



OBJECT ORIENTED PROGRAMMING
DATA STRUCTURES

Extractive Text Summarizer

Evaluating the Effectiveness of HashMap, TreeMap, and
ConcurrentSkipListMap for an Extractive Text Summarizer in Java

16 June 2025

Juan Felix Kusnadi - 2802536386

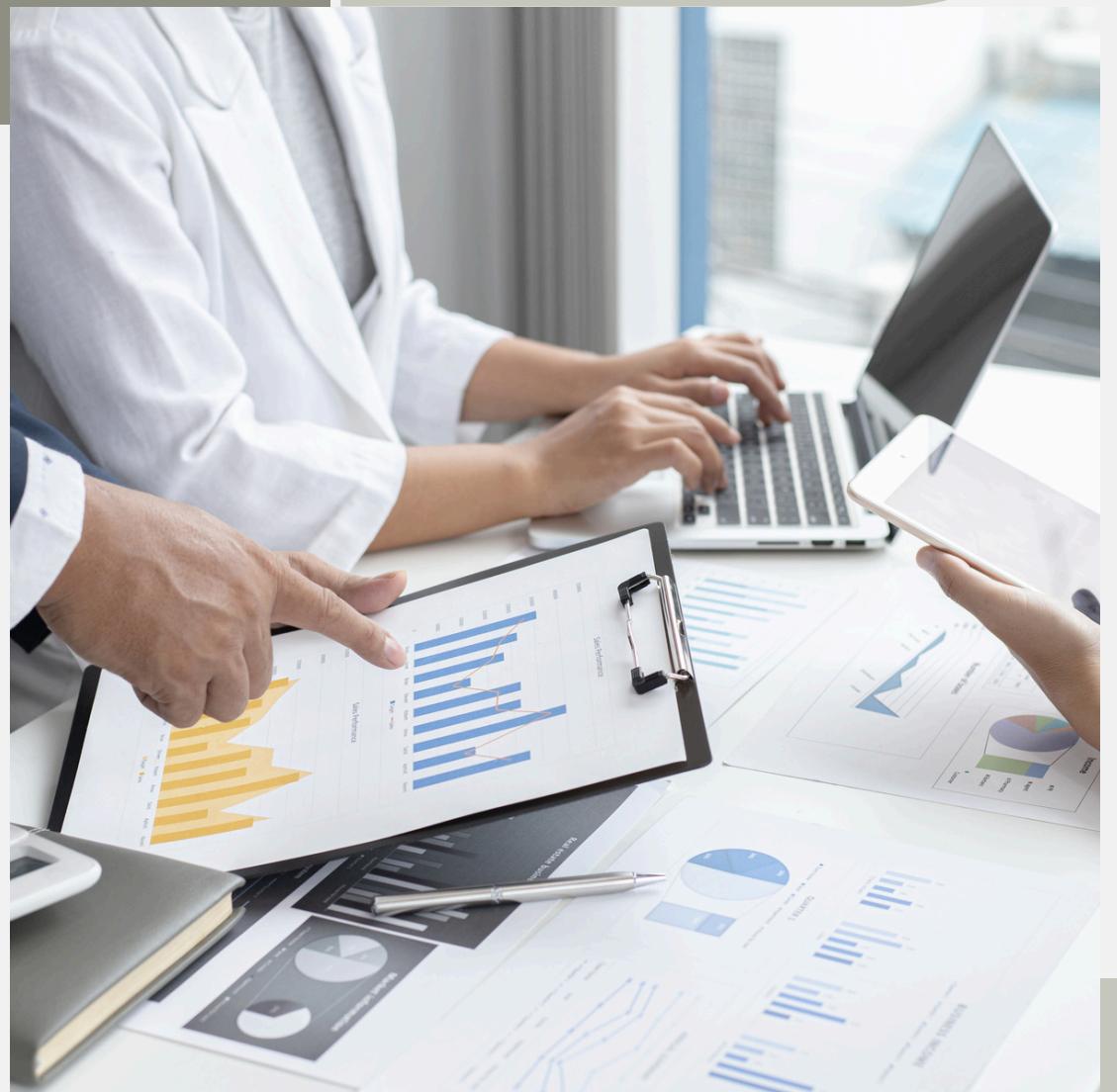
Jonathan Mulyono - 2802537054

Iglesias Sidharta Handjo - 2802530621



Hello!

