

Exploring Test and Reliability Challenges in the Realm of Approximate Circuitry

Sandeep Kumar Mishra,
Assistant Professor,
Dept. of Electronics & Telecommunication Engineering
College of Engineering Bhubaneswar

Abstract—Physical realization of approximate computing (AC) systems is becoming more and more relevant as it becomes increasingly clear that AC is a critical enabler for next-generation computer architectures. This letter explores an important but little-studied topic: the connection between the AC paradigm and hardware malfunctions. We outline the potential for effective testing techniques to raise the yields of the underlying hardware blocks for AC components as well as the related theoretical and practical difficulties. We distinguish between circuits with insufficient provisions and incomplete designs, and we outline the drawbacks of current testing techniques for these two categories of AC block types. We also go over the complex interactions between approximation circuitry and malfunctions that arise over the course of the device's lifetime.

I. INTRODUCTION

VARIOUS types of approximate computing (AC) were shown to offer benefits in applications including telecommunications [6], real-time image processing [2], neural network simulation [3], arithmetic calculations [20], and numerical analysis [17]. Several AC-enabled computer architectures, including Stanford's ERSA [10] and Intel's Runnede [4], have been proposed. AC concepts are being discussed as a core technology for next-generation deep-learning accelerators [21]. With AC getting close to practical application, the physical realization of approximate circuitry shifts into the focus. One key aspect of today's nanoscale technologies is their vulnerability to defects that are introduced in the course of circuit fabrication (chiefly addressed by testing) and to failures during the circuit's lifetime (chiefly addressed by reliability-enhancement techniques, such as online error detection and recovery). As we will show in this letter, the relationship of AC with state-of-the-art test and reliability approaches is insufficiently understood and needs more attention.

AC is defined across abstraction layers and combines concepts on system, software, and hardware levels [1]. However, test and reliability challenges are naturally more pronounced for approximate hardware blocks, also known by the name

“under signed and opportunistic circuitry” [6]. In general, approximate circuit blocks are allowed to deviate from their reference behavior in a well-defined way. Depending on the specific AC variant, this deviation might be due to incomplete design (e.g., an adder that does not function correctly for a few input combinations) [11], under provisioning (voltage or frequency overscaling) [7], [9], defects (error-tolerance) [8] or soft errors (transient error tolerance) [15].

The acceptable behavior is usually formalized by metrics based on error magnitude, error rate, or their combination. *Error magnitude* is the maximal allowed extent of deviation, e.g., maximal numerical difference [8], peak signal to noise ratio [5], or structural dissimilarity [9] in image-processing applications, or frame error rate in wireless telecommunication [19]. *Error rate* is the frequency with which erroneous outputs are produced (which, under simplifying assumptions, corresponds to the probability of error) [18]. In this letter, we discuss the challenges and the opportunities for test and reliability approaches in the context of approximate circuit blocks. We distinguish between two basic types of AC hardware: incomplete designs versus underprovisioned circuitry, because their respective handling from test and reliability perspectives is quite different. The remainder of this letter is organized as follows. Testing of incomplete designs and underprovisioned circuits are discussed in the next two sections, and reliability issues of approximate circuitry are the focus of Section IV. Section V concludes this letter.

II. TESTING INCOMPLETE DESIGNS

In this letter, “incomplete design” refers to an approximate circuit which deviates from its specification in a deterministic and repeatable manner. (In contrast, underprovisioned circuits may fail occasionally, governed by a random process.) One example is an adder which performs correct addition for all possible pairs of operands applied to its input except a small set of such operand pairs. The specification of the incomplete adder may demand that addition works “to some extent” (e.g., with 4-bit instead of 32-bit precision) even for operand pairs from S , or it may allow arbitrary misbehavior for such pairs. For a conventional digital circuit, one single erroneous bit that was observed during testing is a sufficient reason to reject this circuit, as it contains a manufacturing defect. s, testing checks that a manufacturing defect d turned the circuit's regular function f into f_d . Applying a test t can expose the defect through the measurement $f_d(t) / f(t)$. In case of an approximate block, we may opt against rejecting a defective

circuit as long as its function f_d still meets the loosened specification. In fact, one variant of AC, error-tolerance [5], [8], was chiefly motivated by improving manufacturing yield by retaining “good-enough” defective chips about which some guarantees could be made. For instance, such chips could be sold for use in less-critical application, for a lower price than fully functional, defect-free circuits. This differentiation is somewhat related to “speed binning,” where circuits manufactured in the same lot have slightly different performance (attainable clock frequency) due to random process variations. The speed of all circuits from the lot is measured, and faster circuits are used for higher-performance products.

