

# Reliable Computation on Unreliable Hardware: *Can We Have Our Digital Cake and Eat It?*

Vladimir Kiriansky

Saman Amarasinghe

{vlk, saman}@csail.mit.edu

## Abstract

*The digital abstraction allowed us to achieve unprecedented scalability in both hardware and software complexity. Yet, circuit level digital abstraction is becoming increasingly expensive to maintain. We show it's possible to raise the digital abstraction up to software layers and yet provide correctness guarantees.*

*Software modules must expose results to users and to external interfaces with guarantees on confidence and accuracy. A 100% confidence is the basis for the digital abstraction - it is easiest to reason about, and it enables modularity and scalability. We posit that software correctness should be failure-model-oblivious, input-model-oblivious, and usage-model-oblivious.*

*We propose using a combination of reliable and unreliable hardware resources, with a software stack of composable algorithm-specific and generic checkers to ensure correctness. We share preliminary performance results on RSA and SAT checkers, and our anecdotal experience with detected data corruptions on real production hardware.*

*We show we can have the correctness of reliable nodes while enjoying performance and power gains of unreliable nodes to flexibly optimize for throughput, latency, or energy.*

## 1. Introduction

Lessons from analog vs digital circuit design demonstrated the value of the digital abstraction - only logic producing clean binary 0 vs 1 outputs can be iterated infinitely, easily composed vertically, and cheaply scaled horizontally.

Circuit designers strive to ensure digital abstraction over the lifetime of a system in worst-case environment while subject to worst-case workload. Yet, the cost of maintaining the digital abstraction at the circuit level with voltage and frequency guardbands is leading to higher energy and performance loss [18], soon to approach inordinate 50% waste. We believe most efficient use of unreliable execution modes should allow errors to propagate to software, yet without giving up output correctness guarantees.

Real world decisions impacting billions of people and dollars can be effectively guided by computation systems only as well as design assumptions match the real world. Ideally, one would have a good understanding of a system's usage model, input data model, and hardware execution model. Rather than building systems as brittle and complex as our models, we argue system correctness should be designed obliviously to any model.

*Usage-model-oblivious* reusable software components should keep their correctness guarantees even in worst case

scenarios. We'll loosely define correctness, as having 100% confidence in accuracy bounds on outputs. When stakes are small, occasionally not knowing what the system might produce is OK. Some results may have perfectly acceptable accuracy for small input sizes, iterated a few times, on a few computation nodes, with outputs cascading to few other systems. However, composition, layering and replication of these modules with no accuracy bounds will lead to accumulation of errors. Just like tape copying on even the best analog systems suffers unavoidable degeneration, unlike *infinite* digital copying. At real world scale errors may start to dominate in utterly useless results.

*Input-model-oblivious* computation should not depend on input data ranges and patterns for correctness. Otherwise on some inputs the final result may be more affected by errors, or hardware faults may be more likely. Approximate computations that produce acceptable results on inputs drawn from tight uniform or normal distributions (e.g. average person height) may fail dramatically on inputs with wide power-law distributions (e.g. average income). Our programs must work on all real world inputs, not only on easy to synthesize test data.

Interdependence between hardware failure models and input models is even harder to track. We've evaluated vulnerability to delay errors of synthesized ALU pipelines in detailed *analog* SPICE simulations on inputs extracted from actual benchmarks [4]. Visible erroneous bit positions depend on the failing arithmetic logic structure, data ranges, and for static CMOS, depend on tuples of both current and previous inputs. The error rates of a 64-bit adder ranged from 8% to 93% on different benchmarks.

*Failure-model-oblivious* computation must give correctness guarantees without over-relying on an error model of all failure classes and their interdependence. This will be a necessity when using new unreliable hardware operating in high variability regimes. Variations differ by their temporal and spatial locality, e.g. permanent manufacturing defects, aging effects over years, localized thermal events over  $\mu s$ , voltage droop over  $ns$ , and fast and localized capacitive crosstalk events over  $ps$  [15]. It's unlikely system designers would be able to keep up with new hardware and new failure types and predict circuit failure impact on data and control logic and resulting errors visible in architectural state. We'll probably know less and less about the magnitude and frequency of worst case events to fit them in neatly calculable probability distributions.

