

Utilizing Parallelism of TMR to Enhance Power Efficiency of Reliable ASIC Designs

Hagen Sämrow, Claas Cornelius, Jakob Salzmänn, Andreas Tockhorn, Dirk Timmermann

Department of Electrical Engineering, University of Rostock, Rostock, Germany

Email: {hagen.saemrow,claas.cornelius,jakob.salzmänn,andreas.tockhorn,dirk.timmermann}@uni-rostock.de

Abstract — Due to aggressive scaling, reliability issues influence the design process of integrated circuits more and more. A well known technique to tackle these issues represents Triple Modular Redundancy (TMR). It strongly improves reliability of a design at the expense of at least tripled area and power consumption. In this contribution, we propose an enhanced TMR approach that significantly decreases the power overhead of conventional TMR designs. Therefore, the control logic was modified so as to switch between a TMR mode and a parallel mode. This parallel mode allows the circuit to operate with decreased frequency without losing performance by taking advantage of the parallelism offered by the tripled design. Achieved results of investigations on the ISCAS benchmark circuits show power savings of up to 50 % with a small reliability penalty compared to a conventional TMR approach for permanent failures. We also propose strategies how to utilize both operating modes in order to balance the design concerning reliability and power consumption requirements at runtime.

Keywords: *Circuit design, Reliability, Triple Modular Redundancy, Power Consumption, Power-Aware design*

I. INTRODUCTION

Reliability concerns are one of the recent major issues in VLSI design. Besides the rising power density due to scaling effects, transistors and interconnects suffer from an increased sensitivity to different kinds of failures during system operation due to the proceeding miniaturization. What makes the situation even more severe is that reliability issues worsen with non-ideal scaling and increased transistor count as well as power density [1]. These mechanisms lead to increasing current densities and electrical fields as well as rising temperatures and temperature gradients resulting in long-term reliability problems explained in the following [2].

Electric current leads to a transport of wire material alongside the current flow resulting in short connections (hillocks) – especially to ground or supply wires – and increased wire resistances and even opens (voids). This effect is called electromigration and it deteriorates with rising current density and temperature [3]. Further stress mechanisms are self-heating where oxide layers hamper the heat flow [5] and thermal cycling, whereas both effects cause malfunctions of integrated circuits due to fast temperature changes [6]. Another effect is the Time-Dependent Dielectric Breakdown (TDDB), which leads to short connections through transistor gate oxides

and hence, defective transistors caused by tunneling currents [7]. TDDB is negatively affected by high gate voltages and temperatures, though it results rather in delay failures than in immediate errors. However, the timing behavior of a circuit can be disarranged leading to erroneous outputs in the presence of diverse defective gates [8]. Other important mechanisms, which also degrade the gate oxide, are the Negative Bias Temperature Instability (NBTI) that leads to an increased threshold voltage of p-MOSFETs as well as the hot carriers that are trapped electric carriers causing defects in gate oxides [9][10]. Both effects are also dependent on temperature and affect the device characteristics resulting in delay failures and finally, permanent functional errors.

To overcome the challenges of the rising complexity of integrated circuits, one of the key priorities in the future is to support tool assisted insertion of reliability mechanisms [11]. In addition, since full functional system tests have already become unfeasible, on-chip monitoring mechanisms are required during system operation to detect and preferably also to correct errors.

Triple Modular Redundancy (TMR) is an established mechanism to enhance reliability of integrated circuits. Commonly, three identical units of the same function operate with the same inputs. Their results are merged by a voting unit that forwards the final result based on a majority decision [12][13]. A lot of research took place to improve the basic design idea [14][15] and showed the efficiency for reliability improvements [16]. A considerable drawback of TMR is the overhead of area and power consumption. Whereas area recently appears to be a less important design parameter, power concerns arise, similar to reliability problems, with every next technology generation. Besides the obvious impact on rising energy costs or larger accumulators in commercial applications, power consumption leads to heat dissipation that impacts the reliability harmfully and persistently [4]. The common Arrhenius relationship makes imposingly clear that lifetime of integrated circuits decreases exponentially with temperature [17].

