

The Importance of Temporal and Spatial Temperature Gradients in IC Reliability Analysis

UNIV. OF VIRGINIA DEPT. OF COMPUTER SCIENCE TECH. REPORT CS-2004-07

JANUARY 2004

Wei Huang[†], Zhijian Lu[†], Shougata Ghosh[†], John Lach[†], Mircea Stan[†], Kevin Skadron[‡]

Departments of [†]Electrical and Computer Engineering and [‡]Computer Science, University of Virginia
Charlottesville, VA 22904

Abstract

Existing IC reliability models assume a uniform, typically worst-case, operating temperature, but temporal and spatial temperature variations affect expected device lifetime. This paper presents a model that accounts for temperature gradients, dramatically improving interconnect and gate-oxide lifetime prediction accuracy. By modeling expected lifetime as a resource that is consumed over time at a temperature-dependent rate, substantial design margin can be reclaimed and/or less expensive cooling systems may be used. *This report is superseded by TR CS-2004-08.*

Keywords—reliability, electromigration, gate oxide breakdown, thermal management, gradient

1. Introduction

As CMOS technology continues to scale, power density of VLSI circuits has unfortunately been scaling too, and this rapid increase is widely expected to continue. Yet tolerable operating temperatures usually remain fixed from generation to generation, independent of power density. This is because many aging mechanisms in VLSI circuits proceed at a rate that is temperature dependent, and product-lifetime requirements dictate the maximum aging rate. For most applications, maximum temperatures are therefore limited to around 100–120°C. Unfortunately, this means that rising power densities impose rising cooling costs that may eventually become so prohibitive that they limit the development of new products.

Historically, the temperature dependence of many aging processes can be empirically modeled by the Arrhenius Equation

$$\text{MTF} = \text{MTF}_0 \exp\left(\frac{E_a}{kT}\right) \quad (1)$$

where MTF_0 is the mean time to failure at a specified reference temperature, E_a is the activation energy of the failure, and k is the Boltzmann constant. A detailed model based on the physics of interconnect electromigration (EM) and gate-oxide breakdown, two common temperature-dependent IC aging processes, can be found in [1] and [2].

However, these models do not account for dynamic temporal or spatial temperature variations, which have been experimentally shown to have a significant impact on circuit lifetime [3]. Using existing models, circuit designers must assume a constant temperature (usually the worst-case temperature) for the entire circuit, resulting in inaccurate lifetime estimations and excessive design margins.

This paper presents a temperature-dependent reliability model for interconnect EM. The most important contribution in this model is the ability to take into account temporal and spatial temperature gradients and, therefore, more accurately predict circuit lifetime. We also show that the same modeling approach can be used to account for temporal temperature gradients when modeling gate-oxide breakdown. Compared to a standard design methodology that uses a worst-case temperature threshold to achieve an expected lifetime, reliability models that account for temporal and spatial gradients are able to extract higher performance given a particular cost constraint, or to reduce cooling costs while maintaining performance. A particularly attractive approach for taking advantage of information about temperature gradients is to treat circuit lifetime as a resource that is consumed over time by temperature, which fits well with recent techniques for *dynamic*, runtime thermal-management techniques that adapt to specific workload behavior to maximize circuit performance while still meeting reliability requirements. A chip can, in real time, use periods of low temperature to offset brief periods of higher temperature.

2. Interconnect Reliability Analysis With Temperature Gradients

Interconnect EM is the process of self-diffusion due to momentum exchange between electrons and metal atoms. Clement [1] presents a 1-D analytical model of EM-induced interconnect stress build-up that matches empirical measurements. An interconnect failure occurs when the stress reaches a threshold value σ_{th} . The stress build-up caused by atom dislocation can be described by the following equation [1]

$$\frac{\partial \sigma}{\partial t} - \frac{\partial}{\partial x} \left[D_a \left(\frac{B\Omega}{kTL^2\varepsilon} \right) \left(\frac{\partial \sigma}{\partial x} - \frac{q^*LE}{\Omega} \right) \right] = 0 \quad (2)$$

where $\sigma(x, t)$ is the stress function, x is the 1-D position and has been normalized to $[-1, 0]$ for simplicity; D_a is the diffusivity of metal atoms, which has the Arrhenius form of temperature dependency: $D_a = D_{a0} \exp(-Q/kT)$; B , Ω and ε are constants depending on the properties of the metal, the surrounding materials and the interconnect aspect ratio; L is the characteristic length of the interconnect; q^* is the effective charge; and E is the applied electric field, which is equal to ρj , the product of resistivity and current density. The term q^*LE/Ω corresponds to a back-flow flux due to the electric field and determines the steady-state stress solution of Eq.(2).

