

MASARYKOVA UNIVERZITA
FAKULTA INFORMATIKY



Monitorování zátěže a využití výpočetních zdrojů v heterogenním výpočetním prostředí

DIPLOMOVÁ PRÁCA

Juraj Leždík

Brno, jar 2016

Prehlásenie

Prehlasujem, že táto diplomová práca je mojím pôvodným autorským dielom, ktoré som vypracoval samostatne. Všetky zdroje a literatúru, ktoré som pri vypracovaní používal alebo z nich čerpal, v práci riadne citujem s uvedením úplného odkazu na príslušný zdroj.

Vedúci práce: Mgr. Miroslav Ruda

Pod'akovanie

Zhrnutie

Klíčové slova

Contents

Obsah	viii
1 Úvod	1
2 Cloudové technológie	3
2.1 Cloud	3
2.1.1 Kontrola zdrojov cloudu	4
2.2 Virtuálne stroje	4
2.2.1 Hypervízor	5
2.3 Aplikačné kontajnery	5
2.4 MapReduce aplikačné prostredia	6
2.5 Cloudové technológie MetaCentra	6
2.5.1 libvirt/KVM	6
2.5.2 Docker	7
2.5.3 Hadoop	7
3 Zber a uchovávanie monitorovacích dát	8
3.1 Všeobecné problémy monitorovania	9
3.1.1 Monitorovanie v heterogénnom výpočtovom prostredí	11
Problematika intervalov	11
Detekcia nových zdrojov a užívateľov	12
3.2 Časové rady	12
3.2.1 Analýza časových rád	12
3.2.2 Zaobchádzanie s historickými dátami	13
4 Aktuálne monitorovacie riešenia	14
4.1 Linux cgroups	14
4.2 Nagios	15
4.2.1 PluginAPI	15
4.2.2 Timeouty	16
4.3 Zabbix	17
4.3.1 Timeouty	17
4.3.2 Moduly	17
4.4 Icinga	18
4.4.1 Externé pluginy	18
4.4.2 Notifikácie a príkazy udalostí	18

4.4.3	Itegrácie s aplikáciami	18
4.5	<i>Cacti</i>	18
4.6	<i>collectd</i>	18
4.6.1	Obmedzenia	19
4.6.2	Zapisovací plugin Write TSDB	19
4.6.3	Prahy a notifikácie	19
4.6.4	Ostatne	20
4.7	<i>Ganglia</i>	20
4.8	<i>Scout</i>	20
4.9	<i>New Relic</i>	20
4.10	<i>AppDynamics</i>	20
4.11	<i>OpenTSDB</i>	21
4.12	<i>InfluxDB</i>	21
4.12.1	Politiky udržiavania	21
4.12.2	Kontinuálne dotazovanie	22
4.13	<i>RRDTool</i>	22
5	Metriky	23
5.0.1	Význam popisných údajov metrík	23
5.0.2	Periodicita zberania metrík	24
5.0.3	Formát hodnoty metriky	24
5.1	<i>Docker</i>	24
5.1.1	Sieť	24
	Sieť typu most	24
	Prekladaná sieť	25
	Metriky siete	25
5.1.2	Pamäť	25
5.1.3	Procesor	26
5.2	<i>libvirt/KVM</i>	26
5.2.1	Metriky siete	26
5.2.2	Metriky pamäte	26
5.2.3	Metriky zápisu dát	27
5.2.4	Metriky CPU	27
5.3	<i>Hadoop</i>	28
5.3.1	Cluster Metrics API	28
5.3.2	Cluster Application API	29
5.3.3	Node Application API	30
6	Analýza a návrh	31
6.1	<i>MetaCentrum</i>	31
6.2	Požiadavky na aplikáciu	31
6.2.1	Vysoký monitorovací výkon	31

6.2.2	Nízka nadbytočná záťaž	32
6.2.3	Škálovateľnosť	32
6.3	Reakcia na dlhú odozvu modulu	33
6.4	OpenTSDB	33
7	Implementácia	34
7.1	Techniky zbierania metrík	34
7.1.1	Torque	34
7.1.2	Docker	34
7.1.3	libvirt/KVM	35
7.1.4	Hadoop	35
8	Záver	37
	Literatúra	37
A	Kapitola príloha	38

Chapter 1

Úvod

Cloudová infraštruktúra MetaCentra poskytuje výpočetné prostriedky pre mnohé vedecké a výskumné organizácie. Spracovávajú sa v nej veľké objemy dát. Je tvorená mnohými uzlami rozmiestnenými naprieč celou Českou republikou. Podstatou fungovania cloudového modelu je zdieľanie veľkého výpočetného výkonu viacerými subjektami, pre ktoré by zabezpečenie vlastnej infraštruktúry predstavovalo neúmernú ekonomickú, personálnu a prevádzkovú záťaž. Princípom zdieľania je, že prostriedky by mali byť ideálne využívané všetkými rovnako. Nemalo by dochádzať k tomu, že jeden klient vyťaží cloud natoľko, že výpočetné úlohy ostatných dostanú nepomerne malý priestor. To je možné docieľiť monitorovaním využitia výpočtového výkonu a periférií. Ak máme informácie o tom, kto koľko využíval zdroje, je možné tot využívanie účtovať a zároveň do budúcnosti primerane obmedzovať.

Aby bolo možné správne vyvodzovať závery o zaťažení cloudu, je potrebné zbierať údaje o tom periodicky a kontinuálne v čase. Je potrebné sledovať parametre v pravidelných intervaloch. Každému intervalu prináleží hodnota, ktorá vypovedá o využití prostriedkov v danom momente. Takéto dáta sa nazývajú časové rady. Cloud zahŕňa množstvo uzlov, kde na každom môže byť spustených množstvo úloh. O každej je potrebné zbierať viacero parametrov - vyťaženie procesora, pamäte, periférií. To predstavuje veľké množstvo metrických údajov, ktoré je potrebné ukladať a nejakým spôsobom vyhodnocovať. Slúžia na to databázy časových rád. Účtovanie sa tiež deje v určitých pravidelných intervaloch. Je preto potrebné mať možnosť spätne dohľadať údaje o využívaní zdrojov. Uchovávať podrobné metrické dáta má výnraz na určité obdobie. Nie je potrebné vedieť ako bol využívaný cloud úlohou pred piatimi rokmi každých päť sekúnd. Okrem toho by to predstavovalo veľké požiadavky na úložné kapacity, ktoré by rástli lineárne vzhľadom na počet úloh a čas. Je preto žiaduce po nejakej dobe dáta agregovať do väčších celkov a tiež mazať už nazbierané podrobné časové rady.

Aby boli monitorovacie údaje spoľahlivé, je potrebné zabezpečiť pe-

riodický zber metrických údajov. Monitorovanie sa uskutočňuje pravidelným pýtáním sa na vyťaženie daného zdroja. Aplikácia, ktorá zdroj využíva, vytvorí odpoveď a pošle ju späť monitorovacej aplikácii. Problém nastáva, ak odpoveď nie je vytvorená v dostatočne krátkom intervale. Povedzme, že chcem zberať údaje každých päť sekúnd. Aplikácia, ktorá práve rieši úlohy ale môže byť plne zaneprázdnená a nebude odpovedať na výzvy na údaje o vyťaženi. Tento problém je potrebné riešiť. Ak k tomu dôjde, nemali by byť vytvárané ďalšie výzvy na túto zaneprázdnenú aplikáciu v takej miere, ako keď plynulo odpovedala. Monitorovacia aplikácia by si mala zapamätať poslednú nameranú hodnotu a tú odoslať ako aktuálnu. Zároveň by sa mala zaneprázdnenej aplikácii pýtať v menej častých intervaloch. Ak aplikácia začne odpovedať, všetko je v poriadku. Ak nie, je možné reagovať reštartovaním monitorovacej aplikácie. Ak ani táto reakcia problém nevyrieši, komplikácia nastala zrejme v dotazovanej aplikácii.

Cloud predstavuje heterogénnu infraštruktúru. Na riešenie výpočetných úloh sa používajú rôzne technológie a aplikácie. Rôzne aplikácie pristupujú k výpočetným problémom odlišnými postupmi. Všetky úlohy sú ale počítané na akoby jednom veľkom a veľmi výkonnom počítači. Tieto aplikácie teda používajú spoločné zdroje. Metrické dáta, ktoré budem zberať, by takisto mali vypovedať o vyťaženi tých istých druhov zdrojov. Je preto potrebné identifikovať, ktoré metrické dáta z jednej aplikácie je možné porovnať s údajmi z inej aplikácie. Takto môžeme dostať celkový obraz o tom, ako rôzne technológie využívajú jednu skupinu zdrojov a len takto je možné vzájomne tieto technológie zrovnávať a účtovať ich využitie.