Text summarization helps people quickly understand lengthy texts by highlighting key information. This project uses the TF-IDF algorithm to build an extractive summarizer in Java. We compare `HashMap`, `TreeMap`, and `ConcurrentSkipListMap` to evaluate their impact on performance and memory usage.



Problem Description

- The explosion of daily textual data makes it challenging for humans to efficiently extract key information
- Summarization tools condense lengthy texts into concise, meaningful overviews without losing essential content
- Extractive summarization selects and compiles the most important sentences
- This study analyzes how different Java Map implementations affect the performance and behavior of a TF-IDF-based extractive summarizer



Approach

Evaluate the importance of a word within a document relative to a collection of documents (corpus) (Jain, 2025)

Concept I

Words that appear more often in a document are more important (higher Term Frequency)

Concept II

Words that are common across different documents in the corpus are less important as they are less useful for distinguishing different documents (lower Inverse Document Frequency)



Methodology

Preparing the text input

Tokenize, remove stop words, convert lowercase, stem

Calculate Term Frequency (TF)

$TF = w/S$, where

w = number of times word w appears

S = number of words in sentence S

Calculate Inverse Document Frequency (IDF)

$IDF = \log(N/SF(w))$, where

N = number of sentences

$SF(w) =$ sentence frequency containing the word w

Creating a summary

Calculate TF-IDF by $TF \times IDF$

Store this key (word) and value (TF-IDF) pair

Rank top sentences and output top 20%-40% (Demilie, 2022)



