



**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**

BRNO UNIVERSITY OF TECHNOLOGY

**FAKULTA INFORMAČNÍCH TECHNOLOGIÍ**

FACULTY OF INFORMATION TECHNOLOGY

**ÚSTAV INFORMAČNÍCH SYSTÉMŮ**

DEPARTMENT OF INFORMATION SYSTEMS

**DISTRIBUOVANÝ REPOSITÁŘ DIGITÁLNÍCH FORENZ-  
NÍCH DAT**

DISTRIBUTED FORENSIC DIGITAL DATA REPOSITORY

**DIPLOMOVÁ PRÁCE**

MASTER'S THESIS

**AUTOR PRÁCE**

AUTHOR

**Bc. MARTIN JOSEFÍK**

**VEDOUCÍ PRÁCE**

SUPERVISOR

**RNDr. MAREK RYCHLÝ, Ph.D.**

**BRNO 2018**

## **Abstrakt**

Do tohoto odstavce bude zapsán výtah (abstrakt) práce v českém (slovenském) jazyce.

## **Abstract**

Do tohoto odstavce bude zapsán výtah (abstrakt) práce v anglickém jazyce.

## **Klíčová slova**

Sem budou zapsána jednotlivá klíčová slova v českém (slovenském) jazyce, oddělená čárkami.

## **Keywords**

Sem budou zapsána jednotlivá klíčová slova v anglickém jazyce, oddělená čárkami.

## **Citace**

JOSEFÍK, Martin. *Distribučný repositář digitálních forenzních dat*. Brno, 2018. Diplomová práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Vedoucí práce RNDr. Marek Rychlý, Ph.D.

# Distribuovaný repositář digitálních forenzních dat

## Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně pod vedením pana X... Další informace mi poskytli... Uvedl jsem všechny literární prameny a publikace, ze kterých jsem čerpal.

.....

Martin Josefík

13. října 2017

## Poděkování

V této sekci je možno uvést poděkování vedoucímu práce a těm, kteří poskytli odbornou pomoc (externí zadavatel, konzultant, apod.).

# Obsah

<b>1</b>	<b>Úvod</b>	<b>2</b>
<b>2</b>	<b>Digitální forenzní data</b>	<b>3</b>
2.1	Formáty digitálních forenzních dat . . . . .	3
2.2	Způsob uložení . . . . .	3
2.3	Existující systémy . . . . .	3
<b>3</b>	<b>Úložiště pro rozsáhlá strukturovaná i nestrukturovaná data</b>	<b>4</b>
3.1	Big data . . . . .	4
3.2	Distribuované databáze . . . . .	5
3.3	NoSQL, disky, úložiště . . . . .	8
<b>4</b>	<b>Návrh distribuovaného úložiště</b>	<b>9</b>
4.1	Přístup k datům . . . . .	9
4.1.1	Sekvenční, náhodný . . . . .	9
4.1.2	Dotazování . . . . .	9
4.1.3	Big data přístupy . . . . .	9
4.2	Architektura . . . . .	9
4.3	Aplikační rozhraní . . . . .	9
4.4	Technologie . . . . .	9
4.4.1	Docker . . . . .	9
4.4.2	HDFS, Hadoop, Spark . . . . .	9
4.4.3	Cassandra / MongoDB . . . . .	9
4.4.4	Zookeeper . . . . .	9
4.4.5	MQ broker . . . . .	9
<b>5</b>	<b>Implementace</b>	<b>10</b>
5.1	Rozšiřitelnost, znouvupoužitelnost . . . . .	10
<b>6</b>	<b>Testování</b>	<b>11</b>
6.1	Výkon . . . . .	11
<b>7</b>	<b>Závěr</b>	<b>12</b>
	<b>Literatura</b>	<b>13</b>

