Internet Engineering Task Force                          D. Partain (ed)
INTERNET-DRAFT                                            G. Karagiannis
Expires December 2002                                       P. Wallentin
                                                             L. Westberg
                                                                Ericsson
                                                               June 2002

**Resource Reservation Issues in Cellular Radio Access Networks**
**draft-westberg-rmd-cellular-issues-01.txt**

Status of this Memo

Abstract
   This memo describes resource management issues that are relevant to
   the use of IP transport in cellular radio access networks (RANs).
   The document describes the particular characteristics of these kinds
   of networks, the requirements applicable to a resource reservation
   scheme in a cellular RAN, and provides a brief analysis of the
   applicability of existing solutions to this problem space.

Table of Contents

**1**.  **Introduction**

   The rapidly growing popularity of IP and its flexibility make it a
   good candidate to be used for transmission in cellular networks.

   Using IP-based transport on the wired transmission links in the
   cellular networks gives operators an opportunity to upgrade their
   transport network to a packet-based one. When compared with a
   traditional STM-based system, the gain is seen in the statistical
   aggregation of traffic that can be done.  This results in increased
   transmission efficiency and reduced leasing cost for the operator.

   A radio access network (RAN) provides the radio access (e.g., GSM,
   CDMA, or WCDMA) to mobile stations in a cellular network.  To
   accomplish this, radio frames are transported on the wired links
   between different cellular-specific nodes in the RAN. The majority of
   the traffic (up to 100%) is delay-sensitive traffic.

   The cellular user is unaware of the IP-based transport network
   underneath, and the service must work the same way as the user has
   come to expect the cellular services to work in an STM-based
   transport network. In addition to this requirement, the situation is
   further complicated by the fact that the RAN is large in terms of its
   geographic size, the number of inter-connected nodes, and the
   proportion of real-time traffic.

   To satisfy the above requirements, it is absolutely critical that we
   have a simple and scalable bandwidth resource management scheme for
   real-time traffic in this type of network.

   In order for real-time services to function satisfactorily in an IP-
   based RAN, we need to ensure that there are adequate transport
   resources on the links available in the RAN to handle this particular
   instance of the service (e.g., a phone call).  Note that, in rest of
   this draft, whenever the term "resources" is used, it refers to
   bandwidth on the links.

   If the RAN is bandwidth-limited and does not use any mechanism to
   limit the usage of the network resources, congestion might occur and
   degrade the network performance.  For example, speech quality might
   degrade due to packet losses.

**[2](#). Terminology**

The following terminology is used in this memo:

 * BSC:  Base Station Controller

 * RNC:  Radio Network Controller

 * MSC:  Mobile Services Switching Center

 * GGSN: Gateway GPRS Support Node

 * SGSN: Serving GPRS Support Node

 * GPRS: General Packet Radio Service, the packet-switched access
         scheme and service provided in GSM.

 * GSM:  Global System for Mobile Communications

 * UMTS: Universal Mobile Telecommunications System, the third
         generation (3G) mobile system based on WCDMA and GSM,
         specified by 3GPP (third generation partnership project).

 * radio frame: a short data segment coded/decoded and
         transmitted/received by the radio base station.
         It originates from a mobile station or the BSC/RNC.

 * WCDMA: Wideband Code Division Multiple Access, the
         radio transmission technology used in UMTS

**[3](#). Background and motivation**

   The context of the issues described in this document is the cellular
   radio access network (RAN).  This section briefly discusses two
   examples of radio access networks (for GSM - Global System for Mobile
   Communication - and for WCDMA - Wideband Code Division Multiple
   Access) and then outlines the motivation for this memo.

**[3.1](#). IP transport in radio access networks**

   This section introduces the radio access network and its use of IP
   transport. A radio access network (RAN) provides the radio access

(e.g., GSM, CDMA, or WCDMA) to mobile stations for a cellular system.
The boundaries for a RAN are the radio transmission and reception
access points (terminated by base stations) at one end and, at the
other end, the interfaces to the gateways (e.g., MSC and SGSN/GGSN),
which in turn provide connections to the fixed public network.

The radio access network consists of a number of nodes as shown in
the Figure 1 below.

```
                          IP                  IP
                        <--->               <--->

                |---------|     |---------|     |---------|
Mobile     v    |  Base   |     |         |     |  MSC/   |     Fixed
stations   |--| station |-----| BSC/RNC |-----|  SGSN/  |--- public
       ^        |         | ^   |         | ^   |  GGSN   |     network
       |        |---------| |   |---------| |   |---------|
       |                    |               |
    Wireless             Wired           Wired
    interface            interface       interface


             Radio access network
        <------------------------------>
```

        Figure 1:  Typical radio access network and its boundaries

The base station provides the radio channel coding/decoding and
transmission/reception function to and from mobile stations in its
coverage area, which is called a cell.

The BSC/RNC controls a number of base stations including the radio
channels and the connections to mobile stations.  For a WCDMA radio
access network, the BSC/RNC provides soft handover combining and
splitting between streams from different base stations belonging to
the same mobile station.  Furthermore, the BSC/RNC is also
responsible for the allocation of transport resources within the
radio access network.  The transport is either between the base
station and the BSC/RNC, between multiple BSC/RNCs, or between the
BSC/RNC and the MSC/SGSN.


The MSC provides, among others things, support for circuit-switched
services towards mobile stations including mobility management,
access control and call control as well as interworking with external
circuit-switched networks such as the public switched telephony

network (PSTN).  The SGSN/GGSN provide, amongst other things, support
for packet switched services towards mobile stations, including
mobility management, access control and control of packet data
protocol contexts. In addition, the GGSN provides interworking with
external packet-switched networks such as the public Internet.

The radio access network consists of potentially thousands of base
stations and a significant number of BSCs/RNCs.  The traffic volume,
in terms of voice-traffic, generated by these nodes can vary from a
few up to fifty voice calls per base station, and up to several
thousand simultaneous calls (Erlang) per BSC/RNC site.  Therefore, a
router in the network has to handle many thousands of simultaneous
flows.

