

# A RUN-TIME GRAPH-BASED POLYNOMIAL PLACEMENT AND ROUTING ALGORITHM FOR VIRTUAL FPGAS

R. Ferreira<sup>1,2</sup>, L. Rocha<sup>1</sup>, A. Santos<sup>1</sup>, J. Nacif<sup>1\*</sup>

Stephan Wong<sup>2</sup>

Luigi Carro<sup>3</sup>

<sup>1</sup> Departamento de Informática  
Univ. Federal de Viçosa, Brazil  
{ricardo, jnacif}@ufv.br

<sup>2</sup>Computer Engineering  
TU Delft, The Netherlands  
j.s.s.m.wong@tudelft.nl

<sup>3</sup>Instituto de Informática  
UFRGS, Brazil  
carro@ii.ufgrs.br

## ABSTRACT

Dynamic partial reconfiguration enables efficient use of hardware resources by multiplexing system functionality in time. However, many challenges arise from partial reconfiguration implementation. The placement and routing (P&R) of the hardware modules is a computationally intensive task, and the state-of-art algorithms are not suitable to place and route modules at run-time. This paper makes several contributions: (1) *Single Placement at run-time*: we propose a novel P&R algorithm based on greedy heuristic where a single placement is performed at run-time in few milliseconds. (2) *Implicit Graph Model*: the FPGA is modelled as an implicit graph with a direct correspondence to the physical FPGA, and the P&R is performed as a graph mapping problem by exploring the node locality during the depth-first traversal. (3) *Polynomial Placement*: we show that even a single placement can be routed without critical path degradation. (4) *Fragmented Regions*: the graph approach is flexible, and it allows efficient placement even onto fragmented FPGA areas. Compared with the most popular P&R tool running the same benchmark suite our algorithm is on average 864x faster. Moreover, the bitstream for partial reconfiguration is also reduced by a factor of 4.

## 1. INTRODUCTION

Field-Programmable Gate Arrays (FPGAs) can be partially or fully programmed through configuration bitstreams. Partial reconfiguration enables efficient allocation of logic resources by adding more functionality using unused resource or by sharing and/or multiplexing resource during the execution time, depending on the current needs. In dynamic partial reconfiguration (DPaR), the reconfiguration is performed at run-time keeping non-reconfigured FPGA areas running.

In practice, DPaR implementation involves a series of challenges which are being treated by different techniques [1,

2, 3, 4]. In this paper, we address bitstream placement and routing (P&R), one of the most computationally intensive and time-consuming DPaR procedures. When P&R maps a new bitstream at run-time, an available region over the target architecture should be identified. This region should also provide both sufficient hardware resources and minimize the impact to upcoming DPaR procedures.

Traditionally, P&R is performed offline because hardware resource re-arrangement is highly computationally intensive. Traditional approaches [5] include recursive partitioning, analytic, genetic algorithms, and simulated annealing. Unfortunately, because of their high computational requirements none of the previous approaches is suitable for run-time P&R algorithm. Recently, a Just-in-Time (JIT) framework [4, 6] to support DPaR is introduced. This strategy reduces the fragmentation for hardware resources (pre-allocated FPGA area, etc.) in partial reconfiguration context. The configuration bitstream for a virtual FPGA is computed at run-time by performing technology packing, and P&R. However, the P&R still takes few seconds to complete, and hence run-time P&R remained impractical.

In this paper, we propose a novel P&R based on a graph mapping model that performs nearly three orders of magnitude faster in comparison to the state-of-the-art of P&R algorithms [7]. In addition, our approach reduces the bitstream size by a factor of 4. The main contributions of the work presented in this paper are: (1) A novel polynomial P&R greedy heuristic based on graph mapping; (2) A local routing for FPGA by exploring the graph locality and by prioritizing the critical path; (3) An adaptable P&R for fragmented regions; (4) A run-time P&R suitable for dynamic partial reconfiguration frameworks.

This paper is outlined as follows. Section 2 describes our mapping approach. Section 3 presents the graph model used to map the FPGA, and Section 4 describes how the graph model is explored to produce the P&R. In Section 5, our P&R algorithm is evaluated, while we present related work in Section 6. Finally, we conclude our remarks and present future work in Section 7.

\*Financial Support: TU Delft and the Brazilian Inst. and Companies: CNPq & Ciência Sem Fronteiras, FAPEMIG, CAPES, FUNARBE, Gapso

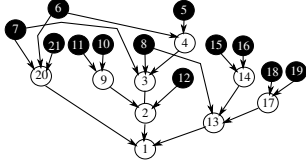


Fig. 1. An input LUT graph example.

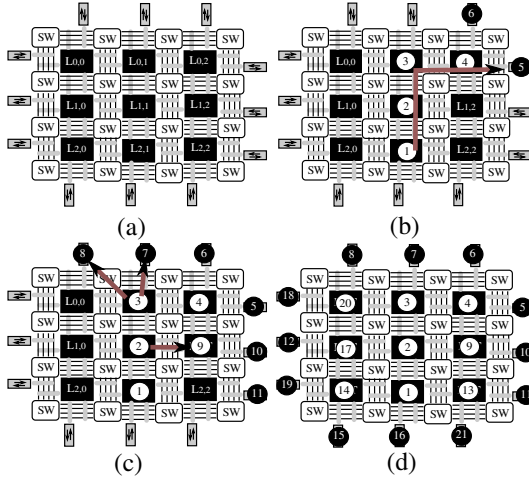


Fig. 2. (a) FPGA (b) First Path (c) Partial (d) Final.

