

DĚLAT  
DOBRÝ SOFTWARE  
NÁS BAVÍ

# PROFINIT

## B0M33BDT – 7. přednáška Architektury a bezpečnost

Marek Sušický  
Milan Kratochvíl

5. prosinec 2018

# Osnova

- › Něco ze života
- › Architektury
  - Hadoop
  - Lambda
  - Kappa
  - Zetta
- › Security a dopady do architektury

# Jak vypadá Hadoop?

- › Yahoo



# Jak vypadá Hadoop?

- › Facebook





# Jak vypadá Hadoop?

- › Google



# Jak vypadá Hadoop?



## Několik otázek

- › Jaká je rychlost světla v optickém kabelu?
- › Jaká je akceptovatelná latence pro telefonní hovor?
- › Kolik událostí za sekundu zvládnou konvenční velké databáze?
- › Kolik stojí malý clusteřík? (5x 2x 10core, 256GB RAM, 10x2TB HDD)



## Několik otázek

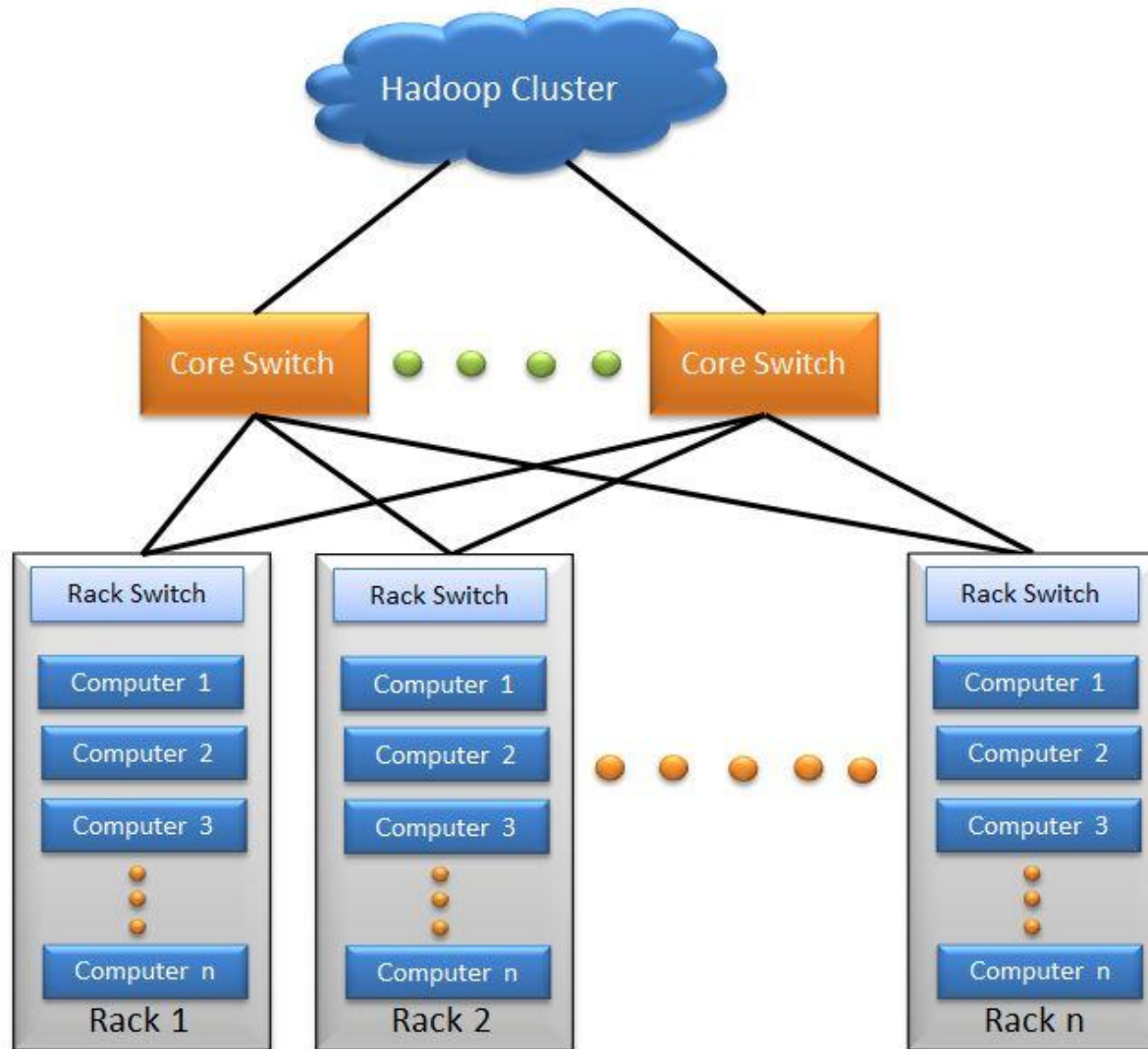
- › Jaká je rychlost světla v optickém kabelu?
  - 200 000km/s
- › Jaká je akceptovatelná latence pro telefonní hovor?
  - 50ms
- › Kolik událostí za sekundu zvládnou konvenční velké databáze?
  - Cca 10 000
- › Kolik stojí malý clusteřík? (5x 2x 10core, 256GB RAM, 10x2TB HDD)
  - Pod 5 M