The transmission between base stations and the BSC/RNC is usually on
leased lines, and this part (due to the wide area coverage of the
cellular network) is usually extremely expensive when compared to the
cost of transmission in the backbone.  Due the number of base
stations, the cost for these leased lines can be quite significant.
Dimensioning using over-provisioning might therefore be prohibitively
expensive, and it is unlikely that the network will be dimensioned
without using the statistical properties of traffic aggregation
(e.g., Erlang trunking).

```
   |--------|
   | Upper  |<---------------------------------------------> towards the
   | layers |                                                MSC/SGSN/GGSN
   |--------|                                    |---------|--------|
   | Radio  |                                    | Radio   |        |
   | Link   |<----------------------------------->| Link    |        |
   | Layer  |                                    | Layer   |        |
   |--------|                                    |---------|        |
   |        |                                    |Radio    |        |<-
   |        |                                    |Physical | Frame  |
   |        |                                    |Layer    | transp.|
   |        |    |-------|--------|              |---------| layer  |
   |        |    |       | Frame  |<-------------->| Frame  |        |
   |        |    |       | transp.|              | transp. |        |
   |Radio   |    | Radio |--------|              |---------|--------|
   |Physical|    | Phys. |  UDP   |<-------------->|  UDP    | UDP    |
   |Layer   |<->| Layer |--------|   |--------|   |---------|--------|
   |        |    |       |  IP    |<->|  IP    |<->|  IP     | IP     |
   |        |    |       |--------|   |--------|   |---------|--------|
   |        |    |       | Link   |   | Link   |<->| Link    | Link   |
   |        |    |       | Layer  |<->| Layer  |   | Layer   | Layer  |
   |        |    |       |--------|   |--------|   |---------|--------|
   |        |    |       |Physical|   |Physical|   |Physical |Physical|
   |        |    |       | Layer  |<->| Layer  |<->| Layer   | Layer  |
   |--------|    |-------|--------|   |--------|   |-----------------|
    Mobile        Base station         Router          BSC/RNC
    Station   ^                   ^              ^
              |                   |              |
        Wireless              Wired          Wired
        interface            interface      interface


                      Radio access network
          <----------------------------------------------------->
```

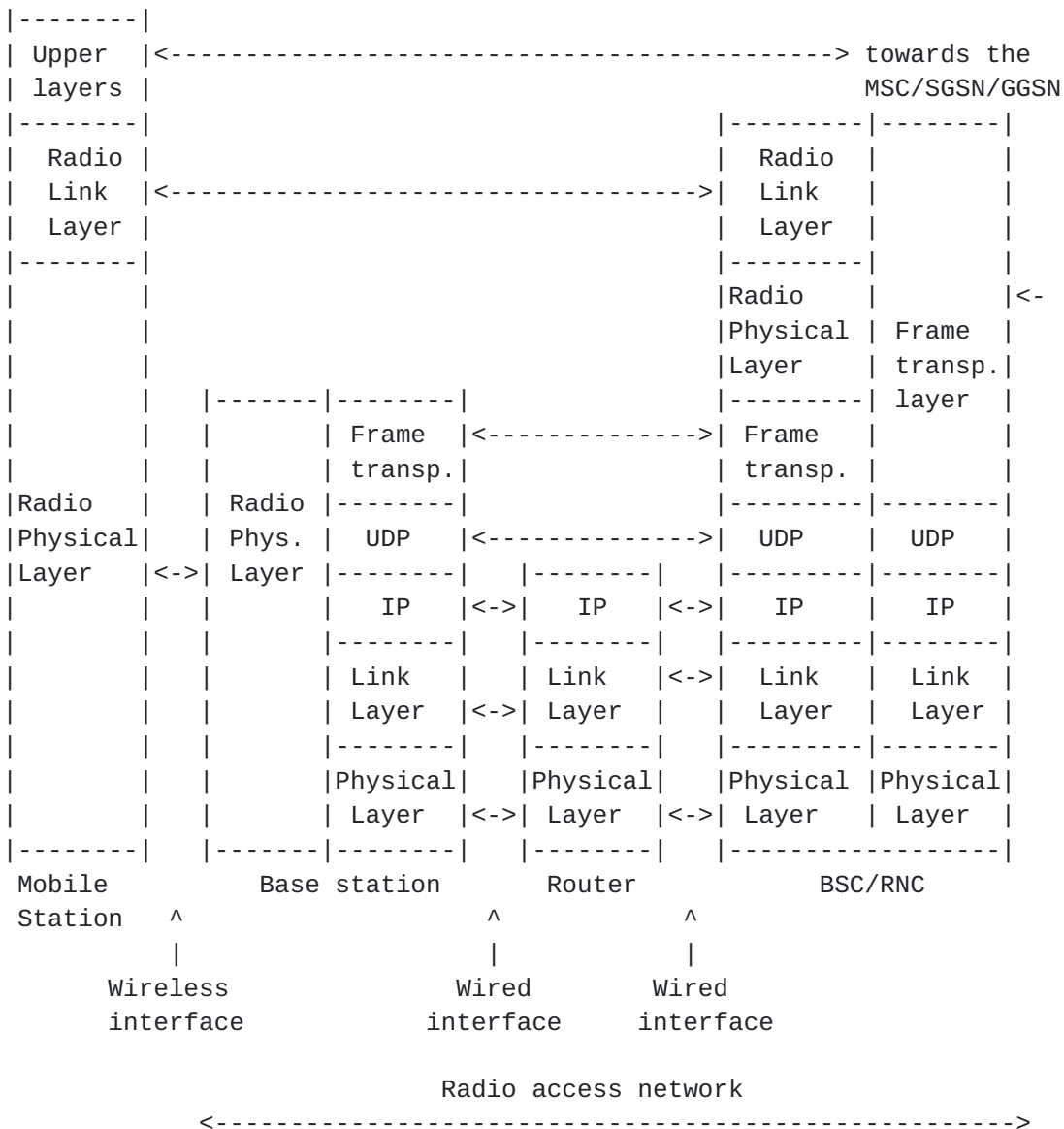             Figure 2:  Example of a protocol stack in the
                     radio access network (simplified)

   Figure 2 shows a simplified example of protocol layering when using
   IP transport in the radio access network.

