# A Dynamic Pipeline Framework implemented in Haskell[*]

Juan Pablo Royo Sales[a], Edelmira Pasarella[a,*], Cristina Zoltan[a], Maria-Esther Vidal[b,c,d]

[a] *Computer Science Department, Universitat Politècnica de Catalunya, C/Jordi Girona, 1-3, 08034, Barcelona, Spain*

[b] *TIB-Leibniz Information Centre of Science and Technology, Hannover, Germany*

[c] *Leibniz University of Hannover, Hannover, Germany*

[d] *L3S Research Centre, 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 suitable for solving problems where the incremental results' generation is critical. In this paradigm, computations are primary entities. Haskell is a purely functional programming language based on solid theoretical foundations. 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 present a Dynamic Pipeline Framework and its implementation in Haskell. As a proof of concept, we implement a solution to the problem of enumerating the weakly connected components of a graph in the proposed Dynamic Pipeline Framework. Moreover, we provide two baseline implementations of the same problem to put in perspective the computational power of the Dynamic Pipeline Framework and its continuous behavior. One baseline is implemented using a graph algorithm Haskell library, while the other one is an *ad hoc* implementation of the problem following the Dynamic Pipeline paradigm. The experimental study is conducted over a benchmark of graphs of various sizes. Two diefficiency metrics measure the continuous behavior of the three approaches, i.e., how good these approaches are while enumerating weakly connected components incrementally. The observed results provide evidence of the suitability of the Dynamic Pipeline Framework for implementing graph stream processing problems where results must be produced continuously.

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

## 1. Introduction

Effective streaming processing of large amounts of data has been studied for several years [7, 15] as a critical factor providing fast and incremental results in big data algorithmic problems. Parallel

techniques that exploit computational power as much as possible represent one of the most explored techniques, regardless of the approach. In that regard, the Dynamic Pipeline Paradigm (DPP) [30] has lately emerged as one of the models that exploit data streaming processing using a dynamic pipeline parallelism approach [15]. This computational model has been designed with a functional focus. The main components of this paradigm are functional stages or pipes that dynamically enlarge and shrink depending on incoming data.

One of the biggest challenges of implementing a Dynamic Pipeline Framework (DPF) is to find a proper set of tools and a 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 is a statically typed pure functional language designed and evolved on solid theoretical foundations where computations are primary entities. At the same time, Haskell makes available a robust set of tools for writing multithreading and parallel programs with optimal performance [18, 21].

In the context of this research, we first assess the suitability of Haskell to implement DPP solutions to stream processing problems. To be concrete, we conduct a proof of concept implementing a Dynamic Pipeline in Haskell for solving a particular and very relevant problem as the computation/enumeration of the weakly connected components of a graph. In particular, the main objective of our proof of concept is to study the critical features required in Haskell for a DPF implementation, the real possibilities of emitting incrementally results, and the performance of such kinds of implementations. Indeed, we explore the basis of an implementation of a DPF in a pure (parallel) functional language as Haskell. This is, we determine the particular features (i.e., versions and libraries) that will allow for an efficient implementation of a DPF. Moreover, we conduct an empirical evaluation to analyze the performance of the Dynamic Pipeline implemented in Haskell for enumerating WCC. To assess the incremental delivery of results, we measure the Diefficiency metrics [1], i.e., the continuous efficiency of the implementation of an algorithm for generating incremental results. Since continuous performance results are encouraging and the programming basis is clearly stated, we develop a general Dynamic Pipeline Framework written in Haskell Programming Language which allows for implementing algorithms under the Dynamic Pipeline Paradigm approach. Then, we conduct, analyze, and report experiments to measure the performance of using this framework to compute weakly connected components incrementally w.r.t. the *ad hoc* Dynamic Pipeline solution. Obtained results satisfy our expectations and encourage us to keep using the Dynamic Pipeline Framework for solving some families of graph stream processing problems where it is critical not to have to wait until the whole results are emitted. Furthermore, the experiments give us insights about the characteristics this family of problems should meet.

**Problem Research and Objective:** The main objective of this work is to design and implement a Dynamic Pipeline Framework using Haskell as programming language. 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 and set the basis of an implementation of a DPF in Haskell. This is, we determine the particular features (i.e., versions and libraries) of this language that will allow for an efficient implementation of the DPF.

**Contributions:** This paper is an extension to our work previously reported by Royo et al. [27]; we present a solution (a.k.a. $DP_{WCC}$) to the problem of enumerating the weakly connected components of a graph as a dynamic pipeline implemented in Haskell. We also empirically show the behavior of $DP_{WCC}$ concerning a solution for enumerating the weakly connected components using Haskell libraries; we name this solution Blocking WCC in Haskell `containers` ($BLH_{WCC}$). Additionally, we present the following novel contributions in this paper: i) the Dynamic Pipeline Framework

(DPF-Haskell) implemented in Haskell. ii) A program for enumerating weakly connected component implemented on top of the Dynamic Pipeline Framework (DPF-Haskell). iii) An empirical evaluation, comparing the performance of $DP_{WCC}$ in DPF-Haskell ($DPFH_{WCC}$) with respect to $BLH_{WCC}$ and $DP_{WCC}$ baseline in Haskell ($DP_{WCC}$). The continuous performance of these implementations is measured in terms of the dieffiency metric dieft [1]. The results of this study suggests that Haskell is a suitable language for implementing DPP and the programming basis for implementing DP solutions is settled. Lastly, it has been shown that implementing $DPFH_{WCC}$ performs well with respect to $DP_{WCC}$. Hence, we envision that the overhead of DPF-Haskell does not degrade the performance of dynamic pipelines implemented on the framework.

The rest of this paper is organized as follows. Section 2 analyzes the related work and positions the Haskell Dynamic Pipeline Framework with respect to existing approaches. Section 3 describes basic notions used in this work, the dynamic pipeline paradigm, the problem of weakly connected components, and the baseline solution $DP_{WCC}$. Section 4 defines the Haskell Dynamic Pipeline Framework and details the most relevant points of its Haskell implementation. Additionally, the implementation of the weakly connected component on top of Haskell Dynamic Pipeline Framework (a.k.a. $DP_{WCC}$ in DPF-Haskell ($DPFH_{WCC}$)) is sketched in section 5. Section 6 reports the results of the empirical evaluation and compares the performance and continuous behavior of the solutions. Finally, section 7 summaries our conclusions and outlooks our future work.

## 2. Related Work

*Streaming Processing.* The development of streaming processing techniques have potentiated areas of massive data processing for data mining algorithms, big data analysis, IoT applications, etc. Data Streaming (DS) has been studied using different approaches ([5, 7, 15] and see [26] for a survey) allowing to process a large amount of data efficiently with an intensive level of parallelization. There are two main different parallelization streaming computational models: Data Parallelism (DAP) and Pipeline Parallelism (PP) (see subsubsection 3.1.2). According to the DAP approach, data are separated and processed in parallel and, all computations taking place in parallel over the different subsets of data are independent of each other. A common model that has been proved successful over the last decade is MapReduce (MR) [4]. Different frameworks or tools like Hadoop [1], Spark [2], etc., support this computational model efficiently. One of the main advantages of this kind of model is the ability to implement stateless operators [26]. Data are treated in different threads or processors without the need for contextual information. However, when it is necessary to be aware of the context, parallelization is penalized, each computational step should be fully calculated before proceeding with the others. This is the case of the `reduce` operation on many frameworks or tools. This problem makes the DAP approach non-viable for delivering results incrementally. Contrary, the dynamic pipeline paradigm exploits pipeline parallelism, allowing for enlarging and shrinking pipeline stages or pipes dynamically. Moreover, each stage produces results as soon as they are ready, enabling, thus, the continuous generations of answers.

*Streaming in Haskell Language.* Streaming computational models have been implemented in Haskell Programming Language during the last ten years. One of the first libraries in the ecosystem was `conduit`[3] in 2011. After that, several efforts on improving streaming processing on the language

---

[1] https://hadoop.apache.org/
[2] https://spark.apache.org/
[3] https://hackage.haskell.org/package/conduit

has been made not only at abstraction level for the user but as well as performance execution improvements like `pipes` [4] and `streamly`[5] lately. 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. `streamly` looks as the most promising library to implement DPP. It provides a `foldrS` combinator that seems to be proper for generating a dynamic pipeline of stages based on the data flow. However, it was not possible to manipulate the channels between the stages to control the data flow. Moreover, 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 crucial motivation to develop our Dynamic Pipeline framework is that we not only want to satisfy our research needs but, as a novel contribution, we aim at providing a DPF to the Haskell community as well. We hope this contribution encourages and helps writing algorithms under the Dynamic Pipeline Paradigm. Another kind of streaming implementation in Haskell is described in [18]. 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 [20]. We do not require reaching that level of parallelization in DDF. In regard to other DP language implementations, a significant contribution to [24] has been done, where a DP implementation in Go Programming Language (Go) for counting triangles of graphs is compared against MapReduce. The reported experiment results show how DP in Go improves the performance in terms of execution time and memory depending on the graph topology. A comparison of different streaming implementations of DP represents a valuable assessment of the power of the paradigm. Nevertheless, this study is out of the scope of this paper, and it is part of our future work.

*Streaming frameworks.* Regarding performance, in general, the primary metrics considered when evaluating stream processing pipelines with massive input data are latency, throughput, and resource utilization [28]. Latency indicates how long the framework takes to deliver a result. Throughput captures the quantity of data processed within a unit of time. A good pipeline framework behavior must report low latency and high throughput. However, there are stream processing problems where obtaining results incrementally is a critical issue, as we said before. This is the case, for instance, of biological, medical, or social systems that need to identify patterns/relationships to make real-time decisions. Consequently, measuring only latency and throughput of stream processing frameworks will not illustrate the goodness of a framework when continuous performance is required. In this regard, diefficiency metrics proposed in [1] are proper tools for this kind of assessment. Regarding resource utilization metrics, one of the significant challenges when deploying a stream processing system on top of a pipeline framework is to estimate the resources consumed during all the time the system is executing. There could be peaks of consumption, but a resource over-estimation means a loss of money and possibly an overhead of the system itself since the over-administration of the resources. According to the DPP approach supported by the DPF, this framework is an elastic parallel stream processing framework. It offers users the possibility of adopting a *pay-as-you-go* policy model [23] use of resources and, hence, the chance of saving money which is a critical issue for any business. Therefore, the resource utilization metric *per se* is not representative enough in front of a stream processing system with peaks and valleys of resource consumption.

