

**Telecommunications and Internet converged Services and
Protocols for Advanced Networking (TISPAN);
Control of Processing Overload;
Stage 2 Requirements**



Reference

DTS/TISPAN-02035-NGN-R2

Keywords

architecture, control, functional

ETSI

650 Route des Lucioles
F-06921 Sophia Antipolis Cedex - FRANCE

Tel.: +33 4 92 94 42 00 Fax: +33 4 93 65 47 16

Siret N° 348 623 562 00017 - NAF 742 C
Association à but non lucratif enregistrée à la
Sous-Préfecture de Grasse (06) N° 7803/88

Important notice

Individual copies of the present document can be downloaded from:

<http://www.etsi.org>

The present document may be made available in more than one electronic version or in print. In any case of existing or perceived difference in contents between such versions, the reference version is the Portable Document Format (PDF). In case of dispute, the reference shall be the printing on ETSI printers of the PDF version kept on a specific network drive within ETSI Secretariat.

Users of the present document should be aware that the document may be subject to revision or change of status. Information on the current status of this and other ETSI documents is available at

<http://portal.etsi.org/tb/status/status.asp>

If you find errors in the present document, please send your comment to one of the following services:

http://portal.etsi.org/chaicor/ETSI_support.asp

Copyright Notification

No part may be reproduced except as authorized by written permission.
The copyright and the foregoing restriction extend to reproduction in all media.

© European Telecommunications Standards Institute 2008.
All rights reserved.

DECTTM, **PLUGTESTS**TM and **UMTS**TM are Trade Marks of ETSI registered for the benefit of its Members.
TIPHONTM and the **TIPHON logo** are Trade Marks currently being registered by ETSI for the benefit of its Members.
3GPPTM is a Trade Mark of ETSI registered for the benefit of its Members and of the 3GPP Organizational Partners.

Contents

Intellectual Property Rights	4
Foreword.....	4
1 Scope	5
2 References	5
2.1 Normative references	5
3 Abbreviations	5
4 TISPAN NGN overload control requirements	6
4.1 High level overload control requirements	6
4.2 General requirements for Nearest Neighbour load control.....	7
4.3 Deployment specific requirements	8
4.4 Application specific requirements	8
Annex A (informative): Comparison between RACS and nearest neighbour overload control	10
History	11

Intellectual Property Rights

IPRs essential or potentially essential to the present document may have been declared to ETSI. The information pertaining to these essential IPRs, if any, is publicly available for **ETSI members and non-members**, and can be found in ETSI SR 000 314: "*Intellectual Property Rights (IPRs); Essential, or potentially Essential, IPRs notified to ETSI in respect of ETSI standards*", which is available from the ETSI Secretariat. Latest updates are available on the ETSI Web server (<http://webapp.etsi.org/IPR/home.asp>).

Pursuant to the ETSI IPR Policy, no investigation, including IPR searches, has been carried out by ETSI. No guarantee can be given as to the existence of other IPRs not referenced in ETSI SR 000 314 (or the updates on the ETSI Web server) which are, or may be, or may become, essential to the present document.

Foreword

This Technical Specification (TS) has been produced by ETSI Technical Committee Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN).

1 Scope

The present document describes the specific TISPAN requirements for controls to manage overload of processing resources in NGNs. In particular, it addresses overload control between nearest neighbours.

2 References

References are either specific (identified by date of publication and/or edition number or version number) or non-specific.

- For a specific reference, subsequent revisions do not apply.
- Non-specific reference may be made only to a complete document or a part thereof and only in the following cases:
 - if it is accepted that it will be possible to use all future changes of the referenced document for the purposes of the referring document;
 - for informative references.

Referenced documents which are not found to be publicly available in the expected location might be found at <http://docbox.etsi.org/Reference>.

For online referenced documents, information sufficient to identify and locate the source shall be provided. Preferably, the primary source of the referenced document should be cited, in order to ensure traceability. Furthermore, the reference should, as far as possible, remain valid for the expected life of the document. The reference shall include the method of access to the referenced document and the full network address, with the same punctuation and use of upper case and lower case letters.

