

Internet Draft  
Expiration: May 1998  
File: [draft-berson-classy-approach-01.txt](#)

S. Berson  
ISI  
S. Vincent  
ISI

## Aggregation of Internet Integrated Services State

November 21, 1997

### Status of Memo

This document is an Internet-Draft. Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

To learn the current status of any Internet-Draft, please check the linebreak "1id-abstracts.txt" listing contained in the Internet-Drafts Shadow Directories on ds.internic.net (US East Coast), nic.nordu.net (Europe), ftp.isi.edu (US West Coast), or munnari.oz.au (Pacific Rim).

### Abstract

The Internet Integrated Services (IIS) architecture[2] has a fundamental scaling problem in that per flow state is maintained at all routers and end-systems supporting a flow. This paper examines the use of aggregation as a technique to reduce the amount of state needed to provide IIS. In our approach, routers at the edge of a region doing aggregation keep detailed IIS state, while in the interior of this region, routers keep a greatly reduced amount of state. Packets will be tagged at the edge with scheduling information that will be used in place of the detailed IIS state. The aggregation scheme described will allow large scale deployment of IIS without overloading routers with state and associated processing.

## 1. Introduction

In order to deploy Internet Integrated Services (IIS), additional resources are needed in the routers to keep state and to process packets. As currently defined [4,2], this state is on a per session basis, where a session is a unicast or multicast destination and optionally a port. IIS state is stored at all routers between the source(s) and destination(s) of a session. Thus, widespread use of IIS would mean a large amount of state at routers in the core of the Internet. This large amount of state and related processing could overwhelm routers, so IIS state models that do not require per-session state at all routers need to be explored.

Several types of per-session state are used with integrated services. First, there is scheduling state that consists of the different traffic service queues for packet forwarding. Second, there is packet classifier state. Classifier state is used to assign each arriving packet to a packet scheduling queue for forwarding. Third, there is the state for the setup protocol, e.g. state carried in RSVP PATH and RESV messages. Finally, current multicast routing algorithms also store state, either per source and session (e.g. DVMRP, PIM-DM) or per session shared-tree (e.g. CBT, PIM-SM). While approaches to aggregation of routing state are similar to aggregation of IIS state, they are beyond the scope of this paper.

Per session state in a router has two harmful effects. First, additional state consumes memory, and second, additional state requires additional CPU cycles to process the state. Depending on the speed with which the memory needs to be accessed, different types of memory can be used. State that is required in order to process each packet (e.g. classifier state and scheduling state) needs to be stored in fast (i.e. cache) memory, while other state (e.g. setup protocol state) can be stored in slower memory. Since fast memory is substantially more expensive than slow memory, and since accesses to fast memory are on the router fast path, our primary goal is to reduce the classifier and scheduler state. Reducing the amount of setup protocol state is an important, but secondary goal.

Any aggregation scheme entails some tradeoffs. Keeping full integrated services state at all routers allows the resources dedicated to the flow of packets from each admitted session to be isolated from the resources for packets from all other sessions. This isolation means that resources reserved and needed for one session will be used only on the reserving session. The disadvantage of doing aggregation is the loss of some portion of this flow isolation, since with aggregation many flows would share the same traffic service class. Our goal is to provide Internet Integrated Services to a network while using a greatly reduced amount of state.

Berson, Vincent

Expiration: May 1998

[Page 2]

We make several assumptions:

1. We assume that an explicitly definable region (e.g. one or more contiguous autonomous systems), called an "aggregating region", will aggregate and will implement supporting mechanisms at its boundary points. We define an "ingress" of an aggregating region as a router where a packet for a specific unicast or multicast session enters the region, and an "egress" as a router where a packet leaves the region. Also, we define an "interior router" as any router in the aggregating region that is on the data path for the session and is not the ingress or egress.
2. We assume that the choice to aggregate and the mechanism used to do so is an intra-domain issue, not an inter-domain issue and therefore there can be variability across regions so long as at the boundaries routers continue to pass along adequate messages to support the setup protocol upstream and downstream.
3. We assume current multicast and unicast routing within the aggregating region; but with no other enhancements.

