



دانشگاه آزاد اسلامی

واحد خمین

## معماری کامپیوتری پیشرفته

دکتر محمد خلیلی درمنی

پائیز ۹۳

## فهرست مطالب

صفحه

عنوان

## معماری کامپیوتر پیشرفته

۴.....	مقدمه.....
	۱- روش های عددی متفاوت
۱۲.....	۱-۱- سیستم اعداد مانده ای (RNS)
۱۴.....	۱-۱-۱- اثبات عمل جمع در مبنای RNS
۱۷.....	۱-۱-۲- برگرداندن اعداد از پیمان به مبنای RNS
۲۰.....	۲-۱- اعداد در مبنای RTIT
۲۱.....	۱-۲-۱- روش به دست آوردن اعداد در RTIT
۲۲.....	۲-۲-۱- روش بدست آوردن نمایش های مختلف اعداد در مبنای RTIT
۲۵.....	۳-۱- اعداد در مبنای BRTIT
	۲- سازمان حافظه
۲۸.....	۱-۲- حافظه تداعی گر
۲۸.....	۱-۱-۲- بررسی مدار سخت افزاری حافظه تداعی گر
	۲-۱-۲- مدار انطباق
۳۱.....	۲-۲- حافظه cache
۳۴.....	۱-۲-۲- محلیت زمانی و مکانی حافظه کش
۳۴.....	۲-۲-۲- کارایی حافظه cache
۳۵.....	۳-۲- نگاشت
۳۶.....	۱-۳-۲- نگاشت مستقیم
۳۷.....	۲-۳-۲- نگاشت تداعی گر
۳۸.....	۳-۳-۲- نگاشت تداعی گر مجموعه ای
۳۹.....	۴-۲- خط لوله
۳۹.....	۱-۴-۲- روش ترتیبی (بدون خط لوله)
۴۰.....	۲-۴-۲- روش با خط لوله

- ۳-۴-۲. مخاطرات (Hazard) خط لوله..... ۴۱
- ۴-۴-۲. انواع پیاده سازی خط لوله..... ۴۵

### ۳- چندپردازنده ها

- ۳-۱- طبقه بندی سیستم ها..... ۵۱
- ۳-۲- ساختار های اتصالات متقابل..... ۵۳
- ۳-۲-۱ سیستم چندپردازنده معمولی..... ۵۳
- ۳-۲-۲ حافظه های چند پورته..... ۵۵
- ۳-۲-۳ حافظه با سوئیچ تقاطعی..... ۵۶
- ۳-۲-۴ شبکه سودیج های چند طبقه (استفاده از سودیج های بانیان)..... ۵۷
- ۳-۲-۵ روش hypercube..... ۵۹
- ۳-۳ Super pipeline..... ۶۲
- ۳-۴ Super scalar..... ۶۳
- ۳-۵ Super pipeline-scalar..... ۶۴

### ۴- ارزیابی کارایی چند پردازنده ها

- ۴-۱- معیار MFLOP..... ۷۱
- ۴-۲- MIPS..... ۷۳
- ۴-۳- معیار SPEC (System Performance Evaluation Cooperative)..... ۷۳
- ۴-۴- قانون امدال (Amdahl's law)..... ۷۴

### ۵- پیوست

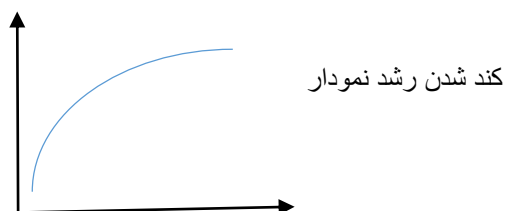
- ۵-۱- روش چینی تسریع شده..... ۷۹

## مقدمه

حافظه ها و CPU در درون IC قرار می گیرند و تکنولوژی IC هم همیشه در حال پیشرفت است و برای داشتن حافظه بیشتر باید تعداد ترانزیستورهای داخل IC را افزایش دهیم اما اندازه IC ها از یک حد بیشتر نمی تواند باشد به همین دلیل اندازه ترانزیستورها کوچکتر و تعداد بیشتری IC در آن قرار می دهیم. با پیشرفت تکنولوژی اندازه ترانزیستورها کوچکتر و حافظه بزرگتر می شود. و بر این اساس شخصی به نام Moore یک قانون که به قانون طلایی Moore معروف است را بیان کرد طبق این قانون تعداد ترانزیستورهای درون IC هر ۱۸ ماه یک بار دو برابر می شود. به طور مثال اگر در یک IC، ۱۰۰۰ ترانزیستور قرار داشته باشد ۱۸ ماه بعد، ۲۰۰۰ ترانزیستور در IC قرار دارند یعنی رشد این پیشرفت به صورت نمایی است. اما وجود دو مشکل باعث شد نمودار از یک حدی به بعد کند پیشرفت کند این مشکلات عبارتند از:

### مشکلات:

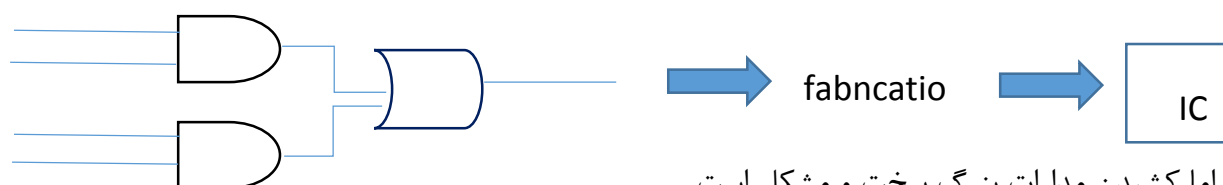
۱. محدودیت سخت افزاری (فیزیکی): وقتی اندازه ترانزیستورها هر ۱۸ ماه یکبار کوچک تر شد تا به ۲۵ نانومتر رسید ولی دیگر نمی توان این سایز را کاهش داد چون ترانزیستور از یک اتم کوچکتر نمی تواند باشد.
۲. محدودیت اقتصادی: ساخت ترانزیستور برای شرکت ها هزینه دارد و بدون هزینه پیشرفتی حاصل نمی شود به همین دلیل این پیشرفت از هر ۱۸ ماه یکبار به ۲۴ ماه یکبار افزایش یافت.



### ارائه راه حل مشکلات:

قبل از بیان راه حل این مشکلات مروری بر نحوه ساخت IC می نمایم. دو راه برای ساخت IC وجود دارد:

راه حل اول: IC ها را یک به یک رسم کرده و خروجی ها را با fabncation می دهیم.



اما کشیدن مدارات بزرگ سخت و مشکل است.

راه حل دوم: استفاده از زبان HDL (زبان توصیف سخت افزار)

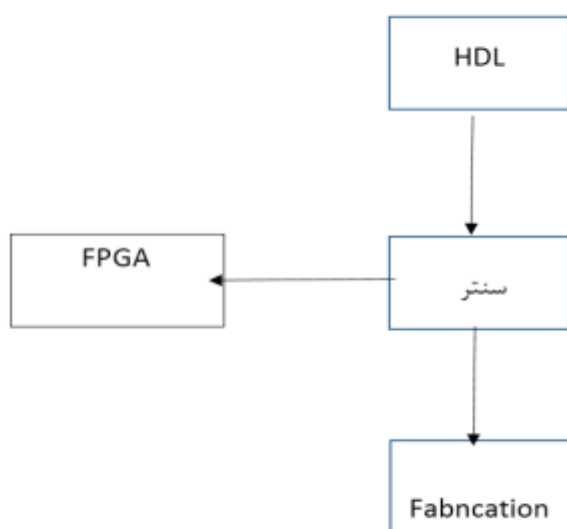
دو زبان معروف برای توصیف سخت افزار وجود دارد :

۱- VHDL

۲- Veriling

با این زبان ها می توان سخت افزار را مدل کرد : به این صورت که مداری که داریم را با این زبان ها می نویسم بعد برنامه ها را به یک ابزار سنتز می دهیم و این ابزار برنامه را به مدار تبدیل کرده و سپس خروجی آن را به کارخانه می دهیم یعنی امروزه برای تولید سخت افزار نیز باید برنامه بنویسیم .

گاهی اوقات هم می توان خروجی سنتز را به FPGA ها داد که FPGA یک IC آماده است که برنامه را روی آن می ریزیم و همان IC مورد نظر ما می شود در این روش هزینه پایین است و برای تعداد کم مناسب است .



راه حل : بیان ایده های جدید به جای روش های قدیمی مثلا به جای ذخیره بیت و صفر و یک از Rbit استفاده کرد و یا از کامپیوتر های کوانتومی استفاده کرد . ما برای هر اتم یک هسته و یک الکترون داریم که دور هسته می چرخند و در آن جهت چرخش الکترون در اتم را یک بیت در نظر می گیرند که به آن QBit می گویند به طور مثال از چپ به راست چرخیدن معادل عدد یک است .

مثال کامپیوتر کوانتومی : در رمز نگاری برای یک کامپیوتر معمولی صد سال طول می کشد تا رمز را تشخیص دهند ولی با کامپیوتر کوانتومی همین کار در یک ساعت انجام می شود و به جای اینکه روی ولتاژ کار کنند روی جریان کار می کنند .

همه این ایده ها سعی می کنند محدودیت ها را از بین ببرند .

تمرین در مورد کامپیوترهای کوانتومی توضیح دهید .

## جواب : تاریخچه کامپیوتر کوانتومی:

نظریه کامپیوتر کوانتومی از سال ۱۹۸۲ مطرح بوده است از زمانی که فیزیک دان مشهور و برنده جایزه نوبل «ریچارد فاینمن» برای نخستین بار، پیشنهاد کرد که باید محاسبات از دنیای دیجیتال وارد دنیای جدیدی به نام کوانتوم شود، همچنین بیان کرد کامپیوتر کوانتومی چگونه ممکن است کار کند. این پیشنهاد تا اوایل سال ۱۹۹۰ مورد توجه جدی قرار نگرفت و به صورت آکادمیک باقی ماند، البته در سال ۱۹۸۵، دویتش متوجه شد که اظهارات فاینمن، می تواند تدریجاً به ساخت کامپیوتر کوانتومی منجر شود و مقاله ای را منتشر کرد مبنی بر اینکه اصولاً هر فرآیند فیزیکی را می توان به خوبی با کامپیوترهای کوانتومی مدل سازی کرد. بالاخره در ۱۹۹۴ «پیتر شور» نخستین گام را برای محقق کردن این آرزو برداشت. وقتی که بعضی از مشکلات کلیدی کامپیوترهای معمولی نشان داده شد، کامپیوترهای کوانتومی در اصل می توانستند خارج از رونوشت های کلاسیکی خود محاسبات را انجام و اجرا نمایند یعنی کارایی بسیار بالاتری را نسبت به کامپیوترهای معمولی از خود نشان می دادند. وی مقاله ای را منتشر نمود که حاوی روشی برای استفاده از کامپیوترهای کوانتومی در حل مشکل پیچیده ای در نظریه اعداد، به نام فاکتورگیری بود. او نشان داد که چگونه یک مجموعه از عملیات ریاضی که منحصراً برای کامپیوترهای کوانتومی طراحی شده اند، می توانند چنین دستگاهی را به انجام فاکتورگیری از اعداد بی شماری با سرعت بالاتر از کامپیوترهای کلاسیک، قادر سازد با این اختراع، محاسبات کوانتومی از یک کنجکاوی به یک توجه جهانی تبدیل شد.

از آن موقع به بعد، گروه های تحقیقاتی در سرتاسر دنیا مسابقه ای را برای پیش قدم شدن در ساخت یک سیستم عملی آغاز نمودند. به این ترتیب ارتباط نوینی بین نظریه ی اطلاعات و فیزیک کوانتومی شروع به شکل گیری کرد که امروزه انرا محاسبات کوانتومی یا محاسبات نانو متری ( Nano Computing ) می نامیم. محاسبات کوانتومی مشکلات گذشته را برطرف می سازد و افق جدیدی را ایجاد می کند.

قدرت خارق العاده کامپیوتر کوانتومی در نتیجه وقوع پدیده ای موسوم به توازی کوانتومی، مکانیزی که انجام و اجرای محاسبات حجیم، زیاد و مکرر را به طور همزمان مقدور می سازد. این یک مقابله سخت و نیرومند و فرق نمایان و بزرگ با کامپیوترهای کلاسیک است که قادرند تنها هر عملیات را فقط به صورت یک عمل در هر بار و البته خیلی سریع انجام دهند.

