

Network Working Group
Request for Comments: 3532
Category: Informational

T. Anderson
Intel Labs
J. Buerkle
Nortel Networks
May 2003

Requirements for the Dynamic Partitioning of Switching Elements

Status of this Memo

This memo provides information for the Internet community. It does not specify an Internet standard of any kind. Distribution of this memo is unlimited.

Copyright Notice

Copyright (C) The Internet Society (2003). All Rights Reserved.

Abstract

This document identifies a set of requirements for the mechanisms used to dynamically reallocate the resources of a switching element (e.g., an ATM switch) to its partitions. These requirements are particularly critical in the case of an operator creating a switch partition and then leasing control of that partition to a third party.

Table of Contents

1. Definitions	2
2. Introduction	3
3. Dynamic Partitioning	6
4. Requirements	7
5. Security Considerations	9
6. Intellectual Property Considerations	9
7. Acknowledgements	9
8. Normative References	10
9. Informative References	10
10. Authors' Addresses	10
11. Full Copyright Statement	11

1. Definitions

In this document, the following definitions will be used.

Switching Element - A device that switches packets (e.g., an ATM switch or MPLS LSR) and whose resources can be divided into partitions, each of which can be independently controlled by a different controller.

Partition - A partition is a set of switching element (SE) resources. Partitions are also referred to as virtual SEs.

Active Partition - An active partition is a partition in which the resources are in use; either under the direct control of a separate controller or under internal policy-based control.

Controller - The entity responsible for controlling the operations of an active partition.

Static Partitioning - In static partitioning, no changes can be made to any active partition's resources without requiring a restart of that partition. Instances of repartitioning in which connections to controllers are disconnected before resources can be reallocated therefore fall into this category.

Dynamic Partitioning - In dynamic partitioning, an active partition's resources can be reapportioned without requiring a restart of the partition.

Frozen Partition - A frozen partition is an active partition that is in the process of being shutdown. A frozen partition's unused resources are relinquished, but all current connections are allowed to remain until removed by the controller. As connections close, the resources are returned to the SE.

Deterministic Partitioning - In deterministic partitioning, each active partition is given an allotted quantity of each resource. The usage of resources in one active partition does not influence the resources available to another active partition. All discussions in these requirements presuppose the use of deterministic partitioning.

Statistical Partitioning - In statistical partitioning, some or all resources are pooled among the active partitions, and allocations may be based on percentages or on some other metric. Discussion of statistical partitions is outside the scope of these requirements.

Proactive Notification - A proactive notification is a message sent from a SE to its controller at the time an event occurs. Specifically, if a SE asynchronously sends the controller a message when it is dynamically partitioned, we say that the SE has proactively notified its controller of the resource reapportionment.

Explicit Reactive Notification - In explicit reactive notification, the SE does not asynchronously send a message when dynamic partitioning occurs. Instead, the SE includes an explicit, resources-reassigned error code in the response to a subsequent request by the controller for an unavailable resource.

Implicit Reactive Notification - This is similar to an Explicit Reactive Notification except that the protocol does not contain any explicit resources-reassigned error codes. In this case, all that the SE can do is to indicate that some general, unknown error or generic resource error (i.e., some resource error problem has occurred but the exact cause is not specified) has occurred when the controller attempts to use unavailable resources. In such cases, the controller may attempt to determine whether a resource shortfall caused the error by using whatever messages are available through the control protocol to query available resources.

