

# ***Half-DRAM: a High-bandwidth and Low-power DRAM Architecture from the Rethinking of Fine- grained Activation***

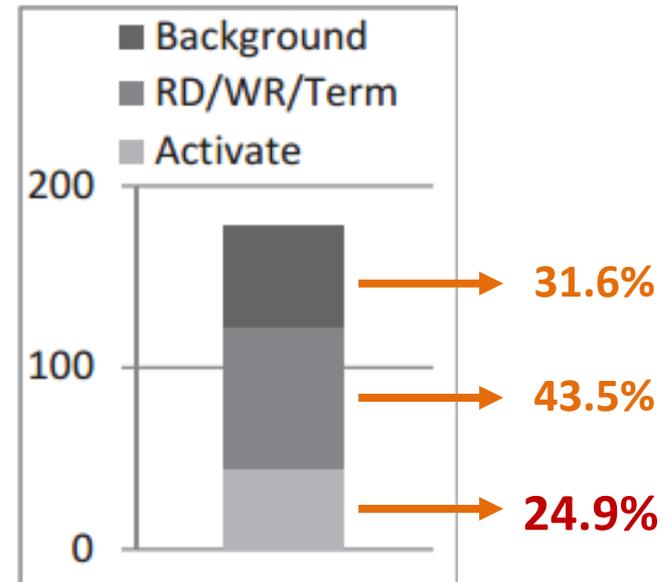
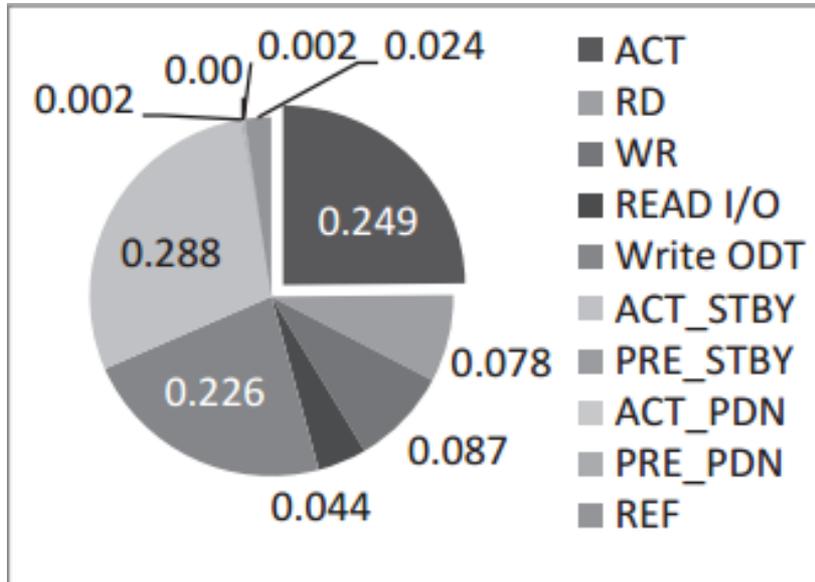
Tao Zhang, Ke Chen, Cong Xu, Guangyu Sun,

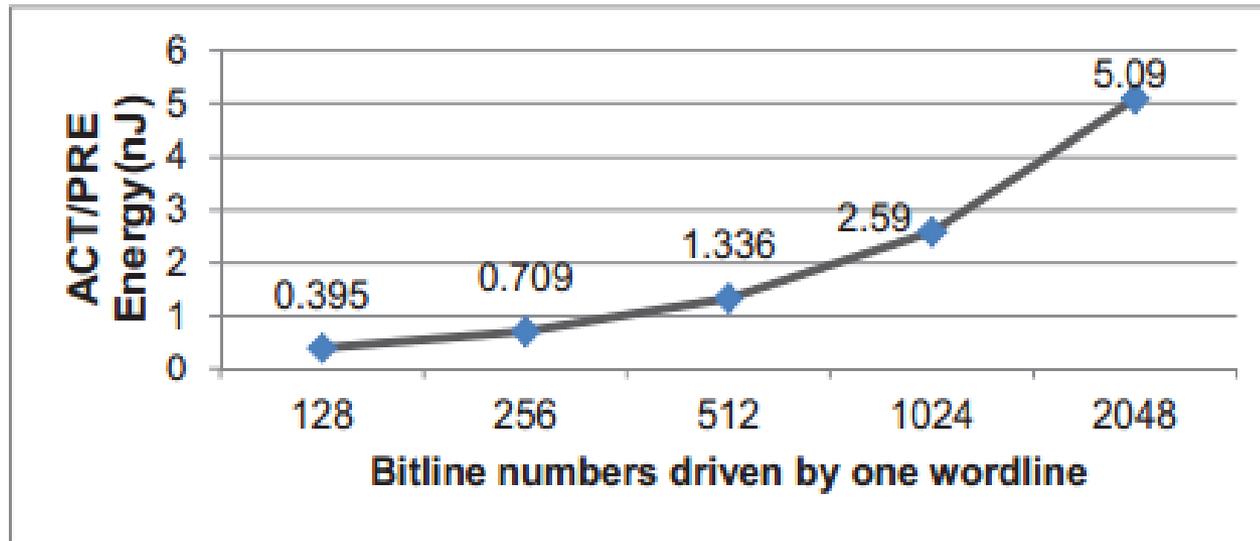
Tao Wang, Yuan Xie

Pennsylvania State University

**ISCA - 2014**

- DRAM can consume **more than 25%** of the total power in a datacenter
  - => improving the power efficiency of DRAM one of the major challenges in the memory architecture design
- Activation and Precharge power can be around **25%** of the total DRAM power



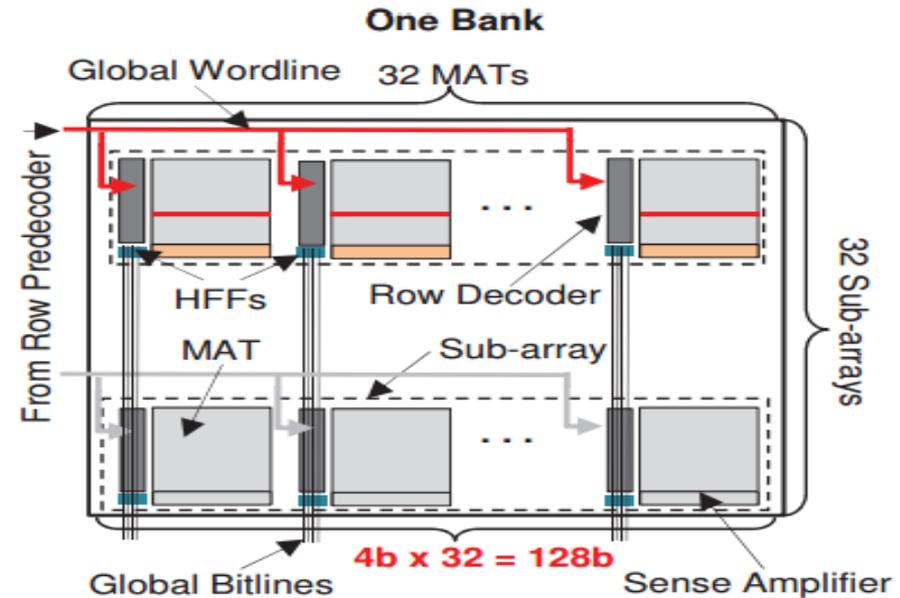
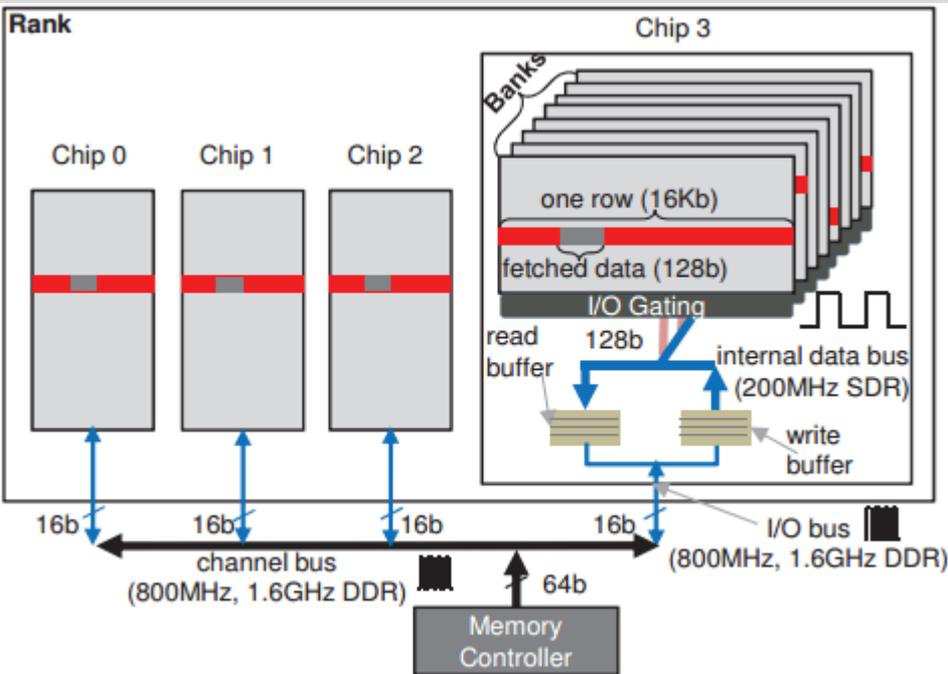


The activation power is *proportional* to the **number of bitlines** being activated during a memory access

- For future memory chips with larger capacity and **more bitlines**, the power efficiency problem of row activation will become **even worse**

*The goal of this work is to minimize the Activation Power*

# DRAM Row Access and Activation Power

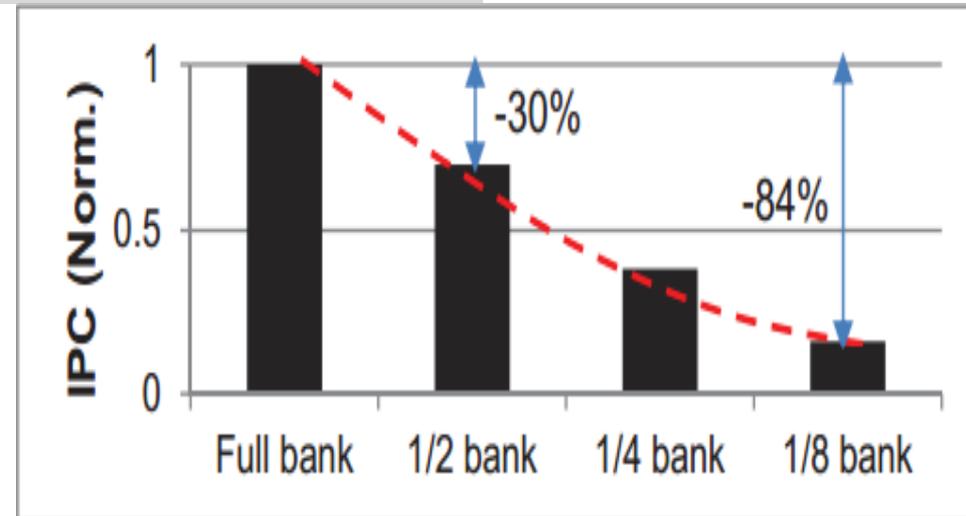
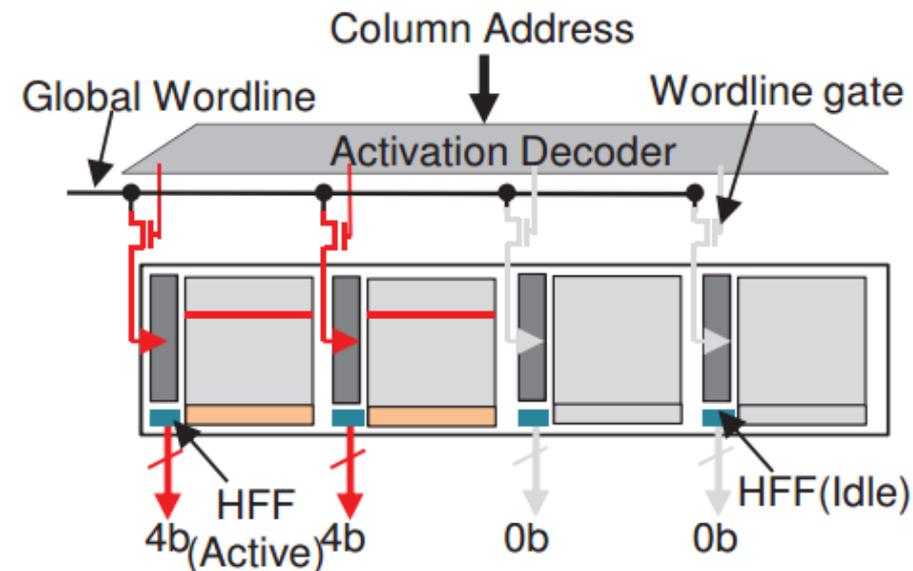


**Row Overfetching:** Entire row to be activated (**red block**) even though only a small portion of data are fetched at a time (**gray block**)

$$I_{ACT} = IDD0 - \frac{IDD3N \times t_{RAS} + IDD2N \times (t_{RC} - t_{RAS})}{t_{RC}}$$

*Reducing activated row size (i.e., number of bitlines) can significantly reduce the activation power*

# Prior Fine-Granule Row Access and the BW dilemma



*HFF: Helper Flip-Flop > latches the selected data and relay them on the internal data bus*

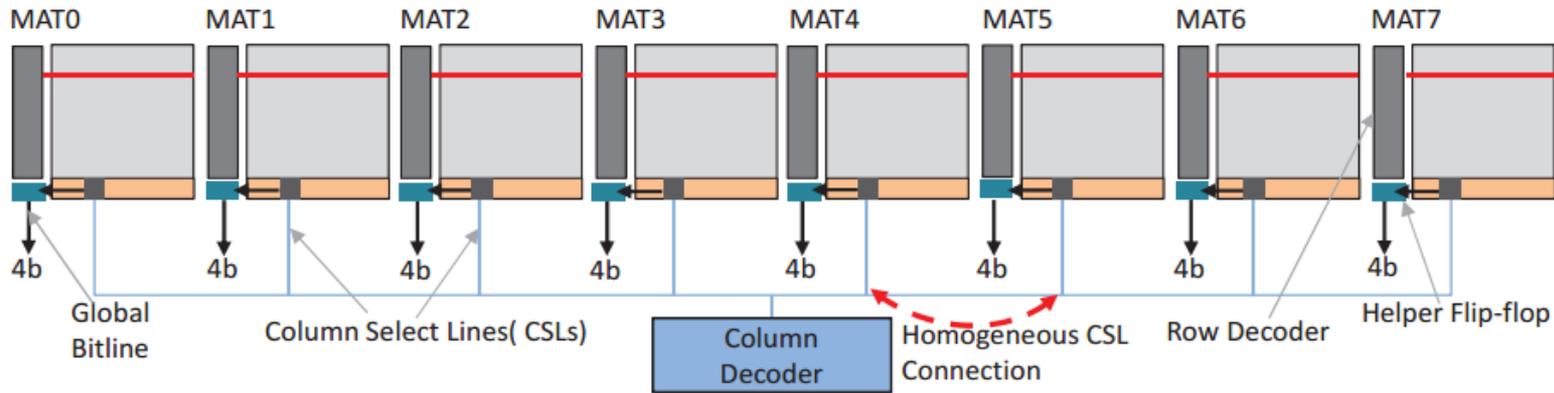
- An activation decoder is introduced to control the number of active MATs

**However...**

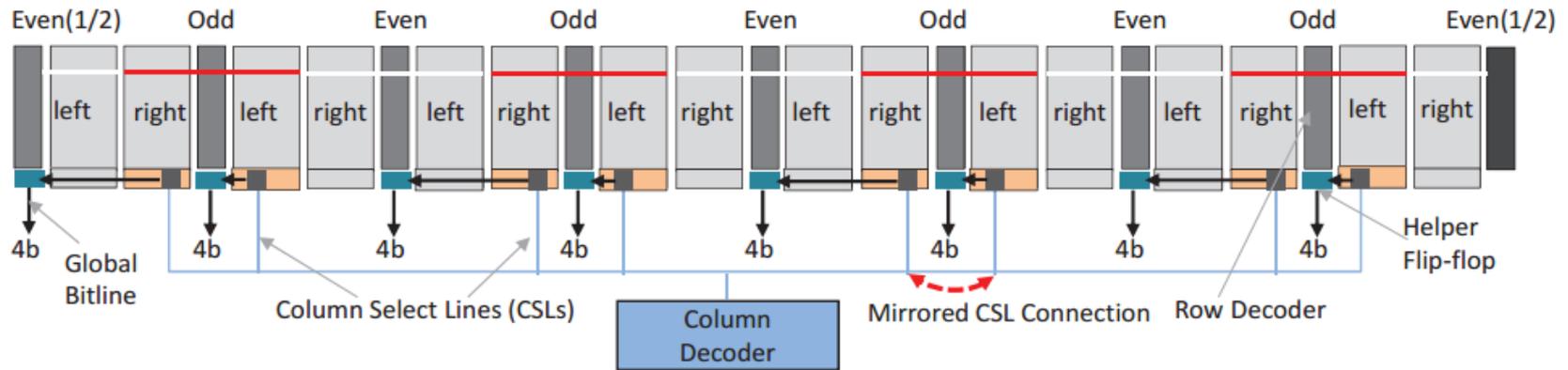
$$B_{DataBandwidth} = W_{DataWidth} \times F_{DataFrequency}$$

**Fine-grained activation techniques can result in significant memory performance degradation**

# Half-DRAM



(a) 1RD-1HFF: traditional structure



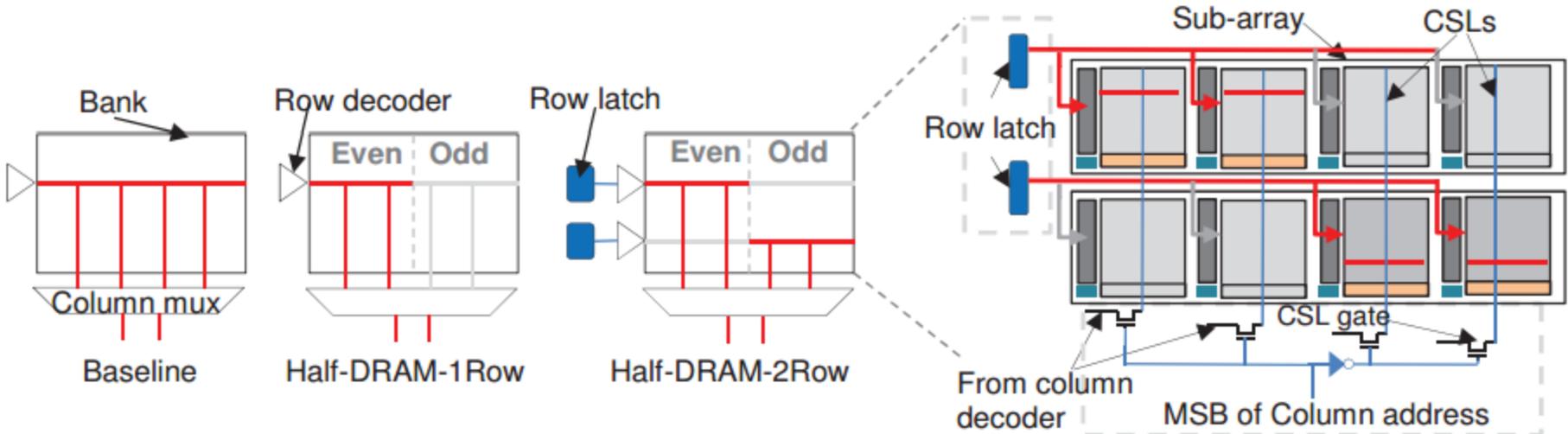
(b) 1RD-2HFF: Half-DRAM structure (proposed)

**Traditional DRAM:** One-to-one relationship exists between row decoder and HFF.

**1RD-1HFF:** One HFF group is dedicated to one row decoder and vice versa

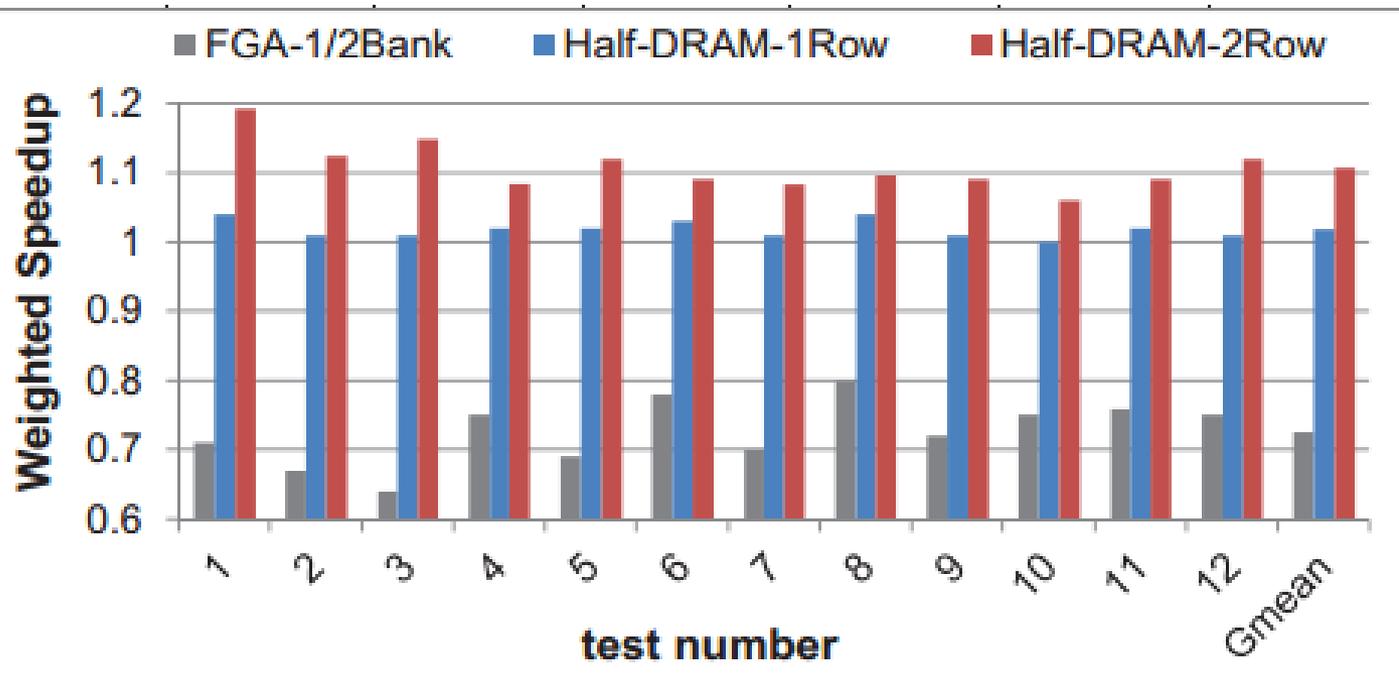
**Half-DRAM:** MAT is split into “left” and “right” block - driven by different row address decoders. The sub-array is further divided into Odd and Even groups

# Half-DRAM



- Here, once a row decoder selects a wordline, both MATs are activated but each with half a row.
- Even if every other row decoder is disabled, all HFFs can still be active
- Though Half-DRAM mitigates the bandwidth reduction problem, it only has half of the row buffer size, which may degrade row buffer hit rate.
- Half-DRAM can be extended to the traditional full bank activation to provide full row buffer size .... *Half-DRAM-2Row*

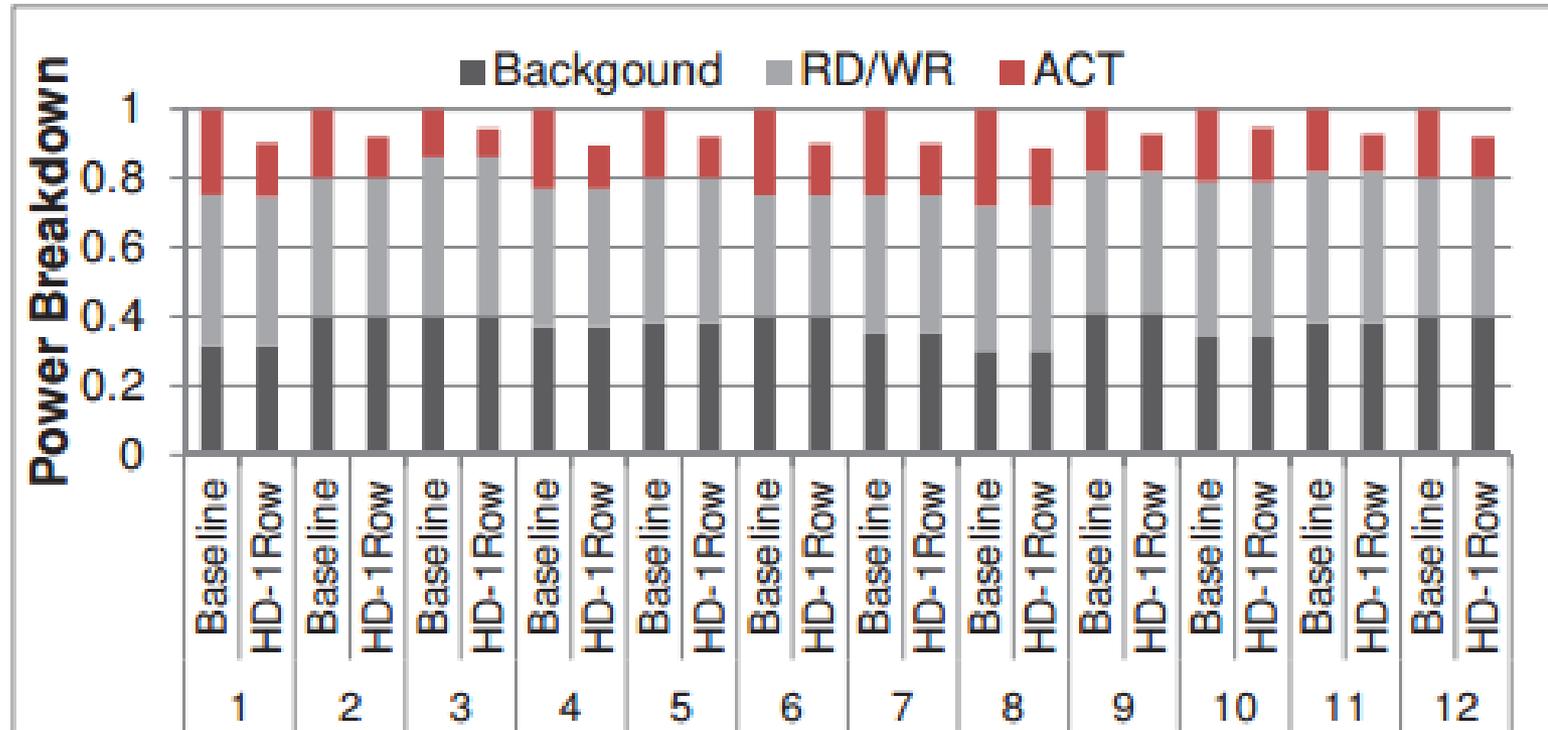
# Performance Evaluation



Avg. performance improvement in Half-DRAM-2Row is **10.7%**

## *test#* Benchmarks (SPEC2006+STRAM)

<sup>1</sup>STREAM×4, <sup>2</sup>bwaves×4, <sup>3</sup>gobmk×4, <sup>4</sup>leslie3d×4,  
<sup>5</sup>libquantum×4, <sup>6</sup>lbm×4, <sup>7</sup>mcf×4, <sup>8</sup>milc×4,  
<sup>9</sup>STREAM-gobmk-lbm-libquantum, <sup>10</sup>bwaves-leslie3d-mcf-milc,  
<sup>11</sup>lbm-libquantum-bwaves-leslie3d, <sup>12</sup>STREAM-gobmk-mcf-milc



On average, Half-DRAM-1Row can achieve **8.4%** improvement on power efficiency over the baseline

**What about Half-DRAM-2Row?**



**Thank you**