



# Power/Capacity Scaling: Energy Savings With Simple Fault-Tolerant Caches

**Mark Gottscho**<sup>1</sup>

Abbas BanaiyanMofrad<sup>2</sup>

Nikil Dutt<sup>2</sup>

Alex Nicolau<sup>2</sup>

Puneet Gupta<sup>1</sup>

<sup>1</sup> NanoCAD Lab  
UCLA Electrical Engineering

[nanocad.ee.ucla.edu](http://nanocad.ee.ucla.edu)

<sup>2</sup> UC Irvine Computer Science

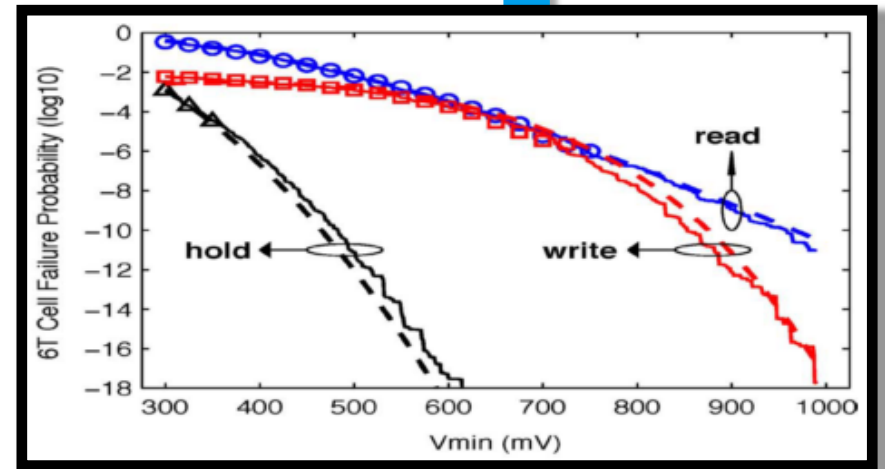
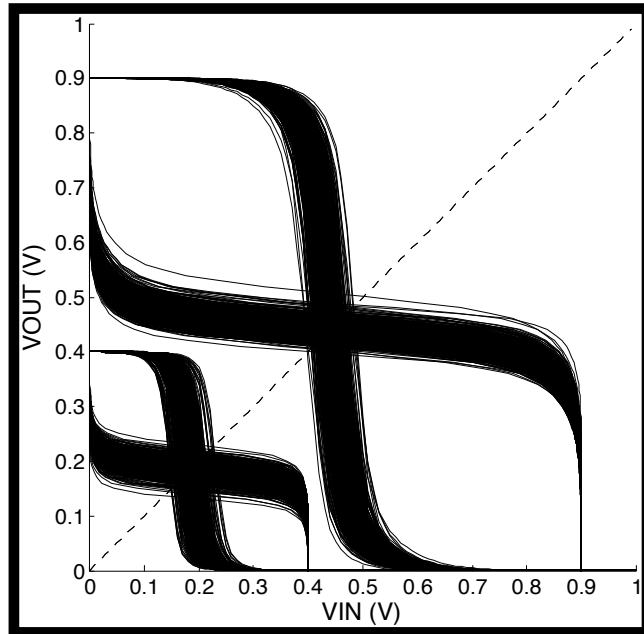
[ics.uci.edu](http://ics.uci.edu)

