# A Gate Level Sensor Network for Integrated Circuits Temperature Monitoring

*Alireza Vahdatpour, Saro Meguerdichian, Miodrag Potkonjak*

Computer Science Department
University of California Los Angeles
{alireza, saro, miodrag}@cs.ucla.edu

*Abstract*— **We present the first sensor network architecture to monitor integrated circuits (IC) thermal and energy activity. The sensor network consists of a set of simple gates, which are superimposed over the actual design of any IC. The sensing network and the actual IC design are completely disjoint in order to enable their simultaneous operation. Since the delay of gates is proportional to their temperature, we can obtain temperature of the network gates, by measuring the delay of the gates in the self-sensing network. Once we measured the delay of the circuit, we use CMOS temperature-delay relation and linear programming formulation to calculate the temperature at any point on the chip. High resolution (spatial and temporal) temperature monitoring allows several run-time optimizations. Protecting shared processors from permanent localized damage through rapid creation of hot spots and efficient accounting of the available energy supply are among two main applications of our IC sensor network.**

## I.    INTRODUCTION

Both long and short term technology trends strongly indicate that exponential growth of microprocessor, graphics and digital signal processors, DRAM and FLASH memory, and field-programmable gate arrays (FPGA) will not just be maintained, but even accelerated. For example, Intel already reported a processor with 80 cores and more than 2 billion transistors. The first Intel graphic processor, Larrabee, will have several thousand processors. Several network processors also have thousands of processors.

One of the key ramifications of shrinking feature size of integrated circuits is continuously increasing density of switching and, therefore, power generation. Increased power consumption directly results into increased temperature, that as its consequence has faster aging of transistors due to NBTI and HCI effects, degradation of interconnect due to electromigration, slowdown of both logic and interconnect, and additional increase in power consumption. The resistance and the delay of interconnect is in particularly strongly affected if its parts are subject to sharply different temperatures. If not addressed, the positive feedback power-temperature inevitably destroys integrated circuits.

In addition, temperature variation makes side-channels attacks easier and more effective (e.g. both electromagnetic (EM) radiation and gate-level power consumption is increased and becomes controllable). In summary, it is important and necessary to monitor the thermal activity of ICs  for detecting and preventing security attacks. In a simple attack, a user can send a request for program execution that creates hot spots, parts of IC with very high temperature. Similar scenario can happen in ASIC where an attacker can feed set up input vectors to the IC, so that the dynamic switching activity of IC increases significantly. In more sophisticated attacks, the observation of changes in operating characteristics of IC can be used for not just local, but also remote extraction of private information. Interesting observation is that it can be also used for hidden information exchange, but that issue is beyond the scope of this paper. Finally, note that collecting information about activity of the circuits can also be used for enforcement of digital right management and the development of better design strategies. For example, the knowledge of which parts of a design are subject to high activity and temperature may results in better placement and routing solutions.  Another solution can be, for each licensed program certain parts of the circuit that are prone to high temperature will require a unique ID to activate otherwise operations will be rescheduled or reassigned. Therefore, it is not surprising that thermal measurements and management has attracted a great deal of research and development attention. However, the current solutions require the addition of thermal sensors that are large, expensive, and slow. Also, they have limited accuracy and cannot address the expressed manufacturing variability of pending silicon technologies.

In this paper, we propose a fundamentally different approach. The basic idea is to add a small network of gates so that the delay measurement of each gate is fast and accurate. The gates are small and do not require separate technology. The measurements are directly taken from the IC and can be done in one or few clock cycles at nanosecond speed. Since there is a well defined dependency of speed and temperature, it is easy to calculate the later in a single clock cycle using look-up table. In addition we impose grid over the IC in such a

way that each its field has the same temperature. Using the information about the change of temperature in each field in a time interval, we can calculate its generated heat and, therefore, its activity. This information can be also used for calculation of temperature at an arbitrary point of the IC using statistical interpolation techniques. We introduce the concept of using simple networks of gates with a small number of simple gates (e.g. smallest, fastest, and the lowest power gate - inverter) to obtain information about speed change at a small set of location on an IC. The change of speed can be used to calculate the temperature and locally generated heat. In longer term, it can be also used for calculation of aging factors of the local gates. The application of linear programming and interpolation techniques is used to calculate change in temperature and generated heat at an arbitrary point of the IC. In the second section, some related work and previous studies are introduced and discussed. Section III consists a brief introduction over the temperature-delay model and delay measuring mechanism that we are leveraging. Section IV introduces the sensing network and section V presents the results. Finally section VI concluded the paper.

## II. RELATED WORK AND BACKGROUND

In the last two decades, energy emerged as a premier design and run-time management metric. Thermal management has attracted a great deal of attention [1]. There is significant interest in modeling and quantifying the impact of high thermal gradients on clock circuitry and interconnect in general [2]. There are recent studies (i.e. [3]) on using digital design blocks for temperature measurement. The measurement circuits in such studies are generally large, and therefore, temperature readings are average values for wide areas in an IC. In addition, approaches by [4] and [5] use processor specific features such as performance counters to estimate processor temperature in the software during the runtime, which is applicable to few processor families supporting the feature and leads to high performance overhead in execution time. Recently, Potkonjak et. al [6] have shown how the post silicon gate-level characterization (in presence of manufacturing variability) is viable using input vector control techniques. We use similar methodology to measure the delay of the sensor gates. While [7] has presented a temperature management technique for leakage minimization, our approach has fundamental differences. Study in [7] assumes the thermal behaviors of the task and processor are known a priori. In addition, it assumes that the leakage energy consumption of the processor is negligible in sleep mode, which requires extra power gating and supply control logic to be implemented in the processor. In addition, up to our knowledge, none of the studies in leakage management considered the spatial variation of the processor temperature.

## III. MODELS

The impact of temperature on delay has been studied widely. Depending on the fabrication process and the circuit technology, several timing-temperature models have been proposed (e.g. [8, 9]). In this study, we used the temperature-delay model introduced by [9]. According to this study, the variation of the delay is about 6%, when temperature varies by 50oC (starting from 50°C). Figure 5 depicts the difference

in the delay-temperature relations for two simple chains, which is due to the differences in the gate size and other characteristics such as capacitance. We use the scheme from Figure 5.a to measure the delay of a logic circuit. Clock cycle period can be dynamically changed in a circuit by the resolution of pico seconds [10]. At the same time, considering gigahertz and megahertz clock frequencies of integrated circuits (where the clock period is larger than hundreds of pico seconds), high-resolution dynamic management of clock frequency is feasible. To accurately measure the propagation delay of a circuit, it is enough to continuously and slightly increase the clock frequency of the measurement circuit (see Figure 2), until the clock cycle period becomes less than the propagation delay of the circuit. Upon reaching this point, the change in the output of the circuit (Output) will not affect FF1 and Out will be different than Output. This clock cycle period can be assumed to be the propagation delay of the logic circuit. Figures 2.b and 2.c depict the conditions where the clock period is slightly more and less than the propagation delay ($t_0$) of the combinational circuit of Figure 2.a. Note that in this example, the Output of the combinational circuit is assumed to become high as the input becomes high. As the clock period becomes closer to the propagation delay, flip-flop metastability may happen as a result of violation of the flip-flops setup and hold times. Since each measurement only takes one clock cycle, measurements can be repeated to reduce the error impact.

## IV. SELF-SENSING STRUCTURE

Figure 1 depicts a sample structure of the proposed temperature monitoring circuit. In general, this network consists of m rows and n columns of inverters; Each row consists of n inverters and n-1 multiplexers and each column is a chain of m inverters. The multiplexer placement is such that 2n-1 combinations of inverters chains can be made on each row by changing the select inputs of the multiplexers (connecting inverters from different rows to each other). As a result, in addition to the n vertical inverter chains, $m.2^{n-1}$ horizontal inverter chains can be constructed dynamically in the network. The depicted network has only one input and two output ports. Flip-flops $FF_{H0}...FF_{Hm-1}$ and $FF_{V0}..FF_{Vn-1}$ are placed on the input side of the cascaded inverters, and the other end of the inverter chains are connected to exclusive-or gates. The chain of flip-flops results in activation of a single horizontal and a single vertical inverter chain in each clock cycle (fed with the new value of the input, considering the input of the circuit is toggled every *max(m, n)* cycles) and cause the change in the value of Output_Horiz and Output_Vert. Note that the output of the exclusive-or gate will change anytime one of its inputs changes. As described earlier, by changing the clock frequency of the flip-flops on both sides of the inverter chains, the propagation delay can be measured. Upon measuring the propagation delay of the inverter chains, we use a linear programming approach to calculate each multiplexer and inverter delay from the measured chain delays. In the proposed linear program, the main variables are high-to-low and low-to-high delay values for each gate (multiplexer and inverter) in the circuit. We also introduce error variables associated with each vertical or

horizontal path to compensate for the effect of possible rounding error in delay measurement and wire delay. To obtain the constraint equations, let us consider the network of Figure 1 (m = 4; n = 3). Consider the input of the circuit is set to high (logical one). The measured delay represents the sum of the delays of the gates that are in the activated path. Therefore, the following linear equations are valid for propagation delays of every vertical inverter chain (for the sake of simplicity, assume that both of the outputs are zero before the input is set to high):

$$e'_{V_i} + d'_{V_{XOR}} + \sum_{j=0}^{\lfloor m/2 \rfloor} d_{V_{2j,i}} + \sum_{j=1}^{\lceil m/2 \rceil} d'_{V_{2j-1,i}} = D_{V_i}, \forall i = 0, 2, 4, \ldots$$

$$e_{V_i} + d_{V_{XOR}} + \sum_{j=0}^{\lfloor m/2 \rfloor} d_{V_{2j,i}} + \sum_{j=1}^{\lceil m/2 \rceil} d'_{V_{2j-1,i}} = D_{V_i}, \forall i = 1, 3, 5, \ldots$$

Here, $d_{V_{j,i}}$ represents the high-to-low delay value of the inverter in vertical inverter chain in the jth row and ith column. $DV_i$ is the measured delay of the ith vertical inverter chain (which is measured when $FF_{V_i}$ has passed the new input to the ith vertical chain). $d_{V_{XOR}}$ is the high-to-low delay of XORV and $e_{V_i}$ is the error variable for the chain, when the output is changing from high to low. Similarly $d_0$ and $e_0$ denote the same gate delay and error when the signal is changing from low to high. For the horizontal inverter chains, assuming that all the select inputs of the multiplexers are set to 0, the following equations are valid:

$$e'_{H_j} + d'_{H_{XOR}} + \sum_{i=0}^{\lfloor n/2 \rfloor} d_{H_{j,2i}} + \sum_{i=1}^{\lceil n/2 \rceil} d'_{H_{j,2i-1}} + \sum_{i=1}^{\lfloor n/2 \rfloor} d_{M_{j,2i}}$$
$$+ \sum_{i=1}^{\lceil n/2 \rceil} d'_{M_{j,2i-1}} = D_{H_j}, \quad \forall j = 0, 2, 4, \ldots$$

$$e_{H_j} + d_{H_{XOR}} + \sum_{i=0}^{\lfloor n/2 \rfloor} d_{H_{j,2i}} + \sum_{i=1}^{\lceil n/2 \rceil} d'_{H_{j,2i-1}} + \sum_{i=1}^{\lfloor n/2 \rfloor} d_{V_{j,2i}}$$
$$+ \sum_{i=1}^{\lceil n/2 \rceil} d'_{V_{j,2i-1}} = D_{H_j}, \quad \forall j = 1, 3, 5, \ldots$$

where $d_{H_{j,i}}$ and $d_{M_{j,i}}$ denote the high-to-low delay value of inverters and multiplexers at the jth row and ith column on the horizontal gate chains, respectively. Similar notation as of the previous equations is used for error and measured path delay values. By changing the select input of the multiplexers, the total number of $m.2n-1$ different combinations of inverters and multiplexers are constructed, and by measuring the path delay, the same number of equations can be derived. In addition, for the case when input is set to low, a series of constraints with similar structure for horizontal and vertical chains are constructed, which we omit for brevity. In order to minimize the effect of manufacturing variability, inverters are designed to be big enough. In addition, in the manufacturing process, inverters of the horizontal and vertical paths at the location j; i will be made with the same size and will be placed in close proximity, therefore their sensed temperature and their delay would be the same. Hence: $d_{V_{j,i}} = d_{H_{j,i}}$

Also, high-to-low and low-to-high propagation delays are linearly related, since they are linear functions of the Load Factor of the gate:

$$d = 0{:}05 + 0{:}017L; \; d^{'} = 0{:}02 + 0{:}038L$$

Finally, the objective function of the linear program is:

$$minimize \sum_{j=0}^{m \cdot 2^{n-1}} (|e_{H_j}| + |e'_{H_j}|) + \sum_{i=0}^{n-1} (|e_{V_i}| + |e'_{V_i}|)$$

The above equations construct constraints for a linear program with m.(2n-1) main variables (gate delays) and 2n+2m.2n-1 constraints. Upon solving the linear equation, estimated delays for each inverter are used to estimate the temperature of the IC at the sensory circuit gate locations. We use the model introduced earlier to calculate the temperature variation using propagation delay values. Thereafter, a two-dimensional linear interpolation is used to extend the temperature estimation for the whole IC surface.

## V. RESULTS

We extensively simulated our IC sensing network using the HotSpot software. Fig. 3 presents a sample temperature map generated from our sensing network, compared to a temperature map generated by HotSpot temperature simulation tool. It is clear that even using a relatively small network of sensing gates, high resolution temperature map can be obtained. Note that the difference between our technique and the compared HotSpot tool is that our technique is a measurement system while HotSpot is a simulation technique. Fig. 4 shows the maximum measurement error using different sizes of sensing network. As presented, the more gates are used in the network, the higher is the accuracy of temperature measurement. In addition, the difference in thermal activity of the applications results in different measurement accuracies. Higher spatial gradient of temperature results in higher measurement inaccuracy.

## VI. CONCLUSION

We have developed the first approach for IC self-sensing using a combination of architectural, design, algorithmic, and statistical techniques. The addition of a small, fast, and low power network enables a procedure for measuring and calculating gate delay and temperature, as well as thermal activity at any arbitrary position in a user specified IC. The approach can be used for addressing numerous IC run-time management tasks and the detection of security attacks. Our simulations indicate the approach's effectiveness in detecting hot spots and potential thermal attacks.

REFERENCES

[1] J Donald, M Martonosi, Techniques for Multicore Thermal Management: Classi_cation and New Exploration, ISCA 2006, pp.78-88.

[2] AH Ajami, K Banerjee and M Pedram, Modeling and analysis of nonuniform substrate temperature e_ects on global ULSI interconnects, IEEE Trans. Comp. Aided Design Integrated Circuits Systems, 24 (6) (2005), pp. 849-861.

[3] P Chen et. al., A Fully Digital Time-Domain Smart Temperature Sensor Realized With 140 FPGA Logic Elements, IEEE Transactions on Circuits and Systems, 54(12): 2661-2668, 2007.

[4] KJ Lee, K Skadron, Using Performance Counters for Runtime Temperature Sensing in High-Performance Processors, IPDPS 2005, page 232.1, DC, USA.

[5] A Merkel, F Bellosa, Balancing power consumption in multiprocessor systems, ACM SIGOPS European Conference on Computer Systems, 2006, pp 403-414.

[6] Miodrag Potkonjak, Ani Nahapetian, Michael Nelson, Tammara Massey, Hardware Trojan horse detection using gate-level characterization, Design Automation Conference DAC, 2009.

[7] L Yuan, S Leventhal, G Qu, Temperature-aware leakage minimization technique for real-time systems. ICCAD 2006, pp. 761-764, CA, USA.

[8] BP Das et. al., Voltage and Temperature Scalable Gate Delay and Slew Models Including Intra-Gate Variations, International Conference on VLSI Design, 2008.

[9] [10] C de Benito, S Bota, JL Rossell_o, J Segura, Temperature impact on multiple-input CMOS gates delay, Proceedings of SPIE, 2007.

[10] J Kalisz, Review of methods for time interval measurements with picosecond resolution, Metrologis, vol. 41, pp. 17-32, 2004
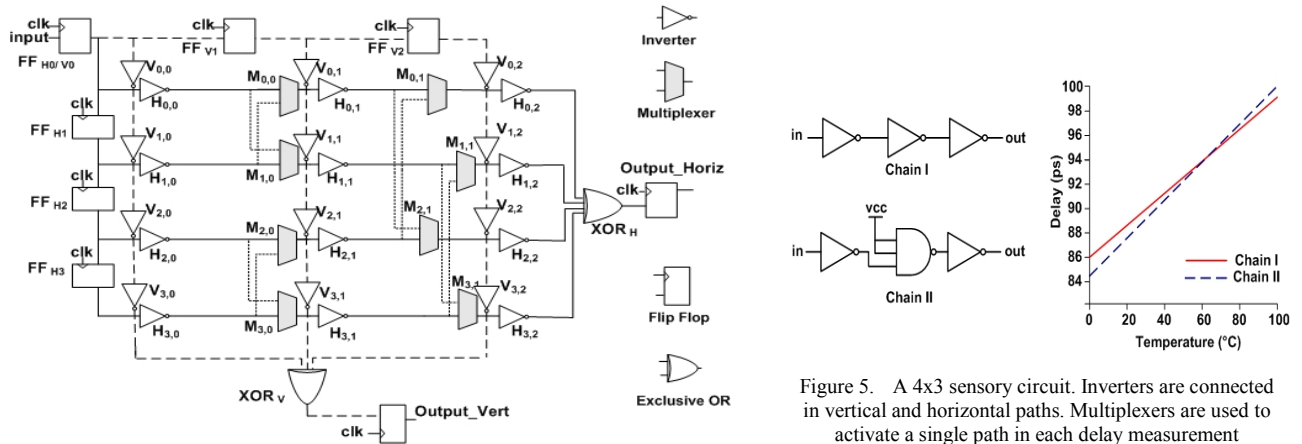
Figure 1. A 4x3 sensory circuit. Inverters are connected in vertical and horizontal paths. Multiplexers are used to activate a single path in each delay measurement



Figure 5. A 4x3 sensory circuit. Inverters are connected in vertical and horizontal paths. Multiplexers are used to activate a single path in each delay measurement
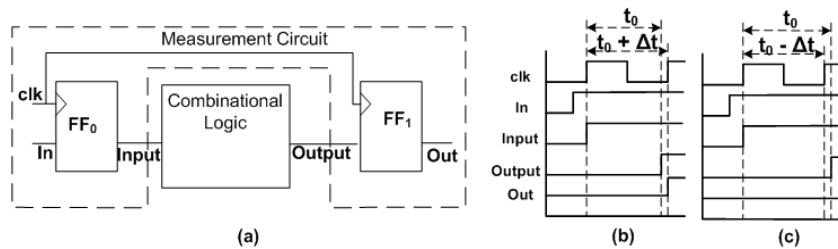


Figure 2. a) A simple mechanism for delay measurement b-c) Circuit behavior when propagation delay is smaller and greater than the clock period
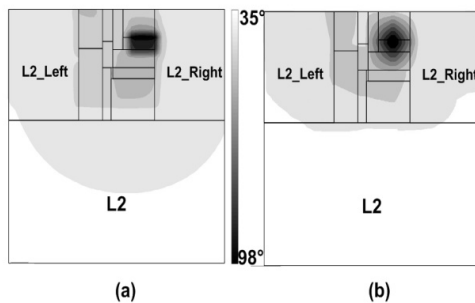


Figure 3. Fig. 3: Temperature map of an Alpha processor executing a benchmark application. (a) High resolution simulation (b) Estimation by an 8x8 sensor network
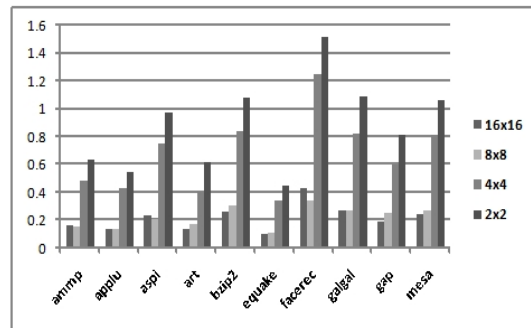


Figure 4. Fig. 4: The maximum measurement error of temperature sensing vs. the size of sensing network, for different SPEC2000 benchmark applications

655