   The radio physical layer performs radio transmission and reception
   functions, including soft handover splitting and combining in case of
   WCDMA.

   The frame transport layer is used to transmit radio frames between

the base station and the BSC/RNC.  A radio frame is a short data
segment coded/decoded and transmitted/received by the radio base
station at a given point in time. The radio frames must be delivered
in a timely fashion with limited delay.  Otherwise, the frames are
discarded by the base station or RNC/BSC. The traffic is therefore
very sensitive to delays.

The radio link layer performs segmentation/re-assembly,
retransmission and multiplexing/scheduling functions as well as radio
resource control.  The adaptation of user data performed by the radio
link layer depends on the type of radio channel and the type of
service. In one case, that very small radio frames might be
transferred, while in other cases the packets are significantly
larger.

Introducing IP in the radio access network implies that an IP QoS-
capable domain, e.g. a Differentiated Services domain, will have to
be introduced and managed in the radio access network. This domain
consists of edge and interior nodes, where the edge nodes are the
nodes located at the boundary of the domain.  All the nodes which are
part of this QoS-capable domain and are not edge nodes are defined as
interior nodes.

An edge node can be defined as an ingress node, or a node that
handles the traffic as it enters the QoS-capable domain.
Alternatively, an edge node might be an egress node, or a node that
handles the traffic as it leaves the QoS-capable domain.  In this
memo, an edge node (ingress or egress) is denoted as the first hop
router that the base station or BSC/RNC is connected to. The first
hop router might be a part of the base station or BSC/RNC.

Furthermore, the base station and BSC/RNC must be able to handle
algorithms used for purposes other than edge node functionality that
are many times more complex than the algorithms required for handling
the edge node functionality. Therefore, the edge node functionality
will only have minimal impact on the complexity of the base station
or BSC/RNC.


**3.2.  Motivation for this memo**


The issues described in this document concern only the cellular radio
access network (RAN).

The architecture of the RAN and the nature of the transported data
mean that the RAN has different characteristics when compared with
other IP-based networks.  However, those differences, which are the
motivation for this memo, are limited to the domain of the RAN and do
not extend into the backbone of the IP network.

In order for the transport within the RAN to function satisfactorily,
even if the transport network is IP-based, we need to ensure that
there are adequate resources in the transport network to meet the
needs of the data flows between the nodes within the RAN.

Based upon the characteristics of the RAN (described in Section 4
below), the current strategies for resource management do not meet
the requirements for an appropriate resource management strategy
within a RAN.  This document seeks to initiate a dialog on how to
correct that situation.

## 4.  Network characteristics of cellular access networks today

Cellular RANs today have a unique set of characteristics compared to
other kinds of IP networks.  These characteristics result in a set of
requirements on any resource reservation scheme that might be used in
the RAN.

### 4.1.  General aspects of the network structure

The network structure for cellular radio access networks can be
described as having the following characteristics:

 * Operator relationship

   The RAN is typically controlled by a single cellular operator
   with full control over the network.  The IP network used to
   transport radio frames might be leased from another operator
   or be built by the same cellular operator.  This network
   could be thought of as an "intranet".

 * Size of network

   RANs can be very large routed networks.  Networks including
   thousands of nodes are certainly within reason.

* Traffic volume

  The traffic from a large number of radio base stations
  needs to be supported by the same transport network. Even
  if a single radio base station generates a modest volume of
  traffic, the total number of flows for radio frame transport
  in the radio access network is very large.

* Transmission sharing

  The network between the BSC/RNC and the base stations is
  built to support transmission sharing between different
  nodes even if they are geographically distributed.
  One transmission link (possibly with redundancy) can
  support more than one base station. In other words, one
  piece of hardware can serve more than one base station and
  therefore can support more cells in one location, such as
  a three sector site. This means that the cells that are
  located at the same location will have to compete for the
  same transmission resources.

* Unicast transport of radio frames

  The transport of radio frames in the radio access network is
  point-to-point transport. Even if the soft handover splitting
  in the BSC/RNC is multicast of radio frames in some sense,
  this is handled above the IP layer by the frame transport
  protocol. For each radio channel in each base station the
  frame transport protocol needs a separate flow. Therefore,
  the frame transport protocol requires unicast transmission
  from the IP layer.

## 4.2. High cost for transmission

The transmission between base stations and the BSC/RNC is usually on
leased lines, and this part (due to the wide area coverage of the
cellular network) is usually extremely expensive when compared to the
cost of transmission in the backbone.  Due to the number of base
stations, the cost for these leased lines can be quite significant.
Even if the cost for the leased line decreases over the years, the
"last mile" to the base station is likely to be expensive due to the
location of the base station.

Cellular RANs are built over a very wide geographic area. There are,

of course, many different networks that cover a wide geographic area
(e.g., across USA, and the world), but the dispersion of nodes over
the area in the RAN case is different. Due to the fact that the base
stations are positioned based on a radio network perspective, i.e.,
radio coverage, and not based on a transmission perspective, a large
proportion of nodes are distributed throughout rural and urban areas,
not close to installed high-capacity transmission hubs.  Even worse,
the the base stations could be positioned far out in the countryside.

The peak bitrate of the multirate radio channels is selected on-
demand.  To utilize the bandwidth of the expensive transmission links
used for radio frame transport efficiently, dimensioning using over-
provisioning might therefore be prohibitively expensive, and it is
unlikely that the network will be dimensioned for peak allocation.
Dynamic allocation and optimization to reduce the cost are therefore
a fundamental requirement.  Resource reservations make it possible to
have high utilization of the network for real-time sensitive traffic
as well as avoiding congestion in the network.

