



INSTITUTO SUPERIOR DE ENGENHARIA DE LISBOA

Área Departamental de Engenharia de Electrónica e Telecomunicações e de Computadores

A Comprehensive Survey of Asynchronous API Approaches in Concurrent I/O Scenarios

Diogo Paulo de Oliveira Rodrigues

Licenciado em Engenharia Informática e de Computadores

Dissertação para obtenção do Grau de Mestre
em Engenharia Informática e de Computadores

Orientador : Doutor Fernando Miguel Gamboa de Carvalho

Júri:

Presidente: [Doutor José Simão]

Vogais: [Doutor Fernando Miguel Gamboa de Carvalho]
[Grau e Nome do segundo vogal]

Abril, 2024



INSTITUTO SUPERIOR DE ENGENHARIA DE LISBOA

Área Departamental de Engenharia de Electrónica e Telecomunicações e de Computadores

A Comprehensive Survey of Asynchronous API Approaches in Concurrent I/O Scenarios

Diogo Paulo de Oliveira Rodrigues

Licenciado em Engenharia Informática e de Computadores

Dissertação para obtenção do Grau de Mestre
em Engenharia Informática e de Computadores

Orientador : Doutor Fernando Miguel Gamboa de Carvalho

Júri:

Presidente: [Doutor José Simão]

Vogais: [Doutor Fernando Miguel Gamboa de Carvalho]
[Grau e Nome do segundo vogal]

Abril, 2024

Aos meus ...

Acknowledgments

I would like to extend special thanks to my advisor, Doctor Miguel de Carvalho, for his persistence and availability. Without this support, I would not have been able to complete this stage. Thank you very much.

I also wish to thank ISEL, the coordinator of the MEIC master's program, and my family, who have always supported me.

To everyone, my deepest thanks.

Abstract

This thesis examines the evolution of data collection methodologies, transitioning from traditional asynchronous I/O systems to modern reactive programming paradigms such as Reactive Programming and Kotlin Flow. These paradigms represent a shift from static, sequential processes to dynamic, responsive data handling strategies, a transformation driven by recent advances in programming techniques.

After presenting a concise historical overview to contextualize the study, the research focuses on a comparative analysis of key technologies like RxJava and .NET AsyncEnumerables, which utilize reactive streams to enhance the efficiency and scalability of large data volume management. The benchmarking methodology used in this research is straightforward, based on time measure of similar operations across the different technologies used in our implementations. The tasks that were used were fundamental and simple, such as identifying the largest word in a dataset and categorizing words by letter count, in a huge data set to simulate a world-like scenario.

This pragmatic approach allows for a direct evaluation of each technology's real-time data processing capabilities, underlining practical benefits and distinguishing their performances.

In conclusion, this work provides a focused comparative analysis to enhance understanding of data collection technologies, while providing key concepts and demonstrating the advancements from traditional to modern tools, making a contribution on how these technologies compare to each other and perform depending of the situation.

Keywords: Keywords (in English): Reactive Programming, Data Collection

Methodologies, Software Paradigms, RxJava, Rx.NET, Async Enumerables, Comparative Analysis, Real-time Data Processing, Programming Paradigms Evaluation, Technology Performance Comparison.

Resumo

Esta tese examina a evolução das metodologias de recolha de dados programaticamente, focando-se na transição de sistemas tradicionais de recolha de dados através de I/O assíncrono para paradigmas modernos de programação reativa. Estes paradigmas representam uma mudança de processos bloqueantes e sequenciais para estratégias de manipulação de dados dinâmicas e responsivas, uma transformação impulsionada pelo aparecimento de novas técnicas e tecnologias de programação assíncrona.

Após apresentar uma visão histórica concisa para contextualizar o estudo, a pesquisa concentra-se numa análise comparativa entre tecnologias chave que permitem programação reativa, como o RxJava e os AsyncEnumerables da Microsoft, que utilizam streams reativos para aumentar a eficiência e escalabilidade na gestão de grandes volumes de dados. A metodologia de benchmarking é direta, baseando-se na recolha de tempos de operação em tarefas fundamentais — como, por exemplo, identificar a maior palavra num conjunto de dados e categorizar palavras pelo número de letras. Para tornar as métricas recolhidas mais representativas, foram utilizadas grandes quantidades de dados de forma a simular situações reais.

Esta abordagem pragmática permite uma avaliação direta das capacidades de processamento de dados em tempo real de cada tecnologia, sublinhando benefícios práticos e distinguindo os seus desempenhos.

Em suma, esta tese oferece uma análise comparativa entre várias tecnologias modernas de programação para a recolha de dados em grande escala, tendo também o objetivo de melhorar a compreensão deste tópico. Apresenta contexto histórico, conceitos-chave e as respostas tecnológicas atuais à leitura de dados em grande escala, contrastando com as práticas do passado.

Palavras-chave: Palavras-chave (em português): Programação Reativa, Metodologias de Coleta de Dados, Paradigmas de Software, RxJava, .NET AsyncEnumerables, Análise Comparativa, Processamento de Dados em Tempo Real, Abstração de Programação, Paralelismo

Contents

List of Figures	xvii
List of Tables	xix
List of Listings	xxi
1 Introduction	1
1.1 Motivation	2
1.2 Goals	2
1.3 Outline	3
2 State of the Art	5
2.1 Background	5
2.2 Asynchronous flow key concepts and design alternatives	9
2.2.1 Synchronous versus Asynchronous	9
2.2.2 Push vs Pull	10
2.2.3 Hot versus Cold	11
2.2.4 Cancelable	13
2.2.5 Error Handling	14
2.2.6 Intrinsic Keywords	15
2.2.7 Back-pressure	17

2.3	State of the Art	18
2.3.1	Java	18
2.3.1.1	RxJava Library	19
2.3.1.2	Reactor Library	20
2.3.1.3	Kotlin Flow	21
2.3.2	.NET C sharp	24
2.3.2.1	RxNet	24
2.3.2.2	Asynchronous Enumerables	25
2.3.3	Javascript	26
2.3.3.1	Javascript Intrinsic Words: async...await and for await...of	26
2.3.3.2	RxJS	27
2.3.3.3	IxJS - Interactive Extensions for JavaScript	28
2.3.4	Technologies comparison	28
3	Benchmarking Use Case Algorithms and Benchmark Setup	29
3.1	Overview of Algorithms	30
3.2	Pipelines	32
3.2.1	Understanding General Pipeline Operations	33
3.2.2	Pipeline Examples	36
4	Detailed Analysis and Comparison of Results	39
4.1	.NET Benchmarking	39
4.1.1	Find the biggest word algorithm results	40
4.1.2	Group Word Results	42
4.2	Java/Kotlin Benchmarking	44
4.2.1	Biggest Word Results	44
4.2.2	Group Word Results	46
4.3	JavaScript Benchmarking	48
4.3.1	Finding the Biggest Word	48
4.3.2	Grouping Words	50

5	Conclusions and Future Directions	51
5.1	Main Conclusions	51
5.1.1	.NET Environment Analysis	52
5.1.2	Java/Kotlin Environment Insights	52
5.1.3	JavaScript Strategy Performance	52
5.1.4	Overarching Conclusions	52
5.2	Future Work and Final Thoughts	52

List of Figures

4.1	Processing times for different strategies for "Find the biggest word".	40
4.2	Processing times for different strategies for "Group Words".	42
4.3	Processing times for different Java/Kotlin strategies for "Biggest Word".	45
4.4	Processing times for different Java/Kotlin strategies for "Group Words".	47
4.5	Processing times for different JavaScript strategies for "Biggest Word"	49
4.6	Processing times for different JavaScript strategies for "Grouping Words"	50

List of Tables

2.1	Your caption	20
-----	------------------------	----

List of Listings

2.1	Asynchronous call example	10
2.2	Synchronous call example	10
2.3	Example of Pull data handling patterns	11
2.4	Example of Pull data handling patterns	11
2.5	Cold flow example	12
2.6	Cold flow example	12
2.7	Intrinsic words .NET C#	15
2.8	Javascript Example with promises	16
2.9	Javascript syntax sugar	16
2.10	Flow builder	22
2.11	Kotlin collect multicall example	23
2.12	Kotlin collect multicall example	24
2.13	Asynchronous Enumerable, C# example	25
2.14	Mozilla's Javascript Asynchronous Enumerables example	27
3.1	Pseudocode for GroupWords function	30
3.2	Pseudocode for FindLargestWord function	32
3.3	C# Pipeline for Parsing Distinct Words into a Dictionary	37
3.4	Java Pipeline for Finding the Largest Word in Files	37



Introduction

As state-of-the-art non-blocking IO programming tools and technologies continue to evolve, new asynchronous idioms emerge across different programming platforms. Examples of these include Dotnet Async Enumerables [1], JS asynchronous generators [4], Kotlin Flow [5], among others.

Simultaneously, some idioms provide pipe lining API's referred for example in [2] that allow developers to build queries over sequences of data elements. This pipeline refers to a series of operations applied to a sequence of elements where each operation in the pipeline takes the output of the previous operation as input. A sequence pipeline with asynchronous data processing may be also denoted as reactive stream pipeline.

The research work that I describe in this dissertation aims to analyze what are the intrinsic overheads behind a reactive stream pipeline processing. While there are advantages regarding their conciseness, expressiveness, and readability, will these pipelines incur any overhead?

To answer this question we will analyze the behavior of different reactive stream pipelines technologies dealing non-blocking IO.

While non-blocking IO pipe lining enabling technologies, allows programmers to exploit concurrency without explicitly managing threads, it introduces new challenges and potential performance bottlenecks. The seamless integration of non-blocking IO with reactive stream pipelines presents a complex interplay of factors that may impact performance and resource utilization.

This dissertation aims, to retrieve performance metrics in several technologies and discuss how these technologies behave in different use cases.

1.1 Motivation

Reactive stream pipelines, it's a concept behind several powerful tools that allows the orchestration of a series of operations on a live feed of sequence of elements in continuous non-blocking IO operations. The ability to perform asynchronous data processing through these pipelines provides developers with a concise, expressive, and readable way to handle complex tasks.

However, with these benefits come potential challenges. The integration of non-blocking IO with reactive stream pipelines introduces a complex interplay of factors that may significantly impact performance and resource utilization. While the industry is rapidly adopting these idioms for their expressiveness and ease of use, there is a growing need to understand the underlying overhead and bottlenecks. Misunderstanding or overlooking these aspects can lead to inefficient code, wasted resources, and ultimately a poor user experience.

Moreover, the variety of technologies offering reactive stream pipeline processing adds another layer of complexity. Different approaches may behave differently in various use cases, leading to uncertainty when selecting the right technology for a specific project. Understanding how these technologies perform in different scenarios is vital for informed decision-making and effective utilization.

