

trap and emulate: <http://www.cs.utexas.edu/~mwalfish/classes/s10-cs372h/lectures/l27.txt>
segmentace! <http://cs.wikipedia.org/wiki/Segmentace>

high availability ~ minimum down-time

VM - virtual machine

VMM - virtual machine monitor - hypervizor, arbitr virtualizovanych PC

host - fyzicky pocitac, na kterem bezi guests

guest - virtualizovany PC

Technika virtualizace

Vývoj:

- mainframe → klientske terminaly pripojene k nejakemu superpocitaci
- jednoduche klient/server prostredi → spousta PC pripojenych pres sit k ruznym serverum
- distrib. prostredi → klienti jsou pripojeni pres web server, ten je pres nejakou branu pripojeny k siti, v siti jsou dalsi klienti, kteri jsou pres rozhrani pripojeni k serverovemu prostredi (nejaka skupina multicomputeru), kde jsou directory server, application server, email server, file server, DB server, management server, content server
 - clustering
 - vykon, stabilita
 - nakladnost
 - virtualizace
 - sdileni prostredku a vykonu, vyvazovani
- pozdeji se klienti pripojuji i pres wifi, VPN, na stanici bezi vice serveru, DB clustery
- cloud → virtualni HW platforma pro provozovani virt. serveru a sluzeb
 - skalovatelnost, elasticita
 - zmena vykonu podle potreb klienta
 - sluzby uctovany podle skutecneho vyuziti
 - vysoka dostupnost a spolehlivost
 - zalozni systemy na ruznych urovnich (servery, infrastruktura, DC)
 - SW, ktery umoznuje za behu rychlou vymenu vadneho clanku
 - automaticke aktualizace
 - sdileni systemu
 - thin provisioning - poskytuje dojem dostupnosti vice prostredku, nez fyzicky existuje - virtualni alokace prostoru na ulozeni dat
 - organizace provozuje jeden fyzicky server a na nem např. čtyři virtuální servery. I když má fyzický server limitovanou velikost operační paměti (např. 8 GB), mohu nadefinovat na virtuálních

serverech v součtu více operační paměti, než kolik je fyzicky dostupné (např. každý virtuální server bude pracovat se 4 GB RAM). Pokročilými systémovými technikami (know-how výrobců virtualizačních technologií) je pak možné sdílet RAM fyzického serveru tak, že opravdu vzniká dojem nezávislých virtuálních serverů s definovanou virtuální velikostí RAM - tyto techniky např. umožňují sdílet podobné bloky dat, využívat nevyužitou paměť jednoho virtuálního serveru jiným apod.

- úspora při investicích do infrastruktury
- multi-tenancy - poskytování služby více klientům
 - virtuální oddíly dat a konfigurace - každý klient pracuje se svou customizovanou verzí
 - virtualizované zdroje jsou logicky odděleny, nelze pracovat s cizími daty
 - v kontrastu s: multi-instancí architektury, kde separátní instance SW pracují s klienty
- druhy cloudu
 - veřejný - poskytování IaaS, PaaS, SaaS třetí stranou - nejčastější
 - zajištěná vysoká škálovatelnost, účtování dle využití zdroje
 - soukromý
 - infrastruktura poskytující služby pouze jedné organizaci
 - schopnost účtování jednotlivým složkám organizace
 - komunitní
 - služby využívány komunitou - spolupracující firmy, projekt
 - hybridní
 - cloud složený z více různých zdrojů
 - cloud interoperability - sky computing

Webhosting

- motivace pro virtualizaci
- postupný vývoj
 - statické stránky
 - 1 web = 1 adresář
 - 1 proces - apache - pro všechny
 - dynamické stránky
 - může obsahovat potenciálně nebezpečný kód
 - 1 web = 1 proces (Apache+PHP)
 - 1 počítač pro všechny
 - uživatelské systémy
 - weby vyžadují odlišné konfigurace
 - 1 web = 1 počítač
 - bez virtualizace je to zbytečně drahé
- jak to zhruba vypadá

- při použití fyzických PC
 - na každém fyzickém PC je jedna doména
 - PC připojeny do internetu
- při použití virtuálních PC
 - na jednom fyzickém PC je spousta virtuálních počítačů, každý spravuje nějakou doménu
 - dale zde bezí hypervizor (virtual machine manager), který spravuje virt. PC
 - fyzický PC připojen do internetu
 - podrobněji:
 - u fyzických PC je u každého počítače NIC (network interface controller)
 - u virtuálních PC je u každého virt. PC virtuální NIC
 - hypervizor je pak napojen na fyzický NIC, který je připojen k internetu a požadavky na jednotlivá virt. PC, která spravuje, organizuje on

Virtualizace

- virtuální
 - very close to being something without being it
- iluze fyzického zařízení, které fyzicky neexistuje
- známe rozhraní fyzického zařízení (stranka paměti, počítač, disk, síťová karta) nebo software zařízení (proces-kernel, HAL - HW abstraction layer) jsou implementovány jinak než obvykle
 - softwarově/jiným HW (SCSI radice)/v kombinaci HW+SW (virtuální paměť, moderní virtualizace procesoru)
- používat tento termín je otázkou zvyku
 - SSH/RDP není nazýváno virtualizace konzole
 - iSCSI není virtuální disk
 - JVM není virtualizace fyzického stroje
- proč?
 - můžeme mít větší počet virtuálních objektů než fyzických
 - virtuální vs fyzická paměť, virtuální počítač
 - virtuální objekty mohou být jiného druhu než fyzické
 - emulace PC jiné architektury
 - virtuální objekty mohou být vzdálené od fyzických
 - vzdálené disky jsou prezentovány jako lokální, vzdálená klávesnice+obrazovka
 - virtuální objekty mohou být implementovány zcela jinak než fyzické
 - disky implementovány souborem
 - virtuální objekty mohou být bez vazby na fyzický svět
 - virtuální síť

- muzeme provadet zasahy do chovani, ktere by bez virtualizace nebyly mozne
 - ladeni, experimenty, mereni
 - šízení všeho druhu - thin provisioning, time sharing apod.
 - migrace, load balancing
- jaka zarizeni?
 - pocitac
 - virtualizovane rozhrani
 - fyzicke rozhrani SW-HW (nazyva se fyzicka virtualizace)
 - SW rozhrani uvnitr OS (paravirtualizace)
 - virtualizace zarizeni uvnitr
 - CPU
 - virtualizovane rozhrani
 - instrukcni sada
 - aplikacni+privilegovane instrukce (HW virtualizace)
 - aplikacni instrukce (paravirtualizace)
 - virtualizovat CPU samo o sobe nema smysl - chybi I/O
 - pamet
 - virtualizovane rozhrani
 - instrukce cteni a zapisu
 - I/O
 - virtualizace
 - na urovni I/O instrukci
 - na SW rozhrani uvnitr OS
- na jake urovni se virtualizuje ~ jak probiha sdileni fyzickeho PC virtualnimi
 - CPU
 - guest OS i hypervizor (ten bezi na host machine) pouzivaji preemptivni multitasking (umi prerusit prave vykonavajici proces bez spoluprace OS)
 - pamet
 - guest OS ma svoji virtualni pamet
 - hypervizor pridava druhou uroven pameti
 - disky
 - virtualni disk je namapovan do spolecneho diskoveho prostoru
 - iSCSI (Internet Small Computer System Interface, IP-based protokol)
 - protokol umoznuje klientum - iniciatorum posilat SCSI prikazy na SCSI disk. zarizeni, ktera jsou ale na vzdalenyh serverech
 - organizace pomoci SAN protokolu
 - SAN protokol (storage area network)
 - dedikovaná (oddělená od LAN, WAN, atd) datová síť, která slouží pro připojení externích zařízení k serverům (disková pole, páskové knihovny a jiná zálohovací zařízení). SAN vznikla hlavně kvůli narůstajícím potřebám na zabezpečení a konsolidaci dat.
 - NAS (network attached storage)
 - označení pro datové úložiště připojené k místní síti LAN. Data toho

úložiště mohou být poskytována různým uživatelům. NAS nemusí mít pouze funkci souborového serveru, ale může mít i jiné specializované funkce. Například klient P2P sítě, webový server a další. Většinou obsahuje nějaký vestavěný počítač, který má za úkol sdílení dat a podporu různých protokolů.

- site
 - trunk mode
 - metoda, která systému poskytuje síťový přístup pro spousty klientů sdílením množiny linek a frekvencí místo toho aby je každému klientovi poskytovala individuálně
 - analogické struktury strumu, kde je jedna hlavní pater a spousta větví k ní připojených
 - NAT (network address translation)
 - síťová maskaráda
 - způsob úpravy síťového provozu přes router přepisem výchozí a/nebo cílové IP adresy, často i změnu čísla TCP/UDP portu u průchozích IP paketů
 - virtuální síť
 - VPN, VLAN (logical LANs based on physical LANs)
- další zařízení
 - resi se exkluzivní přístup/sdílení přístup, vzdálené porty (USB) apod.
- varianty virtualizace
 - aplikační virtualizace
 - technologie, která zapouzdruje (odděluje) aplikační SW od podkladového OS, na kterém běží
 - plně virtualizovaná aplikace není nainstalována v tradičním smyslu, nicméně spouští se, jakoby byla
 - při běhu se aplikace chová jakoby přímo komunikovala s rozhraním originálního OS a všemi prostředky, které OS spravuje, ale může být ve skutečnosti izolována (např. sandboxing - bezpečnostní mechanismus v rámci počítačové bezpečnosti, který slouží pro oddělování procesů běžících se stejným oprávněním - jail, chroot, nebo právě i VM apod.)
 - v tomto kontextu termín virtualizace reflektuje zapouzdření, narozdíl od HW virtualizace, kde termín odrazí abstrakci
 - chroot, WoW (windows on windows), UAC (user account control)
 - skupinám procesů je prezentováno jiné prostředí
 - implementuje se operačním systémem
 - paravirtualizace
 - Xen, MS Hyper-V
 - na fyzickém stroji běží několik upravených OS
 - místo nebezpečných operací volá VMM
 - hypervizor resi alokaci zdroje a serializaci přístupu k zařízením
 - prezentuje SW rozhraní virtuálním strojům, které je podobné, ale není

identické jako rozhraní pro HW, nad kterým leží

- záměr je redukovat čas, který stráví guest OS při provádění náročných operací ve virtuálním prostředí (které bezí jednodušeji na nevirtualizovaném prostředí) - poskytují se speciálně definované hooks, které umožňují guestovi/hostovi požádat/potvrdit o nativní běh těchto operací, které by jinak byly spouštěny ve virt. doméně, což by mělo za následek horší výkon
- vyžaduje aby guest OS umožňoval explicitní podporu para-API (někdy se používají alespoň některé komponenty, které API podporují - např. poskytují balík ovladačů zařízení, které podporují para-API, které se instalují do guest OSu)

- virtualizace

- VMWare, MS Hyper-V
- na fyzickém stroji bezí několik neupravených OS
- hypervizor vytváří každému z nich iluzi fyzického HW
- <http://en.wikipedia.org/wiki/Virtualization>
- HW virtualizace
 - IBM 370
 - VMWare pro 64bit stroje
 - fyzicky proc. vykonává původní instrukce
 - nebezpečné instrukce vyvolávají VMM a jsou jim emulovány - trap and emulate
- SW virtualizace s překladem
 - fyzicky procesor vykonává upravené (přeložené) instrukce
 - instrukce překládá VMM z binárního kódu (BT - binary translation)
 - nebezpečné instrukce jsou při překladu nahrazeny voláním VMM
 - překlad se obvykle týká pouze jádra OS hosta
 - kód určený pro privilegovaný režim procesoru
 - ve virt. stroji bezí v neprivileg. režimu
 - aplikační procesy bez překladu
- často se kombinuje s prvky paravirtualizace
 - ve virtuálním stroji bezí SW (ovládací zařízení) komunikující s VMM
 - omezuje ztráty výkonu při emulaci privilegovaných instrukcí a IO
 - zvyšuje uživatelský komfort

- motivace

- lepší využití CPU
 - většina PC se totiž většinu doby fláka
- lepší využití paměti, diskového prostoru
 - OS neumí bezbolestně přidat nový prostor - velikosti bývají předimenzované
 - virtualizace dokáže prezentovat větší než skutečný prostor
- možnost migrace virtuálních PC

- load-balancing, fault-tolerance
 - vzdalena sprava
 - CD pro instalaci OS lze do virtualni mechaniky vlozit kliknutim mysi
 - checkpointy
 - nepovedene zmeny v konfiguraci lze vratit
 - vyuka uzivatele a spravcu
 - viz. administrace windows, unixu
 - testovani a ladeni OS, siti, aplikaci
 - napr. zkoumani malware
- problémy
 - ztraty vykonu
 - silne zavisi na charakteru aplikaci i technologii virtualizace
 - nekdy jednotky, nekdy desitky procent
 - zmena charakteristik pri migraci
 - ruzna CPU, apod.
 - nespolehlive mereni/ladeni vykonu
 - nepripravenost fyzicke sitove infrastruktury
 - migrace virt. sitovych karet mezi fyzickymi
 - nepripravenost dodavatele SW
 - nevyhodne licencni podminky
 - problémy s individualnimi checkpointy v komunikujicich systemech

Historie virtualizace

- první era
 - 1972 - IBM VM na S/370 mainframes
 - koexistence ruznych OS
 - time-sharing a virtualni pamet nad OS, ten tyto pojmy nezna
 - ladeni OS
 - 1980 - postupny zanik
 - levnejsi architektury - miniPC, PC
 - novy HW nepodporuje virtualizaci
 - nastup Unixu (VM zbytecne komplikuji komunikaci mezi procesy)
- druha era
 - 1999 - VMWare workstation
 - SW virtualizace
 - VMM (hypervizor) jako aplikace Windows NT
 - 2002 - VMWare ESX Server
 - VMM nahrazuje OS hosta
 - 2003 - Xen
 - paravirtualizace - modifikace OS guesta
- architektura x86 je nevhodna
 - dedictvi procesoru 80286
 - pokusy o napravu - rozsireni procesoru: 2005 Intel VT-x, 2006 AMD-V

- naprava není dodnes dokonalá
 - velký podíl SW virtualizace (v 1. řadě byla HW)
 - nezanedbatelná ztráta výkonu - volání jádra, přerušeni, virtualizace virt. paměti
- situace se ale rok od roku zlepšuje

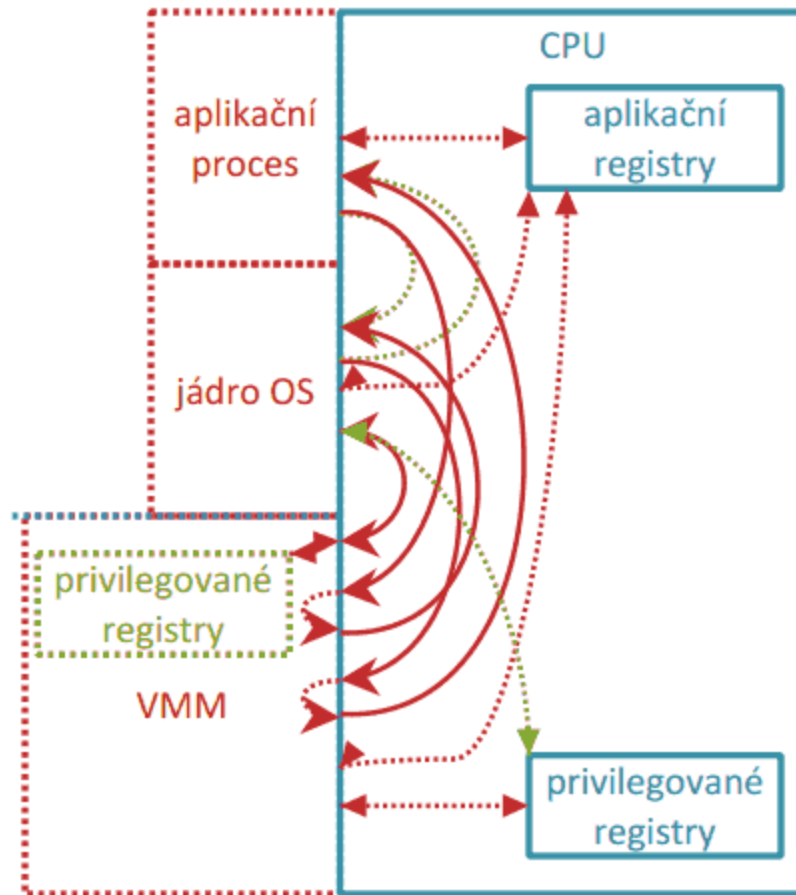
Požadavky na virtualizaci

- equivalence / fidelity (věrnost)
 - program běžící pod VMM se musí chovat v zásadě stejně jakoby běžel na ekvivalentním stroji přímo
- resource control / safety
 - VMM musí mít úplnou kontrolu nad virt. zdroji
- efficiency / performance
 - statisticky převládající část strojových instrukcí musí být prováděna bez zásahu VMM

