

HADOOP

È l'implementazione dello Map Reduce sviluppata da Apache. Deve essere fornito un codice scritto in Java/C++, dove esserci un codice che prende un certo tipo T_{in} in input e che restituisce dopo le fase di reduce il tipo T_{out} in output.

↳ $\langle \text{key}, \text{value} \rangle$

L'obiettivo che ci si pone con Hadoop è suddividere i dati su cui vogliamo lavorare in tante macchine per far in modo di parallelizzare l'accesso ai dati che è la principale bottleneck di un sistema parallelo di questo tipo perché con il tempo performance e velocità di connessione sono migliorate ma l'accesso ai dati no.

Si usa un file system chiamato HDFS che ci permette di avere una replica dei dati in modo che sia garantito sempre l'accesso ai dati ma deve anche garantire la distribuzione di questi dati.

Flume è una possibile ottimizzazione che permette di avere in Hadoop e mi permette di unire due map in una sola map.

STORM

Qui si lavora su uno stream ^{in parallelo} e non su dati che sono già sulla nostra macchina.

Storm ci permette di avere una buona performance nell'esecuzione delle funzioni sullo stream ma è particolarmente adatto solamente per lavori complicati, per le cose semplici non è troppo adatto.

Storm organizza la computazione in una sorta di grafo in cui abbiamo dei nodi di vario tipo:

- Spouts: sono nodi che hanno solo l'output e producono lo stream
- Bolts: nodi che prendono un input, fanno qualcosa e poi producono un output sempre come stream.
- Topologies: sono un'unione di nodi Spouts e di nodi bolts.

Il grafo di Storm ci permette di garantire la fault tolerance perché siamo sicuri che quando viene prodotto uno stream i dati arriveranno fino alla fine dello stream (nodo finale).

SPARK

È una versione più recente di Storm ma permette più o meno di fare la stessa cosa operando comunque su stream.

Permette di fare varie cose:

- SparkSQL ci fa lavorare in un db distribuito usando SQL
- Stream: riusciamo a lavorare su uno stream con delle funzioni più recenti e migliori rispetto a quelli di Storm
- MLlib
- Graphx

In particolare Spark ci fa lavorare sugli stream e ogni volta che arriva una coppia su cui lavorare la memorizziamo anche in un database.

Ci sono vari modi per aggiornare questo database:

- Possiamo rivederlo ogni volta che arrivano nuovi dati
- Possiamo fare un append
- Possiamo fare un update