Chandrakasan showed that parallelized data paths operating with decreased frequency and supply voltage decrease the dynamic power consumption at the cost of increased area and leakage currents [18]. Based on this, our approach combines the principles of redundancy and parallel designs to provide an

enhanced TMR design that consumes less power than the conventional TMR. Thereby, the original integrated circuit will be duplicated twice to construct a basic TMR design, whereas the whole design operates in two modes. On the one hand, a conventional TMR mode detects and masks permanent faults. On the other hand, a power saving mode performs the original function by utilizing the inherent parallelism of the TMR design in order to reduce the operating frequencies. However, the ability of TMR designs to prevent soft errors is lost in this mode. Furthermore, a tool was developed to provide a complete automated design process.

II. DESIGN OF THE ENHANCED TMR

The enhanced TMR design we propose operates in two different modes as long as no fault was detected within the redundant modules. First, every module executes the original circuit function on different data sets creating a threefold parallel system as a whole. Therefore, the operation frequency of the modules during this Non-TMR (NTMR, or parallel mode) mode can be reduced to one-third of the original one without sacrificing data throughput. This phase is regularly interrupted by the TMR mode whereas all modules process the same inputs concurrently with the original frequency. This allows observing if the redundant modules work identically, thus correctly. If a fault is discovered during this TMR mode, a fault_out signal is generated by the error detection unit. When higher level units declare this fault as permanent, the design remains in the TMR mode to perform as a basic TMR. In so doing, physical faults in the modules can be masked on the basis of the majority voting at the outputs.

To allow a controlled switching between the two aspired modes, an accurate assignment of the input signals to the tripled modules as well as correct forwarding of the output signals of the three modules to the design's output have to be ensured. Therefore, a demultiplexer unit at the input and a multiplexer unit at the output have to be designed. These units consume more cells and chip area than in case of the conventional TMR because the mode of parallel operation and the error detection are added. A schematic of the whole design is depicted in figure 1.

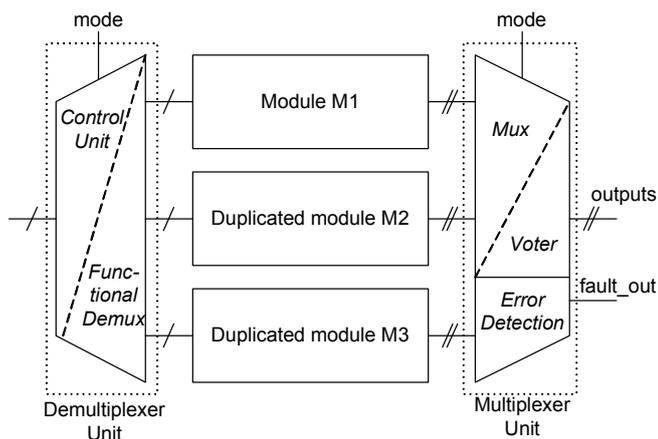


Figure 1 Schematic of the enhanced TMR design with the additional control logic to operate both in a power-saving parallel mode and in a conventional reliable TMR mode

A. Demultiplexer Unit

This unit assigns the data inputs to the appropriate TMR modules by generating enable signals for the input registers. Hence, the control unit interprets the signal 'mode' to decide whether the TMR mode or the NTMR mode is requested. The TMR mode requires enabling the input registers of all three modules concurrently. By contrast, operating in the NTMR mode requires that the registers of only one module are enabled at a time when valid data signals arrive. This is accomplished by switching the enable signals on/off for the appropriate modules. With the help of a small counter, one particular module is chosen by executing a simple round robin. Thus, only every third set of data signals is forwarded to a specific module. This mechanism results in a low-power parallel NTMR mode because the input registers switch only if the enable signal denotes a new set of data for the particular module. Hence, the required operating frequency is just a third of the original one providing the opportunity to use a lower supply voltage for the TMR modules.