## **2. Framework**

Figure 1 shows the architecture of an Internet integrated services device. There are three parts of the device, the scheduler, the setup protocol engine, and the classifier. The scheduler is responsible for enforcing a quality of service for each of the flows (i.e. one flow per scheduling queue). The scheduler achieves this with a scheduling algorithm and a policing algorithm. The scheduling algorithm decides which queue the next packet to be forwarded will come from, while the policing algorithm decides which packets are discarded or put into a different scheduler queue in case of congestion. Typical types of scheduling algorithms are based on fair queueing and/or priority. The classifier is responsible for assigning packets to queues. The classifier accomplishes this by looking at fields (e.g. destination address) in the packet headers. A table in the classifier tells which packet header values assign the packet to each scheduler queue. The setup protocol engine is responsible for setting up the classifier state and scheduler queues. The setup engine receives requests that include a quality of service and a packet classifier entry. Upon receiving a message, the setup engine does admission control to determine if sufficient resources are available to make the reservation. If resources are available, a new scheduling queue is set up in the scheduler for the requested resources, and the classifier entry is installed in the classifier.

Berson, Vincent

Expiration: May 1998

[Page 3]

[FIGURE]

Figure 1: Internet Integrated Services device

## 2.1 Reducing the amount of state

Reducing the amount of state means that some of the basic functions of integrated services must be limited. Since each network service provider will want to make decisions on how to provide integrated services based on their provisioning and traffic patterns, a detailed review of how IIS state is used, is offered.

Scheduler classes are the unit of assignment of network resources. Scheduler state typically takes the form of different traffic classes. These traffic classes are assigned different priorities, link shares, and policing policies to provide a certain level of service. In the basic IIS model, each traffic flow gets its own traffic class, meaning that each flow has certain resources dedicated to it. By reducing the number of traffic classes, a coarser granularity of resource assignment is done in that several flows will be assigned to the same scheduler queue. A consequence of the coarser granularity is that the aggregated traffic gets a certain level of resource usage, but it is impossible inside an aggregating region to isolate flows in the same traffic class from each other. The lack of isolation can be mitigated by adequate provisioning of network resources and by policing at aggregating region borders.

Classifier state is used to assign packets to traffic classes. Each arriving packet is classified into a traffic class at each node according to state established by the setup protocol. In the complete absence of classifier state, no service differentiation is possible. However a packet can be classified at the ingress, and the results of that classification can be encoded in each packet. This reduced classifier state might be as simple as one or more bits segregating different classes of traffic, or alternatively, an encapsulation header. By encoding a small amount of classifier state in the packet, only a very simple classification is needed at interior routers.

Finally setup protocol (e.g. RSVP) state is used to do admission control, provide policing, allow reclassification, and allow merging of reservation messages. The first two of these issues, admission control and policing, are general issues involving both unicast and multicast traffic, while the last two, reclassification and reservation merging, are only issues with multicast traffic.

Berson, Vincent

Expiration: May 1998

[Page 4]

Setup protocol state is used in admission control. Admission control with full IIS state involves computing how much network resources are already committed and whether there are sufficient resources to admit a new flow with certain resource requirements. Lacking setup protocol state, admission control in the aggregating region cannot be based on reservation state per flow. Instead, admission control must be based on aggregated state information at each node. More specifically, state will typically consist only of measurements, and so admission control must be measurement based.

