

به نام خدا

عنوان :

## Speed and Power Scaling of SRAM's

نویسنده :

امیرعلی بلورچیان

کلمات کلیدی : CMOS - SRAM – VLSI - سرعت - توان مصرفی

چکیده : در این مقاله پس از معرفی ساختار کلی SRAM ها به بررسی عوامل مهم در میزان سرعت و توان مصرفی آنها می پردازیم.



معرفی :

SRAM های پربازده و با چگالی زیاد نقش مهمی در ساختار های حافظه سیستم های محاسباتی نوین ایفا میکنند. در این مقاله قصد داریم تا تاخیر و توان SRAM ها را بر حسب تکنولوژی ساخت و اندازه مقایسه کنیم. ساخت و طراحی SRAM ها نیاز به یک حالت بالانس مابین تاخیر ، توان و سطح مقطع مورد نیاز ، دارد. در واقع برای طراحی SRAM ها می بایست داد و ستدی ما بین فاکتور های مهم گفته شده صورت داد. آنالیز مدارات گسترده SRAM با استفاده از نرم افزار های شبه سازی مانند SPICE بسیار زمان بر خواهد بود لذا از آنالیز های تقریبی و خلاصه شده استفاده میکنیم. مدل های آنالیزی و تحلیلی بسیار زیادی برای فاکتور های مورد نظر وجود دارد که می توان از آنها استفاده کرد.

ما در ابتدا ساختار کلی SRAM هارا مطالعه خواهیم کرد و سپس به تاثیر فاکتور های فوق در عملکرد آنها خواهیم پرداخت.

### ساختار کلی SRAM

شکل شماره 1 ساختار تیپیکال SRAM را نشان می دهد. مسیر قابل دسترسی به SRAM ، به دو جزء قابل تقسیم می باشد : یکی دکودر می باشد که قسمتی از آدرس ورودی به word line بوده و دیگری MUX خروجی است که جزعی از سلول ها در خروجی است . در این مقاله ما بر روی زمان خواندن اطلاعات در SRAM توجه بیشتری خواهیم داشت زیرا حساسیت بسیار بالایی دارند.

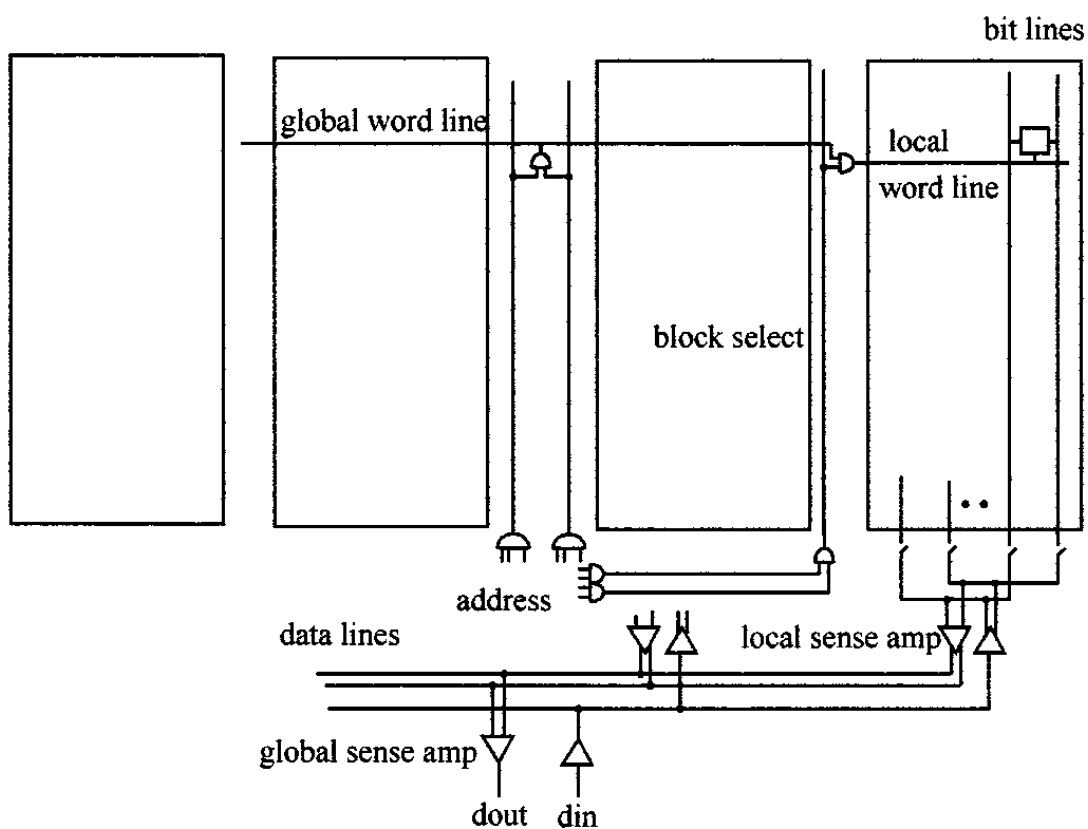
برای توانایی خواندن ورودی می بایستی کد گشایی شود تا بتواند خط word خاصی را فعال نماید. همانطور که در شکل 1 دیده میشود دیکودر معمولا خطوط word مجزایی را (تقسیم شده) به کار می برد به شکلی که قسمتی از آدرس برای فعال کردن خطوط عمودی به کار برده می شود و بقیه آدرس برای خطوط افقی به کار برده میشود. فصل مشترک این دو local word line گفته میشود. سلول هایی که به این word line متصل هستند اطلاعات خود را به bit line ها انتقال می دهند . اطلاعات از زیر بخش bit line ها از طریق MUX ستون ها به تقویت کننده ها فرستاده میشوند و از این طریق این اطلاعات تقویت شده و راهی خطوط دیتا میشوند. این سیگنال ها مجددا توسط تقویت کننده های عمومی تقویت شده و سپس از ساختار خارج میشوند . پراکندگی انرژی در SRAM ها 3 جزء کلی دارد :

1- انرژی دینامیک برای سوچ کردن ظرفیت خازنی در دکودرها، bit lines، خطوط دیتا و سایر سیگنال‌های کنترلی در ساختار.

2- انرژی آمپلیفایرهای سنس شده.

3- انرژی اتلاف شده به دلیل جریان leakage.

به طور تیپیکال ساختار کلی به چند زیر ساختار با اندازه‌های مساوی تقسیم می‌شود. هر کدام از این زیر ساختارها مقداری از accessed word را در خود ذخیره می‌کند که به آن sub word گفته می‌شود.



شکل 1 : SRAM access path

این macroها را میتوان به حالت یک RAM مجزا تصور نمود و تنها می‌بایست که قسمت دیکودر را به طور مشترک استفاده کنند.

## مدلسازی SRAM:

برای بررسی فضای طراحی SRAM های بزرگ ما چند فرض اولیه درباره ویژگی های طراحی انجام می دهیم. ما فرض های اصلی در قسمت بعدی را مشخص می کنیم . لیستی از تمام فرض ها در پیوست موجود می باشد. سپس ما مدل های تحلیلی برای تاخیر، توان و سطح مقطع مورد نیاز اجزای SRAM را گسترش داده و با شبیه سازی مداری HSPICE مقایسه می کنیم.

این مدل ها برای بررسی کارایی رنج وسیعی از SRAM ها در اندازه های متفاوت و تکنولوژی های مختلف برای دستیابی به بهترین وضعیت از لحاظ اندازه و پیکربندی مورد استفاده قرار می گیرد که در قسمت بعدی مورد بحث قرار گرفته است.

### فرضیات:

تکنولوژی اصلی که برای این تحلیل ها بکار می رود پروسه CMOS  $0.25\ \mu\text{m}$  می باشد که جزئیات مناسب در جدول I نشان داده شده است. پروسه مناسب که مستقل از واحد طول نامیده می شود برای تشریح پارامترهای هندسی در این مقاله بکار رفته است که برابر با نصف مینیمم اندازه در هر تکنولوژی می باشد. با فرض اینکه سرعت تمامی دستگاهها و ابعاد آنها رابطه خطی با سایز و اندازه ی طرح دارند. ابعاد منبع و ابعاد سیم بندی از مس و در مقیاس  $0.18\ \mu\text{m}$  به بعد می باشد. ویژگی اصلی که برای تحلیل چهار حالت مختلف در جدول 2 نشان داده شده است.