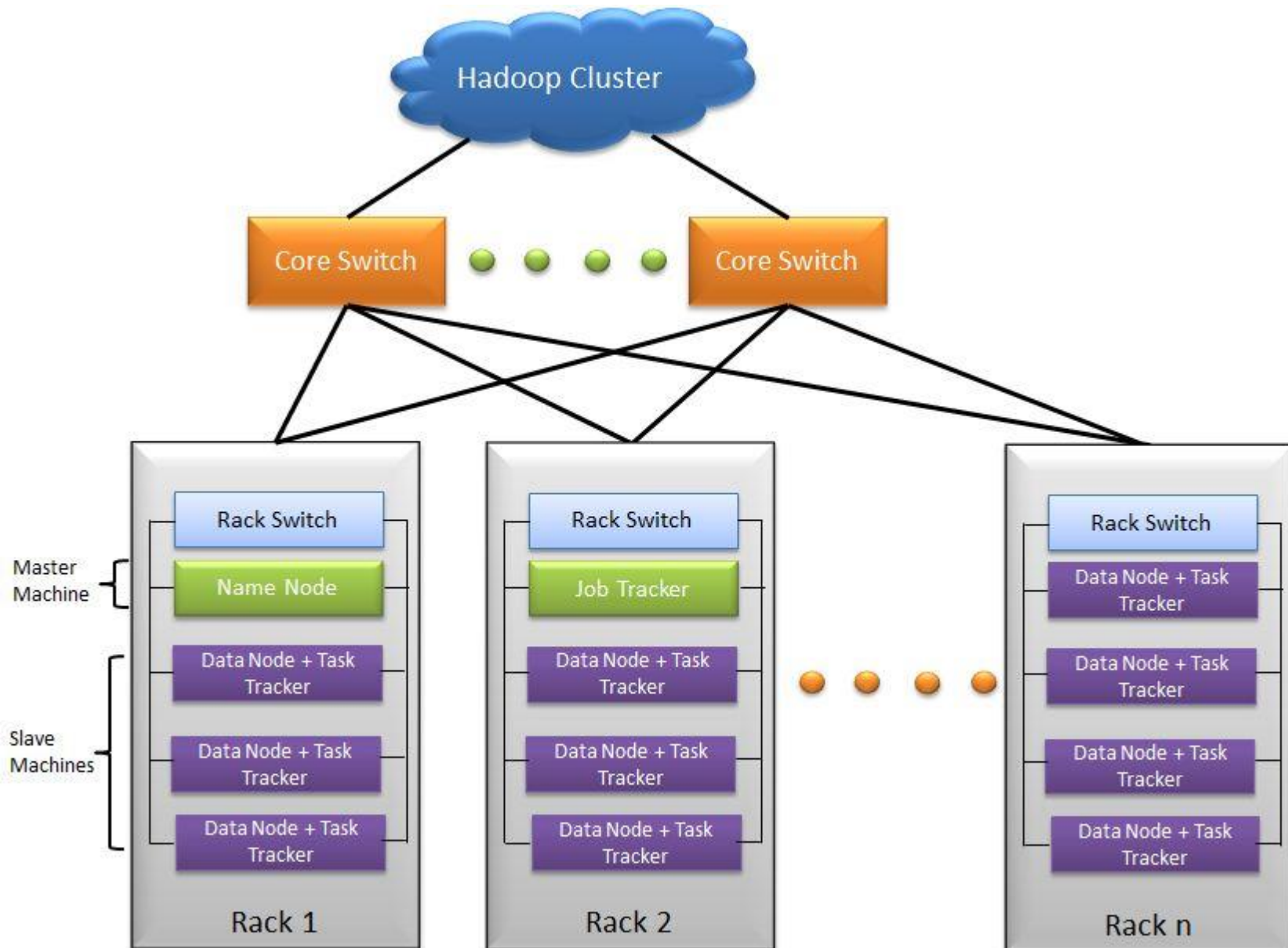
The background consists of a dense field of overlapping, semi-transparent, light gray geometric shapes, primarily polygons and rectangles, scattered across a dark gray gradient. The shapes vary in size and orientation, creating a complex, layered, and crystalline appearance. The overall effect is that of a textured, abstract surface.