# Kapitola 1

## Úvod

## Kapitola 2

# Digitální forenzní data

2.1 Formáty digitálních forenzních dat

2.2 Způsob uložení

2.3 Existující systémy

## Kapitola 3

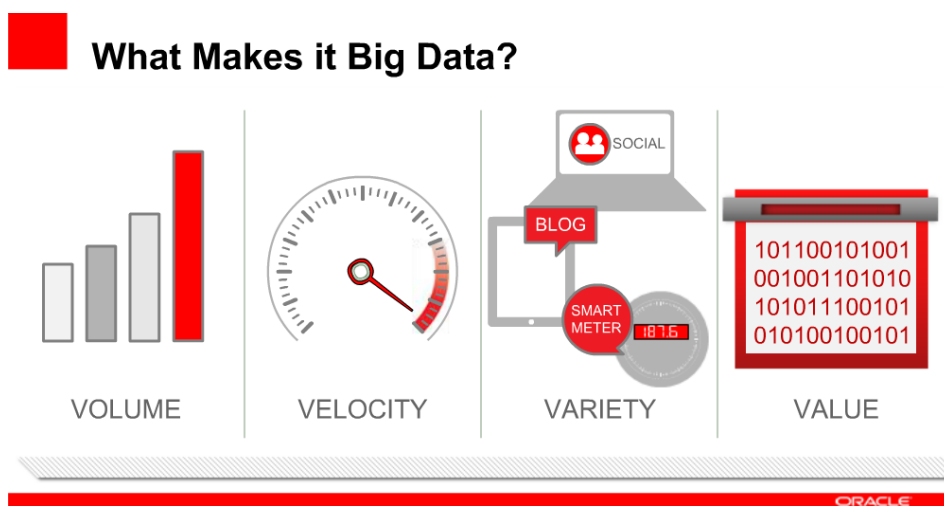
# Úložiště pro rozsáhlá strukturovaná i nestrukturovaná data

V této kapitole budou vysvětleny termíny Big data, distribuované databáze a NoSQL databáze, včetně jejich vlastností, výhod a nevýhod.

### 3.1 Big data

Definicí pro frázi Big data existuje několik. Jedná se o termín použitý na soubory dat, které jsou příliš komplexní z hlediska velikosti a různorodosti, a které je nemožné zpracovávat běžně používanými přístupy a softwarovými nástroji v rozumném čase.

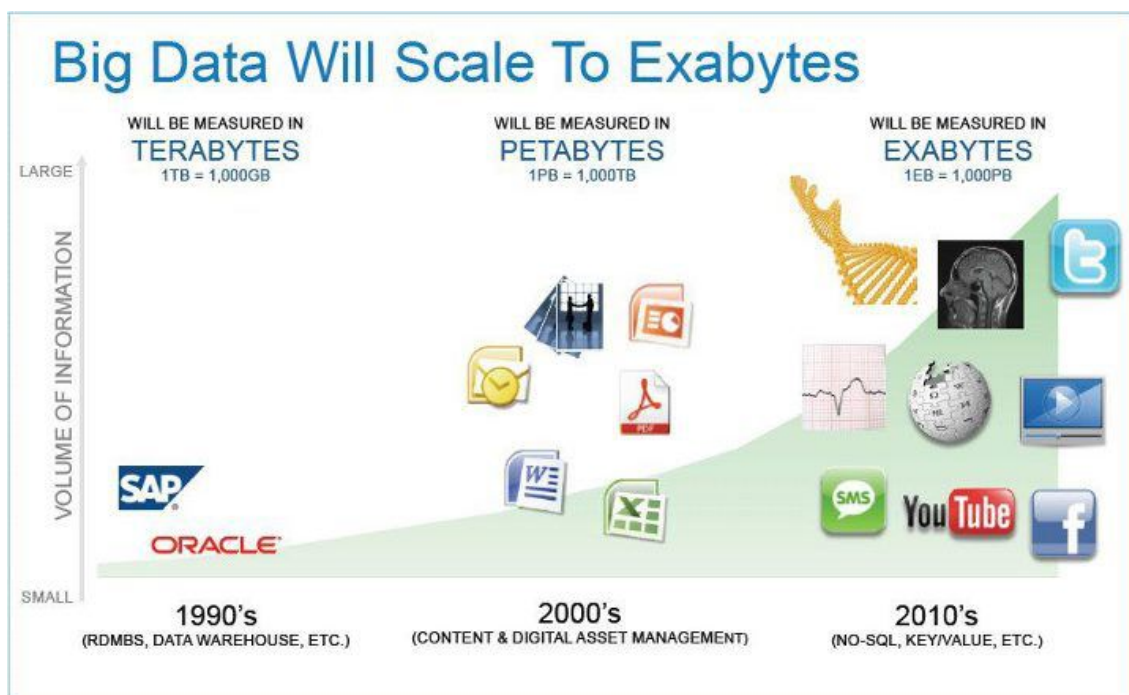
Objem takových dat rychle roste. Vyskytují se v mnoha odvětvích, například sběr informací o počasí, sociální sítě, energetické a telekomunikační společnosti, ekonomie a finančnictví, či data z kamer, měření z různých senzorů apod. Z toho plyne, že se jedná o data různorodých typů, mohou být strukturovaná i nestrukturovaná. Proto je potřeba existence různých technologií pro jejich uložení, zpracování i zobrazení.



