

Solaris Hotplug Framework: Architecture and Design

Revision 1.2

June 5, 2009

Colin.Zou@Sun.COM

Govinda.Tatti@Sun.COM

SHP Project Team

Sun Microsystems, Inc.

Table of Contents

1. Introduction.....	3
1.1 Background & Motivation.....	3
1.2 Goal.....	3
1.3 Terminologies.....	3
1.4 Scope and Road Map.....	4
1.5 References.....	4
2. Architecture Overview.....	5
2.1 Physical Hotplug & Virtual Hotplug.....	5
2.1.1 Concepts.....	5
2.1.2 Background & Goal.....	5
2.2 Overall Architecture.....	6
2.3 Connector & Port.....	7
2.4 Naming of Connection.....	7
2.5 Hotplug States & State Machine.....	8
2.5.1 Overview.....	8
2.5.2 Connector States.....	8
2.5.3 Port and Device States.....	8
2.5.4 State Machine.....	9
2.5.5 A Usage Example in Virtualized Environment.....	11
2.6 Bus Specific Commands Support.....	11
2.7 Events.....	11
2.8 cfgadm shp Plugin.....	12
3. Interfaces.....	12
3.1 modctl Commands.....	12
3.2 Libdevinfo interfaces.....	13
3.3 Hotplug Controller Driver Interfaces.....	13
3.3.1 ndi_hp_register() & ndi_hp_unregister().....	13
3.3.2 ndi_hp_state_change_req().....	14
3.3.3 bus_hp_op.....	15
3.4 Interface Table.....	16

1. Introduction

1.1 Background & Motivation

In Solaris today, every hotplug capable I/O bus technology (e.g., PCI, USB, PCMCIA, Storage) has its own hotplug framework. There is no generic hotplug framework in Solaris kernel. Therefore, each hotplug capable I/O bus needs to invent its own interfaces and mechanism.

In the userland, although there is a popular management software `cfgadm`, it also requires that each bus to have a bus specific plugin.

Therefore, either in Solaris kernel or in userland, there are rooms to improve the current software architecture by having a generic hotplug framework.

If look at the requirements, first, more and more I/O buses are supporting hotplug. That is, we need to develop new bus specific hotplug frameworks from time to time if without a generic Solaris hotplug framework there. And also, we need to maintain more and more bus specific hotplug frameworks as time going on.

Second, the hotplug requirement under virtualized environments (e.g., I/O Virtualization in Ldoms or Xen) raised a new issue here: arbitrary devices in the device tree should be able to be “hotplugged” in order to migrate between domains. This is a significant requirement and one of the major motivations of this project. See section 1.2 and section 2.5.5 for more details.

1.2 Goal

The main goal of this project is to deliver a generic common hotplug framework to Solaris users, a foundation which can support the hotplug functionality for any hotpluggable bus and also, to support device migration functionality in virtualized environment through “virtual hotplug” (see section 1.3 for this term). This new Solaris hotplug framework has a state machine based architecture. It includes generic hotplug states and provides various common hotplug interfaces, which includes:

- Management Interfaces

Generic or bus specific hotplug administration applications including GUIs can be written using these interfaces.

- Hotplug Event Interfaces

Hotplug-aware consumers such as drivers, kernel modules and hotplug aware applications can receive both synchronous and asynchronous hotplug events.

- Device Driver Interfaces

Device drivers can work with the hotplug framework to participate in hotplug activities such as surprise removal, hot replacement, dynamic resource balancing, error handling.

- Hotplug Controller Driver (HPC) Interfaces

HPC drivers will communicate with the hotplug framework through these interfaces. New HPC driver can be written easily for any hotpluggable bus.

This document proposes a new generic Solaris hotplug framework that lays the foundation to address current and future hotplug requirements.

1.3 Terminologies

Physical Hotplug	Hotplug operations on the hardware receptacle.
Virtual Hotplug	Hotplug operations at an arbitrary device node in the Solaris device tree.
Connector	Place where physical hotplug happens.
Port	Place where virtual hotplug happens.
Connection	Place where hotplug (either physical or virtual) happens.
Component	Physical object being hotplugged.
Device	Software object being hotplugged.

Hotplug-object Object being hotplugged (either physical or virtual).

1.4 Scope and Road Map

Currently, there are a lot of issues need to be addressed in hotplug domain. It should be a phase by phase effort to carry out all of them. This case is the first phase and it focuses on the following:

Phase-I:

- Basic hotplug framework features and interfaces
- Virtual hotplug support for PCI/PCIe devices
- Physical hotplug support for PCIe (Native, ACPI) and PCI (SHPC based) devices in new hotplug framework
- Hotplug userland stack (daemon, library and CLI)
- A new cfgadm plugin to support cfgadm(1M) to work with the new hotplug framework

This document focus on the overall architecture and also interfaces exported by kernel modules. Other documents in this case give the hotplug userland stack architecture details and interfaces.

Besides the framework implementation, this phase also ports PCI/PCIe hotplug controller drivers to the new hotplug framework. Any other driver porting jobs are out of scope of this case.