## 4.3.  Transportation of radio frames

The traffic in cellular access networks is dramatically different
from the Internet in general.  The Internet primarily supports best-
effort traffic today, while the traffic on a RAN is (at least today)
largely real-time traffic. The network from the base station to the
BSC/RNC is the part of the network that has the highest volume of
real-time traffic and where delay must be minimized as much as
possible. The reasons for the this are:

 * End-to-end delay for speech traffic consists of delays in the
   mobile stations, the RAN and in the MSC (see Figure 1 in
   Section 3.1). The major portion of delay in the RAN is caused
   by the radio-related functionality (e.g., interleaving and
   coding in the base station and adaptation in the BSC/RNC).
   Therefore, the combined delay in all parts (MSC, radio
   network, and mobile stations) must be minimized as much
   as possible to give the end user proper speech quality.

 * Handover is a major issue.  For GSM, with typically multiple
   handovers per call, excessive delay in the control, e.g.
   radio network control traffic, of the radio network will
   cause a longer handover interruption period. The majority
   of handovers are also made within the radio network.

* The transport of radio frames is very delay-sensitive. In
  the direction from the BSC/RNC to the base station, a radio
  frame is a short segment of data (payload) to be coded and
  transmitted on a given radio channel by the base station at a
  given point in time. In the direction from the base station
  to the BSC/RNC, a radio frame is a short segment of data
  (payload) that was received and decoded by the base station
  at a given point in time and potentially needs to be combined
  in the BSC/RNC with radio frames received by other base
  stations at the same point in time as this particular frame.
  Note that the data segment in a radio frame may contain user
  data but also control signalling information and the same
  type of synchronized frame transport is needed for almost
  all kinds of radio channels and is generally not coupled
  to the type of service. Therefore, even if an end to end
  application is best effort, the transport of the radio frames
  originating from this application might be treated as real-time
  within the radio access network.

The real-time traffic on the RANs is today almost exclusively voice
(up to 90% with 10% signalling), but the cellular systems are
evolving to provide capabilities for other kinds of real-time traffic
(e.g., video).  Nonetheless, voice continues to be one of the most
important sources of revenue in most cellular environments today. The
transport resources are today allocated when the call is accepted,
and the radio frame transport over an IP network has to provide the
same guarantees. If real-time traffic cannot be engineered to work
correctly, the primary revenue stream will disappear.

Some of the sources, such as video-based services and gaming, will be
able to send data at a variable bitrate at higher rates. For radio,
the rates of the radio channels are selected on-demand.  In reality,
the radio network can support a wide range of partitioning of the
radio resources among the different radio channels. A rather large
portion of the transmission resources between the base station and
BSC/RNC will have to be allocated for such services. To be able to
utilize the bandwidth used for radio frame transport efficiently, the
same flexibility is required in assignment of the transport resources
as in the air-interface. Therefore, statically-assigned resources
will induce a cost which is too high for the operator.

**4.4**.  **Mobility aspect of radio frame transport**

   The mobility of the mobile stations imposes strong requirements on
   managing the transmission resources available in the RAN,
   Furthermore, this also implies that there are strong requirements on
   the RAN's internal signaling and not only on transferring of packets
   sent by the mobile station.

   Hard handover is one of the issues. For GSM, with typically multiple
   handovers per call, excessive delay in the control of the radio
   network will cause a longer handover interruption period. Typically,
   most of the handovers will be made between base stations controlled
   by the same BSC and therefore extensive delay between base station
   and BSC will degrade more than delays in the MSC and SGSN/GGSN
   network.

   Moreover, for maximal utilization of radio spectrum in WCDMA (and
   also in CDMA), fast and frequent (soft) handover operations between
   radio channels and radio base stations are required. The frequency of
   handover events is therefore typically higher in WCDMA radio access
   networks than in GSM and means even higher performance requirements
   on the transport solution. If the soft handover cannot  be performed
   fast enough, spectrum cannot be utilized efficiently, which will
   cause degradation of the radio network capacity. At each handover
   event, resource reservation is needed, and therefore resource
   reservation needs to be fast and will be used very frequently.

   The impact of mobility in the radio access network has therefore two
   major differences compared to the fixed network:
    (1) High volume of resource reservation events
    (2) Requirement on short response time for reservations


**5**.  **Requirements on a Resource Reservation Scheme**

   This section outlines what we believe are the fundamental
   requirements placed on any resource reservation scheme in a cellular
   radio access network. Later sections will outline how current schemes
   match these requirements.


**5.1**.  **Main requirement on resource reservation scheme**


   One of the primary requirements that real-time applications impose on

any resource reservation scheme is the provisioning of good QoS
(delay and packet loss) guarantees. This can only be achieved if the
network can be utilized while avoiding congestion and without having
too high packet losses. The level of utilization depends on network
topology, traffic mix, scheduling discipline, delay and packet loss
requirements. The utilization is given by network dimensioning but
should be as high as possible.

The resource reservation scheme must be able to keep the real-time
traffic under a certain pre-defined network utilization in order to
bound congestion. Otherwise, it is impossible to guarantee percentile
bounds on the QOS requirements for real-time traffic.

## 5.2.  IP must provide same service behavior as the transport networks
used today

Today's commercial IP networks are mainly optimized to carry best-
effort traffic. As explained above and also discussed in [WeLi99],
the transport of radio frames in the radio access network puts real-
time requirements on the underlying transport network.  All of these
characteristics are fulfilled by the connection-oriented transport
networks (STM and ATM) used by cellular networks today.  By, at a
minimum, meeting these same requirements, the IP networks will be
capable of providing the same behavior as the transport networks that
are currently used by cellular systems while gaining the advantages
of IP networking. It should be noted that IP networks will be able to
meet these requirements only if the following two constraints are
met:

(1) that service guarantees are percentiles, see Section 5.1
(2) strictly limited to a given operator's IP network, see
    Section 5.9.

