

# Temporal Memoization for Energy-Efficient Timing Error Recovery in GPGPUs

Abbas Rahimi  
CSE, UC San Diego  
La Jolla, CA 92093, USA  
abbas@cs.ucsd.edu

Luca Benini  
UNIBO and ETHZ  
40136 Bologna, Italy  
8092 Zurich, Switzerland  
luca.benini@{unibo.it, iis.ee.ethz.ch}

Rajesh K. Gupta  
CSE, UC San Diego  
La Jolla, CA 92093, USA  
gupta@cs.ucsd.edu

**Abstract**—Manufacturing and environmental variability lead to timing errors in computing systems that are typically corrected by error detection and correction mechanisms at the circuit level. The cost and speed of recovery can be improved by *memoization*-based optimization methods that exploit spatial or temporal parallelisms in suitable computing fabrics such as general-purpose graphics processing units (GPGPUs). We propose here a *temporal memoization* technique for use in floating-point units (FPUs) in GPGPUs that uses value locality inside data-parallel programs. The technique recalls (*memorizes*) the context of error-free execution of an instruction on a FPU. To enable scalable and independent recovery, a single-cycle lookup table (LUT) is tightly coupled to every FPU to maintain contexts of recent error-free executions. The LUT reuses these memorized contexts to exactly, or approximately, correct errant FP instructions based on application needs. In real-world applications, the temporal memoization technique achieves an average energy saving of 8%–28% for a wide range of timing error rates (0%–4%) and outperforms recent advances in resilient architectures. This technique also enhances robustness in the voltage overscaling regime and achieves relative average energy saving of 66% with 11% voltage overscaling.

**Keywords**—Variability; timing errors; error recovery; temporal memoization; value locality; GPGPU

## I. INTRODUCTION

The scaling of physical dimensions in semiconductor circuits opens the way to an astonishing over 7 billion transistors on a 28nm process which gives a grand total of 2,880 CUDA cores in recent GPGPU chips [1]. It is also leading to ever-increasing parametric variations across process, voltage and temperature (PVT) [2]. As designers build circuits operating near-threshold [3] in order to save power, and use voltage overscaling [4] to reach performance targets, the effect of PVT variations are exacerbated.

The most common effect of variation is violation of timing specifications that cause circuit-level timing errors. IC designers commonly use conservative guardbands for the operating frequency or voltage to ensure error-free operation for the worst-case variations. These guardbands have been steadily increasing, thus leaving untapped performance and other costs of overdesign [5]. An alternative to overdesign is to make a design resilient to errors and variations. In this paper, we specifically focus on ‘Design for Resiliency’ (DFR) against timing errors.

Low-voltage DFR applies to both logic and memory blocks. For memory, 8T SRAM arrays utilize a tunable replica bits therefore enables reduction of the minimum operating voltage [7]. Similarly, in non-volatile memory area, resistive

RAM (ReRAM/memristor) is a promising candidate at low power supply voltage [8], [9]. For logic, error-detection sequential (EDS) [6] circuit sensors have been employed to reduce guardbanding in the designed circuits. A common strategy is to detect variability-induced delays by sampling and comparing signals near the clock edge to detect timing errors. The timing errors are corrected by replaying the errant operation with a larger guardband through various adaptation techniques. For instance, a resilient 45nm integer-scalar core [10] places EDS circuits in the critical paths of the pipeline stages. Once a timing error is detected during instruction execution, the core prevents the errant instruction from corrupting the architectural state and an error control unit (ECU) initially flushes the pipeline to resolve any complex bypass register issues. To ensure scalable error recovery, the ECU supports two separate techniques: instruction replay at half frequency, and multiple-issue instruction replay at the same frequency. These techniques impose energy overhead and latency penalty of up to 28 extra recovery cycles per error for the resilient 7-stage integer pipeline [10].

As energy becomes the dominant design metric, aggressive voltage scaling [4] and near-threshold operations [3] increase the rate of timing errors and correspondingly the costs (in energy, performance) of these recovery mechanisms. This cost is exacerbated in FP single-instruction multiple-data (SIMD) pipelined architectures where the pipeline dimensions are expanded both vertically (with wider parallel lanes) and horizontally (with deeper stages). The horizontally expanded deeper pipelines induce higher pipeline latency and higher cost of recovery through flushing and replaying the errant instruction. The FP pipelines consume higher energy-per-instruction than their integer counterparts and typically have high latency for instance over 100 cycles [11] to execute on a GPGPU. Effectively, these energy-hungry high-latency pipelines are prone to inefficiencies under timing errors. Similarly, in vertically expanded pipelines, there is a significant performance drop in a 10-lane SIMD architecture as single-stage-error probabilities increase [12]. In the lock-step execution, any error within any of the lanes will cause a global stall and force recovery of the entire SIMD pipeline.