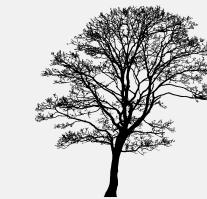
Data Structures

Map Implementations



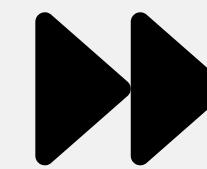
HashMap

Unordered key-value store with hashing



TreeMap

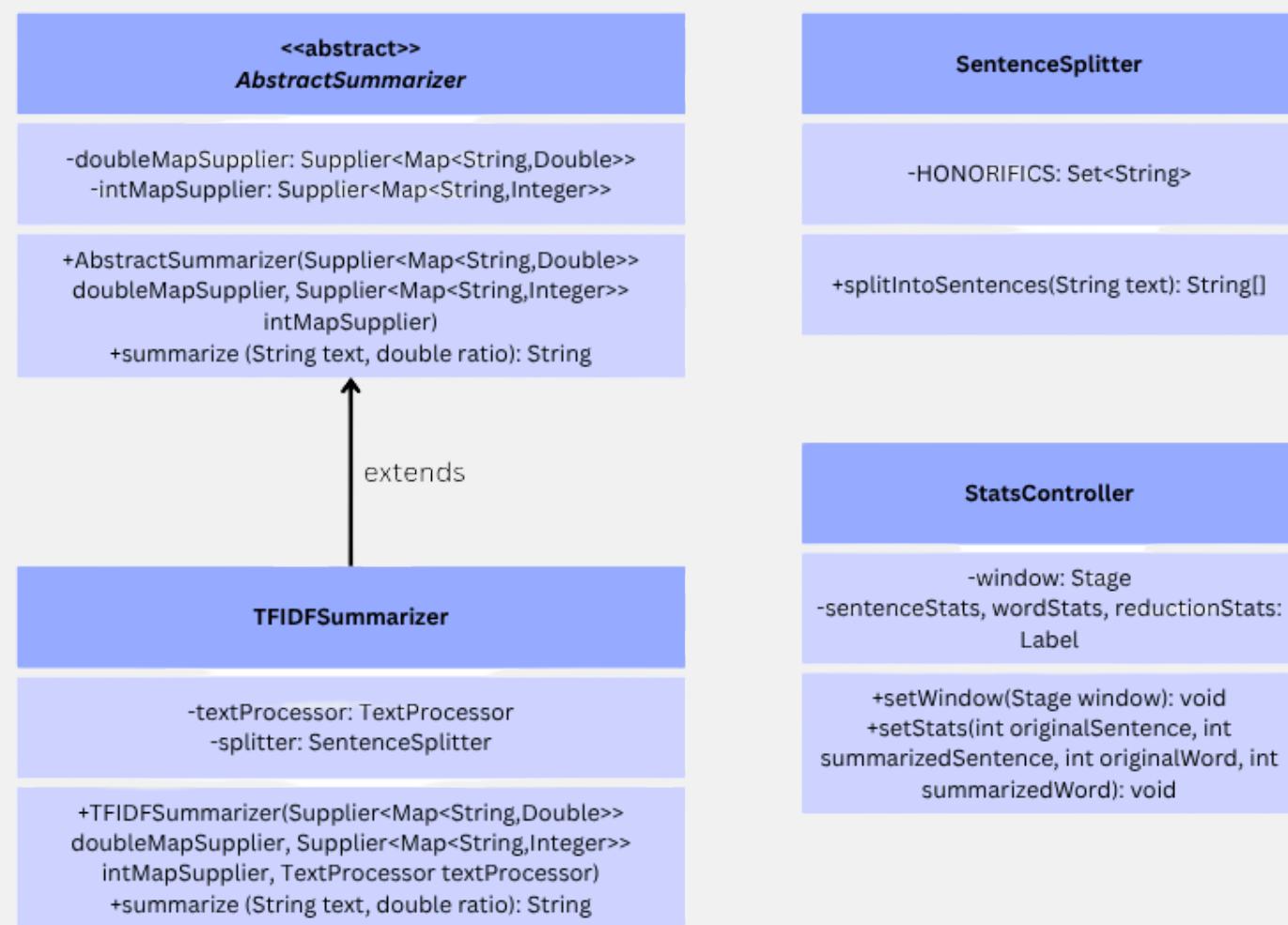
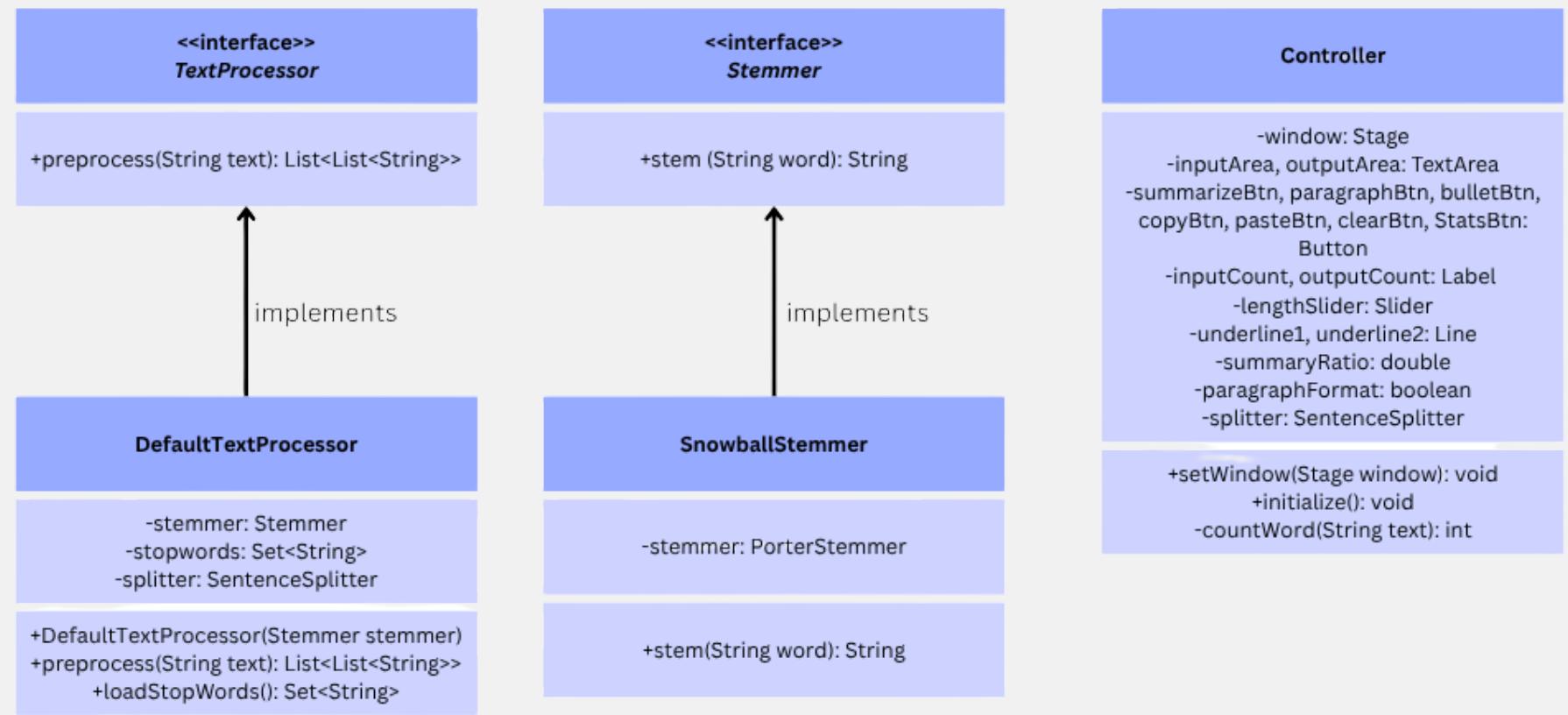
Sorted key-value map using tree