## 2. GRAPH MAPPING

Our approach models the P&R as a graph mapping problem, where an input graph is mapped onto an output graph. The input graph is the circuit to be mapped. It is represented by a  $k$ -input LUT graph. Fig. 1 depicts a simple example. The external inputs are represented by black vertices and the LUTs by white vertices. Since our algorithm is based on a depth-first (DF) graph traversal, the vertex numbers follow a DF order. During the DF traversal, the descendant vertices are ordered by the depth the traversal can reach. For example, vertex 1 has three descendants: 2, 13 and 20. The deepest the traversal can go is through vertex 2, and hence it will be the first to be visited. The goal is to prioritize the critical path during the mapping.

The output graph is the FPGA. For ease of explanation, we start with a simplified view of the FPGA graph as depicted in Fig. 2(a). We use the term vertex for the input graph and the term node for the output graph. The output graph is composed by four node types: LUT, switch box, in/out, and segmented bus. The algorithm maps an LUT vertex onto an LUT node, and an input vertex onto an in/out node. The switch boxes and wiring segments are used for routing the edges. An edge from the input graph is mapped in one or more wiring segments onto the target FPGA graph. For the rest of this paper, when we count a wiring segment or track, it will also include a switch box connection.

We use DF traversal to map the input graph depicted in Fig. 1. First, our algorithm maps the path from 1 to 5 as shown in Fig. 2(b). The LUTs are labeled by the line and column indexes. It is important to note that the FPGA graph is also traversed in DF order. Then, the DF returns to vertex 4, placed at  $LUT_{0,2}$ , and maps the edge  $4 \rightarrow 6$  at an adjacent I/O. Next, the DF returns to vertex 3 to visit 7 and 8. Fig. 2(c) presents a snap-shot when the DF returns to vertex 2, and places vertex 9 at  $LUT_{1,2}$  in order, adjacent to node  $LUT_{1,1}$  where vertex 2 was placed. Fig. 2(d) depicts the final mapping. The time complexity is polynomial as each edge is visited once. Previous work on CGRA has already used graph mapping approach based on depth-first traversal [8, 9]. However, an FPGA graph has nodes with high fanout degrees. Moreover, the routing infrastructure as well as the lower FPGA granularity at the bit level results in a more complex mapping in comparison to CGRAs.

## 3. IMPLICIT FPGA GRAPH

At compile time, the synthesis tool is responsible to generate the input graph, which is stored as an edge list in depth first order. However, the output FPGA graph is much more complex, requiring an efficient and scalable representation. We propose a novel implicit representation, which is different from traditional graph structures like adjacency/incidence list or matrix, which is detailed in the following sections.

### 3.1. Virtual FPGA

Since there are several commercial FPGA families, our approach is based on a virtual FPGA as proposed in [4, 6]. In spite of that, our model can be customized for a specific FPGA family as the mapping is based on a graph approach. Moreover, the virtual FPGA has several advantages: it is independent to the underlying hardware; it can be directly implemented upon a traditional off-the-shelf FPGA device, and it allows partial configuration as shown in [4], even if the target device does not have partial reconfiguration infrastructure.

### 3.2. Logic Block and Track Nodes

Fig. 3 depicts the LUT connections and its graph representation. Let us consider a 4-bit bidirectional channel FPGA. Each  $LUT_{i,j}$  has 4 input multiplexers to connect to the channel tracks. Each LUT will be represented as a LUT node. The directions west, north, east, and south are labeled as 0, 1, 2, and 3, respectively. The LUT output is connected to the south tracks from left (bottom) to right (top) as described in [6]. Each track multiplexer will be represented as a T node and it will be labeled by  $i, j$ , and the track number. Fig. 3(b) depicts 4 track nodes:  $T_{i,j,3}, T_{i,j,2}, T_{i,j,1}, T_{i,j,0}$ .

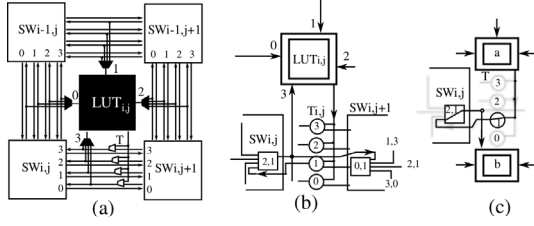


Fig. 3. (a) FPGA LUTs (b) LUT Output (c) Adjacent LUTs.

### 3.3. Switch Box

The switch box (SW) is the key component in the FPGA design. The proposed implementation is based on the Wilton SW [10]. All tracks should traverse at least one SW. The number of in/out depends on the channel width. We propose to represent a SW as a set of small multiplexer SW nodes. Each SW node implements a single switch output, which has a direct correspondence to a physical multiplexer. The multiplexer receives four inputs tracks, one from each SW direction (west, north,...). If there are  $C$  channel tracks, then one switch box will be represented by  $4 \cdot C$  SW nodes. The label  $i, j, d, c$  identifies the SW node, where  $c$  is the track number, and  $d$  is the direction (west,north,...). Fig. 3(b) depicts two SW nodes which are connected to the T node  $T_{i,j,1}$  and the  $LUT_{i,j}$  input 3. The  $SW_{i,j,2,1}$  is from the east side of switch box  $i, j$  and  $SW_{i,j,0,1}$  is from the west part of the switch box  $i, j+1$ . The track  $k = 1$  for the switch  $SW_{i,j+1}$  in Fig. 3(b) is implemented as multiplexer with the following input tracks following Wilton pattern [10]:  $t_{i,j+1,0,k}$ ,  $t_{i,j+1,1,\frac{w-k}{w}}$ ,  $t_{i,j+1,2,k}$ , and  $t_{i,j+1,3,\frac{k-1}{w}}$ .

