

A Dynamic Pipeline Framework implemented in Haskell*

Juan Pablo Royo Sales^a, Edelmira Pasarella^a, Cristina Zoltan^a, Maria-Esther Vidal^b

^a*Universitat Politècnica de Catalunya, 08034, Barcelona, Spain*

^b*TIB/L3S Research Centre at the University of Hannover, Hannover, Germany*

Abstract

Streaming processing has given rise to new computation paradigms to provide effective and efficient data stream processing. The Dynamic Pipeline Paradigm is a computational model for stream processing. In particular, this paradigm is suitable to solving problems where the incremental emission of results is a critical issue. In the implementation of a Dynamic Pipeline Framework computations must correspond to natural primary entities. This fact suggests that a proper implementation language must allow for manipulating functions as first citizens as implementation language. Haskell accomplishes this requirement. It is a pure functional programming language relaying on solid theoretical foundations that provides the possibility of manipulating computations as primary entities. Moreover, from a practical point of view, it has a robust set of tools for writing multithreading and parallel computations with optimal performance. In this work we tackle the problem of specifying and developing a general and algorithm-parametric Dynamic Pipeline Paradigm using Haskell as development language. Additionally, we conduct, analyze and report experiments to measure the performance of computing the weakly connected components of large networks using the Dynamic Pipeline Framework. To assess the incremental delivery of results we measure the Diefficiency metrics, i.e. the continuous efficiency of the implementation of an algorithm for generating incremental results. Obtained results satisfy our expectations and encourage us to keep using the Dynamic Pipeline Framework for solving some graph stream processing problems where is critical not to have to wait until the whole results are emitted.

Keywords: Stream Processing Frameworks, Dynamic Pipeline, Parallelism, Concurrency, Haskell

Contents

1	Introduction	2
2	Preliminaries	3
2.1	Streaming Processing	3
2.2	Dynamic Pipeline Paradigm	4
2.3	Streaming in Haskell Language	5
2.4	Diefficiency Metrics	5

*This work was partially supported by MCIN/ AEI /10.13039/501100011033 [grant number PID2020-112581GB-C21].

3	Dynamic Pipeline Framework in Haskell	6
3.1	Framework Design	6
3.1.1	Background	6
3.1.2	Architectural Design	6
3.2	Implementation	8
3.2.1	DSL Grammar	8
3.2.2	DSL Validation	9
3.2.3	Interpreter of DSL (IDL)	11
3.2.4	Runtime System (RS)	16
3.3	Libraries and Tools	18
3.4	Parallelization	18
3.4.1	Channels	19
4	Enumerating Weakly Connected Component on the DPF	19
4.1	DP _{WCC} Algorithm	19
4.2	DP _{WCC} Implementation	24
4.3	Empirical Evaluation	28
4.4	Running Architecture	28
4.5	Haskell Setup	28
4.6	DataSets	28
4.7	Experiments Definition	29
5	Discussion of Observed Results	29
5.1	Experiment: E1	29
5.2	Experiment: E2	30
5.3	Experiment: E3	32
6	Related Work	34
7	Conclusions and Ongoing Work	35
	Bibliography	35

1. Introduction

Effective streaming processing of large amounts of data has been studied for several years [2, 4, 7] as a key factor providing fast and incremental results in big data algorithmic problems. One of the most explored techniques, regardless of the approach, is the exploitation of parallel techniques to take advantage of the available computational power as much as possible. In that regard, the Dynamic Pipeline Paradigm (DPP) [15] has lately emerged as one of the models that exploit data streaming processing using a dynamic pipeline parallelism approach [7]. This computational model has been designed with a functional focus, where the main components of the paradigm are functional stages or pipes which dynamically enlarge and shrink depending on incoming data. The details of the specific requirements of the Haskell system according to the results of the proof of concept will be presented in section 3.

One of the biggest challenges of implementing a Dynamic Pipeline Framework (DPF) is to find a proper set of tools and programming language which can take advantage of both of its primary aspects: i) *fast parallel* processing and ii) *strong theoretical* foundations that manage computations

as first-class citizens. Haskell Programming Language (Haskell) is a statically typed pure functional language that has been designed and evolved from its birth in 1987, on strong theoretical foundations where computations are primary entities, and at the same time has been providing a powerful set of tools for writing multithreading and parallel programs with optimal performance [9, 10].

Problem Research and Objective: The main objective of this work is to explore the feasibility of using a Functional Programming (FP) language to implement a DPF. In particular, we tackle the problem of establishing the basis of an implementation of a DPF in Haskell, a pure functional language. This is, our aim is to determine the particular features (i.e., versions and libraries) of this language that will allow for an efficient implementation of the DPF. To be concrete, through a particular and very relevant problem as the computation of the Weak Connected Components (WCC) of a graph, we study the critical features required in Haskell for a DPF implementation.

Contributions: A proof of concept of the implementation of a DP for WCC using Haskell; the results of the empirical study suggests that Haskell is a suitable language for implementing DPP. This work also contributes to building the first abstraction approximation of a future framework/library of this computational model in that language.

The rest of this paper is organized as follows. Next section presents the basic notions used through this work. In ?? the proof of concept is analyzed. To be concrete, in this section an algorithm for enumerating WCC using DPP and its implementation using parallel Haskell are introduced. In subsection 4.3, the experiments conducted to assess how Haskell supports the dynamic pipeline implementation for enumerating WCC are described and their results reported. In section 3 the design of the Haskell Dynamic Pipeline Framework and the most relevant details of its implementation are deeply explained. In particular, all the Haskell data types and language techniques used in the implementation are detailed. The justification of the used external libraries for the runtime system are presented at the end of this section. Section 4 the results of the experiments conducted to evaluate the performance of using the DPF are presented. In 6 we present the related work and, finally, conclusions and further work are presented in section 7. Part of the content presented through this work has been published in [14].

2. Preliminaries

Before moving forward with the core of the research, we describe the fundamental concepts that support the different parts of our study. That is, stream processing in general, Dynamic Pipeline Paradigm, streaming processing in the context of Haskell and the metrics we use in the conducted experiments.

2.1. Streaming Processing

The development of streaming processing techniques have potentiated areas as massive data processing for data mining algorithms, big data analysis, IoT applications, etc. Data Streaming (DS) has been studied using different approaches [2? , 7] allowing to process a large amount of data efficiently with an intensive level of parallelization. We can distinguished two different parallelization streaming computational models: Data Parallelism (DAP) and Pipeline Parallelism (PP).

Data Parallelism (DAP). The data is split and processed in parallel and, the computations that perform some action over that subset of data do not have any dependency with other parallel computation. A common model that has been proved successful over the last decade is MapReduce (MR) [?]. Different frameworks or tools like Hadoop [31], Spark [16], etc., support this computational

model efficiently. One of the main advantages of this kind of model is the ability to implement stateless algorithms. Data can be split and treated in different threads or processors without the need for contextual information. On the other hand, when there is a need to be aware of the context, parallelization is penalized, each computational step should be fully calculated before proceeding with the others. For example, this is the case of **reduce** operation on many of the above-mentioned frameworks or tools.

Pipeline Parallelism (PP). It break the computation in a series of sequential stage, where each stage takes the result of the previous stage as an input and downstream its results to the next. Each *pipeline Stage* is parallelized and, it could potentially exist one stage per data item of the stream. The communication between stages takes place through some means, typically channels. One of the main advantages of this model is that the stages are non-blocking, meaning that there is no need to wait to process all data to run the next stage. This kind of paradigm enables computational algorithms that can generate incremental results, preventing the user waits until the end of the whole data stream processing to get a result. On the other hand, the disadvantage over DAP is that although pipeline stages are parallelized, some intensive computation in one stage might delay processing the next stage because of its sequential dependency nature. Therefore, the user must be sure each stage runs extremely fast computations on it.

The nature of our problem requires that results are output incrementally, i.e. Bitriangles are emitted incrementally. Additionally, data need to be aware of the context to compute the BTs. Considering the stream processing models presented above, we have chosen PP. We think it is the model that better fits the requirements of the enumerating incrementally bitriangles problem. We are going to see in the next section what is the specific PP computation model used for that purpose.

2.2. Dynamic Pipeline Paradigm

The *Dynamic Pipeline Paradigm (DPP)* [15] is a PP computational model based on a one-dimensional and unidirectional chain of stages connected by means of channels synchronized by data availability. This chain of stages is a computational structure called *Dynamic Pipeline (DP)*. A DP stretches and shrinks depending on the spawning and the lifetime of its stages, respectively. Modeling an algorithmic solution as a DP corresponds to define a dynamic computational structure in terms of four kinds of stages: *Source* (Sr), *Generator* (G), *Sink* (Sk) and *Filter* (F) stages. In particular, the specific behavior of each stage to solve a particular problem must be defined as well as the number and the type of channels connecting them. Channels are unidirectional according to the flow of the data. The *Generator* stage is in charge of spawning *Filter* stage instances. This particular behavior of the *Generator* gives the elastic capacity to DPPs. *Filter* stage instances are stateful operators in the sense described in [13]. This is, *Filter* instances have a state. The deployment of a DP consists in setting up the initial configuration depicted in Figure 1.

The activation of a DP starts when a stream of data items arrives at the initial configuration of the DP. In particular, when a data stream arrives to the *Source* stage. During the execution, the *Generator* stage spawns *Filter* stage instances according to incoming data and the *Generator* defined behavior. This evolution is illustrated in Figure 2. If the data stream is bounded, the computation finishes when the lifetime of all the stages of DP has finished.

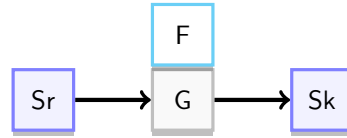


Figure 1: Initial configuration of a Dynamic Pipeline. An initial DP consists of three stages: Sr, G together its filter parameter F, and Sk. These stages are connected through its channels –represented by right arrows– as shown in this figure.

Otherwise, if the stream data is unbounded, the DP remains active and incremental results are output.

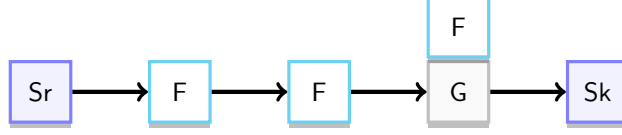


Figure 2: Evolution of a DP. After creating some filter instances (shadow Filter squares) of the filter parameter (light Filter square) in the Generator, the DP has stretched.

2.3. Streaming in Haskell Language

Streaming computational models have been implemented in Haskell Programming Language during the last 10 years. One of the first libraries in the ecosystem was `conduit` [20] in 2011. After that, several efforts on improving streaming processing on the language has been made not only at abstraction level for the user but as well as performance execution improvements like `pipes` [36] and `streamly` [41] lately. Moreover, there is an empirical comparison between those three, where a benchmark analysis has been conducted [18].

Although most of those libraries offer the ability to implement DAP and PP, none of them provide clear abstractions to create DPP models because the setup of the stages should be provided beforehand. In the context of this work, we have done a proof of concept at the beginning, but it was not possible to adapt any of those libraries to implement properly DPP. The closest we have been to implement DPP with some of those libraries was when we explored `streamly`. In this case, there is a `foldrS` combinator that could have been proper to the purpose of generating a dynamic pipeline of stages based on the data flow, but it was not possible to manipulate the channels between the stages to control the flow of the data. It is important to remark that, even though, the library `streamly` implements channels, they are hidden from the end-user and, there is not a clear way to manipulate them.

To the best of our knowledge no similar library under the DPP approach has been written in Haskell Programming Language. One important motivation to develop our own framework is that we not only wanted to satisfy our research needs but, as a novel contribution, we wanted to deliver a DPF to the Haskell community as well. We hope this contribution encourages and helps writing algorithms under the Dynamic Pipeline Paradigm.

