

DĚLAT
DOBRÝ SOFTWARE
NÁS BAVÍ

PROFINIT

B0M33BDT – 2. přednáška

Marek Sušický

3. 10. 2018

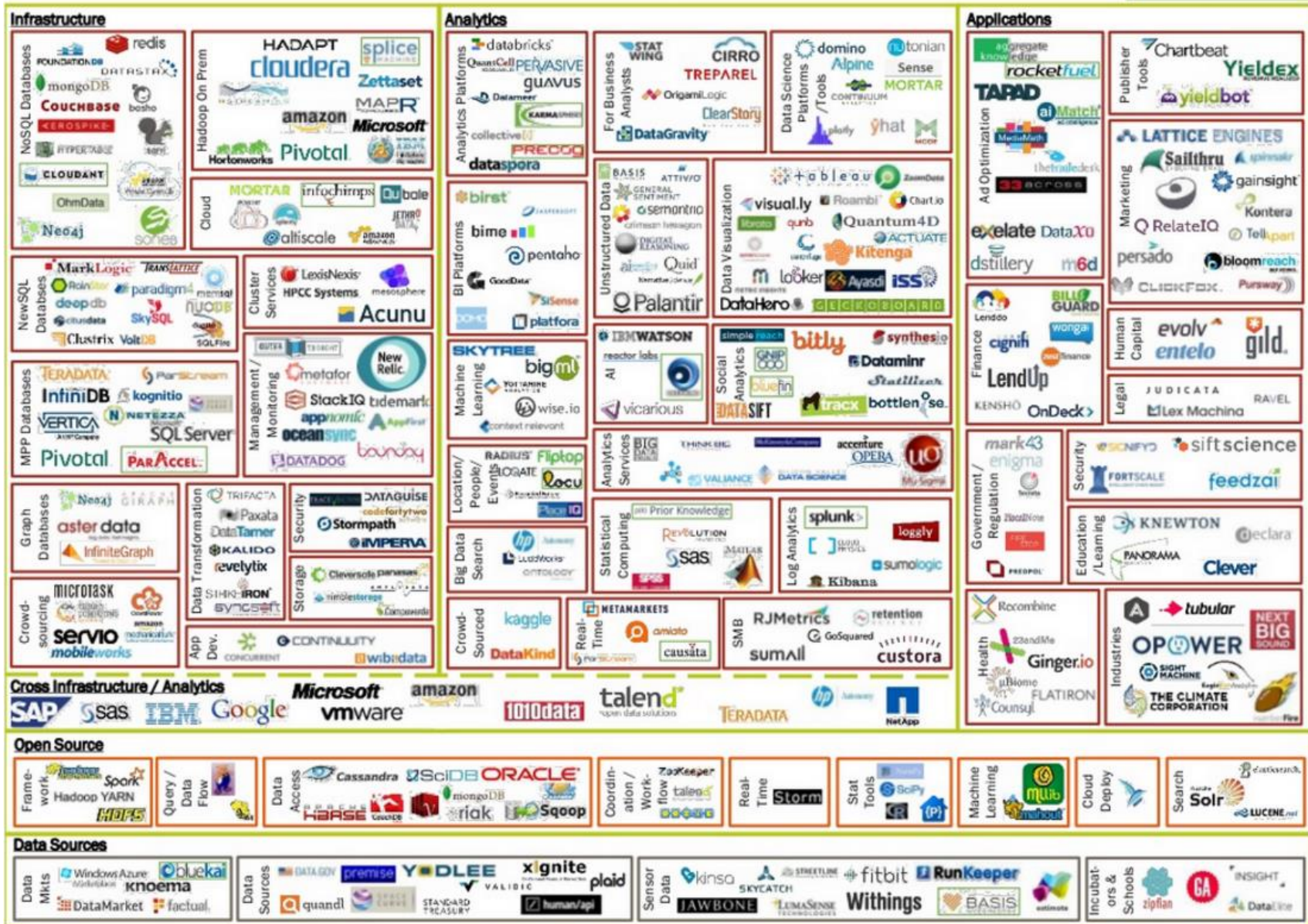
Osnova

- › Big data a Hadoop
- › Na jakém hardware + sizing
- › Jak vypadá cluster - architektura
- › HDFS
- › Distribuce
- › Komponenty
- › YARN, správa zdrojů

Big data neznamená Hadoop

BIG DATA LANDSCAPE, VERSION 3.0

Exited: Acquisition or IPO



PROFIT

Apache Hadoop

- › Wikipedia:
 - Apache Hadoop (pronunciation: /hə'du:p/) is an open-source software framework for **distributed storage** and **distributed processing** of **very large data sets** on computer clusters ***built from commodity hardware***. All the modules in Hadoop are designed with a **fundamental assumption that hardware failures are common** and should be automatically handled by the framework.