While the possible yield improvement is clearly beneficial, there are also strong arguments against approximate circuits having manufacturing defects. First, the test set applied during manufacturing test is rather compact and cannot cover the complete circuit functionality. Even if only acceptable behavior, within the bounds of the loosened circuit specification, was observed during test, there is no guarantee that this will still be the case for arbitrary sequences which will be applied to the circuit in application. This problem does not occur in conventional testing, because a circuit with any defect, however small, is rejected without further examination.

Second, the defect can interfere with critical circuit infrastructure, such as clock distribution network, power grid, or voltage and frequency scaling control. Such interference may cause minimal impact during the rather short test, but may completely disrupt the circuit operation during its actual use. Third, many defects can deteriorate overtime and impact neighboring circuit structures due to factors like electromigration [13] or causing excessive power-supply noise [14]. Detrimental influence of infrastructure and aging also occur for conventional, nonapproximate circuits, but the presence of a tolerated defect magnifies these issues by far.

A. Structural Testing for Magnitude-Oriented Metrics

One promising idea is the extension of structural, model-based testing into the AC domain. This idea, which can be seen as a generalization of threshold testing [8] to more general approximate blocks, is illustrated in Fig. 1. Let the circuit be specified by a reference behavior model (e.g., a circuit without any approximations) and a maximal extent of deviation τ (threshold) according to an *error magnitude metric* Δ (e.g., structural dissimilarity of images). The circuit is acceptable as long as every value produced at its output deviates from the reference behavior by at most τ . Consider a (model of) defect d in the circuit. Create three blocks: one for reference behavior (producing output o), one for circuit with defect d (producing output o_d), and one which calculates the difference $\Delta(o, o_d)$ and compares it with τ .

The construction of Fig. 1 can be reduced to an automatic test pattern generation (ATPG) instance (in this case, the three blocks are represented by digital circuit models) or, equivalently, to a Boolean satisfiability (SAT) formula (then the blocks are mapped to conjunctive normal forms). Running an ATPG algorithm, or a SAT solver,

the faulty circuit exceed the maximum allowed

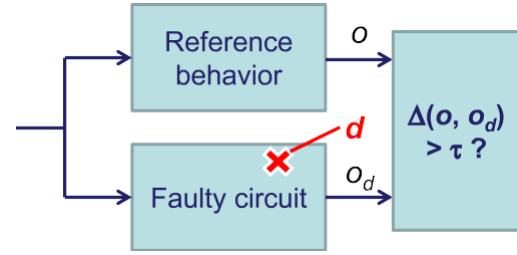


Fig. 1. Model-based test generation in an approximate circuit.

deviation. If the analysis shows that no such t exists, this constitutes a formal proof that the circuit with defect d is still acceptable. Note that threshold testing [8] is a special case, where the outputs are interpreted as numerical values with error magnitude $\Delta(o, o_d)$ set to their difference $|o - o_d|$. The “miter circuit” construction used in conventional circuit testing is an even more restricted variant, with $\Delta(o, o_d)$ being simply 1 if q_{a0} and 0 otherwise.

Two reasons prevent a wide-spread application of structural testing to approximate circuits. First, the need to model elaborated deviation metrics greatly increases the computational complexity of analysis. Already the simple threshold testing [8] was associated with much larger run times than conventional test generation. Incorporating complex error magnitude metrics, such as perceptual image or audio models, will make the test generation intractable. Moreover, the complexity increases further if the approximate block is embedded into a larger system and the test generated for the block needs to be transferred to the system’s top level.

Perhaps even more severe, it is impossible to imply the exact defect (or set of defects) present in a circuit from the measurements performed during manufacturing test. Even if we know for some well-defined defects d_1, \dots, d_n (e.g., some stuck-at or simple bridging faults) that a circuit with these defects is acceptable, we cannot be sure whether a specific circuit instance has one of these defects or some other, perhaps unmodeled, defects. There are diagnostic methods which derive likely defect locations from tester data, but they do not provide guarantees and often have to be complemented by physical failure analysis in practice.