Thus, in FP SIMD pipelines the error rate is multiplied by the wider width while the number of recovery cycles per error increases at least linearly with the pipeline length. This makes the cost of recovery per single error quadratically more expensive relative to scalar functional units. At the same time, parallel execution in the GPGPU architectures – described in Section III– provides an important ability to reuse computation and reduce the cost of recovery from timing errors. This paper, exploits this opportunity to make three

main contributions: First, we propose a *temporal memoization* technique to avoid conventional recovery costs. We observe that the entropy of data-level parallelism is low due to high locality of values. The temporal memoization of recent error-free executions exploits this inherent value locality. We show that the memorized information can be used exactly or approximately depending upon the application needs. Second, we closely integrate a lightweight single-cycle LUT to the FPU to support instruction-level memoization. The design enables a scalable and independent recovery of individual FPUs. Section IV covers these details. Third, we demonstrate the effectiveness of our technique on the Evergreen GPGPU architecture for error-tolerant image processing applications as well as error-intolerant general-purpose applications. Our experimental results in Section V show an average hit rate of 76% on 4-entry LUTs: it leads to an average energy saving of 8% in case of an error-free execution environment (0% timing error); it further reaches to an average energy saving of 28% in the presence of 4% timing error rate thanks for improving recovery cost on the errant instructions. This technique also enhances robustness in the voltage overscaling scenario, up to an average energy saving of 66%. Finally, Section VI concludes this paper.

## II. RELATED WORK

Sodani and Sohi [13] introduced the concept of instruction reuse from the observation that many instructions can be skipped if another instance has already been executed using the same input values. The instruction reuse enables memorization of the outcome of an instruction in hardware tables, therefore a processor can reuse it temporally if the processor performs the same instruction with the same input values within a limited period of time before the entry is overwritten. This temporal memorization technique is fundamentally limited by the latency and the low hit rate of the tables. To improve the hit rates, recent reuse techniques [14], [15] seek to improve association of the entries of the table with similar inputs to the same output. These tolerant techniques rely upon the tolerance in the output precision of multimedia algorithms to achieve high reuse rates, and work at the granularity of the FP instruction [15], or a region of FP instructions [14]. A load value prediction technique also exploits the value locality to predict the register file values being loaded from memory based on previously-seen values [16]. These techniques have been devised for single-core architectures without exploring their potential in combating timing errors in data-level parallel architectures.

Beyond instruction-level reuse, various techniques have been proposed to mitigate the variation-induced timing errors, including adaptive management of guardbanding through ‘predict-and-prevent’ mechanisms [17], [18], [22], [19], [20], and ‘detect-then-correct’ mechanisms [6], [7], [10], [12]. A brief review of the main concepts and their embodiments follows.

*Predict-and-prevent* techniques try to avoid timing errors while reducing guardbands, for instance for individual instructions [17]. The instruction program counter of an out-of-order pipeline is used for an early prediction of an upcoming timing violation by searching in a large predictor table [22]. At higher levels, a *procedure hopping* technique is proposed to avoid voltage droops [18]. In the context of the GPGPU architecture, hierarchically focused guardbanding [19] has been proposed earlier at two levels: fine-grained instruction-level and coarse-grained kernel-level. Rahimi et al. [20] propose a

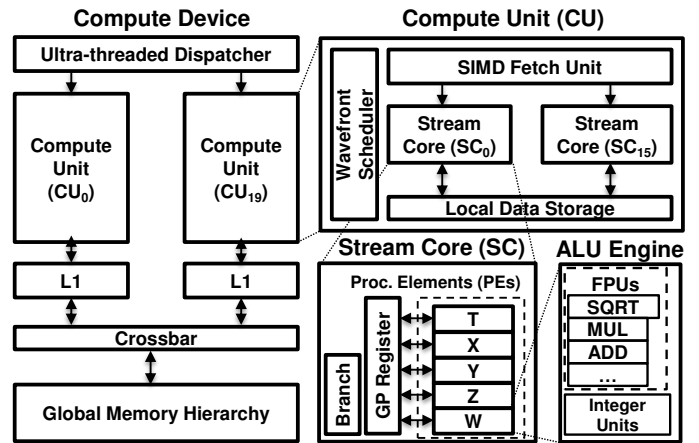


Fig. 1. Block diagram of the Radeon HD 5870 GPGPU.

compiler technique that periodically regenerates healthy codes that reduces the aging-induced performance degradation of the GPGPUs. The predictive techniques cannot eliminate the entire guardbanding to work efficiently at the edge of failure specially so with frequent timing errors in the voltage overscaling and near-threshold regimes.