Clement's reliability analysis was performed assuming a temporally and spatially uniform temperature [1]. In the following two subsections, we extend his model to investigate the effects of temporal and spatial gradients to interconnect lifetime predictions.

A. Temporal temperature gradients

From Eq.(2), if we define

$$\beta(T) = D_a \frac{B\Omega}{kTL^2\varepsilon} = D_{a0} \frac{B\Omega}{kTL^2\varepsilon} \exp(-Q/kT) \quad (3)$$

without the spatial temperature gradients, we have

$$\frac{\partial \sigma}{\partial t} - \beta(T) \frac{\partial}{\partial x} \left(\frac{\partial \sigma}{\partial x} - \frac{q^*\rho j L}{\Omega} \right) = 0 \quad (4)$$

From Eq.(4), we can easily verify that for two different $\beta(T_1)$ and $\beta(T_2)$, we have

$$\sigma_2(x, t) = \sigma_1 \left[x, \left(\frac{\beta(T_2)}{\beta(T_1)} t \right) \right] \quad (5)$$

where σ_1 and σ_2 are stress build-up solutions at temperatures T_1 and T_2 . This shows the advantage of defining β according to Eq.(3)—Eq.(5) implies that $\sigma_1(t_1) = \sigma_2(t_2) = \sigma_3(t_3) = \dots$, as long as $\beta(T_1)t_1 = \beta(T_2)t_2 = \beta(T_3)t_3 = \dots$. β can therefore loosely be considered a rate of aging that is dependent on each wire's characteristics. We can also define a "transformed time" $\varphi_{th} = \beta(T)E(t_f)$, where $E(t_f)$ is the expected time to failure. φ_{th} is an attractive quantity because it is invariant with respect to individual wire characteristics (β).

Figure 1 illustrates the numerical solutions to Eq.(4) by showing the maximum stress build-up along the interconnect for several different values of $\beta(T)$, with the same initial and boundary conditions. Steady-state stress is determined by the term $q^*\rho j L/\Omega$. As dictated by Eq.(5) and seen in Figure 1, when

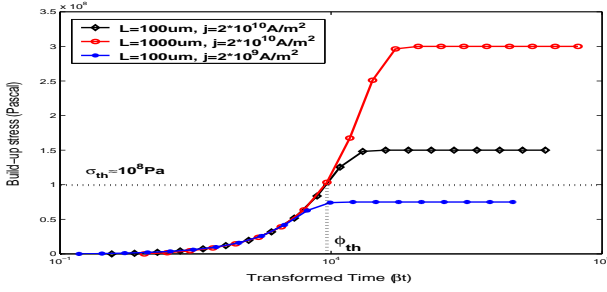


Fig. 1. Numerical solutions of Eq.(4) follow the same track before failure for different values of β .

plotted as a function of transformed time, solutions for different β values follow the same track before reaching their steady states. Also shown in Figure 1 are an example threshold stress $\sigma_{th} = 10^8$ Pa (which is within the range of typical values for critical stress failures) and the corresponding transformed time $\varphi_{th} = \beta(T)E(t_f)$. It follows that $E(t_f) = \varphi_{th}/\beta(T)$ is the true time to failure, and $1/E(t_f) = \beta(T)/\varphi_{th}$ is the failure rate. Notice that φ_{th} is constant for all of the interconnects that will eventually fail. In some cases, if $j \cdot L$ is less than a critical value $(j \cdot L)_c$, the steady-state stress build-up is below σ_{th} , and the interconnect will never fail [4].

If temperature is a function of time, $T(t)$, we can summarize the above analysis by

$$\varphi_{th} = \int_0^{E(t_f)} \beta(T(t)) dt = E(\beta(T(t))) \cdot E(t_f) \quad (6)$$

where $E(\beta(T(t)))$ is the expected value of $\beta(T)$ over time $E(t_f)$. The expected time to failure is therefore

$$E(t_f) = \frac{\varphi_{th}}{E(\beta(T(t)))} \quad (7)$$

A key implication of this new model is the *increased flexibility for dynamic thermal management*. Existing techniques typically specify a fixed temperature threshold based on a defined lifetime goal, and the hardware must ensure that the operating temperature never exceeds that threshold. Our analysis shows that this limitation is overly conservative. Instead, Eq.(6) reveals that interconnect lifetime can be modeled as a resource that is consumed during circuit operation at a rate of $\beta(T)$. Therefore, it is safe to overshoot the fixed temperature threshold for limited periods of time in order to obtain higher circuit performance, and then later compensate for the excess lifetime “consumption” with lower temperatures to ensure the lifetime requirements are still satisfied. Other dynamic thermal management techniques that exploit the flexibility provided by this new model will be explored as part of future work.

In addition, it can be shown that $\beta(T(t))$ is a convex function with respect to temperature within the normal operating temperature range. By applying Jensen’s inequality for convex functions, we have

$$E(\beta(T(t))) \geq \beta(E(T(t)))$$

which, according to Eq.(7), leads to another key observation: *a constant temperature T will always yield a longer expected lifetime than a time-varying temperature with an average of T* . Damping temporal variations can therefore permit a higher operating temperature and performance while maintaining the specified expected lifetime.

B. Spatial temperature gradients

The other aspect of temperature-aware interconnect EM analysis is the effect of temperature variation over the length of a wire, that is, interconnect temperature as a function of position, $T(x)$. This is a concern because atom diffusivity, which is exponentially dependent on temperature as shown in Eq.(3), varies along the wire with temperature. This leads to a different stress build-up distribution than that for a spatially constant wire temperature.

To illustrate the importance of spatial gradients, we consider several simple temperature profiles along a wire: 1) current flows

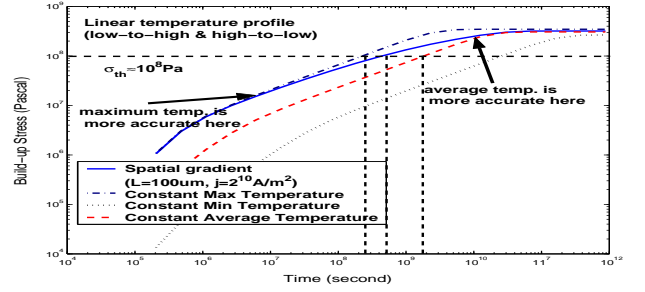


Fig. 2. Solutions for a linear spatial temperature gradient along a wire, together with solutions to constant max, min and average temperatures.

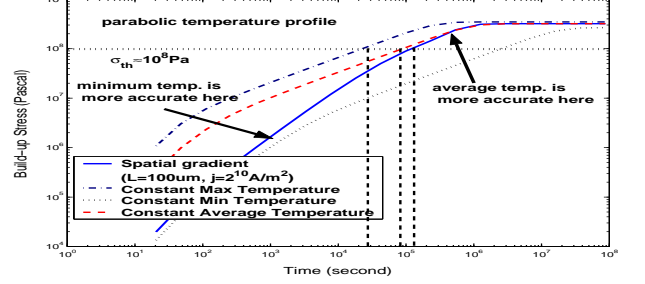


Fig. 3. Solutions for a parabolic spatial temperature gradient along a wire, together with solutions to constant max, min and average temperatures.

from a hot to a cold region of the design with a linear temperature profile; 2) the reverse, with current flowing from cold to hot with a linear temperature profile; and 3) a temperature profile that is parabolic, with the lowest temperatures at the ends and maximum at the center of the wire, as might be found if the wire is within a circuit area of uniform power distribution that is surrounded by a colder region.

By numerically solving Eq.(2) for each of the three temperature profiles, we can plot the maximum stress build-up on the wire over time, as shown in Figures 2–3. These plots also include the stress build-up curves corresponding to three different spatially constant temperatures: maximum temperature $T_{max} = 150^\circ\text{C}$, minimum temperature $T_{min} = 60^\circ\text{C}$ and the average wire temperature for the spatial gradient profiles above). The critical stress for the wire to fail is again set to 10^8 Pa.