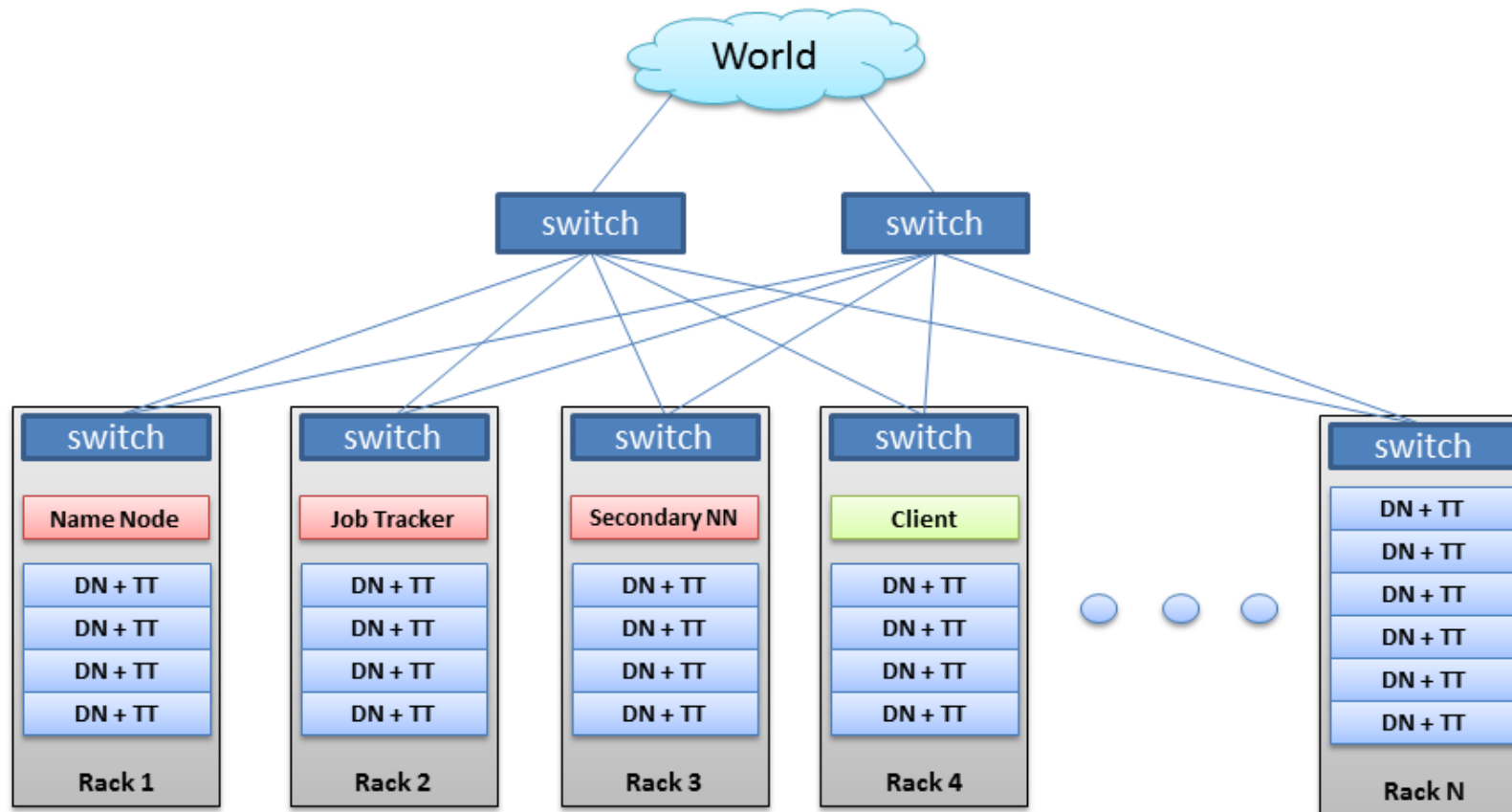
- › Commodity hardware
 - stroje za statisíce CZK (ale ne desítky mil.)
 - 2-4 CPU, každé CPU 10-16 jader
 - 256-512 GB RAM, min. 128GB
 - 10-20 2-4TB HDD
 - rozhodně ne to, co jako server prodává Alza 😊

Jak vypadá levný HW

The screenshot shows the Cloudera Manager web interface in Internet Explorer. The browser address bar shows the URL `http://10.171.64.11:7180/cm/hardware/hosts`. The Cloudera Manager navigation bar includes links for Clusters, Hosts, Diagnostics, Charts, Backup, and Administration. The main content area is titled "All Hosts" and contains several action buttons: Configuration, Add New Hosts to Cluster, Re-run Upgrade Wizard, and Inspect All Hosts. On the left, a "Filters" sidebar shows "STATUS" with a "Good Health" indicator and a count of 5. Below this are expandable sections for CLUSTERS, CORES, COMMISSION STATE, LAST HEARTBEAT, LOAD (1 MINUTE), LOAD (5 MINUTES), LOAD (15 MINUTES), MAINTENANCE MODE, RACK, SERVICES, and HEALTH TESTS. The main table displays a list of hosts with columns for Status, Name, IP, Roles, Last Heartbeat, Load Average, Disk Usage, Physical Memory, and Swap Space. The table contains five rows of host data, all with a "Good Health" status.

