

Kexec Handover and Live Update

Pratyush Yadav <pratyush@kernel.org>

Pratyush Yadav <pratyush@kernel.org>

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

Kexec Handover and Live Update

└─ Agenda

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

Kexec Handover and Live Update
└─What is Live Update?

└─Agenda

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

What is Live Update?

- ▶ It's NOT: live patching, live migration.
- ▶ Updates the kernel or hypervisor with minimal disruption for underlying workloads.
- ▶ Most commonly used for hypervisors.
- ▶ Can also be used by other workloads to reduce kernel patching downtime.
- ▶ Multiple cloud providers working together to upstream it.

2025-10-01

Kexec Handover and Live Update

└─What is Live Update?

└─What is Live Update?

What is Live Update?

- ▶ It's NOT: live patching, live migration.
- ▶ Updates the kernel or hypervisor with minimal disruption for underlying workloads.
- ▶ Most commonly used for hypervisors.
- ▶ Can also be used by other workloads to reduce kernel patching downtime.
- ▶ Multiple cloud providers working together to upstream it.

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

Kexec Handover and Live Update

└─ High level overview

└─ Agenda

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

- ## Kexec Handover and Live Update

└ Live update flow

- ▶ The system is in normal state.
- ▶ The system software starts the live update process.
- ▶ Serializes state keeping VMs active but with limited capabilities.
- ▶ Pauses VMs and does final serialization.
- ▶ Loads and next kernel and hands over the serialized data.
- ▶ Next kernel deserializes the data.
- ▶ Resumes VM, returning normal operation.

What gets preserved?

- ▶ VM metadata
- ▶ VM memory
- ▶ Passthrough devices
- ▶ IOMMU mappings

2025-10-01

Kexec Handover and Live Update

- └ High level overview
 - └ What gets preserved?

What gets preserved?

- ▶ VM metadata
- ▶ VM memory
- ▶ Passthrough devices
- ▶ IOMMU mappings

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

Kexec Handover and Live Update

└ Building blocks

└┬ Agenda

- Complex feature, not easy to do it in one go.
- Upstreaming as a set of building blocks instead.

2025-10-01

- └ Building blocks
 - └ Kexec Handover (KHO)

- Kexec Handover (KHO)

- ▶ Creates a mechanism for kernel-to-kernel communication.
- ▶ Provides mechanism to mark memory as preserved.
- ▶ Makes sure preserved memory does not get used by the next kernel.
- ▶ Passes this information over kexec.

KHO: Memory preservation

```
int kho_preserve_folio(struct folio *folio);  
int kho_unpreserve_folio(struct folio *folio);  
struct folio *kho_restore_folio(phys_addr_t phys);
```

2025-10-01

Kexec Handover and Live Update

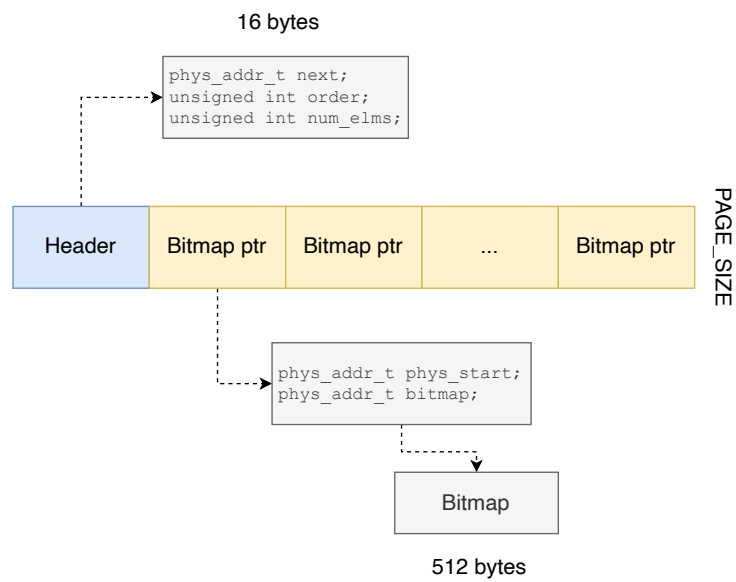
└ Building blocks

└ KHO: Memory preservation

KHO: Memory preservation

```
int kho_preserve_folio(struct folio *folio);  
int kho_unpreserve_folio(struct folio *folio);  
struct folio *kho_restore_folio(phys_addr_t phys);
```