2.4. Diefficiency Metrics

In proof of concept in ?? and in this work on the empirical analysis section 5, we use two important metrics to measure the diefficiency, i.e. continuous efficiency of a program to generate incremental results. The metrics to measure diefficiency are Diefficiency Metric `dief@t` (`dief@t`) and Diefficiency Metric `dief@k` (`dief@k`) [1]. The metric `dief@t` measures the continuous efficiency during the first t time units of execution regarding the results generated by the program. The higher value of the `dief@t` metric, the better the continuous behavior. The metric `dief@k` measures the continuous efficiency while producing the first k answers regarding the results generated by DP-BT-Haskell. The lower the value of the `dief@k` metric, the better the continuous behavior. Both metrics have been measured using `diefpy` Tool (`diefpy`) [23] and traces obtained by the execution of the experiment scenarios (Traces examples are provided in ??). Apart from this, `diefpy` generates two different kind of plots. On the one hand, a bi-dimensional plot containing all the (x, y) points taken from traces like ??, where each x is the t time where the answer was generated and the y is the generated answer

number. This plot is useful to have a visual view of the continuous behavior. On the other hand, a radial plot contains the visual comparison of `dief@t` metric with respects to other non-continuous metrics such as **i)** Completeness (Comp) which is the total number of answers produced by the scenario, **ii)** Time for the first tuple (TFFT) which measure the elapsed time spent by the scenario to produce the first answer, **iii)** Execution Time (ET) which measures the elapsed time spent by the scenario to complete the execution of a query and, **iv)** Throughput (T) which measure the number of total answers produced by the scenario after evaluating a query divided by its execution time ET.

3. Dynamic Pipeline Framework in Haskell

The design and implementation of Haskell Dynamic Pipeline Framework (DPF-Haskell) is a fundamental piece of the present work. A Dynamic Pipeline Framework written in Haskell Programming Language which allow Haskell users to implement any suitable algorithm for Dynamic Pipeline Paradigm. During the process of conducting this research, we have implemented DPF-Haskell [25] and publishes it into The Haskell Package Repository [30]. In this section, we describe the design and implementation details of DPF-Haskell.

3.1. Framework Design

3.1.1. Background

A suitable framework should provide the user the right level of abstraction that removes and hides underlying complexity, allowing the developer to focus on the problem that needs to be solved. There are several design approximations to implement a framework: i) *Configuration Based* where the user only focuses on completing a specific configuration either on a file or a database or both. Once this configuration is completed, the user provides it to the framework's runtime system in order to execute the program. An example of this could be WordPress [?], ii) *Convention over Configuration (CoC)* where the user writes his code and definition following certain patterns in naming or source code location. Using the source code and content-specific information, the framework interprets the execution flow that needs to be executed. This technique has been deeply explored in the last 10 years. One of the first framework that introduce this design paradigm was Ruby on Rails [?]. Other examples are Spring Boot and Cake PHP for example [? ?], iii) *Application Programing Interface (API)* where the framework or library provides a certain amount of functionality implemented in terms of functions or interfaces, and the user needs to compose those functions or implement those abstractions to achieve the desired results. This has been the traditional design paradigm for building any library or framework, and finally iv) *Domain-specific Language (DSL)* [3] where the framework or library provides a new language that represents the domain problem, encoding the solution in terms of that DSL language. An example of this type of design is Hibernate Query Language [?].

There exists two types of Domain-specific Languages [?]: External Domain-specific Language (DSL) and Embedded Domain-specific Language (EDSL). The purpose of DSL, is to create a completely new language with its own semantic, syntax, and interpreter. DSLs are not general-purpose languages, because as their name indicates, they are domain-specific. EDSL are syntactically embedded in the host language of the library, and the user writes in that host language, but restricted by the EDSL abstractions. DPF-Haskell follows a EDSL approach taking advantage of the strong type Haskell system providing correctness at type-level [6].

3.1.2. Architectural Design

In this section we focus on the architectural design of the DPF-Haskell using a EDSL approach. We have built a framework that contains three important components: DSL, IDL and RS.

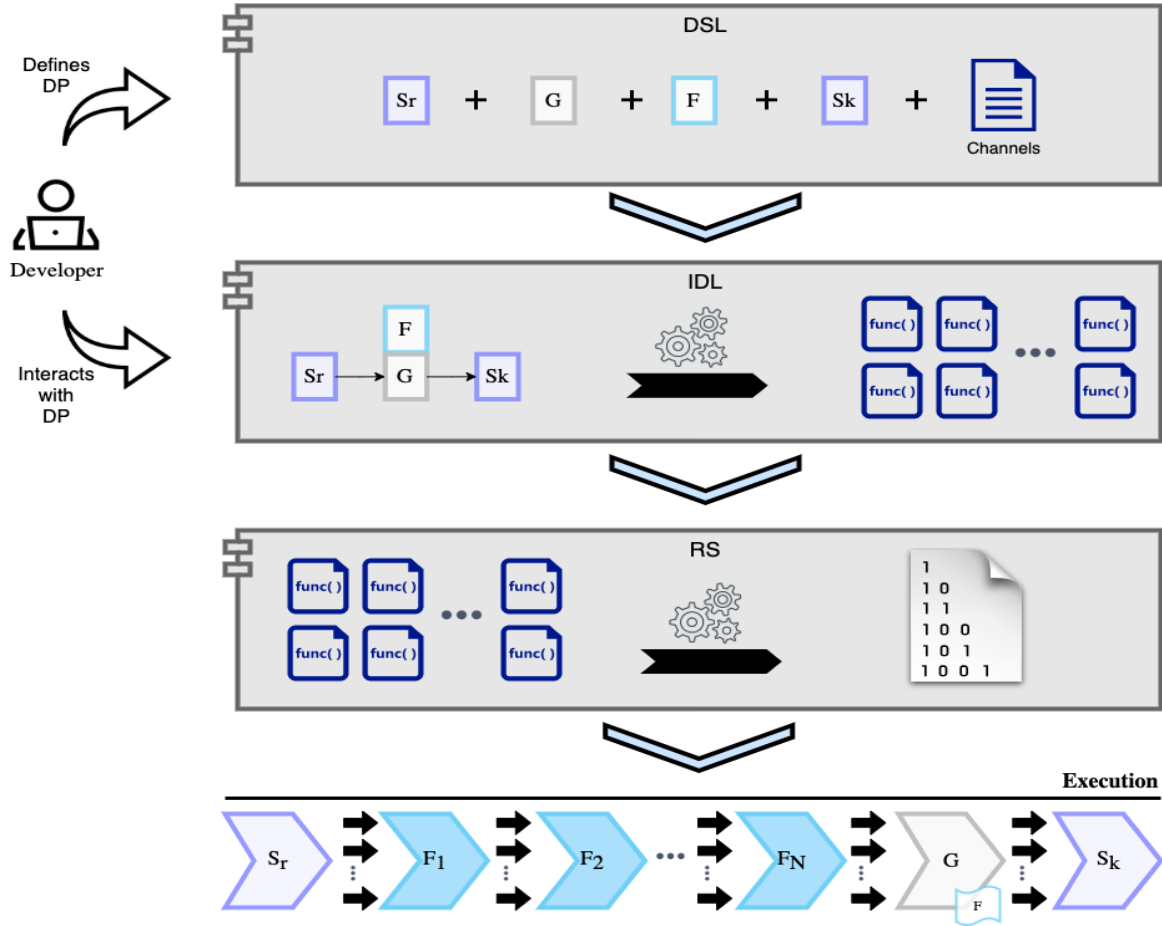


Figure 3: This diagram shows the architectural design of DPF-Haskell. DPF-Haskell is a DSL which is built on three main components: DSL, IDL and RS. In the DSL we can see how the user can compose the main stages of the DPP. IDL is showing how the frameworks is helping the user to transform that definition into real function or computations. Finally RS execute all that definition plus functions. Execution layer indicates an example of a DPP running after being executed.

In Figure 3 we can appreciate the different components mentioned before that are the grey boxes.

DSL. The user interacts with the DSL component where defines how the DPP flow should be. Defining the flow consists on to provide a type-level specification about the channels that communicate each stage of the pipeline, as well as the data types those channels carry. For example, in the case of ?? that we develop the WCC algorithm, the user knows stages S_r , G , and S_k need to be connected with two channels. One of those channels is carrying the edges – **Edge** data type – and the other the accumulated connected components – **ConnectedComp** data type –.

IDL. Based on the definition provided in the DSL, the user interacts with the IDL to build the functions with the algorithms needed for each stage: S_r , G , S_k , F , and actors.

RS. RS is fed with the DPP definition and the functions implementations to finally execute the program.

3.2. Implementation

In this section, we describe the implementation details of each architectural layer: DSL, IDL and RS. As we have explained in ??, this library was published on Hackage [25], the source code is open and can be found on this Github Repository [24].

3.2.1. DSL Grammar

In order to provide correctness verification at compilation level, we define a Context-Free Grammar (CFG) that generates a DPP DSL language. CFG enables the user to define a DPP at type-level.

Definition 1 (DSL CFG). Lets $G_{dsl} = (N, \Sigma, DB, P)$ be a Context-Free Grammar, such that N is the set of non-terminal symbols, Σ the set of terminal symbols, $DP \in N$ is the start symbol and P are the generation rules. Figure 4 shows the formal definition of the grammar.

$$\begin{aligned}
 N &= \{DP, S_r, S_k, G, F_b, CH, CH_s\}, \\
 \Sigma &= \{\text{Source}, \text{Generator}, \text{Sink}, \text{FeedbackChannel}, \text{Type}, \text{Eof}, :=, :<+>\}, \\
 P &= \{ \\
 &\quad DP \rightarrow S_r := G := S_k \mid S_r := G := F_b := S_k, \\
 &\quad S_r \rightarrow \text{Source } CH_s, \\
 &\quad G \rightarrow \text{Generator } CH_s, \\
 &\quad S_k \rightarrow \text{Sink}, \\
 &\quad F_b \rightarrow \text{FeedbackChannel } CH, \\
 &\quad CH_s \rightarrow \text{Channel } CH, \\
 &\quad CH \rightarrow \text{Type } :<+> CH \mid \text{Eof} \}
 \end{aligned}$$

Figure 4: This is the Context-Free Grammar defined for the DSL. In the first box we can see N which is the set of non-terminals symbols of the Grammar. Σ which is the set of the terminal symbols and P the production rules of the grammar.

For encoding G_{dsl} on the Haskell, we use an *Index type* [?] to keep track, at type-level, of the extra information required by the DPP definition such as channels and data types the channels carry.

```

1 data Source (a :: Type)
2 data Generator (a :: Type)
3 data Sink
4 data Eof
5 data Channel (a :: Type)
6 data FeedbackChannel (a :: Type)

```

Source Code 3.1: This code is showing most of the data types that represent the same terminal symbols $\Sigma \in G_{dsl}$. Those types that are indexed by another kind **Type**, allows to store information at type-level needed for interpret the DSL

In Source Code 3.1, there is an *Index type* for each element of Σ encoded in Haskell Types. The

highlighted lines in Source Code 3.1 shows the terminal symbols Σ that are not indexed, because neither **Sink** nor **Eof** are carrying extra type-level information. In the case of **Sink**, since it is the last stage that does not connect further with any other stage, we do not need to indicate any channel information. **Eof** it is just a terminal type to disambiguate the **Channel** ($a :: \text{Type}$) subtree for the full parser tree. **Channel** can carry any type because it needs to be polymorphic to support a different number of channels and data types.

```

1  data chann1 :<+> chann2 = chann1 :<+> chann2
2  deriving (Typeable, Eq, Show, Functor, Traversable, Foldable, Bounded)
3  infixr 5 :<+>
4
5  data a :=> b = a :=> b
6  deriving (Typeable, Eq, Show, Functor, Traversable, Foldable, Bounded)
7  infixr 5 :=>

```

Source Code 3.2: Special terminal symbols $\{<+>, :=>\} \in \Sigma$. This terminal symbols allows to index two types in order to combine several of them and build a chain of stages ($:=>$) and a set of channels ($:<+>$).

