

Introduction to QoS on Janet and in the wider academic community

Intended Audience

This guide provides technical guidance to administrators within Regional Networks and end sites who are considering the deployment and operation of QoS (Quality of Service) in their networks. Site administrators will find this document useful, but will also need to consult with the Regional Network Operator to which they are connected for further information regarding their QoS Policy.

This guide assumes a good knowledge of the networking protocols, including the Internet Protocol (IPv4), and network planning and configuration issues. This is a highly technical subject and the guide is inevitably very acronym heavy. Appendix C provides a glossary of these acronyms.

This Guide and Future Versions

This version of the QoS on JANET Technical Guide is based on work conducted between 2006 and 2008 as part of the second phase of the JANET QoS Development Project. As such, this guide provides a state-of-the-art snapshot of the work conducted and experience gained throughout this period.

This guide reflects current JANET policy towards deployment and use of QoS in the JANET backbone, based on the experiences of JANET(UK) and the other partners during the project. We must make it clear, however, that we do not provide a definitive policy for QoS either in terms of deployment or operation in UK academic networks. This is largely because we have found that there is still no single, unified approach to QoS and at this stage, both within the UK and internationally, work is still ongoing in this area. Instead, this guide presents the findings of the JANET QoS Development Project to help inform network administrators and technical staff who find themselves in a similar position, in the hope of raising awareness of the issues we encountered.

This guide is the final deliverable of the JANET QoS Development Project and, at the time of writing, represents the state of the art in the area of QoS deployment and operation. JANET(UK) reserves the right however to refine the recommendations made in this guide at a later time in order to revise the current QoS Policy if a more definitive picture of QoS deployment emerges.

An Application Perspective on QoS

It is possible to optimise a network intended to provide a single type of service, or even support only a specific application, for the sole benefit of that service or application. The oft-quoted example is the digital telephone network in which circuit-oriented switching and multiplexing is engineered specifically to optimise transmission of the fundamental audio encoding for voice, which generates a stream of bytes at a constant rate. Networks developed to enable data exchange amongst computers adopted a packet-oriented mode of operation, largely to enable statistically efficient multiplexed use of transmission lines in a way that matches well with computer operation, which is essentially bursty. Hence an important property of a traditional data application using the network is that it does not demand that any particular timing constraints should be met by the data transport. Therefore computer packet networks typically do not do this. Such applications are said to have '**Elastic**' requirements for data transmission as they are flexible and willing/able to stretch if necessary. Another aspect of network support for such elastic applications is that, since time is not critical, when a packet is lost through transmission corruption or congestion at a network node, there is time for recovery to be achieved by resending the packet. In contrast to this, in applications that are concerned with conversational audio or video communication, time is crucial. Time-sensitive audio-video applications impose a number of timing constraints on the underlying communications service which include end-to-end delay (one-way and round-trip), variations in this delay and end-to-end synchronization. Such time-sensitive applications are amongst those termed '**real-time**'.

The goal of an integrated services or multiservice network is to support a range of services to meet the transport needs of both elastic and real-time applications. Applications of these two fundamental classes can be further divided into sub-classes according to their specific network requirements. The need to support different classes of traffic in the network led to the concept of being able to allocate resources based on application requirements.

Application Network Requirements

Interactive audio and video communications are well-known examples of communication where timing is important. Whereas in the 1980s JANET was a pure data network supporting elastic applications, during the 1990s the use of videoconferencing became increasingly important, and by 2000 the question of how to expand the range of quality of service provided was part of the SuperJANET4 development programme. Around the same period, other education and research networks were considering the same questions, including the European backbone interconnect, GÉANT, and Internet2 in the USA.

Along with various forms of videoconferencing, applications which may loosely be described as telephony-like have similar timing constraints and make corresponding demands of the underlying network. Although the timing aspect of such 'continuous media' services has been emphasised, these applications also have requirements regarding both transmission capacity and error rate, depending on the quality required and the encoding and compression used. Such requirements are not peculiar to continuous media applications: transfer of a file of a given size within a specified time requires a certain capacity, whatever the nature of the data. Transmission errors must also be guarded against. However, although occasional bit-errors may be tolerable for continuous media, loss of whole packets is often not. Because of the lack of time available to request retransmission, either the application resorts to some form of transmission redundancy (forward error correction) – which may tend to negate the effect of

using compression – or, if the major cause of packet loss is through network congestion, preferential treatment by the network may be sought.

Apart from interactive continuous media, there are other application contexts in which bounds on delay or delay variation (jitter) may be helpful. Within the category of elastic applications are those sometimes referred to as ‘transactional’: one system makes a request of another which responds, typically with data extracted from a collection or perhaps with the result of some action or calculation. In this case, the need for a timely response may originate directly from a person sitting at the requesting system or, in a more complex case, it arises in shared workspace applications, including collaborative use of remote instruments like telescopes, microscopes, accelerators, etc. More recently, there has been considerable development of loosely coupled, wide area, dynamically assembled distributed systems, whether under the guise of Grid-oriented technology in support of science or web-service based applications for commercial or domestic consumption. Such systems have raised the need for timely systems-level transactional support: in this case, typically, to maintain performance in a complex distributed system.

Another area in which timing delay arises, if in a somewhat weaker form, is video or audio streaming. Here, there may be no particular relation to real time to maintain (because the source is a recording), or the source may be near to real time as in some instances of Internet TV or similar. However, even in these latter cases, the requirement to limit the delay between the source and the viewer is usually not stringent and, in order to eliminate jitter to enable smooth playback, the stream can be buffered at the receiver. Reducing the delay introduced by the buffer has to be balanced against reducing the jitter imposed on the stream by the network.

Application-Centric QoS

The term ‘Quality of Service’ (QoS) has a variety of meanings depending on the context. In this guide it relates to techniques of managing network resources to reduce the negative effects of congestion in packet-switched networks and to support the performance constraints of applications by networks carrying a mix of real-time and elastic traffic. The main resources to which QoS techniques are applied are link capacity, router processing capacity and router memory. Commonly, QoS is interpreted to mean the use of packet queuing and scheduling techniques within routers and switches. However, since the onset of congestion depends on loading and network capacity, network provisioning (determining physical resources and their allocation between competing types of traffic) also needs to be included. Inadequate provisioning leads to a slow or broken service but lavish or over-provisioning may be expensive to deploy and maintain. Balancing these for different services and communities tends to be a matter of political or commercial judgement, complicated by the service provisioning failure points for different application traffic types not being the same.

Clearly, if there are no queues on any output port in the network then packet transmission does not involve scheduling, and no delay or packet loss through discard occurs. Although this ideal is not achievable at all times because output port contention is generally unavoidable in packet networks, by increasing provisioning it is possible to reduce the probability of contention and queuing occurring. Thus, attempting to avoid the need for QoS scheduling within the network tends to shift the focus from scheduling to provisioning. Managing QoS scheduling is complex, and hence both expensive and a source of operational

error (whether automated or human). However, increased capacity may also be expensive or even unavailable in some segments of a network. Instead of simply increasing network capacity for all traffic, QoS acts in a more economical way as it takes into account the different tolerances of application classes (elastic and real-time) to packet delays and loss. As such, QoS is based on adequate provisioning of bandwidth to traffic classes (i.e. re-allocation of available bandwidth between traffic classes) according to their needs. As a result traffic of time-sensitive applications is treated in such a way that provides an acceptably low level of delays and packet loss, while traffic of elastic applications experiences greater packet delays/loss which nevertheless are acceptable for this type of application.

In addition to the traditional application areas described above in which QoS may be needed, historically there has always been a demand for 'sequestered' network or link capacity. On the one hand, this may result from a general desire of particular communities or programmes for the equivalent of a guaranteed, dedicated share of the regular network service. On the other hand, it may arise from specific requests for a sequestered share of the underlying bearer service, typically to enable network experiments which would interfere with normal network operation, or because the form of transmission to be used by some activity is incompatible with that in operation for the service network. Collectively, provision of these types of service has come to be referred to as Managed Bandwidth Services (MBS). Providing access to underlying bearer services is beyond the scope of QoS, but QoS techniques can potentially play a role in providing a 'virtualized' share of packet services. As such, MBS is discussed where appropriate within this guide at the national, regional and site level.

QoS is not the only technology to minimise such important network performance characteristics. Some alternative approaches to packet delays and loss are described in sections 2.3 and 2.4.