Obrázek 3.1: Definice Big data podle Oracle. [2]

Big data je často definováno jako 4V z anglických slov Volume, Velocity, Variety a Value. [1]

- Volume – značí množství nebo velikost dat. Big data vyžaduje zpracování vysokých objemů dat neznámých hodnot, například síťový provoz, data sesbírána ze senzorů apod.
- Velocity – vyjadřuje rychlost z hlediska vzniku dat a potřeby jejich analýzy, některá vyžadují zpracování v reálném čase. Nejdůležitější data se zapisují přímo do paměti, a ne na disk, z důvodu co nejrychlejšího zpracování.
- Variety – znamená různorodost typů. Jedná se především o nestrukturovaná data, například text, audio, video, data o geografické poloze a další. Jsou na ně kladeny velmi podobné požadavky jako na data strukturovaná – sumarizace, monitorování, důvěrnost. [1]
- Value – data mají vlastní hodnotu, která musí být analyzována a zjištěna. Nejedná se o jednoduchý proces, je stále potřeba nových metod a technik zpracování.



Obrázek 3.2: S novými technologiemi se masivně zvyšuje růst dat a přibývají nové typy. [3]

Tato práce se zabývá Big daty hlavně typu – PCAP soubory, logy ze síťových zařízení a komunikací. Možnosti uložení Big data budou popsány v následujících podkapitolách.

## 3.2 Distribuované databáze

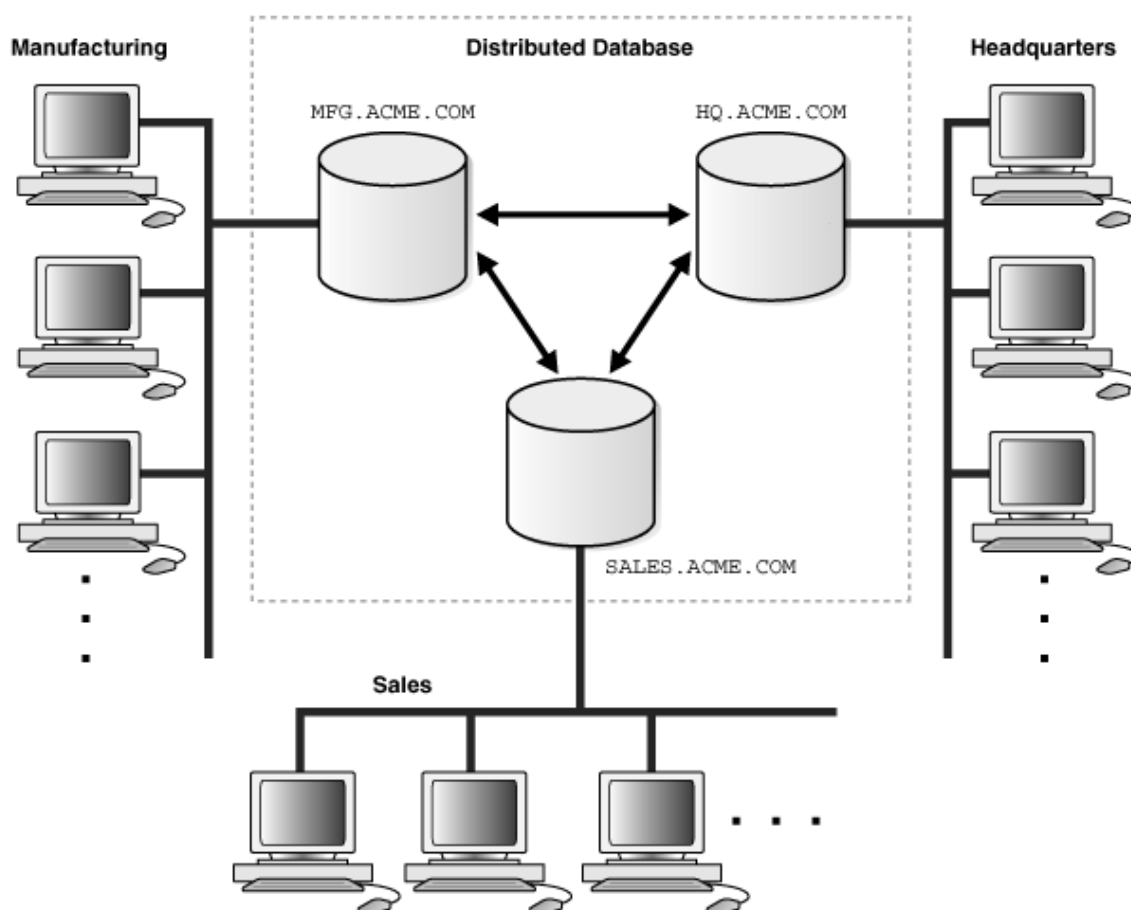
Distribuovaná databáze se skládá z většího počtu samostatných databází, které mohou být geograficky rozmístěny na jiných pozicích. Jednotlivé uzly spolu komunikují přes počítačovou síť. Každý uzel je sám o sobě databázový systém. DSŘBD neboli systém řízení



