```
   1              ------------------------------
   2              | PCI Pass-through in Xen ARM |
   3              ------------------------------
   4              manish.jaggi@caviumnetworks.com
   5              ------------------------------
   6
   7                          Draft-3
   8
   9
  10   -------------------------------------------------------------------------------
  11   Introduction
  12   -------------------------------------------------------------------------------
  13   This document describes the design for the PCI passthrough support in Xen ARM.
  14   The target system is an ARM 64bit Soc with GICv3 and SMMU v2 and PCIe devices.
  15
  16   -------------------------------------------------------------------------------
  17   Revision History
  18   -------------------------------------------------------------------------------
  19   Changes from Draft-1:
  20   --------------------
  21   a) map_mmio hypercall removed from earlier draft
  22   b) device bar mapping into guest not 1:1
  23   c) holes in guest address space 32bit / 64bit for MMIO virtual BARs
  24   d) xenstore device's BAR info addition.
  25
  26   Changes from Draft-2:
  27   --------------------
  28   a) DomU boot information updated with boot-time device assignment and hotplug.
  29   b) SMMU description added
  30   c) Mapping between streamID - bdf - deviceID.
  31   d) assign_device hypercall to include virtual(guest) sbdf.
  32   Toolstack to generate guest sbdf rather than pciback.
  33
  34   -------------------------------------------------------------------------------
  35   Index
  36   -------------------------------------------------------------------------------
  37   (1) Background
  38
  39   (2) Basic PCI Support in Xen ARM
  40   (2.1) pci_hostbridge and pci_hostbridge_ops
  41   (2.2) PHYSDEVOP_HOSTBRIDGE_ADD hypercall
  42
  43   (3) SMMU programming
  44   (3.1) Additions for PCI Passthrough
  45   (3.2) Mapping between streamID - deviceID - pci sbdf
  46
  47   (4) Assignment of PCI device
  48
  49   (4.1) Dom0
  50   (4.1.1) Stage 2 Mapping of GITS_ITRANSLATER space (4k)
  51   (4.1.1.1) For Dom0
  52   (4.1.1.2) For DomU
  53   (4.1.1.2.1) Hypercall Details: XEN_DOMCTL_get_itranslater_space
  54
  55   (4.2) DomU
  56   (4.2.1) Reserved Areas in guest memory space
  57   (4.2.2) New entries in xenstore for device BARs
  58   (4.2.4) Hypercall Modification for bdf mapping notification to xen
  59
  60   (5) DomU FrontEnd Bus Changes
  61   (5.1) Change in Linux PCI FrontEnd - backend driver for MSI/X programming
  62   (5.2) Frontend bus and interrupt parent vITS
  63
  64   (6) NUMA and PCI passthrough
  65   -------------------------------------------------------------------------------
  66
  67   1. Background of PCI passthrough
  68   -----------------------------------
  69   Passthrough refers to assigning a pci device to a guest domain (domU) such that
  70   the guest has full control over the device. The MMIO space and interrupts are
  71   managed by the guest itself, close to how a bare kernel manages a device.
  72
```

```
73   Device's access to guest address space needs to be isolated and protected. SMMU
74   (System MMU - IOMMU in ARM) is programmed by xen hypervisor to allow device
75   access guest memory for data transfer and sending MSI/X interrupts. PCI devices
76   generated message signalled interrupt write are within guest address spaces which
77   are also translated using SMMU.
78   For this reason the GITS (ITS address space) Interrupt Translation Register
79   space is mapped in the guest address space.
80
81   2. Basic PCI Support for ARM
82   --------------------------------
83   The apis to read write from pci configuration space are based on segment:bdf.
84   How the sbdf is mapped to a physical address is under the realm of the pci
85   host controller.
86
87   ARM PCI support in Xen, introduces pci host controller similar to what exists
88   in Linux. Each drivers registers callbacks, which are invoked on matching the
89   compatible property in pci device tree node.
