

Compositionality in Synchronous Data Flow: Modular Code Generation from Hierarchical SDF Graphs*

Stavros Tripakis, Dai Bui, Marc Geilen, Bert Rodiers, and Edward A. Lee

May 7, 2010

Abstract

Hierarchical SDF models are not compositional: a composite SDF actor cannot be represented as an atomic SDF actor without loss of information that can lead to rate inconsistency or deadlock. Motivated by the need for incremental and modular code generation from hierarchical SDF models, we introduce in this paper DSSF profiles. DSSF (Deterministic SDF with Shared FIFOs) forms a compositional abstraction of composite actors that can be used for modular compilation. We provide algorithms for automatic synthesis of non-monolithic DSSF profiles of composite actors given DSSF profiles of their sub-actors. We show how different tradeoffs can be explored when synthesizing such profiles, in terms of modularity (keeping the size of the generated DSSF profile small) versus reusability (maintaining necessary information to preserve rate consistency and deadlock-absence) as well as algorithmic complexity. We show that our method guarantees maximal reusability and report on a prototype implementation.

1 Introduction

Programming languages have been constantly evolving over the years, from assembly, to structural programming, to object-oriented programming, etc. Common to this evolution is the fact that new programming models provide mechanisms and notions that are more *abstract*, that is, remote from the actual implementation, but better suited to the programmer’s intuition. Raising the level of abstraction results in undeniable benefits in productivity. But it is more than just building systems faster or cheaper. It also allows to create systems that could not have been conceived otherwise, simply because of too high complexity.

Modeling languages with built-in concepts of concurrency, time, I/O interaction, and so on, are particularly suitable in the domain of embedded systems. Indeed, languages such as Simulink, UML or SystemC, and corresponding tools, are particularly popular in this domain, for various applications. The tools provide mostly modeling and simulation, but often also code generation and static analysis or verification capabilities, which are increasingly important in an industrial setting. We believe that this tendency will continue, to the point where modeling languages of today will become the programming languages of tomorrow, at least in the embedded software domain.

A widespread model of computation in this domain is Synchronous (or Static) Data Flow (SDF) [14]. SDF is particularly well-suited for signal processing and multimedia applications and has been extensively studied over the years (e.g., see [2, 19]). Recently, languages based on SDF, such as StreamIt [21], have also been applied to multicore programming.

In this paper we consider *hierarchical* SDF models, where an SDF graph can be *encapsulated* into a *composite* SDF actor. The latter can then be connected with other SDF actors, further encapsulated, and

*This work was supported in part by the Center for Hybrid and Embedded Software Systems (CHESS) at UC Berkeley, which receives support from the National Science Foundation (NSF awards #CCR-0225610 (ITR), #0720882 (CSR-EHS: PRET) and #0931843 (ActionWebs)), the U. S. Army Research Office (ARO #W911NF-07-2-0019), the U. S. Air Force Office of Scientific Research (MURI #FA9550-06-0312 and AF-TRUST #FA9550-06-1-0244), the Air Force Research Lab (AFRL), the Multiscale Systems Center (MuSyC) and the following companies: Bosch, National Instruments, Thales, and Toyota.

so on, to form a hierarchy of SDF actors of arbitrary depth. This is essential for *compositional modeling*, which allows to design systems in a modular, scalable way, enhancing readability and allowing to master complexity in order to build larger designs. Hierarchical SDF models are part of a number of modeling environments, including the Ptolemy II framework [8].

The problem we solve in this paper is *modular code generation* for hierarchical SDF models. Modular means that code is generated for a given composite SDF actor P *independently from context*, that is, independently from which graphs P is going to be used in. Moreover, once code is generated for P , then P can be seen as an *atomic* (non-composite) actor, that is, a “black box” without access to its internal structure. Modular code generation can be paralleled to separate compilation, which is available in most standard programming languages: the fact that one does not need to compile an entire program in one shot, but can compile files, classes, or other units, separately, and then combine them (e.g., by *linking*) to a single executable. This is obviously a key capability for a number of reasons, ranging from incremental compilation (compiling only the parts of a large program that have changed), to dealing with IP (intellectual property) concerns (having access to object code only and not to source code). We want to do the same for SDF models. Moreover, in the context of a system like Ptolemy II, in addition to the benefits mentioned above, modular code generation is also useful for speeding-up simulation: replacing entire sub-trees of a large hierarchical model by a single actor for which code has been automatically generated and pre-compiled, removes the overhead of executing all actors in the sub-tree individually.

Our work extends the ideas of modular code generation for *synchronous block diagrams* (SBDs), introduced in [16, 15]. In particular, we borrow their notion of *profile* which characterizes a given actor. Modular code generation then essentially becomes a *profile synthesis* problem: how to synthesize a profile for composite actors, based on the profiles of its internal actors.

In SBDs, profiles are essentially DAGs (directed acyclic graphs) that capture the dependencies between inputs and outputs of a block, at the same synchronous round. In general, not all outputs depend on all inputs, which allows feedback loops with unambiguous semantics to be built. For instance, in a *unit delay* block the output does not depend on the input at the same clock cycle, therefore this block “breaks” dependency cycles when used in feedback loops.

The question is, what is the right model for profiles of SDF graphs. We answer this question in this paper. For SDF graphs, profiles turn out to be more interesting than simple DAGs. SDF profiles are essentially SDF graphs themselves, but with the ability to associate multiple producers and/or consumers with a single FIFO queue. Sharing FIFOs among different actors generally results in non-deterministic models, however, in our case, we can guarantee that actors that share queues are always fired in a deterministic order. We call this model *deterministic SDF with shared FIFOs* (DSSF). DSSF allows, in particular, to decompose the firing of a composite actor into an arbitrary number of *firing functions* that may consume tokens from the same input port or produce tokens to the same output port (an example is shown in Figure 3). Having multiple firing functions allows to decouple firings of different internal actors of the composite actor, so that deadlocks are avoided when the composite actor is embedded in a given context. Our method guarantees *maximal reusability* [16], i.e., the absence of deadlock in any context where the corresponding “flat” (non-hierarchical) SDF graph is deadlock-free, as well as consistency in any context where the flat SDF graph is consistent.

We show how to perform profile synthesis for SDF graphs automatically. This means synthesize for a given composite actor a profile, in the form of a DSSF graph, given the profiles of its internal actors (also DSSF graphs). This process involves multiple steps, among which are the standard *rate analysis* and *deadlock detection* procedures used to check whether a given SDF graph can be executed infinitely often without deadlock and with bounded queues [14]. In addition to these steps, SDF profile synthesis involves *unfolding* a DSSF graph (i.e., replicating actors in the graph according to their relative rates produced by rate analysis) to produce a DAG that captures the dependencies between the different consumptions and productions of tokens *at the same port*.

Reducing the DSSF graph to a DAG is interesting because it allows to apply for our purposes the idea of *DAG clustering* proposed originally for SBDs [16, 15]. As in the SBD case, we use DAG clustering in order to group together firing functions of internal actors and synthesize a small (hopefully minimal) number of

firing functions for the composite actor. These determine precisely the profile of the latter. Keeping the number of firing functions small is essential, because it results in further compositions of the actor being more efficient, thus allowing the process to scale to arbitrary levels of hierarchy.

As shown in [16, 15], there exist different ways to perform DAG clustering, that achieve different tradeoffs, in particular in terms of number of clusters produced vs. reusability of the generated profile. Among the clustering methods proposed for SBDs, of particular interest to us are those that produce *disjoint* clusterings, where clusters do not share nodes. Unfortunately, *optimal disjoint clustering*, that guarantees maximal reusability with a minimal number of clusters, is NP-complete [15]. This motivates us to devise a new clustering algorithm, called *greedy backward disjoint clustering* (GBDC). GBDC guarantees maximal reusability but due to its greedy nature cannot guarantee optimality in terms of number of clusters. On the other hand, GBDC has polynomial complexity.

The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 reviews (hierarchical) SDF graphs. Section 4 reviews rate analysis and deadlock detection for flat SDF graphs. Section 5 reviews modular code generation for SBDs which we build upon. Section 6 introduces DSSF as profiles for SDF graphs. Section 7 elaborates on the profile synthesis procedure. Section 8 details DAG clustering, in particular, the GBDC algorithm. Section 9 presents a prototype implementation. Section 10 presents the conclusions and discusses future work.

2 Related Work

Dataflow models of computation have been extensively studied in the literature. Dataflow models with deterministic actors, such as Kahn Process Networks [12] and their various subclasses, including SDF, are compositional at the semantic level. Indeed, actors can be given semantics as continuous functions on streams, and such functions are closed by composition. (Interestingly, it is much harder to derive a compositional theory of *non-deterministic* dataflow, e.g., see [5, 11, 20].) Our work is at a different, non-semantic level, since we mainly focus on finite representations of the behavior of networks at their interfaces, in particular of the dependencies between inputs and output. We also take a “black-box” view of atomic actors, assuming their internal semantics (e.g., which function they compute) are unknown and unimportant for our purpose of code generation. Finally, we only deal with the particular subclass of SDF models.

Despite extensive work on code generation from SDF models and especially scheduling (e.g., see [2, 19]), there is little existing work that addresses compositional representations and modular code generation for such models. [10] proposes abstraction methods that reduce the size of SDF graphs, thus facilitating throughput and latency analysis. His goal is to have a conservative abstraction in terms of these performance metrics, whereas our goal here is to preserve input-output dependencies to avoid deadlocks during further composition.

Non-compositionality of SDF due to potential deadlocks has been observed in earlier works such as [18], where the objective is to schedule SDF graphs on multiple processors. This is done by partitioning the SDF graph into multiple sub-graphs, each of which is scheduled on a single processor. This partitioning (also called *clustering*, but different from DAG clustering that we use in this paper, see below) may result in deadlocks, and the so-called “SDF composition theorem” given in [18] provides a sufficient condition so that no deadlock is introduced.

More recently, [9] also identify the problem of non-compositionality and propose *Cluster Finite State Machines* (CFSMs) as a representation of composite SDF. They show how to compute a CFSM for a composite SDF actor that contains standard, atomic, SDF sub-actors, however, they do not show how a CFSM can be computed when the sub-actors are themselves represented as CFSMs. This indicates that this approach may not generalize to more than one level of hierarchy. Our approach works for arbitrary depths of hierarchy.