Chapter 2

Cloudové technológie

2.1 Cloud

Cloudové technológie umožňujú využívať veľkú množinu výkonných výpočetných zdrojov mnohým subjektom. Každý užívateľ cloudu využíva časť výpočetného výkonu. Spolu s rozdelením výkonu prichádzajú aj ďalšie výhody. Užívateľ sa nemusí starať o podliehajúci hardvér. Typ fyzického procesora alebo rýchlosť pamäte uzla cloudu sú v určitom zmysle dôležité, ale pre výkon je rozhodujúca veľkosť pamäte, kapacita disku resp. počet procesorových jadier. Tieto požiadavky cloud umožňuje jednoducho spravovať a transparentne meniť. Má to význam, ak sa rozrastú požiadavky užívateľa na výkon, alebo naopak z dôvodu obmedzenej prevádzky či nedostatku výpočetných úloh sa nároky môžu zmenšiť. Klient môže dané prostriedky využívať hneď, bez veľkých počiatočných investícií, ktoré by si buď nemohol dovoliť, alebo ak sú výpočetné úlohy krátkodobejšieho rázu, nákup stroja s požadovaným výkonom by bol nevhodnou investíciou. Užívateľ sa nemusí starať o nákup hardvéru, jeho zostavovanie do funkčných serverov a ich následné rozširovanie a správu. Tieto služby zabezpečuje prevádzkovateľ cloudovej infraštruktúry. Pre neho je zase dôležité mať prehľad o tom, ako sú jeho inštalované kapacity využívané. Či už z pohľadu skvalitňovania vlastných poskytovaných služieb alebo vo vzťahu ku klientovi a k tomu, v akej miere spotrebovávajú poskytovaný výkon. Jednou z charakteristík cloudu je aj heterogenosť v prístupe k zdrojom. Výpočetné prostriedky sú poskytované viacerými spôsobmi:

Infrastructure as a Service (IaaS) - Tento spôsob poskytuje užívateľovi priamo hardvérové prostriedky infraštruktúry. Ten má možnosť určiť, koľko pamäte, procesorov alebo diskovej kapacity požaduje. Je jeho voľbou aký operačný systém použije, aký softvér nainštaluje a aké výpočetné úlohy bude realizovať alebo aké služby bude prevádzkovať.

Platform as a Service (PaaS) - Užívateľ využíva priamo platformu pre-

vádzkovateľa. Jedná sa vrstvu o úroveň vyššie ako pri využití infraštruktúry. Hardvérová konfigurácia je daná, ale užívateľ využíva operačný systém a súbor aplikácií na vývoj a prevádzku vlastných programov. Vlastník cloudu volí platformu a inštaláciu ďalšieho softvéru. Užívateľ ho môže používať a prípadne meniť. Tento spôsob poskytovania zdrojov znižuje záťaž na infraštruktúru tak, že viacerí užívatelia využívajú jedno bežiacie jadro operačného systému.

Software as a Service (SaaS) - V cloudovej infraštruktúre môže byť nainštalovaný softvér, ktorý špecifickým spôsobom zefektívňuje prácu s distribuovanými výpočtovými zdrojmi. Je preto niekedy vhodné poskytovať samotné prostredie jednej aplikácie ako službu. Užívateľ navrhuje výpočetné úlohy alebo aplikácie s využitím knižnic a architektúry konkrétnej aplikácie a táto aplikácia sa zároveň stará o ich vykonávanie.

Momentálne MetaCentrum poskytuje svoj dostupný výkon vo všetkých uvedeých formách. Jednotlivým prístupom zodpovedá konkrétne softvérové riešenie, ktorému sa budem venovať v jednej z nasledujúcich častí.

2.1.1 Kontrola zdrojov cloudu

Bez ohľadu na to, na akej vrstve si klient zvolí prístup k využívaniu zdieľaných výpočetných prostriedkov, jeho výpočty sa v konečnom dôsledku budú musieť realizovať na fyzickom hardvéri poskytovateľa. Aj keď sa daná úloha vypočítava na viacerých uzloch cloudu a rozličnými postupmi, vlastník by mal mať možnosť nejakým jednotným spôsobom určiť ako je reálne celá infraštruktúra vyťažovaná. K tomuto

2.2 Virtuálne stroje

Virtuálne stroje poskytujú úplnú virtualizáciu fyzickej hardvérovej štruktúry. Na jednom hosťujúcom počítači môže byť spustených viacero virtuálnych strojov. Každý má svoj vlastný virtuálny procesor, pamäť, grafický procesor, pevný disk a periférie. Operačný systém spustený vo virtuálnom stroji je izolovaný od hosťovského operačného systému (ak ho hosťovský počítač má). Takéto riešenie má jednu bezpečnostnú výhodu oproti aplikáčnym kontajnerom. Nežiadúce fungovanie jedného virtuálneho stroja neovplyvňuje beh ostatných. V súčasnosti existuje mnoho úrovní virtuálnych strojov. Emulácia inštrukčnej sady, prekladanie programov za behu a

ich optimalizácia, vysokoúrovňové virtuálne stroje (napr. Java) a systémové virtuálne stroje používané ako jednotlivcami tak na serveroch.

2.2.1 Hypervízor

Hypervízor je softvér, ktorý vytvára a zabezpečuje beh virtuálnych strojov. Rozlišujeme 2 typy:

natívny - na hostujúcom počítači nie je nainštalovaný žiadny operačný systém. Hypervízor spravuje hardvér hostujúceho počítača a kontroluje beh operačných systémov, ktoré sa javia ako procesy. Príkladom je VMware ESX/ESXi, Oracle VM Server for x86 alebo Citrix XenServer.

hostovaný - hypervízor je spustený ako bežný program v operačnom systéme hostujúceho počítača. Príkladom je QEMU, VMware Workstation alebo VirtualBox.

[?]

2.3 Aplikačné kontajnery

Kontajnery predstavujú odlišný prístup k virtualizácii ako virtuálne stroje. Tiež ide o snahu spúšťať softvér v prostredí oddelenom od skutočného hardvéru a operačného systému. Na rozdiel od úplných virtuálnych strojov nie je virtualizovaný celý hardvér, ale len softvérové vybavenie nevyhnutné na spustenie programu. Rozdiel v architektúre ilustruje nasledovný obrázok: V kontajneri môže byť spustený nezávislý operačný systém spolu s požadovanými aplikáciami. Hostovský počítač, na ktorom sú kontajnery spustené, má jeden operačný systém a jednu množinu prostriedkov, ktoré tieto kontajnery zdieľajú. Jednotlivé kontajnery dostávajú kontrolovaný prístup k výpočtovému výkonu, pamäti, úložnej kapacite, sieti a prípadne ďalším prostriedkom.

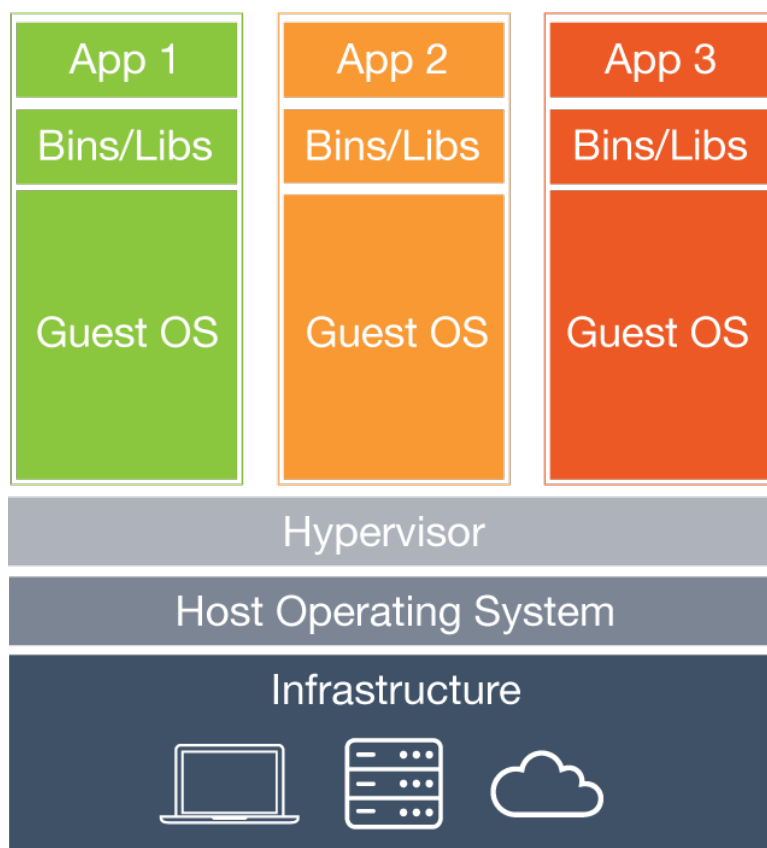


Figure 2.1: Porovnanie architektúry Docker a virtuálnych strojov
[?]

