

FPGA Architecture Optimization Using Geometric Programming

Alastair M. Smith, *Member, IEEE*, George A. Constantinides, *Senior Member, IEEE*,
and Peter Y. K. Cheung, *Senior Member, IEEE*

Abstract—This paper is concerned with the application of geometric programming to the design of homogeneous field programmable gate array (FPGA) architectures. The paper builds on an increasing body of work concerned with modeling reconfigurable architectures, and presents a full area and delay model of an FPGA. We use a geometric programming framework to show how transistor sizing and high-level architecture parameter selection can now be solved as a concurrent optimization problem. We validate the model through the use of simulation program with integrated circuit emphasis (SPICE) models and the versatile place and route (VPR) FPGA architecture simulation tool. Not only does the optimization framework allow architectures to be optimized orders of magnitude faster than previous work, but the combined optimization can lead to different architectural conclusions compared to conventional methods by exploring the coupling between the two sets of optimization variables. Specifically, we show that as delay takes more significance in the objective of the optimization, there should be more lookup tables in a logic block, whereas conventional techniques suggest that there should be fewer lookup tables in an FPGA logic block.

Index Terms—Convex optimization, field programmable gate array (FPGA), geometric programming, reconfigurable architectures.

I. INTRODUCTION

RECENT years have seen considerable evolution in the architecture of field programmable gate arrays (FPGAs). Each generation of commercial FPGAs contains new or refined routing, logic, memory, and embedded block structures. These architectural enhancements are the result of time consuming and expensive experiments, in which FPGA architects (in both industry and academia) use existing or new computer aided design (CAD) tools to map benchmark circuits to the architectures under investigation [1], [2].

Recent work, however, has suggested that this experimental approach can be supplemented by analytical techniques, in which FPGA architectures are modeled by relatively simple equations, and powerful optimization tools are used to “prune” the architecture space, allowing the FPGA architect to invest-

Manuscript received October 22, 2009; revised February 26, 2010. Date of current version July 21, 2010. This work was supported by the Engineering and Physical Sciences Research Council, U.K., under Grants EP/C549481/1 and EP/E00024X/1, and by Xilinx, Inc., San Jose, CA. This paper was recommended by Associate Editor E. K. Bazargan.

The authors are with the Circuits and Systems Research Group, Department of Electrical and Electronic Engineering, Imperial College, London SW7 2BT, U.K. (e-mail: alastair.smith@imperial.ac.uk; g.constantinides@imperial.ac.uk; p.cheung@imperial.ac.uk).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TCAD.2010.2049046

tigate a much wider range of architectures than previously possible [3]–[9].

Much of the effort in analytical modeling of FPGA architectures has been on ascertaining FPGA performance or utilization given a set of high-level parameters describing the logical fabric of the device [3]–[7]. However, it is also possible to model effects due to low-level details of the fabric such as transistor sizing [8], [9]. In the context of FPGA design, the application of transistor-level modeling techniques is particularly interesting, as there are a small number of different resource types which are replicated across the fabric. This means that the entire device architecture can be described by a concise set of low-level parameters, aiding the application of formal optimization.

In this paper, we combine high-level and low-level models, and show how they can be made amenable to geometric programming (GP), a form of convex optimization. This allows many logical parameters and most physical parameters of the device architecture to be optimized concurrently.

Specifically, the contributions of this paper can be summarized as follows.

- 1) Full details of the framework in [5] and [6] that allows concurrent optimization of both high-level (architectural) and low-level (transistor sizing) parameters.
- 2) Formulation of an area-delay model of FPGA fabrics as a geometric program.
- 3) Quantification of the area model accuracy using the versatile place and route (VPR) FPGA simulator and the delay model using the HSPICE simulation environment.

The remainder of this paper is organized as follows. Section II details related work in the field of FPGA architecture modeling and exploration, and provides a brief overview of GP within the domain of digital circuit design. Section III provides details of the FPGA architecture framework studied in this paper. The area and delay models used in this paper are given in Sections IV and V, respectively. These models are then mapped into a GP in Section VI. The GP formulation is studied in Section VII, in which we examine the accuracy of the models used and show how GP can be used to make new conclusions about FPGA architectures. The paper is concluded in Section VIII.

