

Optimalizace správy a zabezpečení virtuálních strojů

Ing. Michal Švamberg, Ing. Jan Krčmář

Západočeská univerzita v Plzni
Centrum informatizace a výpočetní techniky
email: svamberg@civ.zcu.cz, honza801@civ.zcu.cz

14. července 2017

Závěrečná zpráva projektu Fondu rozvoje CESNET, z.s.p.o. pro rok 2017 vedená pod číslem 571R1/2015. Tento projekt je zařazen do oblasti **I.**, tématického okruhu **A.** Řešitel projektu Ing. Michal Švamberg.

Struktura dokumentu

Struktura tohoto dokumentu je v souladu s podklady pro závěrečné oponentní řízení Fondu rozvoje CESNET, z.s.p.o. rozdělena následujícím způsobem:

- Způsob řešení
- Dosažené cíle
- Zdůvodnění změn v projektu
- Konkrétní výstupy
- Přínosy projektu
- Tisková zpráva

Způsob řešení

Projekt byl řešen v souladu se žádostí. Jako první krok jsme provedli vyhledání a srovnání existujících řešení, abychom vybrali vhodného kandidáta pro správu virtualizačního clusteru. Na základě kritérii

- náročnosti implementace,
- nákladů na provoz,
- přehlednosti uživatelského i administrátorského rozhraní,
- využití standardních nástrojů,
- snadnost rozšiřitelnosti,
- aktuálního vývoje

jsme vybrali projekt [WebVirtCloud](https://github.com/honza801/webvirtcloud)¹. V porovnání s ostatními projekty jako je OpenNebula², OpenStack³, oVirt⁴ se nám jevil jako nejvhodnější.

Podle plánu byly nakoupeny tři servery vhodné pro virtualizaci v celkové hodnotě 380.420,- Kč bez DPH.⁵ Servery byly nakoupeny z [výběrového řízení](#)⁶ dle platné legislativy.

Architektura diskového prostoru pro virtualizaci byla upravena do podoby (od nejvyšší vrstvy):

¹<https://github.com/honza801/webvirtcloud>

²Komplexní, určený pro větší nasazení.

³Příliš komplexní pro naše potřeby.

⁴Je silně postavený na Red Hat distribuci a jeho odvozenin.

⁵Zvolili jsme 2 různé konfigurace, původně se počítalo s třemi stroji se 128GB RAM, ale klesající ceny umožnily koupit jeden stroj s 384GB RAM. Oproti žádosti jsme také pořídili vícejádrové procesory, z plánovaných 6 jader na 10 a 12 jaderné. Ostatní parametry se oproti žádosti nezměnily.

⁶https://zakazky.zcu.cz/contract_display_1009.html

1. souborový systém virtuálního stroje, např. `xfs`
2. virtuální zařízení⁷ uvnitř stroje, např. `/dev/vda`
3. soubor `.qcow2` exportovaný jako virtuální zařízení do hosta, např. `/fc-sas/libvirt/storage/ourea22.qcow2`
4. sdílený souborový systém `GFS2`⁸, bude probráno detailněji níže
5. `CLVM`⁹ pro snazší správu diskového prostoru, vznikne např. `/dev/fc-sas/gfs2`
6. `multipath`¹⁰ pro zvýšení dostupnosti¹¹, např. `/dev/dm-1`
7. disková FibreChannel infrastruktura např. `/dev/sda`, `/dev/sdb`

Celé řešení nám takto dává cluster pro virtualizaci. Nejdůležitější je správa souborového systému a Cluster LVM, která musí být přes všechny uzly synchronní. O to se starají nástroje z `Red Hat cluster suite`¹², které byly portovány do mnoha dalších distribucí včetně Debianu.

`GFS2` bylo vybráno z důvodů využít existující FibreChannel diskovou infrastrukturu. V obecném pojetí stačí jakýkoliv souborový systém, například `NFS`. Důležité je, aby byl k dispozici prostor, v němž lze uložit `.qcow2` obrazy virtuálních strojů. Ukládání virtuálních disků do souboru oproti vyhrazeným oddílům¹³ přineslo výhody:

- úspora místa – virtuální stroj nemá žádný vyhrazený oddíl, tudíž nealokuje volné místo
- snadnější manipulace – klonování, vytváření snapshotů či migrace je otázkou vteřin
- snadnější rozšíření prostoru – stačí navýšit limit a zvětšit souborový systém
- jednodušší vytváření snapshotů – podporováno přímo knihovnou `libvirt`
- snadnější migrace virtuálních strojů – opět přímá podpora v `libvirt`

