

# Obrada podataka tehnologijom Apache Spark

## ZAVRŠNI RAD br. 4573

Martin Matak

Mentor: doc. dr. sc. Marko Čupić  
Fakultet elektrotehnike i računarstva  
Sveučilište u Zagrebu

Zagreb, srpanj 2016.

# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 Prvi programi
  - Postavljanje temelja
  - Otporni rasopdijeljeni skup podataka
- 4 Napredno programiranje
  - Algoritam PageRank
  - Skupovi podataka kao uređeni parovi
  - Dijeljene varijable

# Motivacija



# Osnovni gradivni elementi

Spark  
SQL

Spark  
Streaming

MLib  
(strojno  
učenje)

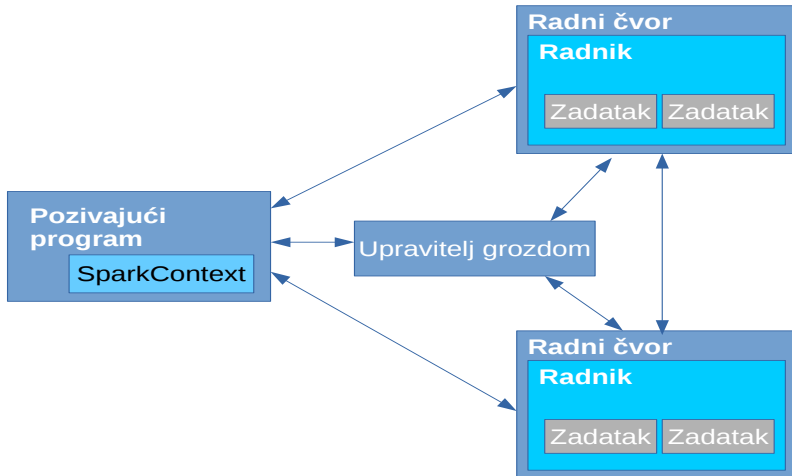
GraphX  
(obrada  
grafova)

Jezgra

# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 **Prvi programi**
  - Postavljanje temelja
  - Otporni rasopdijeljeni skup podataka
- 4 Napredno programiranje
  - Algoritam PageRank
  - Skupovi podataka kao uređeni parovi
  - Dijeljene varijable

# Osnovni elementi aplikacije



# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 Prvi programi**
  - Postavljanje temelja
  - **Otporni rasopdijeljeni skup podataka**
- 4 Napredno programiranje
  - Algoritam PageRank
  - Skupovi podataka kao uređeni parovi
  - Dijeljene varijable

# Otporni raspodijeljeni skup podataka

## Resilient distributed dataset - RDD

- Nepromjenjiva kolekcija podataka



# Otporni raspodijeljeni skup podataka

## Resilient distributed dataset - RDD

- Nepromjenjiva kolekcija podataka

### Transformacije

Operacije koje iz jednog skupa kreiraju drugi, novi skup podataka.  
**Lijena evaluacija** - pokreću ih akcije.

# Otporni raspodijeljeni skup podataka

## Resilient distributed dataset - RDD

- Nepromjenjiva kolekcija podataka

### Transformacije

Operacije koje iz jednog skupa kreiraju drugi, novi skup podataka.  
**Lijena evaluacija** - pokreću ih akcije.

### Akcije

Dohvat jednog ili više elemenata iz nekog skupa.

# Otporni raspodijeljeni skup podataka

Resilient distributed dataset - RDD

- Nepromjenjiva kolekcija podataka

## Transformacije

Operacije koje iz jednog skupa kreiraju drugi, novi skup podataka.

**Lijena evaluacija** - pokreću ih akcije.

## Akcije

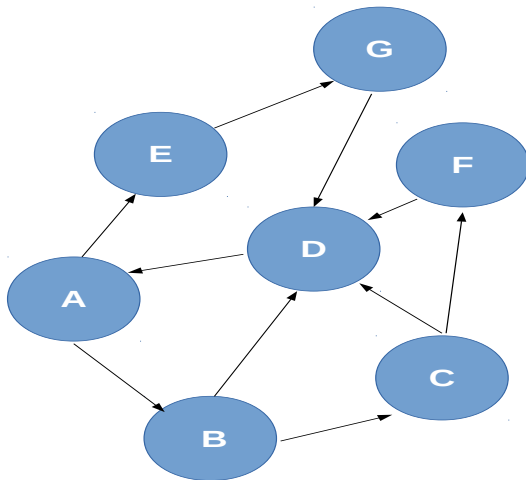
Dohvat jednog ili više elemenata iz nekog skupa.