---

[4] https://hackage.haskell.org/package/pipes
[5] https://hackage.haskell.org/package/streamly

## 3. 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, the metrics used to capture continuous behavior, pipeline parallelism processing in general, and the dynamic pipeline paradigm. Next, we present an implementation of weakly connected components in Haskell following this paradigm.

### 3.1. Basic Concepts

### 3.1.1. Diefficiency Metrics

In this work, we use two relevant metrics to measure the diefficiency, i.e., continuous efficiency of a program to generate incremental results. The metrics to measure diefficiency are `dief@t` and `dief@k` [1]. The metric `dief@t` quantifies 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. The lower the value of the `dief@k` metric, the better the continuous behavior. Both metrics can be computed using Diefficency Metrics (Dm) Tool `diefpy`[6], given the traces of the execution of each of the approaches. Additionally, `diefpy` generates two different kinds of plots from an execution trace: A bi-dimensional plot and a radial plot. In the bi-dimensional plot, the x-axis represents the time when answers are generated and the y-axis represents the number of generated answers. Points $(x, y)$ are taken from traces. The radial plot contains the visual comparison of `dief@t` the metric with respects to other non-continuos metrics, such as, **i)** Completeness (Comp) which is the total number of produced answers. **ii)** Time for the first tuple (TFFT) which measure the elapsed time spent to produce the first answer. **iii)** Execution Time (ET) which measures the elapsed time spent to complete the execution of a query. **iv)** Throughput (T) which measure the number of total answers produced after evaluating a query divided by its execution time ET .

### 3.1.2. Pipeline Parallelism (PP)

A pipeline parallelism approach divides a process computation into sequential stages that are stateful operators [26]. Each stage takes the result of the previous one as an input and downstream its results to the next stage. Each *Pipeline Stage* is parallelized. 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, i.e., 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 users wait until the end of the whole data stream processing to get a result. This feature corresponds with the nature of the Dynamic Pipeline Paradigm. Hence, the PP approach is a proper parallelization streaming computational model for a Dynamic Pipeline Framework implementation. It is worth noting that although pipeline stages are parallelized in the PP approach, unbalanced intensive computation in one stage w.r.t. the other ones might delay the whole processing as a natural consequence of the sequential dependency among stages. As a result, users must be sure each stage runs high-speed computations.

### 3.1.3. Dynamic Pipeline Paradigm

The *Dynamic Pipeline Paradigm* (DPP) [30] 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 DPs. *Filter* stage instances are stateful operators in the sense described in [26], i.e., *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
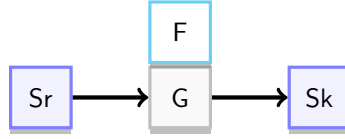
Figure 1: An initial DP consists of three stages: Sr, G together its filter parameter F, and Sk. These stages are connected through its channels –represented by right arrows– as shown in this figure.

items arrives at the initial configuration of the DP. 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; Figure 2 depicts this evolution. If the data stream is bounded, the computation ends when the lifetime of all the stages of DP has finished. Otherwise, if the stream data is unbounded, the DP remains active and results are output incrementally.
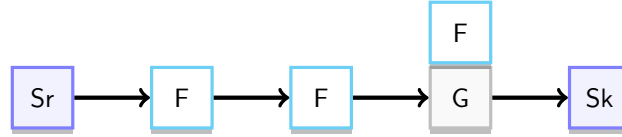
Figure 2: Evolution of a DP. After creating some filter instances (shadow Filter squares) of the filter parameter (light Filter square) in the Generator, the DP has stretched.

### 3.2. The DP$_{WCC}$ Algorithm- Weakly Connected Components of a Graph in the Dynamic Pipeline Paradigm

One of the biggest challenges of implementing a Dynamic Pipeline is to find a programming language with a proper set of tools supporting both of the primary features of the DPP: i) *parallel* processing and ii) *strong theoretical* foundations to manage computations as first-class citizens. We present and illustrate with an example, an implementation for the problem of enumerating weakly connected components using the dynamic pipeline paradigm previously described. More details in [27].

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 3. The input of the Dynamic Pipeline for computing the WCC of a graph, DP$_{WCC}$, is a sequence of edges ending with eof. In the source graph, there are neither isolated vertices nor loops. The connected components are output as soon as they are computed, i.e., they are produced incrementally. Roughly speaking,
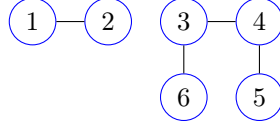
6

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

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, each filter instance receives sets of connected nodes. When in a filter the incoming set of nodes intersects with its set of connected nodes, the set of incoming nodes is joined to the latter and no output is produced. Otherwise, the incoming set of connected nodes is passed to the next stage. $DP_{WCC}$ is defined in terms of the behavior of its four kinds of stages: *Source* ($Sr_{WCC}$), *Generator* ($G_{WCC}$), *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 the 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.
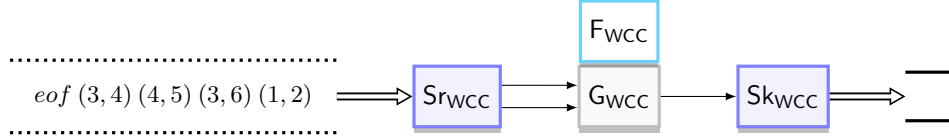


Figure 4: $DP_{WCC}$ Initial setup. Stages Source, Generator, and Sink are represented by the squares labeled by $Sr_{WCC}$, $G_{WCC}$ and $Sk_{WCC}$, respectively. The square $F_{WCC}$ corresponding to the Filter stage template is the parameter of $G_{WCC}$. Arrows $\rightrightarrows$ between represents 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 $\Longrightarrow$ 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.

Let us describe this behavior with the example of the graph shown in Figure 3. Figure 4 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 3 – feeds $Sr_{WCC}$ while $Sk_{WCC}$ emits incrementally the resulting weakly connected components. In what follows, figures 5, 6, 7, 8 and 9 depict the evolution of the $DP_{WCC}$.
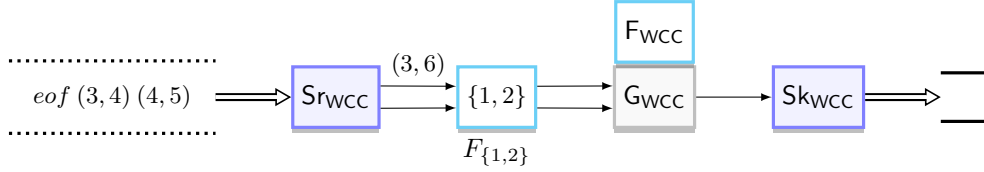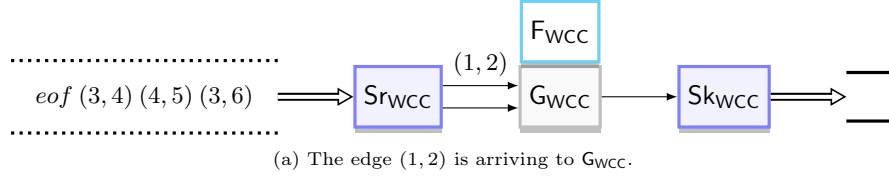
(a) The edge $(1, 2)$ is arriving to $\mathsf{G_{WCC}}$.



(b) When the edge $(1, 2)$ arrives to $\mathsf{G_{WCC}}$, it spawns a new instance of $\mathsf{F_{WCC}}$ before $\mathsf{G_{WCC}}$. Filter instance $F_{\{1,2\}}$ is connected to $\mathsf{G_{WCC}}$ through channels $\mathsf{IC_E}$ and $\mathsf{IC_{set(V)}}$. 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 5: Evolution of the $\mathsf{DP_{WCC}}$: 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 $\mathsf{IC_E}$ to $\mathsf{G_{WCC}}$.



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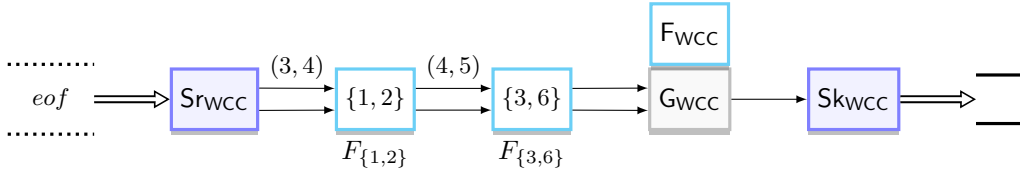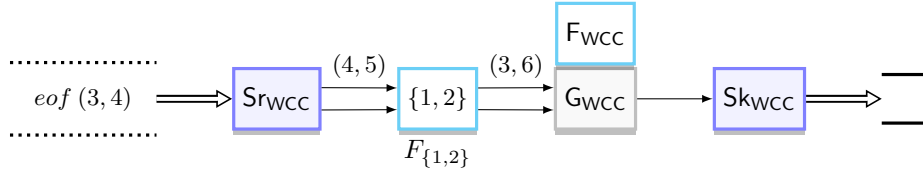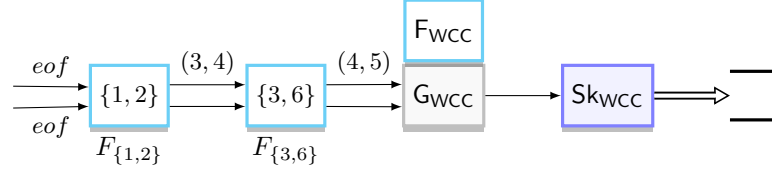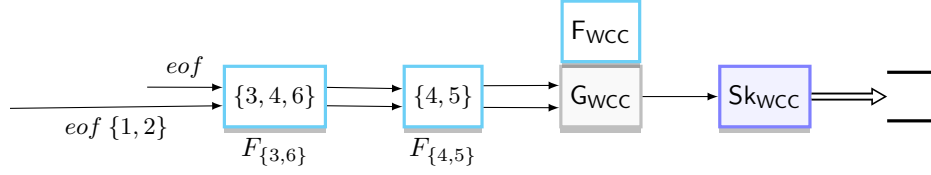
Figure 6: Evolution of the $\mathsf{DP_{WCC}}$: Second state

