

# Apache Pig

**Apache Pig**<sup>[1]</sup> is a high-level platform for creating programs that run on Apache Hadoop. The language for this platform is called **Pig Latin**.<sup>[1]</sup> Pig can execute its Hadoop jobs in MapReduce, Apache Tez, or Apache Spark.<sup>[2]</sup> Pig Latin abstracts the programming from the Java MapReduce idiom into a notation which makes MapReduce programming high level, similar to that of SQL for relational database management systems. Pig Latin can be extended using user-defined functions (UDFs) which the user can write in Java, Python, JavaScript, Ruby or Groovy<sup>[3]</sup> and then call directly from the language.

## Contents

**History**

**Example**

**Pig vs SQL**

**See also**

**References**

**External links**

## History

Apache Pig was originally<sup>[4]</sup> developed at Yahoo Research around 2006 for researchers to have an ad hoc way of creating and executing MapReduce jobs on very large data sets. In 2007,<sup>[5]</sup> it was moved into the Apache Software Foundation.

### Apache Pig



## Apache Pig

<b><u>Developer(s)</u></b>	<u>Apache Software Foundation</u> , <u>Yahoo Research</u>
<b><u>Initial release</u></b>	September 11, 2008
<b><u>Stable release</u></b>	0.17.0 / June 19, 2017
<b><u>Repository</u></b>	<u>svn.apache.org</u> <u>/repos/asf/pig/</u> ( <u>http://svn.apache.org/repos/asf/pig/</u> )
<b><u>Operating system</u></b>	<u>Microsoft Windows</u> , <u>OS X</u> , <u>Linux</u>
<b><u>Type</u></b>	Data analytics
<b><u>License</u></b>	<u>Apache License 2.0</u>
<b><u>Website</u></b>	<u>pig.apache.org</u> ( <u>https://pig.apache.org</u> )

Version	Original release date	Latest version	Release date <sup>[6]</sup>
0.1	2008-09-11	0.1.1	2008-12-05
0.2	2009-04-08	0.2.0	2009-04-08
0.3	2009-06-25	0.3.0	2009-06-25
0.4	2009-08-29	0.4.0	2009-08-29
0.5	2009-09-29	0.5.0	2009-09-29
0.6	2010-03-01	0.6.0	2010-03-01
0.7	2010-05-13	0.7.0	2010-05-13
0.8	2010-12-17	0.8.1	2011-04-24
0.9	2011-07-29	0.9.2	2012-01-22
0.10	2012-01-22	0.10.1	2012-04-25
0.11	2013-02-21	0.11.1	2013-04-01
0.12	2013-10-14	0.12.1	2014-04-14
0.13	2014-07-04	0.13.0	2014-07-04
0.14	2014-11-20	0.14.0	2014-11-20
0.15	2015-06-06	0.15.0	2015-06-06
0.16	2016-06-08	0.16.0	2016-06-08
<b>0.17</b>	2017-06-19	0.17.0	2017-06-19
<b>Legend:</b> <span style="display: inline-block; width: 15px; height: 15px; background-color: #f08080; margin-right: 5px;"></span> Old version <span style="display: inline-block; width: 15px; height: 15px; background-color: #ffff00; margin-right: 5px; margin-left: 10px;"></span> Older version, still maintained <span style="display: inline-block; width: 15px; height: 15px; background-color: #90ee90; margin-right: 5px; margin-left: 10px;"></span> Latest version <span style="display: inline-block; width: 15px; height: 15px; background-color: #ffcc99; margin-left: 10px;"></span> Latest preview version			

# Example

Below is an example of a "Word Count" program in Pig Latin:

```

input_lines = LOAD '/tmp/my-copy-of-all-pages-on-internet' AS (line:chararray);

-- Extract words from each line and put them into a pig bag
-- datatype, then flatten the bag to get one word on each row
words = FOREACH input_lines GENERATE FLATTEN(TOKENIZE(line)) AS word;

-- filter out any words that are just white spaces
filtered_words = FILTER words BY word MATCHES '\\w+';

-- create a group for each word
word_groups = GROUP filtered_words BY word;

-- count the entries in each group
word_count = FOREACH word_groups GENERATE COUNT(filtered_words) AS count, group AS word;

-- order the records by count
ordered_word_count = ORDER word_count BY count DESC;
STORE ordered_word_count INTO '/tmp/number-of-words-on-internet';

```

The above program will generate parallel executable tasks which can be distributed across multiple machines in a Hadoop cluster to count the number of words in a dataset such as all the webpages on the internet.

