# Metaswitch Networks

# White Paper

# A Fresh Look at
# Session Border Control

Martin Taylor, CTO

Metaswitch Networks

## Executive Summary

Over the next decade or so, service providers will be investing more than five billion dollars in purchasing Session Border Controllers (SBCs) to support their build-out of VoIP service access and interconnect – hugely more than has been invested in this technology to date.

The great majority of SBCs that have been deployed in the network to date are compact integrated devices that handle both signaling and media in a monolithic "appliance" form factor. These appliances are dimensioned in terms of numbers of concurrent sessions, where a session is modeled as a voice call of typical average duration.

Deploying SBCs optimized for basic voice telephony may be appropriate for the current mix of services seen in typical consumer or business VoIP environments, but the service mix is expected to change very substantially as service providers introduce new SIP-based capabilities into their networks such as presence, instant messaging and call forking. These services have a dramatic impact on signaling load, driving up SIP message rates by as much as an order of magnitude but having much lower impact on media plane loading.

The current generation of integrated appliance-based SBCs is ill-suited to support this changing service mix, for two reasons. First, these products embody signaling and media capacity in a fixed relationship that is unable to adapt to the changing service mix. And secondly, they offer rather modest overall capacity in the signaling plane, requiring them to be deployed in very large numbers to support the signaling loads that arise with new SIP-based services.

A new solution for session border control is required that enables signaling and media capacity to be scaled independently of one another, and that delivers far greater scalability in the signaling plane. Such next-generation solutions for session border control are finally starting to become available. Service providers would do well to re-evaluate their SBC deployment plans so as to avoid over-investing in architectures that are ill-matched to their medium and longer term SIP services strategy.

## Table of contents

**Meta**switch
Networks

## 1. Introduction

Voice over IP technology has been part of the telecommunications mainstream for more than a decade, but in reality the transition of the public network from circuit-switched to packet-switched technology has barely begun. According to ITU-T worldwide statistics, at the end of 2010 there were about 6.5 billion fixed and mobile phone lines based on circuit-switched technology. By contrast, Infonetics Research estimates that there were about 270 million fixed and mobile VoIP lines in use at this time – just over 4% of the worldwide total. Statistics for inter-carrier traffic exchange using VoIP are hard to come by, but anecdotal evidence suggests that the vast majority of inter-carrier traffic is still handled using traditional TDM connections.

Session Border Controllers (SBCs) are an essential element in the transition of both access and interconnect facilities from circuit-switched to packet-switched technology. They secure the service provider's network against a range of potential threats, support the transit of VoIP traffic across customer firewall and Network Address and Port Translation (NAPT) devices, and promote interoperability with VoIP endpoints. They also represent a very substantial proportion of the cost of deploying VoIP access and interconnect.

The session border control function straddles the signaling and media planes in VoIP networks. In the signaling plane, SBCs perform complex filtering and transformation operations on SIP signaling messages, while in the media plane they perform bandwidth policing, address translation and statistics collection on RTP streams. Most SBCs being deployed today are monolithic devices or "appliances" that perform both the signaling and media functions in one integrated unit.

SBCs are already well established in most service providers' networks at both access and interconnect edges, but total deployed SBC capacity will have to grow massively in the future to address the complete transition of public fixed and mobile voice networks to VoIP. Furthermore, the SBC function will need to evolve very substantially to adapt to the changing SIP services landscape, where SIP-based messaging, presence and call forking will have dramatic impacts on overall SIP signaling load.

This paper explores the impact of SIP growth and changing service mix on the session border control function, and argues that the current generation of integrated appliance-based Session Border Controller products will soon be seen as hopelessly inadequate for the task in hand. The next generation of SBC solutions will be based on a massively scalable distributed architecture, in which signaling and media capacity can be scaled independently – and where appropriate, deployed in the cloud. This evolution of session border control really has to happen before we can realistically contemplate a large-scale transition of public fixed and mobile voice networks to VoIP-based access and interconnect.

## 2. The Evolving SIP Services Landscape

We start by exploring how SIP-based services are expected to evolve beyond basic voice calling, and the impact that this will have in terms of the signaling and media loads experienced by the session border control function in the network.

### 2.1. Current VoIP Services

Today's VoIP networks support three main kinds of services: Consumer VoIP, Hosted PBX and SIP Trunking. These services are invariably delivered over a fixed broadband infrastructure. Wireless operators have yet to deliver VoIP services over mobile broadband in any significant volume – despite the fact that many of them make extensive use of VoIP in their core transport networks.

The characteristics of the signaling and media loads associated with current VoIP services are quite well understood. The vast majority of calls are simple voice calls, with no video content, and they last on average 2-3 minutes. In the signaling plane, each call leg requires between 7 and 16 SIP messages to be exchanged between the endpoint and the network to set up the call and tear it down when completed. The media path may exist between two VoIP terminals on the same network, or between a VoIP terminal and a media gateway. In both cases, the two endpoints almost invariably negotiate a codec that they both support, so that no transcoding is required.

Integrated appliance-based SBCs are engineered to handle this particular mix of signaling and media load. An integrated SBC that is rated to handle 10,000 concurrent calls will typically have just the right amount of processing power to support the minimum 7 messages per call leg, based on an average call hold time of about 2 minutes – about 600 SIP messages per second – and just the right amount of media packet forwarding capacity to support 10,000 typical media streams – for example using G.711 with 20 ms packetization interval.

