

Efficient End-to-End QoS Mechanism Using Egress Node Resource Prediction in NGN Network

Se Youn Ban, Jun Kyun Choi, Hyong-Soon Kim
Information and Communication University

syban@icu.ac.kr, jkchoi@icu.ac.kr, khs@nca.or.kr

Abstract - This paper proposes an efficient end-to-end QoS mechanism using egress node resource prediction via probe method in Next Generation Network (NGN). As we want more smart and intelligent network, NGN is the most important issue to provide users converged and guaranteed services. To provide these services, NGN should support proper end-to-end mechanism to support heterogeneous QoS environment in packet networks. We have studied in IPv6 DiffServ, MPLS DiffServ and other end-to-end QoS with differential service and, there are four ways to provide end-to-end QoS from current best-effort network to NGN; Best-effort QoS in traditional IP network, QoS classification with priority in DiffServ network, QoS aggregation/TE in MPLS network and individual QoS/TE. One of current problems to evolve NGN are there are many legacy equipments which we can not replace at once and various QoS differential service among networks. To address this, there should be proper admission control mechanism to protect core network and support end-to-end QoS in access network. While Planning-based admission is simple but not efficient, and probed-based admission control is efficient but overhead in network, Egress resource prediction-based admission control is efficient and less overhead compared to the previous mechanisms. Egress resource prediction-based admission control mechanism has two parts. First, it checks utilization of egress node by probing. But not like probe-based admission control, it does not send probe as always as there are QoS requests. It handles a bundle of requests with probing, predicts state of egress node and sends probe to handle next bundle of requests. In this mechanism, how to measure current state of egress node and predict its future state.

Keywords — NGN, end-to-end QoS, admission control, differential service.

1. Introduction

Today's hottest issue is a Next Generation Network (NGN) to provide users consistent and ubiquitous services. The NGN discussed in ITU-T has such characteristics; packet-based transfer, broadband capabilities with end-to-end services, interworking with legacy networks, converged services between fixed and mobile network, and so on. To support various user requests, NGN should guarantee various end-to-end QoS [1]. But NGN is a packet-based IP network and congestion which can make network unstable to guarantee end-to-end QoS is the biggest problem. There may be three ways to assure

end-to-end QoS; priority scheduling, resource reservation and admission control. Diffserv IP network shall be core network of NGN and the network resource shall be controlled by RACS in NGN [2]. Admission control is the most efficient way to protect core network to keep its ability to support end-to-end QoS. Actually, Diffserv is not enough to satisfy various users' requests and large amounts of end-to-end connections make pure resource reservation in core network. Here we focus on efficient end-to-end QoS mechanism using admission control in core network of NGN environment.

This paper consists of the following; Section 2 describes the related works in admission control mechanism in packet network. Section 3 proposes an architecture and mechanism of admission control to support end-to-end QoS in NGN. Finally, we conclude in section 4.

2. Related works

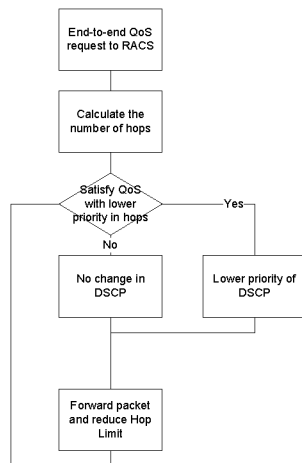
Admission control in IP network can be categorized by two areas; Planning-based admission control like utilization-based admission control [3] uses pre-calculated utilization as admission control parameter. It is very simple and commonly used, but not efficient. Another approach in admission control is Probe-based admission control like RSVP[4] and agent-based selection[5]. It is efficient of using network resource and dynamically adapted to any network situation, but very complex and overhead because it sends a probe every user request over core network.

