✏️

# CACTI-3DD

| ⏱ Created | @October 6, 2024 8:56 AM |
| --- | --- |
| ☰ Class | |

## Get the runnable CACTI-3DD from Sectored DRAM

> https://github.com/CMU-SAFARI/Sectored-DRAM

- Run the areapower.py to generate the runnable makeFile for cacti
- Modify the code containing the sectored DRAM modification, which are the sense amplifier modifications.

> https://github.com/HewlettPackard/cacti

- Replace the mat.cc with the one from CACTI-3DD

## Bug for modifications

- Modification for tsv_length must be modified to get correct timing constraints when inputing the stacking parameter.
- https://github.com/HewlettPackard/cacti/issues/2



- Otherwise the CACTI-3DD cannot be run correctly

## 3D-DRAM Samsung.cfg modification

- https://github.com/HewlettPackard/cacti/issues/6
- Such that you can verify and get similiar results for the parameter of 3D-Samsung model timing constraints.

```
53  Memory Parameters:
54      Total memory size (Gb): 8
55      Stacked die count: 4
56      TSV projection: industrial conservative
57      Number of banks: 16
58      Technology size (nm): 50
59      Page size (bits): 8192
60      Burst depth: 8
61      Chip IO width: 4
62      Best Ndwl: 16
63      Best Ndbl: 32
64      # rows in subarray: 512
65      # columns in subarray: 512
66
67  Results:
68  Timing Components:
69      t_RCD (Row to column command delay): 6.93514 ns
70      t_RAS (Row access strobe latency): 16.6885 ns
71      t_RC (Row cycle): 25.6223 ns
72      t_CAS (Column access strobe latency): 14.1096 ns
73      t_RP (Row precharge latency): 9.67556 ns
74      t_RRD (Row activation to row activation delay): 3.17319 ns
75  Power Components:
76      Activation energy: 1.01003 nJ
77      Read energy: 0.967235 nJ
78      Write energy: 0.967252 nJ
79      Precharge energy: 0.911832 nJ
80  Area Components:
81      DRAM core area: 65.0447 mm2
82      DRAM area per die: 89.3066 mm2
83      Area efficiency: 41.7953%
```

| Validation Targets | Latency (ns) | Real / Model / Err | Power (mW) | Real / Model / Err | Area ($mm^2$) Real / Model / Err |
|---|---|---|---|---|---|
| **Samsung 2Gb 80nm DDR2** | $t_{RCD}$ | 8.8 / 8.39 / -4.6% | | N/A | 195.64 / 189.89 / -2.9% |
| | $t_{RP}$ | 9.0 / 9.04 / 0.4% | | | |
| **Micron 1Gb 78nm DDR3** | $t_{CAS}$ | 15 / 13.9 / -7.3% | $P_{ACT}$ | 90.6 / 93.9 / 3.6% | 100.22 / 90.45 / -9.8% |
| | $t_{RAS}$ | 36 / 32.1 / -10.9% | $P_{RD}$ | 57.1 / 62.4 / 9.4% | |
| | $t_{RC}$ | 51 / 50.6 / -0.9% | $P_{WR}$ | 39.0 / 37.5 / -3.9% | |
| **Samsung 8Gb 3D 60nm DDR3** | $t_{CAS}$ | 15 / 14.9 / -0.9% | | N/A | 98.1×4 / 100.3×4 / 2.2% |

TABLE III

- The CACTI-3DD does not support 60nm tech, thus uses 50nm, can notice that it is actually quite accurate.

# CACTI-IO Technical Report

- https://escholarship.org/content/qt6f6904sc/qt6f6904sc_noSplash_d5a14fc9a279c9b4c6c4b353a7e61ad5.pdf?t=rri97w

# Row, Column array number

- These are for subarrays calculations and output the correct subarray number in parameter.cc

- Senior uses 128×128, adjust it to 128×128. 1Gb per die.
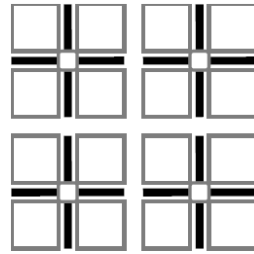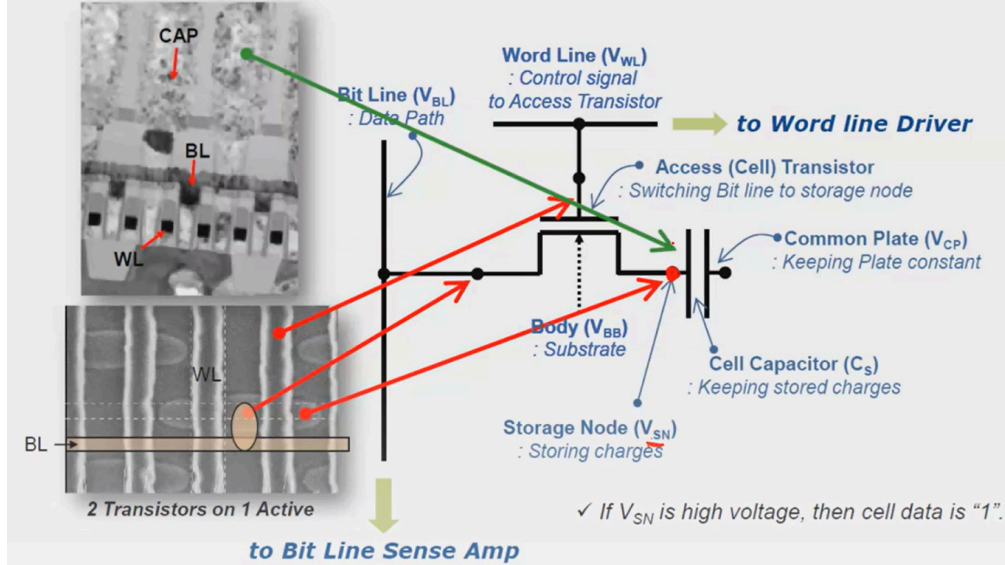
## Memory density parameter(Not used)