ConcurrentSkipListMap

Concurrent sorted map using skip-list

Program Demo

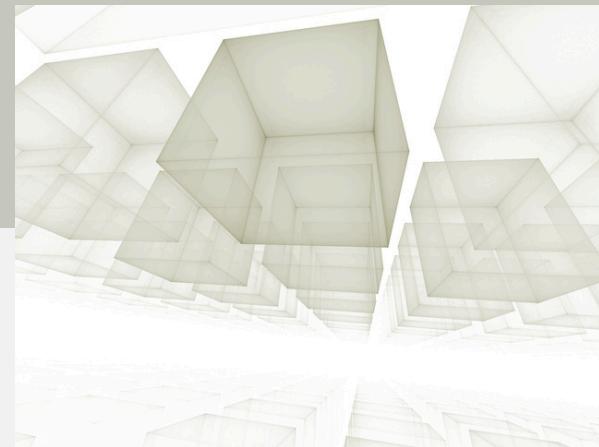


The 4 Principles of OOP



Encapsulation

Classes hide their internal data and helper methods, exposing only defined public methods.



Abstraction

Interfaces and the abstract base class define clear contracts for operations without revealing details.



Inheritance

Reusing common supplier wiring and constructor logic while providing its own summarize() implementation.



Polymorphism

Abstract references allow different summarizer or stemmer implementations to be interchanged at runtime.

Input (AI Generated)

Cybersecurity is essential in today's ... keeping information secure. (100 words)

Artificial intelligence (AI) is revolutionizing the healthcare industry ... determines their success in global health. (500 words)

Renewable energy refers to sources ... high-carbon infrastructure and transition directly. (1000 words)

Social media has become an integral part of modern society ... of social media while minimizing harm. (2000 words)

Coral reefs are among the most diverse and productive ... science to maintain credibility and trust. (5000 words)

Output (medium summary length)

Cybersecurity is essential in today's ... sensitive personal data. (20 words)

Natural language processing (NLP) allows AI ... success in global health. (51 words)

