



# Symmetry Breaking for Multi-Criteria Mapping and Scheduling on Multicores

*Pranav Tendulkar, Peter Poplavko, Oded Maler*

**Verimag Research Report n° TR-2013-3**

11-March-2013

Reports are downloadable at the following address  
<http://www-verimag.imag.fr>

Unité Mixte de Recherche 5104 CNRS - Grenoble INP - UJF

Centre Equation  
2, avenue de VIGNATE  
F-38610 GIERES  
tel : +33 456 52 03 40  
fax : +33 456 52 03 50  
<http://www-verimag.imag.fr>



# Symmetry Breaking for Multi-Criteria Mapping and Scheduling on Multicores<sup>1</sup>

*Pranav Tendulkar, Peter Poplavko, Oded Maler*

11-March-2013

## Abstract

Multiprocessor mapping and scheduling is a long-old difficult problem. In this work we propose a new methodology to perform mapping and scheduling along with buffer memory optimization using an SMT solver. We target split-join graphs, a formalism inspired by synchronous data-flow (SDF) which provides a compact symbolic representation of data-parallelism. Unlike the traditional design flow for SDF which involves splitting of a big problem into smaller heuristic sub-problems, we deal with this problem as a whole and try to compute exact Pareto-optimal solutions for it. We introduce symmetry breaking constraints in order to reduce the run-times of the solver. We have tested our work on a number of SDF graphs and demonstrated the practicality of our method. We validate our models by running an image decoding application on the Tilera multi-core platform.

**Keywords:** synchronous data-flow, multiprocessor, multi-core, mapping, scheduling, SMT, SAT solver

**Reviewers:** Oded Maler

## How to cite this report:

```
@techreport {TR-2013-3,  
  title = {Symmetry Breaking for Multi-Criteria Mapping and Scheduling on Multicores},  
  author = {Pranav Tendulkar, Peter Poplavko, Oded Maler},  
  institution = {{Verimag} Research Report},  
  number = {TR-2013-3},  
  year = {}  
}
```

---

<sup>1</sup>This report is an extension of paper [23]

## 1 Introduction

This work is motivated by a key important problem in contemporary computing: *how to exploit efficiently the resources provided by a multi-core platform while executing application programs*. The problem has many variants depending on the intended use of the platform (general-purpose server or a dedicated accelerator), the specifics of the architecture (memory hierarchy, interconnect), the granularity of parallelism (instruction level, task level), the class of applications and the programming model. We focus on applications such as video, audio and other forms of signal processing which are naturally structured in a *data-flow* style as a network of interconnected software components (actors, filters, tasks). Such a description already exposes the precedence constraints among tasks and hence the task-level parallelism inherent in the application. More specifically we address applications written as *split-join graphs*, which can be viewed as a variant of the Synchronous Data-Flow (SDF) formalism [13, 20], or an abstract semantic models of a subset of a streaming languages such as StreaMIT [25]. Such formalisms, in addition to precedence constraints, also provide a compact *symbolic representation* of data-parallelism, namely, the presence of numerous tasks which are identical in nature and can be executed independently. Once the split-join graph is annotated with execution time figures and the data-parallel tasks have been explicitly expanded we obtain a task graph [3] whose *deployment* on the execution platform is the subject of optimization.

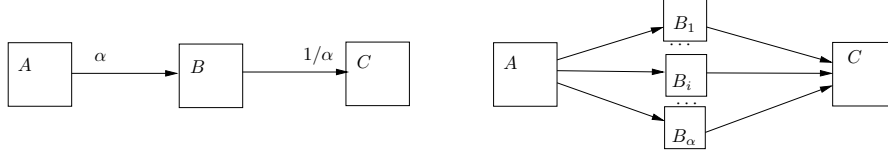
The deployment decisions that we consider and which may affect cost and performance are the following. First we can vary the number of processors used which gives a rough estimation of the cost of the platform (and its static power consumption). On a given configuration it remains to *map* tasks to processors, and to *schedule* the execution order on each processor. The performance measures to evaluate such a deployment are the total execution time (*latency*) and the size of the communication *buffers* which depends on the execution order. This is a multi-criteria (cost, latency and buffer size) optimization problem whose single-criterion version is already intractable. We take advantage of recent progress in SMT (SAT modulo theory) solvers [17, 7] to provide a good approximation of the Pareto front of the problem. We encode the precedence and resource constraints of the problem in the theory of linear arithmetic and, following [15, 14], we submit queries to the solver concerning the existence of solutions whose costs reside in various parts of the multi-dimensional cost space. Based on the answers to these queries we obtain a good approximation of the optimal trade-off between these criteria. The major computational obstacle is the intractability of mapping and scheduling problems aggravated by the exponential blow-up while expanding the graph from symbolic to explicit form. We tackle this problem by introducing “symmetry breaking” constraints among identical processors and identical tasks. For the latter we prove a theorem concerning the optimality of schedules where instances of the same actor are executed according to a fixed *lexicographical* order.

The rest of the paper is organized as follows. In Section 2 we give some background on split-join graphs and their transformation into task graphs and prove a useful property of their optimal schedules. In Section 3 we write down in more detail the constraint-based formulation of deployment and present our multi-criteria cost-space exploration procedure. An experimental evaluation of our approach appears in Section 4, including a validation on the Tilera multi-core platform. We conclude by discussing related and future work.

## 2 Split-Join Graphs

A parallelization factor is any number of the form  $\alpha$  (split) or  $1/\alpha$  (join) for  $\alpha \in \mathbb{N}$ . We use  $\Sigma^*$  to denote the set of sequences over a set  $\Sigma$  and use  $\sqsubset$  for the prefix relation with  $\xi \sqsubset \xi \cdot \xi'$ , where  $\xi \cdot \xi'$  denotes concatenation.

**Definition 2.1 (Split-Join and Task Graphs)** *A split-join graph is  $S = (V, E, d, r)$  where  $(V, E)$  is a directed acyclic graph (DAG), that is, a set  $V$  of nodes, a set  $E \subseteq V \times V$  of edges and no cyclic paths. The function  $d : V \rightarrow \mathbb{R}_+$  defines the execution times of the nodes and  $r : E \rightarrow \mathbb{Q}$  assigns a parallelization factor to every edge. An edge  $e$  is a split, join or neutral edge depending on whether  $r(e) > 1, < 1$  or  $= 1$ . A split-join graph with  $r(e) = 1$  for every  $e$  is called a task-graph and is denoted by  $T = (U, \mathcal{E}, \delta)$ .*



**Figure 1:** A simple split-join graph and its expanded task graph. Actor  $B$  has  $\alpha$  instances.

The decomposability of a task into parallelizable sub-tasks is expressed as a numerical label (parallelization factor) on a precedence edge leading to it. A label  $\alpha$  on the edge from  $A$  to  $B$  means that every executed instance of task  $A$  spawns  $\alpha$  instances of task  $B$ . Likewise, a  $1/\alpha$  label on the edge from  $B$  to  $C$  means that all those instances of  $B$  should terminate and their outputs be joined before executing  $C$  (see Fig. 1). A task graph can thus be viewed as obtained from the split-join graph by making data parallelism *explicit*. To distinguish between these two types of graphs we call the nodes of the split-join graphs *actors* (task types) and those of the task graph *tasks*.

