\* **Apache Pig :→**

        Apache Pig is an abstraction over MapReduce. It is a tool / platform which is used to analyze larger sets of data representing them as data flows. Pig is generally used with Hadoop. We can perform all the data manipulation operations in Hadoop using Apache Pig.

To write data analysis Programs, Pig provides a high - Level language known as Pig Latin. This language provides Various operators using which programmers can develop their own functions for reading, writing and processing data.

\* **features of pig :→**

1. **Rich set of operators :→**

        It provides many operators to perform operations like Join, Sort, filter etc

2. **Easy of Programming :→**

        Pig Latin is similar to SQL and it is easy to write a Pig Script if we are good at SQL.

3. **optimization opportunities :→**

        The tasks in Apache Pig optimize their execution automatically, so the programmers need to focus only on semantics of the language.

4. **Extensibility :→**

        Using the existing operators, users can develop their own functions to read, Process and write data.

**5. Handles all kinds of data :→**

Apache Pig analyzes all kinds of data, both structured as well as unstructured. It stores the results in HDFS.

* **Admiring the Pig Architecture :→**

Pig is made up of two (count 'em, two)

Components :

**1. The language itself :→**

As proof that programmers have a sense of humor the programming language for Pig is known as Pig Latin, a high-level language that allows you to write data processing and analysis programs.

**2. The Pig Latin Compiler :→**

The Pig Latin Compiler Converts the Pig Latin Code into executable Code. The executable Code is either in the form of MapReduce Jobs or it can Spawn a process where a Virtual Hadoop instance is Created to run the Pig Code on a Single node.

The sequence of MapReduce Program enable Pig programs to do data processing and analysis in Parallel, leveraging Hadoop MapReduce and HDFS. Running the Pig Job in Virtual Hadoop instance is a useful strategy for testing out our Pig Scripts.

```
┌─────────────────────────────────────────────────────────────┐
│  Pig                          ┌──────────────────────────┐   │
│                               │   Pig Latin Compiler     │   │
│                               └──────────────────────────┘   │
├─────────────────────────────────────────────────────────────┤
│  Processing framework  [MapReduce V2]  [ Tez ]               │
│                                        ┌────────────────┐    │
│                                        │ MapReduce VI   │    │
│                                        └────────────────┘    │
├─────────────────────────────────────────────────────────────┤
│  Resource Management   [      YARN      ]                     │
├─────────────────────────────────────────────────────────────┤
│  Distributed Storage   [        HDFS        ]                 │
└─────────────────────────────────────────────────────────────┘
```
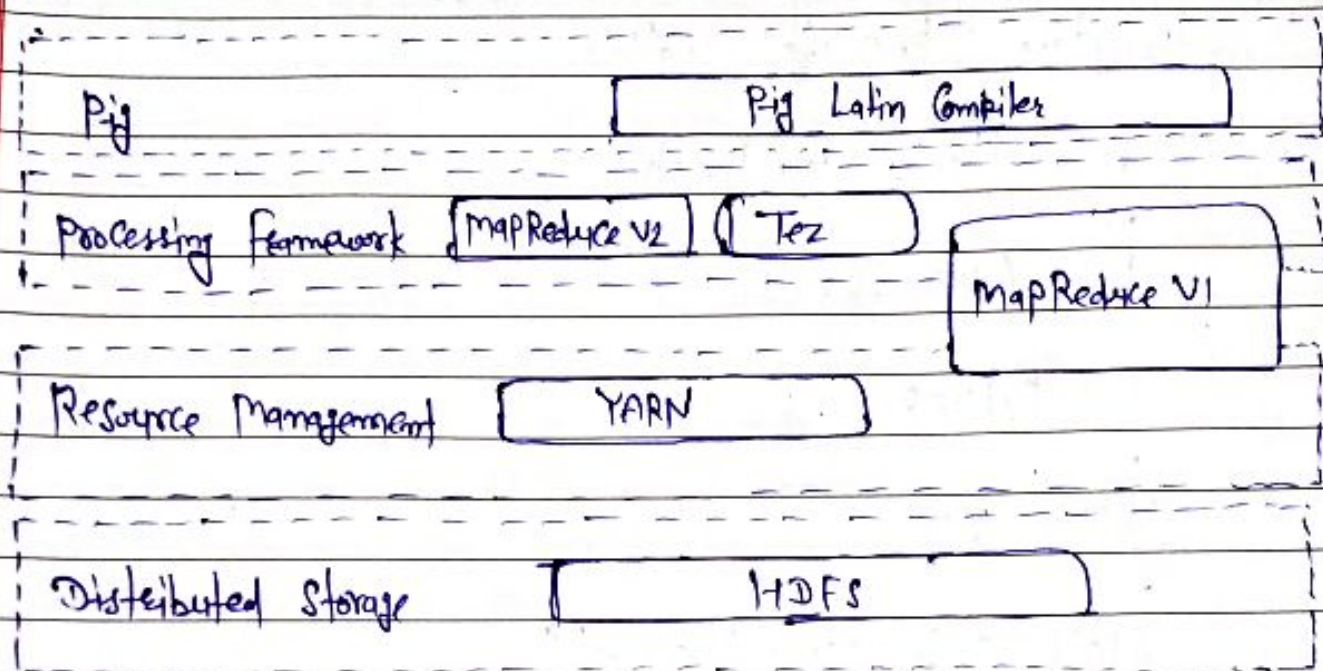
## Figure 1: Pig Architecture

