

Extended Destination ID: Supporting x86 virtual machines with many vCPUs

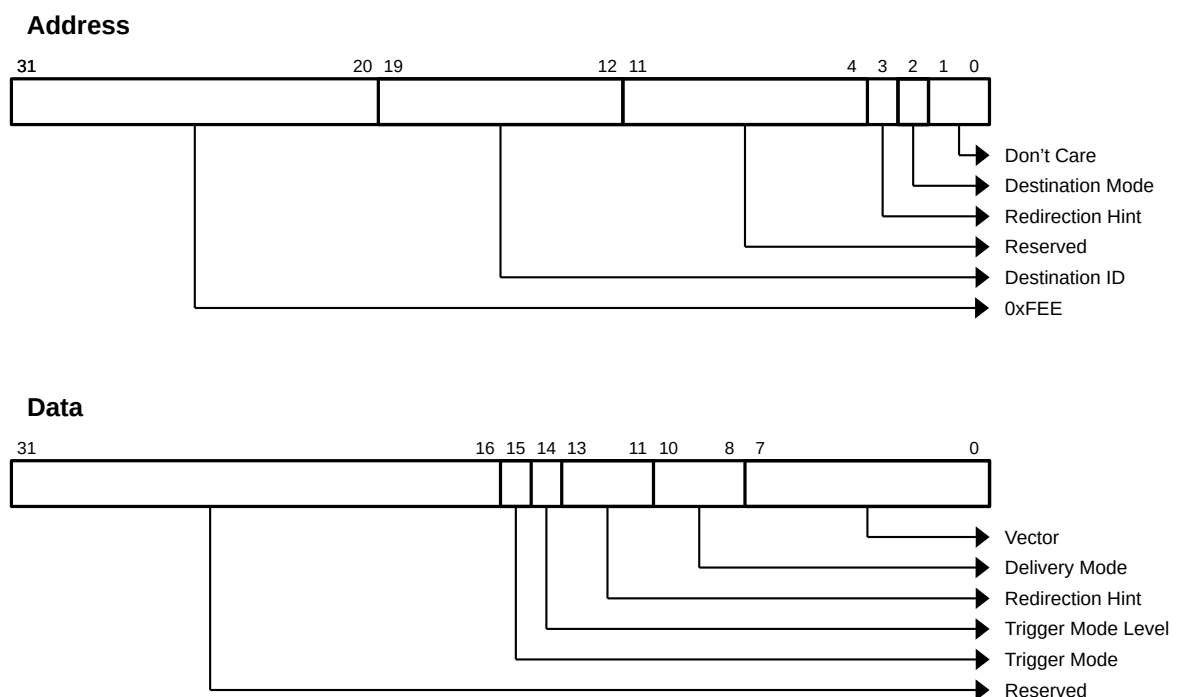
David Woodhouse <dwmw@amazon.co.uk> 2025-08-14

BACKGROUND

Supporting x86 virtual machines with more than 255 vCPUs presents a new challenge. Intel's xAPIC design only supports 8 bits of destination APIC ID for targeting interrupts. The x2APIC design expands this to 32 bits for Inter-Processor Interrupts (IPIs) but not intrinsically for external interrupts.

External interrupts are received as Message Signaled Interrupts (MSI), where the device requesting attention performs a memory write to a physical address range owned by the APIC. The data written, and the low bits of the address, convey all the information about which interrupt vector is to be delivered in which mode to which CPU.

Even with the CPU using x2APIC, the native MSI messages only have space for 8 bits of destination ID. This is the legacy MSI format, now called "Compatibility Format" in Intel VT-d documentation:



The theory of the hardware design is that systems with an APIC ID above 255 shall use Interrupt Remapping, a feature of the IOMMU. With Interrupt Remapping, a new format of MSI messages is used, which merely contains an index into the IOMMU's remapping table — which *can* contain the full 32 bits of destination ID.

DMA translation

Interrupt Remapping was added as an optional feature to Intel and AMD's IOMMU specifications. But DMA translation itself is not an optional feature — that is what the IOMMU was *invented* for in the first place. Thus, there is no architectural way to expose an IOMMU which supports *only* Interrupt Remapping and not DMA translation.

DMA translation is problematic to expose to guests. It is extremely security-sensitive, as an error in configuring DMA translations could expose memory belonging to the host or other guests. Some devices also behave badly in the event of unanticipated DMA faults which could be triggered by a guest controlling its own IOVA space. And in some cases it can also require large amounts of memory to shadow the guest's translation tables, turning the guest's IOVA → GPA translations into IOVA → HPA tables for the IOMMU to use directly.