II. RELATED WORK

Many research works have been concerned with optimization of FPGA architectures. The majority of

an extensive review of GP in the context of circuit design. One of the particularly attractive features of GP is that it has excellent tractability. Interior point methods are used to solve GP with polynomial run-time with respect to the number of variables. By comparison, parameter-sweep methodologies for exploring design spaces have exponential run-time with respect to the number of variables. This means that the GP approach deployed in this paper offers a significant advantage in compute time over traditional methodologies.

Our previous work in [5] made the observation that some high-level models of FPGA fabrics also fit into the GP framework. For example, we made the observation that FPGA routing fabrics consist predominantly of multiplexers, and that the area and number of these multiplexers can be expressed as a GP, leading to modest area savings. Furthermore, in [8], we showed that by employing GP transistor sizing techniques such as those discussed above, area, delay and a combination of the two can be optimized in conjunction with some high-level architectural parameters. In this paper, we present these models in their complete form for the first time, and include comprehensive experimental verification of their accuracy.

III. MODEL FRAMEWORK

The modeling framework we present consists of a number of parts. In Section III-A, we define the target architecture style, which is based on lookup tables as the basic logic element, and with a number of variable parameters. The basic representation of the benchmark circuit is defined in Section III-B, which is used as the input to our model. The model has to use the information about each circuit in conjunction with the information about the architecture to estimate the number of computational resources used by the FPGA architecture, and the number of resources on the critical path. This information is used along with transistor-level details to obtain accurate area information in Section IV and timing information in Section V. Details of how the model is cast as a GP are given in Section VI.

A. Architecture Framework

Throughout this paper we assume an island-style FPGA in which an array of blocks is connected using tracks organized in horizontal and vertical channels with single-driver routing, as represented by VPR 5.0 [1]. There are six high-level architecture parameters that we study in this paper, which are summarized in the top half of Table II and are explained below. The logic blocks in the architecture consist of K -input lookup tables (LUTs) packed into tightly connected configurable logic blocks (CLBs), each with N LUTs and with I external inputs, as shown in Fig. 1(a). Further details of this logic block architecture are available in [2, Sect. 3.1.1]. A K -input LUT can be implemented using a K -level pass transistor multiplexer tree, as shown in Fig. 1. Further details of this structure can be found in [2, Sect. B.1.2].

The logic architecture parameters N , K , and I impact the number of logic blocks required to implement a circuit. For example, a 7-input LUT has a considerably larger logic density per block than a 2-input LUT, hence fewer CLBs will be required in an architecture containing 7-LUTs. Similarly, a

