 1995.
- [9] Crovella M and Bestavros A, "Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes", *IEEE/ACM Transactions on Networking*, 5, 6, pp 835-846, Dec 1997.

- [10] Gibbens R, Sargood SK, Kelly FP, Azmoodeh H, MacFadyen R, MacFadyen N, "An Approach to Service Level Agreements for IP networks with Differentiated Services", Submitted to Phil. Trans. A, Royal society, 2000 and available at <http://www.statslab.cam.ac.uk/~richard/research/papers/sla/>
- [11] Tsybkin YZ, "Adaptation and learning in automatic systems", Mathematics in Science and Engineering Vol 73, Academic press, 1971.
- [12] DiCaro G and Dorigo M, "AntNet: Distributed stigmergic control for communications networks", *J. Artificial Intelligence Research*, 9, pp. 317-365, 1998.
- [13] Fisher DA and Lipson HF, "Emergent algorithms - a new method of enhancing survivability in unbounded systems", *Proc 32nd Hawaii international conference on system sciences*, IEEE, 1999
- [14] Gregory M, White B, Fisch EA and Pooch UW, "Cooperating security managers: A peer based intrusion detection system", *IEEE Network*, 14, 4, pp.68-78, 1996.
- [15] Lewis L, "A case based reasoning approach to the management of faults in telecommunications networks", *Proc. IEEE conf. on Computer Communications (Vol 3)*, pp. 1422-29, San Francisco, 1993.
- [16] Holland JH, "Adaptation in Natural and Artificial Systems" MIT press, 1992.
- [17] Forrest S and Jones T, "Modeling Complex Adaptive Systems with Echo", In *Complex Systems: Mechanisms of Adaptation*, (Ed. R.J. Stonier and X.H. Yu), IOS press pp. 3-21, 1994.
- [18] Burkhart R, "The Swarm Multi-Agent Simulation System", OOPSLA '94 Workshop on "The Object Engine", 7 September 1994.
- [19] Kreft JU, Booth G, and Wimpenny JWT, "BacSim, a simulator for individual-based modelling of bacterial colony growth", *Microbiology*, 144, pp. 3275-3287, 1997.

- [20] Ackley DH and Littman ML, "Interactions between learning and evolution". pp. 487-507 in *Artificial Life II* (ed C.G. Langton, C. Taylor, J.D.Farmer and S. Rasmussen), Addison Wesley, 1993.
- [21] Harvey I, "The Microbial Genetic Algorithm", unpublished work, 1996, available at <ftp://ftp.cogs.susx.ac.uk/pub/users/inmanh/Microbe.ps.gz>
- [22] Sonea S and Panisset M, "A new bacteriology" Jones and Bartlett, 1983.
- [23] Caldwell DE, Wolfaardt GM, Korber DR and Lawrence JR, "Do Bacterial Communities Transcend Darwinism?" *Advances in Microbial Ecology*, Vol 15, p.105-191, 1997.
- [24] Nawa NE, Furuhashi T, Hashiyama T and Uchikawa Y, "A Study on the Relevant Fuzzy Rules Using Pseudo-Bacterial Genetic Algorithm" *Proc IEEE Int Conf. on evolutionary computation* 1997.
- [25] Marshall IW and Roadknight C, "Adaptive management of an active services network", *British Telecom. Technol. J.*, 18, 3, July 2000
- [26] Dawkins R, "The Selfish Gene", Oxford University Press, 1976.
- [27] Kauffman S, "The origins of order", Oxford University Press 1993