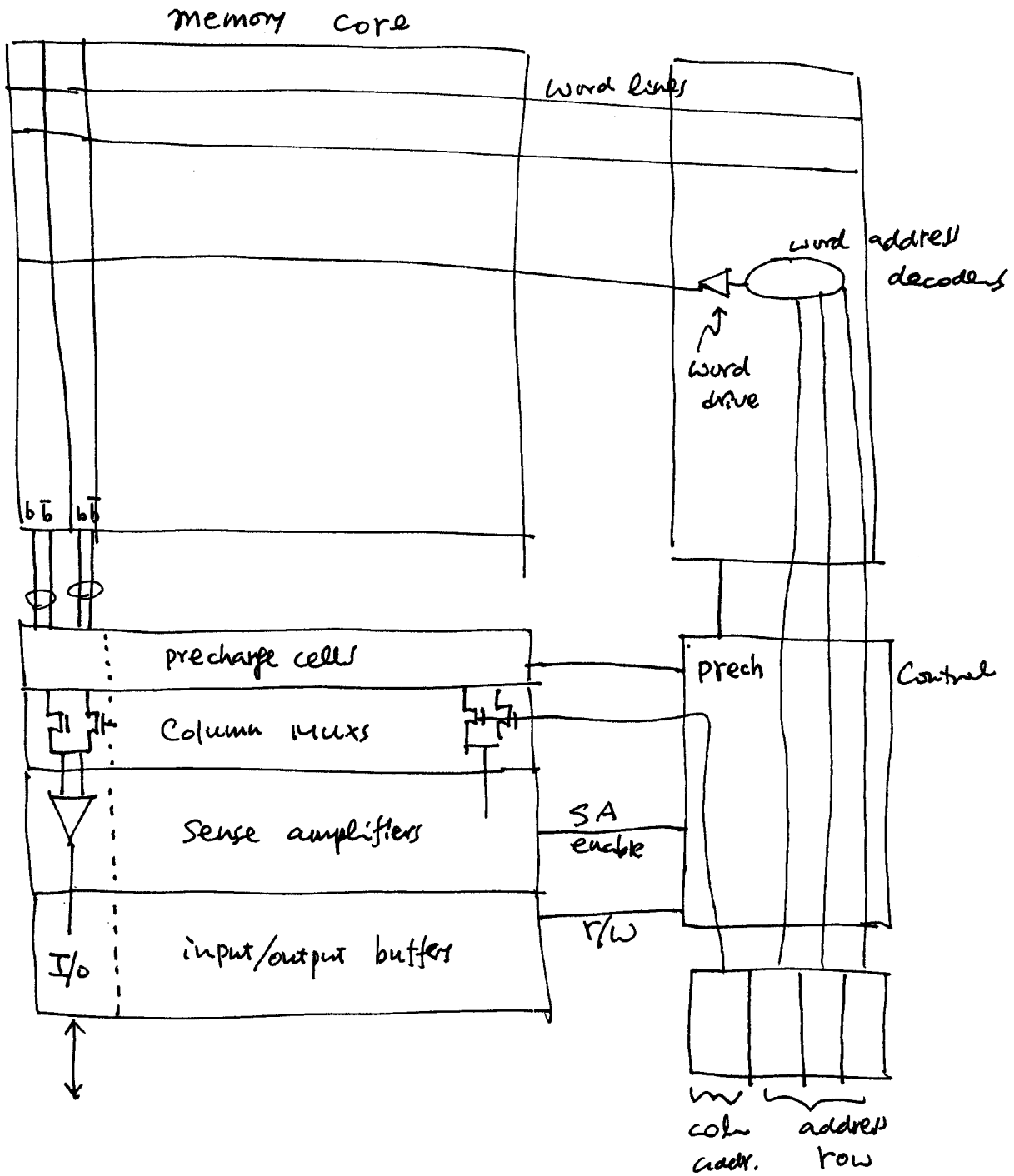
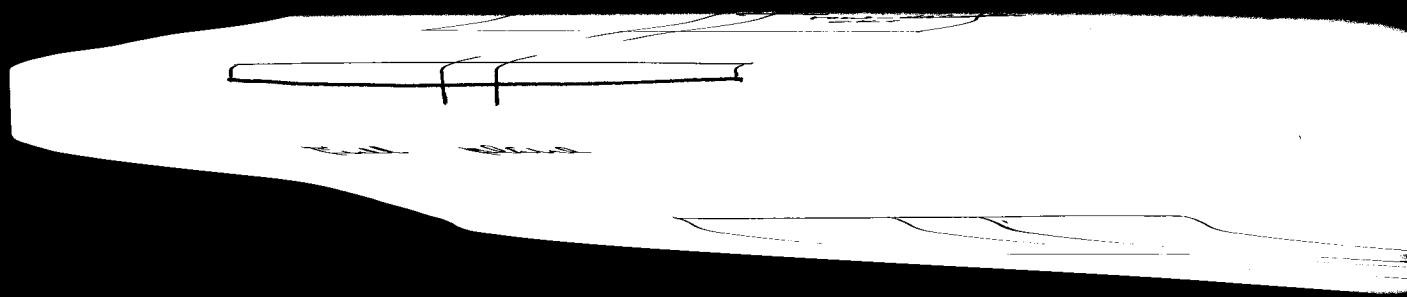
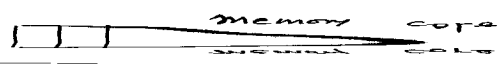


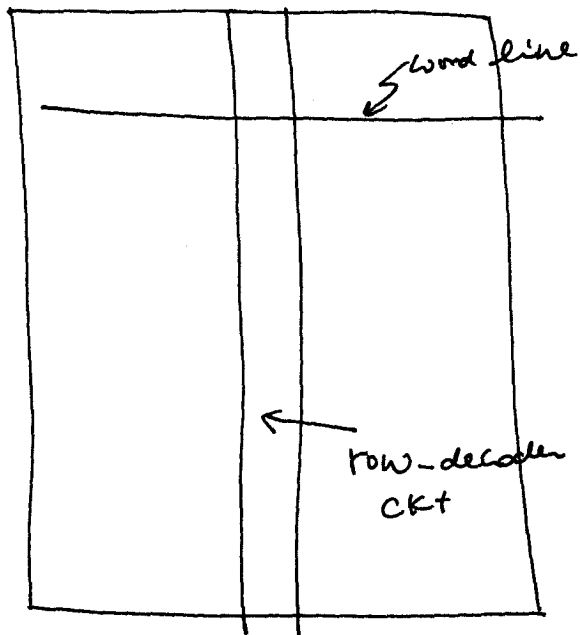
Organization of a RAM.



Organization of a RAM.

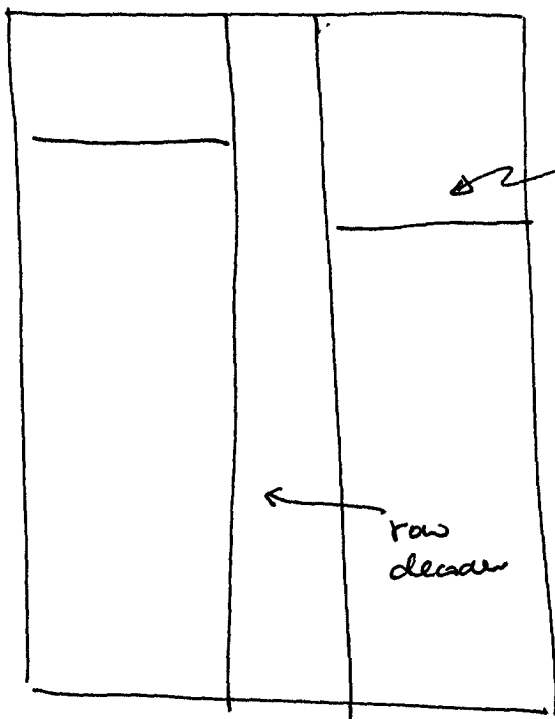


Memory Array Architecture



- Each ^{row} address accesses both sides of the array
- Good for ^{large} capacity array. _{not too}

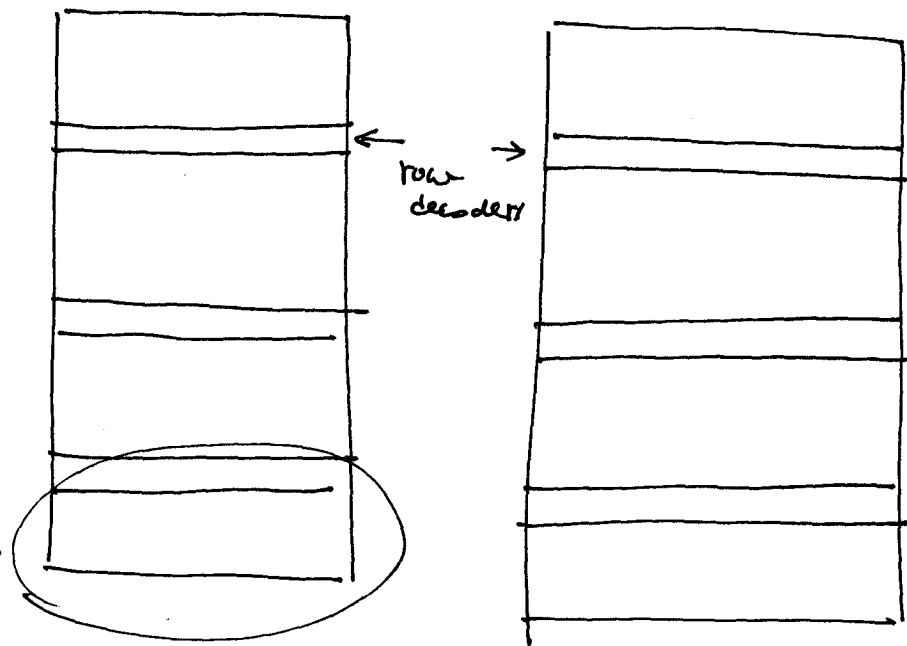
Full plane



- Each row address accesses one side of the array
- Lower power consumption.