3. Architecture and mechanism of admission control to support end-to-end QoS in NGN using prediction

NGN network consists of two parts; access network and core network. It is assumed that IP Diffserv and MPLS Diffserv are supported in NGN core network [6]. RACS supports resource and admission control in NGN and one of its services is to ensure

end-to-end network resource for session service. A problem of RACS is that when it provides like session service, it should check and handle network resource for each session service. This situation causes overload in core network than access network, because all of service request from each access network flows into a core network. The rate to check and handle network resource in core network is much bigger than in access network and it is necessary to reduce complexity in core network. Edge Router (ER) is a gateway between core and access network. The admission control mechanism is located in ER.

NGN is based on IPv6. IPv6 uses Traffic Class of IPv6 header as Differential Service Code Point (DSCP) to support DiffServ. The proposed algorithm uses additional field of Hop Limit of IPv6 header to provide hop-based proportional Diffserv. Hop Limit field has two usages. First, when user request end-to-end connection with QoS, RACS receives the request and Policy Decision Function of RACS determines and provides network parameters to satisfy the request. RACS calculates expected number of hops which packet will proceed and remarks Hop Limit regards with the number of hops. It also match DSCP with user's QoS request and the number of Hop Limit. If the number of Hop Limit is large, it could provide higher priority with DSCP. If the number of Hop Limit is small, it could provide lower priority with DSCP. Second, network uses the remained number of Hop Limit to determine if it changes priority of DSCP. When nodes reduce the number of Hop Limit and the number of Hop is initialized by RACS, the nodes can indicate how much hops remain to reach the destination. If the remains of hops to proceed are small enough to degrade the priority of DSCP, it can provide lower priority of DSCP to increase the utilization of network.



[Fig 4] hop-based proportional Diffserv priority provisioning

5. Numerical Analysis

Let assume user's QoS request is that average minimum bound of packet delay is r ms. The average propagation delay of one hop is p ms. There are two QoS class; High and low. Average queueing delay of High priority per one hop is h ms and that of Low priority is l ms. Let there are two end-to-end connection; longest and shortest routes. One has s hops to reach the destination and the other has d hops to do. In this case, the longest route is a worst case and $(p + h) * d$ ms must be smaller than r ms, user's QoS request. The shortest route has $(p + h) * s$ ms and it is much lower than r ms, user's QoS request. If $(p + l) * s$ ms is smaller than r ms, it can use low priority for the shorted route. It means there are the route which is $(p + h) * (1 - d) / s - p$ hops to get low priority even though the worst case of user's QoS request is bound to high priority. Network can satisfy user's QoS request with lower priority than it expects.

Diffserv network in real world is so complicated, there is a little proper formula to provide total average delay. So I consider Priority Queueing service to bound delay to satisfy quality of priority class and modify Optimal Flow Control Problem with packet network. Final objectives is to find optimal rate of priority class and hops to satisfy priority class delay bound even if it gives lower priority. In this formulation, I try to find optimal rate of priority class in simple case and the later will remain a future study.

Priority queueing model is one node analysis to find average delay and throughput with priority class. There are three kinds of Priority Discipline; Non-preemptive, Preemptive resume and Preemptive non-resume. Non-preemptive is used to formulate the problem.

Average Waiting time for priority class p is below :

$$E[W_p] = E[T_0] + \sum_{k=1}^p E[T_k] + \sum_{k=1}^{p-1} E[T'_k]$$

where T_0 : the completion time of current service

T_k : service time of m_k messages of priority $1, 2, \dots, p$ already waiting

T'_k : service time of $k = 1, 2, \dots, p-1$ high priority message during the waiting time

$$E[W_p] = E[T_0] + \sum_{k=1}^p E[T_k] + \sum_{k=1}^{p-1} E[T'_k] = E[T_0] + \sum_{k=1}^p \rho_k E[W_k] + E[W_p] \sum_{k=1}^{p-1} \rho_k$$

$$E[W_p] = \frac{E[T_0]}{(1 - \sigma_{p-1})(1 - \sigma_p)} \quad \text{where } \sigma_p = \sum_{k=1}^p \rho_k$$