چه تفاوتی میان یک کامپیوتر کوانتومی و یک کامپیوتر کلاسیک وجود دارد؟

بین کامپیوترهای کوانتومی و کامپیوترهای کلاسیک تفاوت های اساسی وجود دارد:

۱. در کامپیوترهای کوانتومی به جای استفاده از ترانزیستورها و مدارهای رایانه ای معمولی از اتم ها و سایر ذرات ریزمانند نانو ذرات نیمه رسانا نقاط کوانتومی Quantum dots برای پردازش اطلاعات استفاده می کنند. یک اتم می تواند به عنوان یک بیت حافظه در رایانه عمل کند و جا به جایی اطلاعات از یک محل به محل دیگر نیز توسط نور امکان می پذیرد. ذخیره اطلاعات در کامپیوترهای نانو به صورت سری هایی از بیت های با حالت های روشن و خاموش صورت می گیرد.
۲. در مقایسه این ۲ نوع کامپیوتر می توان گفت، مسائلی که زمانی تصور می شد در کامپیوترهای کلاسیک غیر قابل حل است، در کامپیوترهای کوانتومی حل خواهد شد و شبیه سازی های صورت گرفته به واقعیت نزدیک تر می شود. حتی ابر کامپیوترها هم در برابر آنها رقیبی محسوب نخواهند شد. به عنوان مثال، به روز رسانی نرم افزار، Email، بانک های آنلاین و تمام قلمرو رمز نگاری عمومی و امضاهای دیجیتال، فقط از دوروش رمز نگاری برای ایمن نگاه داشتن خود استفاده می کنند. RSA و ECC رمزنگاری منحنی بیضی دوروشی هستند که کشف رمز این روش ها، برای کامپیوترهای کلاسیک تا حد زیادی ناشدنی است. ولی یک کامپیوتر کوانتومی به اندازه کافی برای شکستن هر دوی این کدها، قدرتمند است. ECC برای امضاهای دیجیتال استفاده می شود، که اطمینان می دهد یک پیغام واقعا توسط فرستنده مدعی، فرستاده شده است RSA. برای بیشتر سیستم های رمزنگاری کلید عمومی استفاده می شود، که در آن یک پیغام، با یک کلید عمومی مجاز کدگذاری می شود و باید با قوانین ریاضی مبتنی بر کلید سری رمزگشایی شود.
۳. کامپیوترهای کوانتومی از یک خاصیت دیگر هم سود می برند که آنها را از کامپیوترهای کلاسیک مستثنی می کند و آن انتقال از راه دور است. انتقال از راه دور موجب می شود، اطلاعات یک ذره به ذره دیگری منتقل شود. در نتیجه کامپیوترهای کوانتومی برای انتقال بیت در درون و بیرون ساختار خود، نیازمند سیم نیستند.
۴. تفاوت دیگر کامپیوترهای کوانتومی با کامپیوترهای کلاسیک این است که، اندازه ترانزیستورها هر سال کوچکتر می شود. وقتی اندازه ترانزیستورها به ابعاد اتمی نزدیک می شود، دیگر قوانین حاکم بر فیزیک کلاسیک بر رفتار اتم ها حاکم نیست. به طور مثال کسی نمی داند یک الکترون در زمان مشخصی، دقیقاً در کجا قرار دارد یا کسی نمی تواند به درستی تشخیص دهد که الکترون در یک سیم به کجا می رود. یعنی وقتی به ابعاد اتمی نزدیک می شویم، فیزیک کوانتومی رفتار اتم ها را توضیح می دهد و دیگر قوانین فیزیک کلاسیک کاربرد ندارد. در واقع این نوع کامپیوترها با استفاده از فناوری های میکروسکوپی ذره ها کار می کنند.

۵. همان طور که می دانیم در یک کامپیوتر کوانتومی نسبت به کامپیوترهای کلاسیک، اصول حاکم تغییر نموده اند. نه تنها، یک بیت کوانتومی، موسوم به کیوبیت می تواند در حالت های صفر و یک کلاسیک وجود داشته باشد بلکه همچنین می تواند در حالت برهم نهی قرار داشته باشد. هرگاه هر کیوبیت در یک کامپیوتر کوانتومی در حالت برهم نهی واقع شده باشد، آنگاه کامپیوتر می تواند در هر حالت ممکن مجسم کرد که آن کیوبیت می تواند از خود نشان دهند. در واقع کامپیوترهای کوانتومی می تواند در یک زمان چندین حالت داشته باشد و این امکان را ایجاد می کند که میلیون ها بار سریع تر و قدرتمند تر از ابر کامپیوترهای فعلی کار کند. چند حالت پذیری کیوبیت ها همان دلیلی است که باعث می شود کامپیوترهای کوانتومی ذاتاً از پردازش موازی بهره ببرند. پردازش موازی امکان کار کردن بر روی میلیون ها محاسبه در یک لحظه را به این کامپیوترها می دهد در حالی که کامپیوتر شخصی شما فقط یک محاسبه در لحظه انجام می دهد.

### نحوه برقراری ارتباط در کامپیوترهای کوانتومی

با توجه به ماهیت ساختار کامپیوترهای کوانتومی، روش برقراری ارتباط آنها کاملاً متفاوت با کامپیوترهای امروزی است. بدین صورت که پالس های رادیویی نقش صفحه کلید را دارند، که به وسیله آن اطلاعات وارد کامپیوتر می شود و دستگاه تشدید مغناطیسی که شبیه به دستگاه MRI بیمارستان است، نقش صفحه نمایش را ایفا می کند و با ارائه تصویر مغناطیسی از توده مولکولها، کامپیوتر توده محاسبات را به ما می دهد.

از سوی دیگر باید تلاش کرد ترانزیستورهایی از جنس مورد نظر ساخت، زیرا ترانزیستورها، عامل تقویت ولتاژ در مدارهای الکترونیکی هستند و قدرت تقویت کنندگی آنها موجب افزایش سرعت کامپیوترها است.

این ترانزیستورها تاثیر مهمی در تولید کامپیوترهای آینده دارند و در صورتی که در ابتدا یا انتهای ساختار آنها، ترکیبی با دیگر نیمه هادی ها به خصوص طلا ایجاد شود یا حتی روی پوسته آنها نیمه هادی مهمی چون «روی» قرار داده شود، گام مهمی برای تولید قدرتمندترین ترانزیستورها برداشته ایم.

در واقع زمانی که این نیمه هادی های ترکیب شده، به یک باطری متصل می شوند و الکترونیته دریافت می کنند، همچون یک ترانزیستور عمل کرده و موجب تقویت ولتاژ در مدار و همچنین موجب افزایش سرعت کامپیوتر می شوند. از آنجایی که جریان ورودی به این نیمه هادی ها قابل کنترل است، جریان خروجی از



آنها هم قابل کنترل است. این ترکیب ها این قابلیت را دارند که در ساخت نقاط کوانتومی مورد استفاده در کامپیوترهای نسل آینده استفاده شوند.

این نقاط کوانتومی در واقع کریستال هایی از نوع نیمه هادی هستند، که قابلیت ذخیره کردن الکترون ها در آنها فوق العاده بالاست. این نقاط کوانتومی بهترین مکان برای ذخیره سازی اطلاعات در کامپیوترهای پیشرفته هستند. از سوی دیگر در صورتی که بتوان نقاط کوانتومی را با یکدیگر پیوند داد، می توان آنها را به اندازه تنها چند سانتی متر مربع در ساختار سخت افزاری کامپیوترها جای داد. به این ترتیب هر سانتی متر مربع از درایورهای ما می توانند صدها گیگا بایت از اطلاعات را در خود ذخیره کنند.

تا چند سال دیگر، کامپیوترهای کوانتومی از داخل آزمایشگاه های تحقیقاتی دانشمندان علوم رایانه، فیزیک و ریاضی دانان بیرون خواهند آمد و به صورت کاربردی و عملی مورد استفاده قرار خواهند گرفت. آن دسته از مسائل که با محاسبات پیچیده ی خود، کامپیوترهای جبری امروز را به ستوه می آورند، توسط کامپیوترهای کوانتومی به آسانی حل خواهد شد.

### محاسبات کوانتومی

هدف محاسبات کوانتومی یافتن روش هایی برای طراحی مجدد ادوات شناخته شده ی محاسبات (مانند گیت ها و ترانزیستورها) به گونه ای است که بتوانند تحت اثرات کوانتومی، که در محدوده ی ابعاد نانومتری و کوچکتر بروز می کنند کار کنند. ورود به دنیای محاسبات کوانتومی نیازمند دو پیش زمینه مهم است، نخست باید اصول اساسی و برخی تعابیر مهم مکانیک کوانتومی را به طور دقیق بررسی کرد سپس مفهوم اطلاعات در فیزیک نیز، چه به صورت کلاسیک و چه در معنای جدید کوانتومی آن باید درک شود. بنابراین محاسبات کوانتومی را به عنوان یک زمینه و روش جدید و بسیار کارآمد مطرح می کنند. هر سیستم محاسباتی دارای یک پایه اطلاعاتی است که نماینده ی کوچکترین میزان اطلاعات قابل نمایش، چه پردازش شده و چه خام است. همان طور که در قسمت قبل نیز گفتیم در محاسبات کلاسیک، این واحد ساختاری را بیت می نامیم که گزیده واژه « عدد دو دویی » است زیرا می تواند تنها یکی از دو رقم مجاز صفر و یک را در خود نگه دارد به عبارت دیگر هر یک از ارقام یاد شده در محاسبات کلاسیک، کوچکترین میزان اطلاعات قابل نمایش محسوب می شوند. پس سیستم هایی هم که برای این مدل وجود دارند باید بتوانند به نوعی این مفهوم را عرضه کنند. و در محاسبات کوانتومی هم چنین پایه ای معرفی می شود، که آنرا ( QUBIT ) یا بیت کوانتومی می نامیم. اما این تعریف کیوبیت نیست و باید آنرا همراه با مفهوم و نمونه های واقعی و

فیزیکی درک کرد. در ضمن فراموش نمی کنیم که کیوبیت ها سیستم هایی فیزیکی هستند، نه مفاهیمی انتزاعی و اگر از ریاضیات هم برای توصیف آنها کمک می گیریم تنها بدلیل ماهیت کوانتومی آنها است.

در فیزیک کلاسیک برای نگه داری یک بیت از حالت یک سیستم فیزیکی استفاده می شود. در سیستم های کلاسیکی اولیه ( کامپیوترهای مکانیکی ) از موقعیت مکانی دندانها های چند چرخ دنده برای نمایش اطلاعات استفاده می شد. از زمانیکه حساب دودویی برای محاسبات پیشنهاد شد، سیستم های دو حالتی انتخابهای ممکن برای محاسبات عملی شدند. به این معنی که تنها کافی بود تا سیستمی دو حالت یا دو پیکربندی مشخص، متمایز و بدون تغییر داشته باشد تا بتوان از آن برای این منظور استفاده کرد. به همین جهت، از بین تمام کاندیداها، سیستم های الکتریکی و الکترونیکی برای این کار انتخاب شدند. به این شکل، هر بیت، یک مدار الکتریکی است که یا در آن جریان وجود دارد یا ندارد.

هر بیت کوانتومی یا کیوبیت عبارتست از یک سیستم دو دویی که می تواند دو حالت مجزا داشته باشد. به عبارت فنی تر، کیو بیت یک سیستم دو بعدی کوانتومی با دو پایه به شکل  $0 < \text{و} < 1$  است. البته نمایش پایه ها یکتا نیست، به این دلیل که بر خلاف محاسبات کلاسیک در محاسبات کوانتومی از چند سیستم کوانتومی به جای یک سیستم ارجح استفاده می کنیم.

انتخاب ایده ال برای نمایش کیوبیت استفاده از مفهوم اسپین است که معمولا اتم هیدروژن برای آن به کار می رود، چون در یک اتم هیدروژن هم پروتون وهم الکترون، دارای اسپین می باشد. در اندازه گیری اسپین یک الکترون، احتمال بدست آمدن دو نتیجه وجود دارد: یا اسپین رو به بالاست که آنرا با نشان می دهند و معادل  $0 < \text{است و یا رو پایین است که آن را با نشان می دهیم و معادل با } < 1$  است. بالا یا پایین بودن جهت اسپین در یک اندازه گیری از آنجا ناشی می شود که اگر اسپین اندازه گیری شده در جهت محوری باشد که اندازه گیری را در جهت آن انجام داده ایم، آنرا بالا و اگر در خلاف جهت این محور باشد آنرا پائین می نامیم. شاید بتوان گفت مهم ترین تفاوت بیت و کیوبیت در این دانست که بیت کلاسیک فقط می تواند در یکی از دو حالت ممکن خود قرار داشته باشد در حالیکه بیت کوانتومی می تواند به طور بالقوه در بیش از دو حالت وجود داشته باشد. تفاوت دیگر در اینجاست که هرگاه بخواهیم می توانیم مقدار یک بیت را تعیین کنیم اما اینکار را در مورد یک کیوبیت نمی توان انجام داد.

به زبان کوانتومی یک کیوبیت را با عبارت نشان می دهیم. البته اندازه گیری یک کیوبیت حتما یکی از دو نتیجه ممکن را بدست می دهد. از سوی دیگر اندازه گیری روی سیستم های کوانتومی حالت اصلی آنها را تغییر می دهد. کیوبیت در حالت کلی در یک حالت برهم نهاده از دو پایه ممکن قرار دارد.