Half plane

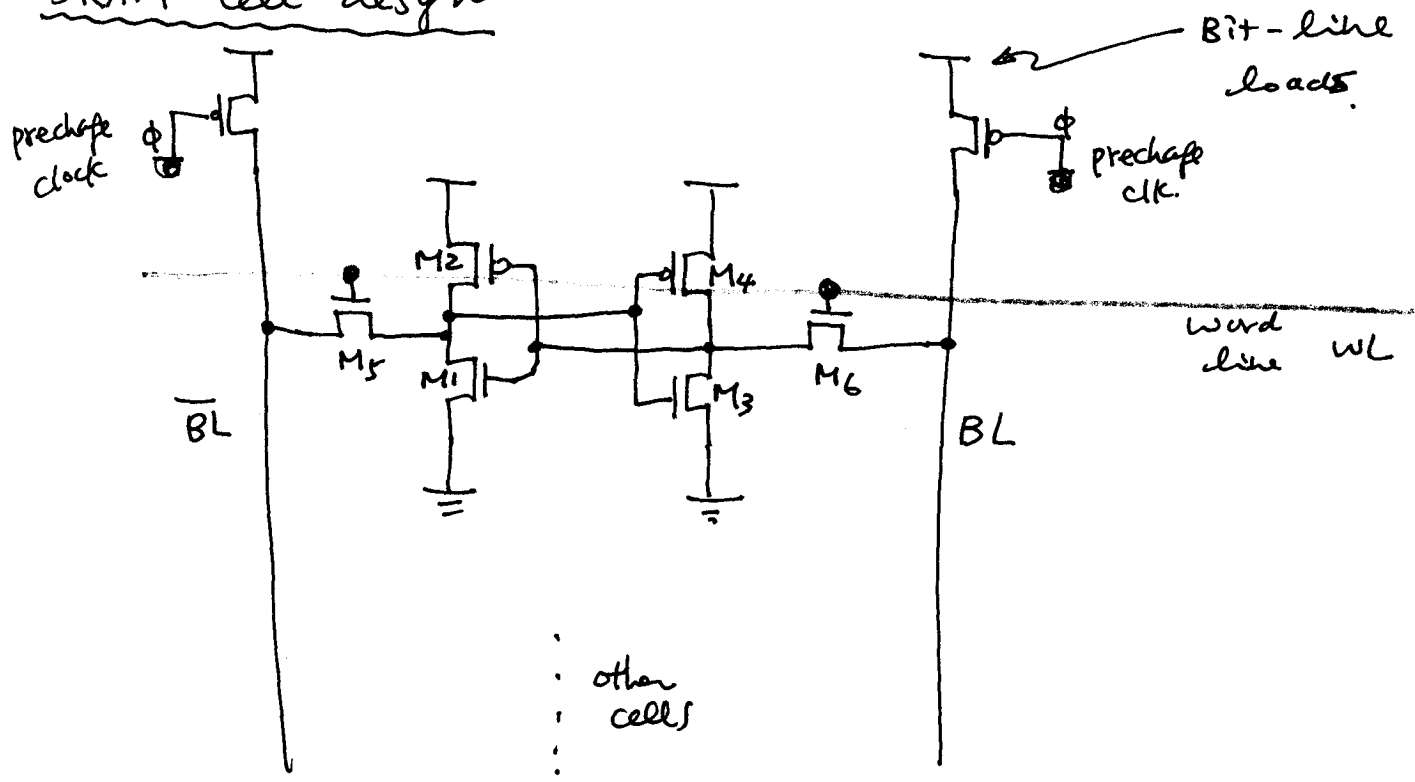
- Block-oriented



• Reduce bit line length

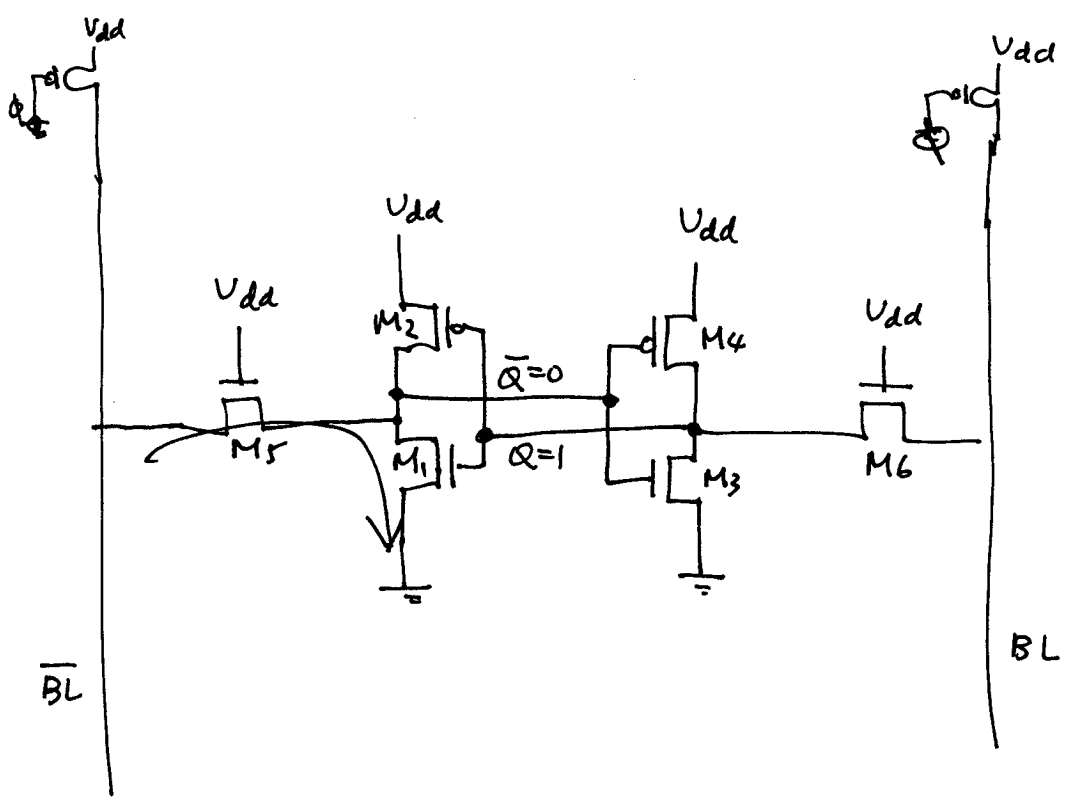
Each block contains its own block selection (to select a specific block) ckt, x-decoder and sense amplifier. (select one of the rows)

• SRAM cell design



• Transistor sizing is very important.

• Example for Read



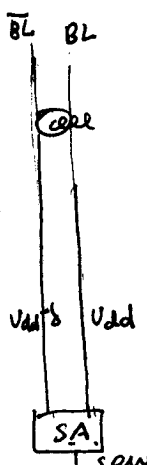
• Assume the cell stores logic 1

$\therefore Q=1, \bar{Q}=0$

• Before read, both bit lines are precharged to 1.

• We expect $BL=1$, but \bar{BL} will be pulled-down

• Capacitances of \bar{BL} and BL are much larger (PF) than transistor capacitances, so the swing of \bar{BL} will not be too large



\therefore needs a sensor amplifier to amplify the signal.
 \hookrightarrow amplify it!!

• $\bar{Q}=0$ cannot be ~~not~~ destroyed.

i.e., the voltage of \bar{Q} cannot be pulled up too much.

$$V_{\bar{Q}} = \left(\frac{V_{dd}}{R_{M5} + R_{M1}} \right) R_{M1} \approx 0.3 V_{dd}$$

$$3.15 \overline{) 1.00} \\ \underline{0.95}$$

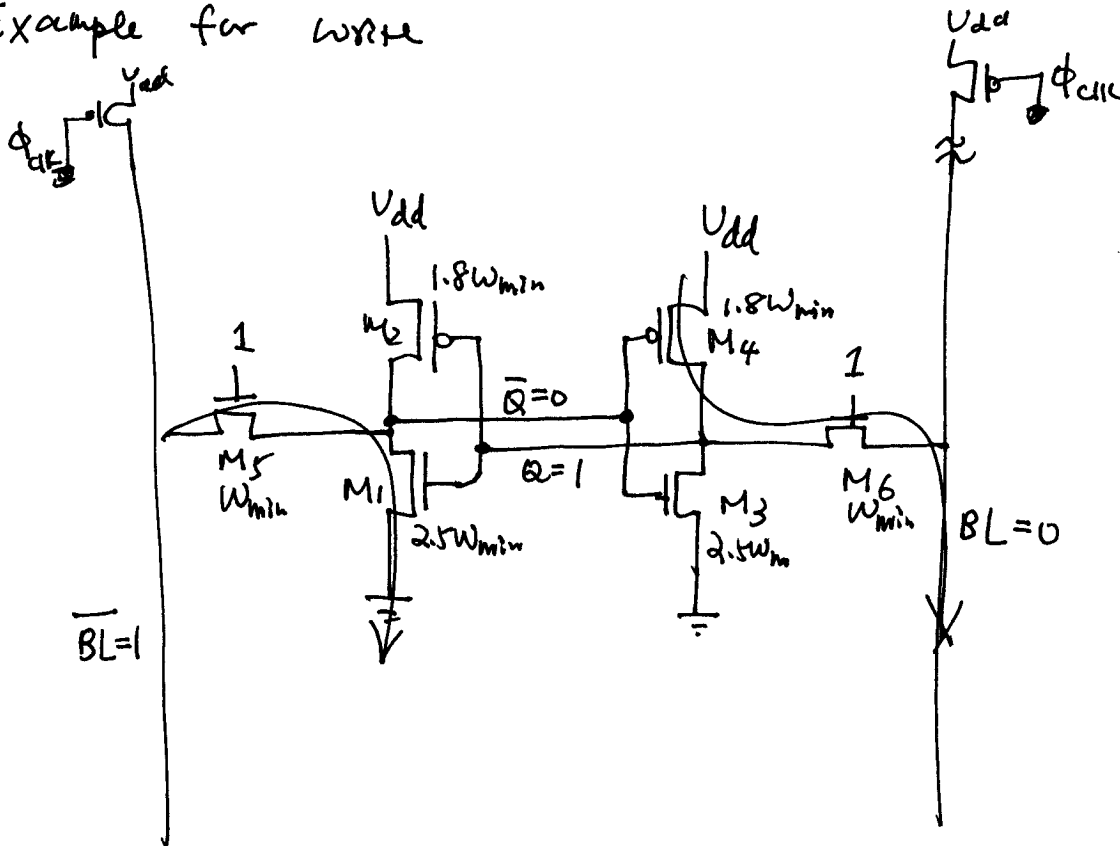
• A reasonable sizing is

$$W_5 = W_{min}, \quad W_1 = 2.5 W_{min}$$

Similarity, $W_6 = W_{min}, \quad W_3 = 2.5 W_{min}$

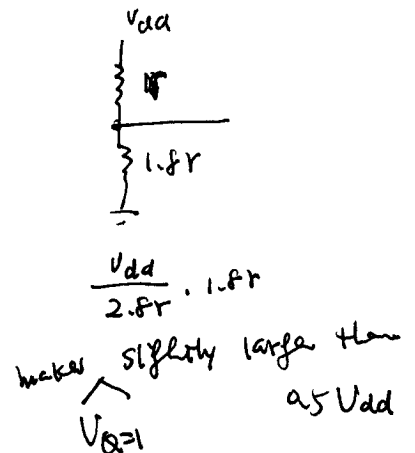
~~...~~
⊗ P on if $V_{gs} < V_{TP}$

• Example for write



• If write 1, no charge!