The following are what can be foreseen for the future phases but out of scope of this case:

Phase-II:

- Hotplug events through Device Contract and LDI frameworks
Device contract and LDI frameworks enhancements to support new hotplug events.
- New hotplug DDI interfaces for device drivers
This task develops DDI interfaces for interaction between hotplug framework and device drivers.
- New PCI resource allocator and a common PCI Configurator with support for dynamic resource re-balancing
This task is to minimum the chance that a PCI hotplug operation would fail due to lacking of resource (e.g., bus numbers, I/O or memory address spaces, etc.) at connector.
- Support hotplug FMA for PCIe devices
- RBAC support for hotplug userland administer tool
In Phase-I, “root” privilege is required to do any operations via hotplug userland administer tool. This task is to deploy RBAC (Role Based Access Control) feature to the tool.

Phase-III:

- Support for GUI administer tool
- Support for ExpressCard
ExpressCard is a industry standard replacing PCMCIA cards. This task is to develop a new hotplug controller driver based on the new hotplug framework interfaces.

In short, this case focuses on the infrastructures, and future phases will add more features based on what are being done here.

1.5 References

- [1] PCI Hot-Plug Specification, Revision 1.0, June 20, 2001
- [2] PCI Standard Hot-Plug Controller Specification, Revision 1.0, June 20, 2001
- [3] PCI Express Base Specification, Revision 2.0, Dec 20, 2006
- [4] Solaris Operating System - Hardware Virtualization Product Architecture, Revision 1.0, November 2007
- [5] Hotplug related ARC cases:

PSARC/1993/687	Hot Plugging
PSARC/1996/285	DR for CPU/Memory Boards
PSARC/1998/327	PCI Hotplug Support
PSARC/1998/460	RCM Framework
PSARC/1999/122	PCI Bus Resource Allocator
PSARC/1999/135	Solaris 8 RAS Enhancements
PSARC/1999/287	Solaris Embedded Fcode Interpreter
PSARC/1999/322	System Event Framework
PSARC/1999/686	Promote PCI Hotplug Interfaces
PSARC/2000/085	Extensions for PCI Hotplug Framework
PSARC/2000/218	RCM Framework Enhancements
PSARC/2000/295	PICL Event Interfaces
PSARC/2000/298	PCI Hotplug slot registration policy
PSARC/2000/532	State Model for Device Information Nodes
PSARC/2002/315	cPCI Autoconfiguration Support
PSARC/2003/193	Solaris Contracts and restart agreements
FWARC/2005/048	1275 Bindings for PCI Express Interconnect
PSARC/2005/375	PCI Hotplug Extensions for PCIe
PSARC/2006/037	PCI Express Hotplug Framework Interrupt Interfaces
FWARC/2006/198	PCI Hotplug Resource Preallocation
PSARC/2007/197	ZFS Hotplug
PSARC/2007/290	Retire Agent for IO devices

2. Architecture Overview

2.1 Physical Hotplug & Virtual Hotplug

2.1.1 Concepts

When a hardware device is plugged to system during running time, in order to configure the newly added device, operating system needs to probe the device, allocate dev_info node and resource (e.g. IO/MEM address spaces) for it, and then tries to attach driver. It is a reversed case when a hardware device is unplugged from system. This is the physical hotplug expected to be supported by an operating system.

If not limit hotplug operations to physical hotpluggable devices, that is, apply the similar operations for an arbitrary dev_info node in the device tree, then this is referred as “virtual hotplug”. For example, hot add/remove an on-board PCI NIC device. Although the device is not physically plug/unplugged to/from the system, from the software view, it is coming to the device tree when virtual plugged and going when virtual unplugged.

2.1.2 Background & Goal

Currently, the hotplug operations supported on Solaris are limited to those devices associated with a hotpluggable physical receptacle (mentioned as “connector” in this document). This is almost fine in the single domain environment in the past, but, nowadays, more and more customers are adopting virtualized environment. Therefore, a requirement is raised that arbitrary device nodes in the device tree should be able to be “hotplugged” and then migrated between domains.

This case tries to provides a uniformed and consistent hotplug model that is capable of adding and removing arbitrary

devices with or without hotpluggable physical receptacles being present, that is, support both physical hotplug and virtual hotplug within one framework.

2.2 Overall Architecture

This section gives an overview of the new Solaris Hotplug Framework. Figure 1 shows the block diagram of software components involved.