اما در اثر اندازه گیری حتما به یکی از پایه ها برگشت می کند. به این ترتیب هر کیوبیت، بیش از اندازه گیری شدن می تواند اطلاعات زیادی را در خود داشته باشد. بر اساس اصل برهم نهی (superposition)، هر سیستم کوانتومی که بیش از یک حالت قابل دسترس دارد، می تواند به طور همزمان در یک ترکیب خاص از آن حالت ها هم قرار داشته باشد. در اصطلاح می گوئیم که سیستم کوانتومی علاوه بر حالت های ناب یک یا چند حالت آمیخته یا برهم نهیده (blend or superposed) نیز دارد. پس اگر یک ساختار حافظه ای  $n$  کیوبیتی داشته باشیم، طبق این اصل، این تعداد می توانند در  $2^n$  پیکربندی متمایز وجود داشته باشند. به این ترتیب یک کامپیوتر کوانتومی این امکان را می یابد که مانند یک کامپیوتر موازی کلاسیک بسیار پر قدرت عمل کند که در یک لحظه روی چندین مسیر اطلاعاتی پردازش می کند. البته مشاهده و متمایز کردن تک تک این محاسبه گرهای کوانتومی غیر ممکن است. چون کامپیوتر کوانتومی با تعداد بسیار زیادی مسیر محاسباتی کار می کند، می توان کاری کرد که این محاسبات با هم تداخل یا برهم تاثیر هم داشته باشند. به عبارتی، محاسباتی که به طور موازی با هم انجام می شوند طبق اصل تداخل می توانند اثر هم را تقویت یا تضعیف کنند. در نتیجه محاسبه ای شبکه ای بوجود می آید که نوعی خاصیت جمعی از تمام محاسبات را نشان می دهد. خاصیت بسیار شگفت انگیز در مکانیک کوانتومی خاصیت درهم تافتگی است. اگر دو یا چند کیوبیت را در برهم کنش با هم قرار دهیم، می توانند برای مدتی در یک حالت کوانتومی مشترک قرار بگیرند به طوری که نتوان آن حالت را به شکل حاصلضربی از حالت های جدا از هم اولیه نشان داد. حالت این واحدهای اطلاعاتی را گنگ یا نادقیق (fuzzy) می نامیم.

این است که یک جفت کیوبیت درهم پیچیده روی یکدیگر درهم تافتگی entanglement یک نتیجه مهم تاثیر همزمانی را می گذارند که به فاصله آن ها از یکدیگر و ماده ای که این فاصله را پرمی کند بستگی ندارد.

## فصل اول

### روش های عددی متفاوت

#### 1-1 RNS (Residue Number System) : سیستم اعداد مانده ای

کامپیوترها امروزه در مبنای ۲ کار می کنند اما در گذشته مبنای کامپیوترها ۱۰ بوده است در حالت کلی هرچه مبنای بالاتر باشد نیاز به هوش بیشتری می باشد .

همه اعمال در کامپیوتر قابل تبدیل به جمع هستند اما در جمع کردن دو عدد باینری وقتی اعداد طولانی می شوند تاخیر انتشار Carry به وجود می آید و راه حل این مشکل RNS می باشد که سرعت پردازنده را بالا می برد . در این روش ما مبحث پیمانه ها را داریم .

قانون پیمانه ها : پیمانه ها باید نسبت به هم اول باشند ( پیمانه ها دو به دو نسبت به هم اول باشند یعنی بزرگترین مقسوم علیه مشترکشان یک باشد ) .

اعداد اول : ۲ و ۳ و ۵ و ۷ و ۱۱ و ۱۳ و .....

بزرگترین عددی که در پیمانه (۹ و ۸ و ۷) می توان نشان داد  $۹ * ۸ * ۷$  می باشد .

در پیمانه (۵ و ۴ و ۳) بزرگترین عدد قابل نمایش  $۵ * ۴ * ۳ = ۶۰$  می باشد

می خواهیم اعداد را در این پیمانه نمایش بدهیم مثلا عدد ۱۲

$$12 = (2, 0, 0)$$

که در آن  $2 = 12 \bmod 5$  و  $0 = 12 \bmod 4$  و  $0 = 12 \bmod 3$  می باشد و رقم نقلی نداریم و جواب باید از پیمانه کوچکتر باشد .

$$14 = (4, 2, 2)$$

$$26 = (1, 2, 2)$$

که در آن ۱ به این صورت محاسبه می شود : جمع ۲ و ۴ عدد ۶ می شود ولی چون پیمانه ۵ است و ۶ بزرگتر از ۵ است پس  $6 - 5 = 1$  می شود چون  $6 \bmod 5 = 1$

نمایش اعداد در پیمانه باعث می شود اعداد کوچکتر شوند مثلا در پیمانه فوق بزرگترین عددی که می توانیم داشته باشیم ۴ است و ۵ هم نمی توانیم داشته باشیم و در ضمن می توانیم محاسبات را موازی انجام دهیم . پس می توان گفت اعداد را به پیمانه می بریم (به حالت RNS) و عمل جمع را انجام داده و دوباره به حالت اولیه می گردانیم

$$\begin{array}{r}
 (5,4,3) \\
 (1,1,0) \quad + \quad 21 \\
 (3,2,0) \quad 18 \\
 \hline
 (4,3,0) \quad 39
 \end{array}$$

$$\begin{array}{r}
 (2,1,2) \quad 17 \\
 (2,1,2) \quad + \quad 17 \\
 \hline
 (4,2,1) \quad 34
 \end{array}$$

در این روش یعنی RNS:

۱- Carry نداریم

۲- اجرای موازی و محاسبات ساده تر

$$\begin{array}{r}
 (.,1,2) \quad * \quad 5 \\
 (1,3,2) \quad 11 \\
 \hline
 (0,3,1) \quad 55
 \end{array}$$

$$\begin{array}{r}
 (3,2,0) \quad - \quad 17 \\
 (0,1,2) \quad 17 \\
 \hline
 (, , -2) \quad 34
 \end{array}$$

برای تبدیل به اندازه پیمانه می توانیم اضافه کنیم  $0+3=3$  ,  $-2+3=1$

$$13=(3,1,1)$$

ضرب و تفریق انجام می شود تقسیم مشکل است .

در تفریق اگر جواب منفی شد به اندازه پیمانه به آن اضافه می کنیم .

اعداد در این پیمانه هدا و در این بازه ها منحصر به فرد هستند .

### ۱-۱-۱ اثبات درستی عمل جمع در RNS:

فرضیات :

$$(m_2, m_1, m_0)$$

$$K = (k_2, k_1, k_0) \begin{cases} k = r_2 m_2 + k_2 \\ K = r_1 m_1 + k_1 \\ K = r_0 m_0 + k_0 \end{cases}$$

$$X = (x_2, x_1, x_0) \begin{cases} x = L_2 m_2 + x_2 \\ X = L_1 m_1 + x_1 \\ X = L_0 m_0 + x_0 \end{cases}$$

حکم قضیه :

$$K+x = (k_2+x_2 \bmod m_2, k_1+x_1 \bmod m_1, k_0+x_0 \bmod m_0)$$

اثبات حکم قضیه :

$$K+x = r_2 m_2 + k_2 + L_2 m_2 + x_2 = (r_2 + L_2) m_2 + k_2 + x_2$$

اگر  $k_2 + x_2$  از  $m_2$  کوچکتر بود که جواب خوش است و اگر  $k_2 + x_2$  از  $m_2$  بزرگتر بود به اندازه پیمانۀ  $m_2$  کم می کنیم و به  $r_2 + L_2$  یک اضافه می کنیم یا همان باقی مانده را حساب می کنیم .

$$K+x = r_1 m_1 + k_1 + L_1 m_1 + x_1 = (r_1 + L_1) m_1 + k_1 + x_1$$

پس وقتی به اندازه پیمانۀ کم می کنیم می شود  $\bmod$

پس اثبات شد که برای جمع دو عدد می توان پیمانۀ ها یا باقی مانده ها را هم جمع کنیم .

تمرین : اثبات انجام شدن ضرب و تفریق - اثبات انجام نشدن تقسیم

پیمانۀ ها را جمع می کنیم، اگر از پیمانۀ کوچکتر بود خودش اما اگر نبود به اندازه  $M$  کم می کنیم به آن طرف ۱ اضافه می کنیم.

تمرین : اثبات ضرب و تقسیم به عنوان تمرین برای شما ( اثبات شدن ضرب و نشدن تقسیم )

۱- اثبات کنید ضرب و تفریق در مبنای RNS درست است و تقسیم درست نیست.

$$x(x_2, x_1, x_0) \rightarrow \begin{cases} x = m_2 k_2 + x_2 \\ x = m_1 k_1 + x_1 \\ x = m_0 k_0 + x_0 \end{cases}$$

$$y(y_2, y_1, y_0) \rightarrow \begin{cases} y = m_2 k'_2 + y_2 \\ y = m_1 k'_1 + y_1 \\ y = m_0 k'_0 + y_0 \end{cases}$$

$$x * y = (m_2 k_2 + x_2) * (m_2 k'_2 + y_2) = m_2^2 k_2 k'_2 + m_2 k_2 y_2 + m_2 k'_2 x_2 + x_2 y_2 = m_2 (m_2 k_2 k'_2 + k_2 y_2 + k'_2 x_2) + x_2 y_2 = m_2 k_2'' + x_2 y_2$$

$$x * y = (m_1 k_1 + x_1) * (m_1 k'_1 + y_1) = m_1^2 k_1 k'_1 + m_1 k_1 y_1 + m_1 k'_1 x_1 + x_1 y_1 = m_1 (m_1 k_1 k'_1 + k_1 y_1 + k'_1 x_1) + x_1 y_1 = m_1 k_1'' + x_1 y_1$$

$$x * y = (m_0 k_0 + x_0) * (m_0 k'_0 + y_0) = m_0^2 k_0 k'_0 + m_0 k_0 y_0 + m_0 k'_0 x_0 + x_0 y_0 = m_0 (m_0 k_0 k'_0 + k_0 y_0 + k'_0 x_0) + x_0 y_0 = m_0 k_0'' + x_0 y_0$$

$$\Rightarrow (x_2 y_2, x_1 y_1, x_0 y_0)$$

$$x(x_2, x_1, x_0) \rightarrow \begin{cases} x = m_2 k_2 + x_2 \\ x = m_1 k_1 + x_1 \\ x = m_0 k_0 + x_0 \end{cases}$$

$$y(y_2, y_1, y_0) \rightarrow \begin{cases} y = m_2 k'_2 + y_2 \\ y = m_1 k'_1 + y_1 \\ y = m_0 k'_0 + y_0 \end{cases}$$

$$x - y = m_2 k_2 + x_2 - (m_2 k'_2 + y_2) = m_2 (k_2 - k'_2) + x_2 - y_2$$

$$x - y = m_1 k_1 + x_1 - (m_1 k'_1 + y_1) = m_1 (k_1 - k'_1) + x_1 - y_1$$

$$x - y = m_0 k_0 + x_0 - (m_0 k'_0 + y_0) = m_0 (k_0 - k'_0) + x_0 - y_0$$

نکته:

در هریک از تفریق های بدست آمده که حاصل عدد منفی شده به اندازه یک پیمانانه به ان اضافه میکنیم.

$$\Rightarrow x - y (x_2 - y_2, x_1 - y_1, x_0 - y_0)$$

$$\frac{x}{y} = (m_2 k_2 + x_2) / (m_2 k'_2 + y_2)$$

این عبارت الزاما برابر با  $\frac{m_2 k_2}{m_2 k'_2} + \frac{x_2}{y_2}$  یا عبارتی که بتوان آن را بصورت  $m_2 k_2'' + \frac{x_2}{y_2}$  نوشت نیست و از طرفی در  $\frac{x_2}{y_2}$  امکان دارد  $y_2$  صفر باشد که حاصل تعریف نشده است و باقیمانده هیچ تقسیمی در ریاضی تعریف نشده نیست بنابراین در دلیل تقسیم در RNS صحیح نیست.

مثال) محدوده (۶۰ تا ۱) یا (۵۹ تا ۰) = ۶۰

(۳ و ۴ و ۵) پیمانه

$$\left. \begin{array}{l} 18 = (3 \text{ و } 2 \text{ و } 0) \\ 17 = (2 \text{ و } 1 \text{ و } 2) \end{array} \right\} \text{forward conversion}$$


---


$$35 \text{ (} 0 \text{ و } 3 \text{ و } 2 \text{)}$$

$$-7 \rightarrow -7 + 30 = 23 \rightarrow (3 \text{ و } 3 \text{ و } 2)$$

$$-8 \rightarrow -8 + 30 = 22 \rightarrow (2 \text{ و } 2 \text{ و } 1)$$

$$-15 \rightarrow \overline{-6} \leftarrow 45 = (0 \text{ و } 1 \text{ و } 0)$$

برای تحت پوشش قرار دادن عدد منفی: هر عدد را +۳۰ می کنیم عدد را یک shift می دهیم، جواب -۳۰ شود تا برگردد به حالت قبلش

(RNS)

عدد را از -۶۰ کم کردیم زیر ۲ تا ۳۰ تا اضافه شده است  $\rightarrow$  -۶۰

راه حل های دیگر:

۱: پیمانه ها زیاد شود اما نسبت به هم اول باشند، (به جای ۳ پیمانه ۴ پیمانه)

۲: اعداد بزرگتر شود.

اگر پیمانه ها زیاد تر شود موازی سازی بیشتر می شود و موازی می توان کار کرد اما: (۱) عدد بردن به پیمانه ها سخت تر می شود (۲) برگرداندن آنها هم سخت تر می شود.