B. Multiplexer Unit

The multiplexer unit has to manage three tasks. Firstly, the outputs of all three TMR modules have to be voted during the TMR mode. That is to say, the consistent results of at least to modules are forwarded to the output. This is executed bit by bit by a modular voter design. Secondly, in the NTMR mode the alternating outputs of the three TMR modules have to be forwarded sequentially to the actual design outputs – i.e. when the parallel operation takes place. This is accomplished based on the signal 'mode'. Two approaches are feasible to switch between both modes. Figure 2 depicts the associated schematics. On the one hand, multiplexers can be used besides

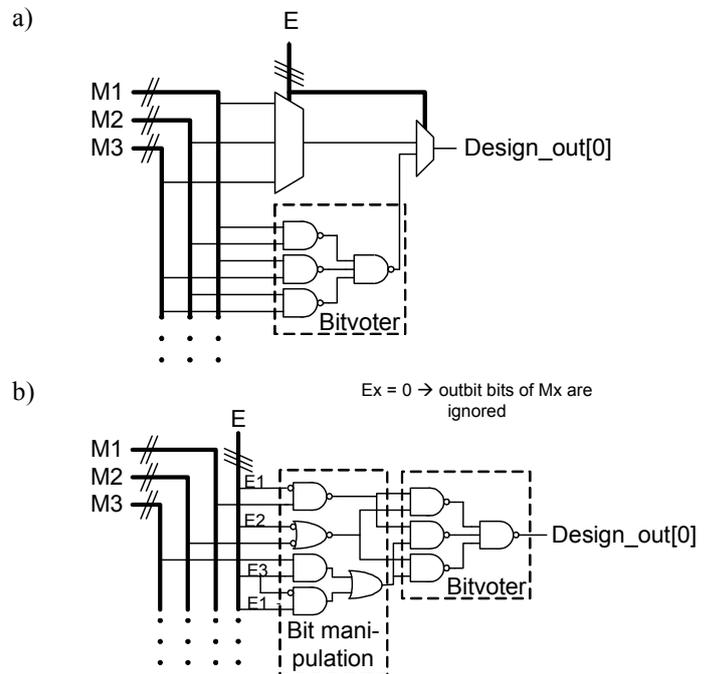


Figure 2 Feasible combinations of the voting and selection unit for one bit: a) using multiplexers b) manipulating the output bits prior to the actual bitvoter

the voter to forward the according output signals (see figure 2 a). On the other hand, the input signals to the voter can be manipulated in such a way that the voter actually forwards the signal of one particular module (see figure 2 b). Thereto, the control signal ‘E’ is used to dynamically forward one of module outputs in case of the NTMR mode or to bypass the original values to the voter in case of the TMR mode. This second implementation results in lower area overhead due to fewer transistors but if the critical path runs through the voting unit the propagation delay of the design can suffer.

Lastly, the multiplexer unit also contains the error detection unit that generates a `fault_out` signal if a fault occurs during the TMR mode. This is essential because the TMR mode is used to check the correct operation of the design. Thus, if an error signal is generated and the failure is declared permanent, the design has to switch to the TMR mode permanently in order to still provide correct outputs. As it is illustrated in figure 3, the generation of the error signal (fault-out) benefits from the TMR approach. This is because when faults occur in the XOR array at least two output bits of the XOR cells switch to a logic one. Therefore, in order to entail an incorrect `fault_out` signal at least two stuck-at-zero faults have to occur at locations belonging together. Hence, this type of implementation provides a certain amount of inherent reliability. However, due to the OR-tree, a large number of output bits can cause the error detection unit to become the critical path because all input values to the OR-tree have to be merged to the single `fault_out` signal which indicates a defective module.

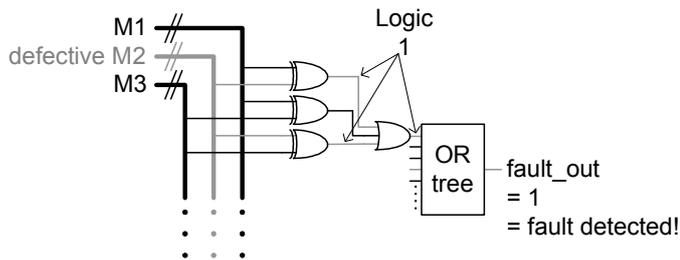


Figure 3 Schematic of the error detection unit that provides some means of inherent reliability due to the array of XOR gates