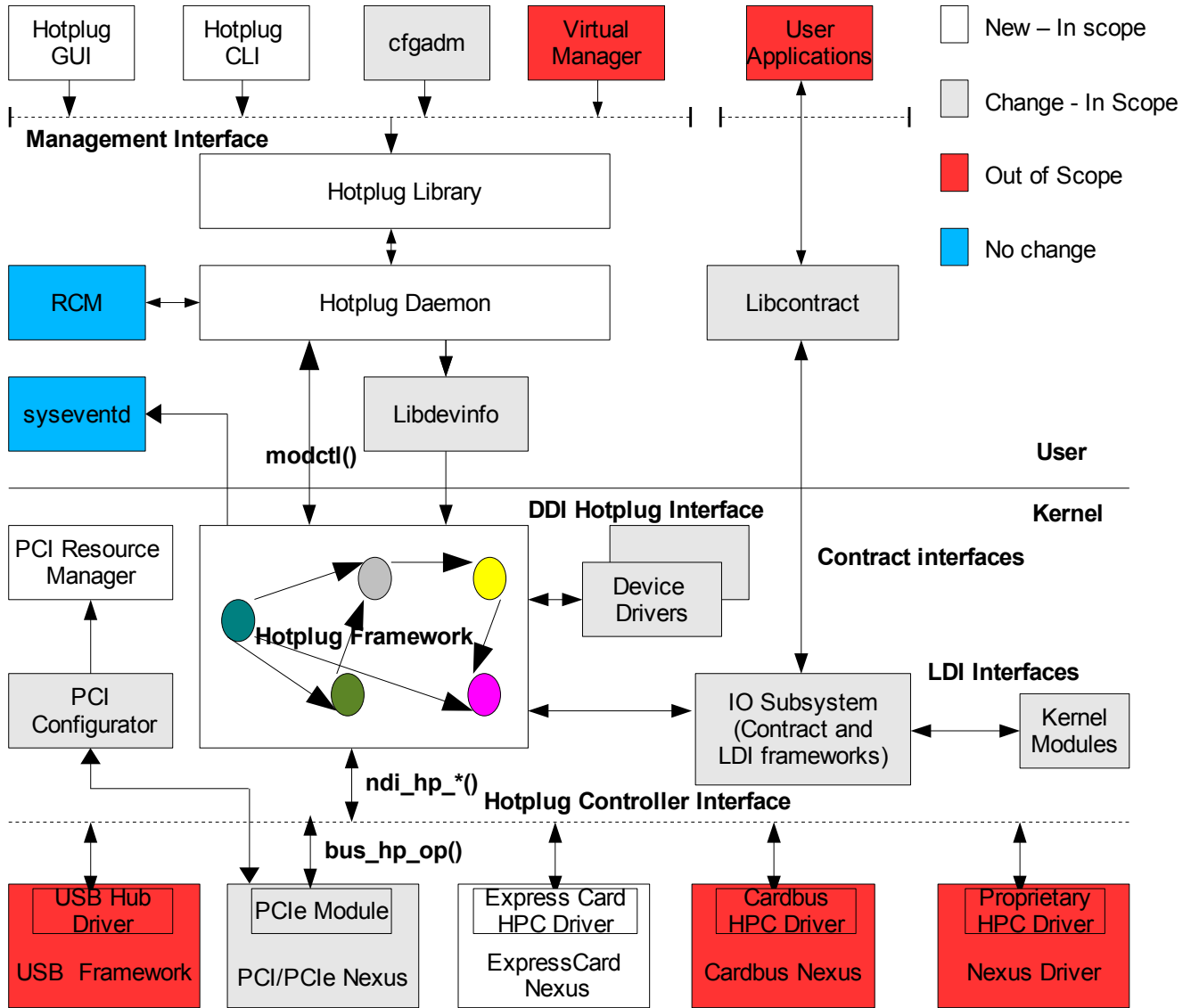


Figure 1. Overall Architecture of Solaris Hotplug Framework

The new Solaris hotplug framework is based on a state machine which is supposed to include generic hotplug states and transitions between states. Refer to section 2.5 for details.

It provides a set of generic and common hotplug interfaces, which includes:

- Management Interfaces

Generic or bus specific hotplug admin applications including GUIs, can be developed using the newly defined management interfaces provided by the hotplug library. Refer to userland design documents of this case for details.

- Hotplug Controller Driver Interfaces

Hotplug controller drivers can be developed for any hotpluggable bus using the newly defined bus_op and NDI interfaces (bus_hp_op() and ndi_hp_*() in Figure 1). Refer to section 3.3 for details.

For the kernel – userland interfaces, a new command is added to the existing modctl system call and some new interfaces are added to libdevinfo. See section 3.1 and 3.2 for details.

From the long term view, cfgadm users are supposed to switch to the hotplug CLI provided by this case. Before that, to make cfgadm running on this framework, so that hotplug controller drivers don't need to maintain cfgadm ioctls any more, a new cfgadm plugin will be developed as the translator from cfgadm commands to Solaris hotplug framework commands. Refer to section 2.8 for details.

And there are plans for some other interfaces to be developed in future phases, such as new hotplug DDI interfaces for device drivers, hotplug events through Device Contract and LDI frameworks, etc. See section 1.4 for details.

2.3 Connector & Port

In this document, either “connector” or “port” is a software object standing for the place where hotplug operations happen.

A connector is corresponding to a hardware receptacle, e.g., a PCI slot, USB port, cardbus slot, SAS port, and fiber channel port, etc. The occupant of the connector is mentioned as “component” in this document. It is corresponding to a PCI card, USB device, cardbus card, etc.

A port is used to specify a “place” under a nexus dev_info node in the device tree in a running Solaris image. For example, “/pci@0,0 pci.1,2” specifies a place under a PCI nexus dev_info node /pci@0,0, and at that place, a dev_info node with attributes “PCI device number 1 and function number 2” is expected. For details of the naming of a connector or port, refer to section 2.4. The occupant of the port is mentioned as “device” in this document. It is corresponding to a dev_info node in device tree. Note, it is forced that the relationship between port and device is 1 : 1. By default, each initialized dev_info node will have a port if the node's bus software (e.g., if PCI dev_info node, its bus software is PCI nexus driver) supports virtual hotplug.