Setup protocol state is used for traffic policing. Policing is done at two points in an IIS network, at a traffic merge point, and at a traffic split point. The traffic merge point is where the traffic from two or more sources for the same session merge. For a shared style reservation, the traffic from all the sources needs to be policed to the reservation parameters at traffic merge points. The second place that policing is needed is at traffic split points. A traffic split point is where multicast traffic has multiple outgoing interfaces. Because of reservation merging and because some outgoing interfaces may have no reservation, traffic at a split point needs to be policed on each outgoing interface to the proper reserved values for that interface. Lacking setup protocol state, no per flow policing can take place inside an aggregating region. This implies that all per flow policing must take place at the edges of the aggregating region.

Setup protocol state is used to allow reclassification of packets from reserved to best effort. This reclassification is necessary at traffic flow split points where there are one (or more) outgoing interfaces with reservations and one (or more) outgoing interfaces without reservations. At these points packets forwarded to the outgoing interfaces with no reservations need to be reclassified as best effort.

Finally, setup protocol state is used to merge reservation messages coming from the egresses to the ingress. With setup protocol state for each session, and hop by hop processing of reservation messages, each router can limit the number of reservation refreshes it sends. Without setup protocol state, each reservation message would be sent to the ingress, requiring a large amount of processing for any large multicast groups.

## 2.2 Basic aggregation architecture

We propose a scheme with a constant amount of scheduler state at all routers, a constant amount of classifier state at interior routers, and no unicast setup protocol state at interior routers.



Berson, Vincent

Expiration: May 1998

[Page 5]

In our approach, a fixed number of traffic service classes are defined and configured at all routers in an aggregating region. Each traffic service class is subject to admission control at the ingress to the region. Traffic policing is done at the ingress for unicast traffic, and at the ingress plus at traffic split points for multicast traffic. Since only a fixed number of traffic classes are available, the scheduler state at all routers is constant.

Data packets are classified and "tagged" on entry to the aggregating region with some identifier indicating to which aggregated traffic class the packet should belong. This identifier can consist of some bits in the packet header (e.g. type of service bits) or the packet could be encapsulated. The tag (or encapsulation) encodes the traffic service class that this packet will receive across the aggregating region. Tagging or encapsulating each packet bounds the amount of classifier state in the interior of the aggregating region.

Full setup protocol state is stored at the borders of the aggregating region, but only multicast setup protocol state is needed in the interior of the region. When a reservation message arrives at an ingress to the region (from an egress), and that reservation passes admission control, then packets from the flow associated with the new reservation are assigned to one of the configured traffic service classes.

Admission control will be measurement-based, and, since there is no setup protocol state in the interior, admission control will be done at the ingress for a session. But the decision will be based on congestion measurements within the aggregating region. Thus each node in the aggregating region, upon receiving an admission control setup protocol message, will make a local congestion decision. This local congestion decision will be made on a threshold basis where new reservations are not admitted if the existing load plus the expected additional load from the new flow is greater than a (possibly dynamic) threshold.

To implement the local congestion decision, each node keeps an estimate of its current load per traffic class. As an admission control message traverses the path between ingress and egress, each interior router must admit that reservation. Each router that provisionally admits the reservation should factor the new reservation into its current traffic estimate. To admit the flow, an estimate of the amount of the reserved traffic due to this reservation must be made, e.g. based on peak rate specified in the admission control message. If some router does not admit the reservation, the admission control message is marked appropriately

Berson, Vincent

Expiration: May 1998

[Page 6]

by the rejecting router and the message is forwarded. When the admission control message arrives at the egress, the egress node checks the message to see if any router rejected the message. If some router rejected the reservation, then the egress will forward the message to the ingress which will reject the reservation and send a reservation error message.

In summary, our approach provides for reduction of packet scheduler state, packet classifier state, and setup protocol (e.g. RSVP) state in the following ways. A constant amount of packet scheduler state at all routers in an aggregating region is achieved by using preconfigured traffic service classes. A constant amount of packet classifier state for unicast traffic at all interior routers in the aggregating region is accomplished by tagging or encapsulating each packet with its traffic service class. No setup protocol state at interior routers for unicast traffic is achieved by doing admission control and policing only at the edges of the aggregating region. For multicast traffic, full RSVP state is stored at all routers, but classifier state is stored only at traffic split points and only on those outgoing interfaces with no reservation.