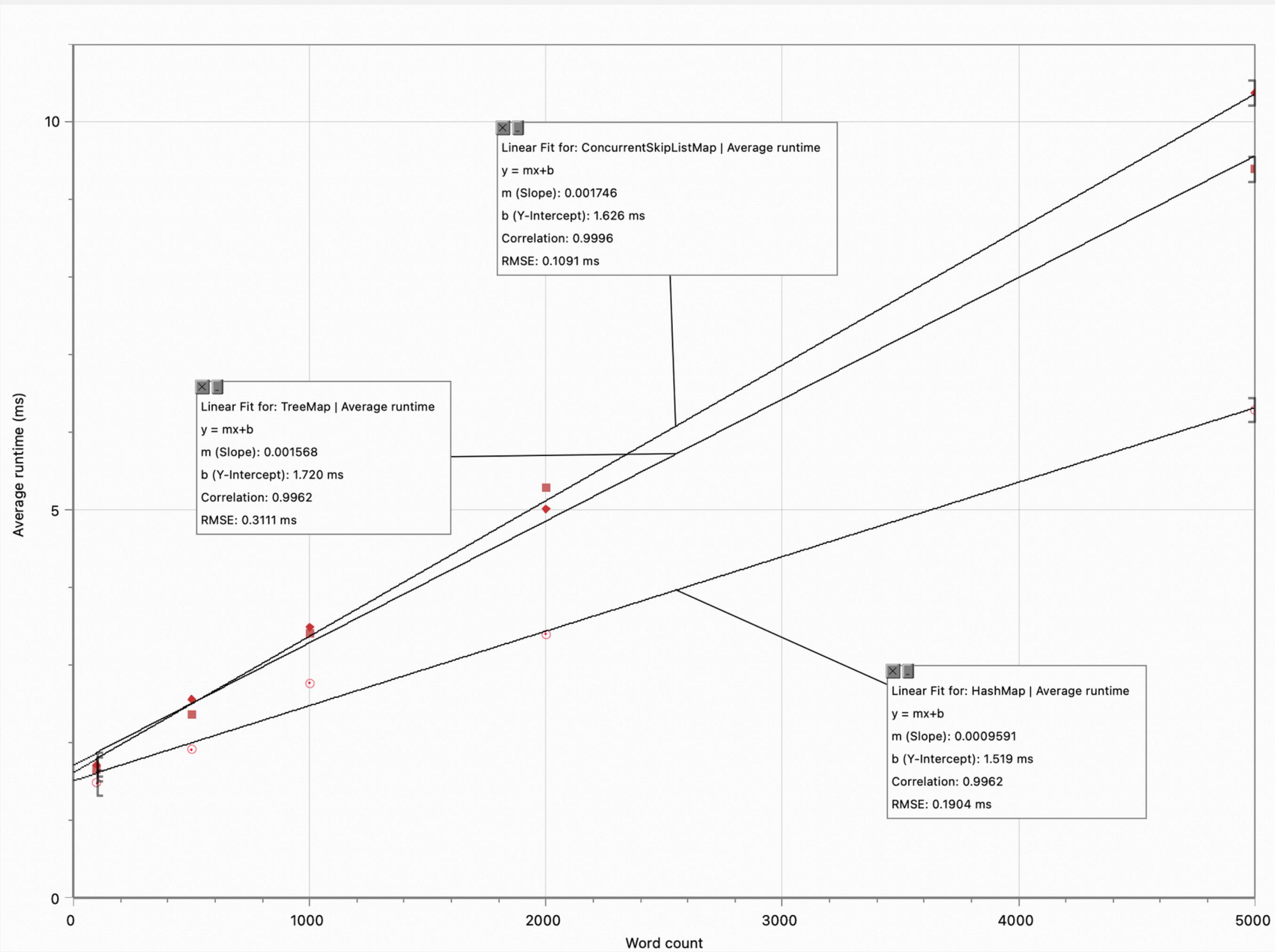
Understanding how these sources work, their advantages ... renewable portfolio standards. (112 words)

