

# Pipeline Architecture

February 28, 2022

Brae Webb

Presented for the Software Architecture course  
at the University of Queensland



THE UNIVERSITY  
OF QUEENSLAND  
AUSTRALIA

# Pipeline Architecture

Software Architecture

February 28, 2022

Brae Webb

## 1 Introduction

The pipeline architecture needs very minimal introduction. Almost every developer will have been exposed to software which implements this architecture. Some notable examples are bash, hadoop, and most functional programming languages. A pipeline architecture consists of modular components which accept input and return output. Data is piped through a sequence of components until the desired output is reached.

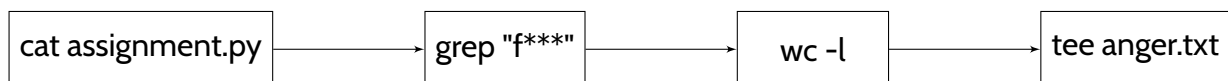


Figure 1: An example of using bash's pipeline architecture to perform statistical analysis.

## 2 Terminology

As illustrated by Figure 2, a pipeline architecture consists of just two components;

**Filters** modular software components, and

**Pipes** the flow of data between filters.



Figure 2: A generic pipeline architecture.

Filters themselves are composed of four major types:

**Producers** Filters where data originates from are called producers, or source filters.

**Transformers** Filters which manipulate input data and output to the outgoing pipe are called transformers.

**Testers** Filters which apply selection to input data, allowing only a subset of input data to progress to the outgoing pipe are called testers.

**Consumers** The final state of a pipeline architecture is a consumer filter, where data is eventually used.

The example in Figure 1 shows how bash's pipeline architecture can be used to manipulate data in unix files. Figure 3 helps to clarify the above terminology by labelling each of the bash processes in the example.

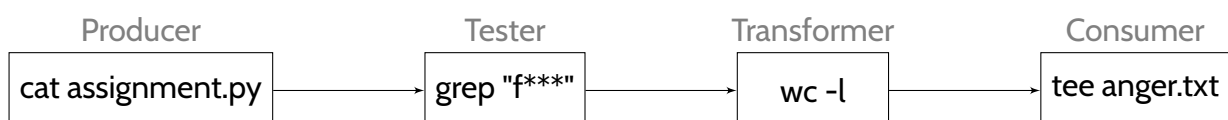


Figure 3: Figure 1 with labelled filter types.

### 3 Case Study: MapReduce

One of the more prevalent uses of the pipeline architecture is the MapReduce pattern. The MapReduce pattern was discovered in 2004 as a solution to the challenges which Google faced managing their search index [1].<sup>1</sup> MapReduce affords impressive parallelism inherent to the programming pattern.

The two key ideas of MapReduce, *map* and *reduce*, come from functional programming.<sup>2</sup> Below are the generic types of the *map* and *reduce* functions in functional programming.

```
1 map : ( $\tau_1 \rightarrow \tau_2$ )  $\rightarrow \tau_1 Seq \rightarrow \tau_2 Seq$ 
2 map f xs
3 reduce : ( $\tau_1 \rightarrow \tau_1 \rightarrow \tau_1$ )  $\rightarrow \tau_1 Seq \rightarrow \tau_1 \rightarrow \tau_1 Seq$ 
4 reduce f xs initial
```

If you're unfamiliar with this notation, the rough English translation is:

**map** The parameters of the *map* function are:

- (a) A function, *f*, which takes a parameter of type  $\tau_1$  and returns a type  $\tau_2$ .
- (b) A sequence of elements of type  $\tau_1$ .

The return type of the *map* function is a sequence of elements of type  $\tau_2$ .

**reduce** The parameters of the *reduce* function are:

- (a) A function, *f*, which takes two parameters both of type  $\tau_1$  and returns a type  $\tau_1$ .
- (b) A sequence of elements of type  $\tau_1$ .
- (c) An initial accumulator value of type  $\tau_1$ .

The return type of the *reduce* function is a sequence of elements of type  $\tau_1$ .

The code snippet below uses the *map* and *reduce* functions to perform the operations of the above bash example. One important thing to note about the example below is the map operation on line 11. Each application of the lambda function within the map operation is completely independent and could, in theory, be executed simultaneously.

```
1 contents = read("assignment.py")
2
3 # filter relevant lines by rebuilding the list
4 contents = reduce( $\lambda$  xs x  $\rightarrow$ 
5                 if x.contains("f***")
6                 then x + xs
7                 else xs,
8                 contents)
9
10 # use map to count occurrences of word
11 contents = map( $\lambda$  line  $\rightarrow$  line.count("f***"), contents)
12
13 # use reduce to sum list of counts
```

<sup>1</sup>Although the pattern was in use prior to their work[2]

<sup>2</sup>I think? Will consult with history textbook (Ian)

```

14 contents = reduce( $\lambda$  total curr  $\rightarrow$  total + curr, contents, 0)
16 write("anger.txt", contents)

```

So by design, code written in this pattern can process data simultaneously. Tools such as [Hadoop](#)<sup>3</sup> are able to take advantage of this to distribute computation automatically.

Using the terminology of a pipeline architecture, what filters do the *map* and *reduce* operators correspond to?

How would you improve the efficiency of the code snippet above?