TABLE I  
FEATURES OF THE BASE  $0.25\text{-}\mu\text{m}$  TECHNOLOGY

Parameter	Value	Comments
$C_g$	$1.29\ \text{fF}/\mu\text{m}$	Unit gate capacitance of an inverter (p-width = $2 \times$ n-width)
$C_j$	$1.72\ \text{fF}/\mu\text{m}$	Unit parasitic junction capacitance of an inverter
$R_g$	$4.4\ \text{K}\Omega\text{-}\mu\text{m}$	Unit output Resistance of an inverter
$\tau_{fo4}$	$90\text{pS}@2.5\text{V}$	Delay of an inverter driving a load four times its size
M1 width	$0.45\ \mu\text{m}$	Minimum width of metal 1
M1 aspect ratio	1.5	Aspect ratio (height/width) of metal 1 cross section.
M1 cap	$0.21\ \text{fF}/\mu\text{m}$	Capacitance of metal 1
M1 res	$0.1\ \Omega/\mu\text{m}$	Resistance of metal 1, (aluminum)

TABLE II  
TECHNOLOGY SCALING OF SOME PARAMETERS

L ( $\mu\text{m}$ )	C <sub>g</sub> (fF/ $\mu\text{m}$ )	$\tau_{\text{fo4}}$ (pS)	Supply (V)	M1 min. Width ( $\mu\text{m}$ )	M1 aspect ratio	M1 cap (fF/ $\mu\text{m}$ )	M1 res, type ( $\Omega/\mu\text{m}$ )
0.25	1.29	90	2.5	0.45	1.5	0.21	0.1, al
0.18	1.29	65	1.8	0.324	1.77	0.23	0.12, cu
0.10	1.29	36	1.2	0.18	2.37	0.29	0.29, cu
0.07	1.29	25	0.9	0.126	2.83	0.33	0.49, cu

برای حذف خطا، 50 mV ضریب تطابق برای thresholds ورودی های تفاضلی در تقویت کننده حسی صرف نظر از پروسه تولید در نظر می گیریم.

تاخیر و انرژی هسته RAM را مدل کرده و از بار خارجی به علت داشتن امپدانس ثابت داخلی صرف نظر می کنیم. با فرض اینکه برای تمامی گیت های RAM طرح مداری استاتیک برای مشخص کردن کاربردهای آن بکار می رود. تاخیر در گیت pMOS به اندازه ی تاخیر در گیت nMOS می باشد بنابراین می توان گیت را به عنوان پارامتر تک سایز مورد بررسی قرار داد. از آنجا که SRAM های پر سرعت، سایز ترانزیستور را در گیت دکودر برای بهبود مسیر بحرانی تغییری دهند لذا ما تاثیر آنرا بر روی آنالیزهای تاخیر مشخص خواهیم کرد. تنوع زیادی در مدارهایی که برای تقویت کننده حسی بکار می رود وجود دارد. در این مقاله ما از تقویت کننده طرح latch شامل یک جفت گین cross-coupled که با کلاک فعال می شوند، استفاده می کنیم. در این ساختار تاخیر تقویت کننده متناسب با لگاریتم ولتاژ مورد نیاز می باشد.

اگر زمان بندی کلاک نمونه برداری به خوبی کنترل شود سرعت اجرا افزایش خواهد یافت. آنها توان کمی مصرف می کنند ما فرض می کنیم که زمان نمونه برداری به خوبی کنترل می شود اما تاثیر ایده آل نبودن کلاک نمونه برداری و ایجاد تاخیر را مشخص می کنیم. زمانیکه تعداد سطوح سیم بندی شده محدود باشد فضا مابین بلوک ها برای خطوط دیتا بکار می رود. که این ناحیه به فضای حافظه اختصاص می یابد. بویژه زمانیکه سایز بلوک کوچک باشد. با توجه به اینکه تعداد سطوح سیم بندی در حال افزایش می باشد لذا در این مقاله فرض بر این است که در صورت زیاد بودن خطوط، خطوط دیتا به صورت عمودی نیز می توانند انتقال یابند. بنابراین فضای خطوط اضافی برای گذرگاه های افقی تنها در انتهای صف مورد نیاز می باشد.

تغییر ابعاد اندازه ترانزیستور روش دیگری برای سبک و سنگین کردن تاخیر، سطح مقطع و توان می باشد. در این مقاله فرض بر این است که اندازه گیت در مسیر انتقال به گونه ای است که باعث ایجاد کمترین تاخیر می شود. بنابراین fanout هرکدام از گیت های منطقی به گونه ای انتخاب شده است که با تاخیر در گنجایش خروجی (fanout) چهار اینورتر بار شده برابر باشد. البته با این فرض، تاخیر کاهش نمی یابد ولی باعث کاهش

انرژی و کاهش سطح مقطع به میزان مطلوب می شود. یک مدل RC ساده که برای گیت‌های منطقی بکار می رود. بنابراین اندازه گیت به گونه ای می باشد که کاهش تاخیر و یا افزایش تاخیر یکسان باشد پارامتر تک سایز  $w$  که سایز ترانزیستور nMOS خواهد بود. اگر  $C_g$ ، خازن ورودی برای پهنای واحد و  $R_g$  مقاومت ورودی برای پهنای واحد اینورتر باشد مقاومت خروجی گیت برابر  $R_o = R_g/w$  و خازن ورودی برابر  $C_{gx} = 3w \times l_c$  می باشد. مقدار تاخیر گیت منطقی  $w$ ، بار  $C_L$  را از طریق مقاومت  $R_w$  و ظرفیت خازنی  $C_w$  درایو می کند. (شکل 4) که در رابطه 1 بر آورد شده است.

$$D = \underbrace{R_o \cdot (C_L + C_w)}_{\text{gate delay}} + p + \underbrace{R_w \cdot \left(\frac{C_w}{2} + C_L\right)}_{\text{wire delay}} \quad (1)$$

تاخیر گیت به علت ظرفیت خازنی نقطه اتصال درین می باشد. در طراحی SRAM پر بازده، توان دینامیکی بر اتلاف توان کل چیره شده و بنابراین انرژی دینامیکی مورد نیاز برای سوئیچ ظرفیت خازنی در دکودر و mux خروجی را مدل می کنیم سپس ما به بررسی جزئیات در دکودر و mux خروجی می پردازیم.

## دکودر:

دکودر از دو جز تشکیل شده است :

دکودر سطر که خطوط کلمه ها را فعال می کند و دکودر ستون که سوئیچ ها را در خطوط بیت و mux خطوط دیتا مرتب می کند. از آنجائیکه دکودر سطر در مسیر دسترسی RAM قرار می گیرد لذا ما تاخیر آنرا مدل می کنیم در حینی که انرژی دکودر سطر و ستون را مدل می کنیم مسیر بحرانی دکودر با رشته ای از سه گیت منطقی که هر کدام شامل گیت NAND و اینورتر می باشد. مسیر دکودر با اینورتر در ورودی با مینیمم سایز  $16\lambda$  برای pMOS و  $8\lambda$  برای pMOS درایو شده است.

تاخیر هر مرحله از مسیر دکودر توسط فرمولی که در رابطه 2 نشان داده شده است محاسبه می شود. هر مرحله تاخیر برای حداکثر گنجایش خروجی 4 اینورتر به منظور کاهش تاخیر اندازه گیری شده است زمانی که مقاومت سیم predecode و خطوط کلمه قابل صرف نظر نباشد (در شکل 5) بافرهای اضافی برای کاهش تاثیر بارهای گیت ورودی مورد نیاز خواهد بود. تعداد مطلوب بافر با محاسبه تاخیر دکودر با تعداد بافرهای مختلف در این دو قسمت قابل محاسبه می باشد.

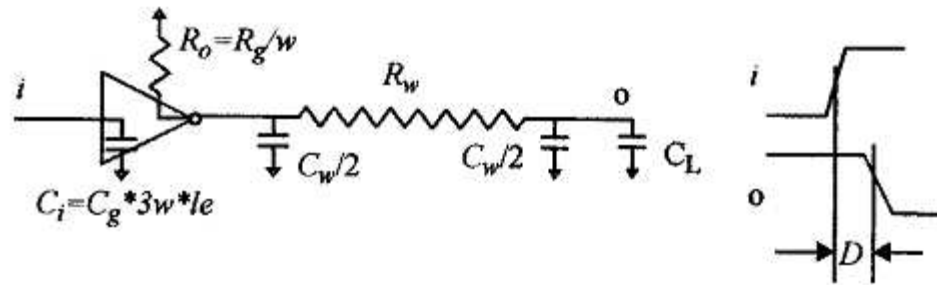


Fig. 4. Delay of a logic gate driving a load through an RC line.

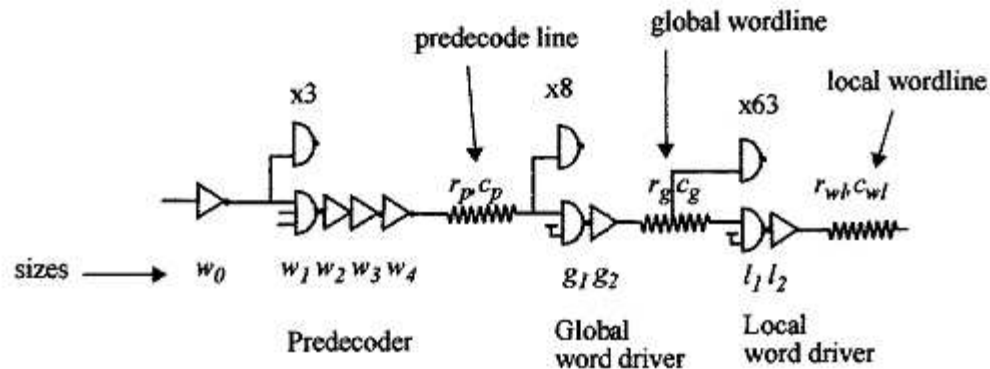


Fig. 5. Model of the critical path of a 12-4096 row decoder.

تاخیر دکودر با حداکثر گنجایش 4 سایز در قسمت 2 بیان شده است .

$$D_{\text{decoder}} = N \cdot t_{foA}x + \sum_{i=0}^N p_i + D_{\text{wire}} \quad (2)$$

دراپورهای کلمه به سیم های predecode متصل شده است. وجود بار به علت تمامی ورودی هایی که به خط کلمه سراسری وصل شده است ، می باشد. در یک SRAM واقعی، دراپورهای کلمه محلی به صورت یکنواخت توزیع شده است .

برای کاهش تاخیر سیم، خط کلمه سراسری را که از مرکز سیم دراپو شده است که دو بخش به صورت موازی دراپو شده بر

روی آن تاثیر می گذارد که هر بخش دارای مقاومت \$r\_g/2\$ و ظرفیت خازنی برابر \$(C\_g+I)/2\$ می باشد. تاخیر در سیم خط کلمه سراسری کلی برابر \$0.5x(C\_g+I)/2x r\_g/2\$ می باشد.

تاخیر سیم خالص برابر با:

$$D_{\text{wire}} = r_p \cdot (C_p/2 + G) + r_g \cdot (C_g + L)/8 + r_{wl} \cdot C_{wl}/8 \quad (3)$$

تاخیر برای دکودر سطر در چهار SRAM مختلف 9٪ تاخیر شبیه سازی در HSPICE می باشد. (شکل 6)

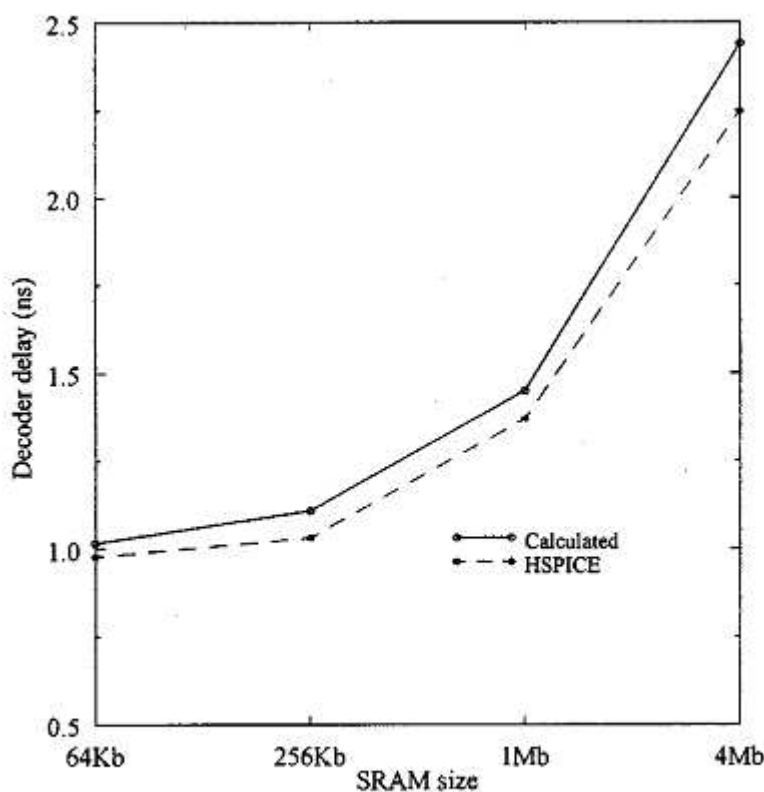


Fig. 6. Comparison of estimated and HSPICE simulated delay for row decoders.

تاخیر خط بیتی به زمان خیز خط کلمه محلی بستگی دارد، لذا  $\text{edge rate}$  در انتهای خط کلمه را محاسبه می کنیم. از شبیه سازی مداری زمان خیز برابر با 1.4 برابر تاخیر مرحله آخر درایور کلمه بدست آمده که به صورت خلاصه در 4 قابل ملاحظه است:

$$\text{rise time} = (\tau_{fo4} + r_{wl} \cdot C_{wl}/8) \times 1.4 \quad (4)$$

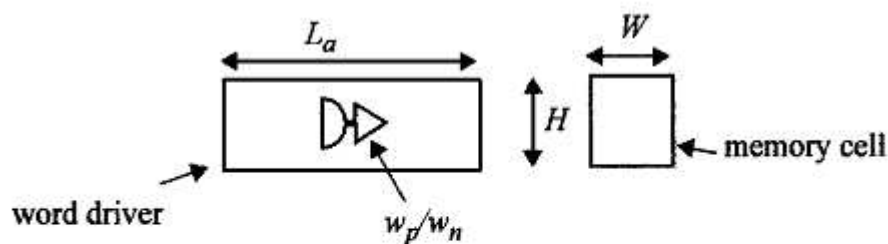
ظرفیت خازنی گیت و سیم در مسیر سیگنال به منظور محاسبه انرژی دکودر اضافه شده است. سطح دکودر با استفاده از

محاسبه سطح درایورهای کلمه محلی و سراسری و سطح مورد نیاز برای سیم های predecode محاسبه می شود. سطح

درایورهای کلمه به عنوان تابع خطی از پهنای کلی دستگاه داخل درایور مدل می شود. (شکل 7) مقدار ثابت برای این تابع



از آرایش شش درایور کلمه مختلف به دست می آید که برابر  $\lambda^2$  نمی باشد که  $\lambda$  نصف مینیمم سائز طرح تکنولوژی می باشد.



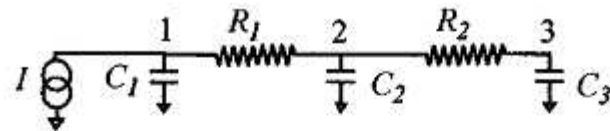
$$L_a * H = 24.05 * 1.25 * (w_p + w_n) + 497 \quad (\text{in } \lambda^2)$$

Fig. 7. Area estimation for the word drivers. The constants have been obtained by an empirical fit on areas from actual layouts.

پهنای کلی دستگاه 1.25 برابر بافر نهائی می باشد. ناحیه کاری predecode عمودی در بلوک سیمی (شکل 1) همواره به ناحیه دکود کلی اضافه می شود. به عنوان مثالی افزایش درپهنای آرایه (صف) SRAM با استفاده از سطح برای درایور کلمه محلی 64، یک درایور کلمه سراسری 64 و سیم بندی عمودی برای 16 سیم predecode و 64 بلوک سیمی انتخابی محاسبه می شود.

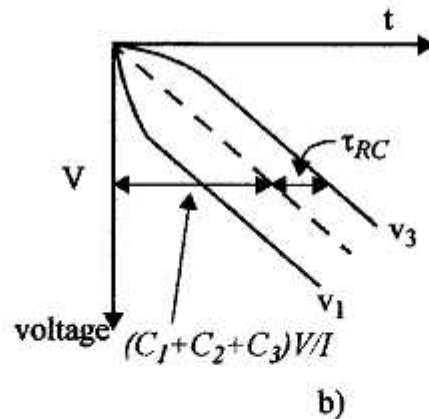
## MUX خروجی

Mux خروجی شامل ماکس خط بیتی می باشد که سلول داده رادر تقویت کننده قرار می دهد و ماکس خطی دیتا، دیتا را از تقویت کننده حسی به خروجی انتقال می دهد. از آنجائیکه سطح سیگنال در هر دوی این ماکس ها کوچک می باشد. ( $\sim 100 \text{ mV}$ ) منبع سیگنال ورودی برای هر دو ماکس منبع جریان ایده آل مورد استفاده قرار گیرد. کاهش تاخیر از طریق شبکه RC برای منبع جریان ورودی متفاوت از منبع ولتاژ ورودی می باشد منبع جریان ایده آل که برای درایو شبکه RC بکار می رود (شکل 8(a)). شکل موج ولتاژ در گره 1 و 3 در شکل 8(b) نشان داده شده است. ثابت زمانی شبکه در رابطه 5 نشان داده شده است.



$$\tau_{RC} = R_1 \frac{C_1(C_2 + C_3)}{(C_1 + C_2 + C_3)} + R_2 \frac{(C_1 + C_2)C_3}{(C_1 + C_2 + C_3)} \quad (5)$$

a)



b)

Fig. 8. (a) Current source driving a RC  $\pi$  network and (b) sketch of the node waveforms.