This work aims to provide a clear picture of how and when to use each technology and strategy, along with an understanding of their respective benefits and drawbacks in comparison to one another.

1.2 Goals

This research is guided by the following central objectives:

1. **Non-blocking IO:** Investigate state-of-the-arts idioms and technologies that enable processing in non-blocking IO use cases.
2. **Baseline Behavior Analysis:** Analyze the behavior of a baseline approach, for example, a piece of code that makes no use of any pipeline of operations,

the logic is made through direct implementation of algorithms using non-abstract programming. Understanding the trade-offs of such an approach, including verbosity, complexity, and maintainability, against performance benefits.

3. **Technology Comparison:** Evaluate and contrast various idioms and technologies that take advantage of non-blocking IO, focusing on how different paradigms (e.g., pull vs. push) concepts and how they perform under various conditions.
4. **Reactive Streams Technologies:** Investigate the practical implementation of reactive stream pipelines using different technologies.

We expect that a baseline approach without the use of any pipeline technology, perform better than any other alternative, despite its verbosity, complexity, and difficult maintainability.

On the other side, we would like to understand which approach that provides an asynchronous pipeline API that performs better, for example in the same language depending on the task.

Lastly, this work aims to have results and conclusions from different programming languages implementations of asynchronous pipelining Technologies. This is crucial to understand, for example, if the same technology e.g. *Reactor* performance differs depending on the programming language; on the other side, its interesting to know too, if any programming environment performs significantly better than the others in this particular set of JAVA, Kotlin, C#, and JavaScript.

We aim to observe how performance behaves in these different scenarios to uncover the distinct advantages and disadvantages of each approach. This in-depth examination is intended to culminate in practical conclusions that can guide informed decision-making in future projects.

1.3 Outline

This dissertation is organized into distinct chapters that systematically guide the reader through the research, concepts, methodologies, and findings:

- **Chapter 1: Introduction** - This initial chapter provides an overview of the research context, the motivation behind the study, the primary goals, and the approach taken.

- **Chapter 2: State-of-the-Art** - Here, we explore the foundational concepts related to reactive stream pipelines and provide a comprehensive review of the existing technologies and strategies in JAVA (including Kotlin), C#, and Javascript. The chapter aims to establish the current landscape and set the stage for our research.
- **Chapter 3: Methodology** - In this chapter, the research design, including the selection of technologies, experimental setup, and methodologies used to gather and analyze data, will be detailed. The rationale behind these choices will also be explained.
- **Chapter 4: Detailed Analysis and Comparison of Results** - This chapter presents the results obtained from the implementations and provides a thorough analysis and discussion of the findings. It includes comparisons between various reactive technologies and other strategies, and how they perform against baseline approaches.
- **Chapter 5: Conclusions and Future Directions** - The final chapter will summarize the main insights derived from the study, highlight the contributions made, discuss the limitations, and propose directions for future research.



State of the Art

This section, aims to show the background and the current state of the art on the subject related with asynchronous pipeline operations involving asynchronous data flows. On the section 2.1, will be made an overview on previously developed work made on this subject, then, on section 2.2, will be made a characterization of the key concepts related to it. By last, on section 2.3, are presented and explained several technologies representative of the state of the Art on how asynchronous operation pipe lining are implemented in different programming realities, e.g. on Kotlin, JAVA and C#.

2.1 Background

From the end of 80's to the beginning of the 2000's, with the acceleration of Moore's Law in hardware and network bandwidth development; the creation of the internet as we know today through the wide spread of use of the HTTP protocol, the necessity of high responsiveness and efficient systems started to grow. This increase in demand of new ways to handle data efficiently and the wide spread of multicore processors created the necessity of the implementation of new software frameworks models that can handle data efficiently taking in account the multithread systems and the asynchronous characteristics of data receiving throw network.

Taking the wave initiated by the Gang of Four in [3], where 23 patterns were compiled to deal with object-oriented problems; a group of researchers wrote a paper [7], where they identify the characteristics that high available systems must have, and how the use of asynchronous IO pipe design patterns can help to achieve these characteristics, as described bellow:

- Concurrency - The server must process multiple client requests simultaneously.
- Efficiency - The software design must be built aiming the use the least hardware resources as possible.
- Simplicity - The code of the solution must be easy to understand, modular and avoid own built design patterns as possible.
- Adaptability - The system must be totally decoupled from client implementations, allowing it to be easily used by any client independently of the underlying technologic realities. To achieve this, may be used standards e.g. [8] or SOAP.

To achieve some of these properties, the authors propose the *Proactor Pattern* [7]. In their opinion, traditional concurrency models, until that date, failed to fully achieve the enumerated properties.

On the other side, before presenting the *Proactor Pattern*, are identified two major existing non-blocking models, namely: *multithreading* and *reactive event dispatching*.

On multithreading, the paper refers that one of the most direct solution of this approach, is the handling of multiple requests by creating a new thread every request. Each request will then be fully processed and the recently created thread is then be disposed after the work is done.

This solution has several serious issues. Firstly, creating a new thread per request is highly costly in terms of computational resources, because are involved context switches between user and kernel modes; secondly, must be taken in account synchronization to maintain data integrity. Then, the authors warn about the fact

that the IO retrieved data is mainly memory-mapped, which rises the question: What happens when the data obtained through IO becomes greater than the system memory can hold? The system stalls until more memory becomes available!? One possible solution is asynchronous pipelines that will be discussed ahead in this dissertation. On last, if the server receives a high demand of requests, the server easily blocks in the process of creating and disposing threads.

To avoid this issue, the authors, recommended the use of dynamic thread pools to process requests, where each request will be linked to a pre-existing thread, avoiding all the overhead of creating and disposing a thread per request; however, issues related with memory-mapping and overhead due to the switching of data between different threads maintains.

Another traditional concurrency model identified by the authors of the paper, is the *Reactive Synchronous Event Dispatching*. In this model, a *Dispatcher*, with a single thread in a loop, is constantly listening requests from clients and sending work requests to an entity named *Handler*. The *Handler*, will then process the IO work Synchronously and request a connection to the client in the *Dispatcher*. When the requested connection is ready to be used, the *Dispatcher* notifies the *Handler*. After the notification, the *Handler* asynchronously sends the data, that is being or has been obtained through IO, to the client.

Although the authors identifying that this approach is positive, because decouples the application logic from the dispatching mechanisms, besides with the low overhead due the use of a single thread, the authors identify several drawbacks with this approach. Firstly, since IO operation are synchronous, the code for this approach is complex because must be set in place mechanisms to avoid IO blocking through hand off mechanisms. Then, if a request processing blocks, the processing of another requests may be impacted.

To keep the positive points but mitigating the identified issues of previous approaches, is suggested the *Proactor Pattern*. This pattern is very similar to the *Reactive Synchronous Event Dispatching*, however, after the requests processed by a single threaded *Completion dispatcher*, the IO work is then dispatched asynchronously to the underlying OS IO subsystems, where multiple requests can be processed simultaneously. For the result to be retrieved, is previously registered a callback in the *Completion Dispatcher* and the OS has the responsibility to queue the finished result in a well known place.

Finally, the *Completion Dispatcher* has the responsibility to dequeue the result placed by the OS and call the correct previously registered callback. With this,

this model creates a platform that provides: decoupling between application and processing mechanisms, offers concurrent work without the issues inherent with the use of threading mechanisms and since IO is managed by the OS subsystems, is avoided code complexity in handling possible blocking and scheduling issues.

The *Proactor Pattern*, creates the ground for several models used by modern platforms that use a reduced number of threads to process client requests and parallel mechanisms to do the heavy work in the background. Some of these technologies, are, for example: *Javascript NODE.JS*, *Spring Webflux*, *vertx* and others.

From what was explained until now, is evident the tendency followed by software architects in terms of asynchronous processing from non-reactive to event driven approaches and trying to make the code simpler and easier to maintain. Initially the systems were non-reactive, where each request had to be processed in a specific thread and that thread blocked until something got ready to go further. The code was usually very complex and hard to maintain. Then, with the asynchronous systems based on events with the introduction of callback systems inspired in patterns like the *Reactor* or *Proactor* the software design started to become more event driven, allowing the servers to be more efficient in responsiveness, resources optimization and code easier to read and maintain by the average programmer.

However, are some limitations in these asynchronous models. For example, if the data to be processed is bigger than the memory available at a given moment or if the data to be calculated is from a source that produces data at a constant rate that must be processed in real time, these models work badly. The traditional models fail to comply these objectives because are mostly eager by design or not comply with the notion of a continuous source of information that requires to be processed in real time. Taken this in account, projects like project Reactor, Asynchronous Enumerable provided by Microsoft or papers like [6], try to deal with these issues, by providing API's that merge the concepts of Fluent API's, functional programming and code syntax that tries to resemble synchronous code, being the complexity inherent with asynchronous models implementations hidden from the programmer since the asynchronous events are treated like, for example, a `Stream` in JAVA would be treated.

2.2 Asynchronous flow key concepts and design alternatives

With the development of several approaches related to asynchronous IO processing, and with the growing necessity to build a name set of properties that simplify the description of asynchronous systems, a dictionary of properties, concepts and design alternatives started to grow by itself. In the following, are discussed several of the concepts related with asynchronous data flow, namely:

2.2.1 Synchronous versus Asynchronous

Before explaining more terms related with asynchronous data flow, it's important to clarify what is synchronous and asynchronous in programming.

Asynchronous in programming, is a call to a function or routine that returns immediately, not blocking the caller until the operation is finished. The operation processing, will be completely independent from the caller execution process and can even be done in another machine. This way, the caller is freed to do more work, even to start N more operations in parallel.

Meanwhile, a call to a synchronous function or routine, blocks the caller until the operation finishes. In this case, the caller has to wait for the completion of the synchronous operation before going forward, which limits the program efficiency if parallelism is applicable. To better visualize what was explained, we have the following examples:

```

1  HttpClient client = HttpClient.
    newHttpClient();
2
3  HttpRequest request = HttpRequest
4      .newBuilder(URI.create("SOME MOVIE API
    "))....
5
6  Task futureResponse = client
7      .sendAsync(request, new
8          JsonBodyHandler<>(DTO.class))
9      .thenAccept(res -> {
10         Console.WriteLine( res.title);
11     });
12 Console.WriteLine("prints something")
13
14 client.Wait();
15
16 //OUTPUT:
17 //prints something
18 //movie title
19

```

Listing 2.1: Asynchronous call example

```

1  HttpClient client = HttpClient.
    newHttpClient();
2
3  HttpRequest request = HttpRequest.
    newBuilder(
4      URI.create("SOME URL"))
5      .header("accept", "application/json")
6      .build();
7
8  HttpResponseMessage<Supplier<DTO>> res =
9      client.send(request, new
10         JsonBodyHandler<>(DTO.class)) //
11         synchronous
12
13 Console.WriteLine(res.title);
14
15 Console.WriteLine("prints something")
16
17 //OUTPUT:
18 //movie title
19 //prints something

```

Listing 2.2: Synchronous call example

As we can see, in the synchronous call example, the operation return only happens after the whole subsequent remote operation is finished, consequently, the caller operation is dependent from several variables to go forward e.g. : HTTP messaging latency, remote server operation speed or bandwidth issues. This causes that the processing of the next code statements to happen only after the synchronous IO call.