There are two important terminal symbols in Σ : $:=>$ and $:<+>$. In Source Code 3.2, the definition shows how $:=>$ and $:<+>$ can combine 2 (two) types. The propose of writing $:=>$ and $:<+>$ as types is to have a syntactic sugar type combinator for writing the DSL according to the CFG. Apart from that, they are different because two distinguishable terminal symbols Σ are needed to separate the encoding of the pipeline stage (Sr, G, Sk) from the encoding of channel composition in the same stage, as we can appreciate in Definition 1.

Now, we can start defining our pipelines at type-level. For example, if we want to generate a DPP that eliminates duplicated elements in a stream, we know that we only need one channel connecting the stages that carries out the type of the element, in this case, **Int** (see Source Code 3.3).

```

1  type DPExample = Source (Channel (Int :<+> Eof))
2                      :=> Generator (Channel (Int :<+> Eof))
3                      :=> Sink

```

Source Code 3.3: This example shows the DSL encoding in DPP of repeated elements problems

3.2.2. DSL Validation

The language generated by the grammar needs to be validated to avoid errors or provide an incorrect DPP definition. Fortunately, Haskell provides several Type-level techniques [?] which allows to verify properties of programs before running them, preventing the users to introduce bugs, reducing errors. This verification done by the compiler establish a Curry-Howard Isomorphism [6], i.e. *Propositions as Types - Programs as Proof*. It is important to remark here that Haskell is not a theorem prover System like Coq [?], but some verifications, as we present in this work, can be done with GHC to ensure correctness on programs. Although Haskell provides tools to build advanced type-level verifications, all these techniques require the addition of *Haskell Language Extensions*.

Once we have the encoded DPP problem in the DSL grammar – see subsection 3.2.1 –, we can proceed on validating that encoded grammar. The implementation of the validation of the DSL CFG at type-level, has been done using *Associated Type Families* [?].

```

1  type family And (a :: Bool) (b :: Bool) :: Bool where
2      And 'True 'True = 'True
3      And a b         = 'False
4
5
6  type family IsDP (dpDefinition :: k) :: Bool where
7      IsDP (Source (Channel inToGen) :=> Generator (Channel genToOut) :=> Sink)
8          = And (IsDP (Source (Channel inToGen))) (IsDP (Generator (Channel
9              ↪ genToOut)))
10     IsDP (Source (Channel inToGen) :=> Generator (Channel genToOut) :=>
11         ↪ FeedbackChannel toSource :=> Sink)
12         = And (IsDP (Source (Channel inToGen))) (IsDP (Generator (Channel
13             ↪ genToOut)))
14     IsDP (Source (Channel (a :<+> more)))
15         = IsDP (Source (Channel more))
16     IsDP (Source (Channel Eof))                = 'True
17     IsDP (Generator (Channel (a :<+> more)))    = IsDP (Generator (Channel more))
18     IsDP (Generator (Channel a))                = 'True
19     IsDP x                                      = 'False
20
21 type family ValidDP (a :: Bool) :: Constraint where
22     ValidDP 'True = ()
23     ValidDP 'False = TypeError
24         ( 'Text "Invalid Semantic for Building DP Program"
25           'Text "Language Grammar:"
26           'Text "DP      -> Source CHANS :=> Generator CHANS
27             ↪      :=> Sink"
28           'Text "DP      -> Source CHANS :=> Generator CHANS
29             ↪      :=> FEEDBACK :=> Sink"
30           'Text "CHANS   -> Channel CH"
31           'Text "FEEDBACK -> FeedbackChannel CH"
32           'Text "CH      -> Type :<+> CH | Eof"
33           'Text "Example: 'Source (Channel (Int :<+> Int)) :=>
34             ↪ Generator (Channel (Int :<+> Int)) :=> Sink'"
35         )

```

Source Code 3.4: Type Families `And`, `IsDP` and `ValidDP` which allows to perform a type-level validation over a DSL CFG definition.

In Source Code 3.4, there are 3(three) Type families that helps to validate the DSL CFG. `IsDP` associated type family is checking the production rules P of the grammar defined in Figure 4, returning a promoted data type `[?]` (not a boolean value) `'True` in case the production rule matches all the generated language, or `'False` otherwise. `ValidDP` is taking the result of `IsDP` type application, associating `'True` promoted boolean type to empty `()` constraint. An empty constraint is an indication of nothing to be restricted, meaning that if `ValidDP` is used as a constraint, and it is fully applied to `()`, it will give the compiler the evidence that there is no error at type-level. `ValidDP` is also associating `'False` to a custom `TypeError` which will appear at compilation time if the DPP DSL definition fully applies to that.

```

1  mkDP :: forall dpDefinition filterState filterParam st.
2      ( ValidDP (IsDP dpDefinition)
3        , DPConstraint dpDefinition filterState st filterParam)
4      => Stage (WithSource dpDefinition (DP st))
5      -> GeneratorStage dpDefinition filterState filterParam st
6      -> Stage (WithSink dpDefinition (DP st))
7      -> DP st ()
8  mkDP = ...
9
10 someFunc = mkDP @DPExample ...

```

Source Code 3.5: Definition of `mkDP` function of the Framework which uses type-level validation of the grammar `ValidDP (IsValid Type)`. Last line of the code is showing that using that function will compile-time check the definition of `DPExample` type.

3.2.3. Interpreter of DSL (IDL)

IDL component takes the DPP definition made on with DSL component to interpret and generate the function definitions that the user needs to fill in for solving a specific problem. In subsection 2.2, we have described what the user needs to provide in a DPP algorithm: `Sr`, `G`, `Sk`, and the `F` with the non-empty set of Actors. The IDL generates the function definitions with an empty implementation to be completed by the user, ensuring that those functions will give "Proof" – in terms of Curry-Howard Correspondence [6] – of the "Propositions" defined on the DSL.

Similar techniques that we used on subsubsection 3.2.2 are also used here. On the first hand, we use *Type-level Defunctionalization* [? ?] to let the compiler generates the signatures of the required functions. On the other hand, we use *Term-level Defunctionalization* to interpret those functions. Moreover, *Indexed Types* [?] and *Heterogeneous List* [?] are used to keep track of the dynamic number and polymorphic types of the functions parameters.

```

1  withSource :: forall (dpDefinition :: Type) st. WithSource dpDefinition (DP st)
2      -> Stage (WithSource dpDefinition (DP st))
3  withSource = mkStage' @(WithSource dpDefinition (DP st))
4
5  withGenerator :: forall (dpDefinition :: Type) (filter :: Type) st. WithGenerator
6      ↪ dpDefinition filter (DP st)
7      -> Stage (WithGenerator dpDefinition filter (DP st))
8  withGenerator = mkStage' @(WithGenerator dpDefinition filter (DP st))
9
10 withSink :: forall (dpDefinition :: Type) st. WithSink dpDefinition (DP st)
11     -> Stage (WithSink dpDefinition (DP st))
12 withSink = mkStage' @(WithSink dpDefinition (DP st))

```

Source Code 3.6: This code is showing the different interpreters combinators to help the user to generate the functions of the principal stages of DPP

In Source Code 3.6 we can appreciate the different combinators of the IDL that helps the user of the framework to interpret the DSL to generate the function definitions. `Stage` data type will be cover in Source Code 3.8, but it is a wrapper type of a pipeline stage – minimal unit of execution –, containing the function to be executed – here is the use *Term-level Defunctionalization* –. `withSource`, `withGenerator`, and `withSink` are aliases of the function `mkStage'` which is the combinator that is applying the *Associated Type* related to that stage. For example `withSource`, is equivalent to `mkStage' @(WithSource dpDefinition (DP st))`. For each *Associated Type Family* defintion, there exist an equivalent term-level definition: `WithSource` type with `withSource` term , `WithGenerator` type with `withGenerator` term, and `WithSink` type with `withSink` term – notice the capital case letter "W" indicating the type and not the term –.

```

1  type family WithSource (dpDefinition :: Type) (monadicAction :: Type -> Type) ::
   ↪ Type where
2  WithSource (Source (Channel inToGen) :=> Generator (Channel genToOut) :=> Sink)
   ↪ monadicAction
3      = WithSource (ChanIn inToGen) monadicAction
4  WithSource (Source (Channel inToGen) :=> Generator (Channel genToOut) :=>
   ↪ FeedbackChannel toSource :=> Sink) monadicAction
5      = WithSource (ChanOutIn toSource inToGen) monadicAction
6  WithSource (ChanIn (dpDefinition :<+> more)) monadicAction
7      = WriteChannel dpDefinition -> WithSource (ChanIn more) monadicAction
8  WithSource (ChanIn Eof) monadicAction
9      = monadicAction ()
10 WithSource (ChanOutIn (dpDefinition :<+> more) ins) monadicAction
11     = ReadChannel dpDefinition -> WithSource (ChanOutIn more ins) monadicAction
12 WithSource (ChanOutIn Eof ins) monadicAction
13     = WithSource (ChanIn ins) monadicAction
14 WithSource dpDefinition _
15     = TypeError
16     ( 'Text "Invalid Semantic for Source Stage"
17       '::$: 'Text "in the DP Definition '"
18       ':<>: 'ShowType dpDefinition
19       ':<>: 'Text "' '"
20       '::$: 'Text "Language Grammar:"
21       '::$: 'Text "DP      -> Source CHANS :=> Generator CHANS :=> Sink"
22       '::$: 'Text "DP      -> Source CHANS :=> Generator CHANS :=> FEEDBACK
   ↪ :=> Sink"
23       '::$: 'Text "CHANS   -> Channel CH"
24       '::$: 'Text "FEEDBACK -> FeedbackChannel CH"
25       '::$: 'Text "CH      -> Type :<+> CH | Eof"
26       '::$: 'Text "Example: 'Source (Channel (Int :<+> Int)) :=> Generator
   ↪ (Channel (Int :<+> Int)) :=> Sink'"
27     )

```

Source Code 3.7: An example of the Associated Type Family `WithSource` that allows to implement *Type-level Defunctionalization* technique that will be the Type-level verification of the term `withSource`

In Source Code 3.7, in the highlighted lines, it can be seen how *Type-level Defunctionalization* is being expanded in a signature function definition with the form `WriteChannel a -> ReadChannel b -> ... -> monadicAct` depending on DPP language definition.

```

1 data Stage a where
2   Stage :: Proxy a -> a -> Stage a
3
4 mkStage' :: forall a. a -> Stage a
5 mkStage' = Stage (Proxy @a)

```

Source Code 3.8: **Stage** data type for implementing *Term-level Defunctionalization* providing evidence to the Type-Level Associated types

In Source Code 3.8, **Stage** data type uses a **Proxy** phantom type. This phantom type allow **Stage** to index the type definition generated by **a**. For example, in Source Code 3.6, when **withSource** interpreter is applied to **WithSource dpDefinition**, the compiler is provided with **dpDefinition** DSL type, it expands the function signature belonging to that DPP definition inside the **Stage**.

Generator and Filter. According to DPP definition in subsection 2.2, G has a F template in order to know how to dynamically interpose a new F during the runtime execution of the program. Let's first study F Data Type in the context of the framework.

```

1 newtype Actor dpDefinition filterState filterParam monadicAction =
2   Actor { unActor :: MonadState filterState monadicAction => Stage (WithFilter
3     ↳ dpDefinition filterParam monadicAction) }
4
5 newtype Filter dpDefinition filterState filterParam st =
6   Filter { unFilter :: NonEmpty (Actor dpDefinition filterState filterParam
7     ↳ (StateT filterState (DP st))) }
8   deriving Generic

```

Source Code 3.9: This code shows the definition of the **Filter** data type which contains a non-empty set of **Actor**. The **Actor** data type is an **Stage** in the Context of the **MonadState** to allow keeping a local memory in the execution context of the filter.