Pig Program Can run on MapReduce 1 or MapReduce 2 without any Code changes, regardless of what mode your cluster is running. However, Pig Scripts Can also run using the Tez API instead.

Going with the Pig Latin Application flow :→
At its Core, Pig Latin is a dataflow language, where you define a data stream and a series of transformations that are applied to the data as it flows through our application. This is in Contrast to a Control flow language ( like c or Java), where we write a series of instructions. In Control flow language, we use Constructs like loop and Conditional logic ( like an if statement) we won't find loops and if statements in Pig Latin.

If we need some convincing that working with pig is a significantly easier than having to write map and reduce programs, start by taking a look at some real pig syntax.
The following listing specifies sample pig code to illustrate the data processing dataflow.

```
A = LOAD 'data-file.txt';
...
B = GROUP ... ;

C = FILTER ...;
...
DUMP B;
...
STORE C INTO 'Results';
```

* working through the ABCs of pig Latin :→

Pig Latin is the language for pig programs. Pig translates the pig Latin script into MapReduce jobs that can be executed within Hadoop cluster. When coming up with pig Latin, the development team followed three key design principles:-

(1) keep it simple :→

Pig latin provides a streamlined method for interacting with Java MapReduce. It's an abstraction, in other words, that simplifiers the creation of parallel programs on the Hadoop cluster for data flows and analysis.
Complex tasks may require a series of interrelated data transformations - such series are encoded as data flow sequences.

writing data transformation and flows as Pig Latin scripts instead of Java MapReduce programs make these programs easier to write, understand and maintain.

(2) Make it Smart :→

You may recall that the Pig Latin Compiler does the work of transforming a Pig Latin program into a series of Java MapReduce Jobs. The trick is to make sure that the Compiler can optimize the execution of these Java MapReduce Jobs automatically, allowing the user to focus on semantics rather than on how to optimize and access the dataset.

(3) Don't limit development :→

Make the Pig extensible so that developers can add functions to address their particular business problems.

* Evaluating Local and Distributed modes of Running Pig Scripts :→

Before you can run your first Pig script, you need to have a handle on how Pig programs can be packaged with the Pig Server.

