

UVOD

MapReduce

- Velik broj podataka
 - o Distribuirani kroz velik broj računala
 - o Unutar razumnog vremena
- Obrada za specijalne svrhe
 - o Indeksirani dokumenti, dnevnik web zahtjeva

Nedostaci:

- Obradni paralelizam
- Distribucija podataka
- Tolerancija na pogrešku

Iz potrebe za rješavanjem ovih problema dolazi do novih software-a

Paralelizam:

- Velika nakupina hardvera
- Spojeni preko Etherneta ili jeftinih switch-eva

Novi software-i:

- Distribuirani datotečni sustav
- MapReduce

Primjer (brojač riječi):

- Ogromni tekstualni dokumenti
- Analiza serverskih dnevnika kako bi se pronašli popularni URL-ovi
- Brojanje koliko puta se svaka izričita riječ pojavila u datoteci

Rješenje:

- Datoteka prevelika za memoriju
- <riječ, brojač> par stane u memoriju
- `words(doc.txt) | sort | uniq -c`

DEFINICIJA

MapReduce – programski model i povezana provedba za obradu i stvaranje velikog skupa podataka

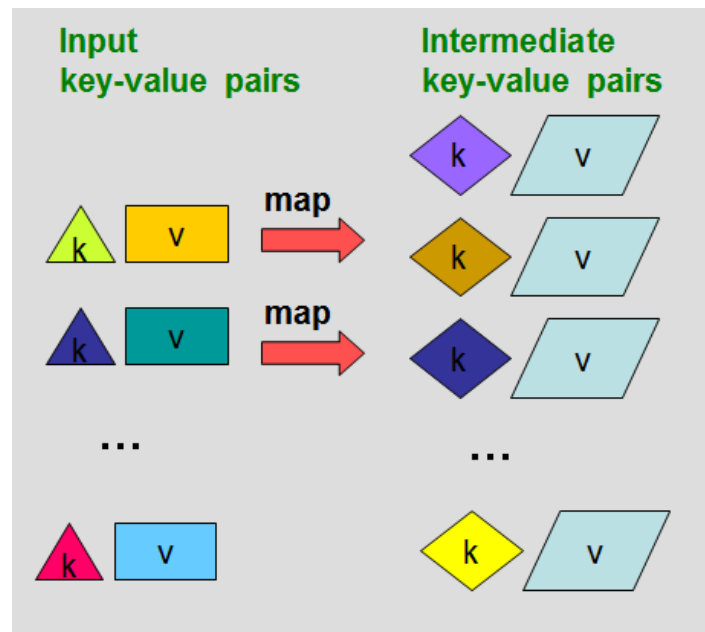
Čitanje podataka – sekvencijalno

Map – izvlačenje vrijedne informacije (<ključ, vrijednost>)