## 5.3.  Efficient network utilization

Due to the high cost of the leased transmission, we must utilize the
network to the highest degree possible, and this must be facilitated
by the resource reservation scheme.

However, in considering a resource reservation scheme, its impact
upon the performance and scalability of the network as a whole must
also be taken into account.  For example, the performance and
scalability impact on the edge and internal routers is a very
important consideration.

5.4.  **Handover performance requirements on resource reservation scheme**

Whatever reservation scheme is used must be highly performant for at
least the following reasons:

  * Handover rates

    In the GSM case, mobility usually generates an average
    of one to two handovers per call. For third generation
    networks (such as WCDMA and cdma2000), where it is
    necessary to keep radio links to several cells
    simultaneously (macro-diversity), the handover rate is
    significantly higher (see for example [KeMc00]).
    Therefore, the admission control process has to cope
    with far more admission requests than call setups alone
    would generate.

  * Fast reservations

    Handover can also cause packet losses. If the processing
    of an admission request causes a delayed handover to the
    new base station, some packets might be discarded, and
    the overall speech quality might be degraded
    significantly.

    Furthermore, a delay in handover may cause degradation
    for other users. This is especially true for radio access
    technologies using macro-diversity, such as WCDMA and CDMA,
    where a handover delay will cause interference for other
    users in the same cell. Furthermore, in the worst case
    scenario, a delay in handover may cause the connection
    to be dropped if the handover occurred due to bad radio
    link quality.

    Therefore, it is critical that an admission control

request for handover be carried out very quickly. Since
the processing of an admission control request is only
one of many tasks performed during handover, the time to
perform admission control should be a fraction of the
time available for handover.

Furthermore, in the situation that the transport network
in the RAN is over-utilised it is preferable to keep
the reservation on already established flows while new
requests might be blocked. Therefore, the handover
requests for resource reservation should be treated
with a higher priority than the new requests for
resource reservation.

## 5.5.  Edge-to-edge reservations, not end-to-end

Real-time applications require a high level of quality of service
(QoS) from the underlying transmission network. This can only be
achieved by accomplishing the QoS management on an end-to-end basis
(i.e., end user to end user), from application to application,
potentially across many domains.

However, this does not mean that the resource reservation protocol
must be applied end-to-end. The end-to-end QoS management
architecture may consist of many interoperable edge-to-edge QoS
management architectures where each of them might use a different
edge-to-edge resource reservation protocol.  In fact, this is far
more likely to be the case than that a global signaling structure
will be available across all different domains in an end-to-end
perspective.  This will increase the flexibility and the openness of
the transmission network since various access networks that are using
the same transmission network and different edge-to-edge QoS
management architectures will be able to interoperate.

It is critical that the appropriate mechanisms for providing the
service guarantees needed in the radio access network be put in place
independently of solving the more difficult problem of end-to-end
QoS.

In our case, the access network is simply an intranet in which we

need to solve a local QoS problem.  This implies that a general
solution which handles the end-to-end QoS problem is unnecessarily
complicated for solving the intranet problem in the cellular access
network.


**5.6**.  **Reservation functionality in edge nodes versus interior routers**

   In our network, it is important that the reservation mechanism be as
   simple as possible to implement in the interior nodes since in most
   cases there might be more interior routers (<= 10 depending on
   network structure) in the path than there are edge nodes.

   Typically, in a RAN there are two edge nodes located in a
   communication path. Moreover, the average number of interior nodes in
   a communication path within a RAN depends on the chosen network
   topology by the RAN operator.  As such, the scheme must be optimized
   for the interior nodes and not for the edge nodes, thus reducing the
   requirements placed on the functionality of the interior routers.
   This means that we can have complicated mapping of traffic parameters
   at the edges and a simplified model in the interior nodes, and that
   the necessary set of parameters required for setting up reservations
   shall be based upon their effect on interior nodes and not on edge
   nodes.

   The edge routers typically have to perform per-session
   management/control, and hence complex per flow handling is not a
   significant burden.


   However, interior nodes do not need to have per flow
   responsibilities. We must therefore optimize for simple QoS
   mechanisms on these interior devices, and use more complex mechanisms
   in the edge devices.


   In our case, edge device functionality is implemented in the first
   hop router that the base station or BSC/RNC is connected to (see
   Section 3.1). In this way we optimize for simple QoS mechanisms on
   the interior devices, while the more complex mechanisms are applied
   on the edge devices, e.g., base station, BSC/RNC.

   This emphasis on simplicity is due to performance requirements listed
   above.  We need to make sure that we understand the minimal level of
   functionality required in the reservation scheme in order to

guarantee the performance of real-time traffic.


**5.7**.  **On-demand and dynamic allocation of resources**

   Real-time services require that a portion of network resources is
   available to them. These resources can be reserved on a static or
   dynamic basis, or potentially based on some kind of measurement of
   network load.

   In the first situation, this may result in an poorly utilized
   network. This is mainly due to the fact that the network resources
   are typically reserved for peak real-time traffic values.  Mobility
   in the network makes static configuration even less desirable as the
   resources will be used even less effectively.

   If using dynamic allocation, this problem will be avoided since the
   resources are reserved on demand.  However, the load from resource
   reservation will be much higher than if static allocation of
   resources is used.  If the dynamic allocation of the resources is
   done on a micro-flow basis, the resulting network load from resource
   reservation might be quite high.

   We might use other methods, such as measurement-based admission
   control, to simplify the reservation protocol, as long as these
   methods can fulfill the requirements (now or in the future).

   What is important is that all of these mechanisms can be used for
   solving part of the network utilization problem, and, as such, any
   reservation scheme must have the flexibility to provide both on-
   demand reservations as well as measurement-based admission control.

   As high bitrate and variable bitrate applications enter the cellular
   space, the need for on-demand reservations of resources will become
   even more acute.