Meanwhile, in the asynchronous operation call, the return happens immediately after the call, however, the processing inherent with that operation will start just when the subsystem that handles the asynchronous function is ready to process that work. For example, when the OS is ready to process the received responses from a remote server that handled the operation. In this case, the statement right after the asynchronous call is processed before the asynchronous operation. Allowing the IO asynchronous operation to be processed outside the program scope and avoiding any block of the main program.

2.2.2 Push vs Pull

Another concept important to understand how asynchronous data flow is handled in programming, is the *Pull* and *Push* processing patterns. In *Pull* pattern,

usually, exists a source of data and the program iterates over that source to operate over each item.

On the other hand, in the *Push* pattern, the items of the data source are "Pushed" to a routine that will operate over that item. To help to assimilate what was just explained, we have the following example:

```

1 Flowable<Long> flow = Flowable
2   .interval(1, TimeUnit.SECONDS);
3
4 Iterator<Long> iterator =
5   flow.blockingIterable().iterator();
6
7 while (iterator.hasNext())
8   System.out.println(iterator.next());
9
10
```

Listing 2.3: Example of Pull data handling patterns

```

1 Flowable<Long> flow = Flowable
2   .interval(1, TimeUnit.SECONDS);
3
4 flow.blockingSubscribe(System.out::println
5   );
6
```

Listing 2.4: Example of Pull data handling patterns

As we can see, in the pull pattern, the items are "pulled" from a data source through an iteration mechanism.

In contrast with that, in the *Push* pattern, items are pushed to a consumer through a supplier.

2.2.3 Hot versus Cold

Another property that must be taken in account when handling with *Reactive Streams* or asynchronous data processing in general, is the nature of the data flow. There are two main adjectives to name a data flow, `Hot` or `Cold`.

A `Cold` data flow, is a flow of information that is produced just when the stream pipeline is subscribed by an observer. In this case, the producer only starts sending/producing data when someone is interested in the data from that source. For example, when program uses a IO mechanism to lazily retrieve a sequence of words from a database, the IO mechanism will only start sending information just when a consumer subscribes that data flow. Usually, the data is sent to the consumer in unicast.

On the other hand, in a `Hot` data flow, the data is produced independently of existing any observer to that information. This mechanism usually work in broadcast and the data is continuously produced and sent to possible observers. In this

case, when an observer subscribes to a publisher, exists the possibility of data items being already lost to that publisher while in the `Cold` flow, the consumer usually receives all items that were produced by the source. In the following examples, a number is produced each 100 milliseconds:

<pre> 1 ConnectableFlowable<Long> hot = Flowable 2 .intervalRange(0, 100, 0, 100, TimeUnit.MILLISECONDS) 3 .publish(); 4 hot.connect(); 5 Thread.sleep(1000); 6 hot.blockingSubscribe(System.out:: println); 7 // output: 8 // 2 9 // 3 10 // 4 11 </pre>	<pre> 1 Flowable<Long> cold = Flowable 2 .interval(100, TimeUnit. MILLISECONDS); 3 Thread.sleep(1000); 4 cold.blockingSubscribe(System.out:: println); 5 // Output: 6 // 1 7 // 2 8 // 3 9 </pre>
--	---

Listing 2.6: Cold flow example

Listing 2.5: Cold flow example

As we can see, in the `Hot` data flow example, the items are emitted from the moment the producer is created, independently of existing any subscriber or observer attached to that publisher. Notice that when a consumer is subscribed to the publisher, 1 seconds after the emission started, the numbers from 0 to 10 were not printed.

In the `Cold` example, the producer only emits data when a subscription is done, and because of that, all the produced numbers were printed, in contrast with what happen in `Hot` stream, where data loss are almost certain.

2.2.4 Cancelable

As already stated above, asynchronous operations may run outside the main program scope. This implies, that the main program loses visibility and control on what happens in the asynchronous operation contexts. Because of that, exists the need to put in place mechanisms of control that allow the main program to maintain control over an asynchronous operation to, for example, cancel the operation or put in place finishing logic that allow, for example: resource disposing, logging, decisions etc...

These mechanisms, are many times done through the concept of *cancelables*. Usually, a cancelable, is an entity that represents an operation that can canceled from an external entity, or, an entity that allows to set logic when an operation finishes by any reason.

In C#, a cancelable is an interface implemented by objects that represent asynchronous operations and provides the means to cancel asynchronous operations, on the fly. This is achieved through a mechanism named: `CancellationToken`, that is used to pass information through different execution threads.

In RxJava, a `Cancelable`, is a functional interface with the method `cancel()`. Then, the `Cancelable` can be associated to a data source representation, the `Observable`, by calling the method `Observable.setCancelable(Cancelable)`. When the `Observable` finishes or is canceled for any reason, the method `Cancelable.cancel()` will be called. This way, proper logic is put in place to handle an asynchronous operation cancelation.

As we saw, these two concepts of cancelable diverge. One, provides the means to cancel an operation on the fly and gives some control over the operation cancelation; the other, provides the means to control an operation cancelation independently of how it was cancelled.

2.2.5 Error Handling

In synchronous environments, usually, when something goes wrong, the way to handle an error in the majority of cases is by throwing an exception and propagate it until the proper code handles it, usually in a try/catch block. However, in cases when exists an asynchronous operation or when a continuous stream of data items are being received, that way of dealing with an error can imply several issues, e.g: exceptions not reaching main program, log losing or asynchronous flow blockage.

Since log losing or blocking a whole operation because of a badly handled error is unacceptable, the best way to deal with errors in asynchronous data flow it is to isolate the error. This way, the flow processing may continue in parallel while the error is properly handled.

The best way to handle this kind of errors, it is to have proper callbacks that are called when an error occurs on the stream item. This way, a function can handle the error properly, without the necessity to blocking any data stream processing, if avoidable and the proper logging and any additional measure to handle it can be put in place.

2.2.6 Intrinsic Keywords

As already stated, asynchronous code is tendentially harder to understand because, in opposition with what happens in synchronous environments, the operations inherent with the sequence of programming statements operations may not happen chronologically ordered. Because of that, many times it is difficult read, debug and sustain asynchronous software.

For that reason, many languages started to add syntax techniques that allow the programmer to build asynchronous code that resembles the synchronous syntax. Under the hoods, the virtual machines that sustain these syntax mechanisms, handle the code bounded with that 'intrinsic words' and builds asynchronous routines that the programmer will not be aware of; being this a way to abstract the programmer from the complexity of handling and sustaining complex asynchronous code.

One example of *intrinsic words* mechanisms, is the `async...await` keywords implemented in *Microsoft's .NET C#* and in *javascript*.

In the next example we can see an example of these keywords being used:

```
1 static async Task Main(string[] args)
2 {
3     IEnumerable<int> enumerable = FetchItems(1000);
4     int i = 0;
5     await foreach (var item in enumerable)
6     {
7         if (i++ == 10){ break;}
8     }
9 }
10
11 static async IEnumerable<int> FetchItems(int delay)
12 {
13     int count = 0;
14     while(true)
15     {
16         await Task.Delay(delay);
17         yield return count++;
18     }
19 }
20
```

Listing 2.7: Intrinsic words .NET C#

Where, for example, we can observe in the line 16, a call to an asynchronous operation, and, by adding the keyword `await`, the next statement although

being a call to an asynchronous operation, the code statement order looks like it is from synchronous set of instructions.

Additionally, the use of `async...await` in .NET, for example, simplifies error handling in asynchronous code. Instead of use a callback to handle an error, by using `async...await` the error can be handled by just using a simple try/catch block.

To better visualize the advantage of using intrinsic words in asynchronous code, on the next example, we can see a code comparison in ECM6 Javascript, with and without the use of intrinsic words in asynchronous code. In the next example, it is possible to observe a decrease of code complexity and increment of readability where is used the "syntax sugar" provided by the 'yield' return.

The example using promises is purposefully made with a "Pyramid of Doom" code to accentuate a difficulty of reading asynchronous if is made without any concern with readability.

<pre> 1 function ourImportantFunction(callback) 2 { 3 task1(function(val1) { 4 task2(val1, function(val2) { 5 task3(val2, callback); 6 }); 7 }); 8 }</pre>	<pre> 1 function ourImportantFunction() { 2 3 var val1 = yield task1(); 4 5 var val2 = yield task2(val1); 6 7 var val3 = yield task3(val2); 8 9 return val3; 10 } 11 12</pre>
--	--

Listing 2.9: Javascript syntax sugar

Listing 2.8: Javascript Example with promises

As we can see, with the use of `yield` keyword, the code that uses a result of several asynchronous operation, instead of being used in a "Matrioska Dool" type of code, with a code made with a chain of callback results; the simple use of a intrinsic keyword like `yield` simplifies the code a lot. Making the code previously hard to read in a easier code to understand and maintain.

2.2.7 Back-pressure

When the *pull* method is used to retrieve items from a source, the producer retrieves only the items it can process in the given time.

However, when the *push* approach is used as data retrieval method from asynchronous flows, the producers have the initiative to push items to its consumers. This can originate situations, where the producer emits items faster than the producers can handle, which can create problems like: unwanted loss of data, lack of responsiveness from consumers, etc...

To resolve these issues, were created strategies and design patterns that are commonly referred as *Backpressure*. There are four main approaches which the majority of *Backpressure* strategies are designed from and can be resumed as:

1. **DROP:** Producer drops items after a retrieving buffer gets full.
2. **Buffer everything:** A buffer, keeps all unprocessed items that are received. Usually, this strategy is used when all received items are critical for the business development and memory management has flexibility to handle the increase of storage needs.
3. **Error:** An error is thrown when the buffer threshold is reached, usually all items received after the threshold is reached are discarded.
4. **Lastest:** Only the last received item in the given moment is kept.
5. **Missing:** No back-pressure strategy it is in place, all items that ca not be processed on arrival, are discarded.