## 4. ALGORITHM

The proposed mapping algorithm is based on DF traversal in both graphs. Therefore, in an optimal scenario, two adjacent vertices  $a \rightarrow b$  should be mapped in two adjacent nodes  $x \rightarrow y$ . However, this assumption is not valid for all vertices since there are physical routing constraints in the FPGA graph. Moreover, two vertices could be placed onto two nodes far away from each other. First, we present the proposed algorithm for adjacent nodes. Then, we discuss our strategy for non-adjacent nodes.

### 4.1. Adjacent LUT Nodes

Initially, let us consider an edge  $a \rightarrow b$ . At least one track and one switch node should be traversed in the FPGA graph as depicted in Fig. 3(c). This mapping example is the best local situation for an edge. However, it could be impossible even for adjacent line and/or column LUT. Therefore, we introduce a novel concept for FPGA adjacent nodes.

Let us consider the node  $L_{i,j}$  at line  $i$  and column  $j$ ,

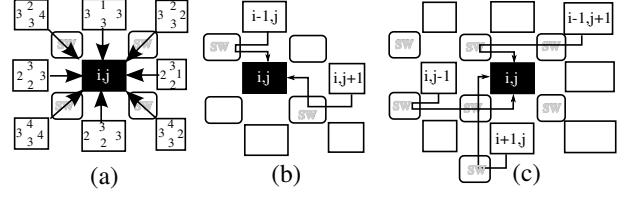


Fig. 4. (a) Routing Cost (b) One SW (c) Two adjacent SW.

Fig. 4(a) depicts the routing costs for an 8-neighborhood. Since  $L_{i,j}$  has four inputs, there are four routing costs. The costs are depicted inside each neighbor. For example, the right node  $L_{i,j+1}$  traverses 2, 3, 1, and 2 SWs to connect to the  $L_{i,j}$  input west, north, east, and south, respectively. Fig. 4(b) depicts the two lowest cost possibilities. The first, from  $L_{i-1,j}$  to  $L_{i,j}$  north input uses only one switch connection because the output LUT port is always in the south side (see Fig. 3). The second case is from  $L_{i,j+1}$  to the east input of  $L_{i,j}$ . Therefore, these are only two routes that use only one switch. It is important to note that the target is a virtual FPGA presented in [4], however, our approach can be modified to different FPGA architectures.

The mapping will first try to P&R the south and east neighbors. Let us again consider the graph example from Fig. 1. For ease of explanation, we only consider the internal nodes. The first path in DF order is  $1 \leftarrow 2 \leftarrow 3 \leftarrow 4$ . The mapping of this path is depicted in Fig. 2(b). Four vertices are placed and routed with minimal routing cost. Edges  $2 \leftarrow 9$  and  $1 \leftarrow 13$  are also placed and routed with the minimal cost. However, vertex 14 could not be placed in an adjacent position next to vertex 13, as shown in Fig. 2(d), since vertex 13 is placed at the bottom-right corner, and there is no free node among its 8 neighbors.

We propose to extend the concept of adjacent node to routing costs greater than one. Fig. 4 depicts three adjacent nodes with routing cost two to connect to three different  $L_{i,j}$  inputs. If the routing cost three is considered, all nodes in Fig. 2(c) are placed and routed as adjacent nodes, even the edges  $14 \rightarrow 13$ ,  $17 \rightarrow 13$ , and  $20 \rightarrow 1$ .

### 4.2. Non-Adjacent LUT Nodes

Fig. 1 depicts an example composed by only 9 vertices. Nevertheless, there are non-adjacent nodes which should still be routed. Let us now consider the input vertices. Some input vertices could be handled as adjacent nodes as such vertices 5 and 6. However, edge  $21 \rightarrow 20$  should be routed by a non-adjacent routing algorithm, since there is no place for the source vertex in the neighborhood of the target node. We find the closest node by using a breadth-first search around  $L_{i,j}$ . Then, the source vertex is placed, and a routing algorithm is performed. Another situation occurs when the source node is already placed, but it has multicast edges such

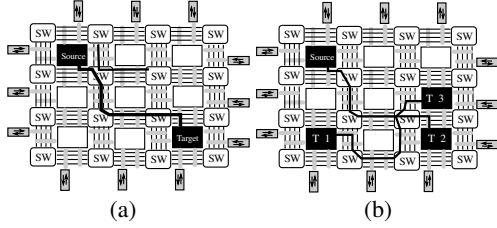


Fig. 5. (a) XY Routing (b) Multicast Routing.

as input vertices 6 and 7 in Fig. 1. Since these vertices have been already placed as adjacent nodes of 4 and 3, when the edge  $6 \rightarrow 20$  is visited in DF order and the routing algorithm handle these connections.

### 4.3. XY Routing and Multicast

For multicast edges, we propose to use a simple and greedy routing algorithm. Our algorithm is based on Network-on-Chip (NoC) XY routing [11, 12]. The routing path starts from the source node, moving through switch boxes until target node final X (line) and Y (column) are reached. We adapted the algorithm to an FPGA architecture considering switch box track channels and LUT side inputs.