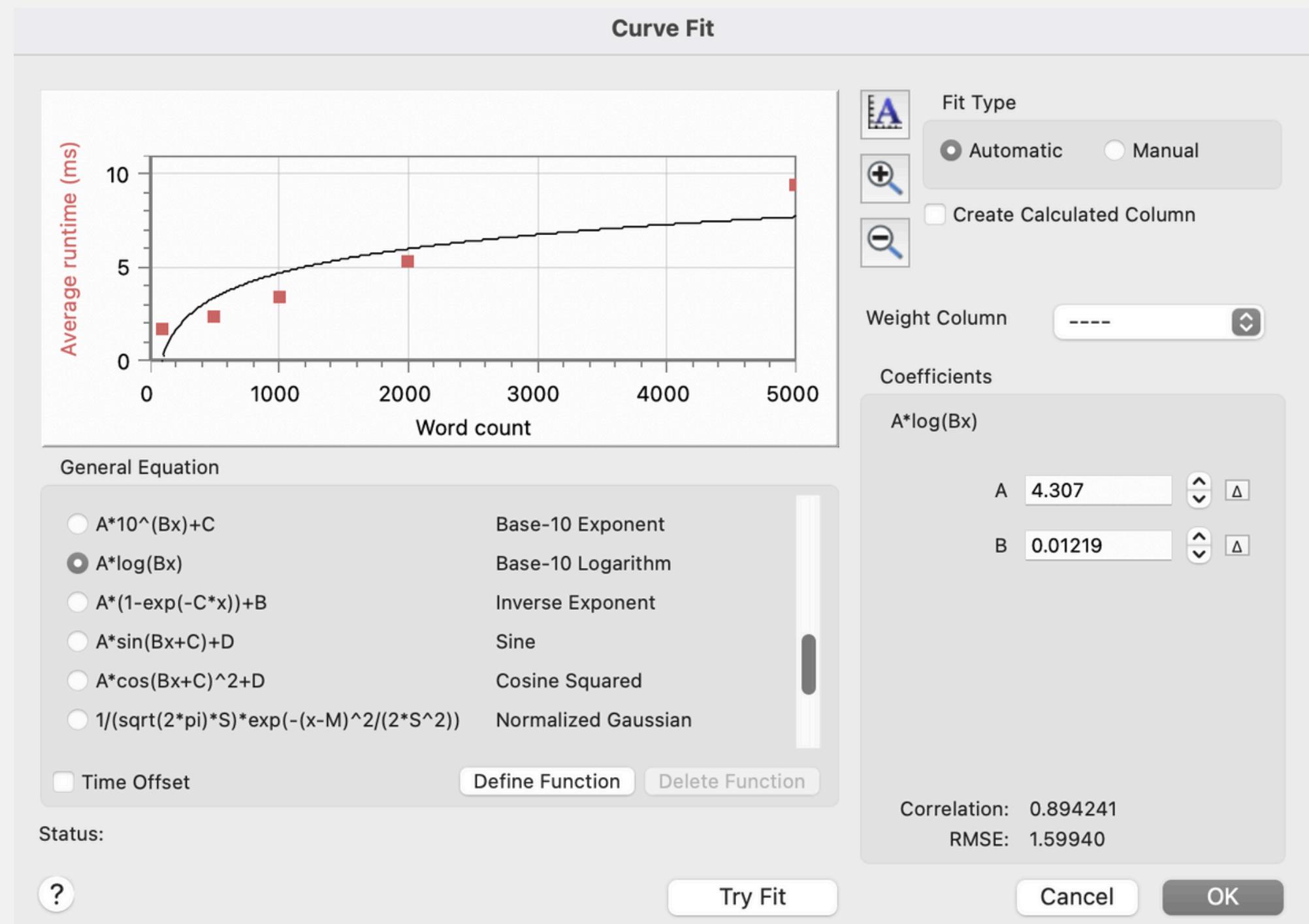
For example, passive use, defined as scrolling ... limits the ability to draw causal inferences. (124 words)