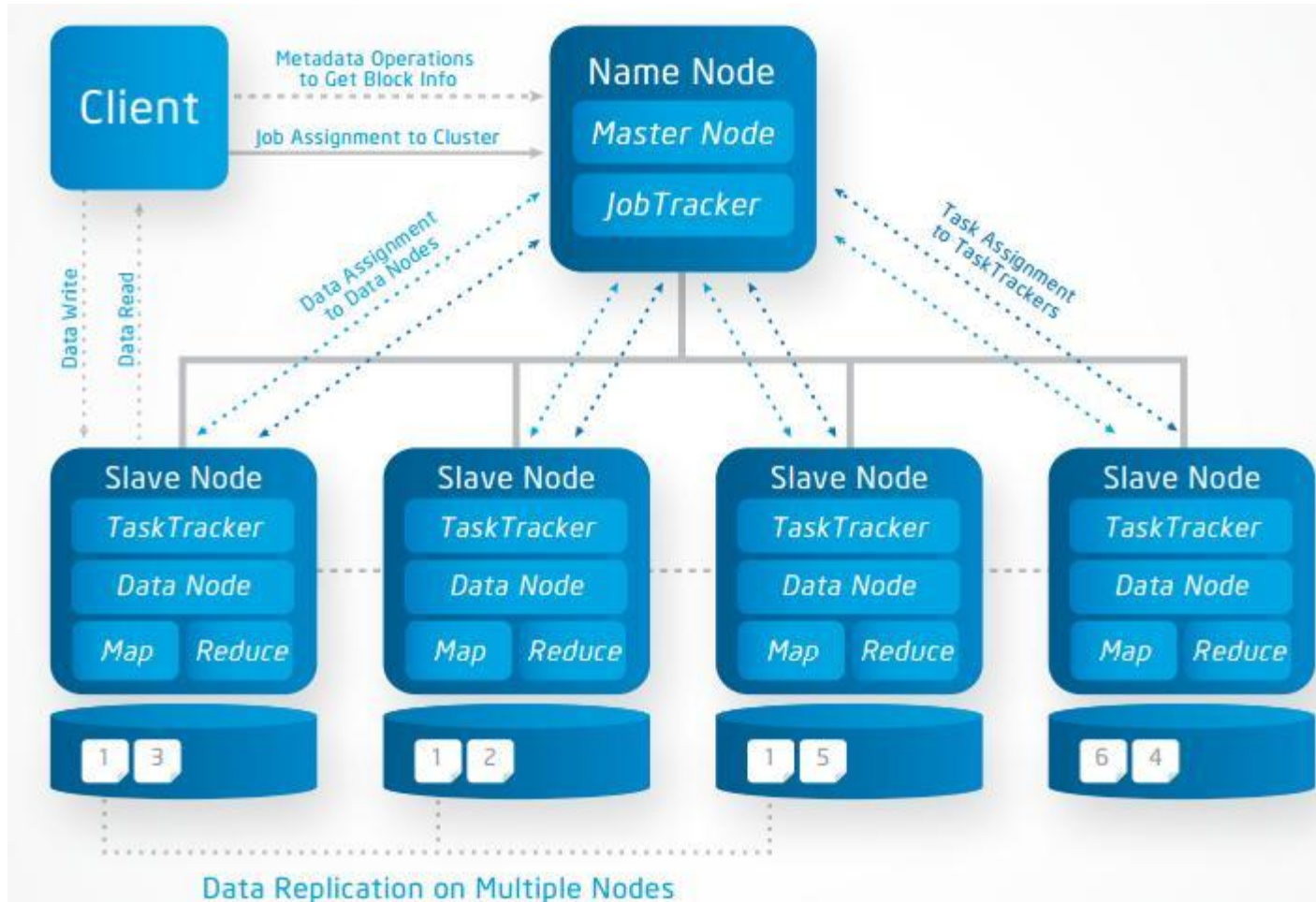
Status	Name	IP	Roles	Last Heartbeat	Load Average	Disk Usage	Physical Memory	Swap Space
Good Health	dhbbdn05	6.93.1.1	5 Role(s)	7.71s ago	0.00 0.02 0.05	13.6 GiB / 22 TiB	8.6 GiB / 251.8 GiB	0 B / 32 GiB
Good Health	dhbbdn06	6.93.4.1	5 Role(s)	11.64s ago	0.00 0.01 0.05	13.7 GiB / 22 TiB	8.8 GiB / 251.8 GiB	0 B / 32 GiB
Good Health	dhbbdn07	6.93.4.1	4 Role(s)	9.77s ago	0.05 0.04 0.05	13.5 GiB / 22 TiB	7.8 GiB / 251.8 GiB	0 B / 32 GiB
Good Health	dhbfe06	6.93.4.1	16 Role(s)	9.09s ago	0.10 0.22 0.38	29.9 GiB / 141.5 GiB	16.3 GiB / 251.8 GiB	0 B / 32 GiB
Good Health	dhbfe07	6.93.4.1	10 Role(s)	7.42s ago	0.12 0.06 0.05	24 GiB / 353.9 GiB	12.4 GiB / 125.8 GiB	0 B / 32 GiB

Hadoop Cluster



BRAD HEDLUND .com

Hadoop – architektura II



Sizing

- › Jak postavit Hadoop
 - Jak si ho objednat
 - HDD parametry
 - Přenosová rychlost
 - RAID
 - 0, 1, 1+0, 5, 6, (2,3,4,7)
 - Síťová rychlost
 - SAN/NAS
 - Paměť
 - CPU jádra
 - Obecná doporučení

Sizing

- › Kalkulace HW požadavků
 - Počet nodů
 - Počet disků
 - HW parametry (CPU, RAM, RAID...) typů nodů

Rychlosti čtení dat

- › RAM
 - DDR4 cca 15 GB/s
- › Síť 10 Gbit
 - 1.25 GB/s
- › SSD disk
 - 200-700 MB/s
 - existují „Enterprise level“, které vydrží (garance 5 let)
 - malé kapacity (max 1TB) a hodně drahé
- › HDD 7.2k
 - latence cca 4ms
 - sekvenční čtení 50-100 MB/s
 - velké kapacity (běžně 4TB-8TB) a relativně levné
 - → Hadoop typicky pracuje s úložištěm

Rychlosti čtení dat – HDD

- › **Sekvenční čtení** – cca 100 MB/s za jedním diskem
- › Random access
 - velikost bloku ext4 bývá 4kB
 - latence, než disk najde blok cca 4ms
 - max. rychlost čistě náhodného čtení $1/0.004 * 4096 = 1 \text{ MB/s}$

Omezující faktory

- › Příklad: 10 nodů, každý node 12 * 2 TB HDD
 - Rychlost čtení v rámci nodu: $12 * 100 \text{ MB/s} = 1.2 \text{ GB/s}$
 - Rychlost čtení v rámci clusteru: 12 GB/s

- › Omezení:
 - rychlost RAM – 10x větší
 - CPU – čtení nezatěžuje
 - síť – na hraně pro jeden node!
 - sběrnice – pozor na počet disků, musí zvládnout



Sizing

- › Ukázkový příklad k zamyšlení
 - 15GB/5 min
 - Historie 30 dní
 - SLA – 5 sekund pro 85% dotazů
 - 10s pro 100% dotazů
 - Při nesplnění vysoké pokuty

Principy

- › Ukládání velkého množství dat
 - mnoho serverů = nodů [desítky až tisíce]
 - každý node mnoho disků [10-20]
- › Zamezení ztráty dat – výpadek nodu
 - replikace (typicky tři kopie každého souboru)
 - 2 repliky ve stejném racku, třetí replika mimo
- › Rychlost čtení
 - data jsou rozložena v celém clusteru – 1 soubor nemusí být celý v jednom nodu!
 - data jsou replikována – lze paralelně číst na několika nodech bez nutnosti přenosu dat přes síť
 - velké soubory – **výhody sekvenčního čtení**
- › Distribuce výpočtů
 - mnoho nodů, přiřazování výpočetního výkonu