2.4 MapReduce aplikačné prostredia

2.5 Cloudové technológie MetaCentra

2.5.1 libvirt/KVM

KVM¹ je plné virtualizačné riešenie pre Linux pre x86 hardvér, obsahujúce virtualizačné rozšírenia (Intel VT or AMD-V). Pozostáva z nahrateľného modulu jadra, `kvm.ko`, ktorý poskytuje základ virtualizačnej infraštruktúry a špecifický modul, `kvm-intel.ko` alebo `kvm-amd.ko`. Je možné virtualizovať obrazy s operačným systémami Linux aj Windows. Každý virtuálny

1. Kernel-based Virtual Machine

stroj má vlastný virtualizovaný hardvér: sieťovú kartu, disk, grafický adaptér, atď. KVM je open-source softvér. Virtualizačný modul jadra sa nachádza v Linuxe od verzie 2.6.20. [?]

libvirt je sada nástrojov na prácu s virtualizačnými schopnosťami Linuxu (a ostatných OS). Je to voľný softvér dostupný pod licenciou GNU LGPL. Obsahuje API v jazyku C a väzby pre bežné programovacie jazyky. [?]

2.5.2 Docker

Docker umožňuje zabaliť aplikáciu so všetkými jej závislosťami do štandardizovanej jednotky (tzv. kontajner) určenej na softvérový vývoj. Kontajner Dockeru obalujú softvér kompletným súborovým systémom, ktorý zahŕňa všetko, čo daný softvér potrebuje na spustenie: kód, nástroje potrebné na beh, systémové nástroje a knižnice. Toto zaručuje, že program bude pracovať rovnako bez ohľadu na prostredie, v ktorom je spustený. [?]

2.5.3 Hadoop

Projekt Apache Hadoop vyvíja open-source softvér na spoľahlivé, škálovateľné, distribuované výpočty. Apache Hadoop je prostredie, ktoré umožňuje distribuované spracovávanie veľkých množstiev dát naprieč clustermi, používajúcimi jednoduché programovacie modely. Je navrhnutý tak, aby bol škálovateľný od jednotlivých serverov po tisíce strojov, kde každý poskytuje lokálny výpočetný výkon a úložný priestor. Nespolieha sa na vysokú dostupnosť hardvérových prostriedkov, ale je navrhnutý, aby detekoval a zvládal chyby na aplikačnej vrstve, takže poskytuje vysoko dostupnú službu nad clusterom počítačov, z ktorých každý je náchylný na chyby.

Projekt pozostáva z týchto modulov:

Hadoop Common: spoločné nástroje, ktoré podporujú ostatné Hadoop moduly

Hadoop Distributed File System (HDFSTM): distribuovaný súborový systém, ktorý poskytuje vysokú priepustnosť

Hadoop YARN: prostredie pre plánovanie úloh a správu zdrojov clusteru

Hadoop MapReduce: systém založený na YARN pre paralelné spracovávanie veľkých množstiev dát, prostredie pre plánovanie úloh a správu zdrojov clusteru

Chapter 3

Zber a uchovávanie monitorovacích dát

Monitorovanie akejkoľvek (nielen počítačovej) prevádzky je dôležitá oblasť, ktorá má význam pre jej správne fungovanie a jej správu, zdokonaľovanie a servis. Na priblíženie uvediem príklad železničnej spoločnosti. Je vhodné vedieť, koľko ľudí prepraví spoločnosť na určitom spoji. Tak môže identifikovať nerentabilné spoje, prípadne spoje, ktoré sú preťažené a je potrebné nejakým spôsobom zväčšiť ich kapacitu. Iným príkladom môže počet vlakov, ktoré prejdú za deň jednou stanicou. Vo všeobecnosti sa teda jedná o vyťaženie zdrojov vlakovkej spoločnosti, či už sú to samotné vlaky alebo stanice, prípadne personál atď. Monitorovanie nám poskytuje dáta, ktoré majú viacero aplikácií:

Prehľad o aktuálnom vyťažení zdrojov Pri zbere dát toto zodpovedá meraniu konkrétnej hodnoty sledovaného parametru. Kvantifikáciou sledovaného zdroja dostaneme predstavu o tom, ako je používaný "teraz", čiže v prítomnosti. Môžeme tak detegovať prípadné preťaženie zdroja, jeho dostupnosť alebo zlyhanie.

Prehľad o vyťažení zdrojov za určitú dobu Vyťaženie zdrojov sa v priebehu času mení. Nielen z krátkodobého hľadiska, kedy napríklad viac ľudí cestuje ráno za prácou a do školy, ale aj zo strednodobého (v období leta ľudia viac cestujú na dovolenky, čiže sa menia vyťažené spoje) a dlhodobého hľadiska (cestujúcich stále pribúda). Uchovávanie monitorovacích dát predstavuje kľúčový požiadavok, aby sme mohli sledovať, aké trendy vo využívaní mali jednotlivé sledované zdroje v priebehu nejakej doby. Z takto nahromadených dát môžeme vypočítavať rôzne štatistiky a ďalej ich analyzovať.

Prehľad o vyťažení zdrojov podľa parametrov Vyťaženie konkrétneho vlaku môže predstavovať jednoduché číslo, ktoré reprezentuje, koľko ľudí sa nachádzalo vo vlaku. Spoločnosť ale prevádzkuje viacero druhov vlakov a množstvo liniek na rôznych. Chceme preto mať možnosť zistiť napr. tom, ako bol využívaný rýchlik na trati Brno-Břeclav

cez prázdniny v utorok ráno. Okrem jednoduchého tvrdenia, že vo vlaku bol nejaký počet cestujúcich, nám dobre nastavené monitorovanie poskytuje aj komplexnejšie informácie.

Prehľad o stave zdrojov z hľadiska správy a plánovania Nezanedbateľný význam má monitorovanie z hľadiska údržby a servisu poskytovaných služieb. Z dostupných dát vieme identifikovať nefunkčný zdroj a nahradiť ho novým. V prípade vysokého vyťaženia zas môžeme vyvodiť záver, že zdroje už nie sú dostačujúce a je potrebné ich nejakým spôsobom rozšíriť prípadne zlepšiť efektivitu ich využívania. Takisto vieme do určitej miery predpovedať, ako budú zdroje v budúcnosti využívané, čo je podstatné pri samotnom plánovaní odstávok, servisnej činnosti alebo dočasného rozšírenia zdrojov.

Účtovanie vyťaženia zdrojov Každá poskytovaná služba má svojho spotrebiteľa. Či už sa jedná o zákazníkov železničnej spoločnosti alebo klientov výpočetného strediska. Dáta o používaní zdrojov sú jediným možným spôsobom ako rozumne stanoviť cenu za používané zdroje. Taktiež na ich základe môžeme sledovať činnosť jednotlivých užívateľov v systémoch a stanoviť im akceptovateľné podmienky na využívanie služieb.

3.1 Všeobecné problémy monitorovania

Problematika monitorovania v sebe zahŕňa viacero aspektov, ktoré možno oddeliť, zistiť pre ne najlepšie riešenia a spojiť ich tak do funkčného celku. V prvom rade je potrebné identifikovať, čo vlastne potrebujeme pre konkrétny systém monitorovať. Jedná sa o určenie zdrojov kľúčových pre dané prostredie. Pre železničnú spoločnosť to môže byť vyťaženie vlakov alebo množstvo spotrebovanej energie. Niektoré údaje, ktoré potrebujeme vedieť, je možné odvodiť z iných už nameraných hodnôt. Takéto meriky budeme nazývať sekundárne. Nakoľko môžu byť dopočítané, nie je vždy nevyhnutné ich v čase sledovať a uchovávať.

Namerané hodnoty metrík je okrem neskoršej štatistickej analýzy potrebné sledovať a analyzovať v reálnom čase. Metrika môže mať svoje kritické hodnoty. Sú to hodnoty, ktoré naznačujú, že daný zdroj sa vymyká bežným očakávaniam o využití. Na toto je vhodné reagovať. Či už sa jedná o zapísanie hlásenia do logu, zobrazenie varovania alebo prípadné zaslanie emailu či sms správy s popisom udalosti.

Ďalším čiastkovým problémom je granularita metrík. Zdroje sú kontin-

uálne využívané v čase. Tieto dáta je potrebné nejakým spôsobom digitalizovať. V určitom momente sa pozrieme na zdroj, kvantifikujeme, do akej miery je využitý a danú hodnotu zobrazíme prípadne uložíme. Meranie je po uplynutí určitého intervalu zopakovať, aby sme získali opäť aktuálne údaje. Rôzne zdroje sa menia v čase rôznou intenzitou. Napr. kapacita vlaku sa v priebehu jazdy mení zriedkavo, zatiaľčo počet cestujúcich sa mení v každej stanici. Je preto dôležité určiť aj periodicitu zberania jednotlivých metrík. Niektoré metriky sa môžu meniť takým spôsobom, že nie je efektívne sledovať ich zmenu pravidelne v nejakých intervaloch. Namiesto toho je efektívnejšie pri zmene daného parametru túto zmenu ohlásiť spolu s novou hodnotou. U vlakov je to napríklad kapacita vlaku, ktorá sa mení len zriedkavo počas jazdy, napríklad pri prepriahaní vozňov z jedného vlaku do druhého.