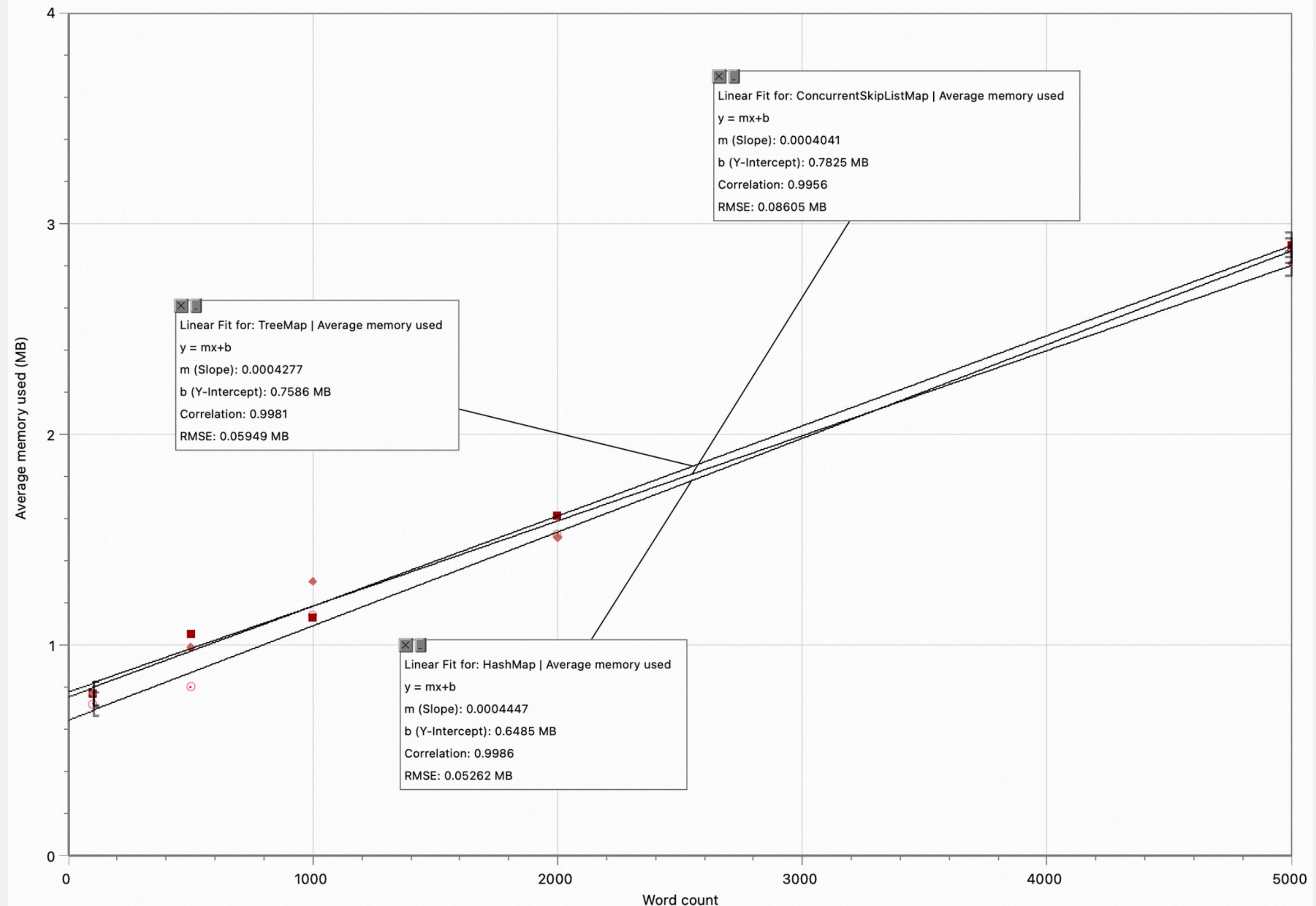
Fertilization occurs externally, and resulting ... emphasizing the economic logic of proactive investment. (192 words)

Average Runtime (ms)	HashMap	TreeMap	ConcurrentSkipListMap
Trial 1	1.481	1.654	1.717
Trial 2	1.917	2.365	2.559
Trial 3	2.764	3.402	3.487
Trial 4	3.399	5.278	5.010
Trial 5	6.284	9.386	10.372

Average Memory (MB)	HashMap	TreeMap	ConcurrentSkipListMap
Trial 1	0.722	0.768	0.772
Trial 2	0.802	1.054	0.990
Trial 3	1.144	1.133	1.302
Trial 4	1.526	1.617	1.512
Trial 5	2.873	2.899	2.812









Conclusion

Performance:

HashMap > TreeMap > ConcurrentSkipListMap

Memory behaviour:

HashMap ≈ TreeMap ≈ ConcurrentSkipListMap

Best data structure: HashMap

1. Lowest time complexity
2. Reasonable space complexity

Evaluation



- Gradient of best fit for HashMap is non-zero (0.0009591)
- Caused by hash code computations and collisions
- String equals() operations may result in growth proportional to N
- Inability to use a logarithmic function for TreeMap and ConcurrentSkipListMap
- $\log(100) = 2$, $\log(5000) \approx 3.70$
- Dominated by the fixed-cost overheads, such as pointer traversal and tree rotation
- 5000 is too small to see asymptotical behaviour



References

- Demilie, W. (2022, September 13). Comparative analysis of automated text summarization techniques: The case of Ethiopian languages. *Wireless Communication and Mobile Computing*, 2022, 13–15.
<https://doi.org/10.1155/2022/3282127>
- Jain, S. (2025, February 7). Understanding TF-IDF (Term Frequency-Inverse Document Frequency). GeeksforGeeks.
<https://www.geeksforgeeks.org/understanding-tf-idf-term-frequency-inverse-document-frequency/>

Thank You

Powered by Canva