

# Aging-Aware Compiler-Directed VLIW Assignment for GPGPU Architectures

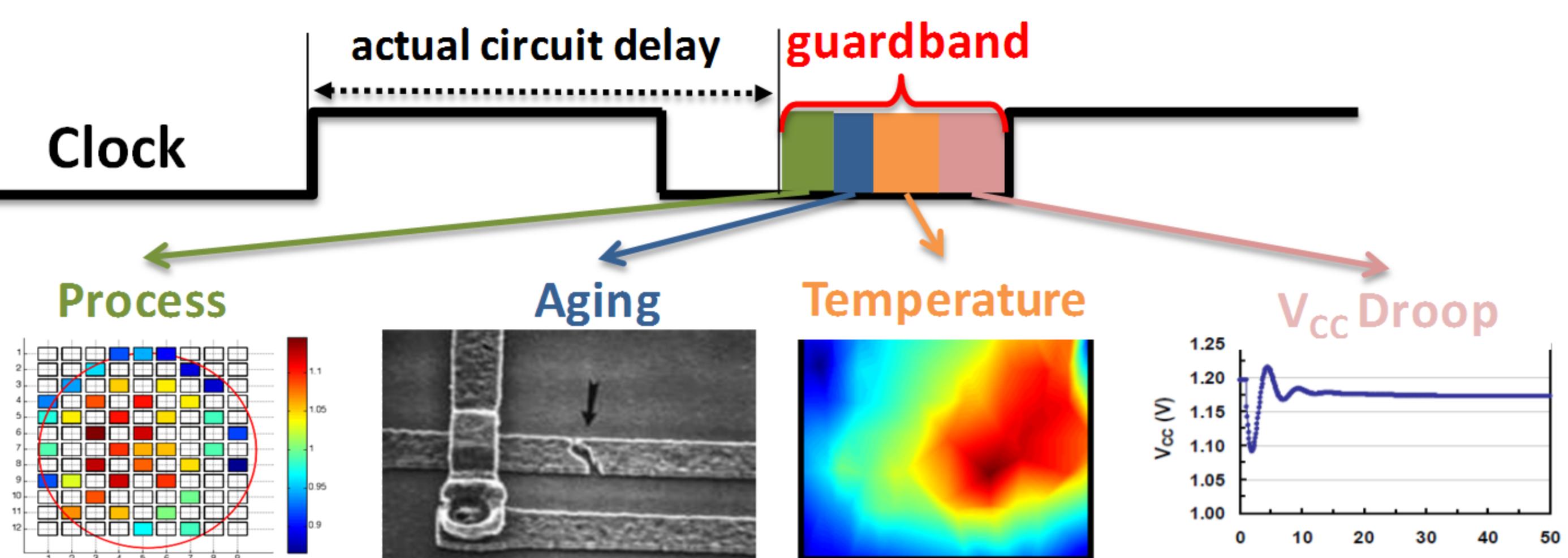


Abbas Rahimi<sup>†</sup>, Luca Benini<sup>‡</sup> and Rajesh K. Gupta<sup>†</sup>