It is important to highlight that during the states shown in figures 5a, 5b, 6a, 6b and 7a the only actor executed in any filter instance is $\mathsf{actor_1}$ (constructing sets of connected vertices). In Figure 7a to Figure 8a although $\mathsf{actor_1}$ is still being executed in some filter instances, there are other ones starting the execution of $\mathsf{actor_2}$ (joining sets of connected vertices, when it applies).
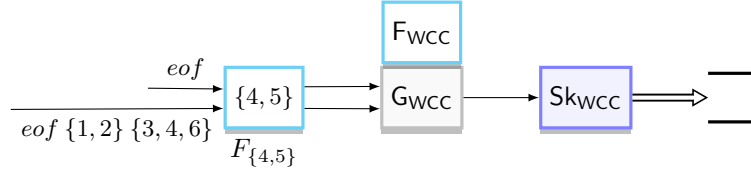
(a) $\mathsf{Sr_{WCC}}$ fed both, $\mathsf{IC_E}$ and $\mathsf{IC_{set(V)}}$, channels with the mark eof received from the input stream in previous state and then, it died. The edge $(3, 4)$ is arriving to $\mathsf{G_{WCC}}$ and the edge $(3, 4)$ is arriving to $F_{\{3,6\}}$.
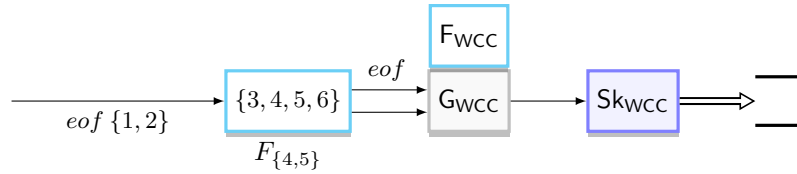


(b) When the edge $(4, 5)$ arrives to $\mathsf{G_{WCC}}$, it spawns the filter instance $F_{\{4,5\}}$ between $F_{\{3,6\}}$ and $\mathsf{G_{WCC}}$. Filter instance $F_{\{3,6\}}$ is connected to the new filter instance $F_{\{4,5\}}$ and this one is connected to $\mathsf{G_{WCC}}$ through channels $\mathsf{IC_E}$ and $\mathsf{IC_{set(V)}}$. 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 $\mathsf{IC_{set(V)}}$, this stage passes its partial set of connected vertices, $\{1, 2\}$, through $\mathsf{IC_{set(V)}}$ and dies. This action will activate $\mathsf{actor_2}$ 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 $\mathsf{IC_{set(V)}}$.

Figure 7: Evolution of the $\mathsf{DP_{WCC}}$: Third state
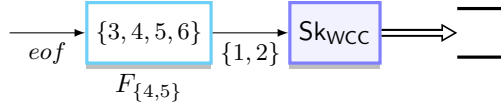


(a) The set of connected vertices $\{3, 4, 6\}$ is arriving to $F_{\{4,5\}}$. The mark eof continues passing to next stages through the channel $\mathsf{IC_E}$.
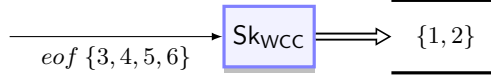


(b) Since the intersection of the set of connected vertices $\{3, 4, 6\}$ arrived to $F_{\{4,5\}}$ and its state is not empty, this state is enlarged to be $\{3, 4, 5, 6\}$. The set of connected vertices $\{1, 2\}$ is arriving to $F_{\{4,5\}}$
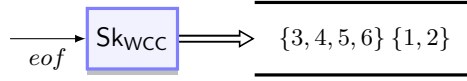
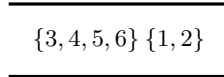Figure 8: Evolution of the $\mathsf{DP_{WCC}}$: Fourth state

9

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



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



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



(d) The weakly connected component of in the graph Figure 3 such as they have been emitted by $\mathsf{DP_{WCC}}$.

Figure 9: Last states in the evolution of the $\mathsf{DP_{WCC}}$

## 4. Dynamic Pipeline Framework in Haskell

This section presents the design and the main features of the implementation of a general Dynamic Pipeline Framework (Haskell). To be concrete, we present the approach that we follow for designing DPF-Haskell, the system architecture and the most relevant implementation details.

### 4.1. Framework Design

A suitable framework should provide users with the proper level of abstraction that hides underlying details and allows developers to focus on the problem to be solved. We follow a Domain-specific Language (DSL) approach [6] to develop the Dynamic Pipeline Framework in Haskell; it provides users with a language to write their solutions. There exists two types of Domain-specific Languages [14]: External Domain-specific Language (DSL) and Embedded Domain-specific Language (EDSL) [11]. The external Domain-Specific Languages (DSL) are completely new languages requiring the development of an interpreter. On the contrary, the Embedded Domain-Specific Languages (EDSL) are languages syntactically embedded in a host language. Accordingly, actually users write their codes in the host language restricted to the EDSL abstractions. In particular, DPF-Haskell follows an EDSL approach where Haskell is the host language. This approach allows for to take advantage of the strong type Haskell system to providing correctness at type-level [10]. In next definition the Dynamic Pipeline Domain Specific Language (DP-EDSL) is formally defined.

**Definition 1.** Syntactical constructions of the Dynamic Pipeline Domain Specific Language (DP-EDSL) are generated by the grammar $(N, \Sigma, DP, P)$ where

$$N = \{DP, S_r, S_k, G, F_b, CH, CH_s\},$$
$$\Sigma = \{\texttt{Source}, \texttt{Generator}, \texttt{Sink}, \texttt{FeedbackChannel}, \texttt{Type}, \texttt{Eof}, \texttt{:=>}, \texttt{:<+>}\},$$

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

The configuration of the initial structure of the DP is stated in terms of their connections through channels and the specification of the types of data these channels carry. This configuration is encoded using the DP-EDSL language. The symbols `:=>` and `:<+>` allow to separate the encoding of the pipeline stages (Sr, G, Sk) from the encoding of channel composition in the same stages, respectively.

**Example 1.** *Consider the problem of defining a DP for removing duplicate elements from a string. The stages are connected by a single channel carrying data of type* `Int`.

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

Source Code 4.1: Encoding in DP-EDSL of a DP for removing duplicate elements in a stream. Note that the DP is defined at type level of the host language (Haskell).

***System Architecture.*** The framework has the three components (Figure 10): EDSL, IDL and RS.



Figure 10: Architecture of the DPF-Haskell. DPF-Haskell is built on three main components: EDSL, IDL and RS. Users interact with the EDSL and the IDL components. EDSL component allows users to enter the initial stages and channels specifications. IDL component supports users to define the Haskell functions corresponding to each stage of the DP. Finally RS execute all that definition plus functions. Execution layer indicates an example of a DPP running after being executed.

*EDSL Component.* This component validates type correctness of the configuration of the dynamic pipeline to be implemented on top of the DPF-Haskell. To do it, it receives from users the DP-EDSL

12

encoding of the initial structure of the DP and a precise specification of the channels connecting the stages of the pipeline.

*IDL Component.* This component is the Interpreter of the DP specification. According to the structure of the dynamic pipeline provided in the DP-EDSL language, it creates skeletons of the stages to allow users define each stage in terms of its functions. Since the Filter stage is a parameter of the Generator stage, its specification is obtained when Generator is defined. Once users have entered the Haskell definition of all the stages, the completely defined DP is passed to the Runtime System Component.

*RS component.* This is the Run System component. It receives a Haskell dynamic pipeline definition, compiles it and launch its execution.

*4.2. Implementation*

This section presents implementation details of each system component: EDSL, IDL and RS.

```
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 4.2: This code is showing most of the data types that represent the same terminal symbols $\Sigma$ in $G_{dsl}$. These types indexed by another kind `Type`, allows us to store information at type-level needed for interpret the DP-EDSL

```
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 4.3: Special terminal symbols $\{$`:<+>`, `:=>`$\} \subset \Sigma$. These terminal symbols allow us to index two types, in order to combine several of them and build a chain of stages (using `:=>`) and a set of channels (using `:<+>`).

*4.2.1. EDSL Validation*

First, let us pay attention to the embedding of the DP-EDSL into Haskell. For this purpose we use an Index type [9] that allows for keeping track –at type-level– of the extra information required by the DP definition, i.e., channels and data types carried by the channels. In Source Code 4.2, an *Index type* for each element of $\Sigma$ an *Index type* is defined to encoded them in Haskell types. Note that `Sink` and `Eof` are not indexed. This is because these symbols do not carry 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` is just a terminal type to disambiguate the `Channel (a :: Type)` of the tree branch. `Channel` can carry any type of data because it is polymorphic to support different number of channels and data types. In Source Code 4.3, the type definition shows how `:=>` and `:<+>` can combine two types. Writing `:=>` and `:<+>` as types is to have a syntactic sugar type combinator in the DP-EDSL.

The DP-EDSL code (with the syntactical constructions from Definition 1) provided by users must be compiled. Fortunately, Haskell provides several type-level techniques [25] that allow for verifying properties of programs before running them. This prevents users to introduce bugs and hence to reduce the occurrence of errors. This verification done by the compiler establish a Curry-Howard Isomorphism [10], i.e., *Propositions as Types - Programs as Proof*. It is important to remark that Haskell is not a theorem prover system like Coq[7], but some verifications, as we present in this work, can be done by the GHC (Glasgow Haskell Compiler) to ensure correctness on programs. Although, Haskell provides tools to build advanced type-level verification, all these techniques require the use of *Haskell Language Extensions*.

