

Titanium Server

Server Group Messaging

Host-Guest Message-Based API & Guest Reference Implementation

Release: 15.12
01/14/2016

WIND RIVER





Copyright Notice

Copyright (c) 2013-2016, Wind River Systems, Inc.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- 1) Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- 2) Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- 3) Neither the name of Wind River Systems nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.



Table of Contents

Introduction.....	3
Host – Guest Server Group Messaging API	4
Message Types and Semantics.....	5
Virtio Serial Device	7
JSON Message Syntax	9
Base JSON Message Layer – Syntax.....	10
Application JSON Message Layer – Syntax.....	11
Examples.....	15
Reference Implementation of Guest Server Group Messaging	17



Introduction

Titanium Server implements a Host-to-Guest Server Group Messaging API to provide a simple low-bandwidth datagram messaging and notification service for servers that are part of the same server group. This messaging channel is available regardless of whether IP networking is functional within the server, and it requires no knowledge within the server about the other members of the group. This document contains the specification for this messaging-based API.

Also included in this document is an overview of the Titanium Server 'Server Group Messaging' SDK Module which provides a Linux-based reference implementation of the Guest-side software for implementing this Messaging-based API in the Guest. The SDK Module provides source code and make/build instructions which can be used strictly as reference or built and included 'as is' in your Guest image. Full build, install and usage instructions can be found in the README files included in the SDK Module. This document simply provides an overview of the reference implementation.



Host – Guest Server Group Messaging API

Titanium Server implements a simple Host-to-Guest Server Group Messaging API to provide a simple low-bandwidth datagram messaging and notification service for servers that are part of the same server group. This messaging channel is available regardless of whether IP networking is functional within the server, and it requires no knowledge within the server about the other members of the group.

The Host-to-Guest Server Group Messaging API is a message-based API using a JSON-formatted application messaging layer on top of a ‘virtio serial device’ between QEMU on the host and the Guest VM. JSON formatting provides a simple, humanly readable messaging format which can be easily parsed and formatted using any high level programming language being used in the Guest VM (e.g. C, Python, Java, etc.). Use of the ‘virtio serial device’ provides a simple, direct communication channel between host and guest which is independent of the Guest’s L2/L3 networking.

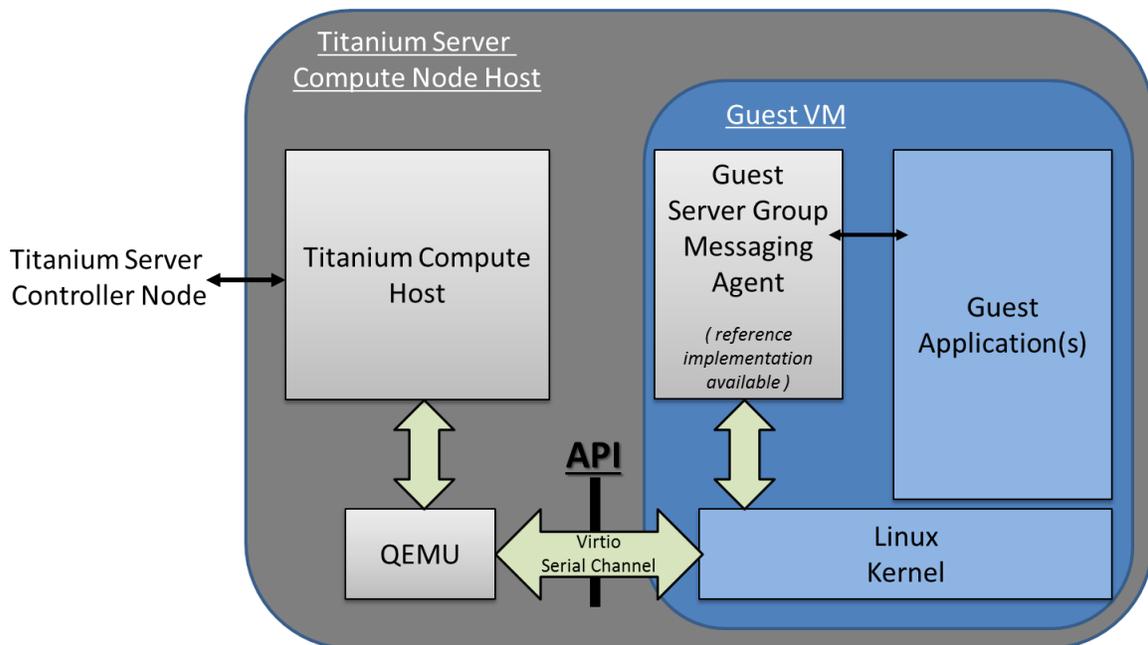


Figure 1 – Titanium Server Host-Guest Server Group Messaging API



Message Types and Semantics

For the Server Group Messaging API, there are four message types; Server Status Query/Response Messages, Asynchronous Server Status Change Notifications, Server Broadcast Messages and a Nack Message.