<sup>†</sup>UC San Diego and <sup>‡</sup>Università di Bologna

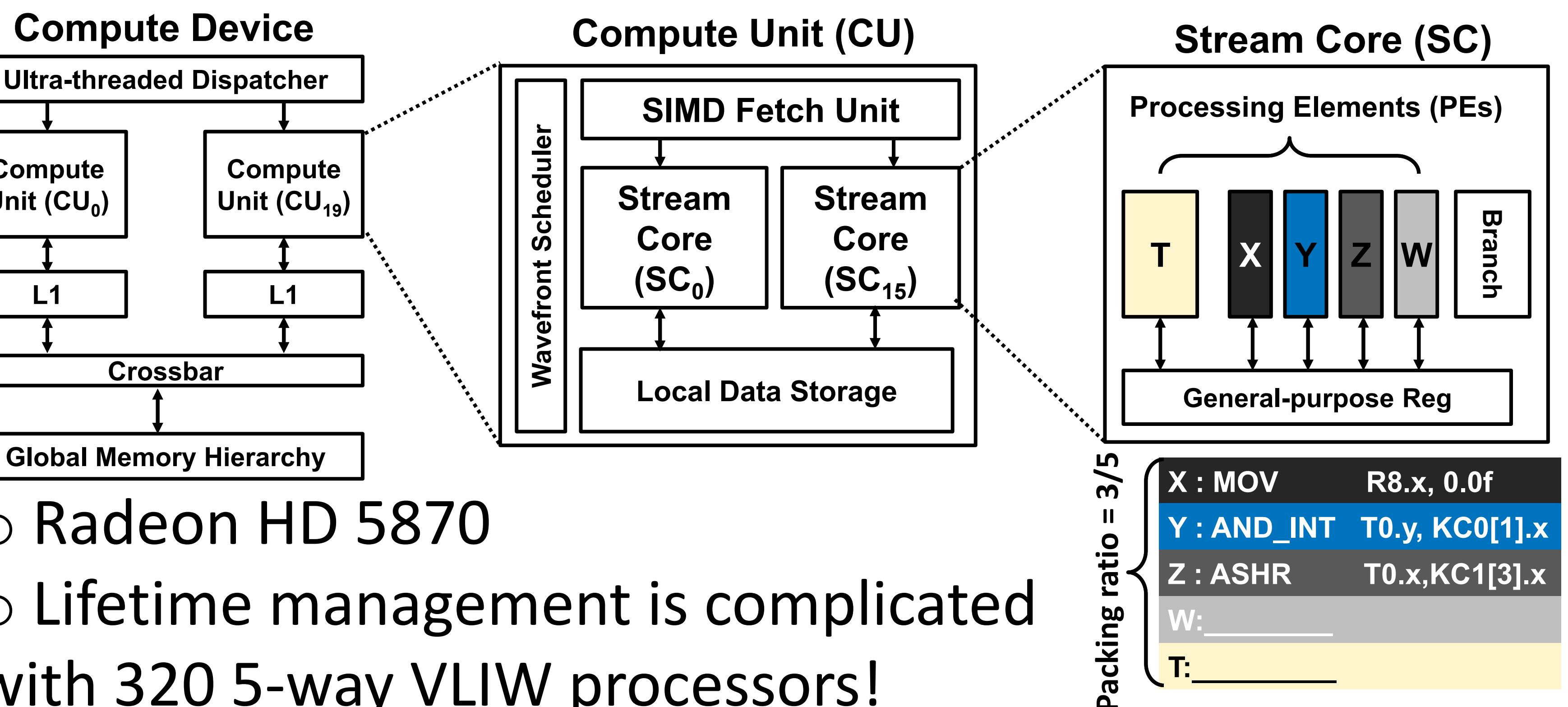


## Variability is about Scale and Cost



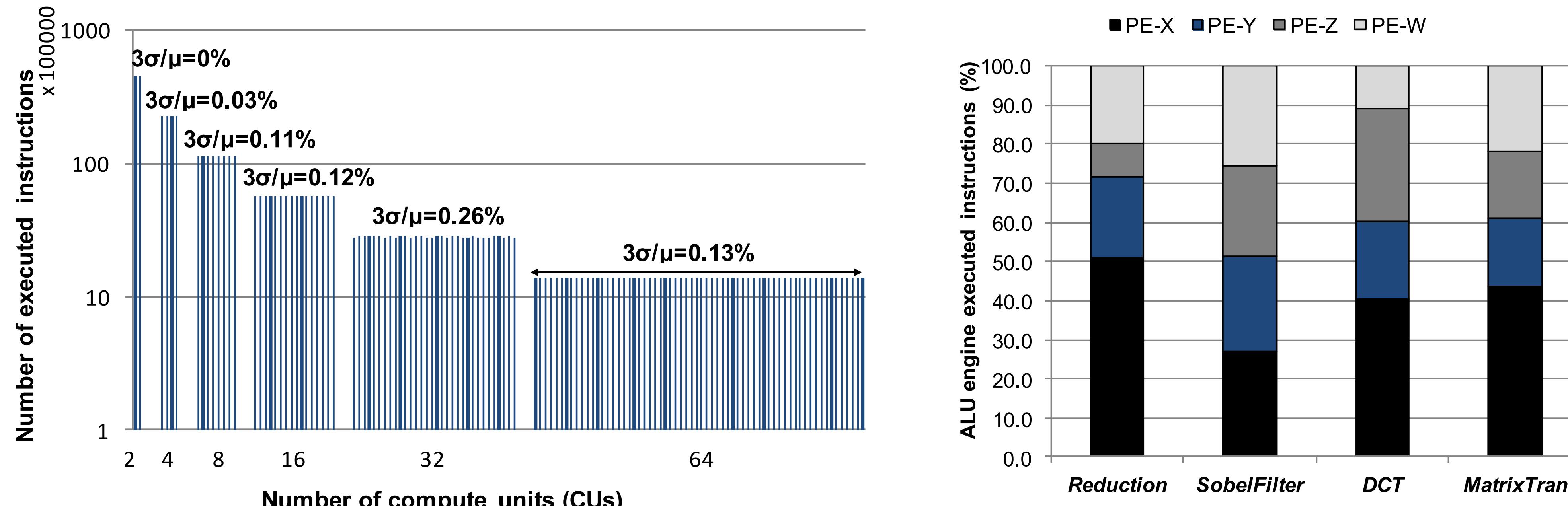
- NBTI-induced performance degradation
- $\Delta V_{TH} = f(\text{Process, Temp, Voltage, Stress})$
- Knobs to tweak: {Power-gating, Duty cycling}