In Source Code 3.9 the definition of the **Filter** data type contains a non-empty set of **Actor**. An **Actor** is a **Stage**, because an actor is the minimal unit of execution of a filter. A **Filter** has a **NonEmpty Actor** – Non-empty List – because a filter is built by a sequence of actors calls. Moreover, **Actor** Stage is defunctionalized with **WithFilter Associated Type Family**. **Filter** runs in an explicit **StateT** monadic context. This is because the F instance should have an state, according to DPP definition in subsection 2.2. For example, in the case of DP_{WCC} , as we have seen in ??, F keeps an updated list of connected components that updates as long as it receives more edges that are connected with the current list of vertices. **Actor** data type – see Source Code 3.9 –, is constrained by **MonadState** which is in the same execution context of the whole **NonEmpty Actor** list of the **Filter**. This means the **StateT** is executed for each **Actor** of that filter, sharing the same state between them.

```

1  mkFilter :: forall dpDefinition filterState filterParam st. WithFilter
   ↳ dpDefinition filterParam (StateT filterState (DP st))
2      -> Filter dpDefinition filterState filterParam st
3  mkFilter = Filter . single
4
5  single :: forall dpDefinition filterState filterParam st. WithFilter dpDefinition
   ↳ filterParam (StateT filterState (DP st))
6      -> NonEmpty (Actor dpDefinition filterState filterParam (StateT filterState
   ↳ (DP st)))
7  single = one . actor
8
9  actor :: forall dpDefinition filterState filterParam st. WithFilter dpDefinition
   ↳ filterParam (StateT filterState (DP st))
10     -> Actor dpDefinition filterState filterParam (StateT filterState (DP st))
11  actor = Actor . mkStage' @(WithFilter dpDefinition filterParam (StateT filterState
   ↳ (DP st)))
12
13  (|>>>) :: forall dpDefinition filterState filterParam st. Actor dpDefinition
   ↳ filterState filterParam (StateT filterState (DP st))
14     -> Filter dpDefinition filterState filterParam st
15     -> Filter dpDefinition filterState filterParam st
16  (|>>>) a f = f & _Wrapped' %~ (a <|)
17  infixr 5 |>>>
18
19  (|>>) :: forall dpDefinition filterState filterParam st. Actor dpDefinition
   ↳ filterState filterParam (StateT filterState (DP st))
20     -> Actor dpDefinition filterState filterParam (StateT filterState (DP st))
21     -> Filter dpDefinition filterState filterParam st
22  (|>>) a1 a2 = Filter (a1 <|one a2)
23  infixr 5 |>>

```

Source Code 3.10: Combinators and small constructor to enable building actors and filter.

Finally, in Source Code 3.10, some combinators and smart constructors are provided in the framework to enable the construction of **Filter** and **Actor**. **mkFilter** is a smart constructor for **Filter** Data Constructor. **single** wraps one actor inside a **Filter**. **actor** is a smart constructor for **Actor** Data Constructor. **(|>>>)** is an appending combinator of an **Actor** to a **Filter**. **(|>>>)** also ensures actor execution order, i.e. the latest actor added is the latest to be executed.

```

1   data GeneratorStage dpDefinition filterState filterParam st = GeneratorStage
2   { _gsGenerator      :: Stage (WithGenerator dpDefinition (Filter dpDefinition
   ↪ filterState filterParam st) (DP st))
3   , _gsFilterTemplate :: Filter dpDefinition filterState filterParam st
4   }

```

Source Code 3.11: **Generator** Data type which contains the **Stage** code of the generator itself, and the **Filter** template that it can be spawned by the **Generator**.

In Source Code 3.11, **G** contains a **F** template and its own stage behavior. **Generator** data type has a field with the **Filter** template that could be spawned by the algorithm defined by the user according to the data received from its input channels. **Generator** has also another field with the behavior of the **G** – a **Stage** –.

3.2.4. Runtime System (RS)

The RS can be divided into two parts: the mechanism to generate stages dynamically in runtime, and the execution entry point of the DPP. Regarding execution entry point, all the stages that we have seen in previous sections are the pieces needed to build an executable **DP st** a monad. This executable monad has an existential type similar to **ST** monad to not escape out from the context on different stages. Once the dynamic pipeline starts to execute, the core of the framework dynamically generates stages between **G** and previous stages, according to the user definition, i.e. an *anamorphism* [11] that creates **F** instances until some condition is met.


```

1  unfoldF :: forall dpDefinition readElem st filterState filterParam l.
   ↳ SpawnFilterConstraint dpDefinition readElem st filterState filterParam l
2  => UnFoldFilter dpDefinition readElem st filterState filterParam l
3  -> DP st (HList l)
4  unfoldF = loopSpawn
5
6  where
7    loopSpawn uf@UnFoldFilter{..} =
8      maybe (pure _ufRsChannels) (loopSpawn <=< doOnElem uf) =<< DP (pull
   ↳ _ufReadChannel)
9
10   doOnElem uf@UnFoldFilter{..} elem' = do
11     _ufOnElem elem'
12     if _ufSpawnIf elem'
13     then do
14       (reads', writes' :: HList l3) <- getFilterChannels <$> DP (makeChansF
   ↳ @(ChansFilter dpDefinition))
15       let hlist = elem' .*. _ufReadChannel .*. (_ufRsChannels `hAppendList`
   ↳ writes')
16       void $ runFilter _ufFilter (_ufInitState elem') hlist (_ufReadChannel .*.
   ↳ (_ufRsChannels `hAppendList` writes'))
17       return $ uf { _ufReadChannel = hHead reads', _ufRsChannels = hTail reads' }
18     else return uf

```

Source Code 3.12: `unfoldF` is the *anamorphism* combinator to spawn new `Filter` types between the `Generator` and previous stages.

In Source Code 3.12, it is presented how is the *anamorphism* mechanism that generates dynamic stages between G and the previous stages. That *anamorphism* is implemented with the function `unfoldF`. That function receives an `UnFoldFilter` Data type, which contains the recipe for controlling that unfold recursive call. In line 12, `_ufSpawnIf` field of `UnFoldFilter`, indicates when to stop the recursion. Inside the conditional, in line 14, new channels are created for the new filter to be spawned. Those new channels connect the new filter with the previous stages and with `Generator`. After that, in line 16 `runFilter` starts the monadic computation, spawning the filter stage with its actors. Finally, the new list of channels are returned for the next recursive step to allow further channel connections.

```

1  mkUnfoldFilter :: (readElem -> Bool)
2      -> (readElem -> DP st ())
3      -> Filter dpDefinition filterState filterParam st
4      -> (readElem -> filterState)
5      -> ReadChannel readElem
6      -> HList l
7      -> UnFoldFilter dpDefinition readElem st filterState filterParam l
8
9
10 mkUnfoldFilterForAll' :: (readElem -> DP st ())
11     -> Filter dpDefinition filterState filterParam st
12     -> (readElem -> filterState)
13     -> ReadChannel readElem
14     -> HList l
15     -> UnFoldFilter dpDefinition readElem st filterState
16     ↪ filterParam l
17
18 mkUnfoldFilterForAll :: Filter dpDefinition filterState filterParam st
19     -> (readElem -> filterState)
20     -> ReadChannel readElem
21     -> HList l
22     -> UnFoldFilter dpDefinition readElem st filterState
23     ↪ filterParam l

```

Source Code 3.13: Combinators for building `UnfoldFilter` types indicating the type of the `unfold` that the user want to achieve.

Several smart constructors are also provided for building `UnfoldFilter` Data Type. In Source Code 3.13 the first combinator is the default smart constructor. i) First field (`readElem -> Bool`) indicate if the a new filter should be spawn or not. ii) Second field (`readElem -> DP st ()`) is a monadic optional computation to do when received a new element, for example logging. iii) Third field `Filter` data type to be spawned. iv) Fourth field (`readElem -> filterState`) is initialization of the `Filter` State. v) Fifth field (`ReadChannel readElem`) that feeds the filter instance. vi) Last field is the *Heterogeneous List* with the rest of the channels to connect with other stages. . The combinator `mkUnfoldFilterForAll` is an smart constructor of `UnfoldFilter` that allows to spawn a new filter for each element received in the G.

3.3. Libraries and Tools

3.4. Parallelization

One of the most important components of the implementation is the selection of concurrency libraries to support an intensive parallelization workload. Parallelization techniques and tools have been intensively studied and implemented in Haskell [10]. Indeed, it is well known that green threads and sparks allow spawning thousands to millions of parallel computations. These parallel computations do not penalize performance when compare with Operative System (OS) level threading [9]. A

straightforward assumption to achieve here, is to use `monad-par` library [32, 10]. Nevertheless, in this work, we have discarded the use of sparks [38] because we can achieve the level of required parallelism spawning green threads only. The next obvious choice is to use `forkIO :: IO () -> IO ThreadId` from `base` library [28]. However, that would imply handling all the threads lifecycles and errors programmatically without any abstraction to facilitate that complex task. Therefore, we choose `async` library [17] which enables to spawn asynchronous computations [9] on Haskell using green threads, and at the same time, it provides useful combinators to managing thread terminations and errors.

3.4.1. Channels

We have several techniques to our disposal to communicate between threads or sparks in Haskell like `MVar` or concurrent safe mechanisms like Software Transactional Memory (STM) [5]. At the same time, in Haskell library ecosystem, we dispose of `Channels` abstractions based on both mentioned communication techniques. In that sense, for conducting the communication between dynamic stages and data flowing in a DPP, we have selected `unagi-chan` library [43] which provides the following advantages to our solution: Firstly, `MVar` channel without using STM reducing overhead. STM is not required in a DPP because one specific stage which is running in a separated thread, can only access to its I/O channels for reading/writing accordingly, and those operations are not concurrently shared by other threads (stages) for the same channels. Second, non-blocking channels. `unagi-chan` library contains blocking and non-blocking channels for reading. This aspect is key to gain speed up on the implementation. Third, the library is optimized for *x86* architectures with use of low-level `fetch-and-add` instructions. Finally, `unagi-chan` is 100x faster [?] on Benchmarking compare with STM and default base `Chan` implementations.

4. Enumerating Weakly Connected Component on the DPF

Having the DPF-Haskell in place as we describe in the previous section 3, we implement the algorithm for enumerating weakly components of a graph using the DPF-Haskell. We describe first the general algorithm in the context of DPP. Secondly, we show the implementation details of that algorithm in the context of the DPF-Haskell and finally we show the experiments and results conducted over that implementation.

4.1. DP_{WCC} Algorithm

Let us consider the problem of computing/enumerating the (weak) connected components of a graph G using DPP. A connected component of a graph is a subgraph in which any two vertices are connected by paths. Thus, finding connected components of an undirected graph implies obtaining the minimal partition of the set of nodes induced by the relationship *connected*, i.e., there is a path between each pair of nodes. An example of that graph can be seen in Figure 5. The input of the Dynamic Pipeline for computing the WCC of a graph, DP_{WCC}, is a sequence of edges ending with `eof`¹. The connected components are output as soon as they are computed, i.e., they are produced incrementally. Roughly speaking the idea of the algorithm is that the weakly connected components are built in two phases. In the first phase filter instance stages receive the edges of the input graph and create sets of connected vertices. During the second phase, these filter instances construct maximal subsets of connected vertices, i.e. the vertices corresponding to (weakly) connected components. DP_{WCC} is defined in terms of the behavior of its four kinds stages: *Source* (Sr_{WCC}), *Generator*

¹Note that there are neither isolated vertices nor loops in the source graph G .