Another difference between the above work and ours is on the representation models, namely, CFSM vs. DSSF. CFSM is a state-machine model, where transitions are annotated with guards checking whether a sufficient number of tokens is available in certain input queues. DSSF, on the other hand, is a data flow model, only slightly more general than SDF. This allows to re-use many of the techniques developed for standard SDF graphs, for instance, rate analysis and deadlock detection, with minimal adaptation.

The same remark applies to other automata-based formalisms, such as I/O automata [17], interface automata [6], and so on. Such formalisms could perhaps be used to represent consumption and production actions of SDF graphs, resulting in compositional representations. These would be at a much lower level than DSSF, however, and for this reason would not admit SDF techniques such as rate analysis, which are more “symbolic”.

To the extent that we propose DSSF profiles as interfaces for composite SDF graphs, our work is related to so-called *component-based design* and *interface theories* [7]. Similarly to that line of research, we propose methods to synthesize interfaces for compositions of components, given interfaces for these components. We do not, however, include notions of refinement in our work. We are also not concerned with how to specify the “glue code” between components, as is done in *connector algebras* [1, 4]. Indeed, in our case, there is only one type of connections, namely, conceptually unbounded FIFOs, defined by the SDF semantics. Moreover, connections of components are themselves specified in the SDF graphs of composite actors, and are given as an input to the profile synthesis algorithm. Finally, we are not concerned with issues of timeliness or distribution, as in [13].

Finally, we should emphasize that our DAG clustering algorithms solve a different problem than the clustering methods used in [18, 9] and other works in the SDF scheduling literature. Our clustering algorithms operate on plain DAGs, as do the clustering algorithms originally developed for SBDs [16, 15]. On the other hand, clustering in [9, 18] is done directly at the SDF level, by grouping SDF actors and replacing them by a single SDF actor (e.g., see Figure 4 in [9]). This, in our terms, corresponds to monolithic clustering, which is not compositional.

3 Hierarchical SDF Graphs

A *synchronous (or static) dataflow (SDF) graph* [14] consists of a set of nodes, called *actors*¹ connected through a set of directed edges. Each actor has a set of *input ports* (possibly zero) and a set of *output ports* (possibly zero). An edge connects an output port y of an actor A to an input port x of an actor B (B can be the same as A). An output port can be connected to a single input port, and vice versa.² Such an edge represents a FIFO (first-in, first-out) queue, that stores *tokens* that the source actor A produces when it fires. The tokens are removed and consumed by the destination actor B when it fires. Queues are of unbounded size in principle. In practice, however, we are interested in SDF graphs that can execute forever using bounded queues.

Actors are either *atomic* or *composite*. A composite actor encapsulates an SDF graph as shown in Figure 1. P is a composite actor while A and B are atomic actors. Composite actors can themselves be encapsulated in new composite actors, thus forming a hierarchical model of arbitrary depth.

Each port of an atomic actor has an associated *token rate*, a positive integer number, which specifies how many *tokens* (i.e., data values) are consumed from or produced to each port every time the actor fires. In the example of Figure 1, A consumes one token from its single input port and produces two tokens to its single output port. B consumes three tokens and produces one token. Each port of an atomic also has a given *data type* (integer, boolean, ...) and connections can only be done among ports with compatible data types, as in a standard typed programming language. Composite actors do not have token rate or data type annotations on their ports. They inherit this information from their internal actors, as we will explain in this paper.

SDF graphs can be *open* or *closed*. A graph is closed if all its input ports are connected; otherwise it is open. The graph shown in Figure 1 is open because the input port of P is not connected. The graphs shown in Figure 2 are closed. These graphs also illustrate another element of SDF notation, namely, *initial tokens*:

¹ It is useful to distinguish between actor *types* and actor *instances*. Indeed, an actor can be used in a given graph multiple times. For example, an actor of type *Adder*, that computes the arithmetic sum of its inputs, can be used multiple times in a given graph. In this case, we say that the *Adder* is *instantiated* multiple times. Each “copy” is an actor *instance*. In the rest of the paper, we often omit to distinguish between type and instance when we refer to an actor, when the meaning is clear from context.

² Implicit *fan-in* or *fan-out* is not allowed, however, it can be implemented explicitly, using actors. For example, an actor that consumes an input token and replicates to each of its output ports models fan-out.

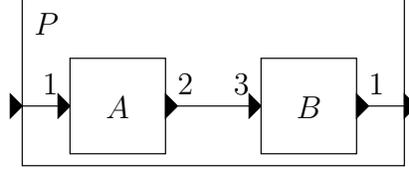


Figure 1: Example of a hierarchical SDF graph.

the edge from the output port of Q to the input port of P is annotated with two black dots which means there are initially two tokens in the queue from Q to P . Likewise, there is one initial token in the queue from P to Q .

An SDF graph is *flat* if it contains only atomic actors. A *flattening* process can be applied to turn a hierarchical SDF graph into a flat graph, by removing composite actors and replacing them with their internal graph, while making sure to re-institute any connections that would be otherwise lost. An example of flattening is shown in Figure 2.

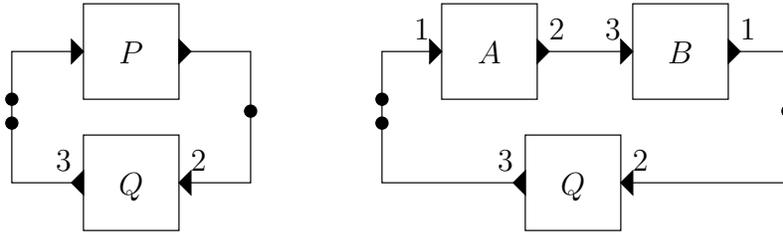


Figure 2: Left: using the composite actor P of Figure 1 in an SDF graph with feedback and initial tokens. Right: the same graph after flattening P .

4 Analysis of SDF Graphs

The SDF analysis methods proposed in [14] allow to check whether a given SDF graph has a *periodic admissible sequential schedule* (PASS). Existence of a PASS guarantees two things: first, that the actors in the graph can fire infinitely often without *deadlock*; and second, that only *bounded queues* are required to store intermediate tokens produced during the firing of these actors. We review these analysis methods here, because we are going to use them for modular code generation (Section 7).

4.1 Rate Analysis

Rate analysis seeks to determine if the token rates in a given SDF graph are *consistent*: if this is not the case, then the graph cannot be executed infinitely often with bounded queues. We illustrate the analysis in the simple example of Figure 1. The reader is referred to [14] for the details.

We wish to analyze the internal graph of P , consisting of actors A and B . This is an open graph, and we can ignore the unconnected ports for the rate analysis. Suppose A is fired r_A times for every r_B times that B is fired. Then, in order for the queue between A and B to remain bounded in repeated execution, it has to be the case that:

$$r_A \cdot 2 = r_B \cdot 3$$

that is, the total number of tokens produced by A equals the total number of tokens consumed by B . The above *balance equation* has a non-trivial (i.e., non-zero) solution: $r_A = 3$ and $r_B = 2$. This means that this

SDF graph is indeed consistent. In general, for larger and more complex graphs, the same analysis can be performed, which results in solving a system of multiple balance equations. If the system has a non-trivial solution then the graph is consistent, otherwise it is not. At the end of rate analysis, if consistent, a *repetition vector* (r_1, \dots, r_n) is produced that specifies the number r_i of times that every actor A_i in the graph fires with respect to other actors. This vector is used in the subsequent step of deadlock analysis.

4.2 Deadlock Analysis

Having a consistent graph is a necessary, but not sufficient condition for infinite execution: the graph might still contain *deadlocks* that arise because of absence of enough initial tokens. Deadlock analysis ensures that this is not the case. An SDF graph is deadlock free if and only if every actor A can fire r_A times, where r_A is the repetition value for A in the repetition vector (i.e., it has a PASS [14]). The method works as follows. For every directed edge e_i in the SDF graph, an integer counter b_i is maintained, representing the number of tokens in the FIFO queue associated with e_i . Counter b_i is initialized to the number of initial tokens present in e_i (zero if no such tokens are present). For every actor A in the SDF graph, an integer counter c_A is maintained, representing the number of remaining times that A should fire to complete the PASS. Counter c_A is initialized to r_A . A tuple consisting of all above counters is called a *configuration* v . A transition from a configuration v to a new configuration v' happens by firing an actor A , provided A is *enabled* at v , i.e., all its input queues have enough tokens, and provided that $c_A > 0$. Then, the queue counters are updated, and counter c_A is decremented by 1. If a configuration is reached where all actor counters are 0, there is no deadlock, otherwise, there is one. Notice that a single path needs to be explored, so this is not a costly method (i.e., not a full-blown reachability analysis). In fact, at most $\sum_{i=1}^n r_i$ steps are required to complete deadlock analysis, where (r_1, \dots, r_n) is the solution to the balance equations.

We illustrate deadlock analysis with an example. Consider the SDF graph shown at the left of Figure 2 and suppose P is an atomic actor, with input/output token rates 3 and 2, respectively. Rate analysis then gives $r_P = r_Q = 1$. Let the edges from P to Q and from Q to P be denoted e_1 and e_2 , respectively. Deadlock analysis then starts with configuration $v_0 = (c_P = 1, c_Q = 1, b_1 = 1, b_2 = 2)$. P is not enabled at v_0 because it needs 3 input tokens but $b_2 = 2$. Q is not enabled at v_0 either because it needs 2 input tokens but $b_1 = 1$. Thus v_0 is a deadlock. Now, suppose instead of 2 initial tokens, edge e_2 had 3 initial tokens. Then, we would have as initial configuration $v_1 = (c_P = 1, c_Q = 1, b_1 = 1, b_2 = 3)$. In this case, deadlock analysis can proceed: $v_1 \xrightarrow{P} (c_P = 0, c_Q = 1, b_1 = 3, b_2 = 0) \xrightarrow{Q} (c_P = 0, c_Q = 0, b_1 = 1, b_2 = 3)$. Since a configuration is reached where $c_P = c_Q = 0$, there is no deadlock.

4.3 Transformation of SDF to Homogeneous SDF

A *homogeneous* SDF (HSDF) graph is an SDF graph where all token rates are equal (and without loss in generality, can be assumed to be equal to 1). Any consistent SDF graph can be transformed to an equivalent HSDF graph using an *unfolding* process [14, 18, 19]. This process consists in replicating each actor in the SDF as many times as specified in the repetition vector. This subsequently allows to identify explicitly the input/output dependencies of different productions and consumptions at the same output or input port. Examples of this process are presented in Section 7.4, where we generalize the process to DSSF graphs.