90
91   2.1 pci_hostbridge and pci_hostbridge_ops
92   ---------------------------------------------
93   The init function in the pci host driver calls to register hostbridge callbacks:
94   int pci_hostbridge_register(pci_hostbridge_t *pcihb);
95
96   struct pci_hostbridge_ops {
97   u32 (*pci_conf_read)(struct pci_hostbridge*, u32 bus, u32 devfn,
98   u32 reg, u32 bytes);
99   void (*pci_conf_write)(struct pci_hostbridge*, u32 bus, u32 devfn,
100  u32 reg, u32 bytes, u32 val);
101  };
102
103  struct pci_hostbridge{
104  u32 segno;
105  paddr_t cfg_base;
106  paddr_t cfg_size;
107  struct dt_device_node *dt_node;
108  struct pci_hostbridge_ops ops;
109  struct list_head list;
110  };
111
112  A pci conf read function would internally be as follows:
113  u32 pcihb_conf_read(u32 seg, u32 bus, u32 devfn,u32 reg, u32 bytes)
114  {
115  pci_hostbridge_t *pcihb;
116  list_for_each_entry(pcihb, &pci_hostbridge_list, list)
117  {
118  if(pcihb->segno == seg)
119  return pcihb->ops.pci_conf_read(pcihb, bus, devfn, reg, bytes);
120  }
121  return -1;
122  }
123
124  2.2 PHYSDEVOP_pci_host_bridge_add hypercall
125  ---------------------------------------------
126  Xen code accesses PCI configuration space based on the sbdf received from the
127  guest. The order in which the pci device tree node appear may not be the same
128  order of device enumeration in dom0. Thus there needs to be a mechanism to bind
129  the segment number assigned by dom0 to the pci host controller. The hypercall
130  is introduced:
131
132  #define PHYSDEVOP_pci_host_bridge_add 44
133  struct physdev_pci_host_bridge_add {
134  /* IN */
135  uint16_t seg;
136  uint64_t cfg_base;
137  uint64_t cfg_size;
138  };
139
140  This hypercall is invoked before dom0 invokes the PHYSDEVOP_pci_device_add
141  hypercall. The handler code invokes to update segment number in pci_hostbridge:
142
143  int pci_hostbridge_setup(uint32_t segno, uint64_t cfg_base, uint64_t cfg_size);
144
```

```
145   Subsequent calls to pci_conf_read/write are completed by the pci_hostbridge_ops
146   of the respective pci_hostbridge.
147
148   2.3 Helper Functions
149   -----------------------
150   a) pci_hostbridge_dt_node(pdev->seg);
151   Returns the device tree node pointer of the pci node from which the pdev got
152   enumerated.
153
154   3. SMMU programming
155   -------------------
156
157   3.1. Additions for PCI Passthrough
158   ----------------------------------
159   3.1.1 - add_device in iommu_ops is implemented.
160
161   This is called when PHYSDEVOP_pci_add_device is called from dom0.
162
163   .add_device = arm_smmu_add_dom0_dev,
164   static int arm_smmu_add_dom0_dev(u8 devfn, struct device *dev)
165   {
166   if (dev_is_pci(dev)) {
167   struct pci_dev *pdev = to_pci_dev(dev);
168   return arm_smmu_assign_dev(pdev->domain, devfn, dev);
169   }
170   return -1;
171   }
172
173   3.1.2 dev_get_dev_node is modified for pci devices.
174   ----------------------------------------------------------------------------
175   The function is modified to return the dt_node of the pci hostbridge from
176   the device tree. This is required as non-dt devices need a way to find on
177   which smmu they are attached.
178
179   static struct arm_smmu_device *find_smmu_for_device(struct device *dev)
180   {
181   struct device_node *dev_node = dev_get_dev_node(dev);
182   ....
183
184   static struct device_node *dev_get_dev_node(struct device *dev)
185   {
186   if (dev_is_pci(dev)) {
187   struct pci_dev *pdev = to_pci_dev(dev);
188   return pci_hostbridge_dt_node(pdev->seg);
189   }
190   ...
191
192
193   3.2. Mapping between streamID - deviceID - pci sbdf - requesterID
194   ----------------------------------------------------------------------
195   For a simpler case all should be equal to BDF. But there are some devices that
196   use the wrong requester ID for DMA transactions. Linux kernel has pci quirks
197   for these. How the same be implemented in Xen or a diffrent approach has to be
198   taken is TODO here.
199   Till that time, for basic implementation it is assumed that all are equal to BDF.
200
201
202   4. Assignment of PCI device
203   -------------------------------
204
205   4.1 Dom0
206   ------------
207   All PCI devices are assigned to dom0 unless hidden by pci-hide bootargs in dom0.
208   Dom0 enumerates the PCI devices. For each device the MMIO space has to be mapped
209   in the Stage2 translation for dom0. For dom0 xen maps the ranges from dt pci
210   nodes in stage 2 translation during boot.
211
212   4.1.1 Stage 2 Mapping of GITS_ITRANSLATER space (64k)
213   -----------------------------------------------------
214
215   GITS_ITRANSLATER space (64k) must be programmed in Stage2 translation so that SMMU
216   can translate MSI(x) from the device using the page table of the domain.
```

```
217
218    4.1.1.1 For Dom0
219    ----------------
220    GITS_ITRANSLATER address space is mapped 1:1 during dom0 boot. For dom0 this
221    mapping is done in the vgic driver. For domU the mapping is done by toolstack.