Množstvo zdrojov predstavuje ďalšiu oblasť, ktorú je potrebné zvážiť. Môžeme disponovať rôznym počtom zdrojov viacerých druhov, rozmiestnenými na viacerých miestach pod správou mnohých ľudí - prípadne oddelení. Metrické dáta vo svojej podstate len kvantifikujú využitie zdroja vo všeobecnosti. Zaobchádzajú s ním v tom zmysle, ako by bol len jeden. Napr. metrika vyťaženie vlaku v počte osôb. Nehovoria nič bližšie o tom, o aký zdroj sa jedná, kde je ho možné nájsť. V rámci širšieho systému je preto dôležitá jednoznačná identifikácia daného zdroja. O danom vlaku preto chceme zistiť napr. jeho názov, linku na ktorej premáva a prípadné ďalšie technické údaje. Takéto údaje sa nazývajú metadáta - čiže dáta o (v tomto prípade o metrických) dátach.

Uchovávanie metrických dát je kľúčovou súčasťou monitorovacieho systému. Doba, po ktorú vlastníci zdrojov uchovávajú metrické dáta, sa líši systém od systému. V niektorých oblastiach to je dokonca upravené zákonmi - napr. telekomunikácie, v iných je to na zvážení majiteľa infraštruktúry. Množstvo dát, ktoré je treba uložiť, závisí od viacerých parametrov. Sú nimi počet druhov zdrojov, počet jednotlivých inštancií zdrojov, počet metrík, ktoré sa o zdrojoch uchovávajú, periodičita zbieraných metrík a do určitej miery aj spôsob identifikácie daného zdroja. Prikladám tabuľku, ktorá ilustruje ako sa mení počet záznamov v závislosti na uvedených parametroch.

Ukladanie metrických dát predstavuje zároveň úzke hrdlo monitorovacieho riešenia. Zistenie hodnôt metrík pre kontrétny predstavuje v zásade serializovanú činnosť. V určitom čase sú všetky zdroje opýtané na to, ako sú vyťažené. Po jednom zistia svoje hodnoty a odošlú ich na uloženie. Množstvo metrík pre jeden zdroj je v bežnej praxi len zlomkom množstva zdrojov. Na časť systému, ktorá je zodpovedná za ukladanie metrík, sú

preto v pravidelných intervaloch kladené pomerne veľké nároky. Je ich už ale možné spracovávať paralelne. Centralizácia úložiska by znamenala prílišnú záťaž a monitorovacie uzly, preto je vhodné, aby takého úložisko bolo distribuované a prispôsobené na paralelné spracovávanie veľkých objemov dát.

3.1.1 Monitorovanie v heterogénnom výpočetnom prostredí

V oblasti cloudového a gridového výpočetného prostredia je pre monitorovanie kľúčové sledovať vyťaženie výpočetných zdrojov. Konkrétne procesor, výpočetná pamäť, trvalý ukladací priestor a sieťová infraštruktúra. Toto je spoločné pre všetky technológie a okrem týchto zdrojov je vhodné sledovať aj parametre špecifické pre jednotlivé oblasti. Tým sa budem konkrétne venovať v kapitole Metriky.

Automatizované monitorovanie výpočetného prostredia so sebou prináša aj špecifické problémy. Odvíjajú sa od toho, že monitorovanie je vykonávané strojovo (tj. pomocou počítača) a monitorované zdroje sú takisto stroje, ktorých stav z pohľadu využitia sa mení veľmi dynamicky.

Rozdielny je aj prístup jednotlivých technológií v tom ako poskytujú výpočetné zdroje. Aplikačné kontajnery a virtuálne stroje sa snažia rozdeliť celú infraštruktúru na menšie viac-menej uzavreté celky, ktoré sú potom sprostredkované užívateľom. Každý tento celok má pridelené zdroje, ktoré potom využíva. Množstvo týchto zdrojov je možné meniť, ale súvisí to s reštartovaním virtuálneho stroja alebo kontajnera. Z pohľadu gridu sú zanedbávané zdroje jedného uzla a na vypočítavanie úloh sa berú do úvahy všetky zdroje celého clusteru.

Problematika intervalov

U vysoko výkonných výpočetných systémov sa jednotlivé operácie vykonávajú vo veľmi krátkych časových intervaloch. Rádovo sú to nanosekundy. Jednotlivé procesy a výpočetné úlohy dostávajú na krátky čas k dispozícii všetky zdroje, čím sa z vyššieho pohľadu zabezpečuje paralelizácia. To spôsobuje vyťaženie množstva inštancií zdrojov mnohými účastníkmi. Celkový stav infraštruktúry z pohľadu vykonaných úloh a prenesených dát sa preto mení veľmi dynamicky a je vhodné zbierať dáta o zdrojoch v rozmedzí jednej až troch sekúnd.

Nezanedbateľnú rolu hrá aj čas potrebný na získanie jedného takéhoto "snímku" využitia zdrojov. Kým v prípade vlakov je napr. dostatok času na spočítanie cestujúcich medzi jednotlivými zastávkami, v prípade het-

erogénneho prostredia môže byť tento čas kritický. Ak sa každé dve sekundy pýtame na metrické údaje, očakávame, že odpoveď príde v kratšom intervale. V najlepšom prípade by táto doba odpovede mala byť len malý zlomok periódy danej metriky. Vopred nevieme, koľko bude trvať, kým príde odpoveď. To je v priamej závislosti od toho, koľko subjektov infraštruktúru využíva a koľko úloh spúšťa. Ak je však doba odpovede dlhšia ako samotný interval metriky, je to potrebné nejakým spôsobom riešiť. Nie je vhodné automaticky paralelne vygenerovať ďalšiu požiadavku na "snímku". Žiaducejším spôsobom riadenia je také, ktoré počká, kým sa daná požiadavka vykoná, a potom v nasledovnom intervale je meranie vykonané znova. Týmto sa predíde nadmernej záťaži celého systému.

To, ako sa systém vysporiadava z takýmto neočakávaným správaním, som zohľadňoval pri jednotlivých dostupných softvérových riešeniach.

Detekcia nových zdrojov a užívateľov

Heterogénna štruktúra v čase mení aj množstvo a kapacitu svojich poskytovaných zdrojov. Nie je možné staticky definovať zoznam procesorov, diskov, alebo virtuálnych strojov ktoré treba monitorovať. Podobne je to aj s užívateľmi. Proces vytvárania nových užívateľov je zautomatizovaný, takisto užívatelia môžu automatizovať vykonávanie svojich výpočetných úloh. Monitorovacie riešenie sa preto musí vedieť vysporiadať s týmito dynamickými zmenami, musí ich vedieť automaticky detegovať a zbierať o nich metrické dáta.

3.2 Časové rady

Monitorovacie dáta vo svojej podstate predstavujú časové rady. Časová rada je sekvencia dát, kde danému časovému okamihu zodpovedá jedna hodnota. Príkladom je zaznamenávanie teplôt v priebehu roka, výšky oceánskeho prílivu alebo množstvo áut, ktoré za určitú dobu prejde jedným bodom diaľnice. Efektívnou metódou vizualizácie dát časových rád sú čiarové grafy. Horizontálna os reprezentuje plynutie času a na vertikálnej osi sú znázornené hodnoty v danom čase.

3.2.1 Analýza časových rád

Analýza časových rád sa primárne zaoberá získavaním štatistík o zozbieraných dátach, napr. priemerná teplota počas celého roka. Medzi ďalšie úlohy patrí:

Exploračná analýza dát

Aproximácia na funkciu

Predpovedanie

Klasifikácia

Mám sa viac o nich rozpísať?

3.2.2 Zaobchádzanie s historickými dátami

Niečo o tom, ako sa v rámci šetrenia priestoru dáta zhľukujú a vypočítavajú sa agregované štatistiky.

Chapter 4

Aktuálne monitorovacie riešenia

Problematike monitorovania softvéru a infraštruktúry sa venuje viacero komerčných alebo open-source aplikácií, prípadne aplikácií zadarmo. Rôznymi technológiami riešia zber dát, ich uchovávanie, vizualizáciu a operácie nad nimi. Najprv sa venujem popisu aplikácií, ktoré súvisia so samotným monitorovaním dát, neskôr rozoberám dostupné riešenia v oblasti uchovávania časových rád.

4.1 Linux cgroups

Linux cgroups je technológia linuxového jadra, ktorá umožňuje limitovať, sledovať a izolovať spotrebu prostriedkov systému jednotlivými procesmi. Zavádza stromovo organizované kontrolné skupiny. Kontrolná skupina obsahuje obmedzenia pre jeden systémový prostriedok, tzv. subsystem. Príklad subsystemov, ktoré poskytuje Red Hat Enterprise Linux:

blkio - this subsystem sets limits on input/output access to and from block devices such as physical drives (disk, solid state, USB, etc.).

cpu - this subsystem uses the scheduler to provide cgroup tasks access to the CPU.

cpuacct - this subsystem generates automatic reports on CPU resources used by tasks in a cgroup.

cpuset - this subsystem assigns individual CPUs (on a multicore system) and memory nodes to tasks in a cgroup.

devices - his subsystem allows or denies access to devices by tasks in a cgroup.

freezer - this subsystem suspends or resumes tasks in a cgroup.

memory - this subsystem sets limits on memory use by tasks in a cgroup, and generates automatic reports on memory resources used by those tasks.

net_cls - this subsystem tags network packets with a class identifier (classid) that allows the Linux traffic controller (tc) to identify packets originating from a particular cgroup task.

net_prio - this subsystem provides a way to dynamically set the priority of network traffic per network interface.

ns - the namespace subsystem.

