

OS Implications for Confidential Computing

Jörg Rödel <jroedel@suse.com>

Confidential Computing Definition

Confidential Computing Definition

- Depending on type of Hypervisor (HV) isolation
- Different definitions for Confidential Computing (CoCo):
 - Software Isolation: HV is isolated from guest by software means
 - Hardware Isolation: HV uses of hardware extensions to isolate guest
- This talk will focus on Confidential Computing using Hardware Isolation
- Newer Hardware supports Trusted Execution Environments (TEEs)
- Attestation is integral part of CoCo

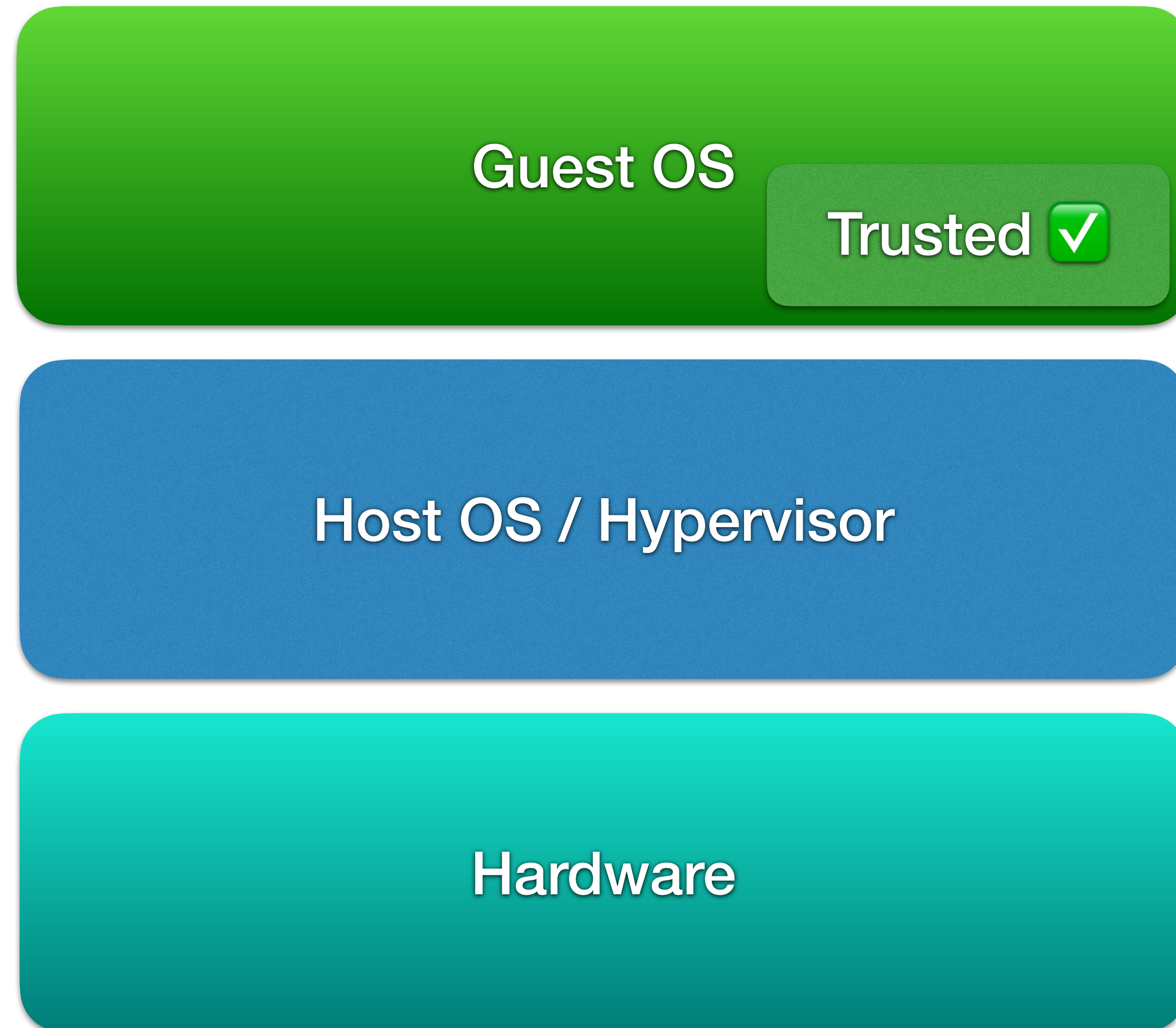
Trusted Execution Base

Guest OS

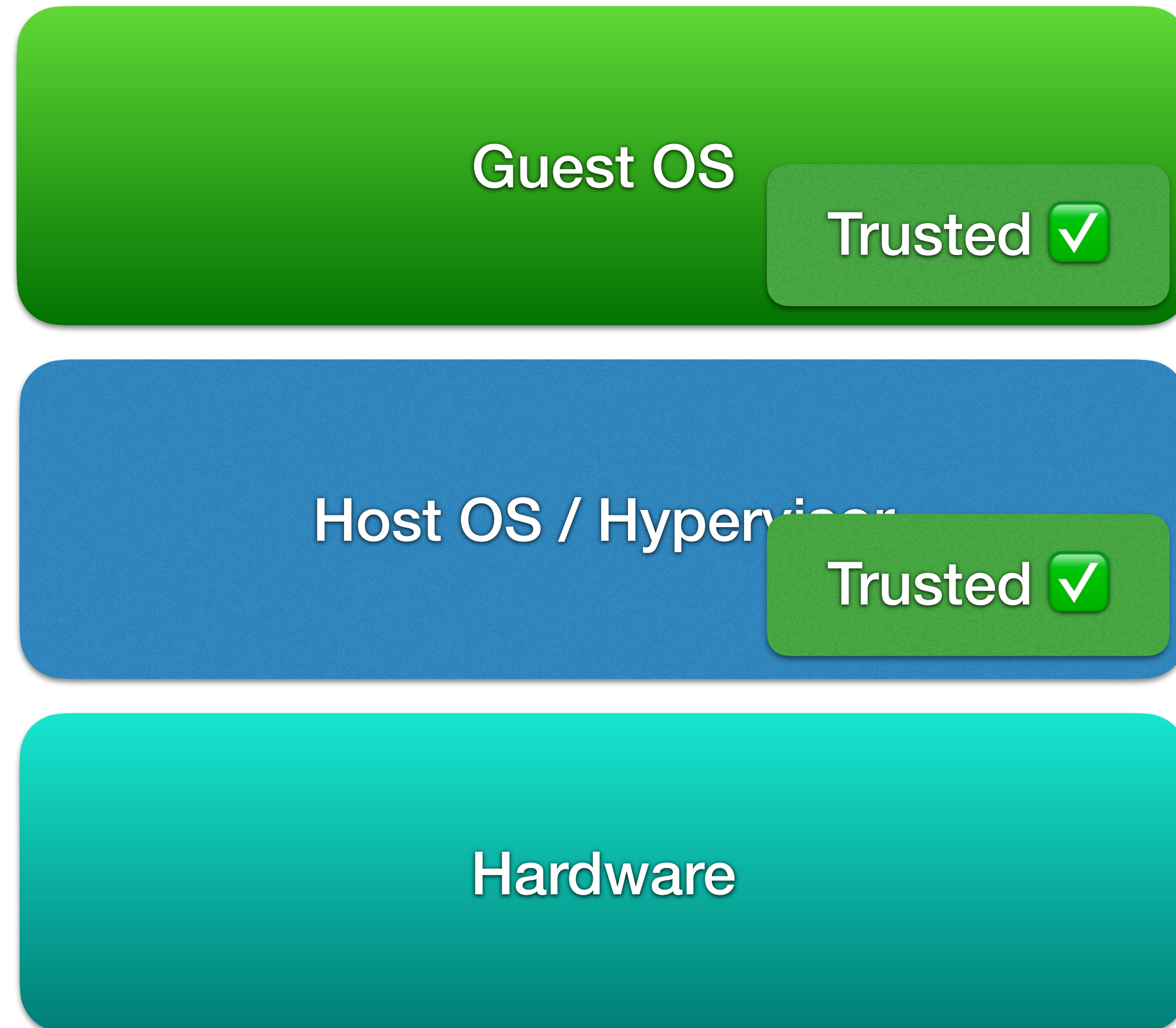
Host OS / Hypervisor

Hardware

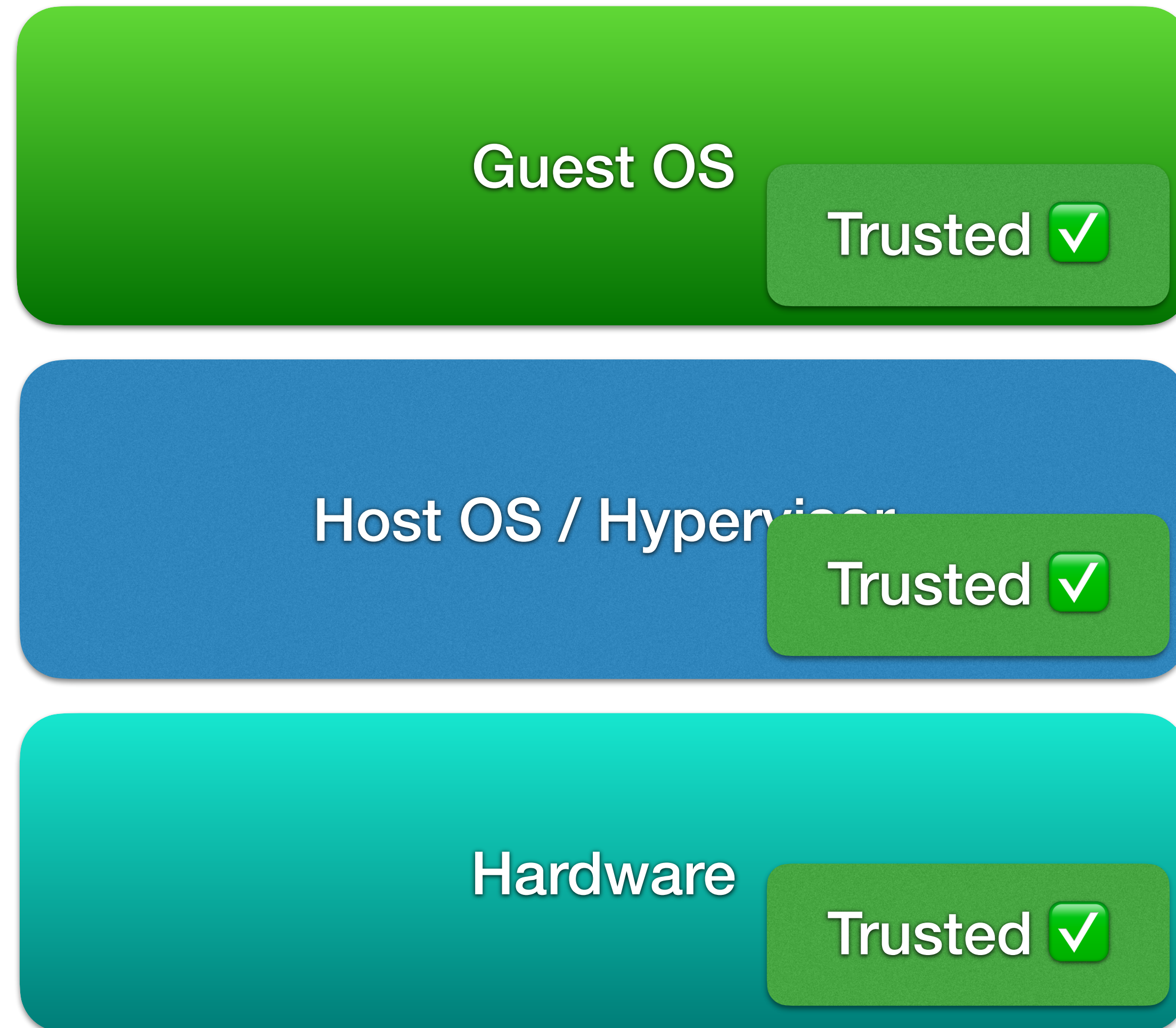
Trusted Execution Base