- Status Query (Guest → Host)
 - This allows a server (Guest) to query the current state of all servers within its server group, including itself,

- Status Response and Status Response Done (Host → Guest)
 - This is the Status Response from the Titanium Server Host containing the current state of all servers within the Guest's server group, including this Guest,
 - This is a multiple message response with each response containing the status of a single server, followed by a final response (status response done) with no data,
 - Each message of the multiple message response has the transaction number (seq) of the status query request that it is related to,

- Notification Message (Host → Guest)
 - This provides the server (Guest) with information about changes to the state of other servers within the server group,
 - Each notification message contains the status of a single server,

- Broadcast Message (Guest → Host) → (Host → Guest)
 - This allows a server (Guest) to send a datagram (thru the Host, with a size of up to 3050 bytes) to all other servers (Guests) within its server group,
 - the payload portion of the message, 'data', can be formatted as desired by the Guest, however it must be a null-terminated string without embedded newlines,
 - the source field of the message is a unique, although opaque, address string representing the server (Guest) that sent the message.



- Nack (Host → Guest)
 - This is a message sent from the Host to the Guest when the Host receives a message with incorrect syntax,
 - It contains the message type of the original (incorrect) message and a log_msg describing the error,
 - This allows the Guest Application developer to debug issues when developing the Guest-side API code.

This service is not intended for high bandwidth or low-latency operations. It is best-effort, not reliable. Applications should do end-to-end acks and retries if they care about reliability.



Virtio Serial Device

The transport layer of the Host-Guest Server Group Messaging API is a 'virtio serial device' (also known as a 'vmchannel') between QEMU (on the host) and the Guest VM. Device emulation in QEMU presents a virtio-pci device to the Guest, and a Guest Driver presents a char device interface to Guest userspace applications. This provides a simple transport mechanism for communication between the host userspace and the guest userspace. I.e. it is completely independent of the networking stack of the Guest, and is available very early in the boot sequence of the Guest.

This is a standard Linux QEMU/KVM feature. The Guest API for interfacing with the 'virtio serial device' can be found at http://www.linux-kvm.org/page/Virtio-serial_API. Examples of Guest code for opening, reading, writing, etc. from/to a 'virtio serial device' can also be found in the source code of the Titanium Server 'Server Group Messaging SDK Module'. This SDK Module provides a Linux-based reference implementation of the Guest-side software for implementing the Guest Serer Group Messaging API. Generally communicating with a 'virtio serial device' is very similar to communicating via a pipe, or a SOCK_STREAM socket.

There are however a few additional considerations to be aware of when using 'virtio serial devices':

- only one process at a time can open the device in the Guest,
- read() returns 0, if the Host is not connected to the device,
- write() blocks or returns -1 with error set to EAGAIN, if the Host is not connected,
- poll() will always set POLLHUP in revents when the Host connection is down.
 - This means that the only way to get event-driven notification of connection is to register for SIGIO. However, then a SIGIO event will occur every time the device becomes readable. The work-around is to selectively block SIGIO as long as the link is up is thought to be up, then unblock it on connection loss so a notification occurs when the link comes back.
- If the Host disconnects the Guest should still process any buffered messages from the device,
- Message boundaries are **not preserved**, the Guest needs to handle message fragment reassembly. Multiple messages can be returned in one read() call, as well as buffers beginning and ending with partial messages. This is hard to get perfect; one can study the `host_guest_msg.c` code in the Titanium Server Guest Server Group Messaging SDK Module for ideas on how this can be handled.



The QEMU/KVM created by Titanium Server in order to host a Guest VM is created with a 'virtio serial device' named:

```
/dev/virtio-port/cgcs.messaging
```

for general Titanium Server Host – to – Guest VM messaging (e.g. Host-Guest Server Group Messaging as well as other Host-Guest Messaging discussed in other Titanium Server – Guest API documents).



JSON Message Syntax

The upper layer messaging format being used is ‘Line Delimited JSON Format’. I.e. a ‘\n’ character is used to identify message boundaries in the stream of data to/from the virtio serial device; specifically a ‘\n’ character is inserted at the start and end of the JSON Object representing a Message.

```
\n{key:value,key:value,...}\n
```

Note that key and values must NOT contain ‘\n’ characters.