Nejzásadnější změna byla tedy v přechodu z vyhrazených LVM oddílů na `.qcow2` formát umístěný na sdíleném souborovém systému `GFS2`. Vyzkoušeli jsme také souborový systém `OCFS2`¹⁴. `OCFS2` jsme ale po úvodních pokusech opustili, občas se stalo, že pádem jednoho uzlu souborového clusteru byly sestřeleny všechny uzly (chyba v jaderném modulu)

⁷Pro linuxové systémy nepoužíváme rozdělování disků, pro operační systém Windows je situace opačná

⁸<https://en.wikipedia.org/wiki/GFS2>

⁹https://access.redhat.com/documentation/en-US/Red_Hat_Enterprise_Linux/6/html/Logical_Volume_Manager_Administration/LVM_Cluster_Overview.html

¹⁰https://en.wikipedia.org/wiki/Linux_DM_Multipath

¹¹Z více cest k jednomu zařízení vytvoří jedno zařízení

¹²https://en.wikipedia.org/wiki/RedHat_cluster_suite

¹³V původním řešení byly vyhrazeny LVM oddíly pro každý virtuální stroj, minimálně kořen a swap.

¹⁴A to na základě zkušeností z provozu Oracle databází.

a servery musely být resetovány. Je třeba říci, že GFS2 není také bez chyby, ale problémy nám přišly snadněji řešitelné a nemávaly tak fatální dopad.

Dále jsme přešli od technologie virtualizace Xen na KVM. Od nasazení KVM si slibujeme lepší podporu pro virtualizaci operačního systému Windows, zvláště pak při aktualizacích virtualizační platformy. Druhou výhodou spatřujeme ve větší komunitě uživatelů (a potažmo vývojarů), navíc podpora KVM je přímou součástí linuxového jádra, což u Xenu donedávna nebylo.

Jak bylo řečeno již dříve, rozhodli jsme se pro správu systémem [WebVirtCloud](#)¹⁵. Jedná se o menší projekt využívající rozhraní [libvirt](#)¹⁶. Je psaný v jazyce Python nad frameworkem [Django](#)¹⁷. WebVirtCloud jsme rozšířili a doplnili, změny byly autorem do [projektu přijaty](#)¹⁸, takže nyní jsou již jeho součástí. Ke dni 16.6.2017 je Jan Krčmář (spoluřešitel) [druhým nejčastějším přispěvatelem](#)¹⁹ do WebVirtCloudu hned po samotném autorovi.

Uživatelé si pochvalují jednoduché a snadné ovládání WebVirtCloudu. Pro nejčastější zásahy je k dispozici administrátorské rozhraní. Pouze pokud přijdou speciální požadavky (nejčastěji týkající se kombinace diskových prostorů) je nutné použít přímo nástroje pro libvirt nebo operaci provést z příkazové řádky hostitele.

Virtuální stroje jsme taktéž chtěli nabídnout běžným uživatelům. Dnes existuje mnoho služeb, kde si takový server lze zřídit, často ze začátku zdarma. Ale my jsme chtěli uživatelům nabídnout server s návazností na existující univerzitní infrastrukturu (jednotné přihlašování, předpřipravený souborový systém AFS, atd.). Zároveň jsme takové stroje chtěli lépe zabezpečit než jsou standardní distribuce linuxu. Další bezpečnostní službou uživatelům měla být aktualizace operačního systému a jeho vybavení v době, kdy vyjdou záplaty chyb.

Shodně s žádostí projektu byla tato (a mnohá další) nastavení zajištěna konfiguračním managementem CFEngine. Konfigurační management je předpřipraven v šablonách (templates) ze kterých vznikají virtuální stroje. Podstatná část projektu je tak zcela oddělená a nebrání přenést znalosti do jiných organizací. Bezpečnostní nastavení pak záleží na přípravě obrazu. My jsme zvolili minimalistický obraz jen s potřebným minimem nástrojů a přednastaveným konfiguračním managementem, který zajistí úpravu operačního systému tak, aby odpovídal aktuálním politikám konfiguračního serveru.

Projekt předpokládal maximálně automatizovat vytváření a rušení serverů včetně návazností na další systémy. Registraci virtuálních strojů jsme vyřešili tak, že máme předgenerované doménové jména a nastaveny k nim statické adresy včetně MAC adres pro DHCP.