2.3 State of the Art

In this section, we will delve into various state-of-the-art frameworks designed for asynchronous IO pipeline processing across multiple programming languages. Initially, we explore the options available in the Java landscape in Section 2.3.1. Here, we begin with an explanation of the multi-language project *ReactiveX.io* and its Java implementation, *RxJava*. Subsequently, we will discuss the Reactor Project, followed by an examination of non-blocking processing in the JVM, specifically through Kotlin's native `Flow` implementation.

Following Java, we shift our focus to Microsoft's .NET C# in Section 2.3.2. A brief mention of ReactiveX's c# implementation precedes a deep dive into how C# natively addresses this challenge, particularly through the implementation of `AsyncEnumerables`.

Next, we explore JavaScript's native solution to this problem, highlighting the use of asynchronous Iterables through the *intrinsic keywords* `FOR` `AWAIT` `...OF` and Promises. Additionally, we refer to ReactiveX's implementation for JavaScript, `RXJS`.

Finally, in Section 2.3.4, we provide a comprehensive overview and draw conclusions about the various technologies discussed, comparing how each approach can be utilized for different problems and objectives. We also make theoretical predictions about the potential performance of each technology under several known circumstances.

2.3.1 Java

In the context of Java, we will explore several libraries and frameworks aimed at simplifying asynchronous data flow handling.

The first is *RxJava*, which is a Java implementation of the *ReactiveX.io* project. This library uses the `Observer` pattern to handle real-time asynchronous processing with and without back-pressure. It is important to note that the *ReactiveX.io* project is a multi-language project, and its Java implementation, *RxJava*, is among the most widely used.

Next, we delve into *Project Reactor*, an integral part of Spring WebFlux's non-blocking web stack. Although it uses a different approach from *ReactiveX*, it still provides the same benefits, like allowing developers to work with a composable API for declarative, event-driven programming.

Lastly, we will discuss the *Kotlin Flow* strategy, which is utilized in the Kotlin language but interoperable with Java. This will provide us with a perspective on how coroutine-based asynchronous data processing is implemented in the JVM environment, contrasting with the Observer pattern used by ReactiveX and Project Reactor.

2.3.1.1 RxJava Library

Before we delve into the specifics of RxJava, it's important to discuss ReactiveX (Reactive Extensions) — the project that gave birth to it. ReactiveX is a multi-language project focusing on combining the observer pattern, iterator pattern, and functional programming techniques to make the handling of asynchronous streams of data manageable and efficient. ReactiveX libraries exist for a variety of programming languages, including JavaScript, Java, C#, and others.

ReactiveX adopts a declarative approach to concurrency, abstracting away the complexities associated with low-level threading, synchronization, thread-safety, concurrent data structures, and non-blocking I/O. Instead, it encourages developers to focus on the composition of asynchronous data streams.

An *observable* sequence in ReactiveX can emit three types of items: "next" notifications (carrying items to observers), "error" notifications (carrying error information), and "complete" notifications (signaling the end of the sequence). Observers subscribe to these observable sequences and react to whatever item they emit. The sequences are lazy; items are not pushed to observers until an observer subscribes.

Now, let's turn our attention to RxJava, the Java implementation of ReactiveX.

RxJava, an open-source project, encapsulates the *Observer pattern*, the *Iterator pattern*, and functional programming techniques to manage asynchronous data flow and control event sequences.

The data source/publisher of an asynchronous event stream in RxJava is represented by the `Observable` and `Flowable` classes. The key distinction is that `Flowable` supports back-pressure through various buffering strategies.

In RxJava, `Observable` and `Flowable` represent a stream of *N* events, while a single event is encapsulated by the `Single` class.

RxJava offers fluent APIs in `Observable` and `Flowable`, enabling operation chaining in pipelines — a feature also found in synchronous environments like

Java's *Stream* fluent API or .NET's Linq framework. Developers can perform stream processing operations like `filter`, `flatMap`, and `distinct` on these asynchronous event sequences, much as they would on synchronous streams in fluent APIs.

Observers/subscribers in RxJava are consumers subscribing to *Observable* through the `Observable.subscribe()` method. These consumers can be implementations of the functional interface `Consumer<T>` or the *Observer* interface. The latter provides enhanced control over stream processing through error handling capabilities, which `Consumer<T>` implementations lack.

The *Observer* can be regarded as an asynchronous counterpart to the Java util interface *Iterator*, as evidenced by the similarity in their interface methods.

1. `onSubscription()`: This method is called immediately after a subscription is made with an *Observable*.
2. `onNext(T item)`: This method is called when an item is emitted by the asynchronous data source.
3. `onError()`: Contrary to the *Iterator*, the RxJava *Observer* is equipped to support error handling. This callback is invoked when an error occurs.
4. `onComplete()`: This method is called when the data source closes or the subscription finishes.

Given the possible relationship between synchronous and asynchronous programming, Table ?? was developed to better visualize the correlation.

	Single Item	Multiple Items
Java	<code>T getData()</code>	<code>Iterable<T> getData()</code>
RxJava	<code>Single<T> getData()</code>	<code>Observable<T> getData()</code>

Table 2.1: Your caption

2.3.1.2 Reactor Library

The *Project Reactor* is another part of the Spring portfolio of projects. It is designed to be a fully non-blocking foundation for Java, compliant with the Reactive Streams specification. It offers efficient demand management (back-pressure) capabilities, making it an ideal choice for scenarios involving live streams of data.

The design of Project Reactor is also based on the *Publisher/Subscriber pattern*. However, instead of using `Observable`, `Flowable`, and `Single`, Project Reactor uses `Flux` and `Mono` to represent asynchronous data streams. `Flux` represents a stream of 0 to N items, while `Mono` represents a stream of 0 or 1 item.

Similar to RxJava, Project Reactor provides a variety of operators that can be used to transform, filter, combine, and manipulate data streams. This allows developers to construct intuitive instruction pipelines.

Just like in RxJava, observers/subscribers in Project Reactor are represented by consumers that are attached to `Flux` and `Mono` via the `subscribe()` method. These consumers can either be implementations of the `Consumer<T>` functional interface or the `Subscription` interface. Here's a brief comparison of how Reactor relates to the conventional synchronous counterparts in Java:

	Single Item	Multiple Items
Java	<code>T getData()</code>	<code>Iterable<T> getData()</code>
Reactor	<code>Mono<T> getData()</code>	<code>Flux<T> getData()</code>

By adhering to the Reactive Streams specification and offering a wide array of operators to handle data, Project Reactor serves as a powerful tool for developing reactive, non-blocking applications in Java.

2.3.1.3 Kotlin Flow

Kotlin, despite being a JVM-based language, differs significantly from Java when it comes to handling asynchronous data flows. Kotlin uses coroutines, a feature natively supported in the language, to simplify asynchronous programming. This feature is used in the implementation of `Flow<T>`, Kotlin's main interface for handling asynchronous data flows.

In comparison to `Observable` and `Flux`, which are based on the Observer pattern, Kotlin's `Flow<T>` is more aligned with the principles of the *Publisher/Subscriber pattern*. This is evident in how the `Flow<T>` interface is implemented. The data source implementation for a `Flow<T>` is done using a builder, and the initiation of the flow sequence is triggered by the `Flow.collect()` method.

Hot flows in Kotlin are represented by the `SharedFlow<out T> : Flow<T>` interface. Unlike `Flow<T>`, which initiates a flow every time `Flow.collect()` is called, `SharedFlow.collect()` emits an unpredictable set of items from an external stream of events initiated before the call to `SharedFlow.collect()`.

Kotlin's `Flow<T>` provides a fluent API of intermediate operators that allow data transformation through operation pipelines, similar to what we have already seen in RxJava and Project Reactor.

A `Flow<T>` data source implementation is done through a builder, like we can see in the next example:

```

1 fun simple(): Flow<Int> = flow { // flow builder
2     for (i in 1..3) {
3         delay(100) // pretend we are doing something useful here
4         emit(i) // emit next value
5     }
6 }
7
8 fun main() = runBlocking<Unit> {
9     launch {
10         for (k in 1..3) {
11             println("I'm not blocked $k")
12             delay(100)
13         }
14     }
15
16     simple().collect { value -> println(value) }
17 }
18
19 //output:
20 //I'm not blocked 1
21 //1
22 //I'm not blocked 2
23 //2
24 //I'm not blocked 3
25 //3
26
27
```

Listing 2.10: Flow builder

A consumer, to start listening a particular data flow has to call the method `Flow.collect()`

. Since `Flow` provides support only to cold flows, calling `collect()` has the particularity of initiating flow the sequence. On the other side, Hot flows in Kotlin are represented by the interface `SharedFlow<out T> : Flow<T>`.

The main difference between the implementation of these two interfaces, is at the result of `collect()` call. While the `Flow.collect()` starts the flow every

time its called, resulting in the limited emission of the same set of items per call; the `SharedFlow.collect()` emits an unpredicted set of items from external stream of events initiated before the call to `SharedFlow.collect()`. On the other side, while the `Flow.collect()` call context is private to the caller, the `ShareFlow.collect()` is shareable by N subscribers, which makes this solution ideal for broadcast mechanisms shared by many users. On the next example, we can see several calls to the `Flow.collect()` that results in retrieving the same set of items; as explained, the call starts a cold flow every time its called:

```
1 fun simple(): Flow<Int> = flow {
2     println("Flow started")
3     for (i in 1..3) {
4         delay(100)
5         emit(i)
6     }
7 }
8
9 fun main() = runBlocking<Unit> {
10     println("Calling simple function...")
11     val flow = simple()
12     println("Calling collect...")
13     flow.collect { value -> println(value) }
14     println("Calling collect again...")
15     flow.collect { value -> println(value) }
16 }
17
18 //Output:
19 //Calling simple function...
20 //Calling collect...
21 //Flow started
22 //1
23 //2
24 //3
25 //Calling collect again...
26 //Flow started
27 //1
28 //2
29 //3
30
31
```

Listing 2.11: Kotlin collect multicall example

As we can see, with the use of the keyword *emit*, it is achieved the same of what we saw in C# with the use of the keyword *yield return*. In this case, the same way the *yield return* returned an item that was part of an asynchronous enumeration of events through *IAsyncEnumerable*, the use of the keyword *emit* will lazily emit data, as it becomes available to be set as event of the Flow.