So Average waiting time for two priority class with exponential distribution:

$$E[W_1] = \frac{(\rho_1/\mu_1) + (\rho_2/\mu_2)}{(1-\rho_1)}$$

$$E[W_2] = \frac{(\rho_1/\mu_1) + (\rho_2/\mu_2)}{(1-\rho_1)(1-\rho_2)}$$

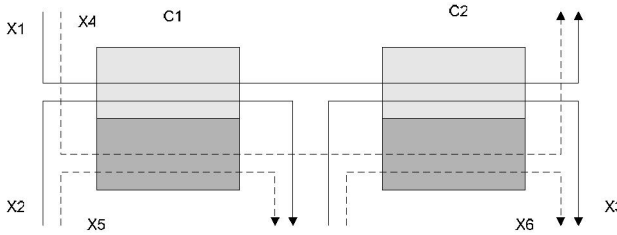
$$E[W] = \frac{(\rho_1/\mu_1) + (\rho_2/\mu_2)}{(1-\rho)} \quad \text{where } \rho = \rho_1 + \rho_2$$

In the case of network, sum of average waiting time in path should be bound to required delay,

$$\sum_i \frac{(\rho_{i1}/\mu_{i1}) + (\rho_{i2}/\mu_{i2})}{(1-\rho_{i1})} \leq D_1 \quad \text{where } \rho_{i1} \text{ is throughput of}$$

high priority in node i, μ_{i1} and μ_{i2} are departure rates of high and low priority. This is a key constraint to formulate the problem.

Based on this facts, let's formulate a simple problem to find optimal proportional rate of high and low priority in Diffserv Network.



[Fig4] Reference Model

This model consists of two nodes which have priority queueing service. There are 6 flows; high priority flows X1, X2 and X3, and low priority flows X4, X5 and

X6. Each flow has his own arrival rate λ_i . The flows of high priority must satisfy the average delay bound D. The flows of low priority is best effort service and has no average delay bound. Each node can receive the same max arrival rate λ_{\max} and the same service rate μ . We have utility function $U(\lambda_i)$ which provides the different utility based on the priority class. Now, we want to find a set of λ_i which satisfy to maximize sum of $U(\lambda_i)$ which satisfy the average delay bound. Let $U(\lambda_i)$ is log and $U(\lambda_i)$ of high priority has twice utility than that of low priority.

1. Determine Decision Variables

For general case

λ_k = arrival rate of priority class

2. Determine Object Function

For general case

$$\text{Maximize utility} = \sum_{\text{high}} U_{\text{high}}(\lambda_i) + \sum_{\text{low}} U_{\text{low}}(\lambda_i)$$

For two priority class case

$$\text{Maximize utility} = 2(\log \lambda_1 + \log \lambda_2 + \log \lambda_3) + (\log \lambda_4 + \log \lambda_5 + \log \lambda_6) \quad \text{---(1)}$$

3. Determine Constraints

For general case

$\sum_{\text{node}} \lambda_i \leq \lambda_{\max}$, sum of arrival rates in a node is bound to max arrival capacity.

$\sum_{\text{priority_flow}} \frac{(\rho_1/\mu_1) + (\rho_2/\mu_2)}{(1-\rho_1)} \leq D_{\text{priority}}$, sum of average delay of priority flow is bound.