2. Introduction

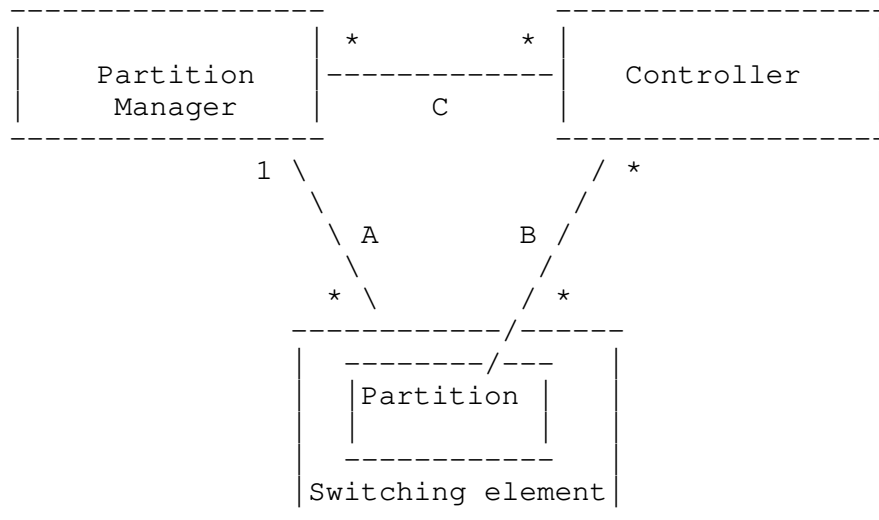
This document identifies the logical entities involved in the partitioning of switching elements. Furthermore, this document provides a set of requirements for the behavior of these logical entities as well as the protocols used by these logical entities to communicate with one another. A primary goal of the requirements specified herein is to allow the resources allocated to a partition to be increased or decreased while the partition is currently active (i.e., it has an active connection with a controller). This document is primarily intended to facilitate the partitioning of GSMP switches. However, while we believe that the logical entities and requirements specified here are necessary for the partitioning of non-GSMP switches and (longest prefix match) forwarders (e.g., routers), we do not believe that these requirements are necessarily sufficient for the partitioning of those devices.

Three logical entities are involved in the partitioning and control of a SE. First, a switching element (for the purposes of this document) is a device that "switches" packets, whose resources can be partitioned and whose partitions can each be controlled by a single controller. This partitioning also implies the ability to enforce this division of resources between competing partitions. Second, the

partition manager (PM) is a management entity that specifies the number of virtual SEs into which the SE should be partitioned and the resources to be allocated to each virtual SE. Lastly, a controller directs the use of the resources of one or more partitions to provide a set of services.

In the rest of this document, we will deal exclusively with logical entities although it is worth noting here that there are many possible mappings of logical entities to physical entities. For example, there may be multiple logical controllers running on a single physical processor (and for convenience we may refer to this processor as a physical controller). Conversely, a single logical controller could consist of processes running on multiple physical processors collaborating to provide proper control. Likewise, there may be multiple partition managers running on a single management workstation. A switching element may consist of one or more whole or fractional physical elements. For example, a SE may be a single whole physical switch (e.g., blade in a chassis), multiple whole physical switches (e.g., two blades in a chassis made to appear as a single logical entity), a single fraction of a physical switch (which would enable nested partitions), or multiple fractions of either the same or different physical switches (e.g., ports 1-3 on blade 1 and ports 2-4 on blade 2). Finally, any combination of these logical entities could theoretically be co-located on the same physical resources.

However, for many reasons, the physical realm often reflects this logical division of functionality. To facilitate this division, several protocols, such as MEGACO [RFC3015] and GSMP [RFC3292], exist that allow control functionality to be physically separated from switching functionality. Recently, some regulatory environments have mandated multi-provider access to a single physical infrastructure. To satisfy these regulations, a common use of partitioning will be for the owner of the SE to partition the SE into several virtual SEs and then to lease these to third parties. In this case, the PM will likely be physically separate from all of the controllers. For locality (and therefore ease) of management, SEs will be remotely configurable and thus the PM will be physically separated from the SE. The following illustration depicts this arrangement. The dashed lines indicate interactions between the entities and are labeled with the cardinality of the relationship between the entities.



Interaction A is one in which the PM partitions the SE and allocates resources to the partitions it creates. There is a one-to-many relationship between PMs and SEs. In order to support dynamic partitioning, this document will place certain requirements on proposed (or new) solutions in this space.