اما اگر اعداد بزرگتر شوند موازی سازی کمتر می شود محدوده اعداد بزرگ می شود محاسبات سخت تر می شود اما forward و backward کم می شوند و باید بر اساس نیاز انتخاب کرد و کار را انجام داد.



## ۱-۱-۲ برگرداندن اعداد از پیمانہ (RNS) به مبنای ۱۰

(روش چینی): backward

زمانی که پیمانہ را داشته باشیم و از ما عدد مربوط به آن را بخواهند از این روش استفاده می کنیم. برای این کار به صورت زیر عمل می کنیم:

اون چه عددی که باقی مانده بر ۵ ← ۲ باقی مانده به ۴ ← ۲ و باقی مانده به ۳ ← ۱ می شود؟

(۳ و ۴ و ۵) پیمانہ

ضرب ۴ و ۳

$$1 = (1 \text{ و } 1 \text{ و } 1)$$

$$? = (1 \text{ و } 0 \text{ و } 0)$$

$$? = (0 \text{ و } 1 \text{ و } 0)$$



$$\left\{ \begin{array}{l} 4 \times 3 = 12 \\ 12 \times 2 = 24 \\ 12 \times 3 = \underline{36} \end{array} \right.$$



$$\left\{ \begin{array}{l} 5 \times 3 = 15 \\ 15 \times 2 = 30 \\ 15 \times 3 = \underline{45} \end{array} \right.$$

$$? = (0 \text{ و } 0 \text{ و } 1)$$



$$\left\{ \begin{array}{l} 5 \times 4 = 20 \\ 20 \times 2 = \underline{40} \end{array} \right.$$

$$(1 \text{ و } 0 \text{ و } 0) = 36$$

$$(0 \text{ و } 1 \text{ و } 0) = 45$$

$$(0 \text{ و } 0 \text{ و } 1) = 40$$

(مثال)

$$19 = (4 \text{ و } 3 \text{ و } 1) = 4 \times (1 \text{ و } 0 \text{ و } 0) + 3 \times (0 \text{ و } 1 \text{ و } 0)$$

$$+ 1 \times (0 \text{ و } 0 \text{ و } 1) \rightarrow 4 \times 36 + 3 \times 45 + 1 \times 40 = 319$$

$$60 = 19 \bmod 319$$

(مثال)

$$60 = 27 \bmod = (2 \text{ و } 3 \text{ و } 0) \rightarrow 2 \times (1 \text{ و } 0 \text{ و } 0) + 3 \times (0 \text{ و } 1 \text{ و } 0) = 2 \times 36 + 3 \times 45 = 207$$

مثال) عدد ( ۱۰۲ و ۹۱ ) را به پیمانه برده و backward انجام دهید؟ ( ۷ و ۸ و ۹ ) پیمانه

$$۹۱ + ۱۰۲$$

$$(۹ \times ۸ \times ۷) = (۹ \text{ و } ۸ \text{ و } ۷) \text{ پیمانه}$$

$$۹۱ (۰ \text{ و } ۳ \text{ و } ۱)$$

$$۷۲ \times ۴ = \underline{۲۸۲}$$

$$\underline{۱۰۲ (۳ \text{ و } ۶ \text{ و } ۴)}$$

$$۱۹۳ (۳ \text{ و } ۶ \text{ و } ۴)$$

$$\left\{ \begin{array}{l} (۰ \text{ و } ۰ \text{ و } ۱) = ۹ \times ۸ = ۷۲ \\ ۷۲ \times ۲ = ۱۴۴ \\ ۷۲ \times ۳ = ۲۱۶ \end{array} \right.$$

$$(۱ \text{ و } ۰ \text{ و } ۰) = ۸ \times ۷ = ۵۶$$

(باقیمانده بر ۹ عدد ۱ شود)

$$\left\{ \begin{array}{l} ۵۶ \times ۲ = ۱۱۲ \\ ۵۶ \times ۳ = ۱۶۸ \\ ۵۶ \times ۴ = ۲۲۴ \\ ۵۶ \times ۵ = \underline{۲۸۰} \end{array} \right.$$

$$(۰ \text{ و } ۱ \text{ و } ۰) = ۹ \times ۷ = ۶۳$$

(باقیمانده بر ۸ عدد ۱ شود)

$$\left\{ \begin{array}{l} ۶۳ \times ۲ = ۱۲۶ \\ ۶۳ \times ۳ = ۱۸۹ \\ ۶۳ \times ۴ = ۲۵۲ \\ ۶۳ \times ۵ = ۳۱۵ \\ ۶۳ \times ۶ = ۳۷۸ \\ ۶۳ \times ۷ = \underline{۴۴۱} \end{array} \right.$$

## مزایا و معایب این روش :

محاسبات در این مبنا بسیار خوب است زیرا آسان می باشد و موازی سازی داریم . وقتی موازی سازی پیش می آید به آن علاقه مند می شویم چون سرعت زیاد می شود و کاری که قبلا در ۱ ثانیه انجام میشد حالا در نیم ثانیه انجام میشود. و سرعت افزایش می یابد.

این محاسبات یکبار برای همیشه انجام می شود و همیشگی ست. پس forward خوب است و محاسبات عالی انجام میشود اما برای backward خوب نیست.

تمرین : برای backward یک روش غیر چینی ارائه دهید. (تمرین خارج از امتحان، جواب در پیوست ۱)

## کاربرد:

در محاسبات کامپیوتری مورد استفاده قرار می گیرد. از این روش برای ارسال داده ها هم میتوان استفاده کرد. داده ها را کوچک کرده (RNS) و ارسال می کنیم.

کاهش توان مصرفی و کاهش خطا را داریم. به دلیل کوچک شدن در شبکه های بیسیم استفاده می شود. اما در multi casting استفاده شده است.

تمرین: اعداد منفی را چگونه به مبنای RNS میبریم؟

برای این کار باید به جای استفاده از سیستم RNS از SNS (signed residue System) استفاده کرد. در این سیستم اعداد دارای علامت هستند که برای ایجاد علامت یک بیت یا صفر می شود یا یک. برای جمع و تفریق اعداد علامت دار باید از روش زیر استفاده کرد.

برای ایجاد سیستم SNS باید ابتدا بازه M بازه اعداد را به دو قسمت تقسیم کنیم:

$$\left\{ \begin{array}{ll} X < M/2 & \text{اعداد مثبت} \\ X \geq M/2 & \text{اعداد منفی} \end{array} \right.$$

که برای بازه منفی از فرمول زیر استفاده می کنیم

$$(X - M) \bmod M = X \bmod M$$

برای محاسبه بازه منفی باید ابتدا Complement عدد را به صورت زیر برای  $X \bmod M$  محاسبه کنیم:

$$X^c = (M - X) \bmod M$$

حال  $X = [x_i]$  و چون  $[x_i] = X \bmod m_i$  پس در نتیجه می توان گفت که Complement عدد X می تواند Complement عدد  $x_i$  هم بشود. پس:

$$x_i^c = (m_i - x_i) \bmod m_i$$

$$[x_i]^c = [x_i^c] = X^c$$

پس عدد  $-x_i$  با عدد  $m_i - x_i$  برابر است

مثال: تفریق دو عدد ۸ و ۹ از هم در مبنای ۵, ۳, ۲ که برای این کار به صورت زیر عمل می کنیم:

$$8[3,2,0]$$

$$9[4,0,1]$$

بدست آوردن complement عدد ۸

$$(8)^c = [2,1,0], 5 - 3, 3 - 2, (2 - 0) \bmod 2$$

برای محاسبه ۸-۹ ابتدا Complement عدد ۹ را محاسبه می کنیم

$$(9)^c = [1,0,1]$$

سپس ۸ را با عدد Complement عدد ۹ جمع می کنیم:

$$\begin{array}{r} 8 = 8 = [3,2,0] \\ -9 = (9)^c = +[1,0,1] \\ \hline -1 \quad \quad \quad [4,2,1] = 29 \text{ or } -1 \end{array}$$

سپس عدد ۲۹ که بدست آمد ابتدا بررسی می کنیم که در چه بازه ای قرار دارد. اگر در بازه دوم است باید آن را از بازه بالا کم کنیم تا جواب صحیح بدست آید. به صورت زیر:

$$29 - 30 = -1$$

جواب صحیح عدد -۱ است.

## ۲-۱ اعداد در مبنای RTIT :

در جلسات قبل یکی از تکنیک های جدیدی که مطرح شده است جهت افزایش کارایی سیستم در حینی که سیستم کامپیوتر ما بر مشکلات قانون moor قائل شود را گفتیم. (RNS)

ما در کامپیوتر مبنای دودویی را یاد گرفتیم. که از صفر و یک تشکیل شده بود. ولی کامپیوترها از ابتدا صفر و یک نبودند بلکه از هزاران لامپ خلا تشکیل شده بودند و در مبنای ۱۰ بودند!!

مبناهای بزرگتر پیچیده تر هم هستند. و هوش بالاتری می خواهند. مثلا مبنای ۲ ساده تر می باشد. در حقیقت ساختن AND, OR در آن ساده تر از مبناهای بالاتر می باشد.

به تازگی به این نتیجه رسیدند که مبنا را از ۲ کمی بالاتر ببرند. یعنی به جای اینکه فقط ۰ و ۱ باشد، مبنای ۳ داشته باشیم، یعنی ۰ و ۱ و ۲. (Ternary Digit).

چون ساختن قطعات بر مبنای ۳ سخت بود، پس مبنای جدیدی به نام RTit به وجود آمد. (۰ و ۱ و ۲). که البته مبنا همان ۲ می باشد و آن عدد اضافه را Redundant میگویند.

سیستم هایی که Redundant دارند، قابلیت اطمینان بالاتری هم دارند. یعنی بیت اضافی که وجود دارد، در صورتی که بیت های اصلی خراب شوند، جایگزین می شود.

مثال:

$$(1000)_{RTit} = 0 \cdot 2^0 + 0 \cdot 2^1 + 0 \cdot 2^2 + 1 \cdot 2^3 = 8$$

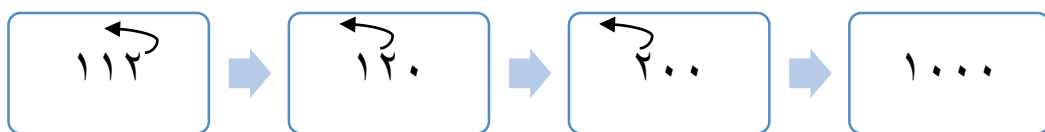
$$(200)_{RTit} = 0 \cdot 2^0 + 0 \cdot 2^1 + 2 \cdot 2^2 = 8$$

دو مساله پیش می آید:

۱. برای عدد ۸ حالت های مختلفی به وجود می آید. و یک مزیت می باشد چون اگر خطایی رخ دهد ما می توانیم نتیجه را به حالت های مختلف ارائه دهیم.

۲. به یک طریقی این اعداد را می توان به دست آورد. یعنی اتفاقی همه ی این ها ۸ نشدند!!

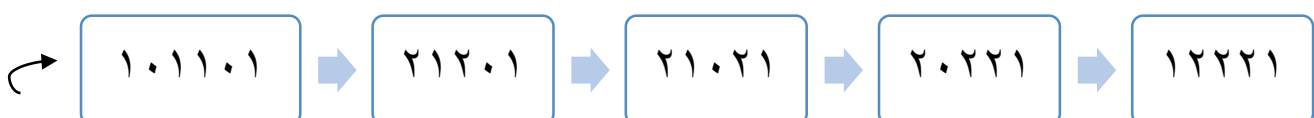
### ۱-۲-۱ روش بدست آوردن اعداد در RTit :



توضیح: عدد ۲ می تواند از رقم پایین تبدیل به ۱ شود. یک مرحله بالا رود.

اگر عدد ۱ از یک مرحله بخواید پایین رود چه می شود؟ ۲ می شود. (یک مرحله بالا رود ۱ و یک مرحله پایین تر رود برعکس می شود).

مثالی دیگر:



نکته: برای تبدیل عددی از مبنای ۱۰ به RTit، عدد باید از مبنای ۱۰ اول به مبنای ۲ تبدیل شود و سپس به RTit تبدیل شود. مبناهای جدید، به تبع، or, and های جدیدی هم دارند. همانطور که می دانید and در مبنای ۲ به شرح زیر است:

Binary digit		and	Or
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	1

در RTit، مینیمم اعداد می شود and و ماکسیمم اعداد می شود or.

RTit		AND	OR
0	0	0	0
0	1	0	1
0	2	0	2
1	0	0	1
1	1	1	1
1	2	1	2
2	0	0	2
2	1	1	2
2	2	2	2

NOT در این مبنا به شکل زیر است:

NOT	M	N	P
0	2	2	2
1	0	1	2
2	0	0	0

در این جدول، P (مثبت) و N (نرمال) و M (منفی) می باشد.

NAND در این مبنا به شکل زیر است: در حقیقت ستون AND را NOT میکنیم. که برای هر بیت ۳ حالت NOT پیش می آید.

## NAND

	RTit	AND	M	N	P
0	0	0	2	2	2
0	1	0	2	2	2
0	2	0	2	2	2
1	0	0	2	2	2
1	1	1	0	1	2
1	2	1	0	1	2
2	0	0	2	2	2
2	1	1	0	1	2
2	2	2	0	0	0

تمرین : NOR, XOR, XNOR را در مبنای RTit بنویسید.