NOTE: While any hyperlinks included in this clause were valid at the time of publication ETSI cannot guarantee their long term validity.

2.1 Normative references

The following referenced documents are indispensable for the application of the present document. For dated references, only the edition cited applies. For non-specific references, the latest edition of the referenced document (including any amendments) applies.

- [1] ITU-T recommendation E.412: "Network management controls".
- [2] ETSI TR 182 015: "Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN); Next Generation Networks; Architecture for Control of Processing Overload".
- [3] ETSI ES 282 003: " Telecommunications and Internet converged Services and Protocols for Advanced Networking (TISPAN); Resource and Admission Control Sub-system (RACS); Functional Architecture".
- [4] ETSI TS 181 005 V2.4.0 (2007-11): "Telecommunications and Internet Converged Services and Protocols for Advanced Networking (TISPAN); Service and Capability Requirements".

3 Abbreviations

For the purposes of the present document, the following abbreviations apply:

IMS	IP Multimedia Subsystem
ISDN	Integrated Service Digital Network
NGN	Next Generation Network

PSTN	Public Switched Telecommunication Network
QoS	Quality of Service
RACS	Resource Admission Subsystem
SCF	Session Control Function
SIP	Session Initiation Protocol
SLA	Service Level Agreement

4 TISPAN NGN overload control requirements

4.1 High level overload control requirements

In TS 181 005 [4] we have the following requirements for overload control.

The NGN shall have mechanisms available to control overload that:

- 1) automatically maximize effective throughput (i.e. admitted service requests/sec) at an overloaded resource;
- 2) achieve this throughout the duration of an overload event, and irrespective of the overloaded resource's capacity or of the number of sources of overload;
- 3) are configurable by the service provider so that, under processing overload, a high proportion of response times at overloaded resources are low enough so as not to cause customers to prematurely abandon service requests;
- 4) should be possible to be applied within a service provider's NGN, and between different service providers' NGNs;
- 5) should be possible to be applied within an NGN subsystem (e.g. IMS, PSTN/ISDN emulation) and between different NGN subsystems.

NOTE: As a general rule, an NGN's call, session and command processing resources can experience prolonged processing overload under the appropriate circumstances (e.g. partial, or full, server failure, high rates of incoming service requests). Consequently, it needs to be equipped with some form of overload detection and control (including expansive controls such as load balancing and resource replication), in order to keep response times just low enough under such processing overload to preclude customers abandoning their service requests prematurely.

Many pieces of equipment will have internal load control, which aims to meet the ideal behaviour described in ITU-T recommendation E.412 [1] and shown graphically in figure 1. The object of these internal controls is to bound the system response time by rejecting some of the workload (because rejecting workload requires less effort than accepting it). As the load increases requests are rejected, but because the rejected requests still consume some processing resource, the rate at which requests can be accepted falls. As a consequence, such internal load control can only protect the physical host against overload to a limited extent. Severe overloads will reduce the rate at which useful work can be done, and very severe overloads may cause the system to operate incorrectly or with unacceptably long response times.

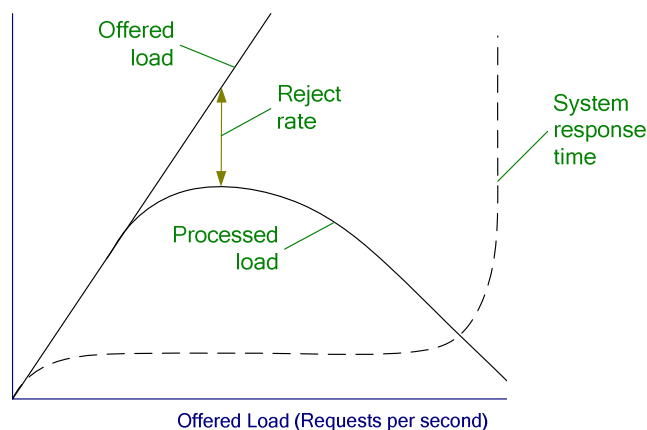


Figure 1: hedoo Typical overload behaviour in a host with internal load control.