Principy Hadoop – syntéza

- › Maximálně využívat sekvenční čtení
- › Pracovat s velkými soubory, které se čtou sekvenčně
 - ušetří se čas na synchronizaci/orchestraci v rámci distribuovaného systému
- › Čím méně dat se načte, tím rychleji se načtou – využití **kompres**e
- › CPU se při čtení fláká, paměť je násobně rychlejší, tj. typicky dekomprese bude rychlejší než IO operace
- › Maximum výpočtů provádět na nodu, kde probíhá čtení dat, přes síť posílat jen co nejvíce agregovaná data
 - síťové rozhraní nebude mít dostatečnou kapacitu
- › **Zásadní omezení – I/O operace**
 - **Tip:** při odhadu, jak dlouho co bude trvat stačí prakticky vždy počítat jen s IO a počtem paralelních čtení/uživatelů (ale neplatí samozřejmě pro ML aplikace)



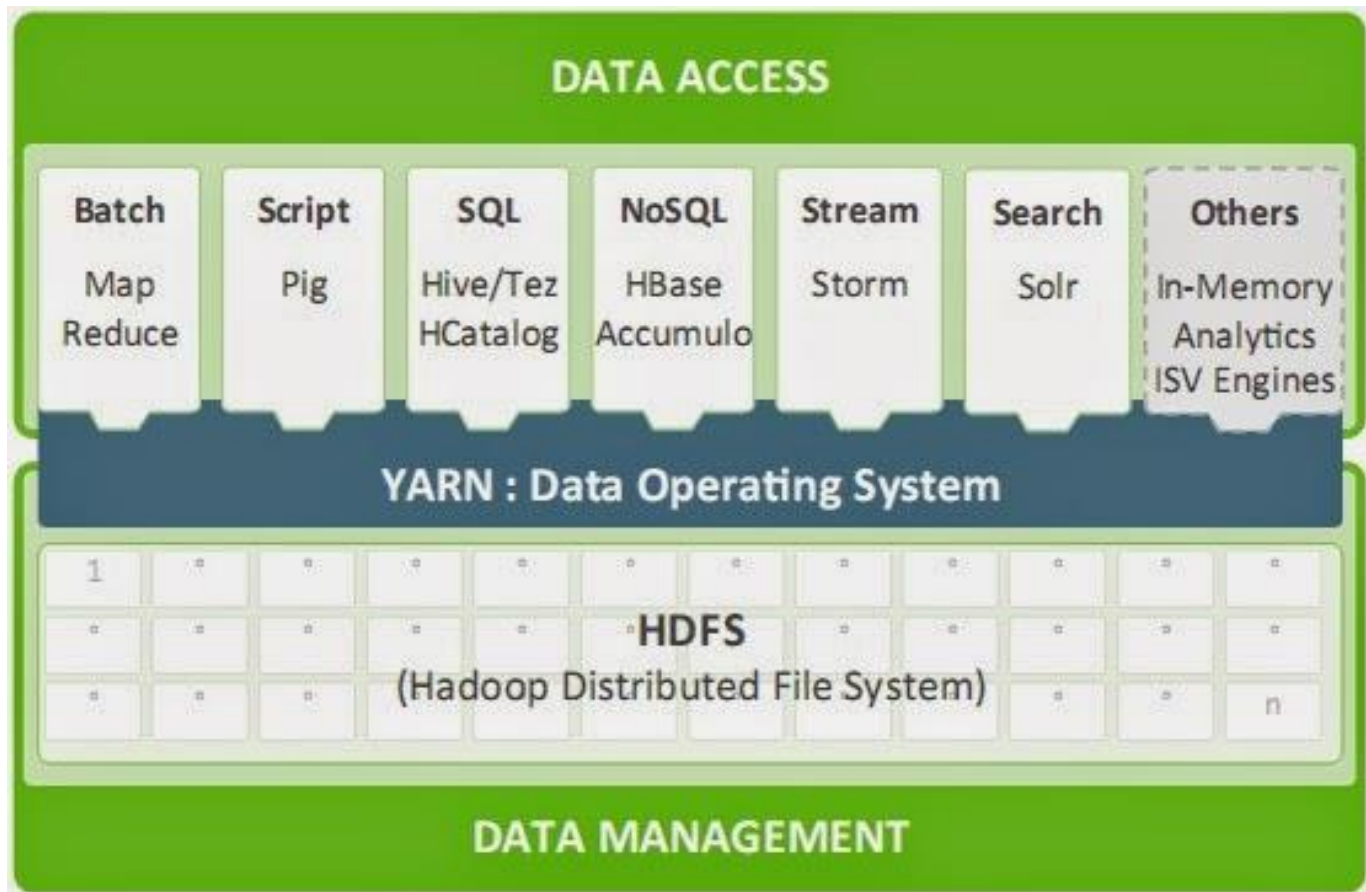
HDFS

HDFS

› HDFS

- NameNode, DataNode
- Replikace
- Operace na souborovém systému
- Bloky, velikost bloku