### 2.3 Multicast

We would like to apply the above approach to multicast traffic. However, there are some fundamental differences between unicast and multicast traffic that results in additional setup protocol state being desirable. Since traffic with multiple destinations is being classified and tagged at the ingress router, a tagged multicast packet will receive reserved service everywhere in the cloud, including branches for which there are no reservations. This all-or-nothing reservation property, caused by tagging at the ingress, complicates aggregating state for multicast sessions. The complications derive from two features of multicast sessions, "heterogeneity" and "dynamic group membership". Heterogeneity refers to the situation where different branches of a multicast distribution tree have different reservations, including the case where some branches have no reservation. Unreserved packets receiving reserved service can adversely affect packet flows that properly have a reservation. Without additional state, there is no way to determine which packets have a legitimate reservation.

The other multicast feature causing complications is dynamic group membership. Dynamic group membership means that the members of a multicast session can be changing over time by having new receivers join, and old receivers leave. Due to the all-or-nothing property, a new best effort receiver joining a multicast session with a reservation in place will be receiving reserved

Berson, Vincent

Expiration: May 1998

[Page 7]

service. Admission control will not have been done for that new receiver, and so those reserved packets to the new receiver can cause congestion for other properly reserved packets. Thus a best effort receiver can disrupt reserved traffic.

Our approach to aggregation and multicast is to keep setup protocol state for multicast sessions. Since setup protocol state is not on the router fast path, this option is reasonable since scheduler state and classifier state are still minimized. Some additional classifier state is needed, but only for sessions with heterogeneous branches, and only at those nodes with heterogeneous outgoing interfaces. This classifier state can be established dynamically from setup protocol state and routing state. When a node discovers that it has heterogeneous outgoing interfaces, classifier entries are established on the best effort branches. The packets that the classifier matches have the reserved markings removed and new markings inserted indicating that the packets are to be forwarded as best effort.

Having setup protocol state for multicast can also solve the problems mentioned in [section 2.1](#) with implosion of reservation and ADREP messages, and with policing. Also note that an additional optimization would be to store setup protocol state for a session only on routers with multiple outgoing interfaces for that session.

#### 2.4 Policing in interior

There is still an issue with state aggregation and shared (e.g. RSVP wildcard or shared explicit) style reservations at traffic merge points. At a merge point, the traffic entering the node may conform to the reservation, but the traffic exiting the node may be non-conforming due to multiple senders. With setup protocol state, the traffic would be policed at the merge node. Since there is no state in the interior of the network, it is impossible to police the traffic of an individual flow in the interior. The traffic will eventually be policed at the next ingress, but may interfere with actual reserved traffic between the merge point and the egress. Since the traffic will eventually be policed, there is no "free bandwidth" for a user trying to exploit this feature. However there is a security problem where a malicious user could tie up reserved bandwidth traffic in a transit network. To protect a network from this sort of abuse, the reservation of an abuser could be cancelled by the egress when excess traffic at an egress is detected.

Berson, Vincent

Expiration: May 1998

[Page 8]

### 3. RSVP operations/extensions

In this section we describe how this aggregation scheme would be applied to RSVP. We assume that a measurement based admission control algorithm is defined, and we describe how to utilize this system to implement aggregation. The basic issues are gathering the measurement data, keeping track of which reservations are new, and describing the appropriate data structures.

Our approach to collecting measurement information for admission control is to originate the state collection from the ingress upon receiving an RSVP reservation message (RESV). The basic exchange of messages is shown in figure 2. Note that no processing is needed in the interior of the aggregating region on the RESV or the preceding PATH messages. But when the ingress receives a new RESV message (i.e. a message for which it has no traffic control state), the ingress initiates admission control. For admission control, the ingress sends an ADmission REQuest (ADREQ) message hop-by-hop through the interior of the aggregating region to the egress (known from the previous RSVP hop field in the RESV message). At each hop through the region, the node attempts to admit the reservation. If admission control succeeds at a node, the ADREQ message is forwarded unchanged to the next hop. If admission control fails at a node, a failure object is appended to the ADREQ message and then the message is forwarded to the next hop. When the egress receives the ADREQ message, it checks if any interior node appended a failure object to the message. If there is no rejection information in the ADREQ message, an ADMISSION REPLY (ADREP) message with an ADMITTED status flag is sent back to the ingress. If there is a failure object in the message, an ADREP message with a REJECTED status is sent back to the ingress. The ingress will treat the ADREP message with a REJECTED status as an admission control failure and a RSVP reservation error will be sent out.

[FIGURE]

Figure 2: Admission control exchange of messages

There are three other possible approaches to aggregated admission control that may be useful in certain types of network environments, but which have serious limitations in general. One approach is to put the admission control request and responses in RSVP reservation messages. This approach is conceptually simpler than the above scheme, but assumes that the reverse route (i.e. egress to ingress) through the aggregating region is known. There are cases where this holds, most notably in link state protocols, but in general this is not the case.



Berson, Vincent

Expiration: May 1998

[Page 9]

The second approach to collecting admission control information is an RSVP path message which travels on the forward route from data ingress to data egress. When this message is processed at each interior node, admission control information is collected. The accumulated admission control information is then included in a corresponding RSVP reservation message from egress to ingress. The main disadvantage with this scheme is that the service class is generally not known only from path information. If a region offers only a small number of services, it would be possible to include admission control information for all of the services in the path message, but this could be a large amount of information in general.

The third approach to collecting admission control information is by out-of-band communication, e.g. as part of routing protocols. In this case, the ingress site can do the admission control based on the out-of-band communication. Further work is needed in this area.

The remainder of this section shows the additions to RSVP to support our approach; using separate ADREQ and ADREP messages in the aggregating region.

### 3.1 RSVP messages

Two new RSVP message types would be needed for RSVP messages, admission request (ADREQ) and admission reply (ADREP). The ADREQ is used by an ingress to set up the reservation across the aggregating region, while the ADREP is used by the egress to report that the reservation has succeeded or not.

In addition to new message types, an admission control object (ADMISSION), and an admission control status (STATUS) object types are needed. The ADMISSION object and zero or more STATUS objects are included as part of a ADREQ or ADREP message. The ADMISSION object contains a common object header, the IP address of the ingress, a handle from the ingress node, and some flags. The handle is included by the ingress node on initiating aggregated admission control and is used to uniquely identify reservation state on the ingress. The STATUS object has a common header, an IP address, and a load status. The IP address is the address of the node supplying the load report, while the load status is the measured load from a node.

The format of the ADREQ and ADREP messages are very simple and are as follows:

```
<ADREQ> ::= <Common Header> <FLOWSPEC> <ADMISSION>
           <status_report>
```

Berson, Vincent

Expiration: May 1998

[Page 10]

```
status_report ::= <empty> | <status_report> <STATUS>
```

```
<ADREP> ::= <Common Header> <ADMISSION> <status_report>
```

When an ingress receives a new RSVP RESV message (and local admission control is successful), the ingress initiates aggregated admission control by sending an ADREQ message. The FLOWSPEC from the RESV message is used in the ADREQ message, as well as an ADMISSION object containing a handle chosen by the ingress.

On receiving an ADREQ message, an interior router will perform admission control based on the measured load and the FLOWSPEC object. If the router is too congested to accept a flow, the rejecting router adds its address and measured load information to a STATUS object which is appended to the ADREQ message. The ADREQ message eventually arrives at the egress with a (possibly empty) list of STATUS objects, containing the address and measured load of nodes that are congested. The egress will set the flags in the ADMISSION object and then return the ADMISSION object and the list of STATUS objects in an ADREP message unicasted to the ingress. The ingress will locate the reservation state block by using the handle in the ADMISSION object. Then the ingress will check the ADREP message for any status objects. If there are any STATUS objects the ingress will reject the reservation and send a reservation error message downstream. If there are no STATUS objects in the ADREP message then the ingress will accept the reservation and the RESV message will be forwarded upstream.