By deploying a distributed overload control, in which systems suppress some service requests before they reach an overloaded neighbour, the overloaded system is protected from more extreme overloads thus enabling it to operate at near optimum load. For the fulfilment of these NGN requirements, the use of such nearest neighbour overload controls is an essential component. The issues regarding overload controls for NGNs have been discussed at some length in [2] which provides some initial requirements for overload controls. These requirements are further elaborated in the present document and allocated into different categories. The requirements are also extended to cover the case where an overloaded target has a non-enumerable set of sources (e.g. SIP user agents overloading an SCF). The overload control requirements for systems loaded by a set of non-enumerable sources may be similar to those for a source loaded by a known, enumerable set of sources, but the approach may be very different. In the enumerable sources case, we may want to control the relative performance seen by each source, offer SLAs to particular sources and have a variety of fairness criteria to use. In the non-enumerable sources case, the issue is one of reducing the aggregate demand, and the fate of individual sources is less predictable. The definition of fairness would perhaps depend on ensuring equitable treatment over a series of overloads, rather than trying to meet fairness criteria between many sources over a single overload event. Those requirements that are not appropriate for non-enumerable sources are labelled as such.

4.2 General requirements for Nearest Neighbour load control

These requirements are general, and any overload control should address these requirements.

- REQ 1 A control shall automatically reduce the load that is sent to an overloaded host.
- REQ 2 A control shall aim to maximize the number of fully processed service requests at an overloaded host, subject to meeting any QoS constraints for the application (e.g. response times).
- REQ 3 The overload control restriction shall apply only to defined events which will depend on the reference point that is being controlled.
- NOTE 1: This information may be thought of as within the domain of the applications on the NGN servers. It is the application that "knows" which requests to offer to the control for possible rejection, the control provides the infrastructure for the restriction to occur.
- REQ 4 A control shall be aware of differing importance levels of service requests, and be configurable to reject lower importance service requests in favour of service requests of higher importance level (for example to protect calls to emergency services).
- REQ 5 When several inter-acting controls are active at an overloaded host at the same time, they shall converge to an acceptable steady-state when the amount of work sent by each of the sources is constant.
- REQ 6 A control shall be configurable by the service provider to enable the enforcement of SLAs i.e. to divide the capacity of an overloaded resource between hosts from which requests are received (or groups of hosts) according to agreed policies.

NOTE 2: It is inappropriate to attempt to manage SLAs for individual sources from a set of non-enumerable sources.

- REQ 7 A control shall enforce fair allocation of overloaded processing resources between competing request sources (a request source in this context means the host from which the overloaded resource receives the request, not necessarily the originating host).
- NOTE 3: It is the SLA that defines the resources allocated to a request source by an overloaded resource. A fair allocation of resources implies that the resource consumed by every source is appropriate given the SLA allocated to that source. It does not necessarily require that the each request source obtains the same amount of processing resource.
- REQ 8 The behaviour of the control shall be fully-specified.
- REQ 9 A control shall allow manual configuration of the overload control's components via a management interface.
- REQ 10 A control shall operate (optionally) without manual configuration of the overload control parameters.
- REQ 11 A control shall output network management data on event occurrence (e.g. control activation/termination) and on demand from network management (e.g. counts of service requests admitted and rejected by the overload control).
- REQ 12 A control shall have adequate security from malicious actions. It shall not be possible for a host to instantiate a restriction able to reject service requests destined for a different host.
- NOTE 4: This implies that a mechanism for the verification of the identity of hosts originating restriction messages is required.
- REQ 13 A control shall be applicable within a network, and between networks.
- REQ 14 It should be possible to apply overload control between different NGN subsystems.
- REQ 15 A control shall react quickly to changes in the workload sent to the host it is protecting.
- REQ 16 A control shall be capable of dealing with work received from hosts which do not support it. Such hosts that are part of a finite set of sources shall not receive disproportionate benefit (i.e. requirements 6 and 7 shall still hold).