Likewise what happens in RxJava, `Flow<T>` provides a fluent API of intermediate operators that allow data transformation through the use of operation pipelines, where is received an upstream flow and the operators return a transformed downstream flow through the traditional push methods like: `filter`, `map()`, `zip()`, `take()` etc... On the next example, we can see an example of an asynchronous data pipeline operation from Kotlin `Flow<T>`, taking advantage of the Fluent API provided by its platform:

```
1 suspend fun performRequest(request: Int): String {
2     delay(1000) // imitate long-running asynchronous work
3     return "response $request"
4 }
5
6 fun main() = runBlocking<Unit> {
7     (1..3).asFlow() // a flow of requests
8         .map { request -> performRequest(request) }
9         .collect { response -> println(response) }
10 }
11
12 //Output:
13 //response 1
14 //response 2
15 //response 3
16
17
```

Listing 2.12: Kotlin collect multicall example

2.3.2 .NET C sharp

2.3.2.1 RxNet

RxNet is the .NET implementation of the ReactiveX project, providing the same powerful programming paradigm from ReactiveX to the .NET ecosystem. This implementation allows developers in the .NET framework to effectively manage asynchronous data flow using the Observer pattern and functional programming techniques, as explained in the ReactiveX and RxJava sections.

In the context of C#, the concepts of Observables and Observers (or Subscribers) are represented by the `IObservable<T>` and `IObserver<T>` interfaces. These are analogous to the `IEnumerable<T>` and `IEnumerator<T>` interfaces in the synchronous realm, and closely resemble the implementations seen in RxJava.

Therefore, the understanding and application of RxNet would follow the same principles and design patterns as observed in RxJava. The key advantage of

RxNet is that it provides .NET developers with an abstracted and simplified approach to asynchronous programming, similar to what the ReactiveX project offers in other languages.

Importantly, RxNet is fully integrated with other .NET asynchronous programming constructs such as Tasks and the `async/await` pattern, offering a robust and comprehensive toolset for addressing various asynchronous and event-based programming scenarios in the .NET framework.

2.3.2.2 Asynchronous Enumerables

In the .NET framework, the concept of an enumerable is represented through the `IEnumerable<T>` interface, which defines a method `GetEnumerator()`. This method returns an `IEnumerator<T>`, enabling iteration over a collection. To extend this concept to the asynchronous world, .NET introduces the `IAsyncEnumerable<T>` interface.

Similar to its synchronous counterpart, `IAsyncEnumerable<T>` returns an `IAsyncEnumerator` but with an asynchronous `MoveNextAsync()` method. This minor but significant modification lets us deal with data sources where data availability is asynchronous, such as real-time feeds, network streams, etc.

Here's a simple example of how asynchronous enumerables can be used in C#:

```
1  static async Task Main(string[] args)
2  {
3      IAsyncEnumerable<int> enumerable = FetchItems(1000);
4      int i = 0;
5      await foreach (int item in enumerable)
6      {
7          if (i++ == 10){ break;}
8          Console.WriteLine(item);
9      }
10 }
11
12 static async IAsyncEnumerable<int> FetchItems(int delay)
13 {
14     int count = 0;
15     while(true)
16     {
17         await Task.Delay(delay);
18         yield return count++;
19     }
20 }
21
22 //
23 //1
24 //1 sec delay
```

```
25 //2
26 //1 sec delay
27 //3
28 //....
29
30
```

Listing 2.13: Asynchronous Enumerable, C# example

This example demonstrates how asynchronous data sources can be worked with in a similar way as synchronous collections, thanks to the use of `IAsyncEnumerable<T>` and the `await foreach` construct.

2.3.3 Javascript

As a functional and dynamically-typed language, JavaScript provides a distinctive approach to managing asynchronous data flow. Enabled by its asynchronous runtime environment, Node.js, JavaScript employs specific intrinsic keywords and libraries to facilitate asynchronous operations. JavaScript's native constructs, namely the `async...await` and `for await...of` keywords, significantly ease the handling of asynchronous tasks. Beyond these built-in facilities, libraries like the Reactive Extensions for JavaScript (RxJS) enhance these capabilities further, enabling more sophisticated operations on asynchronous data streams. This section explores both JavaScript's intrinsic keywords and the RxJS library, emphasizing their respective roles in handling asynchronous flows in JavaScript.

2.3.3.1 Javascript Intrinsic Words: `async...await` and `for await...of`

The JavaScript runtime environment, Node.js, adopts a functional and weakly-typed approach to asynchronous flow processing. The mechanism to process asynchronous streams is somewhat similar to what we've seen in C#, but instead uses intrinsic keywords `async...await` and `for await...of`.

The keyword pair `async...await` provides a syntax that closely resembles synchronous code while making asynchronous calls. The `async` keyword marks a function as asynchronous and enables the use of the `await` keyword within it. The `await` keyword is then used before calling an asynchronous function, indicating that the function should pause and wait for the Promise to resolve or reject.

On the other hand, the keywords `for await...of` provide support for iterating over asynchronous enumerables, as demonstrated in the following example:

```
1      async function* streamAsyncIterableStream {
2      const reader = stream.getReader;
3      try {
4          while true {
5              const { done, value } = await reader.read;
6              if done {
7                  return;
8              }
9              yield value;
10         }
11     } finally {
12         reader.releaseLock;
13     }
14 }
15
16 async function getResponseSizeurl {
17     const response = await fetchurl;
18     let responseSize = 0;
19
20     const iterable = streamAsyncIterableResponse.body;
21
22     for await const chunk of iterable {
23         responseSize += chunk.length;
24     }
25     return responseSize;
26 }
27
```

Listing 2.14: Mozilla's Javascript Asynchronous Enumerables example

This example demonstrates how JavaScript's `async function*` construct can be used to define asynchronous enumerables. These can then be conveniently iterated over using the `for await...of` construct, just like you would with regular collections.

2.3.3.2 RxJS

RxJS represents the JavaScript adaptation of the ReactiveX project, analogous to RxNet in .NET and RxJava in Java. This library furnishes JavaScript developers with the same robust mechanisms and abstractions for handling asynchronous data streams.

While RxJS operates under similar principles as its Java and .NET counterparts, it introduces unique approaches that cater specifically to JavaScript's dynamic and functional nature. As such, the core philosophy remains consistent across

these libraries: simplifying the management of asynchronous data streams, with specific implementations nuanced to suit the distinct characteristics of their respective languages.

2.3.3.3 IxJS - Interactive Extensions for JavaScript

IxJS, also known as Interactive Extensions for JavaScript, is an integral part of the ReactiveX project and aims primary to enable the manipulation of data sequences, similarly to RxJS.

Similarly with another technologies already studied in this work, IxJS focuses on providing powerful abstractions for managing both synchronous and asynchronous sequences of data, while also aiming to be approachable and understandable. It has a robust suite of operators that you can use to write expressive, declarative code. Developers can also create custom operators very easily, extending the core functionality of IxJS to suit specific needs.

2.3.4 Technologies comparison

As we saw, each technology have a set of properties that help to characterize the solution. To help the characterization of each technology documented, we have a relation between the characteristics saw in the chapter 2.2.

	Rx(JAVA/.NET)	FLUX	FLOW	C# async enums	IxJS
Pull		x	x	x	x
Push	x	x			
Cancelable	x	x	x	x	x
Error Handling	x	x	x	x	x
Backpressure	x	x	x		
Intrinsic words			x	x	x



Benchmarking Use Case Algorithms and Benchmark Setup

In order to have metrics to enable the comparison between some of the different technologies we discussed so far, we decided to make software implementations in three programming languages of two key algorithms.

Each selected to simulate two distinct data processing scenarios. These are the "Count each word occurrence" and "Find the biggest word", both of which offer valuable insights into different aspects of data processing.

As a crucial part of this research, the selection of an appropriate and robust data source was imperative. For the empirical analysis and demonstration of the programming paradigms discussed, the "Gutenberg" library was chosen as the primary data source. This library, renowned for its comprehensive dataset, encompasses several thousands of books spanning a multitude of genres and periods. The richness and diversity of the Gutenberg dataset provided an unparalleled opportunity to rigorously test and evaluate the effectiveness of reactive programming techniques in processing and analyzing large volumes of textual data. The expansive nature of the Gutenberg library not only facilitated a broad assessment across various data-intensive scenarios but also underscored the scalability and adaptability of the programming approaches under study.

On the algorithms, in the case of "Group Word" operation, the idea was to have a memory-intensive task. This operation involves analyzing a dataset to identify

and count each word occurrences that fall within a specified size range. Since the algorithm uses data structures, such as dictionaries, for in runtime data storage. These nuances make the "Group Word" operation a curious test case for evaluating the capabilities of asynchronous processing in combination within a memory demanding scenario.

In contrast, the task of finding the largest word, though less memory intensive task, is significant for its simplicity. This operation involves scanning a dataset to find the single longest word, a process that, while straightforward, is critical for understanding the performance of asynchronous processing in more basic computational tasks. It serves as a benchmark for evaluating the efficiency of simpler algorithms and their implementation across different programming languages and environments.

For this work, were chosen these two algorithm as simplistic samples for a simplistic benchmark scenario, but, for future work, can be added a more vast algorithmic variety.

Bellow, we can find the pseudocode of these algorithms, that may help to understand how they function:

3.1 Overview of Algorithms

Pseudocode for the Group Word Operation

```

1 FUNCTION GroupWords(folder, minLength, maxLength)
2   Create an empty map 'wordMap'
3   FOR each file in 'folder' DO
4     Skip the first 14 lines of the file (These are typically
      metadata in Gutenberg project files)
5     FOR each remaining line in 'file' DO
6       IF the line contains "*** END OF" THEN
7         Break (This is the end of the actual content in
          Gutenberg project files)
8       END IF
9       FOR each word in 'line' DO
10        IF length of 'word' is between 'minLength' and '
          maxLength' THEN
11          Increment the count of 'word' in 'wordMap'
12        END IF

```

```
13         END FOR
14     END FOR
15 END FOR
16 RETURN 'wordMap'
17 END FUNCTION
```

Listing 3.1: Pseudocode for GroupWords function

Pseudocode for the Find Biggest Word Operation

```
1 FUNCTION FindLargestWord(folder)
2   Set 'largestWord' as an empty string
3   FOR each file in 'folder' DO
4     Skip the first 14 lines of the file (These are typically
      metadata in Gutenberg project files)
5     FOR each remaining line in 'file' DO
6       IF the line contains "*** END OF" THEN
7         Break (This is the end of the actual content in
          Gutenberg project files)
8       END IF
9       FOR each word in 'line' DO
10        IF length of 'word' is greater than length of '
largestWord' THEN
11          Set 'largestWord' as 'word'
12        END IF
13      END FOR
14    END FOR
15  END FOR
16  RETURN 'largestWord'
17 END FUNCTION
```

Listing 3.2: Pseudocode for FindLargestWord function

To benchmark each algorithm, was chosen to make a baseline implementations, which the code resembles the syntax complexity close to the pseudocode already shown and, on the other side, was used pipeline chains, using fluent API provided, for example, by RxJava or C# in AsyncEnumerable.