C. Switching Strategies

Different strategies have to be evaluated that determine how often and how long the two operation modes are to be executed. Thereby, two parameters are adjustable to define such a switching strategy. First, f_{TMR} represents the frequency how often the TMR mode is switched on during operation time. Second, d_{TMR} declares the duration of the TMR mode during one cycle of TMR and NTMR mode. It is obvious that more failures can be discovered with an increasing frequency f_{TMR} and a longer duration d_{TMR} because the design remains absolutely longer in the reliable TMR mode – that is $t_{TMR} = f_{TMR} \cdot d_{TMR} \cdot t_{OP}$, whereas t_{OP} is the operation time. On the other hand, the power consumption is raised with an increased t_{TMR} because the three modules execute the same

function redundantly. This leads to the conclusion that f_{TMR} as well as d_{TMR} should be chosen as small as possible with the intent of reduced power consumption. Or vice versa, large f_{TMR} and d_{TMR} result in a reliable but power-hungry operation.

The usage of the NTMR mode entails two additional clock cycles of latency for the output signals in comparison to common parallelized circuits. This is because every switching between the two modes requires two clock cycles to compensate for the different frequencies. This results in a performance loss which is linearly dependent on f_{TMR} . To reduce the performance loss to a minimum, f_{TMR} should be chosen as small as possible and d_{TMR} as large as possible to retain a reasonable t_{TMR} .

A further trade-off has to be considered. If the regular data patterns are provided to the inputs during the TMR mode (i.e. during d_{TMR}) no additional performance is lost. However, special test patterns can ease the examination of the redundant modules. In this case, the operation of the regular data signals is delayed, which implies a performance loss. On the other hand, the duration of the TMR mode d_{TMR} can be shortened because the special test patterns cover a better range of critical test cases. Besides, additional storage or the generation of such test patterns has to be provided which is why this option was not implemented here.

III. SIMULATION SETUP

To verify our enhanced TMR designs, we performed different simulations to evaluate reliability and power consumption at gate level. ISCAS benchmark circuits were implemented with cells of an industrial 65 nm gate library. Although delay failures also result from life time reliability issues, we chose stuck-at-one/zero faults to model permanent failures. These faults were induced during the simulations of the various designs. To analyze the behavior of the designs in the presence of faults, we manipulated the netlists by inserting the stuck-at faults to all possible nets of the design. Thereby, the locations of failures were randomly chosen assuming the same failure rate for every net. This leads to an area dependent failure rate for a complete system. More precisely, a small design is less often affected by a fault compared to a large design over a given period of time. In our simulations, the failure rate has been defined according to the enhanced TMR design – termed ENH in the following diagrams. This means that on average one fault is induced per time step in this case. Accordingly, the unprotected design (termed REF) suffers on average only every third time step from a fault because it is about three times smaller. Similar measure hold also true for third design, which is the conventional TMR (termed TMR). These three types of designs were implemented for miscellaneous ISCAS designs. To evaluate the reliability of the designs, the fraction of working units is plotted over time. This kind of reliability $R_{SYS}(t)$ of a system represents the probability to perform as desired until time t . Closely related to the probabilistic term for reliability is the Mean Time To Failure ($MTTF_{SYS}$) which is the average time a system operates until it fails. It is equal to the expected lifetime if the system cannot be repaired. It can be calculated by [19]:

$$MTTF_{SYS} = \int_0^{\infty} R_{SYS}(t) dt \quad (1)$$

Simulating the power consumption has been executed in the absence of faults and with an accurate activity evaluation for every net. Thereby, f_{TMR} and d_{TMR} were adjusted according to the various simulation scenarios. Power consumption values have been derived from simulations based on toggle rates of gate level netlists.

IV. DISCUSSION OF RESULTS

A. Design Overhead

Figure 4 depicts the results for area and delay of the enhanced designs in relation to the conventional TMR circuits. While the simple reference designs are ca. 10 % faster and ca. 65 % smaller than the conventional TMR designs, the results for the area and timing overhead of the modified designs show a diverse picture. Here, the smaller circuits (e.g. c432, c499) as well as the designs with many output bits (e.g. c2670 with 140 output bits, c5315 with 123 output bits) exhibit a larger design overhead. These area costs comprise the overhead due to the control signal generation and the additional area due to the voting circuits and error detection units. The latter units are linearly dependent on the number of output bits because every output bit has to be considered for voting and error detection. Moreover, the propagation delay t_{delay} of fast circuits (e.g. c499, c1908 and c2670) is also more strongly affected than the delay of the slower ones because the additional delay due to error detection and voting is roughly constant. Therefore, the relative delay increase is larger in these cases. More precisely, the additional delay partly depends on the number of outputs. This is because fewer output signals entail smaller OR-trees, which affect the propagation delay less severely. Summarily, although all enhanced TMR designs show a certain design overhead, the unhasty circuits with a small number of output bits (e.g. c6288) come fairly close to the design parameters of the conventional TMR designs.

