

Cactus

Technology mode can't be specified  
10.8 um fixed

(A)

b) enhanced access and cycle time model.  
for on-chip caches.

→ derive relatively simple equations that predict the access/cycle times of cache as a function of various cache parameters, process parameters, and array organization parameters.

→ A Spice model of the cache was used to validate their analytical model.

→ This only shows that model matches Spice model, it does not address the issue of how well the assumed cache model structure reflects a real cache design.

→ Code needs to be compiled using the make file provided.

→ time.c is the model.

→ transistor widths and process parameters are defined in def.h.

→ choose the organisation that gives the smallest access time.

→ Associativity in cache??

$$\begin{array}{r} 32 \rightarrow 68 \\ 32 \times 2 \end{array}$$

Date / /  
Page No.

→ Number of sets  $\rightarrow$  Cache size (C)  
(S)  $\frac{\text{Block size (B)} \times \text{Associativity (A)}}{}$

④ 25  
2  
1024  
32

→ Such an organisation could result in an array that is much longer in one direction causing either bitline or wordlines to be very slow.

→ leads to longer than necessary access time.

→  $N_{dw}$  → indicates how many times the array has been split with vertical cut lines (creating more but shorter wordlines).

→  $N_{dh}$  → indicates how many times the array has been split with horizontal cut lines (causing shorter bitlines).

→ total subarrays  $\rightarrow N_{dw} \times N_{dh}$ .

array  
organizational  
partitioning

→  $N_{hd}$  → how many sets are mapped to a single wordline

→  $N_{tw}, N_{hd}, N_{spd}$  → for tag array

→ Circuit gets decomposed into many equivalent RC circuits, and using simple RC equations to estimate the delay of each stage.

Decade

Sense Amplifiers

Comparators

Column MUX

Wordline

Bitlines

MUX Drivers

Output Drivers

\*Fact 2

↳ An integrated cache timing and power model.

→ Access time, cycle time, power.

→ It added support for full associative caches and for caches with multiple access ports

→ user can specify read/write ports, read only ports and write only ports

\*Fact 1 cache-size block-size assoc tech mode.  
read/write ports read ports write ports.

### Cacti 3

→ timing, power and area model.  
aspect ratio; efficiency.

- fully independent banking of caches.  
This allows users to split the cache into a number of banks with its own set of address and data buses.

### fully associative cache

$$N_{dwl} = N_{wpd} = 1$$

$N_{tbl}$  → can be varied.

parti cache size cache-line assoc tech-mod.  
head/write ports head ports write ports  
m-banks.

### Cacti 4

- leakage power added.
- web interface.

### H-Mac?

- if requested no. of ports is equal to or less than the number of banks, we assume each bank only implements one head/write port.
- can also model memories without tags

## Cacti<sup>5</sup>

- dynamic power, access time, area, and leakage power of caches and other memories.
- cache organisation assumed by cacti is now output graphically to assist users in understanding the output generated by cacti

Date / /  
Page No.

main.cc → takes input file / command line arguments  
 and passes control to carti interface  
 ↗ linked to io.h  
 (uca-org-t  
data)

io.h → linked to const.h, carti-interface.h.

io.cc → io.h, area.h, basic-circuit.h, parameter.h,  
 Ucache.h, muca.h, crossbar.h, arbiter.h.

→ Input Parameter class

↳ parse-cfg func. → parses the cache.cfg  
 file.

↳ display\_ip() → displays the input  
 parameters

+  
 power components } operator  
 powerDef } overloading

uca-org-t      carti\_interface( )      output  
 ↗                  ↗                  ↗ final results.

↳ function overloading  
 (one for file input  
 and other for command  
 input).

uses a 'solve' function to  
 solve the equations.

Δ

error-checking func.

↳ output-data.csv → creates a csv file.

↳ output\_UCA → prints the results.

cacti-interface.h → linked to const.h

↳ defines 3 classes

↳ mem-values +  
↳ memm\_array.  
↳ uca-org-t

→ ~~for~~ defines power components and  
powerDef classes

→ wire types definition

→ defines Input Parameter class

this will change.

uca-org-t class

↳ variables for data and tag array,  
access time, cycle time, area,  
area efficiency, power etc.