For all these reasons, but especially because of the additional IOTLB pressure caused by 4KiB entries, we do not want to have to implement a full virtual IOMMU just to allow a guest to target interrupts at vCPUs above #255.

SOLUTION 1: Intel IOMMU without DMA remapping

Although it does not seem to have been intentional, the Intel IOMMU does actually have a subtle way to indicate that it does not support DMA remapping. There are “Supported Adjusted Guest Address Widths” (SAGAW) bits in the Capability Register which separately indicate support for 3-level, 4-level and 5-level page tables. If *none* of those bits are set, the IOMMU cannot be used for DMA translation.

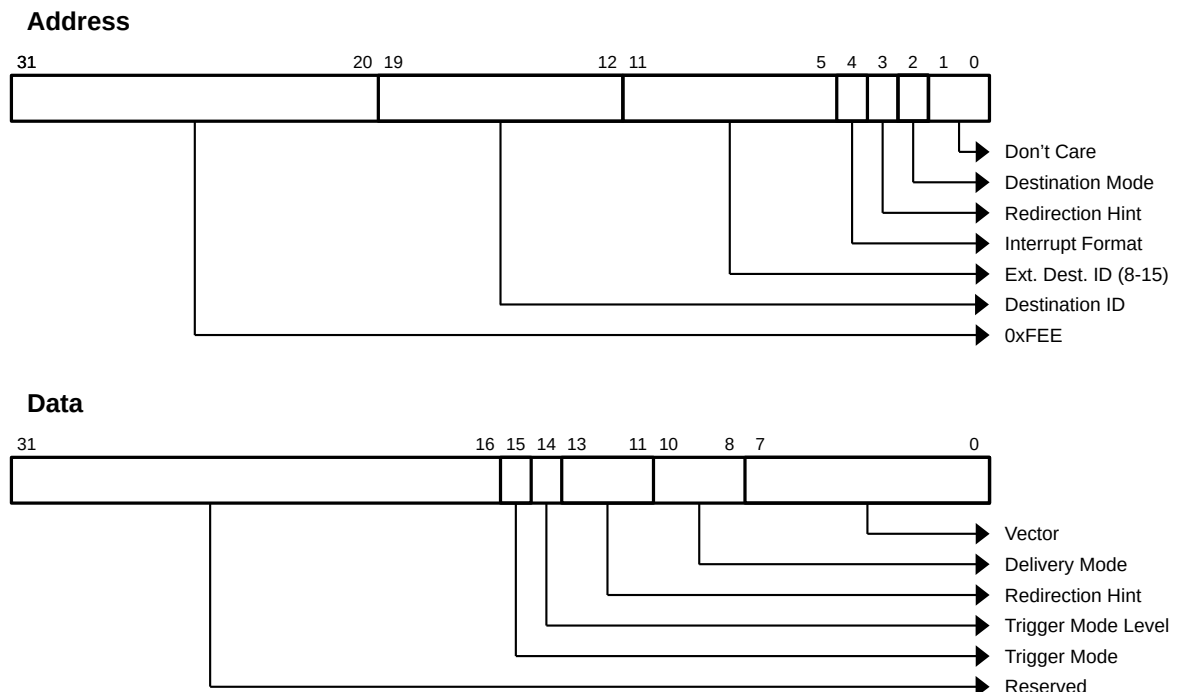
This trick of setting all SAGAW bits to zero is used by QEMU when the ‘`dma-translation=off`’ option is specified to the `intel-iommu` device, and also the EC2 Nitro Hypervisor for instances with high numbers of vCPUs (e.g. [High Memory U7i](#) instances). As a guest, Windows will correctly refrain from attempting to use such an IOMMU for DMA translation, but will use it for Interrupt Remapping. Linux behaves the same, since commit [c40aaac10](#) (present in stable kernels as far back as 4.4.x).

AMD’s IOMMU does not have a similar trick to expose Interrupt Remapping functionality without DMA translation, but that shouldn’t be needed anyway — just as a system may combine an AMD CPU with an Intel network card, it may just as feasibly have an AMD CPU and an Intel IOMMU. This is a combination which is supported by QEMU, and works in Linux and other operating systems. However, current versions of Windows have a bug where they initialize the Intel IOMMU and then feed it MSI messages designed for the AMD IOMMU.

SOLUTION 2: Native 15-bit MSI support using Extended Destination ID field

The simplest option is to forget about Interrupt Remapping completely. There are plenty of extra bits available in the native MSI messages, and what were once marked as ‘Reserved’ in bits 4-11 of the MSI address have even been documented as ‘Extended Destination ID’ in Intel documentation since as far back as ICH4 in 2002 ([82801db-io-controller-hub-4-datasheet.pdf](#) §5.8.5.5, §9.5.10).