It is clear from Figures 2–3 that the actual stress build-up of a wire with spatial temperature gradients cannot be accurately modeled using a spatially constant temperature. With the linear temperature profiles in Figure 2, the actual stress build-up is bounded above by that of the maximum temperature and below by that of the average temperature. Interestingly, when the stress build-up is small (e.g., less than 10^7 Pa), the spatial gradient EM can be well approximated by that of the constant maximum temperature, but when the stress build-up is close to saturation, the spatial gradient plot is similar to that of the constant average temperature.

Although the high-to-low and low-to-high temperature profiles have the same maximal values of stress build-up over time, the direction of the stress in these two cases are opposite. A wire with a low-to-high temperature gradient in the direction of the current is more prone to an “open” failure, because the diffusion of metal atoms (i.e., tensile stress) is more severe at the high-temperature end of the wire. Conversely, “short” failures are more common for wires with high-to-low temperature gradients in the direction of the current, because accumulation of metal atoms (i.e., compressive stress) is more severe at the high-temperature end of the interconnect.

Of even greater interest is the parabolic gradient in Figure 3, where the stress build-up is bounded above by the *average* case, even though some parts of the wire operate at temperatures higher than that average. Detailed analysis reveals that the symmetric temperature distribution along the wire has a damping effect on atom migration, reducing the diffusion speed incurred by the constant maximum temperature.

In the above analysis, we assume that the spatial profile of the temperature is unchanged over time. Future work includes developing an interconnect reliability model that combines both temporal and spatial gradients—in other words, the spatial gradi-

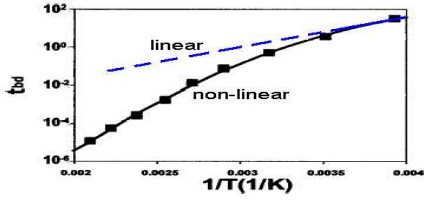


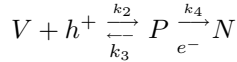
Fig. 4. Arrhenius plot for thin oxide breakdown showing non-linear trend with greater slope at high temperature, which means the breakdown process exacerbates at higher temperatures for thinner gate oxide. (adapted from [2])

ent itself exhibits variations over time. So far we are proposing a conservative approach to solve this problem by finding the worst spatial gradient and applying our analysis to that.

3. Gate-Oxide Reliability Analysis With Temperature Gradients

Generally speaking, CMOS gate-oxide breakdown is a complicated process caused by hot electrons and holes together with other microscopic reactions. Over the years, the thickness of CMOS device gate-oxide has been scaled down to only a few molecular layers as CMOS technology advances. With thinner gate-oxide, Figure 4, adapted from [2], shows that the time-to-break-down follows a nonlinear Arrhenius trend with respect to device operating temperature. This means that gate-oxide breakdown will become worse for future technologies with thinner gate-oxides and higher device operating temperatures. Reliability analysis for gate oxides is therefore growing in importance.

The correct physical model for gate-oxide breakdown is far from being settled. Here we merely wish to observe that the same approach that we described in the previous section will also improve the precision of modeling gate-oxide breakdown. We base our analysis here on work by Cheung [2], who argues that thin oxide breakdown is due to the accumulation of neutral defects. The involved reactions are



where V is a precursor site, P is an intermediate site, and N is a neutral defect. He proposes a model for the calculation of neutral defect concentration, which can well explain the fast temperature acceleration in ultra-thin oxide-breakdown and the nonlinear Arrhenius behavior. The model is

$$[N(t)] = \frac{[V]_0 k_2 J_h k_4 J_e}{(k_3 + k_4 J_e)^2} (1 - e^{-(k_3 + k_4 J_e)t}) \quad (8)$$

where $[N(t)]$ is the time-dependent neutral defect concentration which has a threshold value $[N]_{th}$. If $[N(t)] > [N]_{th}$, the oxide is considered broken-down; $[V]_0$ is the initial concentration of the precursor sites which can be converted to neutral defect sites by first capturing holes/protons followed by capturing electrons. J_h and J_e are the hole and electron current density through the gate oxide. k_2, k_3 and k_4 are the rates for the above microscopic reactions, and they all have the Arrhenius form of temperature dependency (B_2, B_3 and B_4 are constants.)