**5.8**.  **Unicast and not multicast**


   The majority of the traffic in the RAN is point-to-point unicast
   transport of radio frames between the base station and the BSC/RNC.
   As such, the resource reservation scheme need not to be optimized for

multicast.


## 5.9.  A single operator in the RAN


It is realistic to assume that end-to-end communication in IP
networks as well as the end-to-end QoS management architectures will
be managed by more than one operator.

Furthermore, it also realistic to assume that an edge-to-edge
resource reservation protocol can be managed by one single operator.
As such, it is reasonable to limit reservation scheme to a single
operator domain. This will ensure that each operator can optimize the
edge-to-edge QoS management architecture for their needs. Moreover,
this limitation (a single operator domain) means that the reservation
scheme does not need to handle the issues inherent in a multi-
operator domain, thus simplifying the scheme.


## 5.10.  Minimal impact on router performance


The performance of each network node that is used in an end-to-end
communication path has a significant impact on the end-to-end
performance of this communication path. Therefore, the end-to-end
performance of the communication path can be optimized by optimizing
the router performance.  It is absolutely critical that the
introduction of QoS mechanisms and signaling does not overly impact
the performance of the infrastructure.  Obviously, you cannot
introduce new things that need to be done by networking
infrastructure without impacting its performance, but that impact
must be minimized to the greatest extent possible.

One of the factors that can contribute to this optimization is the
minimization of the resource reservation signaling protocol load on
each router. When the dynamic allocation of the resources is on a per
micro-flow basis, the resource reservation signaling protocol could
easily overload a router located in a core network, causing severe
router performance degradation.  Furthermore, any mechanisms defined
must be such that it is reasonable to implement them in hardware
which will increase the scalability of the solution.

**5.11**.  **Scalably Manageable**


   Any strategy for resource management in a RAN must be done in such a
   way that it is easily manageable in a very large network.  This
   implies as little "laying on of hands" as possible and as much
   automation as possible.  In networks made up of many thousands of
   routers, changing of even a single parameter in all routers may be
   prohibitively difficult.  Minimizing the involvement of the operator
   (or the operator's management tools) is therefore an important
   requirement.


**5.12**.  **Bi-directional reservations**


   In current RANs, the BSC/RNC is responsible for initiation of
   reservation of resources in the transport plane.  Therefore, via the
   resource reservation signaling protocol, the BSC/RNC has to support
   the initiation and management of the resource reservations for both
   directions, both to and from the base station, simultaneously. In
   this way a simpler edge-proxy resource reservation functionality will
   be implemented in the base station, decreasing its complexity.


**5.13**.  **Support for non-RAN specific traffic**

   Any strategy for resource management used in a cellular RAN must be
   able to support any type of traffic (RAN-specific or non-RAN
   specific) in the same way, as long as the traffic belongs to the same
   traffic class.

   The RAN-specific traffic is the traffic that is transported through
   the RAN and is generated or used by specific entities belonging to
   the same cellular technology as the one used in the RAN.

   The non-RAN specific traffic is the traffic that is transported
   through the RAN but is neither generated by nor used by any specific
   entity belonging to the same cellular technology as the one used in
   the RAN.

**[6](#)**.  **Evaluation of existing strategies**


   In order to understand whether technology exists today which will
   allow us to manage the resources in cellular networks, we briefly
   look at the protocols that currently exist which address parts or all
   of these requirements.


**[6.1](#)**.  **End-to-end per-flow resource reservation protocol**


   An end-to-end per-flow resource reservation signaling protocol is
   applied in an end-to-end IP communication path, and it can be used by
   an application to make known and reserve its QoS requirements to all
   the network nodes included in this IP communication path.  This type
   of protocol is typically initiated by an application at the beginning
   of a communication session. A communication session is typically
   identified by the combination of the IP destination address,
   transport layer protocol type and the destination port number.  The
   resources reserved by such a protocol for a certain communication
   session will be used for all packets belonging to that particular
   session.  Therefore, all resource reservation signaling packets will
   include details of the session to which they belong.

   The end-to-end per-flow resource reservation signaling protocol most
   widely used today is the Resource Reservation Protocol (RSVP) (see
   [[RFC2210](#)], [[RFC2205](#)]). The main RSVP messages are the PATH and RESV
   messages.  The PATH message is sent by a source that initiates the
   communication session. It installs states on the nodes along a data
   path.  Furthermore, it describes the capabilities of the source. The
   RESV message is issued by the receiver of the communication session,
   and it follows exactly the path that the RSVP PATH message traveled
   back to the communication session source. On its way back to the
   source, the RESV message may install QoS states at each hop. These
   states are associated with the specific QoS resource requirements of
   the destination. The RSVP reservation states are temporary states
   (soft states) that have to be updated regularly. This means that PATH
   and RESV messages will have to be retransmitted periodically. If
   these states are not refreshed then they will be removed.  The RSVP
   protocol uses additional messages either to provide information about
   the QoS state or explicitly to delete the QoS states along the
   communication session path. RSVP uses in total seven types of
   messages:

  * PATH and RESV messages

  * RESV Confirm message

  * PATH Error and RESV Error messages

  * PATH Tear and RESV Tear messages


 An overview of the functionality of the RSVP functionality includes:

  * End-to-end reservation with aggregation of path
    characteristics such as fixed delay.

  * The same type of reservation functionality in all
    routers. Only policy handling separates the edge of the
    domain from other routers.

  * Multicast and unicast reservations with receiver initiated
    reservations. RSVP makes reservations for both unicast and
    many-to-many multicast applications, adapting dynamically
    to changing routes as well as to group membership.

  * Shared reservations for multiple flows.

  * Support for policy handling to handle multi-operator
    situations since more than one operator will be
    responsible for RSVP's operation.