## Ndwl,Ndbl,Nspd & its relationship with subbanks

# A Closer Look at DRAM Cell

$N_{dwl}$ = 4
$N_{dbl}$ = 4
$N_{spd}$ = 1
$N_{subbanks}$ = 2
$N_{mats-in-subbank}$ = 2

$N_{dwl}$ = 8
$N_{dbl}$ = 2
$N_{spd}$ = 2
$N_{subbanks}$ = 1
$N_{mats-in-subbank}$ = 4

$N_{dwl}$ = 8
$N_{dbl}$ = 4
$N_{spd}$ = 1
$N_{subbanks}$ = 2
$N_{mats-in-subbank}$ = 4

Figure 10: Different partitions of a bank.

In order to calculate the optimal organization based on a given objective function, like earlier versions of CACTI [30, 38, 40, 47], each bank is associated with partitioning parameters $N_{dwl}$, $N_{dbl}$ and $N_{spd}$, where $N_{dwl}$ = number of segments in a bank wordline, $N_{dbl}$ = number of segments in a bank bitline, and $N_{spd}$ = number of sets mapped to each bank wordline.

| Parameter Name | Meaning | Parameter Type |
|---|---|---|
| $N_{banks}$ | Number of banks | User input |
| $N_{dwl}$ | Number of divisions in a bank wordline | Degree of freedom |
| $N_{dbl}$ | Number of divisions in a bank bitline | Degree of freedom |
| $N_{spd}$ | Number of sets mapped to a bank wordline | Degree of freedom |
| $D_{bitline-mux}$ | Degree of muxing at bitlines | Degree of freedom |
| $D_{senseamp-mux}$ | Degree of muxing at sense amp outputs | Degree of freedom |
| $N_{subbanks}$ | Number of subbanks | Calculated |
| $N_{mats-in-subbank}$ | Number of mats in a subbank | Calculated |
| $N_{subarr-rows}$ | Number of rows in a subarray | Calculated |
| $N_{subarr-cols}$ | Number of columns in a subarray | Calculated |
| $N_{subarr-senseamps}$ | Number of sense amplifiers in a subarray | Calculated |
| $N_{subarr-out-drivers}$ | Number of output drivers in a subarray | Calculated |
| $N_{bank-addr-bits}$ | Number of address bits to a bank | Calculated |
| $N_{bank-datain-bits}$ | Number of datain bits to a mat | Calculated |
| $N_{bank-dataout-bits}$ | Number of dataout bits from a mat | Calculated |
| $N_{mat-addr-bits}$ | Number of address bits to a mat | Calculated |
| $N_{mat-datain-bits}$ | Number of datain bits to a mat | Calculated |
| $N_{mat-dataout-bits}$ | Number of dataout bits from a mat | Calculated |
| $N_{mat-way-select}$ | Number of way-select bits to a mat (for data array of cache) | Calculated |

Table 1: Organizational parameters of a data array.

## Area Efficiency , Bandwidth & Energy

To calculate the bandwidth of 3D-DRAM (like HBM, HBM2, or HBM3), you can use the following formula:

$$\text{Bandwidth} = \text{Data Rate} \times \text{Bus Width} \times \text{Number of Channels}$$

### Explanation of Terms:

1. **Data Rate**: This is the speed at which data can be transferred per pin, typically measured in gigabits per second (Gbps). For example, HBM2 might have a data rate of 2 Gbps to 3.2 Gbps per pin.

2. **Bus Width**: This refers to the width of the data bus per channel, typically 128 bits (or 16 bytes) for HBM technologies.

3. **Number of Channels**: The number of independent channels available in the memory module. HBM2 usually has 8 channels, while HBM3 might have more.

### Example Calculation:

Suppose you have an HBM2 stack with the following specs:

- **Data Rate**: 2 Gbps
- **Bus Width**: 128 bits (16 bytes)
- **Number of Channels**: 8

$\downarrow$

To calculate the data rate of a customized 3D-DRAM, you can use the following approach:

$$\text{Data Rate} = \frac{\text{Clock Frequency} \times \text{Data Transfers Per Clock Cycle}}{\text{Bus Width}}$$

## Explanation of Terms:

1. **Clock Frequency**: The operating frequency of the DRAM, typically measured in megahertz (MHz) or gigahertz (GHz). This determines how many times data can be transferred per second.

2. **Data Transfers Per Clock Cycle (MT/s)**: This depends on the technology:

   - **DDR (Double Data Rate)**: Transfers data twice per clock cycle (once on the rising edge and once on the falling edge).

   - **QDR (Quad Data Rate)**: Transfers data four times per clock cycle.

   - If it's a customized 3D-DRAM, decide how many transfers you want per clock cycle.

3. **Bus Width**: The width of the data bus per channel, measured in bits. This determines how much data is moved per transfer. Typical 3D-DRAM configurations use 128-bit buses, but you can customize this.

# Number of TSV



- The TSV bit width is determined by the burst depth * io_width, usually burst depth is 2 though, as for HBM
- These information are for a single bank