## AMD Evergreen GPGPU Architecture

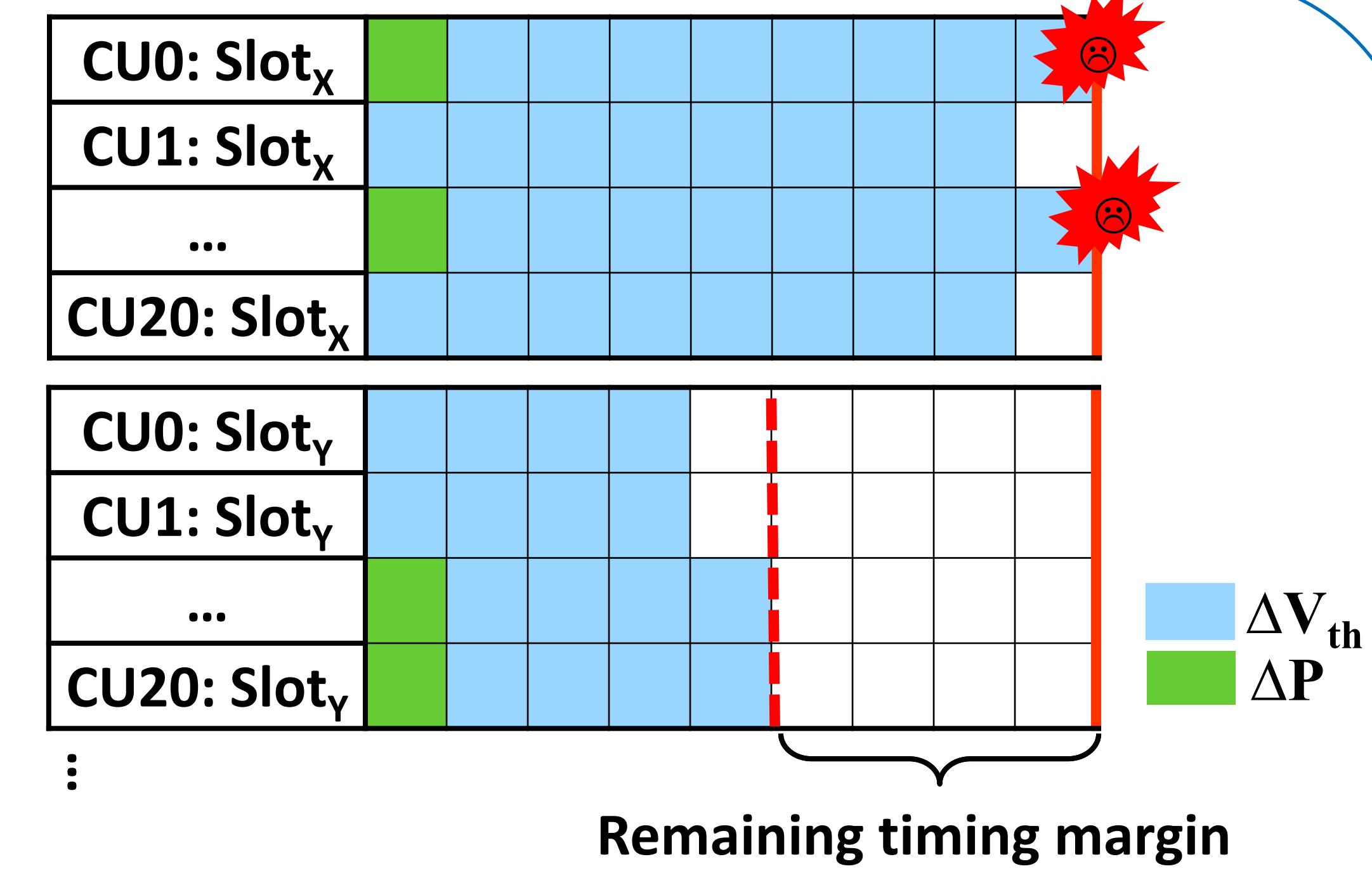


- Radeon HD 5870
- Lifetime management is complicated with 320 5-way VLIW processors!

## GPGPU Workload Variation



- Inter-Compute Units: Uniform 😊
- Inter-Processing Elements: Uneven! 😕
- Inter-Stream Cores: SIMD fashion 😊
- PE<sub>X</sub> executes ~40% of instructions 😕
- Stress consumes timing margin
- The most aged PE limits lifetime