Fig. 5(a) depicts a simple source to target routing example. First, the algorithm tries to route in the X direction. However, since the output track in the X direction is already used, the algorithm routes in the Y direction. The next two steps are accomplished in the X direction, and finally, the routing reaches the target at the north input.

Fig. 5(b) depicts a multicast routing to three target destinations:  $T_1$ ,  $T_2$ , and  $T_3$ . Multicast routing strategy reduces the overall switch cost by sharing routing resources. When using multicast strategy, the routing should observe track directions. For example, south track is connected from left to right. Therefore, since  $T_3$  input is in the south side, our algorithm uses one extra switch box to perform the routing.

Our P&R is based on DF order edge visit. When a source vertex is visited more than two times, our algorithm does not route its edges on the fly. A list of target vertices is created during the traversal for each multiple fanout source node. Finally, when all edges have been visited, all multicast edge lists are processed. In most circuit functions, control signals are connected to several destinations as a multicast signal. The routing cost of these signals strongly contributes to the total routing cost.

## 5. EXPERIMENTAL RESULTS

We compare our results with the VPR tool and benchmark suite [7, 13], which are publicly available. The experiments are performed in an Intel i3 370M, 2.4 GHz, 3MB L2 cache. For a given circuit, depending on the selected P&R options,

Table 1. Execution Time for performing P&R.

Bench	CLB	VPR		DF		gain	Long Wire		
		sec	t	msec	t		1	2	3
fir16	47	0.04	4	0.1	3	400	3	3	3
apex7	153	0.21	10	1.2	14	175	12	12	12
9symml	179	0.19	9	1.29	9	147	9	9	9
alu2	241	0.27	9	1.28	12	210	12	12	12
Alu4	415	0.52	12	0.74	12	702	11	11	11
Apex2	966	1.7	17	1.53	19	1111	15	15	15
Apex4	885	1.4	19	1.18	21	1186	16	12	12
too_large	992	1.62	10	1.58	15	1025	15	15	15
term1	261	0.33	10	0.492	15	670	15	15	15
k2	447	0.64	15	0.72	19	888	14	12	12
Misex3	771	1.2	16	1.05	19	1142	14	12	12
Ex5p	646	1.01	16	0.88	21	1147	16	11	11
Ex1010	964	1.68	20	1.54	23	1090	18	13	13
spla	2730	7.27	21	4.74	37	1533	32	27	22
pdc	2857	7.49	22	4.87	39	1537	34	29	24

the placement operation time can vary from a few milliseconds (msec) to several minutes. Explanation of the VPR configuration options is beyond the scope of this paper. More details can be found in [7].

Regarding the experimental setup, we use a target virtual FPGA similar to the one proposed in [4]. The virtual FPGA is composed by an array of 100x100 slices, each channel containing 50 routing tracks. These architectural definitions do not affect the generality of our solution, since JIT algorithm [4], and our DF algorithm are performed onto the same device. Since our objective is to reduce P&R execution time, we have set VPR to fast mode. The selected parameters are: *-innernum 1 -placealgorithm boundingbox -routechanwidth 50 -routeralgorithm directedsearch*. The *innernum 1* option speeds up the placer by a factor of 10, and the bounding box placement focuses purely on minimizing the bounding box wirelength, not focusing on the critical path delay optimization. The directed search router is routability-driven based on A\* heuristic to improve run-time.

### 5.1. Execution Time

The proposed tool takes only  $2.67\mu\text{sec}$  per CLB while the JIT tool requires about 3.59 msec per CLB, and VPR performance<sup>1</sup> is 26.35 msec per CLB. Table 1 depicts P&R execution time for our DFS approach and VPR tool [7] considering a set of 15 MCNC benchmarks [13]. Our DFS algorithm is 864x faster than VPR in fast mode. Moreover, the run-time improvement can reach up to 1537x for the *Pdc* benchmark.

### 5.2. Simultaneous Mapping

Our proposed approach could also map more than one application at run-time. Fig. 6(a) depicts the FPGA mapping

<sup>1</sup>The VPR parameters are not detailed in [4]

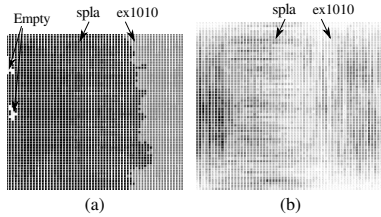


Fig. 6. Spla+ex1010 Occupancy: (a) LUT (b) SW box.

in a 62x62 CLB array of two benchmarks: ex1010 and spla. The white boxes show empty LUTs, the spla is placed on the left black boxes, and the ex1010 is placed on the right grey boxes. The granularity of the mapping at LUT level significantly reduces the fragmentation of the partial reconfiguration, which is also an advantage presented in [4]. There are no rectangular constraints, and the border line is wavy to maximize the occupancy.

Fig. 6(b) depicts the number of tracks per channel occupied for a successful routing. A greyscale is used to show the switch box occupancy, whereas the white color represents low occupancy, and the hotspots are shown in black. The spla benchmark uses more tracks near the inputs while ex1010 is more uniformly distributed.