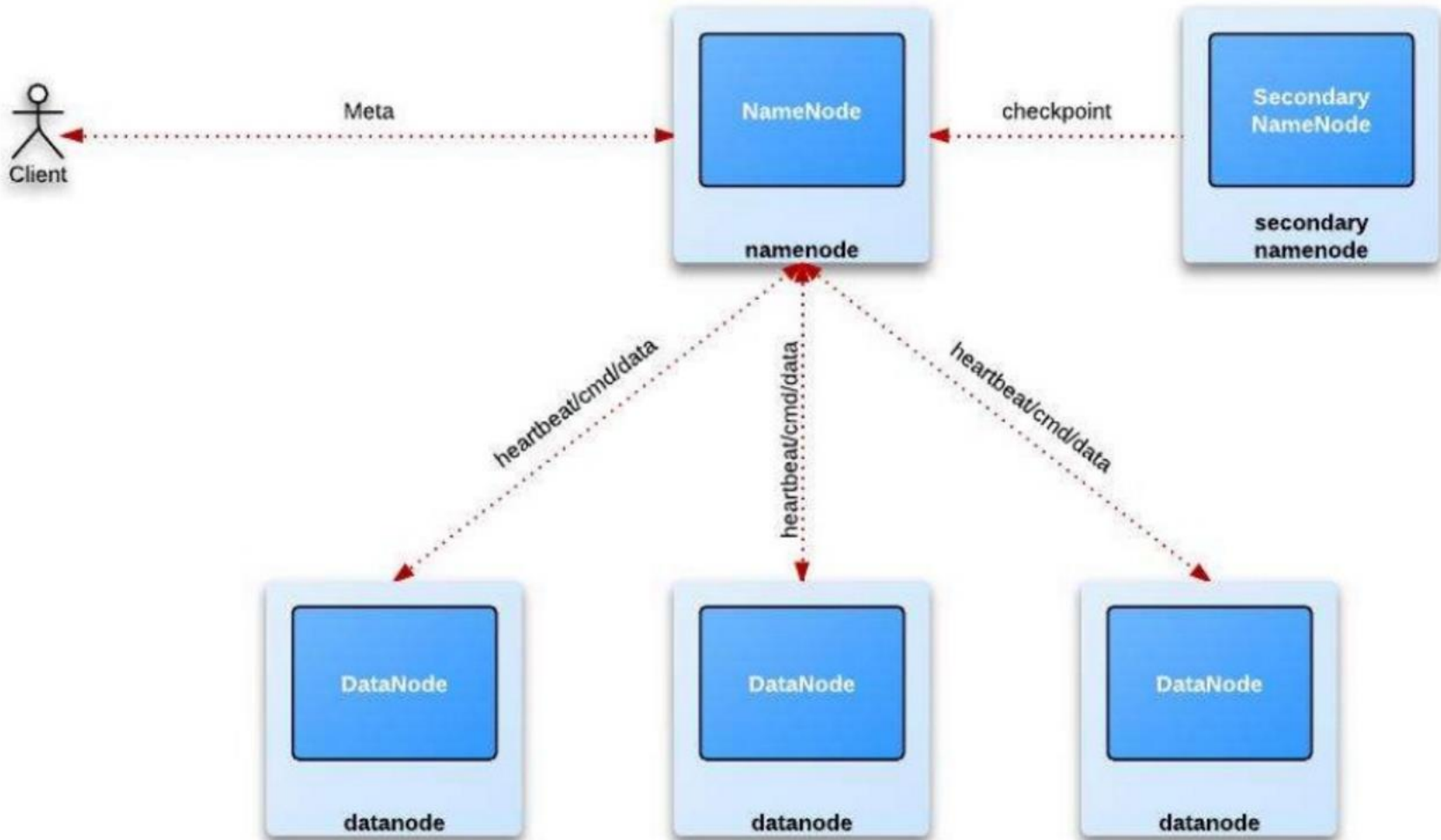
HDFS- architektura



HDFS

- › Hadoop Distributed Filesystem
- › Dobrý pro
 - Velké soubory
 - Streamovaný přístup
- › Špatný pro
 - Spoustu malých souborů
 - Náhodný přístup
 - Nízkolatenční přístup
- › Master-slave design
 - Master – NameNode
 - Slave – DataNode
 - SecondaryNameNode

HDFS



HDFS

- › HDFS soubory jsou rozděleny do bloků
 - Default 64MB/128MB, ale lze změnit
 - Dobré pro velké soubory
 - Děsné pro malé...
- › Replikace
 - Jeden blok může být v nejméně xxx DataNodes
 - Fault tolerant
 - Default hodnota 3

NameNode

- › Metadata filesystemu
 - Kde jsou data
 - V paměti
 - 1GB pro každý milion bloků

- › Ve spojení s
 - Klienty
 - DataNodey
 - SecondaryNameNodey
 - Checkpointing
 - Editlogs a fsimage

DataNode

- › Ukládá Databloky
- › Získává bloky od klientů
- › Získává bloky od ostatních DataNodů
 - Replikace
- › Dostává příkaz delete od NameNode

HDFS filesystem

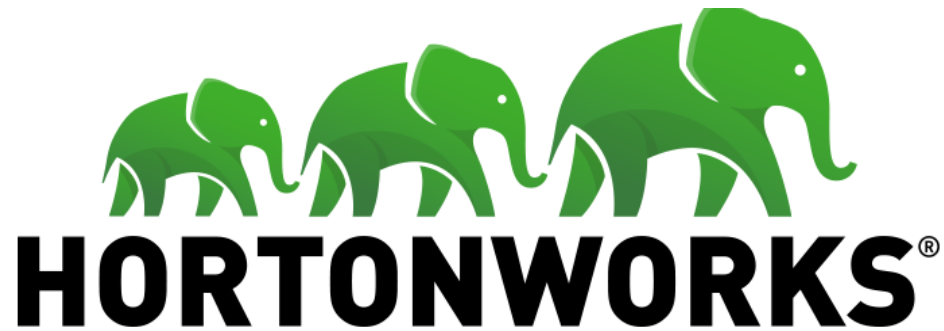
- › put
- › get
- › copyFromLocal
- › ls
- › Rights
 - Chmod
 - Chown
 - Chgrp
- › Další zde
 - <https://hadoop.apache.org/docs/r2.7.1/hadoop-project-dist/hadoop-hdfs/HDFSCommands.html>
- › Vyzkoušíme na cvičení