  * Flexible object definitions. RSVP can transport and maintain
    traffic and policy control parameters that are opaque to
    RSVP. Each RSVP message may contain up to fourteen classes of
    attribute objects. Furthermore, each class of RSVP objects
    may contain multiple types to specify further the format
    of the encapsulated data. Moreover, the signaling load
    generated by RSVP on the routers is directly proportional
    to the flows processed simultaneously by these routers.
    Furthermore, processing of the individual flows in the
    networking infrastructure may impose a significant processing
    burden on the machines, thus hurting throughput. These
    issues make it reasonable to question the scalability and
    performance in a large cellular radio access network.

  * support for uni-directional reservations, not bi-directional.

In the situation where a mobile moves or the connection moves from
one base station to another, it could force the communication path to
change its (source/dest) IP address.  The change of IP address will
require that RSVP establish a new RSVP session through the new path
that interconnects the two end points involved in the RSVP session
and release the RSVP session on the old path.  During this time, the
end-to-end data path connection is incomplete (i.e., QoS disruption)
and it will negatively affect the user performance.

This approach includes much more functionality and complexity than is
required in the cellular RAN. Our problem is significantly simpler to
solve.  The trade-off between performance and functionality is one of
the key issues in the RAN. In our case, the majority of the
functionality in RSVP is not required.  This is true for four
reasons:

 * Unicast reservations are much less complex than multicast.

 * Edge-to-edge with one operator does not require policy
   handling in the interior routers.

 * Path characteristics and flexible traffic parameters and
   QoS definitions could be solved by network dimensioning
   and edge functionality.

 * Per microflow states in intermediate routers cause severe
   scalability problems. Furthermore, receiver-initiated
   reservations impose high complexity in the states due to
   reverse-direction routing of the RESV messages.  A scheme
   based on sender oriented reservation (see e.g., [AhBe99])
   decreases the complexity of the per microflow states due
   to the fact that no reverse-direction routing is
   required.

## 6.2.  Integrated Services over Differentiated Services

The IntServ over DiffServ framework addresses the problem of
providing end-to-end QoS using the IntServ model over heterogeneous
networks. In this scenario, DiffServ is one of these networks
providing edge-to-edge QoS. This is similar to the underlying
architecture for this draft, where the specific network is the
cellular RAN, and where the end-to-end model is unspecified. As such,

the problem addressed by IntServ over DiffServ is similar in nature
to the problem described here, although the specific requirements
(such as network utilization and performance) are different.

The IntServ over DiffServ framework discusses two different possible
deployment strategies. The first is based on statically allocated
resources in the DiffServ domain.  In this strategy the Diffserv
domain is statically provisioned (see Section 6.3). Furthermore, in
this strategy the devices in the Diffserv network region are not RSVP
aware. However, it is considered that each edge node in the customer
network is consisting of two parts. One part of a node is a standard
Intserv that interfaces to the customer's network region and the
other part of the same node interfaces to the Diffserv network
region. Any edge node in the customer network maintains a table that
indicates the capacity provisioned per SLS (Service Level
Specification) at each Diffserv service level. This table is used to
perform admission control decisions on Intserv flows that cross the
Diffserv region.  A disadvantage of this approach is that the edge
nodes in the customer network will not be aware of the traffic load
in the nodes located within the Diffserv domain.  Therefore, a
congestion situation on a communication path within the Diffserv
domain cannot be predicted by any of these edge nodes. Due to the
"Efficient network utilization" requirement explained in Section 5.3,
the RAN is dimensioned such that it may have performance bottlenecks
which are not visible to the edges. More advantages and disadvantages
of this approach are discussed in Section 6.3.


The second possible strategy is based on dynamically allocated
resources in the DiffServ domain. According to [RFC2998], this can be
done using RSVP-aware DiffServ routers.  However, this approach has
most of the drawbacks described in Section 6.1, and per-microflow
state information is kept in the intermediate routers.

Alternatively, resources in the DiffServ domain can be dynamically
allocated using Aggregated RSVP. This will be discussed in Section
6.4.

Other approaches related to the bandwidth broker concept are still
very immature and will not be discussed here.

### [6.3](). **Statically-assigned trunk reservations based on Differentiated** Services

A significant problem in deploying an end-to-end per-flow resource reservation signaling scheme is its scalability. This can be solved by aggregating (trunking) several individual reservations into a common reservation trunk.  The reservation trunks can be either statically or dynamically configured.  When the reservation trunks are statically configured, no signaling protocol is required for performing the reservation of network resources but is likely to be a difficult management problem.  However, due to the different mobility requirements (such as handover) and QoS requirements (such as bandwidth) that the multi-bitrate applications impose on the RAN, it will be difficult to configure the trunked reservations statically and utilize the RAN efficiently.

### [6.4](). **Dynamic trunk reservations with Aggregated RSVP**

The reservation trunks can be dynamically configured by using a signaling protocol that manages various mechanisms for dynamic creation of an aggregate reservation, classification of the traffic to which the aggregate reservation applies, determination of the bandwidth needed to achieve the requirement and recovery of the bandwidth when the sub-reservations are no longer required.

The first router that handles the aggregated reservations could be called an Aggregator, while the last router in the transit domain that handles the reservations could be called a Deaggregator.

The Aggregator and Deaggregator functionality is located in the edge nodes. In particular, an Aggregator is located in an ingress edge node, while a Deaggregator is located in an egress edge node, relative to the traffic source.

The aggregation region consists of a set of aggregation capable network nodes.  The Aggregator can use a policy that can be based on local configuration and local QoS management architectures to identify and mark the packets passing into the aggregated region. For example, the Aggregator may be the base station that aggregates a set of incoming calls and creates an aggregate reservation across the edge-to-edge domain up to the Deaggregator.  In this situation the call signaling is used to establish the end-to-end resource

reservations. Based on policy, the Aggregator and Deaggregator will decide when the Aggregated states will be refreshed or updated.

One example of a protocol that can be used to accomplish QoS dynamic provisioning via trunk reservations is the RSVP Aggregation signaling protocol specified in [BaIt00].