The idea behind this, was to explore readability vs performance in the use of these technologies.

3.2 Pipelines

During the implementation of the benchmarking algorithms, was made a conscious effort to keep the operation pipelines as similar as possible across the different technologies for each algorithm and its syntax closest as possible to what is shown in the algorithms pseudocode. This endeavor aimed to create a fair

and representative evaluation of the behaviors of each technology. By maintaining consistency in pipeline operations, we can more accurately attribute performance differences to the underlying technology, rather than variations in the implemented code. This approach brings us closer to a true comparison of how each technology handles the challenges of asynchronous I/O data retrieval and processing.

In certain instances, particularly with Java, there was a need to incorporate external libraries for non-blocking asynchronous file retrieval. Because, while some environments have native functions that already enable non-blocking IO operations to retrieve data from files, other, like in JAVA, was needed to use external libraries.

3.2.1 Understanding General Pipeline Operations

In the realm of functional programming and modern software development, several key operations are fundamental to process data collections and data flows. These operations, while conceptually similar across different languages, may have different idioms. Below there are explained several key operations that are usually used as operation in several fluent API's to construct operation pipelines.

Sort

The sort operation arranges elements of a collection in a specific order, typically ascending or descending. It's crucial for organizing data in a meaningful way.

In C#: Implemented as `.OrderBy (ascending)` and `.OrderByDescending (descending)`.

In Java: Executed using `.sorted()` with comparators for custom sorting.

In JavaScript: Uses `.sort()`, which can take a comparison function for custom sorting.

Distinct

The distinct operation removes duplicate elements from a collection, ensuring each element is unique.

In C#: Available as `.Distinct()`.

In Java: Achieved using `.distinct()` in the Stream API.

In JavaScript: Typically done using `new Set([...array])` to remove duplicates.

GroupBy

The `groupBy` operation groups elements of a collection based on a specified key. This is useful for categorizing data.

In C#: Done using `.GroupBy`.

In Java: Performed using `.collect(Collectors.groupingBy())`.

In JavaScript: Often implemented with `Array.prototype.reduce()` for custom grouping logic.

Concat and Union

Concat combines two sequences end-to-end, while union merges two sequences and removes duplicates.

In C#: `Concat` for concatenation and `Union` for union operations.

In Java: `.concat` for concatenation and using `.distinct()` after `.concat` for union.

In JavaScript: `.concat()` for concatenation; a combination of `.concat()` and `new Set()` for union.

Any and All

Any checks if any elements in the collection satisfy a condition, while All checks if all elements meet a condition.

In C#: Implemented as `.Any()` and `.All()`.

In Java: `.anyMatch()` and `.allMatch()` in the Stream API.

In JavaScript: `.some()` for any, and `.every()` for all.

Count and Sum

Count returns the number of elements, and Sum calculates the total of the numeric elements in a collection.

In C#: `.Count()` for counting and `.Sum()` for summing.

In Java: `.count()` and using `.mapToDouble()` followed by `.sum()` for summing.

In JavaScript: `.length` for count and `.reduce()` for summing.

Flat and FlatMap

Flat merges all sub-array elements into a new array, while FlatMap first applies a mapping function to each element and then flattens the result into a new array.

In C#: `.SelectMany()` is a close equivalent to FlatMap.

In Java: `.flatMap()` in the Stream API.

In JavaScript: `.flat()` for flat and `.flatMap()` for flatMap.

Zip

Zip combines elements of two collections into pairs or tuples, often based on their position in the collection.

In C#: Available as `.Zip()`.

In Java: Achieved using a combination of `.stream()` and `.map()` with a custom zipper function.

In JavaScript: Implemented manually using methods like `.map()` with additional logic for pairing.

Filter

The filter operation evaluates each element against a predicate (a true/false function) and includes only those elements that satisfy the predicate.

In C#: Implemented as `.Where()`.

In Java: Known as `.filter()`.

In JavaScript: Executed using `.filter()`.

ForEach

The ForEach operation applies a given action to each element in a collection. It's typically used for invoking side effects or operations on each element.

In C#: Available as `.ForEach` in the `List<T>` class or via looping constructs like `foreach`.

In Java: Implemented using `.forEach` in the Stream API or via looping constructs like `for-each`.

In JavaScript: Executed using `.forEach()`, which applies a function to each element of an array or array-like objects.

Map

The map operation applies a function to each element in a collection, transforming them into a new form. It's a cornerstone of functional programming, enabling easy data transformation.

In C#: Named as `.Select`.

In Java: Referred to as `.map`.

In JavaScript: Uses `.map()`, which applies a specified function to each element of an array, returning a new array with the transformed elements.

The consistent use of these operations across different languages underscores a universal shift towards more declarative and expressive programming styles, where the focus is on what needs to be done rather than how to do it. This approach not only enhances code readability and maintainability but also allows for more concise and functional solutions to common programming tasks.

3.2.2 Pipeline Examples

Bellow, we can see the pipeline versions of our key algorithms in C# and Java against their pseudocode representations already shown. The pseudocode, which is closer to a baseline implementation, serves as a point of reference to appreciate the enhanced readability and conciseness offered by pipelines. By comparing these implementations, we can better understand how pipelines abstract complexity and streamline data processing tasks.

C# Pipeline Example

```
1 private Task<string> parseFileDistinctWordsIntoDictionary(string
   filename)
2 {
3     return FileUtils.GetLinesAsyncEnum(filename)
4         .Where(line => line.Length != 0)
5         .Skip(14)
6         .TakeWhile(line => !line.Contains("*** END OF "))
7         .Select(line => Regex.Replace(line, "[^a-zA-Z0-9 -]+", "",
   RegexOptions.Compiled)
8         .Split(' '))
9         .Max(arr => arr) // Find the longest word
10        .AggregateAsync(string.Empty, (biggest, current) => current.
   Length > biggest.Length ? current : biggest)
11        .AsTask();
12 }
```

Listing 3.3: C# Pipeline for Parsing Distinct Words into a Dictionary

Java Pipeline Example

```
1 Files.list(Paths.get("path/to/directory"))
2     .filter(Files::isRegularFile) // Process only regular files
3     .map(file -> EXEC.submit(() -> biggestWord(file))) // Submit a task
   for each file
4     .collect(toList()) // Collect futures into a list
5     .stream()
6     .map(future -> waitForFuture(future).get()) // Process each future
7     .reduce((biggest, curr) -> curr.length() > biggest.length() ? curr
   : biggest) // Reduce to find the largest word
8     .get();
```

Listing 3.4: Java Pipeline for Finding the Largest Word in Files

These C# and Java implementations showcase the power and elegance of pipeline processing in handling tasks that would otherwise require more verbose and complex code.

4

Detailed Analysis and Comparison of Results

In this chapter we present the benchmark results on different technologies and strategies on the two chose algorithms previously mentioned in previous chapter, as detailed in the preceding chapters. Our investigation spanned various strategies and implementations, culminating in a nuanced understanding of the relative performance and characteristics of different approaches across .NET, Java, and other programming environments. Here, we consolidate these findings, reflecting on their broader implications within the field of asynchronous programming. Furthermore, we outline potential avenues for future research that emerge from our study, paving the way for continued exploration and innovation in this domain.

4.1 .NET Benchmarking

In this subsection, we focus on the benchmark using different strategies in .NET programming environment.

4.1.1 Find the biggest word algorithm results

For the find the biggest word algorithm, used i .NET implementations are the following strategies:

- **Baseline:** This strategy serves as the basic approach for finding the biggest word and acts as a baseline for comparison.
- **Asynchronous Baseline :** This strategy uses a single asynchronous task to find the biggest word.
- **Parallel :** An approach that makes the use of parallelization.
- **AsyncEnumerable:** This approach uses asynchronous programming with enumerable collections.
- **RxNet :** This strategy uses the Reactive Extensions (Rx) library with asynchronous file reading operations.

In the following graphic, we have the results in seconds for each strategy:

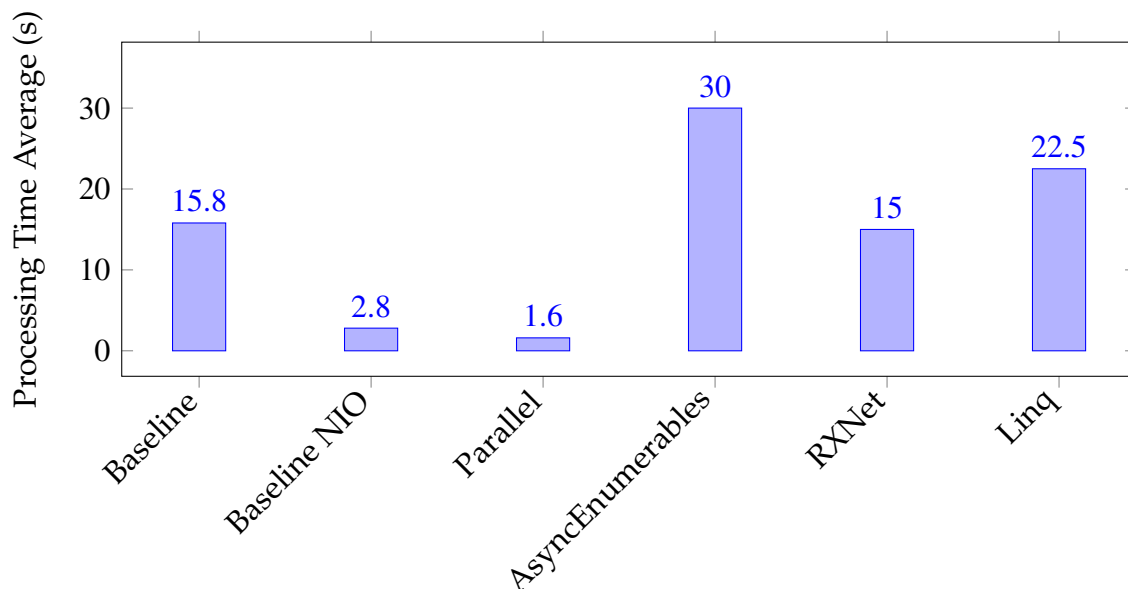


Figure 4.1: Processing times for different strategies for "Find the biggest word".

From these results, it is apparent that the asynchronous baseline approach has an advantage over its blocking counterpart. However, it is also evident that the parallel solution holds an advantage over its pipeline competitors.