### 5.3. Wiring Segments, Tracks and Long Wires

Our placement is a greedy heuristic, and each edge is only visited once. Therefore, we achieve higher speedups in larger circuits. Although the placement could generate a routable solution, the number of wire segments per edge and the maximal track occupation in each channel increase as a function of circuit size. Table 1 depicts the maximal track usage in VPR and in our DFS P&R. For medium size circuits, with up to 1000 4-input LUT or 4000 equivalent gates, the maximal track is close to the VPR solution. However, for larger circuits, our P&R can use up to the double of maximal track usage in smaller circuits. A simple strategy to reduce routing resource usage with no P&R time degradation is to add long wires. In commercial FPGAs, long wires are used to speed up connections and to avoid channel and switch box congesting. They are similar to single wires, except that each one spans two or more LUTs, imposing lower routing delays.

The next experiment evaluates the maximal track usage for the larger MCNC benchmarks in the presence of long wires that span 5 LUTs. The last three columns of Table 1 depict maximal track usage for VPR and for our approach considering the use of 1, 2 or 3 long wire tracks per channel. For medium size circuits the addition of one or two long wire tracks strongly improves P&R resource usage. If more lines are added, the improvements are marginal. However, when we consider larger benchmarks, the addition of three

Table 2. Edge Distribution and Configuration Memory.

Bench	Local %		Internal Multi	Input %		Mem	
	Adj	Near		Single	Multi	DF	Virtual
fir16	67	5	26	2	0	342	2505
apex7	19	5	32	10	33	2166	8640
9symml	26	6	24	2	42	2618	10020
alu2	29	8	30	2	30	3010	13088
Alu4	30	5	39	1	25	5378	22546
Apex2	27	5	45	1	21	13916	52352
Apex4	30	5	57	0.5	6	11528	46012
too_large	27	9	9	1	54	13968	52352
term1	26	3	8	4	58	3930	14775
k2	29	4	43	3	20	5806	24744
Misex3	28	5	46	1	19	10590	40082
Ex5p	31	3	62	0.5	3	8140	13088
Ex1010	27	5	65	0.3	2	13656	52392
spla	29	6	63	0.2	1	35568	143610
pd	30	6	62	0.2	1	37294	149080
<b>Average</b>	30.3	5.3	40.7	1.9	21	10445	40328

long wire tracks promotes a significant reduction maximal track usage.

For the evaluated benchmarks, when VPR timing driven P&R is selected, wire usage increases  $1.15\times$  achieving a  $1.6\times$  critical path performance degradation. Moreover, the P&R is  $10\times$  slower. Our proposed approach requires on average  $3.27\times$  more wiring segments than VPR for the benchmarks shown in Table 1. If 1, 2, and 3 long wires are added, the required number of segments are  $2.49\times$ ,  $2.17\times$ ,  $2.01\times$  bigger than VPR, respectively.

We believe there are three main reasons to explain why our approach requires more wires. First, the placement prioritizes the original critical path, which will result in at least 15% wire length increase, as shown by VPR P&R analysis. Second, although placement is a NP-complete problem, our algorithm tries only one placement based on DF order graph locality. Traditional offline FPGA P&R algorithms evaluate several configurations to reduce wire length. It is important to note that we reduce the P&R time by a factor of  $1000\times$  at the cost of spending  $3\times$  more wires. Nevertheless, those extra wires are already there in the FPGA, so the question that remains is whether this will impact performance, as it will be discussed in Section 5.5.

### 5.4. Edge Distribution

Since placement quality directly impacts on routing cost, we analyze edge distribution in the original graphs to better understand our results. Table 2 depicts three edge types in the original graph. The first two columns show the local edges. As mentioned in Sections 4.1 and 4.2, our placement performs a simultaneous traversal in both original and FPGA graphs. Table 2 shows that 35% of edges are visited, placed and routed directly using less than 3 switch boxes or wire segments. On average less than 2 segments

are used. The main challenge is to route the internal and input edges because they correspond to 63% of the edges. Most of these edges present large fanouts or multicast degrees. Moreover, although the fanout degree is on average 3, few nodes present a very high degree (the internal or input control circuit signals). For the *pd* benchmark instance, 10%, 4%, and 1% of nodes concentrate 48%, 34%, and 18% of the total edges, respectively. Further work include a breadth-first traversal prioritizing these nodes instead the local ones. The third reason of high wiring cost of our approach is the simple routing approach, as explained in Section 4.3. Moreover in Section 5.5, we replace our routing by the VPR routing to show that a better routing could significantly reduce the wiring cost. Further work should be done to include more low cost routing as local expansion wavefront approach of VPR [7]. Despite of that, our DFS approach achieves very low P&R times, around milliseconds.

Table 2 also depicts the configuration memory size. Column DF presents the size to store the original graph in Kbytes. Column Virtual shows the minimal square matrix to store the configuration bits for the virtual FPGA (LUT and SW). Our approach requires on average,  $3.7\times$  less bits, which can reduce the size of the configuration memories in dynamic scenarios.

## 5.5. Critical Path

Table 3 depicts the critical path, the execution time, and the wiring segments for three placements. Our goal is to evaluate the quality of our placement when a better routing is used. We compare our depth-first (DF) placement with VPR. The first VPR placement (FA) optimizes the execution time selecting the parameters: *-innernum 1 bounding-box*. The second VPR placement (CR) optimizes the critical path using the parameters: *-innernum 10 timingdrive*. Column *Ratio Critical* presents the critical path ratio in order to compare our DF placement with VPR. On average, the DF and FA are equivalent, and CR reduces critical path in 26%. Regarding the execution time, our DF placement is on average  $1024\times$  and  $13215\times$  faster than FA and CR placements. The critical path for all three placements is obtained using VPR routing. Column *Routing* depicts the execution time for VPR routing algorithm, which is almost the same for the three placements, except for the last two benchmarks. For the *pd* and *spla*, 1% of the nodes have 15% of edges. These edges will generate more than 40% of the wiring segments. However, our DF approach places these nodes by using only one edge. Further work includes a local search to better evaluate one edge positioning for these hotspot nodes.