The implementation of the validation is done using *Associated Type Families* [2]. In Source Code 4.4, there are three Type families that helps to validate the DP-EDSL code. `IsDP` associated type family is checking the production rules $P$ of the grammar defined in Definition 1, returning a promoted data type [29] (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 no restriction, i.e., if `ValidDP` is used as a constraint. It is fully applied to `()` and gives the compiler the evidence that there is no error at type-level. `ValidDP` is also associating `'False` with a custom `TypeError`; it will appear at compilation time if the DP-EDSL definition fully applies to that – a type checked error –.

```
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 4.5: Definition of `mkDP` function of the Framework which uses type-level validation of DP-EDSL code `ValidDP (IsValid Type)`. Last line of the code is showing that using that function will compile-time check the definition of `DPExample` type (Source Code 4.1)

---

[7]https://coq.inria.fr/

### 4.2.2. Interpreter of DSL (IDL)

IDL component takes the dynamic pipeline specification made on with DP-EDSL component to interpret and generate the function definitions that needs to be implemented for solving a specific

```haskell
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
           ↪  genToOut)))
9      IsDP ( Source (Channel inToGen) :=> Generator (Channel genToOut) :=>
        ↪  FeedbackChannel toSource :=> Sink)
10         = And (IsDP (Source (Channel inToGen))) (IsDP (Generator (Channel
           ↪  genToOut)))
11     IsDP (Source (Channel (a :<+> more)))
12         = IsDP (Source (Channel more))
13     IsDP (Source (Channel Eof))              = 'True
14     IsDP (Generator (Channel (a :<+> more)))  = IsDP (Generator (Channel more))
15     IsDP (Generator (Channel a))             = 'True
16     IsDP x                                   = 'False
17
18 type family ValidDP (a :: Bool) :: Constraint where
19   ValidDP 'True = ()
20   ValidDP 'False = TypeError
21                    ( 'Text "Invalid Semantic for Building DP Program"
22                     ':$$: 'Text "Language Grammar:"
23                     ':$$: 'Text "DP        -> Source CHANS :=> Generator CHANS
                      ↪   :=> Sink"
24                     ':$$: 'Text "DP        -> Source CHANS :=> Generator CHANS
                      ↪   :=> FEEDBACK :=> Sink"
25                     ':$$: 'Text "CHANS     -> Channel CH"
26                     ':$$: 'Text "FEEDBACK -> FeedbackChannel CH"
27                     ':$$: 'Text "CH        -> Type :<+> CH | Eof"
28                     ':$$: 'Text "Example: 'Source (Channel (Int :<+> Int)) :=>
                      ↪   Generator (Channel (Int :<+> Int)) :=> Sink'"
29                    )
```

Source Code 4.4: Type Families `And`, `IsDP` and `ValidDP` which allows to perform a type-level validation over a DP-EDSL CFG definition.

problem. In subsubsection 3.1.3, we describe what should be provided in a DPP algorithm. The IDL generates the function definitions with an empty implementation to be entered by users, ensuring that those functions will give "Proof" – in terms of Curry-Howard Correspondence [10] – of the "Propositions" defined on the DP-EDSL. IDL utilizes techniques similar to the ones used in subsubsection 4.2.1. First, *Type-level Defunctionalization* [3, 13] is used to let the compiler generates the signatures of the required functions. Second, *Term-level Defunctionalization* interprets those functions. Lastly, *Indexed Types* [9] and *Heterogeneous List* [12] 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
   ↪  dpDefinition filter (DP st)
6                  -> Stage (WithGenerator dpDefinition filter (DP st))
7  withGenerator = mkStage' @(WithGenerator dpDefinition filter (DP st))
8
9  withSink :: forall (dpDefinition :: Type) st. WithSink dpDefinition (DP st)
10              -> Stage (WithSink dpDefinition (DP st))
11  withSink = mkStage' @(WithSink dpDefinition (DP st))
```

Source Code 4.6: This code is showing the different interpreters combinators to support users to generate the functions of the stages of the dynamic pipeline

In Source Code 4.6, we can appreciate the different combinators of the IDL that helps to interpret the DSL and generates the function definitions. `Stage` data type will be cover in Source Code 4.8, 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* definition, there is 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 –.

In Source Code 4.7, 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 -> ... -> monadicAction ()` depending on DPP language definition.

In Source Code 4.8, `Stage` data type uses a `Proxy` type. This `Proxy` type allows `Stage` to index the type definition generated by `a`. For example, in Source Code 4.6, when `withSource` interpreter is applied to `WithSource dpDefinition`, the compiler is provided with `dpDefinition` DSL type, expanding the function signature belonging to that DPP definition inside the `Stage`.

*Generator and Filter.* According to DPP definition in subsubsection 3.1.3, the Generator (G) stage has a Filter (F) template as parameter. In order to know how to dynamically interpose a new F

during the runtime execution of the program let us first study F Data Type in the context of the framework.

In Source Code 4.9, the definition of the `Filter` data type contains a non-empty set of `Actor`. 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

```
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 4.7: 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`

```
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 4.8: `Stage` data type for implementing *Term-level Defunctionalization* providing evidence to the Type-Level Associated types

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

Source Code 4.9: This code shows the definition of the `Filter` data type which contains a non-empty set of `Actor`. In the Context of the `MonadState` the data type of `Actor` is `Stage`. This allows for keeping a local memory in the execution context of the filter.

is because the F instance should have an state, according to DPP definition in subsubsection 3.1.3. `Actor` data type – see Source Code 4.9 –, 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.

Finally, in Source Code 4.10, 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.

In Source Code 4.11, 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 algorithms defined by users according to the data received from its input channels. `Generator` has also another field with the behavior of the G – a `Stage` –.

### 4.2.3. Runtime System (RS)

The RS can be divided into two parts: the mechanism to generate stages dynamically at runtime and the execution entry point of the DP. Regarding execution entry point, all the stages defined above 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 users definitions, i.e., an *anamorphism* [22] that creates F instances until some condition is met.

18

```
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 4.10: Combinators and small constructor to enable building actors and filter.

```
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 4.11: `Generator` Data type which contains the `Stage` code of the generator itself, and the `Filter` template that can be spawned by the `Generator`.

```
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 4.12: `unfolF` is the *anamorphism* combinator to spawn new `Filter` types between the `Generator` and previous stages.

In Source Code 4.12, it is presented how is the *anamorphism* mechanim that generates dynamic stages between G and the previous stages. This *anamorphism* is implemented with the function unfoldF. This 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. 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.

Several smart constructors are also provided for building UnfoldFilter Data Type. In Source Code 4.13 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) The 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 a smart constructor of UnfoldFilter that allows to spawn a

```
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
                           ↪   filterParam l
16
17 mkUnfoldFilterForAll :: Filter dpDefinition filterState filterParam st
18                         -> (readElem -> filterState)
19                         -> ReadChannel readElem
20                         -> HList l
21                         -> UnFoldFilter dpDefinition readElem st filterState
                           ↪   filterParam l
```

Source Code 4.13: Combinators for building `UnfoldFilter` types indicating the type of the `unfold` that users want to achieve.

new filter for each element received in the G.

### 4.3. Libraries and Tools

One of the most important task 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 [21]. 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 [18]. A straightforward assumption to achieve, here, is to use `monad-par` library[8]. But, in this first version of the DPF-Haskell we do not use sparks [19, 20] but spawning green threads only. We will consider the use of sparks for future versions. Another choice is to use: `forkIO :: IO () -> IO ThreadId` from `base` library[9]. However, that would imply handling all the threads lifecycles and errors programmatically without any abstraction to facilitate that complex task. Therefore, we choose

---

[8]`https://hackage.haskell.org/package/monad-par`
[9]`https://hackage.haskell.org/package/base-4.15.0.0/docs/Control-Concurrent.html`

`async` library[10] which enables to spawn asynchronous computations [18] on Haskell using green threads, and at the same time, it provides combinators to managing thread terminations and errors.

Regarding channels, there are several techniques to communicate threads or sparks in Haskell like `MVar` or concurrent safe mechanisms like Software Transactional Memory (STM) [8]. At the same time, in Haskell library ecosystem, there are `Channels` abstractions based on previous mentioned communication techniques. In that sense, for conducting the communication between dynamic stages and data flowing in a dynamic pipeline, we have selected `unagi-chan` library [11] which provides the following advantages to our solution: Firstly, `MVar` channel without using STM reducing overhead. STM is not required in a dynamic pipeline because each specific stage running in a separated thread, can only access to its `I/O` channels for reading/writing accordingly, and these 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[12] on Benchmarking compare with STM and default base `Chan` implementations.

## 5. Enumerating Weakly Connected Component on the DPF

This section presents the most relevant details of the implementation of the $DP_{WCC}$ using DPF-Haskell. 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 4, 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. Next, we present the minimal code needed for encoding any DPP using DPF-Haskell. [13]

```
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 5.1: 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 5.1, 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* ($G_{WCC}$)

---

[10]https://hackage.haskell.org/package/async
[11]https://hackage.haskell.org/package/unagi-chan
[12]https://github.com/jberryman/unagi-chan
[13]All the code that we expose here can be accessed publicly in https://github.com/jproyo/dynamic-pipeline/tree/main/examples/Graph

and *Sink* (Sk$_{WCC}$), that matches those declaration. As we can see in Source Code 5.1, highlighted lines 1-3 matches one to one the definition in Figure 4, although in the case of the framework it is not required to provide *Filter* (F$_{WCC}$) definition because DPF-Haskell will deducted from `Generator`. According to this declaration what we need to provide is the correct implementation of `source'`, `generator'` and `sink'` which *Type checked* the DPP type definition[14].

```
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    sink' = withSink @DPConnComp $ \_ cc -> withDP $ foldM_ cc print
9
10   generator' :: GeneratorStage DPConnComp ConnectedComponents Edge st
11   generator' =
12     let gen = withGenerator @DPConnComp genAction
13     in  mkGenerator gen filterTemplate
```

Source Code 5.2: 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 5.2, 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 process is achieved with a DPF-Haskell combinator called `unfoldFile` which is a *catamorphism* of the input data to the stream. `sink'` delivers to the output of the program the upstream connected components received from previous stages. Sk$_{WCC}$ implementation is done 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 3.2. According to what we described in subsubsection 3.1.3, *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, there are two important functions to describe: `genAction` which tells how to spawn a new *Filter* and under what circumstances, and `filterTemplate` which carries the function to be spawn.

