11.2 A 1Mb Embedded NOR Flash Memory with 39µW Program Power for mm-Scale High-Temperature Sensor Nodes

Qing Dong¹, Yejoong Kim¹, Inhee Lee¹, Myungjoon Choi¹, Ziyun Li¹, Jingcheng Wang¹, Kaiyuan Yang¹, Yen-Po Chen¹, Junjie Dong¹, Minchang Cho¹, Gyouho Kim¹, Wei-Keng Chang², Yun-Sheng Chen², Yu-Der Chih², David Blaauw¹, Dennis Sylvester¹

¹University of Michigan, Ann Arbor, MI ²TSMC, Hsinchu, Taiwan

Miniature sensor nodes are ideal for monitoring environmental conditions in emerging applications such as oil exploration. One key requirement for sensor nodes is embedded non-volatile memory for compact and retentive data storage in the event that the sensor power source is exhausted. Non-volatile memory also allows for near-zero standby power modes, which are particularly challenging to achieve at high temperatures when using SRAM in standby due to the exponential rise in leakage with temperature, which rapidly degrades battery life (Fig. 11.2.1). However, traditional NOR flash requires mW-level program and erase power, which cannot be sustained by mm-scale batteries with internal resistances >10k Ω To address this issue, we propose an ultra-low power NOR flash design and demonstrate its integration into a complete sensor system that is specifically designed for environmental monitoring under high temperature conditions: such as when injected into geothermal or oil wells.

The proposed flash design reduces power consumption by using a combined Dickson and ladder pump topology and low-parasitic MIM capacitors to generate 13V with 73% power efficiency, and a cross-sampling current sense amplifier (SA) that doubles the sensing margin. Measured program and erase energy is 49pJ/b and 9.4pJ/b, resulting in a $30\times$ and $22\times$ reduction compared to a standard flash macro. Program power is 39μ W and 82μ W at 25° C and 125° C, enabling a miniature sensor node system to be powered by only two 8μ Ah batteries. The complete measured system is $3.88\times1.70\times1.85$ mm³ and consumes only 190nW of power at 125° C in the flash-enabled deep sleep mode, resulting in a $63\times$ power reduction compared to a conventional sensor design.

Embedded split-gate NOR flash requires >10V for hot-carrier-injection-based program and tunneling-based erase. Circuits to generate these high voltages dominate the write power [1,2]. A cross-coupled Dickson charge pump offers the best power efficiency, but the voltage across the flying capacitor increases linearly in each stage, making it necessary to use high-voltage MOS capacitors. However, these HV MOS capacitors have a parasitic/useful capacitance ratio of 46% (NMOS) and 18% (PMOS), lowering pump efficiency [3]. MIM capacitors offer a very low (1%) parasitic loss, but are limited to <3.6V operation. Using a combined Dickson and Cockcroft-Walton ladder pump [3] (Fig. 11.2.2) allows MIM capacitors to be used while maintaining high power efficiency. A single-stage Dickson pump and a four-stage ladder structure generate an output voltage of 13V. A PMOS body switch is used to avoid diode leakage through the N-well, and a regulation loop is used to stabilize the output voltage. The regulation circuit uses a dual- V_{nn} approach (1.2V and 2.5V) to reduce power by 30%. A V_r-based voltage reference generation circuit [4] provides V_{ref} with sub-nW power consumption. Start-up, erase and program operations need a high VCO frequency (>15MHz) to stabilize V_{out} and hence a high-BW amplifier is required. However, read and standby modes do not require a high-BW amplifier, and therefore amplifier tail current can be lowered in these modes, reducing the total standby power by ~2x. A highresistance diode-chain divider requires low power, but has a long stabilization time. To address this issue, capacitors are placed in parallel with the diode chain to stabilize the loop within 1µs. Figure 11.2.2 shows the 73% peak-power efficiency of the pump loop. The pump itself reduces power by ~4× compared with the baseline design; the whole loop achieves ~7× power reduction using all of the low-power methods mentioned above.

Figure 11.2.3 shows the block diagram of the 1Mb flash macro, which is separated into two banks, each with its own power gating control. When one bank is active, all of the peripherals in the other banks are power gated. Each bank has two arrays with 256×1024 cells. Current SAs, reference current generation, and high-voltage switches are shared by the two arrays. Page-wise erase mode operates on 8kb, whereas only 8b are programmed or read at a time to lower the instantaneous power; this is consistent with common data resolution in sensor systems (e.g., 8b temperature readings). The program power is reduced by ~5×, by using power-gating and a short word length, compared with the baseline design.

At high temperatures, the read current of an erased cell degrades with reduced mobility, whereas that of a programmed cell increases due to its lowered V_{th} . The read current ratio between the two states thus reduces from 8× to 5× as the temperature increases from 25°C to 125°C, complicating SA design at high temperatures. Figure 11.2.3 shows circuit and timing diagrams of a proposed current SA with cross-sampling that doubles the sensing margin. After precharge, (1) S₁ turns on, integrating I_{ref} on C_L (N₀) and I_{cell} on C_R (N₁). (2) S₂ turns on, and S₁ turns off to maintain the gate voltages stored on C_L and C_R and I_{ref} and I_{cell} . When SAEN goes high, I_{cell} - I_{ref} flows out from N₂, while I_{cell} - I_{ref} flows into N₃. The current difference (I_{cell} - I_{ref}) is therefore doubled when the latch-based SA is activated, and mismatch between M₀ and M₁ is cancelled with current sampling instead of voltage sampling [5]. The proposed method halves the current offset compared to a conventional method with same transistor sizing (Fig. 11.2.4 top left). The cross-sampling gate capacitors occupy 16% area of the sense amplifier, which is <0.1% of the whole macro.

The flash macro is fabricated in 90nm embedded ESF3 NOR flash technology. A conventional compiled flash using the same bitcell is also fabricated for baseline comparison. Figure 11.2.4 (top right) shows the measured Shmoo plot of the proposed flash macro achieving 11ns access time at 1.2V and 0.75V read V_{DDmin} . Measured average read V_{DDmin} among 10 dies are 0.739V (Fig. 11.2.4 bottom left). V_{DDmin} across -25°C to 125°C is shown in Fig. 11.2.4 (bottom right). Measured erase and program power at 25°C is 15 μ W and 39 μ W, which represent a 242× and 87× reduction compared to the baseline design. The erase and program energy is 9.4pJ/b and 49pJ/b, which are 30× and 22× lower than the baseline design. At 125°C, the program and erase power is 82µW and 31µW, enabling reliable function of the battery-powered sensor system. At 11ns cycle time, read energy (power) is 2.2pJ/b (1.618mW). Using a read-cycle time suitable for sensor nodes (1us), the design consumes 25µW and shows better power/frequency scaling than the baseline design due to its lower leakage floor. Standby power is 5.4µW even with active charge pump regulation loops: a 4.5× reduction over the baseline design. Figure 11.2.5 compares the measurement results with the baseline design and other work.

The flash macro is incorporated into a high temperature mm-scale sensor system that consists of multiple chip layers: two batteries, flash, PMU, processor, decoupling capacitors, radio, temperature sensor, energy harvester, and a solar cell layer (photo in Fig. 11.2.6). The FSM in the PMU layer turns on the battery switch (SWn=0) once the wake-up timer reaches N_CYCLE specified in the register file (Fig. 11.2.6, top right). Once the battery switch is on, the PMU boots and provides output voltages to the system using switch-capacitor DC-DC converters. The flash layer then automatically programs the processor layer, and the processor executes the sensor program and stores recorded data in flash. It then writes a 'deep sleep' command to register file (SLP_BIT), triggering the FSM to turn off the battery switch (SWn=1) and power gate all of the blocks except for wake-up timer, FSM, and register file. Figure 11.2.6 (top left) shows the simulated power breakdown in deep sleep mode. Measured active system power drawn from the battery is 25µW and 430µW at 25°C and 125°C. Sleep power is 12nW and 190nW at 25°C and 125°C, representing a 3.3× and 63× reduction compared to a conventional system and greatly extending battery lifetime. The stacked system is fully functional at 125°C in stand-alone operation. Figure 11.2.7 shows die and system photos.

Acknowledgements:

This work was supported by TSMC university joint development program and university shuttle program.

References:

[1] M. F. Chang, et al., "A Process Variation Tolerant Embedded Split-Gate Flash Memory Using Pre-Stable Current Sensing Scheme," *IEEE JSSC*, vol. 44, no. 3, pp. 987-994, Mar. 2009.

[2] H. Mitani, et al., "A 90nm Embedded 1T-MONOS Flash Macro for Automotive Applications with 0.07mJ/8kB Rewrite Energy and Endurance Over 100M Cycles Under Tj of 175°C," *ISSCC*, pp. 140-141, Feb. 2016.

[3] T. Ishii, et al., "A 126.6-mm2 AND-Type 512-Mb Flash Memory With 1.8-V Power Supply," *IEEE JSSC*, vol. 36, no. 11, pp. 1707-1712, Nov. 2001.

[4] Q. Dong, et al., "A 114-pW PMOS-Only, Trim-Free Voltage Reference with 0.26% within-Wafer Inaccuracy for nW Systems", *IEEE Symp. VLSI Circuits*, pp. 98-99, June 2016.

[5] M. Jefremow, et al., "Time-Differential Sense Amplifier for Sub-80mV Bitline Voltage Embedded STT-MRAM in 40nm CMOS", *ISSCC*, pp. 216-218, Feb. 2013.



Figure 11.2.1: An SRAM-based sensor system drains a mm-scale battery in minutes at 125°C. A flash-based sensor system can enter deep sleep, reducing power consumption and extending battery life.







Figure 11.2.5: Measurement results comparison with baseline design and other work.



Figure 11.2.2: Combined Dickson and ladder pump topology, and self-adjusted regulation loop, which achieves a 73% peak efficiency.



Figure 11.2.4: The proposed SA halves offset. Measured access time is 11ns at 1.2V. Average read V_{DDmin} among 10 dies is 0.739V. Read V_{DDmin} and write power across temperature shown at bottom right.



Figure 11.2.6: Diagram of PMU layer with wake-up timer and measured stacked system power. Deep sleep power at 25° C and 125° C are 12nW and 190nW, representing a $3.3 \times$ and $63 \times$ power reduction.