*Detect-then-correct* technique for SIMD architectures decouples the lanes through private queues that prevent error events in any single lane from stalling all other lanes [12]. This enables each lane to recover from errors independently. The decoupling queues cause slip between lanes which requires additional architectural mechanisms to ensure correct execution. Further, the decouple queue relies on the recovery based on the global clock-gating which involves stalling the entire lane. This causes one cycle recovery penalty over a two-stage execution unit [12]. However, propagating a global stall signal over a deep GPGPU pipeline [11] is expensive. A *spatial memoization* [21] technique also broadcasts output result of an error-free instruction across all *error-prone* lanes, tightens its scalability. Hence, the cost of scalable recovery (e.g., [10]) per single timing error on these architectures is high limiting their utility to low error rate circumstances. Our present work enhances the scope of ‘detect-then-correct’ approaches in a GPGPU context thanks to an ultra-low cost recovery through memoization, thus offering both scalability and low-cost self-resiliency in the face of high timing error rates.

## III. GPGPU ARCHITECTURE

We focus on the Evergreen family of AMD GPGPUs (a.k.a. Radeon HD 5000 series), that targets general-purpose data-intensive applications. The Radeon HD 5870 GPGPU consists of 20 compute units, a global front-end ultra-thread dispatcher, and a crossbar to connect the memory hierarchy. Each compute unit contains a set of 16 Stream Cores (SCs), i.e., 16 parallel lanes. Within a compute unit, a shared instruction fetch unit provides the same machine instruction for all SCs to execute in a SIMD fashion. Each SC contains five Processing Elements (PEs) – labeled X, Y, Z, W, and T – froming an ALU engine to execute Evergreen machine instructions in a vector-like fashion. Finally, the ALU engine has a pool of pipelined integer and FP units. The block diagram of the architecture is shown in Fig. 1.

The device kernel is written in OpenCL which runs on a GPGPU device. An instance of the OpenCL kernel is called a work-item. Each SC is devoted to the execution of one work-item. In the Radeon HD 5870, a *wavefront* is defined

as the total number of 64 work-items virtually executing at the same time on a compute unit. To manage 64 work-items in a wavefront on 16 SCs of the compute unit, a wavefront is split into *subwavefronts* at the execute stage, where each subwavefront contains as many work-items as available SCs. In other words, SCs execute the instructions from the wavefront mapped to the SIMD unit in a 4-slot time-multiplexed manner using the integer units and FPUs. The time-multiplexing at the cycle granularity relies on the functional units to be fully pipelined.

Evergreen assembly code uses a clause-based format classified in three categories: ALU clause, TEX clause, and control-flow instructions. The control-flow instructions triggering ALU clauses will be placed in the input queue at the ALU engine. There is only one wavefront associated with the ALU engine. After fetch and decode stages, the source operands for each instruction are read that can come from the register file or local memory. For higher throughput, buffers are attached to SCs to read the registers ahead of time. The core stage of a GPGPU is the execute stage, where arithmetic instructions are carried out in each SC. When the source operands for all work-items in the wavefront are ready, the execution stage starts to issue the operations into the SCs. Finally, the result of the computation is written back to the destination operands.

#### IV. TEMPORAL MEMOIZATION

We briefly describe how value locality can increase the resiliency of the GPGPUs followed by a description of the proposed temporal memorization technique. We divide all applications into two classes: error-tolerant image processing applications and error-intolerant general-purpose applications selected from AMD APP SDK 2.5 [23]. For error-tolerant applications, we have examined four image processing filters: Sobel, Gaussian, Box, and URNG. The error-tolerant applications exhibit enhanced error resilience at the application-level when multiple valid output values are permitted, in effect, creating a relation from input values to (multiple) output values. Instead of a single number, the output value is associated with a quality metric that may be within the constraints of application-specific fidelity metrics such as peak signal-to-noise ratio (PSNR) [24]. Therefore, if execution is not 100% numerically correct, the application can still appear to execute correctly from the users perspective [14], [15], [24].

In case of error-intolerant applications that do not have such inherent algorithmic tolerance, even a single bit error could result in unacceptable program execution. In this class, we have examined four applications: Black-Scholes and Binomial models for European-style options in financial engineering; one-dimension Haar wavelet transform; and Eigenvalues of a symmetric matrix.

##### A. FP Instruction-Level Memoization

We focus on the individual FPUs to observe the dispersion of the input operands at the finest granularity. To expose the value locality for each FPU operations, we consider a private FIFO for every individual FPU. These FIFOs have a small depth and keep the distinct sets of the input operands in the order of instruction arrivals. The FIFO matches a set of incoming input operands and the current content of entries of FIFO using a matching constraint. Since the applications exhibit varying degrees of tolerance, two matching constraints are considered: (i) *Exact matching* constraint that enforces full bit-by-bit matching of the input operands of the FPU