Pre každý subsystém existuje jeden strom kontrolných skupín. Ďalej skupina obsahuje zoznam procesov, ktoré podliehajú definovaným obmedzeniam. V strome kontrolných skupín sa proces vyskytuje len raz. V cloudovom prostredí však jednotlivé technológie využívajú a spúšťajú mnoho procesov, ktoré by sa obtiažne prirad'ovalo jednotlivým užívateľom v súvislosti so spustenými aplikáciami, preto tento spôsob monitorovania nie je úplne vhodný.

4.2 Nagios

Nagios je aplikácia, ktorá poskytuje komplexné riešenie na monitorovanie systémov. Poskytuje informácie o kľúčových komponentoch infraštruktúry vrátane aplikácií, služieb, operačného systému, sieťových protokolov, systémových metrík a sieťovej infraštruktúry. [?]

4.2.1 PluginAPI

Scripts and executables must do two things (at a minimum) in order to function as Nagios plugins:

Exit with one of several possible return values Return at least one line of text output to STDOUT The inner workings of your plugin are unimportant to Nagios. Your plugin could check the status of a TCP port, run a database query, check disk free space, or do whatever else it needs to check something. The details will depend on what needs to be checked - that's up to you.

Return Code

Nagios determines the status of a host or service by evaluating the return code from plugins. The following tables shows a list of valid return codes, along with their corresponding service or host states.

Plugin Return Code Service State Host State

0 OK UP

1 WARNING UP or DOWN/UNREACHABLE*

2 CRITICAL DOWN/UNREACHABLE

3 UNKNOWN DOWN/UNREACHABLE

Note Note: If the `use_aggressive_host_checking` option is enabled, return codes of 1 will result in a host state of DOWN or UNREACHABLE. Otherwise return codes of 1 will result in a host state of UP. The process by which Nagios determines whether or not a host is DOWN or UNREACHABLE is discussed here.

Plugin Output Spec

At a minimum, plugins should return at least one of text output. Beginning with Nagios 3, plugins can optionally return multiple lines of output. Plugins may also return optional performance data that can be processed by external applications. The basic format for plugin output is shown below:

TEXT OUTPUT | OPTIONAL PERFDATA

LONG TEXT LINE 1

LONG TEXT LINE 2

...

LONG TEXT LINE N | PERFDATA LINE 2

PERFDATA LINE 3

...

PERFDATA LINE N

The performance data (shown in orange) is optional. If a plugin returns performance data in its output, it must separate the performance data from the other text output using a pipe (|) symbol. Additional lines of long text output (shown in blue) are also optional. [?]

4.2.2 Timeouty

Format: `service_check_timeout=<seconds>` Example: `service_check_timeout=60`

This is the maximum number of seconds that Nagios will allow service checks to run. If checks exceed this limit, they are killed and a CRITICAL state is returned. A timeout error will also be logged.

There is often widespread confusion as to what this option really does. It is meant to be used as a last ditch mechanism to kill off plugins which are misbehaving and not exiting in a timely manner. It should be set to something high (like 60 seconds or more), so that each service check normally finishes executing within this time limit. If a service check runs longer than this limit, Nagios will kill it off thinking it is a runaway processes.

[?]

Nie je mechanizmus pre reštart alebo individuálne timeouty. Má zas-taralý Hadoop plugin, používa staré REST API. Last Release Date 2010-11-05 [?]

4.3 Zabbix

4.3.1 Timeouty

Timeout processing Zabbix will not process a simple check longer than Timeout seconds defined in Zabbix server configuration file. [?]

Zabbix Agent Timeout no 1-30 3 Spend no more than Timeout seconds on processing [?]

Zabbix Server Timeout no 1-30 3 Specifies how long we wait for agent, SNMP device or external check (in seconds). [?]

Neuvádza sa nič o zastavovaní procesov.

4.3.2 Moduly

Loadable modules offer a performance-minded option for extending Zabbix functionality.

There already are ways of extending Zabbix functionality by way of:

user parameters (agent metrics) external checks (agent-less monitoring) system.run[] Zabbix agent item. They work very well, but have one major drawback, namely fork(). Zabbix has to fork a new process every time it handles a user metric, which is not good for performance. It is not a big deal normally, however it could be a serious issue when monitoring embedded systems, having a large number of monitored parameters or heavy scripts with complex logic or long startup time.

Zabbix 2.2 comes with support of loadable modules for extending Zabbix agent, server and proxy without sacrificing performance.

A loadable module is basically a shared library used by Zabbix daemon and loaded on startup. The library should contain certain functions, so that a Zabbix process may detect that the file is indeed a module it can load and work with.

Loadable modules have a number of benefits. Great performance and ability to implement any logic are very important, but perhaps the most important advantage is the ability to develop, use and share Zabbix modules. It contributes to trouble-free maintenance and helps to deliver new functionality easier and independently of the Zabbix core code base. [?]

4.4 Icinga

By default the Icinga 2 daemon is running as icinga user and group using the init script. Using Debian packages the user and group are set to nagios for historical reasons. [?]

4.4.1 Externé pluginy

Icinga determines the status of a host or service by evaluating the return code from plugins. The following tables shows a list of valid return codes, along with their corresponding service or host states [?]

4.4.2 Notifikácie a príkazy udalostí

Unlike notifications, event commands for hosts/services are called on every check execution if one of these conditions match: The host/service is in a soft state The host/service state changes into a hard state The host/service state recovers from a soft or hard state to OK/Up

Ide len o spustenie nejakého systémového príkazu. [?]

4.4.3 Integrácie s aplikáciami

these tiny pure shell+awk plugins for monitoring your hadoop cluster are a enhanced and uptodate version of exchange.nagios.org check_hadoop-dfs.sh [?]

4.5 Cacti

Cacti is a complete network graphing solution designed to harness the power of RRDTool's data storage and graphing functionality. Cacti provides a fast poller, advanced graph templating, multiple data acquisition methods, and user management features out of the box. [?]

4.6 collectd

There are some key differences we think set collectd apart. For one, it's written in C for performance and portability, allowing it to run on systems without scripting language or cron daemon, such as embedded systems. At the same time it includes optimizations and features to handle hundreds of thousands of data sets. It comes with over 90 plugins. It provides powerful networking features and is extensible in numerous ways.

4.6.1 Obmedzenia

It does not generate graphs. It can write to RRD files, but it cannot generate graphs from these files. Monitoring functionality has been added in version 4.3, but is so far limited to simple threshold checking. [?]

4.6.2 Zapisovací plugin Write TSDB

The Write TSDB plugin writes metrics to OpenTSDB, an open-source distributed time-series database based on Apache HBase. [?]

Host Address Hostname or address to connect to. Defaults to localhost.

Port Service Service name or port number to connect to. Defaults to 4242.

HostTags String When set, HostTags is added to the end of the metric. It is intended to be used for name=value pairs that the TSD will tag the metric with. Dots and whitespace are not escaped in this string.

StoreRates false / true If set to true, convert counter values to rates. If set to false (the default) counter values are stored as is, as an increasing integer number.

AlwaysAppendDS false / true If set the true, append the name of the Data Source (DS) to the "metric" identifier. If set to false (the default), this is only done when there is more than one DS.

4.6.3 Prahy a notifikácie

The only action the Threshold plugin takes itself is to generate and dispatch a notification. Every time a value is out of range, notification is dispatched. Also, all values that match a threshold are considered to be relevant or "interesting". As a consequence collectd will issue a notification if they are not received for Timeout iterations. for example, Timeout is set to "2" (the default) and some hosts sends it's CPU statistics to the server every 60 seconds, a notification will be dispatched after about 120 seconds. It may take a little longer because the timeout is checked only once each Interval on the server.

When a value comes within range again or is received after it was missing, an "OKAY-notification" is dispatched. [?]

4.6.4 Ostatne

Zenoss

Munin

4.7 Ganglia

Ganglia is a scalable distributed monitoring system for high-performance computing systems such as clusters and Grids. It is based on a hierarchical design targeted at federations of clusters. It leverages widely used technologies such as XML for data representation, XDR for compact, portable data transport, and RRDtool for data storage and visualization. It uses carefully engineered data structures and algorithms to achieve very low per-node overheads and high concurrency. The implementation is robust, has been ported to an extensive set of operating systems and processor architectures, and is currently in use on thousands of clusters around the world. It has been used to link clusters across university campuses and around the world and can scale to handle clusters with 2000 nodes. Ganglia is a BSD-licensed open-source project that grew out of the University of California, Berkeley Millennium Project [?] Pre Gangliu je dostupný monitorovací plugin pre Hadoop. [?]