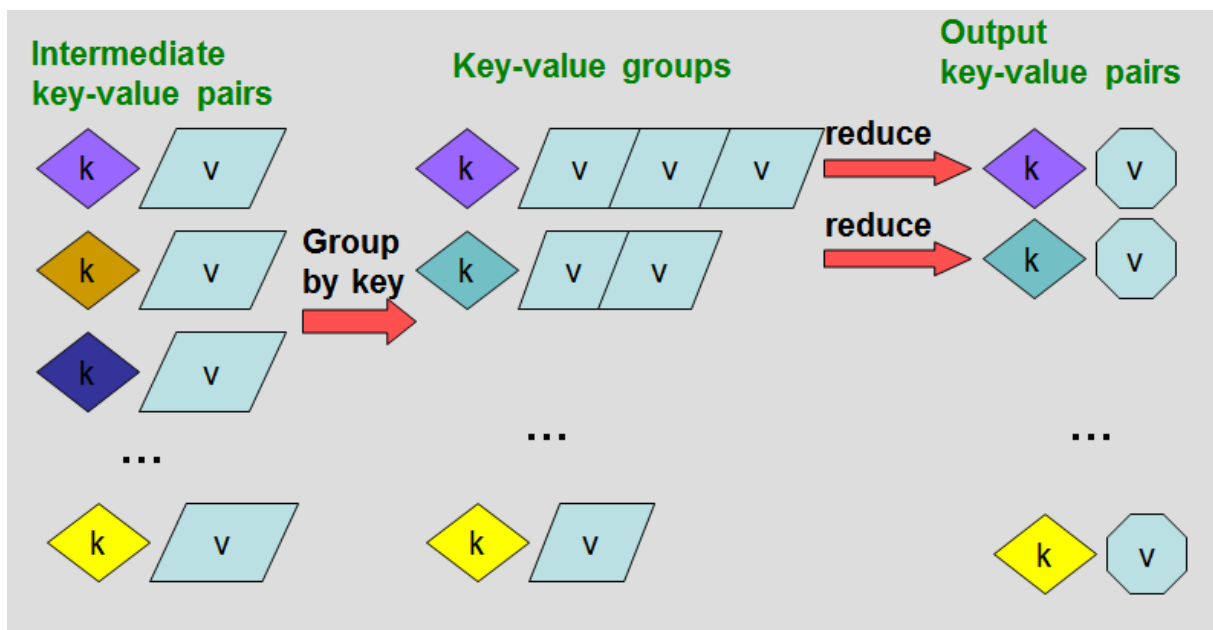
Grupiranje i sortiranje po ključu

Reduciranje – spajanje podataka te generiranje rezultata

MAP



REDUCE



Ulaz je skup parova ključ-vrijednost

$\text{Map}(k, v) \rightarrow \text{list}\langle k' v' \rangle$

- uzima ključ-vrijednost par
- vraća skup parova ključ-vrijednost
- Map poziva za sve parove ključ-vrijednost

$\text{Reduce}(k', \text{list}\langle v' \rangle^*) \rightarrow \langle k'', v'' \rangle^*$

- Vrijednosti v' sa istim ključem k' se reduciraju zajedno i obrađuju se po v' slijedu
- Jedna Reduce funkcija je pozvana po jedinstvenom ključu k'

MapReduce – Broj riječi:

map(String ključ, String vrijednost)

//ključ: naziv dokumenta

//vrijednost: tekst dokumenta

for each word w in value:

emit(w , „1“);

reduce(String ključ, Iterator vrijednosti)

//ključ: riječ;

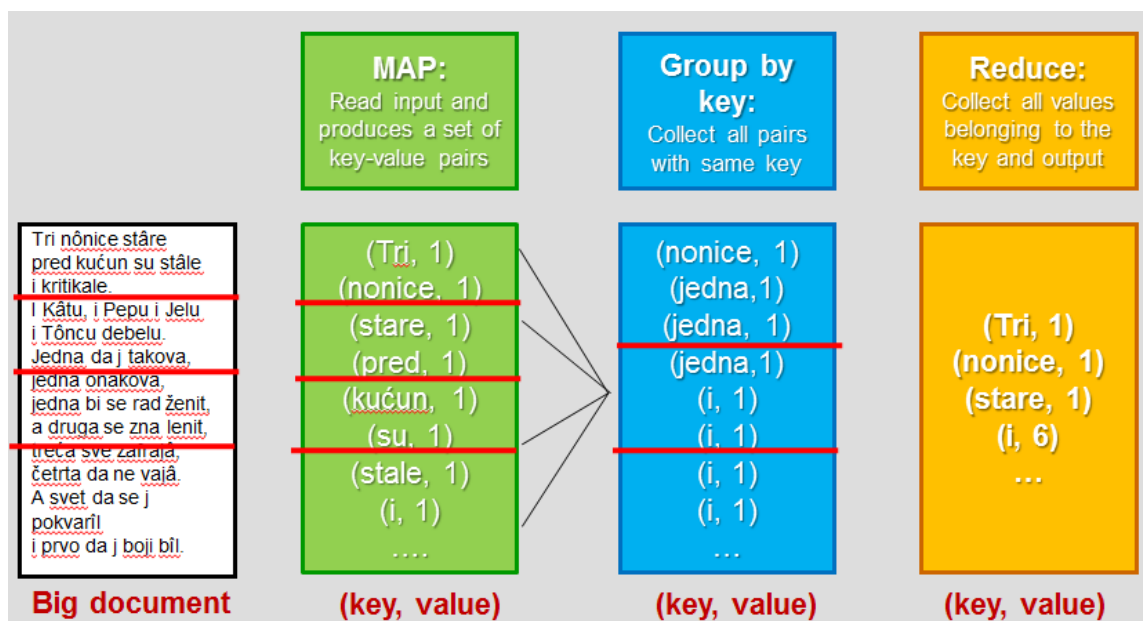
//vrijednost: iterator tijekom brojanja

int rezultat = 0;

for each count v in vrijednosti:

rezultat += Parseint(v);

emit(ključ, rezultat);