Fig. 7 depicts a simple example of a hotspot. The position of vertex *b* shown in Fig. 7(a) will be defined during the first traversal in the DF order when the edge  $a \leftarrow b$  is visited, *b* is placed adjacent to *a*, since the algorithm traverses the graph from the outputs to the inputs. It is important to highlight

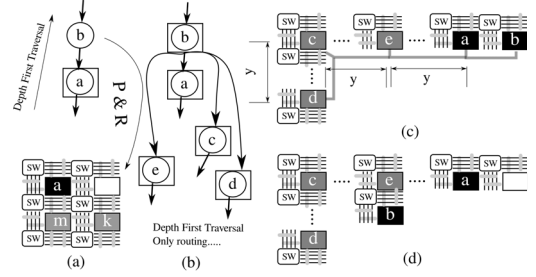


Fig. 7. (a) Place (b) Fanout (c) Routing (d) New Place.

that *b* is not only placed, it is also routed with a minimal local cost. For this reason, we apply the term P&R in this paper. As mentioned before, 35% of the edges are included in this case.

Let us suppose *b* has fanout=4 as shown in Fig. 7(b). When the edges  $c \leftarrow b, d \leftarrow b$ , and  $e \leftarrow b$  are traversed, *b* has been already placed, and these nodes are included in the edge list of *b* as mentioned in Section 4.3. Let us suppose an optimal solution for the routing as depicted in Fig. 7(c). The total wiring segments will be  $2 \cdot y + x + 1$ . However the average distance between *b* and its outputs will have a cost of  $\frac{5 \cdot y + x + 1}{4}$  wiring segments. However, if *b* is move to a new position close to *e*, the average cost is reduced to  $\frac{3 \cdot y + x + 1}{4}$ . Future work will include the investigation of possible improvements on the hotspot placements. Moreover, the critical path could be not affected directly when the connections to *c*, *d*, and *e* are routed from *b*.

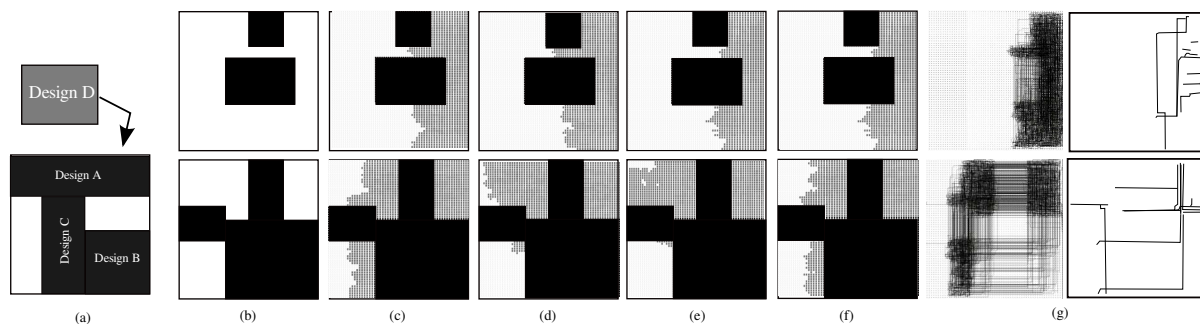
Table 3 also presents the wiring segments. Our DF approach requires on average 46% more wiring segments, since it is based on a greedy heuristic where we evaluate only one single position for each node. Since routing resources in modern FPGAs are not restricted, wiring usage is not a problem. Nevertheless, the DF critical path is on average the same in comparison with FA placement. In fact, it is only 26% worse than the best VPR placement, which tries multiple positions several times and takes  $1222\times$  to  $36318\times$  longer in execution time. Finally, our initial results are very promising since the proposed algorithm tries only once each node connection.

## 5.6. Fragmentation

Considering a function set  $f = A, B, \dots$  to be used during the execution of one application. If the bitstream for each function is generated at compile time, the FPGA could have a poor resource utilization. Moreover, it could be impossible to place all functions, even if there is enough routing and LUT resources. Fig. 8(a) depicts an example where function D can not be placed due to fragmentation. Despite there are two empty areas (white parts) with enough resources it is not possible to place D, since it is generated as a single rectangular region.