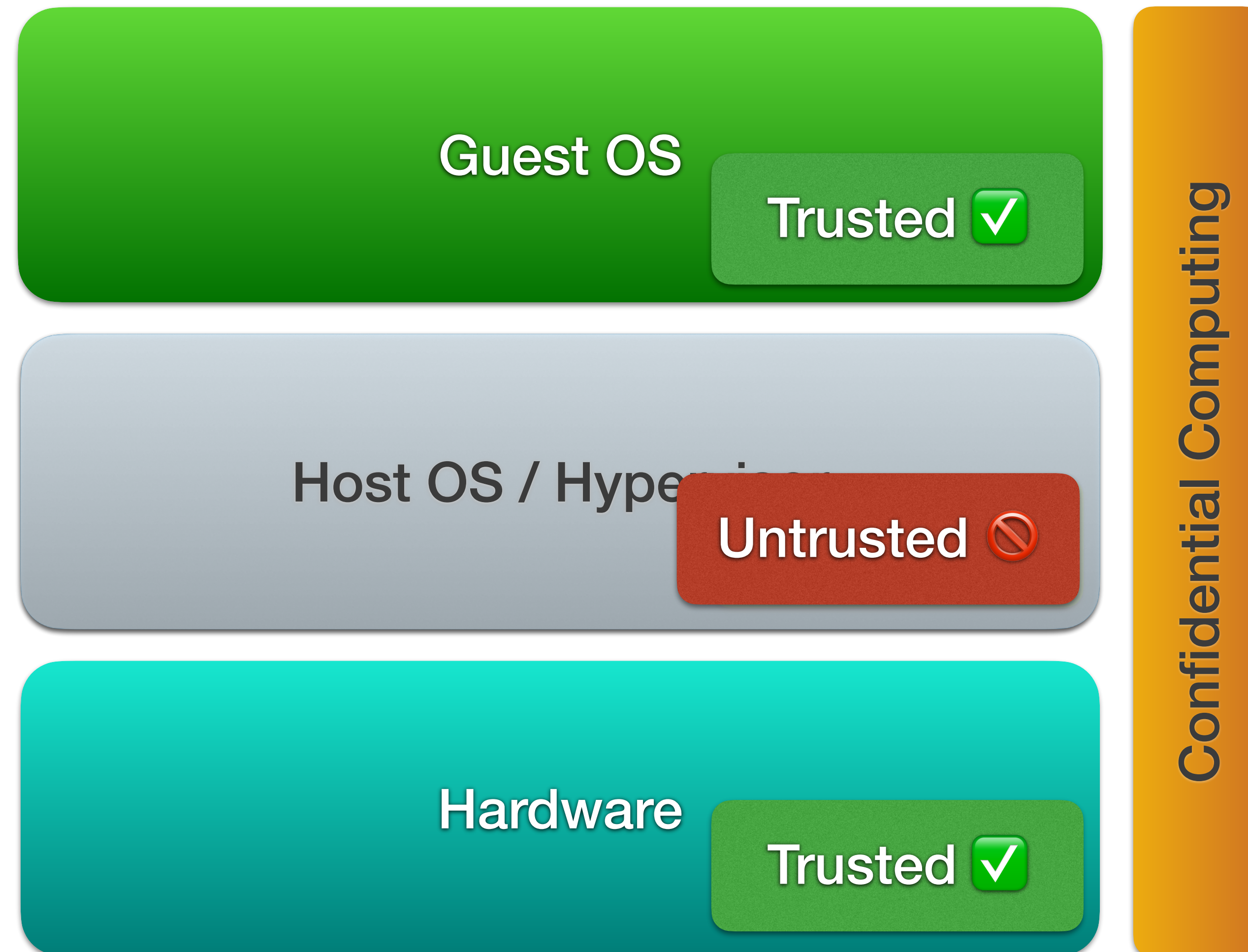
Trusted Execution Base



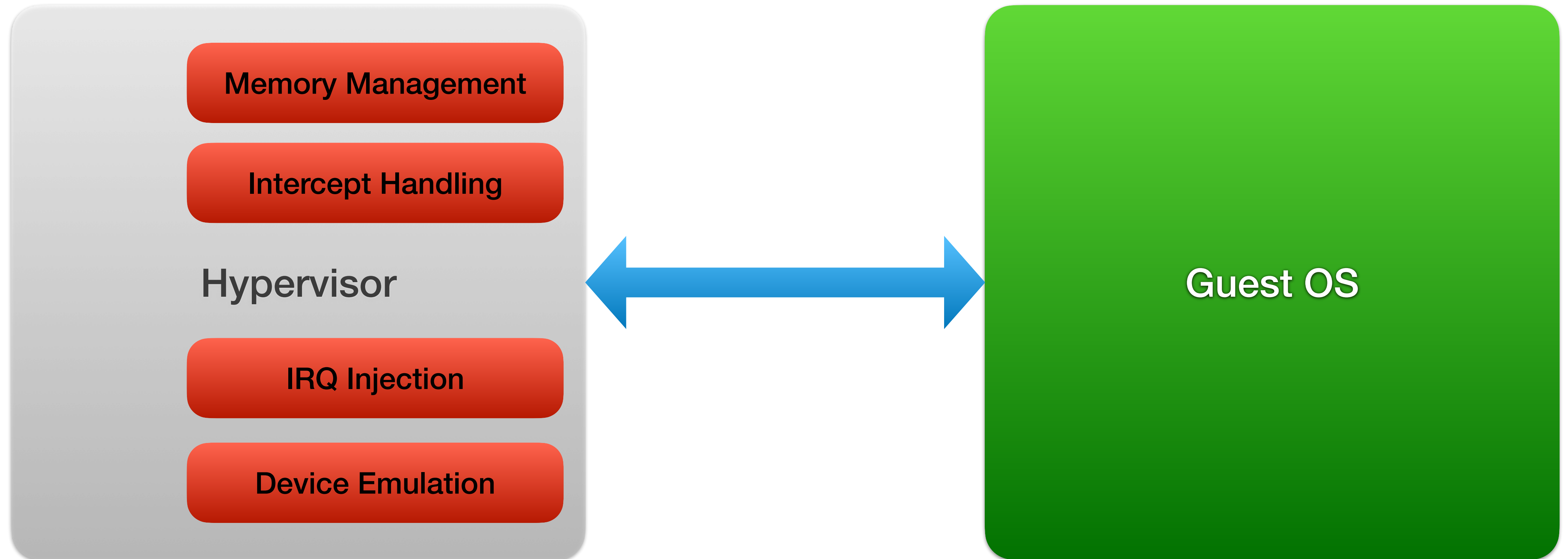
Trusted Execution Base



Trusted Execution Base



Hypervisor Guest Interface



Memory Management

Memory Management

- HV allocates memory for guest usage
- Can read/write arbitrary guest memory at any time
 - Some HV map all guest memory all the time
 - Others use on-demand mappings
- Malicious HV could use that to steal data and/or modify code

Case Study: AMD Secure Encrypted Virtualization

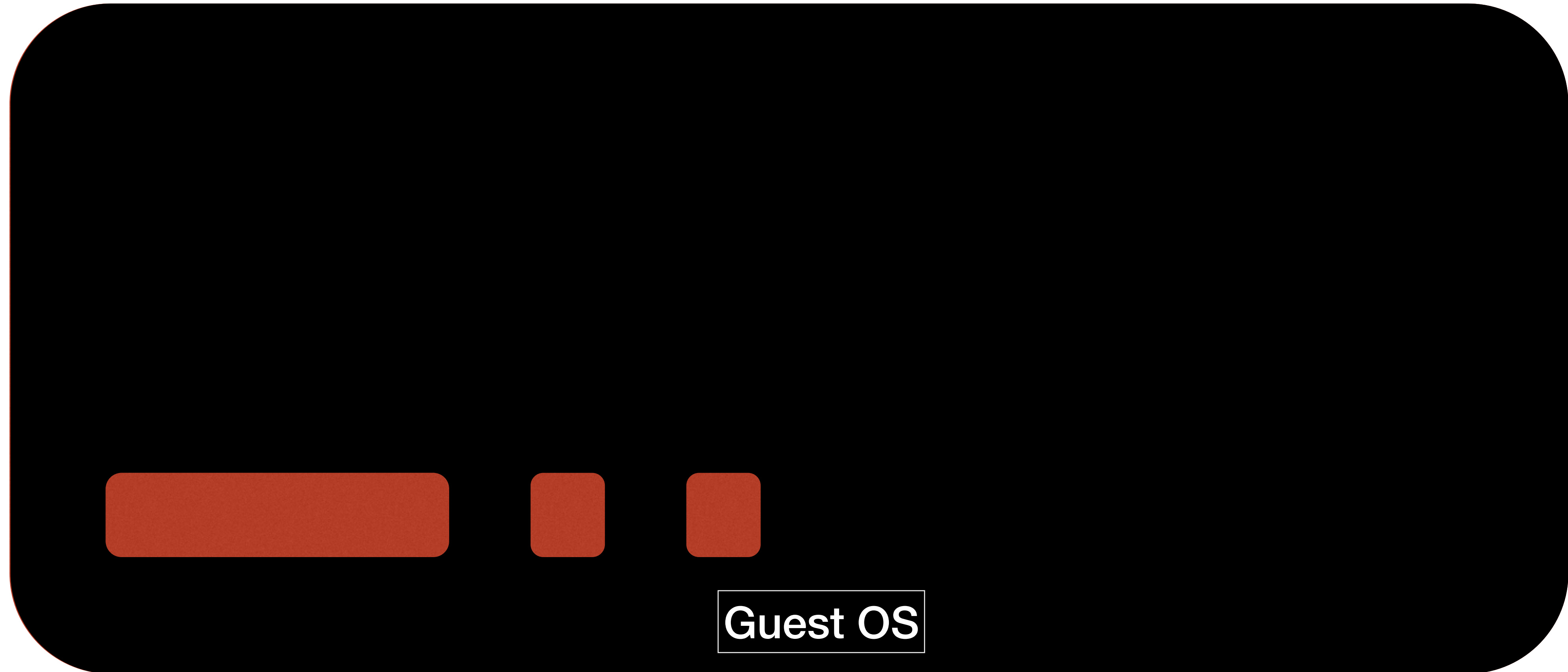
- Confidential Computing in stages - building on each other
 - **Secure Encrypted Virtualization (SEV) - Guest memory encryption**

Guest Memory Encryption

Name: John Doe
Card Number: 1234 5678 9012 3456
Name: Jane Doe
Card Number: 9876 6543 2109 8765
...