The same data structure is defined for both connector and port. And, the hotplug states of connector or port are unified in one state machine. Thus, Solaris hotplug framework provides consistent interfaces for the hotplug operations happened at either connector or port.

2.4 Naming of Connection

A “connection” is either a connector or a port. Any connections in a system are specified in this way:

“Absolute path to parent dev_info node” + “connection name”

A *connection name* is a string which is unique among all of the connections which are immediately under the same parent dev_info node. So that the above “path + name” combination is system wide unique.

A connection's name is decided by bus software (hotplug controller driver). For connector, it is usually derived from hardware properties, e.g., the hardware slot name set by platform manufacturer. For port, it is usually set according to bus specific industry standard.

An example is:

```
/pci@0,0/pci10de,5d@e NEM0
```

“NEM0” is the PCI slot name get from “slot-names” property which is from platform manufacturer.

For virtual hotplug case, the port name is set by hotplug controller drivers (instead of Solaris hotplug framework). Because only the bus specific software has the knowledge of naming a receptacle under the bus node.

An example is:

```
/pci@0,0/pci10de,5d@e pci.1,2
```

“pci.1,2” is the port name which stands for PCI device number 1 and function number 2. It specifies a device place under a PCI bridge node according to PCI industry standard.

Solaris hotplug framework does not define any specific naming rules for connection names, but it assumes that a hotplug controller driver instance will give names which are unique among the immediate children of the dev_info node that the driver instance is being attached.

2.5 Hotplug States & State Machine

2.5.1 Overview

Solaris hotplug framework defines a set of hotplug states and a hotplug state machine, see Figure 2. Here, “hotplug states” are the states of a connection or a hotplug-object. They are exposed to both the hotplug controller drivers and hotplug userland stack. Thus, the hotplug userland stack, framework and controller drivers can then work on the same states and state machine.

The hotplug states are supposed to be generic, that is, it is not aimed to enumerate every possible states in various hotpluggable buses. Usually, the bus specific states are much more detail than the hotplug states defined by the framework. Hotplug controller drivers can run on their own states but they should aware the framework states and have mappings between them when interacting with the framework.

2.5.2 Connector States

To support both physical hotplug and virtual hotplug in one state machine, states for connector, port and device are defined as three sub-sets of states. And we can see the obvious dependencies between them in the state machine.

The first sub-set of states are for connectors:

- **Empty:** No component plugged to the connector.
- **Present:** Component already plugged to the connector.
- **Powered:** The connector is powered.
- **Enabled:** The connector is fully functional.

As mentioned in section 2.3, a connector specifies a physical hotpluggable place. For example, a PCI hotpluggable slot. The following is the example steps of the connector state transition:

1. User plugs a PCI board to an empty slot when system is running.
2. The PCI hotplug controller driver detects there is something plugged in the slot, then, the corresponding connector’s state is turned to “Present” from “Empty”.
3. User issues a command to power the slot. The PCI hotplug controller driver gets the command from hotplug framework and then programs the slot hardware registers to add power to the slot. When this is done, the connector’s state is turned to “Powered”.
4. There might be some more detail programming jobs need to be done for the connector before the plugged PCI board can be accessed by software. User issues another command to “enable” the connector. The connector state is then turned to “Enabled”. Now, the newly plugged PCI board is ready to be probed, that is, OS software can then read the configure space of it and knows what kind of device it is.

Note, because the connectors are managed by nexus drivers (PCI hotplug controller driver in the above example), the power management related issues are handled by the nexus drivers and system PM modules. Hotplug controller drivers should request hotplug framework to update the connector’s state whenever it is changed for some reasons. On the other way, hotplug framework passes the state change commands (e.g., power the connector) to the hotplug controller driver and the latter take care of the rest jobs including PM operations.

The above states and state transitions are what must be experienced before a hotplugged hardware device could be probed by Solaris software. This process is the “physical hotplug” mentioned in section 2.1.

2.5.3 Port and Device States

After the physical device is probed, some dev_info nodes will be allocated for the newly plugged device. The process of linking the dev_info nodes to the right place in the device tree is taken as “virtual hotplug” mentioned in section 2.1.

Let’s continue with the above example. Suppose the newly plugged PCI board has four functions, we know that the existing Solaris PCI configuration software will probe the device’s functions one by one by reading their configuration registers, and then, it allocates a dev_info node for each function and link the four new nodes as children nodes to the PCI nexus node who manages the PCI slot.

According to PCI industry standard, one PCI device (either hotplugged device or on-board device) takes one device number and one or more PCI function numbers. By specifying the PCI nexus node and a device/function number pair, we can

specify a unique “place” in a PCI device tree. Therefore, we know that we can specify a “port” (refer to section 2.3) for a PCI dev_info node by specifying its device number, function number and also the absolute path to its parent node.

It is straightforward that the port has the following two states:

- **Port-Empty:** the port has no device occupying it.
- **Port-Present:** the port is occupied by a device.