→ functions → find-delay, find-energy,  
find-area, find-cyc

(mem-array class → Npd, NdwL, Ndp, area,  
power, delay of  
components)

↳ modify this if my  
component is added.

cacti-interface.cc → linked to area.h,  
basic-circuit.h, component.h,  
const.h, parameter.h,  
cacti-interface.h

↳ defines the functions of uca-org-t  
class.

siml-delay, siml-energy, siml-cyc  
 → return siml values if value of for  
 data and tag arrays are known.

\* const.h

- ↳ defines constants for project
- address bits
- output bits
- max. value of nval, ntbl, ntblc etc
- threshold voltages

*Important* ↳ 6 for transistors.

↳ needs to be changed for more transistors.

\* area.h. → linked to cacti-interface.h,  
 basic\_circuit.h

cell height,  
 width  
 $\text{area} = h * w$

↳ creates 'Areci' class

furn. for  
 getting width,  
 height and  
 returning area

area.cc → linked to area.h, component.h,  
 decoder.h, parameter.h,  
 basic\_circuit.h

→ has no node.

\* basic\_circuit.h → linked to const.h,  
 cacti\_interface.h

basic\_circuit.cc

functions → powers, isbow2, log2,  
 log-two, gate-c, gate-c-pass  
 drain-c, t4\_r\_en, R-to-w

Hoerwitz, simplified pmos-to-mmos size relation.  
 CMOS I leak, mmos leakage, simplified pmos leakage,

$-\log_2 \rightarrow$  log base 2 of num.  
 is power  $\rightarrow$  check if num is power of 2  
 powers  $\rightarrow$  to calculate ~~power~~ base exponent.

gate-c, gate-c pair  $\rightarrow$  same.  
 $\hookrightarrow$  To write gate capacitance in farads.

no-h3  $\ominus$  Technology Parameter class ??

R-to-W  $\rightarrow$  given Resistance, need width  
 Transistor width in um.

p-mos-to-mmos-size ratio.  
 $\hookrightarrow$  sizing ratio.

Hoerwitz  $\rightarrow$  timing models for MOS circuits  
 by Frank Hoerwitz.

CMOS-I leak, simplified-mmos-leakage,  
 simplified-pmos-leakage  
 $\hookrightarrow$  gives off current for CMOS,  
 mmos and pmos respectively.

\* parameter.h

Technology Parameter class

- ↳ parameters which are functions of certain device technology.
- ↳ capacitance, resistance,  $V_{th}$ ,  $I_{on}$ ,  $I_{off}$ ,  $C_{ox}$ ,  $t_{on}$

→ 3 classes inside class

- ↳ Device Type
- ↳ Interconnect Type
- ↳ Memory Type

$w\_comp\_ime\_p1$ ,  $w\_comp\_ime\_p2$ ,  $w\_comp\_ime\_p3$ ,  
 $w\_comp\_ime\_p4$ ,  $w\_comp\_ime\_p5$ ,  
 $w\_comp\_ime\_m3$

↳ & meaning of parameters  
 not clear

time  
64 → 73

Dynamic Parameter class

- ↳  $N_{NpD}$ ,  $N_{NwL}$ ,  $N_{NbL}$

parameter.cc

↳ defines functions for above classes  
 like display for displaying parameter

\* component.h → linked to area.h, parameter.h

Date / /  
Page No.

Classes → Crossbar, Bank, Component

compute - gate area,  
compute - th - width - after -  
folding

height - sense - amplifier.  
logical - effort.

compute - diffusion - width.

component.cc → defines the above functions

logical - effort, compute height - sense - amplifier  
→ need to understand these functions.

\* technology.cc → linked to basic - circuit.h,  
~~parameter.h~~  
parameter.h

→ 3 functions

↳ wire resistance.

↳ wire capacitance.

↳ init - tech - params.

↳ from 376 to 382

→ SRAM cell parameters.