Taking into account the pipeline technologies, we can see that the overall performance is worse compared to strategies that do not use pipelines. In this context, RxNet demonstrates a performance advantage over both Linq and AsyncEnumerables strategies. However, it was expected that, since Linq uses a blocking IO source, it would perform worse than AsyncEnumerables, but that does not occur. This may imply that for this particular algorithm, the AsyncEnumerable strategy incurs a higher overhead than intended.

4.1.2 Group Word Results

In this subsection, we concentrate on the task of grouping words from a file using different strategies. The strategies that we evaluate here include:

- **Baseline:** This strategy serves as the basic approach for word grouping and acts as a baseline for comparison.
- **Linq:** Like in the find word algorithm, this strategy uses Language Integrated Query (LINQ) in a synchronous manner.
- **AsyncEnumerable:** This strategy uses .NET async enumerables to process the asynchronous data .
- **RxNet:** Here, the Reactive Extensions (Rx) library is used to handle data sequences asynchronously and event-based.

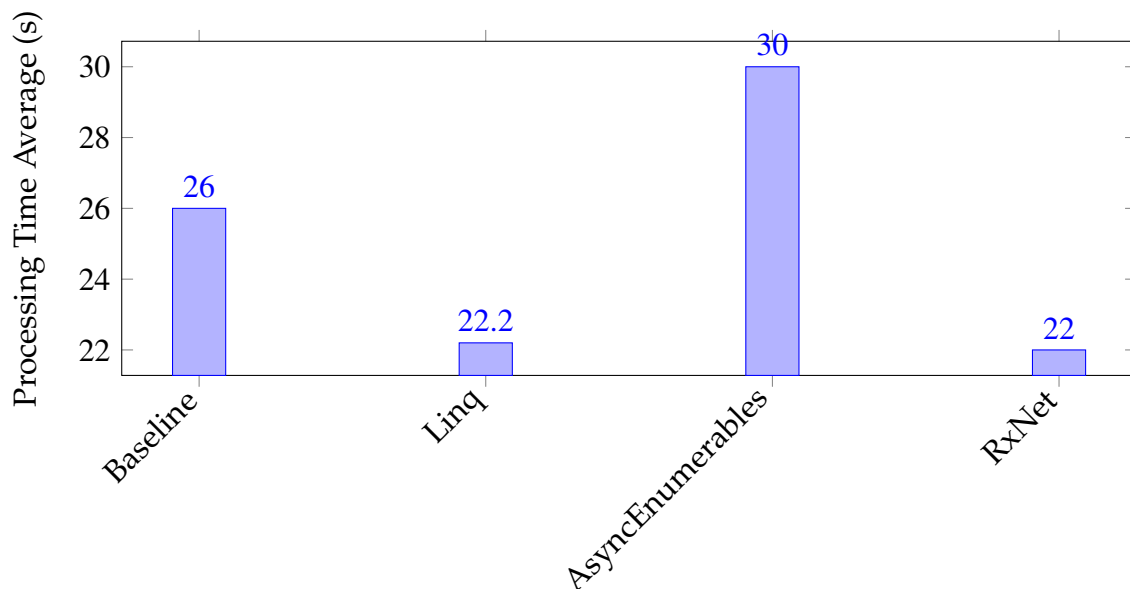


Figure 4.2: Processing times for different strategies for "Group Words".

The baseline strategy serves as a fundamental approach to word grouping and establishes a standard for performance comparison. While straightforward, it does not leverage the advanced features of other approaches, resulting in a moderate processing time. In comparison, the Linq strategy, utilizing Language Integrated Query in a synchronous manner, shows an improvement in performance over

the baseline. This suggests that the streamlined querying capabilities of LINQ, despite being synchronous, can efficiently handle the algorithm's requirements.

On the other hand, the AsyncEnumerable approach, which uses .NET async enumerables to process asynchronous data, exhibits a longer processing time in this context. This indicates that while async enumerables are beneficial for certain applications, their performance in memory-intensive tasks like "Group Words," which heavily relies on dictionaries, might not be optimal.

Furthermore, the RxNet implementation, employing the Reactive Extensions library to handle data sequences asynchronously and event-based, also demonstrates competitive performance. This underscores the potential of RxNet in efficiently managing asynchronous data flows, especially in scenarios where the processing time is crucial.

Overall, while pipeline libraries tend to outperform the baseline approach in memory-intensive algorithms like "Group Words," the choice between blocking and non-blocking IO operations becomes less significant. This is primarily because the overhead from memory-intensive operations tends to overshadow any potential performance gains from using non-blocking IO operations. The results indicate that the optimal strategy for such algorithms would depend more on how they handle memory-intensive operations rather than the differences in IO operation types.

4.2 Java/Kotlin Benchmarking

In this section, we explore and assess diverse strategies applied in Java and Kotlin to process files, and we scrutinize their performances solving the "Find Word" and "Group Word algorithms"

4.2.1 Biggest Word Results

The strategies used in JAVA are:

- **Baseline:** This strategy illustrates a basic non-blocking I/O operation, serving as a comparison baseline.
- **Flux:** These strategies leverage the Reactor Flux model from Java's Project Reactor library. The former follows a standard non-concurrent processing model, while the latter introduces parallelization for improved performance.
- **RXJava:** This strategy employs the RXJava library. They replace the Reactor Flux with Observables, with the distinction being made between non-concurrent and concurrent processing.
- **Streams and parallelization:** Implementation of three strategies that use Java's Streams API and explore handling of blocking operations under three different conditions: standard usage, raw multithreading using threadpools and using parallel method in the streams API.

In the following graphic, we have the results in seconds for each strategy:

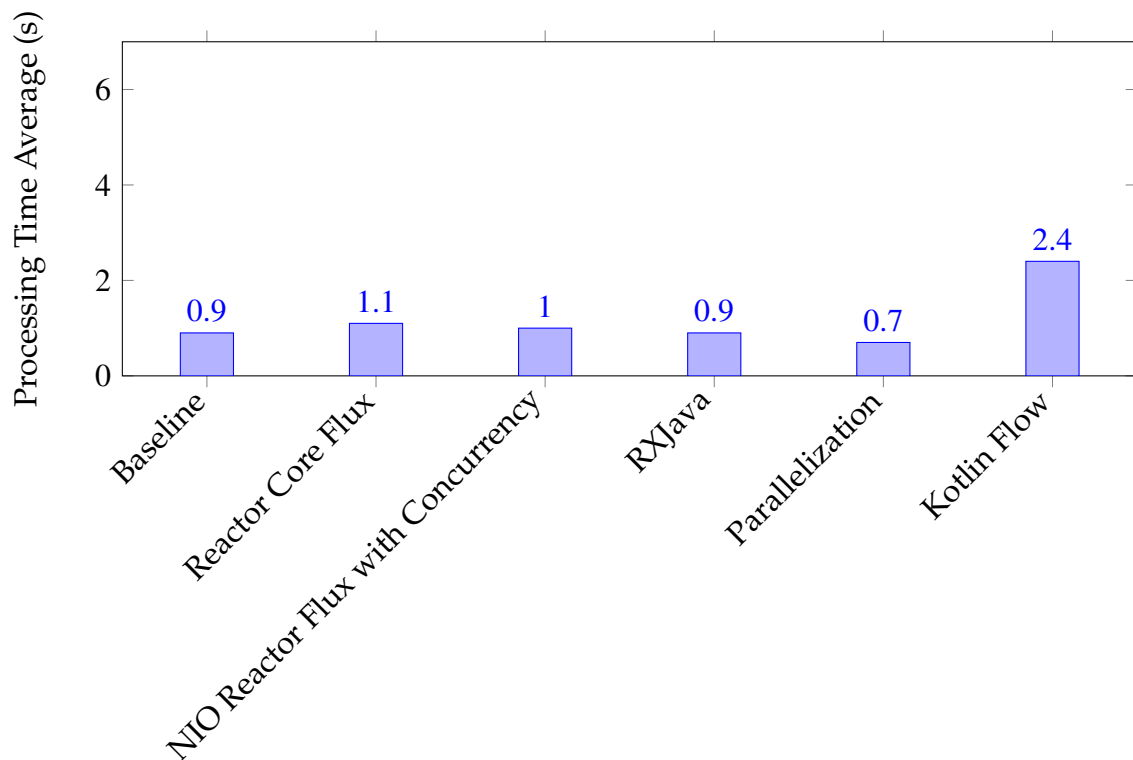


Figure 4.3: Processing times for different Java/Kotlin strategies for "Biggest Word".

For the "Biggest Word" results, the strategies exhibit varied performance levels. The Baseline strategy, illustrating basic non-blocking I/O operations, provides a solid foundation for comparison and demonstrates a competent processing time. Flux strategies, both Reactor Core Flux and its concurrent variant, leverage the Project Reactor library to enhance processing, showing slight variations in efficiency. The RXJava strategy, utilizing the RXJava library, aligns closely with the baseline in terms of performance. A notable distinction is seen in the strategies employing parallelization, particularly those using Java's Streams API, which show the most efficient processing times. This underscores the advantage of parallel processing in optimizing performance.

4.2.2 Group Word Results

The strategies discussed here include:

- **Baseline:** This strategy serves as a baseline for comparison. It illustrates a basic non-blocking I/O operation without the use of any high-level constructs like Reactor Flux or Observable.
- **RXJava:** This strategy employs the RXJava library, popular for building asynchronous and event-driven applications.
- **Reactor Core Flux:** This strategy leverages the Reactor Flux model available in Java's Project Reactor library, providing an efficient approach to handling asynchronous data sequences.
- **Parallel:** This strategy uses Java's Streams API and explores handling of blocking operations with the help of threadpools.
- **Flow (Kotlin):** This strategy utilizes Kotlin's Coroutines and Flow API, which are particularly well-suited for handling multiple values that are emitted sequentially.