One port is for one dev_info node only. In the above example, four ports will be created and the four new nodes will be virtually “plugged” to their ports one by one.

Note, if the bus software supports virtual hotplug functionality, then, for all the nexus nodes of this bus, ports will be created for each immediate child of the nexus node when the child node is initialized, no matter the device is a physical hotplugged device or an on-board device. Thus, arbitrary dev_info nodes can be virtually hotplugged.

When the port state goes to “Port-Present”, user can then manipulate the port's device state. It can be one of the following state:

- **Initialize:** The dev_info node is created and not get to “Probed” state yet.
- **Probed:** The driver's probe(9e) routine has been invoked and it returned success.
- **Attached:** The driver's attach(9e) routine has been invoked and it returned success.
- **Operational:** Post attach processing, including DACF, has been completed, and the device is fully functional.
- **Maintenance:** The device is partly functional or not functional at all. This is either because of a fault or a result of an intentionally administration action. It could has some sub-states, but the definition of the sub-states are out of scope of this case (planned in future phases).

There is a straightforward mapping between the above device states and the DS_* states proposed by PSARC/2000/532.

Table 1. State Mapping

devi_node_state	Port's Device State
DS_PROTO	Initialize
DS_LINKED	Initialize
DS_BOUNDED	Initialize
DS_INITIALIZED	Initialize
DS_PROBED	Probed
DS_ATTACHED	Attached
DS_READY	Operational / Maintenance

Because the relationship between port and device is 1:1 and the device exists after the port is at “Port-Present” state, therefore, to make it easier to tell, this document also mentions the port's device state as port state.

Continue with the above example, corresponding to the four functions in the PCI board, four new dev_info nodes are now in four ports. User can manipulate each of the port's state separately. Note, it is out of the scope of this case to support specific commands to bring a port to “Maintenance” state, but it is planned in future phases.

2.5.4 State Machine

The following is the state machine.

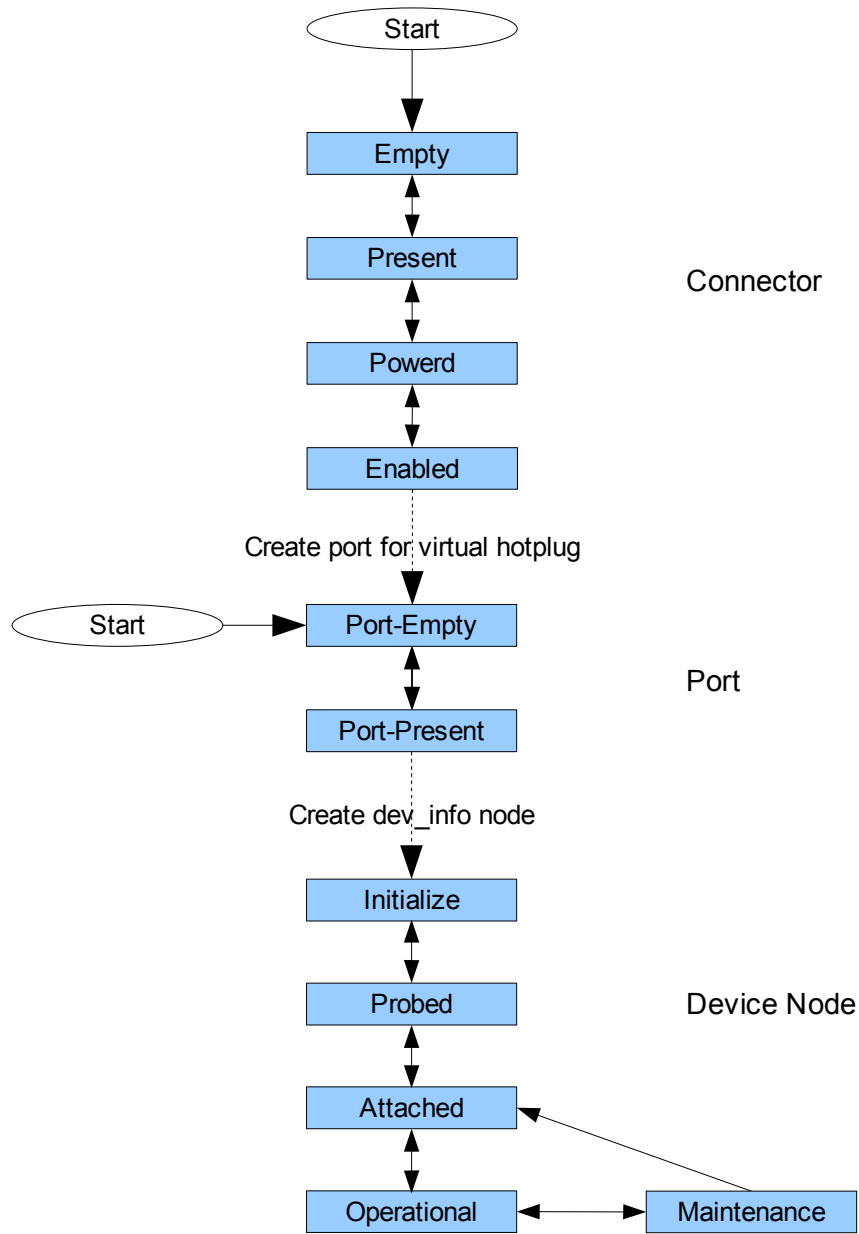


Figure 2. Hotplug State Machine

There are two start points in the state machine. The top one is for the case that a connector is involved. The other one is for the case that no connectors involved, i.e., the on-board device case.

