

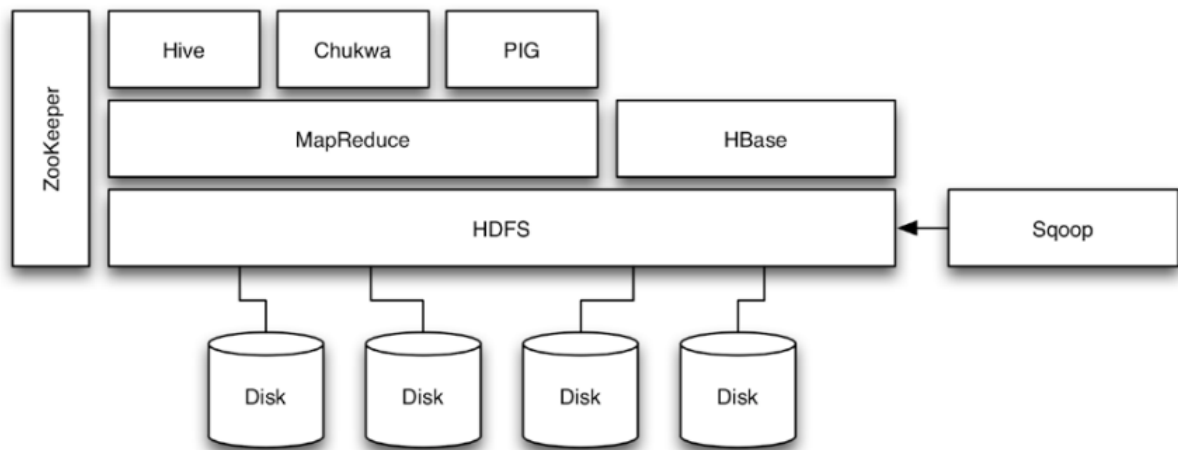


Apache Hadoop

Hădărean Andreea Roxana

Grupa 30242

Apache Hadoop reprezintă un framework open source Java ce are ca principal rol implementarea rețelelor scalabile și de încredere. Acesta include unele subproiecte, precum MapReduce, Pig, ZooKeeper, Hbase și Hive.



Figură 1: Componentele Hadoop [1]

Acest framework se compune din instrumente și funcționalități ce permit serializarea datelor, accesul la sistemul de fișiere și comunicarea între procese. Există două tipuri de configurații: singulare și de tip cluster. Acestea includ în marea majoritate a timpului HDFS, optimizat pentru un debit mare de date de intrare. Se remarcă existența a două componente esențiale: HDFS și MapReduce.

Rolul componentei HDFS este acela de a oferi acces rapid și scalabil la nivelul sistemului distribuit de fișiere. HDFS se bazează pe consistență, o tranzacție neconsiderându-se completă până în momentul în care datele nu sunt scrise în cel puțin două volume configurabile. De asemenea, componenta înglobează toate discurile fizice sau sisteme de fișiere. MapReduce are ca principal rol management-ul nodurilor și



urmarirea modului in care se executa fiecare task. De asemenea, acesta stocheaza toate metodele de mapare si reducere scrise in Java.

Framework-ul Apache Hadoop inglobeaza de asemenea si o alta serie de framework-uri: Chukwa, Hive si HBase. Chukwa are ca principa rol colectarea datelor pentru a monitoriza, a afisa si a analiza inregistrările de la nivelul sistemelor distribuite de dimensiuni mari. Hive se ocupa cu stocarea datelor, permitand de asemenea stocarea lor si transformarea folosind un limbaj de interogare similar cu SQL. Hbase reprezinta baza de date NoSQL care stocheaza in timp real, extrage si permite cautarea in tabele de dimensiuni mari, ruland deasupra HDFS.

Totodata, framework-ul pune la dispozitie si un set de utilitare. Pig reprezinta un set de instrumente pentru a programa procesele de analiza a datelor, oferind un limbaj pentru programare, transformarea datelor si paralelizarea proceselor. Sqoop permite importarea si exportarea datelor din bazele de date relationale in Hadoop sau Hive sau invers. Pentru configurarea, sincronizarea evenimentelor si organizarea nodurilor intr-o retea construita folosind Hadoop este utilizat ZooKeeper.

Pentru sistemele distribuite construite cu ajutorul framework-ului Hadoop se poate realiza deploy-ul in trei moduri: local, pentru implementare si testare, bazat pe thread-uri pe o singura masina virtuala, pseudo-distribuit, pe mai multe masini virtuale, existand un singur nod si distribuit, pentru productie, existand mai multe noduri. Fiecare nod poate executa fie cod MapReduce, fie instructiuni HDFS, executat de un proces ce ruleaza in fundal.

Procesele ce ruleaza in fundal se executa pe trei nivele diferite a sistemului distribuit: master, slave si user application. La nivelul Master sau Name Node, se executa procesul de management al sistemului de fisiere, se monitorizeaza executia task-urilor , se replica blocuri de date si se mentin uniform distribuite. Aceasta componenta poate rula si pe un alt calculator. La al doilea nivel, Slave sau Data Nodes, sunt stocate blocurile de date in sistemul local de fisiere, sunt stocate meta-date, sunt transmise aceste date catre Master. Ultimul nivel, user application, executa implementarea contractelor pentru Java, ofera parametrii specifici executiei aplicatiei si seteaza configuratia parametrilor ce sunt folositi in celelalte doua nivele.



O aplicație construită cu ajutorul Hadoop constă într-unul sau mai multe job-uri. Un job conține un fișier de configurare și una sau mai multe clase Java. Datele trebuie să existe în HDFS. Date de intrare pot avea mai multe formate, precum `KeyValueTextInputFormat` sau `MultiFileInputFormat`. Datele de ieșire se mapează pe formatul datelor de intrare.

În concluzie, Hadoop reprezintă un framework scalabil, flexibil, rapid și tolerant la eșec, util când se dorește stocarea și procesarea datelor. Eficiența acestuia a fost dovedită și prin faptul că este folosit de companii renumite, precum Google și Facebook.

Bibliografie

- [1] E. Ciurana and M. Kalali, "Getting Started with Apache Hadoop," *DZone Refcardz*, pp. 1-6, 2018.