

Peripheral Circuit Design for Field Programmable MCM Systems

Ivo Dobbelaere*, Abbas El Gamal*, Dana How[†] and Bendik Kleveland

Information Systems Laboratory, Stanford University, Stanford, CA 94305-4055

Phone: (415) 723-2525 Fax: (415) 723-8473

Abstract—A field programmable MCM architecture utilizing an array of modified FPGAs is proposed. Interconnections are provided by a fixed wiring network on the MCM, and by programmable interconnection frames on each FPGA. It is shown that full-swing CMOS peripheral circuits are faster than low-swing CMOS circuits. Buffering configurations for the interconnection frame which exploit the MCM performance benefits are selected and optimized. Bidirectional bus implementations using the frame are presented.

1 Overview

We present a new field programmable architecture for prototyping large designs using multiple FPGAs which offers excellent performance and economy while retaining the immediate turnaround of FPGAs at the system level.

For our packaging technology we have selected MCMs[1] because they offer the large pin count per chip necessary for high chip utilization in a partitioned design, and because the off chip delays are smaller than with PCBs. Due to the technological complexities of providing configurability at the MCM level we chose a fixed, statistically-determined wiring pattern on the MCM.

It would be convenient to use commercially available FPGAs (e.g. Actel[4] and Xilinx[5]), but these are not suited to our requirements. First, such FPGAs do not provide the high pin to gate ratio required when partitioning a design among multiple FPGAs. As a result the FPGAs are typically underutilized[6]. Secondly, since the I/O buffers of these chips are designed for general purpose use with PCBs, they do not give the better performance possible with MCMs. In addition, the fixed MCM wiring pattern assumed in our architecture imposes severe routing constraints which we solve by routing many signals through the FPGAs themselves[6]. Since the routing architectures of commercial FPGAs are optimized for local interconnects,

the delays incurred in routing signals from one pin to another through the entire chip are too large. Placing specialized switching chips[9] on the MCM would partially avoid this problem, but such centralization would increase average wirelength and lead to inflexibility in partitioning, placement and routing.

Therefore, we propose to modify the FPGAs to support quick connections from one pin to another, thereby uniformly distributing routing resources across the chips on the MCM. We do so by surrounding the FPGA logic core with an *interconnection frame* (Fig. 1) which supports fast thru-chip as well as chip-to-chip connections. In this paper we focus on the configuration and design of the interconnection frame in CMOS technology.

2 Frame Architecture

We assume the logic core of our modified FPGA to be similar to an existing FPGA (specialized functions could also be used). We also assume that the interconnection frame is programmable via a shift register. Connections leaving the MCM are similar to conventional PCB connections and will not be addressed in this paper. Among inter-chip connections within an MCM, connections between neighboring chips will predominate, so many can be implemented as fixed core-to-core connections. The remainder should be programmable using the interconnection frame.

The frame supports two basic switching patterns:

- Connect a signal from the core of chip A through a pin on chip A, across the MCM and to a pin on chip B. The signal is switched through B's frame and then connects to B's core (Fig. 2a).
- Connect a signal from the core of chip A through a pin on chip A, across the MCM and to a pin on chip B. The signal is switched through B's frame and then connects to another pin on B. The signal is then sent from this pin over the MCM to a pin on chip C which connects to C's core (Fig. 2b). In section 3 we optimize this pattern, since it is more critical than the previous pattern.

*Partially supported by FBI contract J-FBI-89-101

[†]Partially supported by an NSF Graduate Fellowship

The frame is split between the four corner regions of the chip, each containing a switchmatrix (Fig. 3a). Within a switchmatrix, all horizontal lines can connect to all vertical lines (Fig. 3b). Connections between the four switchmatrices and between the switchmatrix and the chip core are also provided.

With MCMs and flip-chip bonding, the direct interconnections between neighboring chips have a much lower capacitance than interconnections between remote chips. Hence, in considering thru-chip interconnection patterns, we will distinguish short interconnections between neighboring chips (5mm and about 1pF for AT&T), and long interconnections between chips that are not close (anything over 2cm); intermediate lengths could be classified either way (Fig. 1). We will take advantage of the MCM performance benefits by optimizing the switchmatrix for the expected mix of attached interconnections. Depending on the length of the connections to the switchmatrix, we have the four thru-chip interconnection cases optimized in section 3: short-to-short, long-to-long, short-to-long, or long-to-short.

At the inputs to the switchmatrix, small buffers may be present to enhance the incoming signal. NMOS transistors, whose gates are driven by 7V to eliminate voltage drop, route the 5V signals inside the switchmatrix. Additionally, output buffers may be needed to drive the off-chip interconnect. Since every signal path must be bidirectional, any output buffer must have an input buffer placed back-to-back; or the bidirectionality may be obtained by simply placing a pair of large programmable bypass transistors across the output buffer (Fig. 3c).

3 Frame Design

3.1 MCM Model

Throughout this section we assume parameters for AT&T's MCM technology[8], with a wiring capacitance of about 1pF/cm, a characteristic impedance between 50 and 100 Ω , and a wire resistance of 20m Ω /cm. We chose 50 μ m \times 50 μ m flip chip bonding pads, with a capacitance of 0.5pF and an inductance of 0.05nH.

3.2 Buffer Circuits

In our application differential lines are unacceptable because of the higher pin count. The advantage of differential lines is mainly noise immunity; the gain in speed is not compelling. In CMOS, there is a choice between low swing ECL-like circuits[3], low swing CMOS[2], and standard full swing CMOS. The low swing ECL-like circuits are not interesting for our

application because of the high DC power consumption. Low swing CMOS circuits consume less power, but simulations show them to be slower, mainly due to the low swing to high swing conversion at the receiver end. Our final circuit choice is full swing CMOS.

3.3 Path Design

We simulated each of the four interconnection cases with HSPICE — Fig. 4 summarizes the signal path we simulated in the short-to-short case; the others were similar. The input buffer is a two stage tristate buffer with adjustable sizing. The output buffer is a three stage tristate buffer with sizing such that MCM transmission line behavior is insignificant (the driver sizes are $W_p=400\mu$ m, $W_n=200\mu$ m). For each case, we considered the placement of either a bypass transistor pair or a buffer at the switchmatrix input, and the placement of a buffer at the matrix output, for switchmatrices varying in transistor widths and in number of short and long paths connected. The propagation delay and the rise and fall times of the output signals were measured for T=70°C. We have measured rise times to the 80% point since 7V NMOS switches degrade the rising edge past 4V, which is not critical for our circuits.

The first simulations indicate that input buffers do not improve the delay for *any* of the interconnection cases. In fact, for different interconnection lengths, input buffers added 0.5ns to the delay even for the fastest input buffers, so we will not use them. However, each path must support signal propagation in both directions, so where there is an output buffer (as determined below) we will place an input bypass transistor pair. Fortunately, a large input bypass transistor pair will not significantly increase the total propagation delay.

Similarly, the best delay with no output buffer in the short-to-short case is at least 20% smaller than the best delay with a buffer, whether or not the signal previously passed through an input bypass transistor pair. The risetime is longer without an output buffer, but not unacceptably so (in the range of 2.5ns). No output buffer should be used in the short-to-short case.

In the case of a long-to-long interconnection, the presence of an output buffer does not significantly affect the best delays. However, since the output buffer does significantly improve the risetime it must be included. In addition, the presence of an output buffer makes small delay times possible without large switch transistor sizes. As previously argued, with the output buffer we must include a large input bypass transistor pair for bidirectionality.

The considerations for the short-to-long case are

| Case | Input | Output |
|----------------|--------------------|------------|
| short-to-short | no bypass | unbuffered |
| short-to-long | no bypass | buffered |
| long-to-short | bypass transistors | unbuffered |
| long-to-long | bypass transistors | buffered |

Table 1: Optimized Path Configurations

very similar to the long-to-long case above. Consequently the short-to-long case will not use an input buffer but it will use an output buffer.

The long-to-short case is the reverse of short-to-long, so there must be an input bypass transistor pair. Omitting the output buffer does decrease the delay (by 10%), but this would increase the risetime by 3ns. However, we have decided that for consistency with the implementations already presented for the other three cases, the increased risetime is acceptable (since the switching delay from propagation and risetime is still less than a long-to-long interconnection delay), and so we will not use an output buffer in the long-to-short case.