Implementace VMM

Virtualizace CPU

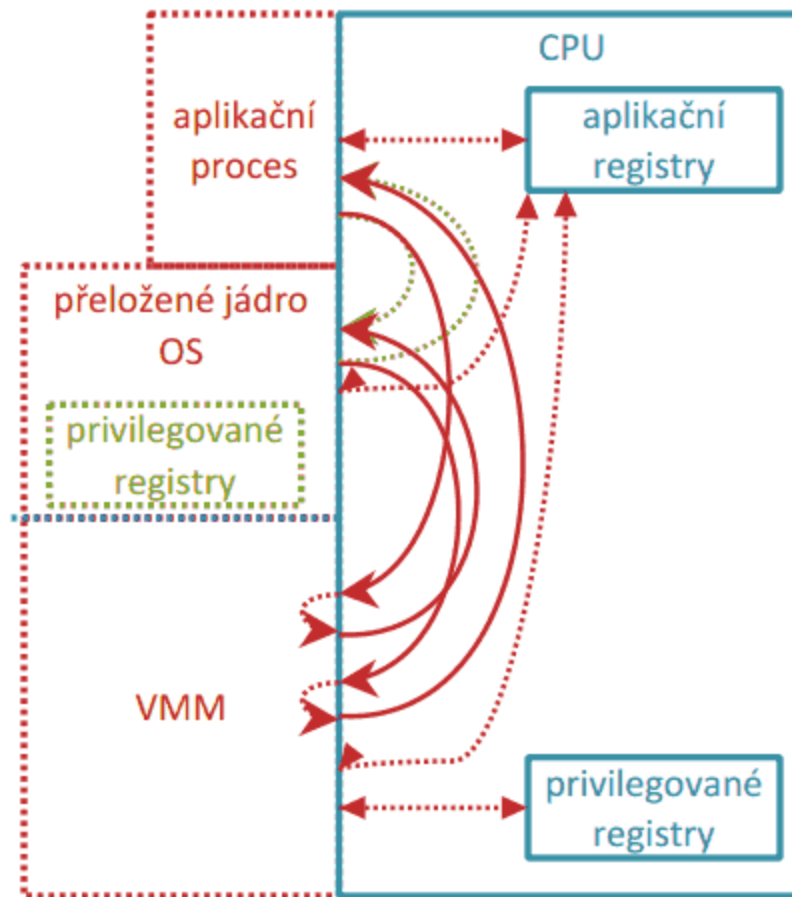
- situace na klasickém fyzickém CPU
 - režim CPU je buď aplikační nebo privilegovaný
 - je to odlišeno příznakem v privilegovaném registru
 - pro vstup do privilegovaného režimu
 - dojde nejprve k přerušeni
 - instrukce pro volání jádra
 - vyvolá se chyba
 - návrat z privilegovaného režimu
 - použití instrukce návratu
 - aplikační proces má svoje aplikační registry, ke kterým přistupuje i jádro OS, jádro má zase privilegované registry, ke kterým má přístup jen ono samo
- situace na virtualizovaném CPU - trap and emulate
 - aplikační proces pracuje normálně
 - jádro OS
 - pracuje v aplikacním režimu
 - privilegované instrukce způsobí SW přerušeni, které obslouží VMM a odemuluje instrukci, která přerušeni způsobila
 - privilegované registry virtuálního CPU jsou uloženy v paměti VMM
 - aplikační proces má opět svoje aplikační registry, ke kterým přistupuje i jádro OS, to ale nyní musí na privilegované registry fyzického CPU přistupovat přes VMM - je tu vložená vrstva navíc



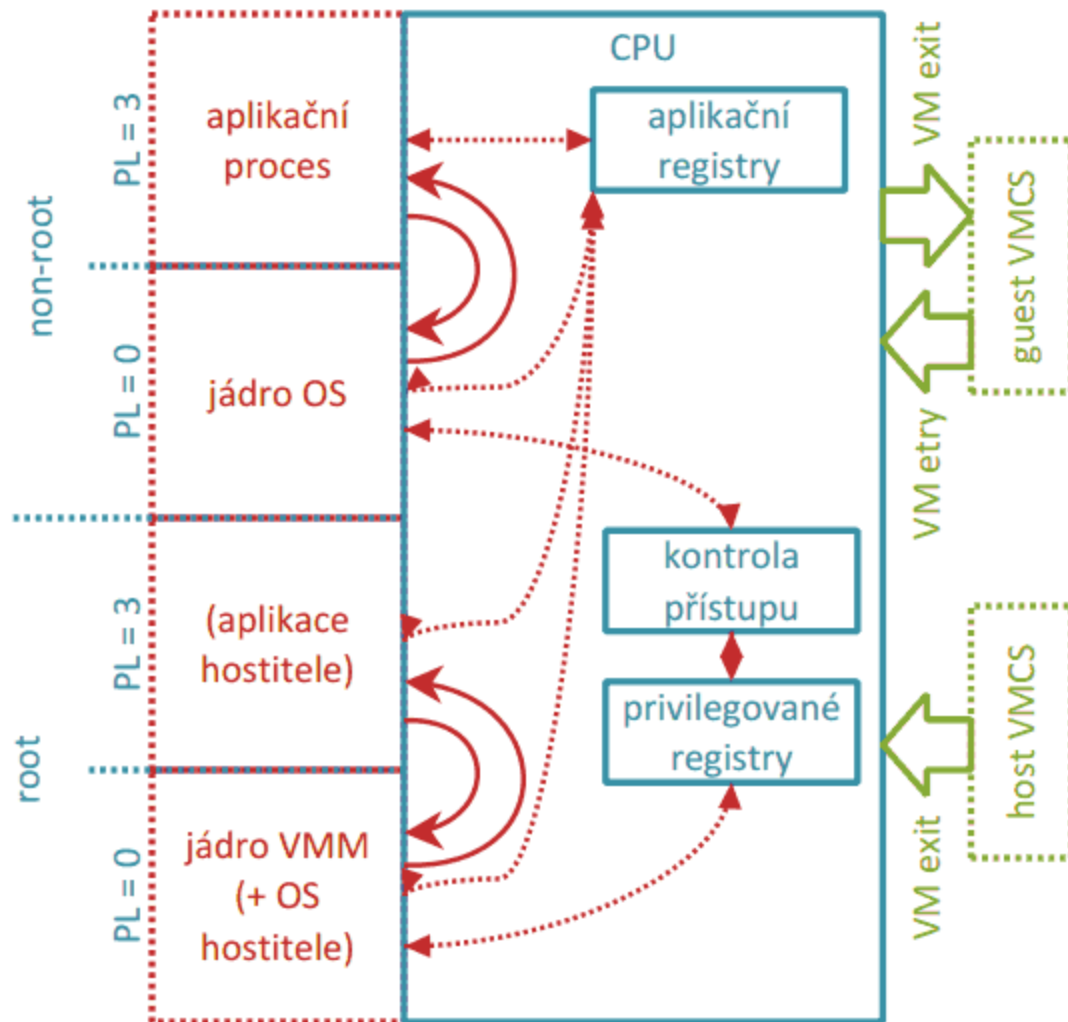
○ problémy

- jádro OS vyvolava hodně privilegovaných instrukcí, počet závisí na dané architektuře CPU, systému a OS
- SW emulace instrukcí je pomalá
 - režie prerušení
 - režie dekodování
 - režie závisí na architektuře CPU
- vhodné u první éry virtualizace - IBM 370
 - I/O řešeny HW kanály → málo privilegovaných instrukcí v OS
 - tenkrát mizivá podpora paralelismu → levné skoky
 - jednoduchá a pravidelná instrukční sada → levné dekodování
 - monolitické aplikace → málo meziprocesové komunikace
- nevhodná pro Intel x86 - typické chyby:
 - část privilegovaných registrů je čitelná neprivilegovanou instrukcí
 - current privilege level (CPL) je uložen ve spodních dvou bitech %cs (code segment - adresa, kde běží program) registru → instrukce typu `mov %cs, %ax` prozradí guest OSu, že je to virtuální mod, což je špatné

- nektore instrukce se v ruznych rezimech chovaji ruzne
 - treba nezpůsobí prerušení
 - příliš mnoho instrukcí jádra OS vyvolává v aplikačním režimu chybu
- jádro virt. OS pracuje na jiné prioritní úrovni než si myslí - “komprese privilegií”
 - nektore instrukce se tak chovají jinak (x86)
 - řeší se to velmi přesně překladem např. u VMWare
- společný adresový prostor
 - CPU nepřepíná adresový prostor při volání jádra
 - ochrana OS řešena privilegovanými stránkami
 - virt. OS nakládá s VAS jako s vlastním
 - nezbyvá místo pro VMM - VMWare řeší segmentaci
- příliš mnoho přechodu VM-VMM
- SW emulace s překladem
 - úprava jádra OS
 - binární kód jádra je překladacovými technikami upraven tak, aby neprováděl privilegované instrukce
 - privilegované registry CPU si upraví kód emuluje sám
 - VMM ale musí v některých případech zasahovat
 - návrat do aplikačního procesu
 - akce s významnými efekty na stránkování apod.
 - I/O operace
 - systém interruptu



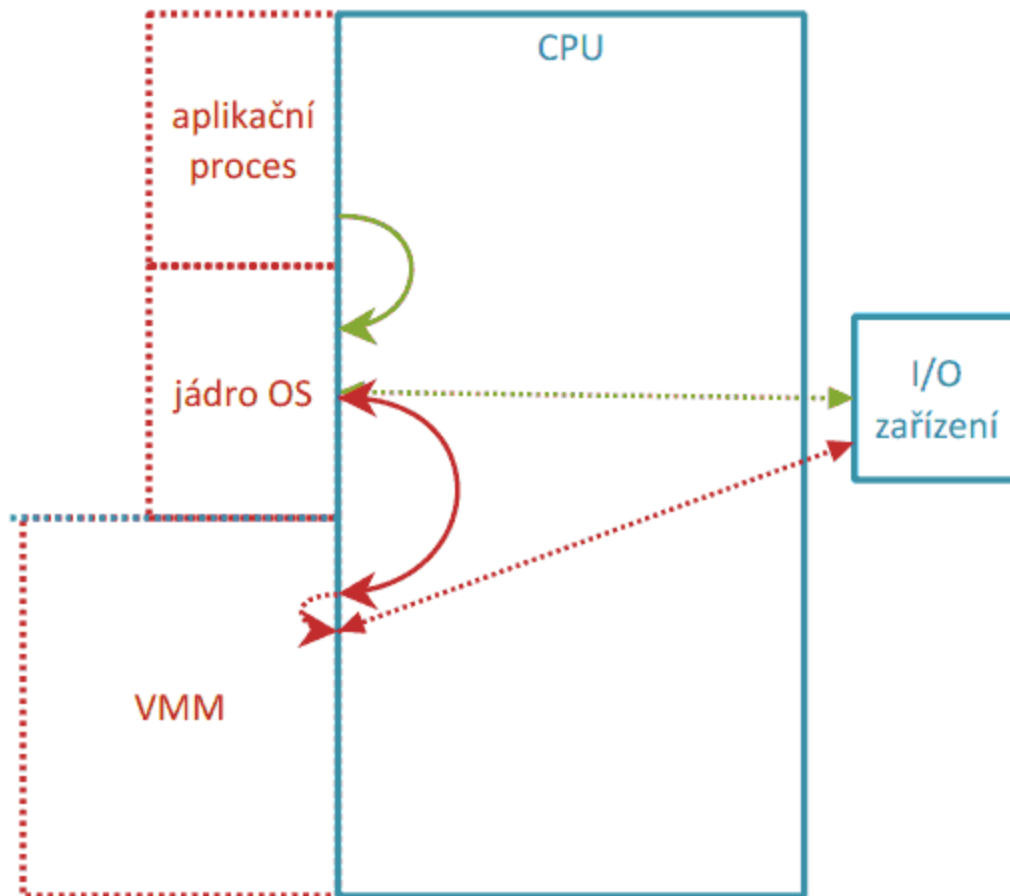
- HW podpora virtualizace
 - nový rozmer privilegovanosti
 - Intel VT-X / AMD-V, provedení se liší
 - VMCS - virtual machine control structures
 - HW-based pametová struktura na Intelu, která zrychluje operace VM entry a VM exit, které slouží pro delegaci vykonávání privilegované instrukce
 - uvažujeme root režim CPU
 - odpovídá CPU bez virtualizace
 - lze využít pro běh host OS
 - non-root režim CPU
 - má omezený přístup k privilegovanému stavu
 - nežádoucí akce způsobují VM exit (vypadek stránky napr.)
 - prepínání režimu
 - kritická část stavu CPU se načítá/ukládá do paměti
 - zahrnuje přepnutí stránkování



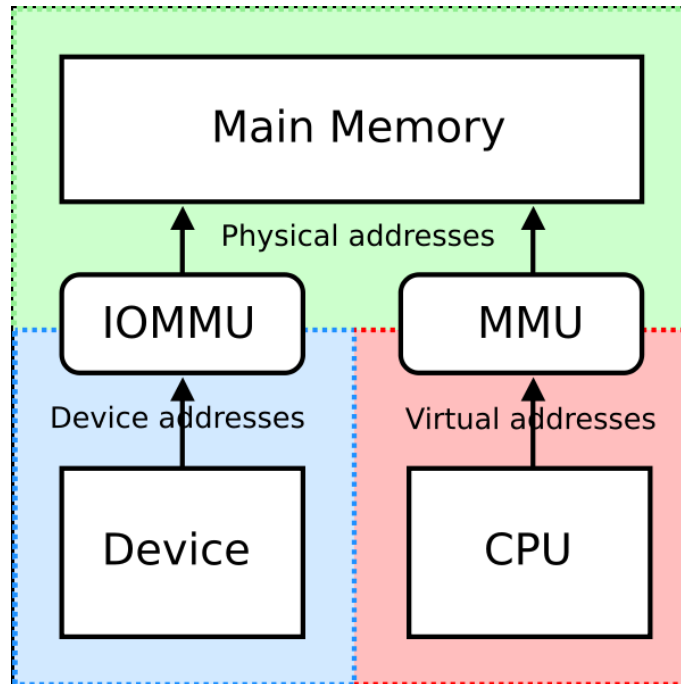
- **vyhody:**
 - odstranena komprese privilegií
 - za podmínky že virt. OSy samotné HW podporu virtualizace nepoužívají
 - nelze virtualizovat virtualizátor
 - na IBM 370 demonstrováno 5 úrovní vnorení virt.
 - prepínání adresového prostoru při VM entry a VM exit
 - ochrana paměti VMM, plná transparence pro virt. OS
 - komplikuje však přístup VMM do paměti VM (emulace I/O apod.)
 - menší počet přechodu VM-VMM
 - lze vyladit konfiguraci HW kontroly přístupu k privilegovanému stavu
 - demonstrováno asi 20x zrychlení některých úloh
 - unix fork & wait benchmark
 - kompilace rozsáhlých projektů s malými moduly

Virtualizace I/O

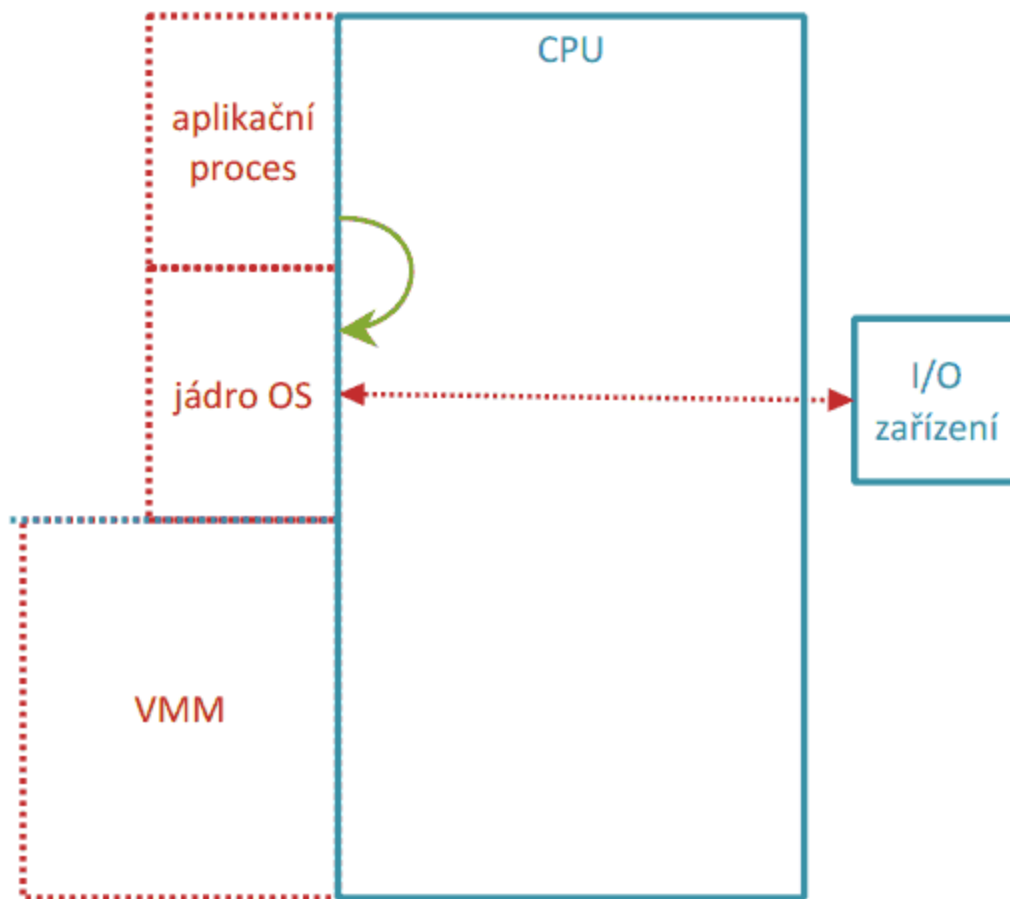
- situace na klasickém fyzickém CPU
 - aplikační procesy realizují všechny I/O voláním OS
 - OS komunikuje s I/O zařízením
 - používá privilegované I/O instrukce nebo
 - paměťové mapované zařízení chráněné stránkovacím mechanismem
- přístup k fyzickému IO na virtuálním CPU
 - privilegované IO instrukce jsou provedeny emulátorem ve VMM
 - paměťové mapované zařízení může být zpřístupněno přímo
 - exklusivní přístup
 - k danému zařízení může přistupovat jen jeden virtuální stroj
 - kromě samotného IO zařízení je třeba zpřístupnit nebo virtualizovat systém přerušování, případně DMA