It's worth noting that there is considerable variation between different network environments in the number of SIP messages required to set up and tear down a call leg. For example, if the network is set up to require SIP authentication on all client requests, early media (for cut-through of ringback tone) and reliable provisional responses (PRACK) then call setup and tear down will require between 9 and 16 messages, depending on the direction of the call attempt and the call release respectively. On a trusted peering point where authentication is not required, but early media and PRACK are used, then between 9 and 12 messages need to be exchanged per call leg. Furthermore, some SIP server systems maintain a "heartbeat" to each SIP endpoint during active calls, to ensure that billing is stopped if an endpoint is unexpectedly disconnected. This heartbeat may be based on an OPTIONS request or a SIP Re-INVITE, requiring 2 or 3 messages per

www.metaswitch.com

heartbeat respectively. If repeated at, say, 30 second intervals, such a heartbeat could require 15 or so additional SIP messages per 3-minute call.

The maximum call processing rate claimed for an SBC product may well be based on an absolute minimum number of SIP messages per call leg, and needs to be discounted – perhaps by more than 50% – to reflect the actual call flows observed in the network.

## 2.2. Emerging Voice and Multimedia Services

While conventional voice calling will continue as the staple offering of SIP-based voice and multimedia networks, a variety of new service capabilities are already starting to become important in these networks – and these new services behave very differently from basic voice calls in the relative amount of load they present to the signaling and media planes respectively.

In the following, we describe the impact of new services on the signaling and media planes. A summary of these service impacts is presented in Table 1 at the end of this section.

### 2.2.1. Messaging

The first standards defining the use of SIP to support messaging services were published as long ago as 2002, and since then we have seen a series of RFCs published that define the usage of SIP for page-mode and session-oriented instant messaging and group text chat.

Despite all the standardization activity, there has been very little real-world deployment of SIP-based messaging to date. Mobile messaging continues to rely on legacy SMS and MMS technology, while in the world of fixed networks the instant messaging scene is dominated by a mix of proprietary protocols and the standard eXtensible Messaging and Presence Protocol (XMPP).

All this is about to change. The 3GPP specifications for Voice over LTE define the use of SIP-based messaging to provide service equivalence and interworking with SMS and MMS – known as "SMS-over-IP". Meanwhile the GSMA has defined standards for the Rich Communications Suite (RCS), which enables mobile network operators to deliver SIP-based session-oriented mobile instant messaging on both 3G and 4G networks. The RCS standards also address fixed network access to support seamless interop of instant messaging and media sharing between and among smartphone, tablet and desktop devices.

As wireless networks evolve to an all-IP architecture with LTE, so SMS messaging will evolve to SMS-over-IP, which makes use of the SIP MESSSAGE method to encapsulate text message content. Each SMS that is sent over IP requires the exchange of two SIP messages – a MESSAGE request followed by a 200 OK response. There is no associated traffic in the media plane.

Many users of mobile phones send many tens or even hundreds of text messages for every voice call that they make. A user that makes 5 voice calls in a day is responsible for generating about 70 SIP messages, while a user who sends 100 text messages in a day is responsible for generating 200 SIP messages, each of which will traverse at least one SBC in the network.

Session-mode messaging services will soon start to be deployed as a complement to SMS on mobile networks, offering a richer user experience including presence sharing and in-session sharing of media such as photos and video clips. The SIP call flows for session-mode messaging more closely resemble those for voice and video calls, in that they start with an INVITE and then make use of the media plane to exchange message content between the users (using the Message Session Relay Protocol, MSRP). The volume of media exchanged during a typical messaging session will vary enormously, from a few tens of bytes of text up to perhaps many megabytes of video files.

As more and more smartphones become capable of session-mode messaging, some of the back-and-forth message exchanges that users would once have conducted over SMS will take place in the context of a messaging session. This will have the beneficial effect of reducing the SIP signaling load imposed by SMS-over-IP, but substantial penetration of session-mode messaging into VoLTE will take time, so it will still be important to plan for the high signaling loads imposed by SMS-over-IP.

SIP-based messaging services will not be confined to mobile networks. Session-mode messaging services using SIP are also starting to be deployed in fixed networks as service providers offering Hosted PBX services expand these offerings to support IM and Presence as key elements of Hosted Unified Communications.

### 2.2.2. Presence

Users of Instant Messaging services typically expect to be able to see the presence of their contacts. Legacy SMS services don't support presence, but both wireless and fixed network operators will want to enhance their next-generation messaging solutions with presence capabilities in the future. Standardized services such as Rich Communication Suite (RCS) and Converged IP Messaging (CPM) already incorporate presence as a key element in their value propositions.

Presence services can impose very substantial additional loading on the SIP signaling network. Experience in the field of social networking suggests that some users may wish to broadcast frequent changes in their presence status to many tens of others. The SIP signaling load associated with such presence updates could easily dwarf that needed to support voice and video calls.

www.metaswitch.com

Some mobile operators have become so concerned about the signaling load that could result from the deployment of the presence service in RCS that they have banded together to define a subset of RCS known as RCS-e, which includes the messaging and media-sharing aspects of RCS but leaves presence out. The cost of deploying sufficient SBC capacity to handle this signaling load may well have figured in their thinking. Sadly, RCS loses a good deal of its appeal when the presence service is taken out. A session border control solution that is far more scalable and cost effective in the signaling plane could potentially transform the economics of presence services, enabling network operators to compete far more effectively in this space with the established social network providers.