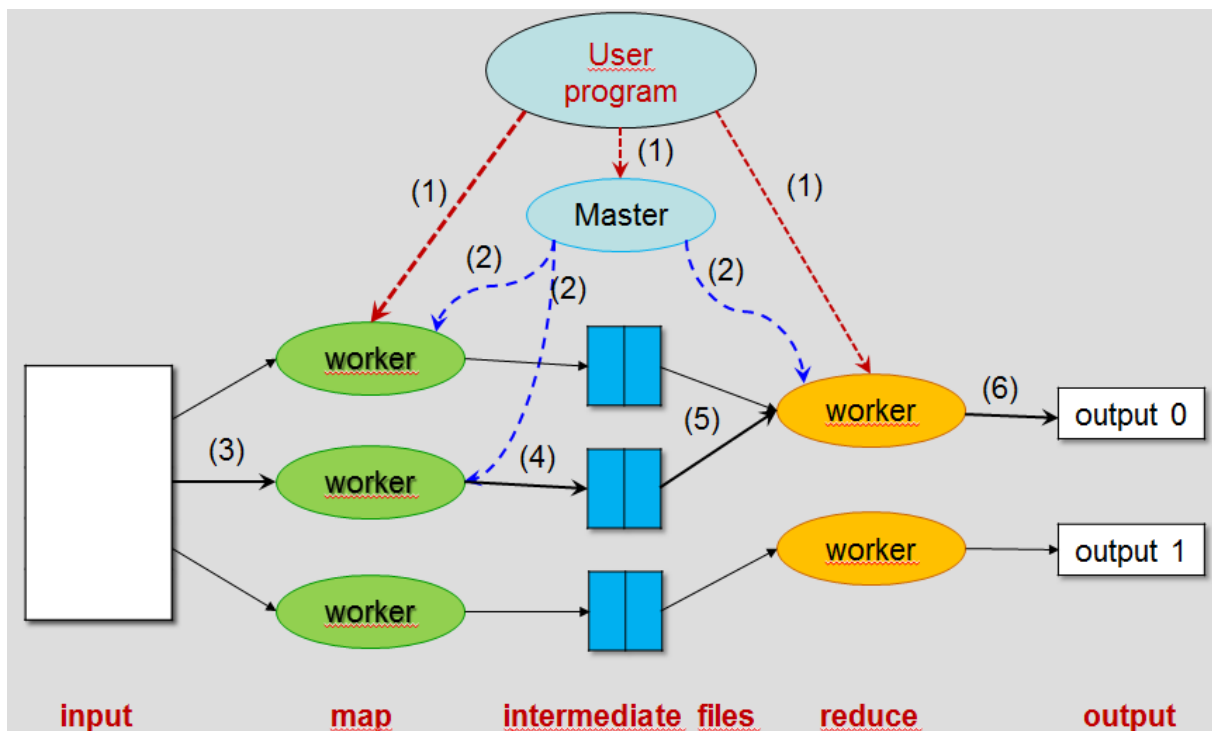


MapReduce – Implementacija

Map-Reduce okruženje se brine o:

- Podjeli ulaznih podataka
- Rasporedu izvršavanja u klasteru
- Grupiranju parova po ključu
- Toleranciji na pogrešku
- Komunikaciji između računala

Korisnik mora napisati funkciju za Map i funkciju za Reduce.