The upper layer messaging format is actually structured as a hierarchical JSON format containing a Base JSON Message Layer and an Application JSON Message Layer:

- the Base Layer provides the ability to multiplex different groups of message types on top of a single ‘virtio serial device’
e.g.
 - resource scaling,
 - server group messaging,
 - etc.

and

- the Application Layer provides the specific message types and fields of a particular group of message types.



Base JSON Message Layer – Syntax

Again, the Base Layer provides the ability to multiplex different groups of message types on top of a single ‘virtio serial device’, e.g. resource scaling versus server group messaging etc.

Host – to – Guest Messages

Key	Value	Optionality*	Example value (for Server Group Messaging)	Description
“version”	integer	M	1	Version of the Base Layer Messaging
“source_addr”	string	M		Opaque string representing the host-side address of the message.
“dest_addr”	string	M	“cgcs.server_grp”	The Guest-side addressing of the message; specifically the Message Group Type
“data”	JSON Formatted String	M	See the following section on Application Layer JSON Message Layer – Syntax for Server Group Messaging.	Application layer JSON message whose schema is dependent on the particular Message Group Type

- M: Mandatory; O: Optional

Guest – to – Host Messages

Guest – to – Host Messages, from a Base Layer perspective, are identical to Host – to – Guest Messages except for swapped semantics of source_addr and dest_addr.



Application JSON Message Layer – Syntax

Again the Application Layer provides the specific message types and fields of a particular group of message types; in this case the messages of Server Group Messaging.

Guest – to – Host Messages

Status Query

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface.
“msg_type”	string	M	“status_query”	Type of the message.
“seq”	integer	M		Transaction number for the query; corresponding status_response and status_response_done messages will have a matching transaction number. This should be incremented on each status_query sent by Guest.

- M: Mandatory; O: Optional; (Condition)

Broadcast Message

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface
“msg_type”	string	M	“broadcast”	Type of the message.
“data”	string	M		Message content; can be formatted as desired by the Guest, however it must be a null-terminated string without embedded newlines.

- M: Mandatory; O: Optional; (Condition)



Host – to – Guest Messages

Status Response

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface.
“msg_type”	string	M	“status_response”	Type of the message.
“seq”	integer	M		Transaction number that the response belongs to.
“data”	string	M	see following info following table	The JSON formatted field containing the same contents as the normal notification that gets sent out by OpenStack’s notification service; see example below.

- M: Mandatory; O: Optional; (Condition)

Example contents of ‘data’ field containing status of a particular server:
(the same contents as the normal notification that gets sent out by OpenStack’s notification service)

```
{
  "state_description": "",
  "availability_zone": null,
  "terminated_at": "",
  "ephemeral_gb": 0,
  "instance_type_id": 10,
  "deleted_at": "",
  "reservation_id": "r-ed4i0c72",
  "instance_id": "4c074ce9-cbde-4040-9fdb-84b36168916b",
  "display_name": "jd_af_vm1",
  "hostname": "jd-af-vm1",
  "state": "active",
  "progress": "",
  "launched_at": "2015-11-26T14:33:03.000000",
  "metadata": {
  },
  "node": "compute-0",
  "ramdisk_id": "",
  "access_ip_v6": null,
  "disk_gb": 1,
  "access_ip_v4": null,
  "kernel_id": "",
  "host": "compute-0",
  "user_id": "369b0103310d4a6bbf43ed389aac211d",
  "image_ref_url": "http://127.0.0.1:9292/images/32b386e1-5a21-47c4-a04a-57910e7b0fc8",
  "cell_name": "",
  "root_gb": 1,
}
```



```

"tenant_id":"98b5838aa73c40728341336852b07772",
"created_at":"2015-11-26 14:32:51.431455+00:00",
"memory_mb":512,
"instance_type":"jdlcpu",
"vcpus":1,
"image_meta":{
  "min_disk":"1",
  "container_format":"bare",
  "min_ram":"0",
  "disk_format":"qcow2",
  "base_image_ref":"32b386e1-5a21-47c4-a04a-57910e7b0fc8"
},
"architecture":null,
"os_type":null,
"instance_flavor_id":"101"
}

```

Status Response Done

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface.
“msg_type”	string	M	“status_response_done”	Type of the message.
“seq”	integer	M		Transaction number that the response belongs to.