**Table 3.** Placement Execution Time and Critical Path by using VPR Routing.

Bench	Critical Path			CPU Time Placement			CPU Time Routing			Wiring Segments			Ratio Critical		Speedup	
	DF	FA	CR	DF msec	FA sec	CR sec	DF sec	FA sec	CR sec	DF	FA	CR	FA/DF	CR/DF	FA/DF	CR/DF
<b>fir16</b>	1.21	1.27	1.18	0.09	0.02	0.11	0.03	0.03	0.04	123	85	81	0.95	1.03	222	1222
apex7	0.74	0.68	0.62	0.31	0.12	1.77	0.16	0.17	0.18	1233	828	799	1.09	1.19	387	5710
9symml	0.97	0.87	0.77	0.39	0.11	1.3	0.17	0.2	0.19	1141	889	878	1.11	1.26	282	3333
<b>alu2</b>	2.31	2.53	2.05	0.53	0.16	1.78	0.23	0.23	0.21	1593	1176	1100	0.91	1.13	302	3358
<b>alu4</b>	2.42	2.58	2.02	0.62	0.31	4.45	0.44	0.36	0.4	3121	2322	2307	0.94	1.20	500	7177
apex2	1.68	1.58	1.16	1.04	1.24	16.6	1.18	0.98	1.18	12429	8184	8032	1.06	1.45	1192	15962
<b>apex4</b>	1.51	1.63	1.14	0.79	1.05	13.01	1.18	0.94	1.02	11249	7980	8151	0.93	1.32	1329	16468
<b>too.large</b>	3.21	3.94	2.95	1.25	1.08	14.5	1.19	1.3	1.2	6522	5357	5197	0.81	1.09	864	11600
term1	0.87	0.89	0.73	0.44	0.22	2.41	0.28	0.24	0.32	1704	1357	1431	0.98	1.19	500	5477
k2	1.65	1.45	1.27	0.54	0.46	5.35	0.57	0.43	0.43	4961	2912	3012	1.14	1.30	852	9907
misex3	1.27	1.17	0.98	0.73	0.86	11.37	0.92	0.81	0.79	8845	5876	5853	1.09	1.30	1178	15575
ex5p	1.55	1.36	1.13	0.63	0.72	9.43	0.75	0.81	0.78	7229	4758	4846	1.14	1.37	1143	14968
<b>ex1010</b>	1.38	1.72	1.1	0.91	1.26	15.96	1.52	1.3	1.47	12450	9318	9976	0.80	1.25	1385	17538
spla	2.68	2.54	1.97	2.18	5.88	73.26	29.9	4.06	4.06	51150	28385	28370	1.06	1.36	2697	33606
pdcc	2.66	2.6	1.79	2.2	6.24	79.9	21.1	4.32	4.4	53045	29897	31684	1.02	1.49	2836	36318
<b>Average</b>											1.46	1.46	1.00	1.26	1045	13215

**Fig. 8.** (a) Fragmentation (b) Initial Scenario (c) Apex2 (d) Apex4 (e) Too\_large (f) Ex1010 (g) Routing and Critical Path

To evaluate the potential of our placement algorithm on fragmented areas, let us consider two scenarios as shown in Fig. 8(b). Let us use the black color for the area which has been already allocated. In the first scenario (top), we have two occupied disjoint areas while in the second scenario (bottom), most part of the area has been already allocated. There are three disjoint free areas in the second scenario. For both scenarios, we map four MCNC benchmark circuits with around 1000 LUTs onto a target FPGA array of 60x60 LUT slices with channel width 50. Fig. 8(c) depicts the area where the function *Apex2* is placed in grey color. Fig. 8(d), 8(e), and 8(f) presents the area for *Apex4*, *Too\_large*, and *Ex1010*.

Our algorithm executes the placement in less than 2 milliseconds, even in fragmented areas. Since it is a graph based approach, it is easily adaptable to fragmented regions. The output of our placement is routed using VPR to obtain the critical path. Fig. 8(g) depicts the VPR [7] screenshot for *Ex1010* benchmark routing and critical path. Most of wiring segments are concentrated on the placement region. More-

over, wiring segments are also used to route through disjoint areas. Table 4 shows the critical path and total wiring segments for both scenarios. Column Ratio depicts critical path length increase compared with optimal critical path. The optimal critical path is obtained by using VPR at timing driven mode mapped on a square contiguous area as shown in Table 3. On average, the critical path increases 41% and 47% for the scenario 1 and 2. On average, the routing resource measured as wiring segments increases 57% and 79% for the scenario 1 and 2.

## 6. RELATED WORK

The state-of-art in P&R algorithm is represented by the VPR framework [7, 14]. The placement approaches are based in three techniques: partitioning, simulated annealing (SA), and analytic placement. In [15], a dynamically adaptive stochastic tunneling algorithm to avoid the freezing problem in SA approach is proposed. However, only a marginal reduction of 18.3% in run-time is obtained over VPR. A par-

**Table 4.** Wiring and Critical Path for Fragmented Regions.

Bench	Scenario 1			Scenario 2		
	Wires	Critical	Ratio	Wires	Critical	Ratio
too_large	7846	3.08	1.04	7328	3.52	1.19
apex2	13328	1.67	1.43	15151	1.72	1.48
apex4	14002	1.62	1.42	17426	1.77	1.55
ex1010	13922	1.91	1.73	17713	1.86	1.69
average			1.41			1.47

allel approach for SA has been recently proposed in [5]. The algorithm evaluates multiple SA moves in parallel, the experimental results achieve an average speedup ranging from  $2\times$  up to  $7\times$  in comparison to VPR. An analytic placement based on a near-linear net model was proposed in [16], which is  $5\times$  faster than VPR fast mode (i.e., with inner\_num=1). Moreover, this placement obtains an 9% reduction in critical-path delay. The routing is performed by using VPR router. The reported CPU time is in the range of seconds for the MCNC benchmarks.

Nowadays, P&R times of high density FPGAs (more than 100K LUTs) can easily reach one day. Unfortunately, the SA approach is not suitable for P&R of this magnitude. An approach to investigate how to reduce the P&R time by employing high-capacity logic blocks has been presented in [17]. In [18], a placement based on a new problem formulation of maximum-bipartite matching (MBP) is presented. The MBP placement is  $30\times$  to  $75\times$  faster than VPR as long run-time is associated with SA for very large circuits (more than 100K CLBs). P&R time for each CLB is around 1 ms.