Interaction B is one in which the controller configures and manages an active partition. Current protocols implementing this interaction include GSMP [RFC3292] and MEGACO [RFC3015]. These protocols allow a many-to-many relationship between controller and partition.

Interaction C is one by which a PM and a controller could communicate to alter the nature of an active partition. There is a many-to-many relationship between PMs and controllers. For example, there are multiple PMs per controller in the case where a controller is managing two partitions from different SEs and there are multiple controllers per PM in the case where a SE has two partitions each managed by a different controller. Possible types of interactions between PM and controller include:

- A controller could request that the resources of one of its active partitions be altered; either increased or decreased.
- The PM could respond to a controller request for altered resource levels.
- The PM could request that a controller release resources currently allocated to one of its active partitions. This could involve the following types of request:

- A request to relinquish allocated, but currently unused resources. That is to put a freeze on additional use of the specified resources.
- A request to relinquish used resources.
- A request to relinquish an active partition. That is a request that a controller release control of an active partition.
- The controller's response to a PM request.

As far as the authors know, no proposed standard solutions currently exist for interactions of type C.

3. Dynamic Partitioning

Static repartitioning of a SE can be a costly and inefficient process. First, before static repartitioning can take place, all existing connections with controllers for the affected partitions must be severed. (This severing must always occur even if the resources to be reapportioned are not currently in use.) When this happens, the SE will typically release all the state configured by the controller. Then, the virtual SE must be placed in the "down" state while the repartitioning takes place. Once the repartitioning is completed, the partitions are placed in the "up" state and the controllers are allowed to reconnect to the partitions. Then, the controllers can reestablish state in those partitions. Thus, static repartitioning results in a period of downtime and a period in which the controllers are reestablishing state for affected partitions. Partitions of a SE that are not affected by a static resource reallocation need not be transitioned to the down state nor would controllers have to reestablish state with unaffected partitions.

Therefore, dynamic partitioning is to be preferred to static partitioning since it avoids the downtime and loss of state associated with static partitioning. However, a different set of potential problems exists for dynamic partitioning. Some questions to be answered include the following:

- How is the controller notified of an increase or decrease in resources?
- What should happen when the PM would like to decrease the resources allocated to a partition but those resources are in use?

4. Requirements

This document does not attempt to answer the preceding questions but instead defines a set of requirements that any solution to these problems MUST satisfy.

1. There MUST be a mechanism by which a PM can create virtual SEs on the SE and allocate SE resources to those virtual SEs.
2. SEs MUST ensure that controllers do not use more resources than those currently allocated to each virtual SE. Therefore, each control protocol MUST provide either an explicit reactive notification or an implicit reactive notification to indicate resource exhaustion.
3. Furthermore, there MUST be a mechanism by which a PM can partition all resources discoverable through GSMP (e.g., label tables). Partitioning of resources used by GSMP indirectly (e.g., CPU), resources used by non-GSMP switches, or resources (e.g., forwarding table entries) used by forwarding-based network elements MAY be supported.
4. If a PM instructs a SE to release resources allocated to an active partition and if any of those resources are currently in use, the SE MUST deny the PM's request. (Requirement #8 addresses the potential starvation issues raised by this requirement.)
5. Subsequent to a resource reallocation failure, the PM SHOULD make use of one or both of the capabilities described in requirements 6 and 7.
6. A PM SHOULD be able to tell a SE to make an active partition into a frozen partition.
7. A PM SHOULD be able to contact the controller to ask it to reduce its resource utilization.
8. The PM MUST be able to exercise "power on/off" type control of the virtual SEs that it has created. When the virtual power to an active partition is turned off, the partition becomes inactive and any controllers associated with that partition are disconnected. This capability allows a PM to resort to static partitioning when a controller is uncooperative about releasing resources. This requirement allows permanent starvation as a result of requirement #4 to be avoided.

