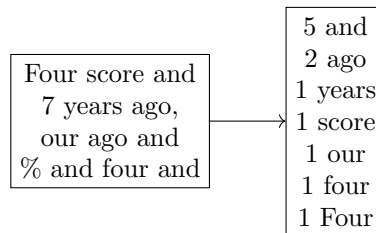We will begin this class with a quiz. We wish to develop a program that takes a sequence of ASCII characters and outputs a frequency-sorted concordance of all words in the input.
We define a word with the following regular expression: [A-Za-z]+
Words must be bounded by white space. An example usage of our program is:

```
                              5 and
                              2 ago
    Four score and            1 years
    7 years ago,         →    1 score
    our ago and               1 our
    % and four and            1 four
                              1 Four
```
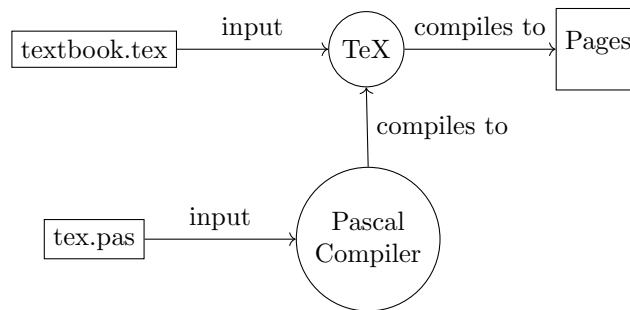
This problem was used by Turing Award winner DE Knuth as press for his book, The Art of Computer Programming
The motivation came from difficulties in creating the initial press for his book;
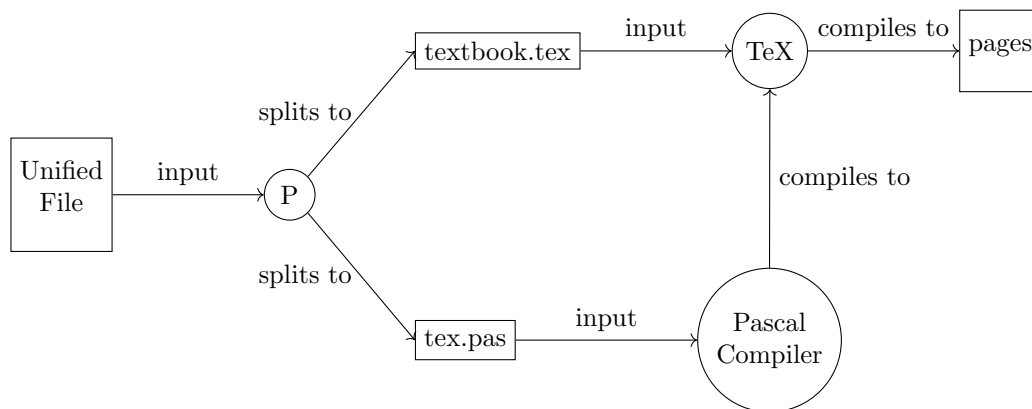    The lead type press could not handle code or equations, which were fundamental!
So Knuth invented a language to produce correct images, called Tex.
He developed a textbook to bring attention to his language, which explained the framework. The framework took the form:

```
                         input            compiles to
    textbook.tex  ──────────────→  TeX  ──────────────→  Pages
                                    ↑
                                    │ compiles to
                                    │
                         input    Pascal
    tex.pas  ──────────────→     Compiler
```

This presents a large (and common) issue – code and documentation must be kept in agreement.
Knuth developed the idea of a unified, interleaved file to combat this.
Thus his program took the following final form:

```
                                            input          compiles to
                     textbook.tex  ──────────────→  TeX  ──────────────→  pages
                   splits to ↗                       ↑
                            /                        │ compiles to
    Unified    input       /                         │
     File  ──────────→  P                            │
                          \                         Pascal
                   splits to ↘          input      Compiler
                             tex.pas  ──────────────→
```

This left him with a nearly 500 page textbook; how would he drum up interest?
Like a good computer scientist, he chose to write a paper. He coined his approach **literate programming**.
This has become standard practice in many disciplines; we can see JavaLibrary for proof.

Knuth approached Doug McIlroy, the manager for the development of UNIX, for help.
McIlroy suggested our problem as an example
Knuth wrote a solution using Hash Tries via literate programming!
Running it through TeX gives pages and through pascal gives the paper!

As impressive as this was, the afterword was what stuck.
McIlroy proposed the following trivial BASH solution (he was the brain behind UNIX pipes after all)
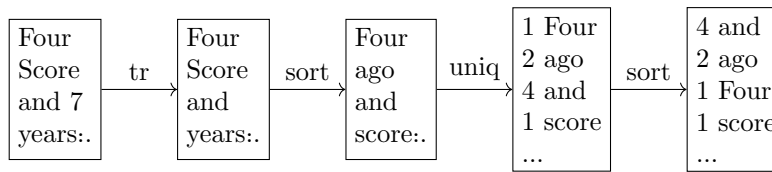
```
        tr -c 'A-Za-z' '\n' | sort | uniq -c | sort -rn
```

The pascal solution was, say, 1000 lines.

   It was faster

   It was more "checked" by compilers

BUT the bash solution is exceedingly simple and documents itself as follows:

| Four Score and 7 years:. | $\xrightarrow{tr}$ | Four Score and years:. | $\xrightarrow{sort}$ | Four ago and score:. | $\xrightarrow{uniq}$ | 1 Four 2 ago 4 and 1 score ... | $\xrightarrow{sort}$ | 4 and 2 ago 1 Four 1 score ... |

This demonstrates an issue fundamental to computer science called **choice of notation**.
There are pros and cons to any choice of notation, and it is up to the engineer to weigh them!
As this is a discussion of languages, consider the following fundamental linguistic hypothesis:

**Sapir-Whorf Hypothesis**.

   a) There is no limit on the structural diversity of language.

   b) The structure of a language determines a native speaker's perception of experience.


This was proposed by linguists for natural languages, but (b) was softened.
Even though it is false, it stuck: see "Eskimos have 19 words for snow".

   Even if false for natural languages, it is true for programming languages!
We explore this in the next lecture