جواب های تمرین به شرح زیر می باشد:

NOR در این مینا به شکل زیر است: در حقیقت ستون OR را NOT می کنیم. که برای هر بیت ۳ حالت NOT پیش می آید.

## NOR

	RTit	OR	M	N	P
0	0	0	2	2	2
0	1	1	0	1	2
0	2	2	0	0	0
1	0	1	0	1	2
1	1	1	0	1	2
1	2	2	0	0	0
2	0	2	0	0	0
2	1	2	0	0	0
2	2	2	0	0	0

فرمول XOR:  $A'.B + A.B'$

A'	B'	AB'	A'B	A'B'	A XOR B
2 2 2	2 2 2	0 0 0	0 0 0	2 2 2 2 2 2 2 2 2	0 0 0 0 0 0 0 0 0
2 2 2	0 1 2	0 0 0	1 1 1	0 1 2 0 1 2 0 1 2	1 1 1 1 1 1 1 1 1
2 2 2	0 0 0	0 0 0	2 2 2	0 0 0 0 0 0 0 0 0	2 2 2 2 2 2 2 2 2
0 1 2	2 2 2	1 1 1	0 0 0	0 0 0 1 1 1 2 2 2	1 1 1 1 1 1 1 1 1
0 1 2	0 1 2	0 1 1	0 1 1	0 0 0 0 1 1 0 1 2	0 1 1 1 1 1 1 1 1
0 1 2	0 0 0	0 0 0	0 1 2	0 0 0 0 0 0 0 0 0	0 1 2 0 1 2 0 1 2
0 0 0	2 2 2	2 2 2	0 0 0	0 0 0 0 0 0 0 0 0	2 2 2 2 2 2 2 2 2
0 0 0	0 1 2	0 1 2	0 0 0	0 0 0 0 0 0 0 0 0	0 0 0 1 1 1 2 2 2
0 0 0	0 0 0	0 0 0	0 0 0	0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0

فرمول XNOR:  $A'.B + A.B'$

در حقیقت جدول XOR، باید NOT شود.

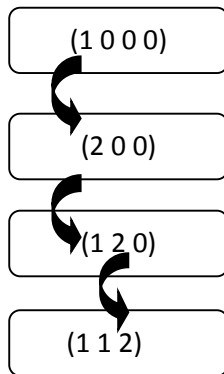
جدول آن به صورت زیر می باشد:

A XOR B	A XNOR B
0 0 0 0 0 0 0 0 0	2 2 2 2 2 2 2 2 2
1 1 1 1 1 1 1 1 1	0 1 2 0 1 2 0 1 2
2 2 2 2 2 2 2 2 2	0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1	0 0 0 1 1 1 2 2 2
0 1 1 1 1 1 1 1 1	1 1 1 1 1 1 1 1 2
0 1 2 0 1 2 0 1 2	1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2	0 0 0 0 0 0 0 0 0
0 0 0 1 1 1 2 2 2	0 1 1 1 1 1 1 1 1
0 0 0 0 0 0 0 0 0	2 2 2 2 2 2 2 2 2



### ۱-۲-۱ روش های به دست آوردن نمایش های مختلف اعداد در مبنای RTIT:

در مبنای RTIT برای یک عدد می توان نمایش های مختلفی بدست آورد. برای این کار از سمت چپ پیمانه یک ارزش پایین تر میرویم، برای مثال اگر عدد ۱ بخواید یک ارزش پایین بیاید به عئی ۲ تبدیل می شود، اگر عدد ۲ بخواید یک ارزش پایین برود، چون عدد ۲ برابر با ۱+۱ است، یکی از یک هایش باقی می ماند و ۱ بعدی یک ارزش پایین می رود و به ۲ تبدیل می شود.



مثال:

### ۱-۳ اعداد در مبنای BR Tit (Balanced RTit)

در این قسمت مبنا به صورت زیر است:  $\{-1, 0, 1\}_2$

اما برای  $-1$ ، از این نمایش استفاده می کنیم:  $\{\bar{1}, 0, 1\}$

جدول نمایش اعداد در مبنای BRTit:

مبنای ۱۰	مبنای BRTit	مبنای ۱۰	مبنای BRTit
7	1 1 1	-7	$\bar{1}\bar{1}\bar{1}$
6	0 1 1	-6	$\bar{1}\bar{1}0$
5	1 0 1 1 1 $\bar{1}$	-5	$\bar{1}0\bar{1}$ $\bar{1}\bar{1}1$
4	1 0 0	-4	$\bar{1}00$
3	0 1 1 1 0 $\bar{1}$ 1 $\bar{1}$ 1	-3	$0\bar{1}\bar{1}$ $\bar{1}01$ $\bar{1}1\bar{1}$
2	0 1 0 1 $\bar{1}$ 0	-2	$0\bar{1}0$ $\bar{1}10$
1	0 0 1 0 1 $\bar{1}$ 1 $\bar{1}$ $\bar{1}$	-1	$00\bar{1}$ $0\bar{1}1$ $\bar{1}11$
0	0 0 0		

تمرین: جدول ۴ بیتی BRTit را بنویسید.

جداول AND, OR, NOT برای اعداد در مبناى BRTit:

	AND	OR
$\bar{1} \bar{1}$	$\bar{1}$	$\bar{1}$
$0 \bar{1}$	$\bar{1}$	0
$1 \bar{1}$	$\bar{1}$	1
$\bar{1} 0$	$\bar{1}$	0
$0 0$	0	0
$0 1$	0	1
$\bar{1} 1$	$\bar{1}$	1
$1 0$	0	1
$1 1$	1	1

جدول NOT:

	M	N	P
$\bar{1}$	1	1	1
0	$\bar{1}$	0	1
1	$\bar{1}$	$\bar{1}$	$\bar{1}$

فواید نمایش اعداد به صورت BRTit:

- نمایش اعداد منفی آسان تر است.
- افزونگی خواهیم داشت.
- محدوده گسترده تر با تعداد خانه کمتر نمایش داده می شود.



## فصل دوم

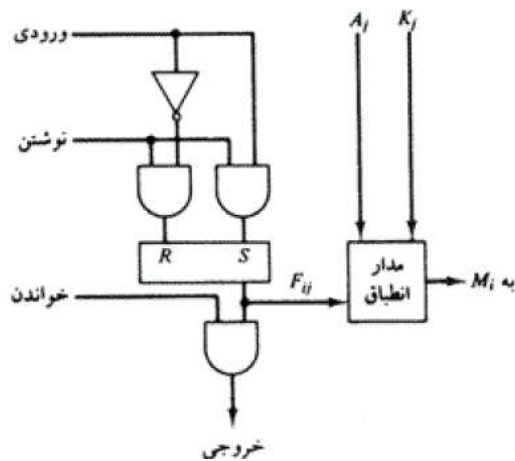
### سازمان حافظه

#### ۲-۱ حافظه تداعیگر:

سرعت بالایی دارد و سرعت جستجو در آن بالاست و به همین دلیل در حافظه های CASH از آنها استفاده می شود و اشکال های آن مدارات پیچیده و بزرگ و هزینه بالای آن است

#### ۲-۱-۲ بررسی مدار سخت افزاری حافظه تداعیگر

مدار هر بیت آن بصورت زیر است:



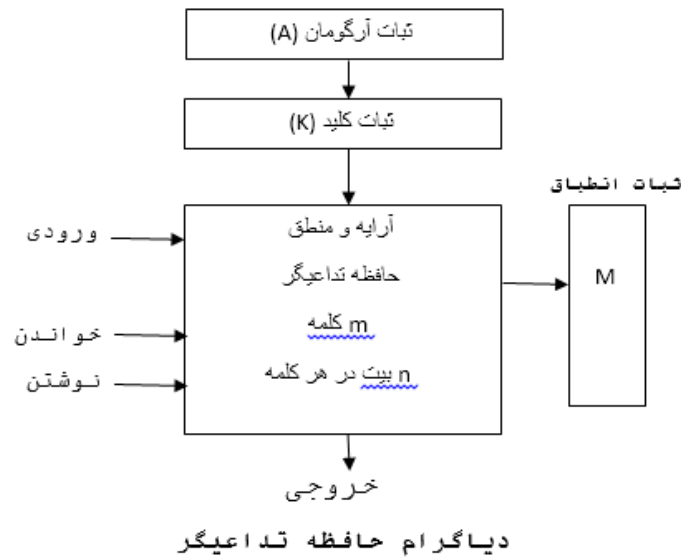
#### ۲-۱-۲ مدار انطباق

خروجی فلیپ فلاپ بصورت زیر است:

S	R	Q
0	0	بدون تغییر
0	1	0
1	0	1
1	1	نامعتبر

برای ساختن حافظه دینامیک از خازن استفاده می کنیم و برای ساختن حافظه استاتیک از لچ و فلیپ فلاپ استفاده می شود و ساختن حافظه استاتیک از نظر منطقی آسانتر ولی هزینه طراحی حافظه دینامیک کمتر است.

## ساختار سخت افزاری:



ثبات  $A$ :  $n$  بیتی است و داخل آن داده مورد جستجو قرار می گیرد.

ثبات  $K$ : یک ثبات  $n$  بیتی است و بیت های متناظر با بیت های این ثبات که برابر یک است در جستجو مهم هستند و بیت های متناظر این ثبات که صفر هستند جستجو نمی شوند در واقع ثبات ماسک می باشد.

مثال عددی: فرض کنید که ثبات آرگومان  $A$  و ثبات کلید  $K$  دارای آرایش بیتی زیر باشند تنها سه بیت سمت چپ  $A$  با کلمات مقایسه زیرا  $K$  در این مکان ها دارای ۱ است.

$A$     ۱۰۱ ۱۱۱۱۰۰

$K$     ۱۱۱ ۰۰۰۰۰۰

انطباق ندارد    ۱۰۰ ۰۰۱۱۱۱    کلمه ۱

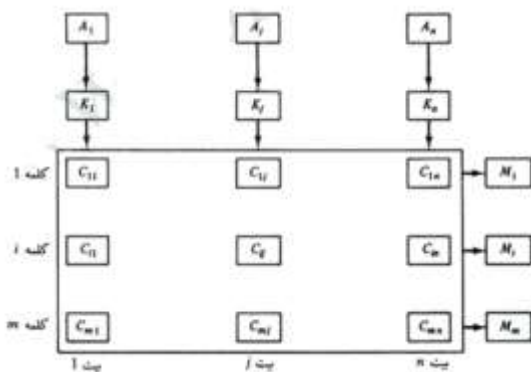
انطباق دارد    ۱۰۱ ۰۰۰۰۰۱    کلمه ۲

کلمه ۲ با میدان آرگومان پوشش نیافته تطبیق دارد زیرا سه بیت سمت چپ آرگومان و کلمه برابرند.

نکته: در حافظه تداعیگر معمولا خطوط آدرس مطرح نمی شود.

حافظه را می توان مانند یک آپارتمان چند طبقه در نظر گرفت که هر طبقه دارای تعدادی مساوی اتاق است که هر کدام از اتاق ها یک بیت هستند اگر تعداد طبقات  $m$  و تعداد اتاق ها در هر طبقه  $n$  بیت

باشد یعنی حافظه دارای  $m$  کلمه  $n$  بیتی است و برای آن یک ثبات انطباق  $m$  بیتی در نظر می گیریم که شکل آن بصورت زیر می باشد:



رابطه بین حافظه و ثبات های خارجی در حافظه تداعیگر در شکل بالا نشان داده شده است سلول های آرایه با حرف  $C$  و دو اندیس مشخص شده اند اولین اندیس شماره کلمه و دومین اندیس مکان بیت را در کلمه مشخص می کند و بیت  $A_j$  در ثبات آرگومان با همه بیت های ستون  $j$  از آرایه مقایسه می شود بشرطی که  $K_j=1$  باشد و این عمل برای همه ستونهای  $j=1,2,\dots,n$  انجام می شود. اگر بین همه بیت های پوشش نیافته آرگومان و بیت های کلمه تطابق وجود داشته باشد بیت متناظر در ثبات انطباق یعنی  $M_i$  برابر ۱ می شود و اگر یک یا چند بیت پوشش نیافته و کلمه مطابقت نداشته باشند  $M_i$  برابر صفر می شود.

ثبات انطباق بصورت زیر طراحی می شود:

$$\begin{cases} x_j & \text{اگر } k_j = 1 \\ 1 & \text{اگر } k_j = 0 \end{cases}$$

رابطه  $M_i$  را بصورت زیر میتوان بدست آورد :

ابتدا بیت های کلید را نادیده میگیریم و آرگومان  $A$  را با بیت های ذخیره شده در سلول کلمه مقایسه می کنیم کلمه  $i$  برابر با آرگومان  $A$  است اگر به ازای  $j=1,2,\dots,n$  داشته باشیم  $A_j=F_{ij}$  و آن را می توان بصورت زیر نشان داد:

$$X_j = A_j F_{ij} + A'_j F'_{ij}$$

$X_j$  زمانی که دو بیت  $A_j$  و  $F_{ij}$  برابر باشند یک می شود و در غیر اینصورت صفر خواهد شد برای اینکه کلمه  $i$  و آرگومان  $A$  برابر باشد باید تمام  $X_j$  ها یک شود این وضعیت شرط یک شدن بیت انطباق  $M_i$  است یعنی عبارت زیر:

$$M_i = X_1 X_2 X_3 \dots X_n$$

که برابر AND کردن جفت بیت های انطباق یافته است.