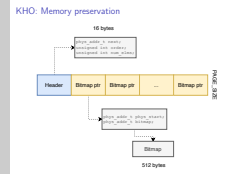
KHO: Memory preservation



2025-10-01

Kexec Handover and Live Update

- Building blocks
- KHO: Memory preservation



KHO: Booting up

- ▶ Pre-reserved scratch area for early boot.
- ▶ Passing KHO metadata: setup data on x86, chosen node in FDT on arm64.

```
struct kho_data {
    __u64 fdt_addr;
    __u64 fdt_size;
    __u64 scratch_addr;
    __u64 scratch_size;
} __attribute__((packed));
```

```
chosen {
    linux,kho-fdt = <...>;
    linux,kho-scratch = <...>;
};
```

2025-10-01

Kexec Handover and Live Update

└ Building blocks

└─KHO: Booting up

- Mention that kexec image and all early boot allocations go in scratch.
- Mention that chosen node gets set at kexec load time.

KHO: Booting up

- ▶ Pre-reserved scratch area for early boot.
- ▶ Passing KHD metadata: setup data on x86, chosen node in FDT on arm64.

```
struct kho_data {
    __u64 fdt_addr;
    __u64 fdt_size;
    __u64 scratch_addr;
    __u64 scratch_size;
} __attribute__((packed));
```

```
chosen {
    linux,kbo-fdt = <...>;
    linux,kbo-scratch = <...>;
};
```

KHO: Booting up

- ▶ On early boot, only allocate from scratch.

```
enum memblock_flags choose_memblock_flags(void)
{
    if (kho_scratch_only)
        return MEMBLOCK_KHO_SCRATCH;
    [...]
}
```

- ▶ After early boot, mark preserved pages as reserved and turn off scratch-only mode
- ▶ Reserved pages don't get released to buddy allocator.

2025-10-01

Kexec Handover and Live Update

└ Building blocks

└ KHO: Booting up

- ▶ On early boot, only allocate from scratch.

```
enum memblock_flags choose_memblock_flags(void)
{
    if (kho_scratch_only)
        return MEMBLOCK_KHO_SCRATCH;
    [...]
}
```

- ▶ After early boot, mark preserved pages as reserved and turn off scratch-only mode

- ▶ Reserved pages don't get released to buddy allocator.

2025-10-01

- ## Kexec Handover and Live Update

- └ Live Update Orchestrator (LUO)

- ### Live Update Orchestrator (LUO)

- ▶ LZO provides a way for userspace to control the live update process.
- ▶ Allows marking which resources to preserve.
- ▶ Provides a state machine to co-ordinate all the components.
- ▶ API is exposed through a set of IOCTLS.

LUO: States

2025-10-01

Kexec Handover and Live Update

- └ Building blocks

- LUO: States

```
struct liveupdate_ioctl_set_event {
    __u32      size;
    __u32      event;
};
```

- ▶ LIVEUPDATE_PREPARE: Normal -> Prepared
- ▶ LIVEUPDATE_FREEZE: Prepared -> Frozen
- ▶ LIVEUPDATE_FINISH: Updated -> Normal
- ▶ LIVEUPDATE_CANCEL: Prepared -> Normal

LUO: States

```
struct liveupdate_ioctl_set_event {
    __u32    size;
    __u32    event;
};
```

- ▶ LIVEUPDATE_PREPARE: Normal -> Prepared
- ▶ LIVEUPDATE_FREEZE: Prepared -> Frozen
- ▶ LIVEUPDATE_FINISH: Updated -> Normal
- ▶ LIVEUPDATE_CANCEL: Prepared -> Normal

LUO: File Descriptors

- ▶ Userspace can pass in supported file descriptors to LUO to mark them for preservation.
- ▶ Not any arbitrary FD, only FDs for supported file types.

```
struct liveupdate_ioctl_fd_preserve {
    __u32          size;
    __s32          fd;
    __aligned_u64  token;
};
```

2025-10-01

Kexec Handover and Live Update

└ Building blocks

└ LUO: File Descriptors

- Give some examples of FDs in Linux: memfd, sockets, VFIO, IOMMUFD, KVM, etc.
- Mention some properties that can change with restore FDs, taking memfd as example.
- Mention that the token can be used to identify the FD after reboot.