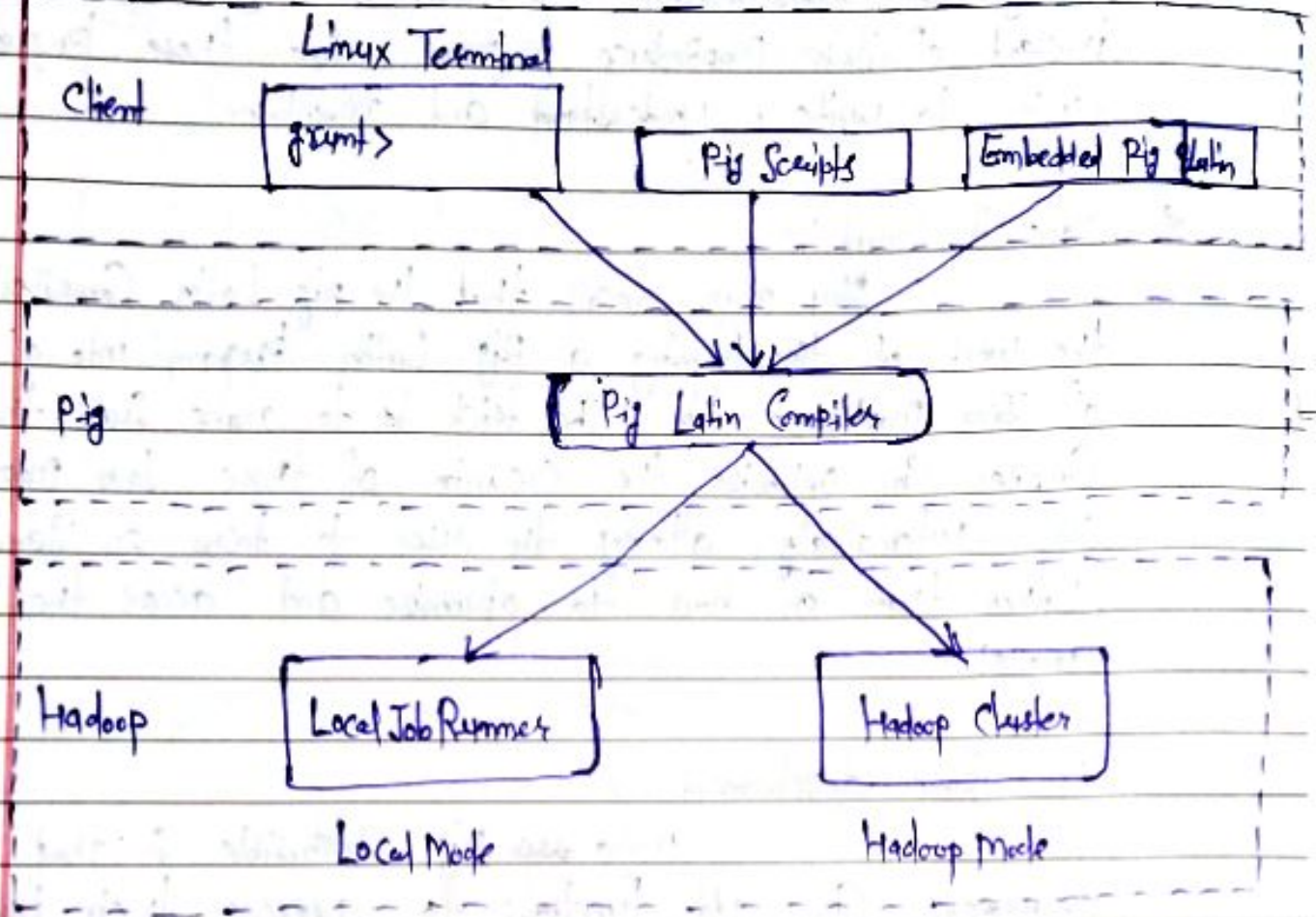Pig has two modes for running Script, as shown in figure :-

Figure : Pig Modes

## Local mode :→

All scripts are run on a single machine without requiring Hadoop MapReduce and HDFS. This can be useful for developing and testing Pig logic. If you're using a small set of data to develope or test your code, then local mode could be faster then going through the MapReduce infrastructure.

Local mode doesn't require Hadoop. When you run in Local mode, the pig program runs in the context of a local Java Virtual machine, and data access is via the local file system of a single machine. Local mode is actually

a local Simulation of MapReduce in Hadoop's LocalJobRunner class.