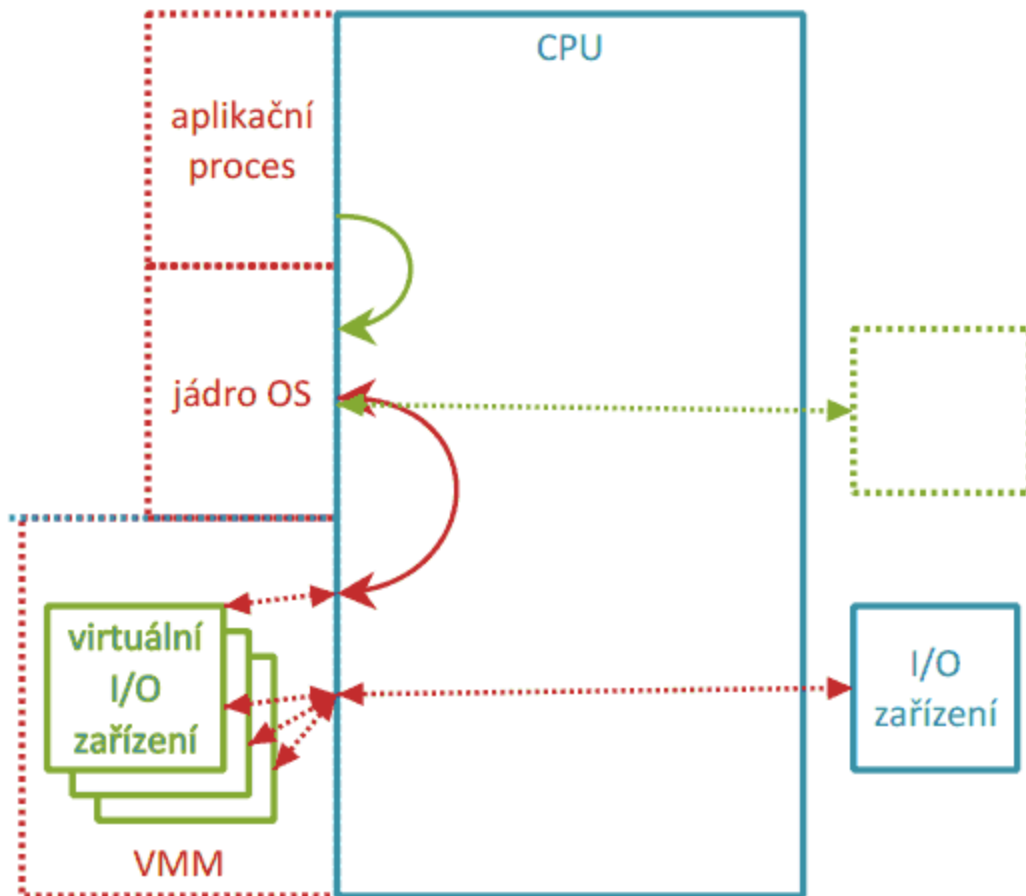
- v případě, že IO instrukce k danému zařízení nejsou privilegované, tak komunikuje jádro virtuálního OS přímo s IO zařízením
 - vyžaduje to konfigurovatelnost HW ochrany IO prostoru - IOMMU



- stejne podminky exkluzivniho pristupu
- moznost sdileneho pristupu
 - IO zarizeni se prezentuje vicekrat v IO adresovem prostoru
 - IO zarizeni ma pro kazdou adresu jednu kopii vnitrich stavovych registru
 - krome samotneho IO zarizeni je treba zpristupnit nebo virtualizovat system preruseni, pripadne DMA - stejne jako u ex. pristupu



- přístup k virtuálnímu IO na virtuálním CPU
 - privilegované IO resp. přístupy na memory-mapped IO jsou emulovány VMM
 - VMM pro každý stroj zvlášť emuluje chování HW
 - daný typ HW nemusí fyzicky existovat
 - sdílený přístup
 - VMM z emulovaného HW extrahuje logické akce
 - ty jsou prováděny fyzickým zařízením



- přístup k virtuálnímu IO s upravou OS
 - zasah do OS
 - paravirtualizace - OS je výrazně upraven
 - při klasické virtualizaci je do OS přidán ovladač virtuálního zařízení
 - výhody:
 - mezi OS a VMM jsou předávány logické příkazy a nikoliv fyzické IO
 - předání nevyžaduje emulaci IO instrukcí
 - logických příkazů je méně
 - serializace příkazů z různých VM je jednodušší

HW podpora virtualizace

- Intel VT-x a VT-d
 - řada rozšíření CPU i podpůrného čipsetu k podpoře virtualizace
 - stále přibývají další
 - jednotlivé úpravy jsou často použitelné nezávisle
 - významné virtualizační SW je využívají téměř všechny
 - Intel spolupracuje s producenty HW
- AMD-V
 - úpravy ve stejném směru podobným směrem

- vetsinou neni kompatibilni s Intellem
- situace znehlednena obchodni politikou
 - ruzne verze CPU maji ruzny uroven podpory
 - obchodni nazvy maskuji podstatu veci
 - nektera rozsireni jsou triviality, jina jsou velmi netrivialni
- Intel VT-x
 - extenze procesoru
 - root/non-root execution - reseni problemu komprese privilegii
 - extended page table (EPT) - reseni problemu virtualizace virtualni pameti
 - separadni mnozina str. tabulek preklada guestovi fyzicke adresy na hostovi fyzicke adresy → guest OS si muze modifikovat sve vlastni str. tabulky a primo obsluhovat vypadky stranek
 - umoznuje VMM vyhnout se VM exitum souvisejicich s virtualizaci str. tabulek, coz je hlavni overhead virtualizace bez EPT
 - flexpriority - virtualizace klicove casti radice preruseni
 - eliminace vetsiny VM exitu zpusobenych pristupy na privilegovane registry
 - flexmigration - virtualizace identifikace CPU a jeho schopnosti
 - moznost vyberu z dostupnych serveru, ktere jsou pro nas nejvhodnejsi - tim ze jsou virtualizovane jejich identifikace a dostupne schopnosti
 - migrace za behu mezi Intel servery
 - virtual processor ID (VPID) - klic zaznamu TLB obsahuje identifikator VM
 - pak neni treba invalidovat celou TLB pri prepinari VM-VMM, VM-VM
 - vguest preemption timer
 - casovac s lepsi granularitou a rychlejsi obsluhou
 - opt. na virtualizaci aplikaci s mirnymi real-time naroky
 - pause-loop exiting
 - HW podpora pro detekci spin-locků způsobující exit (preempci) do VMM
 - zlepšení výkonu
 - pro provoz více virt. procesorů na méně fyzických
 - real-mode support
 - podpora virtualizace pri startu virtualizovaneho OS (early VMM load, guest boot and resume)
- Intel VT-d
 - extenze pro podporu IO virtualizace v chipsetu
 - IOMMU
 - adresový prostor I/O má virtuální a fyzické adresy podobně jako paměť
 - interrupt-remapping support
 - částečná virtualizace řadiče přerušení
 - address translation services support
 - podpora virtualizace MMU při DMA
 - umožnění PCI-E zarizenim cachovat IOTLB položky, ktere se používají při DMA mapování
 - rozšíření standardu sběrnice PCI Express

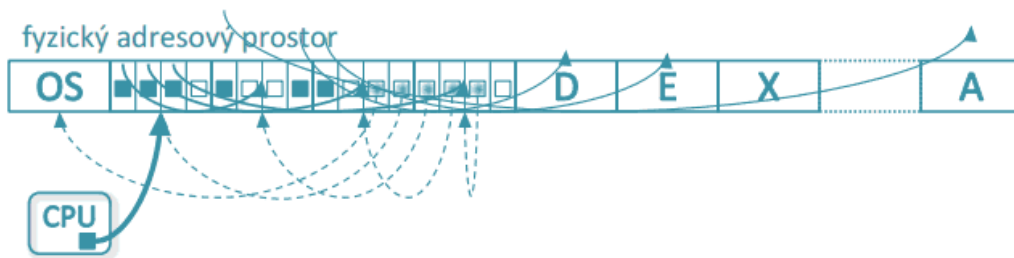
- large Intel VT-d pages
 - podpora vetsich stranek - 2MB, 1GB
 - umozňuje sdílení CPU a DMA verzí stránkových tabulek
- virtual machine device queue
 - zlepšení propustnosti, snížení utilizace CPU
 - uspořádání a grouping paketu na úrovni NIC místo na úrovni VMM
 - network interface card (NIC) s více stavovými prostory pro přímý přístup z VMs
- single-root I/O virtualization (SR-IOV)
 - I/O zařízení deklarují své schopnosti virtualizace (možnost sdílení ostatními)
 - rozšíření standardu PCI Express

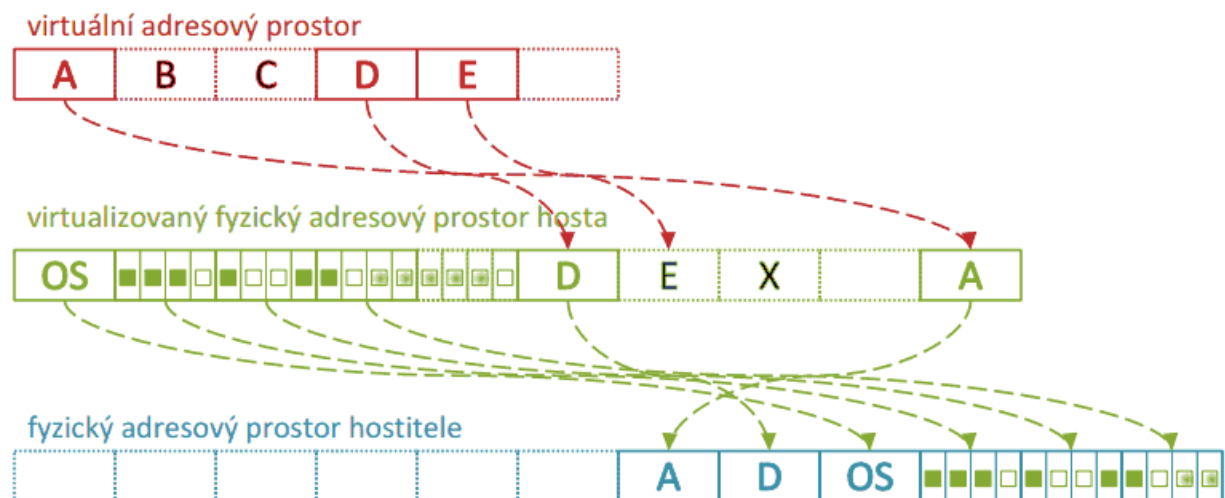
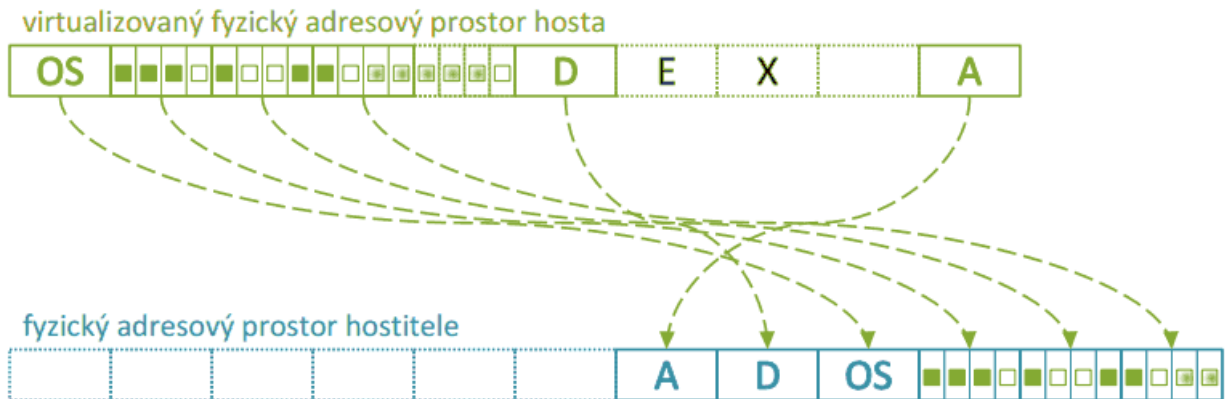
Virtualizace virtuální paměti

- VS ve fyzickém PC
 - adresový prostor z pohledu procesu je jeden nebo více souvislých úseků (rozložení a význam určen dohodou OS a aplikace)
 - dělení na stránky je neviditelné
 - abstraktní pohled na paměť
 - virtuální stránky mapovány na fyzické ramce
 - stránky odložené na disk mapovány nejsou
 - PAS sdílen mnoha procesy
 - 2úrovňové stránkování (x86)
 - při výpadku TLB procesoru se prochází 2 úrovně stránek
 - stránkovací tabulky v PAS
 - fyzická adresa kořene (directory page table) je v registru CPU
 - kód a data OS byvají součástí VAS
 - přepínání stránkování při každém volání OS by bylo neefektivní
 - stránky OS přístupné pouze v privilegovaném režimu procesoru
 - volání OS provedeno spec. instrukcí, která zapíná privilegovaný režim
 - OS plní stránkovací tabulky
 - stránkovací tabulky jsou mapovány podobně jako data OS
 - OS zapisuje do str. tabulek beznými instrukcemi
 - zápis většinou musí být následován privileg. instrukcí TLB flush
- VS ve virtuálním PC
 - VMM poskytuje iluzi PAS
 - mapování PAS guesta na PAS hosta
 - VMM může odkladat stránky na disk podobně jako OS
 - provádění kódu pracuje s VAS guesta
 - potřebné mapování vznikne složením:
 - mapování definovaného OSem guesta
 - mapování definovaného virtualizační PC guesta
 - ve VAS bez:

- aplikacni procesy guesta
 - OS guesta
 - VMM
- potrebujeme 3 urovne privilegii ke strankam
- realizace slozeného mapování
 - pomoci strankovacich tabulek v PAS hosta - obsahují fyzicke adresy v prostoru hosta
 - v systemu jsou dvoji strankovací tabulky
 - strankovací tabulky hosta používane fyzickým CPU
 - obsahují fyzicke adresy v prostoru hosta
 - virtualizované strankovací tabulky guesta používane Osem guesta
 - obsahují virtualizované fyzicke adresy v prostoru guesta
 - VMM počítá výsledné mapování ze strankovacich tabulek guesta
 - OS guesta zapisuje do svých tabulek neprivilegovanými instrukcemi
 - VMM musí zajistit primerenou koherenci fyzických tabulek a tabulek guesta
 - virtualizované tabulky guesta mohou být mapovány read-only a zápisy emulovány
 - koherenci lze udržovat v rámci emulace privilegované instrukce TLB flush
 - slabší VMM neumí odkladání virtualizované paměti na disk
 - mapování virtualizovaných PA na PA je identické
 - VMM pouze kontroluje, zda OS guesta nemapuje nežádoucí PA
 - OS guesta se musí vyrovnat s dírami v PAS
 - používáno převážně při paravirtualizaci (Xen)

virtuální adresový prostor





Data center HW

- motivace
 - standardizace/konsolidace (upevneni)
 - omezit pocet DC v organizaci
 - omezit pocet HW, SW platforem
 - standardizovat vypocetni, sitove a ridici platformy
 - virtualizace
 - konsolidace vice vybaveni DC
 - nizsi kapitalove a provozni vydaje
 - automatizace
 - konfigurace, zaplatovani, vzajemna koordinace, ...

- zabezpečení
 - fyzické, síťové, data, uživatele
- požadavky datových center (serveroven)
 - historicky byly datová centra navrhována bez zavedených standardů, spousta síťových administratorů tak musela celit výběrem použitých technologií a řešením jejich implementace v často omezeném prostoru místnosti
 - 2005 - TIA-942 standard
 - definování norem, které musí centra splňovat z hlediska integrity a funkčnosti výpočetní techniky, která je v nich umístěna
 - definování 4 úrovní data center podle jejich spolehlivosti
 - požadavky pro tier 1
 - jedna neredundantní distribuční cesta pro vedení elektriny a chlazení
 - neredundantní kapacita s dostupností 99.671 %
 - umožňuje aby systém nebyl dostupný max 1729 min. v roce
 - požadavky pro tier 2
 - splňuje požadavky (nebo převyšuje) pro tier 1
 - redundantní kapacita s dostupností 99.741 %
 - umožňuje aby systém nebyl dostupný max 1361 min. v roce
 - požadavky pro tier 3
 - splňuje požadavky (nebo převyšuje) pro tier 1, 2
 - několik nezávislých distribučních cest
 - všechny IT komponenty musí mít napájeny ze dvou zdrojů, navíc plně kompatibilní s topologií infrastruktury site
 - konkurenčně udržitelná infrastruktura s dostupností 99.982 %
 - umožňuje aby systém nebyl dostupný max 95 min. v roce
 - požadavky pro tier 4
 - splňuje požadavky (nebo převyšuje) pro tier 1, 2, 3
 - všechny chladičové komponenty jsou napájeny ze dvou zdrojů
 - infrastruktura s odolností vůči chybám a dostupností 99.995 %
 - umožňuje aby systém nebyl dostupný max 26 min v roce
- problémy
 - design
 - strojírenský design infrastruktury
 - mechanické systémy podílející se na udržování konzistentního vnitřního prostředí
 - topení, ventilace, klimatizace
 - zvlhčování/odvlhčování, tlakování
 - minimalizace prostoru a cen při zachování dostupnosti

- elektrotechnický design infrastruktury
 - distribuce, prepínání, UPS
 - modularne, skalovatelne
- technologická infrastruktura
 - kabelové propojení pro datovou komunikaci, správa počítačů, klávesnice/myši/video snímání technika, ...
- předpoklady dostupnosti
 - pro zajištění větší dostupnosti je nutný větší vstupní kapitál a vyšší provozní náklady
- vyber lokality
 - dostupnost rozvodných sítí, síťové služby, dopravní dostupnost, možnost akutních zásahů
 - klimatické podmínky
- modularita, flexibilita, kontrola teploty, vlhkosti (16-24°C, 40-55% vlhkost), elektrické rozvody (UPS, banky s bateriemi, naftové generátory, ...), ochrana před katastrofami (ohně, záplavy), detektory kouře, automatické ošetrování při požáru, fyzická bezpečnost, ...
- spotřeba energie
 - účinnost využití energie (power usage effectiveness) - PUE = celkový výkon zařízení / výkon veškerého IT vybavení
 - nejmodernější DC mají PUE ~ 1.2
 - analyzování spotřeby energie, chlazení
 - energie činí největší opakující se náklady
 - analýza míst s extrémem - hot spots, over-cooled areas
 - rozmístění vybavení DC
- další aspekty
 - infrastruktura site
 - routery, switchy
 - 2 nebo více poskytovatelů upstream služeb (aby se mohlo vysílat hodně dat)
 - firewally, VPN brány, IDS (intrusion detection system - pro odhalování průniku do systému)
 - řízení DC infrastruktury
 - monitorování apod.
 - aplikace
 - DB, file servery, aplikační servery, zálohy
- blade servery
 - ocesaný server s modulárním designem optimalizovaný pro minimalizaci využití fyzického prostoru a energie
 - narušil od standardního racku, kterému stáčí k funkci napájení a síťový kabel je blade server ocesaný o spoustu komponent pro ušetření místa a minimalizaci spotřeby

- skrin/uzaver (enclosure) se sklada ze spousty blade serveru
 - skrin dostatecne chlazena
 - redundantni zdroje energie
 - chladici system s ridicim systemem starajicim se o regulaci otacek ventilatoru, pouziti kapalinoveho chlazení
 - servery obsahují typicky integrovany NIC pro ethernet nebo radic pro fiber channel rozhrani (pouziti na kroucene dvojlince nebo opticke siti)
 - datova uloziste typicky nejsou umisteny lokálne - pripojeni pres firewire/eSATA/SCSI/iSCSI rozhrani - SAN
- SAN
 - storage area network
 - na zaklade zajmu o oddeleni diskoveho uloziste a procesoroveho vykonu serveru vznikly site SAN
 - narozdil od DAS (directly attached storage), kde je problem se pri poruse serveru jednoduse dostat k datum, spatna rozsiritelnost
 - pro propojeni se pouziva technologie fiber channel, pres opticke kabely, umoznuji propojeni na vzdalenosti desitek az stovek metru
 - propojeni pres FC switch ke spolecnemu datovému ulozisti - typicky diskove pole s vlastní inteligenci (radic RAID)
 - pouzivane technologie:
 - FC
 - vysoke naklady
 - vysoke naroky na obsluhu systemu
 - spis vetsi podniky
 - propustnost 800/1600/3200 MBps
 - ruzne topologie - P2P, switched fabric, arbitrated loop
 - iSCSI
 - posledni dobou se FC nahrazuje iSCSI
 - vychazi ze 2 technologii
 - SCSI rozhrani pro pripojovani disku v serverech
 - technologie TCP/IP
 - ze SCSI si bere pouze protokol, kterým spolu zarizeni komunikují, zcela opousti jeho fyzickou vrstvu (kabely, konektory apod)
 - pri prenosu paketu SCSI se pouzije jejich zapouzdeni do protokolu TCP/IP
 - nepomerne levnejsi nez FC
 - iniciator pripojeni je klient - HW/SW, cilem je prostredek s daty
 - umoznuje network booting
 - FCoE
 - zapouzdeni FC pres Ethernet
 - iSER
 - iSCSI extension over RDMA (remote DMA)

- InfiniBand
- výhody SAN
 - snadné připojení nových diskových polí (u klasických ethernet sítí bez SAN technologií by se musel koupit nový server, což by se prodrazilo)
 - sdílené úložiště (více serverů vidí ten samý logický svazek)
- disková pole
 - diskový úložný systém s několika diskovými zařízeními
 - komponenty
 - ovladač diskových polí
 - cache - RAM, disk
 - skříň/pouzdra na disky
 - zdroje napájení
 - co poskytují
 - dostupnost, udržitelnost, elasticnost
 - redundance dat, hot swap (rychlá výměna disku), RAID
 - kategorie
 - NAS, SAN, hybridní
 - enterprise disková pole
 - vlastnosti navíc
 - automaticky převezme službu při selhání jiného disku
 - výroba snapshotu
 - deduplikace
 - speciální technika komprese dat, která zabraňuje ukládání stejných datových bloků na jednom úložišti
 - replikace
 - tiering
 - technika, která přesouvá datové soubory, oddíly nebo bloky mezi úložišti na základě definované politiky
 - virtuální svazky
 - hot spare
 - disky, které umožňují rychle nahradit funkci poškozeného disku
 - úroveň RAIDu
 - Redundant Array Of Independent Disks
 - proc
 - zvýšení dostupnosti
 - snížení MTBF (mean time between failure)
 - snížení MTTR (mean time to repair)
 - zvýšení výkonu
 - JBOD
 - Just a Bunch Of Disks
 - reálně žádná RAID funkce
 - min # disků = 1

- prostorova efektivita = 1
- odolnost vuci chybam = 0
- mira selhani diskoveho pole = $1 - (1-R)^N$
 - N = # disku
 - R = mira selhani disku
 - priklad: 3 disky, kazdy z disku ma miru selhani 5%,
pouzivaji JBOD, tak mira bude $1 - (1-0.05)^3 \sim 0.14 \sim 1.4\%$
- vykonnost cteni = 1
- vykonnost zapisu = 1

■ RAID 0

- striping - segmentace logicky sekvencne ulozenych dat na ruzna diskova zarizeni
- bez redundance
- min # disku = 2
- prostorova efektivita = 1
- odolnost vuci chybam = 0
- mira selhani diskoveho pole = $1 - (1-R)^N$
- vykonnost cteni = N
- vykonnost zapisu = N
- → pro zvyseni vykonu, kdyz nejde moc o integritu
 - pouziva se u nekterych hernich systemu
 - nicmene ukazalo se ze striping nemusi vzdy pridat vykon, nekdy naopak muze i snizit, navic pri testovani PC her neprineslo moc navyseni, ale u desktop systemu se muze hodit

■ RAID 1

- zrcadleni obsahu
- min # disku = 2
- prostorova efektivita = $1/N$
- odolnost vuci chybam = $N-1$
- mira selhani diskoveho pole = R^N
- vykonnost cteni = N
- vykonnost zapisu = 1
- → pro zvyseni spolehlivosti, rychlost cteni

■ RAID 2

- striping na urovni bitu na zaklade pouziti Hammingova kodu pro zjistení parity
- min # disku = 3
- prostorova efektivita = $1 - 1/n * \log_2(n-1)$
- odolnost vuci chybam = 1
 - vzpamatuje se z kleknuti jednoho disku
- mira selhani diskoveho pole = promenna
- vykonnost cteni = promenna, vykonnost zapisu = promenna

- → uz se nepouziva - zbytocna redundantni kontrola chyb a velka komplexita
- RAID 3
 - striping na urovni bytu na zaklade parity
 - jeden disk je vyhrazen pro ukladani informaci o parite
 - pr. 4 disky, 4. vyhrazen pro ukladani parity, mame urovni bytu → i-ty byte na 4. disku odpovida parite i-tych bytu na discich 1-3
 - min # disku = 3
 - prostorova efektivita = $1 - 1/N$
 - odolnost vuci chybam = 1
 - mira selhani diskoveho pole = $N * (N-1) * R^2$
 - vykonnost cteni = $N-1$, vykonnost zapisu = $N-1$
 - obecne nemuze zpracovavat vice pozadavku naraz, ale hodi se na vysokorychlosti cteni/zapis dlouhych souvislych useku dat - napr. editace nekomprimovaneho videa
- RAID 4
 - striping na urovni bloku na zaklade parity
 - min # disku = 3
 - prostorova efektivita = $1 - 1/N$
 - odolnost vuci chybam = 1
 - mira selhani diskoveho pole = $N * (N-1) * R^2$
 - vykonnost cteni = $N-1$, vykonnost zapisu = $N-1$
 - moc se nepouziva
- RAID 5
 - striping na urovni bloku na zaklade distribuovane parity
 - informace o parite je narozdil od minulych pripadu distribuovana mezi disky
 - min # disku = 3
 - prostorova efektivita = $1 - 1/N$
 - odolnost vuci chybam = 1
 - mira selhani diskoveho pole = $N * (N-1) * R^2$
 - vykonnost cteni = $N-1$, vykonnost zapisu = $N-1$
- RAID 6
 - rozsireni RAID 5 o dodatecny blok s paritou - striping s dvema distribuovanymi paritovymi bloky
 - min # disku = 4
 - prostorova efektivita = $1 - 2/N$
 - odolnost vuci chybam = 2
 - mira selhani diskoveho pole = $N * (N-1) * (N-2) * R^3$
 - vykonnost cteni = $N-2$
 - vykonnost zapisu = $N-2$

- hybridni
 - RAID 0+1 - striped sets in mirrored set
 - RAID 1+0 (RAID 10) - mirrored set in striped sets, kazde zrcadleni muze ztratit disk a nic se nestane
 - RAID 5+0 (RAID 50) - block striping s distr. paritou in striped set
 - odolnost - jeden disk muze spadnou v kazdem RAID5 bloku

Virtualizační infrastruktura (VMWare)

- fyzicky stroj
 - HW
 - CPU, RAM, disky, I/O
 - nevyužitý HW naplno
 - SW
 - jeden aktivni OS
 - OS kontroluje HW
- virtualni stroj
 - abstrakce na HW urovni
 - virtualni HW ~ CPU, RAM, disky, I/O
 - virtualizacni SW
 - oddeluje HW a OS
 - multiplex fyzickeho HW napric nekolika guest VM
 - silne izolovane VM mezi sebou
 - ridi fyzicke prostredky, optimalizuje vyuziti
 - izolace
 - bezpecny multiplexing
 - nekolik VM na jednom fyzickem hostu
 - CPU HW izoluje mezi sebou VM (MMU)
 - silne zaruky
 - SW buggy, pady v ramci jedne VM nemuzou ovlivnit ostatni
 - oddeleni vykonu
 - rozdeleni systemovych prostredku
 - zapouzdeni
 - kazda VM je jeden soubor
 - OS, aplikace, data
 - pamet a stav zarizeni
 - moznost vytvoreni snapshotu, klonu
 - zachytit stav VM za behu a obnovit do nejakeho zachytneho bodu
 - zalohovani, zrcadleni

- jednoduchá distribuce obsahu
 - predkonfigurované aplikace
 - virtuální zařízení
- kompatibilita
 - HW nezávisle
 - fyzický HW skrytý virtualizační vrstvou
 - standardní virtuální HW vystaven každé VM
 - jednou vytvořit, pustit kdekoliv
 - žádné problémy s konfigurací masin
 - migrace VM mezi hosty
 - podpora zastaralých systému
 - můžeme si ho pustit na nové platforme
- správa prostředku (distributed resource scheduler)
 - “boj o prostředky”
 - kontrola alokace prostředku na různých úrovních
 - na úrovni sdílení
 - specifikace relativní důležitosti mezi VM - podíl na výkonosti
 - nároky přímo úměrně podílu
 - VM má nasdíleno dvakrát víc než jiná, je označena, že zere 2x víc
 - abstraktní relativní jednotky - důležitost je pouze poměr
 - low/normal/high
 - 500/1000/2000 shares per virtual CPU
 - na základě rezervování
 - dána garance minimální alokace i když je systém zahlcen
 - máme 2GHz procesoru k dispozici, když nastavíme 1GHz garanci pro proc. 1 a 2, tak oba mají jistotu, že vždy dostanou aspoň 1GHz když chtějí
 - pokud jeden z nich ale využívá jen 500Mhz, druhý může 1.5GHz
 - v konkrétních absolutních jednotkách
 - řízení přístupu: suma rezervací ≤ kapacita
 - default = 0
 - na základě limitu
 - stanovuje horní mez, kolik může VM spotřebovat, i když je systém zahlcen
 - v konkrétních absolutních jednotkách
 - default = inf

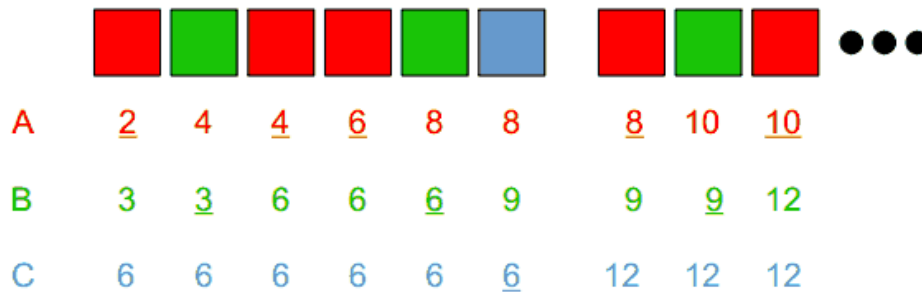
- Proportional-share scheduling

- Simple virtual-time algorithm

- Virtual time = usage / share

- Schedule VM with smallest virtual time

- Example: 3 VM A, B, C with 3:2:1 share ratio



- využití distribuovaného systému

- vyber inicialního hosta když se VM zapne

- předtím, než se začne něco vykonávat nemáme žádný reálný odhad využití prostředku → vybíráme jen na základě nastavení stroje (předpokládá se, že využije vše, co si nastavil)

- migrace bezcíh VM mezi fyzickými hosty

- rychlá migrace, transparentní guest OSu i aplikacím
 - minimální čas nedostupnosti ($\leq 1s$)
 - požadavky:

- sdílené uložení
 - na stejné podsíti
 - kompatibilní CPU

- aby měl CPU, kam se migruje stejnou podmnožinu instrukcí

- nadefinuje administrator
 - instrukce, kterou se tážeme procesoru na jeho instrukční sadu není volána přímo, ale použije se mechanismus trap and emulate