---

[14]The 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    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
       ↪  (readCC .*. HNil)
9      results <- unfoldF unfoldFilter
10     foldM_ (hHead results) (`push` writeCC)
```

Source Code 5.3: 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 4. `genAction` for DP$_{\text{WCC}}$ will use the combinator `mkUnfoldFilterForAll` which will spawn one *Filter* per received edge in the channel, expanding dynamically the stages on runtime. In line 10, 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 5.4: Filter template code composed by 2 Sequential Actors that will calculate the Connected Components and downstream them.

Finally, the *Filter* template code is defined in Source Code 5.4. As we have seen in subsection 3.2, DP$_{WCC}$ *Filter* is composed of 2 Actors. The first actor collect all the possible vertices that are incidence to some vertices edge that was instantiated 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 component, this is `actor2`. At the end of processing, `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.

## 6. Empirical Evaluation

The empirical study aims at evaluating the proposed Dynamic Pipeline Framework. As described previously, we have implemented three solutions to the problem of weakly connected components:

- $BLH_{WCC}$: implemented in `containers` Haskell library [15] using `Data.Graph`.

- $DP_{WCC}$: baseline implementation of the weakly connected components as a dynamic pipeline (subsection 3.2).

- $DPFH_{WCC}$: implementation of the weakly connected components on top of the DPF-Haskell framework (section 5).

Our goal is to answer the following research questions:

**RQ1)** Does $DPFH_{WCC}$ exhibit similar execution time performance compared with $DP_{WCC}$? **RQ2)** Does $DPFH_{WCC}$ enhance the continuous behavior with respect to other approaches, i.e., $BLH_{WCC}$ and $DP_{WCC}$? **RQ3)** Do the proposed approaches handle memory efficiently?

We have configured the following setup to assess our research questions.

*DataSets.* The experiments are executed over networks of the benchmark SNAP [16]. The selected networks correspond to complex undirected graphs with different sizes and connected components; Table 1 reports on the main characteristics of these undirected graphs.

| Network | Nodes | Edges | Diameter | #WCC | #Nodes Largest WCC |
|---|---|---|---|---|---|
| Enron Emails | 36,692 | 183,831 | 11 | 1,065 | 33,696 (0.918) |
| Astro Physics Collaboration Net | 18,772 | 198,110 | 14 | 290 | 17,903 (0.954) |
| Google Web Graph | 875,713 | 5,105,039 | 21 | 2,746 | 855,802 (0.977) |

Table 1: DataSet of Graphs Selected

*Haskell Setup.* We use the following Haskell setup for $DP_{WCC}$ implementation: GHC version 8.10.4, `bytestring` 0.10.12.0[16], `containers` 0.6.2.1[17], `relude` 1.0.0.1[18] and `unagi-chan` 0.4.1.3[19]. The use of `relude` library is because we disabled `Prelude` from the project with the language extension (language options) `NoImplicitPrelude`[20]. 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 [21] as a building tool on top of GHC the compilation command is `stack build`[22]. The setup for $DPFH_{WCC}$ is the same, except that this setup use `dynamic-pipeline` 0.3.2.0[23] library written in the context of this work.

---

[15] https://hackage.haskell.org/package/containers
[16] https://hackage.haskell.org/package/bytestring
[17] https://hackage.haskell.org/package/containers
[18] https://hackage.haskell.org/package/containers
[19] https://github.com/jberryman/unagi-chan
[20] https://downloads.haskell.org/ghc/8.8.4/docs/html/users_guide/
[21] https://docs.haskellstack.org/en/stable/README/
[22] For more information about package.yaml or cabal file, please check https://github.com/jproyo/upc-miri-tfm/tree/main/connected-comp
[23] https://hackage.haskell.org/package/dynamic-pipeline

*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.

*Implementation.* The current version is published on the Github Repository at `https://github.com/jproyo/dynamic-pipeline`. Moreover, DPF-Haskell is available in *Hackage: The Haskell Package Repository*, `https://hackage.haskell.org/package/dynamic-pipeline`.

### 6.1. Experiments Definition

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

*E2: Benchmark Analysis.* This experiment measures *Average Running Time*. The *Average Running Time* is the average running time of 1000 resamples using `criterion` tool. In each sample, the running time is measure from the beginning of the execution of the program until when the last answer is produced. This experiment will help to answer research question [RQ1].

*E3: Continuous Behavior - Dieffiency Metrics.* In this experiment, we conduct two benchmark analysis over execution time comparing $DP_{WCC}$ vs. $BLH_{WCC}$. In the first benchmark analysis, we use `criterion` tool in Haskell which runs over four iterations of each of the algorithms to compute a mean execution time in seconds and compare the results in a plot. In the second benchmark, we use the Dieffiency Metrics (Dm) Tool `diefpy` 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. Based on the reported values of metrics `dief@t` and `dief@k`, we aim at answering research question [RQ2].

*E4: 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 $DP_{WCC}$ is handling parallelization and memory usage. This experiment is performed using two tools, *ThreadScope* Tool[24] for conducting multithreading analysis and *eventlog2html* Tool[25] to assess memory usage analysis. Regarding multithreading analysis, we compute 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. Regarding memory management, we compute 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].

---

[24]`https://wiki.haskell.org/ThreadScope`
[25]`https://mpickering.github.io/eventlog2html/`

## 6.2. Discussion of Observed Results

### 6.2.1. Experiment: E1

Table 2 reports on the execution of $DP_{WCC}$ in the evaluated graphs. 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.

| 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: Total Execution times of each of the networks implemented with $DP_{WCC}$. *MUT Time* is the time of running or executing code and *GC Time* is the time that the program spent doing Garbage collection. *Total Execution time* is the sum of both times

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 output correctness, we have verified that the number of connected components is the same as the metrics already gathered in Table 1. Table 3 represents the *Total execution time* for each of the networks implemented with $DPFH_{WCC}$. We observe similar execution times compared with Table 2 in subsection 3.2. In fact, for *Enron Emails* and *Astro Physics Coll* networks, all times are better than the implementation of the *Proof of Concept*. Regarding *Google Web* network, the time is slightly worse, but there is the fact that DPF-Haskell is adding some overhead over plain code execution. According to these results, we can partially answer [RQ1], because the implementation of $DPFH_{WCC}$ has similar performance compared with $DP_{WCC}$ implementation.

| Network | Exec Param | MUT Time | GC Time | Total Time |
|---|---|---|---|---|
| Enron Emails | `+RTS -N4 -s` | 1.795s | 0.505s | 2.314s |
| Astro Physics Coll Net | `+RTS -N4 -s` | 2.294s | 1.003s | 3.311s |
| Google Web Graph | `+RTS -N8 -s` | 169.381s | 270.784s | 440.176s |

Table 3: Total Execution times of each of the networks implemented with $DPFH_{WCC}$. *MUT Time* is the time of running or executing code and *GC Time* is the time that the program spent doing Garbage collection. *Total Execution time* is the sum of both times