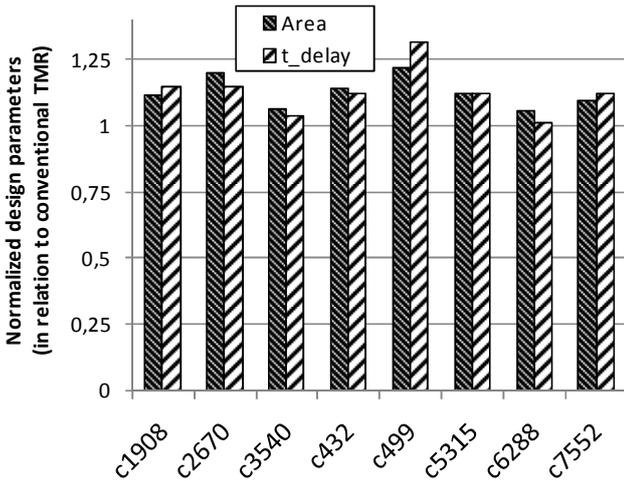


Figure 4 Design overhead of the enhanced designs compared to the conventional TMR approach for different ISCAS circuits

B. Reliability

The initial results on reliability have been obtained by running the designs in the TMR mode for the whole simulation time. This allows evaluating how efficient the different designs still operate in the presence of failures. Figure 5 depicts the reliability $R(t)$ for the circuit c5315 exemplarily. As time goes by (simulated with constant time steps here), more and more faults occur and force erroneous signals at the outputs of the various designs. This means that already after three time steps the unprotected reference design (without any kind of TMR) will only work with a probability of less than 50 %. At this time one fault emerged on average. Interestingly, some of the unprotected designs still work in the presence of faults due to inherent capabilities to mask some of the faults. The conventional as well as the enhanced TMR designs reach the same level of reliability after five time steps. However, due to the larger area of the TMR designs, they suffer at this point already from five faults on average. Beyond that, a further curve (ENH_2) is plotted in the figure as an outlook, which depicts how the reliability is affected when the control logic of the TMR designs is considered fault-free. This is not realistic, but it shows the potential of how beneficial it is to protect the control logic by other means.

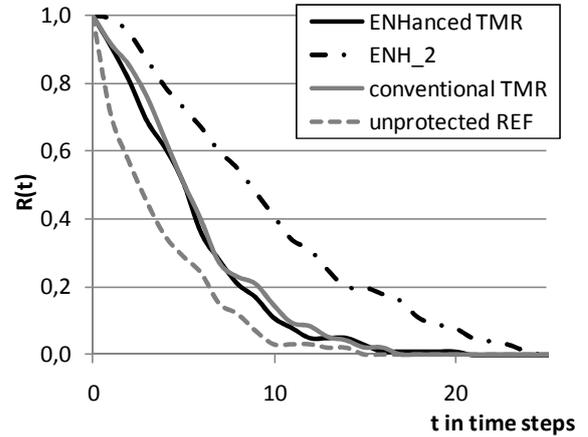


Figure 5 Reliability $R(t)$ for the c5315 circuit implemented as the unprotected REFERENCE, the conventional TMR and the ENHanced, whereas ENH_2 assumes fault-free control logic

In order to compare the miscellaneous reliability results, the MTTF was calculated for all designs based on the reliability according to (1), but with respect to the discrete simulation steps:

$$MTTF \approx \sum_{i=1}^a \frac{(R(i) + R(i-1))}{2} \quad (2)$$

With a being the first time instance where all test cases of a particular design failed, whereas $R(i)$ is the fraction of designs which work correctly at time instance i . Figure 6 plots the MTTF in relation to the unprotected reference design for the diverse ISCAS circuits. As it can be seen, the reliability of the