- dynamický load balancing

- uvažujeme dynamický entitlement - vypočítáváme na základě celkové kapacity clusteru, používaných zdrojů (prostředků) a aktuálního vytížení CPU a paměti na jednotlivých VM

- algoritmus:

- máme snapshot celého clusteru (VMs, hosts)
 - spočítáme normalizovaný entitlement pro každého hosta ~

$N = \sum E_i/C_j$, kde E_i jsou požadavky jednotlivých VM na hostu j a C_j je kapacita hosta j

- tedy napr. host 4GHz, bez tam VMs A = 3GHz, B = 2GHz → entitlement = $5/4 = 1.25$

- spočítáme std. odchylku nároků přes všechny hosty
- pokud je odchylka větší než nějaký threshold (prostě jsou moc zatížené), pokračujeme dál
- pro každou VM v clusteru spočítáme, jak je výhodný přesun na nějakého kompatibilního hosta - jak se zlepší odchylka
- vybereme nejvyhodnější přesun
- takto můžeme vybrat několik přesunů, do nějakého limitu, který je nastaven (konstanta max migrations)

- distribuovaný power management (DPM)

- konsolidace VM na méně hostů a vypnutí hostů, pokud je nízká poptávka
- zapnutí hostů znovu pokud je nutné splnit požadavky vytížení procesoru nebo splnit nastavené podmínky, omezení
- pracuje ve shodě s distribuovaným plánovačem prostředků (DRS)

- distribuovaná vysoká dostupnost

- na základě požadované kapacity k rezervaci restartuje danou VM po selhání jejího hosta na jiném hostu
 - definuje se počet padů hosta, které se ještě tolerují
 - jaké procento kapacity clusteru host zabírá
 - počet hostů, které jsou nachystány bokem pro převzetí kontroly po pádu
- pravidelně a s použitím inteligentních mechanismů monitoruje využití dostupné kapacity
- chytí převzetí funkce poškozeného stroje serverem s nejlepšími dostupnými prostředky
- decentralizovaná detekce poškození hosta
 - hosti v clusteru si navzájem posílají informace o jejich tíkání, když jeden z hostů v nějakém timeoutu neodpoví, spustí se příslušná akce, která restartuje daný stroj
- funguje ve shodě s distribuovaným plánovačem prostředků a distribuovaným power managementem
- DRS/DPM/HA aktivně spravují příslušné volné prostředky

- distribuovaný I/O management

- odolnost vůči chybám

- automatická detekce chyb stroje, okamžité vyvolání triggeru, proběhne neviditelně převzetí funkce poškozeného stroje jiným
- vyvolání vytvoření sekundárního VM na jiném hostu ihned po failoveru pro zajištění nepřetržité ochrany aplikace, tyto pak sdílí prostředky na SAN

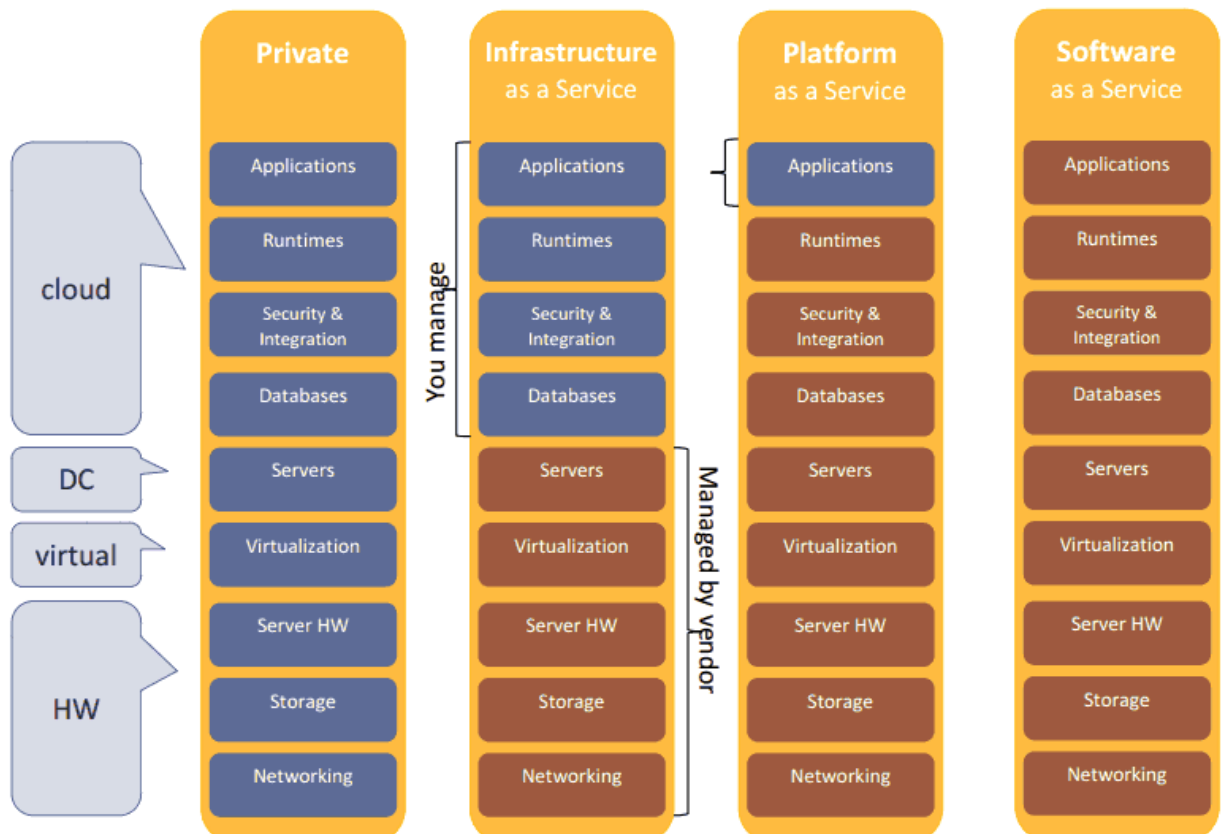
- primarni VM pak funguje normalne a sekundarni provadi pouze read operace
- pokud pak nekdy host pro primarni VM spadne, sekundar ho nahradi bez jakéhokoliv preruseni (~ 5 ms !)
- zajisteni dostupnosti pri zachovani vysoke rychlosti
 - dnesni CPU mohou replikovat proud instrukci
- snapshots/checkpoints
 - dump celeho stavu do stroje do zalozniho souboru
 - VM settings, VM disks content, VM memory content, registry CPU
 - budouci zapisky jsou presmerovany do druheho souboru, který pak funguje systemem copy-on-write → inkrementalni
 - prostě si nejprve udelam cely snapshot a pak si uchovavam jenom nejaky zurnal zmen
- podpora cloud-computingu, multi-tenancy

Cloud computing

Cloud Computing Components

<i>Execution Models</i>	Virtual Machines	Web Sites	Cloud Services
<i>Cloud Storage</i>	SQL Database	Key-Value Tables	Blobs
<i>Data Processing</i>	Map/Reduce	Hadoop	Reporting
<i>Networking</i>	Virtual Network	Connect	Traffic Manager
<i>Messaging</i>	Queues	Service Bus	Service Bus
<i>Caching</i>	Caching	Content Delivery	
<i>Hi-Perf Computing</i>	Scheduler	Load Balancing	
<i>(Multi-)Media</i>	Media Services	Streaming	
	GIS	Searching / Indexing	Marketplace
<i>Other Services</i>	Mobile / E-mail	Language / Translate	Collaboration
	Gaming	Prediction	
<i>SDK</i>	C++ .Net Java PHP Python Node.js ...		

- nabízí se různé úrovně služeb
- IaaS (infrastructure as a service)
 - základní cloudový model
 - poskytovatelé nabízejí počítače - fyzické/virtuální, servery a další prostředky
 - mezi dalšími prostředky patří například úložné soubory, firewally, IP adresy, virtuální LAN site, SW balíky
 - určeno pro síťové architektury
 - Windows Azure VMs
- PaaS (platform as a service)
 - dodání počítačové platformy typicky zahrnující OS, prostředí pro programování, jazyky, databáze, webový server
 - vývojáři aplikací mohou vyvíjet a spouštět své aplikace na cloudové platformě bez toho, aby museli kupovat a nastavovat HW a SW vrstvy, které pod platformou leží
 - Windows Azure, Google App Engine
 - určeno pro vývojáře aplikací
- SaaS (software as a service)
 - přístup k aplikacím SW a databázím
 - "on-demand" SW, typicky se za to platí na bázi "pay-per-use" (monthly/yearly fee per user)
 - určeno pro koncové uživatele
 - Grooveshark, Last.FM, Rapidshare, online hry apod.



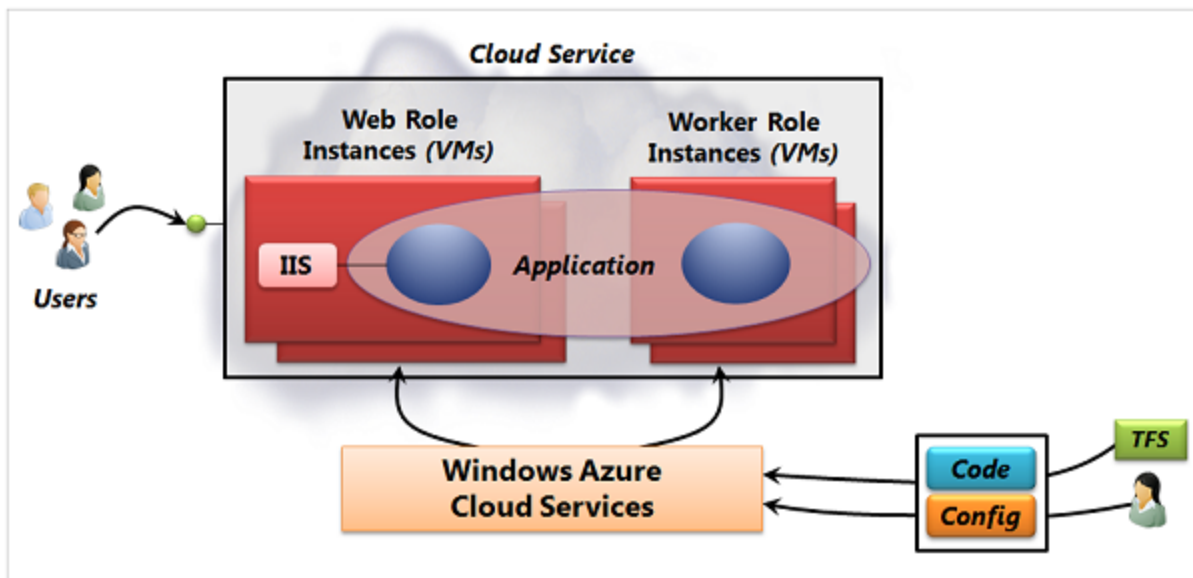
- execution models

- uroven poskytovani sluzeb - nemusim se starat o nic vs muzu cokoliv
- virtualni stroje
 - mame k dispozici obrazy/galerii snych VM
 - dalsi cloudove sluzby - DB, messaging
 - pay-per-use, pay-per-config - memory, processors, disk space
 - administrace pres webovy portal, skriptovaci konzole, API (REST)
 - bezici VM
 - virtualni disky - OS a data
 - typicky BLOBy
 - vyuziti
 - development, test environment
 - beh aplikaci, provoz sluzeb
 - rozsireni datacentra
 - disaster recovery
 - REST (representational state transfer) API a scripting
 - architektura rozhraní, navržená pro distribuované prostředí
 - jednotny pristup ke zdrojum - CRUD (create, read, update, delete)
 - pristup k sluzbam i datum
 - `http://My.tables.MyCloud/<MyTable>(PartitionKey='PK',Ro`

wKey='RK')?\$select=<FirstProperty>

- vsechny zdroje maji vlastni identifikator URI a ke vsem jsou definovane CRUD operace
- bezstavovy protokol typu klient/server
- format vymeny dat - JSON, XML
 - pres HTTP metody get, put, post, delete
- IaaS
 - odolnost proti chybam, monitoring
 - viz virtualizace
 - replikace pres BLOBy
 - v BLOBech ulozeny i disky VM
 - zachovani persistence
 - automaticka aktualizace disku s operacnimi systemy
 - pouze u VHD, ktere poskytuje MS
 - u vlastnich obrazu si aktualizujeme sami
- grouping
 - pri vytvoreni nove VM mame moznost standalone behu nebo ho pripojit do skupiny ostatnich VMs v cloudu
 - kazdy standalone VM ma svoji vlastni verejnou IP narozdil od skupiny, ktera cela vystupuje pod jednou
 - load balancing - rozlozeni pozadavku
 - availability set - rozlozeni na ruzne uzly - ochrana proti havarii
 - komunikace v lokalni siti v ramci skupiny
- webové stránky
 - je to jednodussi a levnejsi nez spravovat virtualni stroje
 - k dispozici je spousta frameworku, knihoven a aplikaci, na kterych se da stavet
 - JAVA, ASP.NET, MySQL, PHP, ...
 - Drupal, WordPress, Dropbox, Joomla, ...
 - nejcastejsi vyuziti cloudovych infrastruktur
 - idealni podminky pro nasazeni v cloudu
 - skalovatelnost
 - lze i pomoci VMs
 - avsak zbytecne slozite
 - nutna vlastni instalace, konfigurace, udrzba
 - PaaS
 - predkonfigurovatelne instalace
 - komplettni framework - OS, DB, web server, knihovny, ...
 - administrace, aktualizace, udrzba vseh komponent
 - dynamicke pridavani instanci, load balancing
 - ruzne urovne izolace - shared / private VM
- cloudové služby
 - poskytovani skalovatelnych SaaS sluzeb

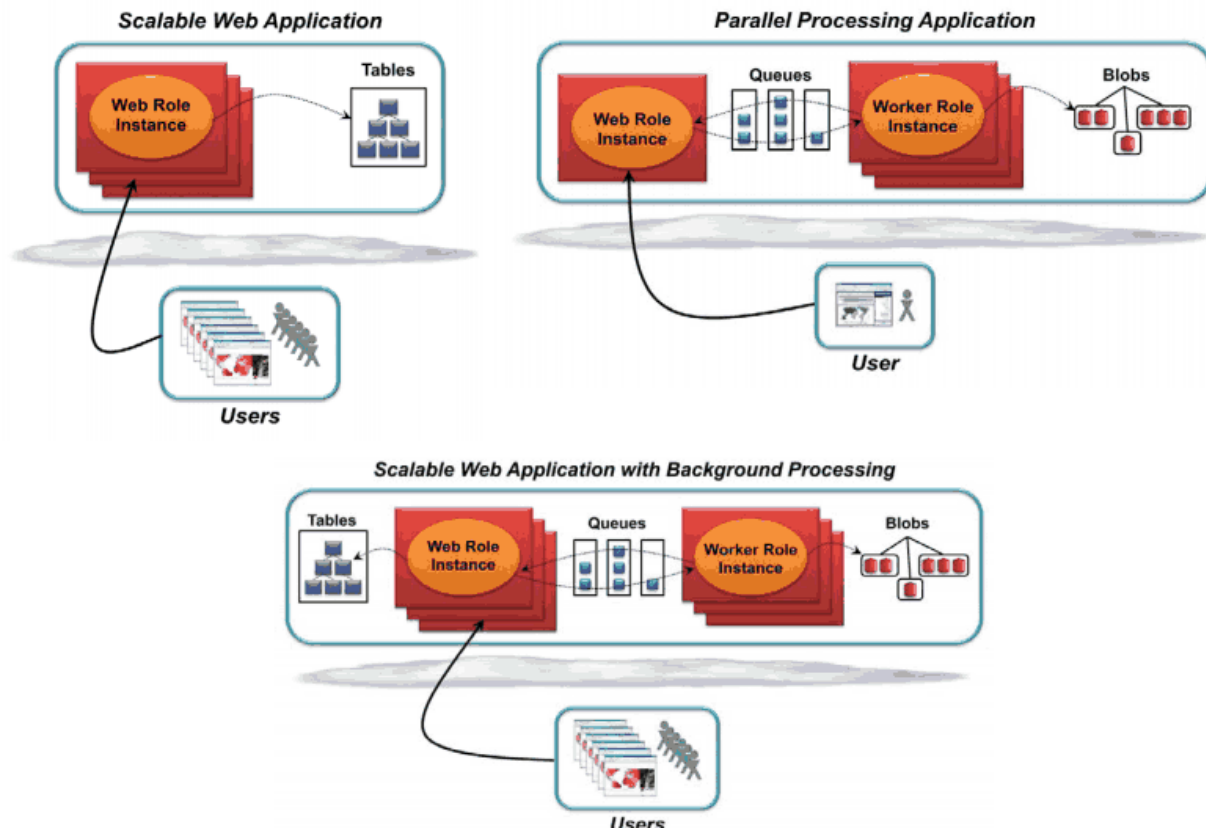
- web sites - omezeny pristup, nelze nainstalovat cokoliv, nejsou dostupne vsechny sluzby
- virtual machines - plny pristup, avsak nutnost administrace a udrzby, samo neskaluje
- PaaS
 - skalovatelnost, dostupnost, spolehlivost, udrzba, ruzne jazyky a platformy
- web roles / worker roles
 - technologie nabizi 2 ruzne pristupy VM
 - web roles - VMs instance, na kterych bezi varianta Windows Serveru s IIS
 - IIS - SW webovy server s velkou skalou rozsirujicich modulu, nejpouzivanejsi po Apache
 - worker roles - VMs instance, kde bezi stejná varianta Windows jako u web roles, ale bez IIS
 - cloudove sluzby se pak sestavaji z nejake kombinace techto