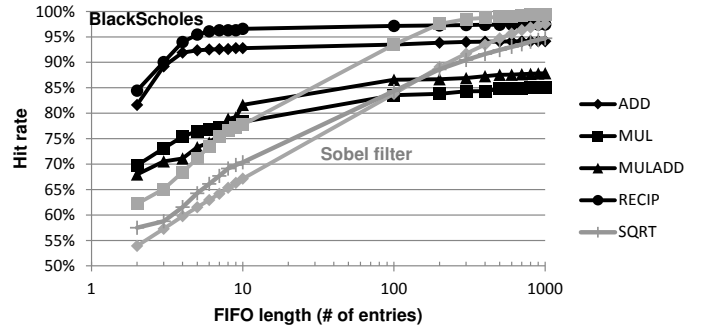


Fig. 2. Hit rate for various FPUs as a function of FIFO length when executing Black-Scholes and Sobel filter.

with the FIFO's entries. This constraint is applicable for the error-intolerant applications; (ii) *Approximate matching* constraint that relaxes the criteria of the exact matching during comparison of the operands by masking the less significant  $N$  bits of the fraction parts. For our image processing filters the approximate matching constraint with masking the less significant 12 bits of the fraction parts guarantees the PSNR of greater than 30dB (this is generally considered acceptable from users perspective in image processing applications).

The matching constraints are programmable and also allow commutativity of the operands where applicable. The FIFO maintains a limited number of recent distinct sets. Therefore, if a set of incoming input operands does not satisfy either matching constraints, the FIFO will be updated by cleaning its last entry and inserting the new incoming operands accordingly. The implementation of the FIFO is detailed in the following subsection.

We measure the overall hit rate of the FIFOs for different types of the FPUs with various FIFO length, ranging from 2 to 1000 entries. Fig. 2 shows the overall hit rate for five types of the FPUs when executing two different applications, Black-Scholes, and Sobel filter. As shown, the hit rate depends on the FPU operations. All FPUs display a hit rate of greater than 50% with a FIFO of 2 entries; and the hit rates increase with the larger FIFOs. For instance, the FP addition units (ADD) with 2 entries FIFOs experience a hit rate of 82% (or 54%) when executing Black-Scholes (or Sobel filter). These hit rates are further increased to 93% (67%) for FIFOs with 10 entries. Clearly, the temporal value locality is a function of both application class and operation type. We also assess the value locality across all types of the FPUs activated during the execution of the two classes of applications. Fig. 3 shows the overall hit rate of all applications when the FPUs utilizing FIFOs with 2, 4, 10, 100, and 1000 entries. As shown, the FIFOs with 2 entries exhibit the lowest hit rate, while the benefit of the large FIFOs is application-dependent and does not improve the hit rate significantly. The hit rate of most applications increases less than 10% when the size of FIFOs is increased by two orders of magnitudes from 10 to 1000. On the other hand, the FIFOs with 4 entries display higher energy efficiency ( $2.8\times$  higher hit rate per power compared to the 10 entries FIFOs), and provide an average hit rate of 76% (up to 97%) for the applications. Therefore, we have used the FIFOs with 4 entries for our proposed temporal memoization technique, and we also measured its energy efficiency in Section V.B.

We note that temporal value locality has been observed earlier in single-core architectures [16]. Our results show that the data-level parallel execution across both application classes

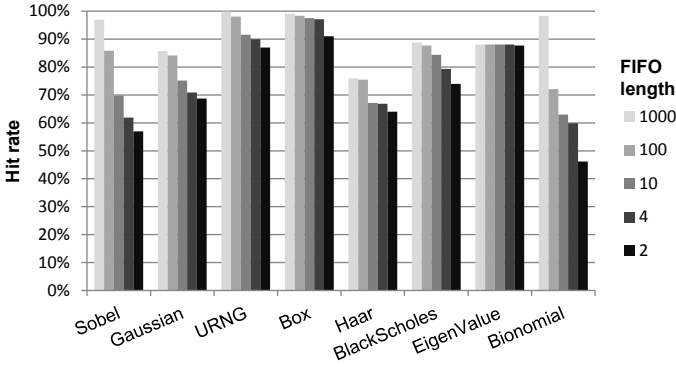


Fig. 3. Overall hit rate of the applications utilizing different FIFOs with 4, 10, 100, 1000 entries.

also displays a significant value locality. This is primarily due to redundant contextual information on the SCs. As a result, the FPU of GPGPUs experience a congested temporal value locality caused by the sub-wavefront time-multiplexing on the SCs that can be exposed by small FIFOs.

### B. Resilient Architecture Utilizing Temporal Memoization

We now describe the design of a resilient GPGPU architecture that utilizes the temporal memoization technique. In Evergreen GPGPUs, every FPU has four execution stages and a throughput of one instruction per cycle. We instrument every FPU pipeline with the error detection and correction mechanisms proposed in [6], [10]. Essentially, every stage uses EDS circuit sensors to detect the timing errors by propagating an error signal toward the end of pipeline that finally reaches to the ECU, the error control unit. The ECU triggers the recovery mechanism through flushing, and issuing the errant instruction multiple times [10]. This forms the baseline ‘detect-then-correct’ mechanism shown as the white components in Fig. 4.

To exploit the value locality, we tightly couple the FPU pipeline with our proposed temporal memoization module (shown as the gray components in Fig. 4). This module has essentially a single-cycle LUT, and a set of flip-flops and buffers to propagate signals through the pipeline. The LUT is shown in the bottom part of Fig. 4 and is composed of two parts: (i) a FIFO with four entries; (ii) a set of combinational comparators. In every entry, the FIFO maintains a set of input operands (e.g., in Fig. 4, two inputs) and the computed result provided by the output of the FPU in the last stage ( $Q_S$ ). The parallel combinational comparators implement the two matching constraints, and are programmable through a 32-bit memory-mapped register as a masking vector. They concurrently make either a full or partial comparison of the input operands with the stored operands in each entry based on the masking vector. The LUT works in parallel with the first stage of the FPU. Therefore, for every set of input operands, the LUT searches the FIFO to find a match between the input operands and the operand values stored in the entries (i.e., whether the matching constraint is satisfied or not). A match directly results in reuse of results computed earlier. Consequently, this affords the temporal memoization module an opportunity to correct an errant instruction with zero cycle penalty.

To enable reuse, the LUT propagates a hit signal alongside with the previously-computed result ( $Q_L$ ) toward the end of pipeline. The LUT raises the hit signal that squashes the remaining stages of the FPU to avoid the redundant computation

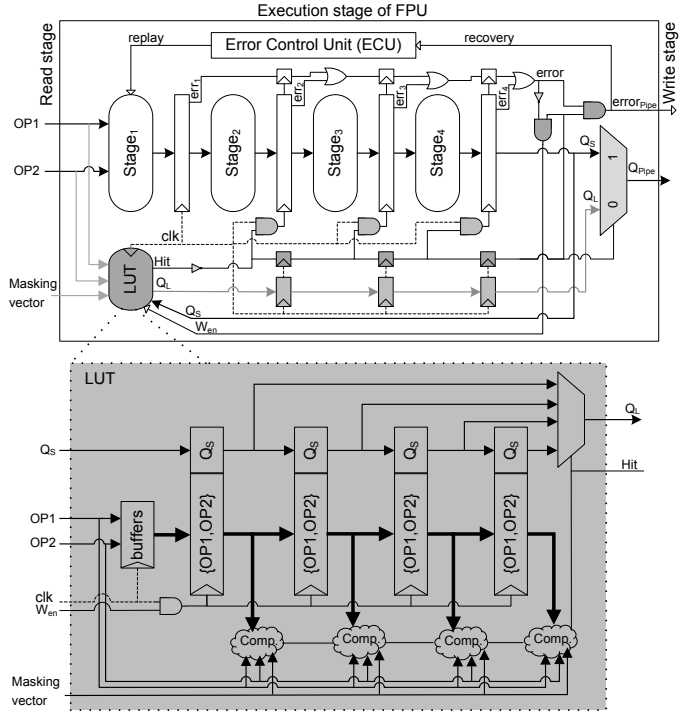


Fig. 4. Execution stage of the resilient FPU with temporal memoization module (in gray); the LUT is shown below.

TABLE I. TIMING ERROR HANDLING WITH TEMPORAL MEMOIZATION MODULE.

Hit	Error	Action	$Q_{Pipe}$
0	0	Normal execution + LUT update	$Q_S$
0	1	Triggering baseline recovery (ECU)	$Q_S$
1	0	LUT output reuse + FPU clock-gating	$Q_L$
1	1	LUT output reuse + FPU clock-gating + masking error	$Q_L$

by clock-gating; the clock-gating signal is forwarded to the rest of stages, cycle by cycle. The stored result is also propagated toward the end of pipeline for the reuse purpose. The hit signal selects the propagated output of the LUT ( $Q_L$ ) as the output of the FPU; it also disables the propagation of timing error signal (if any) to the ECU, thus avoids the costly recovery. Therefore, each hit event reduces energy by locally retrieving the result from the LUT, rather than doing full re-execution by the FPU. In case of a LUT miss, the FIFO is updated to maintain the last recently computed values. It is implemented through a write enable signal ( $W_{en}$ ) that ensures there is no timing error during execution of all stages of the FPU for computing  $Q_S$ . Finally, if a simultaneous timing error and miss occurred, the error signal will be propagated to the ECU that triggers the baseline recovery. Table I summarizes these four states.

Each application has full control over the temporal memoization module as a programmable module through the memory-mapped registers. To determine the matching constraints, the application can set the 32-bit memory-mapped register of the masking vector accordingly. The error-tolerant applications set the masking vector to ignore the differences of the operands in the less significant 12 bits of the fraction part. With the approximate matching constraint, the pair of instructions with two different input operands will have the same output. On the other hand, the exact matching constraint enables the reuse of the previously-computed results stored in the LUT while maintaining the full precision for the error-intolerant applications. Moreover, if an application lacks value locality, it can disable the entire memoization module by power-gating thus avoid any power penalty. Further, compiler-

directed analysis techniques or domain experts with some application knowledge can also store pre-computed values in the LUT to use the most probable or critical results.

## V. EXPERIMENTAL RESULTS

Our methodology uses the AMD Evergreen GPGPUs, but can be applied to other GPGPU architectures as well. Multi2Sim [25], a cycle-accurate CPU-GPU simulation framework, is modified to collect the statistics for computing the temporal value locality. The naive binaries of AMD APP SDK 2.5 [23] kernels are run on the simulator; the input values for the kernels are generated by the default OpenCL host program. We analyzed the effectiveness of proposed technique in the presence of timing errors and voltage overscaling on TSMC 45-nm ASIC flow.

### A. Implementation of Temporal Memoization Module

To keep the focus on the energy-hungry high-latency FP pipelines, we assume that the memory blocks are resilient, for instance by utilizing the tunable replica bits [7]. Since the fetch and decode stages display a low criticality [17], we focus on the execution stage consisting of six frequently exercised functional units: ADD, MUL, SQRT, RECIP, MULADD, FP2FIX. On Evergreen, every ALU functional unit has a latency of four cycles and a throughput of one instruction per cycle [27]. Therefore, VHDL codes of the FPUs are generated and optimized using FloPoCo [26] – an arithmetic synthesizable FP core generator. To achieve a balanced clock frequency across the FP pipelines, the RECIP has a latency of 16 cycles, while the rest of the FPU have four cycles latency.

The temporal memoization module for each FPU operations is described in Verilog synthesizable RTL. To integrate the resilient architecture, the memoization modules are integrated into the FPUs pipelines with the baseline recovery mechanism. Finally, the entire design is synthesized and mapped using the TSMC 45nm technology library. The front-end flow with multi  $V_{TH}$  cells has been performed using *Synopsys Design Compiler* with the topographical features, while *Synopsys IC Compiler* has been used for the back-end. The design has been optimized for timing a signoff frequency of 1GHz at (SS/0.81V/125°C), and then optimized for power using high  $V_{TH}$  cells. The overall area overhead of the temporal memoization modules with four entries FIFO is 0.11% of the total die area of Radeon HD 5870. The power overhead is also negligible and it is entirely paid off by the power saving due to the frequent clock-gating of the FPUs during the hit events that results into even higher energy efficiency detailed in the following subsection. We note that the overhead will be further reduced for deeper pipelines. The memoization module does not limit the clock frequency as it has a positive slack of 14% of the clock period.

### B. Energy Efficiency and Energy Saving

We compare the energy efficiency (GFLOPS/Watt) and energy saving of the proposed architecture with a baseline architecture that also utilizes recent resilient techniques [12], [10]. For the baseline architecture, we consider the decoupling queues technique [12] with error detection of EDS, and the baseline recovery mechanism of the multiple-issue instruction replay [10] adapted for the FPUs to support scalability. Our proposed temporal memoization architecture superposes the temporal memoization modules on the baseline architecture. In our experiments the EDS, ECU and the temporal memoization

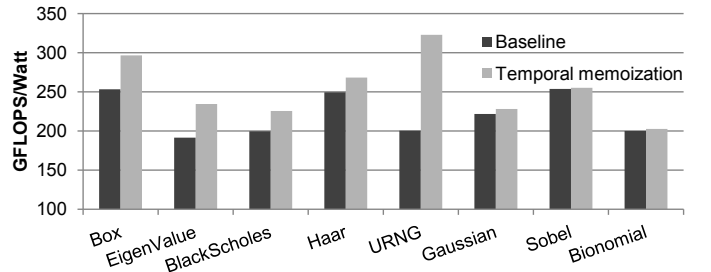


Fig. 5. Energy efficiency of the proposed temporal memoization architecture.

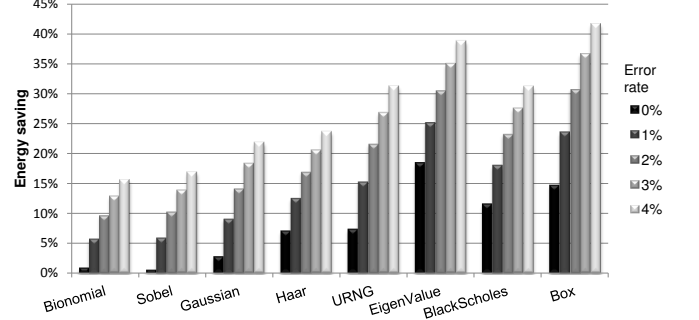


Fig. 6. Energy saving of the proposed temporal memoization architecture while experiencing a range of timing error rates.

modules are always ON; applications have only control on selecting either *approximate* or *exact* matching constraint. Fig. 5 compares the energy efficiency during execution of the applications. As shown, the proposed architecture with temporal memoization technique achieves better energy efficiency across all the applications, thanks to the high hit rates in the LUTs. Further, the technique is still able to maintain the energy efficiency and hide entirely the power cost of the temporal memoization modules in case of a poor locality environment; for example, if the overall hit rate drops by 46% (for EigenValue), 33% (for Black-Scholes), 25% (for Box), 20% (for Haar), and 11% (for URNG). On average, the proposed architecture reaches 16% higher GFLOPS per Watt compared to the baseline architecture.