Table 1 summarizes the path configurations for each case which provide the most effective system solution. From all the arguments we see that we should place an input bypass transistor pair and output buffer combination at a pin if and only if it is attached to a long interconnection. This is the simplest rule which will effectively implement all four thru-chip interconnection cases.

3.4 Transistor Sizing

The switch transistor *widths* can be chosen so that propagation delay and layout area are minimized. For the short-to-long and long-to-long cases, the buffer shields the switchmatrix from the capacitance of the external interconnection and therefore the switch transistor widths have little influence on the delays as long as a certain minimum width is maintained. The paths without an output buffer show a much greater delay dependence on width; the delay minimizing widths are a factor 3 to 4 larger than for the buffered case.

4 Programmable Inter-chip Buses

The configuration proposed in Table 1 provides an efficient way of implementing programmable inter-chip buses. For short interconnections only the I/O circuits at the pin of the core should be bidirectional, because of the absence of directed circuitry in the switchmatrix (Fig. 6a). Buses with long interconnections that are not time critical can be implemented using one interconnection by permanently turning on the input by-

pass transistor pair in the buffer-pass transistor combination at the frame pin. For a critical long bus, a second interconnection is required, to control the direction of the frame I/O circuit (Fig. 6b).

5 Conclusion

We have investigated alternative implementations of an interconnection frame for FPGAs used on an MCM substrate with fixed wiring. We determined that conventional CMOS full swing signals would afford the best inter-chip performance rather than some form of low swing signalling. Surprisingly, I/O buffers were not needed in many cases due to the decreased MCM capacitance and switches controlled by 7V signals. Fig. 5 shows simulated delay versus switchmatrix dimensions of our chosen implementation for the short-to-short case. Although we have yet to determine what mix of short and long connections our prototyping architecture will require, this graph and similar ones for the other cases show that neighboring connection delays of 3.5ns and long-range connection delays of 6ns are attainable when routing a signal through a chip with the interconnection frame. Such speeds will make our prototyping architecture feasible. In addition, the interconnection frame provides a simple implementation of inter-chip programmable buses.

References

- [1] R. Tummala and E. Rymaszewski, *Microelectronics Packaging Handbook*, Van Nostrand Reinhold, 1989.
- [2] T. Knight, Jr. and A. Krymm, "A Self-Terminating Low-Voltage Swing CMOS Output Driver," *IEEE J. Solid State Circuits*, 23-2, April 1988, pp. 457-464.
- [3] B. Chappell *et al*, "Fast CMOS ECL Receivers with 100-mV Worst-Case Sensitivity," *IEEE J. Solid State Circuits*, 23-1, February 1988, pp. 59-67.
- [4] A. El Gamal *et al*, "An Architecture for Electrically Configurable Gate Arrays," *CICC* 1988, pp. 15.4.1-4.
- [5] H. Hsieh *et al*, "A 9000-Gate User Programmable Gate Array," *CICC* 1988, pp. 15.3.1-7.
- [6] S. Walters, "Computer-Aided Prototyping for ASIC-Based Systems," *IEEE Design and Test of Computers*, June 1991, pp. 4-10.
- [7] K. Tai, "Si-on-Si MCM Technology and the Initiation of a University MCM Program," *Multichip Module Workshop*, Santa Cruz, CA, March 1991, pp. 10-13.
- [8] J. Burr *et al*, "System-wide Energy Optimization in the MCM Environment," *Multichip Module Workshop*, Santa Cruz, CA, March 1991, pp. 66-83.
- [9] "Startup Aptix details its FPIC architecture," *Electronic Engineering Times*, January 6, 1992, issue 674.

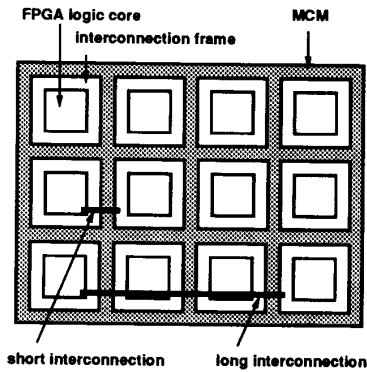


Figure 1: Multiple FPGAs placed on an MCM.

