

DĚLAT  
DOBRÝ SOFTWARE  
NÁS BAVÍ

# PROFINIT

## A4M33BDT

### Technologie pro velká data

## Spark cvičení

Sergii Stamenov  
Hučín Jan

26.4.17

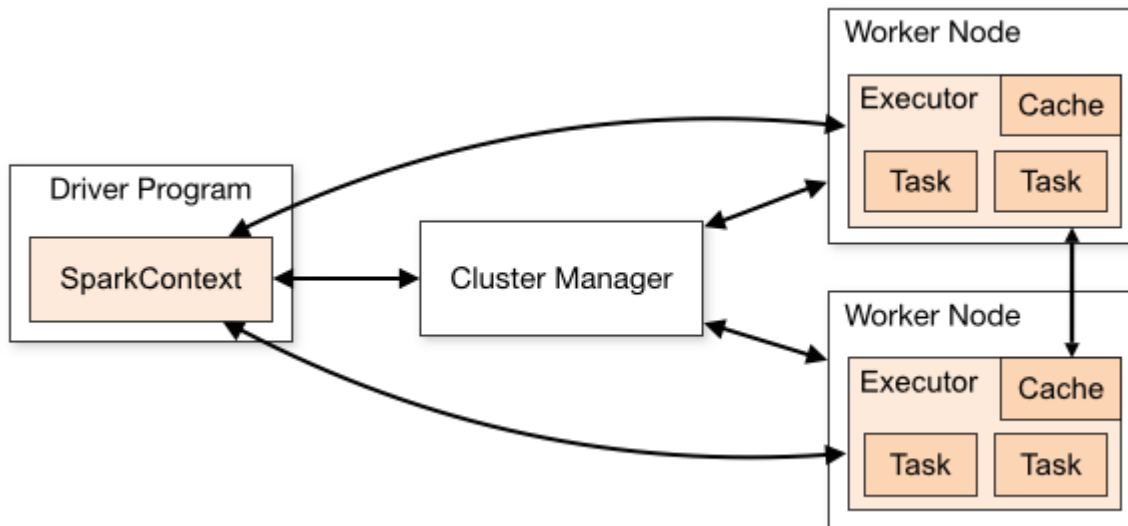


PROFITIT

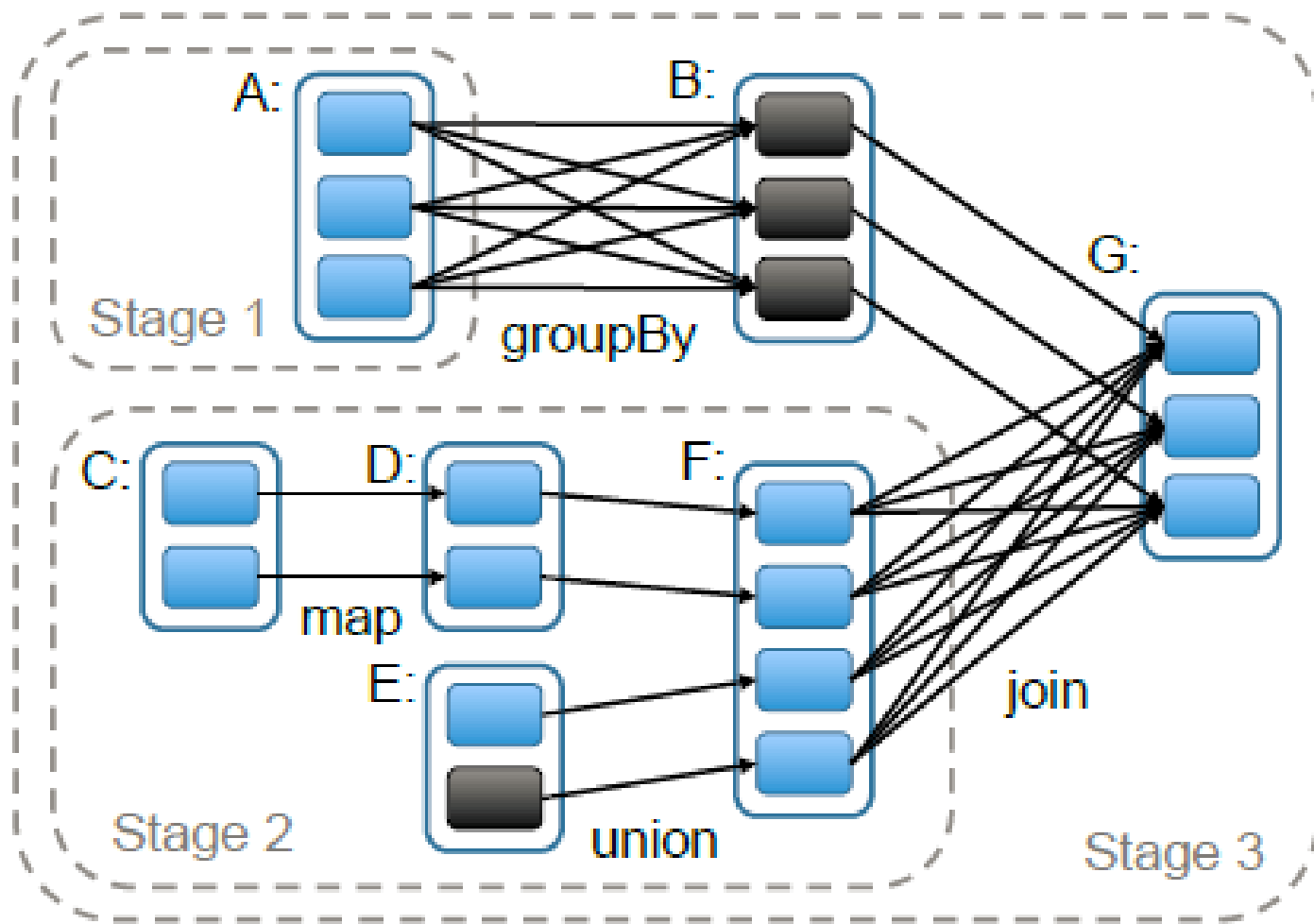


Spark opakování

# Prvky



# Prvky





USE CACHE(), LUKE

# WordCount ukázka

# Úloha #1

## Data

- › 380 000 textu písniček z MetroLyrics
- › Formát: csv
- › Polička: `index, song, year, artist, genre, lyrics`
- › cesta `hdfs:///user/pascep/lyrics_data/`

## Zadáni

- › Pro každou písničku spočítat počet slov
- › Vypsát 10 nejdelších

## Bonus

- › Pomoci akumulátoru spočítat počet chyb (viz Spark **Accumulators**)



# Úloha #2

## Zadání

- › Pro každé slovo spočítat počet písniček ve kterých se to slovo vyskytuje (aka Inverted Index)
- › Vypsát 50 nejčastějších slov

## Bonus

- › Odfiltrovat stop-slova (viz Spark broadcast variable). Seznam stop-slov `/storage/brno2/home/pascepel/SmartStopList.txt`
- › Napočítat počty per genre

## Další zdroje informace

- › Youtube kanal Apache Spark  
<https://www.youtube.com/user/TheApacheSpark>
- › Blogy Cloudera, Databricks

# Děkujeme 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)