## 6.2.2. Experiment: E2

*Criterion Benchmark.* In Figure 11, orange bars report the time taken by `Data.Graph` in $DP_{WCC}$ in Haskell `containers` library[26]. Blue light bars represent the time taken by $DP_{WCC}$ in Haskell.
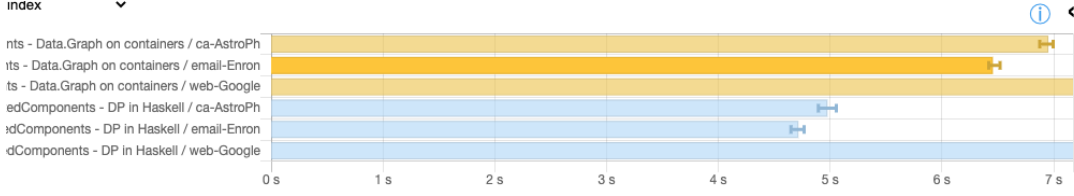


Figure 11: Benchmark of average execution time of $DP_{WCC}$ vs. $BLH_{WCC}$. Average Execution time of running 1000 samples over each of the networks.

Figure 11 shows that $DP_{WCC}$ solution is 1.3 faster compare with $BLH_{WCC}$. Despite this, if we zoom in Figure 11, it can be observed that $DP_{WCC}$ solution is slower compared with $BLH_{WCC}$. Regarding mean execution times for each implementation on each case measure by `criterion` library, we can display the following results:

| Network | $DP_{WCC}$ | $BLH_{WCC}$ | 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 4: Comparison of Mean Execution times between $DP_{WCC}$ vs. $BLH_{WCC}$ for each network

These results allow for answering Question [Q2]. We already had a partial answer with the previous experiment E1 about [Q2] (section 6) where we have seen that the graph topology is affecting the performance and the parallelization, penalizing $DP_{WCC}$ for this particular case. In this benchmark, the solution against $BLH_{WCC}$ confirms the hypothesis.

Moreover, Figure 12, we can appreciate the *Average Execution Time* for each of the networks after running `criterion` tool for the $DPFH_{WCC}$ implementation confirming that after changing the DPP implementation using DPF-Haskell the performance remains stable.
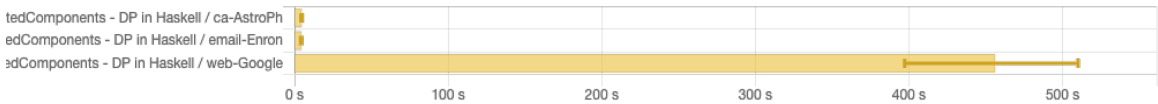


Figure 12: Average Execution time of running 1000 samples over each of the networks with $DPFH_{WCC}$ implementation.

Table 5 reports on the comparison of *Average Execution Times* with the previous obtained in subsection 3.2. All the *Average Execution times* are better for the new implementation with $DPFH_{WCC}$ compared with the $DP_{WCC}$. As it is consistent with the *Total Execution time* of the previous experiment in subsubsection 6.2.1, the only network that perform slightly worse is *Google Web*, but the difference is not significant enough taking into consideration the overhead introduced

---

[26]https://hackage.haskell.org/package/containers

29

by DPF-Haskell. These results support research question [RQ1] confirming that the performance in terms of execution is better for smaller networks and competitive in complex networks.

| Network | DPFH$_{\text{WCC}}$ | DP$_{\text{WCC}}$ | Speed-up |
|---|---|---|---|
| Enron Emails | 4.30s | 4.68s | 0.91 |
| Astro Physics Coll Net | 4.76s | 4.98s | 0.95 |
| Google Web Graph | 456s | 386s | -1.18 |

Table 5: Comparison of Average Execution Times between DPFH$_{\text{WCC}}$ and the implementation of baseline implementation (DP$_{\text{WCC}}$)

*6.2.3. Experiment: E3*

Some considerations are needed before starting to analyze the data gathered with the Dm tool. First, the tool is plotting the results according to the traces generated by the implementation, both DP$_{\text{WCC}}$ and BLH$_{\text{WCC}}$. By the nature of the DPP model, we can register that timestamps as long as the model is generating results. In the case of BLH$_{\text{WCC}}$, 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 BLH$_{\text{WCC}}$ are generated. In those cases, we observe a straight vertical line. 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; it shows how the solution is behaving on the dimensions of Time for the first tuple (TFFT) inverse, Execution Time (ET) inverse, Throughput (T), Completeness (Comp) and Dieficiency Metric `dief@t` (`dief@t`) and how are the tension between them. All these metrics are higher is better. Acosta et al. [1] detail how these metrics are computed and their main properties.
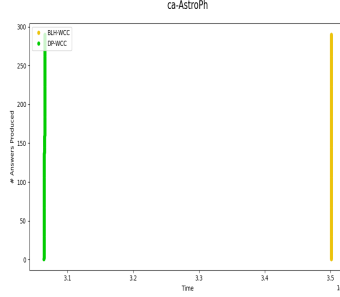
| Network | Implementation | dief@t Metric | dief@k Metric (0.5) |
|---|---|---|---|
| ca-AstroPh | DP$_{\text{WCC}}$ | $1.99 \times 10^5$ | $3.93 \times 10^4$ |
| | BLH$_{\text{WCC}}$ | 0 | 0 |
| email-Enron | DP$_{\text{WCC}}$ | $2.51 \times 10^6$ | $7.06 \times 10^5$ |
| | BLH$_{\text{WCC}}$ | 0 | 0 |
| web-Google | DP$_{\text{WCC}}$ | $1.10 \times 10^7$ | $3.10 \times 10^6$ |
| | BLH$_{\text{WCC}}$ | 0 | 0 |

Table 6: This tables shows the `dief@t` and `dief@k` values gather for DP$_{\text{WCC}}$ vs. BLH$_{\text{WCC}}$. `dief@k` is showing the metric of $k$ values generated at percentile 0.5. We can appreciate that in all cases DP$_{\text{WCC}}$ has a higher value of `dief@t` and a lower value of `dief@k` showing continuos behavior
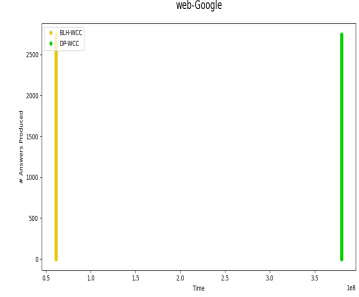
Based on the results shown in Figure 13 and in Table 6, all the solutions in DP$_{\text{WCC}}$ are being generated incrementally. This behavior is measured with the values of `dief@t` and `dief@k` in Table 6. *ca-AstroPh* implemented with DP$_{\text{WCC}}$ has a higher value of `dief@t` and a lower value of `dief@k`, exhibiting more continuous behavior compare with the rest. The *web-Google* using DP$_{\text{WCC}}$ also presents continuous behavior, but its `dief@k` value is not low for the percentile 0.5, meaning that the 50% of the $k$ results are delivered near the end of the execution. We have detected that this behavior appears because *web-Google* concentrates the majority of the vertices in few WCC as we can see in Table 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. Regarding BLH$_{\text{WCC}}$, the `dief@t` and `dief@k` metrics values are 0, indicating, thus, no continuous behavior.
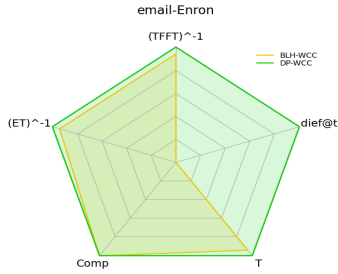
(a) email-Enron Diefficency Metrics DP_WCC vs. BLH_WCC Line Plot

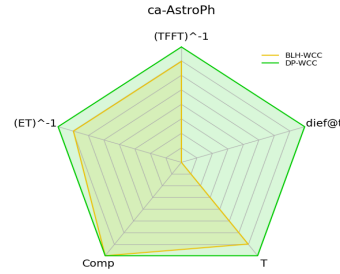(b) ca-AstroPh Diefficency Metrics DP_WCC vs. BLH_WCC Line Plot

(c) web-Google Diefficency Metrics DP_WCC vs. BLH_WCC Line Plot

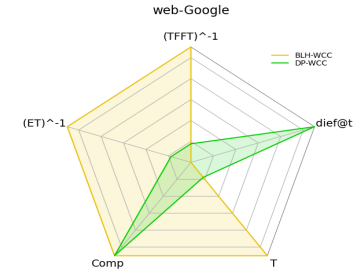Figure 13: These figures show `dief@t` observed results after running all the scenarios for each network with $DP_{WCC}$ (green) vs. $BLH_{WCC}$ (yellow). $y$ axis represents the number of Answers produced and $x$ axis is the $t$ time of the `dief@t` metric describe in subsubsection 3.1.1. The more data points distributed throughout the $x$ axis, the higher, the continuous behavior. The scale in axis $x$ which represents *Time t*, where $t$ is the microsecond difference between the start of the program and the result delivered by it.



(a) email-Enron Diefficency Metrics DP_WCC vs. BLH_WCC Radar Plot

(b) ca-AstroPh Diefficency Metrics DP_WCC vs. BLH_WCC Radar Plot

(c) web-Google Diefficency Metrics DP_WCC vs. BLH_WCC Radar Plot

Figure 14: Radial plots show how the different dimensions values provided by `diefpy` tool such as T, TFFT, `dief@t`, ET and Comp are related each other for each experimental case. These figures show radial plot observed results after running for each network comparing $DP_{WCC}$ (green) vs. $BLH_{WCC}$ (yellow) implementations. `dief@t` is described in subsubsection 3.1.1.

31

As we can appreciate in Figure 14 radar plots, our previous analysis can be confirmed. In all cases, the only approach that has `dief@t` greater than 0 is $DP_{WCC}$ indicating continuous behavior. In the case of $BLH_{WCC}$, we can see that in any of the plots, the yellow region is extending to `dief@t` metric, showing no continuity at all. In the case of *web-Google*, we notice on the plot that the TFFT, T and ET are better in $BLH_{WCC}$, but that only indicates that is faster in execution time as Table 4 already indicates. In spite of this, and as we mentioned before, $BLH_{WCC}$ exhibits no continuous behavior in *web-Google*. Regarding [Q2] (section 6), although $DP_{WCC}$ is faster than $BLH_{WCC}$, the speed-up dimension execution factor is not always the most interest analysis. As we have seen even when, in the case of *web-Google* Graph, $DP_{WCC}$ is slower at execution, it is at least generating incremental results without the need to wait for the rest of the computations.

We also compare $DP_{WCC}$ and $DPFH_{WCC}$, using the same Diefficency Metrics. We aim at confirming that $DPFH_{WCC}$ is still presenting continuous behavior and generates incremental results.
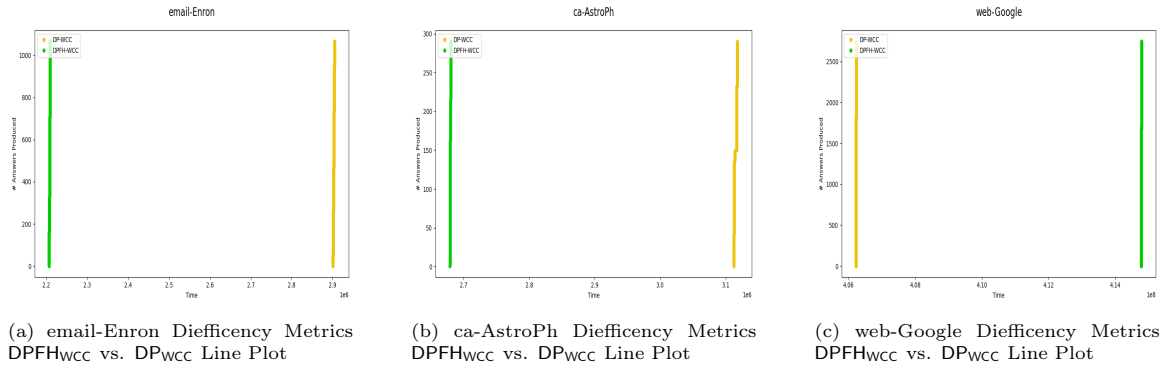


(a) email-Enron Diefficency Metrics $DPFH_{WCC}$ vs. $DP_{WCC}$ Line Plot

(b) ca-AstroPh Diefficency Metrics $DPFH_{WCC}$ vs. $DP_{WCC}$ Line Plot

(c) web-Google Diefficency Metrics $DPFH_{WCC}$ vs. $DP_{WCC}$ Line Plot

Figure 15: These figures show `dief@t` observed results after running all the scenarios for each network with $DPFH_{WCC}$ (green) vs. $DP_{WCC}$ (yellow). $y$ axis represents the number of Answers produced and $x$ axis is the $t$ time of the `dief@t` metric describe in subsubsection 3.1.1. The more data points distributed throughout the $x$ axis, the higher, the continuous behavior. The scale in axis $x$ which represents *Time t*, where $t$ is the microsecond difference between the start of the program and the result delivered by it.

Based on the results shown in Figure 15 and with the support of the metrics values in Table 7; it can be appreciated that both solutions, $DP_{WCC}$ and $DPFH_{WCC}$ show continuous behavior. Moreover, in Table 7, $DPFH_{WCC}$ has a higher value of `dief@t` and a lower value of `dief@k` for all the networks confirming the continuous behavior. In the case of *web-Google* which is the biggest network,

| Network | Implementation | dief@t Metric | dief@k Metric (0.5) |
|---------|----------------|---------------|---------------------|
| ca-AstroPh | $DPFH_{WCC}$ | $1.80 \times 10^5$ | $2.00 \times 10^4$ |
| | $DP_{WCC}$ | $8.77 \times 10^5$ | $1.38 \times 10^5$ |
| email-Enron | $DPFH_{WCC}$ | $1.38 \times 10^6$ | $4.34 \times 10^5$ |
| | $DP_{WCC}$ | $1.98 \times 10^6$ | $5.23 \times 10^5$ |
| web-Google | $DPFH_{WCC}$ | $1.29 \times 10^7$ | $3.08 \times 10^6$ |
| | $DP_{WCC}$ | $1.17 \times 10^7$ | $2.88 \times 10^6$ |

Table 7: This table shows the `dief@t` and `dief@k` values gather for $DP_{WCC}$ vs. $DPFH_{WCC}$. `dief@k` is showing the metric of $k$ values generated at percentile 0.5. We can appreciate that $DP_{WCC}$ has a higher value of `dief@t` for the smaller networks and $DPFH_{WCC}$ has a higher value of `dief@t` for web-Google. In both cases, continuous behavior can be appreciated.

DPFH$_{WCC}$ seems to be more continuous than DPFH$_{WCC}$ according to its `dief@t` and `dief@k` values. In the case of *ca-AstroPh* and *email-Enron*, DP$_{WCC}$ presents slightly bigger `dief@t` and smaller `dief@k` values compared with DPFH$_{WCC}$. That does not mean DPFH$_{WCC}$ is not presenting continuous behavior, because `dief@t` and `dief@k` values for DPFH$_{WCC}$ are still high and low, respectively, showing continuous behavior. As we can appreciate in Figure 16, radar plots also confirm our previous analysis on continuity. On the one hand, we can see in those plots how *ca-AstroPh* and *email-Enron* for DP$_{WCC}$ (yellow) has a bigger `dief@t` value. On the other hand, *web-Google* for DPFH$_{WCC}$ (green) `dief@t` is bigger. Nevertheless, both approaches are continuous and are consistent with Table 7 and Table 5.

In conclusion, we can say that regarding [RQ2], although DP$_{WCC}$ is more continuous in *ca-AstroPh* and *email-Enron* Graphs, DPFH$_{WCC}$ is much faster on those and more continuous in *web-Google* network, showing that it preserves the continuity approach and in some cases it is even faster.
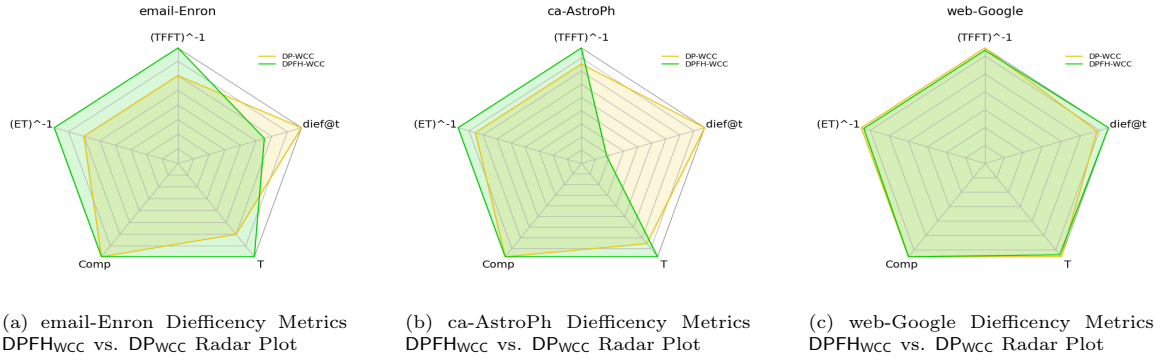


(a) email-Enron Diefficency Metrics DPFH$_{WCC}$ vs. DP$_{WCC}$ Radar Plot

(b) ca-AstroPh Diefficency Metrics DPFH$_{WCC}$ vs. DP$_{WCC}$ Radar Plot

(c) web-Google Diefficency Metrics DPFH$_{WCC}$ vs. DP$_{WCC}$ Radar Plot

Figure 16: Radial plots show how the different dimensions values provided by `diefpy` tool such as T, TFFT, `dief@t`, ET and Comp are related each other for each experimental case. These figures show radial plot observed results after running for each network comparing DP$_{WCC}$ (yellow) vs. DPFH$_{WCC}$ (green) implementations. `dief@t` is described in subsubsection 3.1.1.

*6.2.4. Experiment: E4*
For this type of analysis, our experiment focuses on *email-Enron* network only because profiling data generated by GHC is big enough to conduct the analysis and on the other, and enabling profiling penalize execution time. Moreover, it is important to remark that the analysis was conducted on DP$_{WCC}$ and DPFH$_{WCC}$ but since the results are similar in terms of plots and behavior, we show only the results obtained for DP$_{WCC}$ measurements.

*Multithreading.* For analyzing parallelization and multithreading we have used *ThreadScope* Tool 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.
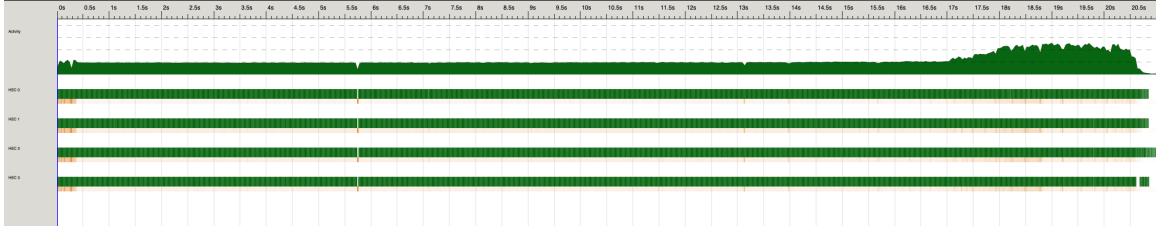
Figure 17: Threadscope Image of General Execution

In Figure 17, 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 17, 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], verifying that Haskell not only supports the required parallelization level but is evenly distributed across the program execution.

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. When the program start, and because of the nature of the DPP model, it is spawning the *Source* stage in



Figure 18: Threadscope Image of Zoomed Fraction

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 18 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* [17]. Having said that, it can be appreciated in Figure 18 our first assumption that the load is evenly distributed because the mean number of executing threads per core is 571.
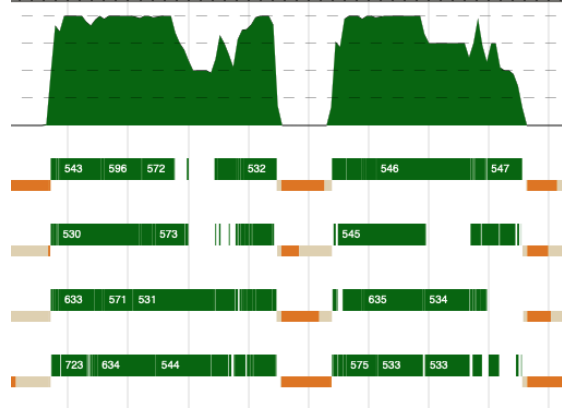
*Memory allocation.* Another important aspect in our case is how the memory is managed to avoid memory leaks or other non-desired behavior that increases memory allocation during the execution time. This is even more important in the particular in $DP_{WCC}$ because it requires to maintain the set of connected components in memory throughout the execution of the program or at least until it outputs the calculated WCC if the last *Filter* is reached, and this WCC will not grow anymore.
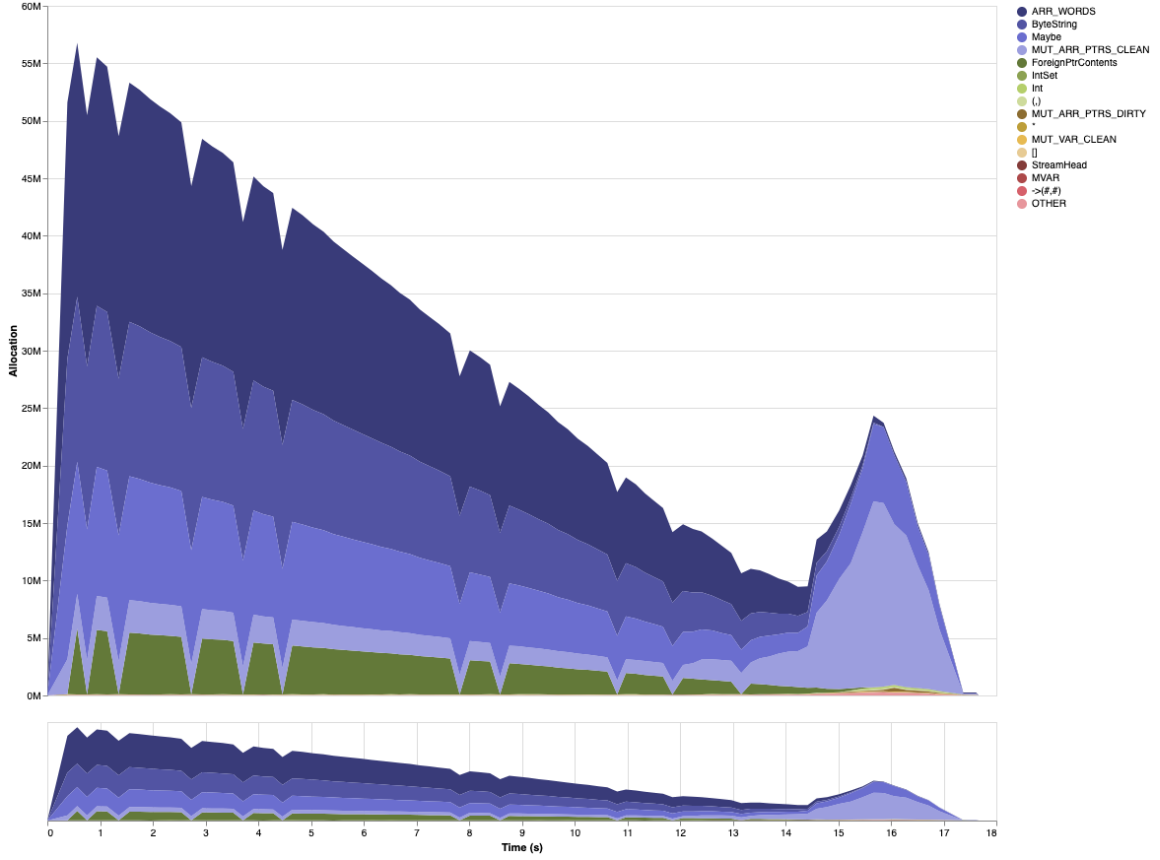
Figure 19: Memory Allocation

In order to verify this, we measure memory allocation with *eventlog2html* which converts generated profiling memory eventlog files into graphical HTML representation. As we can see in Figure 19, DP$_{\text{WCC}}$ does an efficient work on allocating memory, since we are not using more than 57 MB of memory during the whole execution of the program. On the other hand, if we analyze how the memory is allocated during the execution of the program, it can also be appreciated that most of the memory is allocated at the beginning of the program and steadily decrease over time with a small peak at the end that does not overpass even half of the initial peak of 57 MB. The explanation for this behavior is quite straightforward because at the beginning we are reading from the file and transforming a `ByteString` buffer to (`Int, Int`) edges. This is seen in the image in which the dark blue that is on top of the area is `ByteString` allocation. Light blue is allocation of `Maybe` `a` type which is the type that is returned by the *Channels* because it can contain a value or not. Data value `Nothing` is indicating end of the *Channel*.