اکنون بیت  $K_j$  را در مدار اعمال می کنیم که بصورت زیر می باشد:

$$X_j + K'_j =$$

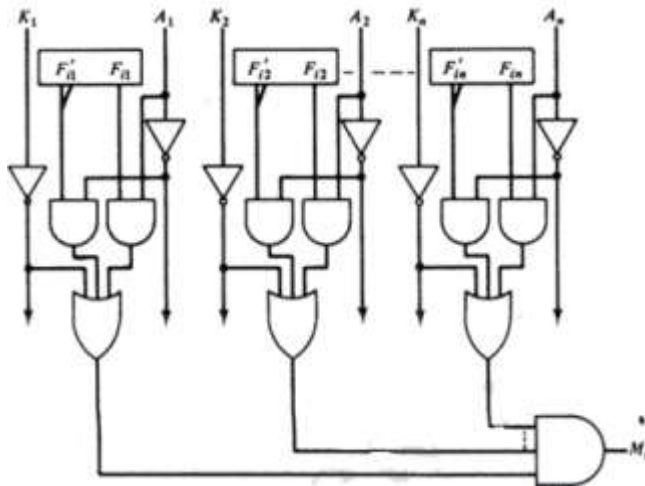
و فرمول بیت انطباق را می توان بصورت زیر نوشت:

$$M_i = (X_1 + K'_1)(X_2 + K'_2) \dots (X_n + K'_n)$$

که با جاگذاری تابع بول آن بصورت زیر درمی آید:

$$M_i = \prod_{j=1}^n (A_j F_{ij} + A'_j F'_{ij} + K'_j)$$

و شکل مدار بصورت زیر است:



قیمت حافظه تداعیگر بیشتر از حافظه معمولی است و سرعت جستجوی آن بالاتر است

## ۲-۲ حافظه Cache

چون سرعت حافظه اصلی نسبت به سرعت CPU پایین تر است از یک حافظه مابین حافظه اصلی و CPU بنام کش قرار می گیرد که سرعت بالایی دارد و هزینه ساخت آن نیز بالا است بنابراین ظرفیت آن پایین و تنها بلوک هایی از حافظه اصلی درون کش قرار می گیرد و CPU در صورت عدم وجود در کش به حافظه اصلی مراجعه می کند و وجود کش به دو دلیل محلیت مکانی و محلیت زمانی است.

### مبانی Caching

Caching یک تکنولوژی استفاده شده برای زیر سیستم های حافظه، در کامپیوتر است. مهمترین هدف یک Cache افزایش سرعت و عملکرد کامپیوتر بدون تحمیل هزینه های اضافی برای تهیه سیستم است. با استفاده از Cache عملیات کاربران با سرعت بیشتری انجام شود.

وقتی پردازنده نیاز دارد که داده ای را بخواند یا بنویسد، ابتدا چک می کند که در Cache موجود است یا نه، این کار به وسیله مقایسه آدرس مکان حافظه با همه تگ های موجود در Cache ((که ممکن است حاوی آدرس باشد)) صورت می پذیرد.

اگر پردازنده آدرس مکان مورد نظر حافظه را در Cache بیابد، می‌گوییم که یک "برخورد Cache" رخ داده، در غیر این صورت گوییم که یک "خطای Cache" روی داده است. در صورت برخورد پردازنده به سرعت داده‌ها را از خط Cache می‌خواند یا می‌نویسد. نسبتی از دسترسی‌ها که منجر به برخورد می‌شود را "نرخ برخورد" گویند، که مقیاسی برای اندازه‌گیری کارایی یک الگوریتم یا برنامه می‌باشد. در صورت بروز خطا، Cache مدخلی دیگر را در نظر می‌گیرد.

### ساختمان مدخل Cache:

مداخل سطری Cache معمولاً ساختاری این چنینی دارند:

Valid Bit	Data Blocks	Tag
-----------	-------------	-----

بلوک‌های داده حاوی داده‌های واکنشی شده از حافظه اصلی می‌باشند. بیت اعتبار مشخص می‌کند که مدخل مذکور، داده معتبر دارد یا نه.

دو مثال ساده:

فرض کنیم که شما هر روز به رستوران می‌روید. هر روز راس ساعت ۵ بعد از ظهر سفارش غذا می‌دهید. هر روز ۴ نوع غذا را به ترتیب خاص سفارش می‌دهید. راس ساعت ۵ همبرگر سفارش می‌دهید، گارسون سفارش شما رو بررسی می‌کند، به آشپزخونه می‌رود، بعد از یک دقیقه همبرگر را برای شما می‌آورد، شما همبرگر را خورده و سفارش سوسیس می‌دهید. مجدداً سفارش توسط گارسون به آشپزخانه منتقل شده، و بعد از یک دقیقه غذا آماده می‌شود. به همین ترتیب شما ۳ غذای دیگه را سفارش داده، و برای هر غذا ۱ دقیقه معطل می‌شوید. خوب شما هرروز همین غذاها را سفارش داده، و برای آماده شدن هر غذا ۱ دقیقه معطل می‌شوید. گارسون با خودش فکر می‌کند که برای اینکه هم خودش کمتر کار کند و هم شما کمتر معطل شوید، ۱ میز آماده دیگر از غذاهای شما را تهیه کند، و بلافاصله بعد از سفارش شما غذا را روی میزتان قرار دهد. در اینجا گارسون "Bus"، آشپزخانه "Ram"، و میز آماده "Cache" در نظر گرفته می‌شوند. بعد از چند روز شما همبرگر را می‌خورید، طبق عادت گارسون برای شما سوسیس می‌آورد، اما شما می‌گویید که امروز پیتزا می‌خواهم، اینجا گارسون مجدداً مجبور می‌شود که ۱ دقیقه شما را در انتظار نگه دارد، تا پیتزا را برایتان بیاورد. در اینجا گارسون میز دومی را تهیه می‌کند که بر اساس انتخاب های دوم شما چیده شده است. بدین ترتیب شما اگر غذایی را سفارش دهید که در میز اول نباشد، اما در میز دوم باشد بلافاصله غذا را میل می‌کنید و معطل نمی‌شوید. میز دوم در اینجا Cache سطح دو (Cache L2) است.

کنتابداری را در نظر بگیرید که در یک کتابخانه مسئول تحویل کتاب به متقاضیان است. فرض کنید در سیستم فوق ((درخواست و تحویل کتاب)) از مفهوم Cache استفاده نمی‌گردد. اولین متقاضی کتابی را



درخواست می نماید ((فرض شده است که متقاضی خود نمی تواند مستقیماً کتاب مورد نظر را از قفسه مربوطه بردارد))، کتابدار، کتاب مورد نظر را از قفسه مربوطه پیدا، و آن را به متقاضی تحویل می دهد. متقاضی پس از ساعاتی مراجعه، و کتاب را تحویل می دهد. کتابدار، کتاب تحویلی را مجدداً در قفسه مربوطه قرار می دهد. پس از لحظاتی یک متقاضی دیگر مراجعه و همان کتاب قبلی را درخواست می نماید، کتابدار مجدداً می بایست به بخش مربوطه در کتابخانه مراجعه، و پس از بازیابی کتاب، آن را در اختیار متقاضی دوم قرار دهد. همانگونه که ملاحظه می گردد، کتابدار مکلف است برای تحویل هر کتاب (حتی کتاب هائی که فرکانس استفاده از آنان توسط متقاضیان زیاد باشد)، به بخش مربوطه مراجعه، و پس از یافتن کتاب آن را در اختیار متقاضیان قرار دهد. فرض کنید بخشی را با ظرفیت حداکثر ده کتاب در مجاورت (نزدیکی) کتابدار آماده نمائیم. کتاب هائی که توسط متقاضیان برگردانده می شود، در بخش فوق ذخیره خواهند شد. مثال فوق را با در نظر گرفتن سیستم Cache ایجاد شده برای کتابدار، مجدداً دنبال می نمائیم. در ابتدای فعالیت روزانه، بخش Cache خالی بوده و هنوز در آن کتابی قرار نگرفته است. اولین متقاضی مراجعه، و کتابی را درخواست می نماید. کتابدار می بایست به بخش مربوطه مراجعه، و کتاب را از قفسه مربوطه برداشته، و در اختیار متقاضی قرار دهد. متقاضی پس از تحویل کتاب، چند ساعت بعد مراجعه، و کتاب را تحویل کتابدار خواهد داد. کتابدار، کتاب تحویلی را در بخش پیش بینی شده، برای Cache قرار می دهد. لحظاتی بعد متقاضی دیگر مراجعه، و درخواست همان کتاب را می نماید. کتابدار در ابتدا بخش مربوط به Cache را جستجو، و در صورت یافتن کتاب، آن را به متقاضی تحویل خواهد داد. در این حالت ضرورتی برای مراجعه کتابدار به بخش و قفسه های مربوطه نخواهد بود. در روش فوق زمان تحویل کتاب به متقاضی، بهبود چشمگیری پیدا خواهد کرد. در صورتی که کتاب درخواستی توسط متقاضی، در بخش Cache کتابخانه نباشد، چه اتفاقی خواهد افتاد؟ در ابتدا مدت زمانی صرف خواهد شد که کتابدار به این اطمینان برسد که کتاب درخواستی در بخش Cache موجود نمی باشد (جستجو)، یکی از چالش های اصلی در رابطه با طراحی Cache، به حداقل رساندن زمان جستجو در Cache است. سخت افزارهای جدید، زمان فوق را به صفر نزدیک کرده اند. پس از حصول اطمینان از عدم وجود کتاب در بخش Cache، کتابدار می بایست با مراجعه به بخش مربوطه آن را انتخاب و در ادامه در اختیار متقاضی قرار دهد.

با توجه به دو مثال فوق، چندین نکته مهم در رابطه با Cache استنباط می گردد:

تکنولوژی Cache، استفاده از حافظه های سریع ولی کوچک، بمنظور افزایش سرعت یک حافظه کند، ولی با حجم بالا است.

زمانیکه از Cache استفاده می گردد، در ابتدا می بایست محتویات آن به منظور یافتن اطلاعات مورد نظر بررسی گردد. فرآیند فوق را Cache Hit می گویند. در صورتی که اطلاعات مورد نظر در Cache موجود

نباشند (Cache miss)، کامپیوتر می بایست در انتظار تامین داده های خود از حافظه اصلی سیستم (حافظه ای کند ولی با حجم بالا) باشد.

اندازه Cache محدود بوده سعی می گردد که ظرفیت فوق حتی المقدور زیاد باشد، ولی به هر حال اندازه آن نسبت به رسانه های ذخیره سازی دیگر بسیار کم است.

این امکان وجود خواهد داشت که از چندین لایه Cache استفاده گردد.

### ۱-۲-۲ محلیت زمانی و مکانی حافظه کش

همانطور که قبلا گفته شده دو دلیل برای استفاده از حافظه وجود دارد؛ محلیت زمانی و محلیت مکانی.

#### ۱) محلیت زمانی (Temporal locality):

داده ها و دستوراتی که اخیرا مورد استفاده قرار گرفته اند، مجددا مورد رجوع قرار خواهند گرفت.

#### ۲) محلیت مکانی (Spatial locality):

داده ها یا دستوراتی که در خانه های نزدیک به هم ((در حافظه)) قرار دارند، با فاصله زمانی اندکی از هم مورد رجوع قرار خواهند گرفت.

### ۲-۲-۲ کارایی cache

کارایی حافظه cache به این جمله بستگی دارد :

پیش بینی محلیت زمانی و مکانی که ما داشتیم، چقدر موفق بوده است؟؟؟

حالا مولفه های زیر را جهت محاسبه کارایی تعریف میکنیم:

متوسط زمان خواندن از حافظه Cache  $T_{Avg} =$

$T_{Cache} =$  متوسط زمان خواندن از Cache

$T_{RAM} =$  متوسط زمان خواندن از RAM

$P =$  (Hit ratio) = احتمال اینکه داده در داخل کاشه باشد.

$(1-P) =$  (Miss ratio) = احتمال اینکه داده در داخل کاشه نباشد.

در ضمن :

**$T_{RAM} > T_{Cache}$**

زیرا حافظه cache سریع تر از RAM می باشد.

$$T_{Avg} = P T_{Cache} + (1-p)(T_{RAM} + T_{Cache})$$

اول به سراغ Cache می رود، اگر داده مورد نظر آن جا بود که می خواند و اگر که نبود، به حافظه RAM مراجعه و از آنجا می خواند و به Cache می دهد و دوباره از آن می خواند. حالا چرا دوباره داخل Cache می گذارد؟؟؟ به دلیل محلیت زمانی و محلیت مکانی.

$$P T_{Cache} + T_{RAM} + T_{Cache} - P T_{RAM} - P T_{Cache} = T_{Cache} + (1-P) T_{RAM}$$

اگر P مقدارش کم باشد، یعنی به صفر نزدیک تر باشد،  $T_{Cache} + (1-P) T_{RAM}$  که در بالا بدست آوردیم، بالا میرود. و اگر P مقدارش بالا رود یعنی به یک نزدیک شود، مقدار تاثیر  $T_{RAM}$  کم می شود.

**نکته** P به اندازه Cache و RAM و اندازه بلوک ها و تکنیک های نگاشت بستگی دارد.

حالا اینکه سیستم Cache داشته باشد یا نداشته باشد، به مقدار P بستگی دارد.

پس حافظه Cache سرعت بالایی دارد. خیلی از قسمت ها در سیستم بافر وجود دارد. اگر قرار باشد که حافظه در داخل خودش فقط اطلاعات را نگه داری کند، همان بافر کافیسست! و اما اگر قرار باشد که جستجویی انجام شود، Cache و حافظه تداعی گر نیاز است. در حقیقت حافظه تداعی گر از Cache استفاده می کند. اما بافر در دل خود از حافظه RAM استفاده می کند. حافظه Cache نیاز دارد که روی خودش هم جستجو انجام شود. بافر یک حافظه میان گیر می باشد. و با سخت افزار بافر متفاوت می باشد. جستجو را انجام نمی دهد. مثل Cache نیست که چیزهای مهم در درون آن ذخیره شوند. حتی توی CPU هم ممکن است وجود داشته باشد.

CPU که می خواهد با RAM ارتباط برقرار کند، همیشه یک کپی از قسمت هایی که مد نظرش هست را به دلیل محلیت زمانی و مکانی در داخل Cache قرار می دهد و اگر این کار به خوبی و درستی انجام شود، سرعت بالا می رود. سوالی که اینجا پیش می آید این است که این عملیات چگونه انجام می شود؟

جواب : با استفاده از نگاشت ها ... که در ادامه با ۳ نوع نگاشت آشنا خواهیم شد ....

### ۳-۲ نگاشت

انواع نگاشت یا mapping از حافظه اصلی به cache :

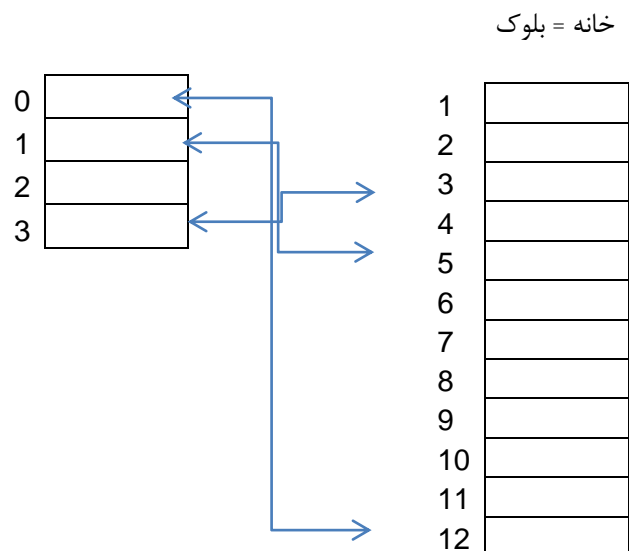
- DIRECT Mapping یا نگاشت مستقیم
- ASSOCIATIVE Mapping یا نگاشت تداعیگر
- SET ASSOCIATIVE Mapping یا نگاشت تداعیگر مجموعه ای

### ۱-۳-۲ نگاشت مستقیم

حافظه تداعیگر در مقایسه با حافظه های دستیابی تصادفی بدلیل مدار اضافی همراه هر سلول گرانتتر است در نگاشت مستقیم، جایی که بلوک باید در آن قرار گیرد مشخص است، و در روش های اشتراکی بلوک در چند محل متفاوت می تواند قرار گیرد در نگاشت مستقیم، جایی که بلوک باید در آن قرار گیرد مشخص است در روش های اشتراکی بلوک در چند محل متفاوت می تواند قرار گیرد در درجه ی اول مکانی انتخاب می شود که بیت اعتبار آن غیر فعال است. در غیر این صورت از بلوکی که کم تر مورد استفاده قرار گرفته است، از حافظه ی نهان خارج می شود.

فرمول نگاشت مستقیم :

شماره بلوک کاشه = تعداد بلوک کاشه mod شماره بلوک حافظه : DIRECT Mapping



شماره ۵ حافظه را بخواهیم ببریم داخل کش  $5 \bmod 4 = 1$  پس باید داخل خانه ۱ کش برود.

شماره ۱۲ حافظه را بخواهیم ببریم داخل کش  $12 \bmod 4 = 0$  پس باید داخل خانه ۰ کش برود.

شماره ۳ حافظه را بخواهیم ببریم داخل کش  $3 \bmod 4 = 3$  پس باید داخل خانه ۳ کش برود.

شماره ۹ حافظه را بخواهیم ببریم داخل کش  $9 \bmod 4 = 1$  پس باید داخل خانه ۱ کش برود، ولی خانه ۱ پر است در اینجا باید خانه یکم را write back کنیم ولی اگر در داخل کش محتوای خانه یکم تغییر کرده و با حافظه اصلی نابرابر است باید Write back کنیم اما اگر مقدار داخل کش تغییر نکرده و با حافظه یکی است دیگر به write back نیازی نیست پاکش می کنیم.

حال اگر فرض کنیم برنامه ما قسمتی در خانه ۱ قسمتی در خانه ۵ و قسمتی در خانه ۹ می باشد خانه ۱ را باید در خانه یکم کش بگذاریم بلوک ۵ را هم در خانه ۱ باید بگذاریم پس باید محتوای خانه ۱ را برگرداند بلوک نهم را در خانه یکم باید بگذاریم باز باید محتوای خانه یک را برگرداند یعنی راندمان پایین می آید چون از کش درست استفاده نکردیم چون با **direct mapping** می گوییم حتما باید در یک خانه خاصی باشد یعنی با **direct mapping** امکانش هست که از کش درست بهره برداری نشود.

مشکل **direct mapping** : **cashe** هنوز فضای خالی دارد ولی باید از یک سری خانه ها بیشتر استفاده شود و همین باعث می شود  $p$  (بهره وری) پایین باید یعنی برخی از خانه ها خالی هستند در حالی که برخی چندین بار خالی و پر می شود .

مزیت **direct mapping** : جستجو در خانه حافظه کش راحت تر است زیرا لازم نیست تمام کش جستجو شود مثلا اگر به بلوک ۱۴ نیاز داریم فقط بلوک ۶ از کش جستجو می شود ( $14 \bmod 8 = 6$ ) و از لحاظ پیاده سازی راحت تر است.

### ۲-۳-۲ نگاهت تداعیگر

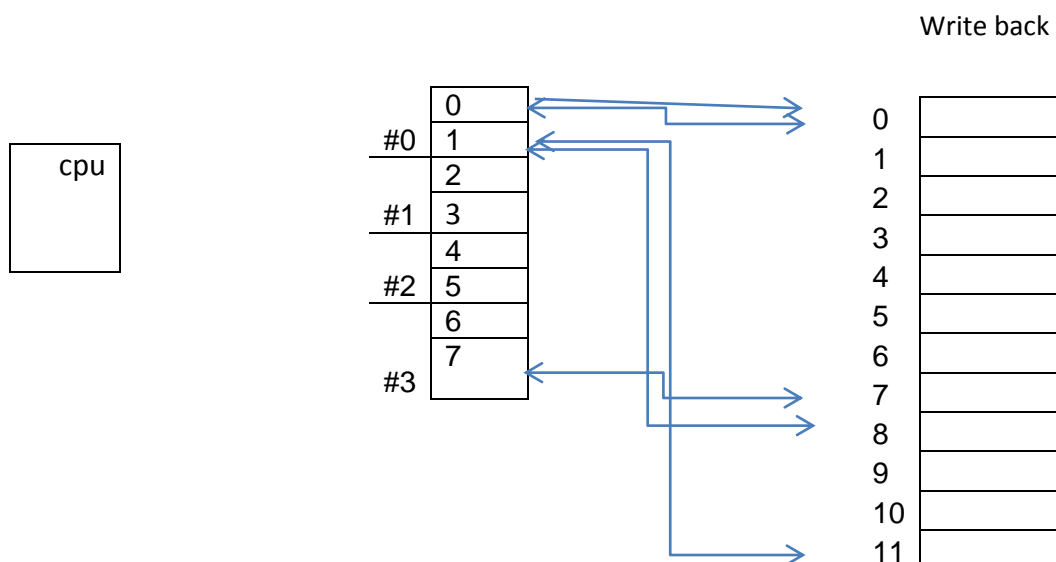
گفتیم باید کل حافظه جستجو شود توی این روش هم شاخص را باید در کش ذخیره کنیم و هم نشانه را باید ذخیره کنیم یعنی در هر کجای کش که خواست بنویسید و برای پیدا کردن داده باید کل حافظه جستجو شود در این روش هم شاخص را باید در کش ذخیره کنیم و هم نشانه را باید ذخیره کنیم یعنی تمام ست های آدرس در داخل حافظه کش ذخیره می گردد ( هر جایی از حافظه می تواند ذخیره شود ) در یک قسمت داده و در قسمت دیگر آدرس و وقتی کش کامل پر شود باید **write back** کنیم و برای **write back** می توانیم از الگوریتم های مختلف مثل **fifo** و غیره استفاده کنیم.

مزیت : استفاده بهینه از حافظه کش

مشکل : چون باید کل خانه های کش جستجو شود سرعت جستجو پایین می آید .

### ۲-۳-۳ نگاهت تداعیگر مجموعه ای

در این روش کش مجموعه مجموعه می شود برای مجموعه بندی باید تعداد بلوک های همه مجموعه ها یکسان باشد.



حافظه کش همیشه باید ۲ به توان X باشد

شماره مجموعه کش = تعداد مجموعه ها mod شماره بلوک حافظه

۷ → #3

۷ → #3

۱۲ → #0

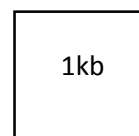
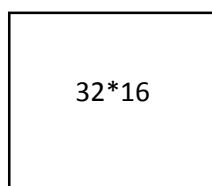
این روش میانه است باید ببینیم کدام بهتر می شود، اگر اندازه مجموعه بزرگ باشد یعنی کل کش یک مجموعه باشد به associative نزدیک می شود و دارای معایب و مزایای associative است. اگر اندازه مجموعه کوچک باشد به روش direct نزدیک می شود و به معایب و مزایای direct است. سرعت جستجوی این روش میانه این دو روش است.

مثال ( حافظه ما یک حافظه 32m دارد و داده ها ۱۶ بیتی است یعنی ۱۶\*۳۲ است

این کامپیوتر 1 kb کش دارد از روش نگاشت مستقیم استفاده می کنیم.

(۱) مکانیسم عملکرد سیستم را توضیح دهید.

(۲) هر بیت نشانه به چند شاخص و هر طبقه کش چقدر بیت نیاز دارد.



در هر طبقه ۱۶ بیت  $= 2^{25} 32 * 2^{20}$

۲۵ خط آدرس

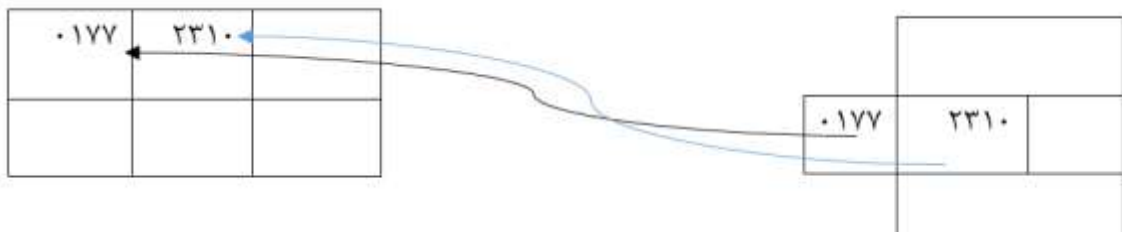
هر طبقه ۳۱ بیت دارد ۱۰ شاخص

۱۵ بیت برای نشانه ۱۵ نشانه

۱۶ بیت برای ذخیره داده

در تداعی گر هم شاخص و هم نشانه را باید ذخیره کنیم

هم آدرس و هم داده کامل ذخیره می شود.



MIPS دومی بیشتر است.

## ۴-۲ خط لوله ( Pipeline )

### ۴-۲-۱ روش ترتیبی (بدون خط لوله)

بطور مثال اجرای هر دستور اسمبلی در یک برنامه دارای ۴ مرحله می باشد (تعداد مراحل بستگی به طراحی سخت افزار دارد)

مرحله اول Fetch یا واکنشی یک دستور از حافظه و قرار دادن درون ثبات است. دومین مرحله از اجرای هر دستور Decode یا خواندن و باز کردن دستور است. سومین مرحله Execute یا همان اجرای کد است. و آخرین مرحله WriteBack یا بازنویسی نتایج اجرای دستور در حافظه اصلی است.

طبق دید معمولی در روال عادی پردازنده تمام این مراحل برای دستور اول انجام می شد و بعد تمام مراحل برای دستور دوم و همینطور دستور سوم و ...