distribuované báze dat (anglicky Distributed Database Database Management System) zajišťuje, že se distribuovaná databáze uživatelům jeví jako jedna jediná databáze. Data jsou fyzicky uložena na různých pozicích. Mohou být spravována rozdílnými SŘBD nezávisle na ostatních pozicích. [5]

Systém řízení distribuované báze dat je centralizovaný systém s těmito vlastnostmi [5]:

- Umí vytvářet, získávat, upravovat a mazat distribuované databáze. Zajišťuje důvěrnost a integritu databází.
- Periodicky synchronizuje databázi a poskytuje mechanismy přístupu tak, aby se databáze uživatelům jevila transparentní.
- Zajišťuje, že změna dat v kterémkoliv uzlu se promítne i v ostatních uzlech.
- Je využíván v aplikacích, kde se předpokládá zpracování velkých objemů dat, ke kterým přistupuje současně mnoho uživatelů.



Obrázek 3.3: Schéma distribuované databáze a současný přístup více zařízení k ní. [4]

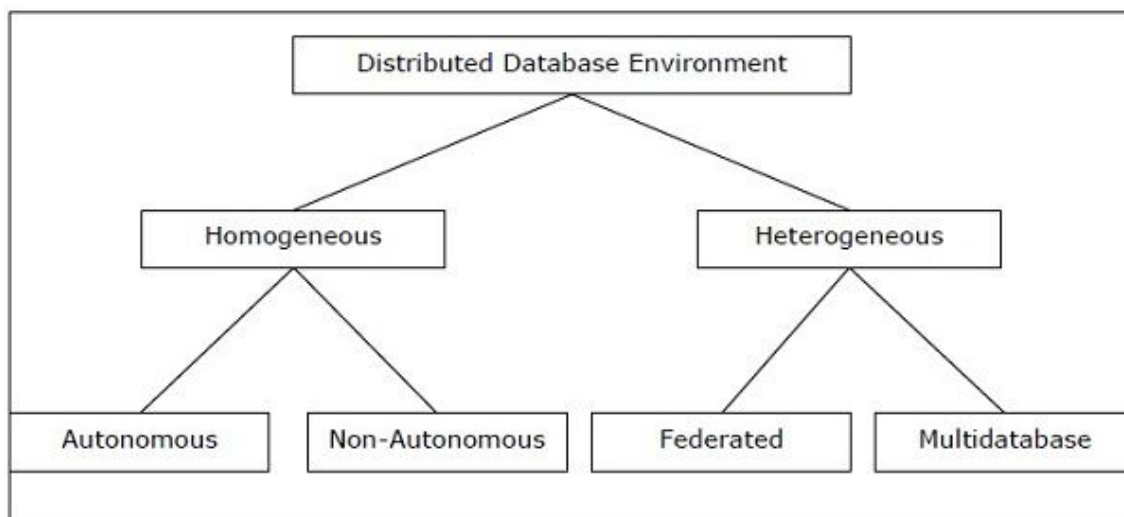
#### Výhody

- Rozšiřitelnost – pokud je potřeba databázový systém rozšířit do nových míst nebo přidat další uzly, stačí přidat nový(é) počítač(e) a lokální data v nové pozici, a na konec je připojit k distribuovanému systému, bez jakéhokoli přerušení funkcionality. Podobný postup je při odebrání uzlu.