در حالت ثابت ( $t \gg \tau_{RC}$ ) گره 1 و 2 و 3 slew rate یکسانی دارند و تاخیر برای نوسان ولتاژ در گره 3 با فرمول 6 به صورت تقریبی قابل محاسبه است که تاخیر را برای زمانی می باشد که مقاومتی جود ندارد. این فرمول برای محاسبه هر دو ماکس خط بیت و خط دیتا بکار می رود

$$D + \frac{(C_1 + C_2 + C_3) \cdot V}{I} + \tau_{RC} \quad (6)$$

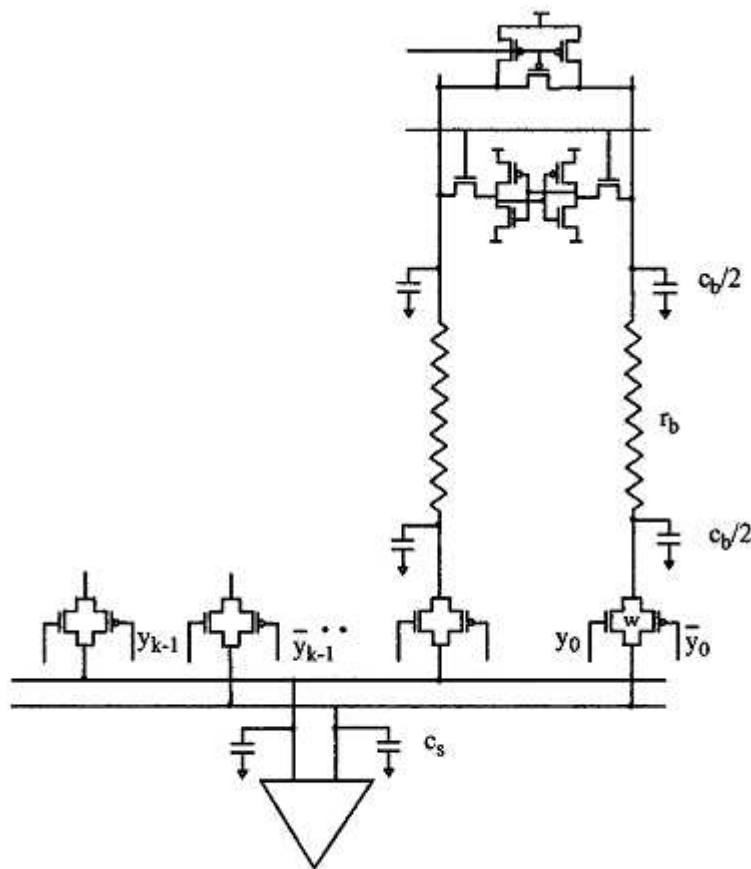


Fig. 9. Schematic of a single-level bit line structure.

یک ماکس خط بیتی تک سطحی در شکل 9 نشان داده شده است و به عنوان یک منبع جریان ایده آل برای درایو شبکه

RC (شکل 8) مدل شده است. سویچ ماکس و سیم خط بیتی محلی و سراسری باعث افزایش مقاومت و ظرفیت خازنی شبکه می شود. تاخیر خط بیتی برای نوسان سیگنال توسط  $\sigma_v$  مجموع تاخیر برای ایجاد نوسان ولتاژ با ثابت زمانی شبکه RC و بدون در نظر گرفتن مقاومت می باشد. خط کلمه محلی طولانی به علت مقاومت خط دارای زمان خیز کندی می باشد از آنجاییکه زمان خیز تاخیر سلولی را تحت تاثیر قرار می دهد لذا باید در مدل آنرا در نظر بگیریم. تاثیر زمان خیز ( $Tr$ ) با اضافه کردن عبارت اضافی به معادله تاخیر که متناسب با آن می باشد، بدست می آید. ثابت تناسبی به نسبت ولتاژ threshold سلول به ولتاژ منبع بستگی دارد و در شبیه سازی مقدار 0.3 را برای رنج وسیع بدست آورده ایم. ثابت زمانی در معادله تاخیر خط بیتی در عبارت 5 محاسبه شده است.

$$D_b = (C_b + j_{mux} \cdot (K + 1) + Cs) \cdot \frac{\delta v}{I_{Cell}} + z \cdot T_r + \tau_{RC} \quad (7)$$

- ظرفیت خارنی خط بیتی  $C_b$
- ظرفیت خارنی اتصال قسمت سویچ ماکس  $j_{mux}$
- تعداد ستون مالتی پلکس شده یه تقویت کننده حسی منفرد  $K$
- ظرفیت خارنی ورودی تقویت کننده حسی  $Cs$
- ولتاژ نوسان در ورودی تقویت کننده حسی  $\sigma_v$
- جریان سلول حافظه  $I_{cell}$
- زمان خیز خط کلمه محلی  $T_r$
- ثابت تناسبی مشخص شده بوسیله HSPICE  $z$
- ثابت زمانی خط بیتی در شبکه RC  $\tau_{RC}$

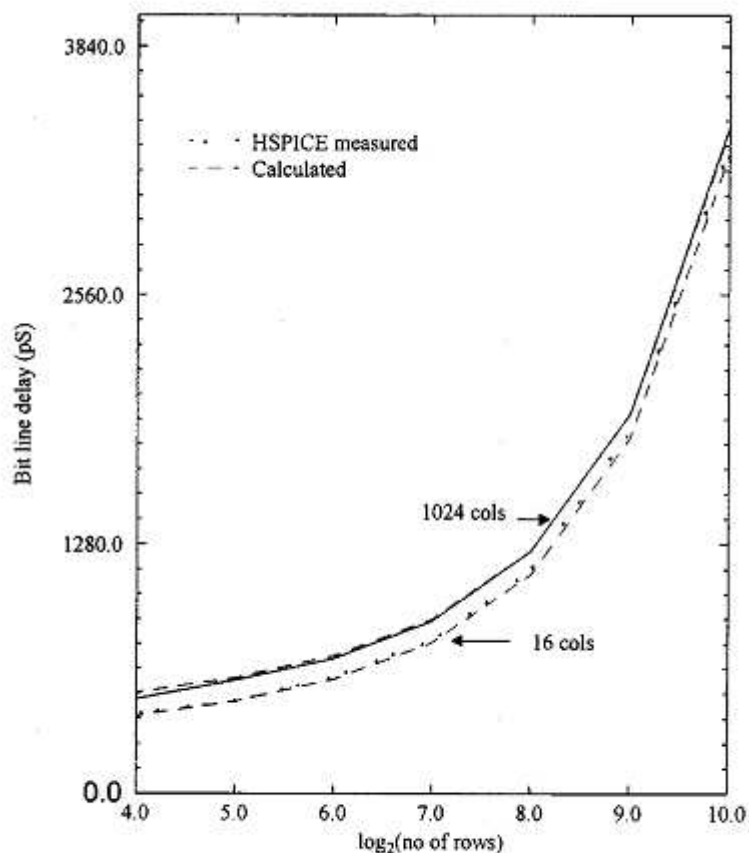


Fig. 10. Bit line delay versus column height;  $0.25 \mu\text{m}$ ,  $1.8 \text{ V}$ , and four columns multiplexing.

شکل 10 نشان دهنده ی تاخیر اندازه گیری شده توسط HSPICE از طریق درایور خط محلی و خط کلمه مقاومتی و خط بیتی ، نسبت به ورودی تقویت کننده حسی می باشد. تاخیر حساب شده برابر  $2.4\%$  تاخیر HSPICE زمانیکه ارتفاع خط بیتی 32 ردیف از هر دو ردیف کوتاه (16 ستون) و بلند (1024 ستون) می باشد .

بافر تقویت کننده حسی که در شکل 11 نشان داده شده است شامل لچ های cross-coupled اصلی که توسط زنجیره ای از اینورترها و یک جفت از درایورهای nMOS به هم وصل شده اند. مبدل لچ سیگنال نویسانی ورودی کوچک را به سیگنال CMOS نویسانی کامل تبدیل کرده و از آن برای هر دو تقویت کننده حسی محلی و سراسری استفاده می کند. در تقویت کننده حسی محلی ، خروجی لچ توسط زنجیره ای از اینورترها و به گیت خروجی nMOS درایو شده است. این ترانزیستورهای nMOS یک سیگنال ولتاژ کوچکی در خروجی آنها با دشارژ کردن خط دیتا بوجود می آورند. تاخیر در ساختار تقویت کننده حسی برابر با مجموع تاخیر تقویت کننده لچ  $\tau_s$  و تاخیر بافر و درایور خروجی می باشد.  $\tau_s$  متناسب با لگاریتم گین ولتاژهای مطلوب و بار خروجی تقویت کننده می باشد. اگر فرض کنیم تمام ترانزیستورهای لچ در همان مقیاس تناسبی باشند، مقاومت خروجی و

ظرفیت خازنی ورودی با تابعی از سایز nMOS که cross-coupled در لچ می باشد، بیان میشود. (شکل 11) درایورهای nMOS به عنوان منبع جریان مدل می شوند که خروجی جریان آنها متناسب با اندازه ی  $w_n$  می باشد. در دکودر سایزهای  $w_s w_1 \dots w_n$  برای کاهش تاخیر ماکس خروجی کلی بکار می رود. معادله 8 نشان دهنده ی مقدار مناسب تاخیر ماکس خروجی برای بهینه سازی می باشد و برابر با مجموع تاخیرهای ماکس خط بیتی، تقویت کننده حسی لچ، بافرها، درایورهای nMOS می باشد.

$$T = D_b(w_s) + \tau_s + \frac{R_s \cdot 3 \cdot C_g \cdot w_1}{w_s} + \frac{R_g \cdot 3 \cdot C_g \cdot w_2}{w_1} + \dots + \frac{R_g \cdot C_g \cdot w_n}{w_{n-1}} + \frac{C \cdot \delta v}{I_n \cdot w_n} + \text{other constants} \quad (8)$$

$$D_b(ws) \approx \frac{Gs \cdot Ws \cdot \delta v}{I_{cell}} \quad (9)$$

$$\tau_s = 2\tau_{f04} \quad \circ$$

$$Gs = 3.8fF/\lambda \quad \circ$$

$$Rs = 36k \quad \circ$$

$$W_s \text{ سایز تقویت کننده حسی} \quad \circ$$

$$Cg, Rg, \text{ مقاومت خروجی و ظرفیت خازنی ورودی} \quad \circ$$

$$I_n = 3.75\mu A/\lambda \quad \circ$$

$$C \text{ ظرفیت خازنی ماکس خط دیتا} \quad \circ$$

برای مشخص کردن روشی برای یافتن مطلوبترین سایز، تاثیر سایز تقویت کننده حسی لچ بر روی ثابت زمانی ماکس خط بیتی نادیده گرفته شده و تنها تاخیر سلولی مورد بررسی قرار گرفته (9) به عبارت ساده تراز تاثیر ظرفیت خازنی اتصال nMOS بر روی ثابت زمانی RC صرف نظر شده است. هر دوی این عاملها تاثیر کمی بر روی بهینه سازی سایز دارند ولی برای محاسبه نتیجه نهایی آنها را در نظر می گیریم. مینیمم تاخیر از طریق ساختار تقویت کننده حسی زمانی صورت می گیرد که هر کدام از عبارتهای در (8) برابر با تاخیر خارجی در 4 اینورتر بار شده باشد. تاخیر در تقویت کننده حسی در طرحی مشابه محاسبه می شود با این تفاوت که در این تحلیل، تاخیر بافر به درایو بار خروجی بررسی نشده است.

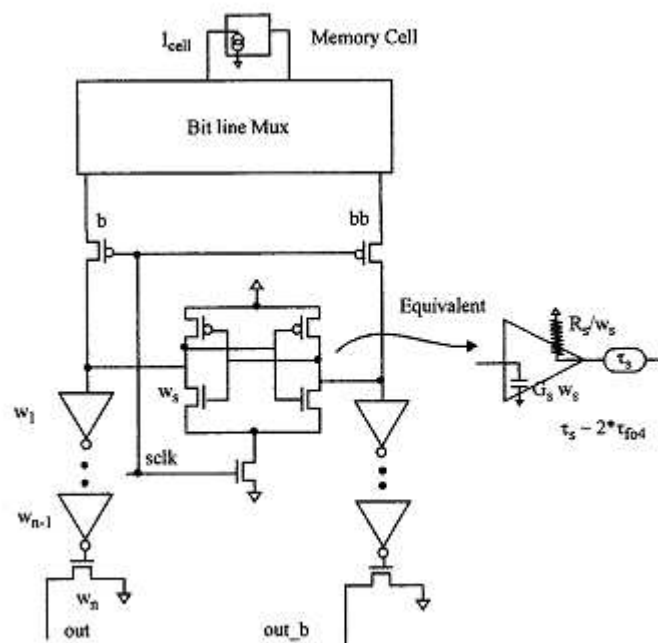


Fig. 11. Local sense amplifier structure.

برای عملیات کم توان، سیگنال در گره ظرفیت خازنی بالا، با نوسان ولتاژ کمی نگه داشته می شود. خط بیت با فرستادن خط کلمه باعث نوسان سیگنال کلی برابر  $2\sigma_v$  می شود. بنابراین انرژی ماکس خط بیت و ماکس خط دیتا که شامل ظرفیت خازنی خط و سیم و اتصالات و ظرفیت خازنی گیت ورودی و ولتاژ منبع می باشد محاسبه می شود. انرژی هر قسمت از شبیه سازی برابر 12 fJ برای 0.25- m پروسه می باشد. سطح سوئیچ در ماکس خط بیتی و تقویت کننده حسی مداری و precharge و درایورهای نوشتاری به قسمت عمودی ردیف SRAM اضافه می شود. (شکل 12) از آنجائیکه درایور نوشتن، precharge و ترانزیستور ماکس بهینه سازی نشده اند بنابراین به ترتیب یک سرجمع ثابت 4، 1 و 2 سلول حافظه اضافه می کنیم. سطح تقویت کننده حسی محلی به عنوان تابع خطی از پهنای دستگاه کلی داخل تقویت کننده حسی می باشد پارامترها برای مدل از 5 طرح مختلف بدست آمده است و در شکل 12 نشان داده شده است. پهنای دستگاه کلی داخل تقویت کننده حسی با پارامترهای  $\omega_p$  و  $\omega_s, \omega_n$  محاسبه می شود. مجموع پهنای دستگاه داخل لچ برابر  $8.7 \times w_s$  که مقدار 8.7 از طرح لچ بدست آمده است. با گنجایش خروجی 4 سائز، سطح فعال بافر بیشتر از 1/3 پهنای درایور نمی باشد. از این رو سطح فعال 2 درایور nMOS و بافرهای آنها برابر با  $1.33 \times 2 \times \omega_n$  می باشد. سپس ما به بررسی نتایجی که با استفاده از این مدلها برای آنالیز ساختار RAM متعدد در سائزهای مختلف و تکنولوژی ساخت متفاوت می پردازیم.

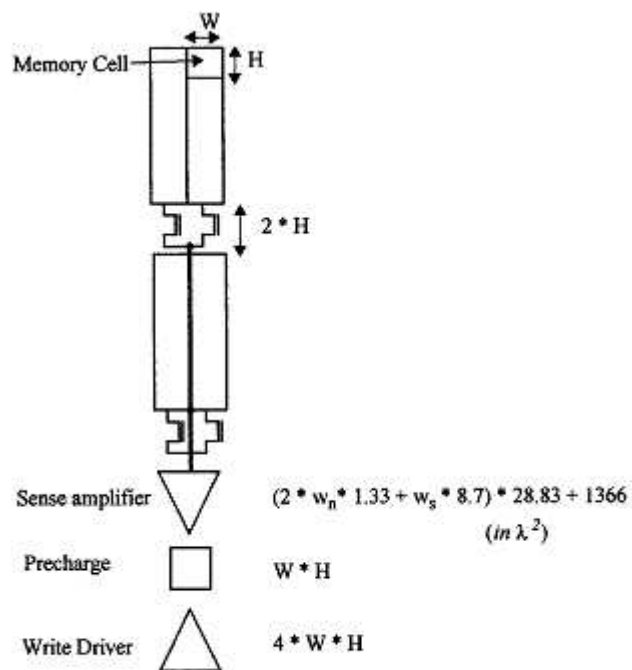


Fig. 12. Area estimation of the output mux.



## نتایج آنالیز:

ما تمامی ساختارهای RAM را یکا یک شمرده و سطح و تاخیر و انرژی هر کدام را با استفاده از مدل ساده ای در قسمت قبلی توضیح دادیم. که مشخص کردن ساختارهای بهینه برای توابع تاخیر، سطح و انرژی را و کمینه ساختن مقادیر آنها را فعال می کند.

$$\min (1 - \beta - \alpha) \cdot \text{Delay} + \alpha \cdot \text{Area} + \beta \cdot \text{Energy}. \quad (10)$$

منحنی تغییرات با عوض کردن مقادیر وزنی  $\alpha$  و  $\beta$  که بین 0 و 1 می باشد بدست می آید. شکل 13 نشان دهنده ی تاخیر ساختار SRAM برای حداقل تاخیر ( $\alpha = \beta = 0$ ) می باشد که به دو صورت با در نظر گرفتن مقاومت سیم و بدون آن ، برای سایز 64 kb تا 16 Mb با پهنای دسترسی 64 بیت در تکنولوژی 0.25-um نشان داده شده است. تاخیر SRAM بدون مقاومت سیم برابر  $15 \tau_{f04}$  برای طرح 64-kb و متناسب با لگاریتم ظرفیت مشاهده شده در می باشد. تاخیر برای هر دو RAM تا  $1.2 \tau_{f04}$  افزایش می یابد. تاخیر برای هر دوی اینها در همان نمودار نشان داده شده و برابر با مقدار بهینه سازی شده برای ساختار SRAM می باشد.

برای دکودر سطر هر بیت آدرس نصف ردیف را انتخاب می کند و بار دیده شده توسط بیت آدرس متناسب با  $S/2$  می باشد. که S تعداد کلی بیت های ردیف می باشد. با دوبرابر کردن در تعداد بیت ها نصف تاخیر  $\tau_{f04}$  اضافه می شود. در مسیر خروجی ، با دوبرابر کردن سایز خازن سیم در ماکس خط دیتا 1.4 برابر افزایش می یابد. که به محیط آرایه متناسب می باشد. بنابراین تاخیر تقویت کننده حسی محلی در حدود  $0.25 \tau_{f04}$  افزایش می یابد. منحنی نهائی در شکل 13 تاخیر با مقاومت سیم را نشان می دهد که سیم های سراسری برای این منحنی پهنای  $7.5 \text{ mm}$  /mm دارند. تاخیر RC سیم با طول سیم افزایش می یابد. تاخیر سیم برای سیم های سراسری در مقیاس SRAM برابر با سایز SRAM می باشد و برای SRAM بزرگ بیشتر نمایان می باشد.

بهینه سازی پهنای سیم به منظور کاهش تاثیر تاخیر داخلی صورت می گیرد. شکل 14 تاخیر کلی را برای 4-Mb SRAM در دو پهنای سیم مختلف در چهار تکنولوژی ساخت مختلف نشان می دهد. فرض بر این است که مس با ضخامت  $0.18 \text{ m}$  می باشد. منحنی پایینی نشان دهنده ی تاخیر با مقاومت صفر برای سیم می باشد. همانطوری که در شکل مشاهده می شود برای  $0.1 \text{ } \mu\text{m}$  تاخیر  $2.2 \tau_{f04}$  و برای  $0.07 \text{ } \mu\text{m}$  تاخیر  $3.6 \tau_{f04}$  افزایش می یابد زمانی که از تاخیر داخلی صرف نظر شده است. منحنی دوم تاخیر سیگنال round-trip را در مسیر دسترسی با فرض سرعت نور متناسب با  $1 \text{ mm}/6.6 \text{ pS}$  و مقدار تاخیر داخلی کمتری را ایجاد می کند. سرعت

تاخیر داخلی نور برابر  $1.75\tau_{f04}$  برای 4-Mb SRAM می باشد که مستقل از تکنولوژی می باشد و برای هر چهار برابر سایز RAM ، دو برابر افزایش می یابد . دو منحنی بالای آن نشاندهنده ی تاخیر با در نظر گرفتن مقاومت غیر صفر برای دو پهنای سیم  $8\lambda$  و  $10\lambda$  می باشد .

کاهش تاخیر سیم زمانی امکانپذیر است که سیم های کلفت برای سیمبندی سراسری بکار رود. از  $0.25\ \mu\text{m}$  با سیم بندی آلومینیومی تا  $0.18\ \mu\text{m}$  با سیم بندی مسی لزوما باعث تاخیر می شود. اما با طرح جمع شدن اضافی (further shrinks)، تاخیر برای تمامی پهنای سیمها کاهش می یابد.

تاخیر RC سیم به خوبی تاخیر گیت مقیاس بندی نشده است. به هر جهت با پهن سازی سیم ، امکان حفظ تاخیر قبلی میسر می شود. پهنای سیم  $10\lambda$  تاخیر با  $\tau_{f04}$  سرعت نور محدود بین  $0.25\text{-}\mu\text{m}$  و  $0.18\text{-}\mu\text{m}$  امکانپذیر می کند. تا زمانی که سیم های پهن تر  $0.1\ \mu\text{m}$  تا  $0.07\ \mu\text{m}$  مورد نیاز می باشد .

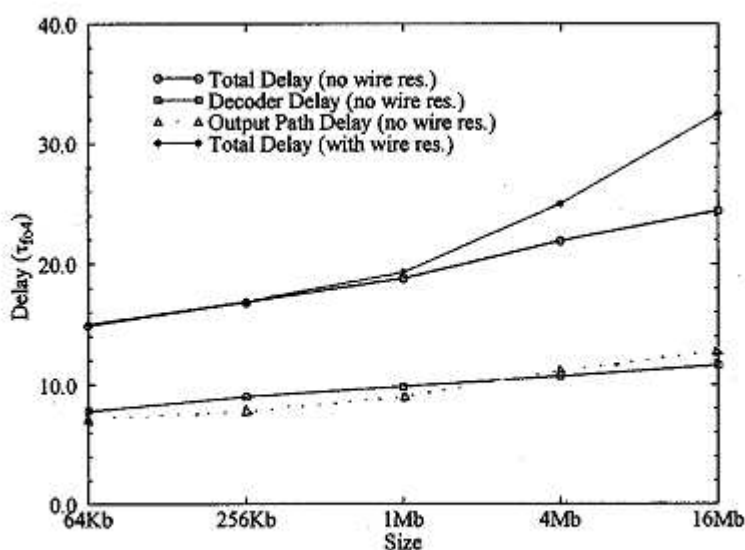


Fig. 13. Delay scaling with size in the  $0.25\text{-}\mu\text{m}$  process.

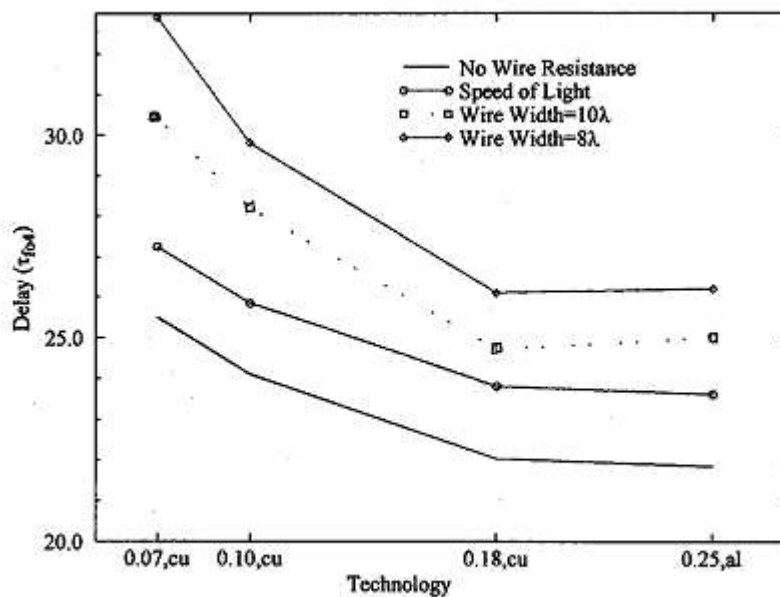


Fig. 14. Delay versus technology for different wire widths for a 4-Mb SRAM.

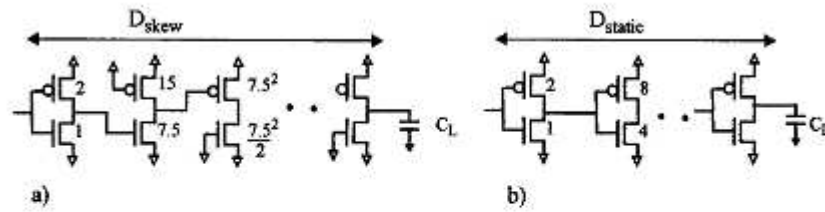
سپس ما روش هایی که عملکرد اجرا و پیاده سازی SRAM واقعی را که ممکن است با منحنی های قبلی متفاوت باشد بررسی می کنیم. معمولاً SRAM بزرگ مدارهای تکراری در مسیر دکود را با هم ترکیب می کند که باعث ایجاد منحنی تاخیر برای شیفت از  $1/2$  تا  $\tau_{f04}$  می شود. SRAM های پر بازده از سبک مداری استاتیک در مسیر دکود استفاده نمی کنند. برای محاسبه افزایش سرعت با skewing، ابتدا زنجیره ای از اینورترها که سیگنال ورودی یا به گیت nMOS یا به pMOS متصل شده باشد را در نظر می گیریم. (شکل 15) با فرض اینکه MOSFET مکمل در دسترس و گیت آن غیر فعال شده باشد. (در راه اندازی واقعی در ری ست کردن فاز، فعال می شود) و آن فقط به گیت با بار خود پرشونده اضافه می شود. با در نظر گرفتن این فرض ها و پارامترهای جدول 1 حداکثر گنجایش خروجی مطلوب برای در حدود 5 تاخیر گیت skewed تقریباً برابر 70٪ گیت بدون skewed می باشد. اگر predecoder و درایور کلمه سراسری بصورت skewed باشند تاخیر دکود برای اجرای استاتیک به جای  $8\tau_{f04}$  به  $6\tau_{f04}$  کاهش می یابد. با دو برابر کردن سایز RAM، تاخیر دکود  $\tau_{f04}$  به  $0.3\tau_{f04}$  به جای  $0.5\tau_{f04}$  افزایش می یابد. سرانجام کلاک برای تقویت کننده حسی معمولاً با تاخیر برای فعال کردن عملکرد در رنج وسیعی صورت می گیرد که موجب تاخیر اضافی تا  $\tau_{f04}$  شود. قسمت بندی کردن برای سبک و سنگین کردن میان تاخیرها، سطح و انرژی بکار می رود. منحنی تغییرات نشان دهنده ی مقادیر مختلف برای  $\alpha\beta$  می باشد که اگر  $\beta$  صفر شود منحنی برای SRAM 4-Mb در پروسه 0.25- m در شکل 16 نشان داده شده است. هر نقطه بر روی این منحنی نشان دهنده ی کمترین سطح قابل

دسترسی از طریق ساختار جدید RAM برای تاخیر مشابه می باشد. با شروع از کمترین تاخیر که به صورت دقیق قسمت بندی شده است، بهبودی مطلوب در سطح با کاهش مقادیر قسمت بندی و قبول تاخیر کوچک امکانپذیر می باشد و کاهش متوالی در قسمت بندی باعث کاهش بهبودی در سطح برای پذیرش تاخیر می شود. پارامترهای قسمت بندی در نقاط A، B و C در شکل نشان داده شده است. نقاط A، B قسمت مرکزی منحنی را تشکیل می دهند که در مقایسه با اجرای سریع، A، 22٪ کندتر و 22٪ سطح کوچکتر و B 14٪ کندتر و 20٪ سطح کوچکتر دارند.

اغلب پارامترهای ساختاری مختلف، تاخیر RAM بسیار حساس به ارتفاع بلوک می باشد و زمان دسترسی سریع با استفاده از ارتفاع کم بلوک میسر می شود. شکل 17 تاخیر و سطح را برای 4-Mb SRAM با ارتفاع بلوک های مختلف نشان می دهد زمانیکه از مقادیر مطلوب برای پارامترها استفاده شده است. ارتفاع کم بلوک تاخیر خط بیتی را کاهش ولی تاخیر سیم سراسری را افزایش می دهد. برای ارتفاع بلوک بلند تاخیر خط بیتی محدود به زمان دسترسی می شود. بنابراین ارتفاع مطلوب برای این مثال برابر 32 ردیف می باشد. افزایش ارتفاع به 128 ردیف منجر به افزایش تاخیر در حدود 8٪ در صورتی که سطح تا 7.6٪ می تواند کاهش یابد.

با قرار دادن  $\alpha = 0$  در معادله 10 میزان انرژی تاخیر از طریق قسمت بندی، بدون در نظر گرفتن محدودیت سطح بدست می آید که در شکل 18 نشان داده شده است. قسمت عمودی سمت چپ میزان انرژی مصرفی برای سوئیچ گیت اینورتر می باشد (Eunit 72 fJ) قسمت بندی امکان مقایسه بین انرژی و تاخیر را فراهم می کند. شکل میزان مطلوب مالتی پلکس ستون (cm) و ارتفاع بلوک (bh) برای پاسخ تاخیر و انرژی برای تعدادی نقاط نشان می دهد.

به این نتیجه رسیدیم که راه حل برای انرژیهای کم استفاده از یک ستون برای مالتی پلکس می باشد. از آنجائیکه ما بهینه سازی برای کاهش تاخیر انجام می دهیم ترانزیستورهای آخری در خروجی تقویت کننده بزرگ می باشند و دارای ظرفیت خازنی بزرگی می باشد. از این رو در طراحی کم انرژی داشتن ارتفاع زیاد بلوک مزیت به شمار می رود که باعث می شود اکثریت muxing در ماکس خط بیتی که ظرفیت خازنی S ترانزیستور سلول دسترسی حافظه بسیار کمتر از ظرفیت خازنی S در ماکس خط دیتا می باشد، صورت گیرد. همچنین مصرف انرژی در بهینه سازی ساختار SRAM به مجموع دو جز بیان می شود که یکی مستقل از ظرفیت می باشد تنها به پهنای دسترسی بستگی دارد و به علت خط کلمه محلی، سیگنال precharge، تقویت کننده محلی و سراسری و... مقیاس قسمت دیگری ریشه دوم ظرفیت می باشد همانطوری که در . مشاهده شد با اتلاف قدرت در سیم های سراسری و دکودرها ارتباط دارد.



$$D_{\text{skew}} = 0.31 * \tau_{f04} + 0.62 * D_{\text{static}} \quad (10)$$

Fig. 15. Optimal sizing for (a) extremely skewed and (b) statically sized inverters.

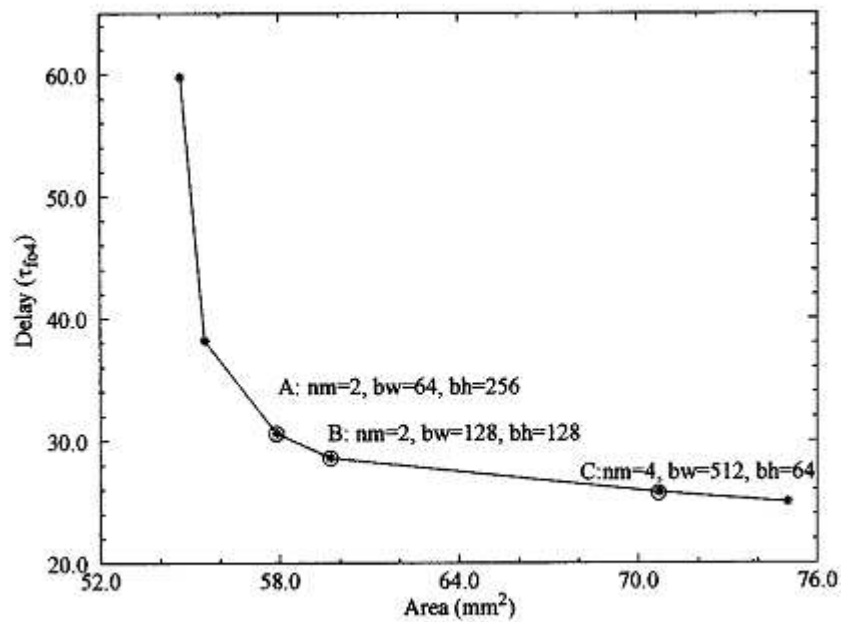


Fig. 16. Delay versus area for a 4-Mb SRAM in the 0.25-μm process.

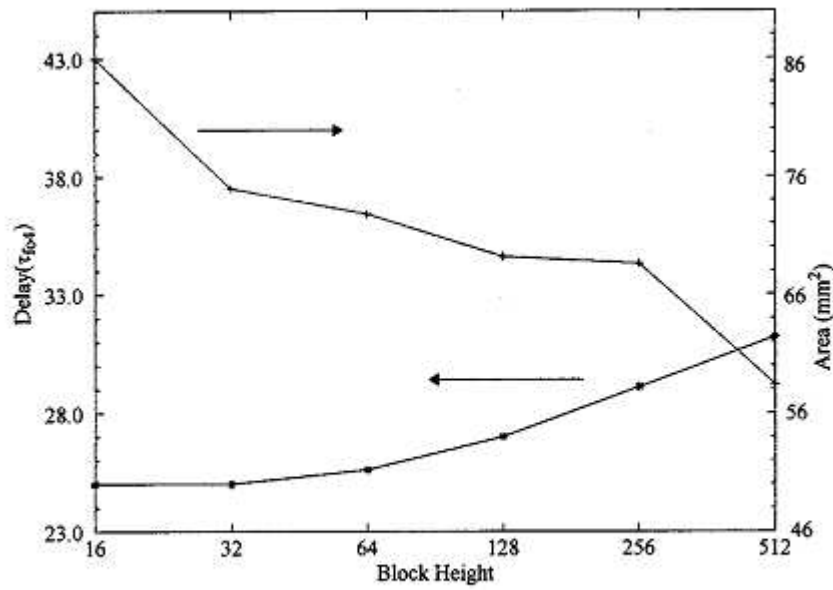


Fig. 17. Delay and area versus block height for a 4-Mb SRAM in a 0.25- $\mu\text{m}$  process.