(1) Račvanje

- Podjela ulazne datoteke u M dijelova
- Pokrene više kopija programa na klasteru

(2) Dodjeljivanje

- Glavni program(Master) dodjeljuje M map i R reduce zadataka

(3) Mapiranje

- Map radnik(worker) čita ulazne komade
- Parsira ključ/vrijednost parove iz ulaza
- Prosljeđuje svaki par u definiranu Map funkciju
- Puni lokalnu memoriju izlazom Map funkcije

(4) Lokalno zapisivanje

- Izlaz mapiranja je zapisan na lokalni disk
- Podijeljeni u R regija od strane funkcije za podjelu
- Lokacije podijeljenih dijelova(particija) su prosljeđene glavnom programu(Master)

(5) Grupiranje i sortiranje

- Master obavještava reduce radnika o lokacijama particija
- Reduce radnik čita prijelazne podatke preko RPC (Remote Procedure Call)
- Sortira prijelazne podatke po prijelaznom ključu

(6) Reduciranje

- Prolaziti preko sortiranih prijelaznih podataka
- Reduce funkcija se poziva za svaki jedinstveni ključ prijelaznih podataka i skupa prijelaznih vrijednosti
- Izlaz reduciranja je dodan u izlaznu datoteku

MapReduce: Master node

Master – struktura podataka

- Stanje zadataka (idealno, u izvođenju, završeno)
- Identitet radnih strojeva
- Veličina i lokacija prijelaznih podataka
 - o R regija prijelaznih podataka nastalih Mapiranjem
 - o Poslani Master-u po završetku Mapiranja

Master – funkcije

- Šalje loakcije prijelaznih podataka zadatku Reduciranja
- Periodički proziva radnike

MapReduce: Fault tolerance

Pogreška radnika:

- Nema odgovora od radnika – greška radnika
- Zadaci Mapiranja dodijeljeni radniku s greškom resetiraju se na idle
 - o Bili oni završeni ili u izvođenju
 - o Preraspodijela po drugim radnicima
- $A(\text{map}) \rightarrow B(\text{map})$
 - o Zadaci za Reduce će biti obaviješteni
 - o Zadaci za Reduce će čitati podatke od radnika B
- Reduce greška – zadaci u izvođenju se resetiraju
 - o Zadatak za Reduciranje se resetira
- Pogreška Mastera
 - o Prekida se MapReduce operacija

MapReduce: Task granularity

M Map i R Reduce zadataka

M mnogo veći od broja radnika

- Poboľjšano dinamičko balansiranje prilikom opterećenja
- Ubrzava obnovu kada dođe do greške radnika
- Jedan DFS dio po mapi (16 -64mb)

R – često ograničen od strane korisnika

- Svaki zadatak reduciranja u odvojenu datoteku
- Mali broj broja radnika

Fizičke granice

- $O(M+R)$ odluka raspoređivanja, $O(M*R)$ stanja

MapReduce: Backup Tasks

„Zaostali“ – spori radnik

- Uzima previše vremena za izvršavanje zadatka
- Produljuje vrijeme posla
 - o Loš disk, obavljanje drugih poslova na mašini

Rješenje:

- Pred kraj MapReduce-a, Master prerasporedi rezervno izvršavanje zadataka u izvođenju
- Zadatak je gotov koji god zadatak prije završi (primarni ili rezervni)
- Značajno smanjuje vrijeme posla
- Ne povisuje resurse drastično

MapReduce: Refinements (preciziranja)

Funkcija podjele

- Funkcija podjele na prijelaznom ključu
- Defaultni – $\text{hash}(\text{key}) \bmod R$
- Zapisi sa istim ključem moraju biti kod istog radnika za Reduce
- Primjer – izlazni ključevi su URL-ovi
 - o Svi ulazi za pojedinog hosta u istoj datoteci
 - o Prilagođene funkcije podjele
 - o $\text{Hash}(\text{Hostname}(\text{urlkey})) \bmod R$

Jamstva redoslijeda

- Unutar dane particije, prijelazni parovi ključ/vrijednost se obrađuju u slijedu povećanja ključa
- Lagano za generiranje sortiranih izlaza po particiji

Preskakanje loših zapisa

- Deterministički pad Map i Reduce
 - o Bug u korisničkom kodu
 - o Prihvatljivo ignoriranje nekih zapisa (na osnovi statističke analize)
- Detektiranje zapisa koji su prouzročili pad i preskakanje istih