5 Modular Code Generation Framework

As mentioned in the introduction, our modular code generation framework for SDF builds upon the work of [16, 15]. A fundamental element of the framework is the notion of *profiles*. Every SDF actor has an associated profile. The profile can be seen as an *interface*, or *summary*, that captures the essential information about the actor. Atomic actors have predefined profiles. Profiles of composite actors are synthesized automatically, as shown in Section 7.

A profile contains, among other things, a set of *firing functions*, that, together, implement firing of an actor. In the simple case, an actor may have a single firing function. For example, actors A, B of Figure 1

may each have a single firing function

```
A.fire(input x[1]) output (y[2]);
B.fire(input x[3]) output (y[1]);
```

The above signatures specify that `A.fire` takes as input 1 token at input port x and produces as output 2 tokens at output port y , and similarly for `B.fire`. In general, however, an actor may have more than one firing function in its profile. This is necessary in order to avoid *monolithic* code, and instead produce code that achieves *maximal reusability*, as is explained in Section 7.3.

The *implementation* of a profile contains, among other things, the implementation of each of the firing functions listed in the profile as a sequential program in a language such as C++ or Java. We will show how to automatically generate such implementations of SDF profiles in Section 7.7.

Modular code generation is then the following process:

- given a composite actor P , its internal graph, and profiles for every internal actor of P ,
- synthesize automatically a profile for P and an implementation of this profile.

Note that a given actor may have multiple profiles, each achieving different tradeoffs, for instance, in terms of *modularity* (size of the profile) and *reusability* (ability to use the profile in as many contexts as possible). We illustrate such tradeoffs in the sequel.

6 DSSF Graphs and SDF Profiles

Deterministic SDF with shared FIFOs, or DSSF, is an extension of SDF in the sense that, whereas *shared FIFOs* are explicitly prohibited in SDF graphs, they are allowed in DSSF graphs, provided determinism is ensured. To see why sharing FIFOs generally results in non-deterministic models, consider two producers A_1, A_2 sharing the same output queue, and a consumer B reading from that queue and producing an external output. Depending on the order of execution of A_1 and A_2 , their outputs will be stored in the shared queue in a different order. Therefore, the output of B will also generally differ.

To guarantee determinism, it suffices to ensure that A_1 and A_2 are always executed in a fixed order. This is the condition we impose on DSSF graphs, namely, that if Q is a queue shared among a set of producers A_1, \dots, A_a and a set of consumers B_1, \dots, B_b , then the graph ensures, by means of other edges, a deterministic way of firing A_1, \dots, A_a , as well as a deterministic way of firing B_1, \dots, B_b . In the context of this paper, we will meet this condition by ensuring that the subgraph restricted to A_1, \dots, A_a is a homogeneous SDF graph that moreover implies a total order among A_1, \dots, A_a , and similarly for B_1, \dots, B_b . Examples of DSSF graphs are provided below (e.g., in Figure 5).

We will use a special type of DSSF graphs to represent profiles of SDF actors, called *SDF profiles*. An SDF profile is a flat DSSF graph where shared queues are only allowed at input or output ports. Moreover, all edges between actors of the profile are such that the number of tokens produced and consumed at each firing by the source and destination actors are equal: this implies that connected actors fire with equal rates. Because shared FIFOs are only allowed in profiles at input or output ports, these FIFOs are called *external*. This is to distinguish them from *internal* shared FIFOs that may arise in other types of DSSF graphs that we use in this paper, in particular, in so-called *internal* DSSF graphs (see Section 7.1).

Figure 3 shows two examples of SDF profiles, namely, two possible profiles for the composite actor P of Figure 1. The left-most profile is a standard SDF actor, whereas the right-most profile is a DSSF graph with two shared FIFOs, depicted as small squares. Dashed-line edges are called *dependency edges* and are distinguished from solid-line edges that are “true” *dataflow edges*. The distinction is made only for reasons of clarity, in order to understand better the way edges are created during profile generation (Section 7.6). Otherwise the distinction plays no role, and dependency edges can be encoded as standard SDF edges with token production and consumption both equal to 1.

We will see how the two profiles of Figure 3 can be synthesized automatically in Section 7. We will also see that these two profiles have different properties. In particular, they represent different pareto points in the modularity vs. reusability tradeoff (Section 7.3).

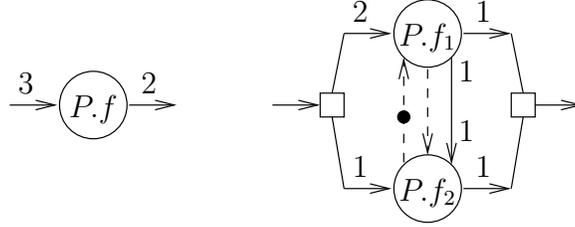


Figure 3: Two profiles for the composite actor P of Figure 1.

The actors of an SDF profile represent firing functions. They are drawn as circles instead of squares to distinguish SDF profiles from SDF models. The left-most profile of Figure 3 contains a single actor, denoted $P.f$, which corresponds to a single firing function $P.\mathbf{fire}$. Profiles that contain a single firing function are called *monolithic*. The right-most profile of Figure 3 contains two actors, $P.f_1$ and $P.f_2$, corresponding to two firing functions, $P.\mathbf{fire1}$ and $P.\mathbf{fire2}$: this is a *non-monolithic* profile. This profile also has two shared FIFOs, one connecting the input port to the two actors, and another from the two actors to the output port. It also has multiple edges, including dependency and dataflow edges. The forward dependency edge from $P.f_1$ to $P.f_2$ represents the fact that $P.\mathbf{fire1}$ must be called *before* $P.\mathbf{fire2}$. Notice that this edge is redundant, since the dataflow edge from $P.f_1$ already encodes this dependency: $P.f_2$ cannot fire before $P.f_1$ fires and produces a token. But the dataflow edge, in addition to a dependency, also encodes a transfer of data between the two functions. The dependency edge from $P.f_2$ to $P.f_1$ encodes the fact that $P.\mathbf{fire1}$ can be called for a second time only *after* the first call to $P.\mathbf{fire2}$. Together these edges impose a total order on the firing of these two functions, which results in deterministic handling of tokens in the shared input and output FIFOs.

Unless explicitly mentioned otherwise, in the examples that follow we assume that atomic blocks have monolithic profiles.

7 Profile Synthesis and Code Generation

As mentioned above, modular code generation takes as input a composite actor P , its internal graph, and profiles for every internal actor of P , and produces as output a profile for P and an implementation of this profile. Profile synthesis refers to the production of a profile for P , while code generation refers to the automatic generation of an implementation of this profile. These two functions are performed a number of steps, detailed below.

7.1 Connecting the SDF Profiles

The first step consists in connecting the SDF profiles of internal actors of P . This is done simply as dictated by the connections found in the internal graph of P . The result is a flat DSSF graph, called the *internal DSSF graph* of P . We illustrate this through an example. Consider the composite actor P shown in Figure 1. Suppose both its internal actors A and B have monolithic profiles, with $A.f$ and $B.f$ representing $A.\mathbf{fire}$ and $B.\mathbf{fire}$, respectively. Then, by connecting these monolithic profiles we obtain the internal DSSF graph shown in Figure 4. In this case, the internal DSSF graph is a standard SDF graph.

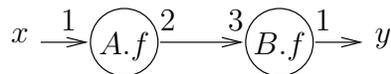


Figure 4: Internal DSSF graph of composite actor P of Figure 1.

Two more examples of connection are shown in Figure 5. There, we connect the profiles of internal actors P and Q of the (closed) graph shown at the left of Figure 2. Actor Q is assumed to have a monolithic profile. Actor P has two possible profiles, shown in Figure 3. The two resulting internal DSSF graphs are shown in Figure 5. The left-most one is a standard SDF graph. The right-most one is a DSSF graph, with two internal shared FIFOs and three standard SDF edges (i.e., non-shared FIFOs).

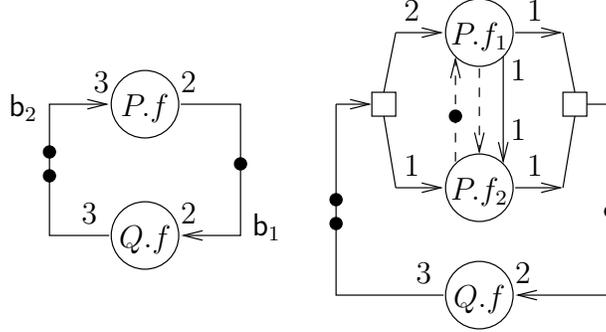


Figure 5: Two internal DSSF graphs, resulting from connecting the two profiles of actor P shown in Figure 3 and a monolithic profile of actor Q , according to the graph at the left of Figure 2.

7.2 Rate Analysis with SDF Profiles

This step is similar to the rate analysis process described in Section 4.1, except that it is performed on the internal DSSF graph produced by the connection step, instead of an SDF graph. This presents no major challenges, however, and the method is essentially the same as the one proposed in [14].

Let us illustrate the process here, for the DSSF graph shown to the right of Figure 5. We associate repetition variables r_p^1, r_p^2 , and r_q , respectively, to $P.f_1, P.f_2$ and $Q.f$. Then, we have the following balance equations:

$$\begin{aligned} r_p^1 \cdot 1 + r_p^2 \cdot 1 &= r_q \cdot 2 \\ r_q \cdot 3 &= r_p^1 \cdot 2 + r_p^2 \cdot 1 \\ r_p^1 \cdot 1 &= r_p^2 \cdot 1 \end{aligned}$$

As this has a non-trivial solution (e.g., $r_p^1 = r_p^2 = r_q = 1$), this graph is consistent, i.e., rate analysis succeeds in this example.

If the rate analysis step fails the graph is rejected. Otherwise, we proceed with the deadlock analysis step.

It is worth noting that rate analysis can sometimes succeed with non-monolithic profiles, whereas it would fail with a monolithic profile. An example is given in Figure 6. A composite actor R is shown to the left of the figure and its non-monolithic profile to the right. If we use R in the diagram shown to the middle of the figure, then rate analysis with the non-monolithic profile succeeds. It would fail, however, with the monolithic profile. This observation also explains why rate analysis must generally be performed on the internal DSSF graph, and not on the internal SDF graph using monolithic profiles for internal actors.

7.3 Deadlock Analysis with SDF Profiles

Success of the rate analysis step is a necessary, but not sufficient, condition in order for a graph to have a PASS. Deadlock analysis is used to ensure that this is the case. Deadlock analysis is performed on the

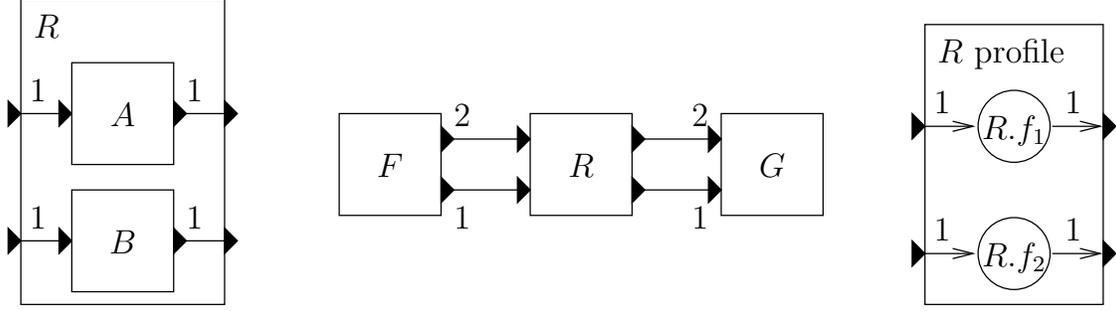


Figure 6: Composite SDF actor R (left); using R (middle); non-monolithic profile of R (right).

internal DSSF graph produced by the connection step. It is done in the same way as the deadlock detection process described in Section 4.2. We illustrate this on the two examples of Figure 5.

Consider first the DSSF graph to the left of Figure 5. There are two queues in this graph: a queue from $P.f$ to $Q.f$, and a queue from $Q.f$ to $P.f$. Denote the former by b_1 and the latter by b_2 . Initially, b_1 has 1 token, whereas b_2 has 2 tokens. $P.f$ needs 3 tokens to fire but only 2 are available in b_2 , thus $P.f$ cannot fire. $Q.f$ needs 2 tokens but only 1 is available in b_1 , thus $Q.f$ cannot fire either. Therefore there is a deadlock already at the initial state, and this graph is rejected.

Now consider the DSSF graph to the right of Figure 5. There are five queues in this graph: a queue from $P.f_1$ and $P.f_2$ to $Q.f$, a queue from $Q.f$ to $P.f_1$ and $P.f_2$, two queues from $P.f_1$ to $P.f_2$, and a queue from $P.f_2$ to $P.f_1$. Denote these queues by b_1, b_2, b_3, b_4, b_5 , respectively. Initially, b_1 has 1 token, b_2 has 2 tokens, b_3 and b_4 are empty and b_5 has 1 token. $P.f_1$ needs 2 tokens to fire and 2 tokens are indeed available in b_2 , thus $P.f_1$ can fire and the initial state is not a deadlock. Deadlock analysis gives:

$$\begin{aligned}
 (c_{p_1} = 1, c_{p_2} = 1, c_q = 1, b_1 = 1, b_2 = 2, b_3 = 0, b_4 = 0, b_5 = 1) & \xrightarrow{P.f_1} \\
 (c_{p_1} = 0, c_{p_2} = 1, c_q = 1, b_1 = 2, b_2 = 0, b_3 = 1, b_4 = 1, b_5 = 0) & \xrightarrow{Q.f} \\
 (c_{p_1} = 0, c_{p_2} = 1, c_q = 0, b_1 = 0, b_2 = 3, b_3 = 1, b_4 = 1, b_5 = 0) & \xrightarrow{P.f_2} \\
 (c_{p_1} = 0, c_{p_2} = 0, c_q = 0, b_1 = 1, b_2 = 2, b_3 = 0, b_4 = 0, b_5 = 1) &
 \end{aligned}$$

Therefore, deadlock analysis succeeds (no deadlock is detected).

This example illustrates the tradeoff between modularity and reusability. For the same composite actor P , two profiles can be generated, as shown in Figure 3. These profiles achieve different tradeoffs. The monolithic profile shown to the left of the figure is more modular (i.e., smaller) than the non-monolithic one shown to the right. The latter is more reusable than the monolithic one, however: indeed, it can be reused in the graph with feedback shown at the left of Figure 2, whereas the monolithic one cannot be used, because it creates a deadlock.

Note that if we flatten the graph as shown in Figure 2, that is, remove composite actor P and replace it with its internal graph of atomic actors A and B , then the resulting graph has a PASS, i.e., exhibits no deadlock. This shows that deadlock is a result of using the monolithic profile, and not a problem with the graph itself. Of course, flattening is not an attractive proposition, because of scalability as well as IP issues, as we explained in the introduction.

If the deadlock analysis step fails then the graph is rejected. Otherwise, we proceed with the unfolding step.

7.4 Unfolding with SDF Profiles

This step takes as input the internal DSSF graph produced by the connection step, as well as the repetition vector produced by the rate analysis step. It produces as output a DAG (directed acyclic graph) that captures

the input-output dependencies of the DSSF. As mentioned in Section 4.3 this step is a generalization of existing transformations from SDF to HSDF. The difference is that we start from a DSSF graph instead of an SDF graph.

The DAG is computed in two steps. First, the DSSF graph is *unfolded*, by replicating each node in it as many times as specified in the repetition vector. These replicas represent the different firings of the corresponding actor. For this reason, the replicas are ordered: dependencies are added between them to represent the fact that the first firing comes before the second firing, the second before the third, and so on. The input and output ports of the actors are also replicated. Finally, for every edge of the original DSSF, a shared FIFO is created in the unfolded graph. The process is illustrated in Figure 7, for the internal DSSF graph of Figure 4. Rate analysis in this case produces the repetition vector $(r_A = 3, r_B = 2)$. Therefore $A.f$ is replicated 3 times and $B.f$ is replicated 2 times. In this example there is a single edge between $A.f$ and $B.f$, therefore, the unfolded graph contains a single shared FIFO. Note that we purposely do not consider the input and output ports to be shared FIFOs at this stage. This is precisely because we want to capture dependencies between each separate production and consumption of tokens at these ports.

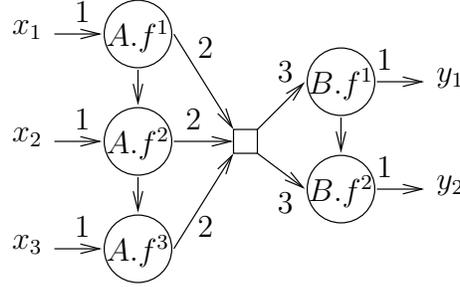


Figure 7: First step of unfolding the DSSF graph of Figure 4: replicating nodes and creating a shared FIFO.

In the second and final step of unfolding, the DAG is produced, by computing dependencies between the replicas. This is done by separately computing dependencies between replicas that are connected to a given FIFO, and repeating the process for every FIFO. We first explain the process for non-shared FIFOs, i.e., standard SDF edges, such as the one between A and B in the DSSF of Figure 4. Let $A \rightarrow B$ be such an edge. Suppose the edge has d initial tokens, A produces k tokens each time it fires, and B consumes n tokens each time it fires. Then the j -th occurrence of B depends on the i -th occurrence of A iff:

$$d + (i - 1) \cdot k < j \cdot n \quad (1)$$

In that case, an edge from $A.f^i$ to $B.f^j$ is added to the DAG. For the example of Figure 4, this gives the DAG shown in Figure 8.

In the general case, a FIFO queue in the internal DSSF of P may be shared by multiple producers and multiple consumers. Consider such a shared FIFO between a set of producers A_1, \dots, A_a and a set of consumers B_1, \dots, B_b . Let k_h be the number of tokens produced by A_h , for $h = 1, \dots, a$. Let n_h be the number of tokens consumed by B_h , for $h = 1, \dots, b$. Let d be the number of initial tokens in the queue. By construction (see Section 7.6) there is a total order $A_1 \rightarrow A_2 \rightarrow \dots \rightarrow A_a$ on the producers and a total order $B_1 \rightarrow B_2 \rightarrow \dots \rightarrow B_b$ on the consumers. As this is encoded with standard SDF edges of the form $A_i \xrightarrow{1} A_{i+1}$, this also implies that during rate analysis the rates of all producers will be found equal, and so will the rates of all consumers. Then, the j -th occurrence of B_u , $1 \leq u \leq b$, depends on the i -th occurrence of A_v , $1 \leq v \leq a$, iff:

$$d + (i - 1) \cdot \sum_{h=1}^a k_h + \sum_{h=1}^{v-1} k_h < (j - 1) \cdot \sum_{h=1}^b n_h + \sum_{h=1}^u n_h \quad (2)$$

Notice that, as should be expected, Equation (2) reduces to Equation (1) in the case $a = b = 1$.

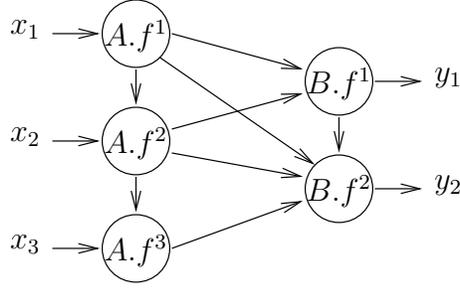


Figure 8: Unfolding the DSSF graph of Figure 4 produces the IODAG shown here.

In the DAG produced by unfolding, input and output port replicas are represented explicitly as special nodes with no predecessors and no successors, respectively. For this reason, we call this DAG an IODAG. Nodes of the IODAG that are neither input nor output are called *internal* nodes.

7.5 DAG Clustering

DAG clustering consists in partitioning the internal nodes of the IODAG produced by the unfolding step into a number of *clusters*. The clustering must be *valid* in the sense that it must not create cyclic dependencies among clusters³. Each of the clusters in the produced clustered graph will result in a firing function in the profile of P , as explained in Section 7.6 that follows. Exactly how DAG clustering is done is discussed in Section 8. There are many possibilities, that explore different tradeoffs, in terms of modularity, reusability, and other metrics. Here, we illustrate the *outcome* of DAG clustering on our running example. Two possible clusterings of the DAG of Figure 8 are shown in Figure 9, enclosed in dashed curves. The left-most clustering contains a single cluster, denoted C_0 . The right-most clustering contains two clusters, denoted C_1 and C_2 .

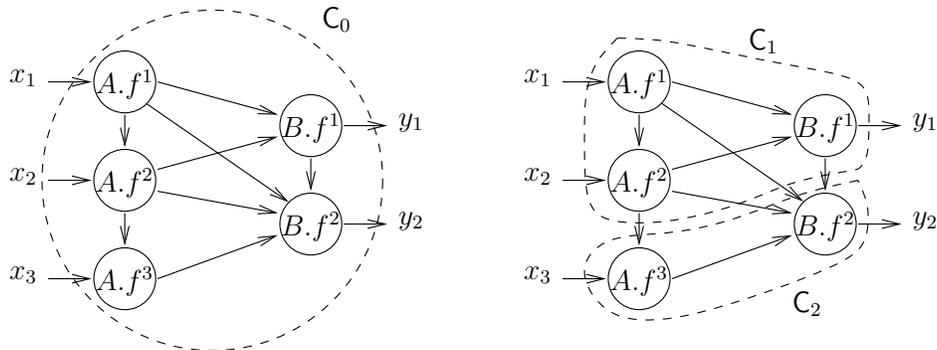


Figure 9: Two possible clusterings of the DAG of Figure 8.

7.6 Profile Generation

Profile generation is the last step in profile synthesis, where the actual profile of composite actor P is produced. The clustered graph, together with the internal diagram of P , completely determine the profile of P . Each cluster C_i is mapped to a firing function fire_i , and also to an atomic node $P.f_i$ in the profile graph of P . For every input (resp. output) port of P , an external, potentially shared, FIFO L is created in

³ A dependency between two distinct clusters exists iff there is a dependency between two nodes from each cluster.

the profile of P . For each cluster C_i , we compute the total number of tokens k_i read from (resp. written to) L by C_i : this can be easily done by summing over all actors in C_i . If $k_i > 0$ then an edge is added from L to $P.f_i$ (resp. from $P.f_i$ to L) annotated with a rate of k_i tokens.

Dependency edges between firing functions are computed as follows. For every pair of distinct clusters C_i and C_j , a dependency edge from $P.f_i$ to $P.f_j$ is added iff there exist nodes v_i in C_i and v_j in C_j such that v_j depends on v_i in the IODAG. Validity of clustering ensures that this set of dependencies results in no cycles. In addition to these dependency edges, we add a set of *backward* dependency edges to ensure a deterministic order of reading from resp. writing to shared FIFOs also across iterations. Let L be a shared FIFO of the profile. L is either an external FIFO, or an internal FIFO of P , like the one shown in Figure 7. Let W_L (resp. R_L) be the set of all clusters writing to (resp. reading from) L . By the fact that different replicas of the same actor that are created during unfolding are totally ordered in the IODAG, all clusters in W_L are totally ordered. Let $C_i, C_j \in W_L$ be the first and last clusters in W_L with respect to this total order, and let $P.f_i$ and $P.f_j$ be the corresponding firing functions. We encode the fact that the $P.f_i$ cannot re-fire before $P.f_j$ has fired, by adding a dependency edge from $P.f_j$ to $P.f_i$, supplied with an initial token. This is a backward dependency edge, $i \leq j$. Note that if $i = j$ then this edge is redundant. Similarly, we add a backward dependency edge, if necessary, among the clusters in R_L , which are also guaranteed to be totally ordered.

To establish the dataflow edges of the profile, we iterate over all internal (shared or non-shared) FIFOs of the internal DSSF of P . Let L be an internal FIFO and suppose it has d initial tokens. Let m be the total number of tokens produced at L by all clusters writing to L : by construction, m is equal to the total number of tokens consumed by all clusters reading from L . For our running example (Figures 4, 8 and 9), we have $d = 0$ and $m = 6$.

Conceptually, we define m output ports denoted z_0, z_1, \dots, z_{m-1} , and m input ports, denoted w_0, w_1, \dots, w_{m-1} . For $i = 0$ to $m - 1$, we connect output port z_i to input port w_j , where $j = (d + i) \div m$, and \div is the modulo operator. Intuitively, this captures the fact that the i -th token produced will be consumed as the $((d + i) \div m)$ -th token of some iteration, because of the initial tokens. Notice that $j = i$ when $d = 0$. We then place $\lfloor \frac{d+m-1-i}{m} \rfloor$ initial tokens at each input port w_i , for $i = 0, \dots, m - 1$, where $\lfloor v \rfloor$ is the integer part of v . Thus, if $d = 4$ and $m = 3$, then w_0 will receive 2 initial tokens, while w_1 and w_2 will each receive 1 initial token.

Finally, we assign the ports to producer and consumer clusters of L , according to their total order. For instance, for the non-monolithic clustering of Figure 9, output ports z_0 to z_3 and input ports w_0 to w_2 are assigned to C_1 , whereas output ports z_4, z_5 and input ports w_3 to w_5 are assigned to C_2 . Together with the port connections, these assignments define dataflow edges between the clusters. Self-loops (edges with same source and destination cluster) without initial tokens are removed. Note that more than one edges may exist between two distinct clusters, but these can always be merged into a single edge.

As an example, the two clusterings shown in Figure 9 give rise, respectively, to the two profiles shown in Figure 3. Another, self-contained example is shown in Figure 10. The two profiles shown at the bottom of the figure are generated from the clustering shown at the top-right, assuming 6 and 17 initial tokens in the queue from A to B , respectively. Notice that if the queue contains 17 tokens then this clustering is not optimal, in the sense that a more coarse-grain clustering exists. However, the clustering is valid, and used here to illustrate the profile generation process.

7.7 Code Generation

Once the profile has been synthesized, its firing functions need to be implemented. This is done in the code generation step. Every firing function corresponds to a cluster produced by the clustering step. The implementation of the firing function consists in calling in a sequential order all firing functions of internal actors that are included in the cluster. This sequential order can be arbitrary, provided it respects the dependencies of nodes in the cluster. We illustrate the process on our running example (Figures 1, 3 and 9).

Consider first the clustering shown to the left of Figure 9. This will result in a single firing function for P , namely, **P.fire**. Its implementation is shown below in pseudo-code:

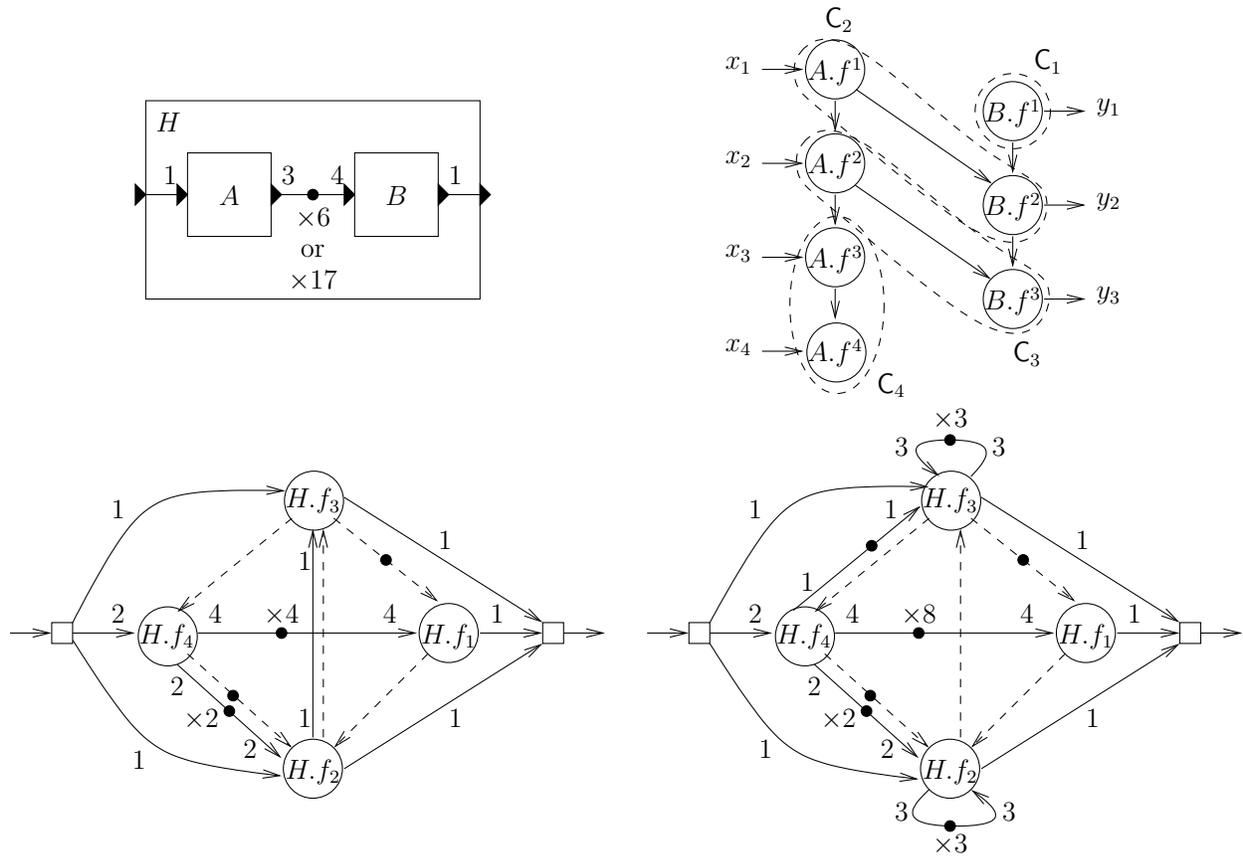


Figure 10: Composite SDF actor H (top-left); possible clustering produced by unfolding (top-right); SDF profiles generated for H , assuming 6 initial tokens in the queue from A to B (bottom-left); assuming 17 initial tokens (bottom-right). Firing functions $H.f_1, H.f_2, H.f_3, H.f_4$ correspond to clusters C_1, C_2, C_3, C_4 , respectively.

```

P.fire(input x[3]) output y[2]
{
  local tmp[4];
  tmp <- A.fire(x);
  tmp <- A.fire(x);
  y <- B.fire(tmp);
  tmp <- A.fire(x);
  y <- B.fire(tmp);
}

```

In the above pseudo-code, `tmp` is a local FIFO queue of length 4. Such a local queue is assumed to be empty when initialized. A statement such as `tmp <- A.fire(x)` corresponds to a call to firing function `A.fire`, providing as input the queue `x` and as output the queue `tmp`. `A.fire` will consume 1 token from `x` and will produce 2 tokens into `tmp`. When all statements of `P.fire` are executed, 3 tokens are consumed from the input queue `x` and 2 tokens are added to the output queue `y`, as indicated in the signature of `P.fire`.

Now let us turn to the clustering shown to the right of Figure 9. This clustering contains two clusters, therefore, it results in two firing functions for `P`, namely, `P.fire1` and `P.fire2`. Their implementation is shown below:

```

persistent local tmp[N]; /* N is a parameter */
assumption: N >= 4;

P.fire1(input x[2])          P.fire2(input x[1], tmp[1])
output y[1], tmp[1]         output y[1]
{
  tmp <- A.fire(x);          {
  tmp <- A.fire(x);          tmp <- A.fire(x);
  tmp <- A.fire(x);          y <- B.fire(tmp);
  y <- B.fire(tmp);         }
}

```

In this case `tmp` is declared to be a *persistent* local variable, which means its contents “survive” across calls to `P.fire1` and `P.fire2`. In particular, of the 4 tokens produced and added to `tmp` by the two calls of `A.fire` within the execution of `P.fire1`, only the first 3 are consumed by the call to `B.fire`. The remaining 1 token is consumed during the execution of `P.fire2`. This is why `P.fire1` declares to produce at its output `tmp[1]` (which means it produces a total of 1 token at queue `tmp` when it executes), and similarly, `P.fire2` declares to consume at its input 1 token from `tmp`.

Dependency edges are not implemented in the code, since they carry no useful data, and only serve to encode dependencies between firing function calls. These dependencies must be satisfied by construction in any correct usage of the profile. Therefore they do not need to be enforced in the implementation of the profile.

7.8 Discussion: Inadequacy of Cyclo-Static Data Flow

It may seem that *cyclo-static data flow* (CSDF) [3] can be used as an alternative representation of profiles. Indeed, this works on our running example: we could capture the composite actor `P` of Figure 1 using the CSDF actor shown in Figure 11. This CSDF actor specifies that `P` will iterate between two “firing modes”. In the first mode, it consumes 2 tokens from its input and produces 1 token at its output; in the second mode, it consumes 1 token and produces 1 token; the process is then repeated. This indeed works for this example: embedding `P` as shown in Figure 2 results in no deadlock, if the CSDF model for `P` is used.

In general, however, CSDF models are not expressive enough to be used as profiles. We illustrate this by two examples⁴, shown in Figures 6 and 12. Actors `R` and `W` are two composite actors shown in these

⁴ [9] also observe that CSDF is not a sufficient abstraction of composite SDF models, however, the example they use embeds a composite SDF graph into a dynamic data flow model. Therefore the overall model is not strictly SDF. The examples we provide are much simpler, in fact, the models are homogeneous SDF models.

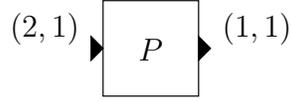


Figure 11: CSDF actor for composite actor P of Figure 1.

figures. All graphs are homogeneous (tokens rates are omitted in Figure 12, they are implicitly all equal to 1; dependency edges are also omitted from the profile, since they are redundant). Our method generates the profiles shown to the right of the two figures. The profile for R contains two completely independent firing functions, $R.f_1$ and $R.f_2$. CSDF cannot express this independence, since it requires a fixed order of firing modes to be specified statically. Although two separate CSDF models could be used to capture this example, this is not sufficient for composite actor W , which features both internal dependencies and independencies.

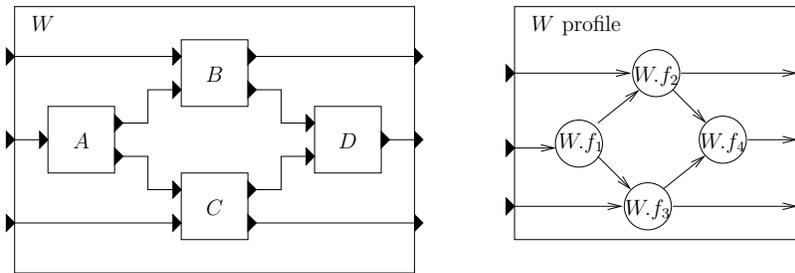


Figure 12: An example that cannot be captured by CSDF.

8 DAG Clustering

DAG clustering is at the heart of our modular code generation framework, since it determines the profile that is to be generated for a given composite actor. As mentioned above, different tradeoffs can be explored during DAG clustering, in particular, in terms of modularity and reusability: the more *fine-grain* the clustering is, the more reusable the profile and generated code will be; the more *coarse-grain* the clustering is, the more modular the code is, but also less reusable in general.

DAG clustering takes as input the IODAG produced by the unfolding step. A trivial way to perform DAG clustering is to produce a single cluster, that groups together all internal nodes in this DAG. This is called *monolithic* DAG clustering and results in monolithic profiles that have a single firing function. This clustering achieves maximal modularity, but it results in non-reusable code in general, as the discussion of Section 7.3 demonstrates.

In this section we describe clustering methods that achieve *maximal reusability*. This means that the generated profile (and code) can be used in *any* context where the corresponding flattened graph could be used. Therefore, the profile results in no information loss as far as reusing the graph in a given context is concerned. At the same time, the profile may be much smaller than the internal graph.

To achieve maximal reusability, we follow the ideas proposed in [16, 15] for SBDs. In particular, we present a clustering method that is guaranteed not to introduce *false input-output dependencies*. These dependencies are “false” in the sense that they are not induced by the original SDF graph, but only by the clustering method.

To illustrate this, consider the monolithic clustering shown to the left of Figure 9. This clustering introduces a false input-output dependency between the third token consumed at input x (represented by node x_3 in the DAG) and the first token produced at output y (represented by node y_1). Indeed, in order

to produce the first token at output y , only 2 tokens at input x are needed: these tokens are consumed respectively by the first two invocations of `A.fire`. The third invocation of `A.fire` is only necessary in order to produce the *second* token at y , but not the first one. The monolithic clustering shown to the left of Figure 9 loses this information. As a result, it produces a profile which is not reusable in the context of Figure 2, as demonstrated in Section 7.3. On the other hand, the non-monolithic clustering shown to the right of Figure 9 preserves the input-output dependency information, that is, does not introduce false dependencies. Because of this, it results in a maximally reusable profile.

The above discussion also helps to explain the reason for the unfolding step. Unfolding makes explicit the dependencies between different productions and consumptions of tokens *at the same ports*. In the example of actor P (Figure 4), even though there is a single (open) input port x and a single (open) output port y in the original DSSF graph, there are three copies of x and two copies of y in the unfolded DAG, corresponding to the three consumptions from x and two productions to y that occur within a PASS.

Unfolding is also important because it allows us to re-use the clustering techniques proposed for SBDs, which work on plain DAGs [16, 15]. In particular, we can use the so-called *optimal disjoint clustering* (ODC) method which is guaranteed not to introduce false IO dependencies, produces a set of pairwise *disjoint* clusters (clusters that do not share any nodes), and is *optimal* in the sense that it produces a minimal number of clusters with the above properties. Unfortunately, the ODC problem is shown to be NP-complete in [15]. This motivated us to develop a “greedy” DAG clustering algorithm, which is one of the contributions of this paper. Our algorithm is not optimal, i.e., it may produce more clusters than needed to achieve maximal reusability. On the other hand, the algorithm has polynomial complexity. The greedy DAG clustering algorithm that we present below is “backward” in the sense that it proceeds from outputs to inputs. A similar “forward” algorithm can be used, that proceeds from inputs to outputs.

8.1 Greedy Backward Disjoint Clustering

The greedy backward disjoint clustering (GBDC) algorithm is shown in Figure 13. GBDC takes as input an IODAG (the result of the unfolding step) $G = (V, E)$ where V is a finite set of nodes and E is a set of directed edges. V is partitioned in three disjoint sets: $V = V_{\text{in}} \cup V_{\text{out}} \cup V_{\text{int}}$, the sets of input, output and internal nodes, respectively. GBDC returns a partition of V_{int} into a set of disjoint sets, called clusters. The partition (i.e., the set of clusters) is denoted \mathcal{C} . The problem is non-trivial when all $V_{\text{in}}, V_{\text{out}}$ and V_{int} are non-empty (otherwise a single cluster suffices). In the sequel, we assume that this is the case.

\mathcal{C} defines a new graph, called the *quotient graph*, $G_{\mathcal{C}} = (V_{\mathcal{C}}, E_{\mathcal{C}})$. $G_{\mathcal{C}}$ contains clusters instead of internal nodes, and has an edge between two clusters (or a cluster and an input or output node) if the clusters contain nodes that have an edge in the original graph G . Formally, $V_{\mathcal{C}} = V_{\text{in}} \cup V_{\text{out}} \cup \mathcal{C}$, and $E_{\mathcal{C}} = \{(x, \mathcal{C}) \mid x \in V_{\text{in}}, \mathcal{C} \in \mathcal{C}, \exists f \in \mathcal{C} : (x, f) \in E\} \cup \{(\mathcal{C}, y) \mid \mathcal{C} \in \mathcal{C}, y \in V_{\text{out}}, \exists f \in \mathcal{C} : (f, y) \in E\} \cup \{(\mathcal{C}, \mathcal{C}') \mid \mathcal{C}, \mathcal{C}' \in \mathcal{C}, \mathcal{C} \neq \mathcal{C}', \exists f \in \mathcal{C}, f' \in \mathcal{C}', (f, f') \in E\}$. Notice that $E_{\mathcal{C}}$ does not contain self-loops (i.e., edges of the form (f, f)).

The steps of GBDC are explained below. E^* denotes the transitive closure of relation E : $(v, v') \in E^*$ iff there exists a path from v to v' , i.e., v' depends on v .

Identify input-output dependencies Given a node $v \in V$, let $\text{ins}(v)$ be the set of input nodes that v depends upon: $\text{ins}(v) := \{x \in V_{\text{in}} \mid (x, v) \in E^*\}$. Similarly, let $\text{outs}(v)$ be the set of output nodes that depend on v : $\text{outs}(v) := \{y \in V_{\text{out}} \mid (v, y) \in E^*\}$. For a set of nodes F , $\text{ins}(F)$ denotes $\bigcup_{f \in F} \text{ins}(f)$, and similarly for $\text{outs}(F)$.

Lines 1-3 of GBDC compute $\text{ins}(v)$ and $\text{outs}(v)$ for every node v of the DAG. We can compute these by starting from the output and following the dependencies backward. For example, consider the DAG of Figure 8. There are three input nodes, x_1, x_2, x_3 and two output nodes, y_1 and y_2 . We have: $\text{ins}(y_1) = \text{ins}(B.f^1) = \text{ins}(A.f^2) = \{x_1, x_2\}$, and $\text{ins}(y_2) = \text{ins}(B.f^2) = \{x_1, x_2, x_3\}$. Similarly: $\text{outs}(x_1) = \text{outs}(A.f^1) = \{y_1, y_2\}$, and $\text{outs}(x_3) = \text{outs}(A.f^3) = \{y_2\}$.

Line 4 of GBDC initializes \mathcal{C} to the empty set. Line 5 initializes Out as the set of internal nodes that have an output y as an immediate successor. These nodes will be used as “seeds” for creating new clusters (Lines 7-8).

Input: An IODAG $G = (V, E)$. $V = V_{\text{in}} \cup V_{\text{out}} \cup V_{\text{int}}$.
Output: A partition \mathcal{C} of the set of internal nodes V_{int} .

```

1 foreach  $v \in V$  do
2   | compute  $\text{ins}(v)$  and  $\text{outs}(v)$ ;
3 end
4  $\mathcal{C} := \emptyset$ ;
5  $\text{Out} := \{f \in V_{\text{int}} \mid \exists y \in V_{\text{out}} : (f, y) \in E\}$ ;
6 while  $\bigcup \mathcal{C} \neq V_{\text{int}}$  do
7   | partition  $\text{Out}$  into  $C_1, \dots, C_k$  such that two nodes  $f, f'$  are grouped in the same set  $C_i$  iff
8     |  $\text{ins}(f) = \text{ins}(f')$ ;
9   |  $\mathcal{C} := \mathcal{C} \cup \{C_1, \dots, C_k\}$ ;
10  | for  $i = 1$  to  $k$  do
11    | while  $\exists f \in C_i, f' \in V_{\text{int}} \setminus \bigcup \mathcal{C} : (f', f) \in E \wedge \forall x \in \text{ins}(C_i), y \in \text{outs}(f') : (x, y) \in E^*$  do
12      |  $C_i := C_i \cup \{f'\}$ ;
13    | end
14  | end
15  |  $\text{Out} := \{f \in V_{\text{int}} \setminus \bigcup \mathcal{C} \mid \neg \exists f' \in V_{\text{int}} \setminus \bigcup \mathcal{C} : (f, f') \in E\}$ ;
16 end
17 while quotient graph  $G_{\mathcal{C}}$  contains cycles do
18   | pick a cycle  $C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_k \rightarrow C_1$ ;
19   |  $\mathcal{C} := (\mathcal{C} \setminus \{C_1, \dots, C_k\}) \cup \{\bigcup_{i=1}^k C_i\}$ ;
20 end

```

Figure 13: The GBDC algorithm.

Then the algorithm enters the while-loop at Line 6. $\bigcup \mathcal{C}$ is the union of all sets in \mathcal{C} , i.e., the set of all nodes clustered so far. When $\bigcup \mathcal{C} = V_{\text{int}}$ all internal nodes have been added to some cluster, and the loop exits. The body of the loop consists in the following steps:

Partition seed nodes with respect to input dependencies Line 7 partitions Out into a set of clusters, such that two nodes are put into the same cluster iff they depend on the same inputs. Line 8 adds these newly created clusters to \mathcal{C} . In the example of Figure 8, this gives an initial $\mathcal{C} = \{\{B.f^1\}, \{B.f^2\}\}$.

Create a cluster for each group of seed nodes The for-loop starting at Line 9 iterates over all clusters newly created in the previous step and attempts to add as many nodes as possible to each of these clusters, going backward, and making sure no false input-output dependencies are created in the process. In particular, for each cluster C_i , we proceed backward, attempting to add unclustered predecessors f' of nodes f already in C_i (while-loop at Line 10). Such a node f' is a candidate to be added to C_i , but this happens only if an additional condition is satisfied: namely $\forall x \in \text{ins}(C_i), y \in \text{outs}(f') : (x, y) \in E^*$. This condition is violated if there exist an input node x that some node in C_i depends upon, and an output node y that depends on f' but not on x . In that case, adding f' to C_i would create a false dependency from x to y . Otherwise, it is safe to add f' , and this is done in Line 11.

In the example of Figure 8, executing the while-loop at Line 10 results in adding nodes $A.f^1$ and $A.f^2$ in the cluster $\{B.f^1\}$, and node $A.f^3$ in the cluster $\{B.f^2\}$, thereby obtaining the final clustering, shown to the right of Figure 9.

In general, more than one iteration may be required to cluster all the nodes. This is done by repeating the process, starting with a new Out set. In particular, Line 14 recomputes Out as the set of all unclustered nodes that have no unclustered successors.

Removing cycles The above process is not guaranteed to produce an acyclic quotient graph. Lines 16-19 remove cycles by repeatedly *merging* all clusters in a cycle into a single cluster. This process is guaranteed not to introduce false input-output dependencies, as shown in Lemma 5 of [15].

8.1.1 Termination and Complexity

Theorem 1 *Provided the set of nodes V is finite, GBDC always terminates.*

Proof: G is acyclic, therefore the set `Out` computed in Lines 5 and 14 is guaranteed to be non-empty. Therefore, at least one new cluster is added at every iteration of the while-loop at Line 6, which means the number of unclustered nodes decreases at every iteration. The for-loop and foreach-loop inside this while-loop obviously terminate, therefore, the body of the while-loop terminates. The second while-loop (Lines 16-19) terminates because the number of cycles is reduced by at least one at every iteration of the loops, and there can only be a finite number of cycles. ■

Theorem 2 *GBDC is polynomial in the number of nodes in G .*

Proof: We provide only a rough and largely pessimistic complexity analysis. A more accurate analysis is beyond the scope of this paper.

Let $n = |V|$ be the number of nodes in G . Computing sets `ins` and `outs` can be done in $O(n^2)$ time (perform forward and backward reachability from every node). Computing `Out` can also be done in $O(n^2)$ time (Line 5 or 14). The while-loop at Line 6 is executed at most n times. Partitioning `Out` (Line 7) can be done in $O(n^3)$ time and this results in $k \leq n$ clusters. The while-loop at Lines 10-12 is iterated no more than n^2 times and the *safe-to-add- f'* condition can be checked in $O(n^3)$ time. The quotient graph produced by the while-loop at Line 6 contains at most n nodes. Checking the condition at Line 16 can be done in $O(n)$ time, and this process also returns a cycle, if one exists. Executing Line 18 can also be done in $O(n)$ time. The loop at Line 16 can be executed at most n times, since at least one cluster is removed every time. ■

8.1.2 Correctness

GBDC is correct, in the sense that, first, it produces disjoint clusters and clusters all internal nodes, second, the resulting clustered graph is acyclic, and third, the resulting graph contains no input-output dependencies that were not already present in the input graph.

Theorem 3 *GBDC produces disjoint clusters and clusters all internal nodes.*

Proof: Disjointness is ensured by the fact that only unclustered nodes (i.e., nodes in $V_{\text{int}} \setminus \bigcup \mathcal{C}$) are added to the set `Out` (Lines 5 and 14) or to a newly created cluster C_i (Line 11). That all internal nodes are clustered is ensured by the fact that the while-loop at Line 6 does not terminate until all internal nodes are clustered. ■

Theorem 4 *GBDC results in an acyclic quotient graph.*

Proof: This is ensured by the fact that all potential cycles are removed in Lines 16-19. ■

Theorem 5 *GBDC produces a quotient graph G_C that has the same input-output dependencies as the original graph G .*

Proof: We need to prove that: $\forall x \in V_{\text{in}}, y \in V_{\text{out}} : (x, y) \in E^* \iff (x, y) \in E_C^*$. We will show that this holds for the quotient graph produced when the while-loop of Lines 6-15 terminates. The fact that Lines 16-19 preserve IO dependencies is shown in Lemma 5 of [15].

The \Rightarrow direction is trivial by construction of the quotient graph. There are two places where false IO dependencies can potentially be introduced in the while-loop of Lines 6-15: at Lines 7-8, where a new set of clusters is created and added to \mathcal{C} ; or at Line 11, where a new node is added to an existing cluster. We examine each of these cases separately.

Consider first Lines 7-8: A certain number $k \geq 1$ of new clusters are created here, each containing one or more nodes. This can be seen as a sequence of operations: first, create cluster C_1 with a single node $f \in \text{Out}$, then add to C_1 a node $f' \in \text{Out}$ such that $\text{ins}(f) = \text{ins}(f')$ (if such an f' exists), and so on, until C_1 is complete; then create cluster C_2 with a single node, and so on, until all clusters C_1, \dots, C_k are complete. It suffices to show that no such creation or addition results in false IO dependencies.

Regarding creation, note that a cluster that contains a single node cannot add false IO dependencies, by definition of the quotient graph. Regarding addition, we claim that if a cluster C is such that $\forall f, f' \in C : \text{ins}(f) = \text{ins}(f')$, then adding a node f'' such that $\text{ins}(f'') = \text{ins}(f)$, where $f \in C$, results in no false IO dependencies. To see why the claim is true, let $y \in \text{outs}(f'')$. Then $\text{ins}(f'') \subseteq \text{ins}(y)$. Since $\text{ins}(f'') = \text{ins}(f)$, for any $x \in \text{ins}(f)$, we have $(x, y) \in E^*$. Similarly, for any $y \in \text{outs}(f)$ and any $x \in \text{ins}(f'')$, we have $(x, y) \in E^*$.

Consider next Line 11: The fact that f' is chosen to be a predecessor of some node $f \in C_i$ implies that $\text{ins}(f') \subseteq \text{ins}(f) \subseteq \text{ins}(C_i)$. There are two cases where a new dependency can be introduced: Case 2(a): either between some input $x \in \text{ins}(C_i)$ and some output $y \in \text{outs}(f')$; Case 2(b): or between some input $x' \in \text{ins}(f')$ and some output $y' \in \text{outs}(C_i)$. In Case 2(a), the safe-to-add- f' condition at Line 10 ensures that if such x and y exist, then y already depends on x , otherwise, f' is not added to C_i . In Case 2(b), $\text{ins}(f') \subseteq \text{ins}(C_i)$ implies $x' \in \text{ins}(C_i)$. This and $y' \in \text{outs}(C_i)$ imply that $(x', y') \in E^*$: indeed, if this is not the case, then cluster C_i already contains a false IO dependency before the addition of f' . ■

8.2 Clustering for Closed Models

It is worth discussing the special case where clustering is applied to a closed model, that is, a model where all input ports are connected. This in particular happens with top-level models used for simulation, which contain source actors that provide the input data. By definition, the IODAG produced by the unfolding step for such a model contains no input ports. In this case, a monolithic clustering that groups all nodes into a single cluster suffices and the GBDC algorithm produces the monolithic clustering for such a graph. Such a clustering will automatically give rise to a single firing function. Simulating the model then consists in calling this function repeatedly.

9 Implementation

We have built a preliminary implementation of the SDF modular code generation described above in the open-source Ptolemy II framework [8] (<http://ptolemy.org/>). The implementation uses a specialized class to describe composite SDF actors for which DSSF profiles can be generated. These profiles are captured in Java, and can be loaded when the composite actor is used within another composite. For debugging and documentation purposes, the tool also generates in the GraphViz format DOT (<http://www.graphviz.org/>) the graphs produced by the unfolding and clustering steps.

Using our tool, we can, for instance, generate automatically a profile for the Ptolemy II model depicted in Figure 14. This model captures the SDF graph given in Figure 3 of [9]. Actor A2 is a composite actor designed so as to consume 2 tokens on each of its input ports and produce 2 tokens on each of its output ports each time it fires. For this, it uses the DownSample and UpSample internal actors: DownSample consumes 2 tokens at its input and produces 1 token at its output; UpSample consumes 1 token at its input and produces 2 tokens at its output. Actors A1 and A3 are homogeneous. The SampleDelay actor models an initial token in the queue from A2 to A3. All other queues are initially empty.

Assuming a monolithic profile for A2, GBDC generates for the top-level Ptolemy model the clustering shown to the left of Figure 15. This graph is automatically generated by DOT from the textual output

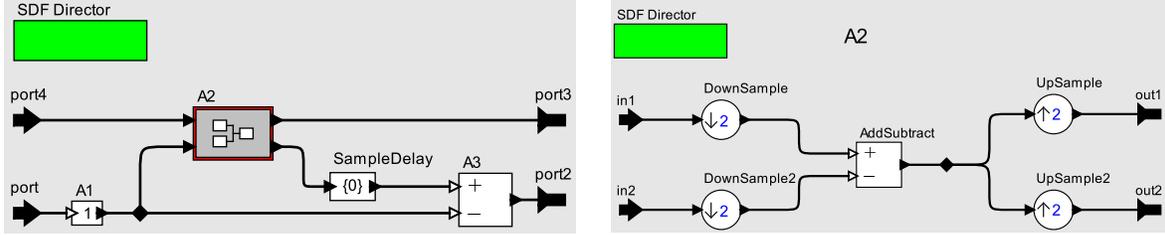


Figure 14: A hierarchical SDF model in Ptolemy II. The internal diagram of composite actor A2 is shown to the right.

automatically generated by our tool. The two replicas of A1 are denoted A1_1.0 and A1_2.0, respectively, and similarly for A2 and A3. Two clusters are generated, giving rise to the profile shown to the right of the figure. It is worth noting that there are 4 identical backward dependency edges generated for this profile (only one is shown). Moreover, all dependency edges are redundant in this case, thus can be removed. Finally, notice that the profile contains only two nodes, despite the fact that the Ptolemy model contains 9 actors overall.

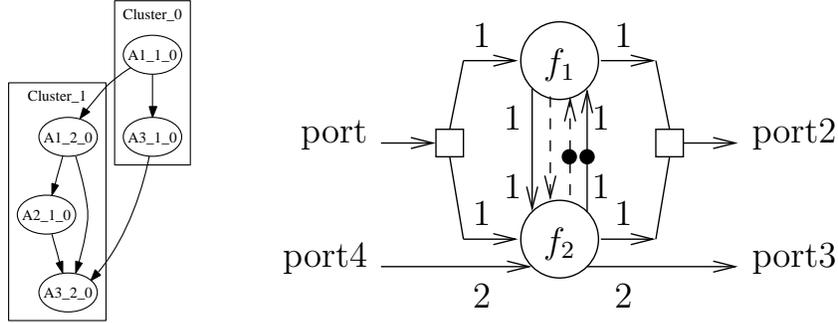


Figure 15: Clustering (left) and DSSF profile (right) of the model of Figure 14.

10 Conclusions and Perspectives

Hierarchical SDF models are not compositional: a composite SDF actor cannot be represented as an atomic SDF actor without loss of information that can lead to deadlocks. Extensions such as CSDF are not compositional either. In this paper we introduced DSSF profiles as a compositional representation of composite actors and showed how this representation can be used for modular code generation. In particular, we provided algorithms for automatic synthesis of DSSF profiles of composite actors given DSSF profiles of their sub-actors. This allows to handle hierarchical models of arbitrary depth. We showed that different trade-offs can be explored when synthesizing profiles, in terms of modularity (keeping the size of the generated DSSF profile minimal) versus reusability (preserving information necessary to avoid deadlocks) as well as algorithmic complexity. We provided a heuristic DAG clustering method that has polynomial complexity and ensures maximal reusability.

In the future, we plan to examine how other DAG clustering algorithms could be used in the SDF context. This includes the clustering algorithm proposed in [16], which may produce overlapping clusters, with nodes shared among multiple clusters. This algorithm is interesting because it guarantees an upper bound on the number of generated clusters, namely, $n + 1$, where n is the number of outputs in the DAG. Overlapping

clusters result in complications during profile generation that need to be resolved.

Another important problem is efficiency of the generated code. Different efficiency goals may be desirable, such as buffer size, code length, and so on. Problems of code optimization in the SDF context have been extensively studied in the literature, see, for instance [2, 19]. One direction of research is to adapt existing methods to the modular SDF framework proposed here.

We would also like to study possible applications of DSSF to contexts other than modular code generation, for instance, compositional performance analysis, such as throughput or latency computation. Finally, we plan to study possible extensions to dynamic data flow models.

Acknowledgments

We would like to thank Jörn Janneck and Maarten Wiggers for their valuable input.

References

- [1] F. Arbab. Abstract behavior types: a foundation model for components and their composition. *Sci. Comput. Program.*, 55(1-3):3–52, 2005.
- [2] S. Bhattacharyya, E. Lee, and P. Murthy. *Software Synthesis from Dataflow Graphs*. Kluwer, 1996.
- [3] G. Bilsen, M. Engels, R. Lauwereins, and J. A. Peperstraete. Cyclo-static data flow. In *IEEE Int. Conf. ASSP*, pages 3255–3258, May 1995.
- [4] S. Bliudze and J. Sifakis. The algebra of connectors: structuring interaction in bip. In *Embedded Software (EMSOFT’07)*, pages 11–20, New York, NY, USA, 2007. ACM.
- [5] J.D. Brock and W.B. Ackerman. Scenarios: A model of non-determinate computation. In *Proc. Intl. Colloq. on Formalization of Programming Concepts*, pages 252–259, London, UK, 1981. Springer-Verlag.
- [6] L. de Alfaro and T. Henzinger. Interface automata. In *Foundations of Software Engineering (FSE)*. ACM Press, 2001.
- [7] L. de Alfaro and T. Henzinger. Interface theories for component-based design. In *Embedded Software (EMSOFT’01)*. Springer, LNCS 2211, 2001.
- [8] J. Eker, J. Janneck, E. Lee, J. Liu, X. Liu, J. Ludvig, S. Neuendorffer, S. Sachs, and Y. Xiong. Taming heterogeneity – the Ptolemy approach. *Proc. IEEE*, 91(1), January 2003.
- [9] J. Falk, J. Keinert, C. Haubelt, J. Teich, and S. Bhattacharyya. A generalized static data flow clustering algorithm for mpso scheduling of multimedia applications. In *Embedded Software – EMSOFT’08*, pages 189–198. ACM, 2008.
- [10] M.C.W. Geilen. Reduction of Synchronous Dataflow Graphs. In *Design Automation Conference, DAC 2009*. ACM, 2009.
- [11] B. Jonsson. A fully abstract trace model for dataflow and asynchronous networks. *Distrib. Comput.*, 7(4):197–212, 1994.
- [12] G. Kahn. The semantics of a simple language for parallel programming. In *Information Processing 74, Proceedings of IFIP Congress 74*. North-Holland, 1974.
- [13] H. Kopetz. Elementary versus composite interfaces in distributed real-time systems. In *ISADS’99: 4th Intl. Symp. Autonomous Decentralized Systems*, pages 1–8. IEEE, 1999.
- [14] E.A. Lee and D.G. Messerschmitt. Static scheduling of synchronous data flow programs for digital signal processing. *IEEE Trans. Comput.*, 36(1):24–35, 1987.

- [15] R. Lubliner, C. Szegedy, and S. Tripakis. Modular code generation from synchronous block diagrams: modularity vs. code size. In *Principles of Programming Languages – POPL’09*, pages 78–89. ACM, January 2009.
- [16] R. Lubliner and S. Tripakis. Modularity vs. reusability: code generation from synchronous block diagrams. In *Design, Automation and Test in Europe – DATE’08*, pages 1504–1509. ACM, March 2008.
- [17] N. Lynch and M. Tuttle. Hierarchical correctness proofs for distributed algorithms. In *PODC’87: Proc. 6-th ACM Symp. on Principles of Distributed Computing*, pages 137–151, New York, NY, USA, 1987. ACM.
- [18] J.L. Pino, S.S. Bhattacharyya, and E.A. Lee. A Hierarchical Multiprocessor Scheduling Framework for Synchronous Dataflow Graphs. Technical Report UCB/ERL M95/36, EECS Department, University of California, Berkeley, 1995.
- [19] S. Sriram and S. Bhattacharyya. *Embedded Multiprocessors: Scheduling and Synchronization – 2nd ed.* CRC Press, 2009.
- [20] Eugene W. Stark. An algebra of dataflow networks. *Fundam. Inform.*, 22(1/2):167–185, 1995.
- [21] W. Thies, M. Karczmarek, and S. Amarasinghe. StreamIt: A language for streaming applications. In *11th Intl. Conf. on Compiler Construction, CC’02*, volume LNCS 2304. Springer, 2002.