• So, let us try to write 0.



- \overline{BL} cannot charge the state of \overline{Q} as we discussed.
- $Q=1$ must be charged first by discharging through $M6 \rightsquigarrow BL$ line.
- $Q=1$ must drop below $V_{dd} + U_{TP0}$ such that $M2$ can be on. So, there is a sizing problem between $M4$ and $M6$.
- Once Q drops to below $V_{dd} + U_{TP0}$, $M2$ begins to be on and charge \overline{Q} to 1 and charge the state of the cell. (Get stabilized)!!
- $M4$ must be ~~large enough~~ ^{small}, such that Q can drop to below $V_{dd} + U_{TP0}$.
 small (in size)
 i.e., its R must be large !!

$$\frac{5 - 0.7}{4/3}$$

③

$$\left(\frac{W}{L}\right)_4 \leq 1.8 \left(\frac{W}{L}\right)_6$$

This makes $V_{dd} + U_{TP0}$

$$\Rightarrow W_4 = W_2 = 1.8 W_{min}$$

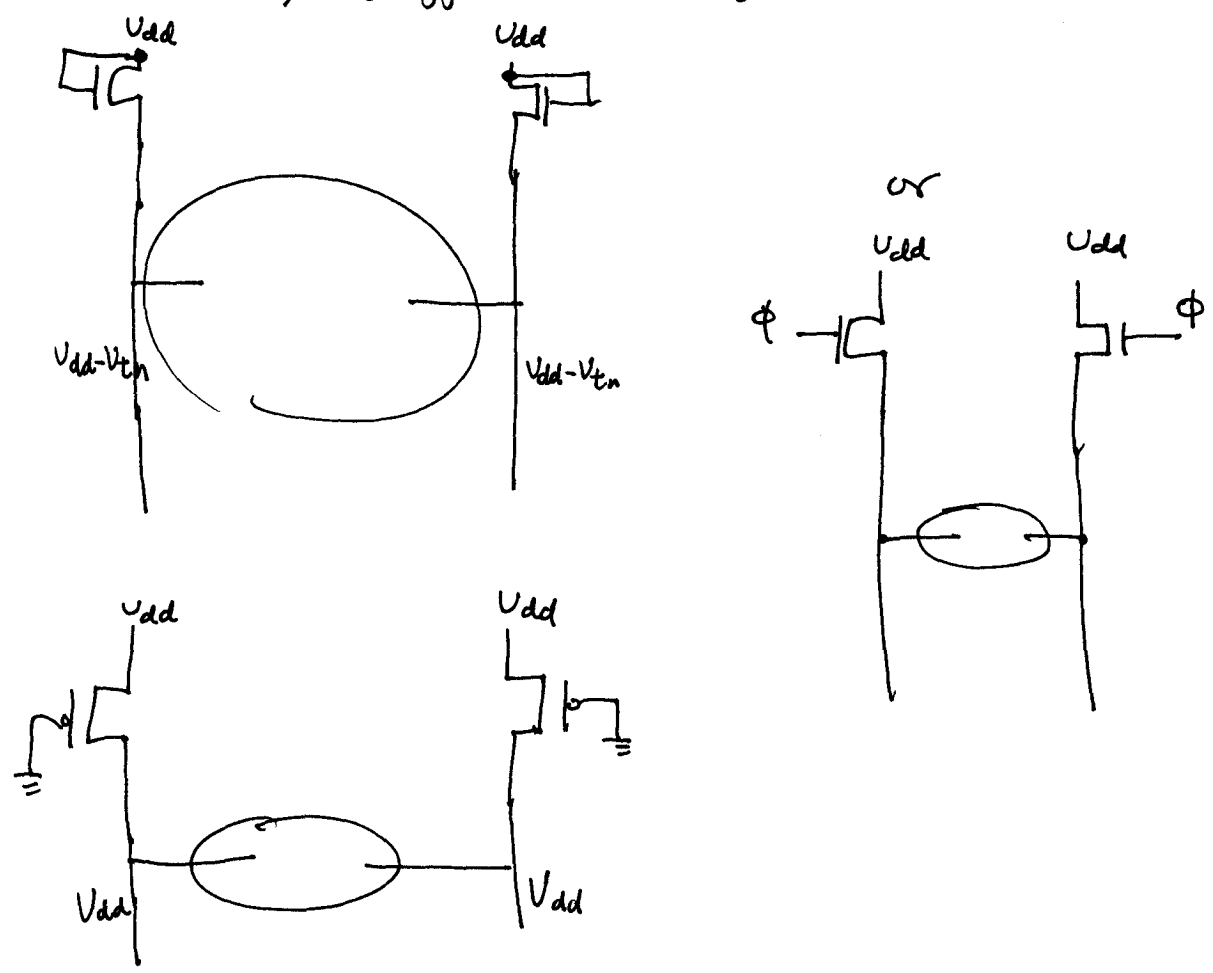
• Bit line Load.

• In fact, the bit lines can be charged to lower than V_{dd} to save power

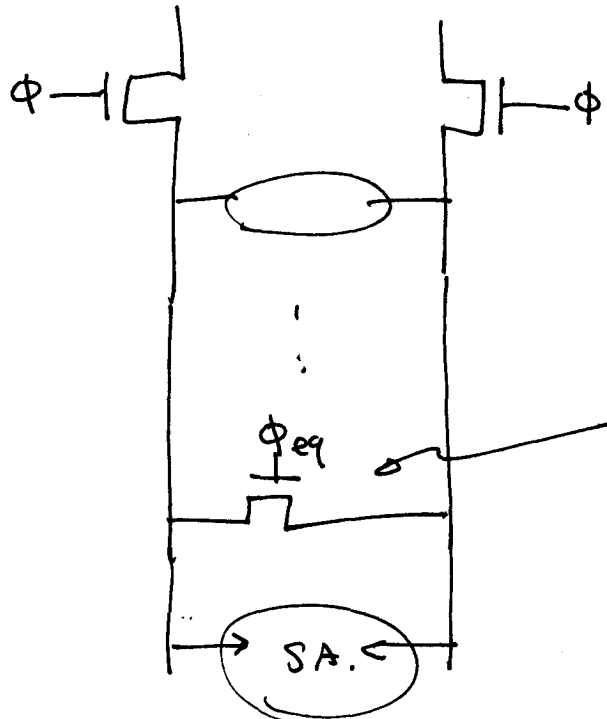
• However, if BL (in P.184) is too low, $Q=1$ cannot be maintained, and the state may be changed.

∴ Bit line voltage $\uparrow \Rightarrow$ more stable memory state.

• can have many different designs.

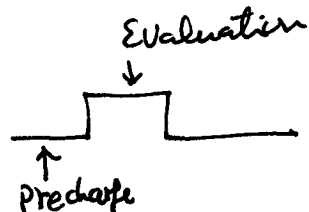
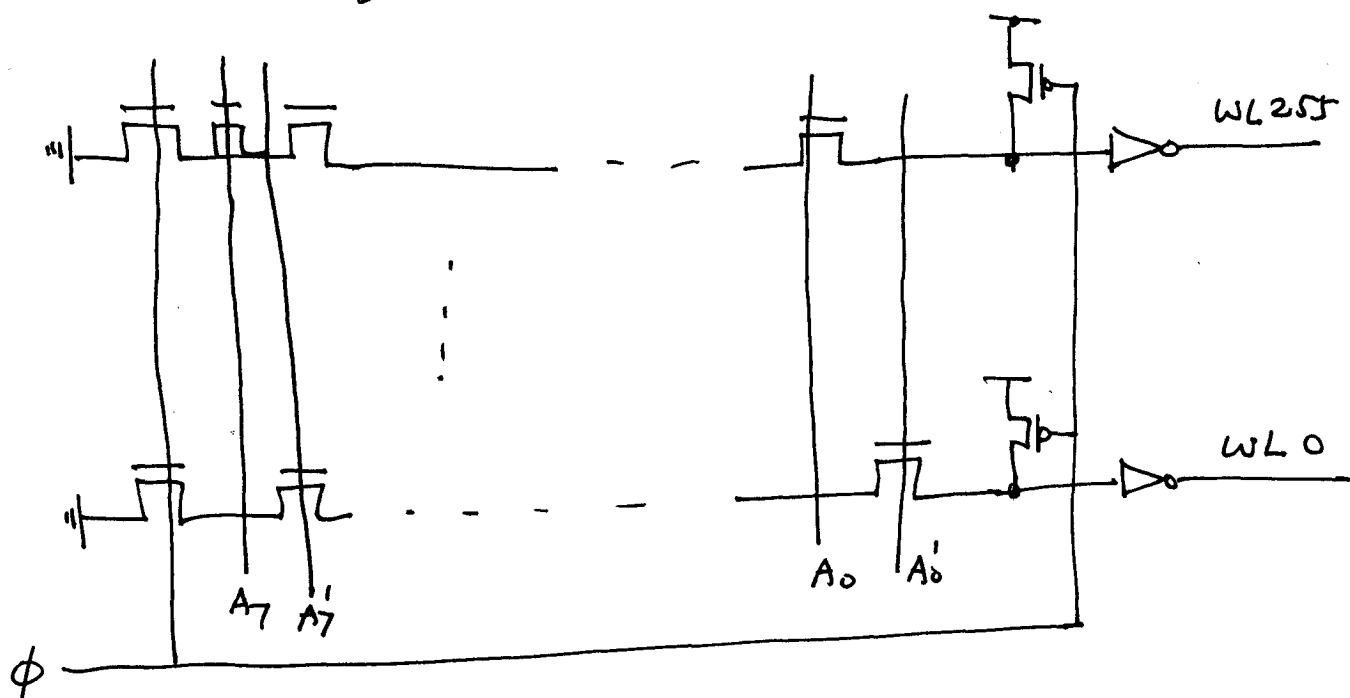


• no matter what kind of pull-up ckt, equalizing device must be added.



After each r/w operation, this must be activated for a while to equalize the voltage.

• Row Decoder $8 \rightarrow 256$ mapping (Example) - NAND-type



• Dynamic CKT

• $\phi = 0$, all outputs are zero.