Also, note that while RCS-e does not make use of presence, it does make extensive use of the SIP OPTIONS method for capabilities exchange between RCS-e endpoints. The RCS-e specifications call for the sending of a SIP OPTIONS message (and receiving a response to it) prior to each communications attempt, e.g. voice call or SMS send. This clearly adds very significantly to total SIP signaling load on the network.

### 2.2.3. Forking

SIP-based voice and video services offer the useful characteristic that two or more devices can register with the same address, enabling users to make calls from or receive calls on multiple different devices with the same phone number. Such devices may include smartphones, PC-based softphones, tablets, analog telephone adapters or dedicated SIP phone devices.

When a call attempt is addressed to a user who has multiple SIP devices registered, the network will send the SIP INVITE request to each device – a process known as "forking". The first device that answers will return a 200 OK response, and the network will then send a CANCEL request to all of the other devices. A media session is only established to the device that answers. Forking therefore has the effect of substantially increasing the signaling load for a voice or video call without impacting the load on the media plane.

Note that forking applies not just to voice and video calls, but also to attempts to establish instant messaging sessions.

### 2.2.4. Video Calling

SIP-based video calling is still in its infancy, but has already become an established element of many Hosted PBX service offerings. With the rapid penetration of smartphones equipped with front-facing cameras and inexpensive mobile broadband services, and with the growing proportion of PCs and laptops equipped with built-in cameras, demand for video calling has started to ramp up rapidly.

A SIP-based video call presents the same signaling load as a voice call – requiring between 7 and 16 SIP messages per call leg – but presents vastly greater load to the media plane. While a voice call may require only 16 kbps of data bandwidth, a video call typically requires at least 300 kbps, and preferably 500 kbps or more to deliver reasonable quality.

Growth in video calling will clearly drive demand for additional SBC media handling capacity. However, despite the far greater demands on the media plane associated with a video call, the total forecast demand for video calling represents a tiny amount of media load compared with voice. For example, in October 2010 In-Stat forecast that mobile video calling would drive 9 petabytes of data demand in North America in 2015 [ref http://www.instat.com/newmk.asp?ID=2898]. If we take FCC statistics for mobile voice published in 2010, which indicated a total mobile subscriber base in North America of 270M using an average of about 700 voice minutes per month, and we assume a voice data bandwidth of 16 kbps, then the total amount of data movement associated with voice would be over 540 petabytes per year.

Video calling presents new kinds of media processing requirements to support interoperation between devices with different video capabilities. Video media processing may include transcoding, trans-rating and aspect ratio adjustment. Such processing is most cost-effectively handled by specialized hardware, and simple economics suggests that these resources – which may only be required on a minority of video calls – should be deployed centrally in the network rather than being associated with local SBCs.

### 2.2.5. Wideband Voice Calling

Traditional circuit-switched voice calling uses the G.711 codec which is limited to a 3 kHz audio bandwidth. VoIP networks are not limited in this way, and support the transport of voice and audio streams using any arbitrary codec. Many VoIP endpoints take advantage of this capability to offer superior quality, for example by supporting the G.722 codec which supports an audio bandwidth of 6 kHz. This provides a substantial improvement in the perceived quality of voice communications – and is often marketed as "HD Voice". Skype offers wideband audio with its proprietary SILK codec, and credits at least some of its success in the market to the improved audio quality this offers.

There are several different wideband audio codecs that have broad market acceptance. Many SIP business phones and PC-based softphones support G.722, which is a relatively simple technology and royalty-free. We've already mentioned SILK, which Skype licenses at no cost. For next-gen mobile phones supporting Voice over LTE, 3GPP has standardized on the use of G.722.2, also known as Adaptive Multi-Rate Wideband (AMR-WB).

**Meta**switch
Networks

www.metaswitch.com

These codecs are all different and non-interoperable, so to support wideband audio between, say, a mobile phone using G.722.2 and a PC-based softphone using G.722, the network needs to insert a transcoder in the path of the call.

Audio transcoding resources can be integrated into the media plane functions of a Session Border Controller, or they can be deployed as separate elements in the core of the voice network. As with video transcoding, this function is provided most cost-effectively with specialized hardware. The difficulty with planning the deployment model for audio transcoding is the degree of uncertainty over how much total transcoding capacity will be required. This will depend on numerous factors including the relative market penetration of the different wideband codecs in both fixed and mobile endpoints, user expectations for interoperability of wideband audio, and service provider policy with respect to wideband audio services.

Currently, the prevailing view among service providers seems to be that audio transcoding is better deployed as a separate function in the core of the network, rather than integrated into the media plane function of SBCs.

2.3.    Hostile Network Activity

From time to time, individuals or groups with a particular agenda choose to attack Web sites via the Internet. There is plenty of freely available software technology for distributing viruses to PCs so as to create "botnets", and one of the uses of such botnets is to orchestrate various kinds of attack, such as distributed denial-of-service attacks which flood a site with so much traffic that it becomes unresponsive.

The exact same technology can be used to build botnets that attack voice and multimedia network services based on SIP. A range of attack techniques have already been witnessed in the field, ranging from simple brute force attacks that bombard the network with thousands of spurious SIP registration attempts to more insidious techniques that take over authenticated endpoints and send malformed or un-routable SIP requests.

Such attacks on SIP networks are a commonplace occurrence, and they are typically handled reasonably well by today's appliance-based SBCs. But these are mostly rather small-scale attacks, and are probably the result of individuals experimenting with

| Service Type | Signaling Plane Impact | Media Plane Impact | Comment |
|---|---|---|---|
| Page-mode messaging (e.g. SMS over IP) | At least 2 SIP messages per SMS leg (more if read report requested) | None | Required to support text messaging in LTE networks |
| Session-mode messaging (e.g. fixed or mobile IM) | From 5 to 10 SIP messages per session leg (depending on authentication model) | Varies widely according to session content – from tens of bytes per session for text only to megabytes for image or video sharing | Required to support Rich Communications Suite or Converged IP Messaging services |
| Presence | 2 SIP messages per status change per presence watcher | None | Required for RCS but not for RCS-e |
| Capability exchange using OPTIONS | 2 SIP messages per call leg for each voice or video call, instant messaging session or SMS message | None | Required for RCS-e |
| Forking | Typically 7 SIP messages per call per additional device to which voice calls are forked | None | Applies to voice calls, video calls and attempts to establish messaging sessions |
| Video calling | Same load as for voice calls | 20 to 200 times the data volume of a voice call, and may also require media processing e.g. transcoding | |
| Wideband voice calling | Same load as for voice calls | Similar loading to narrowband voice calls, but may require transcoding | Wideband codec required for VoLTE (AMR-WB) is not generally supported by desktop clients |

*Table 1 – Summary of signaling and media plane impacts of new service capabilities*

www.metaswitch.com

malicious software rather than making a concerted effort to bring a SIP network down. In any case, even if they were to bring a SIP network down, the overall impact would be limited since there are still relatively few subscribers being served by such networks, compared with the subscriber population as a whole.

As the number of subscribers served by SIP networks grows to many millions, these networks become a more interesting target both for organized cyber-terrorists and for individuals who bear some ill-will towards a particular service provider. And successful attacks on such large-scale SIP networks will have a much more widespread impact.

SIP networks that support fixed broadband access, and therefore permit devices such as PC-based softphones to access SIP services are obviously vulnerable to the same kind of botnet-based attack as the Web. These networks clearly need protecting well against such attacks.

The situation with mobile networks is less clear. Some mobile operators take the view that the security mechanisms in place for authentication of mobile IMS endpoints are sufficient to protect against the risk of botnet-based attacks. But history has shown that supposedly secure systems are viewed by certain types of individuals as a challenge that, sooner or later, human ingenuity will find a way to penetrate. The rapid rise to prominence of open source mobile software platforms such as Android would seem to provide an ideal culture medium for the kinds of virus development that have the potential to bring down mobile SIP-based voice and messaging services. Furthermore, many mobile phones will make use of WiFi as well as 4G mobile networks to access their IMS-based services, which means that these services will be exposed to the public Internet.

From the point of view of session border control, protection against hostile activity is almost entirely a challenge for the signaling plane. As the threat of large scale attacks grows, so the processing power required to handle legitimate SIP traffic in the face of hostile signaling activity needs to grow very substantially – but with no requirement for equivalent growth of media processing capacity.

## 3. A Distributed Model for Session Border Control

The evolving SIP services landscape described above is not well served by the current generation of integrated Session Border Controller products, because they embody signaling and media handling capacity in a fixed relationship that has been optimized for voice calling. As the service mix changes to include messaging, presence and forking services that dramatically drive up signaling load, and as the threat of malicious signaling activity grows, these SBCs will have to be severely down-rated in terms of the number of subscribers they can support. This will result in the stranding of a great deal of the media handling capacity in those SBCs. It will also mean a huge increase in the unit numbers of SBCs that need to be deployed in the network.

Given this, it is clear that an alternative architecture for the session border control function is urgently called for – one which supports independent scaling of signaling and media capacity. The distributed SBC architecture meets this requirement.

### 3.1.    The Distributed SBC Architecture

The distributed SBC architecture separates the signaling plane functions of session border control from the media plane functions, and enables the signaling and media planes to be implemented on physically separate elements that communicate with each other by means of a control protocol.

This is exactly the same principle as the softswitch architecture, which has dominated the transition of the PSTN to VoIP for well over a decade. A softswitch system separates signaling plane functions, such as signaling protocol interworking and call control, from media plane functions such as TDM to VoIP interworking. In the softswitch world, this separation permits the independent scaling of call control and media capacity, and also allows the signaling and media elements to be geographically distributed so as to benefit from centralization of complex configuration in the signaling plane – such as call routing – while avoiding the need for inefficient media backhaul.

The distributed SBC architecture is so similar in concept to the softswitch architecture that it is able to leverage the same standard protocol for the control connection between signaling plane and media plane functions, namely H.248.
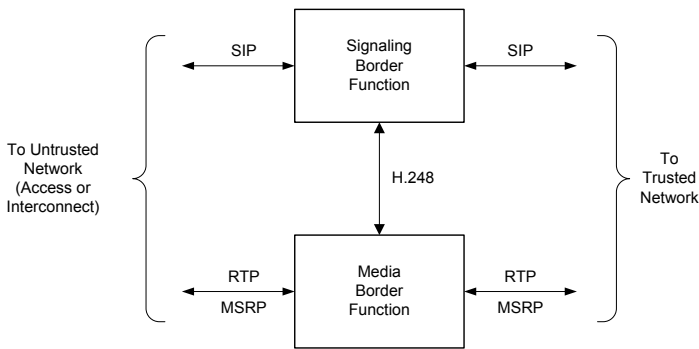
**Meta switch** Networks

*Figure 1 - Distributed SBC Architecture*

The 3GPP specifications for IMS describe the distributed SBC architecture and the usage of H.248 between signaling and media plane functions in detail. These specifications use terminology specific to IMS – for example, the signaling plane function for network interconnect is called the Interconnect Border Control Function (IBCF) – but the concepts are exactly as described in this white paper.