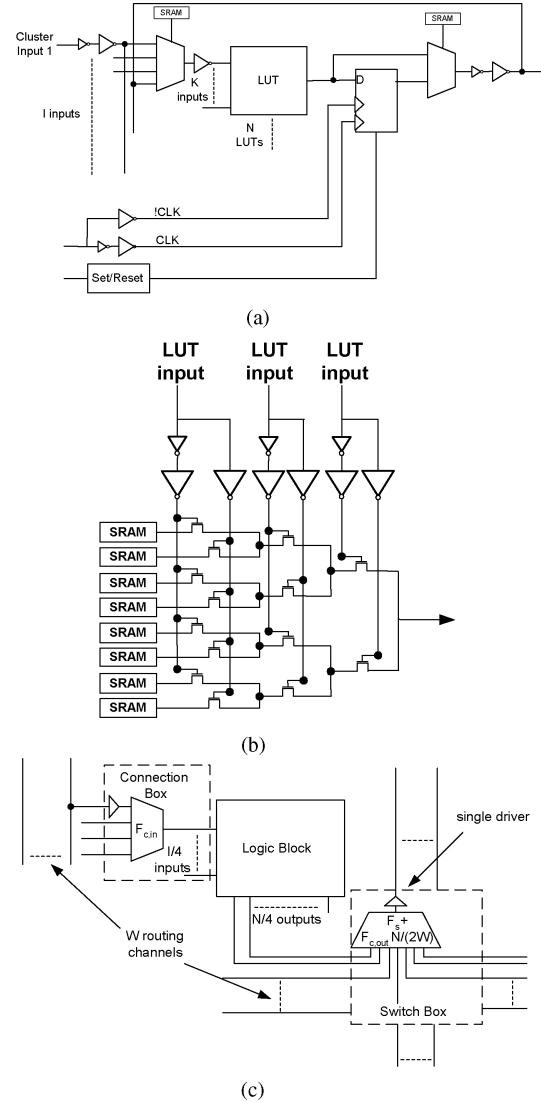


Fig. 1. Detailed view of FPGA architecture. (a) Structure of a CLB. (b) A 3-input LUT multiplexer tree. (c) Routing fabric of an FPGA.

TABLE II
MODEL PARAMETERS

High-level Architectural Parameters:	
K	Number of inputs per lookup table
N	Number of lookup tables per logic block
I	Number of inputs per logic block
$F_{c,in}$	Number of tracks that connect to each logic input pin
$F_{c,out}$	Number of tracks each logic block can connect to
F_s	Number of track end-points that connect to each track driver
Circuit Parameters:	
p	Rent parameter of a given circuit
n_2	Number of 2-LUTs in a given circuit
d_2	Depth of circuit netlist in number of 2-LUTs

large value of N implies fewer architecture blocks are required, as the CLBs have increased capacity. Due to the internal structure of a CLB, the tradeoff between these three parameters is not straightforward; the most recent study of this tradeoff is in [1].

The routing architecture is used to connect signals between logic resources, as shown in Fig. 1(c). Routing architectures in FPGAs have changed considerably since [2], and modern commercial FPGAs use single-driver routing [20]. These