Regarding the routing step, a new approach is presented in [19]. When combined with a low-cost architecture change, this new approach results in a 34% reduction in router run-time, at the cost of a 3% area overhead. The routing time is around one second for the MCNC benchmarks.

Our work focus on P&R time reduction for partial reconfiguration. A similar work [4] proposes a JIT P&R that is 7.34 faster than VPR. Moreover, our approach as well the JIT approach leads to significant lower fragmentation of hardware resources at LUT level. However, while JIT P&R time is around seconds, our DF based algorithm executes in few milliseconds, allowing its usage in run-time applications.

## 7. CONCLUSIONS AND FUTURE WORK

We proposed a novel run-time placement and routing (P&R) algorithm for FPGA-based designs. Based on experimental results, our P&R performance is on average  $864\times$  faster, compared with the state-of-art VPR [7] tool. We perform each P&R process in milliseconds, making run-time partial reconfiguration feasible in practice. Moreover, we reduced the configuration memory by a factor of 4. We also proposed an implicit graph model to represent the FPGA. The

data structures have a direct correspondence to the physical and/or virtual configuration memories. Our approach allows partial reconfiguration at LUT granularity level. This feature reduces the fragmentation during the mapping and remapping of several applications on the fly.

Apart from the virtual FPGA used in our experiments, the proposed graph based approach is generic, and it will be adapted for a specific physical FPGA architecture as future work. On average, 35% of LUT connections are routed with minimal routing resources during the DF traversal performed by our algorithm.

Although our algorithm presents some penalties in terms of the maximum number of tracks per channel and the total of wire segments, additional improvements are feasible. Further work includes local optimization to replace a subset of nodes to reduce the required routing resources. In addition, better routing strategies for multicast hotspot signals should also be developed.

## 8. REFERENCES

- [1] K. Papadimitriou, A. Dollas, and S. Hauck, "Performance of partial reconfiguration in fpga systems: A survey and a cost model," *ACM Trans. Reconf. Technol. Syst.*, vol. 4, Dec. 2011.
- [2] M. Gericota, G. Alves, M. Silva, and J. Ferreira, "Run-time management of logic resources on reconfigurable systems," in *Design, Automation and Test in Europe Conference and Exhibition, 2003*, 2003, pp. 974–979.
- [3] M. Handa and R. Vemuri, "An efficient algorithm for finding empty space for online fpga placement," in *Proceedings of DAC*, 2004.
- [4] H. Sidiropoulos, K. Siozios, P. Figuli, D. Soudris, and M. Hubner, "On supporting efficient partial reconfiguration with just-in-time compilation," in *PhD Forum (IPDPSW)*, IEEE, 2012.
- [5] A. Ludwin and V. Betz, "Efficient and deterministic parallel placement for fpgas," *ACM Trans. Des. Autom. Electron. Syst.*, vol. 16, no. 3, 2011.
- [6] M. Hubner, P. Figuli, R. Girardey, D. Soudris, K. Siozios, and J. Becker, "A heterogeneous multicore system on chip with run-time reconfigurable virtual fpga architecture," in *Workshops and PhD Forum (IPDPSW)*, 2011.
- [7] VPR. (2009). [Online]. Available: <http://www.eecg.toronto.edu/vpr/>
- [8] R. Ferreira, A. Garcia, T. Teixeira, and J. Cardoso, "A polynomial placement algorithm for data driven coarse-grained reconfigurable architectures," in *ISVLSI*, 2007.
- [9] R. Ferreira, J. Cardoso, A. Damiani, J. Vendramini, and T. Teixeira, "Fast placement and routing by extending coarse-grained reconfigurable arrays with omega networks," *Journal of Systems Architecture*, vol. 57, no. 8, 2011.
- [10] S. J. E. Wilton, "Architectures and algorithms for field-programmable gate arrays with embedded memory," Ph.D. dissertation, University of Toronto, 1997.
- [11] M. Dehyadgari, M. Nickray, A. Afzali-Kusha, and Z. Navabi, "Evaluation of pseudo adaptive xy routing using an object oriented model for noc," in *Microelectronics, International Conference on*, 2005.
- [12] X. Lin, P. McKinley, and L. Ni, "Deadlock-free multicast wormhole routing in 2-d mesh multicomputers," *Parallel and Dist. Syst, IEEE Trans on*, vol. 5, 1994.
- [13] MCNC. (2010). [Online]. Available: <http://cadlab.cs.ucla.edu/kirill/>
- [14] V. Betz and J. Rose, "Vpr: a new packing, placement and routing tool for fpga research," in *FPL*, 1997.
- [15] M. Lin and J. Wawrzynek, "Improving fpga placement with dynamically adaptive stochastic tunneling," *CAD of Integrated Circuits and Systems, IEEE Transactions on*, vol. 29, no. 12, 2010.
- [16] M. Xu, G. Grewal, and S. Areibi, "Starplace: A new analytic method for fpga placement," *Integration, the VLSI Journal*, vol. 44, no. 3, pp. 192 – 204, 2011.
- [17] S. Y. Chin and S. J. Wilton, "Towards scalable fpga cad through architecture," in *Proceedings of FPGA*, 2011.
- [18] H. Bian, A. C. Ling, A. Choong, and J. Zhu, "Towards scalable placement for fpgas," in *Proceedings of FPGA*, 2010.
- [19] M. Gort and J. Anderson, "Reducing fpga router run-time through algorithm and architecture," in *FPL*, 2011.