

Hardware-supported virtualization (AMD/Intel)

▶ Intel VT-x a VT-d

- ▶ Řada rozšíření CPU i podpůrného chipsetu k podpoře virtualizace
 - Neustále přibývají další
- ▶ Jednotlivé úpravy jsou často použitelné nezávisle
- ▶ Významné virtualizační softwary je využívají téměř všechny
 - Intel spolupracuje s producenty software

▶ AMD-V

- ▶ Úpravy ve stejném čase (2006) podobným směrem
 - Většina není kompatibilní s Intelem

▶ Situace znepráhledněna obchodní politikou

- ▶ Různé verze CPU mají různý stupeň podpory
- ▶ Obchodní názvy maskují podstatu věci
 - Některá rozšíření jsou triviality, jiná jsou velmi netriviální

- ▶ Root/non-root execution (2005)
 - ▶ Řešení problému komprese privilegií
- ▶ Extended Page Table (EPT) (2008)
 - ▶ Řešení problému virtualizace virtuální paměti
- ▶ VMCS Shadowing (2013)
 - ▶ Podpora rekurzivní virtualizace

- ▶ Odstraněna komprese privilegií
 - ▶ Podmínka: Virtualizované OS samy HW podporu virtualizace nevyužívají
 - VMCS Shadowing (Intel 2013): Rekurzivní virtualizace je možná
 - Na IBM VM/370 bylo demonstrováno 5 úrovní vnoření virtualizace
- ▶ Přepínání adresového prostoru při VM entry/exit
 - ▶ Ochrana paměti VMM, plná transparence pro virtualizovaný OS
 - Komplikuje přístup VMM do paměti VM (při emulaci I/O apod)
- ▶ Menší počet přechodů VM-VMM
 - ▶ Lze vyladit konfigurací HW kontroly přístupu k privilegovanému stavu
 - ▶ Demonstrováno cca. dvojnásobné zrychlení některých úloh
 - Unix fork and wait benchmark
 - Kompilace rozsáhlých projektů s malými moduly

- ▶ FlexPriority
 - ▶ Virtualizace klíčové části řadiče přerušení (APIC)
- ▶ Pause-loop exiting
 - ▶ Detekce spin-locků způsobující exit do VMM
 - Pro provoz více virtuálních procesorů na méně fyzických
- ▶ VGuest Preemption Timer
 - ▶ Časovač s lepší granularitou a rychlejší obsluhou
 - Pro virtualizaci aplikací s mírnými real-time nároky
- ▶ FlexMigration
 - ▶ Virtualizace identifikace CPU a jeho schopností
- ▶ Virtual Processor ID (VPID)
 - ▶ Klíč záznamu v TLB obsahuje identifikátor VM
 - Není třeba invalidovat celou TLB při přepínání VM-VMM a VM-VM
- ▶ Real-mode support
 - ▶ Podpora virtualizace při startu virtualizovaného OS

▶ IOMMU

- ▶ I/O zařízení přistupují k paměti přes MMU podobně jako CPU
- ▶ Address Translation Services (ATS) support
 - Rozšíření standardu sběrnice PCI Express
- ▶ Large Intel VT-d Pages
 - Umožňuje sdílení CPU a DMA verzí stránkových tabulek

▶ Interrupt-remapping support

- ▶ Částečná virtualizace řadiče přerušení

▶ Virtual Machine Device Queue

- ▶ Network Interface Card s více stavovými prostory pro přímý přístup z VM

▶ Single-Root I/O Virtualization (SR-IOV)

- ▶ I/O zařízení deklarují své schopnosti virtualizace
- ▶ Rozšíření standardu PCI Express

▶ Graphics Virtualization Technology

- ▶ Využití výpočetní síly (Intel) GPU ve virtuálních strojích
 - Exkluzivní přístup (GVT-D)
 - Sdílený přístup (GVT-S) – vyžaduje přizusobení ovladačů ve virtuálních strojích
 - Časový multiplex (GVT-G)

▶ Data Direct I/O Technology (DDIO)

- ▶ Zpřístupnění CPU cache pro DMA – snížení latence síťové komunikace

▶ Virtualization Extensions to the x86 Instruction Set

- ▶ Enables software to more efficiently create VMs so that multiple operating systems and their applications can run simultaneously on the same computer

▶ Tagged TLB

- ▶ Hardware features that facilitate efficient switching between VMs for better application responsiveness

▶ Rapid Virtualization Indexing (RVI)

- ▶ Helps accelerate the performance of many virtualized applications by enabling hardware-based VM memory management

▶ AMD-V Extended Migration

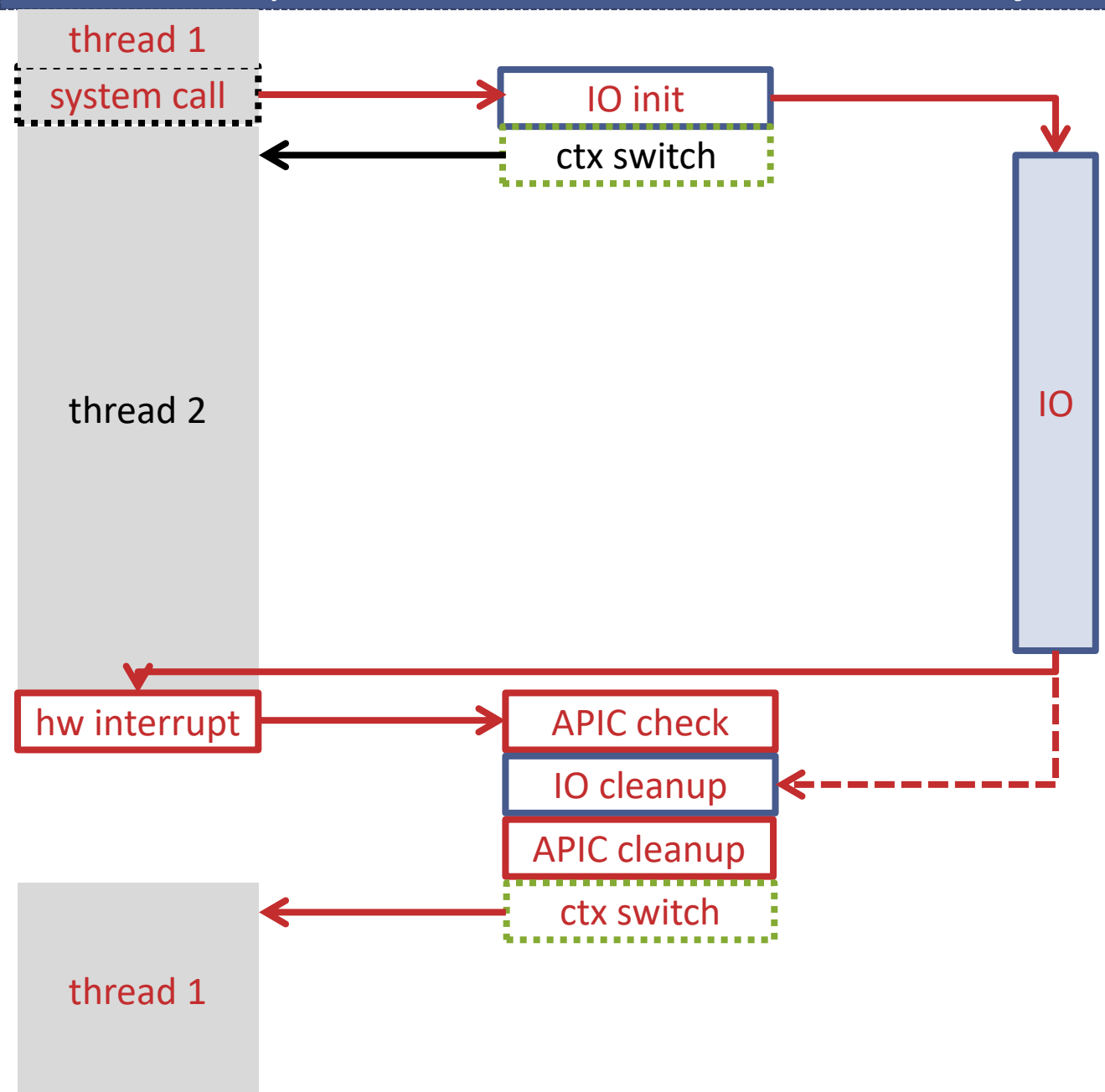
- ▶ Helps virtualization software with live migrations of VMs between all available AMD Opteron processor generations

▶ I/O Virtualization

- ▶ Enables direct device access by a VM, bypassing the hypervisor for improved application performance and improved isolation of VMs for increased integrity and security

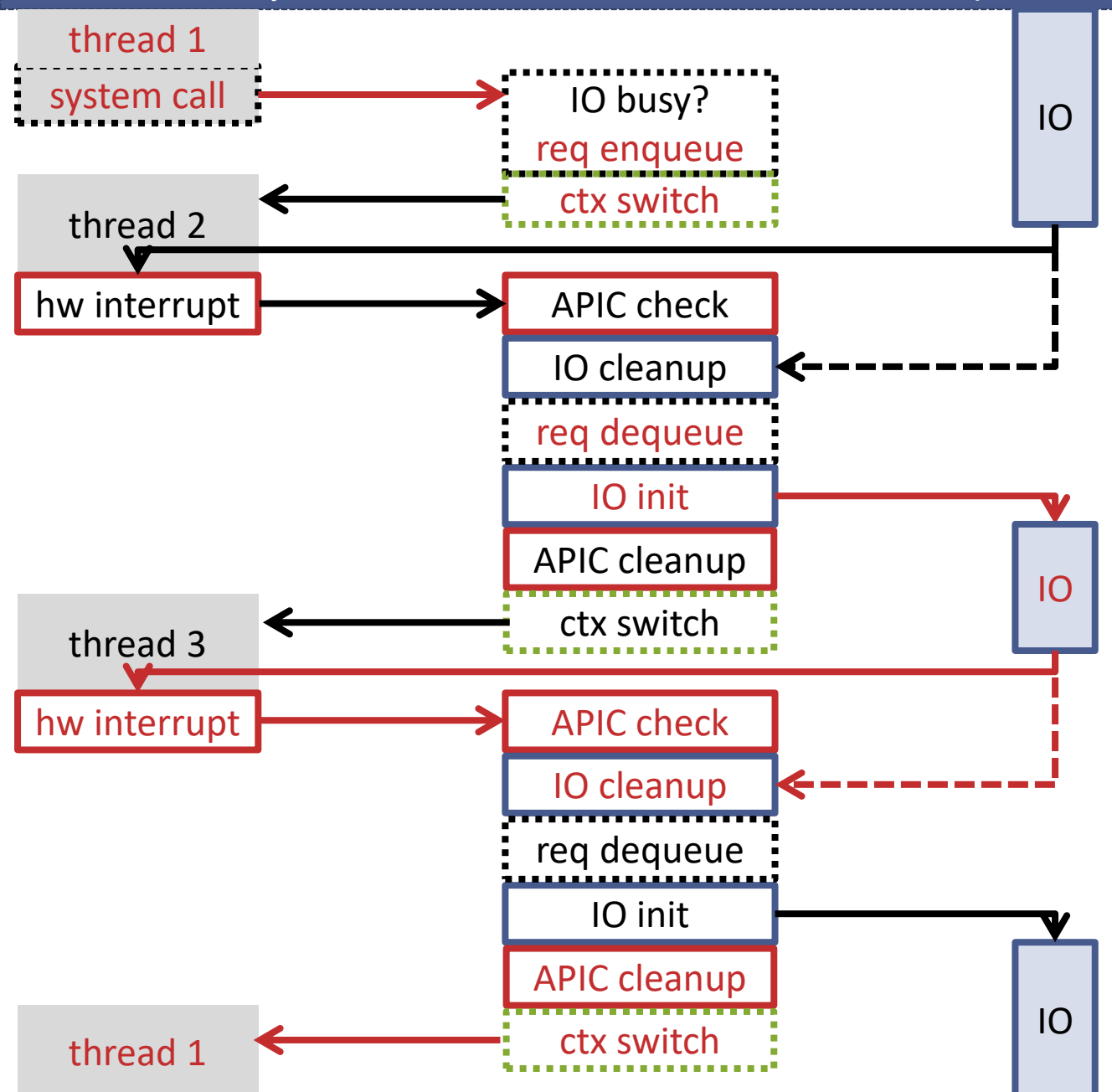
Virtualizace I/O

Obsluha IO požadavku v OS bez virtualizace - zjednodušeno



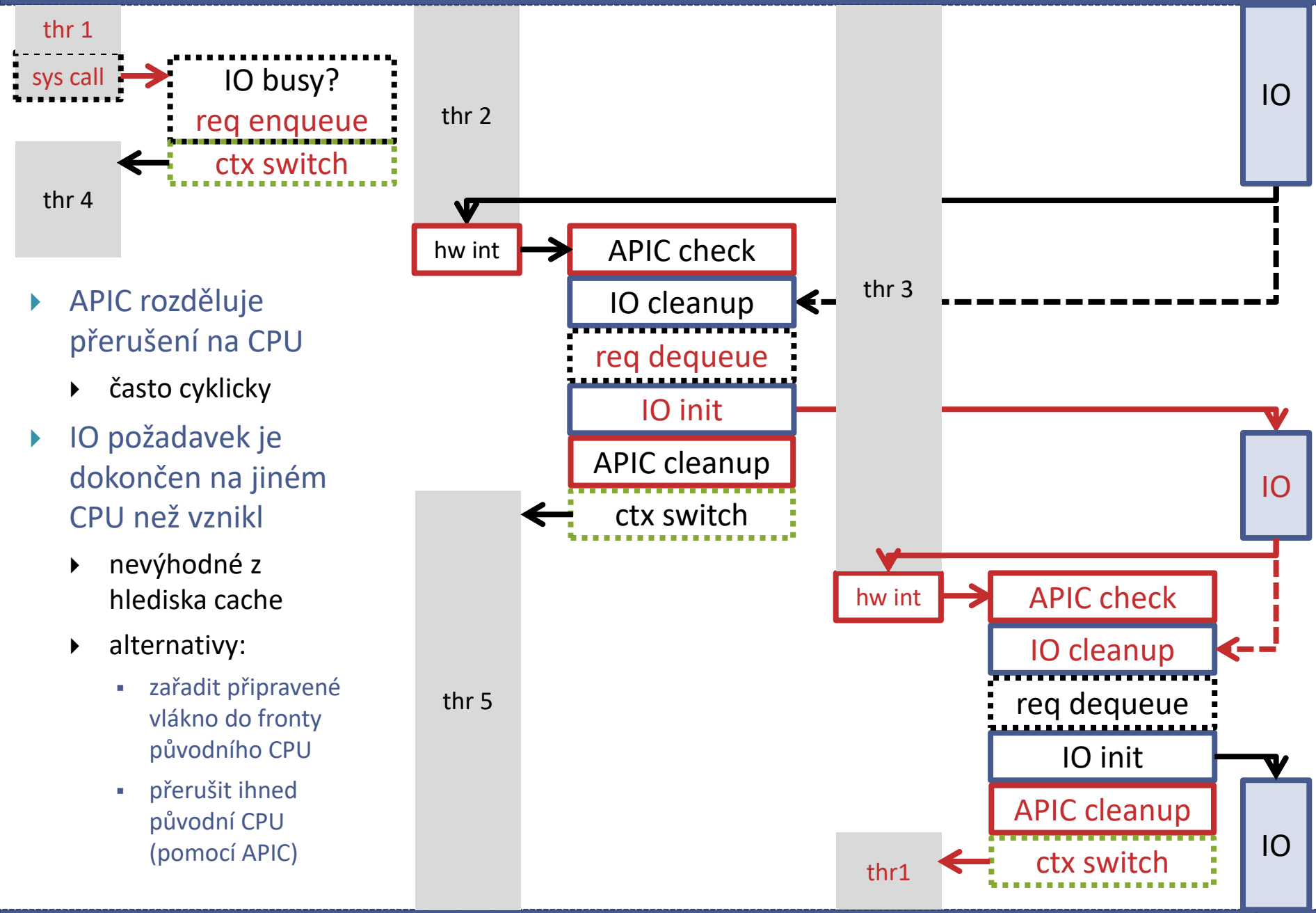
- ▶ Aplikační proces volá jádro
- ▶ Jádro nastartuje IO operaci
 - ▶ I/O instrukce
- ▶ Po dobu čekání běží jiné aplikační vlákno
- ▶ Dokončení operace je signalizováno přerušením
 - ▶ To putuje od IO zařízení přes APIC
 - (terminologie x86: advanced programmable interrupt controller)
- ▶ Obsluha přerušení začíná zkoumáním důvodu
 - ▶ APIC sdružuje různé zdroje přerušení
- ▶ Jádro testuje úspěšnost operace
 - ▶ I/O instrukce
 - ▶ Samotný výsledek I/O bývá v paměti
- ▶ Ukončení obsluhy přerušení se hlásí APICu
 - ▶ Ochrana před rekurzí