| | Access Border | Interconnect Border |
|---|---|---|
| IMS terminology for signaling plane function | IMS Application Layer Gateway (IMS-ALG) embedded in Proxy Call Session Control Function (P-CSCF) | Interconnect Border Control Function (IBCF) |
| IMS terminology for media plane function | IMS Access Gateway | Transition Gate Way (TrGW) |
| Designation for control interface reference point | Iq | Ix |
| 3GPP specification for H.248 interface | TS 29.334 | TS 29.238 |

*Table 2 - Distributed SBC architecture as described in 3GPP specifications for IMS*

In a distributed SBC, the signaling function handles all of the processing of SIP messages including authentication, topology hiding, denial-of-service attack detection and prevention, traversal of network address and port translation (NAPT), SIP header and SDP manipulation, session routing, billing event reporting and so on. When it determines that a session setup request has progressed to the point where a media path needs to be established, the signaling function sends an H.248 command to the media function to establish a "gate" for the relaying of RTP media. This command includes details of the bandwidth permitted

for the gate, address or port number translations needed to support NAPT traversal by the RTP media streams, and details of any media transcoding function that needs to be inserted in the path. When the session is torn down by the endpoints, the signaling function uses H.248 commands to collect statistics about the RTP session (e.g. packet loss, jitter and delay) from the relevant gate before closing it.

### 3.2. Benefits of the Distributed SBC Architecture

The distributed SBC architecture offers a wide range of compelling advantages over traditional appliance-based integrated Session Border Controllers. These advantages translate into lower capital costs for service providers who are faced with the need to massively scale up their SBC deployments to support a rapidly growing population of SIP endpoints, as well as lower operational costs arising from the presence of fewer, more centralized session border signaling functions to manage. The rest of this section explains how these benefits arise.

### 3.2.1. Independent scalability of signaling and media

We have discussed at length the factors that are leading to vastly different trends in signaling load versus media load in real-world SIP networks. The single most overriding advantage of the distributed SBC architecture is that it permits service providers to grow signaling border control, media border control and transcoding capacity in their networks according to the actual respective demand for each, rather than in the fixed ratio that is imposed by integrated appliance-based SBCs. Since each of these functions represents a very substantial cost element in the deployment of VoIP networks, service providers cannot afford to pursue SBC strategies based on fixed ratios of signaling, media and transcoding that make highly inefficient use of the capacity of any of these elements.

### 3.2.2. Geo-distributed deployment of signaling and media

In the distributed SBC architecture, the signaling border control and media border control functions are logically separate and communicate via a standard protocol. This allows these two functions to be placed in geographically separated locations if appropriate.

The flexibility to deploy signaling border control and media border control in separate locations is a major advantage of the distributed architecture. It permits the signaling border control function to be concentrated near the core of the network in a small number of powerful systems, which control a larger number of media border control functions that are deployed close to the edge of the network. This simplifies the task of managing the complex configuration of the session border control function – which exists entirely in the signaling border controllers – while permitting media paths to be optimized and media backhaul to be minimized.

www.metaswitch.com

Geographically distributed deployment of session border control is typically most appropriate for the access SBC function, because of the highly distributed nature of the network access edge. For the interconnect SBC function, where large numbers of signaling and media connections are concentrated at points of interconnect between carriers, co-location of the signaling and media border control functions is likely to make more sense.

One argument that is sometimes advanced against the idea of centralizing the signaling border function is that it allows malicious signaling traffic to penetrate deeper into a service provider's network than would be the case with appliance-based SBCs deployed at the edge. While this is undoubtedly true, it should be remembered that SIP signaling consumes, on average, two orders of magnitude less bandwidth than the RTP media with which it is associated. The extra network resources that might be consumed by malicious traffic in the event of an attack on a signaling border function in the centralized case are therefore relatively insignificant.

### 3.2.3. Media path optimization

SIP endpoints such as IP business phones, PC-based softphones and mobile SIP clients are typically provisioned with the SIP URI of the access SBC to which they must send SIP signaling to obtain their services. With an integrated SBC, the media from a given endpoint is forced always to traverse the media border function that is co-resident with the signaling border function whose identity they have been provisioned with. Where a SIP endpoint is nomadic and is accessing the network from multiple different locations over time, this may mean that the media path is sometimes routed very inefficiently.