Další systém, který jsme chtěli napojit na virtuální stroje byl monitoring. Po prozkoumání možností jsme zjistili, že nejsnazší způsob by byl přístup přes REST API do monitoringu a stroj se sám nastavil. Bohužel aktuálně provozovaný monitoring²⁰ nám toto

¹⁵<https://github.com/honza801/webvirtcloud>

¹⁶<https://libvirt.org/>

¹⁷<https://www.djangoproject.com/>

¹⁸<https://github.com/retspen/webvirtcloud/commits?author=honza801>

¹⁹<https://github.com/retspen/webvirtcloud/graphs/contributors>

²⁰Nagios a správa konfigurace přes webové rozhraní NConf.

| Name | Description | Host User | Status | VCPU | Memory | Actions |
|---------|----------------------|-------------------------|--------|------|---------|--|
| themis | pajovo testovani idm | nesoi1 paja | Active | 8 | 8192 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| styx | Guacamole server | nesoi1 dextor | Active | 1 | 2048 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| loco | pajovo kolobezky | nesoi1 paja | Active | 1 | 2048 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| koleje | kolejni server | nesoi1 svamberg (+1) | Active | 1 | 2048 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| shib-sp | pajovo shibboleth | nesoi1 paja | Active | 1 | 2048 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| themis3 | pajovo testovani idm | nesoi1 paja | Active | 2 | 8192 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| themis2 | pajovo testovani idm | nesoi1 paja | Off | 2 | 4096 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| wxkDC1 | wxk DC1 | nesoi1 mlavicka | Active | 4 | 4096 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| ourea1 | ourea1 | nesoi1 svamberg | Off | 1 | 2048 MB | [Play] [Pause] [Power] [Refresh] [Eye] |
| atlas | atlas | nesoi1 svamberg (+1) | Active | 4 | 6144 MB | [Play] [Pause] [Power] [Refresh] [Eye] |

Obrázek 1: Přehledová stránka všech spravovaných strojů

neumožňuje. Navíc připravujeme zcela nově postavený monitoring na [Icinga 2](https://www.icinga.com/)²¹, která má REST API vestavěné. Prozatím jsme vyřešili dohled tak, že CFEngine zná konfiguraci Nagiosu a hlídá, zda odpovídá provozovaným službám.²² Toto řešení není dokonalé, ale do doby než bude připraveno nové monitorovací prostředí, je dostačující.

Jak již bylo naznačeno dříve, instalaci jsme se rozhodli řešit formou předpřipravených obrazů, které se po startu přizpůsobují (roztáhnou souborový systém, nastaví hostname, inicializují konfigurační management, ...). Dříve zvažovaný systém [FAI](http://fai-project.org/)²³ jsme opustili a přešli k nástroji [DebianInstaller s preseed](https://wiki.debian.org/DebianInstaller/Preseed)²⁴ nastavením. Uživatelé pro instalaci používají

²¹<https://www.icinga.com/>

²²Pokud se na stroj nainstaluje databáze a v Nagiovi není sonda na její kontrolu, tak upozorní na tuto nesrovnalost správce. Ten ve webovém rozhraní NConf sondy danému hostu přidá.

²³<http://fai-project.org/>

²⁴<https://wiki.debian.org/DebianInstaller/Preseed>

klonování obrazů, které je rychlé a v cloudovém prostředí běžné. Chtěli jsme ale instalaci virtálních a fyzických strojů co nejvíce sblížit, tedy instalovat jen potřebné minimum. Přípravu instalace fyzických strojů tak nyní můžeme otestovat snadno na virtálním stroji, což nám přináší zjednodušení a urychlení práce.

The screenshot shows the WebVirtCloud management interface. At the top, there is a navigation bar with 'WebVirtCloud', 'Instances', 'Computes', 'Users', and 'Logs'. The user 'svamberg' is logged in. Below the navigation bar, the instance 'icarus (icarus)' is shown as 'Active' with 1 Vcpu, 2048 MB Ram, 10.0 GB Disk, and 512.0 MB Disk. A toolbar contains icons for Power, Access, Resize (highlighted), Snapshots, Settings, Graphs, and Destroy. The 'Resize Instance' dialog box is open, showing the following configuration:

- Logical host CPUs: 24**
 - Current allocation: 1
 - Maximum allocation: 1
- Total host memory: 378.6 GB**
 - Current allocation (MB): 2048
 - Custom value
 - Maximum allocation (MB): 2048
 - Custom value
- Disk allocation (B):**
 - Current allocation (vda): 10.0 GB
 - Current allocation (vdb): 512.0 MB

A green 'Resize' button is located at the bottom right of the dialog box.

Obrázek 2: Možnost změnit parametry virtuálního stroje

Dosažené cíle

Byla vytvořena zcela nová virtualizační platforma, která zjednodušuje práci administrátorům a uživatelům nabízí přehledné a jednoduché rozhraní. Uživatelé si mohou vytvořit

vlastní virtuální stroje a to do limitů stanovených správci. Otázka bezpečnosti byla zahrnuta jako součást předpřipravených obrazů, zároveň konfigurační management běží také na všech serverech virtualizační platformy. Konfigurační management na strojích (hostitelích i hostech) hlídá aktuálnost nainstalovaných balíčků, nastavení firewallu, hesel a mnohé další volby, které ovlivňují bezpečnost. Další možnosti přidáváme podle toho, jak přicházejí nové požadavky.

Novým virtualizačním systémem se podařilo odstranit největší bolesti toho starého. Šlo zejména o rozdělování disků a přiřazování diskového prostoru, chybějící uživatelské rozhraní, komplikovanou instalaci a inicializaci virtuálního stroje a nízká úroveň bezpečnosti. To byly hlavní důvody, proč jsme nemohli virtuální servery nabídnout všem uživatelům univerzity.

Všechna bezpečnostní opatření popsaná výše závisí na tom, zda bude aktivní konfigurační management na daném virtuálním stroji. U strojů, které jsou ve správě, tuto vlastnost ohlídá monitorovací systém. Problém nastává u „rychloobrátkových“²⁵ virtuálních strojů nebo těch, kde práva administrátora má také uživatel, který stroj vytvořil. Administrátor může totiž konfigurační management vypnout, čímž také vypne všechna bezpečnostní opatření. Tomu jsme se rozhodli zabránit tak, že na serveru konfiguračního managementu se kontroluje časová značka, kdy naposledy klient konfiguračního managementu kontaktoval server. Pokud se stane tato značka neaktuální, je vlastníkově poslána výstraha emailem a virtuální stroj je vypnut. Uživatel si jej může z webového rozhraní opět nastartovat a provést nápravu, pokud tak opět neučiní, server se zase vypne.

Bohužel se nám nepodařilo najít žádného studenta, který by byl ochotný na projektu hlouběji spolupracovat. Není mnoho studentů, kteří se chtějí věnovat správě linuxových operačních systémů. Alespoň pro testování jsme získali Marka Zimmermanna, který nám již v rané fázi přípravy přinesl zpětnou vazbu.

Zdůvodnění změn v projektu

Oproti původnímu plánu jsme provedli několik drobných změn v projektu. Místo původní konfigurace serverů jsme nakoupili jeden v lepší konfiguraci. Důvodem byla klesající cena serverů. Pro stanovení kvóra clusteru jsme potřebovali minimálně 3 stroje. Ke konci projektu jsme k přidali další 2 stroje, které do hospodaření projektu nezapočítáváme. Jsou to stroje, které navazují na projekt a rozšiřují jej.

V implementační části projektu jsme se rozhodli nenapojovat systém těsněji na monitoring než je v současnosti. Zjistili jsme, že vhodné řešení by vyžadovalo významnější zásahy do monitorovacího systému. Proto jsme s kolegy začali diskutovat o možné změně monitoringu²⁶ tak, aby to vyhovalo i dynamickému vzniku či zániku serverů v cloudovém prostředí. Současné řešení nás upozorňuje na změny, ale ty se musí zanášet do systému ručně, což není úplně dle představ, ale do nového monitoringu s tím vydržíme.

²⁵Pokud potřebuji stroj jen na krátkodobý pokus. Stroj může být používán v jednotkách hodin nebo několika dní.

²⁶Aktuálně se testují možnosti přechodu z Nagios 3.5.1 na Icinga 2.

V projektu se nám nepodařilo úplně dobře čerpat z cestovného a konferenčního. Důvodem byl fakt, že původně plánované výdaje na konferenci OpenAlt se rozhodl zaplatit organizátor konference sám.

V žádosti jsme uvedli, že bude třeba oslovit uživatele a toto řešení jim nabídnout. Zatím řešení virtualizace nabízíme, pokud přijde požadavek na zřízení serveru. Hromadněji chceme uživatele oslovit až po aktualizaci operačního systému, pravděpodobně se začátkem nového akademického roku 2017/18.

Konkrétní výstupy

K 10.7.2017 jsme na celkem 5 hostovských serverech provozovali 104 virtuálních strojů. Požadovaný diskový prostor virtuálních strojů je 3,5TB, ale reálně zabráno 1,7TB. To znamená, že oproti původnímu Xen řešení šetříme 52% diskového prostoru.

Během projektu jsme uspořádali několik seminářů, na kterých jsme představili naše řešení, ukázali jeho výhody a předvedli používání.