Intel now uses the lowest of those Extended Destination ID bits as the Interrupt Format bit which indicates Remappable Format MSIs, so the remaining 7 (bits 5-11) are available to provide an additional 7 bits of destination APIC ID, allowing for up to 32768 vCPUs to be targeted by interrupts.



Guest Operating System support

- **Linux:** has [supported](#) this enlightenment since the v5.11 kernel, released in March 2020.
- **Windows:** Microsoft have been engaged to support this enlightenment (*and fix the IOMMU bug above*).
- **FreeBSD:** Support requested in https://bugs.freebsd.org/bugzilla/show_bug.cgi?id=288122

Hypervisor advertisement

This enlightenment is common to Hyper-V, Xen, KVM and FreeBSD/Bhyve. Each exposes the feature through CPUID:

- **Hyper-V:** sets the `HYPERV_VS_PROPERTIES_EAX_EXTENDED_IOAPIC_RTE` bit in the `HYPERV_CPUID_VIRT_STACK_PROPERTIES` leaf.
(Ref: [Linux commit d981059e13ff](#) from Microsoft as *Hyper-V itself is poorly documented*).
- **Xen:** sets the `XEN_HVM_CPUID_EXT_DEST_ID` bits in the HVM-specific features leaf.
(Ref: [Xen commit f592322062f](#)).
- **Linux/KVM:** sets the `KVM_FEATURE_MSI_EXT_DEST_ID` in the `KVM_CPUID_FEATURES` leaf.
(Ref: [Linux commit 5a169bf04cd2](#)).
- **FreeBSD/Bhyve:** FreeBSD 15.0 adds `CPUID BHYVE_FEATURES` leaf and `CPUID BHYVE_FEAT_EXT_DEST_ID` bit.
(Ref: [commits 313a68ea20b4](#), [4322d597453d](#)).

TECHNICAL NOTES

I/O APIC vs. MSI

The I/O APIC is a device for turning pin-based interrupts into MSIs. Each pin is configured with a Redirection Table Entry which contains all the same bits as the MSI message described above, just in a different order. The I/O APIC has also documented the “Extended Destination ID” bits since the ICH4 specification referenced above.

The 7 bits used for the Extended Destination ID in this specification correspond to bits 49-55 of the I/O APIC RTE.

Given a sufficiently clean software design, no changes are needed in I/O APIC implementation in either host or guest side in order to implement this specification. QEMU’s [ioapic_entry_parse](#) function, for example, merely swizzles bitfields of the RTE into the resulting MSI address/data fields.

Since commit [5d5a97133887](#), Linux just does the converse of the same bit-swizzling, allowing the ‘upstream’ target (*either IOMMU or native APIC*) to compose the MSI message and then just shifting the bits around.

This detail of the I/O APIC is why only 7 additional bits from the existing Extended Destination ID are used. It would be theoretically possible to use more bits of the MSI message, giving a full 32 bits of destination APIC ID — but that would not be addressable through the I/O APIC without an extension to the RTE format and corresponding modifications to I/O APIC support code in guests and hypervisors. This version is even supportable in hardware.

For more detail about MSI message formats and the I/O APIC, see

<http://david.woodhou.se/more-than-you-ever-wanted-to-know-about-x86-msis.txt>

Physical vs. Logical destination mode

Although most operating systems don’t use it for external interrupts, the logical addressing mode allows an interrupt to be target at more than one CPU within a ‘cluster’. With x2APIC enabled, the 32-bit target APIC ID is split into two 16-bit fields, the top half being the cluster ID and the bottom half being a bitmask of CPUs within that cluster to address.

Since this enlightenment provides only 15 bits and not the full 32 bits, that would only be enough to address the first 15 CPUs (in cluster #0). Operating systems are therefore expected to use the physical addressing mode, which allows them to address CPUs up to #32767.

Hypervisors may set the `FORCE_APIC_PHYSICAL_DESTINATION_MODE` bit in the ACPI FADT table as a hint to guest operating systems, or guest operating systems may work that requirement out for themselves. The logical vs. physical addressing mode is a per-interrupt setting, so it is perfectly feasible to use logical addressing for Inter-Processor Interrupts (allowing multicast/broadcast) while using physical addressing for external interrupts.

xAPIC broadcast

In the xAPIC mode, the target 0xFF is considered a broadcast, and all CPUs will respond to such an interrupt. This requires no special handling or compatibility with the 15-bit extension because the interpretation is done by each (potential) target CPU. Just as with real hardware, if a target APIC ID has all of the low bits set (e.g. 0x1FF), any CPU in xAPIC mode will interpret that as a broadcast and receive it accordingly.