With regards to aggregated RSVP, even if the reservation is based on aggregated traffic, the number of re-negotiations of the allocated resources due to mobility (handover) does not decrease and each re-negotiation of resources has the same performance requirements as the per-flow reservation procedure.

Note that the aggregated RSVP solution may use a policy to maintain the amount of bandwidth required on a given aggregate reservation by taking account of the sum of the underlying end to end reservations, while endeavoring to change it infrequently. However, such solutions (policies) are very useful assuming that the cost of the overprovisioned bandwidth is not significant, since this implies the need for a certain "slop factor" in bandwidth needs. In a RAN, where overprovisioning is not preferable due to high costs of transmission links, a more dynamic QoS provisioning solution is needed.

Furthermore, the aggregated RSVP scheme is receiver initiated and cannot support bi-directional reservations.

In the aggregated RSVP scheme the resource reservation states stored in all the RSVP aware Edge and Interior nodes represent aggregated RSVP sessions (i.e., trunks of RSVP sessions). Therefore, the number of the resource reservation states in the aggregated RSVP scheme compared to the (per-flow) RSVP scheme, is decreased.  However, in a Diffserv based RAN the number of the aggregated RSVP sessions depends on:

  *  the number of Aggregators/Deaggregators; Considering
     that each base station and each BSC/RNC is used as
     Aggregator/Deaggregator, the total number of
     Aggregators/Deaggregators within the RAN is
     significantly high. This is due to the fact
     that the number of BSCs/RNCs is significantly
     high and the number of base stations in a RAN is
     in the range of thousands, see Section 3.1.
  *  the network topology used; The communication between
     RNCs is performed in a meshed way, i.e., all to all
     communication. This will imply that many communication

       paths will have to be maintained by the RAN
       simultaneously.
   *   the number of Diffserv Code Points (DSCPs) used; More
       than one traffic classes will be supported
       within the RAN. Therefore, the number of the Diffserv
       Code Points (DSCPs) used within the RAN will probably
       be higher than one.

   Therefore, the number of the aggregated RSVP reservation states
   within a Diffserv based RAN will be significantly large.


7.  **Conclusion**


   Cellular radio access networks and coming wireless applications
   impose different requirements on reservation strategies than typical
   Internet conditions.

   Firstly, the reservation solution does not need to have the same
   level of complexity:

     * Edge-to-edge not end-to-end: The IP traffic is generated
       in the network and is only transported as far as the
       cellular-specific nodes (such as the base station and
       BSC/RNC).

     * Single operator domain and no inter-domain transport: The
       transport is owned and managed by a single operator.

     * Only unicast not multicast: The end-to-end payload is
       transported between nodes. This transport only requires
       a unicast reservation.

   Furthermore, a cellular radio access network has much higher
   performance requirements on the reservation strategy:

     * Efficient usage of the transmission network: The transport
       between the BSC/RNC and the base station represents a significant
       cost for the cellular operator, and efficient usage of
       the transmission network is therefore critical from a cost
       point of view. The network should allow dynamic allocation
       of resources to allow efficient statistical aggregation

of traffic without causing congestion.

* A wide-area network with significant volume of real-time
  traffic: Real-time traffic levels up 100% must be
  supported.

* The resource reservation process has to handle a
  significantly higher volume of requests, and the process
  has to be fast enough to avoid packet losses in the
  air-interface during handover.

* The scheme must be optimal for interior nodes and not
  for the edge nodes. In this way the necessary set of
  parameters required for setting up reservations should be
  based upon their effect on interior nodes and not on edge
  nodes. This reduces the complexity on the interior routers.

Given these requirements, we believe that appropriate standardization
should take place to create the necessary protocols for edge-to-edge
resource management for a single operator domain.

## 8. References

[AhBe99]  Ahlard, D., Bergkvist, J., Cselenyi, I., "Boomerang
          Protocol Specification", Internet draft, Work in progress.

[BaIt00]  Baker, F., Iturralde, C. Le Faucher, F., Davie, B.,
          "Aggregation of RSVP for Ipv4 and Ipv6 Reservations",
          IETF RFC 3175, 2001.

[KeMc00]  Kempf, J., McCann, P., Roberts, P., " IP Mobility
          and the CDMA Radio Access Network: Applicability
          Statement for Soft Handoff", Internet draft, Work
          in progress.

[RFC2205] Braden, R., Zhang, L., Berson, S., Herzog, A., Jamin, S.,
          "Resource ReSerVation Protocol (RSVP) -- Version 1
          Functional Specification", IETF RFC 2205, 1997.

[RFC2210] Wroclawski, J., "The use of RSVP with IETF Integrated
          Services", IETF RFC 2210, 1997.

[WeLi99]  Westberg, L., Lindqvist, M., "Realtime Traffic over

Cellular Access Networks", Internet draft, Work in
progress (expired).

[RFC2998] Bernet, Y., Yavatkar, R., Ford, P., baker, F., Zhang, L.,
         Speer, M., Braden, R., Davie, B., "Felstaine, E.,
         "Framework for Integrated Services operation over
         Diffserv Networks", IETF RFC 2998, 2000.

## 9.  Acknowledgements

## 10.  Authors' Addresses

David Partain
Ericsson Radio Systems AB
P.O. Box 1248
SE-581 12  Linkoping
Sweden
EMail: David.Partain@ericsson.com

Georgios Karagiannis
Ericsson EuroLab Netherlands B.V.
Institutenweg 25
P.O.Box 645
7500 AP Enschede
The Netherlands
EMail: Georgios.Karagiannis@eln.ericsson.se

Pontus Wallentin
Ericsson Radio Systems AB
P.O. Box 1248
SE-581 12  Linkoping
Sweden
EMail: Pontus.Wallentin@era.ericsson.se

Lars Westberg

Ericsson Research
Torshamnsgatan 23
SE-164 80 Stockholm
Sweden
EMail: Lars.Westberg@era-t.ericsson.se