- Stanje u memoriji?

# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 Prvi programi
  - Postavljanje temelja
  - Otporni rasopdijeljeni skup podataka
- 4 Napredno programiranje**
  - Algoritam PageRank**
  - Skupovi podataka kao uređeni parovi
  - Dijeljene varijable

## Koja stranica je najvažnija?



# Algoritam PageRank

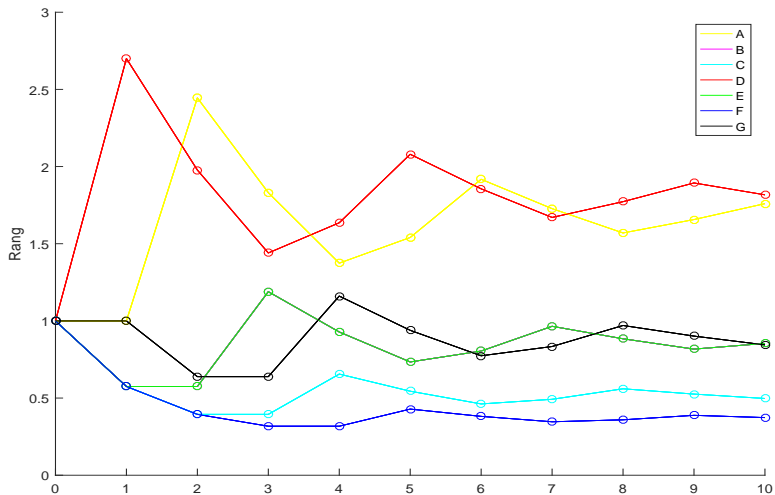
---

## Algoritam 1 Algoritam Pagerank

---

- 1: Početni rang svake stranice postavi se na 1.0
  - 2: **for**  $i = 1$  to  $P$  **do**
  - 3:   Svaka stranica  $n$  šalje svojim susjedima doprinos  
       $\text{rang}(n) / \text{brojSusjeda}(n)$
  - 4:   Postavi ukupni rang stranice prema formuli:  
       $0.15 + 0.85 * \text{ukupan primljeni doprinos}$
  - 5: **end for**
-

# Rezultati



# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 Prvi programi
  - Postavljanje temelja
  - Otporni rasopdijeljeni skup podataka
- 4 **Napredno programiranje**
  - Algoritam PageRank
  - **Skupovi podataka kao uređeni parovi**
  - Dijeljene varijable



# Uređeni parovi te spremanje i čitanje podataka

- Podatci u obliku ključ - vrijednost

# Uređeni parovi te spremanje i čitanje podataka

- Podatci u obliku ključ - vrijednost
- Posebne **transformacije** i **akcije**

## Uređeni parovi te spremanje i čitanje podataka

- Podatci u obliku ključ - vrijednost
- Posebne **transformacije** i **akcije**
- Mogućnosti spremanja i čitanja: baze podataka, tekstualne datoteke (JSON, CSV, TSV) ...

# Sadržaj

- 1 Uvod
- 2 Osnovni gradivni elementi
- 3 Prvi programi
  - Postavljanje temelja
  - Otporni rasopdijeljeni skup podataka
- 4 Napredno programiranje**
  - Algoritam PageRank
  - Skupovi podataka kao uređeni parovi
  - Dijeljene varijable**

## Dijeljene varijable: odašiljači i akumulatori

- **Problem:** Svaki zadatak na grozdu ima svoju kopiju varijabli

## Dijeljene varijable: odašiljatelji i akumulatori

- **Problem:** Svaki zadatak na grozdu ima svoju kopiju varijabli

### Odašiljatelj

Nepromjenjiva varijabla koja zauzima malo memorije čiju vrijednost je moguće dohvatiti na cijelom grozdu.

## Dijeljene varijable: odašiljatelji i akumulatori

- **Problem:** Svaki zadatak na grozdu ima svoju kopiju varijabli

### Odašiljatelj

Nepromjenjiva varijabla koja zauzima malo memorije čiju vrijednost je moguće dohvatiti na cijelom grozdu.

### Akumulator

Globalna varijabla čiju vrijednost je moguće mijenjati iz cijelog grozda.

# Sažetak

- Tehnologija za obradu velike količine podataka.
- Algoritam PageRank.
- Dijeljene varijable.
- Još bi trebalo:
  - Postaviti i isprobati Apache Spark na grozdu.
  - Razraditi svaku od komponenata - SparkSQL, MLib, Spark Streaming i GraphX.