B. Testing for Error-Rate Oriented Metrics

An approximate circuit is acceptable with respect to an *error-rate bound* $\epsilon \in [0, 1]$ if the fraction of its produced outputs that deviate from the reference behavior is ϵ or less. Existing error-rate test approaches [18] make the assumption that the input distribution is uniform; then, the circuit with input space I is acceptable if the input subspace S for which it deviates from the reference behavior fulfills $|S|/|I| \leq \epsilon$. Error-rate test generation [18] aims at deciding, for a given defect d , whether a circuit affected by d will still satisfy the error-bound. This is done by generating a test set T that is indicative for the error-rate of the entire input space I (which is normally exponential in the number of circuit inputs). In other words, $|S \cap T|/|T|$ should approximate $|S|/|I|$ as tight

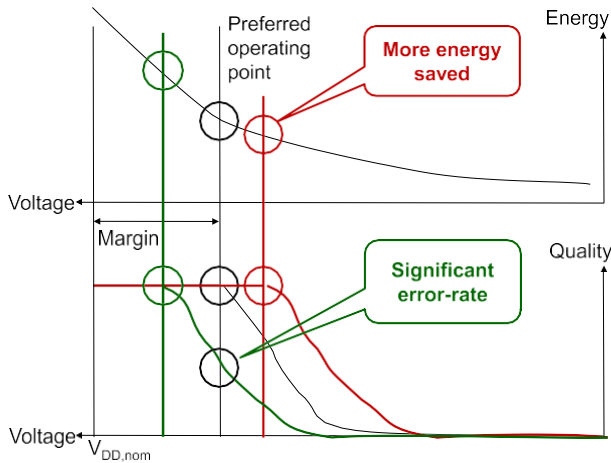


Fig.2.Effects of voltage scaling on energy consumption and output quality for a fast (red), nominal (black), and slow (green) circuit instance.

as possible. Such test sets tend to be larger than conventional detection tests, and they are valid only with a confidence.

An alternative approach to the error-rate testing problem is to generate a set of acceptable defects $D_{acc} = \{d_1, \dots, d_n\}$

and to perform diagnosis for a specific failing circuit. If the circuit is as one of the defects from D_{acc} , it can be considered acceptable despite being defective. Problems listed for magnitude-oriented metrics (insufficient diagnostic resolution, unmodeled defects, and reliability impact) still apply in this case.

III. TESTING UNDERPROVISIONED CIRCUITS

Underprovisioned circuits are operated at voltage-frequency points for which their correct functioning is not guaranteed. This is illustrated in Fig. 2 for the case of supply voltage V_{DD} . Suppose that the circuit's V_{DD} is gradually reduced, i.e., moved to the right of its nominal value $V_{DD,nom}$, while its clock frequency (cycle duration) is kept constant. The power and the energy consumed by the circuit is reduced roughly quadratically with V_{DD} decrease, as indicated in the upper part of Fig. 2. At the same time, the switching speed of logic gates reduces, and at some point one or multiple paths through the circuit will become slower than the cycle time and errors will occur on circuit outputs.

Circuits are usually designed with a margin, i.e., the difference between the cycle duration and the delay of the longest path in the circuits exceeds zero. Because of natural process variations, the margins are different for specific manufactured circuit instances. It is advantageous to operate the circuit with just the minimal V_{DD} before failure, because this minimizes the energy consumption. Fig. 2 shows the preferred operating points for three circuit instances of different speed. It is important to select individual operating points rather than one precalculated for the nominal circuit, because otherwise faster

circuits cannot realize maximal energy benefits, and slower circuits enter the regime beyond their margin and produce excessive error-rate.

Note that temperature changes and V_{DD} fluctuations may shift the curves from Fig. 2 dynamically. This can be addressed by performing introspection at run-time by detecting errors and adapting (increase voltage or reduce frequency) upon their occurrence [9]. As an alternative, it is possible to avoid adaptation and continue operation in the beyond-the-margin regime [7]. However, in many practical circuits no guarantees about error magnitudes or rates in this regime can be made.

Self-adapting underprovisioned approximate circuits pose conceptual challenges to testing. One key question is whether the self-adaptation features should be enabled or disabled during testing. If self-adaptation is enabled on a circuit which has a small-delay defect just outside the normal extent of variability, it may compensate this defect by, e.g., increasing V_{DD} to the maximum possible level, such that the circuit passes the test. This circuit is then operated but has no resilience reserves and will fail during the next temperature increase or occurrence of power-supply noise. If self-adaptation is disabled, the tester may confuse "regular" transient failures, from which it is designed to recover, with defects and reject circuits which could be normally operated, thus increasing the yield loss.