- vhodne pro vicevrstve skalovatelne aplikace
- jedna adresa, vice web/worker roli
 - automaticky load balancing a availability set
- vyvojove a produkcní prostredi
 - jednoduchy vyvoj/testovani a prechod na novou verzi
- monitorovani (nejen) virtualniho HW
 - agent uvnitr web/worker role
- cloudove sluzby vs webove stranky
 - administrativni pristup do VM - instalace libovolneho potrebného SW
 - web/worker role - vicevrstve aplikace - vlastní VM pro aplikacni logiku
 - možnost propojeni cloudove aplikace s privatnimi uzly

- remote desktop pro primy pristup k aplikacni VM

Cloud Services - scénáře

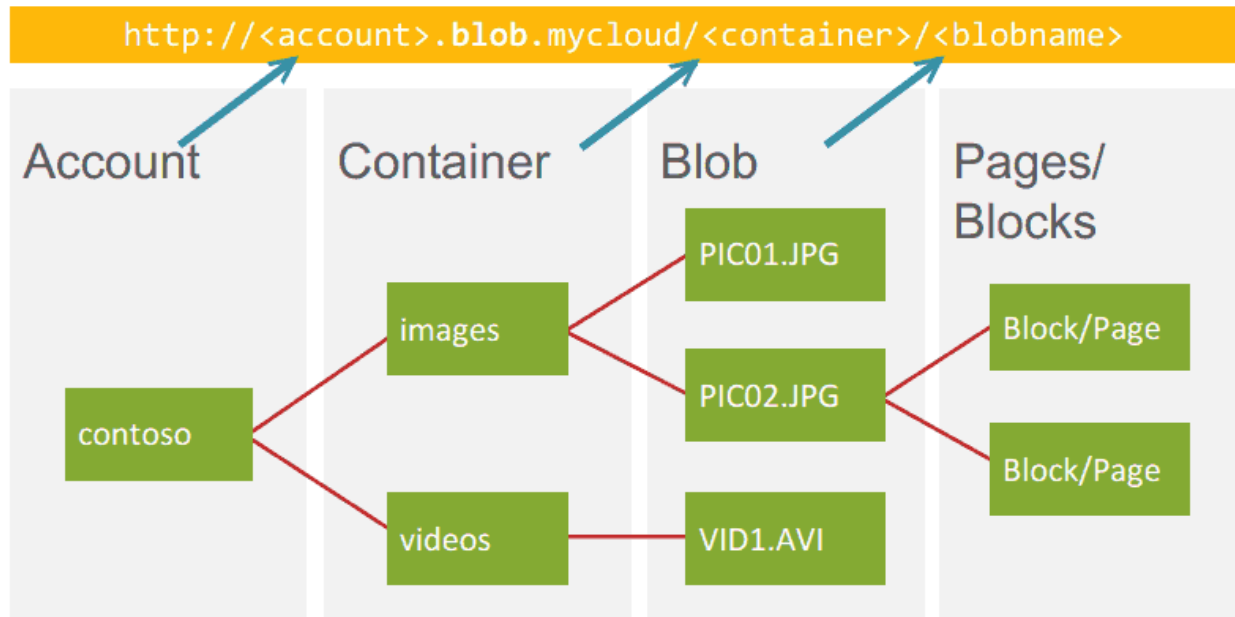


Windows Azure use cases

- high-performance computing
 - hlavní uzel přijima pracovní requesty a rozesílá je worker uzlům
 - laboratorní simulace, výzkum leku apod.
 - neočekávaná / periodická potřeba výkonu
 - pravidelně chceme něco zaznamenávat
 - neočekávané chceme spočítat velký náklad dat apod.
 - investice do provozu
 - paralelní problémy
 - cluster - privatní / cloud / kombinovaný
 - používá se job scheduler (DRS), cluster manager (nějaký GUI pro konfiguraci služeb clusteru)
 - problémy
 - privatní a citlivá data
 - závislost na externím SW
 - přenos velkých dat
 - data management

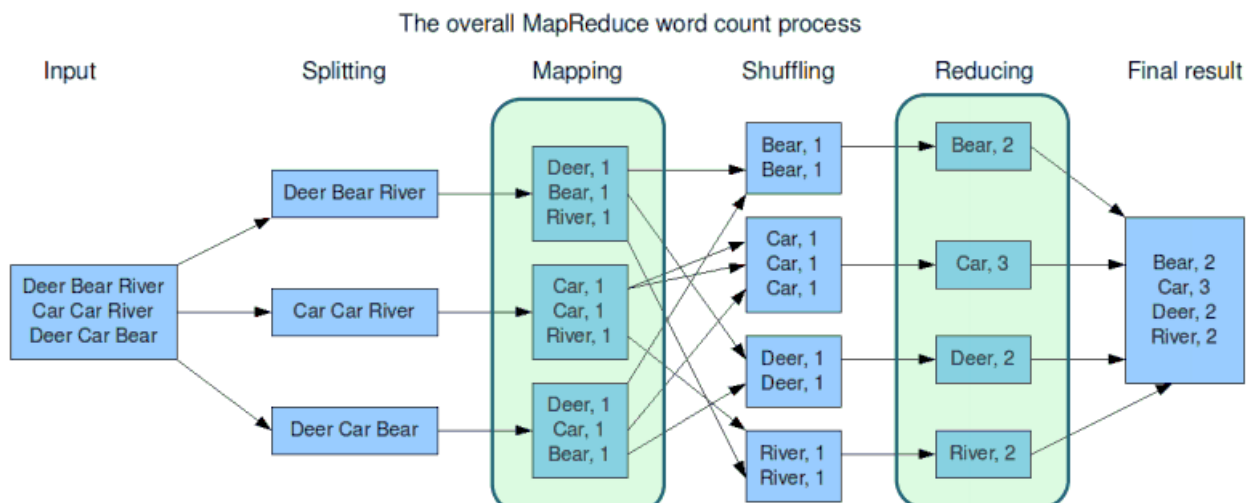
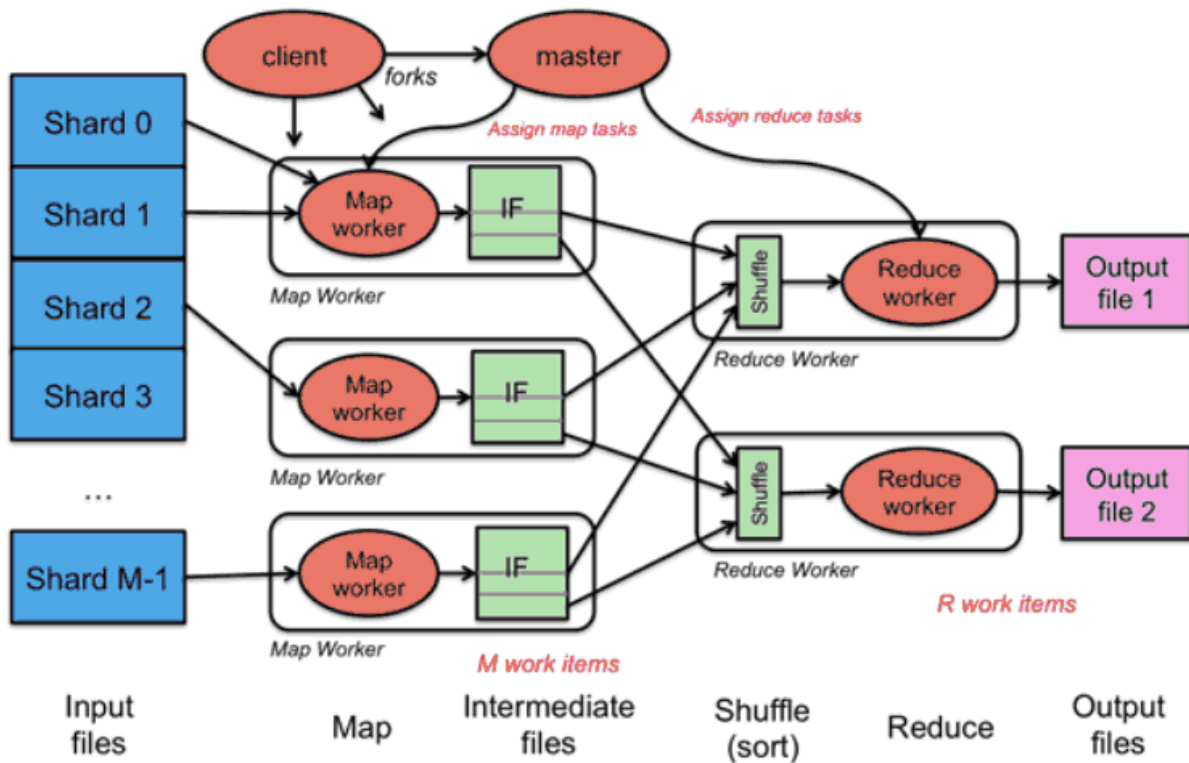
- SQL
 - různé implementace a knihovny - JDBC, ADO.NET, ...
 - není to jen DBMS v cloudu - je to PaaS
 - DB se stará o administrativní údržbu - řízení HW infrastruktury, automatická aktualizace DB a OS SW, distribuce dat na více serveru - replikace
 - přístupné pro cloud i externí aplikace
 - typicky distribuované řešení - poskytuje availability a scalability
 - primární / sekundární replika
 - load balancing
 - cluster index
- NoSQL
 - když nepotřebujeme komplexní SQL dotazování
 - nejde o relační databázi
 - tabulky klic-hodnota
 - velmi rychlý škálovatelný přístup k velkým (řád TB) typovaným datům
 - identifikace - key, partition key/row key
 - nejdou složité dotazy a třídění dle více klíčů
 - struktura
 - account, table, partition, entity, property (ta má pak třeba jméno, typ, hodnotu)
 - table
 - množina řádek
 - jednoznačné id - partition key/row key
 - partition key - rozlišení entit mezi různě replikovanými uložišti
 - row key - jednoznačné id v rámci partition
 - třídění jen dle PK/RK
 - entity
 - množina properties (atributů)
 - vždy PK/RK + timestamp
 - další atributy aplikačně definované
 - bez pevného schématu
 - různé entity → různé atributy
 - partitions
 - operace na jedné partition jsou atomické
 - rozložení dat mezi uzly řízeno aplikací
 - volba PK/RK
 - související entity na jednom uzlu - shodný PK
 - rovnoměrné rozložení entit mezi partitions - různé PK
 - efektivní dotazy - filtr přes PK
 - příklad: blogovací server
 - témata, příspěvky v tématech
 - partition key ~ téma

- row key ~ datum, čas príspevku
 - ďalšie atribúty ~ text, hodnotenie
 - efektívny dotaz je napr. nájdenie najaktualnejších príspevkov v tematu
- príklad: verzované uloženie dokumentu
 - PK: meno súboru
 - RK: verzia súboru
 - ďalšie atribúty: dátum, ID, poznámka, ...
 - vlastné údaje typicky v BLOB uložení
- column stores
 - sloupce orientovaná architektúra
 - riadky - veľmi mnoho hodnôt - riadky milióny
 - adresácia - column family / row / column → value
 - zadné schéma - možné rôzne počty a typy hodnôt
 - verzovanie hodnôt, timestamps, možnosť prístupu ke starším verziam
 - replikácia, sharding (horizontálne partitioning na jednotlivé shardy), rýchly prístup, bez transakcií, rôzne úrovne konzistencie dát
 - využitie v dátových skladoch, CRM (customer relationship management)
 - Cassandra (Apache open source distributed DBMS), C-store
 - ako sa udáva sloupceová architektúra
 - rozpad bežnej riadkovej tabuľky na sloupce
 - potom máme pre každý sloupec tabuľku ID, value
 - urobíme ešte RLE kompresiu, aby sa v slupci neopakovali hodnoty - stačí nám vedieť, pre aký rozsah ID sú rovnaké
 - potom sa to serializuje do súboru
- BLOBs
 - neštruktúrovaná, veľká (rády TB) binárna data
 - sdružená v replikovaných kontajneroch - meno kontajneru ako PK
 - integrácia s file systemom, persistence
 - ukladanie videí, virtuálnych diskov s VM, backup → je to lacné
 - typy
 - block blob
 - sekvencie blokov - až milióny
 - optimalizované pre sekvencný prístup - streamy
 - rýchly zápis (append) a sekvencný čítanie
 - page blob
 - stránkové orientované
 - náhodný read/write prístup
 - vyššia rýchlosť
 - použitie s RESTful API - putblob, getblob, deleteblob, copyblob, ...



- business analytics - SQL reporting
 - reporty, analiza dat spolecnosti
 - vydaje, prijmy v kvartale, vytizeni pracovniku, ...
 - data mining nad SQL / NoSQL
 - pouziti rozhodovacich stromu, lesu
 - regresivni analiza (statisticke metody, kdy odhadujeme hodnoty nahodne veliciny na zaklade znalosti jinych velicin)
 - clustering
 - ruzne typy vystupu - grafy, sestavy, agregace, XML, PDF, ...
- Hadoop - široky system propojenych produktu ("ZOO")
 - BigData, nezpracovatelná rozumne jednim DBMS
 - nejsou treba ani relacni data - logy, raw data apod.
 - paradigma Map / Reduce - paralelizace a distribuovanost vypoctu
 - dalsi soucasti a souvisejici moduly
 - HDFS - hadoop distributed file system
 - Pig, PigLatin - paralelizacni platforma, high-level jazyk
 - Hive, HiveQL - data warehouse & querying, analiza dat
 -
 - map/reduce
 - paradigma pro paralelizaci a distribuovanost
 - map: input → (key,value)
 - reduce: (key, list of values[]) → output
 - sjednoceni hodnot pro kazdy klic
 - zpracovani, vysledek
 - mezivysledky mohou byt duplicitni
 - zpracovani velmi velkych objemu dat

- odolnosť proti chybam
- aplikácie - frekvencie slov
 - spočítat celkovú frekvenciu slov v mnohých dokumentoch
 - $\text{map}(\text{file}, \text{text}) \rightarrow (\text{word}, \text{count})$
 - $\text{reduce}(\text{word}, \text{list of counts}) \rightarrow \text{count}$



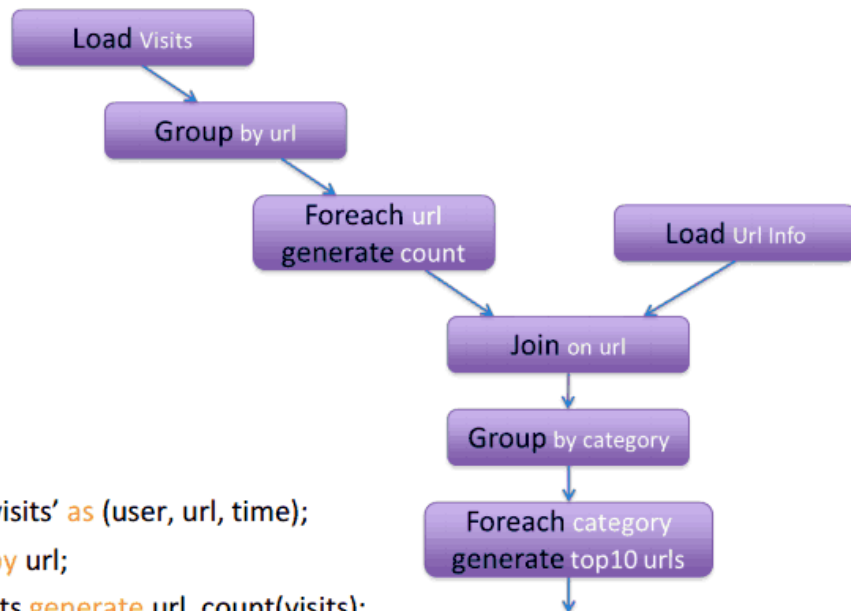
- master/slave architektúra
 - jobtracker - single master server