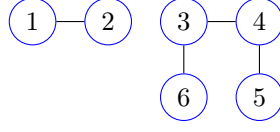


Figure 5: Example of a graph with two weakly connected components: $\{1, 2\}$ and $\{3, 4, 5, 6\}$

(G_{WCC}), *Sink* (Sk_{WCC}), and *Filter* (F_{WCC}) stages. Additionally, the channels connecting these stages must be defined. In DP_{WCC} , stages are connected linearly and unidirectionally through the channels IC_E and $IC_{set(V)}$. Channel IC_E carries edges while channel $IC_{set(V)}$ conveys sets of connected vertices. Both channels end by the *eof* mark. The behavior of F_{WCC} is given by a sequence of two actors (scripts). Each actor corresponds to a phase of the algorithm. In what follows, we denote these actors by $actor_1$ and $actor_2$, respectively. The script $actor_1$ keeps a set of connected vertices (CV) in the state of the F_{WCC} instance. When an edge e arrives, if an endpoint of e is present in the state, then the other endpoint of e is added to CV . Edges without incident endpoints are passed to the next stage. When *eof* arrives at channel IC_E , it is passed to the next stage, and the script $actor_2$ starts its execution. If script $actor_2$ receives a set of connected vertices CV in $IC_{set(V)}$, it determines if the intersection between CV and the nodes in its state is not empty. If so, it adds the nodes in CV to its state. Otherwise, the CV is passed to the next stage. Whenever *eof* is received, $actor_2$ passes-through $IC_{set(V)}$ the set of vertices in its state and the *eof* mark to the next stage; then, it dies. The behavior of Sr_{WCC} corresponds to the identity transformation over the data stream of edges. As edges arrive, they are passed through IC_E to the next stage. When receiving *eof* on IC_E , this mark is put on both channels. Then, Sr_{WCC} dies.

Let us describe this behavior with the example of the graph shown in Figure 5.

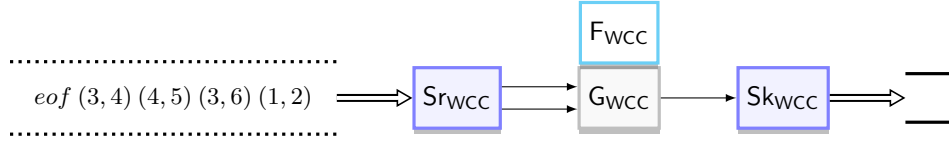
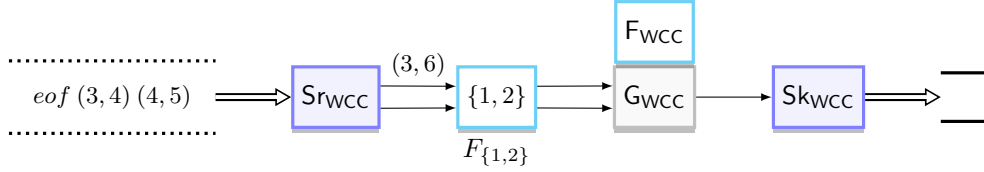
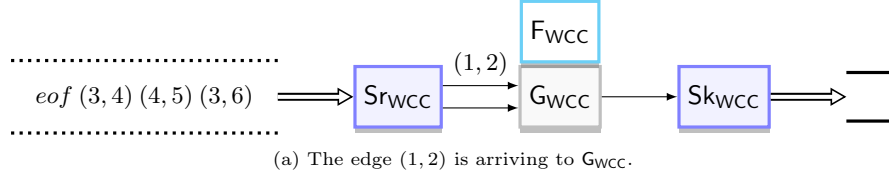


Figure 6: DP_{WCC} Initial setup. Stages Source, Generator, and Sink are represented by the squares labeled by Sr_{WCC} , G_{WCC} and Sk_{WCC} , respectively. The square F_{WCC} corresponding to the Filter stage template is the parameter of G_{WCC} . Arrows \Rightarrow between represents the connection of stages through two channels, IC_E , and $IC_{set(V)}$. The arrow \rightarrow represents the channel $IC_{set(V)}$ connecting the stages G_{WCC} and Sk_{WCC} . The arrow \Rightarrow stands for I/O data flow. Finally, the input stream comes between the dotted lines on the left and the WCC computed incrementally will be placed between the solid lines on the right.

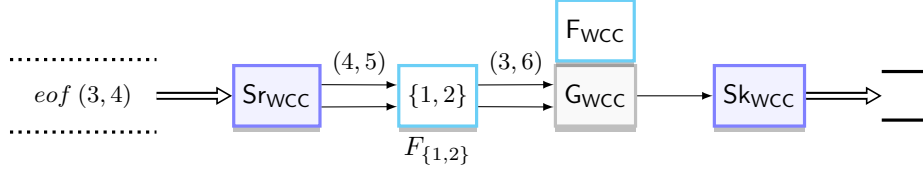
Figure 6 depicts the initial configuration of DP_{WCC} . The interaction of DP_{WCC} with the "external" world is done through the stages Sr_{WCC} and Sk_{WCC} . Indeed, once activated the initial DP_{WCC} , the input stream – consisting of a sequence containing all the edges in the graph in Figure 5 – feeds Sr_{WCC} while Sk_{WCC} emits incrementally the resulting weakly connected components. In what follows Figure 7, Figure 8, Figure 9, Figure 10 and Figure 11 depict the evolution of the DP_{WCC} .

It is important to highlight that during the states shown in Figure 7a, Figure 7b, Figure 8a, Figure 8b

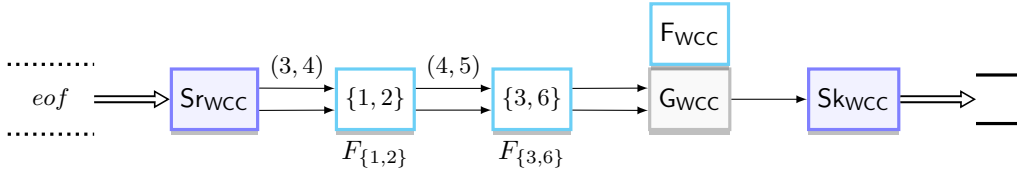


(b) When the edge (1, 2) arrives to G_{WCC} , it spawns a new instance of F_{WCC} before G_{WCC} . Filter instance $F_{\{1,2\}}$ is connected to G_{WCC} through channels IC_E and $IC_{set(V)}$. The state of the new filter instance $F_{\{1,2\}}$ is initialized with the set of vertices $\{1, 2\}$. The edge (3, 6) arrives to the new filter instance $F_{\{1,2\}}$.

Figure 7: Evolution of the DP_{WCC} : First state



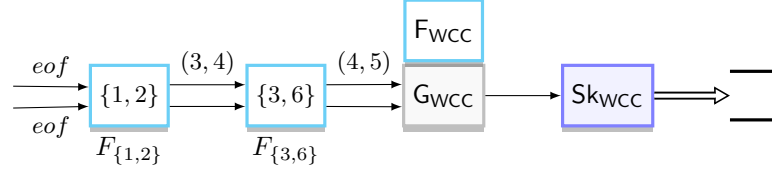
(a) None of the vertices in the edge (3, 6) is in the set of vertices $\{1, 2\}$ in the state of $F_{\{1,2\}}$, hence it is passed through IC_E to G_{WCC} .



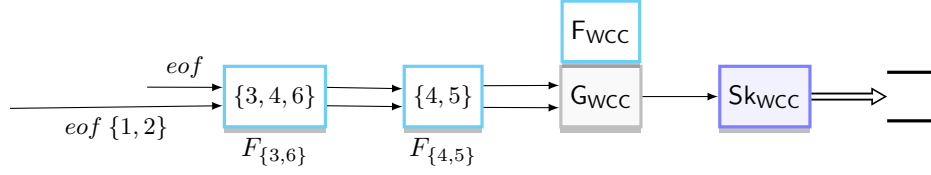
(b) When the edge (3, 6) arrives to G_{WCC} , it spawns the filter instance $F_{\{3,6\}}$ between $F_{\{1,2\}}$ and G_{WCC} . Filter instance $F_{\{1,2\}}$ is connected to the new filter instance $F_{\{3,6\}}$ and this one is connected to G_{WCC} through channels IC_E and $IC_{set(V)}$. The state of the new filter instance $F_{\{3,6\}}$ is initialized with the set of vertices $\{3, 6\}$. The edge (3, 4) arrives to $F_{\{1,2\}}$ and Sr_{WCC} is fed with the mark *eof*. Edges (3, 4) and (4, 5) remain passing through IC_E .

Figure 8: Evolution of the DP_{WCC} : Second state

and Figure 9a the only actor executed in any filter instance is $actor_1$ (constructing sets of connected vertices). Afterwards, although $actor_1$ can continue being executed in some filter instances, there are some instances that start executing $actor_2$ (constructing sets of maximal connected vertices). This is shown from Figure 9a to Figure 11a.

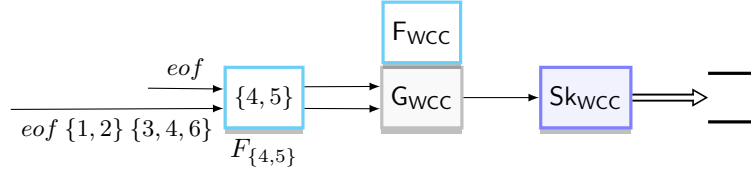


(a) Sr_{WCC} fed both, IC_E and $IC_{set(V)}$, channels with the mark *eof* received from the input stream in previous state and then, it died. The edge (4,5) is arriving to G_{WCC} and the edge (3,4) is arriving to $F_{3,6}$.

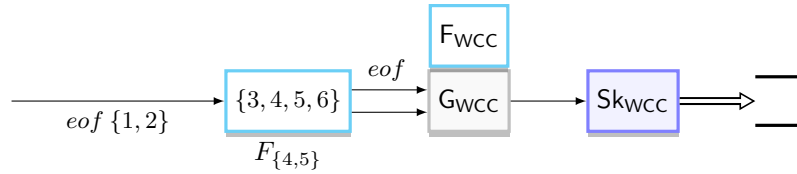


(b) When the edge (4,5) arrives to G_{WCC} , it spawns the filter instance $F_{4,5}$ between $F_{3,6}$ and G_{WCC} . Filter instance $F_{3,6}$ is connected to the new filter instance $F_{4,5}$ and this one is connected to G_{WCC} through channels IC_E and $IC_{set(V)}$. Since the edge (3,4) arrived to $F_{3,6}$ at the same time and vertex 3 belongs to the set of connected vertices of the filter $F_{3,6}$, the vertex 4 is added to the state of $F_{3,6}$. Now, the state of $F_{3,6}$ is the connected set of vertices $\{3,4,6\}$. When the mark *eof* arrives to the first filter instance, $F_{1,2}$, through $IC_{set(V)}$, this stage passes its partial set of connected vertices, $\{1,2\}$, through $IC_{set(V)}$ and dies. This action will activate *actor*₂ in next filter instances to start building maximal connected components. In this example, the state in $F_{3,6}$, $\{3,4,6\}$, and the arriving set $\{1,2\}$ do not intersect and, hence, both sets of vertices, $\{1,2\}$ and $\{3,4,6\}$ will be passed to the next filter instance through $IC_{set(V)}$.

Figure 9: Evolution of the DP_{WCC}: Third state

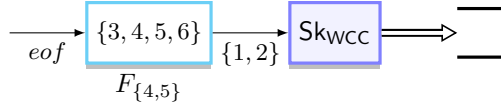


(a) The set of connected vertices $\{3,4,6\}$ is arriving to $F_{4,5}$. The mark *eof* continues passing to next stages through the channel IC_E .

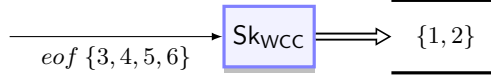


(b) Since the intersection of the set of connected vertices $\{3,4,6\}$ arrived to $F_{4,5}$ and its state is not empty, this state is enlarged to be $\{3,4,5,6\}$. The set of connected vertices $\{1,2\}$ is arriving to $F_{4,5}$.

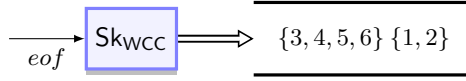
Figure 10: Evolution of the DP_{WCC}: Fourth state



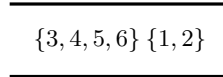
(a) $F_{\{4,5\}}$ has passed the set of connected vertices $\{1, 2\}$ and it is arriving to Sk_{WCC} . The mark *eof* is arriving to $F_{\{4,5\}}$ through $IC_{set(V)}$.



(b) Since the mark *eof* arrived to $F_{\{4,5\}}$ through $IC_{set(V)}$, it passes its state, the set $\{3, 4, 5, 6\}$ through $IC_{set(V)}$ to next stages and died. The set of connected vertices $\{1, 2\}$ arrived to Sk_{WCC} and this implies that $\{1, 2\}$ is a maximal set of connected vertices, i.e. a connected component of the input graph. Hence, Sk_{WCC} output this first weakly connected component.