# Motivation: Increasing Process Variability Limits SRAM Voltage Scaling

Process variability  $\uparrow$



$V_T$  variations  $\uparrow$   
SRAM  $\sigma_{SNM}$   $\uparrow$



J. Wang and B. Calhoun  
TVLSI 2011



VDD  $\downarrow$



SRAM BER  $\uparrow$   
exponentially

# Question

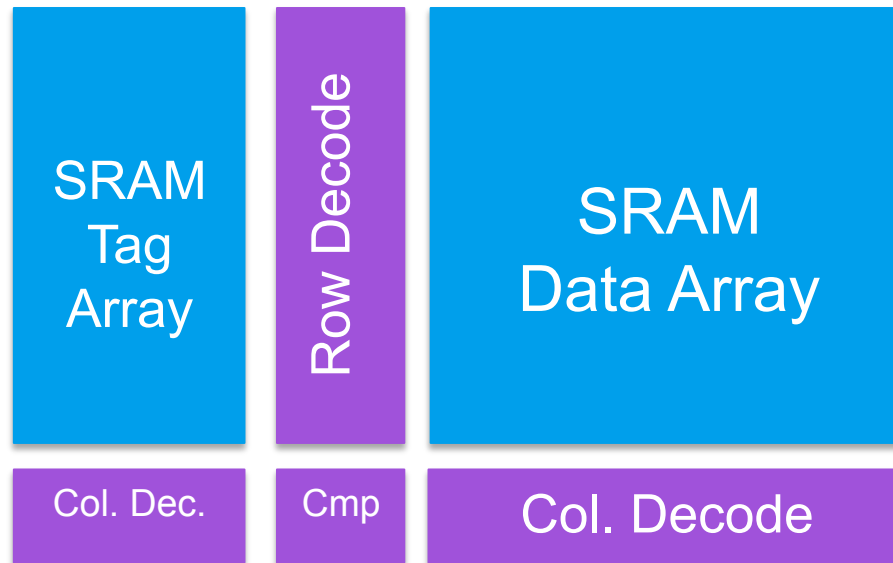
---

**How to optimize SRAM  
for the “best” *system-level* tradeoffs in  
energy, reliability, performance, & area?**

# Using Fault Tolerance to Achieve Lower min-VDD

Many fault-tolerant, voltage-scalable (FTVS) approaches lower min-VDD using sophisticated fault tolerance methods

Baseline Cache @ Nominal VDD – No Fault Tolerance



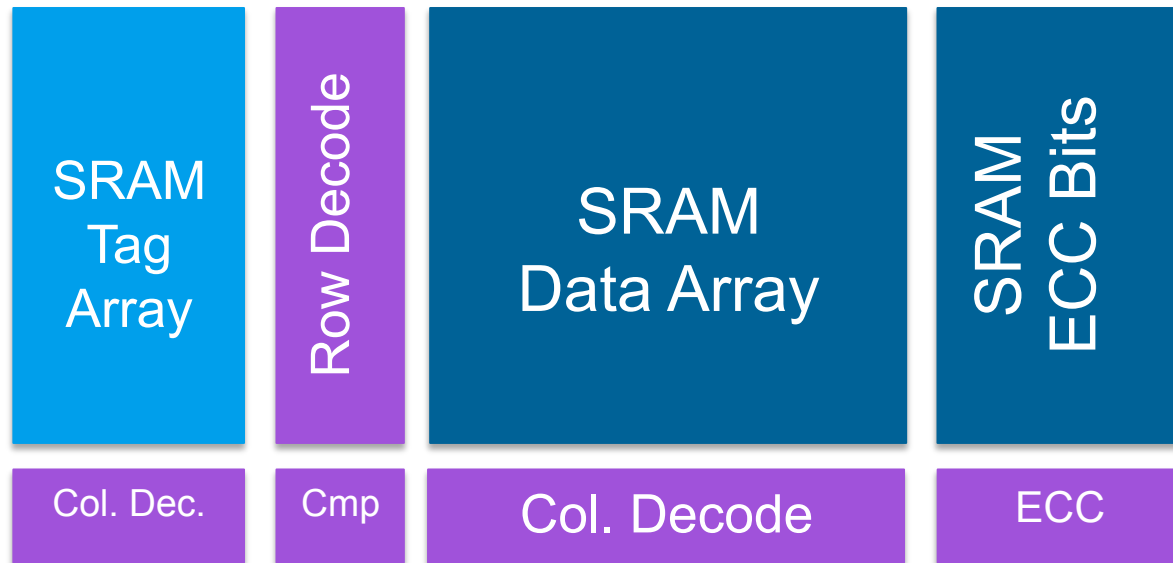
**PURPLE** = periphery  
@ full VDD

**BLUE** = SRAM cells  
(bright is higher  
VDD)

# Using Fault Tolerance to Achieve Lower min-VDD

Many fault-tolerant, voltage-scalable (FTVS) approaches lower min-VDD using sophisticated fault tolerance methods