- Spolehlivost – když nějaký z připojených uzlů selže, nepřestane distribuovaná databáze fungovat, sníží se maximálně výkon.
- Ochrana (záloha) dat – při zničení jednoho uzlu a smazání dat z něj, mohou být stejná data zálohována i na jiných uzlech.
- Výkonnost – pokud jsou data efektivně distribuována, může být uživatelův požadavek uspokojen rychleji. Transakce mohou být také distribuované a provedeny rychleji.

#### Nevýhody

- Integrita dat – data musí být průběžně synchronizována na více uzlech, aby na stejné dotazy nebyly z různých uzlů vráceny rozdílné odpovědi.
- Komunikační režie – i zdánlivě jednoduchá operace může vyžadovat spoustu zbytečné komunikace.
- Cena – DSŘDB vyžaduje drahý a složitý software ke koordinaci uzlu a zajištění transparentnosti. [5]
- Mezi další patří – složitost, zabezpečení, řízení souběžného přístupu k datům.



Obrázek 3.4: Distribuované databáze můžeme rozdělit na homogenní a heterogenní, a tyto ještě dále dělit. [5]

Homogenní – všechny uzly používají identické SŘBD a operační systémy. Uzly mají informace o ostatních uzlech a spolupracují při zpracování uživatelských požadavků. Homogenní distribuovaná databáze se navenek jeví uživateli jako jeden systém. Je jednodušší jej navrhnout a spravovat.

Heterogenní – uzly mohou mít rozdílné operační systémy a SŘBD, které nejsou kompatibilní. Mohou také využívat rozdílná schémata (relační, objektově orientované, hierarchické, ...). Rozdílnost schématu je hlavním problémem při zpracování dotazu a transakcí. Kvůli tomu je také složité dotazování. [6]

Architekturami distribuovaných databází jsou centrální architektura, klient-server, peer-to-peer, multi-databázová architektura.

### 3.3 NoSQL, disky, úložiště

## Kapitola 4

# Návrh distribuovaného úložiště

### 4.1 Přístup k datům

#### 4.1.1 Sekvenční, náhodný

#### 4.1.2 Dotazování

#### 4.1.3 Big data přístupy

### 4.2 Architektura

### 4.3 Aplikační rozhraní

### 4.4 Technologie

#### 4.4.1 Docker

#### 4.4.2 HDFS, Hadoop, Spark

#### 4.4.3 Cassandra / MongoDB

#### 4.4.4 Zookeeper

#### 4.4.5 MQ broker

## Kapitola 5

# Implementace

### 5.1 Rozšiřitelnost, znovupoužitelnost

## Kapitola 6

# Testování

### 6.1 Výkon

## Kapitola 7

## Závěr

# Literatura

- [1] Heller, P.; Piziak, D.; Stackowiak, R.; aj.: *An Enterprise Architect's Guide to Big Data*. [Online; navštíveno 26.09.2017].  
URL <http://www.oracle.com/technetwork/topics/entarch/articles/oea-big-data-guide-1522052.pdf>
- [2] Louwers, J.: *Big Data is sometimes Fast Data*. [Online; navštíveno 27.09.2017].  
URL <http://johanlouwers.blogspot.cz/2013/01/big-data-is-sometimes-fast-data.html>
- [3] Nambiar, R.: *What is Big Data ?*. [Online; navštíveno 27.09.2017].  
URL <http://rrnamb.blogspot.cz/2012/09/what-is-big-data.html>
- [4] Oracle Help Center: *Distributed Database Architecture*. [Online; navštíveno 29.09.2017].  
URL [https://docs.oracle.com/cd/B28359\\_01/server.111/b28310/ds\\_concepts001.htm](https://docs.oracle.com/cd/B28359_01/server.111/b28310/ds_concepts001.htm)
- [5] Tutorials Point (I) Pvt. Ltd.: *Distributed DBMS Tutorial*. [Online; navštíveno 29.09.2017].  
URL [https://www.tutorialspoint.com/distributed\\_dbms/](https://www.tutorialspoint.com/distributed_dbms/)
- [6] Wikipedia: *Distributed database*. [Online; navštíveno 02.10.2017].  
URL [https://en.wikipedia.org/wiki/Distributed\\_database](https://en.wikipedia.org/wiki/Distributed_database)