When a connector is created, it starts from “Empty” state. For example, the PCI hotplug slot case or USB port case. The connector is created usually when hotplug controller driver discovers a physical hotpluggable receptacle in the system.

When a port is created, it starts from “Port-Empty” state. It is created for receiving a dev_info node. The port can either depend or not depend on a connector. That is, either a physical hotpluggable device or an on-board device can have a port. For a physical hotpluggable device, only if its connector is at “Enabled” state, the corresponding port and device node would exist.

The port is created mainly in two scenarios:

- When a dev_info node is being initialized by bus software (nexus driver), for example, during device enumeration

when boot or at the end of a physical hotplug process. In this case, the port state will be updated to “Port-Present” from “Port-Empty” immediately after the port is created.

- When a userland command is received to create a specific port. The port state will be kept as “Port-Empty” until a virtual hotplug operation happens.

When it is known that a port's device exist, the port's state goes to “Port-Present”. After the dev_info node of the device is allocated, the port state goes through “Initialize”, “Probed”, “Attached” and “Operational”. Now the device is fully functional. And then, if something happened (except driver detach or predetach) which made the device not fully functional, then it goes to “Maintenance” state.

2.5.5 A Usage Example in Virtualized Environment

As mentioned in section 2.1.2, one goal of this case is to support device migration in virtualized environment. The following is a usage example to elaborate how the virtual hotplug functionality can be used to do a device migration between domains.

Suppose there is a PCI NIC device in a system, and it has two functions. Each function has a corresponding network node in system device tree. The system administrator want to move the second function from root domain to an I/O domain so that the I/O domain can have direct access to that NIC function. The following are the example steps to perform this action:

1. Login the root domain. Administrator issues a command to list the ports in the system and find the port name corresponding to the second function of the PCI NIC device.
2. By specifying the port name and the absolute path to the parent node, administrator issues a command to turn the port state to “Initialize” in root domain. The command should include a option to indicate that a flag should be marked to the corresponding dev_info node that the device's driver can not be automatically re-attached.
3. Login the target I/O domain. By specifying the same port name and path, administrator issues a command to create the port in I/O domain. The port is in “Port-Empty” state.
4. Still in the I/O domain, administrator issues another command to the port to turn the port state to “Port-Present”. This will cause the hotplug framework and the bus software in I/O domain to do a read-only probe. The NIC function will then be discovered.
5. In the I/O domain, administrator issues the last command to the port to bring the port's device state to “Operational”. Done.

Note, the above example steps are for the purpose of elaborating the details of what will happen during a device migration. The actual steps can be much simpler for a system administrator. For example, the hotplug administer tool can combine step 3, 4 and 5 into one command which includes creating port and upgrading port state to “Operational” from “Port-Empty”.

2.6 Bus Specific Commands Support

Besides the generic hotplug state machine and commands, Solaris hotplug framework also provides a way to support bus specific commands related with hotplug operations, e.g., turn on/off the LED of a PCI hotplug slot. The general idea is, the userland stack and Solaris hotplug framework just pack and pass the command from user to hotplug controller drivers, and then hotplug controller drivers do all the interpretation job and execute the command. The command results (if any) will also be just passed up without any interpretation. Thus, to add a new bus's hotplug support, the only things to do are to update the userland hotplug CLI man page for new bus specific commands and add a new hotplug controller driver. There is no need to update the code of userland stack and Solaris hotplug framework.

The initial idea is to use name-value pairs to pack user commands and return command results. More details TBD.

2.7 Events

From the view of the consumers of hotplug sysevents, nothing is changed by this case. Specifically, the following two sysevents are mentioned here:

- ESC_DR_AP_STATE_CHANGE
- ESC_DR_REQ

The hotplug controller drivers continue to generate ESC_DR_REQ. For example, when a PCI attention button is pressed, the PCI hotplug controller driver gets the interrupt and then generate an ESC_DR_REQ event as before, and the existing userland consumers will also just work as before.

For ESC_DR_AP_STATE_CHANGE, there are two cases:

- Userland issues a state change request to hotplug framework, then the framework will generate this sysevent when the state is changed. For example, a user issues a command to enable a connector (that is, change the connector state to “Enabled”) via hotplug administration tool.
- The hotplug controller driver submit a state change request to hotplug framework by calling `ndi_hp_state_change_req()` (see section 3.3.2 for this interface), then the framework will generate this sysevent after it completes the state change job. For example, when a PCI power fault interrupt is received by the PCI hotplug controller driver's interrupt handler, the handler will submit a state change request to the framework that the PCI slot's state should be changed to a state less than “Powered” (that is, power off the slot); then, the framework will generate an ESC_DR_AP_STATE_CHANGE sysevent when the state change job is done.