In the following table and graphic, we have the results in seconds for each strategy:

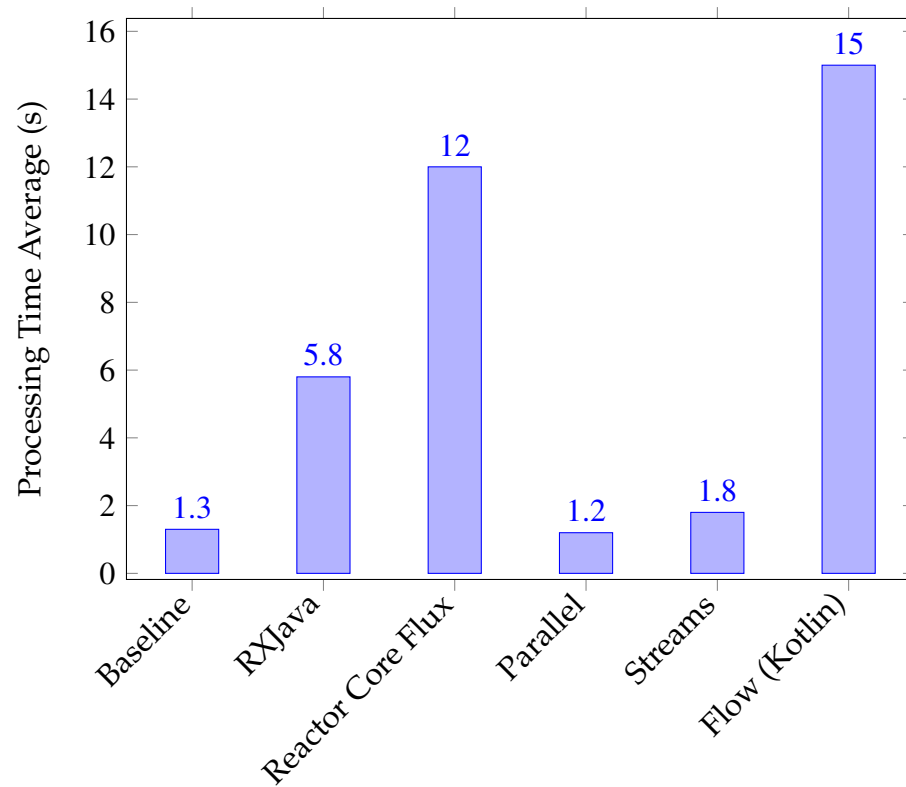


Figure 4.4: Processing times for different Java/Kotlin strategies for "Group Words".

On "Group Word" benchmark results, the strategies again present a range of performances. The Baseline strategy, as before, sets a standard for comparison, showing moderate efficiency. RxJava, known for its asynchronous and event-driven capabilities, demonstrates a longer processing time in this context, indicating potential overheads in handling complex data sequences. Reactor Core Flux, while efficient in handling asynchronous sequences, shows a significant processing time, possibly due to the complexity of the Group Word algorithm. Strategies involving parallel processing, whether through direct use of threadpools or the Streams API, exhibit improved efficiency, highlighting the benefits of parallelism in handling computationally and memory intensive tasks. Notably, Kotlin's Flow API, designed for handling sequential emissions in a coroutine-based environment, shows a longer processing time, which could reflect the overheads associated with coroutine management in complex data processing scenarios.

Overall to JAVA, these results illustrate the trade-offs between different programming paradigms and libraries in Java and Kotlin. While some strategies excel in simple tasks, others are more suited to complex algorithms, indicating the importance of choosing the right approach based on the specific requirements of the algorithm and the nature of the data being processed.

4.3 JavaScript Benchmarking

4.3.1 Finding the Biggest Word

Here, we investigate the following three strategies:

- **Baseline:** This strategy serves as the basic JavaScript approach for finding the biggest word, acting as a baseline for comparison.
- **Stream:** This strategy uses JavaScript streams, which provide a way to handle reading/writing files, network communications, or any kind of end-to-end information exchange in an efficient manner.
- **RxJS:** This strategy leverages the Reactive Extensions for JavaScript (RxJS) library, which offers a set of methods for dealing with asynchronous data sequences in an effective way.

In the following graphic, we have the results in seconds for each strategy:

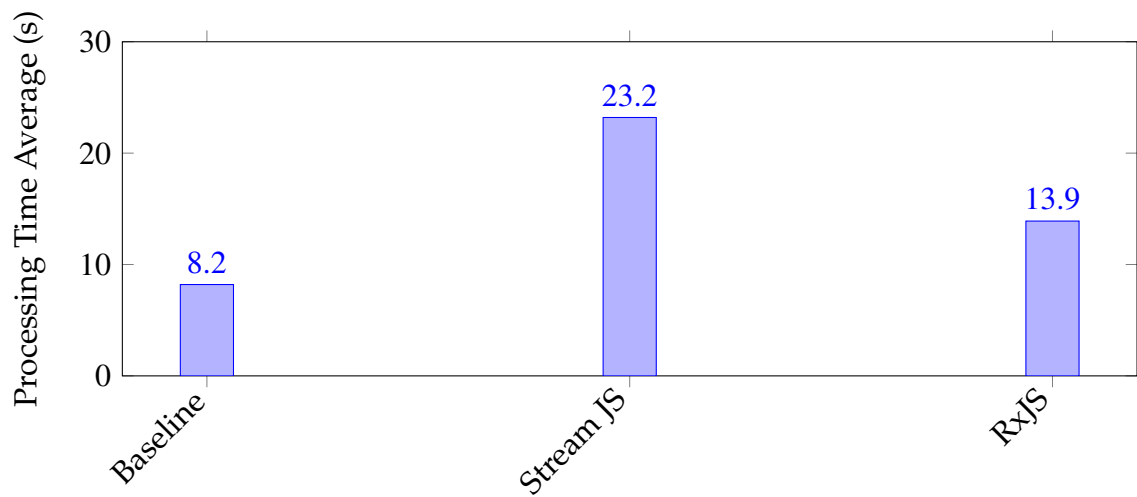


Figure 4.5: Processing times for different JavaScript strategies for "Biggest Word"

As we can see in this javascript implementation, the performance for find biggest word algorithm of the baseline behaves as expected, having the best performance among its alternatives. On the other side, the algorithm implemented using RXJS library behaves slightly better than the alternative that makes use of JS streams API.

4.3.2 Grouping Words

In this subsection, we evaluate the following strategies for grouping words:

- **Baseline:** This strategy serves as the basic JavaScript approach for grouping words, acting as a baseline for comparison.
- **Stream JS:** This strategy uses JavaScript streams, which provide a way to handle reading/writing files, network communications, or any kind of end-to-end information exchange in an efficient manner.
- **RxJS:** This strategy leverages the Reactive Extensions for JavaScript (RxJS) library, which offers a set of methods for dealing with asynchronous data sequences in an effective way.

In the following graphic, we present the results in seconds for each strategy:

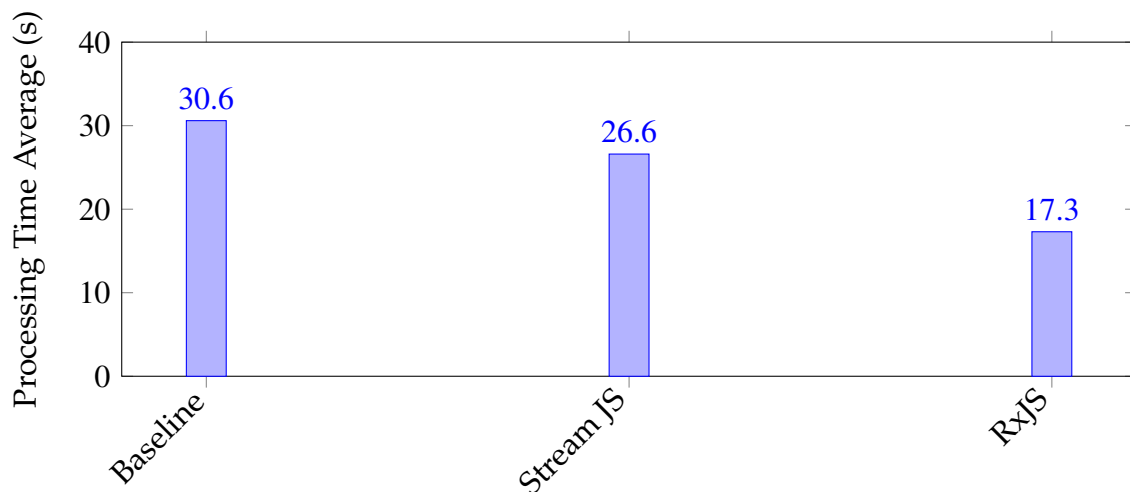


Figure 4.6: Processing times for different JavaScript strategies for "Grouping Words"

In this case, the performance of the algorithms change radically from what we saw in the "Find the biggest word" algorithm. In this case, the baseline has the worst performance among the three strategies, while the RxStrategy has the best. One explanation for these results is that the pipeline instructions of the pipeline instructions are optimized while the baseline is not.

Conclusions and Future Directions

5.1 Main Conclusions

The analysis of reactive streams programming and pipelining technologies throughout this thesis highlights their substantial benefits in terms of user-friendliness and maintenance simplicity. Nonetheless, these technologies also present certain drawbacks which must be carefully considered in the context of specific application requirements and desired outcomes.

For systems demanding high performance and robustness, a baseline-like approach, which allows for custom optimization, might be ideal. Conversely, for rapid development with satisfactory performance and easier maintenance, tools like RxJava and .NET's AsyncEnumerable provide valuable solutions. The choice largely depends on the project's long-term objectives and operational context. For example, high-frequency trading systems require the raw speed and low latency that finely tuned baseline systems provide. In scenarios where time-to-market is crucial, the productivity benefits of higher-level frameworks may outweigh their performance costs.

Another significant finding is the scalability offered by reactive programming technologies, which is essential for applications expected to handle growing data volumes or increased throughput over time.

For each language we took, some conclusions:

5.1.1 .NET Environment Analysis

In the .NET ecosystem, strategies utilizing parallel computing clearly outperformed others, emphasizing the importance of leveraging modern CPU architectures to enhance processing capabilities.

5.1.2 Java/Kotlin Environment Insights

In Java and Kotlin environments, parallelization strategies proved most effective, highlighting the necessity of concurrent processing for handling complex or large data sets efficiently.

5.1.3 JavaScript Strategy Performance

The RxJS library demonstrated its efficacy in JavaScript, particularly for tasks like "Grouping Words," showcasing its capability to manage data streams efficiently—a crucial aspect for web applications.

5.1.4 Overarching Conclusions

While baseline approaches generally performed better in specific algorithms, the practicality of using advanced frameworks like RxJava or .NET's AsyncEnumerable in enterprise environments cannot be overlooked. They facilitate rapid development and easier maintenance, which are often prioritized over minor performance gains in many professional settings.

Moreover, the selection of the appropriate technology or framework is invariably influenced by the specific demands of the project, underscoring the need for a tailored approach in software development.

5.2 Future Work and Final Thoughts

This thesis underscores that asynchronous I/O technologies serve not only as technical tools but also as catalysts for broader transformations in software engineering. They hold the potential to redefine the paradigms of software development and maintenance.

As we advance into an increasingly digital era, the strategic use of these technologies—especially in managing complex data flows in real time already is critical. They are poised to enhance system robustness and reduce failure rates significantly and today its a critical concept that any IT engineer must take in account.