(c) Finally, the set of connected vertices $\{3, 4, 5, 6\}$ arrived to Sk_{WCC} and was output as a new weakly connected component. Besides, the mark *eof* also arrived to Sk_{WCC} through $IC_{set(V)}$ and thus, it dies.



(d) The weakly connected component of in the graph Figure 5 such as they have been emitted by DP_{WCC} .

Figure 11: Last states in the evolution of the DP_{WCC}

4.2. DP_{WCC} Implementation

As we said before, the DP_{WCC} implementation has been made as a proof of concept to understand and explore the limitations and challenges that we could find in the development of a future DPF in Haskell. In Section 3 we emphasize that the focus of DPF in Haskell is on the IDL component. Hence, the development of the DP_{WCC} is as general as possible using most of the constructs and abstractions required by the IDL. Lets introduce the minimal code needed for encoding any DPP using DPF-Haskell. ²

```
1  type DPConnComp = Source (Channel (Edge :<+> ConnectedComponents :<+> Eof))
2                      :=> Generator (Channel (Edge :<+> ConnectedComponents :<+> Eof))
3                      :=> Sink
4
5  program :: FilePath -> IO ()
6  program file = runDP $ mkDP @DPConnComp (source' file) generator' sink'
```

Source Code 4.1: In this code we can appreciate the main construct of our DP_{WCC} which is a combination of Sr_{WCC} , G_{WCC} and Sk_{WCC}

In Source Code 4.1 there are two important declarations. First, the *Type Level* declaration of the DP_{WCC} to indicate DPF-Haskell how our stages are going be connected, and using that *Type Level* construct, we use the IDL to allow the framework interpret the type representation of our DPP and ensuring at compilation time that we provide the correct stages, *Source* (Sr_{WCC}), *Generator* (G_{WCC}) and *Sink* (Sk_{WCC}), that matches those declaration. According to this declaration what we need to provide is the correct implementation of `source'`, `generator'` and `sink'` which *Type checked* the DPP type definition³.

²All the code that we expose here can be accessed publicly in <https://github.com/jproyo/dynamic-pipeline/tree/main/examples/Graph>

³The names of the functions are completely choosen by the user of the framework and it should not be confused with the internal framework combinators.

```

1  source' :: FilePath
2         -> Stage
3         (WriteChannel Edge -> WriteChannel ConnectedComponents -> DP st ())
4  source' filePath = withSource @DPConnComp
5                      $ \edgeOut _ -> unfoldFile filePath edgeOut (toEdge . decodeUtf8)
6
7  sink' :: Stage (ReadChannel Edge -> ReadChannel ConnectedComponents -> DP st
8             ↪ ())
9
10 sink' = withSink @DPConnComp $ \_ cc -> withDP $ foldM_ cc print
11
12 generator' :: GeneratorStage DPConnComp ConnectedComponents Edge st
13 generator' =
14     let gen = withGenerator @DPConnComp genAction
15     in mkGenerator gen filterTemplate

```

Source Code 4.2: In this code we can appreciate the Sr_{WCC} , G_{WCC} and Sk_{WCC} functions that matches the type level definition of the DP. Sr_{WCC} and Sk_{WCC} are completely trivial but G_{WCC} will be analyzed later due to its internal complexity.

As we appreciate in Source Code 4.2, Sr_{WCC} and Sk_{WCC} are trivial. In the case of `source'` the only work it needs to do is to read the input data edge by edge and downstream to the next stages. That's achieved with a DPF-Haskell combinator called `unfoldFile` which is a catamorphism of the input data to the stream. In the case of Sk_{WCC} it is also a simple implementation but doing the opposite as Sr_{WCC} using an anamorphism combinator provided by the framework as well, which is `foldM_`. The G_{WCC} Stage is a little more complex because it contains the core of the algorithm explained in subsection 4.1. According to what we described in subsection 2.2, *Generator* stage spawns a *Filter* on each received edge in our case of DP_{WCC} . Therefore, it needs to contain that recipe on how to generate a new *Filter* instance – in our case of Haskell it is a defunctionalized Data Type or Function –. Then, we have two things `genAction` which tells how to spawn a new *Filter* and under what circumstances, and `filterTemplate` with the function to be spawn.

```

1  genAction :: Filter DPConnComp ConnectedComponents Edge st
2              -> ReadChannel Edge
3              -> ReadChannel ConnectedComponents
4              -> WriteChannel Edge
5              -> WriteChannel ConnectedComponents
6              -> DP st ()
7  genAction filter' readEdge readCC _ writeCC = do
8      let unfoldFilter = mkUnfoldFilterForAll filter' toConnectedComp readEdge
9          ↪ (readCC .*. HNil)
10     results <- unfoldF unfoldFilter
11     foldM_ (hHead results) (`push` writeCC)

```

Source Code 4.3: In this code we can appreciate the Generator Action code which will expand all the filters in runtime in front of it and downstream all the connected components calculated for those, to the Sink

DPF-Haskell provides several combinators to help the user with the *Generator* code, in particular with the spawning process as it has been describe in section 3. `genAction` for DP_{WCC} will use the combinator `mkUnfoldFilterForAll` which will spawn one *Filter* per received edge in the channel, expanding dynamically the stages on runtime. In the third highlighted line we can appreciate how after expanding the filters, the generator will downstream to the *Sink*, the received Connected Components calculated from previous filters.

```

1  filterTemplate :: Filter DPConnComp ConnectedComponents Edge st
2  filterTemplate = actor actor1 |>> actor actor2
3
4  actor1 :: Edge
5      -> ReadChannel Edge
6      -> ReadChannel ConnectedComponents
7      -> WriteChannel Edge
8      -> WriteChannel ConnectedComponents
9      -> StateT ConnectedComponents (DP st) ()
10 actor1 _ readEdge _ writeEdge _ =
11     foldM_ readEdge $ \e -> get >>= doActor e
12 where
13     doActor v conn
14         | toConnectedComp v `intersect` conn = modify' (toConnectedComp v <>)
15         | otherwise = push v writeEdge
16
17 actor2 :: Edge
18     -> ReadChannel Edge
19     -> ReadChannel ConnectedComponents
20     -> WriteChannel Edge
21     -> WriteChannel ConnectedComponents
22     -> StateT ConnectedComponents (DP st) ()
23 actor2 _ _ readCC _ writeCC = do
24     foldWithM_ readCC pushMemory $ \e -> get >>= doActor e
25
26 where
27     pushMemory = get >>= flip push writeCC
28
29     doActor cc conn
30         | cc `intersect` conn = modify' (cc <>)
31         | otherwise = push cc writeCC

```

Source Code 4.4: Filter template code composed by 2 Sequential Actors that will calculate the Connected Components and downstream them.

Finally, in Source Code 4.4 the *Filter* template code is defined. As we have seen in subsection 4.1, DP_{WCC} *Filter* is composed of 2 Actors. The first actor collect all the possible vertices that are incidence some of the vertices edge that was instantiate with. Once it does not receive any more edges, it starts downstream it set of vertices to the following filters in order to build a maximal connected components, that is **actor2**. At the end **actor2** will downstream its connected component to the following stages. As we show, with the help of the Haskell Dynamic Pipeline Framework, building a DPP algorithm like WCC enumeration consist in few lines of codes with the *Type Safety* that Haskell provides.

4.3. Empirical Evaluation

The empirical study aims at evaluating the performance of DP_{WCC} when implemented in Haskell. Our goal is to answer the following research questions:

- RQ1)** Does DP_{WCC} in Haskell support the dynamic parallelization level that DP_{WCC} requires?
RQ2) Is DP_{WCC} in Haskell competitive compared with default implementations on base libraries for the same problem? **RQ3)** Does DP_{WCC} in Haskell handle memory efficiently?

We have conducted different kinds of experiments to test our assumptions and verify the correctness of the implementation. First, we have performed an *Implementation Analysis* in which we have selected some graphs from Stanford Network Data Set Collection (SNAP) [40] and analyze how the implementation behaves under real-world graphs if it timeouts or not and if it is producing correct results in terms of the amount of WCC that we know beforehand. We have also tested the implementation doing a *Benchmark Analysis* where we focus on two different types of benchmarks. On the one hand, using `criterion` library [22], we have evaluated a benchmark between our solution and WCC algorithm implemented in `containers` Haskell library [21] using `Data.Graph`. On the other hand, we have compared if the results are being generated incrementally in both cases and how that is done during the pipeline execution time. This last analysis has been conducted using `diefpy` tool [1, 23]. Finally, we have executed a *Performance Analysis* in which we have to gather profiling data from Glasgow Haskell Compiler (GHC) for one of the real-world graphs, to measure how the program performs regarding multithreading and memory allocation.

4.4. Running Architecture

All the experiments have been executed in a *x86 64 bits* architecture with a *6-Core Intel Core i7* processor of 2,2 GHz which can emulate up to 12 virtual cores. This processor has *hyper-threading* enable. Regarding memory, the machine has *32GB DDR4* of RAM, *256 KB* of L2 cache memory, and *9 MB* of L3 cache.

4.5. Haskell Setup

Regarding specific libraries and compilations flags used on Haskell, we have used GHC version 8.10.4. We have also used the following set of libraries: `bytestring 0.10.12.0` [19], `containers 0.6.2.1` [21], `relude 1.0.0.1` [37] and `unagi-chan 0.4.1.3` [43]. The use of `relude` library is because we disabled `Prelude` from the project with the language extension `NoImplicitPrelude` [27]. Regarding compilation flags (GHC options) we have compiled our program with `-threaded`, `-O3`, `-rtsopts`, `-with-rtsopts=-N`. Since we have used `stack` version 2.5.1 [39] as a building tool on top of GHC the compilation command is `stack build`⁴.

4.6. DataSets

For all the experiments, we have used the following networks taken from SNAP [40]. In this particular experiment setup, we have selected the following specific data sets that can be found here [34, 33, 35]

⁴For more information about `package.yaml` or `cabal` file please check <https://github.com/jproyo/upc-mir-tfm/tree/main/connected-comp>

Network	Nodes	Edges	Diameter	#WCC	#Nodes Largest WCC
Enron Emails	36692	183831	11	1065	33696 (0.918)
Astro Physics Collaboration Net	18772	198110	14	290	17903 (0.954)
Google Web Graph	875713	5105039	21	2746	855802 (0.977)

Table 1: DataSet of Graphs Selected

The criteria for selecting the networks have been followed the idea of testing the solution in more complex graphs, in which all of them are undirected but with different sizes concerning its number of nodes as we can see in Table 1.

4.7. Experiments Definition

E1: Implementation Analysis. In this experiment, we measure GHC statistics running time enabling `+RTS -s` flags. The metrics that we measure are *MUT Time* which is the amount of time in seconds GHC is running computations and *GC Time* which is the number of seconds that GHC is running garbage collector. *Total execution time* is the sum of both in seconds. At the same time, we are going to check the correctness of the output counting the number of WCC generated by the algorithm against the already known topology of it in subsection 4.6. The experiment’s primary goal is to help answer the research question [RQ2].

E2: Benchmark Analysis. In this experiment, we conduct two benchmark analysis over execution time comparing DP-WCC-Haskell with Haskell `containers` default implementation. In the first benchmark analysis, we use `criterion` [22] tool in Haskell which runs over four iterations of each of the algorithms to get a mean execution time in seconds and compare the results in a plot. In the second benchmark, we use Diefficiency Metrics (Dm) Tool *diefpy* [23] in order to measure with the ability of DPP model to generate results incrementally [1]. This is one of the strongest feature of DPP Paradigm since it allows process and generate results without no need of waiting for processing until the last element of the data source. This kind of aspect is essential not only for big data inputs where perhaps the requirements allow for processing until some point of the time having partial results but at the same time is important to process unbounded streams. The experiment’s primary goal is to help answer the research question [RQ2] as well.

E3: Performance Analysis. In this experiment, we measure internal parallelism in GHC and memory usage during the execution of one of the example networks. The motivation of this is to verify empirically how DP-WCC-Haskell is handling parallelization and memory usage. This experiment is conducted using two tools, *ThreadScope* [42] for conducting multithreading analysis and *eventlog2html* [26] to conduct memory usage analysis. Regarding multithreading analysis the metrics that we measure are the distribution of threads among processors over execution time which is how many processors are executing running threads over the whole execution; and the mean number of running threads per time slot which is calculated by zooming in 8 time slots and taking the mean number of threads per processor to see if it is equally distributed among them. In regards to memory management, the metric that we measure is the amount of memory in *MB* consumed per data type during the whole execution time. The experiment helps to answer the research questions [RQ1,RQ3].

5. Discussion of Observed Results

5.1. Experiment: E1

The following represents the execution for running these graphs on our DPP implementation.

Network	Exec Param	MUT Time	GC Time	Total Time
Enron Emails	+RTS -N4 -s	2.797s	0.942s	3.746s
Astro Physics Coll Net	+RTS -N4 -s	2.607s	1.392s	4.014s
Google Web Graph	+RTS -N8 -s	137.127s	218.913s	356.058s

Table 2: Execution times

It is important to point out that since the first two networks are smaller in the number of edges compared with *web-Google*, executing those with 8 cores as the `-N` parameters indicates, does not affect the final speed-up since GHC is not distributing threads on extra cores because it handles the load with 4 cores only.

As we can see in Table 2, we are obtaining remarkable execution times for the first two graphs and it seems not to be the case for *web-Google*. Doing a deeper analysis on the topology of this last graph, we can see according to Table 1 that the number of *Nodes in the largest WCC* is the highest one. This means that there is a WCC which contains 97.7% of the nodes. Moreover, we can confirm that if we analyze even deeper how is the structure of that WCC with the output of the algorithm, we can notice that the largest WCC is the last one on being processed. Having that into consideration we can state that due to the nature of our algorithm which needs to wait for collecting all the vertices in the `actor2` filter stage it penalizes our execution time for that particular case. A more elaborated technique for implementing the actors is required to speed up execution.

Regarding the correctness of the output, we have verified with the outputs that the number of connected components is the same as the metrics already gathered in Table 1.

5.2. Experiment: E2

Criterion Benchmark. In Figure 12, orange bars report the time taken by `Data.Graph` in Haskell `containers` library [21]. Blue light bars represent the time taken by DP-WCC-Haskell.

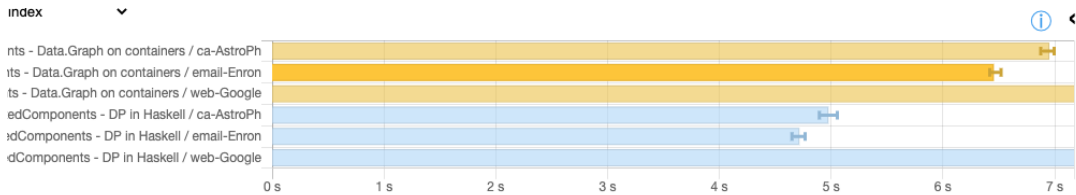


Figure 12: Benchmark 1 - DP in Haskell vs. `Data.Graph` Haskell

Figure 12 shows that DP-WCC-Haskell solution is 1.3 faster compare with Haskell `containers` library. Despite this, if we zoom in Figure 12, it can be observed that DP-WCC-Haskell solution is slower compared with Haskell `containers`; the reasons behind this have been explained in subsection 5.1.

Regarding mean execution times for each implementation on each case measure by `criterion` library [22], we can display the following results:

Network	DP-WCC-Haskell	Haskell containers	Speed-up
Enron Emails	4.68s	6.46s	1.38
Astro Physics Coll Net	4.98s	6.95s	1.39
Google Web Graph	386s	106s	-3.64

Table 3: Mean Execution times

These results allow for answering Question [Q2]. We already had a partial answer with the previous experiment E1 about [Q2] (subsection 4.3) where we have seen that the graph topology is affecting the performance and the parallelization, penalizing DP-WCC-Haskell for this particular case. In this benchmark, the solution against a non-parallel `containers Data.Graph` confirms the hypothesis.

Diefficiency Metrics. Some considerations are needed before starting to analyze the data gathered with Dm tool. Firstly, the tool is plotting the results according to the traces generated by the implementation, both DP-WCC-Haskell and Haskell *containers*. By the nature of DPP model, we can gather or register that timestamps as long as the model is generating results. In the case of Haskell *containers*, this is not possible since it calculates WCC at once. This is not an issue and we still can check at what point in time all WCC in Haskell *containers* are generated. In those cases, we are going to see a straight vertical line.

It is important to remark that we needed to scale the timestamps because we have taken the time in nanoseconds. After all, the incremental generation between one WCC and the other is very small but significant enough to be taken into consideration. Thus, if we left the time scale in integer milliseconds, microseconds, or nanoseconds integer part it cannot be appreciated. In case of escalation, we are discounting the nanosecond integer of the first generated results resulting in a time scale that starts close to 0. This does not mean that the first result is generated at 0 time, but we are discarding the previous time to focus on how the results are incrementally generated.

Having said that, we can see the results of Dm which are presented in two types of plots. The first one is regular line graphs in where the x axis shows the time escalated when the result was generated and the y axis is showing the component number that was generated at that time. The second type of plot is a radar plot in which shows how the solution is behaving on the dimensions of Time for the first tuple (TFFT), Execution Time (ET), Throughput (T), Completeness (Comp) and Diefficiency Metric `dief@t` (`dief@t`) and how are the tension between them; all these metrics are higher is better. All the details about these metrics are explained here [1].

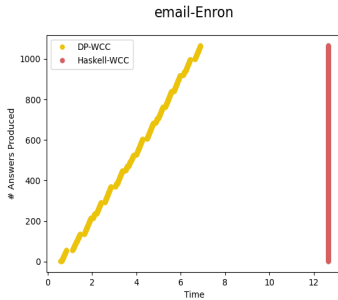


Figure 13: email-Enron Dm

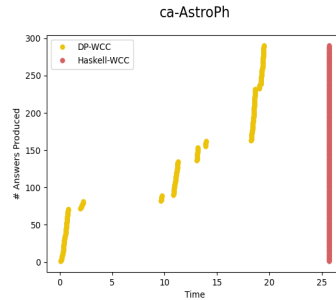


Figure 14: ca-AstroPh Dm

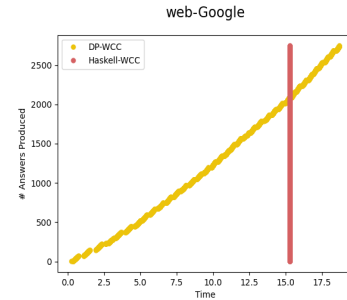


Figure 15: web-Google Dm

Based on the results shown in all the figures above, all the solutions in DP-WCC-Haskell are being

generated incrementally, but there is some difference that we would like to remark. In the case of *email-Enron* and *ca-AstroPh* graphs as we can see in Figure 13 and Figure 14, there seems to be a more incremental generation of results. This behavior is measured with the values of Diefficiency Metric dief@t (dief@t). *ca-AstroPh* as it can be seen in Figure 14, is even more incremental showing a clear separation between some results and others. The *web-Google* network which is shown in Figure 15, is a little more linear and that is because all the results are being generated with very little difference in time between them. This is due to the fact of the explained reasons in subsection 5.1. Having the biggest WCC at the end of *web-Google* DPP algorithm it is retaining results until the biggest WCC can be solved, which takes longer.

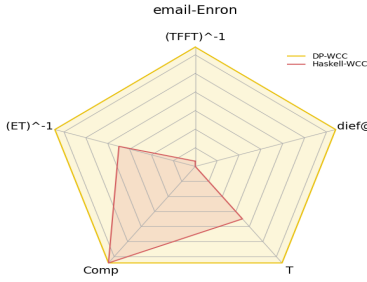


Figure 16: email-Enron Dm

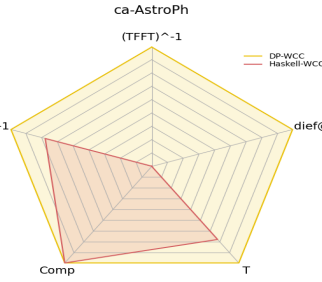


Figure 17: ca-AstroPh Dm

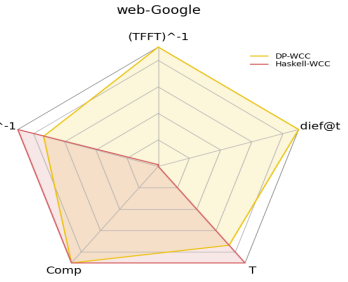


Figure 18: web-Google Dm

As we can appreciate in the above radar plots our previous analysis can be confirmed. We can see for example that the Throughput of *web-Google* in Figure 18, in the case of DP-WCC-Haskell is worse than Haskell **containers**, which is not happening for the others.

In conclusion, we can say that regarding [Q2] (subsection 4.3) although DP-WCC-Haskell is faster than the traditional approach, the speed-up dimension execution factor is not always the most interest analysis that we can have, because as we have seen even when in the case of *web-Google* Graph DP-WCC-Haskell is slower at execution, it is at least generating incremental results without the need to wait for the rest of the computations.

5.3. Experiment: E3

For this type of analysis, our experiment focuses on *email-Enron* network [34] only because profiling data generated by GHC is big enough to conduct the analysis and on the other, and enabling profiling penalize execution time.

Multithreading. For analyzing parallelization and multithreading we have used *ThreadScope* [42] which allows us to see how the parallelization is taking place on GHC at a fine grained level and how the threads are distributed throughout the different cores requested with the `-N` execution `ghc-option` flag.

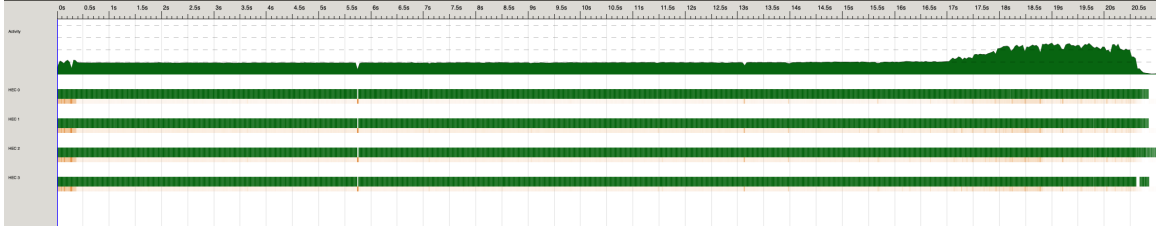


Figure 19: Threadscope Image of General Execution

In Figure 19, we can see that the parallelization is being distributed evenly among the 4 Cores that we have set for this execution. The distribution of the load is more intensive at the end of the execution, where `actor2` filter stage of the algorithm is taking place and different filters are reaching execution of that second actor.

Another important aspect shown in Figure 19, is that this work is not so significant for GHC and the threads and distribution of the work keeps between 1 or 2 cores during the execution time of the `actor1`. However, the usages increase on the second actor as pointed out before. In this regard, we can answer research questions [Q1] and [Q3] (subsection 4.3), verifying that Haskell not only supports the required parallelization level but is evenly distributed across the program execution too.

Finally, it can also be appreciated that there is no sequential execution on any part of the program because the 4 cores have *CPU* activity during the whole execution time. This is because as long the program start, and because of the nature of the DPP model, it is spawning the *Source* stage in a separated thread. This is a clear advantage for the model and the processing of the data since the program does not need to wait to do some sequential processing like reading a file, before start computing the rest of the stages.