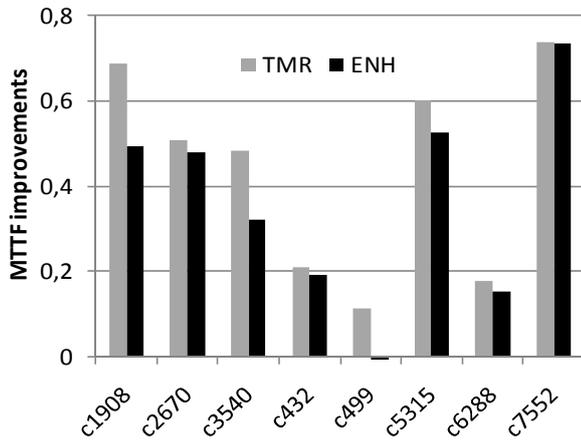


Figure 6 MTTF improvements of the pure and enhanced TMR designs compared to the references

enhanced TMR designs is in fact slightly smaller compared to the conventional TMR designs due to the larger control overhead. However, the differences between the conventional and the enhanced TMR designs range only between 1 % and 11 %. In case of the c499 design, the MTTF even deteriorates. This is because the area overhead is very large in comparison to the actual design and there are plenty of output bits. Consequently, the chance of a fault to occur increases similar to the area increase and exceeds the possible gains of the TMR. Summarizing the findings, with the focus on reliability in these analyses, it could be shown that larger designs suffer relatively less from the overhead due to the TMR. Especially the largest design (c7552) shows significant reliability enhancements of more than 75 %. Furthermore, the MTTFs of the enhanced TMR designs were slightly smaller compared to the conventional TMR designs.

C. Power analyses

Figure 7 depicts the power consumption of the enhanced TMR designs in relation to the power of the conventional TMR. These results are given as a function of the TMR rate – that is t_{TMR} / t_{OP} which is the ratio of time that the designs are in TMR mode to the whole simulation time. For instance, a TMR rate of zero means that the design solely operates in the NTMR (or parallel) mode. By contrast, a TMR rate of one denotes that the design runs exclusively in the TMR mode. The power savings are linearly dependent on t_{TMR} with a minimum at the lowest possible t_{TMR} . Similar to the results for area and propagation delay, the larger designs with few output bits benefit considerably more than the smaller designs – such as designs c6288 and c3540. The reasons are again that few output signals require less control overhead and that the necessary overhead is rather small compared to the large designs. Strictly speaking, the proposed modified TMR does not pay off for small designs or designs with many outputs bits. Note the designs c499, c432 and c2670 in figure 7.

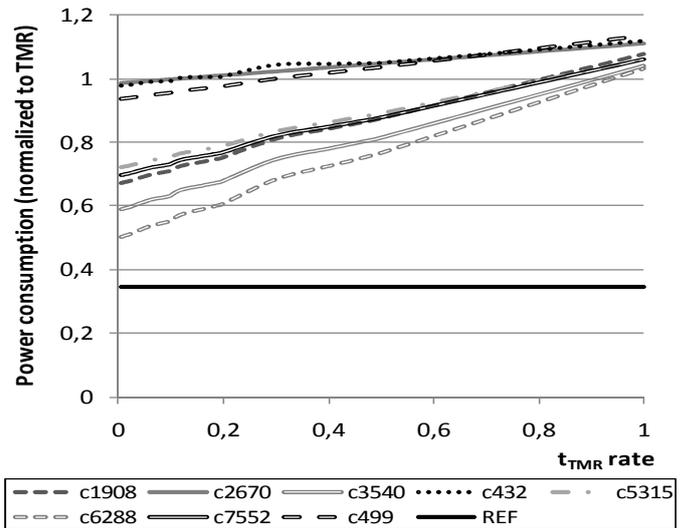


Figure 7 Power consumption depending on the t_{TMR} rate for the enhanced TMR designs normalized to the conventional TMR

As a result, these designs consume more power than the conventional TMR designs as soon as the TMR rate is increased. However, all the other designs allow significant power savings compared to the conventional TMR as long as the TMR mode is not operated longer than 75 % of the time. For a very small t_{TMR} , power savings of up to 50 % were achieved for the circuit c6288 which is a 16x16-multiplier. This amount of power dissipation is 45 % higher compared to the unprotected reference design (i.e. without TMR), which is also depicted as a constant curve in figure 7 (see curve REF). Hence, the conventional TMR necessitates about three times more power than the unprotected reference design, which attributes to the tripled redundancy.