Distribuce

Distribuce

cloudera



Hotové distribuce

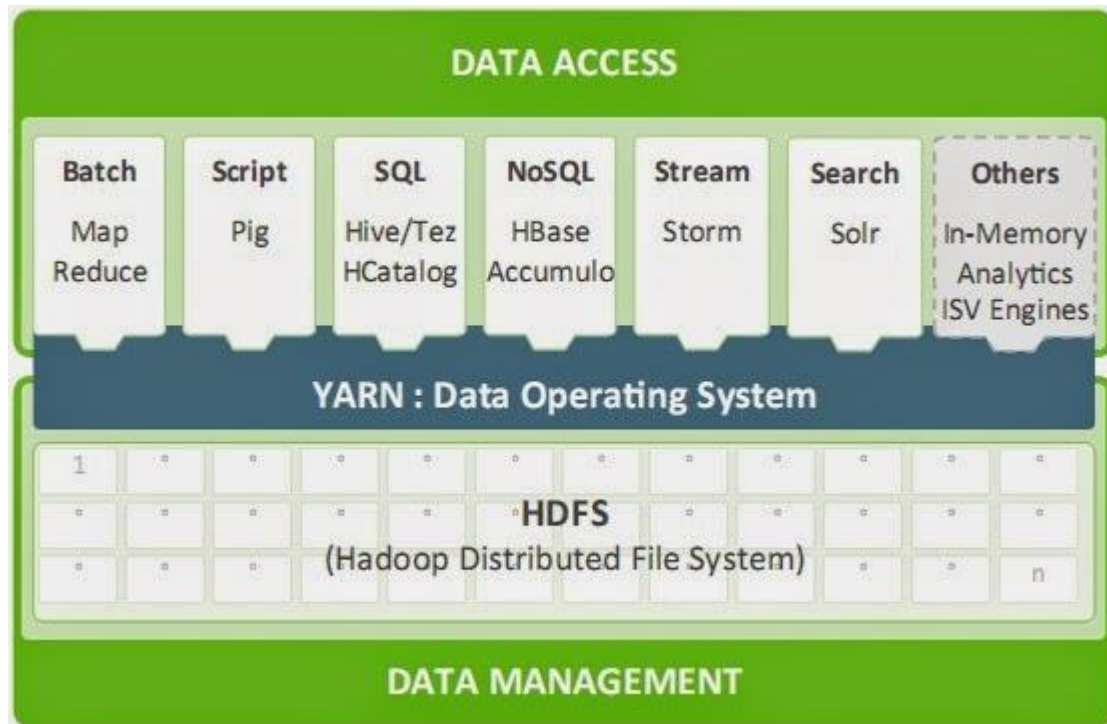
- › Řeší peklo závislostí – jeden update může vyvolat řetězovou vlnu
- › Nabízejí komerční podporu
- › Rychlejší reakce ?
- › Co je zadarmo, je špatné ?!
- › Proč znovu vymýšlet kolo