ECC Cache, Data Array @ 0.7 VDD



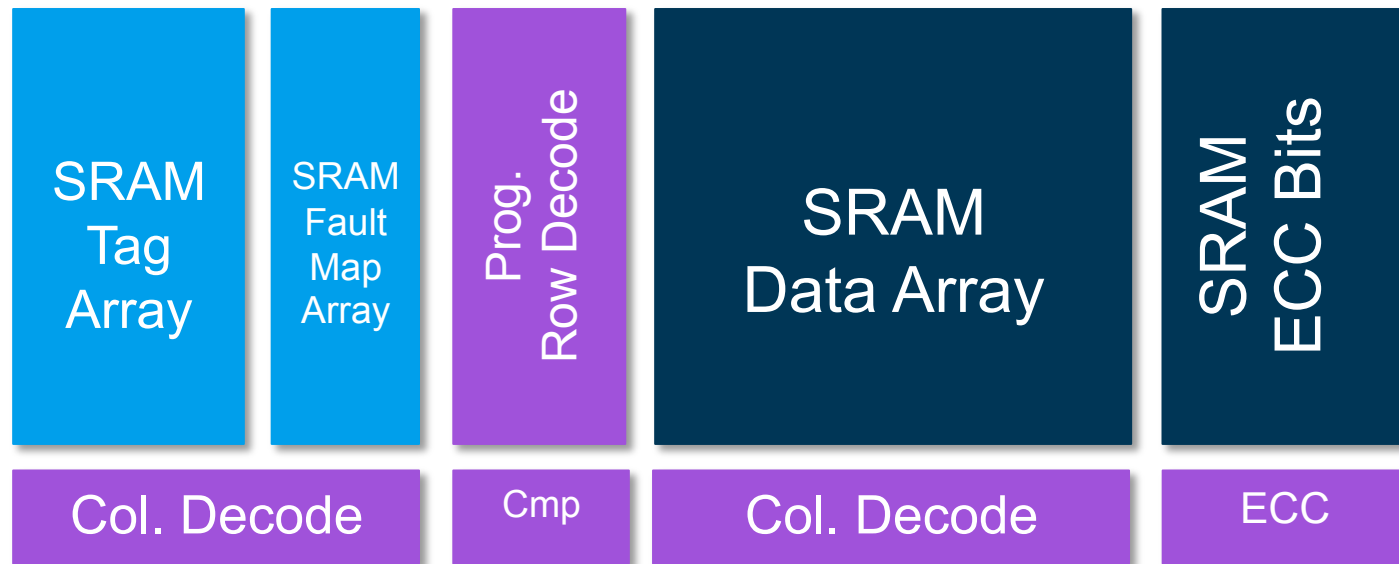
**PURPLE** = periphery  
@ full VDD

**BLUE** = SRAM cells  
(bright is higher  
VDD)

# Using Fault Tolerance to Achieve Lower min-VDD

Many fault-tolerant, voltage-scalable (FTVS) approaches lower min-VDD using sophisticated fault tolerance methods

ECC + Faulty Set Remapping Cache, Data Array @ 0.5 VDD



**PURPLE** = periphery @ full VDD

**BLUE** = SRAM cells (bright is higher VDD)

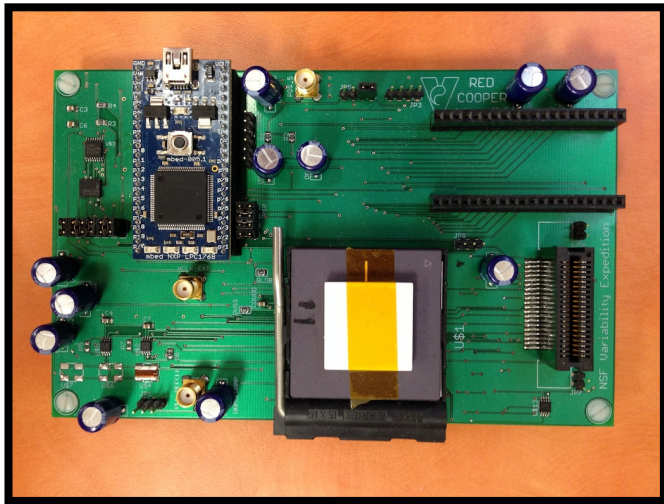
# This Work

---

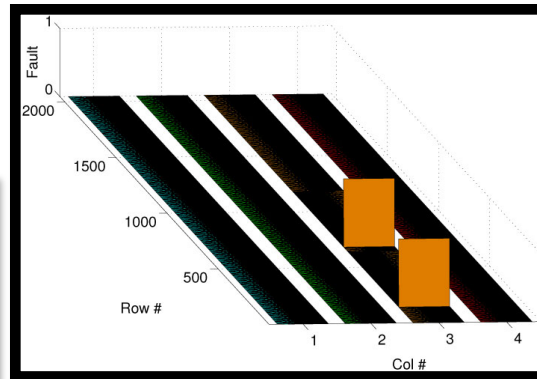
Save energy via  
*simple & low-overhead*  
fault-tolerant, voltage-scalable (FTVS)  
SRAM cache architecture