The DAG structure naturally induces a partial-order relation  $\angle$  over the actors such that  $v \angle v'$  if there is a path from  $v$  to  $v'$ . The set of *minimal* elements with respect to  $\angle$  is  $V_\bullet \subseteq V$  consisting of nodes with no incoming edges. Likewise, the maximal elements  $V^\bullet$  are those without outgoing edges. An *initialized path* in a DAG is an alternating sequence of nodes and edges  $\pi = v_1 \cdot e_1 \cdot v_2 \cdots v_k$  starting from some  $v_1 \in V_\bullet$ . Such a path is *complete* if  $v_k \in V^\bullet$ . With any such path we associate the *multiplicity signature*

$$\xi(\pi) = (v_1, \alpha_1) \cdot (v_2, \alpha_2) \cdots (v_{k-1}, \alpha_{k-1})$$

where  $\alpha_i = r(e_i)$ . We will also abuse  $\xi$  to denote the projection of the signature on the multiplication factors, that is  $\xi(\pi) = \alpha_1 \cdot \alpha_2 \cdots \alpha_{k-1}$ .

To ensure that different instances of the same actor communicate with the matching instances of other actors and that such instances are joined together properly, we need an *indexing scheme* similar to indices of multi-dimensional arrays accessed inside nested loops. Because an actor may have several ancestral paths, we need to ensure that its indices via different paths agree. This will be guaranteed by a well-formedness condition that we impose on the multiplicity signatures along paths.

**Definition 2.2 (Parenthesis Alphabet)** Let  $\Sigma = \{1\} \cup \Sigma_l \cup \Sigma_j$  be any set of symbols consisting of a special symbol  $1$  and two finite sets  $\Sigma_l$  and  $\Sigma_j$  admitting a bijection which maps every  $\alpha \in \Sigma_l$  to  $\alpha' \in \Sigma_j$ .

Intuitively  $\alpha$  and  $\alpha'$  correspond to a matching pair consisting of a split  $\alpha$  and its inverse join  $1/\alpha$ . These can be viewed also as a pair of (typed) left and right *parentheses*.

**Definition 2.3 (Canonical Form)** The canonical form of a sequence  $\xi$  over a parentheses alphabet  $\Sigma$  is the sequence  $\bar{\xi}$  obtained from  $\xi$  by erasing occurrences of the neutral element  $1$  as well as matching pairs of the form  $\alpha \cdot \alpha'$ .

For example, the canonical form of  $\xi = 5 \cdot 1 \cdot 3 \cdot 1 \cdot 1 \cdot 1 / 3 \cdot 1 \cdot 2$  is  $\bar{\xi} = 5 \cdot 2$ . Note that the (arithmetic) products of the factors of  $\xi$  and of  $\bar{\xi}$  are equal and we denote this value by  $c(\xi)$ . A sequence  $\xi$  is *well-parenthesized* if  $\bar{\xi} \in \Sigma_l^*$ , namely its canonical form consists only of left parentheses. Note that this notion refers also to signature *prefixes* that can be *extended* to well-balanced sequences, namely, sequences with no violation of being well-parenthesized by a join not compatible with the *last* open split.

**Definition 2.4 (Well Formedness)** A split-join graph is well formed if:

1. Any complete path  $\pi$  satisfies  $c(\xi(\pi)) = 1$ ;
2. The signatures of all initialized paths are well parenthesized.

The first condition ensures that the graph is meaningful (all splits are joined) and that the multiplicity signatures of any two paths leading to the same actor  $v$  satisfy  $c(\xi) = c(\xi')$ . We can thus associate

unambiguously this number with the actor itself and denote it by  $c(v)$ . This *execution count* is the number of instances of actor  $v$  that should be executed.

The second condition forbids, for example, sequences of the form  $2 \cdot 3 \cdot 1/2 \cdot 1/3$ . It implies an additional property: every two initialized paths  $\pi$  and  $\pi'$  leading to the same actor satisfy  $\bar{\xi}(\pi) = \bar{\xi}(\pi')$ . Otherwise, if two paths would reach the same actor with different canonical signatures, there will be no way to close their parentheses by the same path suffix. Although split-join graphs *not* satisfying Condition 2 can make sense for certain computations, they require more complicated mappings between tasks and they will not be considered here. For well-formed graphs, a *unique canonical signature*, denoted by  $\bar{\xi}(v)$ , is associated with every actor.

**Definition 2.5 (Indexing Alphabet and Order)** *An actor  $v$  with  $\bar{\xi}(v) = \alpha_1 \cdots \alpha_k$  defines an indexing alphabet  $A_v$  consisting of all  $k$ -digit sequences  $h = a_1 \cdots a_k$  such that  $0 \leq a_i \leq \alpha_i - 1$ . This alphabet can be mapped into  $\{0, \dots, c(v) - 1\}$  via the following recursive rule:*

$$\mathcal{N}(\varepsilon) = 0 \quad \text{and} \quad \mathcal{N}(h \cdot a_j) = \alpha_j \cdot \mathcal{N}(h) + a_j$$

We use  $\ll_v$  to denote the lexicographic total order over  $A_v$  which coincides with the numerical order over  $\mathcal{N}(A_v)$ .

Every instance of actor  $v$  will be indexed by some  $h \in A_v$  and will be denoted as  $v_h$ . We use notation  $h$  and  $A_v$  to refer both to strings and to their numerical interpretation via  $\mathcal{N}$ . In the latter case  $v_h$  will refer to the task in position  $h$  according to the lexicographic order  $\ll_v$ .

Algorithm 1 scans the split-join graph forward and generates incrementally the corresponding task graph. If  $v$  is a parent of  $v'$  with factor  $\alpha$ , every instance of  $v$  will yield  $\alpha$  instances of  $v'$  which will later be merged when  $\alpha'$  is encountered. We assume implicitly that inserting an edge  $(u, u')$  into  $\mathcal{E}$  involves the insertion of  $u'$  to  $U$  in case it is not there already. The duration of a task is inherited from its actor, that is,  $\delta(v_h) = d(v)$ . The tasks can be partitioned naturally according to their actors, letting  $U = \bigcup_{v \in V} U_v$  and  $U_v = \{v_h : h \in A_v\}$ . A permutation  $\omega : U \rightarrow U$  is *actor-preserving* if it can be written as  $\omega = \bigcup_{v \in V} \omega_v$  and each  $\omega_v$  is a permutation on  $U_v$ .

**Definition 2.6 (Deployment)** *A deployment for a task graph  $T = (U, \mathcal{E}, \delta)$  on an execution platform with a finite set  $M$  of processors consists of a mapping function  $\mu : U \rightarrow M$  and a scheduling function  $s : U \rightarrow \mathbb{R}_+$  indicating the start time of each task and satisfying precedence and mutual exclusion constraints, namely, for each pair of tasks we have:*

$$\text{Precedence:} \quad (u, u') \in \mathcal{E} \Rightarrow s(u') - s(u) \geq \delta(u)$$

$$\text{Mutual exclusion:} \quad \mu(u) = \mu(u') \Rightarrow [(s(u') - s(u) \geq \delta(u)) \vee (s(u) - s(u') \geq \delta(u'))]$$

Note that  $\mu(u)$  and  $s(u)$  are decision variables while  $\delta(u)$  is a constant. The latency of the deployment is the termination time of the last task,  $\max_{u \in U} (s(u) + \delta(u))$ .