We will show how to continue building robust and scalable software systems harvesting power and performance gains of

unreliable hardware, while the correctness safety net depends on adequately reliable hardware and software. Our high level system architecture is based on efficient software *checkers* running on reliable hardware nodes or modes of execution. Current production hardware is considered adequately resilient and can be hardened further by derating from specifications to a lower performance but much higher reliability mode. For example, an insightful study [14] of a million PCs showed CPUs underclocked below factory ratings to be 39% more reliable.

We’ll start with our high level design overview, show early experiments on existing hardware, and then discuss observed errors and related techniques.

## 2. Design Overview

We believe the abstraction boundary can be redrawn to relax the circuit-level and micro-architecture contracts, while software layers take charge of correctness guarantees.

We are working on a generic checker framework for error detection in any program, as well as efficient algorithm-specific checkers. We’ll cover in future full publications the design and performance characteristics of our generic checkers, and composability of nested checkers. We’ll detail below design patterns for algorithm-specific checkers.

### 2.1. Design Criteria

For the rest of this paper, our default design choices are for general purpose systems that give correctness guarantees with 100% confidence in 100%- $\epsilon$  accuracy and worst case 100% latency, limited by maximum power and available resources. Any uncertainty is shifted solely to performance metrics over short time scales (e.g. latency of 99.9% of requests in one minute), while we try to minimize long term average latency, resource utilization, and energy.

We assume that confidentiality is either not an issue, e.g. data and code are not critical if exposed, or that unreliable nodes have no egress path to public networks. We maintain data and service availability within worst case service SLA – recovery must always be possible from data available on reliable nodes (from memory, or durable storage). Replication among unreliable nodes is useful only to increase throughput and latency for common case independent failures.

We assume reliable nodes use a practically reliable configuration with a negligible failure probability. Reliable node computation, communication, and storage resources are insulated from unreliable nodes. Unreliable nodes do not have to be insulated from each other as their resources are considered completely unreliable.

Design criteria for all layers below the new abstraction line diverge between the two types of nodes. Many existing techniques discussed in Section 4 can be added everywhere, but trade-offs of performance vs reliability change dramatically. For example, a more resilient OS design to survive more crashes might offer higher availability on unreliable nodes at

low cost. On unreliable systems we only need a *low* error rate with minimal overhead for best-effort correctness and availability, while expensive reliability techniques can be used on the fewer reliable nodes with little performance impact.

### 2.2. Algorithm-specific Checkers

In our survey of classic algorithms we’ve found many important computations for which checking a result is much cheaper than solving the problem. We’ll call *solver* any original program whether solving a decision problem, search problem, or an optimization problem. An acceptable output is either an exact output in discrete problems, or an output within tolerance for accuracy. We’ll call a *checker* a program that dynamically checks whether a solver output is acceptable for the inputs of a given problem instance.

Algorithm-specific checkers are often asymptotically faster, and for realistic input sizes need 10–1000x less resources than solvers. Overall system reliability is determined by the reliability of checkers, while performance and energy is dominated by unreliable solvers.

Program checking for many important problems has been investigated over the years [24], though often as “a theoretical curiosity” [22] guided by different motivations than our design principles. In our settings only asymmetric checkers faster than a solver are useful, yet a fast checker with low communication overhead can be used to check the integrity of any solver, even if running unreliable software on unreliable hardware in untrusted cloud settings.

A checker must be *sound* – it should never accept an invalid output. However, checkers don’t have to be *complete* with respect to the original solver – some valid outputs may be too hard to check. Thus a *constrained solver* may produce only easy to check acceptable outputs. An *augmented solver* dynamically generates a certificate as a proof for a solver execution which can be efficiently checked. It’s usually better to shift work to augmented solvers which would execute on unreliable nodes more efficiently, while keeping checkers simple and fast.

Little domain expertise is needed to understand what makes an algorithm result acceptable, feasible, and/or optimal. Usually the theoretical groundwork proving the original solver sound and complete is easy to reuse. A verification tool may be useful to verify correctness of simple checkers, for example, an automated equivalence checker [16] can be used to verify that for all inputs a checker accepts only the outputs of a reliable solver.

Probabilistic checkers can often be constructed, but to keep things simple (for statistically challenged users), and for easy composability for developers, our default checker choices are simple and deterministic. For example, our sorting checker is deterministic linear instead of the known probabilistic  $O(n)$  checker by hashing [24, 21].

Algorithm-specific checkers are useful even if worst case asymptotic time is no better than best solver time, when aver-

age case time or hidden constants are better. We are therefore considering language support for expressing common patterns of checking and batch checking, eager use of speculative results, and lazy checking. Checking of intermediate results should be deferred as much as possible, as it’s often ultimately unnecessary when final outputs can be verified cheaply or with just a few intermediate results.

Algorithm extensions and framework support can also help offloading and overlapping communication, efficient check-pointing, robust forward error recovery, progress and progress rate tracking.

### 3. Experiments

We show preliminary results of a couple of our algorithm-specific checkers from our growing algorithm-specific collection. Rather than simulating unreliable execution, we *stimulated* realistic hardware errors on production processors by configuring them out of specifications.

#### 3.1. RSA Decryption Checker in `openssl`

The Secure Sockets Layer (SSL) cryptographic protocol is a cornerstone of secure Internet communication, critical for Web commerce and increasingly used with rising privacy concerns. RSA decryption is the most computationally intensive server component of a new session handshake, especially after the recent transition to 2048-bit keys [3].

Our RSA checker takes advantage of a performance asymmetry between RSA decryption and encryption. Even though both depend on modular exponentiation, the exponent  $e$  of an RSA public key  $(e, n)$  is usually a small odd integer, while its inverse  $d$  in an RSA private key  $(d, n)$  has the full  $\beta$ -bitlength.<sup>1</sup> Modular exponentiation with an RSA public key  $M^e$  takes only  $O(\beta^2)$  operations, while applying a private key  $M^d$  takes  $O(\beta^3)$  [8].

We’ve implemented an RSA accelerator in `openssl-1.0.1e` offloading the expensive private key operations to unreliable hosts over the network and checking the results back on the reliable host with much cheaper encryption. Whenever encryption of decrypted ciphertext matches the original ciphertext, we’d mark unreliable decryption successful. Traditional RSA correctness theorems prove that for a given message  $M$ , its corresponding ciphertext  $C = M^e \pmod n$  can be decrypted with  $C^d = M^{ed} \equiv M \pmod n$ . If we plug in  $C$  above, we see that  $C^{de} \equiv C^{ed} \equiv C \pmod n$ .

On our ARM1176 processors rated at 700MHz, public key checking is 35x faster than private key exponentiation for 2048-bit keys, and 65x faster for 4096-bit keys, as seen in Table 1. After taking into account CPU overhead for communication, a single reliable checker can check the results of 25–50 unreliable solvers. We introduced unreliability by over-clocking nodes to 990MHz, which correspondingly increased

Key length bits	Decrypt/sec $M^d$	Encrypt/sec $M^e$	Speedup Ratio
512	7085 ops	74206 ops	10x
1024	1496 ops	26400 ops	18x
2048	238 ops	8337 ops	35x
4096	36 ops	2331 ops	65x

Table 1: RSA throughput on ARM1176.

decryption throughput by 40%, and reduced latency even over 100Mbps Ethernet.