In a distributed SBC architecture, signaling border controllers and media border controllers may have a many-to-many control relationship. This means that, while a nomadic SIP endpoint will typically always send its signaling to a particular signaling border controller, that signaling border controller is free to choose a media border controller which is physically close to the current location of the SIP endpoint. By this means, the media path can be optimized on a call-by-call basis so that it minimizes usage of transport resources and maximizes quality of experience for the user.
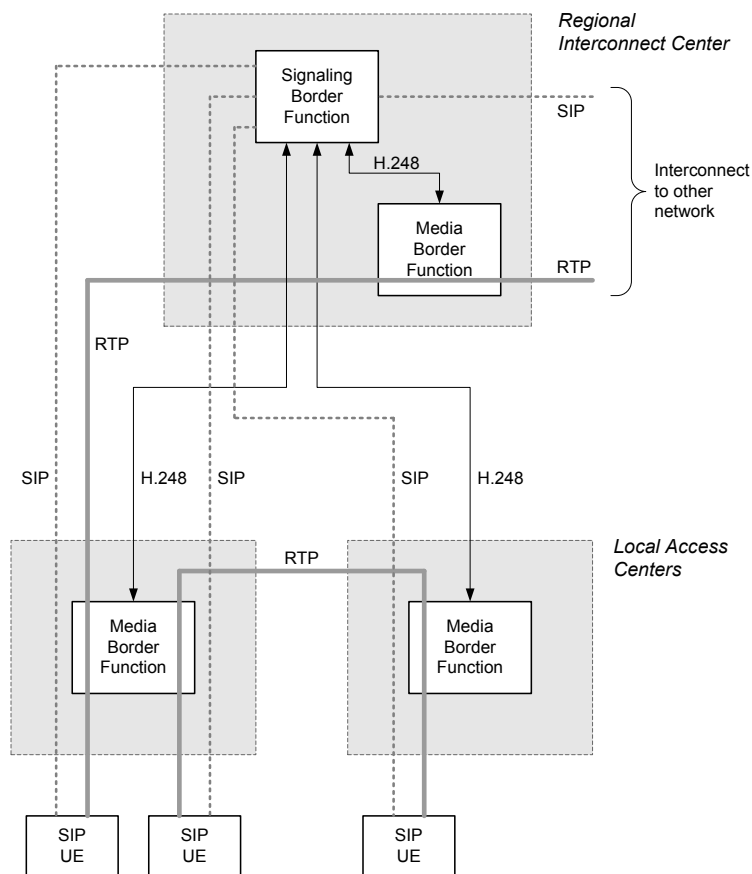


Figure 2 - Example of distributed session border control deployment, showing on-net and off-net calls
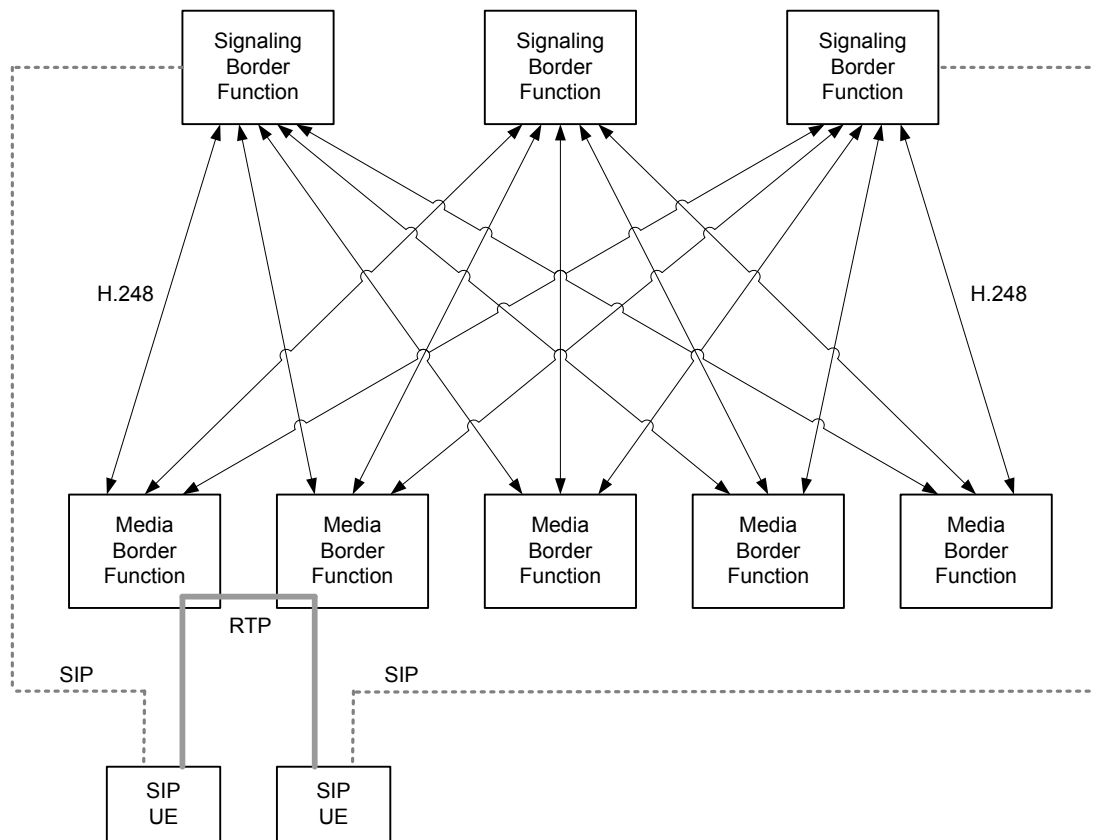
www.metaswitch.com

*Figure 3 - Media path optimization in a large scale session border control deployment*

3.2.4.    De-duplication of media border function

In the general case, in a large network, a signaling session between two SIP endpoints traverses two signaling controllers.  This is unavoidable because each SIP endpoint maintains a signaling relationship with the network via a given signaling border controller instance.

With integrated appliance-based SBCs, each signaling border control function is tightly coupled (and co-resident) with its associated media border control function.  As a result, the media path always transits the network via the same two SBC devices as the signaling path.

The media border control function is responsible for policing the bandwidth characteristics of media streams, performing network address and port number manipulation to support NAPT traversal, and collecting statistics about the quality of the media session.  If the media path between two endpoints transits via two integrated SBC devices, then the media border control function is applied to the media path twice.  This is actually wasteful and unnecessary.

On any given media path between two SIP endpoints, a single media border controller is, in principle, able to perform all of the bandwidth policing, network address and port number translation, and statistics collection that is required.   There is no fundamental reason why the media border control function should be performed twice – the fact that integrated appliance-based SBCs do so is simply a limitation of their architecture.