- M: Mandatory; O: Optional; (Condition)

Notification Message

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface.
“msg_type”	string	M	“notification”	Type of the message.
“data”	string	M	see contents of ‘data’ field documented for status_response	The JSON formatted output of the response to the Compute API GET /<version>/<tenant_id>/servers/<server_id> (see Compute API documentation for exact contents of response)

- M: Mandatory; O: Optional; (Condition)



Broadcast Message

Key	Value	Optionality*	Example value	Description
“version”	integer	M	1	Version of the interface
“msg_type”	string	M	“broadcast”	Type of the message.
“source_instance”	string	M		The unique, although opaque, address string representing the server (Guest) that sent the message.
“data”	string	M		Message content; can be formatted as desired by the Guest, however it must be a null-terminated string without embedded newlines.

- M: Mandatory; O: Optional; (Condition)

Nack

Key	Value	Optionality*	Example value	Description
“version”	integer	M	2	Version of the interface
“msg_type”	“nack”	M	“nack”	The type of message.
“orig_msg_type”	string	M	“broadcast”	The type of message that host previous received from guest.
“log_msg”	string	M	“failed to parse version”	Error message

- M: Mandatory; O: Optional; (Condition)



Examples

Examples of ‘full’ Server Group Messaging JSON messages, containing the Application JSON Message Layer encapsulated inside the Base JSON Messaging Layer.

Status Query:

Guest sends a query to TiS for status of all servers in Guest’s Server Group:

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server_grp",
"data":{"version":1,"msg_type":"status_query","seq":1}}\n
```

TiS responds with the status of a server in the Guest’s Server Group; one or more messages, each containing the status of one server in the Guest’s Server Group:

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server
_grp","data":{"version":1,"msg_type":"status_response","seq":1,"data":{"
state_description": "", "availability_zone": null, "terminated_at":
"", "ephemeral_gb": 0, "instance_type_id": 10, "deleted_at": "",
"reservation_id": "r-ed4i0c72", "instance_id": "4c074ce9-cbde-4040-
9fdb-84b36168916b", "display_name": "jd_af_vm1", "hostname": "jd-af-
vm1", "state": "active", "progress": "", "launched_at": "2015-11-
26T14:33:03.000000", "metadata": { }, "node": "compute-0",
"ramdisk_id": "", "access_ip_v6": null, "disk_gb": 1, "access_ip_v4":
null, "kernel_id": "", "host": "compute-0", "user_id":
"369b0103310d4a6bbf43ed389aac211d", "image_ref_url":
"http://127.0.0.1:9292/images/32b386e1-5a21-47c4-a04a-
57910e7b0fc8", "cell_name": "", "root_gb": 1, "tenant_id":
"98b5838aa73c40728341336852b07772", "created_at": "2015-11-26
14:32:51.431455+00:00", "memory_mb": 512, "instance_type": "jd1cpu",
"vcpus": 1, "image_meta": { "min_disk": "1", "container_format":
"bare", "min_ram": "0", "disk_format": "qcow2", "base_image_ref":
"32b386e1-5a21-47c4-a04a-57910e7b0fc8" }, "architecture": null,
"os_type": null, "instance_flavor_id": "101" }}}\n
```

TiS responds with response done for the current outstanding query request; with no data:

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server
_grp","data":{"version":1,"msg_type":"status_response_done","seq":1}}\n
```



Notification:

A notification of a server state change from TiS:

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server_grp",  
"data":{"version":1,"msg_type":"notification","data":{"state_description":  
"", "availability_zone": null, "terminated_at": "", "ephemeral_gb": 0,  
"instance_type_id": 10, "deleted_at": "", "reservation_id": "r-ed4i0c72",  
"instance_id": "4c074ce9-cbde-4040-9fdb-84b36168916b", "display_name":  
"jd_af_vml", "hostname": "jd-af-vm1", "state": "active", "progress": "",  
"launched_at": "2015-11-26T14:33:03.000000", "metadata": { }, "node":  
"compute-0", "ramdisk_id": "", "access_ip_v6": null, "disk_gb": 1,  
"access_ip_v4": null, "kernel_id": "", "host": "compute-0", "user_id":  
"369b0103310d4a6bbf43ed389aac211d", "image_ref_url":  
"http://127.0.0.1:9292/images/32b386e1-5a21-47c4-a04a-57910e7b0fc8",  
"cell_name": "", "root_gb": 1, "tenant_id":  
"98b5838aa73c40728341336852b07772", "created_at": "2015-11-26  
14:32:51.431455+00:00", "memory_mb": 512, "instance_type": "jdlcpu", "vcpus":  
1, "image_meta": { "min_disk": "1", "container_format": "bare", "min_ram":  
"0", "disk_format": "qcow2", "base_image_ref": "32b386e1-5a21-47c4-a04a-  
57910e7b0fc8" }, "architecture": null, "os_type": null, "instance_flavor_id":  
"101" }}\n
```