4.8 Scout

Scout runs within Docker containers without any special configuration. [?]

4.9 New Relic

<http://newrelic.com/docker>

4.10 AppDynamics

<https://www.appdynamics.com/community/exchange/extension/docker-monitoring-extension/>

This extension works only with the standalone machine agent.

Metrics include:

Hadoop Resource Manager App status and progress: submitted, pending, running, completed, killed, and failed app count
Memory size, memory usage
Allocated containers, container count in different states
Node sta-

tus, count of nodes in different states Scheduler capacity, app and container count Ambari Individual host metrics including CPU, disk, memory, JVM, load, network, process, and RPC metrics Service component metrics including CPU, disk, memory, JVM, load, network, process, RPC, and component-specific metrics [?]

4.11 OpenTSDB

OpenTSDB je databáza na uchovávanie a sprístupňovanie veľkých objemov časových dát. Pozostáva z Time Series Daemon (TSD) a z utilít pre príkazový riadok. Interakcia s OpenTSDB je primárne realizovaná cez jedného alebo viacerých TSD. Každý TSD je nezávislý. Neexistuje žiadny riadiaci proces, žiadny zdieľaný stav, takže je možné spustiť toľko TSD, koľko je potrebné na zvládnutie požadovanej záťaže. Každý TSD používa open-source databázu HBase na ukladanie a vyberanie dát časových rád. HBase schéma je vysoko optimalizovaná na rýchlu agregáciu podobných časových rád, aby minimalizovala požiadavky na úložný priestor. Používatelia TSD nemusia pristupovať do HBase priamo. S TSD je možné komunikovať cez jednoduchý protokol podobný Telnetu, cez HTTP API alebo cez jednoduché GUI. Všetka komunikácia sa deje na tom istom porte (TSD odhadne protokol klienta pohľadom na prvých niekoľko bajtov, ktoré obdrží). [?]

Na vizualizáciu dát existuje nástroj Metrilyx. Je to open-source webový engine, ktorý vytvára grafy zo zhromaždených dát. Je možné meniť časové rozpätie, za ktoré sa majú grafy metrík zobrazíť.

4.12 InfluxDB

InfluxDB je platforma na zbieranie, uchovávanie, vizualizáciu a správu časových dát. Užívateľ môže vytvoriť viacero databáz. Dáta sa zapisujú a čítajú pomocou rozhrania príkazového riadka, rôznych klientskych knižníc, alebo pomocou HTTP API. Na vizualizáciu dát používa modul *chronograf*.

4.12.1 Politiky udržiavania

Každá databáza obsahuje pravidlá, ktoré definujú, po akú dobu majú byť ukladané dáta časových rád a koľko kópií dát má byť vytvorených. Jedna databáza môže mať niekoľko takýchto politík. Pri zápise do nej je možné špecifikovať, ktorá politika sa má pre zápis použiť. Pri vytvorení databázy je automaticky vytvorená jedna politika

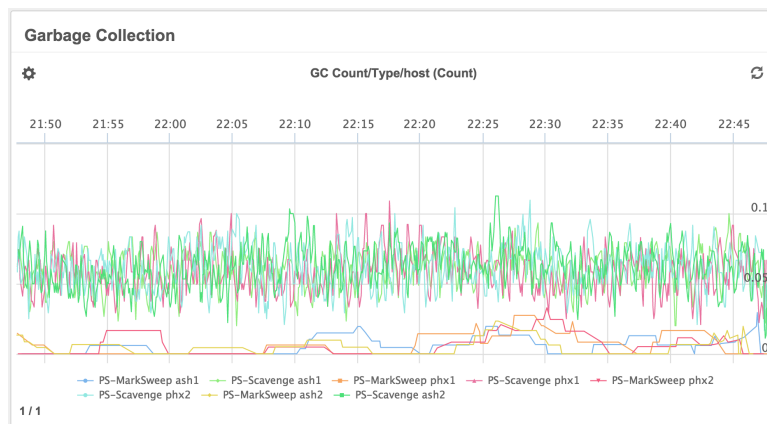


Figure 4.1: Vizualizácia časových rád OpenTSDB pomocou Metrilyx
[?]

4.12.2 Kontinuálne dotazovanie

Databáza je schopná periodicky vykonávať požiadavky na dáta. Cieľom je zmenšovanie objemu dát. Ide o zhlukovanie dát s vysokou frekvenciou zberu, čím vzniknú dáta s menšou hustotou zberu. Táto hodnota je potom zvyčajne uložená do inej databázy.

4.13 RRDTool

RRDTool je nástroj na uchovávanie, spravovanie a vizualizáciu časových dát. Využíva round-robin databázu. Je to databáza s dopredu danou maximálnou veľkosťou. V prípade, že príde požiadavka na zápis hodnoty a databáza je už plná, dôjde k prepisu najstaršej hodnoty. Tento nástroj takisto obsahuje funkcie na konsolidáciu dát. Konsolidovaná hodnota je typicky priemer, minimum alebo maximum z viacerých hodnôt zozbieraných za dlhší časový úsek. Tieto hodnoty sú ukladané do round-robin archívu.

Chapter 5

Metriky

Každá z technológií cloudového výpočetného strediska je zdrojom mnohých dát o vyťažení zdrojov. Je potrebné určiť, ktoré zdroje monitorovať a ktoré metriky zbierať. Je vhodné mať také metriky pre všetky využívané cloudové technológie, ktoré je možné nejakým spôsobom porovnať medzi sebou. V prípade procesora sa jedná o jeho aktuálne vyťaženie, ale zároveň je zaujímavým údajom aj prepočítaný procesorový čas. Táto metrika môže byť efektívnym nástrojom pre následné účtovanie jednotlivým užívateľom. Keďže majú ale tieto technológie principiálny rozdiel vo svojom určení, nie je možné vždy o každom zdroji zbierať rovnaké dáta. O distribuovaných výpočtoch napríklad nie je možné efektívne zistiť aktuálne vyťaženie procesora. Takéto úlohy sa počítajú na viacerých uzloch clusteru. O poradí a rozmiestnení požiadaviek na zdroje rozhoduje riadiaca aplikácia, takže sa nejedná o jeden procesor. Táto aplikácia má ale prehľad o tom, koľko času strávil celý cluster počítaním danej úlohy. Takže v určitom ohľade zrovnateľná s virtuálnym strojom a jeho spotrebou procesorového času.

Podobná situácia nastáva aj u monitorovania sieťových rozhraní. Distribuované výpočty využívajú sieť svojším spôsobom a len na účel výpočítania komplexnejšieho problému. Jedná sa o prepojenie uzlov v rámci clusteru. Je to jednoúčelová vysokorýchlostná sieť. Z hľadiska poskytovania výpočetných kapacít dáva väčší zmysel orientácia na využívanie konektivity smerom do internetu.

5.0.1 Význam popisných údajov metrík

Metrické dáta vypovedajú o využití zdrojov. Vieme určiť ako dlhý čas a v akej miere bol využívaný výpočtový výkon. Z nich samotných nevieme presne určiť, kto zdroje využíval. Pre to aby mali tieto dáta zmysel pre neskoršie účtovanie, je potrebné mať možnosť ich priradiť k užívateľom, či už sa jedná o vlastníka virtuálneho stroja alebo užívateľa, ktorý spustil konkrétnu úlohu. Na základe týchto ďalších popisných údajov je potom možné zisťovať, kto zdroje vyťažoval za časové obdobie najviac a podľa toho

stanoviť prípadnú cenu za používanie zdrojov. Okrem identity je vhodné metriky popisovať aj ďalšími údajmi, ako je miesto, kde sa zdroj nachádza, názov stroja, ktorý poskytuje svoje kapacity, a ďalšie špecifické údaje, ktoré sa týkajú jednotlivých technológií, ktoré budem popisovať konkrétnejšie v ďalších častiach.

5.0.2 Periodicita zberania metrík

Zberané metriky sa líšia tým, ako veľmi sa v čase menia. Kým záťaž procesora, vstupno-výstupné operácie alebo množstvo prenesených dát sa mení v čase pomerne rýchlo, veľkosť clustera v počte poskytnutých procesorov alebo virtuálnej pamäte sa nemení tak často. Má preto zmysel uvažovať o kontrole intervalu zbierania jednotlivých metrík. Nie len na úrovni zhluku metrík pre konkrétnu technológiu ale aj pre jednotlivé metriky samostatne.

5.0.3 Formát hodnoty metriky

Metriky ako záznamy obsahujú čas, kedy bola hodnota nameraná, samotnú hodnotu nejakého sledovaného javu a popisné dáta. Na metrické hodnoty sa môžeme pozerieť ako na prírastky alebo ako na absolútne hodnoty. Prírastky hovoria o rozdielne aktuálnej nameranej hodnoty a poslednej hodnoty. Absolútne hodnoty predstavujú aktuálnu nameranú hodnotu využitia zdroja.

5.1 Docker

5.1.1 Sieť

Aby mohli medzi sebou jednotlivé kontajnery komunikovať, Docker im poskytuje sieťové rozhrania. Každé rozhranie má nakonfigurovanú sieť, do ktorej patrí. Na to, aby kontajnery spolu mohli komunikovať, musia byť členmi rovnakej siete. Komunikácia naprieč sieťami nie je možná. Užívatelia si môžu definovať vlastné siete. Docker na vytvorenie týchto sietí poskytuje dva ovládače.

Sieť typu most

Je to jednoduchý typ siete určený pre malé siete. Je ju možné vytvoriť príkazom

```
$ docker network create --driver bridge NÁZOV_SIETE
```


Po vytvorení siete je možné spustiť kontajnery v tejto sieti príkazom

```
$ docker run --net=NÁZOV_SIETE --name=NÁZOV_KONTAJNERA
```

Prekladaná sieť

Docker umožňuje vytvoriť aj sieť, v ktorej sa nachádza viacero host'ujúcich počítačov zároveň. To umožňuje komunikovať medzi sebou aj kontajnerom, ktoré sú spustené v rozličných sieťach, prípadne na inom host'ujúcom počítači..

Metriky siete

Pre jednotlivé sieťové rozhrania je možné zbierať tieto metriky:

rx_bytes - počet prijatých bajtov

rx_dropped - počet prichádzajúcich zahodených bajtov

rx_error - počet chybných bajtov

rx_packets - počet prijatých paketov

tx_bytes - počet odoslaných bajtov

tx_dropped - počet zahodených bajtov pri pokuse o odoslanie

tx_errors - počet odoslaných chybných bajtov

tx_packets - počet odoslaných paketov

5.1.2 Pamäť

Cez API Dockeru je možné získať nasledovné metriky pamäte.

usage - spotreba pamäte

failcnt - počet chýb

Docker neposkytuje údaj o tom, koľko pamäte poskytuje host'ujúci počítač. Zistiť tento údaj však v implementácii nepredstavuje problém, a preto je tiež zberaná táto metrika.

5.1.3 Procesor

Procesor predstavuje jeden z najdôležitejších údajov o vyťažení zdrojov. Docker poskytuje viaceré metriky o procesore, ktorých hodnoty predstavujú výpočetný čas strávený na procesore:

percpu_usage - využitie jednotlivých jadier procesora

usage_in_usermode -

total_usage -

usage_in_kernelmode -

system_cpu_usage -

5.2 libvirt/KVM

5.2.1 Metriky siete

Pre jednotlivé sieťové rozhrania je možné zbierať tieto metriky:

rx_bytes - počet prijatých bajtov

rx_dropped - počet prichádzajúcich zahodených bajtov

rx_error - počet chybných bajtov

rx_packets - počet prijatých paketov

tx_bytes - počet odoslaných bajtov

tx_dropped - počet zahodených bajtov pri pokuse o odoslanie

tx_errors - počet odoslaných chybných bajtov

tx_packets - počet odoslaných paketov

5.2.2 Metriky pamäte

VIR_DOMAIN_MEMORY_STAT_SWAP_IN - The total amount of memory written out to swap space (in kB).

VIR_DOMAIN_MEMORY_STAT_SWAP_OUT - Page faults occur when a process makes a valid access to virtual memory that is not available. When servicing the page fault, if disk IO is required, it is considered a major fault. If not, it is a minor fault. These are expressed as the number of faults that have occurred.

VIR_DOMAIN_MEMORY_STAT_MAJOR_FAULT - počet chybných bajtov

VIR_DOMAIN_MEMORY_STAT_MINOR_FAULT - počet prijatých paketov

VIR_DOMAIN_MEMORY_STAT_UNUSED - The amount of memory left completely unused by the system. Memory that is available but used for reclaimable caches should NOT be reported as free. This value is expressed in kB.

VIR_DOMAIN_MEMORY_STAT_AVAILABLE - The total amount of usable memory as seen by the domain. This value may be less than the amount of memory assigned to the domain if a balloon driver is in use or if the guest OS does not initialize all assigned pages. This value is expressed in kB.

VIR_DOMAIN_MEMORY_STAT_ACTUAL_BALLOON - Current balloon value (in KB).

VIR_DOMAIN_MEMORY_STAT_RSS - Resident Set Size of the process running the domain. This value is in kB

VIR_DOMAIN_MEMORY_STAT_NR - The number of statistics supported by this version of the interface. To add new statistics, add them to the enum and increase this value.

5.2.3 Metriky zápisu dát

rd_req - number of read requests

rd_bytes - number of read bytes

wr_req - number of write requests

wr_bytes - number of written bytes

errs - In Xen this returns the mysterious 'oo_req'.

5.2.4 Metriky CPU

Libvirt poskytuje o procesore údaj o tom, koľko výpočetného času strávil daný virtuálny stroj na procesore hostujúceho počítača.

5.3 Hadoop

5.3.1 Cluster Metrics API

Toto API poskytuje metriky o celom clusteri.

appsSubmitted - The number of applications submitted

appsCompleted - The number of applications completed

appsPending - The number of applications pending

appsRunning - The number of applications running

appsFailed - The number of applications failed

appsKilled - The number of applications killed

reservedMB - The amount of memory reserved in MB

availableMB - The amount of memory available in MB

allocatedMB - The amount of memory allocated in MB

totalMB - The amount of total memory in MB

reservedVirtualCores - The number of reserved virtual cores

availableVirtualCores - The number of available virtual cores

allocatedVirtualCores - The number of allocated virtual cores

totalVirtualCores - The total number of virtual cores

containersAllocated - The number of containers allocated

containersReserved - The number of containers reserved

containersPending - The number of containers pending

totalNodes - The total number of nodes

activeNodes - The number of active nodes

lostNodes - The number of lost nodes

unhealthyNodes - The number of unhealthy nodes

decommissionedNodes - The number of nodes decommissioned

rebootedNodes - The number of nodes rebooted

5.3.2 Cluster Application API

Tento koncový bod API poskytuje informácie o jednotlivých aplikáciách, ktoré boli alebo sú spustené na clusteri.

id - The application id

user - The user who started the application

name - The application name

Application Type - The application type

queue - The queue the application was submitted to

state - The application state according to the ResourceManager - valid values are members of the YarnApplicationState enum: NEW, NEW_SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED

finalStatus - The final status of the application if finished - reported by the application itself - valid values are: UNDEFINED, SUCCEEDED, FAILED, KILLED

progress - The progress of the application as a percent

trackingUI - Where the tracking url is currently pointing - History (for history server) or ApplicationMaster

trackingUrl - The web URL that can be used to track the application

diagnostics - Detailed diagnostics information

clusterId - The cluster id

startedTime - The time in which application started (in ms since epoch)

finishedTime - The time in which the application finished (in ms since epoch)

elapsedTime - The elapsed time since the application started (in ms)

amContainerLogs - The URL of the application master container logs

amHostHttpAddress - The nodes http address of the application master

allocatedMB - The sum of memory in MB allocated to the application's running containers

allocatedVCores - The sum of virtual cores allocated to the application's running containers

runningContainers - The number of containers currently running for the application

memorySeconds - The amount of memory the application has allocated (megabyte-seconds)

vcCoreSeconds - The amount of CPU resources the application has allocated (virtual core-seconds)

cit

5.3.3 Node Application API

Toto API poskytuje metriky o využití jednotlivých uzlov clusteru.

rack - The rack location of this node

state - State of the node - valid values are: NEW, RUNNING, UNHEALTHY, DECOMMISSIONED, LOST, REBOOTED

id - The node id

nodeHostName - The host name of the node

nodeHTTPAddress - The nodes HTTP address

healthStatus - The health status of the node - Healthy or Unhealthy

healthReport - A detailed health report

lastHealthUpdate - The last time the node reported its health (in ms since epoch)

usedMemoryMB - The total amount of memory currently used on the node (in MB)

availMemoryMB - The total amount of memory currently available on the node (in MB)

usedVirtualCores - The total number of vCores currently used on the node

availableVirtualCores - The total number of vCores available on the node

numContainers - The total number of containers currently running on the node

Chapter 6

Analýza a návrh

6.1 MetaCentrum

Projekt MetaCentrum vznikol v roku 1996 a od roku 1999 je jeho činnosť zastrešovaná organizáciou CESNET. Zaoberá sa budovaním národnej gridovej infraštruktúry a prepojením s podobnými projektami za hranicami Českej republiky. Projekt je oficiálnou súčasťou Európskej gridovej iniciatívy (EGI). Úlohou MetaCentra je predovšetkým koordinácia a rozširovanie infraštruktúry či už o vlastné zdroje alebo prostredníctvom partnerov, ktorý poskytujú výpočetný výkon svojich clusterov. Jedná sa hlavne o akademickú spoluprácu. MetaCentrum spravuje výpočetné prostriedky a dátové úložiská AV, JČU, MU, MZLU, UK, VUT, ZČU. V súčasnosti disponuje (stav k 30.7. 2010) 1500 jadrami CPU, 100 TB využiteľnej diskovej kapacity v podobe poľa a 400 TB kapacity v podobe pások. Služby využíva 385 registrovaných aktívnych užívateľov, ktorí spolu na 750 tisíc úlohách využili 7 miliónov hodín procesorového času.