## MapReduce Mode / Hadoop Mode :→

MapReduce Mode is also known as Hadoop Mode. Pig is executed on the Hadoop cluster. In this case, the Pig Script gets Converted into a Series of MapReduce Jobs that are then run on the Hadoop Cluster. If you have a teradayte of data that you want to perform operations on and you want to interactively develop a program, you may soon find things slowing down Considerably and you may start growing your storage.

✱   Checking out the Pig Script Interfaces :→

The Pig Programming language is designed to handle any kind of data tossed its ways - Structured, Semi-structured, unstructured data, Pig Programs Can be packaged in three different ways :-

(1) Script :→

This method is nothing more than a file Containing Pig Latin Commands, identified by the Pig suffix (FlightData.pig). Ending your Pig program with the .pig extension is a Convention but not required. The Commands are interpreted by the Pig Latin Compiler and executed in the order determined by the Pig optimizer.

## (2) Grunt :→

Grunt Grunt acts as a Command interpreter where we can interactively enter Pig Latin at the Grunt Command line and immediately see the response. This method is helpful for prototyping during initial development and with what-If scenarios.

## (3) Embedded :→

Pig Latin statements can be executed within Java, Python or JavaScript programs.

Pig scripts, Grunt shell Pig Commands, and embedded Pig programs can run in either Local mode or MapReduce mode. The Grunt shell provides an interactive shell to submit Pig Commands or run pig scripts.

To start the Grunt shell shell is executed is in Interactive mode, Just Submit the Command Pig at your shell. To specify whether a script or Grunt shell is executed locally or in Hadoop mode Just specify it in the -X flag to the Pig Command.

The following is an example of how you'd specify running your pig Script in local mode.

- Pig -X local miles Per Carrier. Pig

Here's how you'd run the Pig script in Hadoop mode, which is the default if you don't specify the flag.

- Pig -X mapreduce miles Per Carrier. Pig

By default, when we specify the Pig Command without any parameters, it starts the Grunt shell in Hadoop mode. If we want to start the Grunt shell in local mode just add the -x local flag to the Command.
Here is an example :-

Pig -x local


**\* Scripting with Pig Latin :→**

Hadoop is a rich and quickly evolving ecosystem with a growing set of new applications. Rather than try to keep up with all the requirements for new capabilities, Pig is designed to be extensible via user-defined functions, also known as UDFs. UDFs Can be written in a number of programming languages, including Java, Python and Javascript. Developers are also posting and sharing a growing collection of UDFs online.

( Look for Piggy Bank and DataFu, to name just two examples of such online collections)
Some of the Pig UDFs that are part of these repositories are LOAD/STORE function (XML, for example), date time, functions, text, math, and stats functions.


Pig Can also be embedded in host languages such as Java, Python and Javascript, which allows you to integrate Pig with your existing applications.