- rozhrani pro klienty
 - fronta jobů, zpracování FIFO, přiřazení jobů tasktrackerům
 - tasktrackers - slave servers, jeden na každý uzel v clusteru
 - vykonávání úkolu dle jobtrackeru
 - přesuny dat mezi fázemi map, reduce
- výhody - velká nestrukturovaná data, batch/offline režim
- nevýhody - nevhodné pro malá data, komunikace a synchronizace mezi uzly, vyšší latence, transakčnost
- co zahrnuje ZOO
 - HDFS - optimalizovaný pro Hadoop
 - HBase - NoSQL DB, distributed, column-oriented, cílem je ukládat biliony řádků, miliony sloupců
 - Pig, PigLatin - vysokourovnňový jazyk pro analýzu dat, překlad na map/reduce joby
 - HIVE, HIVEQL - data warehouse infrastruktura, SQL-like jazyk pro dotazy nad HDFS
 - Mahout - framework pro strojové učení a data mining
 - Hama - distribuovaný (grid) výpočetní framework
 - zookeeper - high-performance koordinátor distribuovaných aplikací
 - naming, konsensus, členství ve skupinách, výběr koordinátora, fronty, event notification, workflow & cluster management, ...
- HDFS
 - rádově 10 000 uzlů, milion souborů, PB dat
 - uložení strukturovaných i nestrukturovaných dat
 - vlastnosti
 - odolnost proti chybám HW
 - replikace, detekce chyb, zotavení
 - optimalizace pro velké množství levného HW
 - streamování souborů
 - optimalizace propustnosti, nízké přístupové doby
 - orientace na velké soubory
 - velká data uložena v menším množství velkých souborů
 - write-once, read-many
 - jednou zapsaný soubor již není možné změnit ani appendovat k němu data
 - podpora pro append už sice existuje, ale jen ve spec. případech
 - heslo: moving computation is cheaper than moving data
 - architektura
 - NameNode

- master, spravuje metadata
 - obsahuje seznamy souboru, adresaru, jejich mapovani na bloky a umisteni bloku
- DataNode
 - slave, ulozi bloky, operace zapisu a cteni
 - nevi nic o souborech
 - NameNode mu muze naridit blok zreplikovat nebo smazat
- operace se souborovym systemem
 - cteni: NameNode zjist umisteni bloku souboru, klient komunikuje primo s DataNode
 - zapis: NameNode zalozi soubor a rozhodne, jaky DataNode pouzit
- replikace
 - v atributu souboru se uchovava info o # pozadovanych replik
 - na jednom stroji vzdy jedna replika souboru
 - odolnost proti vypadku disku i celeho stroje
 - zapis pomoci odlozene replikace
- implementace ruznych FS na zaklade rozhrani HDFS/HDFS(hadoop cluster FS)
 - Azure BLOB storage, CassandraFS, Symantec cluster FS, ...
- Pig + PigLatin
 - nevyhody Map/Reduce: prilis low-level, nepodporuje slozitejsi data-flow
 - Pig je behove prostredi (platforma) pro webově skalovane zpracovavani dat
 - provadi transformace (kompilace) na soustavu MapReduce programů/jobů
 - BigData - davkove zpracovani
 - pouziva se vysokourovnovy jazyk PigLatin pro datovou analyzu
 - paralelni zpracovavani dat
 - operace pro manipulaci s relacnimi daty
 - imperativni styl programovani (vzor Java)
 - pigs eat anything
 - umi operovat nad daty obsahujici/neobsahujici metadata, nad relacnimi/nested/nestrukturovanymi daty
 - pigs live anywhere
 - jazyk pro paralelni zpracovani dat, neni spjatý s zadnym konkrétnim frameworkem
 - pigs are domestic animals
 - designovane pro jednoduche rizeni a modifikaci uzivateli

- podpora integrace uzivatelskych kodu
 - snadna uzivatelska rozsireni, ktere mohou uzivatele psat v Jave, Pythonu, Javascriptu, Ruby, Groovy a pak volat primo z PigLatinu
- pigs fly
 - rychle zpracovani dat
 - trvale zlepsovani vykonnosti

Pig Data Flow & Pig Latin



```

visits      = load '/data/visits' as (user, url, time);
gVisits     = group visits by url;
visitCounts = foreach gVisits generate url, count(visits);
urlInfo     = load '/data/urlInfo' as (url, category, pRank);
visitCounts = join visitCounts by url, urlInfo by url;
gCategories = group visitCounts by category;
topUrls     = foreach gCategories generate top(visitCounts,10);
store topUrls into '/data/topUrls';
  
```

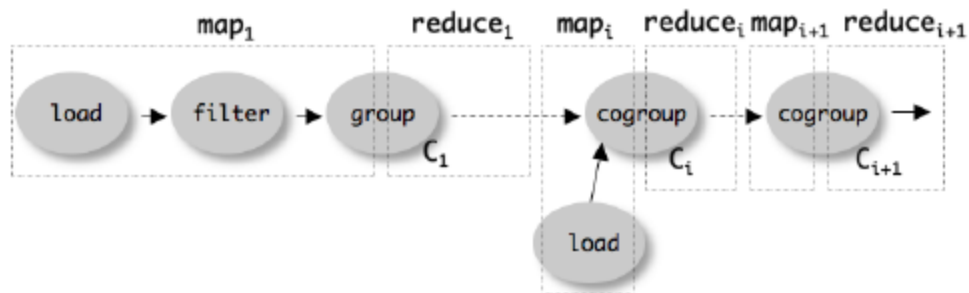
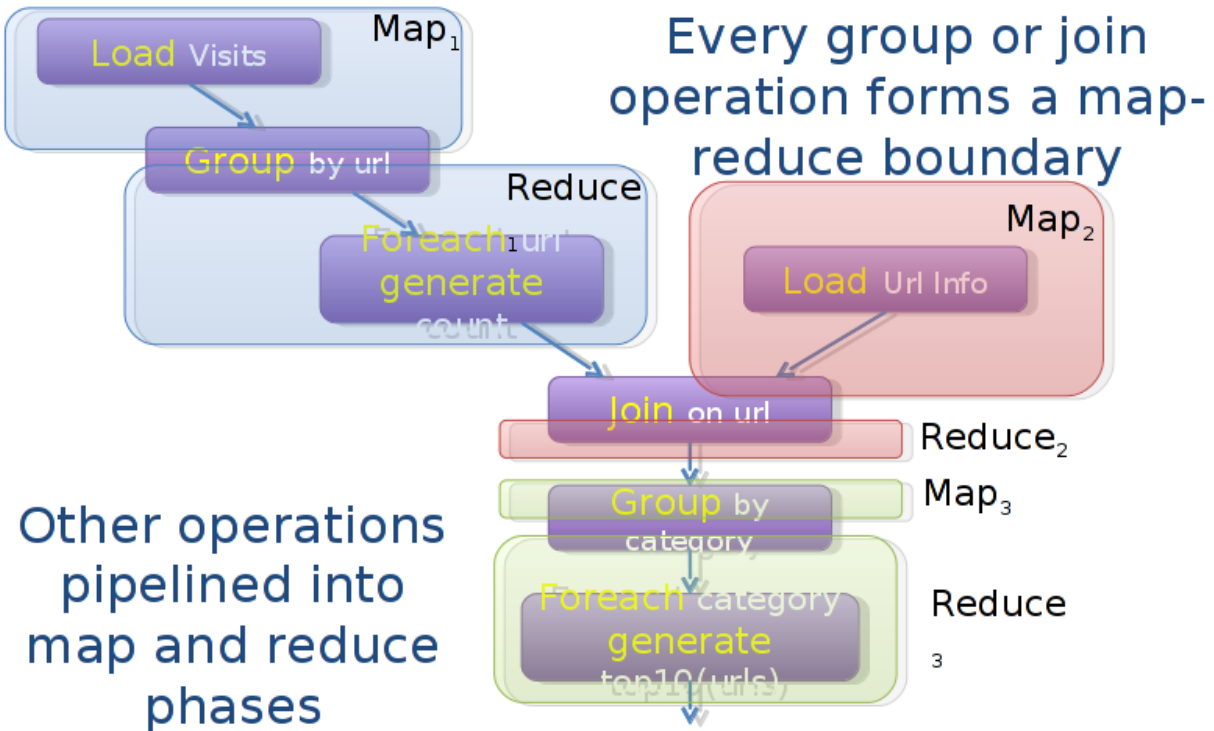
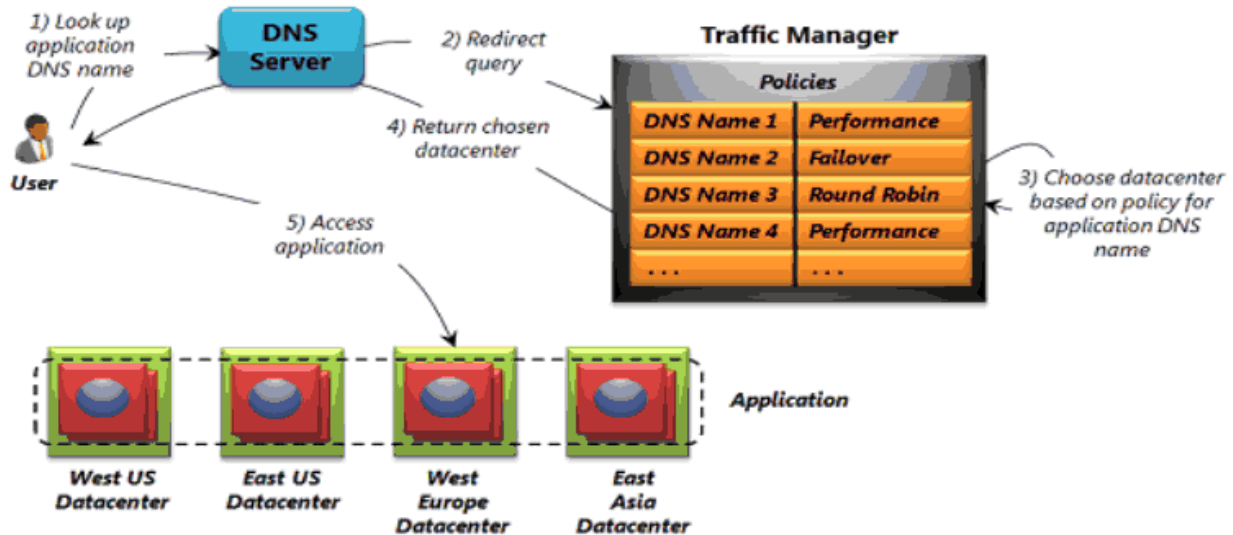


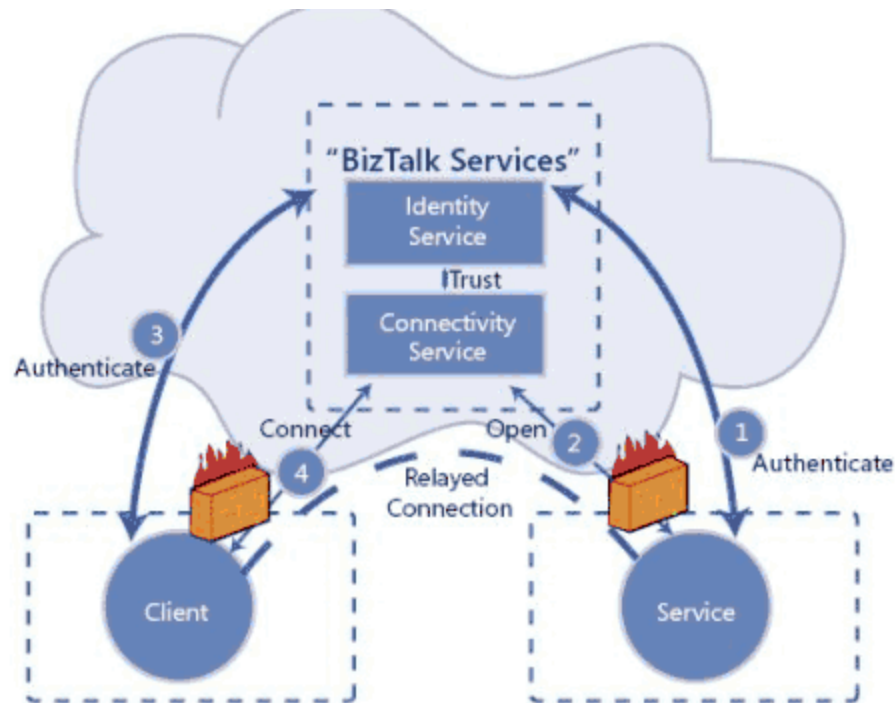
Figure 3: Map-reduce compilation of Pig Latin.

- Hive + HiveQL
 - data warehousing infrastructure
 - nástroje pro umožnění jednoduché extrakce/transformace/náctení dat
 - mechanismus pro uložení struktury na různou škálu datových formátů
 - umožňuje přístup k souborům uloženým na HDFS nebo v systémech datových úložišť jako je HBase
 - spuštění dotazu přes MapReduce
 - HiveQL je SQL-like jazyk sloužící pro dotazy nad HDFS
 - rozšiřitelný o vlastní MapReduce mapovace a redukovače
- cloud networking
 - Azure běží dnes na spoustě datových centrech v USA, Evropě a Asii
 - když chceme spustit aplikaci nebo uložit data můžeme si vybrat z těchto datových center, které chceme použít
 - existuje také několik způsobů, jak se k těmto centřům připojit
 - virtuální síť - propojení vlastní lokální sítě s definovanou množinou Azure VMs (s nějakým cloudem)
 - rozšíření vlastní sítě
 - přístup k službám lokální sítě z cloudových modulů
 - a naopak přístup k službám cloudu z lok. sítě
 - také propojení několika data center - geografická škálovatelnost
 - connect, přímé připojení - propojení cloudové sítě přímo s lokálním PC
 - nutnost VPN gateway, admin sítě
 - klientský SW
 - přístup ke cloudovým službám z lokálních PC
 - bez nutnosti síťování
 - cloudové služby se jeví jako v lokální síti
 - možnost individuálního nastavení apod.
 - traffic manager
 - pokud aplikace běží na více datacentrech, můžeme použít traffic manager, který používá inteligentní směrování uživatelských požadavků mezi instance aplikace
 - load balancing požadavků mezi data centra
 - definovaná pravidla směrování - od vývojáře aplikace:
 - typ performance - prostě nejbližší DC
 - typ failover - prioritní
 - typ round robin - kruhová fronta
 - vhodné pro velké aplikace, škálovatelnost



- cloud messaging
 - často je potřeba interakce aplikací
 - Azure nabízí několik možností, jak komunikaci řešit
 - service bus
 - založeno na mechanismu publish-subscribe
 - odesílatele (publishers) posílají zprávy pro celou třídu, bez toho, aniž by věděli, kdo do ní patří
 - příjemci (subscribers) se zapisují do třídy, o které mají zájem a přijímají tyto zprávy, aniž by věděli, kdo je odesílatelem
 - volně spájená komunikace
 - skupinová komunikace 1:N, více příjemců
 - API mimo cloud - letecká společnost, zprístupnění rezervací letenek a oznamování změn
 - mechanismy queues/topics/relays
 - queues
 - spolehlivý asynchronní jednosměrný kanál mezi rolími
 - web role přijme soubor s videem
 - uloží jej do BLOBu, pošle zprávu do message queue
 - worker role zprávu vyzvedne a provede konverzi formátu
 - strukturovaná/nestrukturovaná část zprávy
 - binární data v body, key/values properties
 - škálovatelnost - nezávislý počet read/write rolí
 - odesílatel i příjemce může být více
 - automatický load balancing
 - vyzvednutí zprávy
 - receive/delete režim
 - service bus přijme požadavek, označí jej, že se zpracovává a smaze zprávu z fronty
 - pokud aplikace, která má přijmout zprávu spadne, je

- zprava ztracena
 - idempotentni sluzby
 - jednodussi, mozna ztrata pri vypadku
- peek/lock rezim
 - dvoufazovy protokol pro dorucovani zprav
 - kdyz zprava dorazi do fronty, je uvedena do zamceneho stavu a setrvava v nem, dokud se stav nezmeni nebo dojde k timeoutu, kdy dojde zase ke zviditelneni zpravy ve fronte
 - kdyz je zprava v zamcenem stavu, muze aplikace vykonat pouze jedinou z techto operaci:
 - complete - zprava uspesne zpracovana
 - abandon - dorucovani selhalo (ale ne na zaklade obsahu zpravy), zprava je odemknuta a muze se opet dorucovat
 - defer - odlozene zpracovani
 - deadletter - odlozene zpracovani, kdyz dorucovani selhalo na zaklade obsahu zpravy
 - odolnost proti chybam
- topics
 - sekvence zprav tykajici se nejakeho tematu
 - predplatne (subscribing) je definovano jmenem a filtrem
 - filtr na zaklade vlastnosti zpravy (WHERE klauzule)
 - kazde predplatne ma svoji virtualni kopii zpravy
 - jedna zprava muze byt zpracovana vice prijemci
 - receive/delete, peek/lock
- relays
 - obousmerne propojeni na rozdil od queues/topics
 - motivace
 - obe aplikace za firewallem - blokovane prichozi porty
 - NAT - promenniva IP adresa
 - reseni
 - komunikace pres service bus relay service
 - obe aplikace se prihlasi
 - registr relay IDs



- caching
 - aplikace casto pristupuji na ta sama data
 - chceme uchovat kopii takovych dat co nejbliž k aplikaci, která s nimi pracuje
 - Azure in-memory caching - data se zachovají přímo v operační paměti VM dane aplikace
 - key/value store, kde value je nějaký chunk dat, možno i strukturované
 - různé způsoby - intra-app, inter-app, inter-cloud (v rámci app, mezi app, mezi cloudy)
 - CDN (content delivery network)
 - mnoho uzlů rozmístěných po celé zemi, galaxii
 - serve content to end-users with high performance & high availability
 - první přístup - kopie dat do uzlu CDN
 - další přístup - jen k uzlu CDN
- multimedia
 - video služby, media služby
 - uploading BLOBu, encoding, konverze formátu, packaging, streaming, protection (komponenta pro digital rights management), advertising
- GIS (geographic information systems)
 - mapy - free/private/special
 - vazby na relační data
 - API, web
- marketplace
 - aplikace, služby

- prodej MP3, e-knihy
- vydavatele - pricing, updates, ranking, evaluations, ...
- customers - identity, subscriptions, push, updates, reviews, rankings, ...
- publishing/subscription API
- mobilní služby
 - mobilní klienti
 - autentizace
 - přístup k DB
 - SMS, email, MMS
 - MightyText
 - další služby: hry, překladače, google docs, google drive

Security

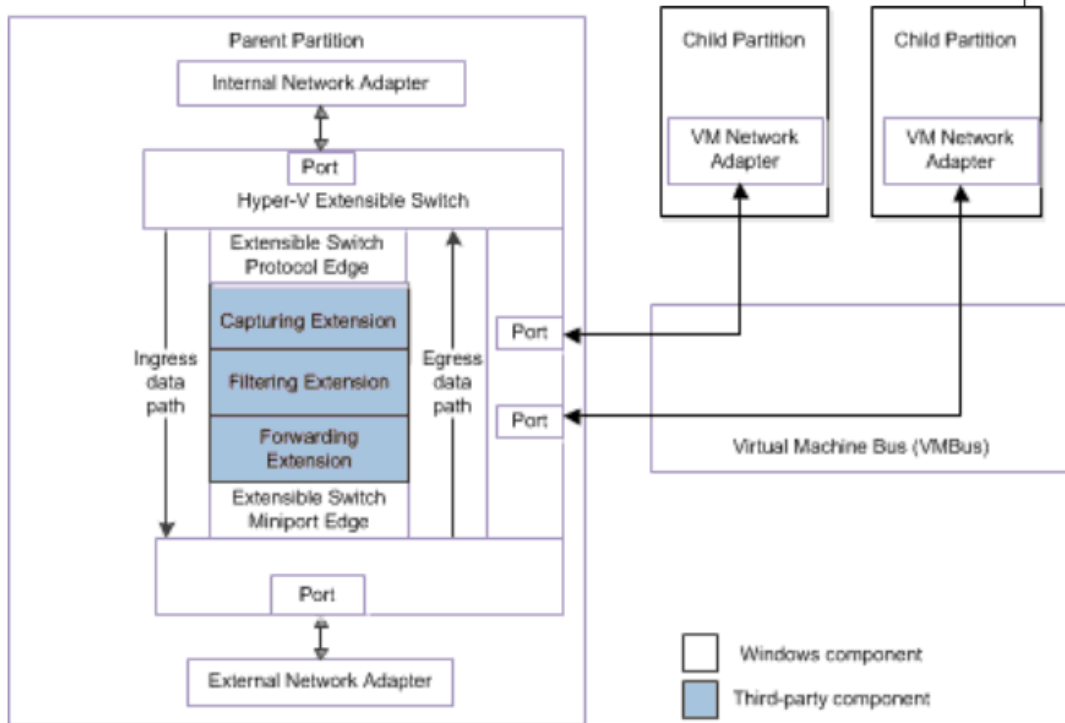
Bezpečnostní rizika virtualizace

- blue pill attack
 - chyti do pastí instancí bezcího OS tím, že nainstaluje hypervizor a virtualizuje stroj, na kterém běží
 - potřebuje však oprávnění ke spuštění privilegovaných instrukcí
 - předchozí OS bude stále udržovat své existující reference na otevřené soubory a zařízení, ale skoro cokoliv jiného se může zachytit - skoro libovolná zpráva - a posílat falešné odpovědi (HW interrupty, pozdravy na data, systémový čas apod.)
 - předvedla Rutkowska 2006, označeno za 100% nedetekovatelné
 - zveřejnila i kód, který zjistí, že systém běží virtualizovaný (jsme v Matrixu) - red pill
 - idea je, že na virtualizovaném systému, kde takto běží hypervizor a původní OS se musí sdílet prostředky, VMM musí relokovat strukturu s vektorem přerušení, aby nekolidovala s hostovou strukturou
 - na x86 existuje instrukce v neprivilegovaném módu, která ukládá do operandu pravou adresu tabulky
 - můžeme otestovat, zda je tabulka na svém místě nebo ne
 - další kód detekující běh ve virtualizovaném prostředí: timing attack
 - trap and emulate trvá mnohem déle než provedení nativní instrukce
 - je potřeba používat externí zdroj času (NTP), protože lokální čas může být podvržen
 - programy ale detekují pouze přítomnost virtualizace, ne blue pill malware
 - spousta programů používá virtualizaci - velký počet falešných odhalení
 - navíc jsou detekce ne úplně přesné
- zranitelnost VMM
 - problém je, že útokem na VMM může útočník napadnout více serverů naráz
- fyzické, virtuální firewally
 - při použití virtualizace servery z různých trust zone většinou sdílí stejné fyzické

prostředky (paměť, síťové karty, apod)

- webové rozhraní pro přístup k DC
 - spousta datacenter mohou být administrovány přes webové rozhraní - přes nějakou centrální konzoli, ta pak musí být dobře zabezpečena
- další nebezpečí
 - spolehani na tradiční bezpečnostní bariéry
 - fyzické firewally, apod
 - nefungují efektivně při dynamické povaze virtuálních instancí
 - aktuálně rychle se rozšiřující poskytování těchto služeb klientům
 - bezpečnost přenechána netradičním zaměstnancům
 - ...
- řešení
 - virtuální firewally
 - běží plně ve virtualizovaném prostředí
 - nabízejí službu jako fyzický firewall
 - dodávaný např. jako tradiční SW na guest VM, jako ucelově navržený systém pro virtuální síť, virtuální switch s dalšími bezpečnostními schopnostmi nebo řízený kernelový proces běžící v rámci hypervizoru hosta
 - podpora migrace za běhu
 - použití u stretched clusters
 - 2 a více fyzických hostů nainstalovaných na separátních místech méně než 100 kilometrů od sebe, které jsou součástí stejného vSphere clusteru
 - agentless antivirus
 - standardní agent-based bezpečnost je pro virtuální prostředí a jeho dynamickou povahu slabá
 - takto se dají ochránit všechny VM na hostovi bez nutnosti instalace agenta na každé VM
 - je to nějaký virtuální host, který se stará o detekci
 - s agenty má detekce velký dopad na výkon systému - AV storm, tedy se tímto budeme předcházet
 - extensible switches
 - virtuální switche umožňují komunikaci mezi VM, inteligentní směrování paketů a jejich inspekci
 - extensible switche umožňují přidávání pluginů - doplňující funkcionality od vývojáře
 - přizpůsobení switche pro dané prostředí, ve kterém pracuje

Extensible Switch



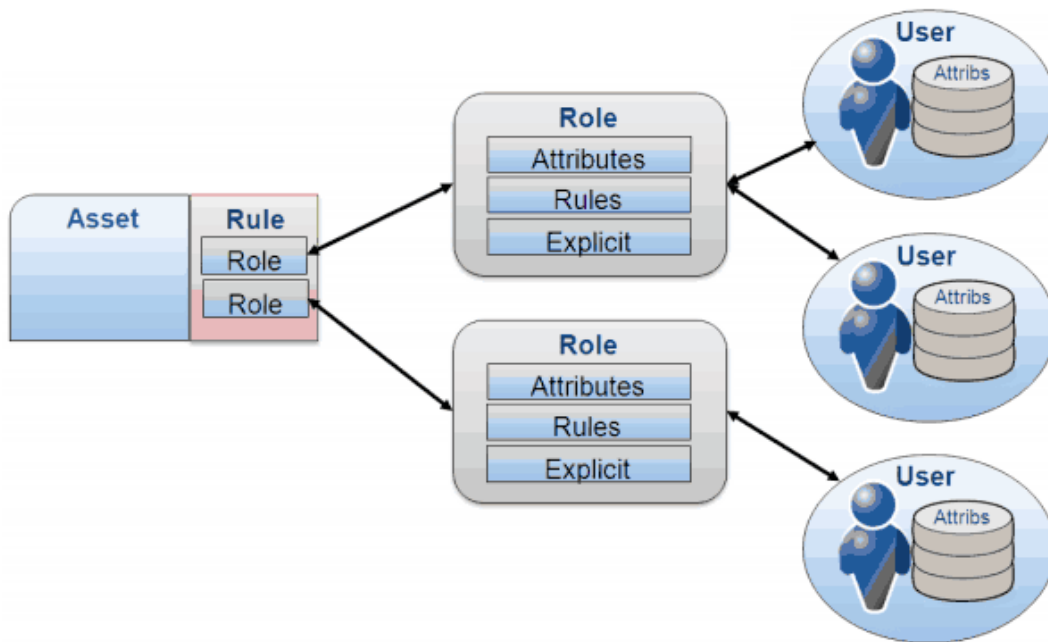
Bezpečnostní rizika cloud computingu

- cizí přístup k našim datům
- dostatečné fyzické zabezpečení
- nepoužívané, poškozené pevné disky sesrotovat v drtíčkách
- nejisté lokace uložených dat
- dodržování právních předpisů
- nedostatek investigativních mechanismů
- zotavení po havárii
- dlouhodobá životaschopnost

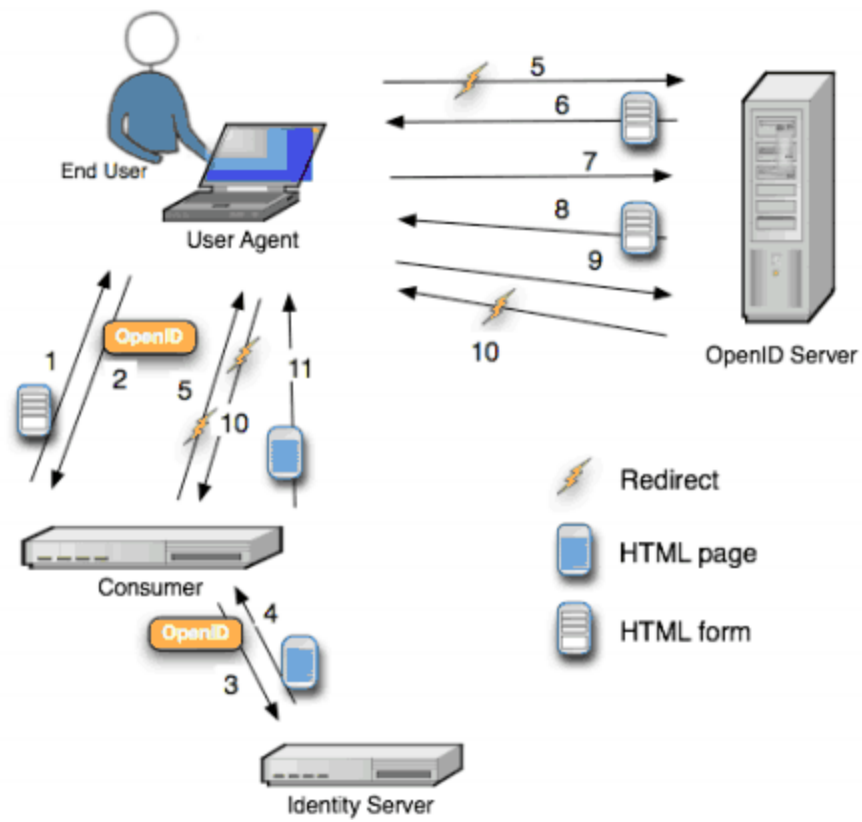
Identity management

- jde nám o autentizaci a autorizaci
 - prokázání, zda jde o správného uživatele
 - prokázání, zda má daný uživatel pravomoce pro přístup k datům
- externí uživatelské databáze
 - active directory / LDAP protokol použitý jako centrální autentizační mechanismus
 - záznamy ve stromové struktuře, informace o uživateli, do jakých skupin patří, jaké má privilegia
- dvoufázová autentizace
 - dvě různé ověřovací metody pro zajištění dostatečné bezpečnosti

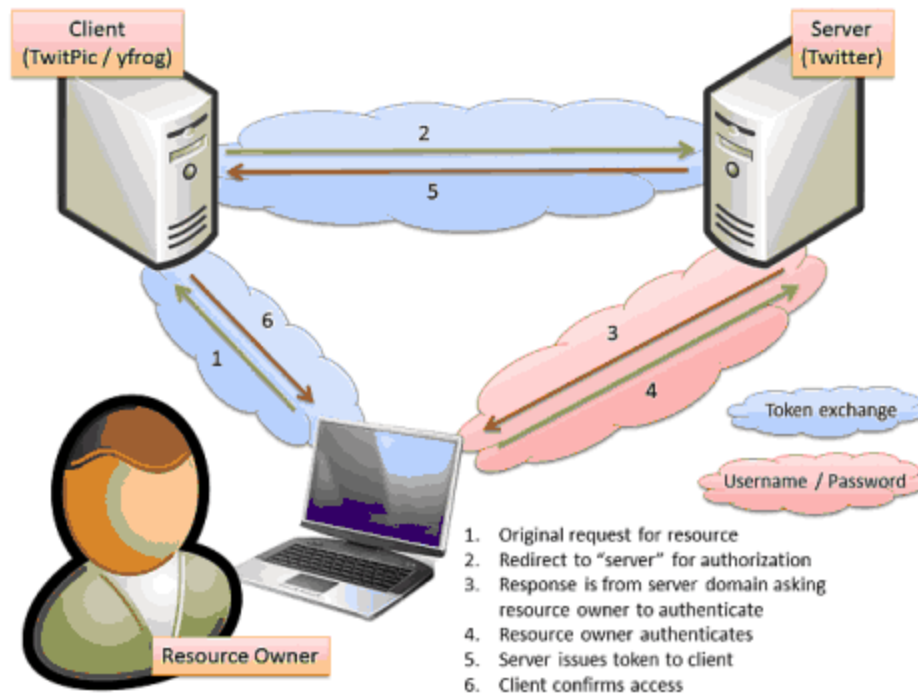
- přístupová karta + USB flash
- PIN + overovací kód z mobilu
- ...
- dnes již bezne - google, facebook, dropbox, paypal, ...
- role-based access control
 - motivací je specifikovat a vynutit bezpečnostní politiku specifickou pro celý systém VMs, která přirozeně reflektuje strukturu uživatelského systému
 - pro různé funkce v rámci firmy jsou vytvořeny různé role
 - povolení na provedení specifické akce se pak řídí podle role, kterou uživatel má/vykonává
 - přiřazení rolí k uživatelským účtům



- OpenID
 - otevřený standard popisující decentralizovaný způsob autentizace uživatele, který odstraňuje potřebu na straně provozovatele služby poskytovat a vyvíjet vlastní systémy pro autentizaci a který rovněž samotným uživatelům služby umožňuje konsolidaci jejich digitálních identit
 - má tvar unikátního URL, ke kterému je přiřazeno heslo
 - služba, která uživatelům autentizaci pomocí OpenID nabízí, při přihlašování uživatele přesměruje požadavek na ověření identity na správce daného OpenID účtu a ten vrátí informaci o povolení či zamítnutí žádosti o autentizaci
 - používá Google, IBM, seznam.cz



- OAuth
 - pouziva se delegovani uzivatelske autorizace treth strane



- SAML
 - podobne jako OpenID, ale cilene na podniky
 - security assertion markup language
 - XML-based
 - podpora single sign-on (SSO)
 - uzivatel se autentizuje na serveru A, pak se chce prihlasit na server B, B zjisti se vstupuje ze serveru A, zepta se serveru, jestli uz se neautentizoval u nej a muze rovnou povolit vstup
 - nutnost vzajemne duvery mezi identity providerem a service providerem
 - podpora mapovani protokolu na jine nez HTTP
 - MS Active Directory Federation services
 - komponenta na Windows serveru umoznujici SSO → SAML-based
 - CAS (central authentication service)
 - SSO protokol
 - postaveny na open-source projektu Shibboleth
- Eduroam
 - pouziva na autentizaci hierarchickou strukturu RADIUS serveru