With the distributed SBC architecture, it is possible for the two signaling border controllers involved in any given session setup request to co-operate in such a way as to ensure that only one media border controller is inserted in the media path within the service provider's network.  This "single hop media relay" technique effectively halves the total amount of media border controller capacity that is required in the network.   Given that the media plane function typically represents at least two-thirds of the cost of session border control, a 50% reduction in required media plane capacity translates into at least 30% reduction in the total cost of SBC deployments.
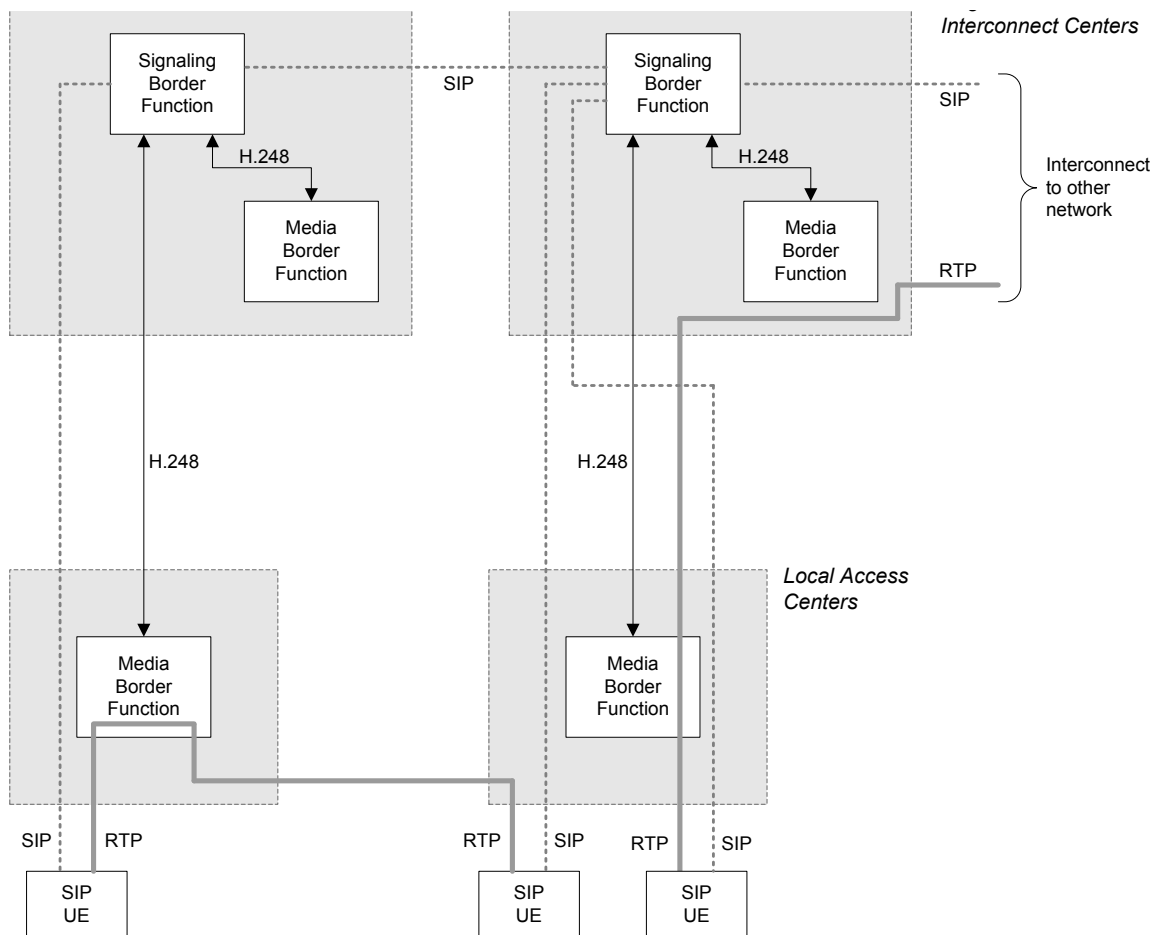
www.metaswitch.com

*Figure 4 - De-duplication of Media Border Function, showing both on-net and off-net calls*

### 3.2.5.    Cost reduction of high-availability configurations

Session Border Controllers are usually deployed in a high-availability configuration to protect against hardware or software failures, which would otherwise result in loss of SIP-based services.  With integrated appliance-based SBCs, a 1:1 model is invariably used to provide high availability.  In this model, signaling and session state is replicated between the active and passive members of the pair, so that hitless failover can be achieved in both signaling and media planes.

The 1:1 architecture doubles the hardware cost of the SBC solution.  Depending on how the SBC software is licensed, it may also double the overall cost of the SBC solution.  High availability can therefore be an expensive option for SBCs, but most service providers consider it essential.

In reality, hitless failover is more important in the signaling plane than it is in the media plane.  Users are generally tolerant of occasional unexpected disconnections of a session.  But they expect that "dial tone" is always available – so if a call is dropped, they can re-connect immediately by just re-dialing.

The distributed SBC architecture enables a more cost-effective model for high availability which is particularly appropriate in larger networks.  In this HA model, the signaling border control function is protected on a 1:1 basis while the media border control function is protected on an N:1 basis.  If an individual media border controller instance fails, all calls that are being handled by that instance will be dropped.  But when users attempt to re-connect, the signaling border control function will establish the session via one of the remaining active media border controller instances, so service is still available.