We also compare the energy saving for a range of timing error rates [0%, 4%] on the execution state of the FPUs. Our implementation excludes the fact that the temporal memoization module may produce an erroneous result, because the module has a positive slack of 14% of the clock period. Therefore it is unlikely to face the timing errors. When there are no timing errors (0%), the proposed architecture has an average energy saving of 8% compared to the baseline architecture for all applications. This scenario is similar to the value prediction techniques as proposed in [16], that is extended to GPGPU architectures. Moreover, as shown in Fig 6, the temporal memoization technique has a great potential of energy saving in the high error rate circumstances. On average 14%, 20%, 24%, 28% lower total energy is achieved compared to the baseline at the error rates of 1%, 2%, 3%, 4%. This is accomplished through the efficient memoization-based timing error recovery that does not impose any latency penalty as opposed to the baseline recovery.

### C. Memoization-Based Voltage Overscaling

We also assess the efficacy of the proposed architecture in the voltage overscaling regime, while maintaining constant clock frequency. We scale down the voltage of the FPUs in the range of 0.9V–0.8V. To ensure always correct functionality

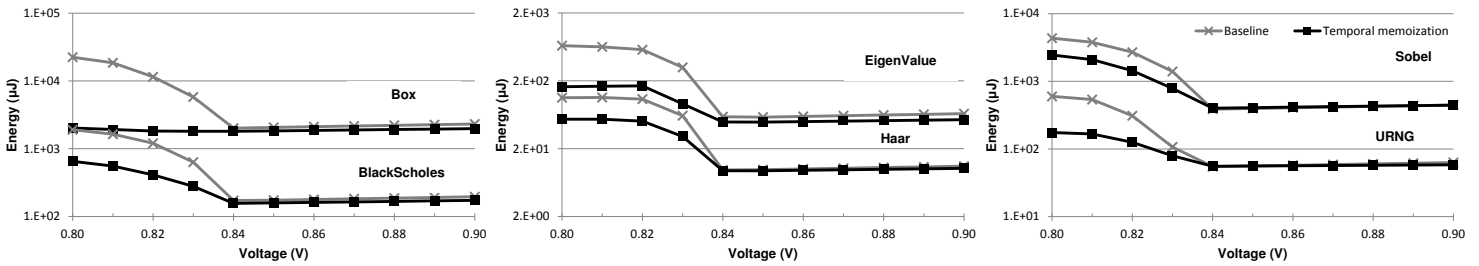


Fig. 7. Total energy consumption of the proposed temporal memoization architecture in comparison with the baseline architecture under voltage overscaling.

of the temporal memoization modules, we maintain their operating voltage at the fixed nominal 0.9V. Voltage scaling feature of *Synopsys PrimeTime* is employed to analyze the delay and power variations under the voltage overscaling. Then, the voltage overscaling-induced delay is back annotated to the post-layout simulation which is coupled with Multi2Sim simulator to quantify the timing error rate. The baseline architecture triggers the recovery mechanism when any voltage overscaling-induced timing error occurs, while our proposed architecture does it in case of simultaneous events of the miss and the error.

Fig. 7 illustrates the energy consumption of the two architectures at different voltage overscaling points for each application. The proposed memoization architecture exhibits a great potential of survival in the voltage overscaling. On average, for different applications: (i) The proposed architecture achieves 8% average energy saving at the nominal voltage of 0.9V. (ii) Scaling down the voltage to 0.84V, reduces the gain of our energy saving to 6% since the FPU's of the baseline are nicely reduced their total power as consequence of negligible error rate, while we cannot proportionally scale down the power of the temporal memoization modules. (iii) Beyond 0.84V, our technique surpasses the baseline architecture due to the abrupt increasing of the error rate and therefore frequent recoveries. At voltage of 0.8V, the proposed technique reaches an average energy saving of 66%.

## VI. CONCLUSION

We exploit value locality in improving timing error correction in GPGPUs. A fast lightweight temporal memoization module independently stores recent error-free executions of a FPU which is sufficient enough to temporarily protect individual FPU's against timing errors. To efficiently reuse computations, the technique supports both exact and approximate error corrections for error-intolerant general-purpose and error-tolerant image processing applications, respectively. These real-world applications exhibit a low entropy, that is high contextual information, yielding an average hit rate of 76% on the 4-entry LUTs. This avoids costly recovery, therefore improves the energy efficiency and reduces the total energy by average savings of 8% (for 0% timing error rate) to 28% (for 4% timing error rate). This technique also surpasses the baseline architecture by enhancing robustness and energy saving in the voltage overscaling regime. Our ongoing work is focused on realizing ultra low-power temporal memorization circuits using low-voltage memories, or ReRAM that could further enhance the energy savings.