Alternatives to QoS in JANET

In principle, some or all of the application requirements discussed in the previous section could potentially be met by a combination of dedicated switched physical links, dedicated switched analogue frequency or wavelength channels, dedicated switched digital time division channels, or varying levels of traffic engineering applied to packet-level services. During the period of SuperJANET3, ATM (asynchronous transfer mode) was used as the bearer for the IP service, and dedicated switched virtual channels were available on request as what became known as MBS. Policing of the link capacity used by these channels within the ATM infrastructure was enforced to protect other users of the links, including the service network.

The advent of SuperJANET4, in which ATM was no longer used, led to interest in the use of QoS techniques in pursuit of such capacity sharing. The requirement is superficially similar to that for link sharing for a particular application or class of applications. However, in the MBS case, the requirement typically relates to a (virtual) community rather than a class of applications.

From the beginning of 2007, in common with many other networks, JANET adopted the approach of having multiple bearer services and correspondingly making their use available to the community. This has currently moved the focus away from the use of QoS techniques for MBS-like purposes.

Overview of Activity in the Academic Community – a World-wide

Perspective

This section presents an overview of ongoing QoS development work to supplement and provide context to the work done within this project. In this summary, we describe the QoS projects of other network operators, before discussing ongoing standardisation work and related QoS research projects in the wider academic community.

Related QoS Projects

An overview of the QoS projects of other network operators can be found on the JANET(UK) website [QoSProj] and is summarised below:

- **GÉANT.** GÉANT is the pan-European interconnect network for research and education which is planned, built and managed by DANTE (Delivery of Advanced Network Technology to Europe) [Dante]. Current GÉANT QoS services include IP Premium, Best Effort and Less-than-Best Effort. The QoS services were established on the basis of several earlier projects, the most important of which was SEQUIN.
- **GÉANT2.** GÉANT2 started on 1 September 2004 with the aim of developing and deploying the seventh generation of the pan-European research and education network, the successor to GÉANT. The project will run for four years and JANET(UK) participates in GÉANT2 Service Activity 3 (GN2 SA3), aiming to implement an end-to-end IP Premium QoS service across GÉANT2 and NRENs which are interested in this activity [GÉANTPIP].
- **MB-NG.** The Managed Bandwidth Next Generation (MB-NG) project [MB-NG] addressed the issues which arise in the sector of high performance inter-Grid networking, including sustained and reliable high performance data replication, end-to-end advanced network services and QoS. The project finished in September 2004.
- **Internet2 QoS Working Group.** The Internet2 QoS Working Group [Internet2QoS] is the developer of the Scavenger service which is also known as the Less than Best Effort service. The group also explored the possibilities of the IP Premium service and published useful material on Bandwidth Brokering.

Standardisation Work

The IETF was responsible for developing the key QoS standards and protocols in the Integrated Services Working Group [IntServWG] and Differentiated Services Working

Group [DiffServWG]. Although these working groups were concluded in 2000 and 2003 respectively, they defined many of the underlying protocols that our work is based on such as IntServ and DiffServ. Although formal standardisation work in this area has now concluded, related work is still being done in the IP Performance Metrics [IPPMWG] and Datagram Congestion Control Protocol [DCCPWG] working groups.

While a thorough discussion of the standards and protocols developed here is beyond the scope of this guide, the key RFCs will be identified and referenced later in this guide as appropriate.

QoS Research Projects

QoS related research has been actively conducted since the mid 1990s and has been important in steering the development of QoS. While a complete history of QoS related projects cannot be included here, a selection of the major UK and European-based efforts is presented below:

- Traffic Engineering for Quality of Service in the Internet, at Large Scale (TEQUILA) [TEQUILA]
- Multi Service Access Everywhere (Muse) [Muse]
- Policy Analysis for Quality of Service Management (PAQMAN) [PAQMAN]
- Management of End-to-end Quality of Service Across the Internet at Large (MESCAL) [MESCAL]
- The Moby Dick Project [MobyDick]
- The ENTHRONE Project [ENTHRONE].

Source URL: <https://community-stg.jisc.ac.uk/library/janet-services-documentation/introduction-qos-janet-and-wider-academic-community>