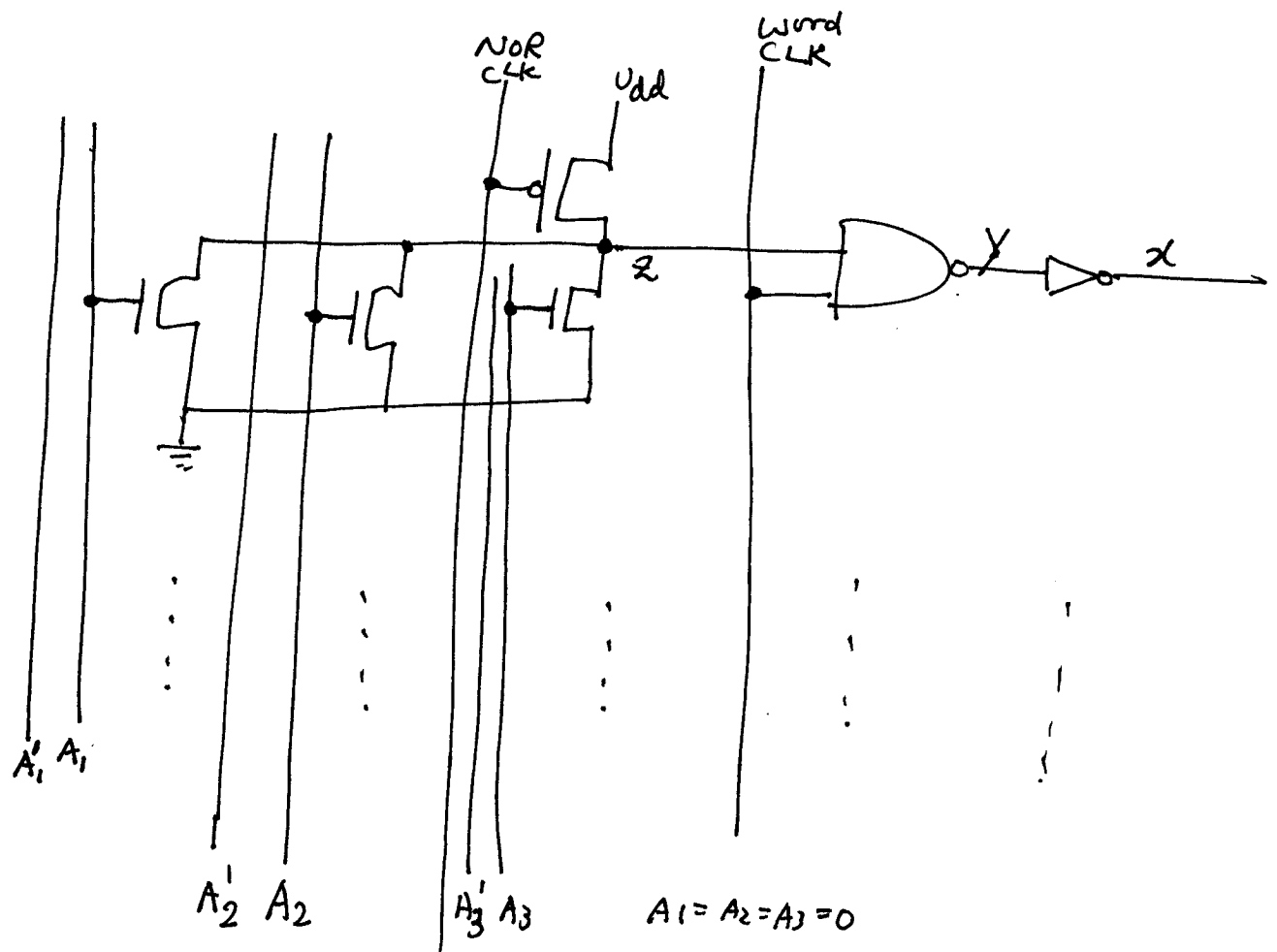
$\phi = 1$, only one word line would be 1

• Very power-saving, slower (\because larger height of evaluation NMOS stack)

• must add inverter buffer, since there is no power source to drive the high capacitive load.

• NOR-type

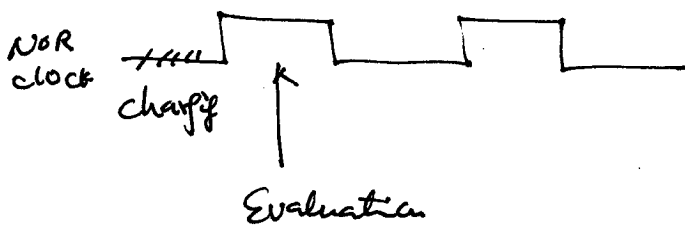
3 \rightarrow 8 mapping



$A_1 = A_2 = A_3 = 0$
 $\Rightarrow Z = 1,$
 $\Rightarrow Y = 0$
 $X = 1$

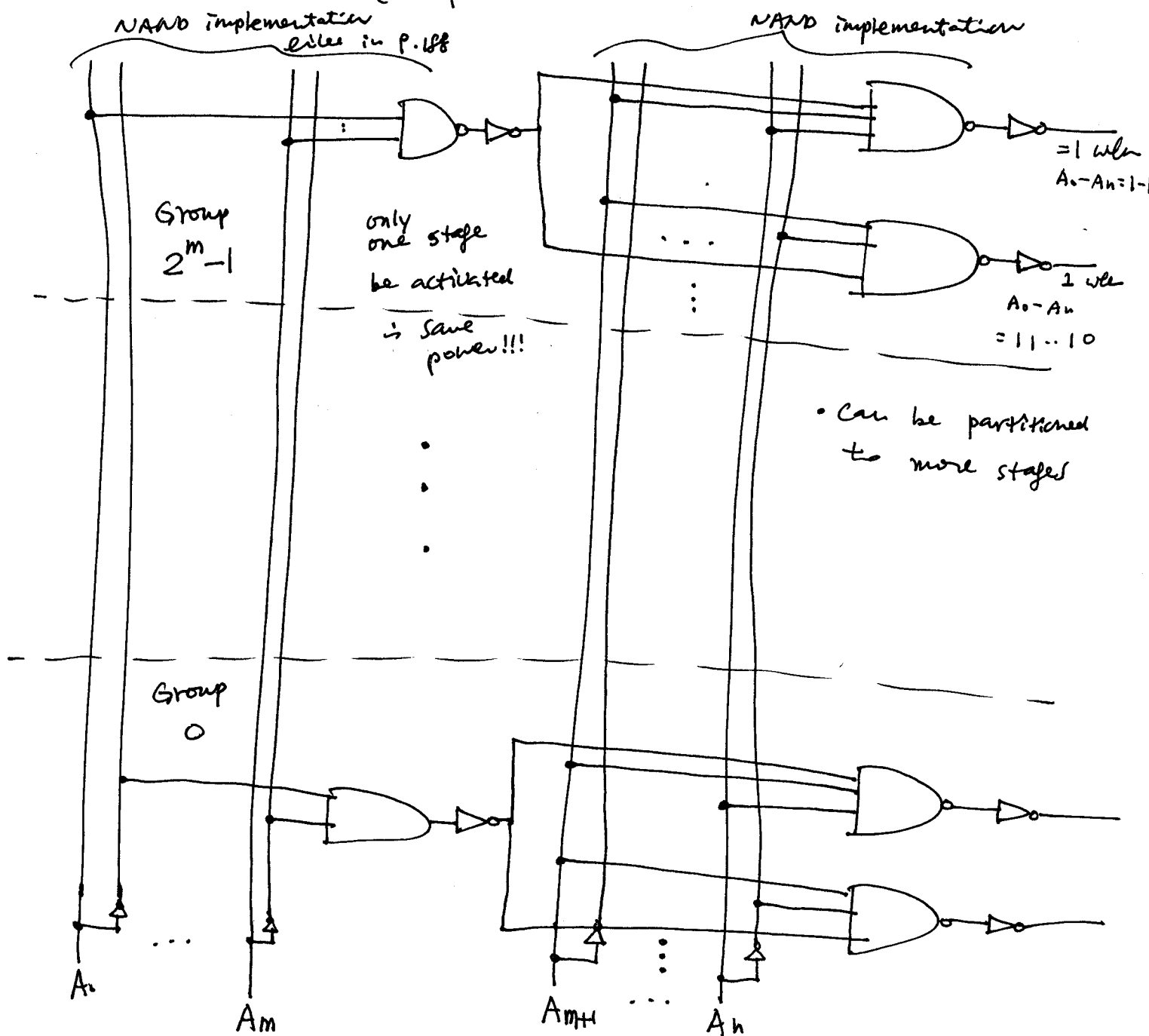
• faster, but consumes more power

∴ N-1 rows are discharging.

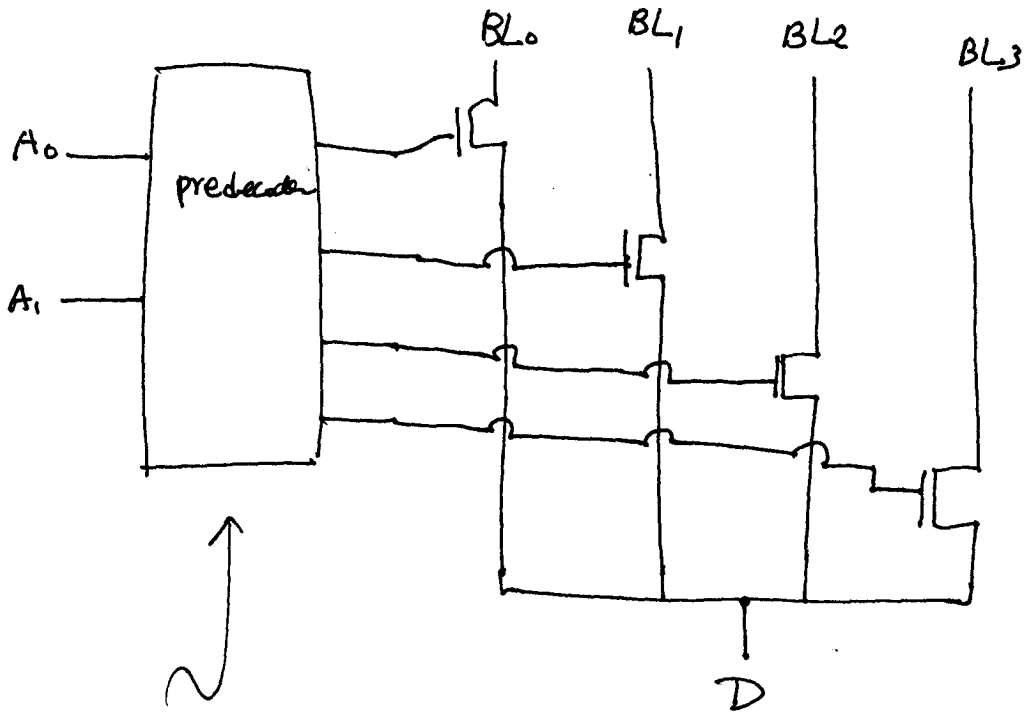


(or called predecoder)