D. Switching strategies

Since the TMR mode is not run all the time in order to save power, it is possible that faults remain undetected. Therefore, this section investigates the impact of the switching strategy on the number of undetected faults.

Figure 8 shows the exemplary results of according simulations for the c6288 design. As a start, one physical fault is injected into the design. Following, the number of logical errors is counted before the control logic detects the error and compensates for the physical fault by running the TMR mode. Thus, the numbers of undetected errors are depicted in the figure against the duration that the TMR mode is executed. Furthermore, various functions are given that refer to different TMR rates – that is the ratio of TMR time in comparison to the execution time. All simulations were started in the TMR mode which resembles an initial self-test. Two important things can be determined. First, if the TMR rate is high enough – for instance, larger than 0.1 = 10 % – faults in the modules can quickly be detected during the operation. Second, an appropriately high duration of the TMR mode d_{TMR} – here, more than 60 clock cycles – enables the

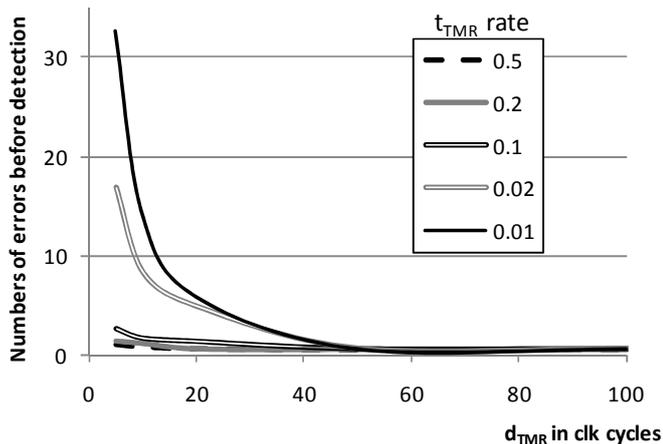


Figure 8 Average number of design errors which occur in an enhanced c6288 design before the errors have been detected from the circuit

enhanced designs to detect faults, even when the TMR rate t_{TMR} is very small. This observation matches very well with the requirements for performance that requires a small f_{TMR} and a large d_{TMR} (see chapter II. C). Moreover, a small TMR rate also advantages the power consumption as shown in the previous section.

V. FURTHER DESIGN IMPROVEMENTS

Further enhancements are conceivable. First, as mentioned in chapter 1, parallelism allows reducing the supply voltage in combination with the frequency adjustments. Although, such an approach seems promising, the design of an efficient level converter that shifts between the different voltage levels as the operating modes switch is a challenging task. Second, the control mechanisms – (de-)multiplexer, voter, error detection – should also be made more reliable because these design units are the weak points. For instance, doubling the demultiplexer unit necessitates little additional area and would allow to at least detect errors there. However, to harden the multiplexer/voting unit is rather difficult or costly as the area needs are linearly dependent on the number of output bits.

Apart from that, implementing other modes than the TMR and NTMR could also be useful. For instance, if the module which generates an erroneous output is detected, it is beneficial to completely turn off that module. Consequently, another mode would operate then similar to Duplication with Comparison (DwC) [20]. That is to say, the two remaining modules operate in parallel with halved frequency to save power or they work at full speed on the same data set to detect errors.

VI. CONCLUSION

This contribution identified the needs for improvements of lifetime reliability. Therefore, an enhanced TMR design was introduced which operates in two modes. First, the NTMR mode exploits the modular redundancy to run the three modules concurrently with reduced frequency and thus decreased power consumption. Second, the TMR mode works

just as a conventional TMR design with a voter and serves as a self-test to monitor the tripled modules. Due to the additional control logic of the enhanced TMR design, area needs are slightly higher and an additional latency arises. However, power savings of up to 50 % compared to a conventional TMR design can be achieved while errors can still be detected and masked. In a nutshell, the proposed enhanced TMR design performs better for designs with larger area and few outputs.