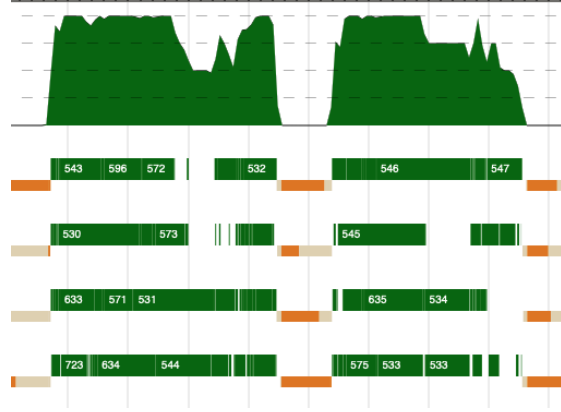


Figure 20: Threadscope Image of Zoomed Fraction

Figure 20 zooms in on *ThreadScope* output in a particular moment, approximately in the middle of the execution. We can appreciate how many threads are being spawned and by the tool and if they are evenly distributed among cores. The numbers inside green bars represent the number of threads that are being executed on that particular core (horizontal line) at that execution slot. Thus, the number of threads varies among slot execution times because as it is already known, GHC implements *Preemptive Scheduling* [8].

Having said that, it can be appreciated in Figure 20 our first assumption that the load is evenly distributed because the mean number of executing threads per core is 571.

Memory allocation. Another important aspect in our case is how the memory is being managed to avoid memory leaks or other non-desired behavior that increases memory allocation during the execution time. This is even more important in the particular implementation of WCC using DPP

Another kind of streaming implementation in Haskell is described in [9]. In that work, the author describes how to encode pipeline parallelism with `Par Monad`. Although this could have been a suitable alternative for implementing DP, the parallelization level used by `Par Monad` is sparks [38]. As we have explained in section ??, we do not require to reach that level of parallelization in our current model.

In regards to other DP language implementations, a significant contribution on [12] has been done, where a DP implementation in Go Programming Language (Go) for counting triangles of graphs is compared against MapReduce. Those experiment results have shown how DP in Go improves the performance in terms of execution time and memory depending on the graph topology. It would be interesting and a matter of future work, to compare different language implementations of DPs, taking into consideration those promising results and the ones presented in this article.

7. Conclusions and Ongoing Work

The empirical evaluation of the DP-WCC-Haskell implementation to compute weakly connected components of a graph, evidence suitability, and robustness to provide a Dynamic Pipeline Framework in that language. Measuring using `dief@t` metrics in section 5.2 reveals some advantageous capability of DP_{WCC} implementation to deliver incremental results compared with default containers library implementation. Regarding the main aspects where DPP is strong, i.e. pipeline parallelism and time processing, the DP_{WCC} performance shows that Haskell can deal with the requirements for the WCC problem without penalizing neither execution time nor memory allocation. In particular, the DP_{WCC} implementation outperforms in those cases where the topology of the graph is more sparse and where the number of vertices in the largest WCC is not big enough. We think this work has gathered enough evidence to show that the implementation of Dynamic Pipeline in Haskell Programming Language is feasible. This fact opens a wide range of algorithms to be explored using the Dynamic Pipeline Paradigm, supported by purely functional programming language. As we mentioned in section ??, we are addressing the design and the definition of the DSL in Haskell, taking into account the knowledge obtained in this work. The complete DPF's implementation will contain the *Type-Level* DSL allowing the user to define algorithms in terms of DP and the Interpreter of DSL (IDL) that will be mainly based on what has been presented here.

Bibliography

- [1] M. Acosta, M.-E. Vidal, and Y. Sure-Vetter. Diefficiency metrics: measuring the continuous efficiency of query processing approaches. In *International Semantic Web Conference*, pages 3–19. Springer, 2017.
- [2] F. N. Afrati, D. Fotakis, and J. D. Ullman. Enumerating subgraph instances using map-reduce. In *29th IEEE International Conference on Data Engineering, ICDE 2013, Brisbane, Australia, April 8-12, 2013*, pages 62–73, 2013. doi: 10.1109/ICDE.2013.6544814. URL <http://dx.doi.org/10.1109/ICDE.2013.6544814>.
- [3] M. Fowler. *Domain-Specific Languages*. Addison-Wesley, Upper Saddle River, NJ, 2010. ISBN 978-0-321-71294-3. URL <https://www.safaribooksonline.com/library/view/domain-specific-languages/9780132107549/>.
- [4] M. I. Gordon, W. Thies, and S. Amarasinghe. Exploiting coarse-grained task, data, and pipeline parallelism in stream programs. *ACM SIGOPS Operating Systems Review*, 40(5):151–162, 2006.

- [5] T. Harris, S. Marlow, and S. Peyton Jones. Composable memory transactions. In *PPoPP '05: Proceedings of the tenth ACM SIGPLAN symposium on Principles and practice of parallel programming*, pages 48–60. ACM Press, January 2005. ISBN 1-59593-080-9. URL <https://www.microsoft.com/en-us/research/publication/composable-memory-transactions/>.
- [6] W. A. Howard. The formulae-as-types notion of construction. In H. Curry, H. B., S. J. Roger, and P. Jonathan, editors, *To H. B. Curry: Essays on Combinatory Logic, Lambda Calculus, and Formalism*. Academic Press, 1980.
- [7] I. Lee, T. Angelina, C. E. Leiserson, T. B. Schardl, Z. Zhang, and J. Sukha. On-the-fly pipeline parallelism. *ACM Transactions on Parallel Computing*, 2(3):17, 2015.
- [8] P. Li, S. Marlow, S. Peyton Jones, and A. Tolmach. Lightweight concurrency primitives for ghc. In *Proceedings of the ACM SIGPLAN Workshop on Haskell Workshop*, Haskell '07, page 107–118. Association for Computing Machinery, 2007. ISBN 9781595936745. doi: 10.1145/1291201.1291217.
- [9] S. Marlow. Parallel and concurrent programming in Haskell. In V. Zsóok, Z. Horváth, and R. Plasmeijer, editors, *CEFP 2011*, volume 7241 of *LNCS*, pages 339–401. O'Reilly Media, Inc., 2012.
- [10] S. Marlow, R. Newton, and S. Peyton Jones. A monad for deterministic parallelism. *SIGPLAN Not.*, 46(12):71–82, 2011. ISSN 0362-1340. doi: 10.1145/2096148.2034685.
- [11] E. Meijer, M. Fokkinga, and R. Paterson. Functional programming with bananas, lenses, envelopes and barbed wire. In J. Hughes, editor, *Functional Programming Languages and Computer Architecture*, pages 124–144, Berlin, Heidelberg, 1991. Springer Berlin Heidelberg. ISBN 978-3-540-47599-6.
- [12] E. Pasarella, M.-E. Vidal, and C. Zoltan. Comparing mapreduce and pipeline implementations for counting triangles. *Electronic proceedings in theoretical computer science*, 237:20–33, 2017.
- [13] H. Röger and R. Mayer. A comprehensive survey on parallelization and elasticity in stream processing. *ACM Computing Surveys (CSUR)*, 52:1 – 37, 2019.
- [14] J. P. R. Sales, E. Pasarella, C. Zoltan, and M.-E. Vidal. Towards a dynamic pipeline framework implemented in (parallel) haskell. In *PROLE2021*. SISTEDES, 2021. URL <http://hdl.handle.net/11705/PROLE/2021/017>.
- [15] C. Zoltan, E. Pasarella, J. Araoz, and M.-E. Vidal. The Dynamic Pipeline Paradigm. In *PROLE2019*. SISTEDES, 2019. URL <http://hdl.handle.net/11705/PROLE/2019/017>.
- [16] Apache Foundation. Apache spark. <https://spark.apache.org/>. Accessed: 2021-04-27.
- [17] Simon Marlow. Haskell async library. <https://hackage.haskell.org/package/async>. Accessed: 2021-04-17.
- [18] Harendra Kumar. Haskell stream libraries benchmarks. <https://github.com/composewell/streaming-benchmarks>. Accessed: 2021-05-09.
- [19] D. Stewart, D. Coutts. Haskell bytestring library. <https://hackage.haskell.org/package/bytestring>. Accessed: 2021-04-17.

- [20] Michael Snoyman. Haskell conduit types. <https://hackage.haskell.org/package/conduit>. Accessed: 2021-05-09.
- [21] Haskell.org. Haskell containers library. <https://hackage.haskell.org/package/containers>. Accessed: 2021-04-17.
- [22] Bryan O’Sullivan. Haskell criterion tool. <https://hackage.haskell.org/package/criterion>. Accessed: 2021-04-17.
- [23] Maribel Acosta. Python diefficiency metric tool. <https://github.com/SDM-TIB/dieffpy/>. Accessed: 2021-05-03.
- [24] Juan Pablo Royo Sales. Haskell dynamic pipeline library - github repository. <https://github.com/jproyo/dynamic-pipeline>. Accessed: 2021-09-02.
- [25] Juan Pablo Royo Sales. Haskell dynamic pipeline library. <https://hackage.haskell.org/package/dynamic-pipeline>. Accessed: 2021-09-07.
- [26] M. Pickering, D. Binder, C. Heiland-Allen. Haskell eventlog2html tool. <https://mpickering.github.io/eventlog2html/>. Accessed: 2021-04-17.
- [27] Glasgow Haskell Compiler. Glasgow haskell compiler user’s guide. https://downloads.haskell.org/ghc/8.8.4/docs/html/users_guide/glasgow_exts.html#language-options. Accessed: 2021-05-09.
- [28] Haskell.org. Haskell base library. <https://hackage.haskell.org/package/base-4.15.0.0/docs/Control-Concurrent.html#v:forkIO>. Accessed: 2021-04-26.
- [29] Haskell.org. Haskell ghc-heap types. <https://downloads.haskell.org/~ghc/8.10.4/docs/html/libraries/ghc-heap-8.10.4/GHC-Exts-Heap-ClosureTypes.html>. Accessed: 2021-05-09.
- [30] haskell.org. Hackage: The haskell package repository. <https://hackage.haskell.org>. Accessed: 2021-05-09.
- [31] Apache Foundation. Apache hadoop. <https://hadoop.apache.org/>. Accessed: 2021-04-30.
- [32] S. Marlow, R. Newton. Haskell monad-par library. <https://hackage.haskell.org/package/monad-par>. Accessed: 2021-05-11.
- [33] J. Leskovec, J. Kleinberg and C. Faloutsos. Snap datasets: Stanford large network dataset collection. <https://snap.stanford.edu/data/ca-AstroPh.html>, 2007.
- [34] William Cohen. Snap datasets: Stanford large network dataset collection. <https://snap.stanford.edu/data/email-Enron.html>, 2004.
- [35] Google Inc. Snap datasets: Stanford large network dataset collection. <https://snap.stanford.edu/data/web-Google.html>, 2002.
- [36] Gabriel Gonzalez. Haskell pipes library. <https://hackage.haskell.org/package/pipes>. Accessed: 2021-05-09.
- [37] D. Kovanikov, V. Romashkina, S. Diehl, Serokell. Haskell relude library. <https://hackage.haskell.org/package/containers>. Accessed: 2021-04-17.

- [38] Glasgow Haskell Compiler. Glasgow haskell compiler user's guide. https://downloads.haskell.org/ghc/8.8.4/docs/html/users_guide/parallel.html#annotating-pure-code-for-parallelism. Accessed: 2021-04-28.
- [39] FP Complete. Haskell stack tool. <https://docs.haskellstack.org/en/stable/README/>. Accessed: 2021-05-11.
- [40] Stanford University. Snap datasets: Stanford large network dataset collection. <https://snap.stanford.edu/data/index.html>. Accessed: 2021-04-15.
- [41] Harendra Kumar. Haskell streamly types. <https://hackage.haskell.org/package/streamly>. Accessed: 2021-05-09.
- [42] S. Singh, S. Marlow, D. Jones, D. Coutts, M. Konarski, N. Wu, E. Kow. Haskell threadscope tool. <https://wiki.haskell.org/ThreadScope>. Accessed: 2021-04-17.
- [43] Brandon Simmons. Haskell unagi-chan library. <https://hackage.haskell.org/package/unagi-chan>. Accessed: 2021-04-17.