Further issues in testing underprovisioned approximate circuits are parameter-dependent defects and defects in the adaptation infrastructure itself. Parameter-dependent defects are defects which manifest themselves only under specific voltage/frequency conditions (they may be seen as irregularities on what test engineers know as "Shmoo plots"). The most obvious technique to address them is to test the circuit under a number of different voltage/frequency combinations, which, however, increases test time, complexity, and cost. Defects in self-adaptation circuitry are difficult to observe during test, because they are disconnected from the functional part of the circuit. Reproducing all possible self-adaptation scenarios during test is infeasible, and special design-for-testability features will be typically required. This problem generally arises for conventional circuits as well, but the comparatively large margins of such circuits make them less prone to self-adaptation related failures than underprovisioned circuits and soften their requirements on test procedures.

It is logical to apply error-rate metrics discussed above to underprovisioned circuits. However, their nondeterministic and complex, intermittent nature of failure patterns poses known error-rate testing concepts inapplicable to them. Any error-rate statements for such circuits should perhaps be obtained by more elaborate probabilistic characterization procedures, such as tomographic testing originally developed for quantum and other probabilistic circuits [12].

APPROXIMATE CIRCUITS

Like their conventional counterparts, approximate circuits are potentially affected by failures that occur after manufacturing and during the device's lifetime. Early life failures and wearout mechanisms could, e.g., slowly increase the circuit's error rate from an acceptable value to one over the threshold ϵ . For example, negative-bias temperature instability, or hot-carrier effects may increase the logic gate delays within the

circuit, and electromigration could do the same with interconnect delays [13]. As a consequence, the circuit will exceed its cycle duration more often.

Some of the specific aspects of approximate circuits make them more and some less vulnerable to early life failures and wearout in comparison with regular circuits. Approximate circuits often have some degree of error-resilience or graceful degradation, i.e., capability to operate (at least to some extent) in presence of errors. This feature increases their robustness and, potentially, their lifetime, but may hinder error-detection which would miss the buildup of many small, uncritical errors to a large and critical failure.

The underprovisioned species of approximate circuits tend to be operated under more relaxed operating points than functionally identical conventional designs, reducing the probability of over stress that lead to early life failures and excessive aging. At the same time, they are operated in a marginal regime, with little or no noise margins. The smallest disturbance will lead to a failure, and an unexpectedly large disturbance may exceed the capability of the self-adaptation infrastructure to recover the circuit from the failure. It appears essential that self-adaptation strategies and in particular introspection routines account for the possibility of aging and deterioration during the circuit's lifetime. For this purpose, they may have to be combined with built-in self-test and characterization procedures which not only decide whether the circuit is faulty but also distinguish intermittent errors (which the circuit is designed to tolerate) from permanent wearout-induced failures [16].

V. CONCLUSION

Approximate circuits must be manufactured, tested, and their post-production quality must be regulated if they are to be used. In the event of incomplete designs, conventional test methods might not fully grasp the potential of yield optimization, and they might not be accurate for circuitry that is underprovisioned. There are methods for testing circuits with looser specifications (concerning error magnitude or error rate), but they are often imprecise, predicated on faulty premises, and not very scalable. Instead of using simulations, there seems to be an urgent need for data obtained on real circuits produced in large volumes. What proportion of the flaws that were found were classified as "benign"? Did the tests designed for a basic fault model, such as stuck-at faults, accurately depict the real failures? In the context of reliability, the lack of data is even more significant. Approximate circuits need to be efficiently tested, which requires the solution of several unresolved issues. Automated proofs for

approximation boundaries under the assumption of a complicated acceptability metric and test generation under such a metric are of essential interest. Better failure isolation and diagnosis during post-manufacturing testing must be made possible by design-for-testability techniques, and on-chip introspection circuitry must offer fine-grained circuit health monitoring spontaneously. These characteristics already advantageous for conventional circuits, but they are crucial for approximation circuits as well. The additional expense of these solutions must be carefully weighed against the efficiency advantages brought about by the AC paradigm. Vendors of electrical design automation tools must undertake an economic analysis and a legal evaluation of the hazards associated with delivering circuits that are known to be defective in tandem with these investigations.