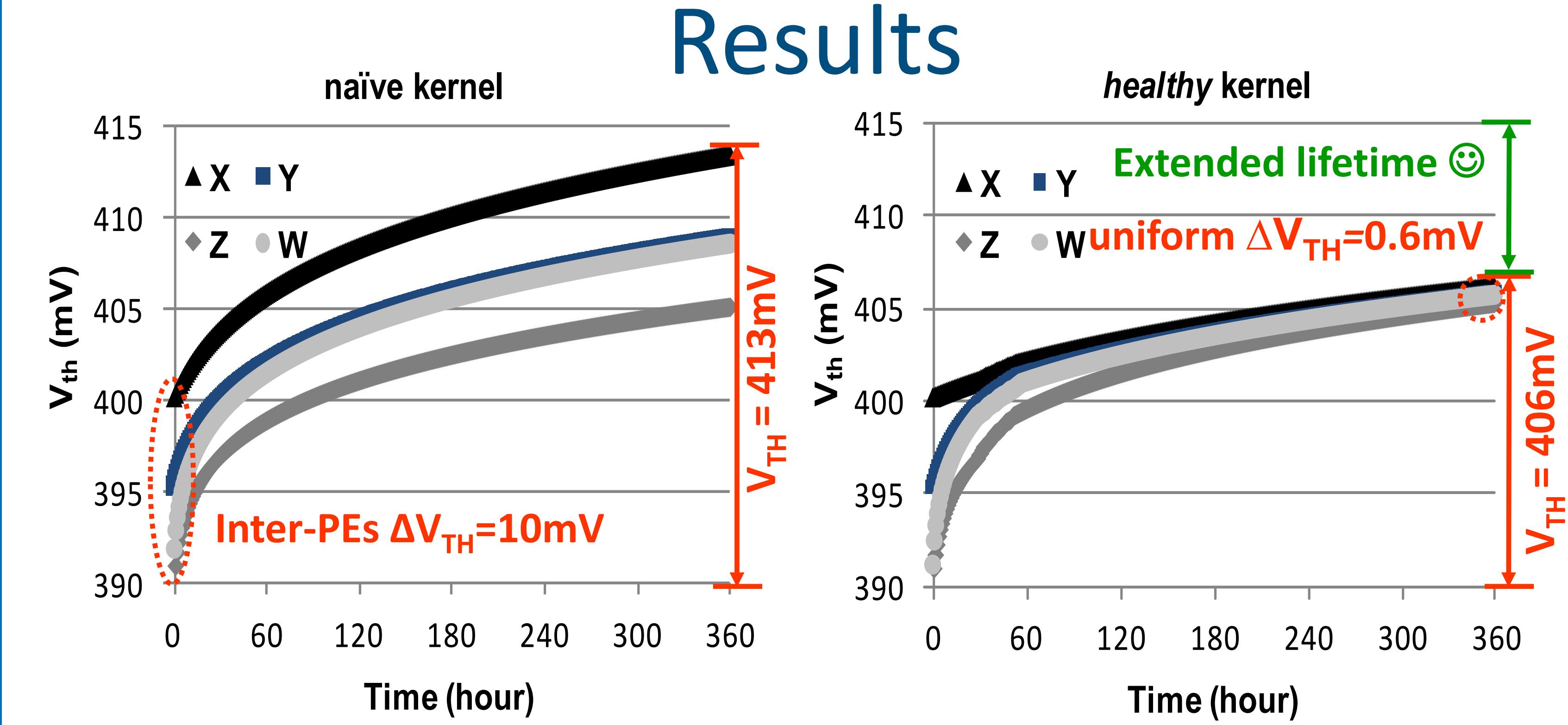