The problem of optimal scheduling of a task-graph is already NP-hard due to the non-convex mutual exclusion constraints. This situation is aggravated by the fact that the task-graph will typically be exponential in the size of the split-join graph. On the other hand, it admits many tasks which are identical in their duration and isomorphic in their precedence constraints. In the sequel we exploit this symmetry by showing that all tasks that correspond to the same actor can be executed according to a *lexicographic order* without compromising latency.

**Definition 2.7 (Ordering Scheme)** *An ordering scheme for a task-graph  $T = (U, \mathcal{E}, \delta)$  derived from a split-join graph  $G = (V, E, r, d)$  is a relation  $\prec = \bigcup_{v \in V} \prec_v$  where each  $\prec_v$  is a total order relation on  $U_v$ . An ordering scheme is *prefix-consistent* if for every  $(v, v') \in E$  and every  $h \neq h' \in A_v$ , we have*

$$(r(v, v') = 1) \Rightarrow [(v_h \prec v_{h'}) \iff (v'_h \prec v'_{h'})]$$

and

$$[(r(v, v') = \alpha) \vee (r(v', v) = 1/\alpha)] \Rightarrow [(v_h \prec v_{h'}) \iff (\forall a, a' v'_{h \cdot a} \prec v'_{h' \cdot a'})]$$

**Algorithm 1** Transforming a split-join graph into a task graph

---

**Input:** a well-formed split-join graph  $S = (V, E, d, r)$   
**Output:** An equivalent task graph  $T = (U, \mathcal{E}, \delta)$

**Initialization**  
 $F := V_\bullet$  % frontier  
 $U := \{v_\varepsilon : v \in V_\bullet\}$  % empty index for minimal nodes that execute once  
 $\mathcal{E} := \emptyset$  % no edges initially

**while**  $F \neq \emptyset$  **do**  
  pick  $v \in F$  % take an actor in the frontier  
  **for every**  $e = (v, v') \in E$  **do**  
     $F := F \cup \{v'\}$  % update frontier  
    **case**  
       $r(e) = 1$ : % neutral  
      **for each**  $v_h \in U$   
         $\mathcal{E} := \mathcal{E} \cup (v_h, v'_h)$   
       $r(e) = \alpha > 1$ : % split  
      **for each**  $v_h \in U$   
        **for**  $a = 1.. \alpha$   
           $\mathcal{E} := \mathcal{E} \cup (v_h, v'_{h.a})$   
       $r(e) = 1/\alpha$ : % join  
      **for each**  $v_{h.a} \in U$   
         $\mathcal{E} := \mathcal{E} \cup (v_{h.a}, v'_h)$   
    **endcase**  
  **endfor**  
   $F := F - \{v\}$  % update frontier  
**endwhile**

---

The lexicographic ordering scheme  $\ll$  is prefix consistent. We say that a schedule  $s$  is *compatible* with an ordering scheme  $\prec$  if  $v_h \prec v_{h'}$  implies  $s(v_h) \leq s(v_{h'})$ . We denote such an ordering scheme by  $\prec^s$  and use notation  $v[h]$  for the task occupying position  $h$  in  $\prec_v^s$ .

**Lemma 2.1** *Let  $s$  be a feasible schedule and let  $v$  and  $v'$  be two actors such that  $(v, v') \in E$ . Then*

1. *If  $r(v, v') = \alpha \geq 1$ , then for every  $h \in [0, c(v) - 1]$  and every  $a \in [0, \alpha - 1]$  we have*

$$s(v'[\alpha h + a]) - s(v[h]) \geq d(v).$$

2. *If  $r(v, v') = 1/\alpha$  then for every  $h \in [0, c(v) - 1]$  and every  $a \in [0, \alpha - 1]$  we have*

$$s(v'[h]) - s(v[\alpha h + a]) \geq d(v).$$

**Proof 2.1** *Since each instance of  $v$  is a predecessor of  $\alpha$  instances of  $v'$  (Case 1), the execution of  $v'[j]$  must occur after the termination of at least  $\lceil j/\alpha \rceil$  instances of  $v$ . By definition, this is not earlier than the termination of the first  $\lceil j/\alpha \rceil$  instances of  $v$  to occur in schedule  $s$ . A similar argument holds for Case 2.*

■

**Theorem 2.1 (Lexicographic Ordering)** *Every feasible schedule  $s$  can be transformed into a latency-equivalent schedule  $s'$  compatible with the lexicographic order  $\ll$ .*

**Proof 2.2** *Let  $\omega_s$  be an actor-preserving permutation on  $U$  defined as  $\omega_s(v_h) = v[h]$ . In other words,  $\omega_s$  maps the task in position  $h$  according to  $\ll$  to the task occupying that position according to  $\prec^s$ . The new*

deployment is defined as

$$\mu'(v_h) = \mu(\omega_s(v_h)) \quad \text{and} \quad s'(v_h) = s(\omega_s(v_h)).$$

Permuting tasks of the same duration does not influence latency nor the satisfaction of resource constraints. All that remains to show is that  $s'$  satisfies precedence constraints. Each  $v_h$  is mapped into  $v[h]$  and each of its  $v'$  sons (resp. parents) is mapped into  $v'[\alpha h + a]$ ,  $0 \leq a \leq \alpha - 1$ . Hence a precedence constraints for  $s'$  of the form

$$s'(v_{h \cdot a}) - s'(v_h) \geq d(v)$$

is equivalent to

$$s(v[\alpha h + a]) - s(v[h]) \geq d(v)$$

which holds by virtue of Lemma 2.1 and the feasibility of  $s$ .  $\blacksquare$

The implication of this result, which holds for any prefix-consistent order, is that we can introduce additional lexicographic constraints to the formulation of the scheduling problem without losing optimality and thus significantly reduce the search space.

### 3 Constraint-Based Feasible Cost-Space Exploration

In this section we formulate multi-core deployment for split-join graphs as a set of constraints in the theory of linear arithmetics. Expressing scheduling problems using constraints is fairly standard [1, 2, 28, 15] and various formulations may differ in the assumptions they make about the application and the architecture and the aspects of the problem they choose to capture.

We target shared-memory multi-core architectures such as [16, 26, 11]. To avoid data copying overhead, such architectures provide groups of multiple processors – so-called *clusters* – with an instantaneous access to a dedicated shared memory. We assume that the whole application fits into the local memory of such a cluster and hence the communication between the tasks is modeled as instantaneous. The access to the shared memory (including contentions) is taken into account in the task execution times  $\delta$ . We assign a buffer to each edge (channel)  $(v, v')$  of the split-join graph so that tasks associated with the same actor read from and write to the same buffers.