- ▶ Userspace can pass in supported file descriptors to LUO to mark them for preservation.
- ▶ Not any arbitrary FD, only FDs for supported file types.

```
struct liveupdate_ioctl_fd_preserve {
    __u32          size;
    __u32          fd;
    __aligned_u64  token;
};
```


LUO: Subsystems

- ▶ For things that can't be described by a FD.
- ▶ Examples: PCI, NVME, ftrace, etc.

2025-10-01

Kexec Handover and Live Update

- └ Building blocks

- LUO: Subsystems

- Mention that not much work done on this so use cases and usage model still unclear.

- ▶ For things that can't be described by a FD.
- ▶ Examples: PCI, NVME, ftrace, etc.

Memory File Descriptor (memfd)

- ▶ memfd attaches a file descriptor to anonymous memory.
- ▶ State preserved: memory contents, size and position.
- ▶ After preserve, cannot add or remove pages from the memfd.
- ▶ Limitations: no sparseness, no swap.

2025-10-01

Kexec Handover and Live Update

└ Building blocks

└ Memory File Descriptor (memfd)

- Mention that memfd is the first user of LUO.
- Mention that pages are pinned and holes are filled.

Memory File Descriptor (memfd)

- ▶ memfd attaches a file descriptor to anonymous memory.
- ▶ State preserved: memory contents, size and position.
- ▶ After preserve, cannot add or remove pages from the memfd.
- ▶ Limitations: no sparseness, no swap.

memfd: preservation format

```

/ {
    pos = <0x...>;
    size = <0x...>;
    folios = [array of memfd_luo_preserved_folio]
};

```

```
struct memfd_luo_preserved_folio {
    u64 foliodesc;
    u64 index;
};
```

- ▶ Foliodesc: bottom 12 bits for flags, rest for PFN.

2025-10-01

Kexec Handover and Live Update

└ Building blocks

- └ memfd: preservation format

memfd: preservation format

```

/ {
    pos = <0x...>;
    size = <0x...>;
    folios = [array of memfd_luo_preserved_folio]
};

struct memfd_luo_preserved_folio {
    u64 foliodesc;
    u64 index;
};

```

- ▶ Foliodesc: bottom 12 bits for flags, rest for PFN

- Explain why we use FDT.

VFIO, PCI, IOMMU, etc...

2025-10-01

Kexec Handover and Live Update

└─ Building blocks

└─ VFIO, PCI, IOMMU, etc. . .

VFIO, PCI, IOMMU, etc...

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

Kexec Handover and Live Update
└ Upstream status

└ Agenda

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

Upstream status

- ▶ KHO is in mainline. See `kernel/kexec_handover.c` and `include/linux/kexec_handover.h`.
- ▶ LUO v4 sent out few days ago. [Patch posting](#). It is starting to stabilize and is on path to upstream soon.
- ▶ memfd support will get merged with the LUO patches.
- ▶ RFCs for PCI, VFIO, IOMMU out.

2025-10-01

Kexec Handover and Live Update

└ Upstream status

└ Upstream status

Upstream status

- ▶ KHO is in mainline. See `kernel/kexec_handover.c` and `include/linux/kexec_handover.h`.
- ▶ LUO v4 sent out few days ago. [Patch posting](#). It is starting to stabilize and is on path to upstream soon.
- ▶ `memfd` support will get merged with the LUO patches.
- ▶ RFCs for PCL, VFIO, IOMMU out.

Agenda

What is Live Update?

High level overview

Building blocks

Upstream status

Future work

2025-10-01

- Kexec Handover and Live Update
 - Future work
 - Agenda

Agenda

- What is Live Update?
- High level overview
- Building blocks
- Upstream status
- Future work

Future work

- ▶ Supporting more subsystems: huge pages, VFIO, IOMMU, PCI, etc.
- ▶ Implementing luod.
- ▶ Improving performance for reboots.
- ▶ Defining a mechanism for kernels to negotiate versions to enable rollback and roll forward to a wider set of kernels.
- ▶ Testing and validation.

2025-10-01

Kexec Handover and Live Update

Future work

└ Future work

Future work

- ▶ Supporting more subsystems: huge pages, VFIO, IOMMU, PCI, etc.
- ▶ Implementing load.
- ▶ Improving performance for reboots.
- ▶ Defining a mechanism for kernels to negotiate versions to enable rollback and roll forward to a wider set of kernels.
- ▶ Testing and validation.

2025-10-01

Kexec Handover and Live Update

- Future work

Thank you for attending the talk!

Thank you for attending the
talk!

2025-10-01

Kexec Handover and Live Update
└ Future work