Funkcija kombiniranja

- Značajno ponavljanje prijelaznih ključeva nastalih funkcijom Mapiranja $(k, v_1)(k, v_2)$ za isti ključ k .
 - o Reduciranje je komutativno i asocijativno
- Funkcija kombiniranja – spajanje podataka particija prije slanja preko mreže
 - o $\text{Combine}(k, \text{list}(v_1)) \rightarrow v_2$
 - o Izvršeno na svakoj mašini koja izvodi Mapiranje
- Izlaz Reduca \rightarrow izlazna datoteka
- Izlaz kombinacije \rightarrow ulaz Reduce-a

Brojači

- Brojati pojave različitih događaja(Map)
- Master prikuplja vrijednosti brojača i vraća ih korisničkom kodu

Ulaz i izlaz

- Tekstualni način unosa – svaka linija je par ključ/vrijednost
- Svaka implementacija zna kako podijeliti podatke
- Korisnik mora implementirati jednostavno sučelje za čitanje
- Izlaz funkcionira na sličan način

Problemi pogodni za MapReduce

MapReduce nije rješenje za sve probleme

- Mnogo čvorova u paraleli
- Velike i rijetko ažurirane datoteke
- Primjer:
 - o Online maloprodaja (Amazon) – nije prikladan
 - o Računanje ranka stranice (Google) – prikladan

Brojanje učestalosti pristupa URL-u

- Funkcija mapiranja obrađuje zapise zahtjeva web stranice
 - o Izlaz (URL, 1)
- Funkcija Reduca dodaje zajedno sve vrijednosti za isti URL
 - o Emits (URL, total count)

Obrnuti Web-Link graf

- Map: (cilj, izvor) za svaku poveznicu(link) prema ciljanom URL-u, koji je pronađen na stranici izvora
- Reduce: konkatencija svih izvora URL-a povezanih sa ciljanim URL-om, (target, list(source))

Distribuirani grep

- Map: emitiraju liniju ako odgovara uzorku
- Reduce: kopira podržane podatke u izlaz

Distribuirano sortiranje

- Map: izvlačenje ključa iz svakog zapisa
 - o Emits(key, record)
- Reduce: emitira sve nepromijenjene parove
- Garantira slijed i funkcija podijele

Pojam-Vector po domaćinu

- Sažima najvažnije riječi u dokumentu kao popis parova (riječ, učestalost)
- Map: emits (hostname, termvector) za svaki dokument
- Reduce: po dokumentu izrazvektor za određenog domaćina
 - o Dodaje izrazvektore zajedno
 - o Ignorira rijetke izraze
 - o Emitira parove (domaćin, izrazvektor)

Matrix – Vector multiplication

- $n \times n$ matrica M
 - o $m_{i,j}$ – element u i -to retku i j -tom stupcu
- Vektor v duljine n sa j -tim elementom v_j
- (vektor) $x = M$ (vektorski) v
- $x_i = \sum_{j=1}^n m_{ij} * v_j$
- M i v su spremljeni u DFS
 - o n je velik, ali (vektor) v može biti spremljen u memoriju
 - o pozicije $m_{i,j}$ i v_j su vidljivi

Funkcija za Map

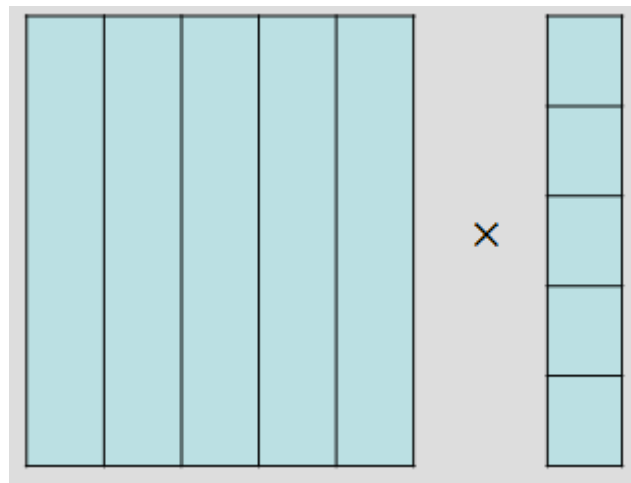
- primjeni se na jednom elementu od M
- (vektor) v se čita prvi, da bude dostupan svim Map radnicima
- Za svaki m_{ij} map emits($i, m_{ij} * v_j$)