## 4 Case Study: Compilers

An interesting case study of the pipeline architecture is a compiler.<sup>4</sup> As a foundational technology, compilers have undergone rigorous refinement and are perhaps the most well studied type of software. Modern compilers have well-defined modular phases as illustrated by Figure 4, each phase of a compiler transforms the representation of the program until the target program is produced.

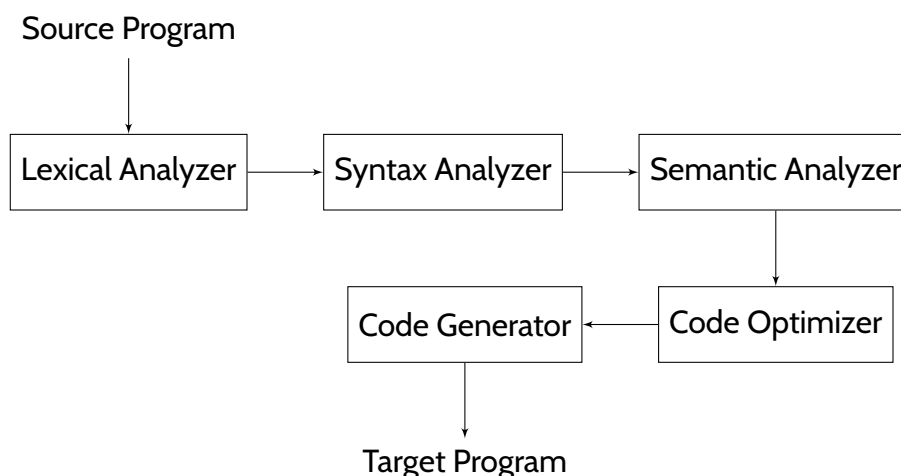


Figure 4: Typical phases of a compiler.

However, a compiler is not well suited to use a pipeline architecture. In general, the modules of a pipeline architecture should be independent of their input source. This is not the case in compilers, as each phase relies on the completion of the previous phase. As a result, the input dependencies of a compiler make it too restrictive for a true pipeline architecture.

Instead, compilers are often built as a hybrid of a pipeline architecture and the *Blackboard Architecture*. The blackboard architecture consists of;

- a knowledge base, the ‘blackboard’,

<sup>3</sup><https://hadoop.apache.org/>

<sup>4</sup>You don’t need to understand the phases of a compiler — two data structures, the Symbol Table and AST, are transformed in each compiler phase.

- knowledge sources which use and update the knowledge base, and
- a control component to coordinate the operation of knowledge sources.

In modern compilers, the data which would be passed through pipes, the Symbol Table and AST, are used and updated by each phase. They are subsequently used as the 'blackboard'. Each phase is considered a knowledge source which uses the knowledge base to transform and update the knowledge base. Finally, in this hybrid, the control component is not required as the sequence of phase execution in a pipeline coordinates operation. Figure 5 illustrates this proposed architecture. Of course, there are many compilers out there, many of them deviate from this architectural hybrid.

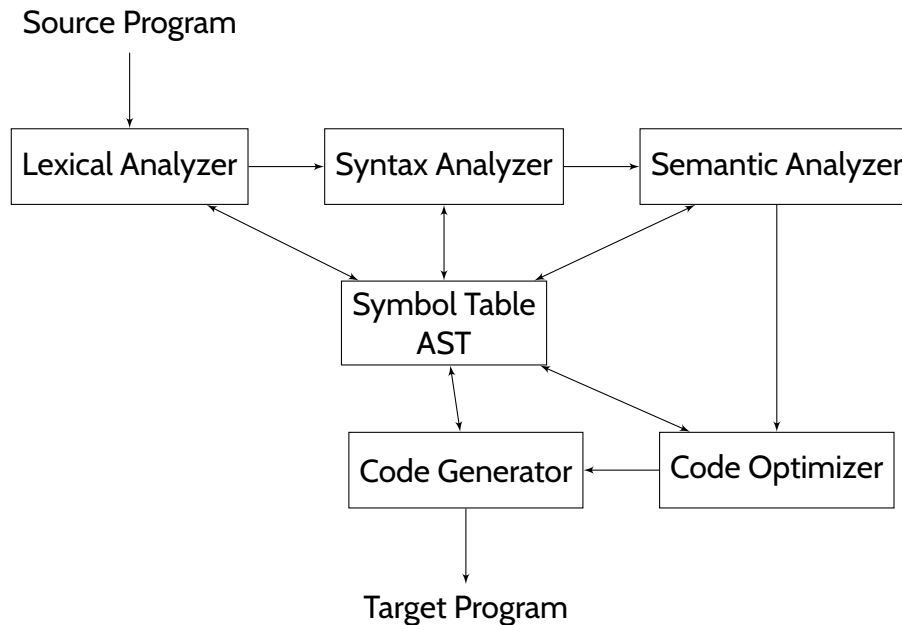


Figure 5: Modern phases of a compiler.

## References

- [1] J. Dean and S. Ghemawat, "Mapreduce: Simplified data processing on large clusters," in *OSDI'04: Sixth Symposium on Operating System Design and Implementation*, (San Francisco, CA), pp. 137–150, 2004.
- [2] D. J. DeWitt and M. Stonebraker, "Mapreduce: A major step backwards." <https://dsf.berkeley.edu/cs286/papers/backwards-vertica2008.pdf>, January 2008.