# * HDFS Architecture :→

HDFS holds very large amount of data and provides easier access. To store such huge data, the files are stored across multiple machines. These files are stored in redundant fashion to rescure the system from possible data losses in case of failure. HDFS also makes applications available to parallel processing.
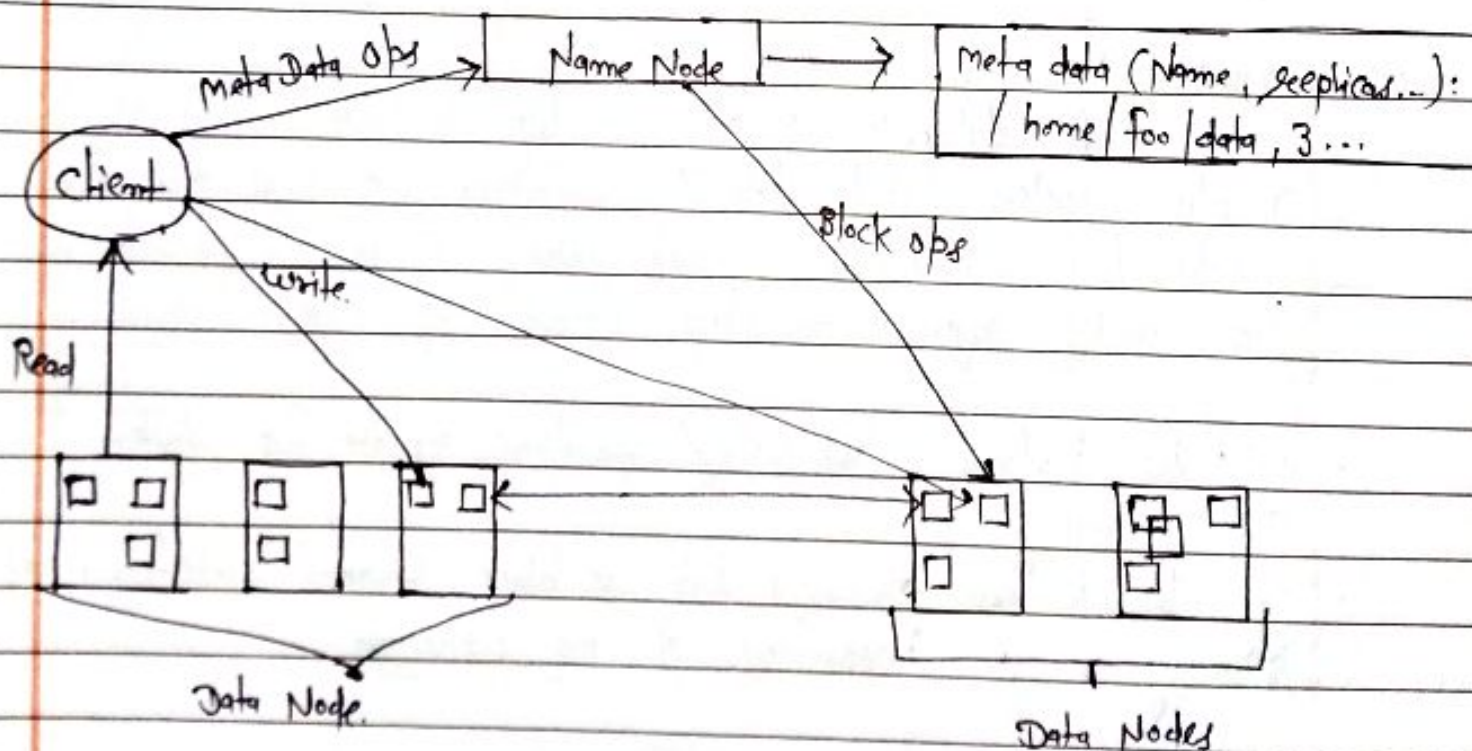
## HDFS Architecture.



figure: HDFS Architecture.

→ **Name Node :→**

The NameNode is the Commodity hardware that Contains the GNU/LINUX Operating System and the NameNode Software. It is a Software that can be run on Commodity hardware. The system having the NameNode acts as the master Server and it does the following tasks —

- Manages the files
- Store Meta Data
- It also executes file System operations Such as renaming, closing, and opening files and directories.

→ **Data Node :→**

The datanode is a Commodity hardware having the GNU/Linux Operating System and dataNode Software. for every Node (Commodity hardware/System) in a Cluster, there will be a datanode. These nodes manage the data Storage of their System.

- DataNodes perform read-write operations on the file System, as per Client request.
- They also perform operations such as block Creation, deletion, and replication according to the Instructions of the NameNode.

→ **Name Node :→**

The NameNode is the Commodity hardware that Contains the GNU/LINUX operating system and the NameNode Software. It is a Software that can be run in Commodity hardware. The System having the NameNode acts as the master Server and it does the following tasks —

- Manages the files
- Store Meta Data
- It also executes file System operations Such as renaming, closing, and opening files and directories.

→ **Data Node :→**

The datanode is a Commodity hardware having the GNU/Linux operating System and dataNode Software. For every Node (Commodity hardware/System) In a Cluster, there will be a datanode. These nodes manage the data Storage of their System.

- DataNodes perform read-write operations on the file System, as per Client request.
- They also perform operations Such as block Creation, deletion, and replication according to the Instructions of the NameNode.