• Multi-stage decoder (Compromise between NAND type & NOR type)



• ~~Row~~ Decoder : show by 2 → 4 mapping
Column

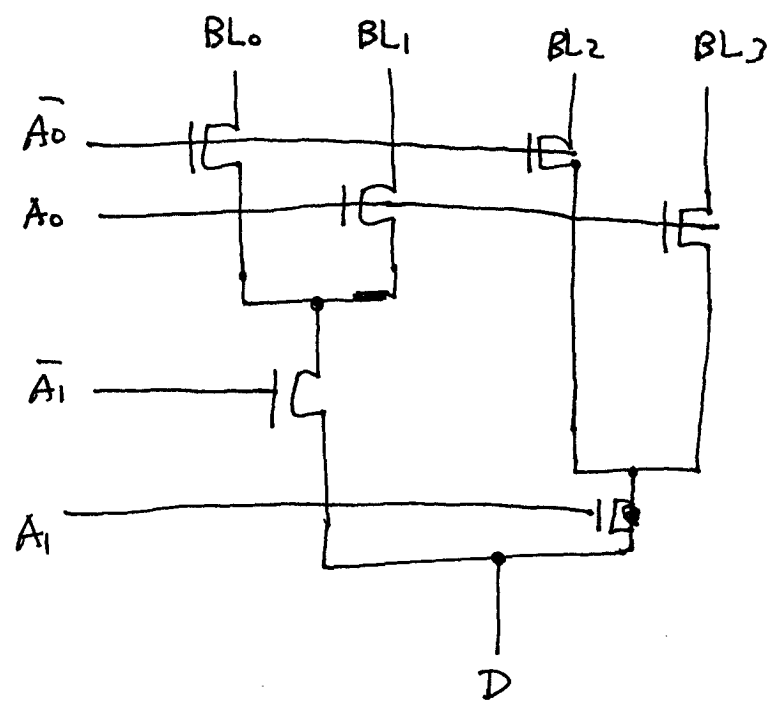


predecoder +
pass transistor

Don't need to be fast!

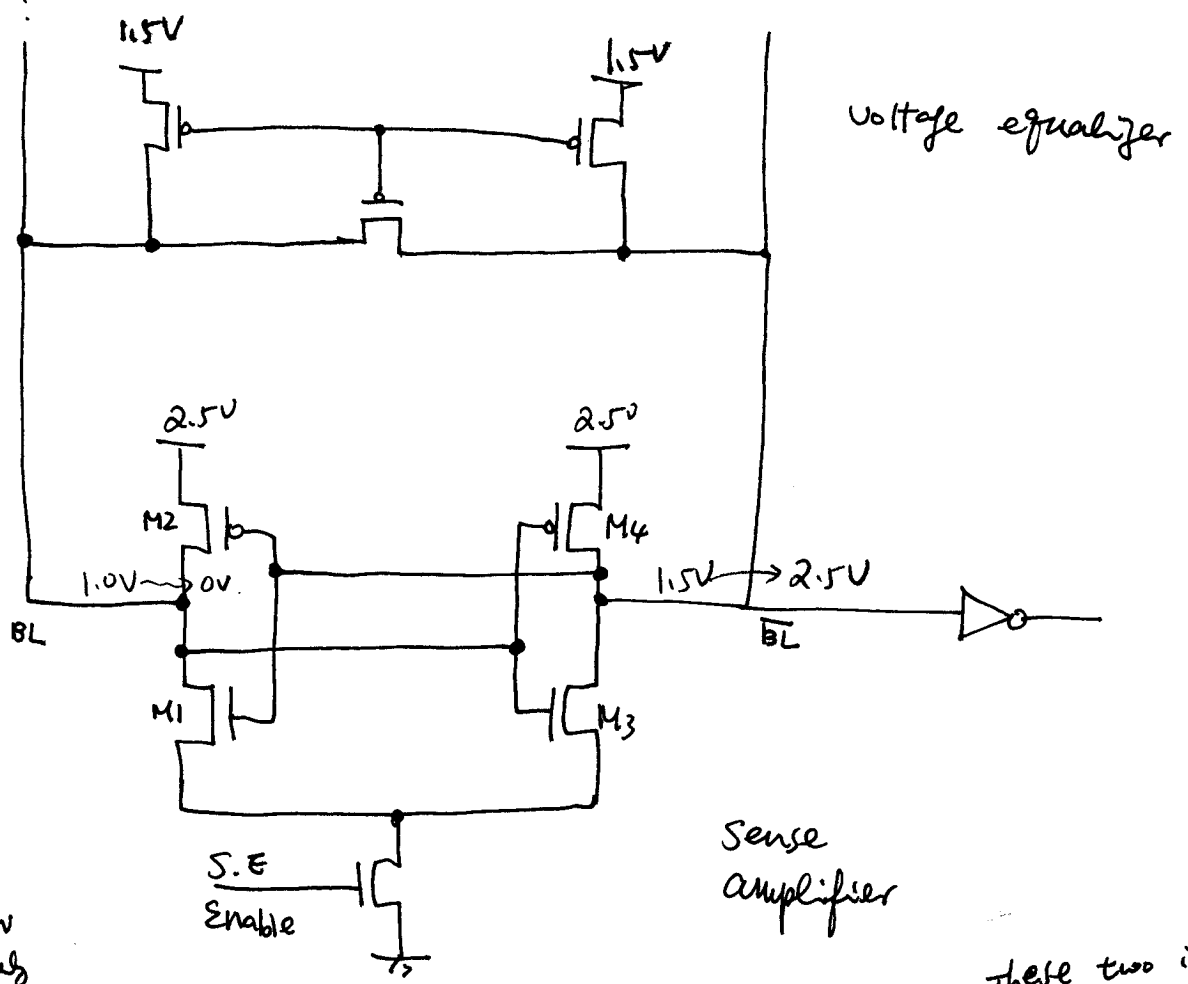
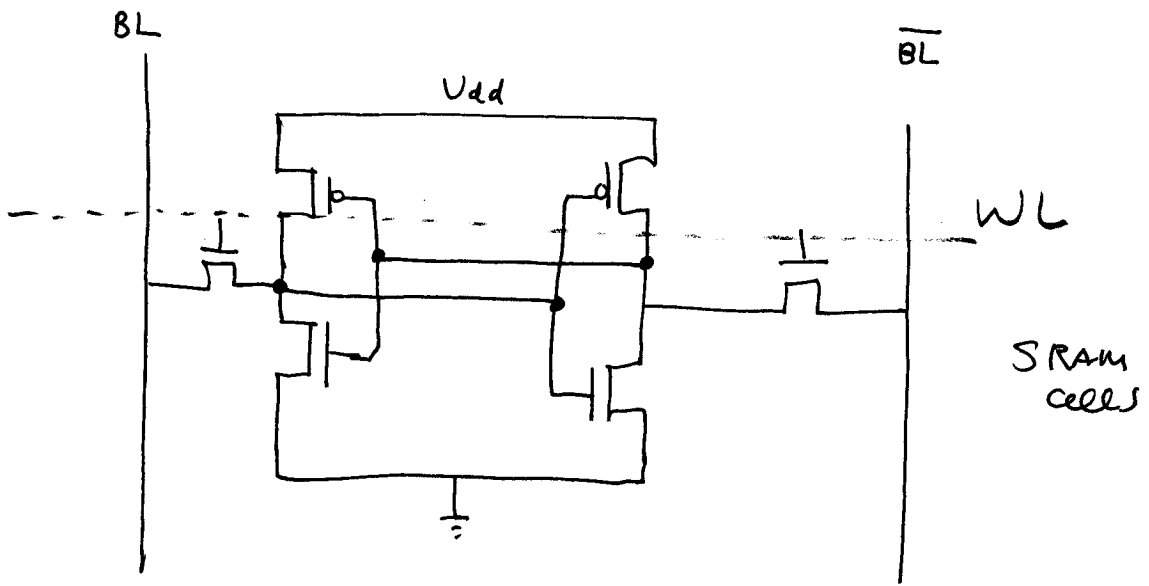
~~this is near the end of
read operation.~~

Tree Decoder



slow!
may need to pad
many transistors

• Sense Amplifier



Step 1:
BL changes 1.5V \rightarrow 1.0V

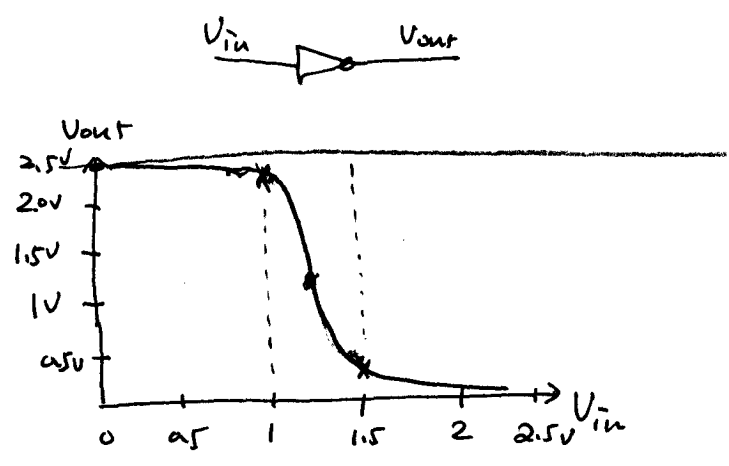
Step 2:
 \overline{BL} changes 1.5V \rightarrow 2.5V quickly

Step 3:
This quickly makes BL change from 1.0V \rightarrow 0V.

Step 4: This feedback and quickly makes \overline{BL} change 1.5V \rightarrow 2.5V.

• These two inverters push each other !!!

• Inverter is a high-gain device



- Change in input 0.5V \implies output change about 2.0V
- \therefore 4 times of gain.

• Initially, both BL & \overline{BL} are precharged to ~~2.5~~ 1.5V.

• If BL is discharged to 1.0V, \overline{BL} remains at 1.5V $\xrightarrow{\text{increase to } 2.5V}$ 2.5V
 \overline{BL} remains at 1.5V \longrightarrow BL dropped to 0.5V $\xrightarrow{\text{eventually to } 0V}$ 0V.

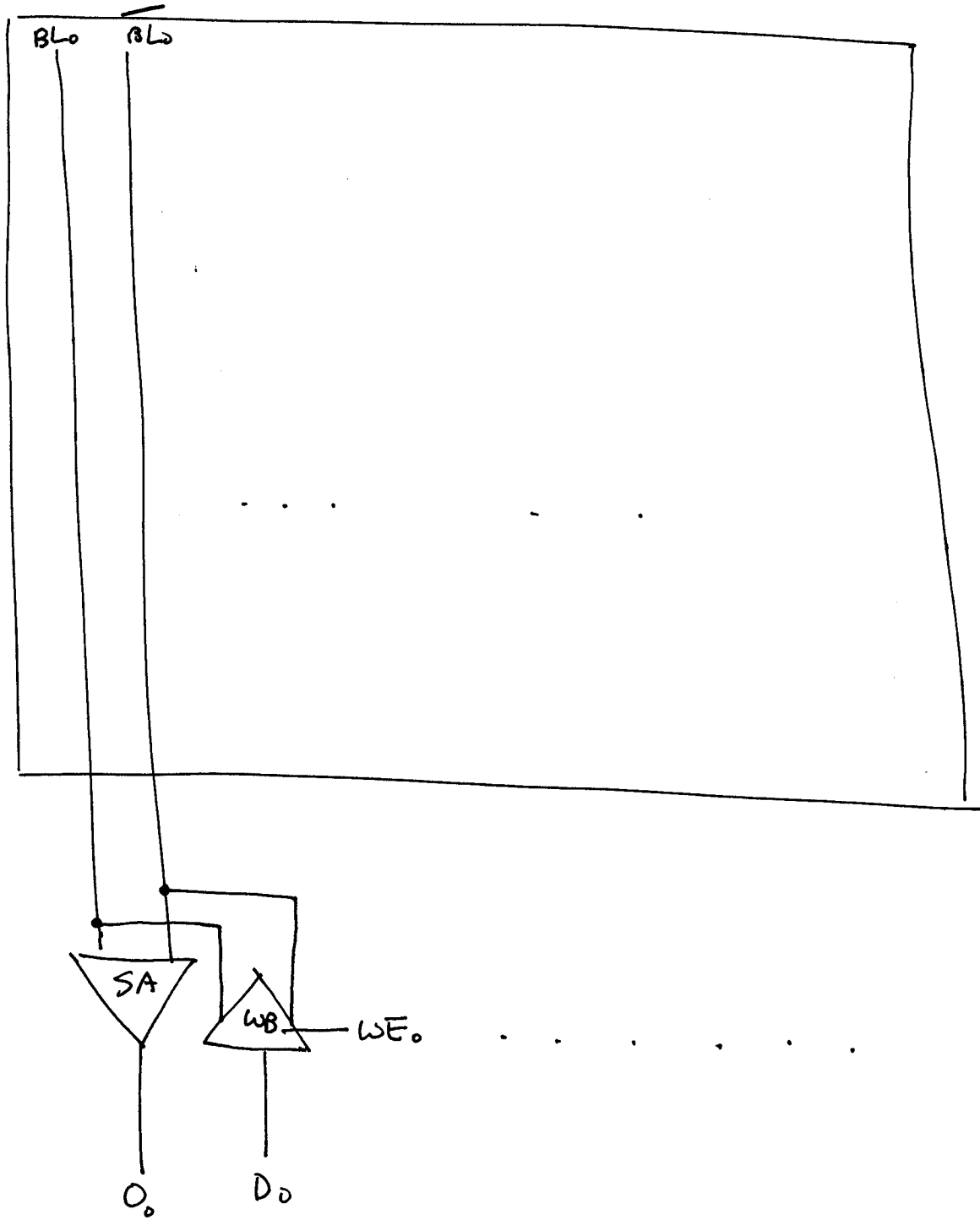
• If \overline{BL} is discharged to 1.0V, \longrightarrow BL remains at 1.5V $\xrightarrow{\text{eventually to } 2.5V}$ 2.5V
 BL remains at 1.5V \longrightarrow \overline{BL} dropped to 0.5V $\xrightarrow{\text{eventually to } 0V}$ 0V.

• This is most commonly used in real designs and consume much less power than current mirror based S.A.

• The sense amplifier enable signal is generated using a self-timed approach to reduce the power consumption.

That is $\left. \begin{array}{l} \text{BL will not charged to } 2.5V \\ \overline{BL} \text{ will not be discharged to } 0V \end{array} \right\} \text{ or vice versa}$

• Read/write operation



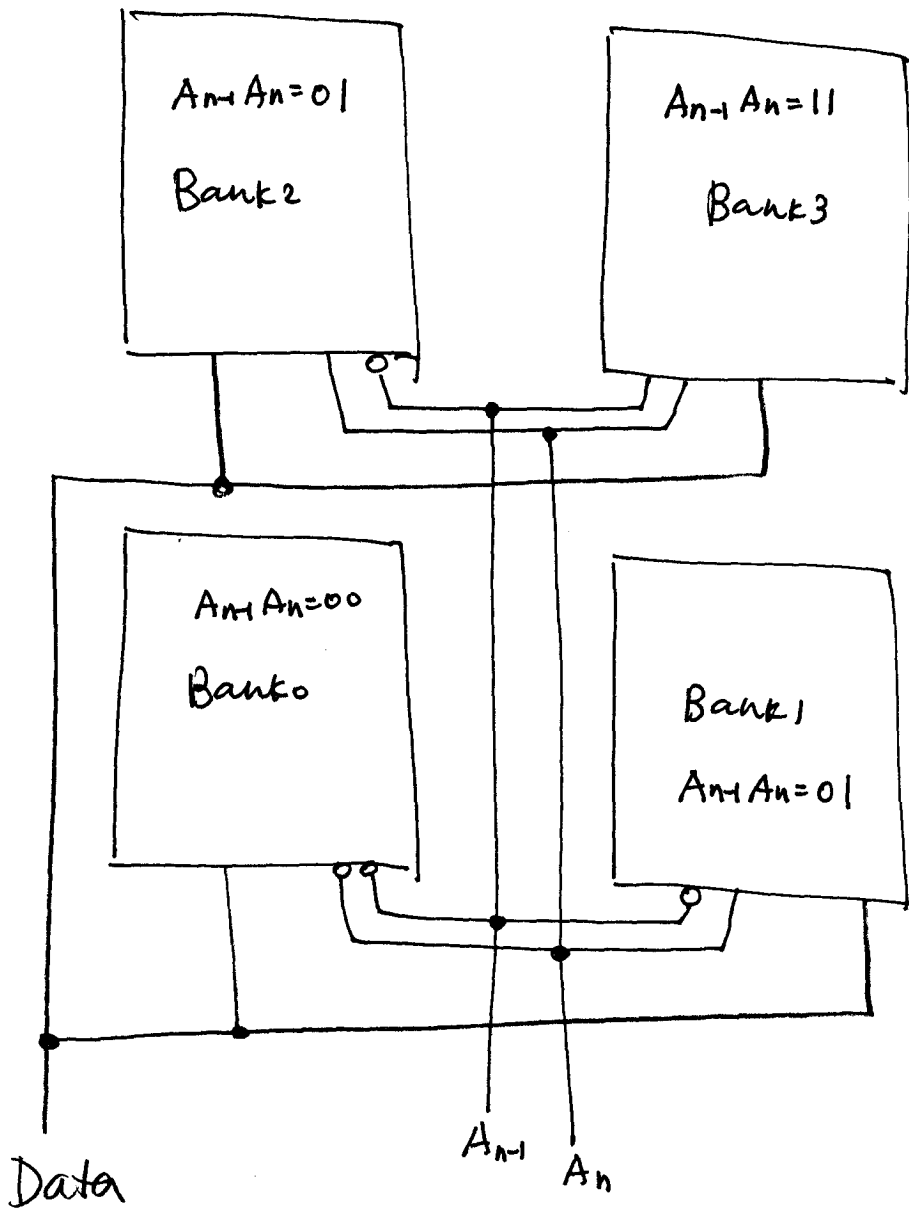
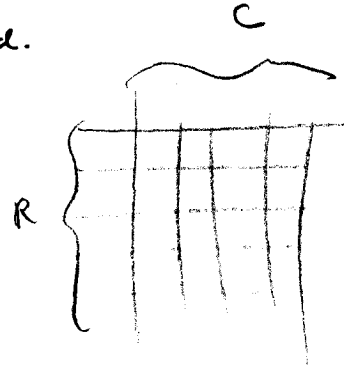
• Low-power memory design

① Memory Sub-banking

• R rows, C columns, C_{cell} : C (bit line switch/cell)

↳ totally, $R \cdot C \cdot C_{cell}$ of cap. is switched.

• Try to partition the memory

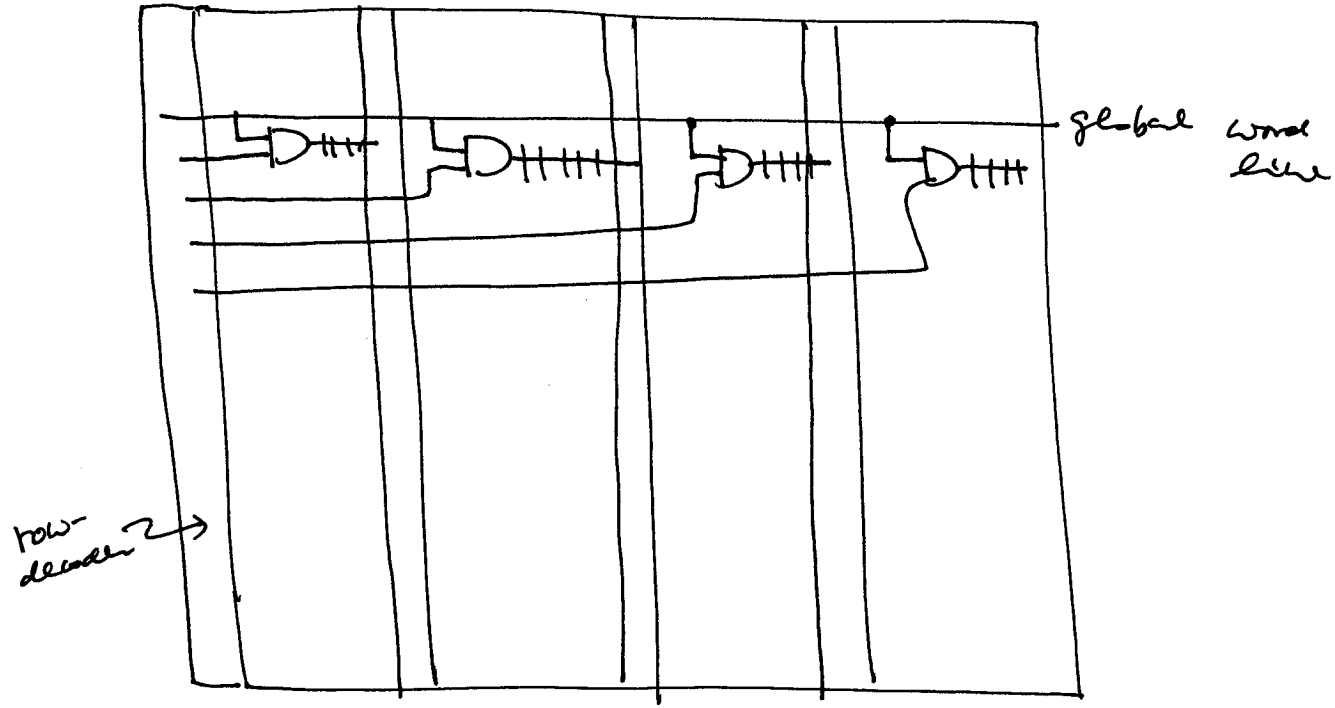


Make it B banks

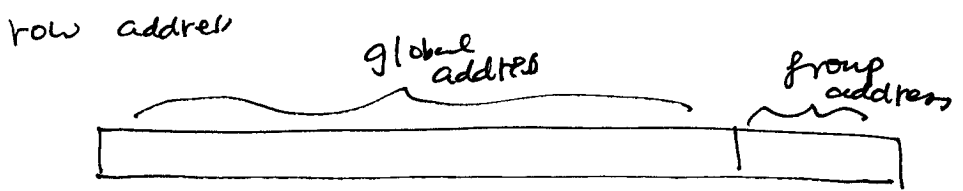
$$\Rightarrow \frac{R \cdot C \cdot C_{cell}}{B}$$

- if interconnect is not considered.
- Can find an optimal solution.

bit bit

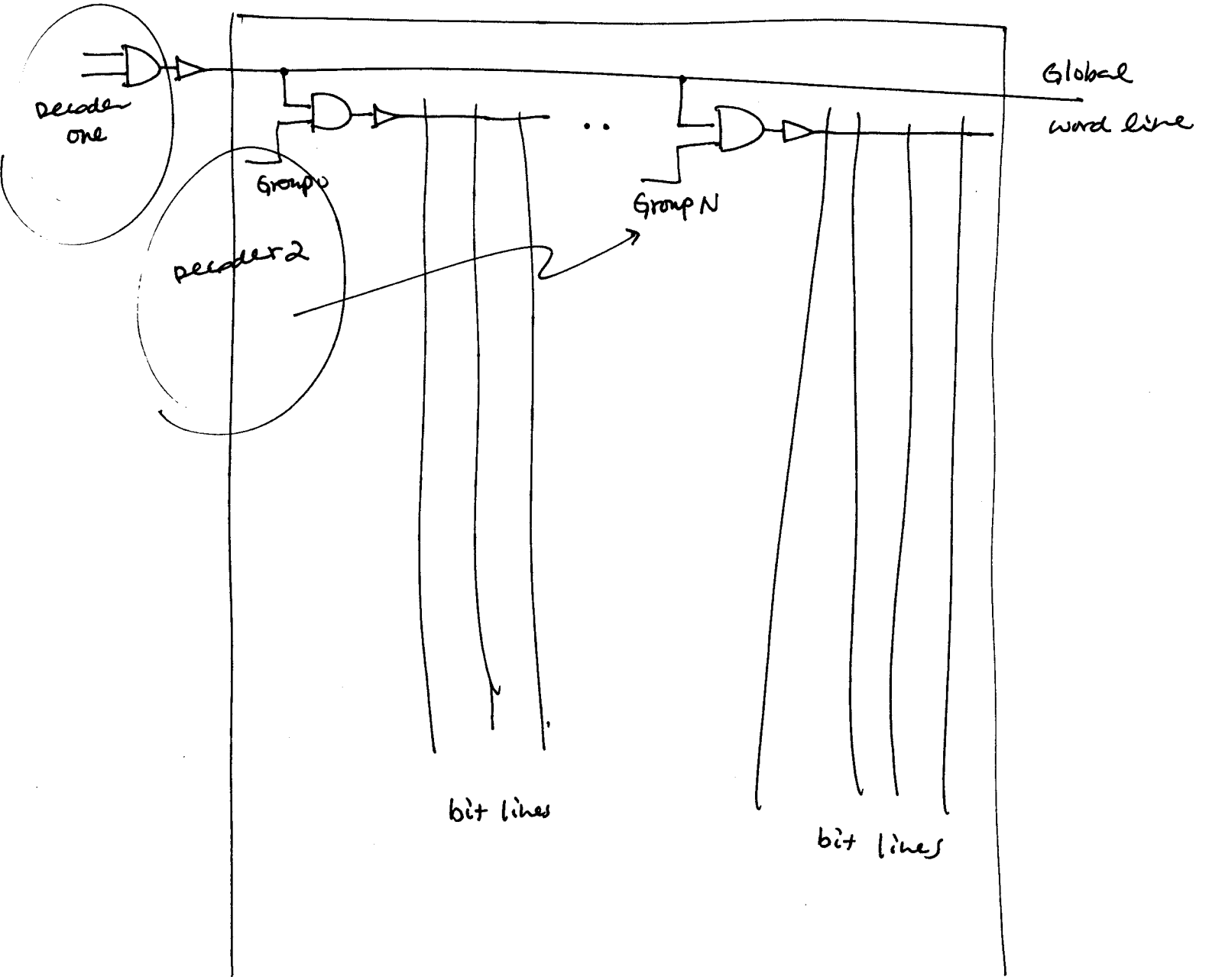


Divided word line



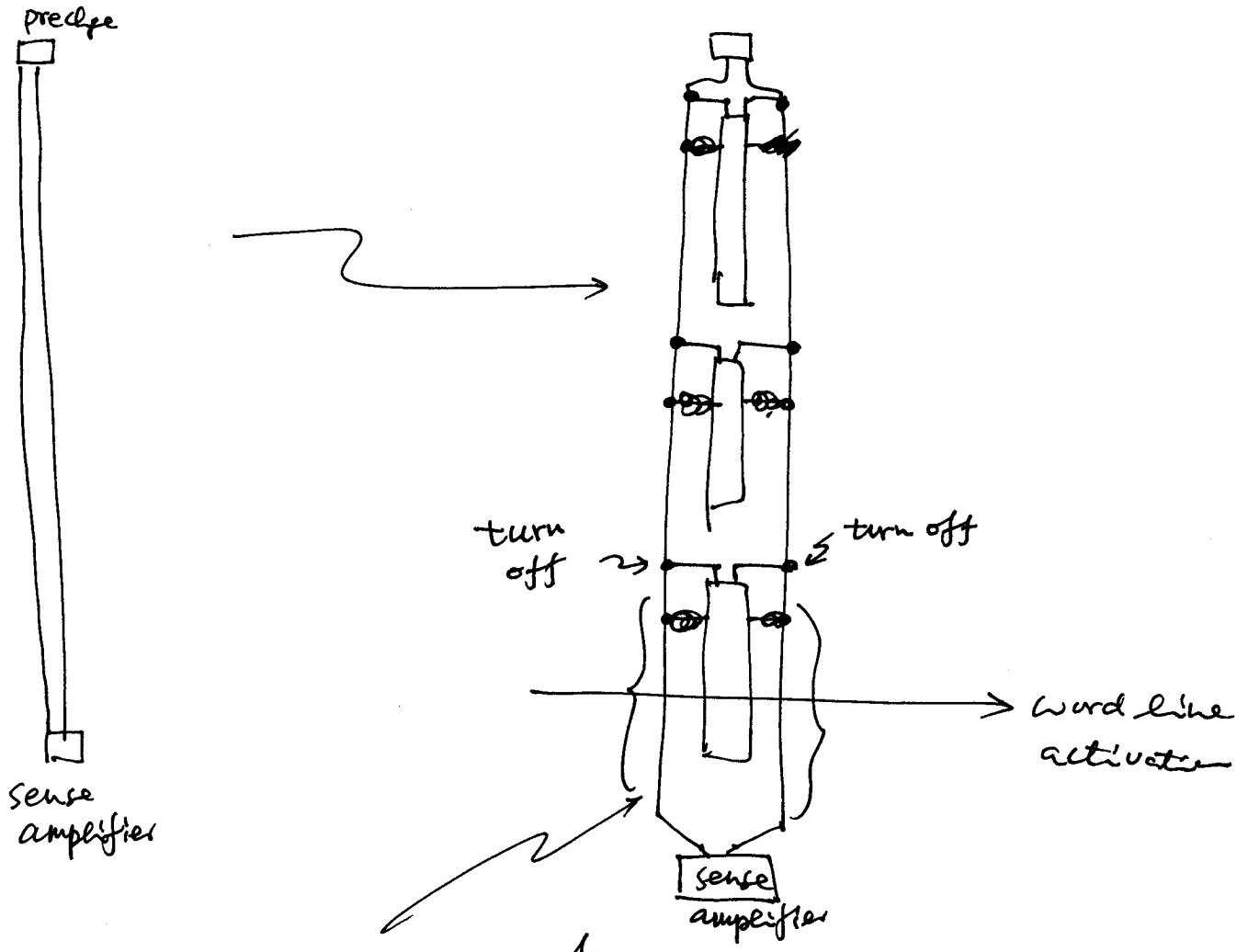
- Needs at least 2 metal lines
- local word line: polysilicon
- global word lines: metal 1
- bit lines: metal 2.

• Divided Word Line Architecture



- Reduce power consumption in Bit Lines.

• Bit line segmentation



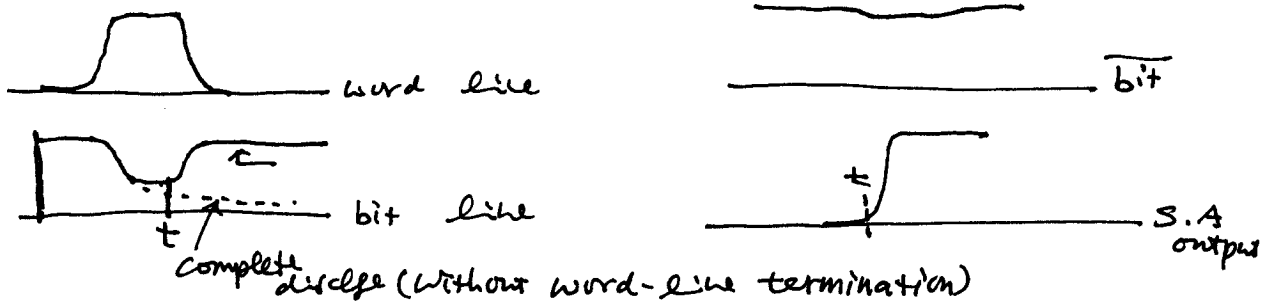
Only either of these two segments to be discharged!!