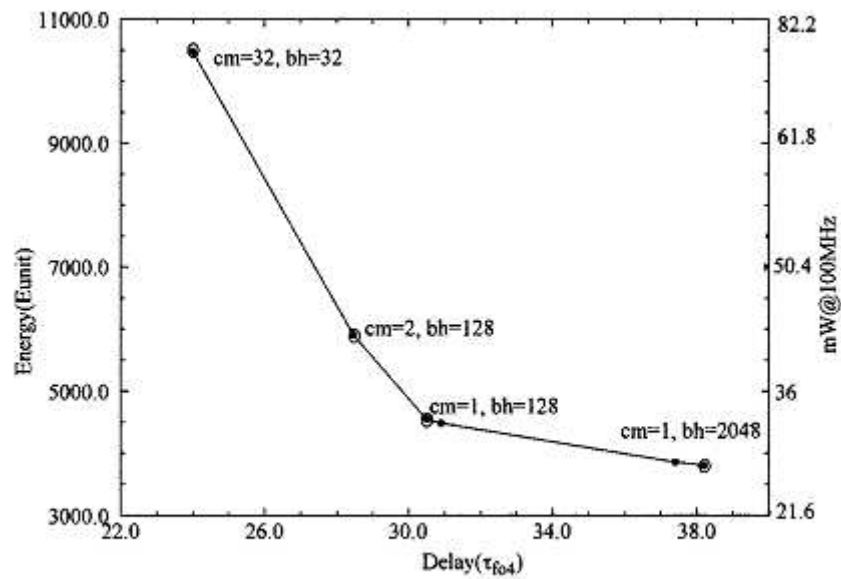


Fig. 18. Energy versus delay for a 4-Mb SRAM in a 0.25- $\mu\text{m}$  process.

## نتیجه گیری :

مدل های تحلیلی برای تاخیر، سطح و انرژی ، طراحی در محدوده کوچک زمانی را امکانپذیر می سازد. این مدل ها برای بررسی تاثیر قسمت بندی SRAM بکار می رود. تغییرات قابل توجه بین سطح، تاخیر و انرژی از طریق انتخاب ساختارهای SRAM قابل مشاهده است. این مدل ها برای پیش بینی میزان تاخیر با ظرفیت و تکنولوژی پروسه بکار می رود. تاخیر SRAM به دو قسمت قابل تفکیک می باشد. یکی به علت ترانزیستورها ی بکار رفته در تکنولوژی می باشد. (تاخیر گیت) و دیگری به خاطر اتصال درونی (تاخیر سیم ) می باشد. تاخیر گیت به ازای دو برابر کردن سایز RAM،  $1.2\tau_{f04}$  افزایش می یابد. زمانیکه از طرح مداری استاتیک برای طراحی دکودر استفاده می شود برای 64-kb RAM از  $15\tau_{f04}$  شروع می شود. برای بهینه سازی ساختار 4-Mb SRAM، افزایش تاخیر در حدود  $2\tau_{f04}$  برای  $0.1\text{-}\mu\text{m}$  و در حدود  $3.6\tau_{f04}$  برای  $0.07\text{-}\mu\text{m}$  می باشد که در ساختار دیگر کمتر می باشد. برای اکثر ساختارهای SRAM با استفاده از طرح های مرتبه ای برای مسیرهای خط بیتی و خط دیتا و با استفاده از تکنیکهای جبران آفت، می توان از افزایش تاخیر جلوگیری کرد. تاخیر سیم برای RAM های بیشتر از 1-Mb تاثیر گذار می باشد. با استفاده از روش shrink، تاخیر سیم کاهش می یابد و طراحی دوباره سیم برای حفظ تاخیر سیم نسبت به تاخیر گیت با همان نسبت قبلی باید انجام گیرد. تقسیم ساختار خط کلمه برای دکودرها و ماکس ستون برای مسیر خط بیتی، فضای کافی برای استفاده از سیمهای کلفت تر ایجاد می کند که باعث کنترل تاخیر برای 4-Mb و طراحی کوچکتر می باشد. تاخیر سیم دارای حد کمتری بوسیله سرعت نور می باشد که برای 4-b SRAM در حدود  $1.75\tau_{f04}$  می باشد و با چهار برابر کردن ظرفیت، دو برابر می شود. در سطح های بالا تر، معماری RAM نیاز به تغییر برای استفاده از مسیر یابی آدرس و دیتا دارد. (مثال 14 را مشاهده کنید). تاخیر سیم متناسب با سطح سلول می باشد و طراحی سلول با سطح کمتر برای RAM بزرگ مزیت به شمار می رود. حتی اگر سلول ضعیفتر باشد. بنا براین سلول DRAM، سلول چند ارز multivalued ، سلول TFT-based و طراحی سلول های جدید دیگر برای کارایی زیاد و ظرفیت زیاد RAM مورد نیاز می باشد.

## ضمیمه :

ما تمام فرض هایی را که در قسمت مدل سازی بیان نشد را در این قسمت لیست کرده ایم :

## تکنولوژی :

تکنولوژی اصلی در مقیاس  $0.25\mu\text{m}$  CMOS فرض شده است و توضیحات مناسب در جدول 1 نشان داده شده است. ویژگیهای اساسی برای 4 ساختار مختلف در جدول 2 بیان شده است. فلز کاری مس با فرض تولید  $0.18\mu\text{m}$  به بعد انجام شده است. فلز سطح بالاتر به حالت ضخیم تر طراحی شده است: ارتفاع و پهنای آنها در مقیاسی طراحی شده است که دارای قسمت عبوری بزرگتری باشند اما افزایش ارتفاع آن تنها با ریشه دوم فاکتور افزایش پهنای امکانپذیر می باشد. برای مثال لایه فلزی سطح بالا تر با دو برابر پهنای مینیمم لایه فلز 1 دارای ارتفاع 1.4 برابر ارتفاع فلز 1 می باشد بنابراین مقاومت فاکتور 3 کوچکتر از مقاومت فلز 1 می باشد. فرض ما بر این است که pitch سیمبندی دوبرابر پهنای سیم برای سیم های سراسری می باشد.

## معماری :

SRAM سنکرون می باشد برای مثال کلاک امکان دسترسی را شروع می کند بنابراین به آسانی نتیجه می شود که با اضافه کردن SRAM آسنکرون توان و تاخیر برای ردیابی انتقال آدرس (ATD) بکار می رود. ساختار SRAM جا سازی شده مفروض است. تمام بیت های دیتا کلمه ی دسترسی از هسته حافظه در مجاورت فیزیکی یکدیگر حاصل می شود (شکل 1) برخلاف SRAM مستقل، محل قرارگیری پورت IO دیتا برای ارتباط خارجی بهینه سازی شده است. این بهینه سازی یک مقدار ثابت آفست به تاخیر و توان هسته SRAM اضافه می کند که نتیجه این مطالعه قابل اجرا برای SRAM مستقل می باشد. اندازه سلول برای تحلیل در نشان داده شده و سطح سلول معمولاً از طرح بندی 6 ترانزیستور CMOS با چگالی بالا می باشد.

## ساختار مدار :

RAM برای کارکرد پرسرعت با توان کم طراحی شده است که اتلاف انرژی را بدون تأثیر بر سرعت کاهش می دهد. خط کلمه محلی برای کنترل نوسان خط بیتی بکار می رود که نوسانهای کم برای کاهش توان در خط دیتا بکار می رود. از آنجایی که این روش سرعت RAM را تحت تأثیر قرار نمی دهد نتایج تحلیل مامربوط به میزان تاخیر قابل اجرا برای طراحی SRAM با سرعت مطلوب می باشد. یک تقویت کننده حسی طرح لچ با کنترل زمانی کامل برای تقویت کننده فرض شده است که توان بسیار کمی مصرف می کند و کارایی سریع دارد. بنابراین نتایج تحلیل ما منجر به SRAM پرسرعت و کم توان می شود. برای پروسه  $0.25\mu\text{m}$  نوسان



ورودی مطلوب که تاخیر تقویت کننده را کاهش می دهد از شبیه سازی مقدار  $100\text{ mV}$  که  $50\text{ mV}$  آفست ورودی می باشد. ترانزیستورها در ماکس خط بیتی سائز ثابتی دارند و در ماکس خط دیتا برای مشخص کردن تحلیل پهن تر می باشند. شبیه سازی مداری تاخیر RAM را نشان می دهد که حساسیت کمی به سائز این ترانزیستورها دارند.

#### مدل انرژی :

نوسان در خط بیتی و خط IO به عملیات کم توان محدود می شود. زمانی که دقیقا به مقدار مطلوب توسط تقویت کننده حسی محدود می شود، در طراحی عملی، متغیرهای کمکی برای کنترل مناسب آن وجود دارد و فرض می کنیم که دوبرابر نوسان سیگنال مطلوب باشد. بنابراین برای پروسه  $0.25\text{-}\mu\text{m}$ ، نوسان برابر می باشد که  $200\text{ mV}$  می باشد که نوسان مطلوب برای تقویت کننده برابر  $100\text{ mV}$  می باشد.