## VII. ACKNOWLEDGMENTS

This work was supported by the NSF's Variability Expedition under award n. 1029783, ERC-AdG MultiTherman

(291125), and FP7 Virtical (288574).

## REFERENCES

- [1] *Whitepaper NVIDIA's Next Generation CUDA Compute Architecture: Kepler GK110*, 2012.
- [2] ITRS [Online]. Available: <http://public.itrs.net>
- [3] M.R. Kakoe, et al., *Variation-Tolerant Architecture for Ultra Low Power Shared-L1 Processor Clusters*, IEEE TCAS II, 59(12) (2012).
- [4] D. Jeon, et al., *Design Methodology for Voltage-Overscaled Ultra-Low-Power Systems*, IEEE TCAS II, 59(12) (2012) pp. 952-956.
- [5] P. Gupta, et al., *Underdesigned and Opportunistic Computing in Presence of Hardware Variability* IEEE TCAD, 32(1) (2013) pp. 489-499.
- [6] K.A. Bowman, et al., *Energy-Efficient and Metastability-Immune Resilient Circuits for Dynamic Variation Tolerance*, IEEE JSSC, 2009.
- [7] A. Raychowdhury, et al., *Tunable Replica Bits for Dynamic Variation Tolerance in 8T SRAM Arrays*, IEEE JSSC, 2011.
- [8] I. Kazi, et al., *A ReRAM-based non-volatile flip-flop with sub-VT read and CMOS voltage-compatible write*, Proc. IEEE NEWCAS, 2013.
- [9] A. Ghofrani, et al., *Towards Data Reliable Crossbar-based Memristive Memories*, Proc. IEEE ITC, pp. 6-13, 2013.
- [10] K.A. Bowman, et al., *A 45 nm Resilient Microprocessor Core for Dynamic Variation Tolerance*, IEEE JSSC, 2011.
- [11] M.-M. Papadopoulou, et al., *Micro-benchmarking the GT200 GPU* Technical report, Computer Group, ECE, University of Toronto, 2009.
- [12] R. Pawlowski, et al., *A 530mV 10-lane SIMD Processor with Variation Resiliency in 45nm SOI*, Proc. IEEE ISSCC, 2012, pp. 492-494.
- [13] A. Sodani and G.S. Sohi, *Dynamic Instruction Reuse*, Proc. IEEE/ACM ISCA, 1997, pp. 492-494.
- [14] C. A. Martinez, et al., *Dynamic Tolerance Region Computing for Multimedia*, IEEE Transactions on Computers, 61(5) (2012).
- [15] C. Alvarez, et al., *Fuzzy Memoization for Floating-point Multimedia Applications*, IEEE Transactions on Computers, 54(7) (2005).
- [16] M. -H. Lipasti, et al., *Value Locality and Load Value Prediction*, Proc. ACM ASPLOS, 1996, pp. 138-147.
- [17] A. Rahimi, et al. *Analysis of Instruction-level Vulnerability to Dynamic Voltage and Temperature Variations*, Proc. IEEE/ACM DATE, 2012.
- [18] A. Rahimi, et al., *Procedure hopping: A low overhead solution to mitigate variability in shared-L1 processor clusters*, Proc. ACM/IEEE ISLPED, 2012, pp. 415-420.
- [19] A. Rahimi, et al., *Hierarchically Focused Guardbanding: An Adaptive Approach to Mitigate PVT Variations and Aging*, Proc. ACM/IEEE DATE, 2013.
- [20] A. Rahimi, et al., *Aging-Aware Compiler-Directed VLIW Assignment for GPU Architectures*, Proc. ACM/IEEE DAC, 2013.
- [21] A. Rahimi, et al., *Spatial Memoization: Concurrent Instruction Reuse to Correct Timing Errors in SIMD Architectures*, IEEE TCAS II, 2013.
- [22] K. Chakraborty, et al., *Efficiently Tolerating Timing Violations in Pipelined Microprocessors*, Proc. ACM/IEEE DAC, 2013.
- [23] AMD APP SDK 2.5 [Online]. Available: [www.amd.com/stream](http://www.amd.com/stream)
- [24] M. A. Breuer, *Multi-media Applications and Imprecise Computation*, Proc. IEEE DSD, 2005.
- [25] Multi2Sim [Online]. Available: <http://www.multi2sim.org/>
- [26] FloPoCo [Online]. Available: <http://flopoco.gforge.inria.fr/>
- [27] AMD APP OpenCL Programming Guide, Chapter 6.6.1, pp. 157, 2012.