Guest OS

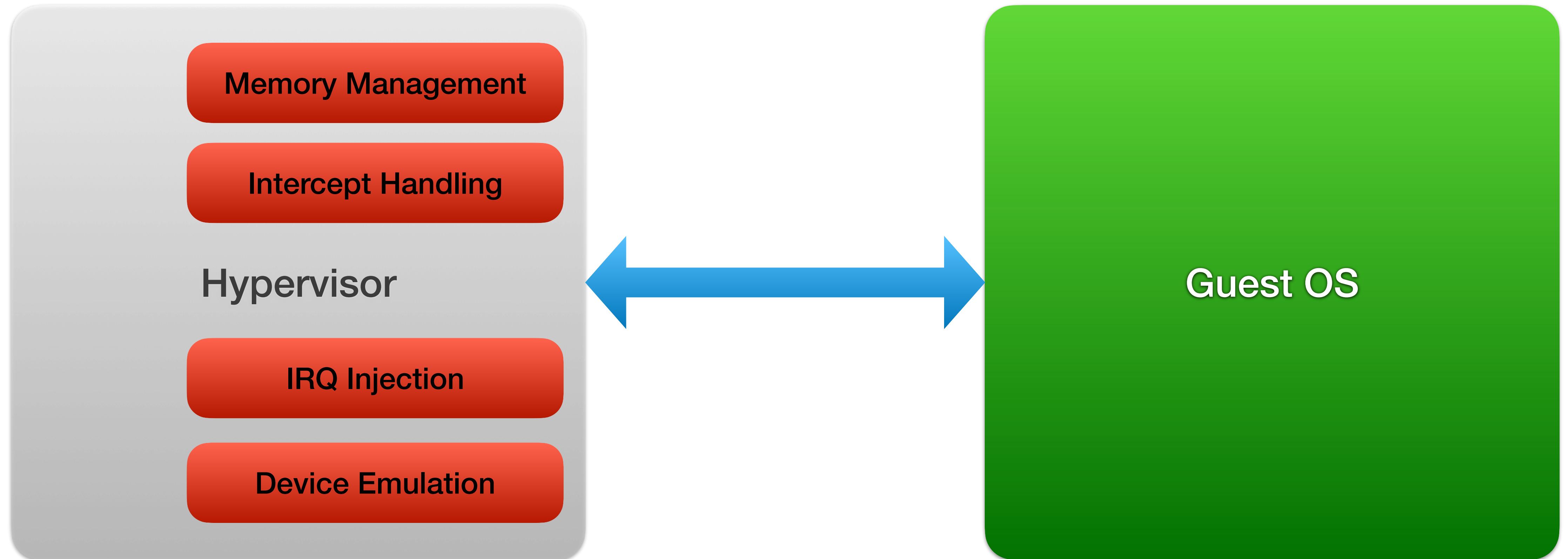
Guest Memory Encryption



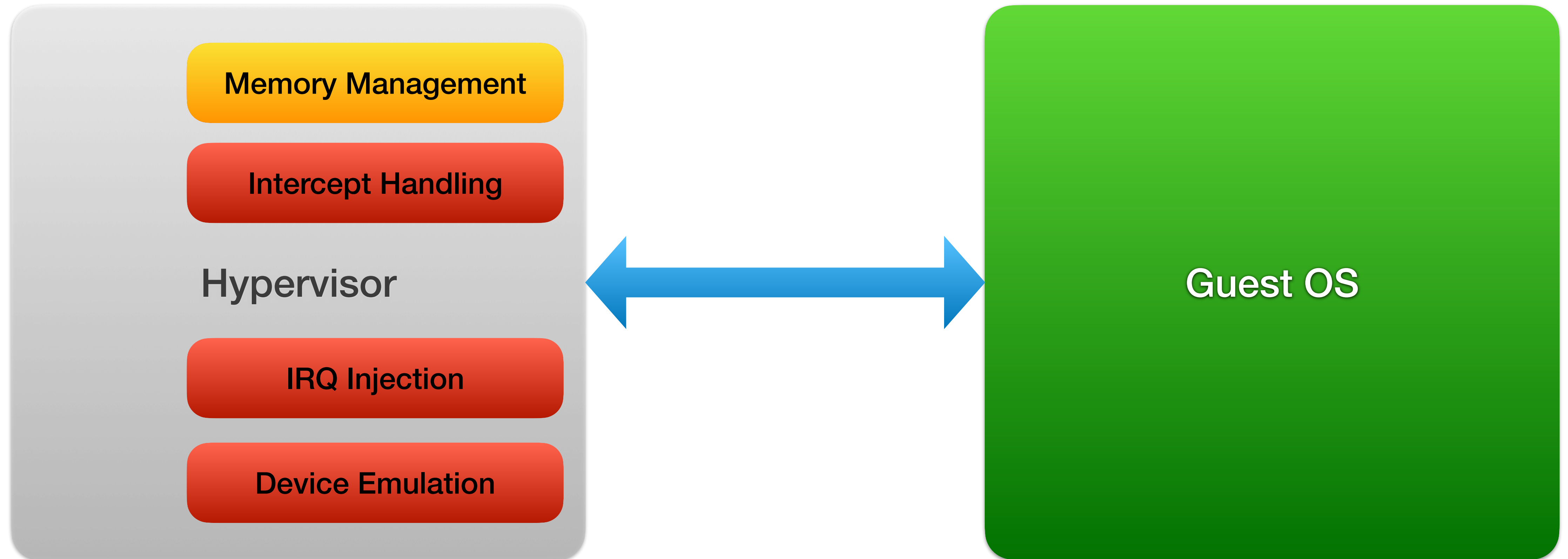
OS Implications

- OS needs to be aware of shared and private (encrypted) memory
- Can decide which memory is shared with the HV
- Sharing implementations:
 - Page-table based: Shared/Private is a flag in the PTE
 - GPA space partitioning: Shared/Private parts of the GPA space
- OS needs changes in Page-table, MMIO, and DMA memory allocation code
- HV can still do memory replay and remapping attacks

Hypervisor Guest Interface



Hypervisor Guest Interface



Intercept Handling

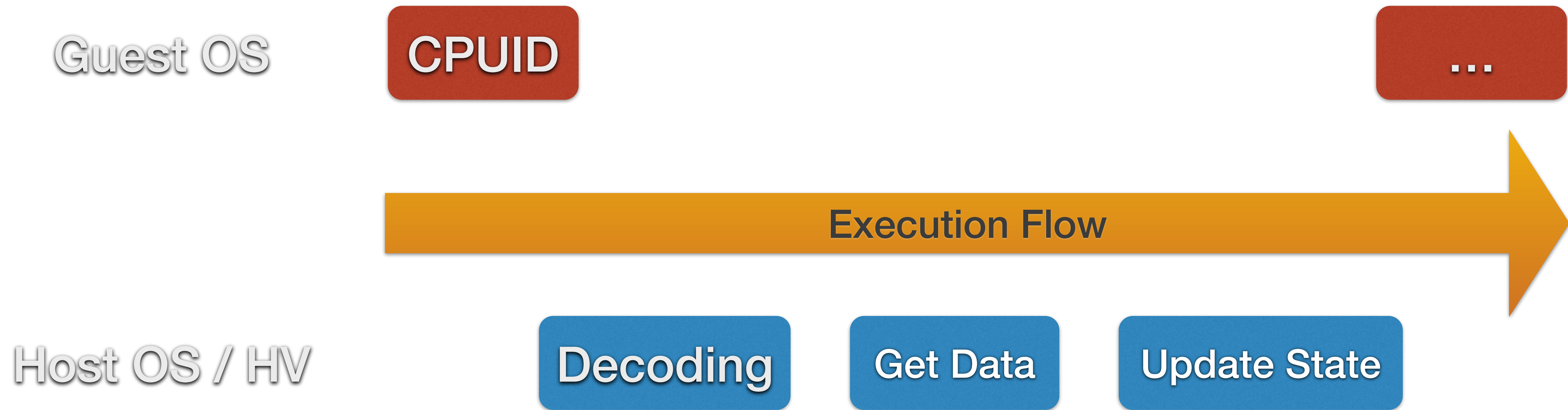
Intercept Handling

- For intercept handling HV needs access to guest registers
- Registers can be implicit or explicit instruction parameters
- For instruction emulation the HV needs to update guest registers
- Malicious HV could attack the guest:
 - Changing control flow by changing IP or SP
 - Steal secret data, e.g. AES keys stored in FPU registers

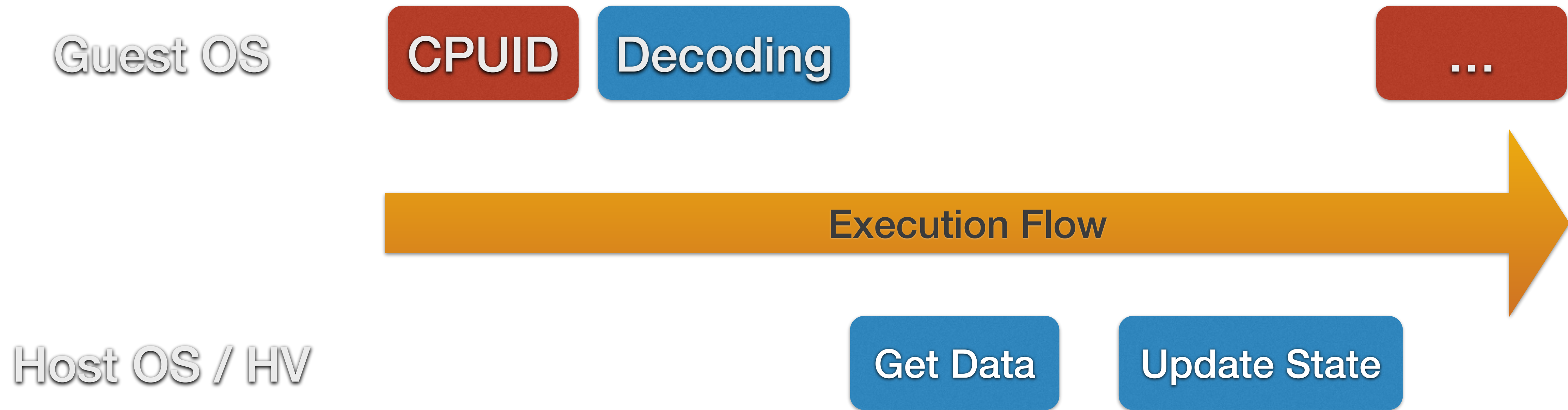
Case Study: AMD Secure Encrypted Virtualization