$$k_{2,3,4} = B_{2,3,4} \exp\left(\frac{-Q_{2,3,4}}{kT}\right)$$

Thus, if we define

$$\gamma(T) = k_3 + k_4 J_e = B_3 \exp\left(\frac{-Q_3}{kT}\right) + J_e B_4 \exp\left(\frac{-Q_4}{kT}\right) \quad (9)$$

and consider the term $\alpha(T) = ([V]_0 k_2 J_h k_4 J_e) / (k_3 + k_4 J_e)^2$ in Eq.(8) as a constant because $\alpha(T)$ is weakly dependent on temperature compared to the term $e^{-\gamma(T)t}$, which is super-exponential, we have

$$[N(t)] \approx \alpha(1 - e^{-\gamma(T)t}) \quad (10)$$

We find that we can use exactly the same method we used in Section 2(A) when analyzing the effect of temporal temperature

gradients on interconnect reliability: define $\psi_{th} = \gamma(T)t_{bd}$ as the transformed time at which the defect concentration reaches $[N]_{th}$. Therefore, $t_{bd} = \psi_{th}/\gamma(T)$ is the time-to-breakdown.

If temperature varies over time (i.e., there is temporal temperature gradient $T(t)$ at a particular device) we have

$$\psi_{th} = \int_0^{t_{bd}} \gamma(T(t)) dt$$

Thus, $\gamma(T(t))$ is the rate at which the gate-oxide's lifetime is consumed. Similar to Section 2(A), we can argue that analysis based on a uniform worst-case temperature is overly conservative. We have only considered temporal temperature gradients in this section because CMOS devices are so tiny that spatial temperature gradients do not apply to individual devices for gate-oxide reliability analysis.

4. Temperature-Aware Design Using Temperature Gradients

In this section, we present an interconnect reliability case study showing the benefits of temperature-aware reliability analysis. The effects of temporal and spatial temperature gradients on EM failure are investigated for a simulated $0.13\mu\text{m}$ microprocessor. Interconnect temperatures are obtained from a new, validated compact thermal model we have developed that divides the chip area into a fine grid of equal-area cells. We briefly introduce the thermal model and its utility to the VLSI design community, and then present and discuss the simulation results. In this paper, we do not perform a case study for gate-oxide reliability analysis due to the fact that the appropriate values of k_3 and k_4 were not discovered in the literature. Such a study will be part of future work, but based on the similarity of the analysis, we expect that accounting for temporal temperature variation in gate-oxide reliability would yield similar benefits as for interconnect EM.

A. A grid-like compact thermal model

The accuracy of reliability analyses using the models presented here depends greatly on accurate and detailed temperature estimations. A useful model must be able to simulate transient, not just steady-state, temperatures. It should also be parameterized so that a correct model is generated regardless of the materials, layout, or thermal package. Finally, different levels of granularity should be supported, so that thermal analysis can be a part of the design process from early architectural studies through different circuit-design stages, including detailed design, placement, and routing.

We have extended our prior HotSpot model [5], which meets the first two requirements but was only designed to model temperature at a microarchitecture level of granularity. Instead of the ad-hoc structure of that original model, we have developed a grid-based model that accommodates the required range of granularities. The grid-based model is also necessary for obtaining data for spatial temperature gradients. For example, the grid can be set to the granularity of individual fundamental circuit structures, standard cells, functional units, etc. The grid-based approach is therefore more general. Its adjustable grid size enables designers to perform temperature-related analyses at any level of granularity, including the microarchitecture level targeted by the original model. It is also useful to note that the new model can be easily extended to a multi-resolution grid with fine resolution in critical portions of the circuit and coarse resolution elsewhere.

Figure 5 shows the structure of the model. It consists of a grid of cells, with each cell consisting of thermal resistors and capacitors representing silicon, heat spreader, thermal interface material and heat sink. Validation shows that an essential component of the model (in addition to the detailed model of the die and the various components of the thermal package) is the thermal interface material between the layers of the thermal package. (After modeling the interface material, the spatial temperature gradient across the die increases to a range of 30 to 50 degrees, which is more accurate than the gradient of about 10 degrees predicted in [5].) This extended thermal model is also able to predict interconnect temperatures by adding layers of thermal resistors and capacitors representing interconnect metal layers together with inter-layer dielectrics and vias that are above each of the silicon grids. The numerical solution for the grid-based compact model follows the same algorithm used in [5].