Over one day of execution we detected 11 mismatches in RSA decryption/encryption pairs. These would otherwise have been silent data corruptions! Our unreliable node also suffered an order of magnitude more Linux kernel and application crashes. The reduced availability resulted simply in throughput loss as our requests are short-lived. Tight timeouts would eliminate latency losses since the expected computation and network time are well known by the caller. More interestingly, in a week-long study of 250+ kernel crash back-traces, we observed bursty crashes highly skewed to a few locations. We have yet to analyze causes and effects for these highly correlated errors: whether faults are due to self-induced voltage droops that can be smoothed [17]; whether highly vulnerable locations use data patterns that make errors more likely; whether circuit failures result in datapath or control logic errors; whether catastrophic system failures are due to affecting software control flow.

We will discuss our further lessons from these observations in Section 4.

#### 3.2. SAT and UNSAT satisfiability Checkers

Proving satisfiability of a 3-CNF boolean formula is a classic NP-complete problem. Checking whether an assignment is satisfiable takes linear time, while finding an assignment (SAT) or proving that no such assignment exists (UNSAT) are  $\Omega(2^n)$  problems. Yet many real world instances can be efficiently solved thus boolean formulas are a very useful modeling framework. A SAT solver balances the two tasks of quickly finding a satisfiable assignment, or learning quickly a set of lemmas that can prove unsatisfiability.

We implemented an offline SAT checker in the `minisat` [11] SAT solver. We evaluated the checker on an Intel i7-2600K using datasets and solvers from the SAT Challenge 2012 [1]. Checking satisfiable instance results was 4–40,000x faster than solving them with `minisat-2.2.0` and the best overall solver `glucose-2.1`. We have not yet detected data corruption errors with this workload on unreliable systems.

Checking unsatisfiable instances (UNSAT) is harder. Although many modern SAT solvers are able to produce a proof of unsatisfiability with low overhead[13], proof checking sometimes is as hard as solving the problem. Even harder is checking indeterminate instances where a SAT solver times out before finding a satisfiable assignment or a proof of un-

<sup>1</sup>Typically  $e = 65537$ , thus  $\lg e = O(1)$ , while  $\lg d \leq \beta$ ,  $\lg n \leq \beta$ .

satisfiability. We’re working on efficient on-line checkers for unsatisfiable and indeterminate instances.

## 4. Discussion and Related Techniques

We draw three main conclusions from these preliminary experiments. First, data corruption errors are real, as not all errors are catastrophic! Second, fail-stop errors still may be a lot more frequent than data corruptions. Third, errors are likely following a power law distribution in both time and space. (We extrapolate from the fail-stop errors, that data corruptions also follow a power law distribution). We then discuss applicability of known reliability and availability techniques for correctness and performance in these settings.

### 4.1. Failure Models

There’s little convincing data to show what real errors look like, and our real hardware and analog experiments are no exception. Many fault tolerant design evaluations assume uniform random bit-flip errors. Even if memory soft errors were so well behaved, a fault model for delay errors is not readily available. A good fault model would have to capture the probability distribution of error structure and magnitude, temporal and spatial locality, different input sequences vulnerability, and be adjusted for aging and environmental conditions.

In our anecdotal experience, errors are rare but with large magnitude. We suspect both temporal and spatial locality in errors, highly correlated to both workload and environment conditions, thus we can’t make guarantees about their independence. For correctness, we prefer to not make assumptions that we know anything for certain about the error model, e.g. frequency of errors in a given interval, magnitude, burstiness.

Software may need to react dynamically to the observed error rate to determine the best error avoidance, error detection, and error recovery technique. Graceful performance degradation if failures become worse than expected is desirable. Prediction would be better if we can fit each device observations to a failure model with further dependence on system configuration, workload, and environment.

### 4.2. Reliability Techniques

Practical software error detection techniques are needed to handle data corruptions to guarantee correctness.

We can’t put a bound on the number of possible errors before an error detection opportunity. Venerable techniques like dual or triple modular redundancy [23] are rendered ineffective when errors are correlated, and there is no upper bound on errors in any interval.

Our algorithm-specific checkers don’t make any correctness assumptions based on predicted error distribution.

Techniques like Rely [6] can be used as reliability amplification to increase the likelihood of unreliable computations being correct. We can still use approximate failure models to optimize performance and power in the common case. Overall

efficiency results would improve as our models for each device get better, even if we accept that our error models may not cover worst case scenarios.