222
223    4.1.1.2 For DomU
224    ----------------
225    For domU, while creating the domain, the toolstack reads the IPA from the
226    macro GITS_ITRANSLATER_SPACE from xen/include/public/arch-arm.h. The PA is
227    read from a new hypercall which returns the PA of the GITS_ITRANSLATER_SPACE.
228    Subsequently the toolstack sends a hypercall to create a stage 2 mapping.
229
230    Hypercall Details: XEN_DOMCTL_get_itranslater_space
231
232    /* XEN_DOMCTL_get_itranslater_space */
233    struct xen_domctl_get_itranslater_space {
234    /* OUT variables. */
235    uint64_aligned_t start_addr;
236    uint64_aligned_t size;
237    };
238    typedef struct xen_domctl_get_itranslater_space xen_domctl_get_itranslater_space;
239    DEFINE_XEN_GUEST_HANDLE(xen_domctl_get_itranslater_space;
240
241    4.2 DomU
242    ------------
243    There are two ways a device is assigned
244    In the flow of pci-attach device, the toolstack will read the pci configuration
245    space BAR registers. The toolstack has the guest memory map and the information
246    of the MMIO holes.
247
248    When the first pci device is assigned to domU, toolstack allocates a virtual
249    BAR region from the MMIO hole area. toolstack then sends domctl
250    xc_domain_memory_mapping to map in stage2 translation.
251
252    4.2.1 Reserved Areas in guest memory space
253    -------------------------------------------
254    Parts of the guest address space is reserved for mapping assigned pci device's
255    BAR regions. Toolstack is responsible for allocating ranges from this area and
256    creating stage 2 mapping for the domain.
257
258    /* For 32bit */
259    GUEST_MMIO_BAR_BASE_32, GUEST_MMIO_BAR_SIZE_32
260
261    /* For 64bit */
262
263    GUEST_MMIO_BAR_BASE_64, GUEST_MMIO_BAR_SIZE_64
264
265    Note: For 64bit systems, PCI BAR regions should be mapped from
266    GUEST_MMIO_BAR_BASE_64.
267
268    IPA is allocated from the {GUEST_MMIO_BAR_BASE_64, GUEST_MMIO_BAR_SIZE_64}
269    range and PA is the values read from the BAR registers.
270
271    4.2.2 New entries in xenstore for device BARs
272    ---------------------------------------------
273    toolstack also updates the xenstore information for the device
274    (virtualbar:physical bar).This information is read by xenpciback and returned
275    to the pcifront driver configuration space reads for BAR.
276
277    Entries created are as follows:
278    /local/domain/0/backend/pci/1/0
279    vdev-N
280    BDF = ""
281    BAR-0-IPA = ""
282    BAR-0-PA = ""
283    BAR-0-SIZE = ""
284    ...
285    BAR-M-IPA = ""
286    BAR-M-PA = ""
287    BAR-M-SIZE = ""
288
```

```
289   Note: Is BAR M SIZE is 0, it is not a valied entry.
290
291   4.2.4 Hypercall Modification for bdf mapping notification to xen
292   ---------------------------------------------------------------
293   Guest devfn generation currently done by xen-pciback to be done by toolstack
294   only. Guest devfn is generated at the time of domain creation (if pci devices
295   are specified in cfg file) or using xl pci-attach call.
296
297   5. DomU FrontEnd Bus Changes
298   ---------------------------------------------------------------------------------
299   5.1 Change in Linux PCI ForntEnd - backend driver for MSI/X programming
300   -----------------------------------------------------------------------------
301   FrontEnd backend communication for MSI is removed in XEN ARM. It would be
302   handled by the gic-its driver in guest kernel and trapped in xen.
303
304   5.2 Frontend bus and interrupt parent vITS
305   -------------------------------------------
306   On the Pci frontend bus msi-parent gicv3-its is added. As there is a single
307   virtual its for a domU, as there is only a single virtual pci bus in domU. This
308   ensures that the config_msi calls are handled by the gicv3 its driver in domU
309   kernel and not utilising frontend-backend communication between dom0-domU.
310
311   It is required to have a gicv3-its node in guest device tree.
312
313   6. NUMA domU and vITS
314   -------------------------
315   a) On NUMA systems domU still have a single its node.
316   b) How can xen identify the ITS on which a device is connected.
317   - Using segment number query using api which gives pci host controllers
318   device node
319
320   struct dt_device_node* pci_hostbridge_dt_node(uint32_t segno)
321
322   c) Query the interrupt parent of the pci device node to find out the its.
```