Another important improvement over our prior work is physical validation of the thermal model, which we have performed by comparing the estimated silicon surface temperatures from the model against those obtained with a commercial thermal test

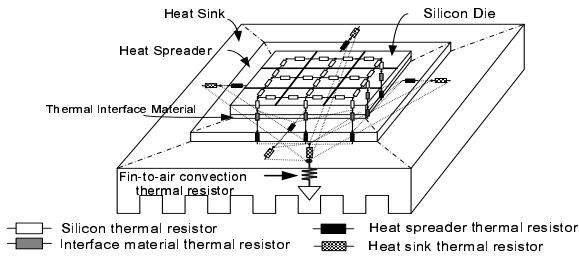


Fig. 5. Example structure of the grid-like compact thermal model with 3x3 grids, thermal interface material, heat spreader and heat sink. Thermal capacitors and heat sources are not drawn for clarity.

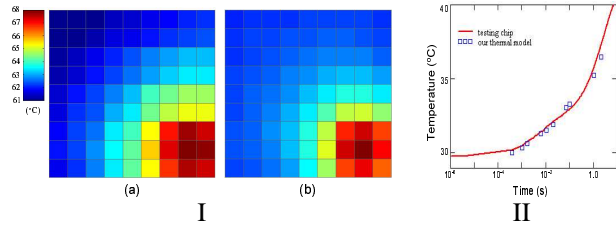


Fig. 6. (I)—Steady-state validation of the compact thermal model: (a) Testing chip measurements (b) Results from the model with errors less than 5%. (II)—Transient validation of the compact thermal model. Percentage error is less than 7%. (Transient temperature response of one power dissipater is shown here.)

chip [6]. The test chip has a 9x9 grid of power dissipating resistors, which can be turned on or off individually, and a corresponding 9x9 grid of temperature sensors; and can measure both steady-state and transient temperatures for each of the power dissipators. We build the same 9x9 grid-like chip structure in our thermal model, turn on specific sets of power dissipators in the test chip, and assign exactly the same power values at the same locations in the grid-based thermal model. Figure 6 (I) compares steady-state thermal plots for one particular pattern; and Figure 6 (II) compares transient temperature measurements. Power density in this experiment is $0.5\text{W}/\text{mm}^2$ in the heat-dissipating area. As can be seen, our compact thermal model is quite accurate, with the worst case error for steady-state and transient temperatures less than 5% and 7%, respectively. Similar results were obtained for other power-dissipation patterns.

B. A case study: Impact of temperature gradients for typical program behavior

To demonstrate the benefits of accounting for temporal and spatial temperature gradients, we present a case study using temperature values obtained from simulating a microprocessor with characteristics similar to a $0.13\mu\text{m}$ Alpha 21364. Using the described grid-based compact thermal model, we can obtain the steady-state and transient temperature responses of the devices and interconnect. For example, the transient thermal behavior of interconnect above the floating point register is shown in Figure 7 for the SPEC2000 benchmark program *applu*, revealing obvious temporal temperature gradients. The thermal-package characteristics we used in simulating *applu* were derived so that the constant temperature value that, according to Section 2A, yields the same expected lifetime as the pattern in Figure 7 is 110°C (a common limit). This temperature has also been plotted in Figure 7 as a straight line.

These results illustrate the potential benefits of accounting for temporal variation. If the lifetime budget is used to dictate only some fixed worst-case temperature (e.g., 110°C), then a more expensive cooling solution is required to bring *applu*'s actual behavior within specification while achieving the same performance. The alternative is that the voltage and clock speed must be reduced, or a dynamic thermal management technique must be engaged to reduce processor activity and enforce the 110° limit whenever the operating temperature exceeds the threshold. Using microarchitecture simulation techniques described in [5], we estimate that selecting a lower design point for voltage and frequency would require a 13% reduction in clock frequency; and that dynamic thermal management would reduce performance by about 10% using dynamic voltage scaling and 50% using fine-grained fetch gating. If temporal temperature variation is taken into account, none of these costly solutions are needed!