9. During dynamic repartitioning, a SE MUST maintain all existing state associated with the partitions being modified.
10. Control protocols SHOULD NOT include any mechanism by which a SE can ask its controller to reduce its resource usage.
11. Control protocols MAY contain proactive resource notification messages by which a SE could instantaneously inform the controller of an increase or decrease in resources. (We do not specifically require control protocols to contain proactive notifications because all control protocols must already have explicit or implicit reactive notifications as mentioned in requirement #2).
12. A PM MAY directly inform a controller of a change in virtual SE resources rather than rely on the implicit resource exhaustion mechanism of the control protocol.
13. SEs MAY inform the PM of resource exhaustion on a particular partition.
14. A controller MAY ask the PM for further resources or a reduction in existing resources.
15. To support the automation of interaction between the PM and attached controllers, the PM MUST be able to determine from the SE the addresses of the controllers that are currently attached to a virtual SE. Additionally, the SE MAY allow the PM to determine which control protocol (and version thereof) is currently managing each active partition.
16. A SE MAY support the ability to have one virtual SE provide a service to another virtual SE within the same physical SE. For example, a SE may be configured to provide a virtual link between two virtual SEs. Furthermore:
 - a. There MUST be a mechanism by which the SE can inform the PM which of these partition-to-partition services are provided by the SE.
 - b. There MUST be a mechanism by which the PM can configure the available partition-to-partition services.
 - c. If the configuration of a partition-to-partition service results in a virtual port being added/removed from a virtual SE, the SE MUST notify all controllers attached to that virtual SE (assuming that the corresponding control protocol supports such notifications).

17. There MUST be a mechanism by which a PM can query a SE to determine the resources of that SE, the partitions currently configured on that SE and the resources allocated to each partition.

5. Security Considerations

Only authorized PMs MUST be allowed to dynamically repartition a SE. Therefore, SEs MUST use a secure process by which an authorized entity may instruct the SE as to which PM should control it. This instruction MAY specify the PM explicitly or MAY specify the use of a (discovery) protocol to dynamically locate the PM. Similarly, only the PM (or an authorized agent of the PM) that is authorized to partition a SE MUST be allowed to contact controllers to request that they decrease their resources or inform them that their resources have been increased. Likewise, the PM MUST verify and authenticate that any requests for additional/fewer resources for a virtual SE have come from a controller authorized to control the specified virtual SE.

6. Intellectual Property Considerations

The IETF takes no position regarding the validity or scope of any intellectual property or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; neither does it represent that it has made any effort to identify any such rights. Information on the IETF's procedures with respect to rights in standards-track and standards-related documentation can be found in RFC 2026. Copies of claims of rights made available for publication and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementors or users of this specification can be obtained from the IETF Secretariat.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this standard. Please address the information to the IETF Executive Director.

7. Acknowledgements

The authors would like to acknowledge the contributions of Avri Doria and Jonathan Sadler to this document.

8. Normative References

[RFC2119] Bradner, S. "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.

[RFC3292] Doria, A., Hellstrand, F., Sundell, K. and T. Worster, "General Switch Management Protocol (GSMP) V3", RFC 3292, June 2002.

9. Informative References

[RFC3015] Cuervo, F., Greene, N., Rayhan, A., Huitema, C., Rosem, B. and J. Segers, "Megaco Protocol 1.0," RFC 3015, November 2000.

10. Authors' Addresses

Todd A. Anderson
Intel Labs
JF2-60
2111 NE 25th Avenue
Hillsboro, OR 97124 USA

Phone: +1 503 712 1760
EMail: todd.a.anderson@intel.com

Joachim Buerkle
Nortel Networks Germany GmbH & Co. KG
Hahnstrasse 37-39
60528 Frankfurt

Phone: ++49 (0)69 6697 3281
EMail: joachim.buerkle@nortelnetworks.com

11. Full Copyright Statement

Copyright (C) The Internet Society (2003). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Acknowledgement

Funding for the RFC Editor function is currently provided by the Internet Society.