Moreover, prospective improvements were suggested that incorporate further operating modes and conditions. Since, the first results herein are promising, further works will aim at considering additional failure models and application scenarios.

REFERENCES

- [1] G Srinivasan, J., Adve, S., Bose, P. and Rivers, J., "The Impact of Technology Scaling on Lifetime Reliability", In Proc. of DSN, 2004.
- [2] Mc Pherson, J.W.: "Reliability trends with advanced CMOS scaling and the implications for design", In Proc. of Custom Integrated Circuits Conference, 2008.
- [3] Srinivasan, J.; Adve, S.; Bose, P.; Rivers, J.: "The Case for Lifetime Reliability-Aware Microprocessors", In Proc. of 31st International Symposium on Computer Architecture (ISCA), 2004.
- [4] O. Semonov, et al., "Impact of Self-Heating Effect on Long-Term Reliability and Performance Degradation in CMOS Circuits", in IEEE Trans. on Device and Materials Reliability, 2006.
- [5] Tenbroek, B., Lee M., Redman-White, W., Bunyan, R., Uren, M.: "Self-heating effects in SOI MOSFETs and their measurement by small signal conductance techniques", In IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 1996.
- [6] Gardner, D., Flinn, P.: "Mechanical stress as a function of temperature in aluminum films", In IEEE Transactions on Electron Devices, 1988.
- [7] Vogel, E. et al., "Reliability of Ultra-Thin Silicon Dioxide Under Combined Substrate Hot Electron and Constant Voltage Tunneling Stress", Electron Devices, 2000.
- [8] Kaczer, B. et al., "GOB in FET devices and circuits: From nanoscale physics to system-level reliability", Elsevier, 2007.
- [9] Schröder, D.; Babcock, J.: Negative Bias Temperature Instability: Road to Cross in Deep Submicron Silicon Semiconductor Manufacturing. Journal of Applied Physics, vol. 94, 2003.
- [10] Li, P.-C., Stamoulis, G., Hajj, I.: "Probabilistic timing approach to hot-carrier effect estimation", In Trans. On Computer-Aided Design of Integrated Circuits and Systems, 1994.
- [11] Semiconductor Industry Association (SIA), "International Technology Roadmap for Semiconductors", Release 2009.
- [12] Von Neumann, J.: "Probabilistic logics and the synthesis of reliable organisms from unreliable components", In Automata Studies, Princeton University Press, 1956.
- [13] Siewiorek, D.: "Architecture of fault-tolerant computers: An historical perspective", In Proc. of the IEEE 79 (12), 1991
- [14] Mitra, S., Mc Cluskey, E.: "Word-voter: a new voter design for triple modular redundant systems", In Proc. of the IEEE VLSI Test Symposium, 2000.
- [15] Gaitanis, N.: "Design of totally self-checking TMR fault-tolerant systems", In IEEE Transactions on Computers, 1988.
- [16] Mitra, S., Mc Cluskey, E.: "Design of redundant systems protected against common-mode failures", In Proc. of IEEE VLSI Test Symposium, 2001.
- [17] J. Srinivasan, et al.: "RAMP: A Model for Reliability Aware Microprocessor Design", In IBM Research Report, RC23048, 2003.
- [18] Chandrakasan, A., Sheng, S., Brodersen, R.: "Low-power CMOS digital design", In IEEE Journal of Solid-State Circuits, 1992.
- [19] Koren, I.; Krishna, M.: Fault-Tolerant Systems. San Francisco: Morgan-Kaufmann, 2007.
- [20] Johnson, B.; Aylor, J.; Hana: "Efficient Use of Time and Hardware Redundancy for Concurrent Error Detection in a 32-bit VLSI Adder", In IEEE Journal of Solid-State Circuits, vol. 23, no. 1, 1988.
- [21] Zhu, D., Melhem, R., Mossé, D., Elnozayh, E.: "Analysis of an energy efficient optimistic TMR scheme", In Proc. of ICPADS, 2010