↳ for 90nm.

↳ for 65nm → 574 - 580

↳ for 45nm → 776 - 780

↳ for 32nm → 976 - 980

defines  
technology  
parameters

Vdd, gate length,  
Cost, width etc.

\* arbiter.h → linked to basic\_circuit.h,  
 cacti-interface.h, component.h,  
 parameter.h, mat.h, wire.h.

Arbiter class → inherited from component class.

private variables → NTm<sup>1</sup>, PTm<sup>1</sup>, NTm<sup>2</sup>, PTm<sup>2</sup>,  
 R, PTi, NTi, flit-size,  
 NT<sub>H</sub>, PT<sub>H</sub>, T<sub>hi</sub>S<sub>1</sub>, T<sub>hi</sub>S<sub>2</sub>  
 → ??

arbiter.cc  
 ↴ line 56-63 → ??

→ width of m-mos and pmos transistors.

↳ used in leakage power computation.

\* crossbar.h → linked to basic\_circuit.h,  
 cacti-interface.h, component.h,  
 parameter.h, mat.h, wire.h.

Crossbar class → inherited from component class.

Crossbar.cc → output-buffer, compute-power functions.

prints → flit size, width and height  
of array, dynamic and  
leakage power.

\* Ucache.h → linked to area.h, router.h,  
mca.h.

min-values\_t class → update-min-values  
function.

struct solution.

calculate-time, solve, init-tech-params  
functions.

Ucache.cc → linked to area.h, bank.h,  
basic-circuit.h, component.h,  
const.h, decoder.h, parameter.h,  
subarray.h, uca.h.

a) uses UCA class from uca.h.

functions → check-uca-org, check-mem-org,  
find-optimal-uca, filter-tag-arr,  
filter-data-arr

\* uca.h. → linked to area.h, bank.h, component.h, parameter.h, htrec2.h.

class UCA → inherited from Component class.

compute-delay(),  
returns extrisetime.  
↓  
defined in  
uca.cc.

compute-power-energy()  
power for a bank  
of an array.

uses functions from bank.h and  
bank.cc.  
and from htrec2.h.  
for delays and power.

\* bank.h. → linked to component.h, decoder.h, prot.h, htrec2.h.

Bank class → ~~is~~ inherited from Component class

functions → compute-delay, compute-power-energy

Bank.cc

Bank constructor → Head-write ports,  
exclusive Head and write  
ports.

address bits, data input bits, data output bits.

compute delays → uses mat. compute-delays()  
 ↳ in mat.h

compute-power-energy() → uses  
 mat. compute-power-energy()  
 ↳ computes dynamic and leakage  
 power for read operation.

\* mat.h. → linked to component.h, decoder.h,  
 wire.h, subarray.h.  
 class)

Mat class → inherited from component class.

functions → compute-delays,  
 compute-power-energy.

ols  
 mat.cc  
 ↳ lime 92 → 2 transistors per cell??

lime 364 → delay-fa-tag ()  
 ↳ definition??  
 ↳ lime 490.

functions → width-write-driver-on-write-max  
 compute-comparators-height,  
 compute-bitline-dddy., compute-sa-delay.

compute - subarray - out - arr  
compute - comparator - delay

compute - power - energy

\* Subarray.h → linked to area.h,  
Component.h, parameter.h

Subarray class → inherited from Component class

variables → num-Hous, num-cells,



cell

↳ (instance of Area class)

CWL, C-bl → capacitances.

func. → get\_total\_cell\_area

compute-C()

↳ for bitline and wordline  
capacitance.

cell.get\_area() \* num-Hous \* num-cells

\* Htree2.h → linked to basic-circuit.h,  
component.h, parameter.h,  
assert.h, subarray.h,  
ctrl-interface.h, wire.h

Htree2 class → inherited from component class.

functions → in-htree(), out-htree(),  
limited-in-htree(), limited-out-htree(),  
input-command, output-buffer

intree2.cc

Date: / /  
Page No.:

8(1)

input - mand → mand gate sizing calculation.

output - buffer → tristate buffer model  
consisting of not, mand, nor  
and driver transistors.

say!

\* decoder.h. → linked to area.h,  
component.h, parameter.h.

1.h,

Decoder class → inherited from component  
class.

func. → compute\_widths, compute\_area,  
compute\_delays  
↳ returns otherwise time.

nl

rl

l

PreddecBlk class → inherited from component  
class.

??

PreddecBlk Dv → inherited from component  
class

Preddec class → inherited from component  
class.

-max  
U,

-lay.

Driver class

Date / /  
Page No.

\* decoder.cc

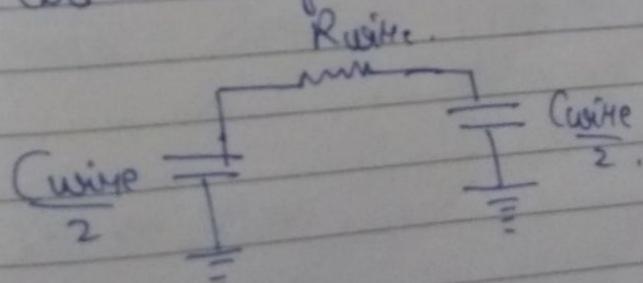
\* wire.h. → linked to basic-circuit.h, component.h,  
parameter.h, cacti-interface.h.

wire class → inherited from component  
class.

## \* Cacti 5.0 Tech Report

Date / /  
Page No.

- dynamic power, access time, area and leakage power of caches and other memories.
- uses ITRS ~~for~~ technology projections 90 nm, 65 nm, 45 nm, 32 nm.
- optimization steps →
  - find all solutions with area within a certain percentage of area of solution with best area efficiency.
    - max area constraint.
  - find all solutions with access time that is within a certain percentage of best access time.  
solution in this reduced subset of solutions
    - max acc time constraint. ?
  - apply opt. function.
- wire delay model.



Ecc → Error correction code.

- a variable to incorporate + specify the number of data bits per Ecc bit
- default value = 8

### Data Array Organisation

Highest Level → Independent Banks

$N_{\text{banks}}$ .

↓  
multiple identical subbanks.  
 $N_{\text{subbanks}}$  (1 subbank activated per access).

↓  
multiple identical mats

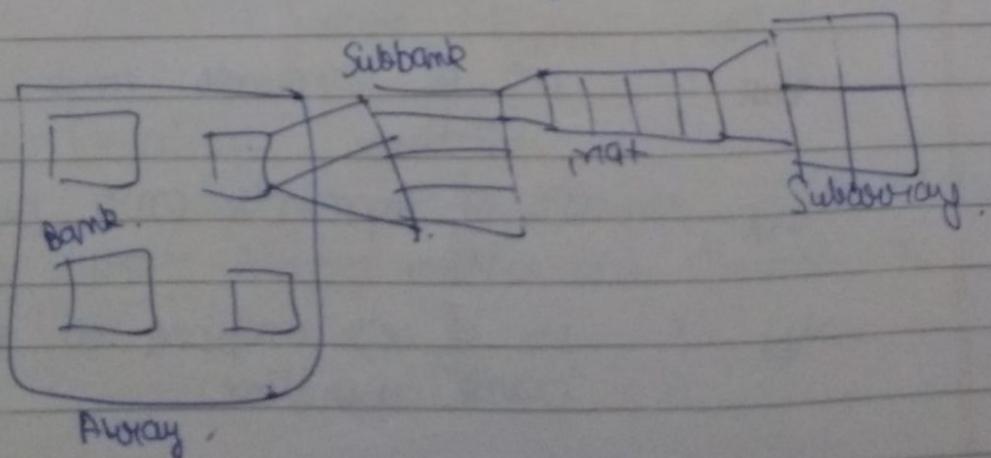
$N_{\text{mats-in-subbank}}$ .

(All mats are activated per subbank access).

↓  
4 identical subarrays

and associated predecoding logic.

Subarray will be the part where cell and its peripheral circuitry will be defined.



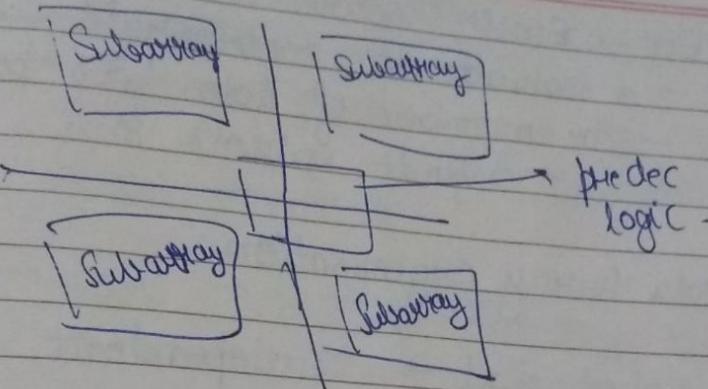
+ Mat

Date / /  
Page No.

(C)

guy

Subarray



Sub array

Precharge and  
Equilibration.

2D array  
of memory cells.

Wordline  
drivers

How  
decode gates.

Bitline MUX  
Sense Amp  
Sense Amp MUX  
Subarray Out Drivers  
Write MUX and Drivers

$N_{\text{row}}$  = no. of segments in a bank wordline.

$N_{\text{tbl}}$  = no. of segments in a bank bitline.

$N_{\text{slot}}$  = no. of sets mapped to each bank wordline.

- Each predecode block has two levels.
- the first level is composed of one 2-4 decode unit and one 3-8 decode unit.
- at second level, 4 outputs from 2-4 decoder and 8 outputs from 3-8 decoder are combined using 32 NAND<sub>2</sub> gates.

$$R_{\text{Predec-output-wire}} = \frac{L_{\text{Predec-output-wire}}}{\text{Unit-length}_w, h}$$

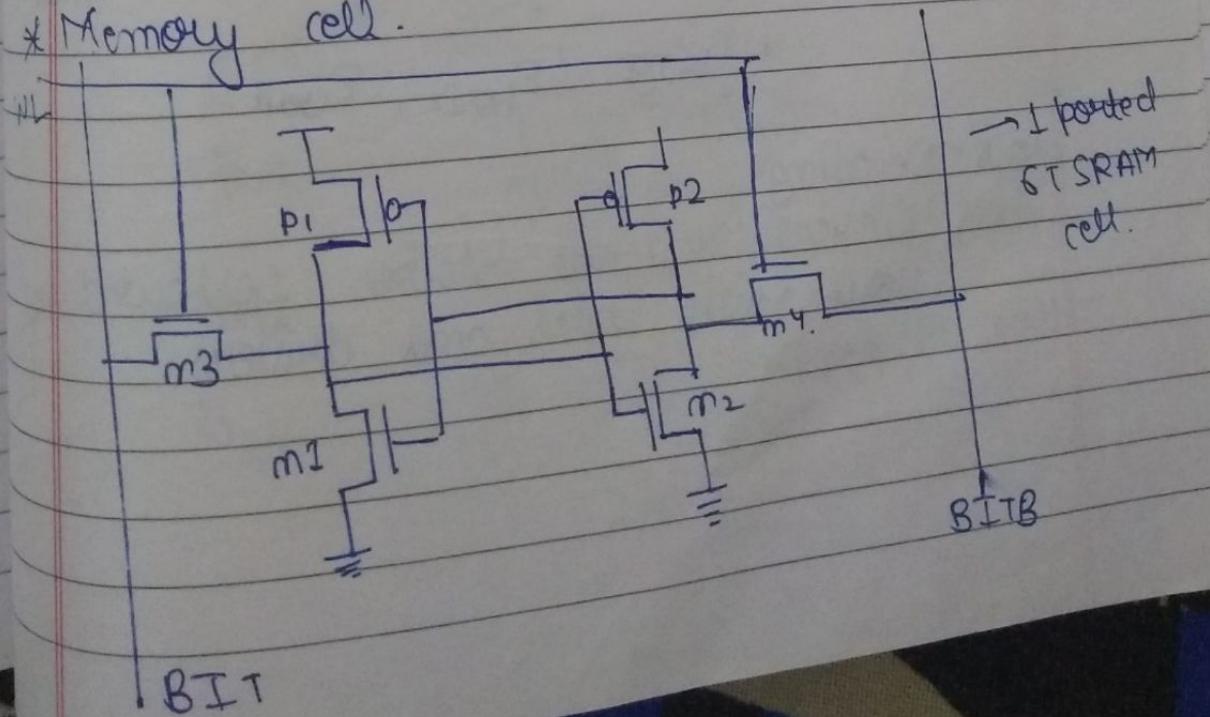
$$C_{\text{Predec-output-wire}} = \frac{L_{\text{Predec-output-wire}}}{\text{Unit-length}_w, h}$$

↓

max. length amongst  
lengths of predecode  
output wires.

→ sizing of gates in each circuit path is calculated using the method of logical effort.

### \* Memory cell.



$$N_{\text{subbands}} = \frac{N_{\text{dbl}}}{2}$$

$$N_{\text{proto-im-subbands}} = \frac{N_{\text{dbl}}}{2}$$

### Wire Modelling

$$R_{\text{wire}} = L_{\text{wire}} R_{\text{unit-length-wire}}$$

$$C_{\text{wire}} = L_{\text{wire}} C_{\text{unit-length-wire}}$$

$L_{\text{wire}} \rightarrow$  length of wire.

$$R_{\text{unit-length-wire}} = \alpha_{\text{scatter}} \frac{P}{(\text{thickness} - \text{barrier} - \text{dishing}) * (\text{width} - 2 * \text{barrier})}$$

$$C_{\text{unit-length-wire}} = \epsilon_0 (2 M \epsilon_{\text{horiz}} \frac{\text{thickness}}{\text{spacing}} + 2 \epsilon_{\text{vert}} \frac{\text{width}}{I D_{\text{thick}}})$$

+ fringe ( $\epsilon_{\text{horiz}}, \epsilon_{\text{vert}}$ )

### Row-Decoding

↳ two row pre-decode blocks followed by row-decode gates and drivers.

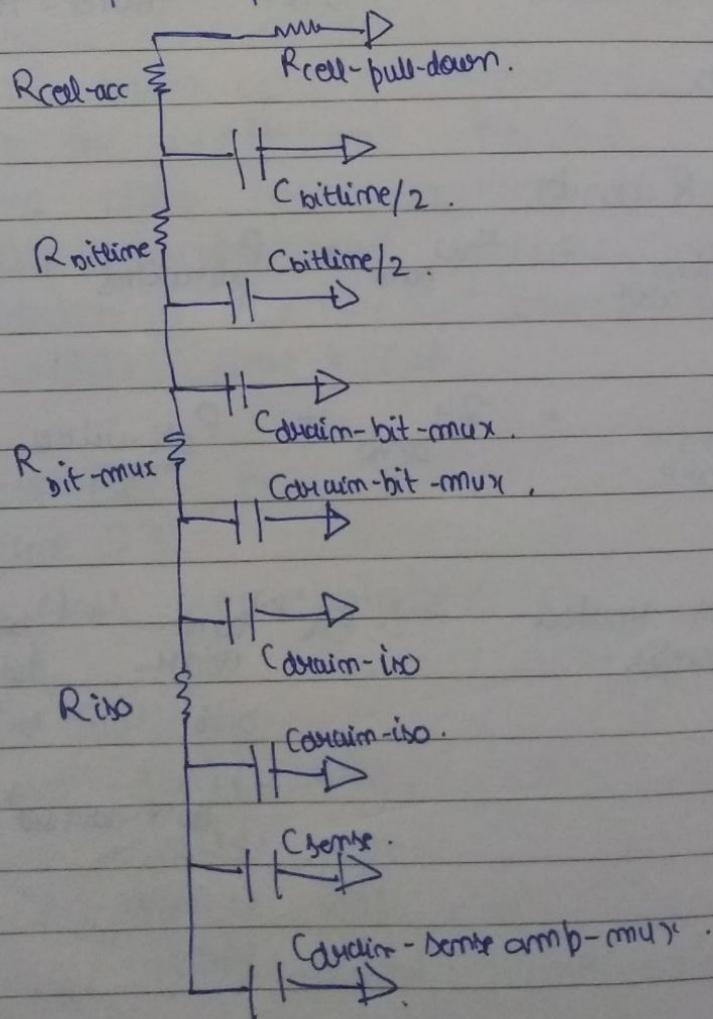
Read [7] → for SRAM cell size and width.