Delay errors in current ALU designs affect disproportionately more-significant bits, unlike some approximate computing technique assumptions. New problem-specific ALU designs [12] may be more resilient. Generic solutions at circuit [5, 9] and architectural [20, 2] levels focused on correctness of every pipeline stage or instruction may reduce guardbands but are expensive in the absence of errors, increase design complexity, increase buffering [5], or fail in the presence of control logic errors [9]. At low error rate we believe it’s simpler and better to leave software to handle detection and recovery.

The ERSA [7] prototype evaluation demonstrates several probabilistic applications are inherently error-resilient and even at high rates of injected errors converge timely to 90% or better accuracy. We share many of the same design criteria but want to address most deterministic algorithms and give 100% accuracy guarantee.

### 4.3. Availability Techniques

The frequency of fail-stop errors is dominating in our experiments. More typical deployment configurations would preferably not be exposed to as high error rates. More resilient OSes, fast reboot, and tight health tracking would increase the system availability. Failure-oblivious [19] methods can be used to further increase availability – we don’t need the OS to be fail-safe. By the end-to-end principle, it’s only the application results and data that matter, and all externally visible effects are gated by reliable nodes.

Applications can similarly be fault-tolerant to fail-stop errors in the traditional meaning in distributed systems; failure-oblivious when detection is timely and results are partially or less likely to be impacted; or *fail-fast* when errors are rare and recovery fast.

### 4.4. Recovery Techniques

Efficient backward or forward recovery techniques become critical for both fail-stop (erasures) and corruption errors at higher error rates.

Computations organized as small re-executable tasks with many short-lived computations are naturally easy to recover. Large problem instances need to be solved either using larger number of systems and/or longer time for the execution of each algorithmic step. The probability of experiencing an error increases, and algorithm convergence rate or termination may be impacted. Generic recovery techniques can be used, or more efficient algorithm-specific error recovery techniques.

Thanks to reliability and availability techniques, both a wrong response and no response would get treated as a slow response. When a late result is no better than no result, large scale distributed systems are already deploying tail-tolerant [10] techniques to keep latency in check.

## 5. Conclusion

We've shown our new digital cake recipe - we embrace approximate computing with a base of unreliable hardware, add unreliable software layers, but we top it with a *soft digital* icing. We then have always correct results, while enjoying performance and energy efficiency; making simple or no application modifications. Not only we'll get to have the cake and eat it; it's (mostly) free too.

## 6. Acknowledgements

We thank George Kurian for helping with our analog SPICE simulations, and Vijay Ganesh for sharing his expertise on SAT and UNSAT challenges.

