# Design Decisions

The code has been structured to be simple and clear, making it easy for anyone to understand and modify. The function *find_anagrams()* does exactly what it promises, it finds anagrams. This clear structure ensures the code is both intuitive and easy to maintain. Error handling has also been integrated to improve user-friendliness. If the program cannot locate the provided file, it won't crash. Instead, it will display a helpful error message, which makes debugging easier. Additionally, the use of *os.path.join()* ensures that the program works across different operating systems, such as Windows, Linux, or macOS, by handling file paths flexibly.

The program is designed to be scalable, able to handle a large number of words. Sorting each word alphabetically is a straightforward but effective method for grouping anagrams. While sorting may take some time, it is far more efficient than comparing every word with every other word in the list. The use of *defaultdict(list)* helps store the anagram groups efficiently, ensuring that as the number of words grows, the program can continue processing without significant delays. Processing the file line by line, rather than loading the entire file into memory, helps maintain low memory usage, which is especially important when dealing with large files.

Performance is a critical consideration, particularly when handling large input files. Sorting the words, while time-consuming, is the most efficient way to group anagrams. Without sorting, the program would require a slower method for checking if two words are anagrams. The defaultdict allows for efficient group lookups, speeding up the process. The program also uses UTF-8 encoding, ensuring compatibility with various characters and types of input text files without errors.

The *os* library, a built-in Python module, handles file paths effectively, ensuring portability across different operating systems without the need for special adjustments. The use of defaultdict from the collections module simplifies the code by automatically managing key existence checks, making the program cleaner and more efficient. Overall, the design balances simplicity and performance.

# Scalability Considerations

To handle 10 million words, the current solution would still work, but the process would take some time. Sorting each word alphabetically to group anagrams is efficient, but when you reach millions of words, sorting every word one by one can become slow. The memory usage would also increase since we need to store all the words and their corresponding anagram groups in memory. However, because the program processes the file line by line, memory usage would remain manageable, and it could handle 10 million words without overwhelming most modern systems. It would just take a few minutes to process, depending on the hardware.

When dealing with 100 billion words, the solution would need some significant changes. First, storing 100 billion words and their anagram groups in memory would be impractical. The program would quickly run out of memory, so it would need to be restructured to use much less memory. Instead of loading everything into memory at once, we would need to process the data in smaller chunks, storing intermediate results on disk or in a database.

For processing 100 billion words, distributed computing becomes essential. Using a framework like Apache Spark or Dask would allow the data to be split across multiple machines, with each machine processing part of the data in parallel. This would drastically reduce the time needed to process the words. Additionally, instead of sorting every word, we could use a hashing technique to group words by their character frequencies, which would be faster and more memory-efficient.

Another strategy would be to use cloud services. These platforms can scale automatically, so they would handle the large volume of data and the heavy processing load.