Broadcast:

A broadcast message to/from another server:

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server_grp",  
"data":{"version":1,"msg_type":"broadcast","source_instance":"instance-  
00000001","data":"Hello World"}}\n
```

Nack:

A Nack from TiS for an invalid broadcast message sent from Guest.

```
\n{"version":1,"source_addr":"cgcs.server_grp","dest_addr":"cgcs.server_grp",  
"data":{"version":1,"msg_type":"nack","orig_msg_type":"broadcast","log_msg":  
"failed to parse version"}}\n
```



Reference Implementation of Guest Server Group Messaging

This section provides an overview of the Linux-based reference implementation of the Guest-side software for implementing this Host-to-Guest Server Group Messaging API in the Guest.

This reference implementation can be found in the Titanium Server Guest Server Group Messaging SDK Module. This Module provides source code and make/build instructions which can be used strictly as reference or built and included ‘as is’ in your Guest image. Full build, install and usage instructions can be found in the README files included in the SDK Module. This section simply provides an overview of the reference implementation.

The diagram below provides the architecture diagram of the reference implementation:

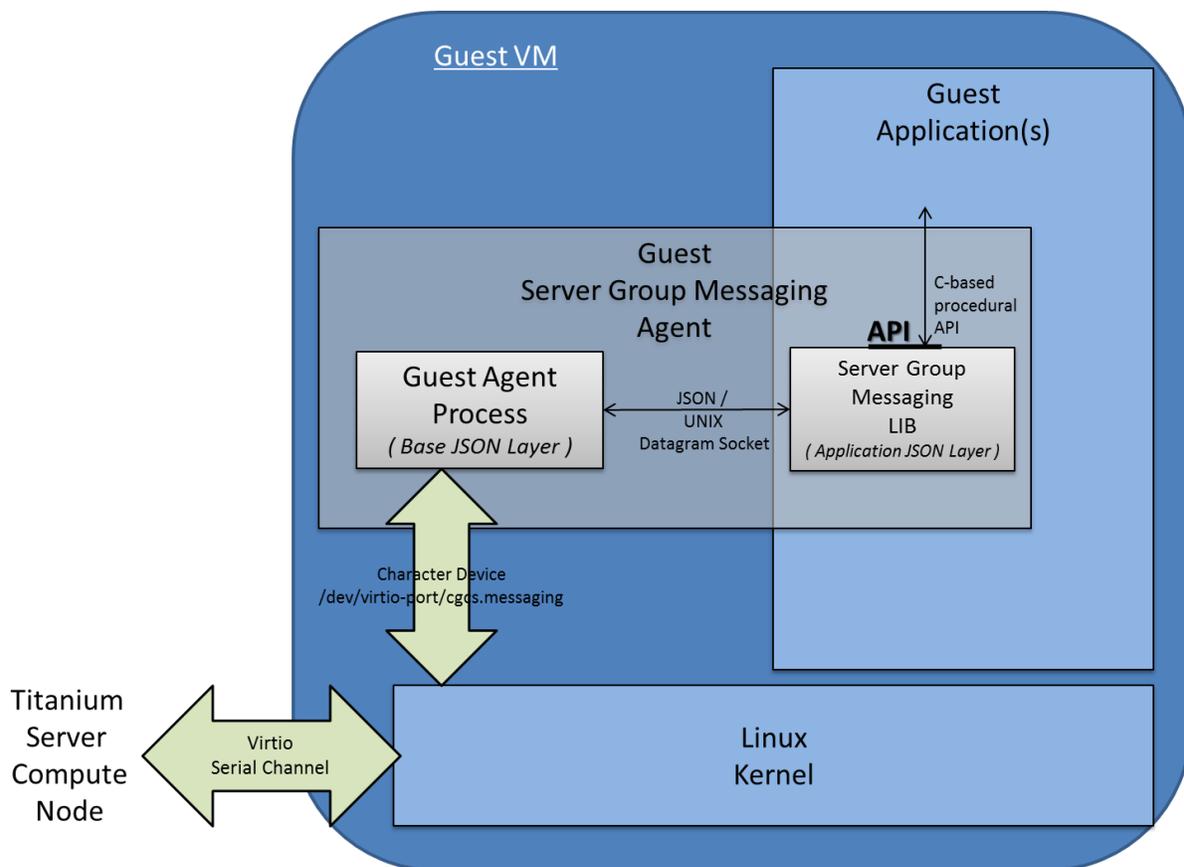


Figure 2 – Reference Implementation Architecture for Guest Server Group Messaging



Where:

- A Guest Agent Process implements the Base JSON Messaging Layer.
This includes:
 - opening/reading,/writing and general management of the virtio serial device between the Guest and the Host,
 - parsing/processing/formatting of the Base JSON Messaging Layer of the Guest-Host interface, where processing of the messages involves:
 - the multiplexing/de-multiplexing of Application Layer messages to/from registered Guest Application Layer Agents; in this particular case a Guest Application Process responsible for handling Server Group Messaging for the Guest,
 - the interface between the Guest Agent Process and the Guest Application Process responsible for Server Group Messaging for the Guest:
 - is a message-based interface;
 - specifically a JSON Messaging Layer over a UNIX Datagram socket,
 - where the UNIX Socket Address is the Message Group Type (cgcs.server_grp in this particular case) specified within the Base JSON Messaging Layer and
 - where the JSON Message consists of the ‘data’ field contents specified within the Base JSON Messaging Layer.
 - NOTE
 - The implementation files for the Guest Agent Process within the Titanium Server Guest Server Group Messaging SDK Module are:
 - misc.h, guest_host_msg.h, host_guest_msg_type.h,
 - guest_agent.c, host_guest_msg.c, lib_guest_host_msg.c
- A Server Group Messaging Lib which provides a C-based procedural API for a Guest Application Process to interface with the Guest Agent Process.

Specifically this library implements:

- the interface described above; a JSON Messaging Layer over a UNIX Datagram socket.
 - where the UNIX Socket Address is the Message Group Type (cgcs.server_grp in this particular case) specified within the Base JSON Messaging Layer and
 - the JSON Message consists of the ‘data’ field contents specified within the Base JSON Messaging Layer.
- with a C-based procedural API.



○ NOTE

- the definition and implementation of the Server Group Messaging Lib within the Titanium Server Guest Server Group Messaging SDK Module are:

- server_group.h and server_group.c
- server_group_app.c (*a sample usage of the API*)