Another important aspect is the green area which represents `IntSet` allocation, which in the case of our program is the data structure that we use to gather the set of vertices that represents a WCC. This means that the amount of memory used for gathering the WCC itself is minimum, and it is decreasing over time, which is another empirical indication that we are incrementally releasing results to the user. It can be seen as well that as long the green area reduces the lighter

35

blue (`MUT_ARR_PTRS_CLEAN`[27]) increases at the same time indicating that the computations for the output (releasing results) is taking place.

Finally, according to what we have stated above, we can answer the question [Q3] showing that not only memory management was efficient, but at the same time, the memory was not leaking or increasing across the running execution program.

## 7. Conclusions and Future Work

The empirical evaluation of the $DP_{WCC}$ implementation to compute weakly connected components of a graph, evidence suitability, and robustness to provide a Dynamic Pipeline Framework in Haskell. Measuring using `dief@t` and `dief@k` metrics provide evidence of the continuous behavior of $DPFH_{WCC}$ and $DP_{WCC}$ without adding overhead over the baseline implementation using the Haskell library $BLH_{WCC}$. Regarding the main aspects where DPP is strong, i.e., pipeline parallelism and time processing, the $DPFH_{WCC}$ and $DP_{WCC}$ performance show that Haskell can deal with the requirements for the WCC problem without penalizing neither execution time nor memory allocation. In particular, the $DPFH_{WCC}$ and $DP_{WCC}$ implementations outperform $BLH_{WCC}$ in those cases where the topology of the graph is more sparse, and the largest WCCs do not include a large percentage of vertices of the original graph. Moreover, the reported results indicate that implementing Dynamic Pipeline in Haskell Programming Language is feasible. More importantly, these results put in perspective the goodness of the proposed Dynamic Pipeline Framework to implement a wide range of algorithms supported by purely functional programming language.

In the future, we plan to apply DPF to other problems of subgraph enumeration, e.g., bi-triangle or graphlets, essential for efficient network analysis. We also plan to compare different language implementations of DP, taking into consideration the promising results reported in the literature [24] and the experimental outcomes presented in this article. Moreover, as a result of assessing the DPF potential, we plan to develop a new version according to [18] to reach a higher parallelization level using sparks (`Par` *Monad*). We envision that using this new approach will allow the DPF to scale. As a result, the DPF will scale up to practical problems with huge input data streams.

## 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] M. M. Chakravarty, G. Keller, and S. P. Jones. Associated type synonyms. *ACM SIGPLAN Notices*, 40(9):241–253, 2005.

[3] O. Danvy and L. R. Nielsen. Defunctionalization at work. In *Proceedings of the 3rd ACM SIGPLAN international conference on Principles and practice of declarative programming*, pages 162–174, 2001.

[4] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.

---

[27]GHC-Exts-Heap-ClosureTypes.html

[5] T. Dunning and E. Friedman. *Streaming architecture: new designs using Apache Kafka and MapR streams.* " O'Reilly Media, Inc.", 2016.

[6] M. Fowler. *Domain-specific languages.* Pearson Education, 2010.

[7] 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.

[8] T. Harris, S. Marlow, S. Peyton-Jones, and M. Herlihy. Composable memory transactions. In *Proceedings of the tenth ACM SIGPLAN symposium on Principles and practice of parallel programming*, pages 48–60, 2005.

[9] R. Hinze, J. Jeuring, and A. Löh. Type-indexed data types. *Science of Computer Programming*, 51(1-2):117–151, 2004.

[10] 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.

[11] P. Hudak. Building domain-specific embedded languages. *ACM Computing Surveys (CSUR)*, 28(4es):196–es, 1996.

[12] O. Kiselyov, R. Lämmel, and K. Schupke. Strongly typed heterogeneous collections. In *Proceedings of the 2004 ACM SIGPLAN Workshop on Haskell*, pages 96–107, 2004.

[13] O. Kiselyov, S. P. Jones, and C.-c. Shan. Fun with type functions. In *Reflections on the Work of CAR Hoare*, pages 301–331. Springer, 2010.

[14] T. Kosar, S. Bohra, and M. Mernik. Domain-specific languages: A systematic mapping study. *Information and Software Technology*, 71:77–91, 2016.

[15] I.-T. A. Lee, C. E. Leiserson, T. B. Schardl, Z. Zhang, and J. Sukha. On-the-fly pipeline parallelism. *ACM Transactions on Parallel Computing (TOPC)*, 2(3):1–42, 2015.

[16] J. Leskovec and A. Krevl. SNAP Datasets: Stanford large network dataset collection. `http://snap.stanford.edu/data`, June 2014.

[17] 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*, pages 107–118, 2007.

[18] S. Marlow. Parallel and concurrent programming in Haskell. In V. Zsók, Z. Horváth, and R. Plasmeijer, editors, *CEFP 2011*, volume 7241 of *LNCS*, pages 339–401. O'Reilly Media, Inc., 2012.

[19] S. Marlow, S. Peyton Jones, and S. Singh. Runtime support for multicore haskell. In *Proceedings of the 14th ACM SIGPLAN international conference on Functional programming*, pages 65–78, 2009.

[20] S. Marlow, P. Maier, H.-W. Loidl, M. K. Aswad, and P. Trinder. Seq no more: better strategies for parallel haskell. *ACM Sigplan Notices*, 45(11):91–102, 2010.

[21] S. Marlow, R. Newton, and S. Peyton Jones. A monad for deterministic parallelism. *ACM SIGPLAN Notices*, 46(12):71–82, 2011.

[22] 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.

[23] T. T. Nguyen, M. Weidlich, H. Yin, B. Zheng, Q. H. Nguyen, and Q. V. H. Nguyen. Factcatch: Incremental pay-as-you-go fact checking with minimal user effort. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2165–2168, 2020.

[24] 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.

[25] S. Peyton-Jones and E. Meijer. Henk: a typed intermediate language. *TIC*, 97:10, 1997.

[26] H. Röger and R. Mayer. A comprehensive survey on parallelization and elasticity in stream processing. *ACM Computing Surveys (CSUR)*, 52(2):1–37, 2019.

[27] 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`.

[28] G. Van Dongen and D. Van den Poel. Evaluation of stream processing frameworks. *IEEE Transactions on Parallel and Distributed Systems*, 31(8):1845–1858, 2020.

[29] B. A. Yorgey, S. Weirich, J. Cretin, S. Peyton Jones, D. Vytiniotis, and J. P. Magalhães. Giving haskell a promotion. In *Proceedings of the 8th ACM SIGPLAN Workshop on Types in Language Design and Implementation*, pages 53–66, 2012.

[30] 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`.

# Appendix A. Summary of URL references

| |
|---|
| `bytestring 0.10.12.0` |
| `https://hackage.haskell.org/package/bytestring` |
| `criterion` library |
| `https://hackage.haskell.org/package/criterion` |
| `Data.Graph`-containers `0.6.2.1` |
| `https://hackage.haskell.org/package/containers` |
| Diefficency Metrics (Dm) Tool `diefpy` |
| `https://github.com/SDM-TIB/diefpy/` |
| `dynamic-pipeline 0.3.2.0` |
| `https://hackage.haskell.org/package/dynamic-pipeline` |
| *eventlog2html* Tool |
| `https://mpickering.github.io/eventlog2html/` |
| `MUT_ARR_PTRS_CLEAN` |
| GHC-Exts-Heap-ClosureTypes.html |
| `NoImplicitPrelude` |
| `https://downloads.haskell.org/ghc/8.8.4/docs/html/users_guide/` |
| stack version 2.5.1 |
| `https://docs.haskellstack.org/en/stable/README/` |
| `relude 1.0.0.1` |
| `https://hackage.haskell.org/package/relude` |
| `stack build` |
| `https://github.com/jproyo/upc-miri-tfm/tree/main/connected-comp` |
| *ThreadScope* Tool |
| `https://wiki.haskell.org/ThreadScope` |
| `unagi-chan 0.4.1.3` |
| `https://github.com/jberryman/unagi-chan` |

Table A.8: Summary of libraries and tools used in this work