# Pig vs SQL

In comparison to SQL, Pig

1. has a nested relational model,
2. uses lazy evaluation,
3. uses extract, transform, load (ETL),
4. is able to store data at any point during a pipeline,
5. declares execution plans,
6. supports pipeline splits, thus allowing workflows to proceed along DAGs instead of strictly sequential pipelines.

On the other hand, it has been argued DBMSs are substantially faster than the MapReduce system once the data is loaded, but that loading the data takes considerably longer in the database systems. It has also been argued RDBMSs offer out of the box support for column-storage, working with compressed data, indexes for efficient random data access, and transaction-level fault tolerance.<sup>[7]</sup>

Pig Latin is procedural and fits very naturally in the pipeline paradigm while SQL is instead declarative. In SQL users can specify that data from two tables must be joined, but not what join implementation to use (You can specify the implementation of JOIN in SQL, thus "... for many SQL applications the query writer may not have enough knowledge of the data or enough expertise to specify an appropriate join algorithm."). Pig Latin allows users to specify an implementation or aspects of an implementation to be used in executing a script in several ways.<sup>[8]</sup> In effect, Pig Latin programming is similar to specifying a query execution plan, making it easier for programmers to explicitly control the flow of their data processing task.<sup>[9]</sup>

SQL is oriented around queries that produce a single result. SQL handles trees naturally, but has no built in mechanism for splitting a data processing stream and applying different operators to each sub-stream. Pig Latin script describes a directed acyclic graph (DAG) rather than a pipeline.<sup>[8]</sup>

Pig Latin's ability to include user code at any point in the pipeline is useful for pipeline development. If SQL is used, data must first be imported into the database, and then the cleansing and transformation process can begin.<sup>[8]</sup>

## See also

---

- Apache Hive
- Sawzall — similar tool from Google

## References

---

1. "Hadoop: Apache Pig" (<http://pig.apache.org/>). Retrieved Sep 2, 2011.
2. "[PIG-4167] Initial implementation of Pig on Spark - ASF JIRA" (<https://issues.apache.org/jira/browse/PIG-4167>). *issues.apache.org*. Retrieved 2018-12-29.
3. "Pig user defined functions" (<http://pig.apache.org/docs/r0.11.1/udf.html>). Retrieved May 3, 2013.
4. "Yahoo Blog:Pig – The Road to an Efficient High-level language for Hadoop" (<https://web.archive.org/web/20160203181220/https://developer.yahoo.com/blogs/hadoop/pig-road-efficient-high-level-language-hadoop-413.html>). Archived from the original (<https://developer.yahoo.com/blogs/hadoop/pig-road-efficient-high-level-language-hadoop-413.html>) on February 3, 2016. Retrieved May 23, 2015.
5. "Pig into Incubation at the Apache Software Foundation" (<https://web.archive.org/web/20160203162733/https://developer.yahoo.com/blogs/hadoop/pig-incubation-apache-software-foundation-393.html>). Archived from the original (<https://developer.yahoo.com/blogs/hadoop/pig-incubation-apache-software-foundation-393.html>) on February 3, 2016. Retrieved May 23, 2015.

6. "Apache Pig Releases" (<https://pig.apache.org/releases.html>). *Apache*. Retrieved 2019-03-13.
7. "Communications of the ACM: MapReduce and Parallel DBMSs: Friends or Foes?" (<https://web.archive.org/web/20150701205317/http://database.cs.brown.edu/papers/stonebraker-cacm2010.pdf>) (PDF). Archived from the original (<http://database.cs.brown.edu/papers/stonebraker-cacm2010.pdf>) (PDF) on July 1, 2015. Retrieved May 23, 2015.
8. "Yahoo Pig Development Team: Comparing Pig Latin and SQL for Constructing Data Processing Pipelines" (<https://web.archive.org/web/20150530103839/https://developer.yahoo.com/blogs/hadoop/comparing-pig-latin-sql-constructing-data-processing-pipelines-444.html>). Archived from the original (<https://developer.yahoo.com/blogs/hadoop/comparing-pig-latin-sql-constructing-data-processing-pipelines-444.html>) on May 30, 2015. Retrieved May 23, 2015.
9. "ACM SigMod 08: Pig Latin: A Not-So-Foreign Language for Data Processing" (<http://infolab.stanford.edu/~olston/publications/sigmod08.pdf>) (PDF). Retrieved May 23, 2015.

## External links

---

- [Official website \(https://pig.apache.org\)](https://pig.apache.org)
- 

Retrieved from "[https://en.wikipedia.org/w/index.php?title=Apache\\_Pig&oldid=1015404722](https://en.wikipedia.org/w/index.php?title=Apache_Pig&oldid=1015404722)"

---

This page was last edited on 1 April 2021, at 08:36 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.