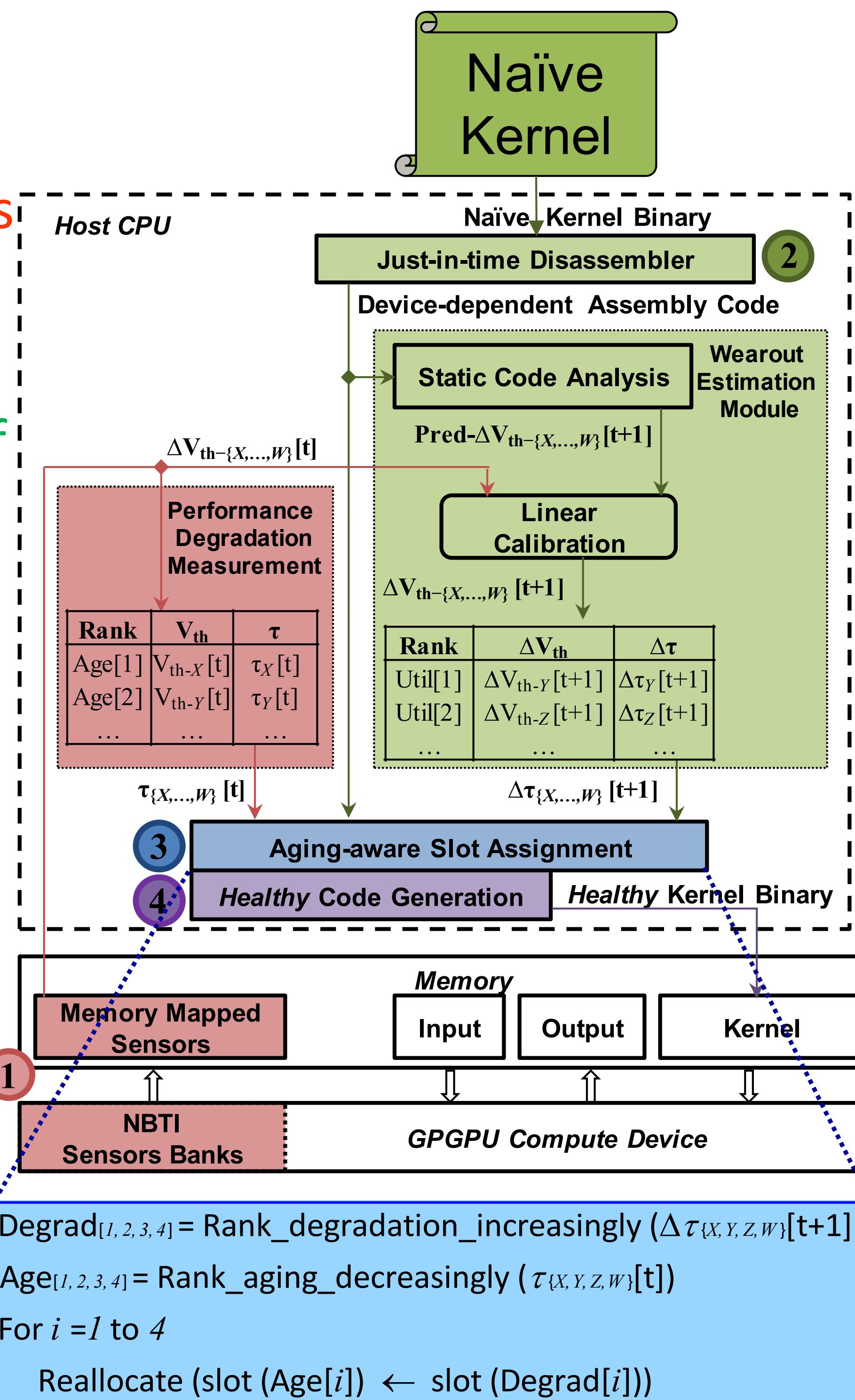
## Aging-Aware Compilation