To take buffer storage into account, we enrich the split-join graph model to become  $G = (V, E, d, w, r)$  with  $w(v, v')$  assigning to any edge in  $E$  the amount of data (in bytes) communicated from an instance of  $v$  to an instance of  $v'$  (this is called *token size* in the SDF literature). The corresponding task graph is  $T = (U, \mathcal{E}, \delta, w^\uparrow, w^\downarrow)$  where  $w_{v,v'}^\uparrow(u)$  is the amount of data *written* by  $u$  on the channel and  $w_{v,v'}^\downarrow(u')$  is the amount *read* by  $u'$ , where  $u \in U_v$  and  $u' \in U_{v'}$ . We assume that  $u$  allocates this memory space while starting and that  $u'$  releases it upon termination. The relation between  $w$ ,  $w^\uparrow$  and  $w^\downarrow$  depends on the type of the edge: for a split edge  $w_{v,v'}^\uparrow(u) = \alpha w(v, v')$  and  $w_{v,v'}^\downarrow(u') = w(v, v')$  while for join edges it is the other way around.

In the following we write down the constraints that define a feasible schedule and its cost in terms of latency, number of processors and buffer size.

- **Completion time and precedence:**  $e(u)$  is the time when task  $u$  terminates and a task cannot start before its predecessors terminate.

$$\bigwedge_{u \in U} e(u) = s(u) + \delta(u) \quad \wedge \quad \bigwedge_{(u, u') \in \mathcal{E}} e(u) \leq s(u')$$

- **Mutual exclusion:** tasks running on same processor should not overlap in time.

$$\bigwedge_{u \neq u' \in U} (\mu(u) = \mu(u')) \Rightarrow (e(u) \leq s(u') \vee (e(u') \leq s(u)))$$

- **Buffer:** these constraints compute the buffer size of every channel  $(v, v') \in E$ . They are based on the observation that buffer utilization is piecewise-constant over time, with jumps occurring upon initiation of writers and termination of readers. Hence the peak value of memory utilization can be found among one out of finitely-many starting points.

The first constraint defines  $W_{v,v'}^\uparrow(u, u_*)$ , the contribution of writer  $u \in U_v$  to the filling of buffer  $(v, v')$  observed at the start of a writer  $u_* \in U_v$ :

$$\bigwedge_{(v,v') \in E} \bigwedge_{u \in U_v} \bigwedge_{u_* \in U_v} (s(u) > s(u_*)) \wedge (W_{v,v'}^\uparrow(u, u_*) = 0) \vee (s(u) \leq s(u_*)) \wedge (W_{v,v'}^\uparrow(u, u_*) = w_{v,v'}^\uparrow(u))$$

Likewise the value  $W_{v,v'}^\downarrow(u', u_*)$  is the (negative) contribution of reader  $u' \in U_{v'}$  to buffer  $(v, v')$  observed at the start of writer  $u_*$ :

$$\bigwedge_{(v,v') \in E} \bigwedge_{u' \in U_{v'}} \bigwedge_{u_* \in U_v} ((e(u') > s(u_*)) \wedge (W_{v,v'}^\downarrow(u', u_*) = 0)) \vee ((e(u') \leq s(u_*)) \wedge (W_{v,v'}^\downarrow(u', u_*) = w_{v,v'}^\downarrow(u')))$$

The total amount of data in buffer  $(v, v')$  at the start of task  $u_* \in U_v$ , denoted by  $R_{v,v'}(u_*)$ , is the sum of contributions of all readers and writers already executed:

$$\bigwedge_{(v,v') \in E} \bigwedge_{u_* \in U_v} R_{v,v'}(u_*) = \sum_{u \in U_v} W_{v,v'}^\uparrow(u, u_*) - \sum_{u' \in U_{v'}} W_{v,v'}^\downarrow(u', u_*)$$

The buffer size for  $(v, v')$ , denoted by  $B_{v,v'}$  is the maximum over all the start times of tasks in  $U_v$ :

$$\bigwedge_{(v,v') \in E} \bigwedge_{u_* \in U_v} R_{v,v'}(u_*) \leq B_{v,v'}$$

- **Costs:** The following constraints define the cost vector associated with a given deployment, which is  $C = (C_l, C_n, C_b)$ , where the costs indicate, respectively, *latency* (termination of last task), number of *processors* used and total *buffer size*.

**Latency**  $C_l$  is the maximal completion time, defined by:

$$\bigwedge_{u \in U} e(u) \leq C_l.$$

**Mapping Cost**  $C_n$  defines the number of processors required for the schedule. The constraints impose processor utilization which is monotonic in the processor *index* (a processor  $m$  is used only if  $m - 1$  is used) and hence  $C_n$  is the largest processor index used:

$$\bigwedge_{m \in M} (\bigvee_{u \in U} \mu(u) = m) \Rightarrow ((\bigwedge_{m' < m} \bigvee_{u \in U} \mu(u) = m') \wedge C_n \geq m)$$

**Buffer Size Cost**  $C_b$  defines the total buffer usage of the schedule.

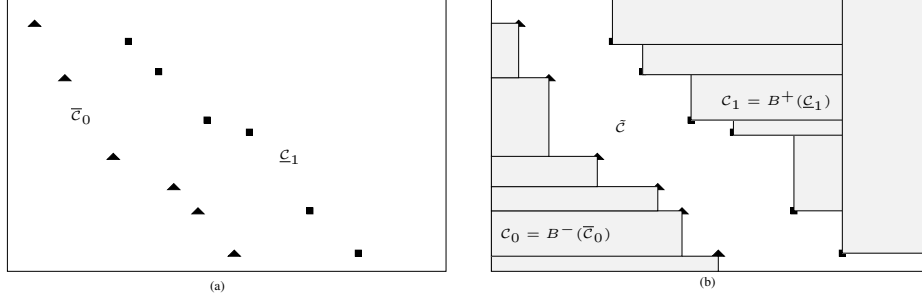
$$C_b = \sum_{(v,v') \in E} B_{v,v'}$$

In Appendix B we give an example of encoding all the constraints for the Z3 solver.

We refer to the totality of these constraints as  $\varphi(\mu, s, C)$  which are satisfied by any feasible deployment  $(\mu, s)$  whose cost is  $C$ . We sometimes add to them two types of symmetry-breaking constraints that do not change optimal costs: the lexicographic task ordering constraints described previously (henceforth: task symmetry) and fairly standard constraints to exploit processor symmetry: processor 1 runs task 1, processor 2 runs the lowest index task not running on processor 2, etc.

SAT and SMT solvers were designed for deciding satisfiability, not for optimization. However, such solvers can be used to find optimal costs by submitting queries concerning the existence of solutions with





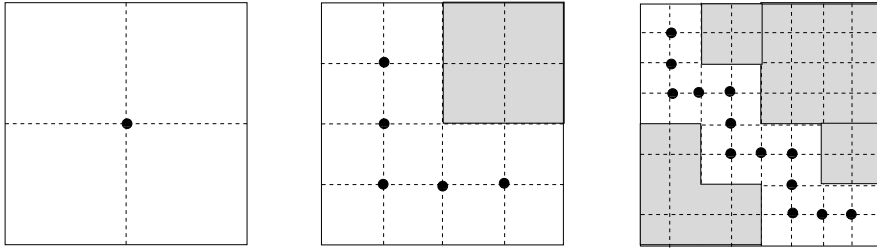
**Figure 2:** (a) Sets  $\mathcal{C}_0$  (**unsat**) and  $\mathcal{C}_1$  (**sat**) represented by their extremal points  $\bar{c}_0$  and  $\underline{c}_1$ ; (b) The state of our knowledge at this point as captured by  $\mathcal{C}_0$  (infeasible costs)  $\mathcal{C}_1$  (feasible costs) and  $\tilde{C}$  (unknown). The actual Pareto front is contained in the closure of  $\tilde{C}$ .

specific costs, which can be viewed as a binary search in the cost space with the solver acting as an oracle. We focus on *multi-criteria* optimization problems where we seek to find optimal trade-offs between latency  $C_l$ , number of processors  $C_n$  and buffer storage  $C_b$ . Such problems [8] do not admit a unique optimal solution but rather a set of *efficient* Pareto solutions [18] that cannot be improved in one cost dimension without being worsened in others. The set of such solutions, known as the *Pareto front*, represents the optimal trade-offs between the conflicting criteria. Following [14] we use queries to an SMT solver to find an approximation of the Pareto front. We summarize below the essence of the exploration methodology of [14], which can be viewed as a multi-dimensional generalization of binary search. Other approaches for multi-criteria optimization can be found in [8, 29, 6].

Let  $Q(c)$  be a shorthand for the satisfiability query  $\exists\mu\exists s\exists C$  s.t.  $\varphi(\mu, s, C) \wedge C \leq c$  which asks whether there is a feasible deployment whose cost vector is smaller or equal to  $C$ . If the solver answers affirmatively with some cost  $c$  we have a solution and may also conclude any cost in *forward cone* of  $c$  set  $B^+(c) = \{c' \mid c' \geq c\}$  is feasible. If the answer is negative we can conclude that any cost in the *backward cone*  $B^-(c) = \{c' \mid c' \leq c\}$  is infeasible. After submitting any number of queries with different values of  $c$  we face a situation illustrated in Fig. 2. The sets  $\bar{C}_0$  and  $\underline{C}_1$  are, respectively, the maximal costs known to be infeasible (**unsat**) and minimal costs found (**sat**). Sets  $\mathcal{C}_0$  and  $\mathcal{C}_1$  are defined as the sets of all points known to be **unsat** and **sat**, they are equal to the forward/backward cone of the extremal points. The feasibility of costs which are outside  $\mathcal{C}_0 \cup \mathcal{C}_1$  is unknown. The set  $\underline{C}_1$  constitutes an approximation of the Pareto front and its quality, defined as a kind of Hausdorff distance to the actual front, is bounded by its distance to the boundary of the backward cone of  $\bar{C}_0$ .

Before we apply the exploration procedure we need to bound the cost space. For latency  $C_l$ , a lower bound is the size of the the longest path (in terms of  $\delta$ ) through the task graph. The upper bound is the total amount of work (sum of  $\delta$  over all tasks). The bounds on buffers size are obtained by the two extreme scenarios. The lower bound is when each buffer is filled by the writer(s) to the minimal level required by the reader(s) to execute, that is,  $B_{v,v'} = \alpha w(v, v')$  for an edge with multiplicity  $\alpha$  or  $1/\alpha$ . The upper bound should cover the execution of all instances of  $v$  before any instance of  $v'$ ,  $B_{v,v'} = w(v, v') \cdot \max(c(v), c(v'))$ . The number of processors ranges trivially between 1 and the maximal number of processors on the platform. The width of the task-graph, when smaller than the number of processors, can serve as a tighter upper-bound as it limits the number of tasks that can execute in parallel.

Unlike the distance-oriented algorithm of [14], we use here a simpler exploration algorithm based on grid refinement. At every stage of the algorithm we refine the grid defined on the cost space and ask  $Q(c)$ -queries with  $c$  ranging over those newly-added grid points which are outside  $\mathcal{C}_0 \cup \mathcal{C}_1$ . Note that not all these new points will necessarily be queried because each query increases the size of  $\mathcal{C}_0 \cup \mathcal{C}_1$  so as to include some of these points. The description so far was based on the assumption that all queries terminate. However it is well-known that as  $c$  gets closer to the boundary between **sat** and **unsat**, the computation time may grow prohibitively and the solver can get stuck. To tackle this problem we bound the time budget per query and when this bound is reached we abort the query and interpret the result as **unsat**. Choosing the appropriate value for this time-out bound is a matter of trial and error.

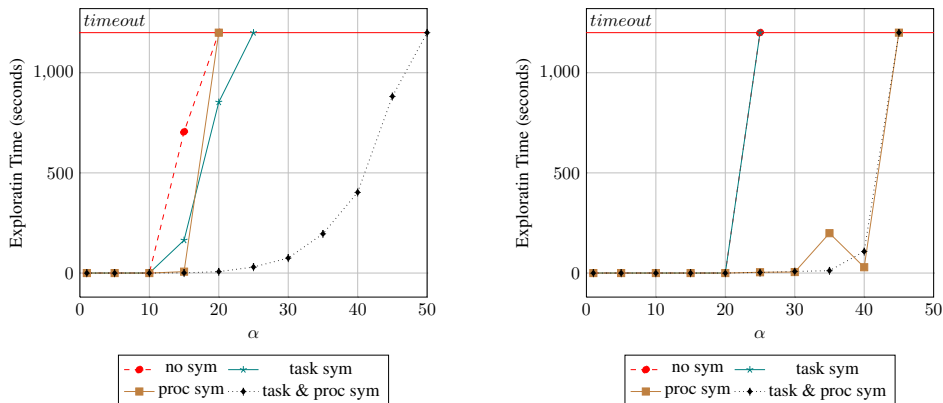


**Figure 3:** Exploring the cost space via grid refinement. The dark points indicate the new candidates for exploration after each refinement.

## 4 Experiments

In this section we investigate the performance of the cost-space exploration algorithm. First, we assess the contribution of the symmetry reduction constraints on the execution time and the quality of solutions for a synthetic example. Then we explore the cost space for a split-join graph derived from a real video application. These experiments use version 4.1 of the Z3 Solver [17] running on a Linux machine with *Intel Core i7* processor at 1.73 GHz with 4 GB of memory. Finally, we validate the model used to derive the solution by deploying a JPEG decoder on the Tiler platform [26] according to the derived schedule. The measured performance is very close to the predicted model.

**Finding Optimal Latency:** We use the split-join graph of Fig. 1 with various values of  $\alpha$  to explore the effect of the symmetry reduction constraints on the performance of the solver. We start with a single cost version of the problem and perform binary search to find the minimum latency that can be achieved for a fixed number of processors. We solve the same problem using four variations of the constraints: without symmetry reduction, with processor symmetry, with task symmetry and with both. Figure 4 depicts the computation times for finding the optimal latency as a function of  $\alpha$  on platforms with 5 and 20 processors. We use time-out per query of 20 minutes, which is much larger than the one minute we typically use because we want to find the *exact* optimum in order to compare the effects of different symmetry constraints. Scheduling problems are known to be easy when the number of processors approaches the number of tasks. For the difficult case of 5 processors, task symmetry starts dominating beyond 10 tasks and the combination of both gives the best results. It increases the size of graphs whose optimal latency can be found (with no query executing more than 20 minutes) from  $\alpha = 12$  to  $\alpha = 48$ . Not surprisingly, for 20 processors, the relative importance of processor symmetry grows.



**Figure 4:** Time to find optimal latency as a function of the number of tasks for 5 and 20 processors.

**Processor-Latency Trade-offs:** To demonstrate the effect of symmetry reductions on the Pareto front exploration we fix  $\alpha = 30$  and seek trade-offs between latency and the number of processors. We use a time budget of one minute per query. Fig. 5 depicts the results obtained with and without symmetry

breaking constraints. The square points show the **unsat** points whereas the circle are the **sat** points. The black curve is the approximation of the Pareto front, connecting all the minimal **sat** points. Points whose queries took long time to answer are surrounded by a dark halo whose intensity is proportional to the time (the darkest areas are around the **timeout** points). As one can see from the figure, symmetry constraints reduce significantly the number of time-outs with processor symmetry doing the job on the upper-left part of the curve while task symmetry is useful around the middle. The total time to find the minimal latency for each and every value of  $C_n$  is 42 minutes without symmetry, and 16 minutes with both types of symmetry constraints.

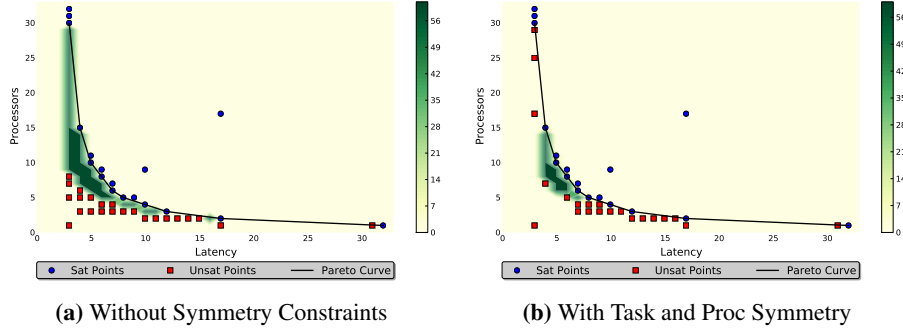


Figure 5: Pareto Exploration Result

**Video Decoder:** Next we perform a 3-dimensional cost exploration for a model of a video decoder taken from [10] and described in more detail in the appendix A.1. The application admits 11 actors expanding to 122 tasks. Without any symmetry constraints the solver quickly times out for most queries of interest. Symmetry constraints do not completely eliminate time-outs but reduce them significantly and quality of the Pareto front is much better, as shown in Fig. 6. Note that for a sequential implementation ( $C_n = 1$ ) the constraints improve buffer size from 276 to 182 and for the most parallel deployment ( $C_n = 122$ ) they reduce latency from 10 to 7 and buffer size from 333 to 229. The Pareto point (14, 333, 62) found without symmetry reduction is strongly dominated by the point (10, 229, 31) found with symmetry breaking. This solution improves the latency and buffer usage by roughly a third while using half of the processors. We believe it is a promising indication of the applicability of our approach and of the potential performance gains in treating the optimization problem globally.

**Jpeg Decoder:** Finally we validate our model by deploying a JPEG decoder on the Tiler platform [26] which is a 64-core symmetric multi-core platform running at 862.5 MHz. The theoretical scheduling problem that we solve is *deterministic* where task durations are assumed to be precisely known. The obtained schedule is time-triggered, given in terms of the *exact* start time function  $s$ . In reality, there are variations in execution times and in our implementation we use a static order schedules, preserving only the *order* of task execution on each processor. This is a common way to generate schedules for task graphs and SDF, see for example [20]. When task durations agree with the nominal values used in the optimization problem, this scheduling policy coincides with  $s$ . Unlike the traditional work on dataflow mapping, we support mappings where the writers and readers of the same buffer storage can be spread over more than two different processors. Our experience confirms that this dynamic scheduling policy can be implemented with a reasonable amount of additional synchronization between the cores. Note also that when the schedule is compatible with lexicographical task order (justified by Theorem 2.1), the accesses to the channels automatically become FIFO and this facilitates the implementation of cyclic buffers.

The split-join graph of the decoder is shown in the appendix A.2. It has three main actors: variable length decoding (*vld*), inverse quantization and inverse discrete cosine transform (*iq/idct*) combined and color conversion (*color*). To measure execution times we run the decoder several times on a single processor and measure the execution time of each actor. To mitigate cache effects, we consider the average execution time rather than worst case, which occurs only in the first execution due to cache misses. We use these average execution time in the model we submit to the solver. We then deploy the decoder on the platform and run it 100 times (again to dampen cache effects). The relation between the average latency (in  $\mu s$ ) observed experimentally and the Pareto points computed by the SMT solver is depicted below and

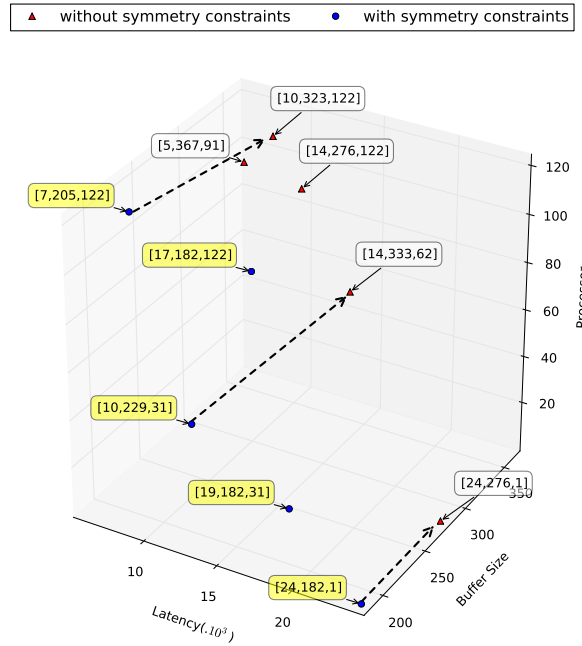


Figure 6: Video Decoder Pareto Points

the deviation is typically smaller than 15%.

no. proc	1	3	4	6	8	12
predicted	506	314	278	261	243	226
measured	461	309	299	307	300	351

## 5 Discussion

The deployment of programs on parallel machine is a very old problem whose parameters change with the evolution of computer architecture. The problem exists in both software [12] and hardware [5] and in the latter it is viewed as an instance of *high-level synthesis*. Due to problem complexity the problem is often solved using heuristics such as list scheduling and/or decomposed into separate phases, for example, optimizing latency and buffers separately [21]. Recent advances in SAT and SMT solvers and other constraint propagation techniques suggests an opportunity to formulate and solve the problem in a monolithic way, avoiding the sub-optimality of decomposed solutions. For example, [15] exploit SMT solvers to combine multiple deployment sub-problems: the task-to-processor assignment, the ordering of tasks on each processor and the assignment of scalable voltage per processor. For SDF graphs, [2] and [28] combine multiple phases using a *constraint programming* (CP) engine. In the context of high-level synthesis, the tool FACTS (see [9] for references) uses branch-and-bound approach combined with constraint analysis, whereas [5] discusses various ILP formulations. In [27] a quantitative model checking engine is developed using a variant of timed automata for combined scheduling and buffer storage optimization of SDF graphs.