In a large scale SBC deployment, this approach can reduce the total cost of the media border function by almost 50%.  The N:1 approach for high availability in the media border function can be combined with the single hop media relay approach described in the previous section to achieve an overall reduction of almost 75% in the total cost of the media border function.

www.metaswitch.com

3.2.6.    Open interoperability between signaling & media elements

Just as the softswitch architecture supports open interoperability between call agents and media gateways from different vendors, so the distributed SBC architecture supports open interoperability between signaling border controllers and media border controllers from different vendors.

While this might appear to be more of a theoretical advantage than a real one, it is worth pointing out that certain router vendors have integrated a media border controller function in some of their router products, supporting the standard H.248 control interface. With a separate signaling border controller driving the embedded media border control function in an edge router, service providers can leverage the excellent packet processing hardware capabilities of their routers while reducing the number of separate elements they need in their VoIP networks.

The distributed SBC architecture therefore supports additional deployment options that enable service providers to better optimize the session border control function in their specific network environments.

## 4.    Scalability of Signaling Capacity

A typical integrated SBC appliance of the kind that is widely deployed in networks today can support at most 100,000 subscribers, assuming a signaling load that arises entirely from voice calling, and also assuming a relatively modest busy hour calling load.  When this is adjusted to allow for the extra signaling load associated with SIP-based messaging, presence and call forking, and to provide adequate head-room to cope with concerted denial of service attacks, an absolute upper limit of 25,000 subscribers per unit should be considered much more realistic. Given that integrated SBCs need to be deployed in pairs to provide a high-availability solution, this implies that a network with 10M VoIP subscribers would need 800 SBC appliances to handle access alone.  Depending on the percentage of on-net traffic, a comparable (though smaller) additional quantity would be required to support interconnect.

Dealing with the logistics of installing, configuring, administering and performing software upgrades on a device population of this size is a truly daunting prospect.  The sheer quantity of devices also introduces some network design challenges, requiring for example the introduction of large-scale SIP-aware load-balancing functions in order to hide the existence of so many different IP addresses from SIP endpoints.

An analysis of the hardware on which this type of device is based reveals the reason for the modest performance:  the typical integrated SBC appliance is equipped with a CPU comparable in power to that in a budget laptop.

Carrier-class hardware solutions such as ATCA-based server blades are available with at least an order of magnitude more processing power than these appliances.  To achieve such high throughput, these processors typically incorporate 12 or more CPU cores. Note that running SBC software designed for a single or dual-core processor environment on these massively parallel systems won't generally deliver substantial performance gains.  The software has to be designed from the ground up to take full advantage of a multi-core processor architecture.

The carrier-class SBC solutions of the future will not only be based on a distributed architecture, but they will also leverage state-of-the-art multi-CPU commodity processing power.  This is the only rational way to achieve the levels of signaling scalability that are necessary to cope with anticipated signaling traffic loads with a reasonably manageable number of network elements.

## 5.    Conclusion

The separation of signaling and media functions is a well-accepted principle of all next-generation voice and multimedia networks, up to and including IMS.   This principle recognizes the very important differences between the types of work performed by network elements in the signaling path and media path respectively, and the value of concentrating signaling elements centrally in the network (for ease of management) while permitting media to follow the shortest path between the two endpoints of a call.

As we have seen, a rational analysis of the requirements for session border control, informed by an appreciation of clearly apparent trends in evolving SIP service mix, leads to one inevitable conclusion:  a distributed architecture for session border control is the only approach that really makes sense.

The current popularity of integrated appliance-based Session Border Controllers appears to fly in the face of this principle, and we can only speculate about the reasons for this apparent anomaly. Session border control is a very complex area of technology, and there are few products on the market that do a really good job of it. Perhaps the vendors who have been most successful in the market just happen to have promoted the integrated appliance-based architecture – and their products have succeeded not because they have the best architecture, but because they have implemented the most complete feature set and delivered the best software quality.

As the SBC market matures, service providers will benefit from access to a wider range of products and solutions that meet all their functional and quality requirements for session border control.  In these circumstances, SBC architecture and scalability become critical differentiators – and the benefits of a distributed and highly scalable architecture for session border control are so overwhelming that no service provider can afford to overlook them.

www.metaswitch.com

## Glossary

3GPP    Third Generation Partnership Project

AMR-WB Adaptive Multi-Rate Wide Band

ATCA    Advanced Telecommunications Computing Architecture

CPM    Converged IP Messaging

IMS    IP Multimedia Subsystem

IP    Internet Protocol

LTE    Long Term Evolution

MMS    Multimedia Message Service

MSRP    Message Session Relay Protocol

NAPT    Network Address and Port Translation

PBX    Private Branch eXchange

RCS    Rich Communications Suite

RFC    Request for Comments

RTP    Real Time Protocol

SBC    Session Border Controller

SIP    Session Initiation Protocol

SMS    Short Message Service

TDM    Time-Division Multiplexed

VoIP    Voice over Internet Protocol

URI    Uniform Resource Identifier

VoLTE    Voice over Long Term Evolution

XMPP    eXtensible Messaging and Presence Protocol

Martin Taylor has spent over 20 years in the telecom and network equipment industries, with diverse experience in product marketing, engineering, technology planning and business development. Since joining Metaswitch in 2004, he has been responsible for developing the company's IMS technology strategy and has also led key technology initiatives in hosted VoIP services, Web Services interfaces and Web portal evolution.

**Meta**switch Networks

www.metaswitch.com