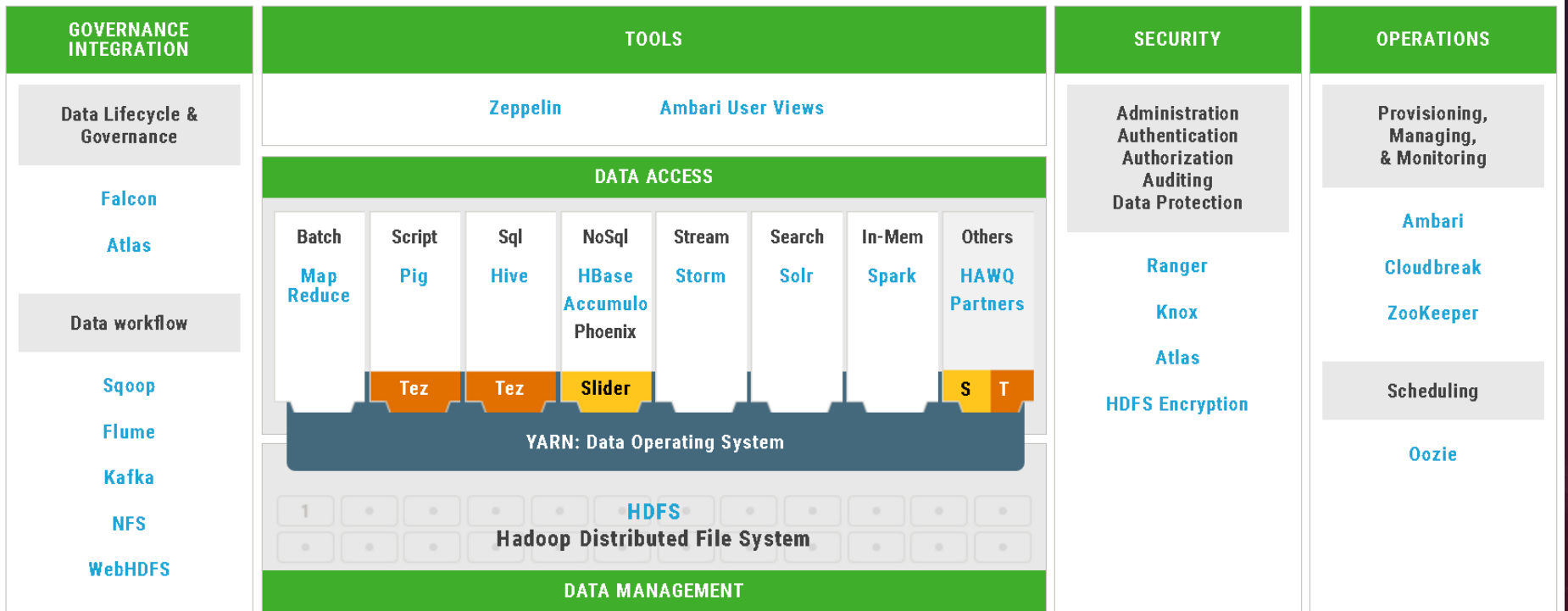
- › „Musí se k tomu dospět“

Komponenty

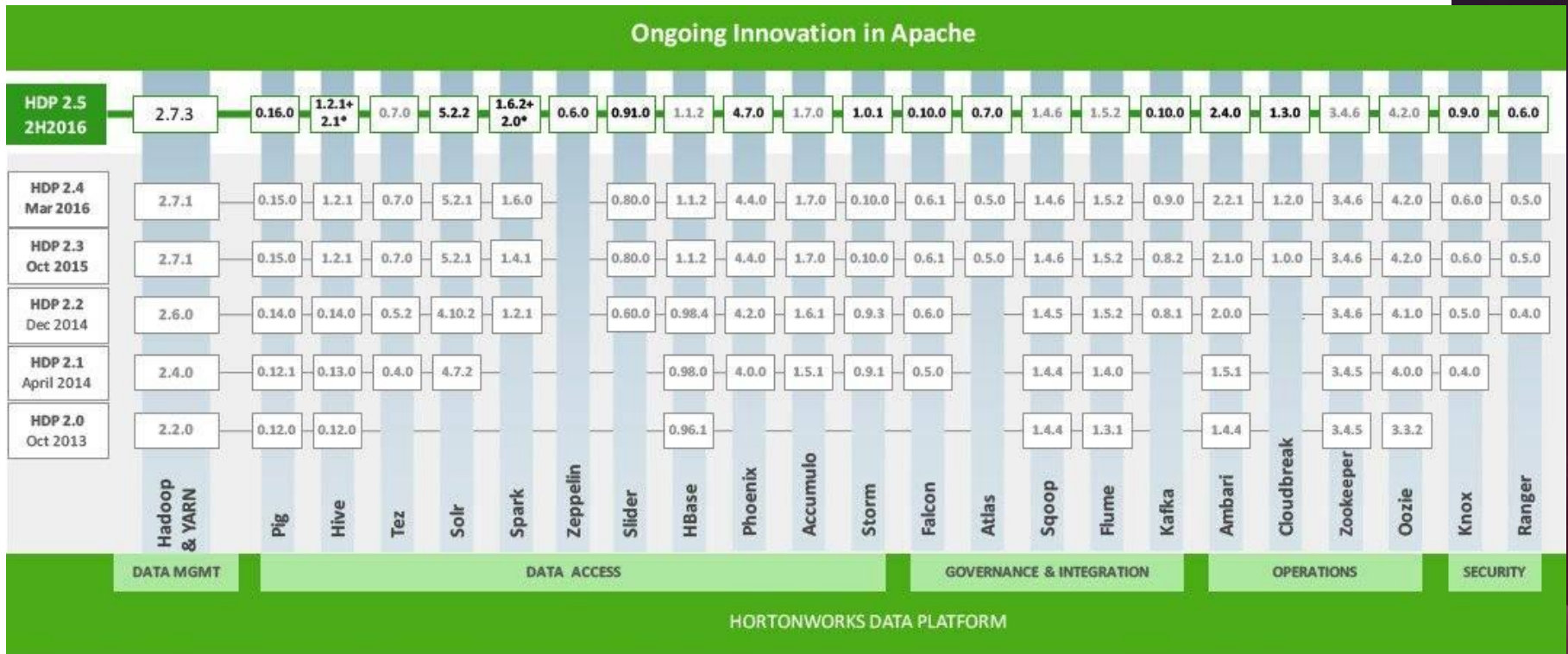
The background of the slide is a dark gray color, overlaid with a complex pattern of numerous overlapping, semi-transparent polygons. These polygons vary in size and orientation, creating a layered, crystalline effect. The overall aesthetic is modern and technical.

Zvěřinec - zjednodušeně



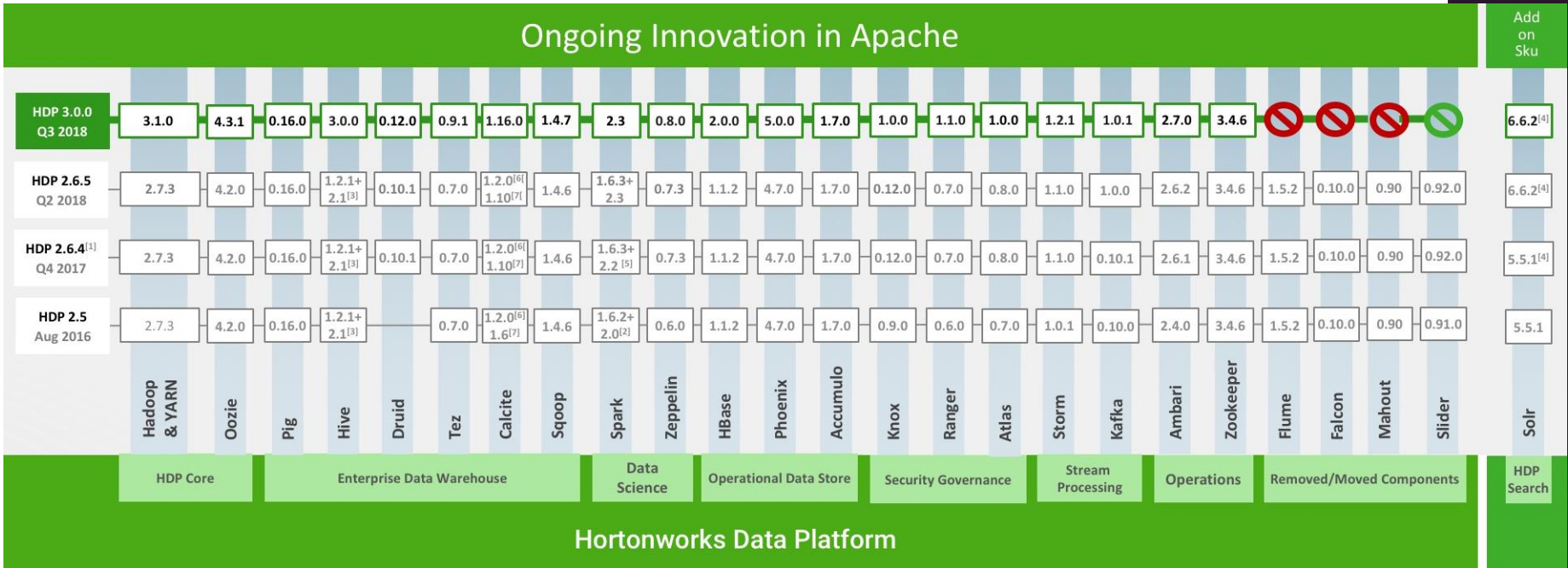


Verze mezi releases



* Spark 1.6.2+ Spark 2.0 – HDP 2.5 support installation of both Spark 1.6.2 and Spark 2.0. Spark 2.0 is Technical Preview within HDP 2.5.
Hive 1.2.1+ Hive 2.1 – Hive 2.1 is Technical Preview within HDP 2.5.