### Peripheral Circuitry

- Sense Amplifier. → clock-latch based sense amplifier.
- Bitline and Sense Amp. MUXes.
  - ↳ NMOS based MUX.
- Precharge and Equalization Circuitry.
  - ↳ fixed widths for NMOS.

~~fixed widths for NMOS from [20]~~

Circuit Model of blue memory cell and Sense Amp. Input.



[26] → analysis of latch-based sense amplifiers  
and equations for sensing delay under different assumptions. (1)

### \* Area Model

total diffusion width → eq<sup>m</sup> 14  
page 27.

$$\text{Area of data array} = H_{\text{data}} \cdot W_{\text{data}} \cdot h,$$

$$P_{\text{all-wires}} = P_{\text{wires}} \cdot N_{\text{wires - Mouted-to-banks}}$$

pitch

for 8 banks.

$$W_{\text{data}} = 4W_{\text{bank}} + P_{\text{all-wires}} + 2 \frac{P_{\text{all-wires}}}{4}$$

$$W_{\text{data}} = 24_{\text{bank}} + \frac{P_{\text{all-wires}}}{2}$$

$$N_{\text{wires - Mouted}} = 8 \left( N_{\text{bank-addy}} + N_{\text{bank-data}} + N_{\text{data-select}} \right)$$

$$+ N_{\text{bank-data out}} + N_{\text{way-select signals}}$$

max  
V,  
4.

$N_{\text{banks}} \text{ in horizontal dim}^m$ .

$\Rightarrow 2 \times N_{\text{banks}} \text{ in vertical dim}^m$

Page 30 -

Area of Mat equations.

eq<sup>m</sup> (33) and (3)

\* \* eq<sup>m</sup> 33 and 34 uses height and width of memory cell

eq<sup>m</sup> 38  $\rightarrow$  H<sub>subarr-bitline</sub>  $\rightarrow$  will be changed.

\* area.h, area.cc  $\rightarrow$  returns area given the width and height,

Area class has functions to set width, height and getArea() for returning the area based on saved width and height.

↓  
Nothing has to be changed in these files.

\* parameter.h.

line 123

class MemoryType &

{

public :

double b\_w;

double b\_h;

double cell\_a\_w;

double cell\_pmos\_w;

double cell\_mmos\_w;

double vbitbar;

will probably will be changed.

} memory cell related a

line 164-173

Date / /  
Page No.

(j)

parameter cc → has a function definition  
for display() → to display the 3x11  
parameters.

\* Delay Modeling → access time,  
→ random cycle time.

Access time

↳  $T_{\text{access}}$  → includes time taken by  
row-decoder path, bit-line  
decoder path, sense amp decoder  
path.

$$T_{\text{bitline}} = \left\{ \begin{array}{l} \sqrt{2T_{\text{step}} \frac{V_{DD} - V_{TH}}{m}}, \quad T_{\text{step}} \leq 0.5 \frac{V_{DD} - V_{TH}}{m} \\ . \quad T_{\text{step}} + \frac{V_{DD} - V_{TH}}{2m}, \quad T_{\text{step}} > 0.5 \frac{V_{DD} - V_{TH}}{m}. \end{array} \right. ?$$

from [28]

$T_{\text{step}}$  → uses  $R_{\text{cell}}$  pull-down,  $R_{\text{cell}}$  acc.

↳ will change

$$m = \frac{V_{DD} - V_{TH}}{2T_{\text{step}}}$$

### Random Cycle Time

- ↳ limited by wordline and bitline delays.
- ↳ will not change except  $V_{bitline}$ -swing

### \* Power Modeling

$$V_{bitline\_swing} = 2V_{sense}.$$

↳ this may change.

- ↳ dynamic energy
- ↳ leakage power.

eqn 75 and 76 may change as  $V_{bitline}$ -swing may change.

↳  $W_{line-mmols}$ ,  $W_{line-pmol}$  → width of transistors

eqn 91, 92 → related to power consumed by cells.

$$P_{mem\_cell} = V_{DD} I_{mem\_cell}$$

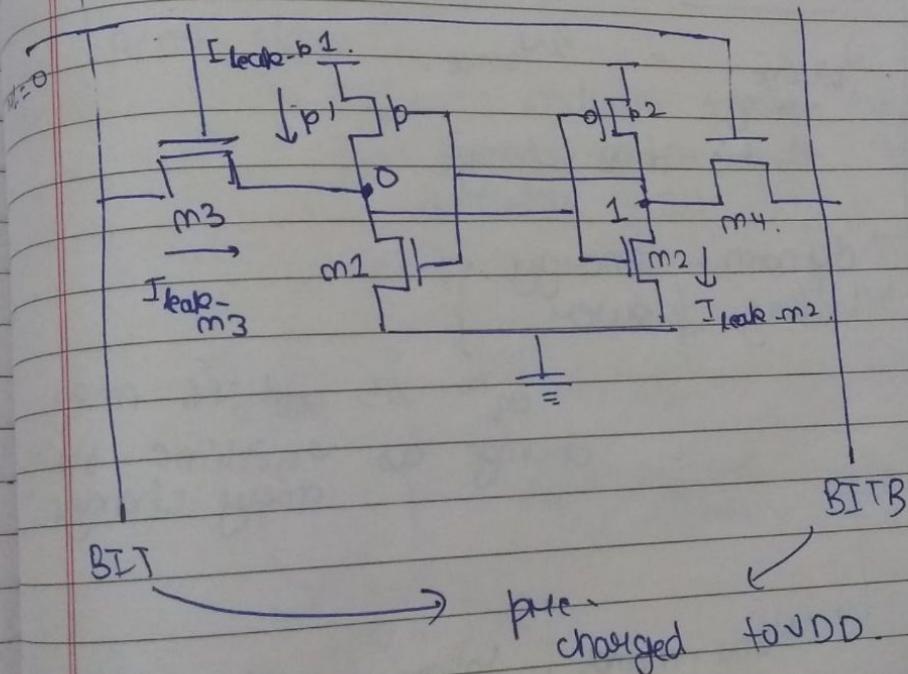
$$I_{mem\_cell} = I_{P_1} + I_{m_2} + I_{m_3}$$

$$I_{P_1} = w_{P_1} I_{off\_pmos}$$

$$I_{m_2} = w_{m_2} I_{off\_mmos}$$

$$I_{m_3} = w_{m_3} I_{off\_mmos}$$

These  
will  
change.



## Technology Modelling

→ Technology projections from ITRS

→ 180, 90, 65, 45, 32 nm → nodes.

Three Device Types

- High performance (HP)
- Low Standby Power (LSP)
- Low Operating Power (LOP).

can be set for data array cell-type,  
data-array peripheral type, tag array  
cell-type and tag array peripheral  
type in configuration file.

Table - 5 page 39

Line 575 - 579 0.155 mm mcr in technology.cc	Technology data - for SRAM cell.
	Area of SRAM cell
	$w_{SRAM-cell-acc}$ width of SRAM cell access transistor.
	$w_{SRAM-cell-pd}$ width of SRAM cell pull-down transistor.
	$w_{SRAM-cell-pu}$ Pull up transistor.
	AR <sub>SRAM-cell</sub> Aspect ratio of cell.

→ TRIASTAR can be used to obtain device data for scaling models and assumptions that are different from those of ITRS.