- Confidential Computing in stages - building on each other
 - Secure Encrypted Virtualization (SEV) - Guest memory encryption
 - **SEV Encrypted State (SEV-ES) - Guest register encryption**

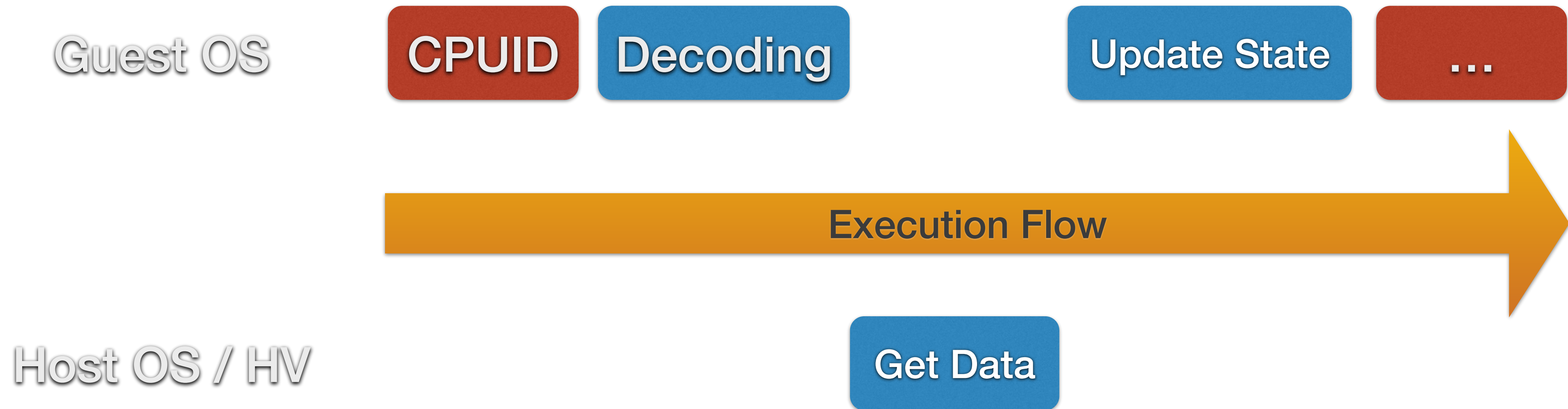
SEV-ES Guest-Host Flow



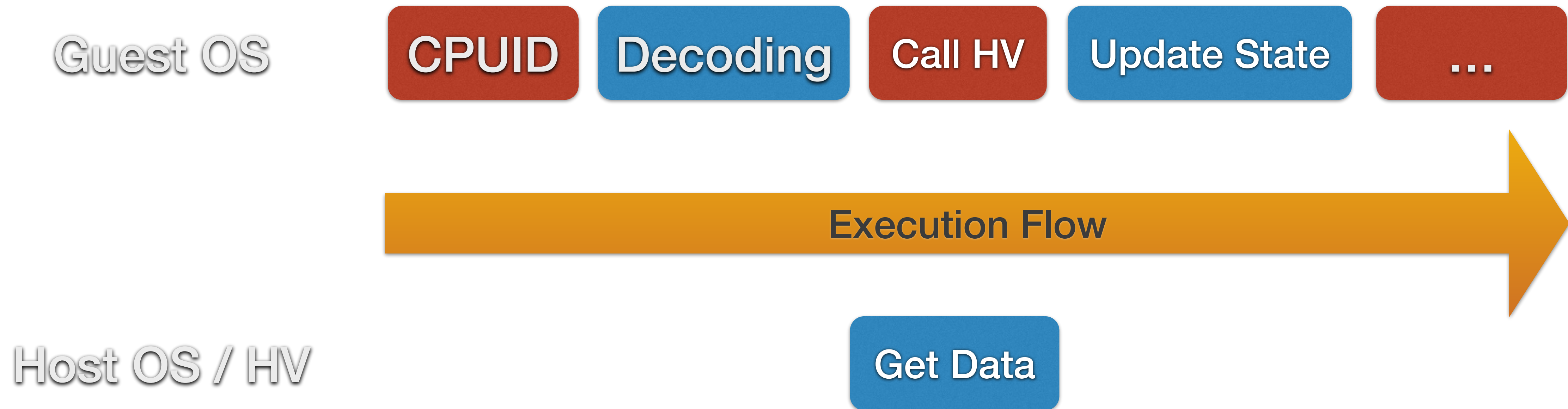
SEV-ES Guest-Host Flow



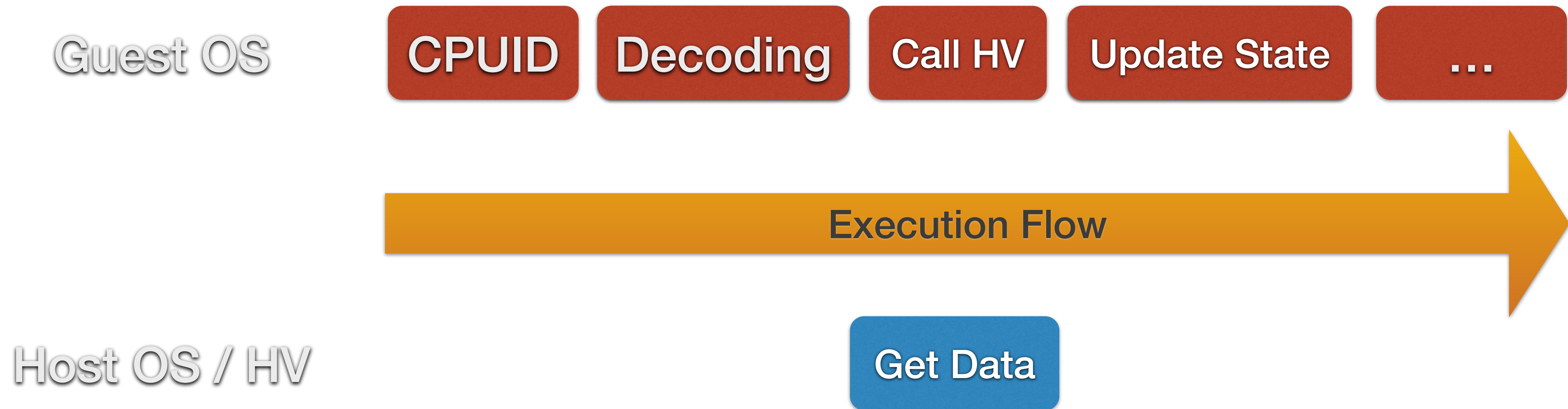
SEV-ES Guest-Host Flow



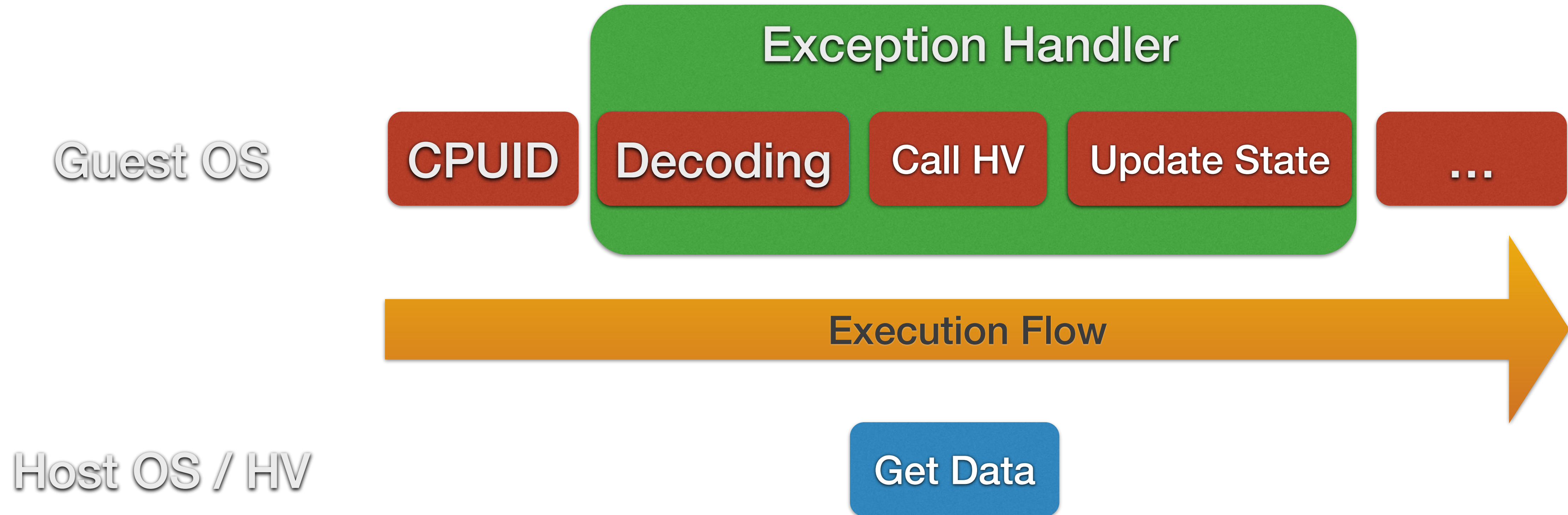
SEV-ES Guest-Host Flow



SEV-ES Guest-Host Flow



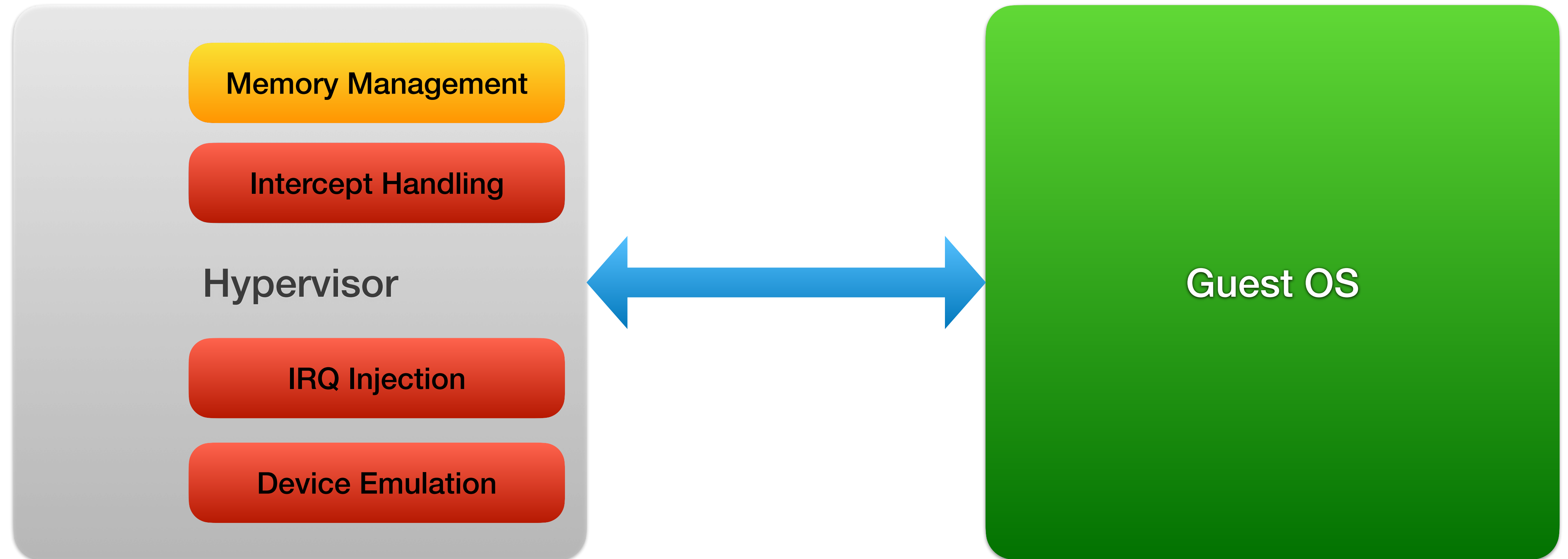
SEV-ES Guest-Host Flow



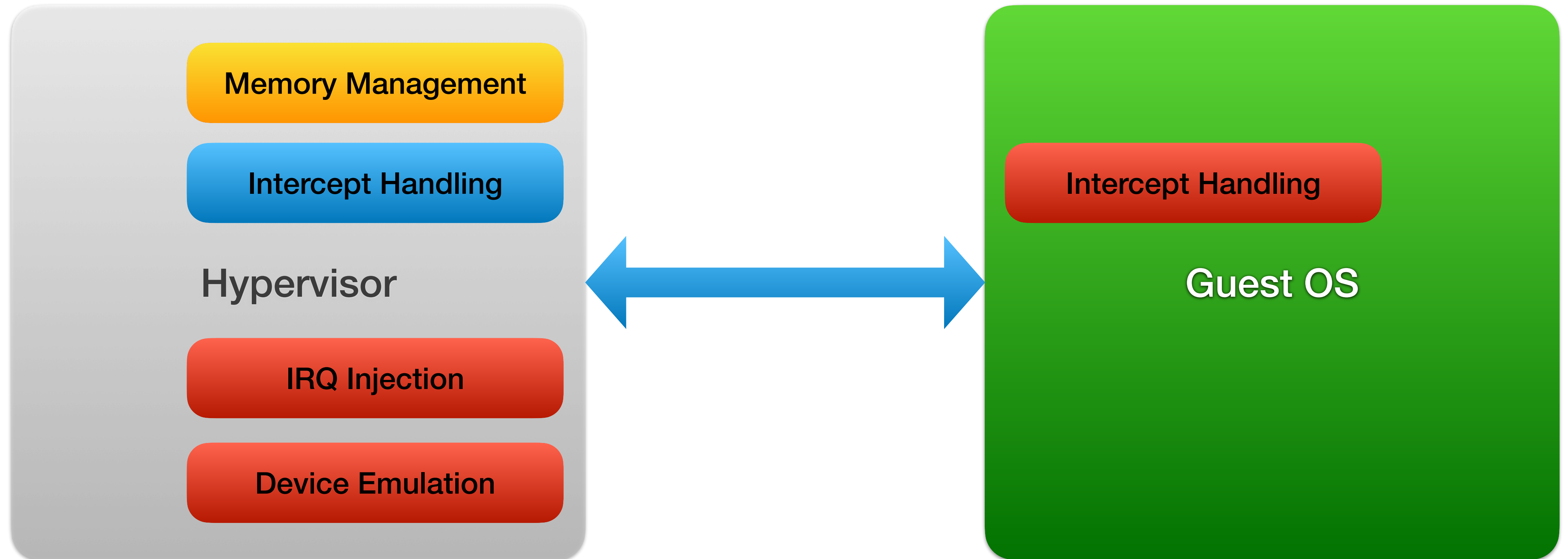
OS Implications

- OS needs to implement new exception handler
 - Will handle intercepts in trusted guest context