4.3 Deployment specific requirements

The following requirements relate to different deployment options of the protected host and its traffic sources. Particular controls will need to address these deployments, so any overload control architecture must support or enable these requirements.

- REQ 17 Controls are required for scenarios where hosts receive service requests from an uncountable number of unknown sources as well as scenarios where hosts receive service requests from a known set of sources. A particular overload target may have a mixture of known, enumerable sources as well as an unknown, non-enumerable set of sources. In those cases, it shall be possible for the controls to inter-work, such that the SLAs for the enumerable sources are protected and the non-enumerable sources are managed on the basis of the aggregate workload from them all.
- REQ 18 A control shall inter-work with proxies, load balancing and load forking.

4.4 Application specific requirements

The general requirements above relate to overload controls that act between hosts that are nearest neighbours, i.e. host load control. Application specific requirements arise from the need for some applications to perform functions closely related to host load control. Many of these requirements are not naturally part of host load control, as they are specific to a particular protocol/application, but they may be implemented by an application using the basic infrastructure used for host overload control. An important feature of these application level controls is that they have an application defined granularity, whereas host load control only distinguished between flows on the basis of the host from which the request is received.

- REQ 19 A control shall enforce fair allocation of an overloaded processing resource between competing controlled application layer level streams of service requests, where such streams need not have a 1:1 correspondence to nearest neighbour hosts.
- REQ 20 A control shall be configurable by the service provider to enable the enforcement of SLAs for application level flows (to divide the capacity of an overloaded resource between competing application level requests flows according to agreed policies). This would allow a service provider to use the nearest neighbour load control infrastructure to protect well behaved flows when processing congestion is being caused by "badly behaved" flows, even though those flows are from the same nearest neighbour.
- REQ 21 A control shall facilitate an application to automatically limit ineffective service requests by detecting specific destination application layer names/addresses that are attracting a high reject rate and selectively controlling demand to them.

Annex A (informative): Comparison between RACS and nearest neighbour overload control

The Resource Admission Control Sub-system (RACS) [3] is responsible for regulating access to NGN resources. Initially, one might think that nearest neighbour load control is simply a specific subset of the functionality of RACS as its role is to regulate access to computational resource. In that case, it might seem paradoxical that nearest neighbour load control is additional to the architecture while RACS is an integrated part of the architecture. The key to resolving this apparent paradox is to understand the fundamentally different objectives of the two.

RACS is specifically designed to manage access to resources in the transport stratum that deliver service to end users, i.e. it regulates access to bandwidth to ensure the QoS targets of the user session are maintained. It is an integral part of the service delivery and is designed to allow end to end admission control to be achieved. The RACS infrastructure enables the resource for a particular end user session to be secured on an end to end basis before the session is admitted.

Nearest neighbour load control, on the other hand, is not integrated into the service logic. There is no end to end co-ordination of the nearest neighbour load control admission decisions, they are all local (independent) decisions extending only as far as nearest neighbours - the admission decisions are not service based - nearest neighbour load control is not accepting or rejecting a user session rather it is accepting or rejecting the request processing on that particular node.

To clearly demonstrate the difference, consider the fact that nearest neighbour load control may be deployed between servers that implement RACS. The differences between RACS and nearest neighbour load control are summarized in table 1.

Table 1: A comparison between RACS and nearest neighbour load control

RACS	Nearest neighbour load control
Request acceptance by RACS implies that resource is reserved for that request	Acceptance of request by nearest neighbour load control does not imply a reservation of any processing resource
Provides complex distributed functionality	Only affects nearest neighbours
Objective is service QoS control	Objective is processing infrastructure protection
May be a concrete instantiation - one can touch a physical system that only implements RACS functions	Only exists as a component in physical systems that implement NGN functions

History

Document history		
V2.0.0	January 2008	Publication