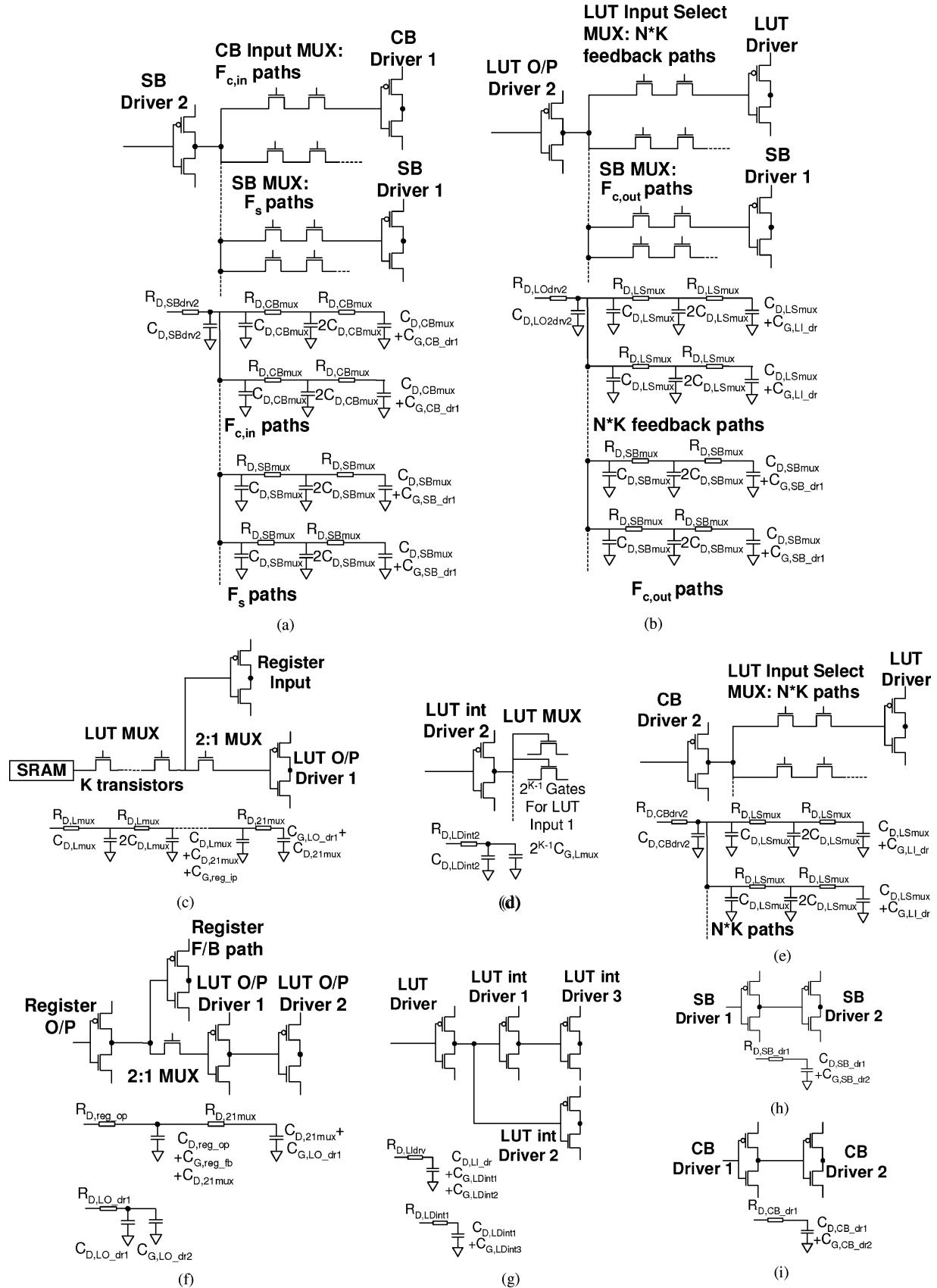


Fig. 4. RC delay models for circuit path, transistor variables are summarized in Table V. (a) RC model of paths from a switch box driver through a connection box multiplexer. (b) RC model of paths from a LUT output driver through the LUT input (feedback) multiplexer. (c) RC model of the path through a LUT to the register and output driver. (d) RC model of the path from a LUT input driver to the LUT gates. (e) RC model of the paths from a connection block through the LUT input multiplexer. (f) RC model of the paths emanating from the register on the output of a LUT. (g) RC model of delay paths in the LUT input driver. (h) RC model of the path between inverters in the switch box driver. (i) RC model of the path between inverters in the connection box driver.

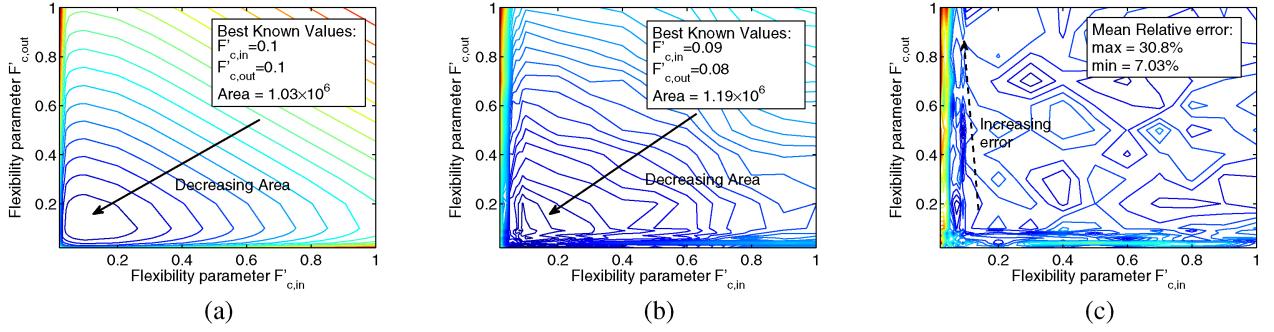


Fig. 5. Comparison of routing area model and experiment. (a) Area model implemented as a GP. (b) VPR experimental method. Each contour represents a difference of 30k transistors. (c) Mean relative error graph, where each contour represents a 1% difference.