 - Uses hyper-calls to exchange data with HV
- Exception handler needed very early in OS boot
 - Before first intercepted instruction is executed
- Exception handler usually includes a full instruction decoder

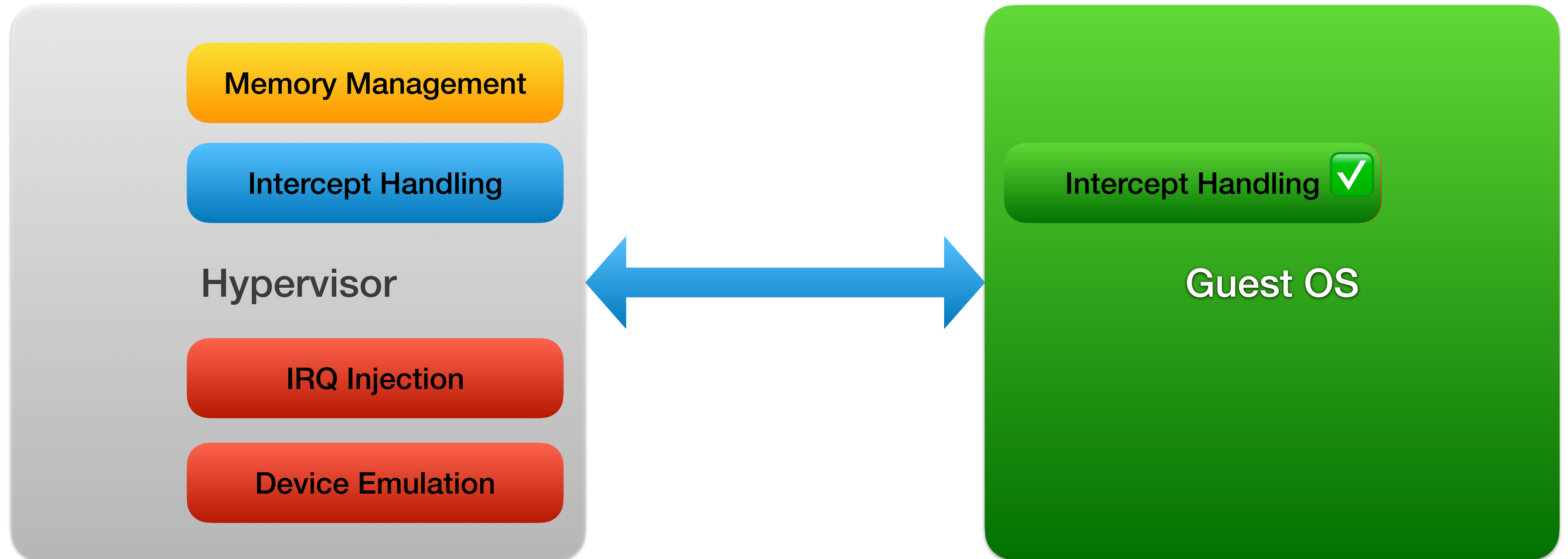
Hypervisor Guest Interface



Hypervisor Guest Interface



Hypervisor Guest Interface



Memory Management Revisited

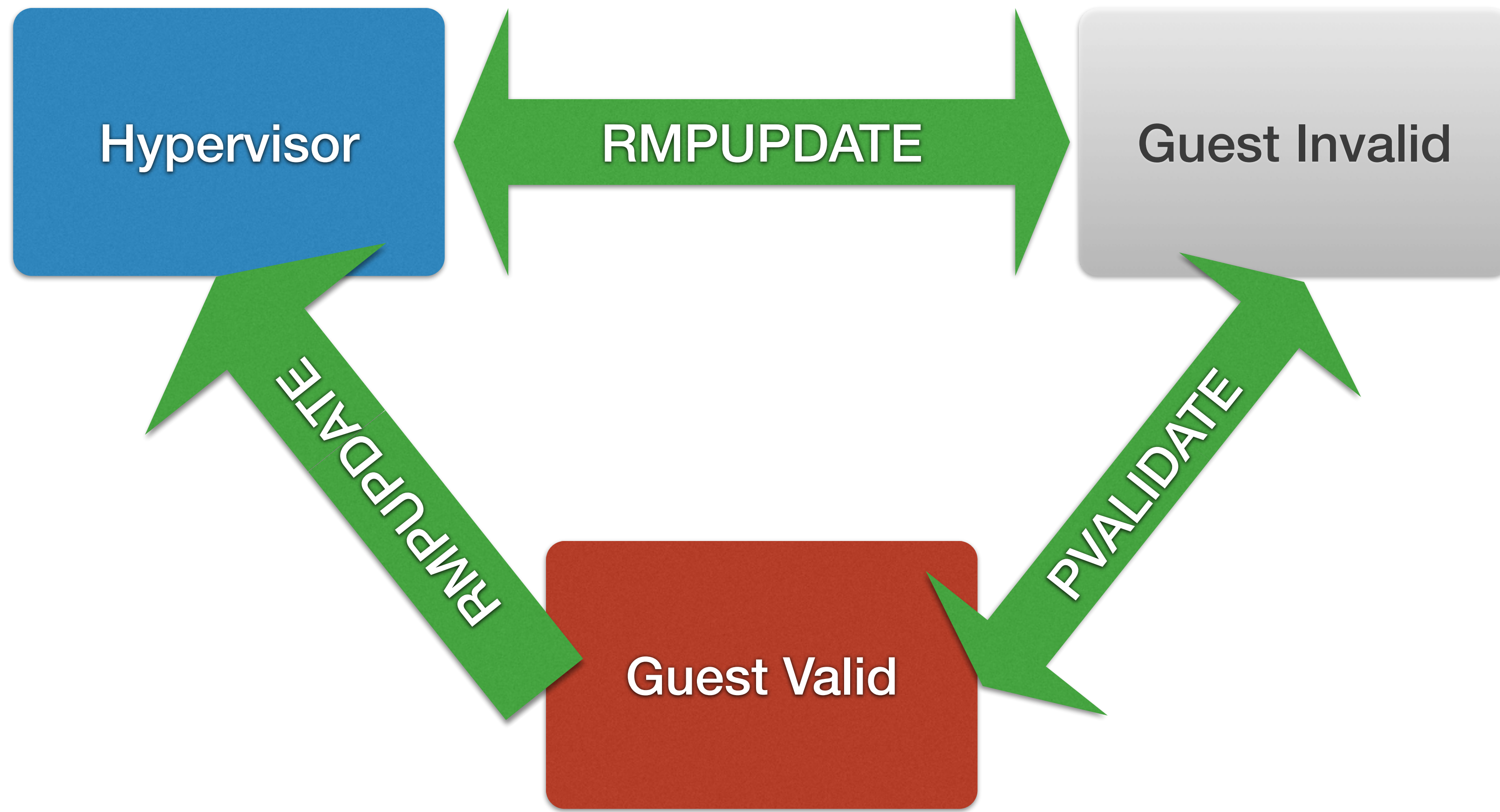
Memory Management Revisited

- Using only memory encryption leaves attack vectors open
 - Memory replay: HV replays an old version of an encrypted page
 - Memory remapping: HV maps encrypted page at different GPA
- Not possible to mitigate
- But become detectable via hardware extension