### 1. Observability: Reading NBTI sensors measurements

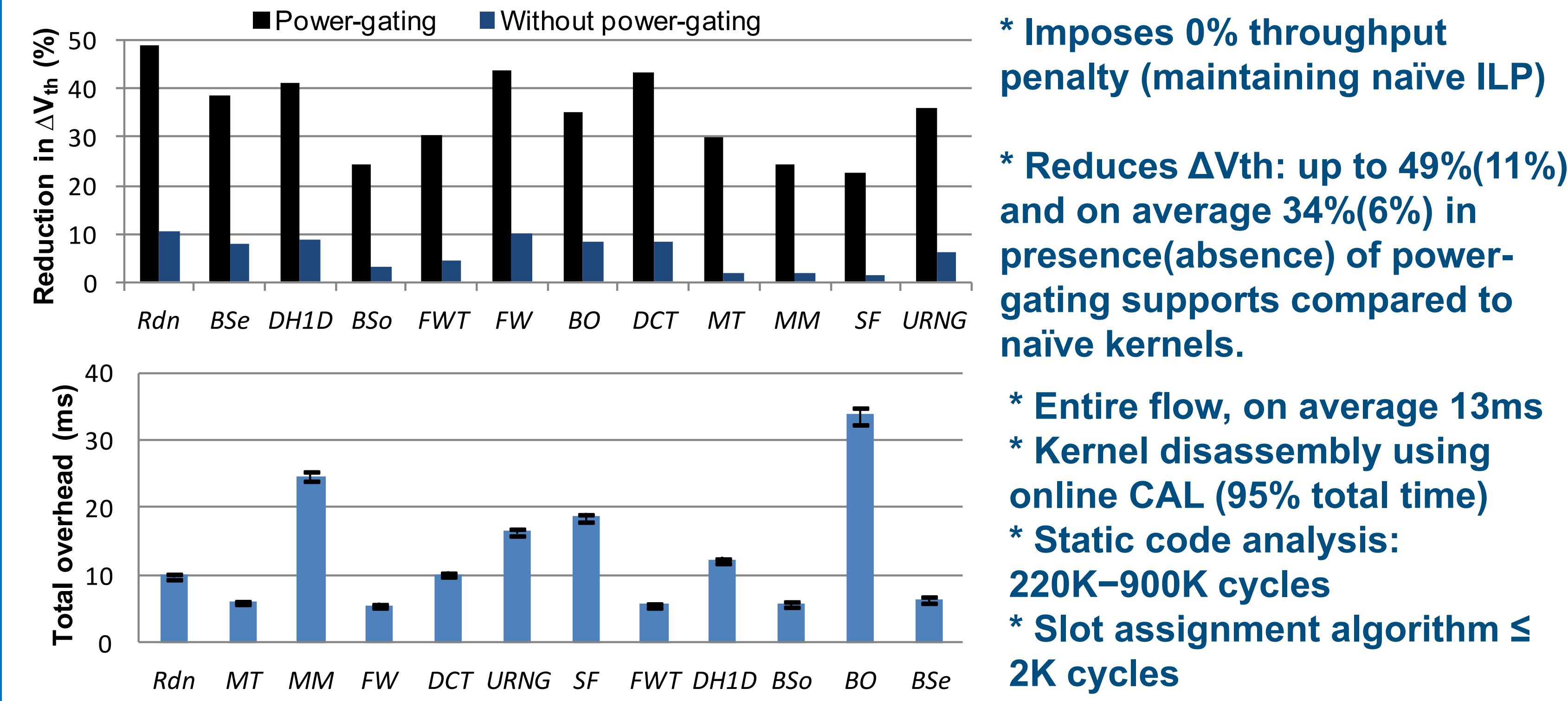
2. Prediction: Static code analysis technique estimates the percentage of instructions that will carry out on every PE (a linear calibration module later fits the predicted  $\Delta V_{TH}$  shift to the observed  $\Delta V_{TH}$  shift).

3. Controllability: Uniform slot assignment assigns fewer/more instructions to higher/lower stressed slots.

4. Periodic healthy kernel generation 😊



Adapting kernels periodically leads to a uniform  $\Delta V_{TH}$  among PEs (without power-gating)



Work in progress: Memory subsystems, reducing  $\Delta V_{TH}$  by up to 43% for register files.

NSF Expedition in Computing, Variability-Aware Software for Efficient Computing with Nanoscale Devices <http://variability.org>

UC San Diego

UCLA

MICHIGAN

STANFORD  
UNIVERSITY

UCIRVINE

ILLINOIS  
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN