

Optimal FPGA Module Placement with Temporal Precedence Constraints

Sándor P. Fekete*

Ekkehard Köhler†

Jürgen Teich‡

Abstract

We consider the optimal placement of hardware modules in space and time for FPGA architectures with re-configuration capabilities, where modules are modeled as three-dimensional boxes in space and time. Using a graph-theoretic characterization of feasible packings, we are able to solve the following problems:

- Find the minimal execution time of the given problem on an FPGA of fixed size,
- Find the FPGA of minimal size to accomplish the tasks within a fixed time limit.

Furthermore, our approach is perfectly suited for the treatment of precedence constraints for the sequence of tasks, which are present in virtually all practical instances. Additional mathematical structures are developed that lead to a powerful framework for computing optimal solutions. The usefulness is illustrated by computational results.

1 Introduction

A Field-Programmable Gate Array (FPGA) typically consists of a regular rectangular grid of equal configurable cells (logic blocks) that allow the prototyping of simple logic functions together with simple registers and with special routing resources (see Figure 1). A particular design is realized by customizing a configuration: In traditional SRAM-based chips, this can be done at power-up by loading a configuration bit-stream serially into the chip. These chips may only be reconfigured as a whole with typical re-configuration times ranging in the order of milliseconds.

Today, new generations of FPGAs have become partitionable and dynamically reconfigurable, even partially. These chips (see e.g. [1, 24]) may support several independent or interdependent tasks and designs at a time, and parts of the chip can be reconfigured quickly during run-time.

For a start, the reader may think of architectures similar to the Xilinx 6200 FPGA [24] architecture, where col-

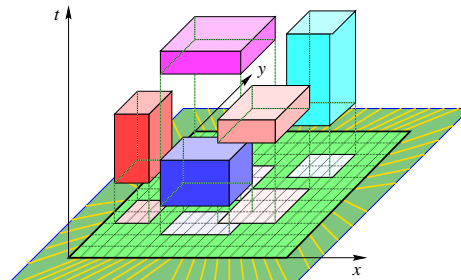


Figure 1. An FPGA and a set of five modules (tasks), shown in ordinary two-dimensional space and in three-dimensional space-time. Modules must be placed inside the chip and must not overlap if executed simultaneously on the chip.

umn read-ins and read-outs of flip-flop contents may be performed during run-time without interfering with other configured parts of the chip. (In Section 3, we describe how to model additional overhead for devices such as a Xilinx Virtex.) Under these assumptions, a task, or module, may be represented by a cuboid, with two spatial dimensions, and one representing the time of computation, see Fig. 2.

However, even if the configuration time is short, the compilation time for constructing the configuration stream for a task is still rather long. This diminishes the results that have been reported recently on on-line strategies for compiling and reconfiguring such devices. Important examples include speeding up computational problems in hardware by task compaction on hypercubes [16], or approaches to dynamic allocation of a sequence of tasks on an FPGA of given size by using heuristics to compact tasks in execution on the chip during run-time [3, 4].

Here we consider *statically defined problems where a task set is given*. In previous work reported in [22, 23] we have described how these problems can be understood by virtue of an easy graph-theoretic characterization of feasible packings. In this paper, we show how to extend this approach in order to deal with very restrictive constraints that are present in virtually all practical instances: Typically, there are *temporal precedence constraints* imposed on the set of computing modules, since the output of one task may be needed as input for another task. Such a set of

*Department of Mathematics, TU Berlin, Berlin, Germany, email: fekete@math.tu-berlin.de.

†Department of Mathematics, TU Berlin, Berlin, Germany, email: ekoe@math.tu-berlin.de.

‡Computer Engineering Laboratory, University of Paderborn, Paderborn, Germany, email: teich@date.upb.de.

precedence constraints may be described by a dependency graph, see Figure 2. For a problem instance of this type, we are interested in finding exact solutions to the following problems. (In the following, a spatial placement is called *feasible* if the locations occupied by each pair of tasks that have overlapping execution intervals are disjoint and fit into the available space and time. In the presence of precedence constraints for the tasks, *feasibility of a schedule* implies that all of these constraints are met.)

- Find the chip of smallest size to accommodate all tasks such that a given maximum total execution-time is satisfied (*MinA&FindS*) together with a feasible schedule. A subproblem called *MinA&FixedS* arises when the precise starting times of all tasks are already given.
- Check whether for a chip of given size and given maximum execution time, there is a feasible placement and a feasible schedule that accommodates a set of tasks (*FeasAT&FindS*).
- Find the smallest execution time of the set of tasks for a chip of fixed size (*MinT&FindS*).
- Check whether a chip of given size and a given feasible schedule allow a feasible placement (*FeasA&FixedS*).

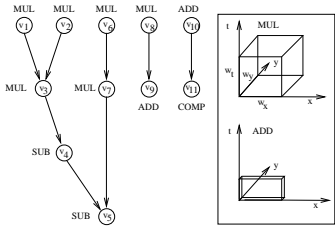


Figure 2. Dependency graph of tasks and shape of modules (3D boxes) with the spatial dimensions x and y and the temporal dimension t (execution time)

Our results reported in [23] deal with solving the above *FixedS*-problems. For these cases, the problems may be modeled by three-dimensional orthogonal packing problems (OPP) without precedence constraints. Note that already the one-dimensional counterparts without precedence constraints are known to be NP-complete in the strict sense [13]. Existing ILP formulations for similar problems such as [2] fail to solve technical problems of interesting size to optimality because they use a grid decomposition and model the placement of a module at location (x, y) and time t by a 0-1-variable, requiring $\mathcal{O}(|V||X||Y|)$ 0-1 variables and $\mathcal{O}(|X||Y||T|)$ constraints where $|X|$, $|Y|$, $|T|$ are the dimensions of the underlying grid in the x -, y -, and t -direction. The largest two-dimensional packing problems that have been solved with this technique place about 20 rectangles on a 30×30 grid [2, 15]. Solving a three-dimensional problem with about 100 nodes is hopeless if these standard solution techniques are used.

A breakthrough to solve these problems to optimality was due to the introduction of so-called *packing classes* [7, 21] that drastically reduce the search space for feasible packings. The application of this idea to *FixedS*-problems in the context of FPGA module reconfiguration has been elaborated in [22, 23]. Here, we develop a framework to solve also problems with precedence constraints such as given by data dependencies.

After introducing some basic mathematical terminology in Section 2, we describe in Section 3 how to extend our approach to a new line of FPGAs with limited reconfiguration capabilities. In Section 4, we sketch the mathematical concepts of packing classes and a solution to packing problems without precedence constraints, which are indispensable for understanding how to deal with precedence constraints. In Section 5, we introduce precedence constraints, describe the mathematical foundations for incorporating them into the search, and explain how to implement the resulting algorithms. Finally, we present computational results for two realistic benchmarks in Section 6.

2 Mathematical Modeling

Problem instances. We assume that a problem instance is given by a *set of tasks* V . Each task has a *spatial requirement* in the x - and y -direction, denoted by $w_x(v)$ and $w_y(v)$, and a *duration*, denoted by a size $w_t(v)$ along the time axis. The reconfigurable chip H consists of an array of $h_x \times h_y$ cells. In addition, there may be an overall allowable time h_t for all tasks to be completed. A *schedule* is given by a start time $p_t(v)$ for each task. A schedule is *feasible*, if all tasks can be carried out without overlap of computation tasks in time or space, such that all tasks are within spatial and temporal bounds.

Graphs. Some of our descriptions make use of a number of different graph classes. An (undirected) graph $G = (V, E)$ is given by a set of vertices V , and a set of edges E ; each edge describes the adjacency of a pair of vertices, and we write $\{u, w\}$ for an edge between vertices u and w . For a graph G , we obtain the *complement graph* \bar{G} by exchanging the set E of edges with the set \bar{E} of non-edges. In a directed graph $D = (V, A)$, edges are oriented, and we write (u, w) to denote an edge directed from u to w .

Precedence constraints. There may be a temporal precedence requirement between some of the tasks, since some tasks need to be finished before others can get started. Mathematically, this means that we are given a partial order on V , which can be described by a directed acyclic graph $D = (V, A)$, where A is the set of directed arcs. In the presence of such a partial order (denoted A), a feasible schedule also needs to satisfy these additional constraints.

Packing problems. In the following, we treat tasks as three-dimensional boxes and feasible schedules as arrangements of boxes that satisfy all side constraints. This is

implied by the term of a *feasible packing*. As described in the introduction, there are different types of objectives, corresponding to different types of packing problems. The *Orthogonal Packing Problem* (OPP) is to decide whether a given set of boxes can be placed within a given “container” of size $h_x \times h_y \times h_t$. The *Base Minimization Problem* (BMP) is to minimize the size h_x for a fixed h_t such that all boxes fit into a container $h_x \times h_x \times h_t$ with quadratic base. This corresponds to minimizing chip size to carry out a set of computations within a given time – called *MinA&FindS* in the introduction. The *Strip Packing Problem* (SPP) is to minimize the size h_t for a given base size $h_x \times h_y$, such that all boxes fit into the container $h_x \times h_y \times h_t$. This corresponds to minimizing the time to carry out a set of computations on a given chip – called *MinT&FindS* in the introduction.

3 Technical Aspects and Overheads

In previous work [23], we made the assumption that communication and reconfiguration can be described by fixed time periods, so they are not subject to optimization. Here we quantify these overheads for a recent Xilinx Virtex device, the XCV 1000. Therefore, we describe in more technical detail how our approach can also be applied to this kind of FPGAs. After computing the overheads for communication and reconfiguration for this device, Section 6.2 presents an example that takes these overheads into account.

Communication overhead. The Xilinx Virtex devices are column-oriented in the sense that reconfiguration and read-in/out of flip-flop contents may only be specified for a full column of CLBs of the chip. In Fig. 3, one column of CLBs (configurable logic blocks) is shown. Now, our model of communication assumes that a task finishing its execution communicates with another task (data-dependence) in the following way: All flip-flops of CLBs containing result values to be read by the direct successor task are read-out using a bus interface and stored to external memory. Subsequently, assuming the receiving module has been loaded at the location as specified by the packing, this module performs a read-in of the results of the sending module into its input flip-flops, not necessarily at adjacent CLB columns.

Unfortunately, for a Virtex device, only complete column read-outs and read-ins can be performed and the bus width is only 8 bit. For a Virtex XCV 1000 device, reading one column of flip-flops corresponds to addressing one frame containing 1248 bits. This takes approximately 180 clock cycles (156 for reading the column in chunks of 8 bits with the other, 24 cycles overhead for configuration of the interface). The same overhead is necessary to perform a write of a column of flip-flops specified by another frame.

Based on this overhead for reading/writing a column of flip-flops, we can bound the communication overhead from above by 2 times the overhead of 180 clock cycles time the

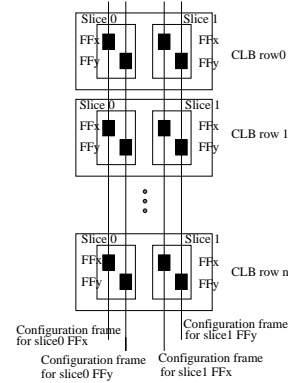


Figure 3. One-column slice of a Xilinx Virtex FPGA with CLBs containing two slices each

number of columns containing data that must be communicated. The problem of having to update a complete column of flip-flops during a read-in may be addressed by using a particular column of flip-flops (such as FFy in Slice 0) for read-ins only and designing the modules such that FFy flip-flops of Slice 0 are only used during the first cycle of module execution.

As we will see in an example (presented in Section 6.2), this overhead may be accounted for by adding the overhead to the execution time of a task.

Reconfiguration overhead. Consider reconfiguring a hardware task v occupying a rectangle of $w_x(v) \times w_y(v)$ CLBs. In case of a Virtex FPGA, we have to reconfigure $w_y(v)$ columns of CLBs in this case. Each CLB column consists of 48 frames to be configured, each responsible for configuring a certain type of resource. When using auto-increment of addresses within the 48 frames of a CLB column, each column to reconfigure requires $48 \times 39 \times 4$ (1248 bit = $39 * 32$ bit, 32 bit requires 4 cycles ($32 = 4 * 8$)) cycles leading to a complete reconfiguration overhead of 7512 cycles per column reconfiguration, for a total of $w_y(v) \times 7512$ clock cycles. Again we treat this reconfiguration overhead as a constant that we add to the execution time of a hardware task.

4 Solving Unconstrained Problems

4.1 A Framework for Optimal Solutions

Before discussing precedence constraints, we describe a number of fundamental mathematical insights and resulting computational methods for unconstrained packing problems. Mathematical details can be found in our previous papers [7, 8, 9, 10, 11, 21, 22, 23].

If we have an efficient method for solving OPPs, we can also solve BMPs and SPPs by using a binary search. However, deciding the existence of a feasible more-dimensional packing is a hard problem in higher dimensions, and proposed methods suggested by other authors [2, 15] have been of limited success.

Our framework uses a combination of different approaches to overcome these problems:

1. Try to disprove the existence of a packing by fast and good classes of lower bounds on the necessary size.
2. In case of failure, try to find a feasible packing by using fast heuristics.
3. If the existence of a packing is still unsettled, start an enumeration scheme in form of a branch-and-bound tree search.

By developing good new bounds for the first stage, we have been able to achieve a considerable reduction of the number of cases where a tree search needs to be performed. (Mathematical details for this step are described in [8, 10] and are omitted from this short paper.) However, it is clear that the efficiency of the third stage is crucial for the overall running time when considering difficult problems. Using a purely geometric enumeration scheme for this step by trying to build a partial arrangement of boxes is easily seen to be immensely time-consuming. In the following, we describe a purely combinatorial characterization of feasible packings that allows to perform this step more efficiently.

4.2 Packing Classes

If we consider a feasible packing in d -dimensional space, we can extract some partial information by considering the relative arrangement of coordinate intervals. More precisely, we can consider the projections of the boxes onto the three coordinate axes, and thus reduce the one d -dimensional arrangement to d one-dimensional ones. (See Figure 4 for an example in $d = 2$.) In a second step, we can disregard the exact coordinates of the resulting intervals in direction i and only consider the *component graph* $G_i = (V, E_i)$: Two boxes u and v are connected by an edge in G_i , iff they have overlapping i coordinates. Mathematically, a graph with this characterization of edges is called an *interval graph*. These graphs have been studied intensively in graph theory (see [14, 20]), and they have a number of very useful algorithmic properties.

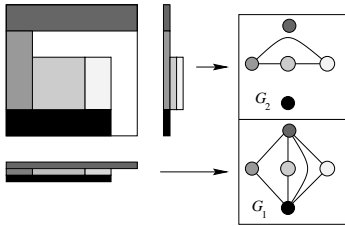


Figure 4. The projections of the boxes onto the coordinate axes define interval graphs (here in 2D: G_1 and G_2).

Considering sets of d component graphs G_i instead of complicated geometric arrangements has some clear advantages. (Algorithmic implications for our specific purposes

will be discussed further down.) It is not hard to check that the following three conditions must be satisfied by all d -tuples of graphs G_i that are constructed from a feasible packing:

C1: G_i is an interval graph, $\forall i \in \{1, \dots, d\}$.

C2: Any independent set S of G_i is i -admissible, $\forall i \in \{1, \dots, d\}$, i.e., $w_i(S) = \sum_{v \in S} w_i(v) \leq h_i$, since all boxes in S must fit into the container in the i th dimension.

C3: $\bigcap_{i=1}^d E_i = \emptyset$. In other words, there must be at least one dimension in which the corresponding boxes do not overlap.

A d -tuple of component graphs satisfying these necessary conditions is called a *packing class*. The remarkable property (proven in [21, 9]) is that these three conditions are also sufficient for the existence of a feasible packing.

Theorem 1 *A d -tuple of graphs $G_i = (V, E_i)$ corresponds to a feasible packing, iff it is a packing class, i. e., if it satisfies the conditions C1, C2, C3.*

This allows it to consider only packing classes in order to decide the existence of a feasible packing, and disregard most of the geometric information.

4.3 Solving OPPs

Our search procedure works on packing classes, i.e., triples of component graphs with the properties C1, C2, C3. Since each packing class represents not only a single packing but a whole family of equivalent packings, we are effectively dealing with more than one possible candidate for an optimal packing at a time. (The reader may check for the example in Figure 4 that there are 36 different feasible packings that correspond to the same packing class.)

The search tree is traversed by Depth First Search, see [11, 21] for details. Branching is done by fixing an edge $\{b, c\} \in E_i$ or $\{b, c\} \notin E_i$. After each branching step, it is checked if one of the three conditions C1, C2, C3 is violated, or whether a violation can only be avoided by fixing further edges. This is easy for two of the conditions: enforcing C3 is obvious; property C2 is hereditary, so adding edges to E_i later will keep it satisfied. (Note that computing maximum weighted cliques on comparability graphs can be done efficiently, see [14].) In order to ensure that property C1 is not violated, we use a number of graph-theoretic characterizations of interval graphs and comparability graphs. These characterizations are based on two forbidden substructures (again, see [14] for details). In particular, this means that the following configurations have to be avoided:

1. induced chordless cycles of length 4 in E_i ;
2. so-called 2-chordless odd cycles in the set $\overline{E_i}$ of edges excluded from E_i (see [11, 14] for details);
3. infeasible stable sets in E_i .

Each time we detect such a fixed subgraph, we can abandon the search on this node. Furthermore, if we detect a fixed subgraph, except for one unfixed edge, we can fix this edge, such that the forbidden subgraph is avoided.

Our experience shows that these conditions are already useful when only small subsets of edges have been fixed, since by excluding small sub-configurations, like induced chordless cycles of length 4, each branching step triggers a cascade of more fixed edges.

5 Solving Optimization Problems with Precedence Constraints

For most practical instances, we have to satisfy additional constraints for the temporal placement, i.e., for the start times of tasks. It should be stressed that for standard approaches, adding constraints makes the three-dimensional packing problems much harder. This is significantly different from our approach, where the nature of the data structures simplifies these problems from three-dimensional to purely two-dimensional ones: If the whole schedule is given, all edges E_t in one of the graphs are determined, so we only need to construct the edge sets E_x and E_y of the other two graphs¹. As we have worked out in detail in [22, 23], this allows it to solve the resulting FixedS-Problems quite efficiently.

A more realistic, but also more involved situation arises if only a set of precedence constraints is given, but not the full schedule. The following describes how to convert packing classes into ordered arrangements, and how to deal with order constraints.

5.1 Packing Classes and Interval Orders

Any edge in a graph G_i corresponds to an overlap between the corresponding intervals. This means that the complement graph \overline{G}_i given by the complement \overline{E}_i of the edge set E_i consists of all pairs of coordinate intervals that are “comparable”: Either the first interval is “to the left” of the second, or vice versa. Any (undirected) graph of this type is a so-called *comparability graph* [14]. By orienting edges to point from “left” to “right” intervals, we get a partial order of the set V of vertices, a so-called *interval order* [20]. Obviously, this order relation is transitive, i.e., $e \prec f$ and $f \prec g$ imply $e \prec g$, which is the reason why we also speak of a *transitive orientation* of the undirected comparability graph G_i . See Figure 5 for a (two-dimensional) example of a packing class, the corresponding comparability graph, a transitive orientation, and the packing corresponding to the transitive orientation.

Now consider a situation where we need to satisfy a partial order A of precedence constraints in the time dimension.

¹To emphasize the motivation of temporal precedence constraints, we write E_t to suggest that the time coordinate is constrained, and E_x and E_y to imply that the space coordinates are unrestricted. Clearly, our approach works the same way when dealing with spatial restrictions.

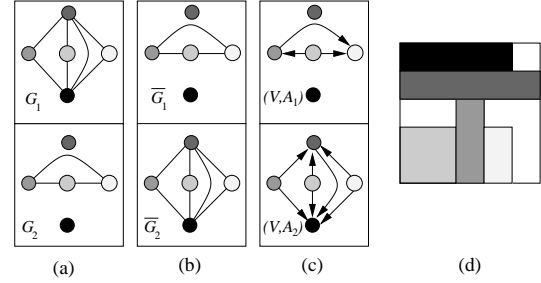


Figure 5. (a) A two-dimensional packing class. (b) The corresponding comparability graphs. (c) A transitive orientation. (d) A feasible packing corresponding to the orientation.

It follows that each arc $a = (u, w) \in A$ in this partial order forces the corresponding undirected edge $e = \{u, w\}$ to be excluded from E_t . Thus, we can simply initialize our algorithm for constructing packing classes by fixing all undirected edges corresponding to A to be contained in \overline{E}_t . After running the original algorithm, we may get additional comparability edges. As the example in Figure 6 shows, this causes an additional problem: Even if we know that the graph \overline{G}_t has a transitive orientation, and all arcs $a = (u, w)$ of the precedence order (V, A) are contained in \overline{E}_t as $e = \{u, w\}$, it is not clear that there is a transitive orientation that contains all arcs of A .

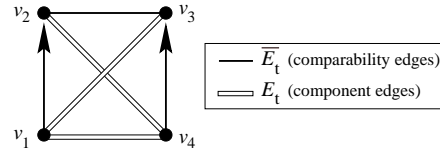


Figure 6. A comparability graph $\overline{G}_t = (V, \overline{E}_t)$ with a partial order A contained in \overline{E}_t , such that there is no transitive orientation of \overline{G}_t that extends A .

5.2 Finding Feasible Transitive Orientations

Consider a comparability graph \overline{G} that is the complement of an interval graph G . Deciding whether \overline{G} has a transitive orientation that extends a given partial order A is a problem that has been studied in the context of scheduling. Korte and Möhring [18] give a linear-time algorithm for determining a solution, or deciding that none exists.

In principle, it is possible to solve our more-dimensional packing problems with precedence constraints by adding this algorithm as a black box to test the leaves of our search tree for packing classes: In case of failure, backtrack in the tree. However, the resulting method cannot be expected to be reasonably efficient: During the course of our tree search, we are not dealing with one fixed comparability graph, but only build it while exploring the search tree. This means that we have to expect spending a considerable

amount of time testing similar leaves in the search tree, i.e., comparability graphs that share most of their graph structure. It may be that already a very small part of this structure that is fixed very “high” in the search tree constitutes an obstruction that prevents a feasible orientation of all graphs constructed below it. So a “deep” search may take a long time to get rid of this obstruction. This makes it desirable to use more structural properties of comparability graphs and their orientations to make use of obstructions already “high” in the search tree.

5.3 Implied Orientations

As in the basic packing class approach, we consider the component graphs G_i and their complements, the comparability graphs \overline{G}_i . This means that we continue to have three basic states for any edge: (1) edges that have been fixed to be in E_i , i.e., component edges; (2) edges that have been fixed to be in \overline{E}_i , i.e., comparability edges; (3) unassigned edges.

In order to deal with precedence constraints, we also consider orientations of the comparability edges. This means that during the course of our tree search, we can have three different possible states for each comparability edge: (2a) one possible orientation; (2b) the opposite possible orientation; (2c) no assigned orientation.

A stepping stone for this approach arises from considering the following two configurations – see Figure 7:

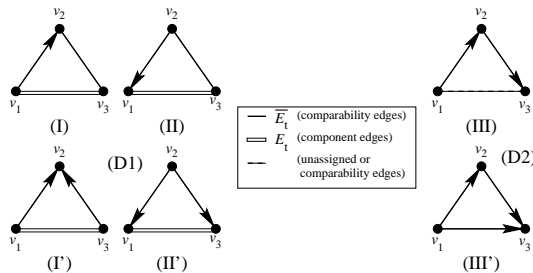


Figure 7. Implications for edges and their orientations: Above are path implications (D1, left) and transitivity implications (D2, right); below the forced orientations of edges.

The first configuration consists of two comparability edges $\{v_1, v_2\}, \{v_2, v_3\} \in \overline{E}_i$, such that the third edge $\{v_1, v_3\}$ has been fixed to be an edge from the component graph E_i . Now any orientation of one of the comparability edges forces the orientation of the other comparability edge, as shown in the figure. Since this configuration corresponds to an induced path in \overline{G}_i , we call this arrangement a *path implication*.

The second configuration consists of two directed comparability edges $(v_1, v_2), (v_2, v_3)$. In this case we know that the edge $\{v_1, v_3\}$ must also be a comparability edge, with an orientation of (v_1, v_3) . Since this configuration corresponds

to a triangle in \overline{G}_i , we call this arrangement a *transitivity implication*.

Clearly, any implication arising from one of the above configurations can induce further implications. Considering sequences of path implications leads to the following partition of comparability edges into *path implication classes*: Two comparability edges are in the same implication class, iff there is a sequence of path implications, such that orienting one edge forces the orientation of the other edge. For an example, consider the arrangement in Figure 6. Here, all three comparability edges $\{v_1, v_2\}, \{v_2, v_3\}$, and $\{v_3, v_4\}$ are in the same path implication class. Now the orientation of (v_1, v_2) implies the orientation (v_3, v_2) , which in turn implies the orientation (v_3, v_4) , contradicting the orientation of $\{v_3, v_4\}$ in the given partial order A . We call this type of contradiction a *path conflict* on a path implication class.

It is not hard to see that the path implication classes form a partition of the comparability edges, since we are dealing with an equivalence relation.

Similar to possible orientation conflicts for path implication classes, we may get a violation of transitivity implications, as a sequence of implications may force a directed cycle. (An example can be found in our mathematical report [6].) This type of violation we call a *transitivity conflict*.

Thus, we have the following necessary conditions for the existence of a transitive orientation that extends a given partial order A :

- D1:** Any path implication can be carried out without a conflict.
- D2:** Any transitivity implication can be carried out without a conflict.

These necessary conditions are also sufficient:

Theorem 2 (Fekete, Köhler, Teich) *Consider a partial order A with arc set contained in the edge set of a given comparability graph G . A can be extended to a transitive orientation of G , iff all arising path implications and transitivity implications can be carried out without creating a path conflict or a transitivity conflict.*

A proof and further mathematical details² are described in our forthcoming mathematical paper [5].

5.4 Solving OPPs with Precedence Constraints

We start by fixing for all arcs $(u, v) \in A$ the edge $\{u, v\}$ as an edge in the comparability graph \overline{G}_i , and we also fix its orientation to be (u, v) . In addition to the tests for enforcing the conditions for unoriented packing classes (C1, C2, C3),

²The interested reader may take note that we are extending previous work by Gallai [12], who extensively studied implication classes of comparability graphs. See Kelly [17], Möhring [20] for informative surveys on this topic, and Krämer [19] for an application in scheduling theory.

we employ the tests suggested by path implications and triangle implications. Like for packing classes, we can again get cascades of fixed edge orientations. If we get an orientation conflict or a cycle conflict, we can abandon the search on this tree node. The correctness of the overall algorithm follows from Theorem 2.

6 Computational Experiments

The first example is a numerical method for solving a differential equation (DE) with 11 nodes. The node operations are either multiplications or ALU-type operations. In this example, we treat overheads as constants, so they are not part of the optimization. In a second example, a video-codec using the H.261 norm is optimized. For this example, we consider the various possible overheads derived in Section 3.

6.1 DE Benchmark

In this benchmark, we assume a module library containing two hardware modules (box types): an array-multiplier and a module of type ALU that realizes all other node operations (comparison, addition, subtraction). For a word-length of $n=16$ bits, we assume a module geometry of 16×1 cells for the ALU module, and of 16×16 cells for the multiplier. Furthermore, the execution time of an ALU node takes one clock cycle, while a multiplication requires 2 clock cycles on our target chip.

The dependency graph is shown in Fig. 2. First, we compute the transitive closure of all data dependencies to allow our algorithm to find contradictions to feasible packings already in the input.

Next, we solve several instances of the BMP problem for different values of h_t reported in Table 1. Each h_t listed yields a test case for which the container size is minimized ($MinA$), assuming $h_x = h_y$. Also shown is the CPU-time needed for finding a solution.

Table 1. Computational results for optimizing reconfigurations for the DE benchmark

test	container sizes			CPU-time (s)
	h_t	h_x	h_y	
1	6	32	32	6.54 s
2	13	17	17	0.03 s
3	14	16	16	0.02 s

The reported optimization times were measured as the CPU-times on a SUN-Ultra 30 architecture.

For the DE benchmark, it turns out that a chip of 32×32 freely programmable cells is necessary to obtain a latency between 6 and 12 clock cycles. As the longest path in the graph has length 6, there does not exist any faster schedule. For 12 and 13 cycles, a chip of size 17×17 is necessary, for $h_t \geq 14$, a chip of size 16×16 cells is sufficient which is

the smallest chip possible to implement the problem as one multiplication by itself uses the full chip.

Similarly, the SPP is solved. The tradeoff between area size and necessary time is visualized in Fig. 8, where the Pareto-optimal points are shown. The figure also shows the Pareto-points for the case where no partial order needs to be satisfied (shown dashed).

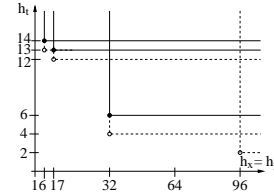


Figure 8. Pareto-optimal points for minimizing chip area and processing time for the DE benchmark. (a) Including partial order constraints (solid lines). (b) Without consideration of partial order constraints (dashed lines).

6.2 Video-Codec

In the following, we describe an optimal schedule for a hybrid image sequence coder/decoder. (See Figure 9.) Its purpose is to compress and decode video images using the H.261 standard.

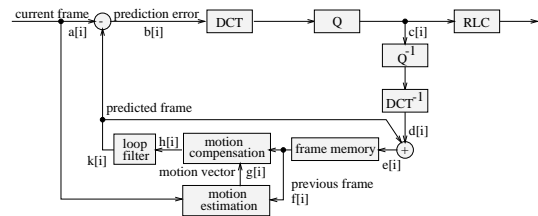


Figure 9. Block diagram of a video-codec (H.261)

For realizing the device, we have a library of three different modules. One is a simple processor core with a (normalized) area requirement of 625 units (25×25 cells, normalized to other modules in order to obtain a coarser grid) called PUM. Secondly, there are two dedicated special-purpose modules: a block matching module (BMM) that is used for motion estimation and requires $64 \times 64 = 4096$ cells; and a module DCTM for computing DCT/IDCT-computations, requiring $16 \times 16 = 256$ cells.

To illustrate the impact of various overheads (as discussed in Section 3) and of precedence constraints, we give three sets of computational results, as shown in Table 2. Tests “A” are without, “B” with precedence constraints. Test 1 uses a freely configurable device without any overhead. Test 2 assumes additional communication overhead. Test

3 considers communication overhead *and* expensive reconfiguration overheads for a Xilinx Virtex device. h_t reported corresponds to nanoseconds for a clock rate of 100 MHz.

Table 2. Computational results for optimizing reconfigurations for the Video-Codec

<i>test</i>	container sizes			CPU-time (s)
	h_t	h_x	h_y	
1 A	88,000	88	89	0.33 s
1 B	129,000	88	89	0.43 s
2 A	91,600	88	89	0.35 s
2 B	161,400	88	89	0.53 s
3 A	6,782,880	88	89	191 s
3 B	18,630,900	88	89	12.97 s

7 Conclusion

We have presented results for an idealized scenario without communication and reconfiguration overheads, as well as a case study with overheads for a Xilinx Virtex chip. It should be stressed that the huge overheads do not limit the feasibility of our approach, but rather indicate that this kind of chip is not very well suited for dynamically reconfigurable tasks. Enhancements of this chip should be a) higher bandwidth of the reconfiguration interface, and b) partial reconfiguration capabilities as in former generations such as the Xilinx 6200 series.

Finally, for the current technology, we conclude that reconfigurability may only be exploited with gain for task execution times that are much larger than the communication and reconfiguration times that are in the order of milliseconds today.

Acknowledgment

We are very grateful to Jörg Schepers for letting us use and extend the code for more-dimensional packing that he started as part of his thesis [21], to Marcus Bednara for support in estimating communication and reconfiguration overheads, and to several anonymous reviewers for helpful and encouraging suggestions.

References

- [1] Atmel. *AT6000 FPGA configuration guide*. Atmel Inc.
- [2] J. E. Beasley. An exact two-dimensional non-guillotine cutting tree search procedure. *Op. Research*, 33:49–64, 1985.
- [3] O. Diessel and H. ElGhindy. Partial FPGA rearrangement by local repacking. Report 97-08, Dept. of Comp. Sci and Software Eng., Univ. of Newcastle, Australia, 1997.
- [4] O. Diessel and H. ElGhindy. Run-time compaction of FPGA designs. *Proc. of FPL'97—the 7th Int. Workshop on field-programmable logic and applications*, pages 131–140, Berlin, 1997.
- [5] S. P. Fekete, E. Köhler, and J. Teich. Extending partial suborders and implication classes. Report 697-2000, TU Berlin.
- [6] S. P. Fekete, E. Köhler, and J. Teich. More-dimensional packing with order constraints. Report 698-2000, TU Berlin.
- [7] S. P. Fekete and J. Schepers. A new exact algorithm for general orthogonal d-dimensional knapsack problems. In *Algorithms – ESA '97*, volume 1284, pages 144–156, Springer LNCS, 1997.
- [8] S. P. Fekete and J. Schepers. New classes of lower bounds for bin packing problems. In *Proc. Integer Programming and Combinatorial Optimization (IPCO'98)*, volume 1412, pages 257–270, Springer LNCS, 1998.
- [9] S. P. Fekete and J. Schepers. On more-dimensional packing I: Modeling. Report 97-288, Center for Applied Computer Science, Universität zu Köln, Available at <http://www.zpr.uni-koeln.de/ABS/~papers>, 1997.
- [10] S. P. Fekete and J. Schepers. On more-dimensional packing II: Bounds. Report 97-289, Universität zu Köln, 1997.
- [11] S. P. Fekete and J. Schepers. On more-dimensional packing III: Exact algorithms. Report 97-290, Universität zu Köln, 1997.
- [12] T. Gallai. Transitiv orientierbare Graphen. *Acta Math. Acad. Sci. Hungar.*, 18:25–66, 1967.
- [13] M. Garey and D. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman, New York, 1979.
- [14] M. Golumbic. *Algorithmic graph theory and perfect graphs*. Academic Press, New York, 1980.
- [15] E. Hadjiconstantinou and N. Christofides. An exact algorithm for general, orthogonal, two-dimensional knapsack problems. *European Journal of Operations Research*, 83:39–56, 1995.
- [16] C.-H. Huang and J.-Y. Juang. A partial compaction scheme for processor allocation in hypercube multiprocessors. In *Proc. of 1990 Int. Conf. on Parallel Proc.*, pages 211–217, 1990.
- [17] D. Kelly. Comparability graphs. In I. Rival, editor, *Graphs and Order*, pages 3–40. D. Reidel Publishing, Dordrecht, 1985.
- [18] N. Korte and R. Möhring. Transitive orientation of graphs with side constraints. In H. Noltemeier, editor, *Proceedings of WG'85*, pages 143–160. Trauner Verlag, 1985.
- [19] A. Krämer. *Scheduling Multiprocessor Tasks on Dedicated Processors*. Doctoral thesis, Fachbereich Mathematik und Informatik, Universität Osnabrück, 1995.
- [20] R. H. Möhring. Algorithmic aspects of comparability graphs and interval graphs. In I. Rival, editor, *Graphs and Order*, pages 41–101. D. Reidel Publishing Company, Dordrecht, 1985.
- [21] J. Schepers. Exakte Algorithmen für orthogonale Packungsprobleme. Doctoral thesis, Universität Köln, 1997, available as Report 97-302.
- [22] J. Teich, S. Fekete, and J. Schepers. Compile-time optimization of dynamic hardware reconfigurations. In *Proc. Int. Conf. on Parallel and Distributed Processing Techniques and Applications (PDPTA'99)*, pages 1097–1103, Las Vegas, U.S.A., June 1999.
- [23] J. Teich, S. Fekete, and J. Schepers. Optimization of dynamic hardware reconfigurations. *J. of Supercomputing*, to appear, 2001.
- [24] Xilinx. XC6200 field programmable gate arrays. Technical report, Xilinx, Inc., October 1996.