Obsluha IO požadavku v OS bez virtualizace (1 CPU)



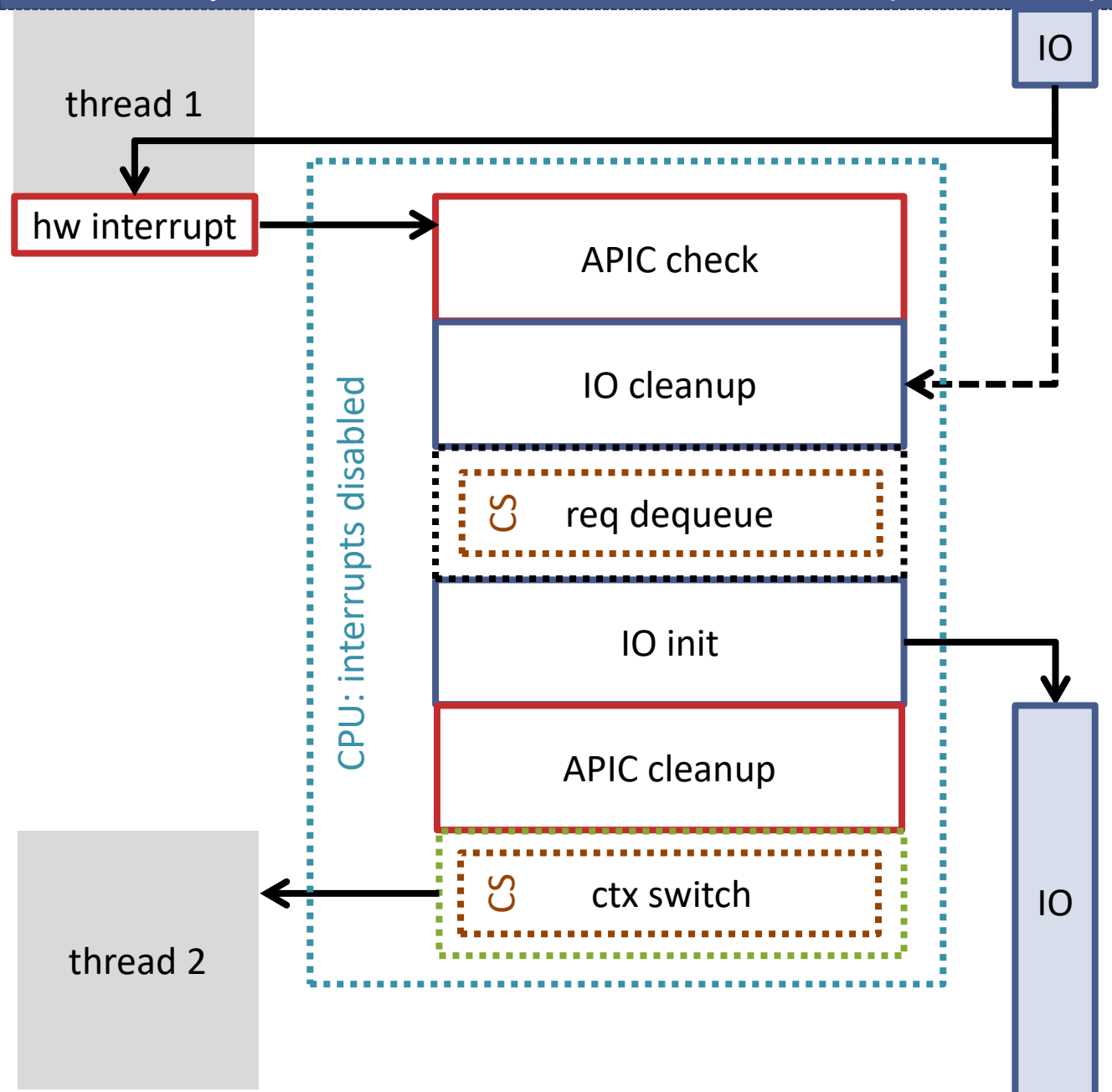
- ▶ V reálném případě je IO zařízení často obsazeno
 - ▶ Požadavky čekají ve frontě organizované jádrem OS
 - ▶ Obsluha přerušení typicky dokončuje starý požadavek a startuje nový
- ▶ Zdrojů přerušení je víc než signálů, které vedou k CPU
 - ▶ Obsluha jednoho přerušení občas řeší více zdrojů přerušení najednou

Obsluha IO požadavku v OS bez virtualizace (více CPU)



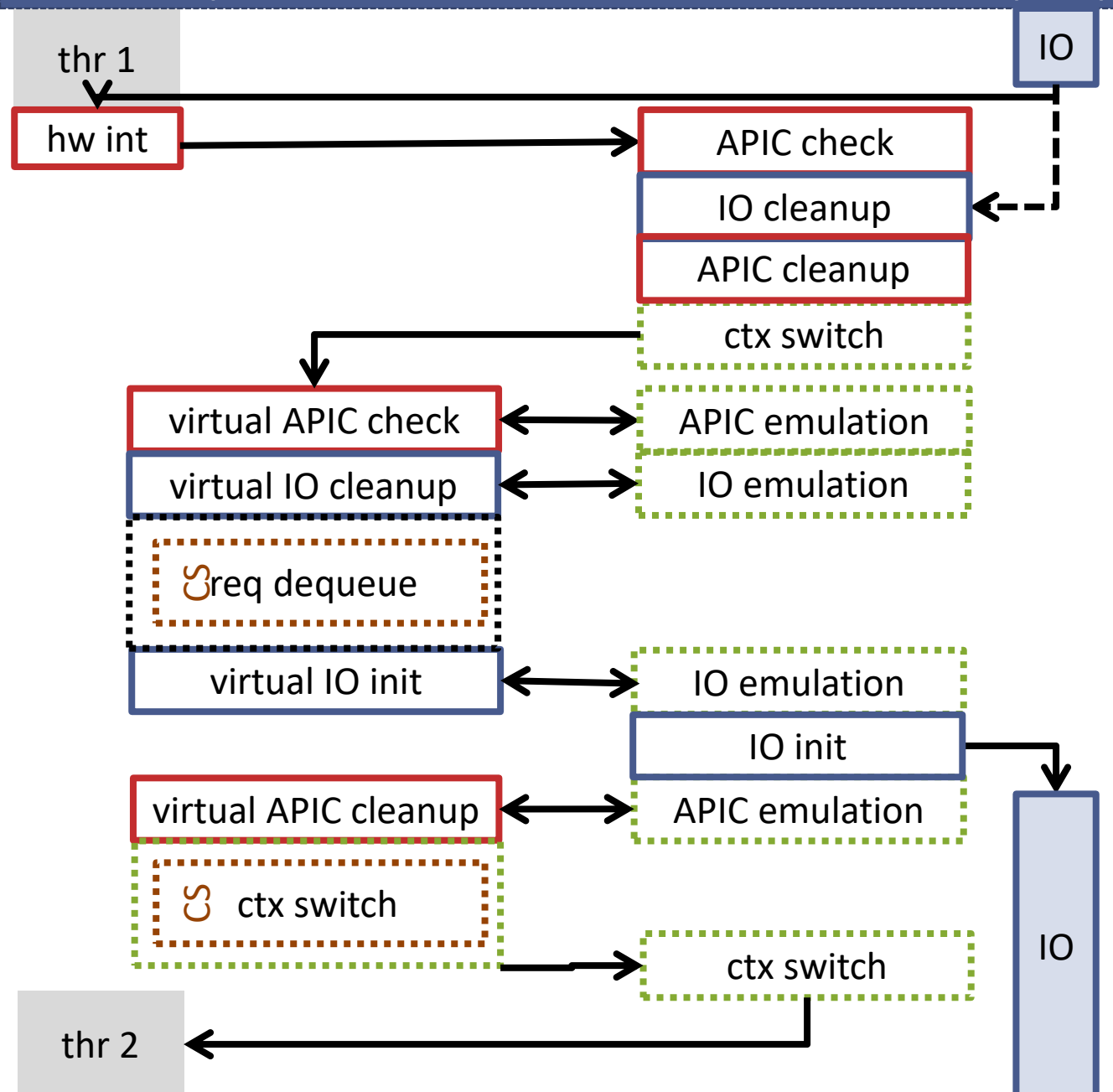
- ▶ APIC rozděluje přerušení na CPU
 - ▶ často cyklicky
- ▶ IO požadavek je dokončen na jiném CPU než vznikl
 - ▶ nevýhodné z hlediska cache
 - ▶ alternativy:
 - zařadit připravené vlákno do fronty původního CPU
 - přerušit ihned původní CPU (pomocí APIC)

Obsluha přerušení v OS bez virtualizace (více CPU)



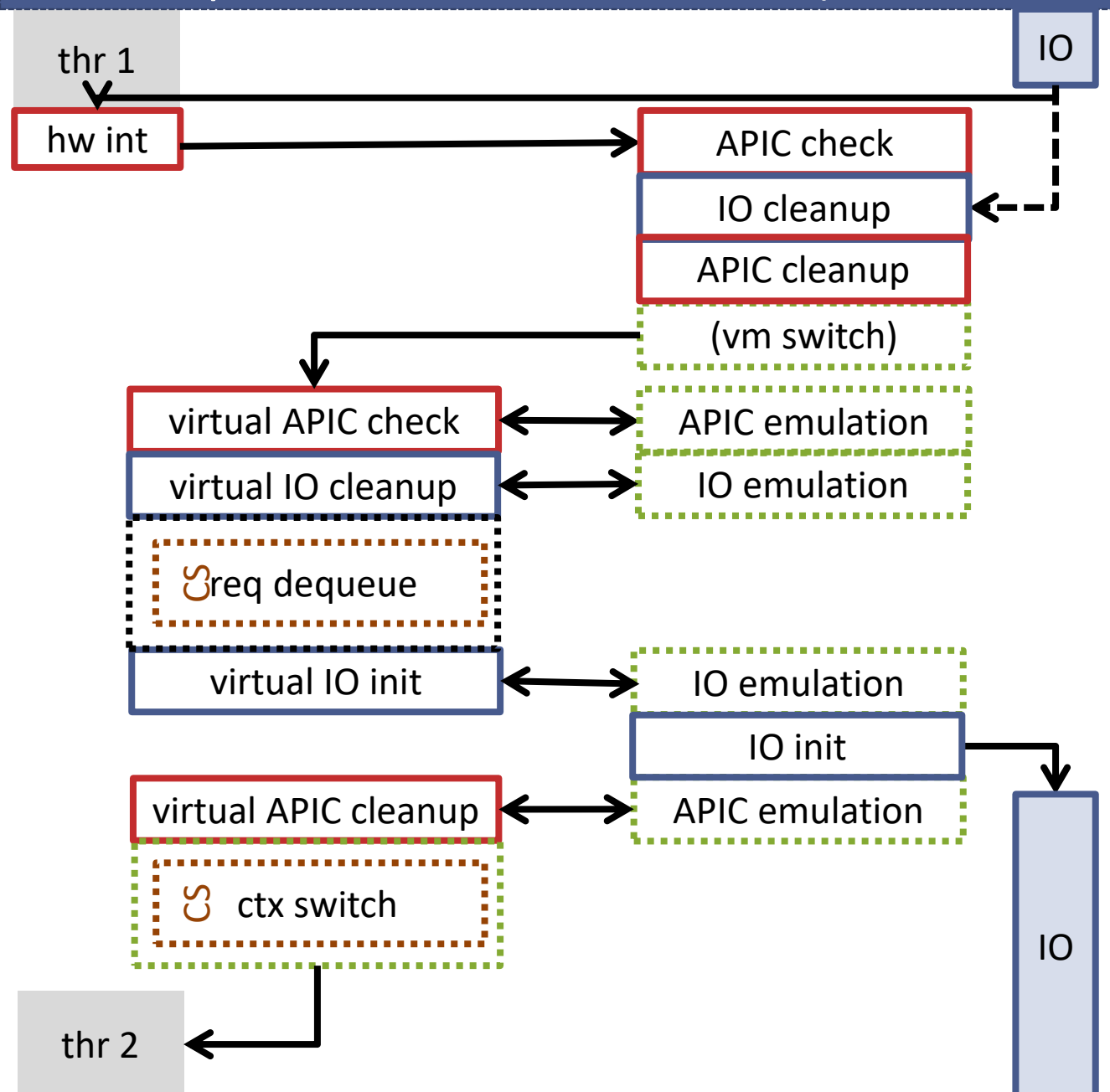
- ▶ CPU obvykle automaticky zakáže vnější přerušení při vstupu do jeho obsluhy
 - ▶ Na jiných CPU ale přerušení zakázána nejsou
- ▶ Datové struktury jádra musejí být chráněny
 - ▶ CS - Spinlock

Obsluha přerušení v OS s virtualizací (bez HW podpory virtualizace)



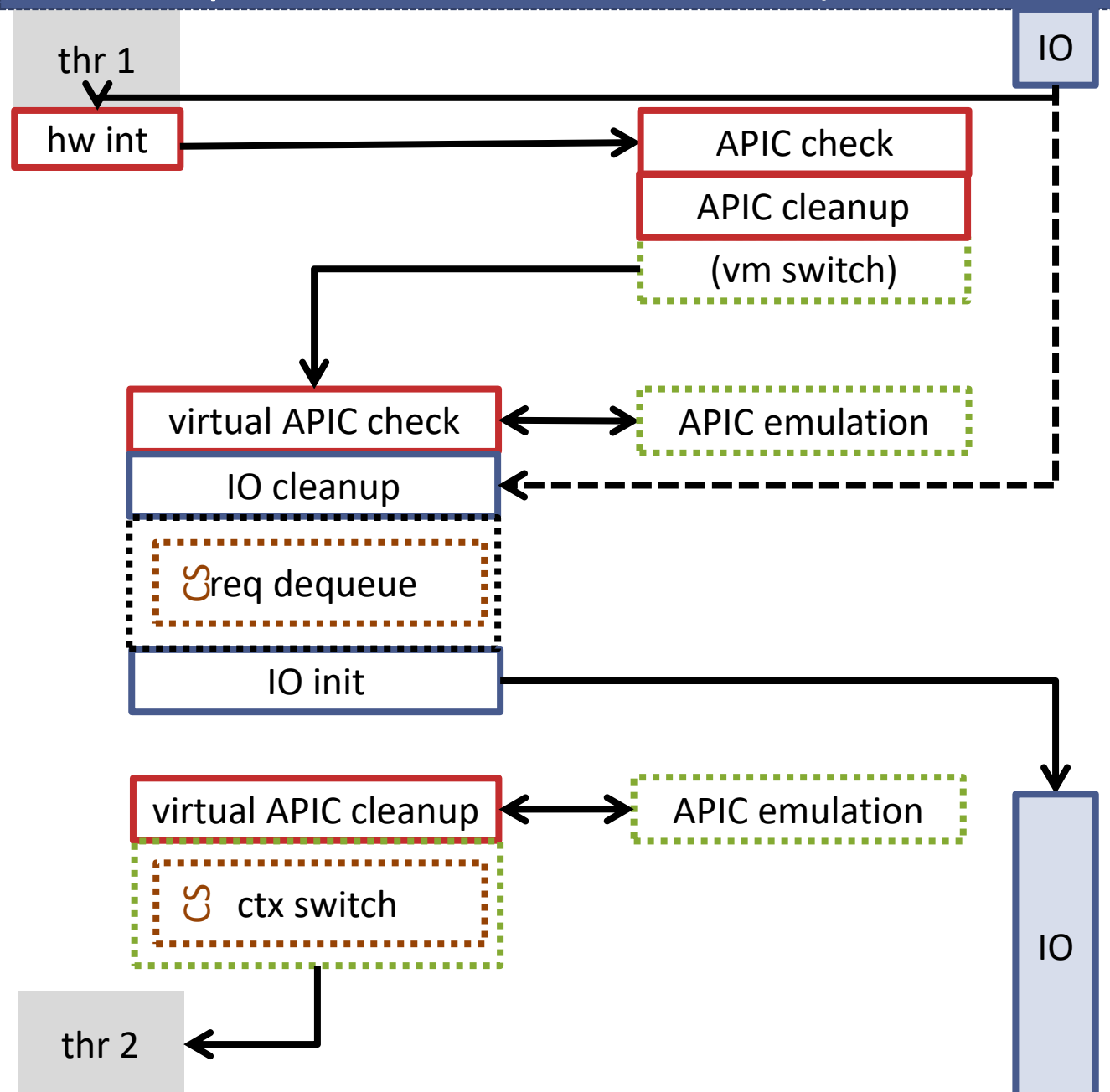
- ▶ Jádro OS běží v aplikačním režimu
 - ▶ Přerušení jsou fyzicky povolena
 - ▶ VMM je musí odložit na příпустný okamžik
- ▶ Interakce s IO a APIC musí být virtualizována
- ▶ Kritické sekce mohou být přerušeny fyzickým přerušením a přeplánovány
 - ▶ Trvají nepřipustně dlouho
 - ▶ Ostatní virtuální CPU jsou zablokovány ve spinlocku - aktivně!

Obsluha přerušení v OS s virtualizací (s root/non-root režimy)



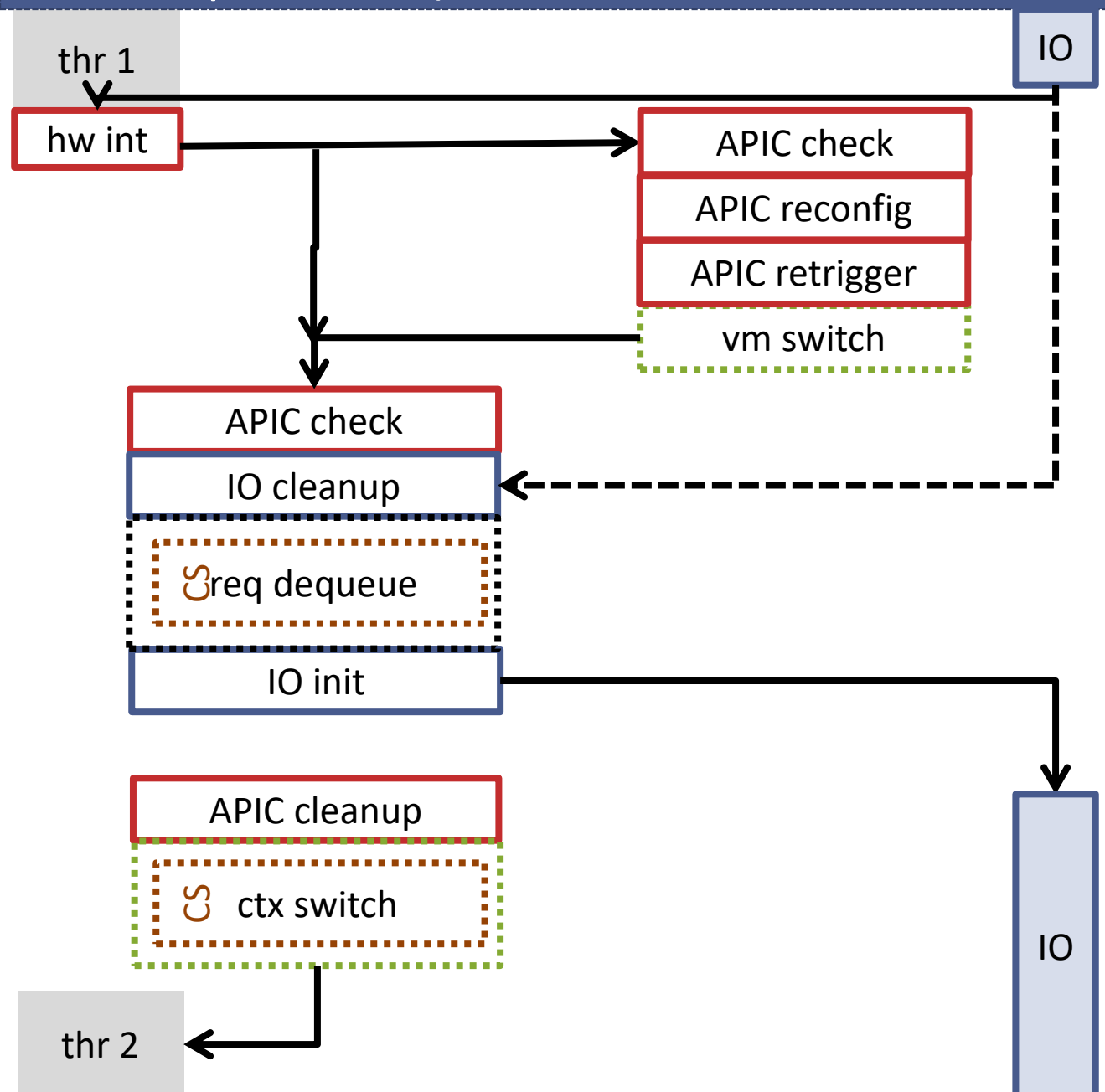
- ▶ Jádro OS běží v non-root režimu
 - ▶ Přerušení jsou fyzicky povolena, CPU je řeší jako VM-Exit
 - ▶ VMM je musí poslat správnému VM
- ▶ Interakce s IO a APIC musí být virtualizována

Obsluha přerušení v OS s virtualizací (s root/non-root a VMDQ)



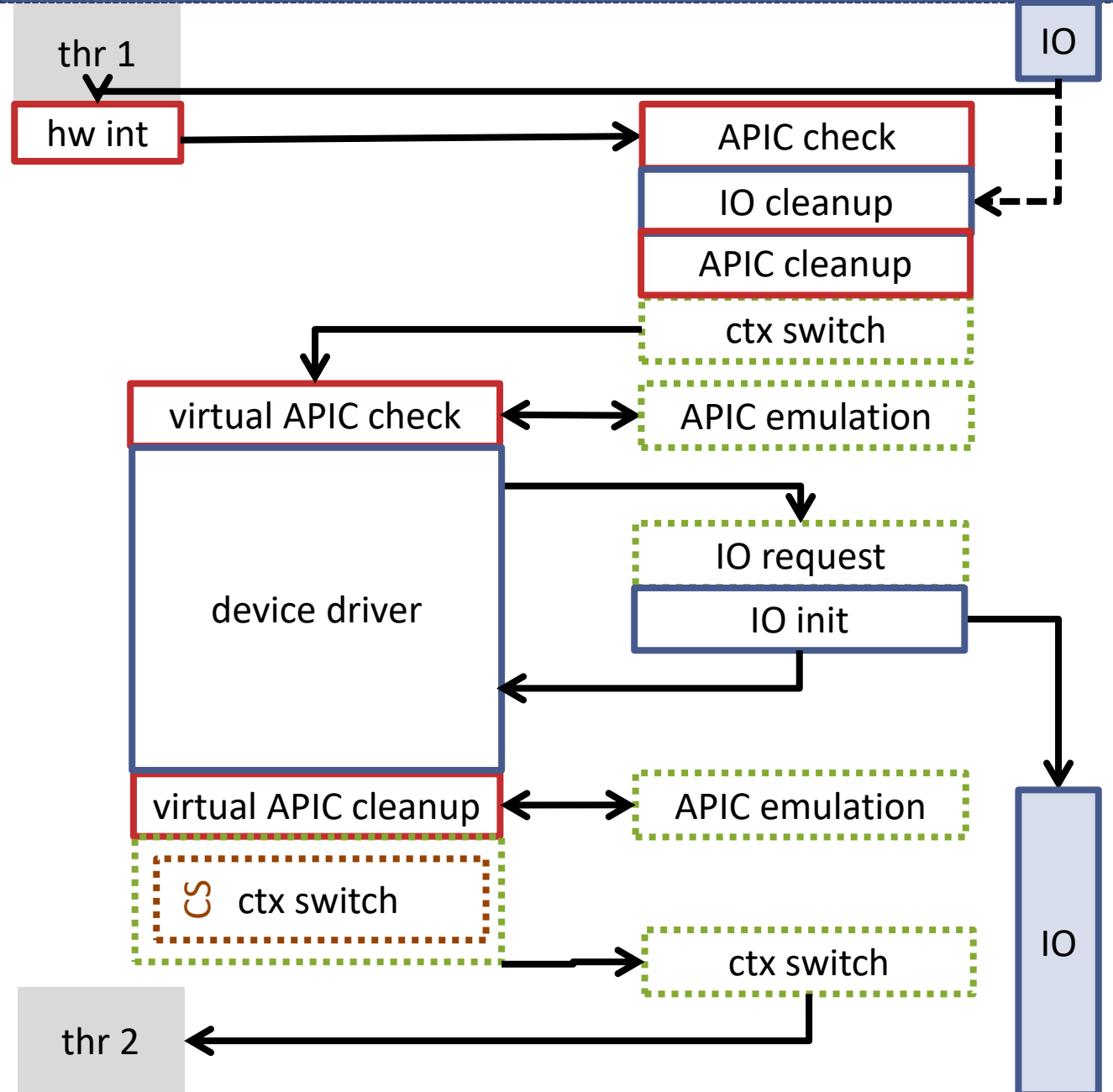
- ▶ Jádro OS běží v non-root režimu
 - ▶ Přerušení jsou fyzicky povolena, CPU je řeší jako VM-Exit
 - ▶ VMM je musí poslat správnému VM
- ▶ IO je přímo zpřístupněno v exklusivním režimu (VMDQ)
- ▶ APIC stále musí být virtualizován

Obsluha přerušení (s root/non-root, VMDQ a FlexPriority)



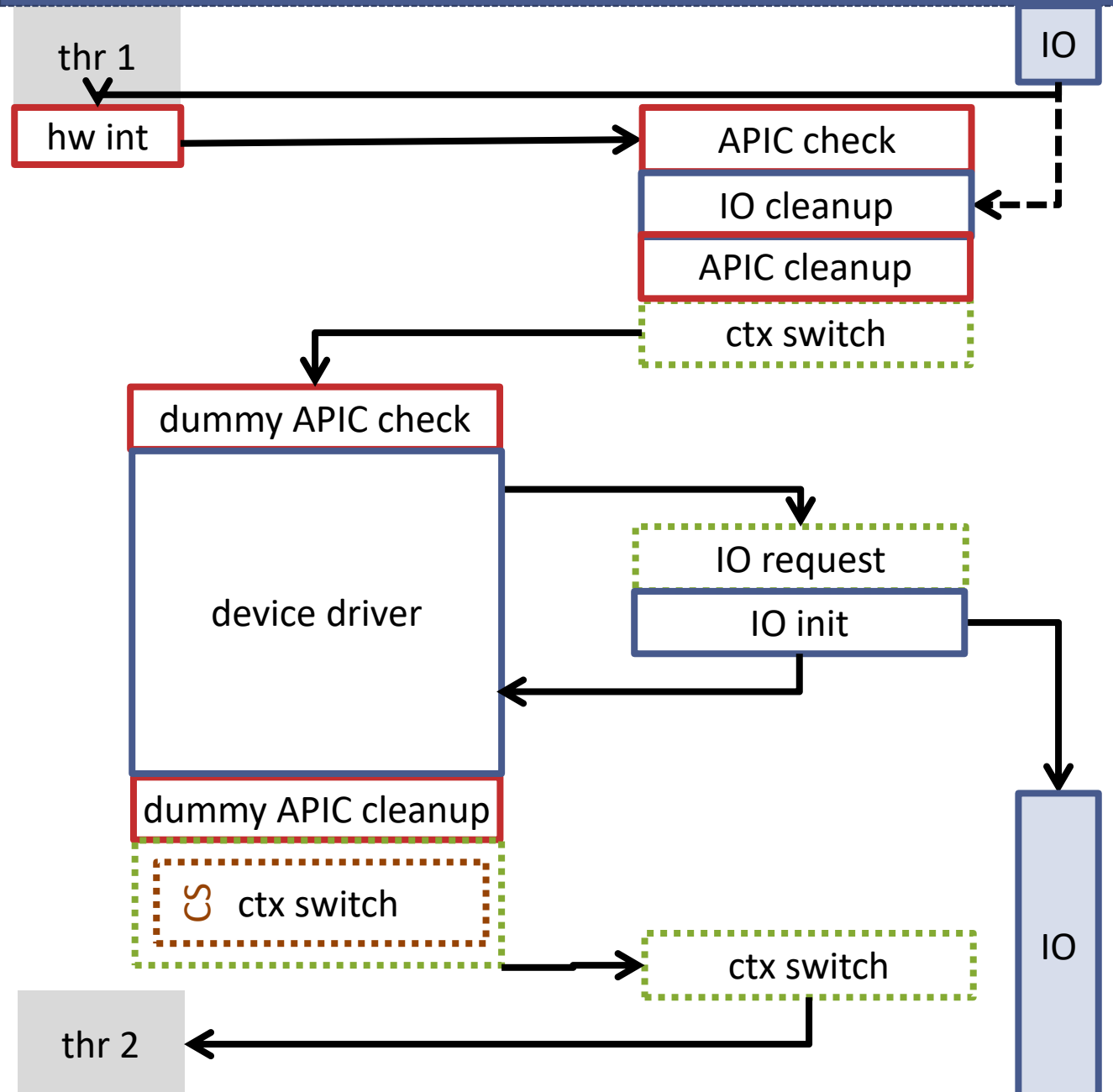
- ▶ IO je přímo zpřístupněno v exkluzivním režimu (VMDQ)
- ▶ Virtualizace APIC je řešena hardwarově
 - ▶ Přerušení od zařízení obsluhovaných exkluzivně právě běžící VM jsou v APIC/CPU konfigurovány tak, že nevyvolávají VM exit – obsluhuje je jádro OS
 - ▶ Ostatní přerušení VM exit vyvolávají
 - ▶ Při přepnutí VM musí VMM rekonfigurovat APIC

Obsluha přerušení v OS s ovladačem virtuálního zařízení



- ▶ Instalace speciálního ovladače virtuálního zařízení do OS
 - ▶ Vyřeší emulaci samotného IO zařízení
 - ▶ Nevyřeší emulaci APIC a režii přepínání kontextu
 - Tyto části bývají společné všem zařízením a nelze je vyměnit bez zásahu do OS
 - ▶ Problém spinlocku zůstává

Obsluha přerušení v OS s ovladačem virtuálního zařízení a FlexPriority



- ▶ Instalace speciálního ovladače virtuálního zařízení do OS
 - ▶ Vyřeší emulaci samotného IO zařízení
 - ▶ Emulaci fiktivního APIC zařídí přímo HW

- ▶ Každé asynchronní (I/O) i synchronní (syscall) přerušení může způsobit oživení vlákna - co s ním?
 - ▶ Přerušené vlákno může také zůstat živé
 - ▶ A) přerušené vlákno má přednost před oživeným
 - Počítající vlákna prakticky znemožní běh komunikujících
 - 1 obrátka komunikace za časové kvantum preemptivního plánovače
 - ▶ B) oživené vlákno dostane přednost před přerušeným
 - Vlákna intenzivně provádějící I/O potlačí běh nekomunikujících
 - ▶ C) něco mezi tím
 - Unix: Dynamické priority
- ▶ Jádro OS předpokládá známý počet stále běžících CPU
 - ▶ To při virtualizaci neplatí
 - ▶ Iluze většího počtu CPU je vytvářena preemptivním plánovačem ve VMM
 - Perioda plánovače rozhoduje o efektivitě komunikujících virtuálních CPU

Starší metody virtualizace

Virtualizace CPU – starší metody

- ▶ Rozhraní software-hardware v klasickém CPU (s dvěma módy)
 - ▶ Neprivilegovaný mód (aplikační procesy)
 - Instrukční sada (neprivilegované instrukce)
 - Stav CPU
 - Skutečný obsah neprivilegovaných registrů
 - Efekt stavu privilegovaných registrů
 - Paměťový prostor procesu
 - R/W operace
 - Manipulace prostřednictvím služeb OS
 - Mechanismus volání jádra OS
 - Speciální instrukce nebo (úmyslná) chyba
 - ▶ Privilegovaný mód (jádro OS) - navíc
 - Privilegované instrukce
 - Stav CPU
 - Skutečný obsah privilegovaných registrů
 - Paměťový prostor procesu
 - Manipulace s HW částí stránkování (TLB)
 - I/O operace
 - Obsluha přerušení synchronních i asynchroních

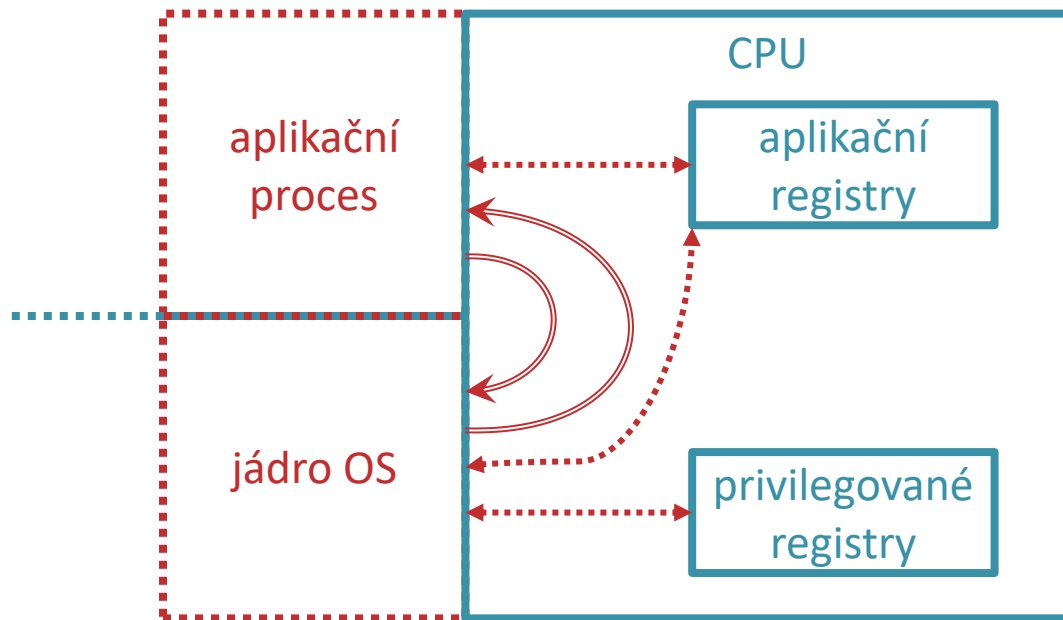
- ▶ Cíl: implementovat "jinak" rozhraní software-hardware
 - ▶ Nepřiznat sdílení fyzického CPU s jinými virtuálními stroji
 - ▶ Situace v procesorech bez HW podpory - s dvěma módy
- ▶ Neprivilegovaný mód (aplikační procesy)
 - ▶ Instrukční sada (neprivilegované instrukce)
 - Implementováno bez úprav fyzickým CPU - jinak to kvůli výkonu nelze
 - Fyzické CPU je občas odebráno preemptivním multitaskingem hypervizoru
 - ▶ Stav CPU
 - Skutečný obsah neprivilegovaných registrů
 - Preemptivní multitasking hypervizoru vyměňuje obsah registrů
 - Efekt stavu privilegovaných registrů
 - Fyzický stav je nutně jiný, ale odchylka nesmí být viditelná
 - ▶ Paměťový prostor procesu
 - Mechanismus virtuální paměti musí kombinovat sdílení paměti mezi VM a sdílení mezi procesy jednoho VM
 - ▶ Mechanismus volání jádra OS
 - Nelze přímo připustit původní semantiku (přepnutí do privilegovaného módu a skok do OS)

- ▶ Cíl: implementovat "jinak" rozhraní software-hardware
 - ▶ Nepřiznat sdílení fyzického CPU s jinými virtuálními stroji
- ▶ Situace v procesorech bez HW podpory - s dvěma módy
- ▶ Privilegovaný mód (jádro OS) - navíc
 - ▶ Privilegované instrukce
 - Většinou nelze připustit přímé provedení fyzickým CPU
 - ▶ Stav CPU
 - Fyzický obsah většiny privilegovaných registrů nelze přiznat
 - ▶ Paměťový prostor procesu
 - Manipulace s HW částí stránkování (TLB) musí být emulována
 - Nutnost zaznamenat i manipulaci prováděnou neprivilegovanými instrukcemi
 - ▶ I/O operace
 - U sdílených zařízení nelze připustit přímé provedení fyzickým CPU
 - ▶ Obsluha přerušení synchronních i asynchronních
 - Přerušení obvykle přepíná CPU do privilegovaného režimu
 - Začátek obsluhy přerušení musí být ve VMM

- ▶ Aplikační kód je vždy vykonáván fyzickým CPU
- ▶ Kód jádra OS lze vykonávat různými způsoby

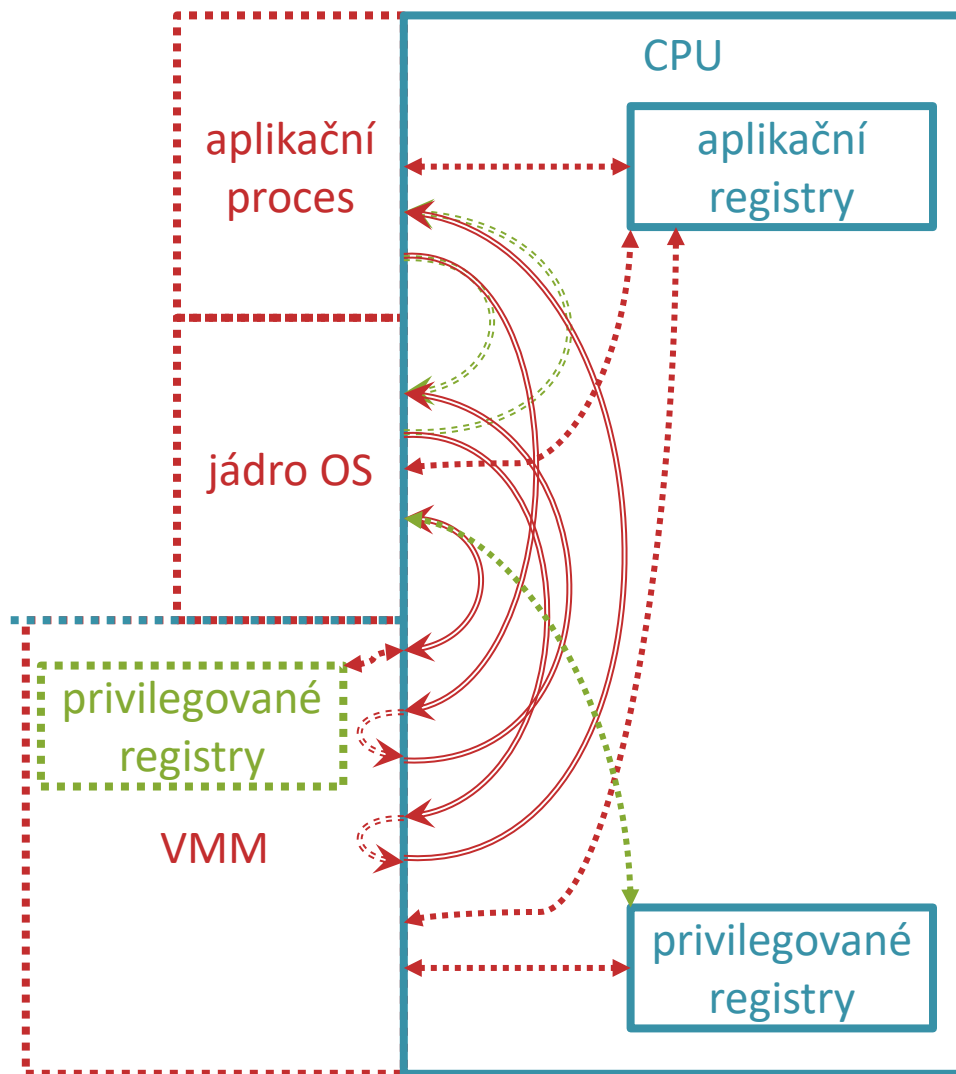
- ▶ Situace v procesorech bez HW podpory
 - ▶ Trap and Emulate - Jádro OS vykonáváno fyzickým CPU běžícím v aplikačním režimu [první éra virtualizace]
 - Privilegované instrukce způsobí výjimku a jsou emulovány
 - Časová penalizace každého přechodu aplikace-jádro
 - Neprivilegované instrukce musejí běžet identicky s privilegovaným režimem
 - Nepřiznat změněný stav procesoru neprivilegovanou instrukcí
 - ▶ Binary translation - přeložení do "bezpečného" kódu [VMWare x86]
 - Instrukce manipulující s privilegovanou částí stavu jsou upraveny
 - Vyžaduje čas a prostor pro překlad (on demand)
 - Bezpečný kód vykonáván fyzickým CPU [v privilegovaném režimu ?]
 - Šetří čas na přepínání režimu CPU
- ▶ Hardwarová podpora virtualizace [VMWare x64]
 - ▶ Více úrovní privilegovanosti ve stavu procesoru
 - Na stroji s N úrovněmi lze emulovat virtuální stroj s N-1 úrovněmi
 - ▶ Oddělené (stínové) registry pro fyzický a virtuální stav CPU

- ▶ Aplikační kód je vždy vykonáván fyzickým CPU
- ▶ Kód jádra OS lze vykonávat různými způsoby
 - ▶ Fyzickým CPU běžícím v aplikačním režimu [první éra virtualizace]
 - ▶ Binary translation - přeložení do "bezpečného" kódu [VMWare x86]
 - Bezpečný kód vykonáván fyzickým CPU v privilegovaném [?] režimu
 - ▶ Hardwarová podpora virtualizace [VMWare x64]
 - Více úrovní privilegovanosti ve stavu procesoru
 - Na stroji s N úrovněmi lze emulovat virtuální stroj s N-1 úrovněmi
- ▶ Paravirtualizace [Xen, některé varianty Hyper-V]
 - Upravený OS nepoužívá privilegované instrukce
 - Ani jinak nezasahuje přímo do privilegovaného stavu (stránkování)
 - OS je vykonáván fyzickým CPU v aplikačním režimu
 - Ochrana proti chybám v OS
 - Privilegované akce nahrazeny voláním hypervizoru
 - Výrazně menší režie než při emulaci rozhraní SW-HW
 - Zůstává režie přepínání aplikace-OS přes hypervizor

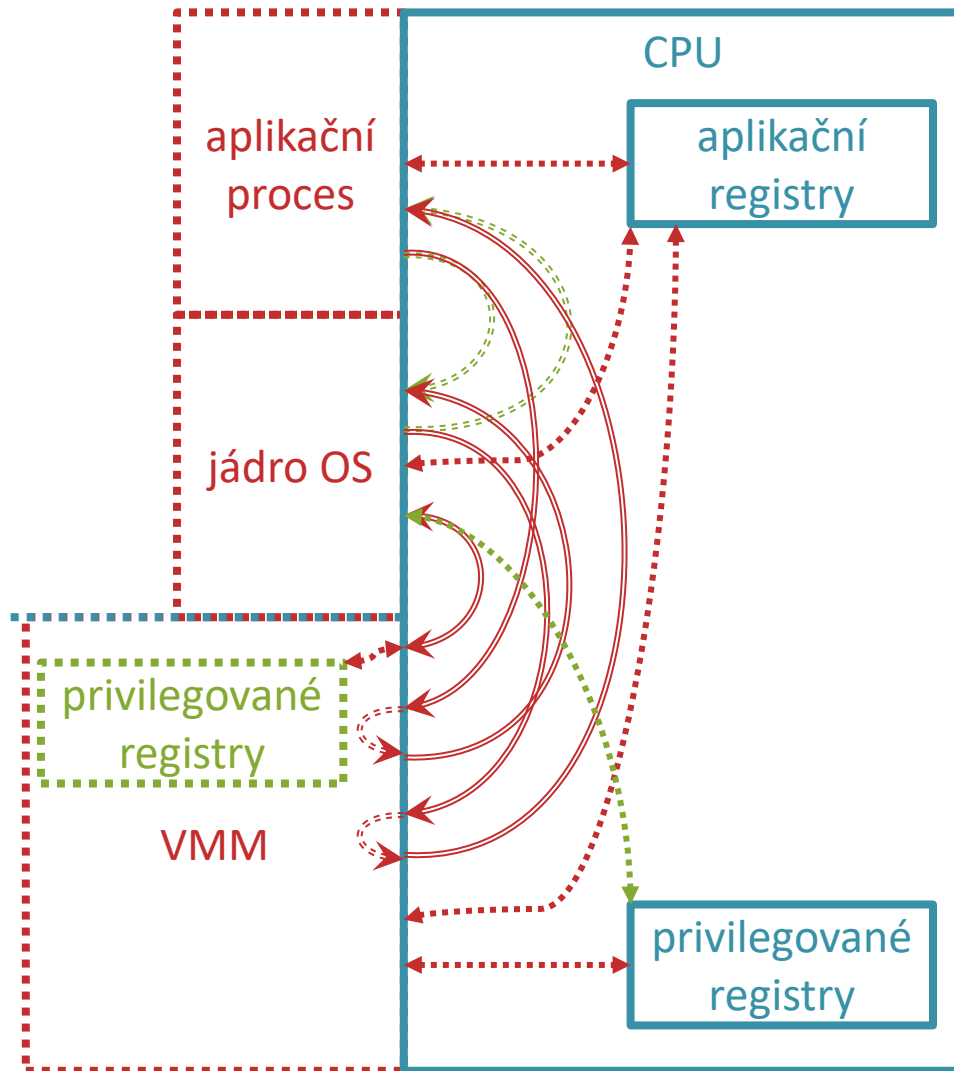


- ▶ Režim CPU
 - ▶ Aplikační
 - ▶ Privilegovaný
 - ▶ Odlišeno příznakem v privilegovaném registru
- ▶ Vstup do privilegovaného režimu
 - ▶ Přerušení
 - ▶ Instrukce pro volání jádra (SYSCALL)
 - ▶ Chyba
- ▶ Návrat z privilegovaného režimu
 - ▶ Instrukce návratu (IRET, SYSRET)

Operační systém na virtualizovaném CPU - Trap and Emulate



- ▶ Aplikační proces
 - ▶ Pracuje normálně
- ▶ Jádro OS
 - ▶ Pracuje v aplikačním režimu
 - ▶ Privilegovaná instrukce způsobí softwarové přerušení
 - ▶ VMM toto přerušení obsluhuje - emuluje instrukci, která ji způsobila
- ▶ Privilegované registry virtuálního CPU
 - ▶ uloženy v paměti VMM



- ▶ Jádro OS vyvolává hodně privilegovaných instrukcí
 - ▶ Počet závisí na architektuře CPU, systému a OS
- ▶ Softwarová emulace instrukcí je pomalá
 - ▶ Režie přerušení
 - ▶ Režie dekódování
 - ▶ Režie závisí na architektuře CPU

▶ První éra virtualizace

- ▶ IBM 370
 - I/O řešeno HW kanály = málo privilegovaných instrukcí v OS
 - Mizivý paralelismus = levné skoky
 - Jednoduchá a pravidelná instrukční sada = levné dekódování
 - Monolitické aplikace = málo meziprocesové komunikace
- ▶ Trap and Emulate byla vhodná technika

▶ Dnes

- ▶ Jádro OS vyvolává hodně privilegovaných instrukcí
 - Intenzivní komunikace mezi procesy a s I/O zařízeními
- ▶ Softwarová emulace instrukcí je pomalá
 - Režie přerušování
 - Režie dekódování
 - Režie závisí na architektuře CPU

▶ *Kompresa privilegií*

- ▶ Jádru virtualizovaného OS pracuje na jiné prioritní úrovni, než si myslí
 - Některé instrukce se chovají jinak (intel x86)
 - Řešeno velmi pracně překladem (VMWare)

▶ *Společný adresový prostor*

- ▶ CPU nepřepíná adresový prostor při volání jádra
 - Ochrana OS řešena privilegovanými stránkami
- ▶ Virtualizovaný OS nakládá s virtuálním adresovým prostorem jako s vlastním
- ▶ Nezbyvá místo pro VMM
 - VMWare: řešeno segmentací (dostupná pouze v 32-bit režimu)

▶ *Příliš mnoho přechodů VM-VMM*

Virtualizace virtuální paměti – starší metody

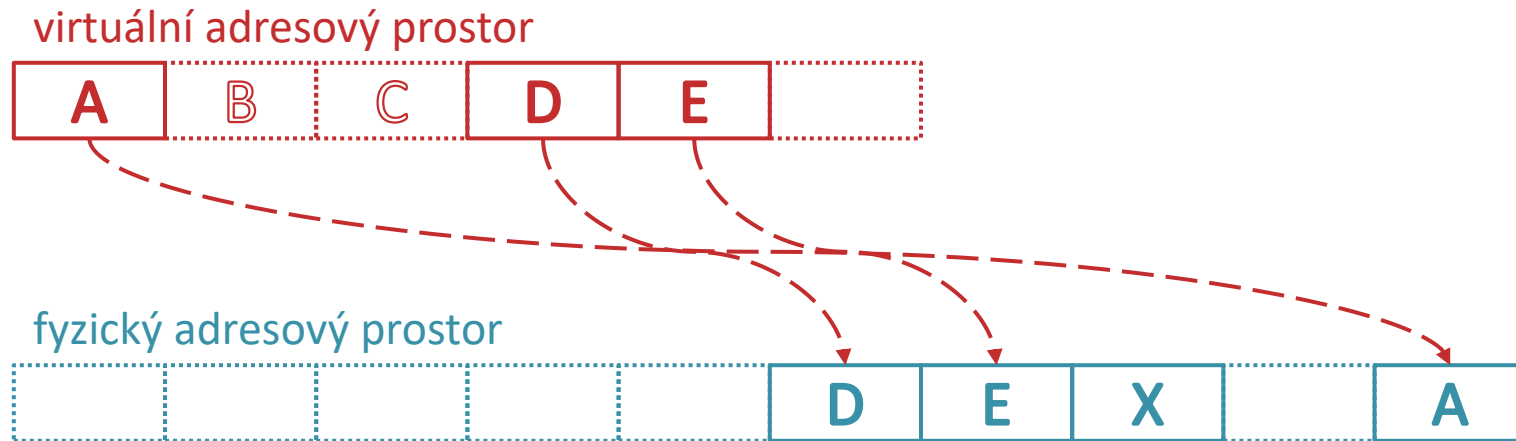
adresový prostor z pohledu procesu



▶ Adresový prostor z pohledu procesu

- ▶ Jeden nebo více souvislých úseků
 - Rozložení a význam určen dohodou aplikace a OS
- ▶ Dělení na stránky je neviditelné

Virtuální paměť ve fyzickém počítači



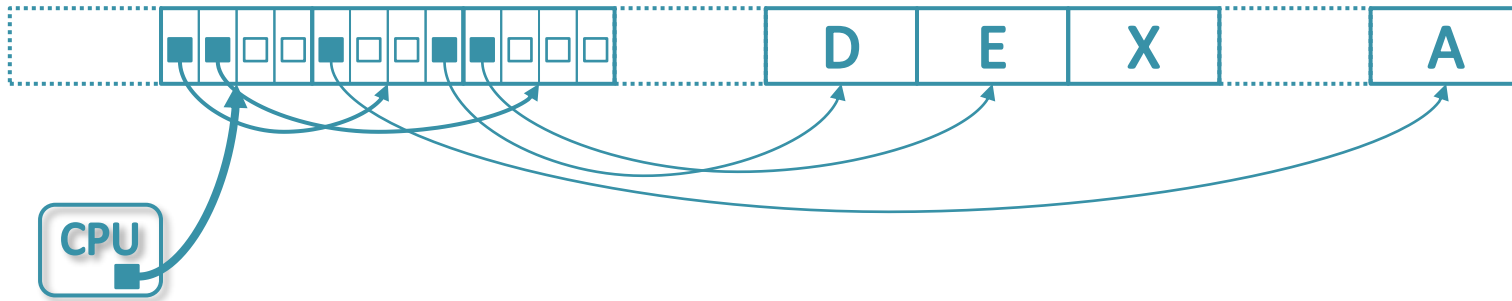
- ▶ Abstraktní pohled na virtuální paměť
 - ▶ Virtuální stránky jsou mapovány na fyzické rámce
 - ▶ Stránky odložené na disk mapovány nejsou
 - ▶ Fyzický adresový prostor je sdílen mnoha procesy

Virtuální paměť ve fyzickém počítači

virtuální adresový prostor



fyzický adresový prostor



► Realizace dvouúrovňovým stránkováním (x86)

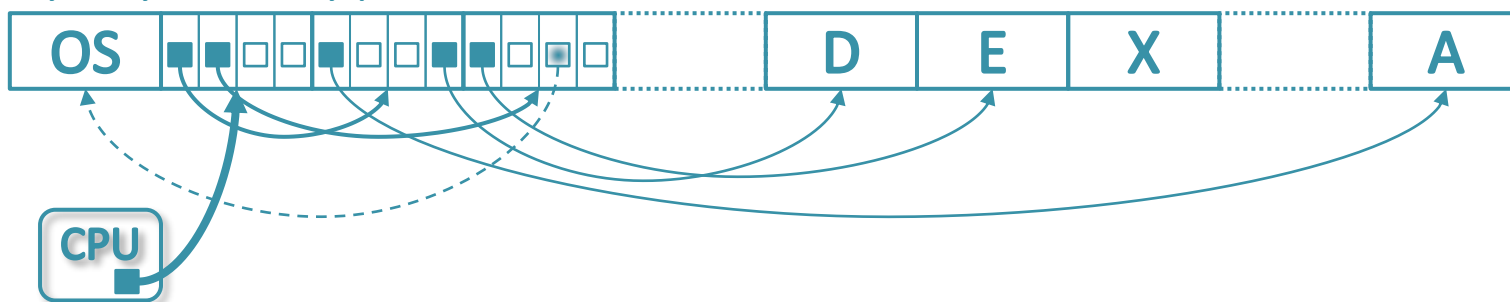
- Při výpadku TLB procesor prochází dvě úrovně stránkových tabulek
- Stránkové tabulky uloženy ve fyzickém adresovém prostoru
- Fyzická adresa kořene uložena v registru CPU

Virtuální paměť ve fyzickém počítači

virtuální adresový prostor



fyzický adresový prostor



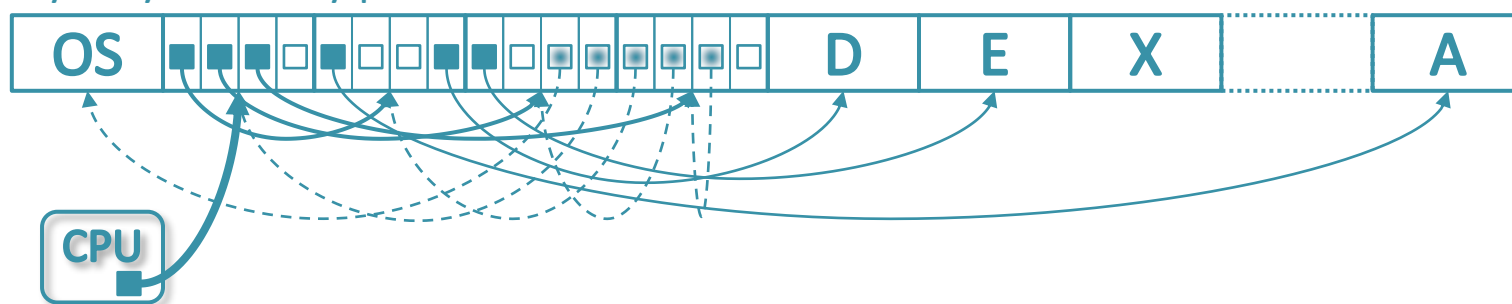
- ▶ Kód a data OS bývají součástí virtuálního adresového prostoru
 - ▶ Přepínání stránkování při každém volání OS by bylo neefektivní
 - ▶ Stránky OS přístupné pouze v privilegovaném režimu procesoru
 - ▶ Volání OS provedeno speciální instrukcí, která zapíná privilegovaný režim

Virtuální paměť ve fyzickém počítači

virtuální adresový prostor



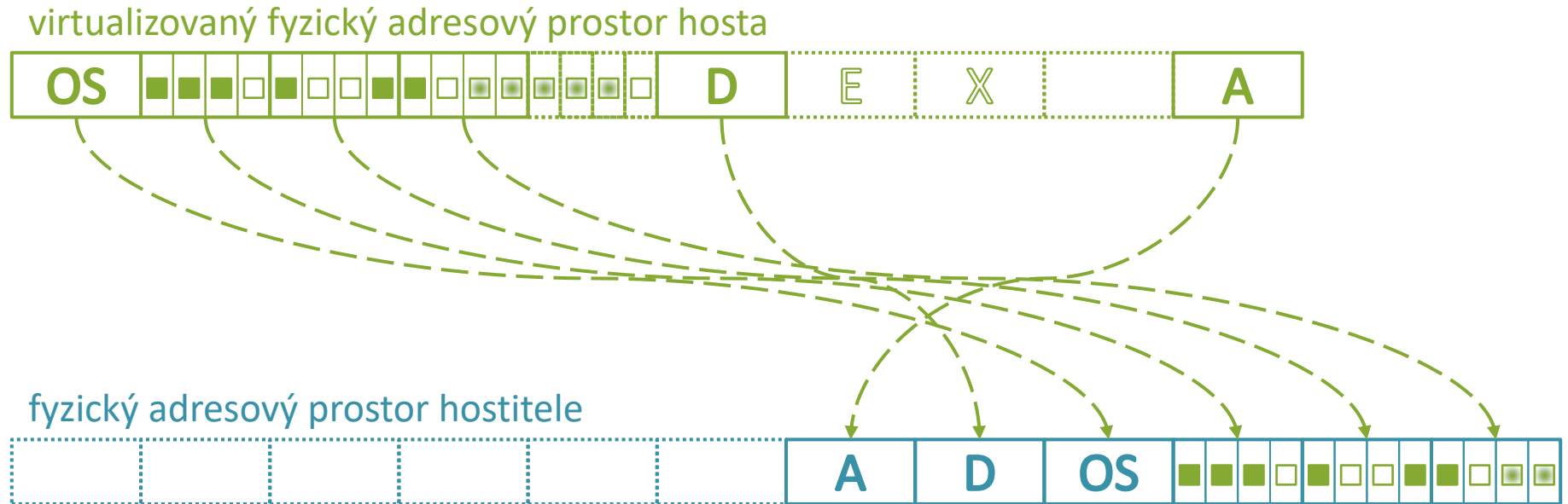
fyzický adresový prostor



▶ Operační systém plní stránkovací tabulky

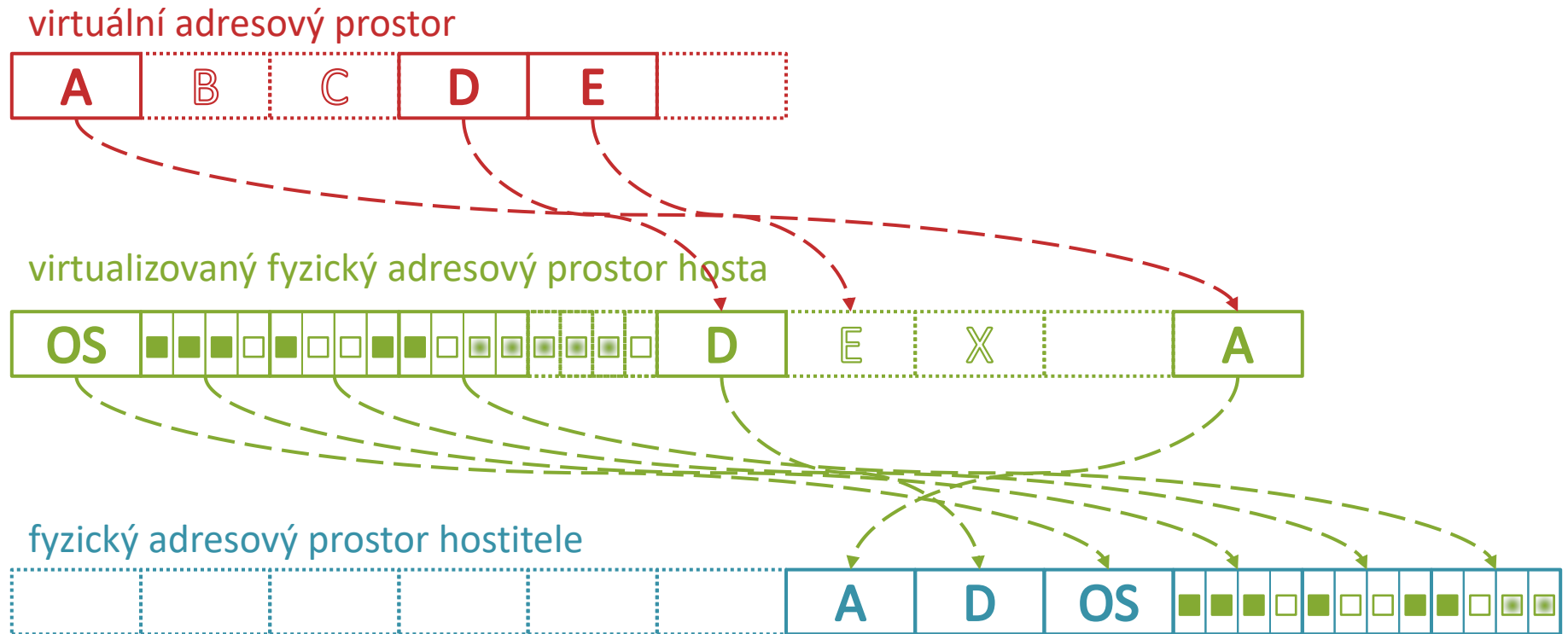
- ▶ Stránkovací tabulky jsou mapovány podobně jako data OS
- ▶ OS zapisuje do stránkovacích tabulek běžnými instrukcemi
- ▶ Zápis většinou musí být následován privilegovanou instrukcí "TLB flush"

Virtuální paměť ve virtuálním počítači



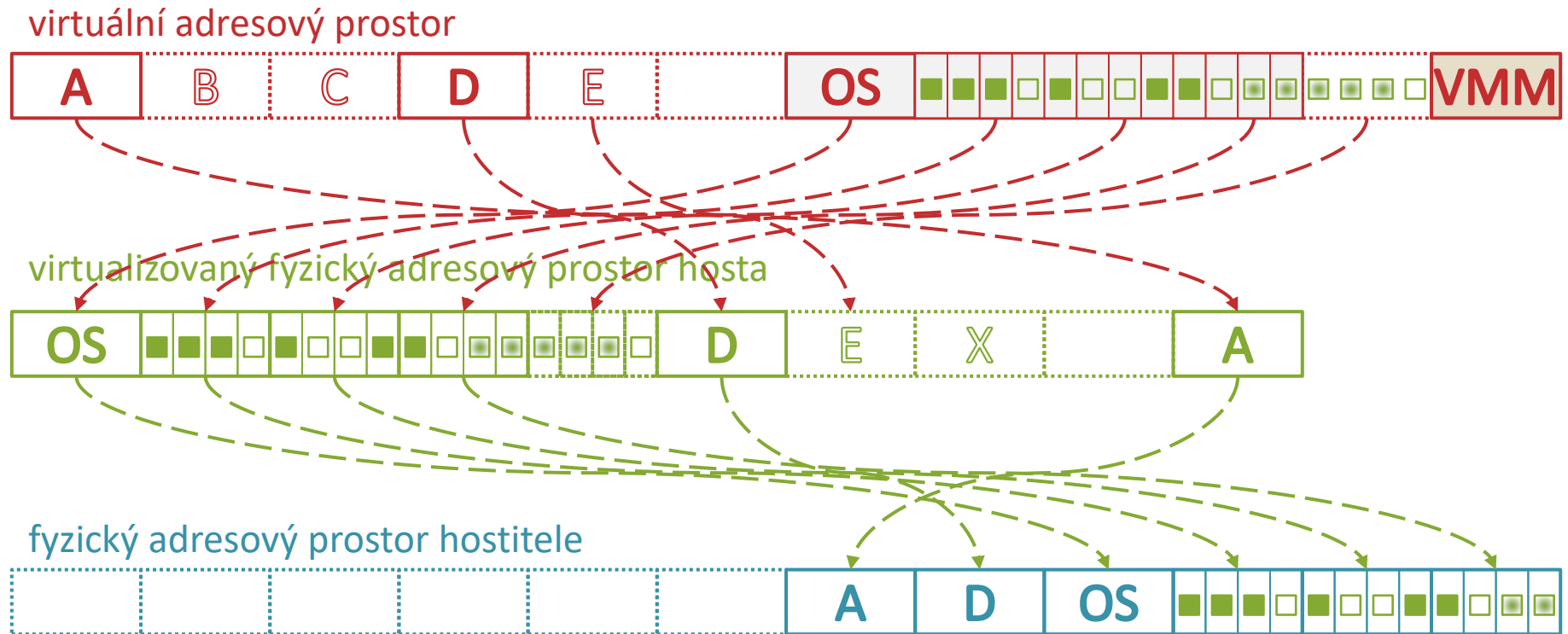
- ▶ VMM poskytuje iluzi fyzického adresového prostoru
 - ▶ Mapování fyzického prostoru hosta na fyzický prostor hostitele
 - ▶ VMM může odkládat stránky na disk podobně jako OS

Virtuální paměť ve virtuálním počítači



- ▶ Prováděný kód pracuje s virtuálním adresovým prostorem hosta
- ▶ Potřebné mapování vznikne složením
 - ▶ mapování definovaného operačním systémem hosta
 - ▶ mapování definovaného virtualizací počítače hosta

Virtuální paměť ve virtuálním počítači

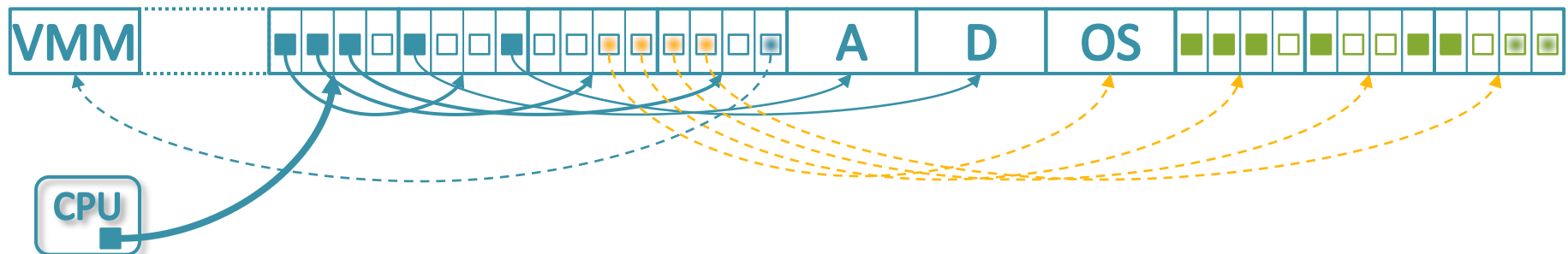


- ▶ Ve virtuálním adresovém prostoru běží
 - ▶ Aplikační procesy hosta
 - ▶ OS hosta
 - ▶ VMM
- ▶ Potřebujeme 3 úrovně privilegií ke stránkám

virtuální adresový prostor



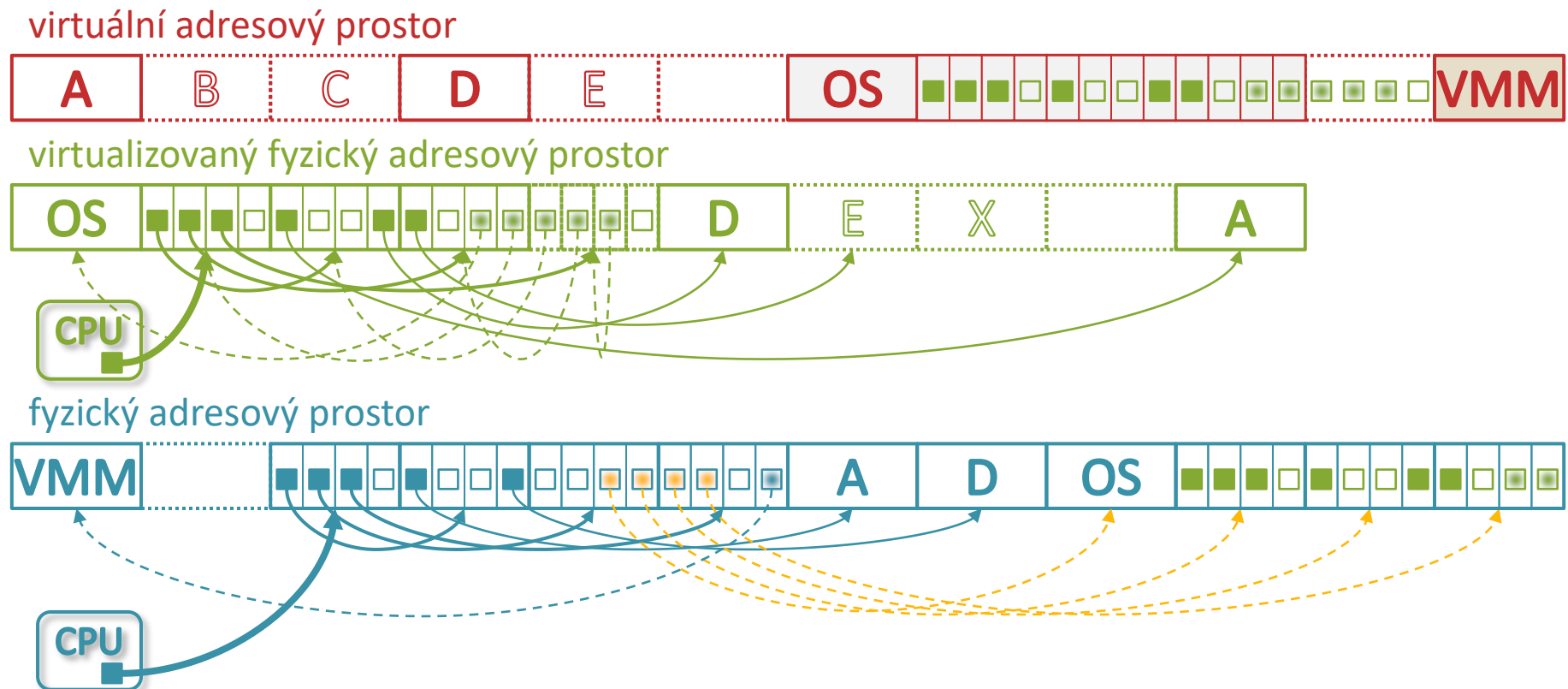
fyzický adresový prostor



► Realizace složeného mapování

- Stránkovacími tabulkami ve fyzickém prostoru hostitele
 - Obsahují fyzické adresy v prostoru hostitele

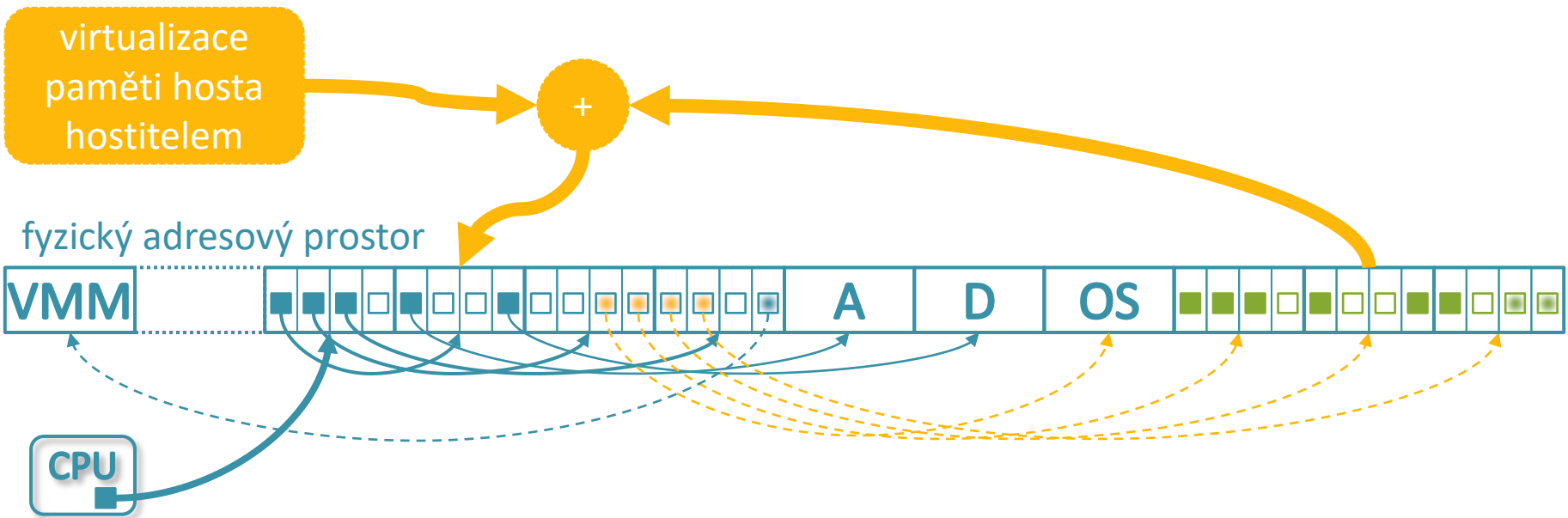
Virtuální paměť ve virtuálním počítači - bez HW podpory



- ▶ V systému jsou dvojí stránkovací tabulky
 - ▶ Stránkovací tabulky hostitele používané fyzickým CPU
 - Obsahují fyzické adresy v prostoru hostitele
 - ▶ Virtualizované stránkovací tabulky hosta používané OS hosta
 - Obsahují virtualizované fyzické adresy v prostoru hosta

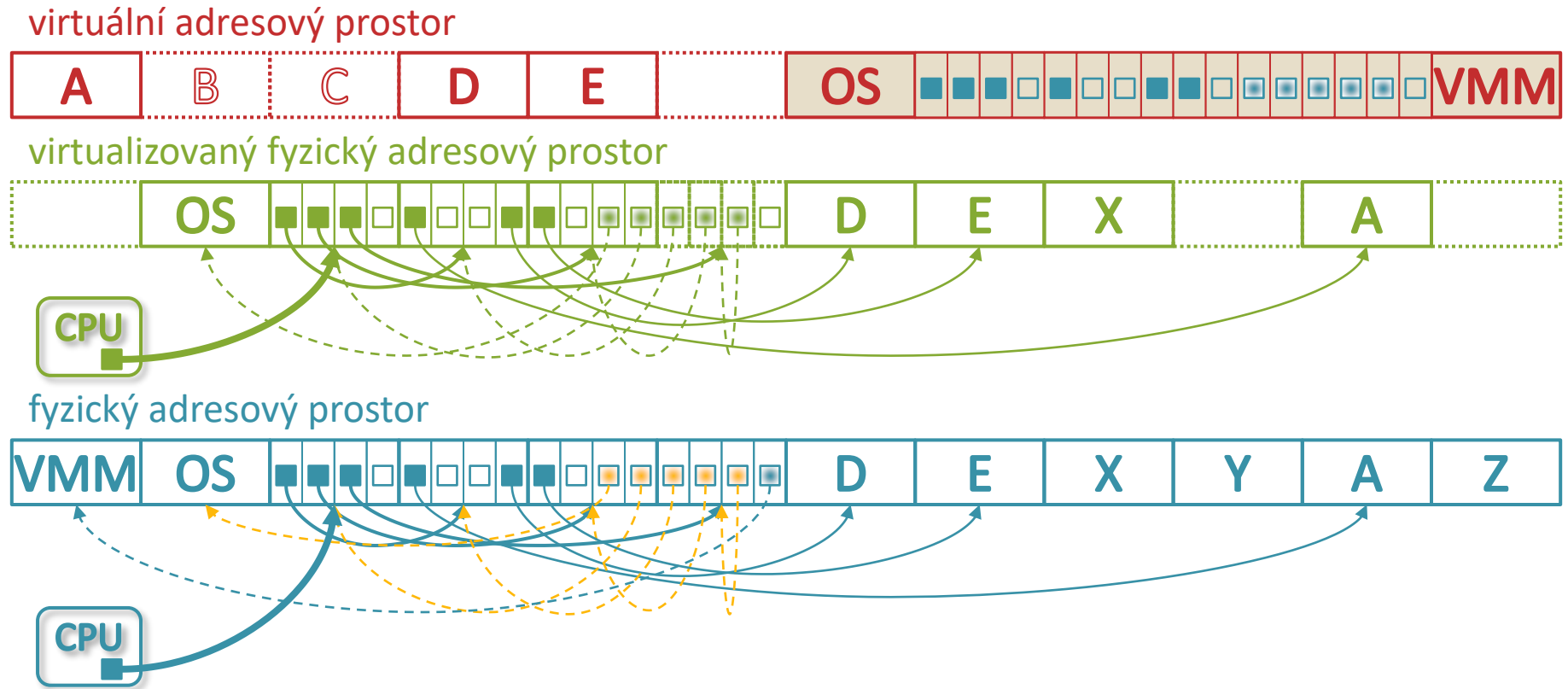
Virtuální paměť ve virtuálním počítači - bez HW podpory

virtuální adresový prostor

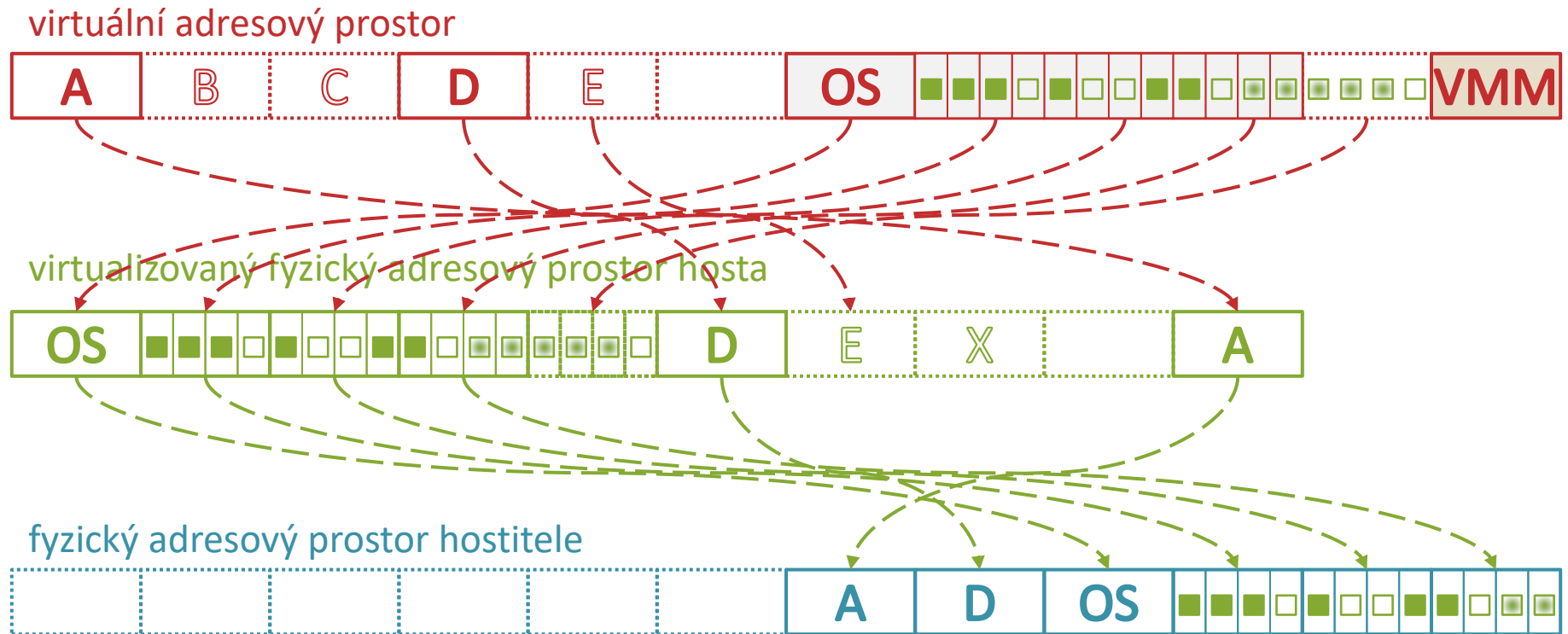


- ▶ VMM počítá výsledné mapování ze stránkovacích tabulek hosta
 - ▶ OS hosta zapisuje do svých tabulek neprivilegovanými instrukcemi
 - ▶ VMM musí zajistit přiměřenou koherenci fyzických tabulek a tabulek hosta
 - Virtualizované tabulky hosta mohou být mapovány read-only a zápisy emulovány
 - Koherenci lze udržovat v rámci emulace privilegované instrukce "TLB flush"

Virtuální paměť ve virtuálním počítači - bez druhé virtualizace



- ▶ Slabší VMM neumí odkládání virtualizované paměti na disk
 - ▶ Mapování virtualizovaných fyzických adres na fyzické je identita
 - VMM pouze kontroluje, zda OS hosta nemapuje nežádoucí fyzické adresy
 - ▶ OS hosta se musí vyrovnat s dírami ve fyzickém adresovém prostoru
 - Používáno převážně při paravirtualizaci (Xen)



► Extended (Intel) / Nested (AMD) Page Tables

- Skládání dvojího mapování provádí procesor sám
 - Týká se mechanismu hardware page-walk
 - TLB obsahuje složené mapování
 - Řeší i chybné naplnění stránkových tabulek

UNUSED SLIDES

Virtualizace - klady a zápory

- ▶ Lepší využití CPU
 - ▶ Většina počítačů se většinu doby fláká
- ▶ Lepší využití paměti, diskového prostoru, ...
 - ▶ OS neumí bezbolestně přidat nový prostor - velikosti bývají předimenzovány
 - ▶ Virtualizace dokáže prezentovat větší než skutečný prostor
- ▶ Možnost migrace virtuálních počítačů
 - ▶ Load-balancing, fault-tolerance
- ▶ Vzdálená správa
 - ▶ CD pro instalaci OS lze do virtuální mechaniky vložit kliknutím myši
- ▶ Checkpointy
 - ▶ Nepovedené změny v konfiguraci lze vrátit
- ▶ Výuka uživatelů/správce
- ▶ Testování a ladění OS, sítí i aplikací
 - ▶ Zkoumání malware

- ▶ Účel virtualizace (technický pohled)
 - ▶ Větší počet virtuálních objektů než fyzických
 - virtuální vs. fyzická paměť, virtuální počítač
 - ▶ Virtuální objekty jiného druhu než fyzické
 - emulace počítače jiné architektury
 - ▶ Virtuální objekty vzdálené od fyzických
 - vzdálené disky prezentované jako lokální, vzdálená klávesnice+obrazovka
 - ▶ Virtuální objekty implementované zcela jinak, než fyzické
 - disky implementované souborem
 - ▶ Virtuální objekty bez vazby na fyzický svět
 - virtuální síť
- ▶ Zásahy do chování, které by bez virtualizace nebyly možné
 - ladění, experimenty, měření
 - šízení všeho druhu (time sharing, thin provisioning)
 - migrace, load balancing

- ▶ Ztráta výkonu
 - ▶ Silně závisí na charakteru aplikací i technologii virtualizace
 - ▶ Někdy jednotky, někdy desítky procent
- ▶ Změna charakteristik při migraci
 - ▶ Různá CPU
- ▶ Nespolehlivé měření/ladění výkonu
- ▶ Nepřipravenost fyzické síťové infrastruktury
 - ▶ Migrace virtuálních síťových karet mezi fyzickými
- ▶ Nepřipravenost dodavatelů software
 - ▶ Nevýhodné licenční podmínky
 - ▶ Problémy s individuálními checkpointy v komunikujících systémech

Virtualizace – technické principy

▶ Virtualizovaná zařízení

▶ Počítač

- Virtualizované rozhraní
 - fyzické rozhraní software-hardware (fyzická virtualizace)
 - softwarové rozhraní uvnitř OS (paravirtualizace)
- Zahrnuje virtualizaci zařízení uvnitř počítače

▶ CPU

- Virtualizované rozhraní = instrukční sada
 - Aplikační + privilegované instrukce (hardwarová virtualizace)
 - Aplikační instrukce (paravirtualizace)
- Samo o sobě nema smysl - CPU nemá vstup/výstup

▶ Paměť

- Virtualizované rozhraní = instrukce čtení a zápisu

▶ I/O zařízení

- Virtualizováno
 - na úrovni I/O instrukcí
 - na softwarovém rozhraní uvnitř OS

▶ Sdílení fyzického počítače virtuálními

▶ CPU

- guest OS i hypervizor používají preemptivní multitasking

▶ Paměť

- guest OS už má svou virtuální paměť
- hypervizor přidává druhou úroveň

▶ Disky

- virtuální disk mapován do společného diskového prostoru
- iSCSI, SAN, NAS,...

▶ Sítě

- trunk mode, NAT, virtuální sítě,...

▶ Další zařízení

- exkluzivní přístup
- sdílený přístup
- vzdálené USB apod.