# Architektury

# Hadoop



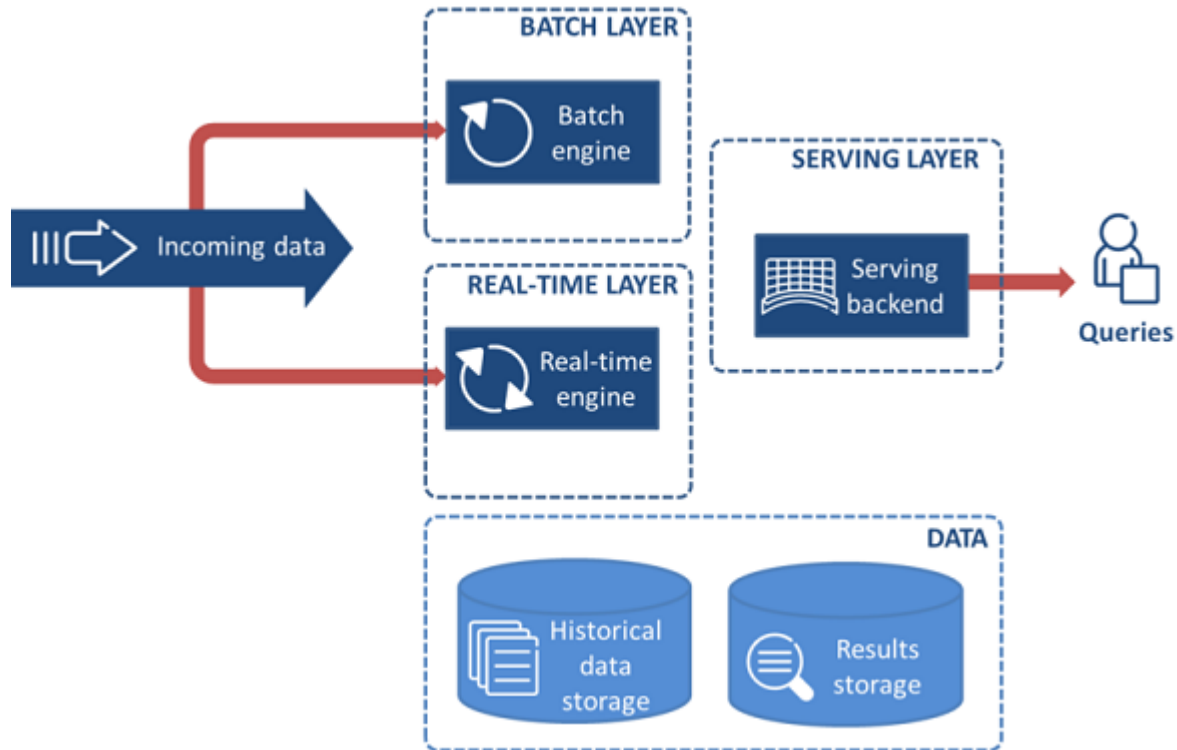
# Hadoop





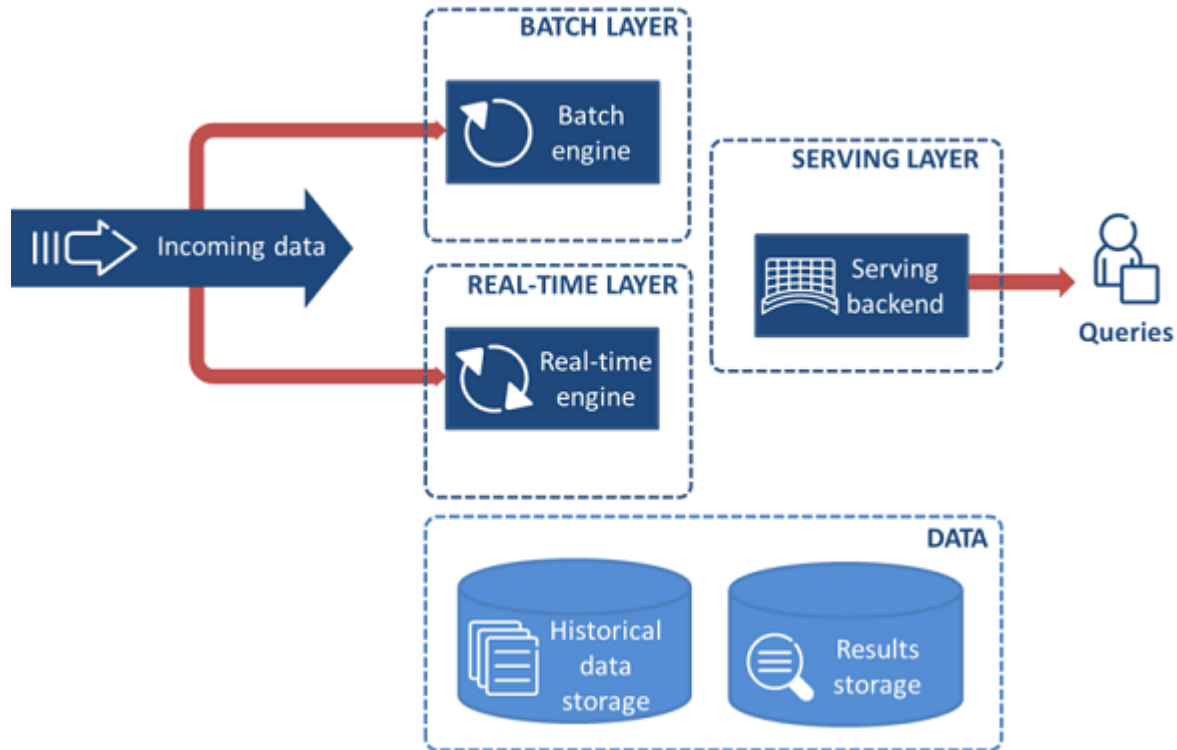
# Lambda

- › Z Apache Storm
- › Nathan Marz, 2011
- › <http://nathanmarz.com/blog/how-to-beat-the-cap-theorem.html>
- › Yahoo, Netflix



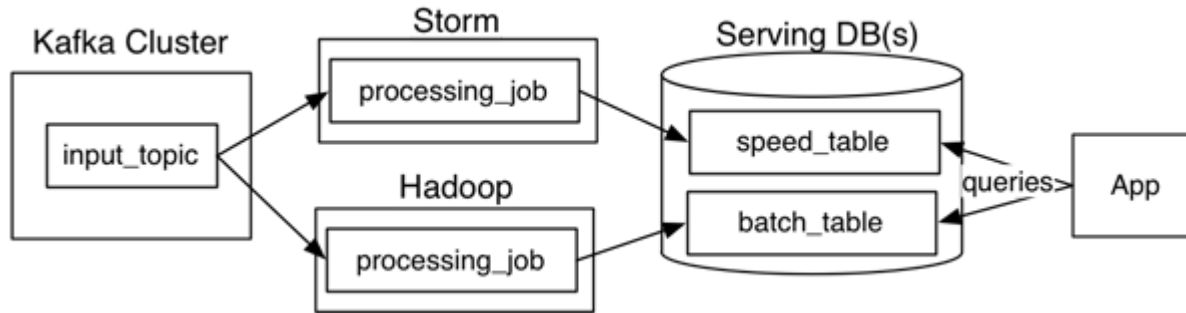
# Lambda

- › 4 vrstvy



# Lambda

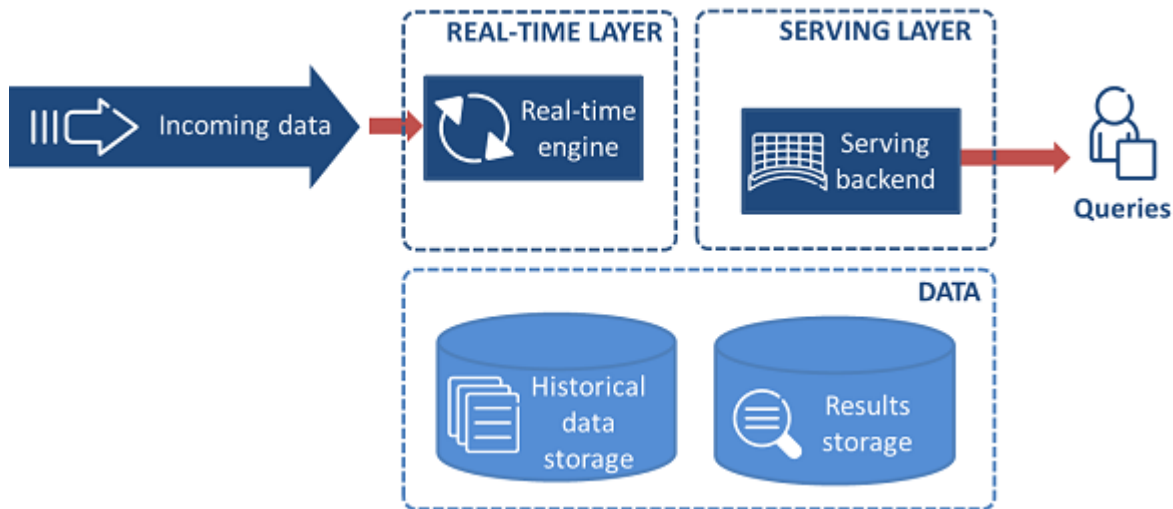
- › Konkrétní technologie





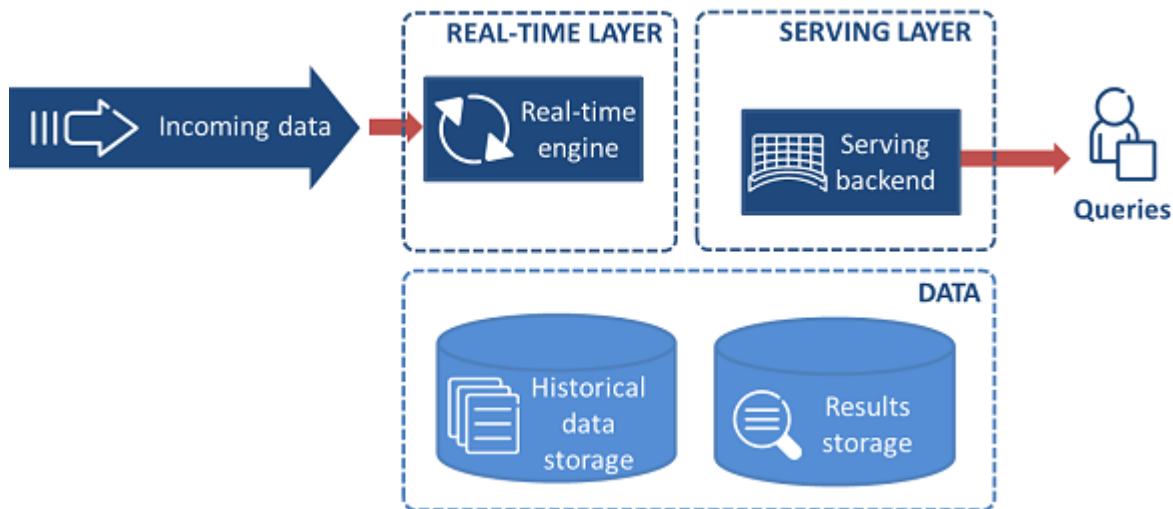
# Kappa

- › 2014 Jay Kreps – LinkedIn
- › <https://www.oreilly.com/ideas/questioning-the-lambda-architecture>



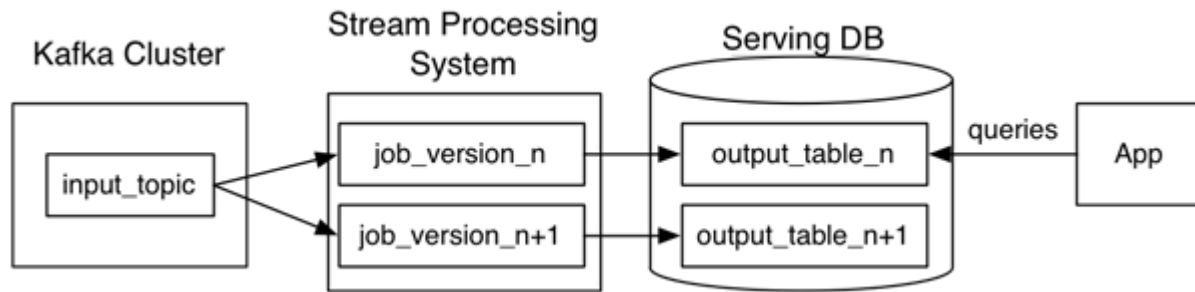
# Kappa

- › 3 vrstvy – odstranění batch vrstvy
- › Lze použít dlouhou retenci
- › Problém se stavem – microbatche?



# Kappa

- › Konkrétní technologie



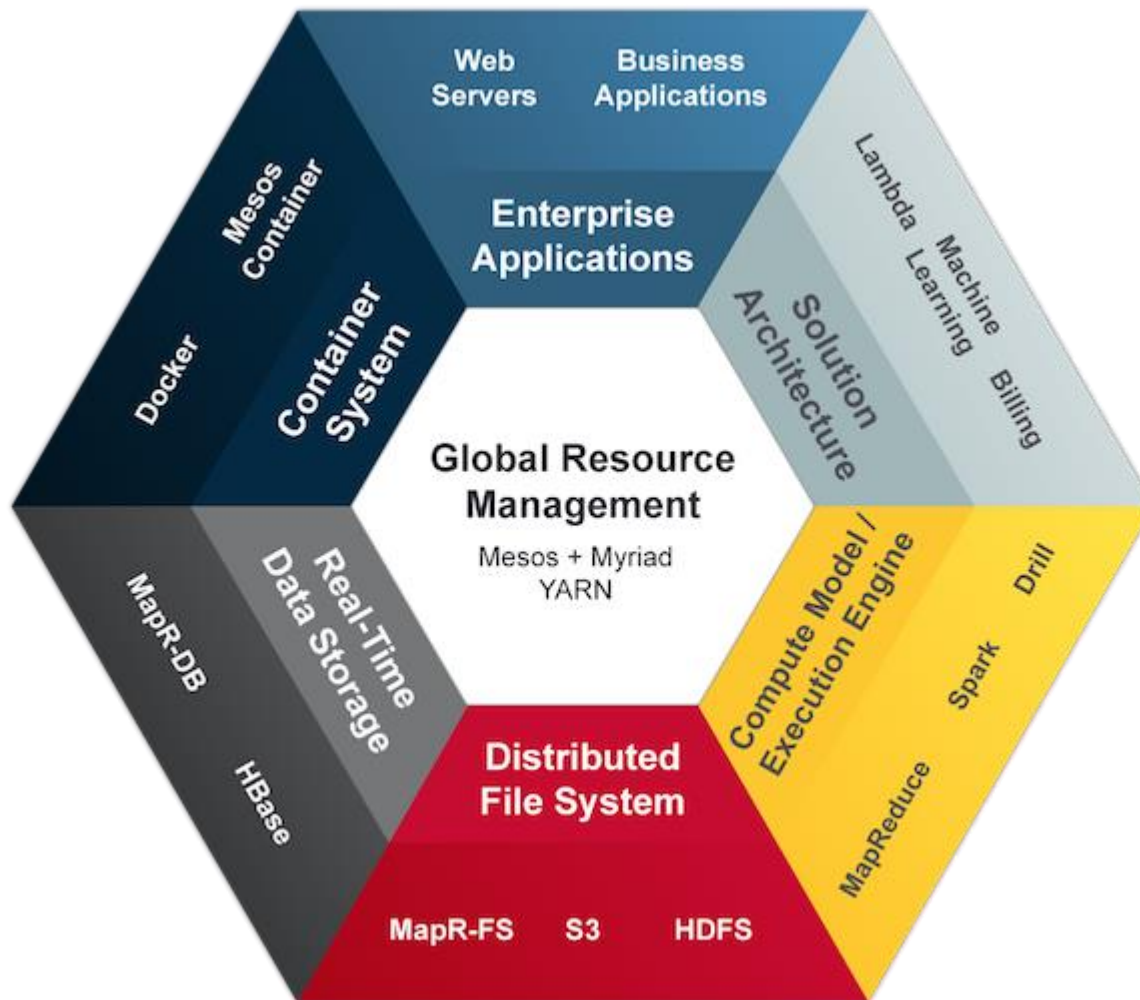


# Kappa

- › Tooty
- › Log data store
  - Kafka
- › Streaming computation systems
  - Samza
  - Storm
  - Kafka Streams
  - Flink

# Zetta

- › Jim Scott - MapR
- › (Zetta je 6 číslo řecké abecedy), data-centric



# Zetta

- › Google Zetta





# Zetta

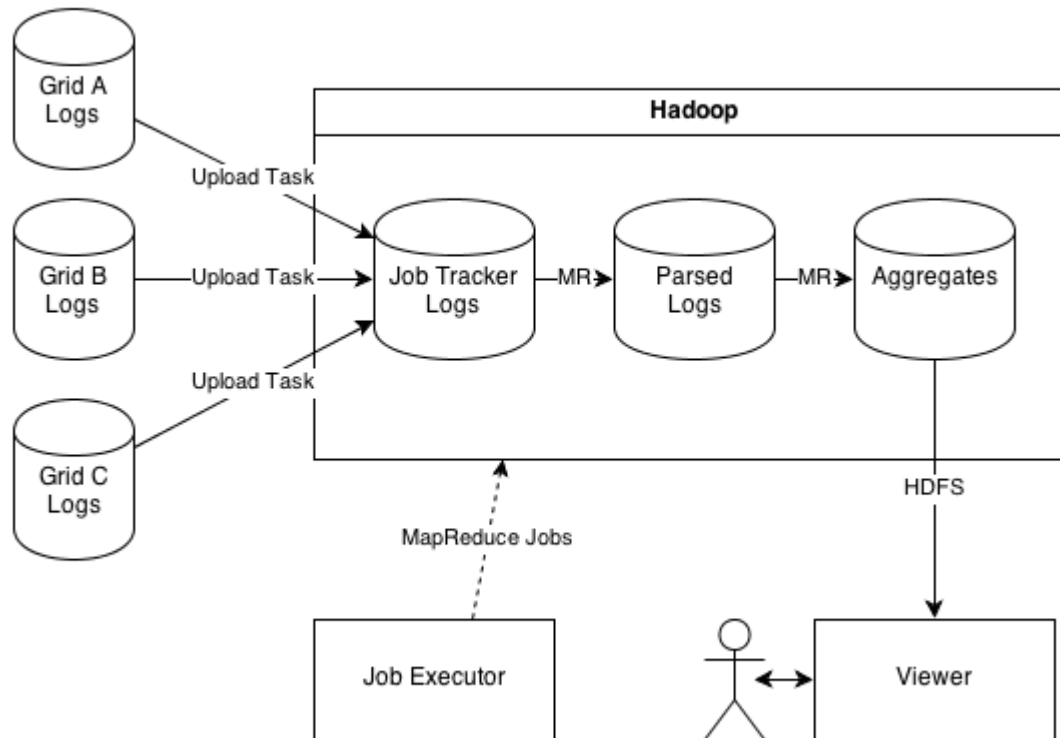
- › Co to znamená?
- › Všechno na Mesos
- › Dynamická alokace zdrojů
- › Omezení přesunů dat
- › Zatím spíš okrajová

## Několik ukázek ze života

- › Sběr logů
- › Reklamní platforma
- › DWH Offloading
- › Analytické pískoviště

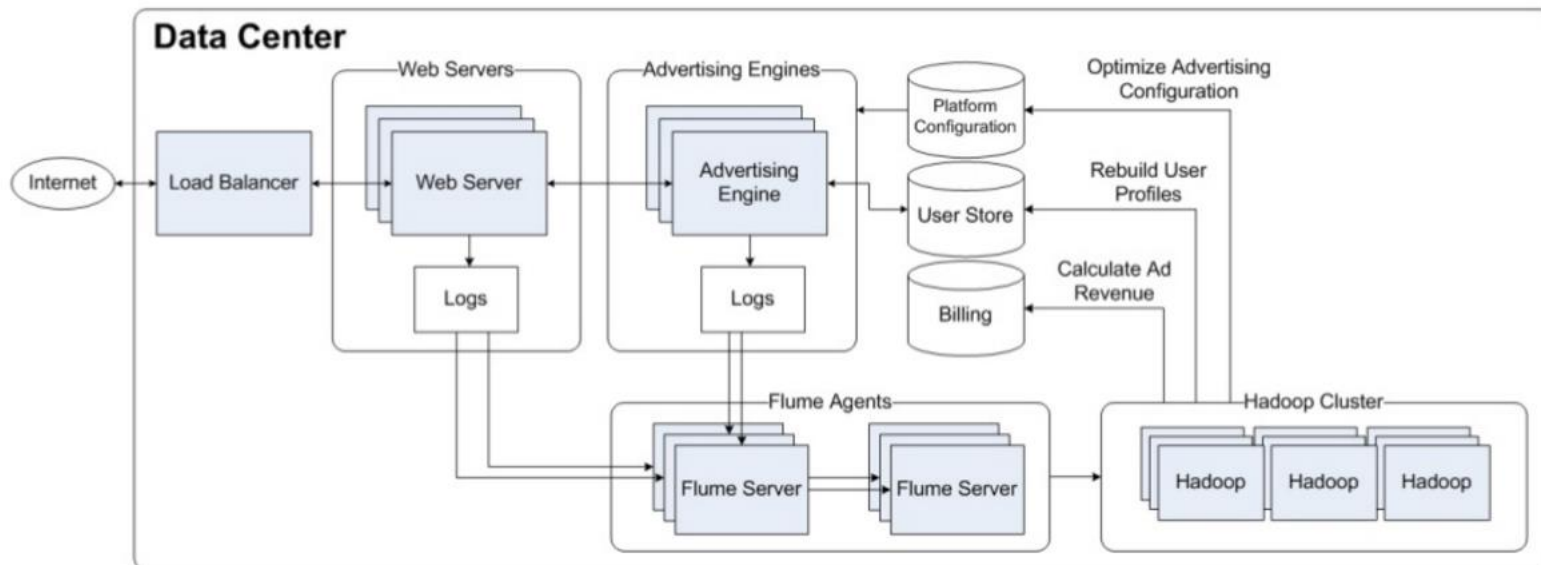
# Sběr logů

- › Web server tvoří logy
- › Ty se ukládají na disk – rotace
- › Pak se posílají na jiné servery
- › Logy se zpracovávají