REFERENCES

- [1] A. Agrawal *et al.*, "Approximate computing: Challenges and opportunities," in *Proc. IEEE Int. Conf. Rebooting Comput.*, San Diego, CA, USA, 2016, pp. 1–8.
- [2] A. Alaghi, C. Li, and J. P. Hayes, "Stochastic circuits for real-time image processing applications," in *Proc. DAC*, 2013, Art. no. 136.
- [3] B. D. Brown and H. C. Card, "Stochastic neural computation. I. Computational elements," *IEEE Trans. Comput.*, vol. 50, no. 9, pp. 891–905, Sep. 2001.
- [4] N. P. Carter *et al.*, "Runnemed: An architecture for ubiquitous high-performance computing," in *Proc. IEEE HPCA*, Shenzhen, China, 2013, pp. 198–209.
- [5] I. S. Chong and A. Ortega, "Hardware testing for error tolerant multimedia compression based on linear transforms," in *Proc. IEEE DFTS*, Monterey, CA, USA, 2005, pp. 523–531.
- [6] P. Gupta *et al.*, "Underdesigned and opportunistic computing in presence of hardware variability," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 32, no. 1, pp. 8–23, Jan. 2013.
- [7] R. Hegde and N. R. Shanbhag, "Soft digital signal processing," *IEEE Trans. Very Large Scale Integr. (VLSI) Syst.*, vol. 9, no. 6, pp. 813–823, Dec. 2001.
- [8] Z. Jiang and S. K. Gupta, "Threshold testing: Improving yield for nanoscale VLSI," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 28, no. 12, pp. 1883–1895, Dec. 2009.
- [9] P. K. Krause and I. Polian, "Adaptive voltage over-scaling for resilient applications," in *Proc. DATE*, Grenoble, France, 2011, pp. 1–6.
- [10] L. Leem, H. Cho, J. Bau, Q. A. Jacobson, and S. Mitra, "ERSA: Error resilient system architecture for probabilistic applications," in *Proc. DATE*, Dresden, Germany, 2010, pp. 1560–1565.
- [11] J. Miao, A. Gerstlauer, and M. Orshansky, "Multi-level approximate logic synthesis under general error constraints," in *Proc. ICCAD*, San Jose, CA, USA, 2014, pp. 504–510.
- [12] A. Paller, A. Alaghi, I. Polian, and J. P. Hayes, "Tomographic testing and validation of probabilistic circuit," in *Proc. IEEE ETS*, Trondheim, Norway, 2011, pp. 63–68.
- [13] M. Pecht, R. Radojic, and G. K. Rao, *Guidebook for Managing Silicon Chip Reliability*. Boca Raton, FL, USA: CRC Press, 1999.
- [14] I. Polian, "Power supply noise: Causes, effects, and testing," *ASPJ. Low Power Electron.*, vol. 6, no. 2, pp. 326–338, 2010.
- [15] I. Polian, J. P. Hayes, S. Kundu, and B. Becker, "Transient fault characterization in dynamic noisy environments," in *Proc. ITC*, Austin, TX, USA, 2005, pp. 1048–1057.
- [16] L. R. Gómez *et al.*, "Adaptive Bayesian diagnosis of intermittent faults," *J. Electron. Test. Theory Appl.*, vol. 30, no. 5, pp. 527–540, 2014.
- [17] A. Schöll, C. Braun, and H.-J. Wunderlich, "Applying efficient fault tolerance to enable the preconditioned conjugate gradient solver on approximate computing hardware," in *Proc. IEEE DFTS*, 2016, pp. 21–26.
- [18] S. Shahidi and S. K. Gupta, "ERTG: A test generator for error rate testing," in *Proc. ITC*, Santa Clara, CA, USA, 2007, pp. 1–10.
- [19] V. Tomashevich, C. Gimmler-Dumont, N. When, and I. Polian, "Reliability analysis of MIMO channel preprocessing by fault injection," in *Proc. IEEE WISSE Conf.*, Noordwijk, The Netherlands, 2014, pp. 1–6.
- [20] R. T. Uppu, R. K. Uppu, A. D. Singh, and A. Chatterjee, "A high throughput multiplier design exploiting input based statistical distribution in completion delays," in *Proc. Int. Conf. VLSI Design*, Pune, India, 2013, pp. 109–114.
- [21] J. Yoshida, *Race for AI Chips Begins*. EE Times, San Francisco, CA, USA, 2016. [Online]. Available: http://www.eetimes.com/document.asp?doc_id=1331052