## References

- [1] SAT Challenge 2012. <http://baldur.itl.kit.edu/SAT-Challenge-2012/>.
- [2] T.M. Austin. DIVA: a reliable substrate for deep submicron microarchitecture design. In *Microarchitecture. MICRO-32. Proceedings 32nd Annual International Symposium on*, pages 196–207, 1999.
- [3] Elaine Barker and Allen Roginsky. Transitions: Recommendation for transitioning the use of cryptographic algorithms and key lengths. *NIST Special Publication*, 800:131A, 2011.
- [4] Christian Bienia. *Benchmarking Modern Multiprocessors*. PhD thesis, Princeton University, 2011.
- [5] K.A. Bowman, J.W. Tschanz, S.L. Lu, P.A. Aseron, M.M. Khellah, A. Raychowdhury, B.M. Geuskens, C. Tokunaga, C.B. Wilkerson, T. Karnik, and V.K. De. A 45 nm resilient microprocessor core for dynamic variation tolerance. *Solid-State Circuits, IEEE Journal of*, 46(1):194–208, Jan. 2011.
- [6] Michael Carbin, Sasa Misailovic, and Martin C. Rinard. Verifying quantitative reliability for programs that execute on unreliable hardware. In *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages and Applications, OOPSLA '13*, pages 33–52, New York, NY, USA, 2013.
- [7] Hyungmin Cho, Larkhoon Leem, and Subhasish Mitra. ERSA: Error resilient system architecture for probabilistic applications. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 31(4):546–558, April 2012.
- [8] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms, 3rd ed.* MIT Press, Cambridge, MA, 2009.
- [9] S. Das, C. Tokunaga, S. Pant, W.-H. Ma, S. Kalaiselvan, K. Lai, D.M. Bull, and D.T. Blaauw. RazorII: In situ error detection and correction for pvt and ser tolerance. *Solid-State Circuits, IEEE Journal of*, 44(1):32–48, jan. 2009.
- [10] Jeffrey Dean and Luiz André Barroso. The tail at scale. *Communications of the ACM*, 56(2):74–80, 2013.
- [11] Niklas Eén and Niklas Sörensson. An extensible SAT-solver. In *Theory and applications of satisfiability testing*, pages 502–518. Springer, 2004.
- [12] Jie Han and Michael Orshansky. Approximate computing: An emerging paradigm for energy-efficient design. In *IEEE ETS*, May 2013.
- [13] Marijn JH Heule, Warren A Hunt Jr, and Nathan Wetzler. Verifying refutations with extended resolution. In *Automated Deduction-CADE-24*, pages 345–359. Springer, 2013.
- [14] Edmund B Nightingale, John R Douceur, and Vince Orgovan. Cycles, cells and platters: an empirical analysis of hardware failures on a million consumer PCs. In *Proceedings of the Sixth conference on Computer Systems*, pages 343–356. ACM, 2011.
- [15] Jan M. Rabaey, Anantha Chandrakasan, and Borivoje Nikolić. *Digital integrated circuits : a design perspective*. Prentice Hall electronics and VLSI series. Prentice-Hall, Upper Saddle River, NJ, 2003.
- [16] David A Ramos and Dawson R Engler. Practical, low-effort equivalence verification of real code. In *Computer Aided Verification*, pages 669–685. Springer, 2011.
- [17] V.J. Reddi, S. Kanev, Wonyoung Kim, S. Campanoni, M.D. Smith, Gu-Yeon Wei, and D. Brooks. Voltage smoothing: Characterizing and mitigating voltage noise in production processors via software-guided thread scheduling. In *Microarchitecture (MICRO), 2010 43rd Annual IEEE/ACM International Symposium on*, pages 77–88, Dec. 2010.
- [18] V.J. Reddi, S. Kanev, Wonyoung Kim, S. Campanoni, M.D. Smith, Gu-Yeon Wei, and D. Brooks. Voltage noise in production processors. *Micro, IEEE*, 31(1):20–28, Jan.-Feb. 2011.
- [19] Martin Rinard, Cristian Cadar, Daniel Dumitran, Daniel M. Roy, Tudor Leu, and William S. Beebe, Jr. Enhancing server availability and security through failure-oblivious computing. In *Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation - Volume 6, OSDI'04*, pages 21–21, Berkeley, CA, USA, 2004. USENIX Association.
- [20] E. Rotenberg. AR-SMT: a microarchitectural approach to fault tolerance in microprocessors. In *Fault-Tolerant Computing, 1999. Digest of Papers. Twenty-Ninth Annual International Symposium on*, pages 84–91, 1999.
- [21] Nirmal R. Saxena and Edward J. McCluskey. Linear complexity assertions for sorting. *Software Engineering, IEEE Transactions on*, 20(6):424–431, 1994.
- [22] Justin Thaler, Mike Roberts, Michael Mitzenmacher, and Hanspeter Pfister. Verifiable computation with massively parallel interactive proofs. *CoRR*, abs/1202.1350, 2012.
- [23] John Von Neumann. Probabilistic logics and the synthesis of reliable organisms from unreliable components. *Automata studies*, 34:43–98, 1956.
- [24] Hal Wasserman and Manuel Blum. Software reliability via run-time result-checking. *J. ACM*, 44(6):826–849, November 1997.