Case Study: AMD Secure Encrypted Virtualization

- Confidential Computing in stages - building on each other
 - Secure Encrypted Virtualization (SEV) - Guest memory encryption
 - SEV Encrypted State (SEV-ES) - Guest register encryption
 - **SEV Secure Nested Paging (SEV-SNP) – Introducing page states**

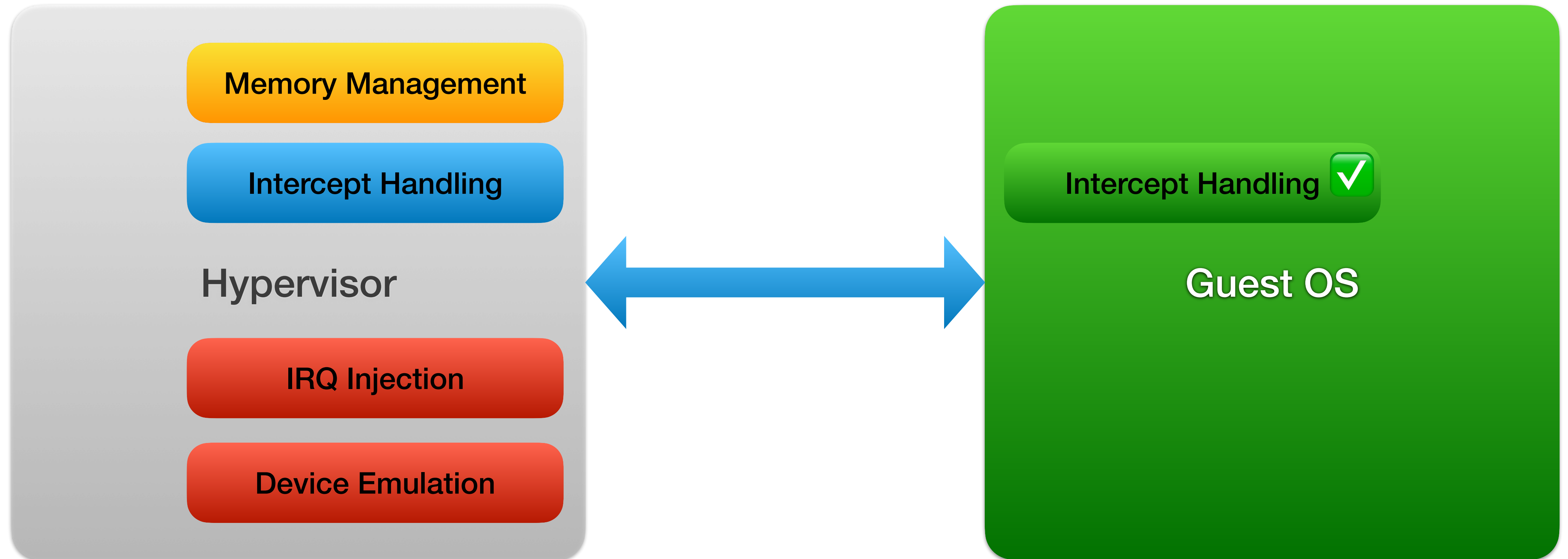
Page States in AMD SEV-SNP



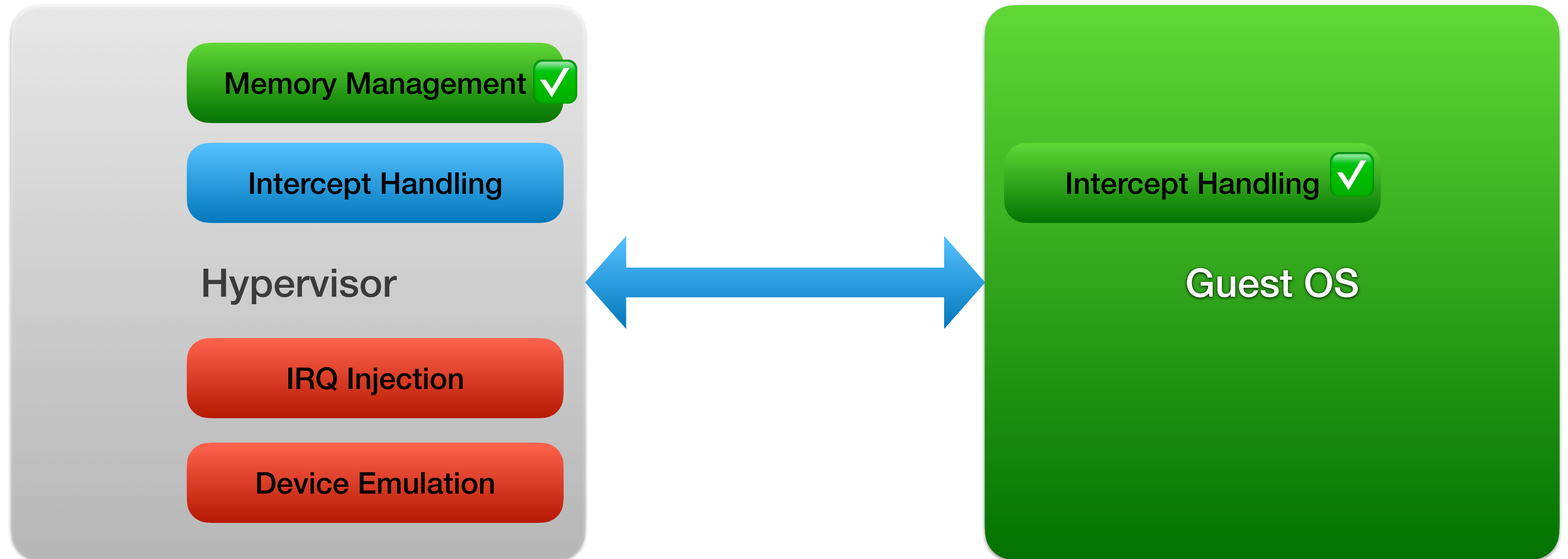
OS Implications

- OS needs to keep track of page states
 - Hypervisor vs. Guest-invalid vs. Guest-valid
- Keeping track allows to reliably detect malicious HV behavior
- OS also needs a new paravirtual interface to HV
 - Some page-state changes need to be coordinated with HV
 - OS needs to tell HV when pages switch between HV and Guest-Invalid

Hypervisor Guest Interface



Hypervisor Guest Interface



Interrupt Injection

Interrupt Injection

- HV can inject interrupts at any time
- Usually the HV tracks when the guest is ready for IRQs
 - IRQs enabled
 - No Interrupt shadow
- OSes disable IRQs when not able to handle them
- Malicious HV could inject IRQs while the guest is in a critical state
 - Could be used to change control flow in guest and reveal secrets

Case Study: AMD Secure Encrypted Virtualization

- Confidential Computing in stages - building on each other
 - Secure Encrypted Virtualization (SEV) - Guest memory encryption
 - SEV Encrypted State (SEV-ES) - Guest register encryption
 - SEV Secure Nested Paging (SEV-SNP) - Introducing page states
 - **SEV-SNP Secure Interrupt Injection - Guest is responsible for delivering IRQs**

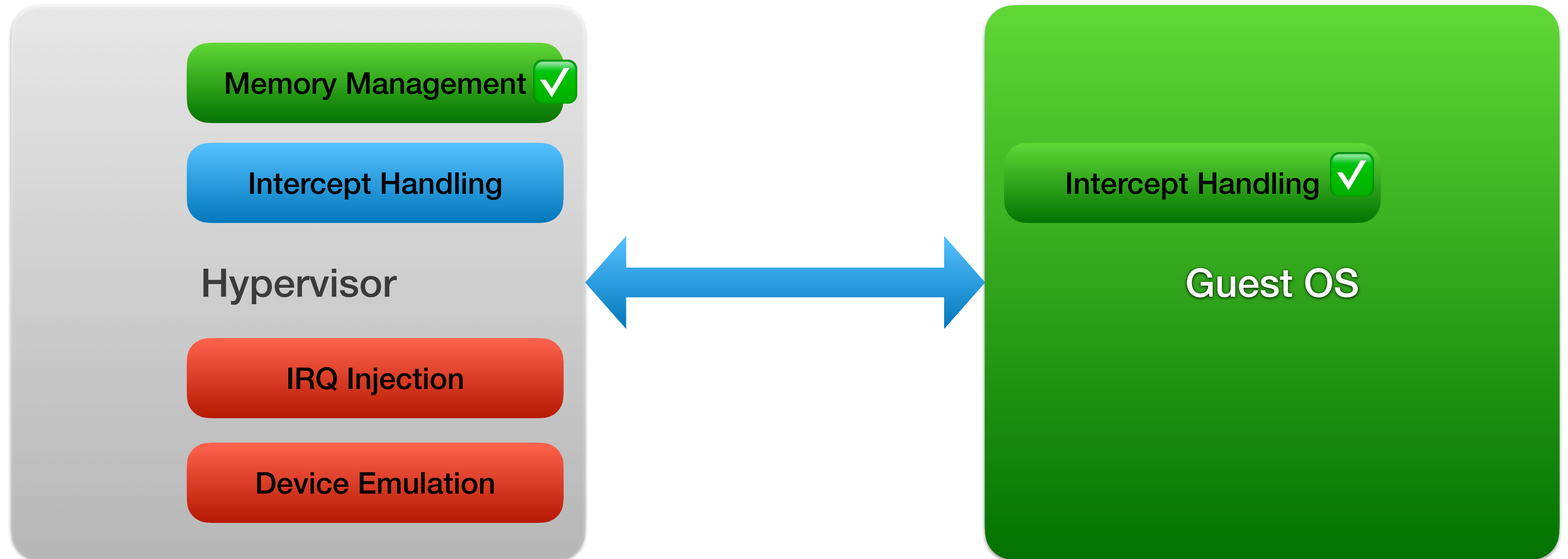
SEV-SNP Secure Interrupt Injection

- When enabled HV can only inject one event: #HV
- Guest OS needs to be prepared to receive this event at any time
 - Needs to be an IST vector
- PV protocol used to communicate which event is delivered
 - Also involves software blocking of new #HV events

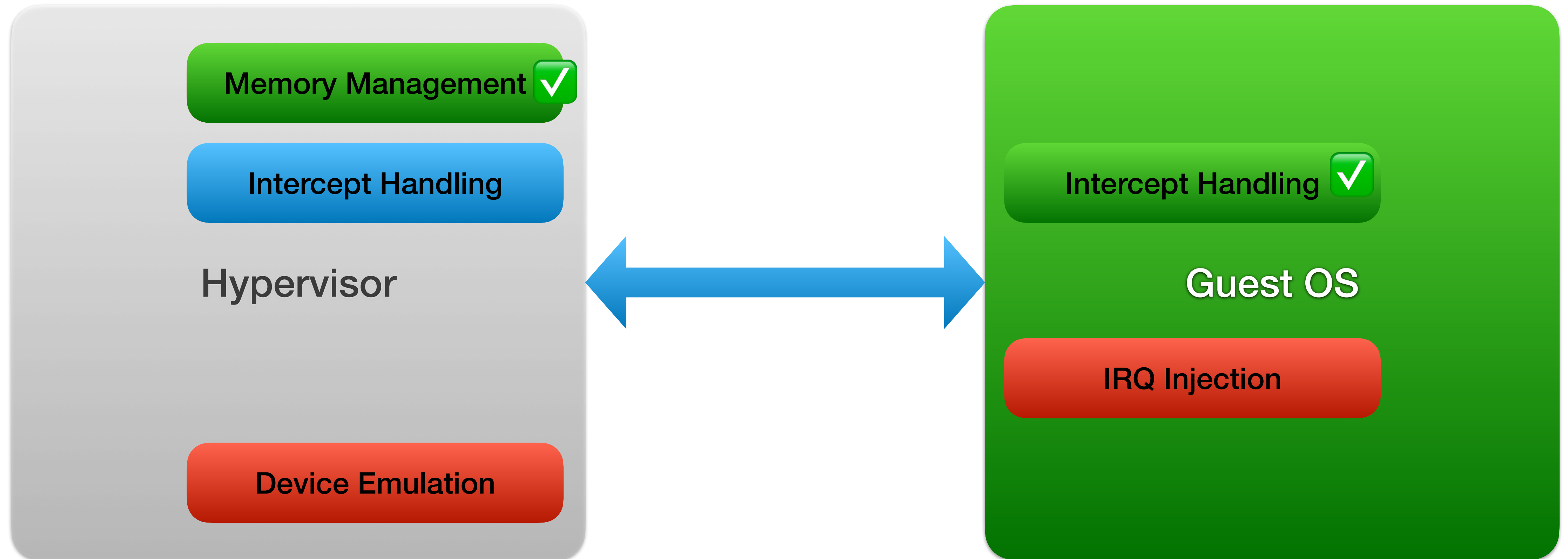
OS Implications

- OS needs to implement handler for #HV and harden it
 - Must be IST
 - Must support nesting and detection for malicious injection
- OS needs to track when IRQ handlers can run and deliver events itself
 - Instrumentation of IRQ enable/disable events
 - Manual IRQ handler launching with stack switching

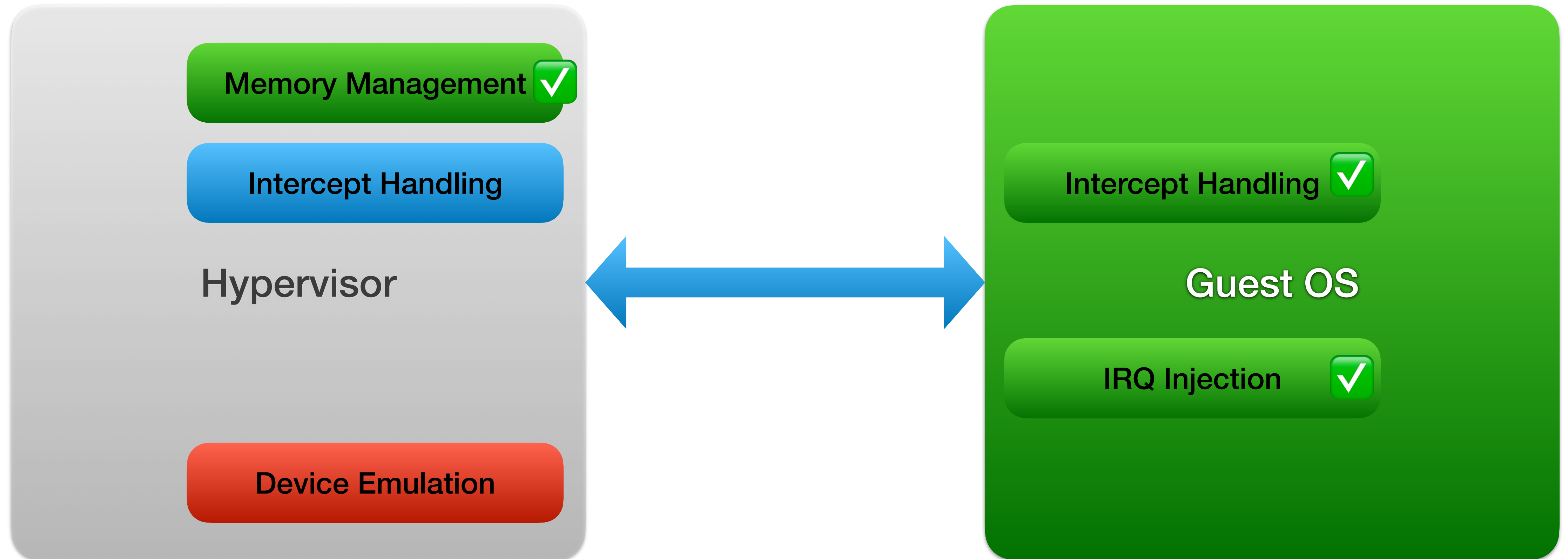
Hypervisor Guest Interface



Hypervisor Guest Interface



Hypervisor Guest Interface



Device Emulation

Device Emulation

- All HV-emulated are treated as insecure
- Some devices carry security sensitive state: TPM
- Two approaches:
 - Device driver hardening
 - In-guest Paravisor

Device Driver Hardening

- Often done via device driver fuzzing
- Implemented using HV-side fuzzers
 - Had some success in finding bugs in device drivers
 - Patches are sent to the Linux kernel for device driver hardening
- But, overall, a difficult approach which is never finished

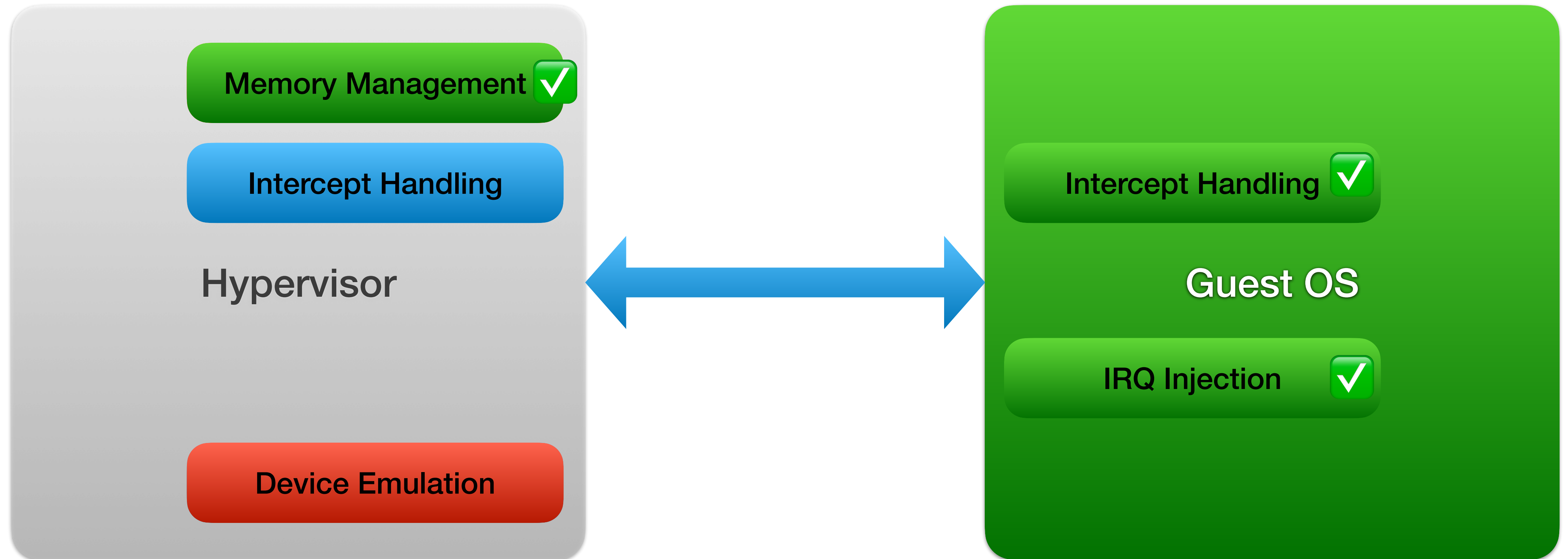
In-Guest Paravisor

- Uses hardware capabilities for isolation within TEE
 - AMD SEV-SNP VM Privilege Levels (VMPLs)
- VMPLs allow memory separation within an SEV-SNP guest
- 4 Levels - each with its own CPU state
- Use cases:
 - Make guest memory inaccessible to OS
 - Securely move HV functionality into TEE

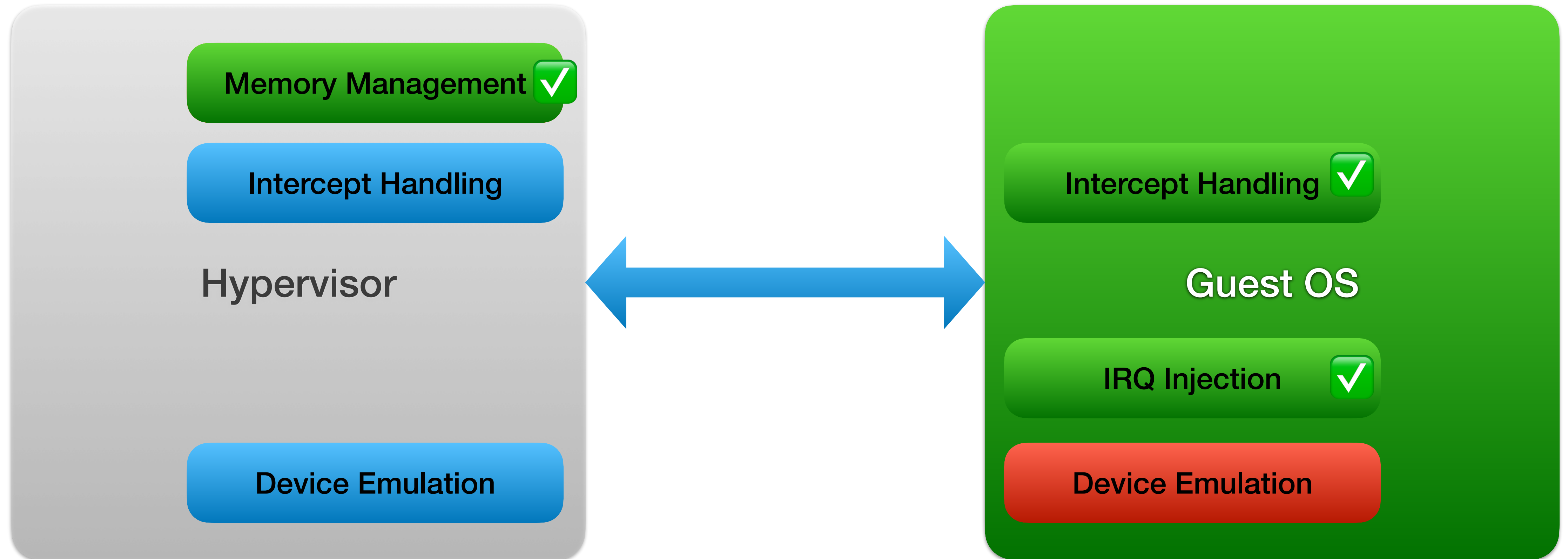
OS Implications

- OS needs to harden device drivers
 - Involves continuous fuzzing and code review
- For paravisor support OS needs implement PV calls to paravisor
 - Specific protocol using shared memory
 - Additional protocols for emulated devices

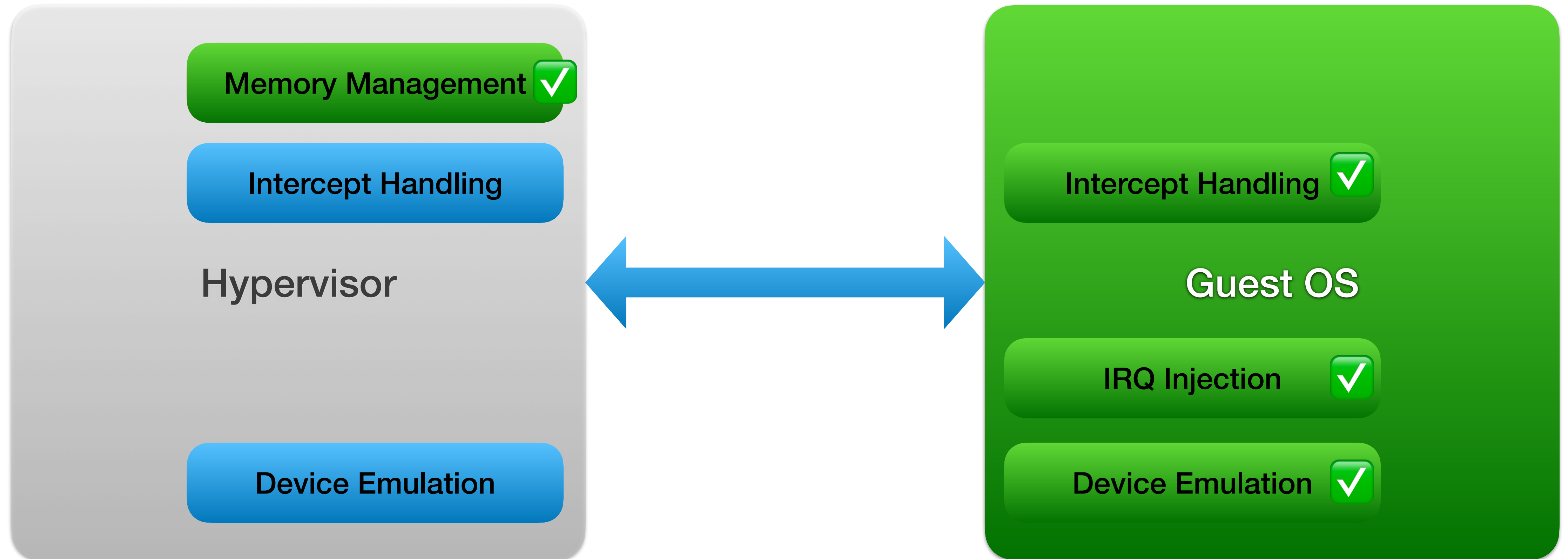
Hypervisor Guest Interface



Hypervisor Guest Interface



Hypervisor Guest Interface



COCONUT Secure VM Service Module

- COCONUT Secure VM Service Module (SVSM) currently under development
- OS-level project written in stable Rust
- Will support unprivileged execution mode (CPL3)
- First use-case: Secure TPM 2.0 emulation for CoCo guests
 - Needed for attestation
- Can emulate more devices in the future
- Allows to move some CoCo specifics from OS into Paravisor

COCONUT Secure VM Service Module

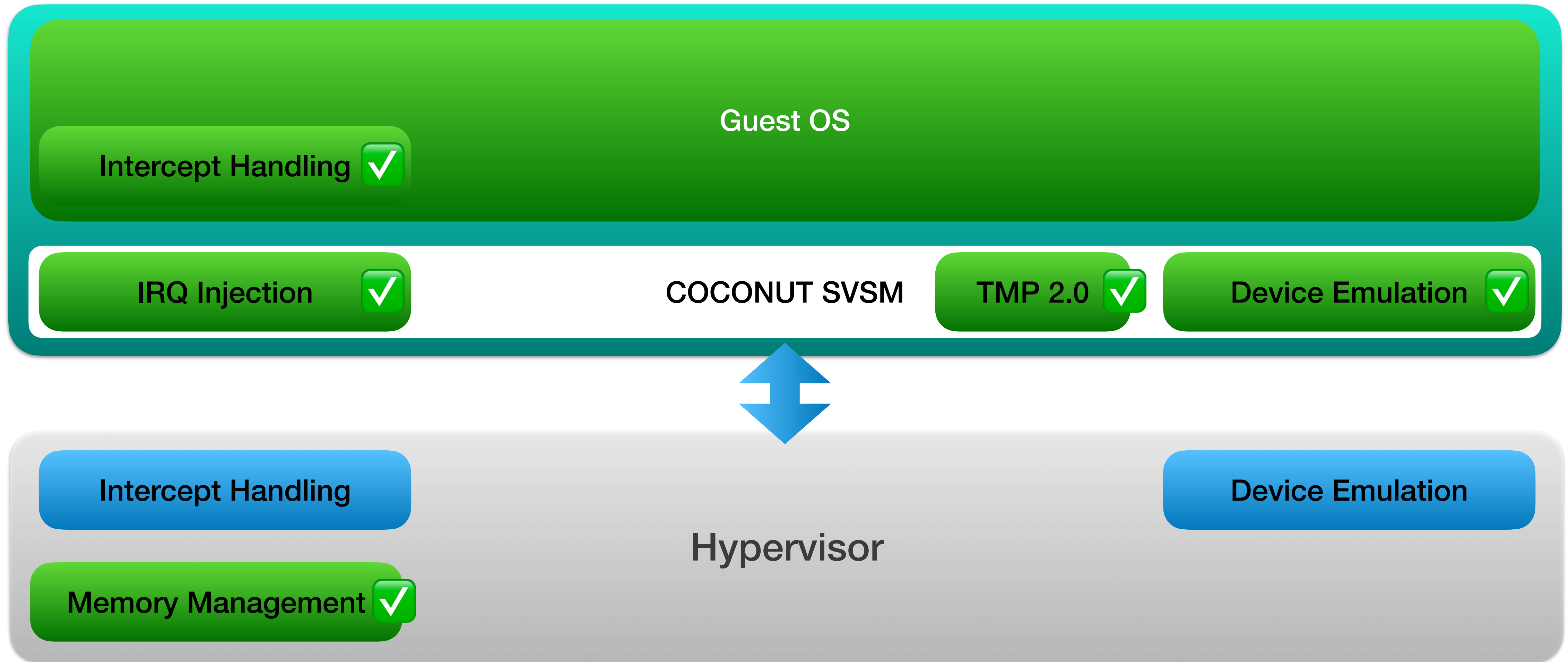
The screenshot shows the GitHub repository page for `coconut-svsm/svsm`. The repository is public and has 17 forks, 36 stars, and 11 watchers. The main branch is selected, and there are 1 branch and 0 tags. The repository contains several files and folders, including `.cargo`, `.github/workflows`, `scripts`, `src`, `stage1`, `utils`, `.gitignore`, `.mailmap`, `CONTRIBUTING.md`, `Cargo.lock`, `Cargo.toml`, and `INSTALL.md`. The repository is managed by `joergroedel` and has 682 commits. The repository is licensed under MIT and has no releases published.

File/Folder	Commit Message	Commit Date
<code>.cargo</code>	build: fix tests for stable Rust	last month
<code>.github/workflows</code>	Merge pull request #81 from 00xc/build/stablerust	last month
<code>scripts</code>	Merge pull request #84 from osteffenrh/container-build-with-podman	3 weeks ago
<code>src</code>	Merge pull request #98 from 00xc/sev/pvalidate	yesterday
<code>stage1</code>	stage1: Fix rounding error in kernel fs bin size	4 months ago
<code>utils</code>	Change License and update license headers	6 months ago
<code>.gitignore</code>	gitignore: add .idea directory	7 months ago
<code>.mailmap</code>	Add a .mailmap file	last year
<code>CONTRIBUTING.md</code>	docs: mention development mailing list	5 days ago
<code>Cargo.lock</code>	Update <code>bitflags</code> crate to 2.4.0	last month
<code>Cargo.toml</code>	Update <code>bitflags</code> crate to 2.4.0	last month
<code>INSTALL.md</code>	INSTALL.md: update references to build target	last month

<https://github.com/coconut-svsm/svsm/>



CoCo Virtualization Stack Vision



Thank You! Questions?