# Reklamní platforma

- › Web logy a informace o zobrazování reklamy
- › Logy Flumem do HDFS
- › Pak počítáme a vracíme zpět na znovuzpracování



# DWH offloading

- › Aktivní archiv
- › Levnější úložiště dat
  
- › Typicky Sqoop
- › Flume
- › ETL řešené v Hive, Sparku, nebo přes nástroje třetích stran



# DWH offloading

- › Je možné dělat vrstvy jako v normálním DWH
- › Tzn. L0, L1, L2
- › ETL řízené např. pomocí Oozie
- › Často ale komerční nástroje – Talend ETL, Informatica BDM, Oracle ODI
- › Ne vždy to je ale výhra

# Analytické pískoviště

- › Data nahrávána většinou ad-hoc
- › Standardní přísun dat přes Flume/Sqoop/scp
- › Velkou roli má R, python a Spark (pySpark)
- › Využití toolů jako Zeppelin, Jupyter, Hue Notebook, či Cloudera Workbench
- › Většinou se moc neřeší bezpečnost



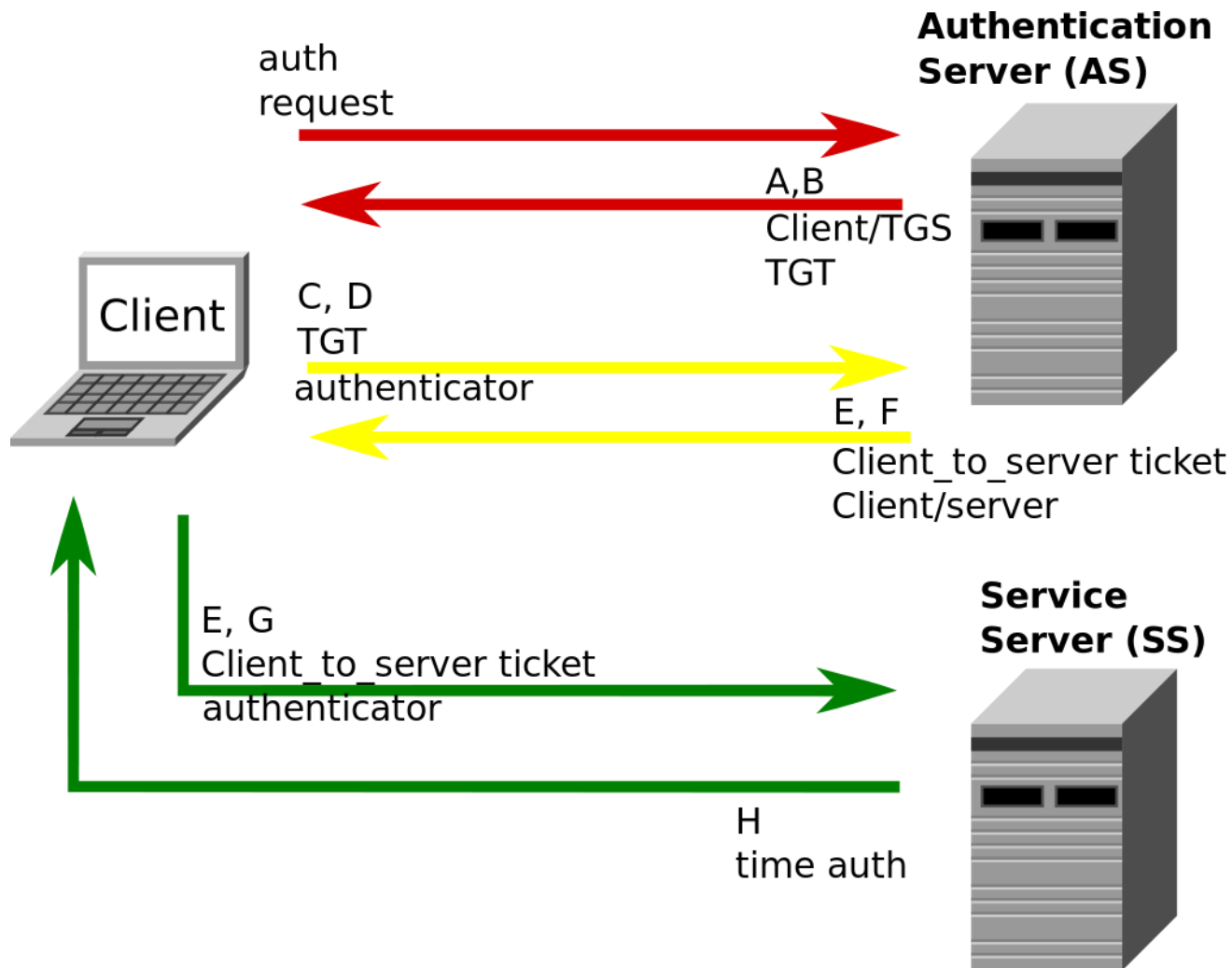
# Security a vliv na architekturu

# Kerberos

- › Lze rozdělit na tři části
  - KDC – Kerberos Distribution Center
  - Server – poskytuje služby
  - Klienti – uživatelé, počítače, služby
- › KDC nabízí
  - AS – Autentizační server
  - TGS - Ticket Granting Service
- › Pojmy
  - TGT – Ticket Granting Ticket
  - Service Server