- Seminář pro IT techniky CIV, jaro 2016, Plzeň (neveřejné)
Sloužilo administrátorům pro seznámení s novým prostředím virtualizace a jejich možnostmi. Cílem byla diskuse o použitých technologiích, navázání na infrastrukturu a hlavně získání zpětné vazby po tom, co jej začali používat.
- OpenAlt, 5.11.2016, Brno
[Přednáška o řešení projektu](#)²⁷ s krátkou ukázkou používání. K dispozici je také [video záznam](#)²⁸.
- EurOpen, 16.5.2017, Myslovice
[Prezentace stavu virtualizace](#)²⁹ a jeho řešení. Součástí je [příspěvek do sborníku konference](#)³⁰. V něm lze najít porovnání nejčastějších cloudových rozhraní i konkrétní technické řešení tohoto projektu.

Vzhledem k velmi silné pozitivní vazbě od správců IT, kteří virtuální prostředí začali aktivně používat, jsme se rozhodli informaci o možnosti zřídit si vlastní virtuální stroj uvolňovat trochu obezřetněji. A to minimálně do doby než virtualizační infrastrukturu posílíme³¹. Na druhou stranu, pokud je příležitost, tak systém nabídneme. Máme tak již virtuální servery pro katedry, projekty i studenty.

²⁷<https://openalt.cz/2016/data/Svamberg%20-%20webvirtcloud.pdf>

²⁸<https://www.superlectures.com/openalt2016/jednoduchy-cloud-vhodny-i-pro-univerzitu>

²⁹http://europen.cz/Proceedings/50/Svamberg_europen_webvirtcloud.pdf

³⁰<http://europen.cz/Anot/50/eo-1-17.pdf>

³¹Aktuálně je uzavřené výběrové řízení na dodávku serverů mezi kterými je také posila virtualizační platformy. Dodání serveru očekáváme v průběhu léta 2017.

Z hlediska správy operačního systému se nám podařilo sjednotit management serverů včetně zohlednění bezpečnostních hrozeb a aktualizací. Nyní můžeme upravit hromadně všechny servery do 10 minut³².

Vznikla [uživatelská dokumentace s návodem](#)³³, jak si vytvořit virtuální stroj. Vytvořena [dokumentace s poznatky o provozu clusteru](#)³⁴ pro virtualizaci v Debianu.

Cluster se aktuálně sestává z 5 nodů³⁵ s celkovou pamětí 2T a počtem procesorů 192. Na této infrastruktuře aktuálně provozujeme 104 virtuálních strojů. Navýšení bylo potřebné z toho důvodu, že připravujeme převod původních virtuálních strojů z technologie Xen na KVM a také aktualizovat infrastrukturu z Debian Jessie na Stretch³⁶.

Přínosy projektu

Zjednodušení přípravy a samotné správy virtuálních strojů je hlavním přínosem pro administrátory operačních systémů. Pro uživatele pak snadný způsob vytvoření, ovládání, přístupu i zrušení stroje. O oblíbenosti tohoto rozhraní hovoří čísla. V systému je k 11.7.2017 aktivních 45 uživatelů, z toho 40 jsou běžní uživatelé, tj. nemají aktivovanou roli administrátora.

Přechodem na virtualizaci KVM očekáváme lepší podporu pro nelinearové operační systémy, zvláště pak pro MS Windows. Zde byly potíže při aktualizaci hypervizorů.

Velká úspora je pak na diskových polích, kde je prostor velmi drahý a je ho omezené množství. Nyní je úspora místa 52% oproti původnímu řešení.

Tisková zpráva

Projekt na Západočeské univerzitě zvýšil dostupnost a snadnost vytvoření vlastního virtuálního stroje. Umožnil lepší využití hardware a významně šetří diskový prostor oproti původnímu řešení. Byla vylepšena automatizace životního cyklu virtuálního stroje a napojení na návazné informační systémy. Pro správu virtuálních strojů byl nasazen projekt WebVirtCloud. Získané znalosti je možno využít i v dalších organizacích.

³²Synchronizace master serveru s GIT repozitářem se provádí v 5 minutových intervalech, klienti kontrolují konfiguraci ze serveru každých 5 minut.

³³http://support.zcu.cz/index.php/Cloudov%C3%A9_sl%C5%BEby

³⁴http://support.zcu.cz/index.php/Public:Honz801/debian_cluster

³⁵Oproti původně plánovaným třem jsme museli již během projektu počet navýšit.

³⁶Debian Stretch byl jako stabilní distribuce vydán 17.6.2017.