```
server_group.h:

/* Function signature for the server group broadcast messaging callback function.
 * source_instance is a null-terminated string of the form "instance-xxxxxxx".
 * The message contents are entirely up to the sender of the message.
 */
typedef void (*sg_broadcast_msg_handler_t)(const char *source_instance,
                                           const char *msg, unsigned short msglen);

/* Function signature for the server group notification callback function. The
 * message is basically the notification as sent out by nova with some information
 * removed as not relevant. The message is not null-terminated, though it is
 * a JSON representation of a python dictionary.
 */
typedef void (*sg_notification_msg_handler_t)(const char *msg, unsigned short msglen);

/* Function signature for the server group status callback function. The
 * message is a JSON representation of a list of dictionaries, each of which
 * corresponds to a single server. The message is not null-terminated.
 */
typedef void (*sg_status_msg_handler_t)(const char *msg, unsigned short msglen);

/* Get error message from most recent library call that returned an error. */
char *sg_get_error();

/* Allocate socket, set up callbacks, etc. This must be called once before
 * any other API calls.
 *
 * Returns a socket that must be monitored for activity using select/poll/etc.
 * A negative return value indicates an error of some kind.
 */
int init_sg(sg_broadcast_msg_handler_t broadcast_handler,
            sg_notification_msg_handler_t notification_handler,
            sg_status_msg_handler_t status_handler);

/* This should be called when the socket becomes readable. This may result in
 * callbacks being called. Returns 0 on success.
 * A negative return value indicates an error of some kind.
 */
int process_sg_msg();

/* max msg length for a broadcast message */
#define MAX_MSG_DATA_LEN 3050

/* Send a server group broadcast message. Returns 0 on success.
```



```
 * A negative return value indicates an error of some kind.
 */
int sg_msg_broadcast(const char *msg);

/* Request a status update for all servers in the group.
 * Returns 0 if the request was successfully sent.
 * A negative return value indicates an error of some kind.
 *
 * A successful response will cause the status_handler callback
 * to be called.
 *
 * If a status update has been requested but the callback has not yet
 * been called this may result in the previous request being cancelled.
 */
int sg_request_status();
```