• Reduce voltage swing on Bit Lines

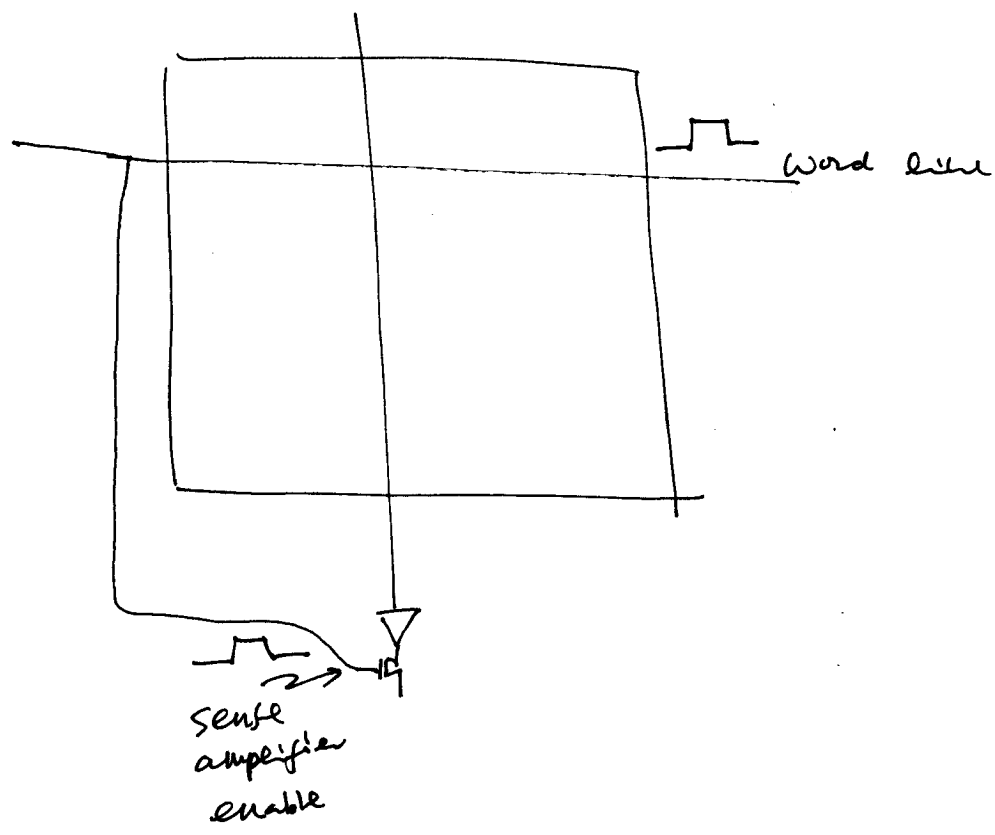
* Not quite feasible

∴ affect noise sensitivity, Complexity of sense amplifier,

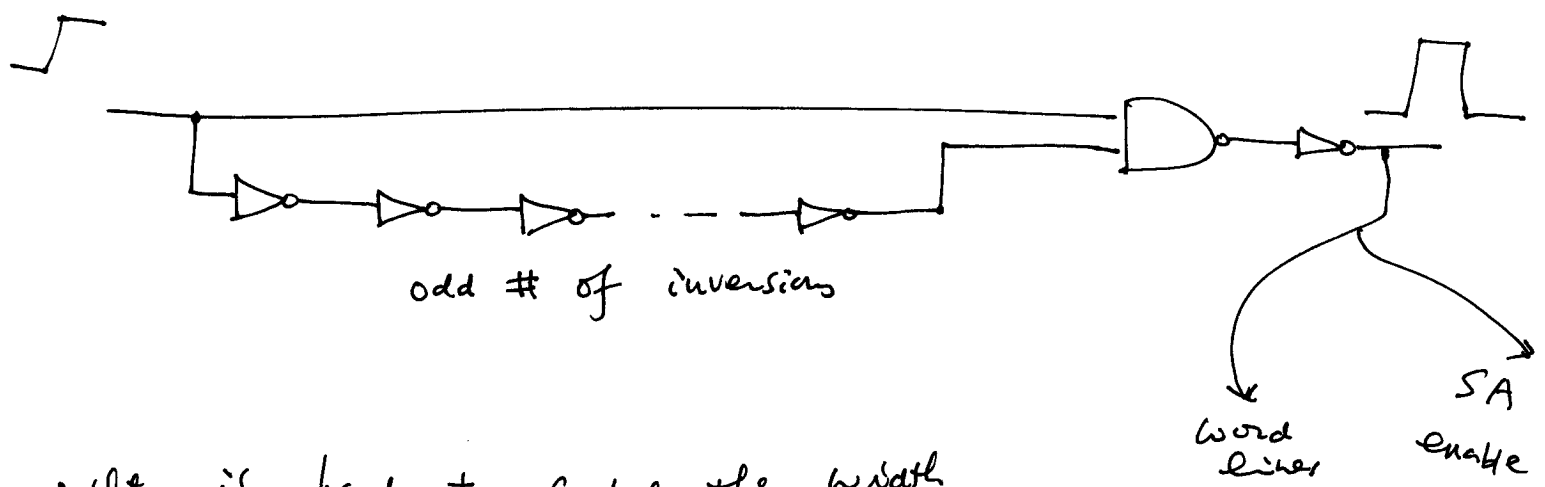
performance of the RAM Core.



Pulsed word Lines



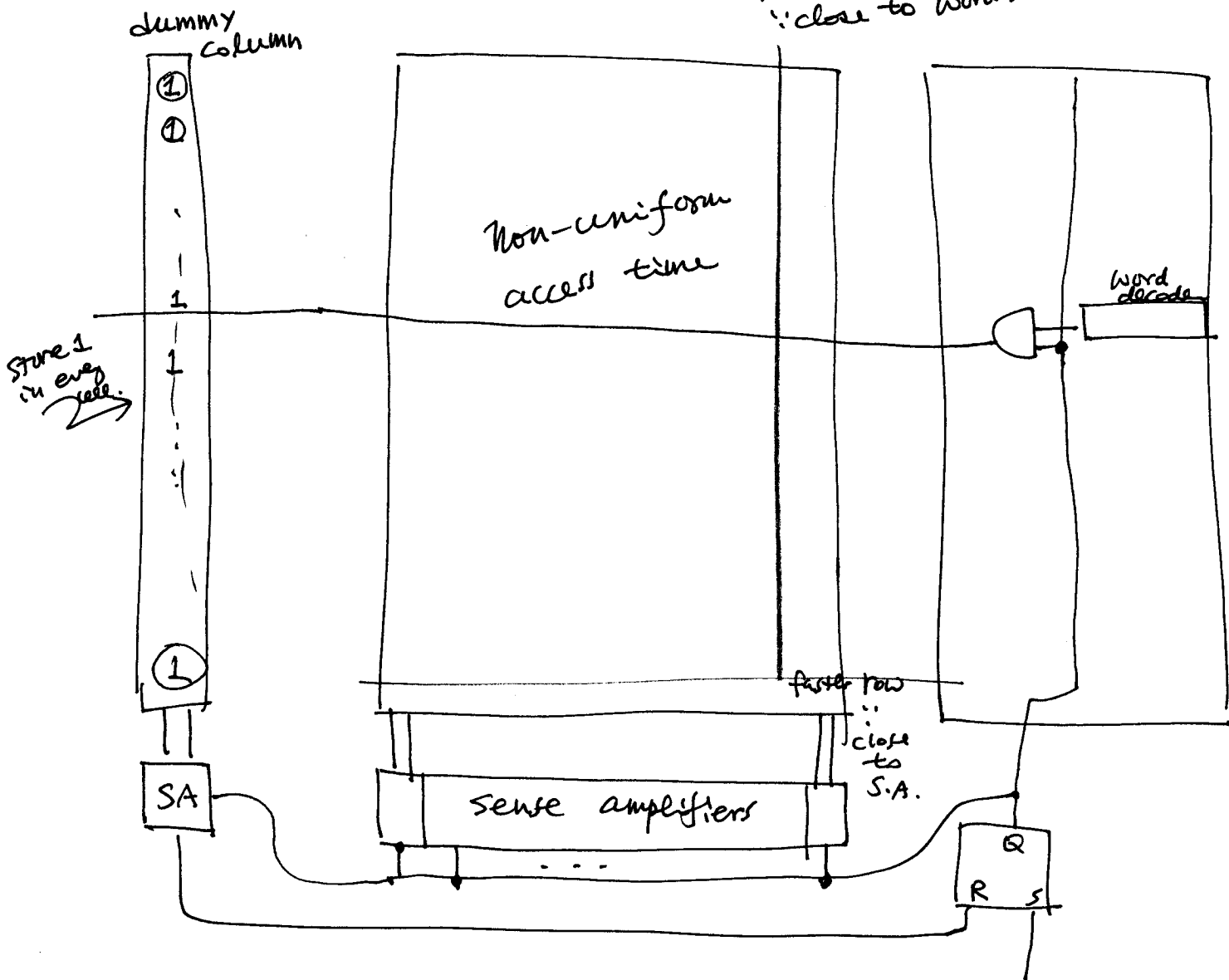
- Try to limit the bit line voltage discharge.
- Try to enable the word lines for precisely the time needed to develop the bit cell voltage discharge.



- It is hard to control the width
- try to use self-timed word lines.

- Self-timed word lines.

fast column
"close to word line driver;



① SR flip-flop is set, and the word line is triggered.
(SA)

② When the dummy SA (slowest) generate 1, the rest of the columns would have been sensed.

③ The high signal reset the flip flop and turn off the word line,

(SA)

Summary: power optimization for SRAM

(P200)

Reduce total capacitance
switched

- ① Banked organization
- ② word-line division
- ③ Multi-stage decoder
- ④ Bit line segmentation

Reduce the voltage swing
across switched capacitances

- ① Self-timed RAM

turn off word line &
SA enable, s.t. the
voltage swing in the bit
lines can be minimized.

- ② Reduce bit line voltage
swing.