MetaCentrum primárne poskytuje svoj výpočetný výkon a úložnú kapacitu. Taktiež sprístupňuje svoje programové vybavenie a vývojové prostredie a hlavne množstvo aplikácií využívaných na výskumné účely, ako napr. Ansys, Gaussian, Matlab, Mathematica. Taktiež sa venuje vývoju v oblasti gridového a cloudového počítania, napr. v oblasti plánovania, gridového middleware, optimalizácie a paralelizácie výpočtov a virtualizácie infraštruktúry. Dôležitou funkciou je účasť na medzinárodných projektoch, využívanie medzinárodnej výpočetnej infraštruktúry a využívanie skúseností na rozvoj v domácom prostredí.

6.2 Požiadavky na aplikáciu

6.2.1 Vysoký monitorovací výkon

Cloudová infraštruktúra MetaCentra pozostáva z mnohých výpočetných uzlov. Sú prepojené sieťami s veľkou prenosovou kapacitou a vysokou

prieupustnosťou. Je potrebné zbierať dáta o využití množstva zapojených clusterov, na ktorých je tiež spustených mnoho výpočetných úloh či virtuálnych strojov. Metriky sú zbierané periodicky v určitých intervaloch z jednotlivých uzlov. Veľká záťaž je kladená na databázu, do ktorej sú tieto metrické dáta pravidelne odosielané.

6.2.2 Nízka nadbytočná záťaž

Primárnou úlohou cloudu je poskytovanie svojho výpočetného výkonu a prostriedkov. Je žiadúce, aby monitorovacia aplikácia predstavovala čo najmenšiu záťaž pre systémy, na ktorých je spustená. Ak by monitorovanie samotné spotrebovávalo príliš veľa zdrojov, takto získané metriky by nemali požadovanú presnosť. Je to možné doceliť výberom efektívneho programovacieho jazyka, napr. C, C++, prípadne rýchle skriptovacie jazyky ako Python, Ruby či Go. Nie je príliš vhodné používať jazyky, ktoré na svoj beh potrebujú ďalšiu vrstvu v podobe virtuálneho stroja, ako napr. Java.

6.2.3 Škálovateľnosť

Monitorovacia aplikácia by mala poskytovať zrovnateľné výsledky v oblasti rýchlosti odozvy v prípade, že bude spustená na jednom uzle, ale aj v prípade, že jej úlohou bude monitorovať desiatky až stovky výpočetných uzlov s množstvom spustených aplikácií. To je možné doceliť vhodným paralelizovaním dotazov na metriky

Monitorovacia aplikácia bude pozostávať z dvoch častí.

démon - jeho úlohou je zber metrík a odosielanie do databázy

zásúvné moduly - ich úlohou je zisťovanie metrických dát a odosielanie démonovi

Démon je program, ktorý je spustený raz a beží v systéme na pozadí. Podľa konfigurácie pri štarte zistí, ktoré moduly sa budú používať. Takisto dôjde ku konfigurácii jednotlivých modulov, napr. nastavenie potrebných ciest k požadovaným súborom. Následne dôjde k inicializácii jednotlivých modulov. V tejto fáze moduly inicializujú prostriedky, ktoré potrebujú v priebehu zberu metrík. Napr. pripojenie na správcu kontajnerov alebo hypervízora. Nie je efektívne, aby boli tieto prostriedky inicializované pri každej požiadavke na metriku, pretože by to spomaľovalo proces samotného zberu dát. Potom nasleduje fáza behu. Démon periodicky spúšťa jednotlivé moduly, ktoré zisťujú metrické dáta. Tie následne vracajú ako

odpoveď démonovi. V prípade, že je potrebné démona ukončiť, dôjde najprv k ukončeniu jednotlivých modulov. V tejto fáze moduly uvoľnia všetky prostriedky, ktoré mali naalokované.

6.3 Reakcia na dlhú odozvu modulu

Ak je hodnota nejakej metriky mimo určitý rozsah, je generované hlásenie. Na to je možné reagovať. Na situáciu, keď časť zodpovedná za zbieranie dát neodpovedá, je ale možné reagovať len reštartovaním celej monitorovacej aplikácie. Nie je možné jednotlivé pluginy ovládať nezávisle.

V princípe nejde nijako odlišiť, či daný modul čaká na údaje alebo došlo k chybe a modul neodpovedá. Preto bude potrebné vytvoriť niektoré pluginy tak, aby jedna časť bola neustále dostupná a reagovala na výzvy od riadiacej aplikácie. Bude definovaný časový interval na vrátenie hodnoty. V prípade ak plugin úspešne v časovom intervale zistil dané metrické dáta, vráti ich riadiacej aplikácii. V prípade, že v danom intervale plugin neobdržal metrické dáta, vráti poslednú hodnotu. Zároveň sa nebudú vytvárať nové požiadavky na tento údaj. Tento časový interval si bude môcť užívateľ nastaviť pre každú sondu. Predmetom testovania bude zistiť, aký interval by bol vhodný pre tú ktorú sondu.

Ďalšou prahovou hodnotou bude počet opakovaní, pri ktorých plugin vracia poslednú hodnotu danej metriky. Ak dôjde k prekročeniu tejto hodnoty, bude reštartovaná celá monitorovacia aplikácia.

6.4 OpenTSDB

Ako databázu na uchovávanie časových dát som si zvolil OpenTSDB. Dôvodom je používanie databázy HBase. Je to distribuovaná databáza určená pre veľké objemy dát v rádoch stoviek miliónov a milárd záznamov. Je typom NoSQL databázy. Oproti SQL databázam je linárne škálovateľná. Ak dôjde k zdvojnásobeniu výpočetných zdrojov, dôjde aj k zdvojnásobeniu výkonu databázy. To je dôležité pri zbere časových dát z mnohých uzlov, ktoré sa v gridovej infraštruktúre MetaCentra nachádzajú. Ďalším dôvodom je, že v MetaCentre je aktuálne databáza HBase využívaná.

Chapter 7

Implementácia

Využijem existujúcu aplikáciu na zbieranie metrických dát *collectd*. Na zbieranie jednotlivých metrík som vytvoril moduly pre tento program. Tieto údaje následne bude odosielať do databázy OpenTSDB pomocou modulu WriteTSDB.

7.1 Techniky zbierania metrík

7.1.1 Torque

Na zbieranie metrík som použil príkaz *qstat -j*.

Prints either for all pending jobs or the jobs contained in *job_list* various information. The *job_list* can contain *job_ids*, *job_names*, or wildcard expression *sge_types(1)*.

For jobs in E(rror) state the error reason is displayed. For jobs that could not be dispatched during in the last scheduling interval the obstacles are shown, if 'schedd_job_info' in *sched_conf(5)* is configured accordingly.

For running jobs available information on resource utilization is shown about consumed cpu time in seconds, integral memory usage in Gbytes seconds, amount of data transferred in io operations, current virtual memory utilization in Mbytes, and maximum virtual memory utilization in Mbytes. This information is not available if resource utilization retrieval is not supported for the OS platform where the job is hosted.

7.1.2 Docker

Budem zberať metriky o spotrebovávaných zdrojoch pre každý kontajner. Na komunikáciu s Dockerom je možné využiť:

príkazy aplikácie v príkazovom riadku

Remote API

Používanie príkazov aplikácie môže byť o čosi rýchlejšie, ale následne by bolo potrebné analyzovať textový výstup programu. Rozhodol som sa použiť Remote API. Toto API funguje pre účely monitorovania na princípe REST a odpovede vracia vo formáte JSON, čo predstavuje zjednodušenie spracovania výstupu. Démon Dockeru "počúva" na lokálnom sockete, čo by nemalo spôsobovať výrazné oneskorenie odpovede.

O každom kontajneri spustenom v Dockeri budem sledovať využívanie týchto zdrojov :

procesor

pevný disk

sieť

pamäť

Dáta je možné zberať buď ako prúd, alebo po jednorazových žiadostiach. To ešte nemám veľmi preštudované, neviem aké sú tam intervaly, tak sa rozhodnem až neskôr. Predpokladám ale, že z hľadiska rýchlosti a alokácie bude cesta asi ten stream.

7.1.3 libvirt/KVM

Na monitorovanie virtuálnych strojov použijem existujúcu implementáciu sond v jazyku Python. Je šírená pod licenciou open-source a vznikla v rámci organizácie Cesnet.

7.1.4 Hadoop

The Hadoop YARN web service REST APIs are a set of URI resources that give access to the cluster, nodes, applications, and application historical information. The URI resources are grouped into APIs based on the type of information returned. Some URI resources return collections while others return singletons. HTTP Requests

To invoke a REST API, your application calls an HTTP operation on the URI associated with a resource.

Summary of HTTP operations

Currently only GET is supported. It retrieves information about the resource specified.

Security

The web service REST API's go through the same security as the web UI. [?]

Kerberos

Chapter 8

Záver

Appendix A

Kapitola priloha