Funkcija za Reduce

- Zbroji sve vrijednosti povezane sa ključem i
- Rezultat je par (i, x_i)

Kada vektor v ne može stati u memoriju

- Podjela matrice M i vektora v na pruge
- Map: komad M pruge i cijelu v prugu



Upiti nad bazom podataka

- Upiti preveliki za zajedničke relacijske baze podataka
- SQL, relacija R , n -torke
- Relacija R može biti spremljena u DFS

Relacijsko-algebarske operacije

- Selekcija – stanje $C, \sigma_C(R)$
- Projekcija – podskup $S, \pi_S(R)$
- Unija, presjek, razlika
- Nature join $R \bowtie S$
- Grupiranje i agregacija

Selekcija

- Stanje C , $\sigma_C(R)$
- Map – ako n -torka $t \in R$ zadovoljava uvjet C , emit (t, t)
- Reduce – identitet

Projekcija

- Podskup S , $\pi_S(R)$
- Map
 - $\forall t \in R$ izgradi n -torku t' bez komponenti koje nisu u S
 - Emit (t', t')
- Reduce – eliminacija duplikata
 - Reduce $(t', [t', t', \dots, t'])$ u (t', t')

Natural join $R(A,B) \bowtie S(B,C)$

A	B		B	C		A	C
a ₁	b ₁		b ₂	c ₁		a ₃	c ₁
a ₂	b ₁		b ₂	c ₂		a ₃	c ₂
a ₃	b ₂		b ₃	c ₃		a ₄	c ₃
a ₄	b ₃						
R		\bowtie	S		=	$R \bowtie S$	

- Map: za ulaz $R(a, b)$ emit $(b, (a, R))$
 - Za ulaz $S(b, c)$ emit $(b, (c, S))$
- Reduce: spoji $(b, (a, R))$ sa $(b, (c, S))$
 - Emit (a, b, c)

Množenje matrica

- $P = M \times N$
- $p_{ik} = \sum_j m_{ij} \cdot n_{jk}$
- Map
 - Za svaki m_{ij} emit $(j, (M, i, m_{ij}))$
 - Za svaki n_{jk} emit $(j, (N, k, n_{jk}))$
- Reduce
 - Za svaki ključ jistraži listu vrijednosti
 - Za svaku vrijednost iz (M, i, m_{ij}) i iz (N, k, n_{jk}) napravi par ključ-vrijednost $((i, k), (m_{ij} \cdot n_{jk}))$
- Map2:
 - Identitet - $((i, k), v) \rightarrow ((i, k), v)$
- Reduce2:
 - Za svaki ključ (i, k) napravi sumu vrijednosti sa tim ključem
 - Rezultat: par $((i, k), v)$ gdje je $v = P_{ik}$