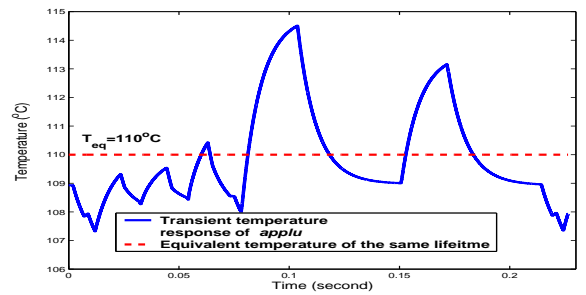


Fig. 7. Transient temperature response of floating point register interconnects. Constant operating temperature for the same interconnect lifetime is also shown.

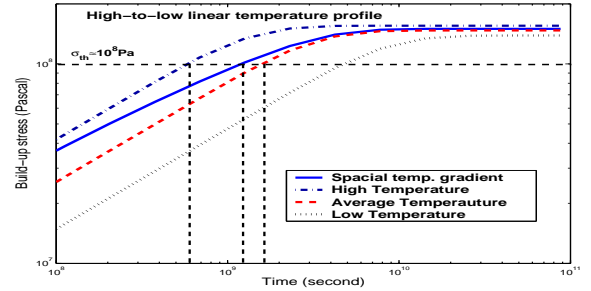


Fig. 8. Lifetime prediction comparison for a wire with and without a spatial temperature gradient.

Now consider spatial temperature gradients along the interconnect. From the compact thermal model, for a typical point in time we find that the temperature of the interconnect above the floating point register is about 110°C , while the temperature for the interconnect above the adjacent floating point queue is about 73°C , a 37°C gradient for interconnect from the floating point register to the floating point queue. From our temperature-aware analysis shown in Figure 8, which plots the stress build-up over time for this particular spatial-gradient pattern, we can see that assuming a uniform average temperature along the interconnect overestimates expected lifetime by 30% compared to the actual case, and using a uniform maximum underestimates expected lifetime by 80%. This means that the typical approach of worst-case analysis will yield drastically lower temperature specifications than necessary, which again would require an unnecessarily expensive cooling solution or unnecessary performance sacrifices.

5. Conclusions

In this paper, we have presented a new approach to interconnect and device gate-oxide reliability analysis that accounts for temporal and spatial variations in temperature. We have also developed and validated a dynamic compact thermal model for simulating, at various levels of detail, the time-dependent evolution of on-chip temperatures across an IC. Reliability analysis using temporal and spatial gradient values obtained from a real application on a simulated processor show the importance of accounting for temperature gradients. Worst-case analysis can drastically underestimate expected lifetime, requiring either unnecessarily aggressive and costly cooling solutions or else reductions in power dissipation that incur unnecessary sacrifices in IC performance. We propose that, instead of designing for a maximum tolerated temperature based on a worst-case analysis, expected lifetime should be viewed as a resource that is consumed over time at a temperature-dependent rate. This dynamic, reliability-driven approach to managing operating temperature fits particularly well with the recent advent of dynamic thermal management techniques; and this paper shows that lifetime requirements are the proper objective function rather than fixed temperature thresholds. This modeling approach will be extended to other failure mechanisms and integrated with techniques for real-time thermal management.

Acknowledgments

This work is supported in part by the National Science Foundation under grant nos. CCR-0105626, CCR-0133634, two grants from Intel MRL, and a grant from the University of Virginia Fund for Excellence in Science and Technology. The au-

thors would also like to thank Karthik Sankaranarayanan for his help with the thermal model simulations.

References

- [1] J.J. Clement, *J. Applied Physics*, vol. 82, pp. 5991–6000, Dec. 1997.
- [2] K.P. Cheung, *Applied Physics Let.*, vol. 83, pp. 2399–2401, Sep. 2003.
- [3] C.J.M. Lasance, *Microelectronics and Reliability*, vol. 43, pp. 1969–74, Dec. 2003.
- [4] J.J. Clement, *IEEE Trans. on Device and Materials Reliability*, vol. 1, pp. 33–42, Mar. 2001.
- [5] K. Skadron, M.R. Stan, W. Huang, S. Velusamy, K. Sankaranarayanan and D. Tarjan, *Proc. of the 30th Int'l Symp. on Computer Architecture*, pp. 2–13, June 2003.
- [6] V. Székely, C. Márta, M. Renze, G. Végh, Z. Benedek and S. Török, *IEEE Trans. on Components, Packaging, and Manufacturing Technology–Part A*, vol. 21, pp. 399–405, Sep. 1998.