Although the hotplug controller drivers are not generating ESC_DR_AP_STATE_CHANGE directly as before, the sysevent consumers does not see any changes because they just get the same events as before.

In the future phases, it is possible that we would replace the existing events and also modify the events consumers in order to support new functionality.

2.8 cfgadm shp Plugin

As mentioned in Figure 1, a new `cfgadm shp` plugin will be developed as the translator between `cfgadm` commands and Solaris hotplug framework commands. This phase will port PCI/PCIe hotplug controller drivers to Solaris hotplug framework interfaces. The current `cfgadm ioctl` implementation in these drivers will be removed. But all the existing `cfgadm PCI` functionality will be supported via the new `cfgadm shp` plugin. In the future (not in this case scope), when some other hotplug controller drivers are also ported to Solaris hotplug framework, the `cfgadm shp` plugin can also be used to supported those drivers.

More details TBD.

3. Interfaces

This section focus on interfaces exported by Solaris hotplug framework in the kernel. For interfaces in userland stack, refer to userland documents of this case.

3.1 modctl Commands

`modctl` is a system call in Solaris. This document adds a new command to it. Userland programs call `modctl()` to issue the new command for hotplug operations.

- New `modctl` command and sub commands :

```
#define MODHPOPS 45
```

The above is the new `modctl` command to be added to `uts/common/sys/modctl.h` for hotplug operations. The following are three sub commands of it:

```
#define MODHPOPS_CHANGE_STATE 0
```

```
#define MODHPOPS_CREATE_PORT 1
```

```
#define MODHPOPS_REMOVE_PORT 2
```

```
#define MODHPOPS_BUS_CMD 3
```

- `modctl` arguments corresponding to different sub commands :

```
int modctl(MODHPOPS, MODHPOPS_CHANGE_STATE, char *path, char *cn_name, int state);
```

```
int modctl(MODHPOPS, MODHPOPS_CREATE_PORT, char *path, char *port_name);
```

```
int modctl(MODHPOPS, MODHPOPS_REMOVE_PORT, char *path, char *port_name);
```

```
int modctl(MODHPOPS, MODHPOPS_BUS_CMD, char *path, char *cn_name, TBD);
```

Solaris hotplug framework implements the above interfaces.

When `MODHPOPS_CHANGE_STATE` command is received, the framework tries to operate the connection specified by `path` and `cn_name`, and change connection state to the target `state`. The `modctl()` returns success (0) if the target state is

reached. Otherwise, it returns an errors code.

When MODHPOPS_CREATE_PORT command is received, the framework creates a new port under dev_node specified by path with name port_name. The modctl() returns success (0) if the port is successfully created. Otherwise, it returns an errors code.

When MODHPOPS_REMOVE_PORT command is received, the framework removes the port specified by path and port_name. The modctl() returns success (0) if the port is successfully removed. Otherwise, it returns an errors code.

When MODHPOPS_BUS_CMD command is received, the framework just pass the bus specific command to the hotplug controller driver. The latter parses the command and operate the connector specified by path and cn_name. And then return the bus specific results to hotplug framework. The framework then pass the result on and eventually the modctl() caller will get the result. The modctl() returns success (0) if the bus specific command is successfully passed to the driver and also get the result data from the driver. Otherwise, it returns an errors code.

3.2 Libdevinfo interfaces

To get information of the connections in a system (including connection name, current state, etc.), new libdevinfo interfaces are added. More details see userland documents in this case.

3.3 Hotplug Controller Driver Interfaces

Solaris hotplug framework defines some interfaces for hotplug controller drivers. Via these interfaces, the framework and drivers are able to interact with each other.

3.3.1 ndi_hp_register() & ndi_hp_unregister()

Hotplug controller drivers call these interfaces which are implemented by Solaris hotplug framework. Also refer to Figure 1.

ndi_hp_register()

Synopsis

```
int ndi_hp_register(dev_info_t *dip, ddi_hp_cn_info_t *cn_info);
```

Arguments

dip parent node of the connection.
cn_info pointer to the structure of connection information.

Return Values

NDI_SUCCESS operation succeed.
NDI_FAILURE operation failed.
And other NDI_* values to specify an error type on operation failure case.

Context

This function can be called from user and kernel context.

Hotplug controller drivers call this interface to register a new connection to the framework, that is, to create the connection data structure on a nexus node.

The following structure is input as cn_info argument. It specifies the connection to be registered.

```
typedef struct ddi_hp_cn_info {  
    uint16_t        cn_ver;           /* DDI_HP_CN_INFO_VERSION */  
    char            cn_name[DDI_HP_CN_NAME_LEN];   /* Connection name */  
    int             cn_connector_num; /* Connector number */  
    int             cn_port_num;    /* Port number */  
    char            *cn_type_str;    /* The string for the type of the connection. */  
    dev_info_t      *cn_child;       /* The device occupying this port; It is kept NULL for a connector. */  
};
```

```
} ddi_hp_cn_info_t;
```

cn_name is a string name of a connection, refer to section 2.4 for details.