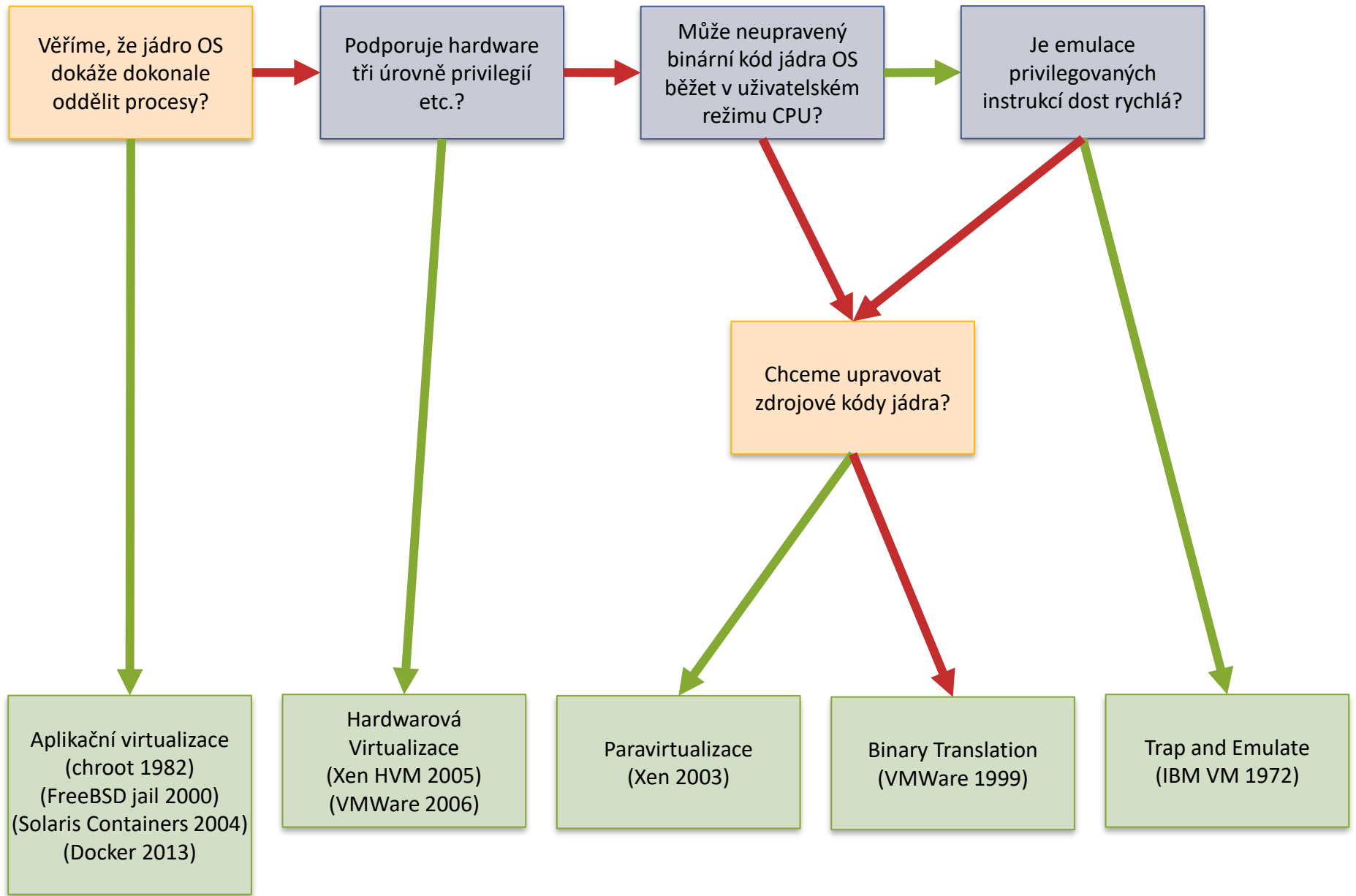
Virtualizace – abstraktní pohled

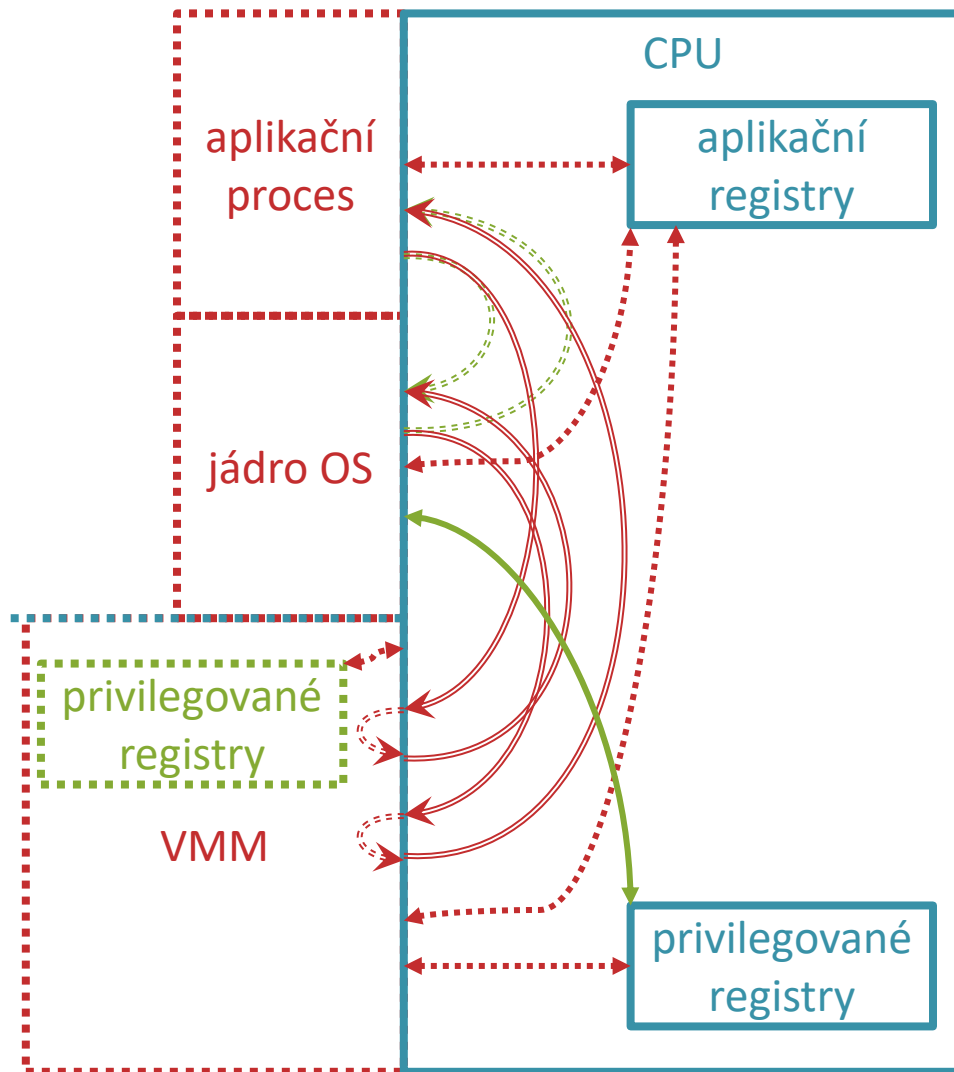
- ▶ Virtualizace se týká běhu software
- ▶ Cílem virtualizace je
 - ▶ sdílet hardware pro běh několika (nesouvisejících) kusů software
 - ▶ přesouvat běžící nebo pozastavené kusy software jinam
 - ▶ pozorovat chování software apod.
- ▶ Software je
 - ▶ FORTRAN, C, ... – přeložený program obsahující instrukce fyzického CPU
 - ▶ Java, C#, ... – napůl přeložený program obsahující instrukce virtuálního stroje
 - ▶ Python, PHP, ... – zdrojové kódy
- ▶ Běh software vždy zahrnuje běh knihovných funkcí
 - ▶ Část knihoven je vždy v podobě instrukcí fyzického CPU (typicky přeloženo z C)
 - V mnoha případech se v nativních knihovnách odehrává většina běhu
 - Vyšší jazyk a nativní knihovny typicky komunikují sdílením paměti
 - Oddělení běhu nativních knihoven od vyšších vrstev je prakticky nerealizovatelné
 - ▶ I v případě vyšších jazyků je vhodné chápat běh software jako provádění instrukcí fyzického CPU (knihovna, kód generovaný JIT překladem, interpreter)

- ▶ Software je program prováděný jako posloupnost nativních instrukcí
- ▶ Software dnes obvykle neběží jako jeden samostatný program
 - ▶ Z hlediska balancování zátěže a migrace je potřeba řešit skupinu programů společně
 - ▶ Jejich spolupráce je zajištěna službami OS – tyto služby je nutné také zahrnout do balancovaného/migrovaného softwarového balíku
- ▶ Základní otázkou virtualizace je místo, kde je balík software odříznut od okolí
 - ▶ Kontejnery – na rozhraní aplikace-OS
 - ▶ Paravirtualizace – uprostřed OS
 - ▶ Skutečná virtualizace – na rozhraní OS-HW

- ▶ Virtualizaci lze dělat na mnoha úrovních
 - ▶ Aplikační virtualizace
 - chroot, WoW, UAC, kontejnery, bash.exe
 - skupinám procesů je prezentováno jiné prostředí
 - implementováno operačním systémem
 - ▶ Paravirtualizace
 - Xen, Microsoft Hyper-V
 - na fyzickém stroji běží několik upravených operačních systémů
 - hypervizor řeší alokaci zdrojů a serializaci přístupu k zařízením
 - ▶ (Hardwarová/Klasická) Virtualizace
 - VMWare, Microsoft Hyper-V, Xen HVM
 - na fyzickém stroji běží několik neupravených operačních systémů
 - hypervizor vytváří každému z nich iluzi fyzického hardware

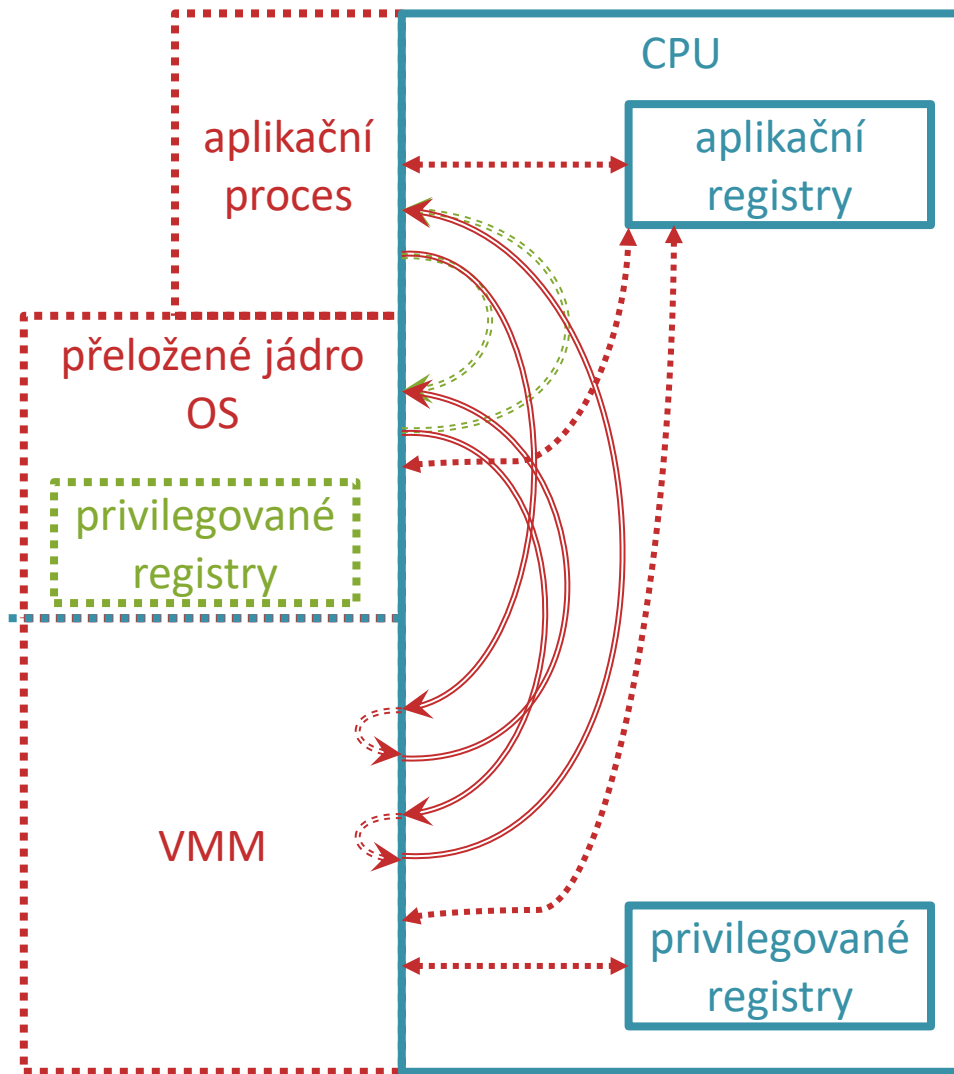
Volba přístupu k virtualizaci



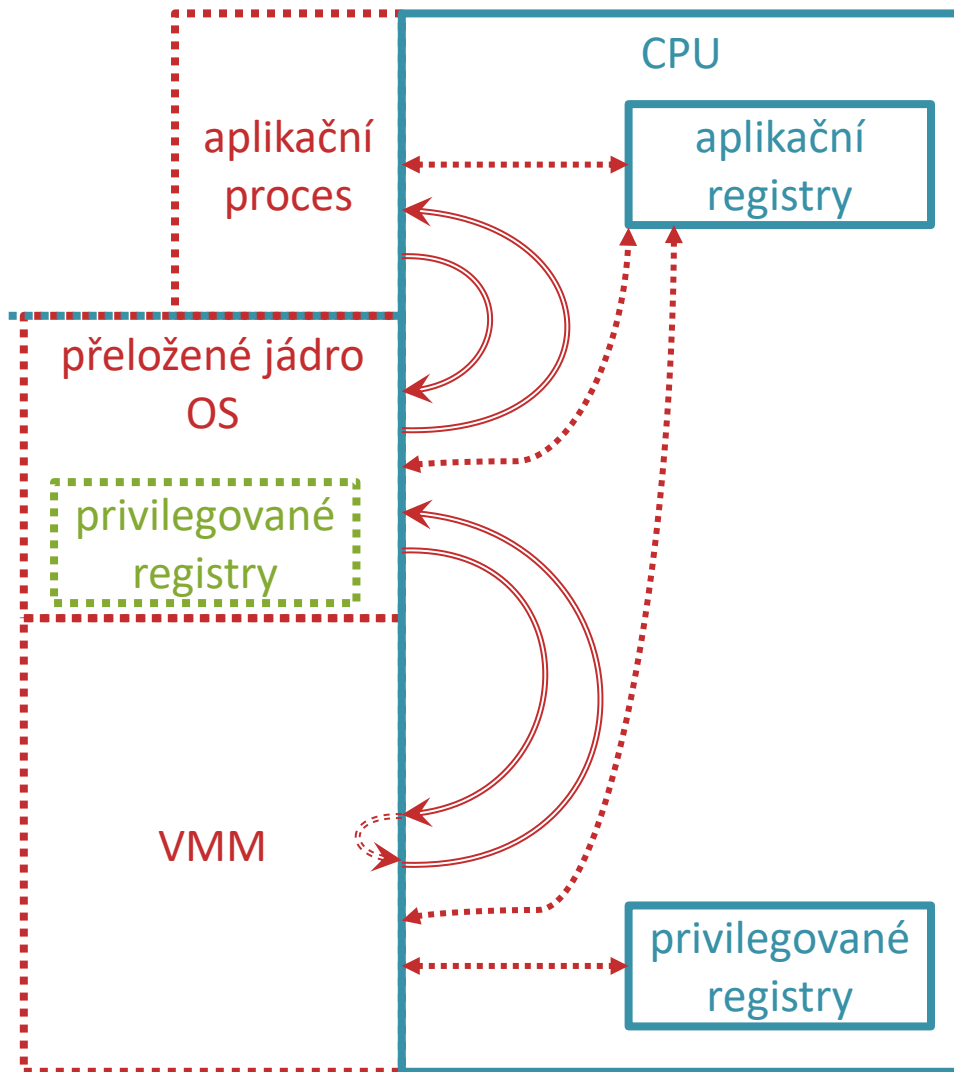


- ▶ Nevhodná architektura procesoru neumožňuje použití Trap and Emulate
 - ▶ Intel x86
- ▶ Typické chyby
 - ▶ Část privilegovaných registrů je čitelná nepriviligovanou instrukcí
 - ▶ Některé instrukce se v různých režimech chovají různě
 - ▶ Příliš mnoho instrukcí jádra OS vyvolává v aplikačním režimu chybu

Binary Translation



- ▶ Úprava jádra OS
 - ▶ Binární kód jádra je překladačovými technikami upraven tak, aby neprováděl privilegované operace
 - ▶ Privilegované registry CPU si upravený kód emuluje sám
 - ▶ Zásah VMM nutný pro:
 - Přechody aplikace-jádro
 - Akce s významným efektem (např. na stránkování)
 - I/O operace
 - Systém přerušení



- ▶ Úprava jádra OS
 - ▶ Binární kód jádra je překladačovými technikami upraven tak, aby neprováděl privilegované operace
 - ▶ Privilegované registry CPU si upravený kód emuluje sám
 - ▶ Jádro je překladem fakticky donuceno dobrovolně spolupracovat s VMM
 - Přechody aplikace-jádro jsou stejně drahé jako bez VM
 - Přechody jádro-VMM jsou levné
 - ▶ Nutná důvěra v mechanismus překladu