Note that if either an ADREQ or an ADREP message is lost, then the ingress behaves as a router that lost an RSVP reservation message. The ingress will wait until the reservation message is refreshed and then send a new ADREQ. It would also be possible to resend an ADREQ after a configurable amount of time.

```
struct {
    Object_header  header;
    struct in_addr ingress;
    Handle         handle;
    int            flags;
} ADMISSION;
```

```
#define ADMITTED      1
#define REJECTED      0
```

```
struct {
    Object_header  header;
    struct in_addr sess;
    Load          load;
```

Berson, Vincent

Expiration: May 1998

[Page 11]

```
} STATUS;  
  
#define RSVP_ADREQ      11  
#define RSVP_ADREP     12
```

#### 4. Reserved traffic congestion

One last issue is how to deal with traffic congestion in reserved traffic classes. Congestion is typically detected in an ADREQ/ADREP message. Similar to the load threshold for admission, there can be a load threshold for congestion which is higher.

There are two ways that admitted load can cause high levels of congestion on a router, measurement based admission control failure and link failure. Measurement based admission control failure occurs due to past load not being representative of the future. For example, if many idle reserved sources become very active at the same time, some links may become overloaded. Link failure is caused by a link being declared down and a new set of routes are generated. The traffic that was traversing the failed link may now be added to an already heavily loaded link. In both of these cases, links may become overcommitted.

For small amounts of congestion, reserved traffic can preempt best effort traffic, i.e. best effort traffic packets are dropped, but new admission control requests are rejected. For larger amounts of congestion, reserved packets need to be dropped or reservations need to be preempted. For link failure in a per-flow IS state network, the setup protocol would typically use reservation preempting (while admission control mistakes are expected not to happen). With aggregated IS state, the point at which the routing changes may not have setup protocol state. Thus a node typically cannot distinguish an overload caused by a link failure from one caused by transient heavy congestion. The response is the same for either cause of congestion; Ingress nodes in this case will need to preempt reservations.

#### 5. Acknowledgements

This Internet Draft is the result of discussions with many people particularly Bob Braden, Bob Lindell, Deborah Estrin, and Daniel Zappala.

#### REFERENCES

- [1] Braden, R., Zhang, L., Berson, S., Herzog, S., and Jamin, S., "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification," [RFC 2205](#), September 1997.

Berson, Vincent

Expiration: May 1998

[Page 12]

- [2] Braden, R., Clark, D., and Shenker, S., "Integrated Services in the Internet Architecture: an Overview," [RFC 1633](#), June 1994.
- [3] Floyd, S., and Jacobson, V., "Link-sharing and Resource Management Models for Packet Networks," IEEE/ACM Transactions on Networking, Vol. 3 No. 4, pp. 365-386, August 1995.
- [4] Jamin, S., Shenker, S., and Danzig, P., "Comparison of Measurement-based Admission Control Algorithms for Controlled-Load Service," Infocomm '97.
- [5] Rampal, S., "Flow Grouping for Reducing Reservation Requirements for Guaranteed Delay Service," Internet Draft, December 1996.

#### Security Considerations

Security considerations have not been addressed in this draft.

#### Author's Address

Steven Berson  
USC Information Sciences Institute  
4676 Admiralty Way  
Marina del Rey, CA 90292

Phone: +1 310 822 1511  
EMail: berson@isi.edu

Subramaniam Vincent  
USC Information Sciences Institute  
4676 Admiralty Way  
Marina del Rey, CA 90292

Phone: +1 310 822 1511  
EMail: svincent@isi.edu