# Kerberos

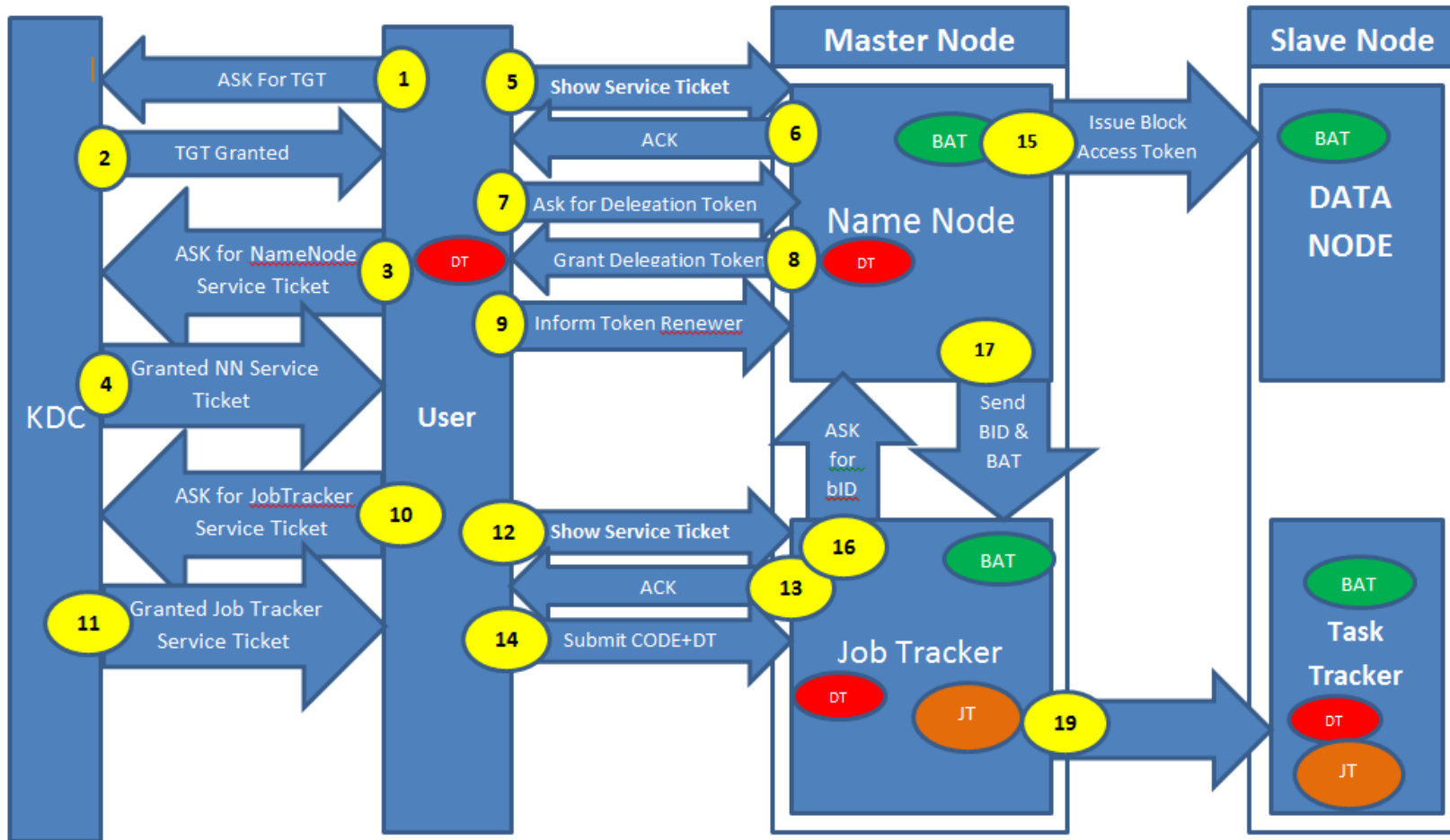
- [https://en.wikipedia.org/wiki/Kerberos\\_\(protocol\)](https://en.wikipedia.org/wiki/Kerberos_(protocol))





# Kerberos

› Jak to funguje



# Security

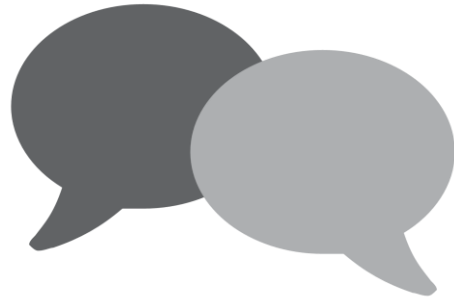
- › HDFS encryption
- › End to end encryption
- › Security komponenty
  - Sentry
  - Ranger

# Security

- › Data locality
- › Data privacy
- › Data labeling
- › Expirace dat
  
- › GDPR
  - Anonymizace
  - Pseudonymizace

# Data masking

- › Jak získat data pro testovací prostředí?
- › Syntetická data?
- › Jak zajistit byznys relevantnost těchto dat?
- › Jak to udělat výkonné a škálovatelné?





# Díky za pozornost

PROFINIT

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



Telefon  
+ 420 224 316 016



Web  
[www.profinit.eu](http://www.profinit.eu)



LinkedIn  
[linkedin.com/company/profinit](https://linkedin.com/company/profinit)



Twitter  
[twitter.com/Profinit\\_EU](https://twitter.com/Profinit_EU)