For this case

$$\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5 \leq \lambda_{\max} \quad \text{----- (2)}$$

$$\lambda_1 + \lambda_3 + \lambda_4 + \lambda_6 \leq \lambda_{\max} \quad \text{----- (3)}$$

$$\frac{\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5}{1 - \frac{\lambda_1 + \lambda_2}{u}} \leq D \quad \text{----- (4)}$$

$$\frac{\lambda_1 + \lambda_3 + \lambda_4 + \lambda_6}{1 - \frac{\lambda_1 + \lambda_3}{u}} \leq D \quad \text{----- (5)}$$

$$\frac{\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5}{1 - \frac{\lambda_1 + \lambda_2}{u}} + \frac{\lambda_1 + \lambda_3 + \lambda_4 + \lambda_6}{1 - \frac{\lambda_1 + \lambda_3}{u}} \leq D \quad \text{----- (6)}$$

Because (6) constraint is not convex set, we assume that all nodes are identical and average delay constraints change like :

$$\frac{\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5}{1 - \frac{\lambda_1 + \lambda_2}{10}} \leq 1.5 \quad \text{----- (7)}$$

$$\frac{\lambda_1 + \lambda_2 + \lambda_4 + \lambda_6}{1 - \frac{\lambda_1 + \lambda_3}{10}} \leq 1.5 \quad \text{----- (8)}$$

(7) and (8) constraints are convex set and the problem is a convex problem.

Now I increase the high priority average delay bound with the original object function:

	λ_1	λ_2	λ_3	λ_4	λ_5	λ_6	Node Arrival rate
D=4	0	1.2000	1.2000	0	0.8000	0.8000	2.0000
D=5	0	1.3846	1.3846	0	0.9231	0.9231	2.3091
D=6	0.5128	1.0256	1.0256	0.3333	0.6667	0.6667	2.5384
D=7	0.5761	1.1523	1.1523	0.3889	0.7778	0.7778	2.8954
D=19	1.2281	2.4561	2.4561	0.7719	1.5439	1.5439	5.3842
D=20	1.3333	2.6667	2.6667	0.6667	1.3333	1.3333	6
D=	1.3333	2.6667	2.6667	0.6667	1.3333	1.3333	6

There is a point between D=5 and D=6 where long path flow has a part of optimal solution. As delay bound increases, the sum of optimal flows rate at node increases.

After D=19, there are same results. Because increasing D infinite is removing the delay bound constraint, the sum of optimal flow rate at node is as same as maximum flow rate capacity at node. I guess there may be co-relationship between rate capacity and the saturation point - It means we may calculate node capacity given delay bound in Diffserv network. As we see the result, long path flow causes network performance low. It is unexpected result that to optimize the network utility is decreasing long path flow of low priority. We can find the optimal flow rate given delay-bound and node capacity. If we know the optimal flow rate, we can modify Diffserv traffic with the optimal flow rate - like increasing or decreasing high priority flow to the optimal flow. We shows the optimal flow and our algorithm is useful to increase network utility.

4. Conclusion and Future works

In this paper, we have present category of admission control in IP network in NGN, architecture and mechanism of efficient admission control using usage prediction. The admission control mechanism reduces network resource reservation complex in unit of log scale. We have also proposed segment reservation to support fast network resource reservation in egress node.

REFERENCES

- [1] Y.NGN-FRA, "Functional Requirements and Architecture of the NG", ITU-T, May 2004
- [2] TR.RACS, "Functional Requirements and Architecture for Resource and Admission Control in Next Generation Networks", ITU-T, May 2004
- [3] D. Xuan, C. Li, R. Bettati, J. Chen, W. Zhao, "Utilization- Based Admission Control for Real-Time Applications," *The IEEE International Conference on Parallel Processing*, Canada, Aug. 2000
- [4] L. Zhang, S. Deering, D. Estrin, S. Shenker and D. Zappala, "RSVP: a new resource reservation protocol," *IEEE Networks Magazine*, vol. 31, No. 9, pp. 8-18, September 1993.
- [5] G. Papaioannou, S. Sartzetakis, and G.D. Stamoulis. Efficient agent-based selection of DiffServ SLAs over MPLS networks within the ASP service model. Journal of Network and Systems Management, Special Issue on Management of Converged Networks, Spring 2002
- [6] Y.e2eqos.2, "An end-to-end QoS architecture based on centralized resource control for IP networks supporting NGN services", ITU-T, May 2004