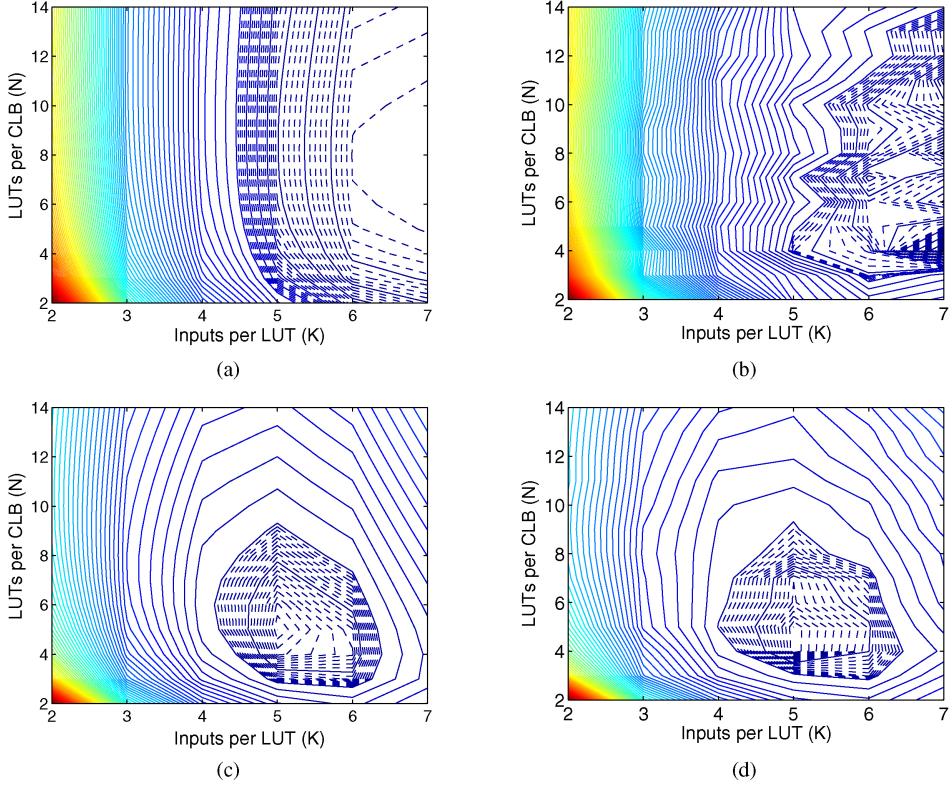


Fig. 5(c) shows the difference between the model and the VPR experimentation. The model is particularly inaccurate for extremely low values of $F'_{c,out}$. This is because the channel width model from [4] breaks down in this region: VPR must architect the routing multiplexers that connect to a very small number of tracks, for example one in a hundred, and this discretization of $F'_{c,out}$ is difficult to make consistent across the array. Nevertheless, the mathematical model demonstrates that the trend tracking is sufficient to detect the correct region of the design space for determining the best routing architecture.

The GP framework has a significant advantage in terms of compute time required. Each point in the VPR experimentation can take between 2 min and 4 h to obtain for a single benchmark dependent on benchmark size and interconnect complexity. Conversely, the GP framework takes approximately one minute to run, independent of the benchmark size, a speedup of up to 240 \times . Moreover, the optimal values for F_c in the

model can be found in a single run of the GP rather than a parameter sweep, reducing the execution time of any sweeping approach by two orders of magnitude.

B. Delay Model Verification

In this paper, we focus on verifying the circuit delay model. This is done through the use of HSPICE [27]. The SPICE model is constructed by specifying the transistors and connections in the critical path. The transistor sizings are specified by the results of the GP optimization. The optimally derived high-level architecture parameters such as channel width and the flexibility parameters must be quantized for the SPICE model in order to obtain multiplexers with an integral number of inputs. Certain parts of the critical path are replicated, such as the switch box to switch box delay, and are only specified once. Each of the constituent parts of the SPICE model are extracted so that they can be scaled up by the factors in (33) and com-