cn_connector_num is a digital number of a connector. It is unique among the sibling connectors. It is given by bus software (hotplug controller driver), usually derived from hardware properties, e.g., the hardware slot number set by platform manufacturer. If a *ddi_hp_cn_info_t* structure is for a connector, its *cn_connector_num* is set to the connector's number. If a *ddi_hp_cn_info_t* structure is for a port, its *cn_connector_num* is set to the port's dependent connector's number; if the port has no dependent connector, then *cn_connector_num* is set to -1.

cn_port_num is a digital number of a port. It is unique among the sibling ports. It is given by bus software (hotplug controller driver) according to bus industry standards. For example, for a port under a PCI nexus node, the port number is encoded from PCI device number and function number. If a *ddi_hp_cn_info_t* structure is for a port, its *cn_port_num* is set to the port's number. If a *ddi_hp_cn_info_t* structure is for a connector, its *cn_port_num* is set to -1.

cn_type_str is for the connection types, e.g., "PCI Slot", "PCI Express Slot", "Port", etc. It is supposed that each hotplug controller driver will register their own types to the framework. These type string is exposed to userland and end users.

cn_child points to the child of a port. It is NULL if the port is at state "Port-Empty" or the *ddi_hp_cn_info_t* structure is for a connector.

ndi_hp_unregister()

Synopsis

```
int ndi_hp_unregister(dev_info_t *dip, char *cn_name);
```

Arguments

dip parent node of the connection.
cn_name name of the connection.

Return Values

NDI_SUCCESS operation succeed.
NDI_FAILURE operation failed.
And other NDI_* values to specify an error type on operation failure case.

Context

This function can be called from user and kernel context.

Hotplug controller drivers call this interface to un-register an existing connection from the framework, that is, to remove the connection data structure from a nexus node.

3.3.2 ndi_hp_state_change_req()

Hotplug controller drivers call this interface to submit a state change request to hotplug framework.

Synopsis

```
int ndi_hp_state_change_req(dev_info_t *dip, char *cn_name, ddi_hp_cn_state_t state, int flag);
```

Arguments

dip parent node of the connection.
cn_name name of the connection.
State target state
flag DDI_HP_REQ_SYNC, return until the state change is done.
 DDI_HP_REQ_ASYNC, return once the state change request is submitted.

Return Values

NDI_SUCCESS operation succeed.
NDI_FAILURE operation failed.

And other `NDI_*` values to specify an error type on operation failure case.

Context

This function can be called from user and kernel context. If with flag `DDI_HP_REQ_ASYNC`, it can also be called from interrupt context.

See section 2.7 for more about the usage of this interface.

3.3.3 bus_hp_op

`bus_hp_op` is a new `bus_op` added to `struct bus_ops` by this case. Solaris hotplug framework calls this interface which is implemented by hotplug controller drivers. Also refer to Figure 1.

```
struct bus_ops {
    ...
+   /*
+   * NOTE: the following busop entrypoint is available with version
+   * 10 or greater.
+   */
+   int      (*bus_hp_op)(dev_info_t *dip, char *cn_name,
+                        ddi_hp_op_t op, void *arg, void *result);
};
```

Arguments

<code>dip</code>	parent node of the connection.
<code>cn_name</code>	name of the connection.
<code>op</code>	operations to be done for the connection.
<code>arg</code>	additional data to specify the operation.
<code>result</code>	operation result.

Return Values

`DDI_SUCCESS` operation succeed.

`DDI_FAILURE` operation failed.

And other `DDI_*` values to specify an error type on operation failure case.

The operations specified by `op` are defined as following:

```
/*
* Typedef for Hotplug OPS commands used with bus_hp_op()
*/
typedef enum {
    DDI_HPOP_CHECK_HP_SUPPORT = 1, /* 01 Check hotplug support */
    DDI_HPOP_CN_GET_STATE,        /* 02 Get the current connection state */
    DDI_HPOP_CN_CHANGE_STATE,    /* 03 Change connection state */
    DDI_HPOP_CN_PROBE,           /* 04 Probe connection */
    DDI_HPOP_CN_UNPROBE,         /* 05 Unprobe connection */
    DDI_HPOP_CN_ONLINE_CHILDREN, /* 06 Online all the children of connection */
};
```

```

        DDI_HPOP_CN_OFFLINE_CHILDREN, /* 07 Offline all the children of connection */
    } ddi_hp_op_t;

```

3.4 Interface Table

The following is a summarize of all the interfaces mentioned in section 3.

Table 2. Interfaces

Interface	Stability	Comments
MODHPOPS	Project private	Modctl command for hotplug operations
MODHPOPS_CHANGE_STATE	Project private	Modctl sub-command for changing hotplug state
MODHPOPS_CREATE_PORT	Project private	Modctl sub-command for creating a port
MODHPOPS_REMOVE_PORT	Project private	Modctl sub-command for removing a port
ndi_hp_register()	Consolidation private	Register a connection to hotplug framework
ndi_hp_unregister()	Consolidation private	Unregister a connection from hotplug framework
ndi_hp_state_change_req()	Consolidation private	Submit a state change request to hotplug framework
ddi_hp_cn_info_t	Consolidation private	Structure for the information of a connection
bus_hp_op()	Consolidation private	bus_ops entry for hotplug operations
ddi_hp_op_t	Consolidation private	Hotplug operations