Various approaches to facilitate the task of the solver by additional symmetry breaking constraints have been tried, for example [19] for graph coloring or an automated method for discovering graph automorphism [4] which can lead to significant improvements [9]. However, our deployment problem does not require complex detection of isomorphic subgraphs. Instead we exploit the knowledge about the structure of the task graphs coming from the original split-join graph. Theorem 1 provides the necessary compact symmetry breaking constraints that do the job. As for the restrictions that we imposed on the split-join graph compared to more general SDF graphs admitting non-divisible token production and consumption rate, let us first remark that Theorem 1 can be extended, somewhat less elegantly, to this more general

case. Moreover, the extensive study of StreamIT benchmarks found in [24] reports that most actors in most applications, fall into the category of well-formed split-join graphs that we treat.

The contribution of the paper can be summarized as follows. We provide a framework for multi-criteria optimization and cost-space exploration, not based on heuristic sub-optimal decomposition. Using symmetry reduction justified by Theorem 1, we could conduct a 3-dimensional cost-space exploration for a non-trivial problem with 122 tasks. The theorem itself generalizes the result of [9] which proves optimality of lexicographic order for one level of nesting. We prove the result for arbitrary nesting depth and give a simpler proof.

In future, we plan to extend this work in several directions. First we will employ more refined models of data communication where different mappings imply different data transfer costs. Secondly we will consider *pipelined* executions as was done in [15, 2, 28, 27], using *e.g.*, a finite unfolding. This will increase the number of tasks but will reduce the effect of precedence constraints. Thirdly we should adapt the methodology to a more significant variability in task duration and this will require an implementation of scheduling under uncertainty that can deviate from the task execution order provided by *s*. Finally we will seek ways for a more direct exploitation of the symbolic representation of data-parallel tasks and a tighter interaction between the cost exploration algorithm and the solver.

## References

- [1] Baptiste, P., Le Pape, C., Nuijten, W.: Constraint-Based Scheduling. Kluwer international series in engineering and computer science: VLSI, computer architecture, and digital signal processing, Kluwer (2001) 3
- [2] Bonfietti, A., Benini, L., Lombardi, M., Milano, M.: An efficient and complete approach for throughput-maximal sdf allocation and scheduling on multi-core platforms. In: Design, Automation Test in Europe Conference Exhibition (DATE), 2010. pp. 897–902 (2010) 3, 5
- [3] Coffman, E.G.: Computer and job-shop scheduling theory. Wiley (1976) 1
- [4] Darga, P.T., Sakallah, K.A., Markov, I.L.: Faster symmetry discovery using sparsity of symmetries. In: Proceedings of the 45th annual Design Automation Conference. pp. 149–154. DAC '08, ACM, New York, NY, USA (2008), <http://doi.acm.org/10.1145/1391469.1391509> 5
- [5] De Micheli, G.: Synthesis and optimization of digital circuits. Electrical and Computer Engineering Series, McGraw-Hill Higher Education (1994), <http://books.google.fr/books?id=FESDQgAACAAJ> 5
- [6] Deb, K.: Multi-Objective Optimization Using Evolutionary Algorithms. Wiley paperback series, Wiley (2009), <http://books.google.fr/books?id=U0dnPwAACAAJ> 3
- [7] Dutertre, B., Moura, L.: A fast linear-arithmetic solver for dpll(t). In: Ball, T., Jones, R. (eds.) Computer Aided Verification, Lecture Notes in Computer Science, vol. 4144, pp. 81–94. Springer Berlin Heidelberg (2006), [http://dx.doi.org/10.1007/11817963\\_11](http://dx.doi.org/10.1007/11817963_11) 1
- [8] Ehrgott, M.: Multicriteria Optimization. Springer Berlin · Heidelberg (2005), [http://books.google.fr/books?id=AwRjo6iP4\\_UC](http://books.google.fr/books?id=AwRjo6iP4_UC) 3
- [9] van Eijk, C.A.J., Jacobs, E.T.A.F., Mesman, B., Timmer, A.H.: Identification and exploitation of symmetries in dsp algorithms. In: Proceedings of the conference on Design, automation and test in Europe. DATE '99, ACM, New York, NY, USA (1999), <http://doi.acm.org/10.1145/307418.307572> 5
- [10] Fradet, P., Girault, A., Poplavko, P.: Spdf: A schedulable parametric data-flow moc. In: Design, Automation Test in Europe Conference Exhibition (DATE), 2012. pp. 769–774 (2012) 4
- [11] Kalray: Kalray MPPA 256, <http://www.kalray.eu/> 3

- [12] Kwok, Y.K., Ahmad, I.: Static scheduling algorithms for allocating directed task graphs to multiprocessors. *ACM Comput. Surv.* 31(4), 406–471 (Dec 1999), <http://doi.acm.org/10.1145/344588.344618> 5
- [13] Lee, E., Messerschmitt, D.: Synchronous data flow. *IEEE* 75(9), 1235 – 1245 (1987) 1
- [14] Legriel, J., Guernic, C., Cotton, S., Maler, O.: Approximating the pareto front of multi-criteria optimization problems. In: Esparza, J., Majumdar, R. (eds.) *Tools and Algorithms for the Construction and Analysis of Systems*, *Lecture Notes in Computer Science*, vol. 6015, pp. 69–83. Springer Berlin Heidelberg (2010), [http://dx.doi.org/10.1007/978-3-642-12002-2\\_6](http://dx.doi.org/10.1007/978-3-642-12002-2_6) 1, 3, 3
- [15] Legriel, J., Maler, O.: Meeting deadlines cheaply. In: *ECRTS*. pp. 185–194. *IEEE* (2011) 1, 3, 5
- [16] Melpignano, D., Benini, L., Flamand, E., Jegou, B., Lepley, T., Haugou, G., Clermidy, F., Dutoit, D.: Platform 2012, a many-core computing accelerator for embedded socs: performance evaluation of visual analytics applications. In: *Proceedings of the 49th Annual Design Automation Conference*. pp. 1137–1142. *DAC '12*, ACM, New York, USA (2012), <http://doi.acm.org/10.1145/2228360.2228568> 3
- [17] Moura, L., Bjorner, N.: Z3: An efficient smt solver. In: Ramakrishnan, C., Rehof, J. (eds.) *Tools and Algorithms for the Construction and Analysis of Systems*, *Lecture Notes in Computer Science*, vol. 4963, pp. 337–340. Springer (2008), [http://dx.doi.org/10.1007/978-3-540-78800-3\\_24](http://dx.doi.org/10.1007/978-3-540-78800-3_24) 1, 4
- [18] Pareto, V.: Manuel d'économie politique. *Bull. Amer. Math. Soc.* 18, 462–474 (1912) 3
- [19] Ramani, A., Aloul, F., Markov, I., Sakallah, K.: Breaking instance-independent symmetries in exact graph coloring. In: *DATE*. vol. 1, pp. 324–329 Vol.1 (2004) 5
- [20] Sriram, S., Bhattacharyya, S.: *Embedded Multiprocessors: Scheduling and Synchronization*, Second Edition. *Signal Processing and Communications*, Taylor & Francis (2009), <http://books.google.fr/books?id=v13bnBCKJLEC> 1, 4
- [21] Stuijk, S., Geilen, M., Basten, T.: Exploring trade-offs in buffer requirements and throughput constraints for synchronous dataflow graphs. In: *43rd annual Design Automation Conference*. pp. 899–904. *DAC '06*, ACM, New York, NY, USA (2006), <http://doi.acm.org/10.1145/1146909.1147138> 5
- [22] Tendulkar, P.: Z3 Sample Constraints. <http://www-verimag.imag.fr/~tendulka/files/formatsPaperExampleConstraints.z3> (2013) B
- [23] Tendulkar, P., Poplavko, P., Maler, O.: Symmetry breaking for multi-criteria mapping and scheduling on multicores. In: *Formal Modelling and Analysis of Timed Systems (FORMATS)*, 2013 (2013) 1
- [24] Thies, W., Amarasinghe, S.: An empirical characterization of stream programs and its implications for language and compiler design. In: *international conference on Parallel architectures and compilation techniques*. pp. 365–376. *PACT '10*, ACM, NY, USA (2010), <http://doi.acm.org/10.1145/1854273.1854319> 5
- [25] Thies, W., Karczmarek, M., Amarasinghe, S.: Streamit: A language for streaming applications. In: Horspool, R. (ed.) *Compiler Construction*, *Lecture Notes in Computer Science*, vol. 2304, pp. 179–196. Springer (2002), [http://dx.doi.org/10.1007/3-540-45937-5\\_14](http://dx.doi.org/10.1007/3-540-45937-5_14) 1
- [26] Tilera, LTD: Tilera TILE64 processor, <http://www.tilera.com/> 3, 4, 4
- [27] Yang, Y., Geilen, M., Basten, T., Stuijk, S., Corporaal, H.: Exploring trade-offs between performance and resource requirements for synchronous dataflow graphs. In: *ESTImedia*. pp. 96–105 (2009) 5

- [28] Zhu, J., Sander, I., Jantsch, A.: Buffer minimization of real-time streaming applications scheduling on hybrid cpu/fpga architectures. pp. 1506–1511. DATE '09, European Design and Automation Association, Leuven, Belgium (2009), <http://dl.acm.org/citation.cfm?id=1874620.1874980> 3, 5
- [29] Zitzler, E., Thiele, L.: Multiobjective evolutionary algorithms: A comparative case study and the strength Pareto approach. IEEE transactions on Evolutionary Computation 3(4), 257–271 (1999) 3

## A Application Example Graphs

### A.1 Video Decoder

In our video decoder programming model, a complete video frame can be processed by a repeated execution of the split-join graph shown in Fig. 7. The *vld* actor parses the input bitstream, extracting the subsequent macroblocks. The parameter  $x$  selects the number of macroblocks processed per one graph execution. The larger  $x$  the more macroblocks can be processed in parallel, but the more difficult it is to generate an optimal schedule. The actors indexed by ‘L’ process the 4 luminance blocks per macroblock and the actors indexed by ‘C’ process the 2 chrominance blocks. The *Color* converts the frame into the RGB format. The *Fetch* actors fetch the reference blocks from the previous frame for motion compensation, *MC*. *Upscale* scales 2 chrominance blocks into 8. The weights  $w$  of the channels (in blocks) are depicted in parentheses in the figure. We perform the exploration for  $x = 5$  which yields a task-graph with 122 tasks.

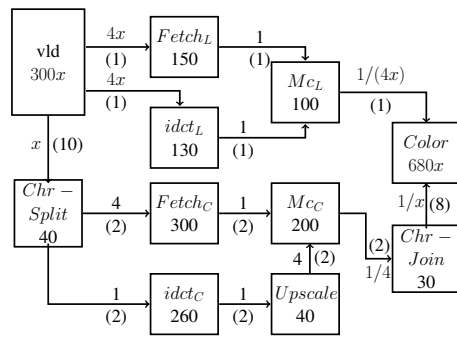
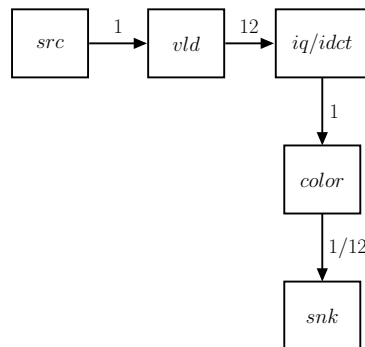


Figure 7: Video Decoder Example

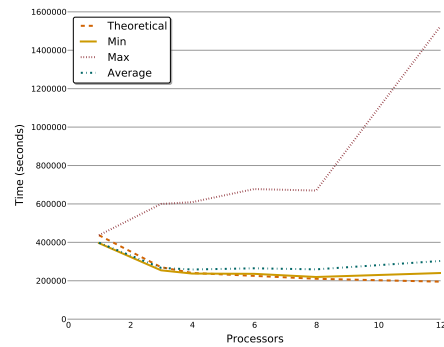
Actor	Execution Time
<i>vld</i>	$300 \cdot x$
<i>Fetch<sub>L</sub></i>	150
<i>idct<sub>L</sub></i>	130
<i>Mc<sub>L</sub></i>	100
<i>Color</i>	$680 \cdot x$
<i>Chr - Split</i>	40
<i>Fetch<sub>C</sub></i>	300
<i>idct<sub>C</sub></i>	260
<i>Upscale</i>	40
<i>Mc<sub>C</sub></i>	200
<i>Chr - Join</i>	30
<i>Color</i>	$680 \cdot x$

Table 1: Execution Times

### A.2 Jpeg Decoder



(a) Jpeg Decoder Example

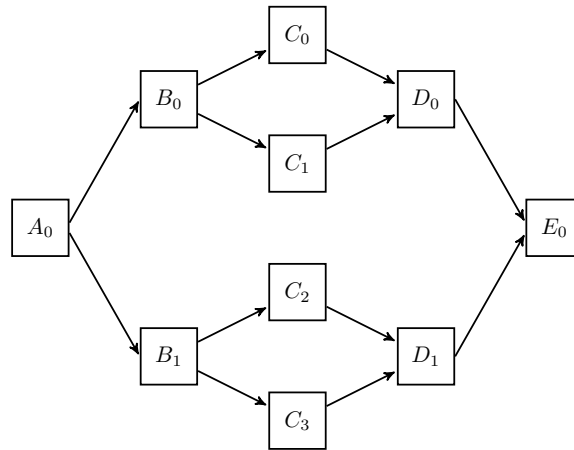


(b) Execution time on Tileria Platform

Figure 8: Jpeg Decoder



## B Z3 SAT Solver Constraints



**Figure 9:** Task Graph Example to demonstrate Z3 Constraints

The Z3 Constraints for the task graph in figure 9 are available at [22].