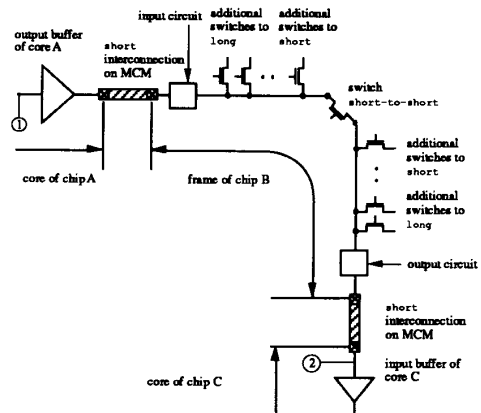


Figure 4: short-to-short signal path from 1 to 2.

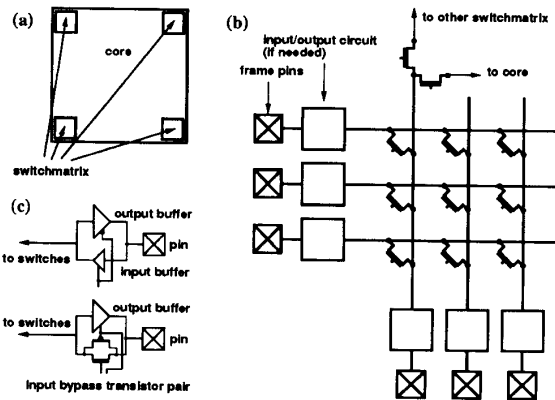


Figure 2: (a) Switchmatrix placement. (b) Switchmatrix structure. (c) I/O circuits.

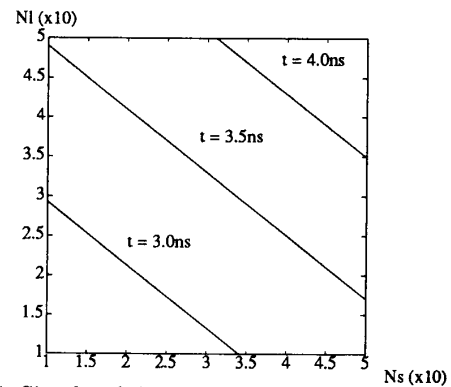


Figure 5: Simulated short-to-short delay versus switchmatrix sizes N_s and N_l .

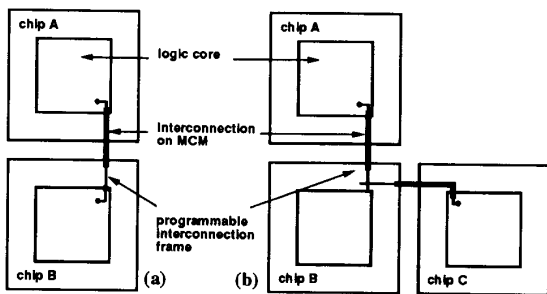


Figure 3: The frame supports two main switching patterns: (a) Connection from core A to core B, with switching in frame B. (b) *Thru-chip* interconnection from core A to core C, with switching in frame B.

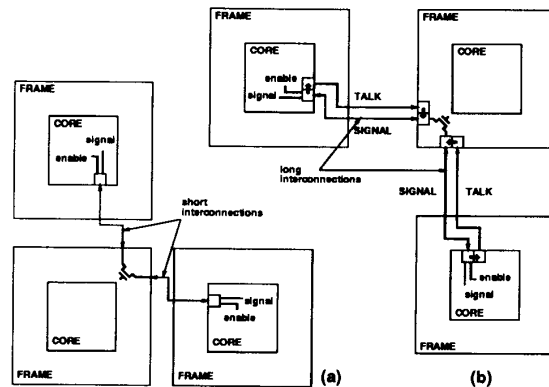


Figure 6: Bus implementation: (a) short interconnections. (b) long interconnections, fast implementation.