# SRAM “Fault Inclusion Property”

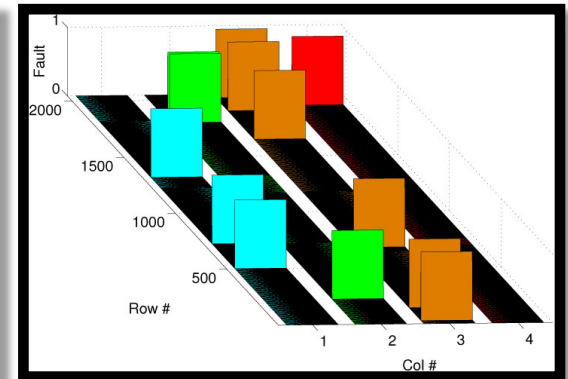
NSF Variability Expedition  
“Red Cooper” test chips<sup>1</sup>  
based on ARM Cortex M3



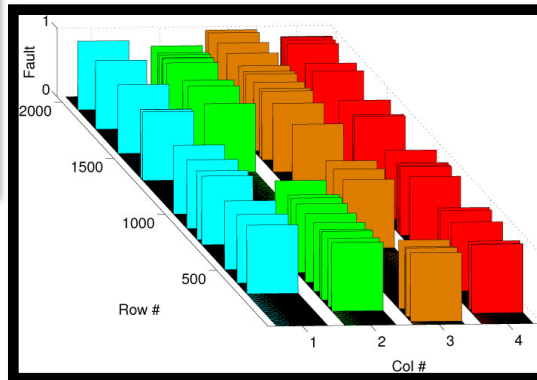
<sup>1</sup>L. Lai ASP-DAC'14



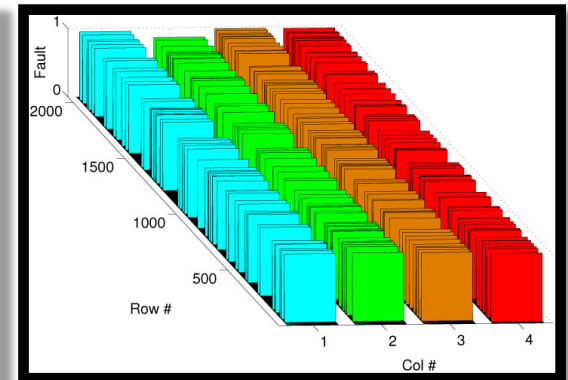
550 mV



525 mV



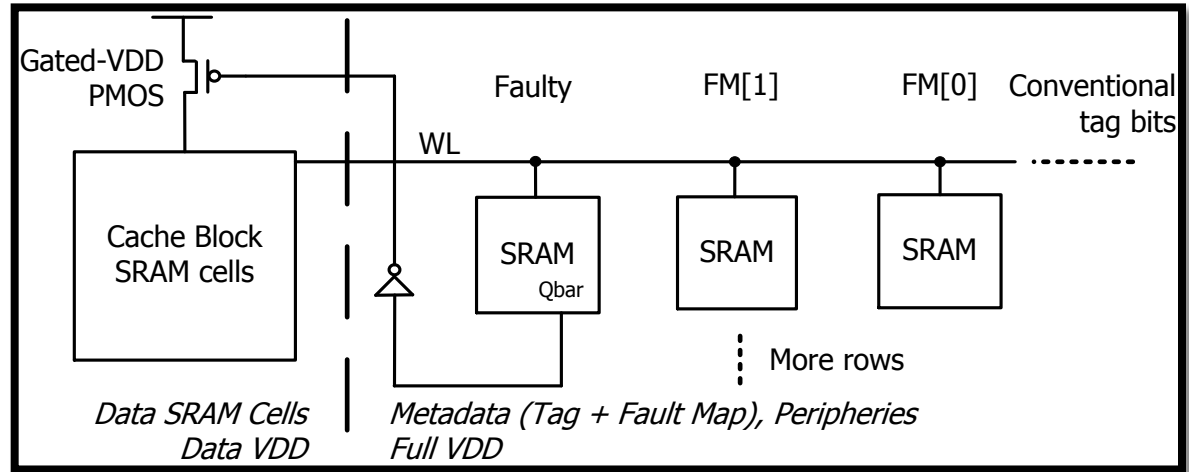
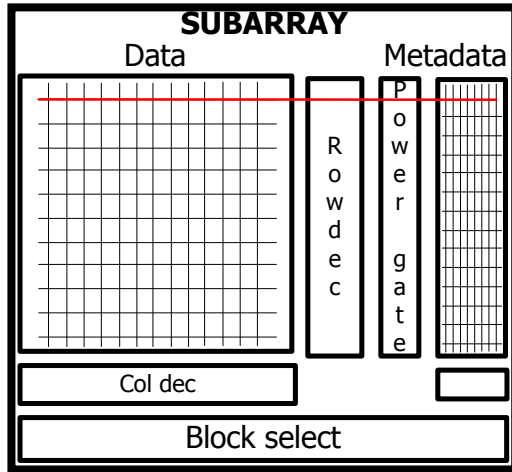
500 mV



475 mV



# Architectural Mechanism



- No redundancy – just sacrifice faulty blocks as VDD scales
  - # good blocks fall off a “cliff” anyway
    - Redundancy can only do so much
  - Negligible area overhead

# Power/Capacity Scaling

---

- To adjust data array VDD:
  - Temporarily stall accesses
  - Cache controller finds the blocks that will become faulty at next VDD using *FM* bits
    - Flush those blocks that are also *Valid & Dirty*
    - Then set *Faulty* bits, power gating them
  - Adjust VDD, wait for voltage to settle
  - Resume operations

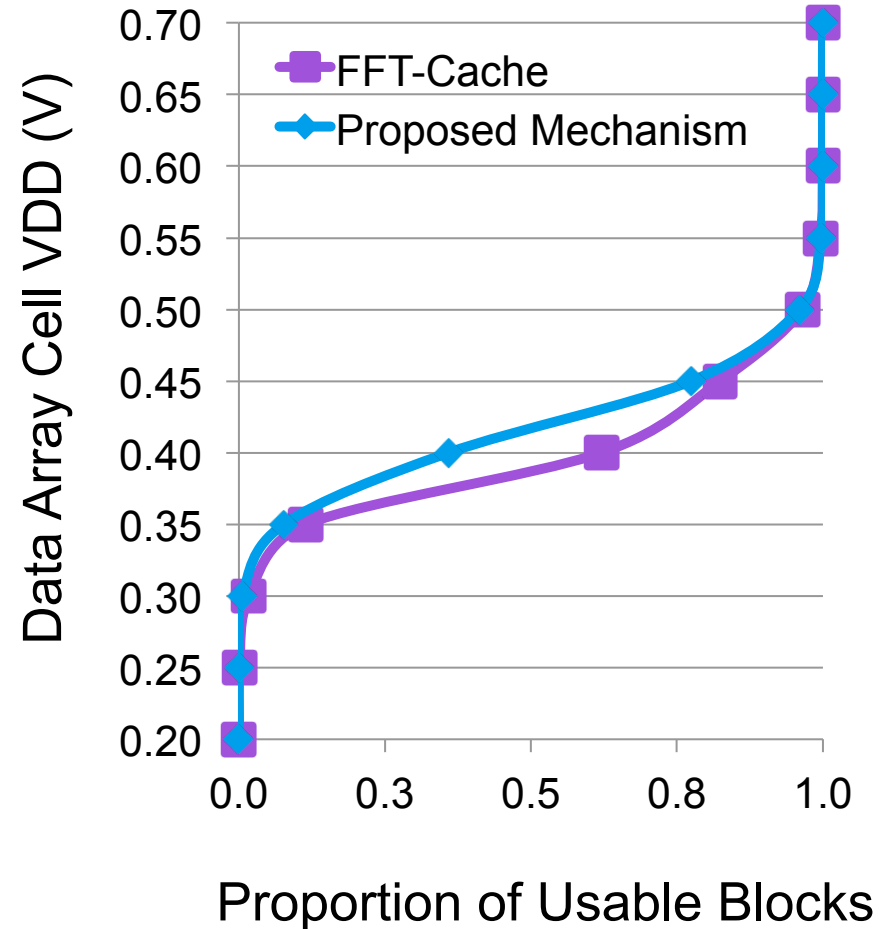
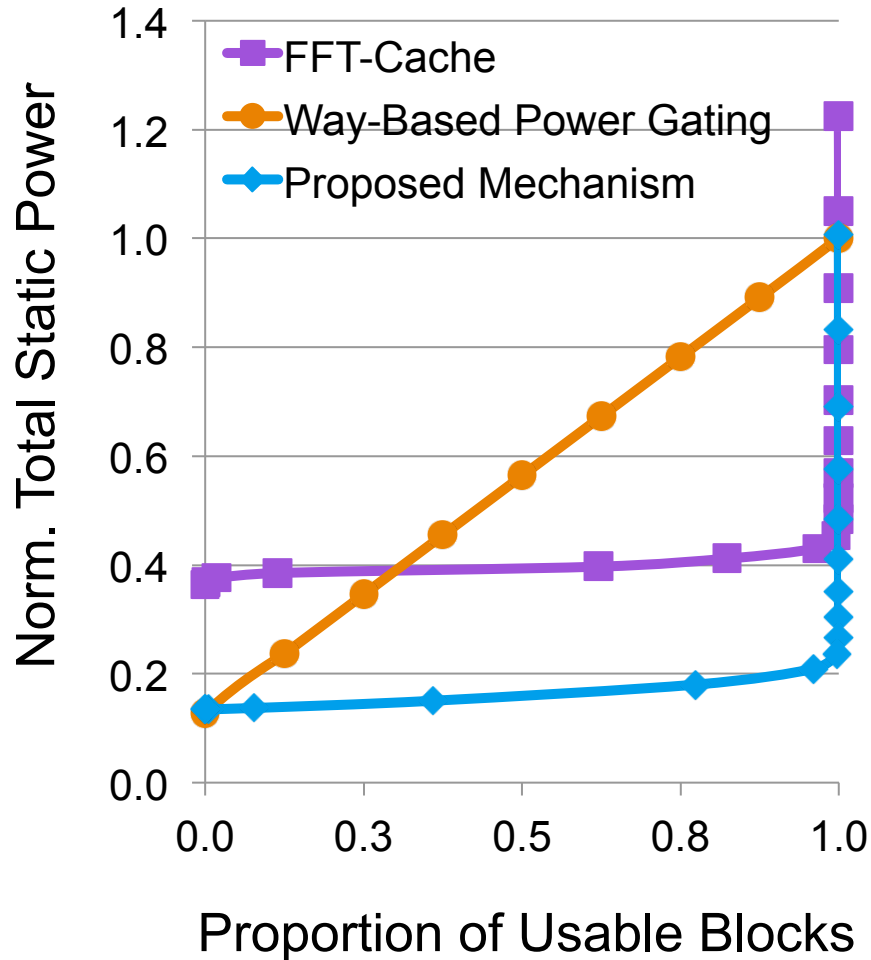
# Static & Dynamic Power/Capacity Scaling Policies

---

- Static (SPCS) Policy: Choose single optimal VDD at design, test, or boot time
- Dynamic (DPSCS) Policy: Decide on optimal VDD periodically during runtime



# Sample Analytical Results: L2 Cache Config. 1



# Simulation Results Summary

Against baseline 6T SRAM cache @ 1V:

SPCS: Up to 55% total cache energy savings

DPCS: Up to 69% total cache energy savings

DPCS: < 4% performance overhead

< 5% area overhead

(g) Normalized Total Cache Energy, Config. A

(h) Normalized Total Cache Energy, Config. B

Figure 4: Simulation Results for SPCS and DPCS

# Summary

---

Power vs. effective capacity

Fault Inclusion Property

Keep it simple

# Directions for Future Work

---

Optimality study of DPCS policies

DPCS for shared memory MPSoCs

DPCS-aware software strategies



# Q&A

**This work was supported by the NSF Variability Expedition  
Grant Numbers CCF-1029783 and CCF-1029030**

**[www.variability.org](http://www.variability.org)**