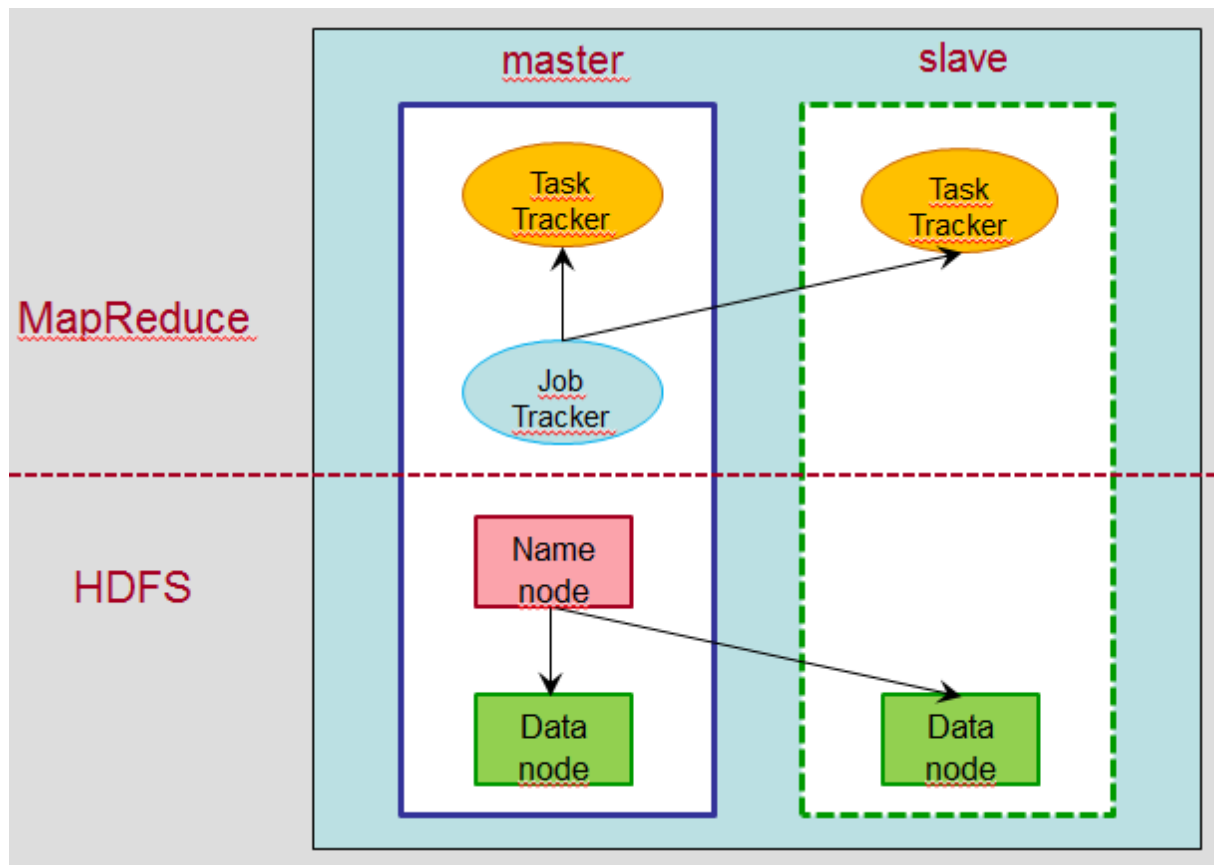
HADOOP – implementacija MapReduce-a otvorenog koda

- Napisan u Javi
- Koristi HDFS za spremanje
- Hbase – distribuirana baza podataka

HDFS

- Distribuirani, skalabilni i portabilni datotečni sustav
- 128 Mb blokovi
- Replikacija – 3 kopije svakog bloka
 - Dvije na istoj, jedna na različitoj polici
- Čvor naziva i čvor podataka

Konfiguracija klastera



Glavni čvor – JobTracker

- MapReduce program (job)
 - o Job – JAR arhiva, konfiguracijska datoteka i particije
- Pruža JobClient s jedinstvenom identifikacijskom oznakom za posao
- Pridodjeli zadatke radnicima
 - o Zadržati posao što je bliže moguće podacima
 - o Svijest o policama – JobTracker zna koji čvor sadrži podatke i koje druge mašine su u blizini
- Izvršava zadatke konkurentno na radnicima
 - o 4 utora (2 map, 2 reduce) na svakom radniku

JobClient

- Podijela i distribucija podataka preko HDFS-a

Job queue

- Poslovi dostupni u HDFS

Job scheduler – inicijalizacija zadaća (Map i Reduce)

Čvor radnika – TaskTracker

- Izvršava zadatke
- Periodički obavještava Mastera o svom stanju
 - o Prazni Map i Reduce utore
- Stvara nove instance JVM- a i izvršava JAR zadatke
 - o Spriječava rušenje TaskTracker-a ako posao sruši JVM

Tolerancija na pogreške

- Ako master čvor dobiva poruke o greška stanja od radnika
 - o Master prerasporedi zadatke
 - o Izbjegava preraspodjelu na istom čvoru
- Ako master ne dobiva poruke od radnika
 - o Prerasporedi zadatke svih radnika
 - o Ukloni radnika

Ograničenja

- Opterećenja sustava i dostupnosti dodijeljenog stroja
- Problem sporih radnika

Alternativna raspodijela

- Pravedan raspoređivač
- Raspoređivač po kapacitetu