Verze mezi releasesy



[1] HDP 2.6 – Shows current Apache branches being used. Final component version subject to change based on Apache release process.

[2] Spark 1.6.3+ Spark 2.1 – HDP 2.6 supports both Spark 1.6.3 and Spark 2.1 as GA.

[3] Hive 2.1 is GA within HDP 2.6.

[4] Apache Solr is available as an add-on product HDP Search.

[5] Spark 2.2 is GA

Koloběh technologií

- › Nadšení
 - Našel jsem skvělou technologii ! Vyřeší naše problémy.
- › Realita – o několik hodin později
 - Kde je dokumentace?
- › Vystřízlivění – podle povahy o několik hodin, až jednotek dní později
 - Ono to asi vážně nefunguje.
- › Zklamání – o několik zoufalých dnů později
 - Nefungují ani příklady na webu ! Kdo probůh tohle vyvíjí ?
 - Porozhlédneme se tedy jinde...



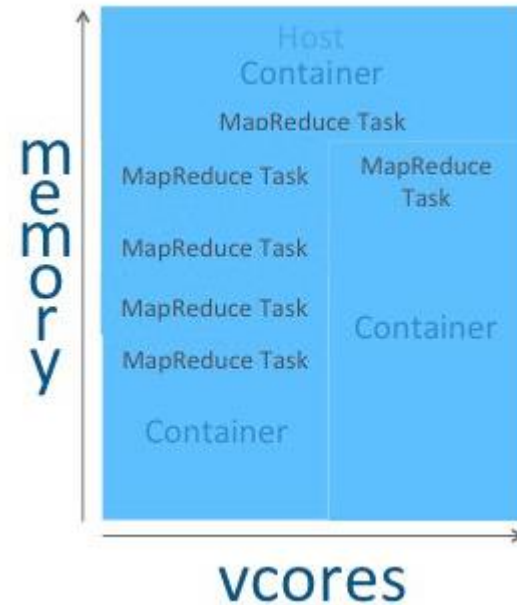
YARN

YARN

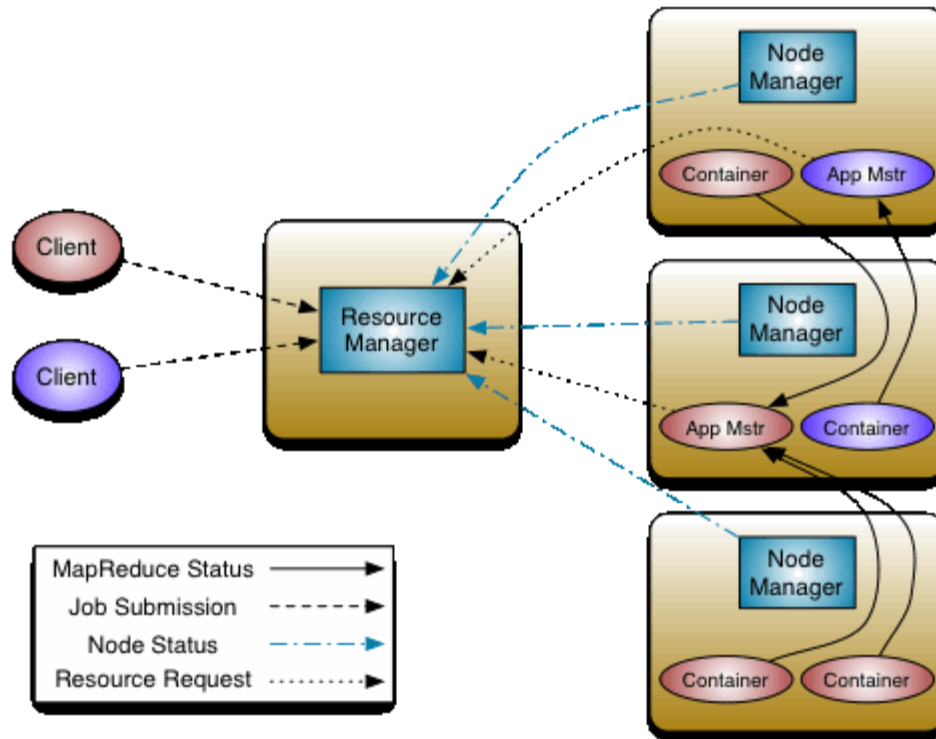
- › **Yet Another Resource Negotiator**
- › tj. plánovač a alokátor zdrojů
 - paměť
 - CPU
 - počet vláken
 - síť...
- › Většinou se využívá transparentně, uživatel o něm moc neví, záleží ale hodně na konfiguraci
- › Ne všechny aplikace YARN využívají, např. Impala má vlastní plánovač
 - každý alokátor by tak měl mít výhradní zdroje...

YARN

- › Application – klientská aplikace
- › Container – zdroje přiřazené aplikaci na konkrétním nodu
- › Resource Manager – globální správce zdrojů pro cluster
- › Node Manager – podřízený správce zdrojů na nodu



YARN



Díky za pozornost

PROFINIT

Profinit, s.r.o.
Tychonova 2, 160 00 Praha 6



Telefon
+ 420 224 316 016



Web
www.profinit.eu



LinkedIn
linkedin.com/company/profinit



Twitter
twitter.com/Profinit_EU