F <sub>1</sub>	D <sub>1</sub>	EX <sub>1</sub>	WB <sub>1</sub>	F <sub>2</sub>	D <sub>2</sub>	EX <sub>2</sub>	WB <sub>2</sub> .....
----------------	----------------	-----------------	-----------------	----------------	----------------	-----------------	-----------------------

اگر زمان اجرای هر دستور را T<sub>n</sub> و زمان اجرای هر مرحله از دستور را T<sub>k</sub> در نظر بگیریم در نتیجه زمان اجرای هر دستور برابر است با kT<sub>k</sub>

(در اینجا تعداد مراحل 4 در نظر گرفته شده است یعنی  $k=4$ ) و زمان اجرای  $n$  دستور برابر است با  $nkT_k$

## ۲-۴-۲ روش با خط لوله

پایپلین کردن یا pipelining تکنیکی است که در آن اجرای چند دستور باهم همپوشانی پیدا می‌کنند. امروزه پایپلین کردن یک امر حیاتی برای افزایش سرعت پردازنده‌ها به حساب می‌آید. در این بخش برای توضیح اصطلاحات و مفاهیم پایپلین کردن از مقایسه با دنیای بیرون استفاده خواهیم کرد.

کسانی که در اتاق رختشویی کار میکنند و تعداد زیادی لباس را شستشو می‌کنند، به طور ناخودآگاه از تکنیک پایپلین کردن استفاده می‌کنند. یک روش غیر پایپلین برای شستن این تعداد لباس به صورت زیر است:

۱. تعدادی لباس کثیف را داخل شوینده قرار (washer) دهید .

۲. بعد از تمام شدن عملیات شستن، لباس های خیس را داخل خشک کننده (dryer) قرار دهید.

۳. وقتی که کار خشک کننده تمام شد، لباسهای خشک را بر روی یک میز قرار داده و تا کنید.

۴. وقتی که کار تا کردن تمام شد، لباس های تا شده را داخل کمد قرار دهید و یا از همکارتان بخواهید که

آنها را بیرون ببرد.

وقتی که لباس های تا شده را داخل کمد قرار دادید، دوباره یک کار جدید را با شستن دسته دیگری از لباس های کثیف، شروع کنید.

در پردازنده‌ها نیز اجرای دستورات می‌تواند به صورت پایپلین شده انجام شود. به این صورت که به علت بهره‌وری پردازنده و بیکار نماندن آن در زمانه‌ای بین اجرای هر مرحله از اجرای دستور، هر مرحله بعد از اتمام کار به سراغ دستور بعد رفته و مرحله مربوط به خود را در دستور بعد اجرا کند.

یعنی مرحله Fetch پس از اتمام کار خود به سراغ دستور بعد رفته و دستور بعد را Fetch کند. همین‌طور مرحله Decode و مراحل Execute و WriteBack.

F <sub>1</sub>	D <sub>1</sub>	EX <sub>1</sub>	WB <sub>1</sub>		
	F <sub>2</sub>	D <sub>2</sub>	EX <sub>2</sub>	WB <sub>2</sub>	
		F <sub>3</sub>	D <sub>3</sub>	EX <sub>3</sub>	WB <sub>3</sub>
			F <sub>4</sub>	D <sub>4</sub>	EX <sub>4</sub> WB <sub>4</sub>



در اینجا همانند روش ترتیبی دستورات پشت سر هم اجرا می‌شوند ولی از مکانیزم دیگری برای این کار استفاده می‌شود.

در این حالت زمان اجرای  $n$  دستور برابر است با  $kT_k + (n-1) T_k$  چون از دستور دوم به بعد به  $T_k$  زمان نیاز است.

### مقایسه کارایی در دو حالت فوق:

ما می‌توانیم میزان افزایش سرعت پردازنده پایپلاین شده به پردازنده غیر پایپلاین را که در بالا صحبت کردیم در قالب فرمول بیان کنیم. اگر همه مراحل کاملاً بالانس باشند (تأخیرشان یکسان باشد)، در این صورت زمان بین دستورات در پردازنده پایپلاین شده (با فرض شرایط ایدآل) برابر خواهد بود با:

$$\frac{\text{کارایی با خط لوله}}{\text{کارایی بدون خط لوله}} = \frac{\text{زمان اجرای بدون خط لوله}}{\text{زمان اجرای با خط لوله}} = \frac{nkT_k}{[k + (n-1)]T_k} = \frac{nk}{k+n-1}$$

$$= \lim_{n \rightarrow \infty} \frac{nk}{k+n-1} = k$$

در شرایط ایده آل میزان افزایش سرعت پایپلاین کردن مساوی تعداد مراحل پایپلاین خواهد بود.

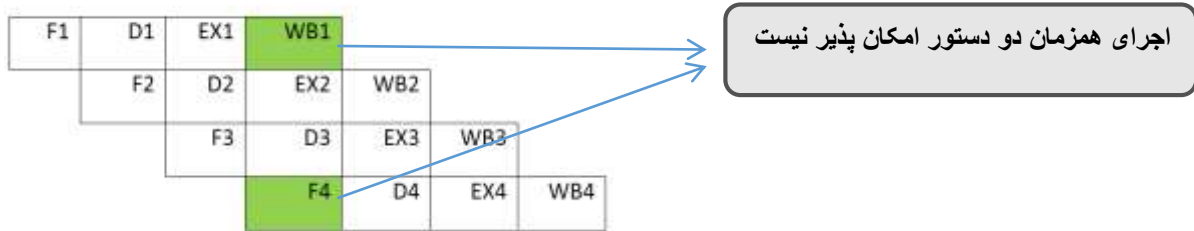
طبق مقایسه بالا کارایی سیستم در حالت دوم  $k$  برابر بیشتر از حالت اول است. یعنی دستورات  $k$  برابر سریع‌تر انجام می‌شوند.

در اینجا چون  $k$  را ۴ فرض کردیم پس سرعت انجام دستورات ۴ برابر سریع‌تر انجام می‌شوند.

### ۳-۴-۲ Hazard های خط لوله

هر چقدر  $K$  بیشتر شود هازارد هم بیشتر می‌شود و هر کدام مشکلات و مسائل خاص خودش را دارد و باعث می‌شود سرعت ما  $K$  برابر نشود و کمتر شود. و وقتی هازارد بیشتر شود مشکلات هم بیشتر می‌شود. پس عملاً اولین مسئله این است که هر چقدر  $K$  را بیشتر کنیم افزایش کارایی مان بیشتر می‌شود ولی به مثابه آن هازاردها هم بیشتر می‌شود.

## اولین hazard در خط لوله:



برای مثال زمانی که ما fetch می کنیم یعنی دستوراتی را می خواهیم از حافظه واکنشی کنیم، اگر همان زمان دستور wb (یعنی نتیجه ی اجرای دستورات را روی حافظه بنویسیم ) بخواند انجام شود یعنی همان زمان یک دستور می خواهد بخواند و یک دستوردیگرمی خواهد بنویسد که این امکان برای حافظه های معمولی امکان پذیر نیست.

به طور مثال در شکل زیر هر دو دستور wb1 و f4 می خواهند هم زمان اجرا شوند اما چون f4 باید بخواند و wb1 می خواهد بنویسد امکان خواندن و نوشتن هم زمان روی حافظه معمولی نیست که به این مشکل HAZARD می گویند.

این هازارد ها باعث می شود کارایی خط لوله به اون حد K برابری خودش نرسد و هر چه تعداد مراحل را هم بیشتر کنیم این هازارد ها هم بیشتر می شوند

راه حل : استفاده از حافظه های DUALPORT که هم زمان هم می توان خواند و هم می توان نوشت در این صورت مشکل هازارد اول حل می شود اما این نوع حافظه گران تر است.

## دومین hazard خط لوله :

اگر فرض کنیم دستور اول، دستور jump باشد (یعنی پرش کردن به نقطه ای) زمانی که مرحله ی اول fetch است در این مرحله یعنی fetch کردن مشخص نمی شود که دستور اول jump است چون فقط وظیفه واکنشی دارد مثلا اگر دستور اول پرش به دستور پنجم باشد، در زمان واکنشی دستور اول این مشخص نمی شود پس بنابراین وقتی می خواهد دستور بعدی fetch شود دستور دوم یعنی add واکنشی می کند و دستور ۵ را fetch نمی کند .

Ex:

دستورات Jump ( پرش به دستور پنجم )

Add

Sub

Mul

دستور پنجم ۵:

در مرحله بعد یعنی مرحله ی اجرا هنوز هم هیچ دیدگاهی نسبت به نوع دستور وجود ندارد فقط زمانی که اجرا تمام شود نوع دستور مشخص می شود پس زمانی که اجرای دستور اول تمام می شود تازه متوجه می شود که دستور اول jump بوده و دو دستور بعدی ، بی خودی fetch شده اند و نیاز به واکنش کردن آنها نبوده است و این خودش یک hazard است.

راه حل :

یک مدار به نام prefetch وجود دارد که کارش بررسی دستورات است مثلاً در مثال بالا بررسی می کند ببیند دستور jump است یا نه اما ممکن است باعث افزایش هزینه یا افزایش زمان اجرا شوند.

**تمرین: چهار مورد از هازادهای (مخاطره های) خط لوله را بنویسید ؟**

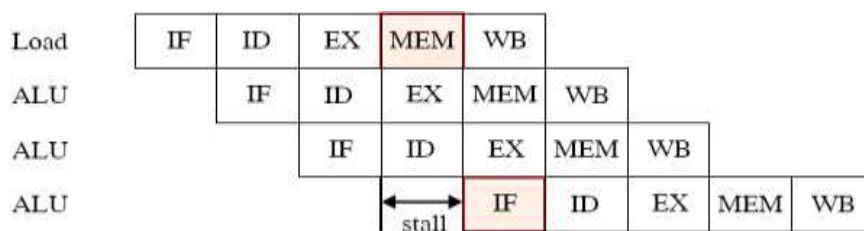
در معماری های خط لوله بزرگترین مانع وجود موانع یا در اصطلاح stall های خط لوله است. مانع هنگامی رخ می دهد که مجموعه جدید ورودی را نتوان در خط لوله بواسطه شرایطی به نام مخاطره یا hazard به اجرا در آورد. هنسی و پاترسون مخاطرات را به سه دسته تقسیم می کنند:

۱. مخاطرات ساختاری، بدین واسطه رخ می دهند که منابع ماشین نتوانند کلیه ترکیبهای همپوشانی دستورالعمل را پشتیبانی کنند.

۲. مخاطرات داده ای، هنگامی که نتیجه حاصل از یک دستورالعمل مورد نیاز دستورالعمل بعدی باشد.

۳. مخاطرات کنترلی، هنگام پردازش دستورالعملهای انشعاب ایجاد می شود.

مخاطره ساختاری می تواند در صورت عدم وجود منبع رخ دهد. برای نمونه اگر ماشینی صرفاً یک درگاه برای حافظه داشته باشد و نتوان داده و دستورالعمل را به طور همزمان از حافظه خواند این مشکل می تواند رخ دهد. برای نمونه به شکل ذیل توجه نمایید.



شکل: شکل هازارد ساختاری در درگاه حافظه ماشین های DLX

همانگونه که در بالا مشخص است نمی توان چهارمین fetch از حافظه را انجام داد زیرا دستورالعمل اولی در این زمان نیاز به دسترسی به حافظه دارد. برای ممانعت از مواجه شدن با همچنین وضعیتی در هنگام تولید کد سعی به استفاده بیشتر از ثباتها می شود. برای نمونه به دستورالعملهای ذیل توجه نمایید:

ADD R1,R2,R3

SUB R4,R1,R5

این دو دستورالعمل بدون هیچگونه تاخیری قابل اجرا هستند زیرا نتیجه دستورالعمل اول در دومی بلافاصله استفاده می شود. اما، در مورد دو دستورالعمل زیر نمی توان از ایجاد مانع یا در اصطلاح *stall* جلوگیری نمود.

LW R1,0(R2)

ADD R3,R1,R4

علت بوجود آمدن مانع در شکل ذیل مشخص شده است:

LW R1,0(R2)	IF	ID	EX	MEM	WB		
ADD R3,R1,R4		IF	ID	stall	EX	MEM	WB

شکل: هازاد داده ای

می توان با استفاده از امکان زمانبندی دستورالعملها یا در اصطلاح *Instruction Scheduling* یا *Instruction Reordering* این مشکل را حل نمود. برای این منظور با در نظر گرفتن وابستگی های داده ای و کنترلی امکان جابجایی دستورالعملهایی که ثبات R1 را استفاده نمی کنند و مشکل *stall* را ایجاد نمی کنند، را به میان دو دستورالعمل فوق بررسی باید نمود. در صورت یافتن دستورالعمل مناسب آنرا به بین دو دستورالعمل می توان انتقال داد.

مشکل فوق در اجرای عبارتی مثل  $A+B+C$  نیز امکان پذیر است. برای نمونه به شکل ذیل توجه نمایید. در اینجا چنانچه اجرا از چپ به راست انجام شود دستورالعمل جمع دومی می بایست برای یک سیکل در انتظار اجرا باقی بماند.

ADDD R3,R1,R2	IF	ID	EX1	EX2	MEM	WB	
ADDD R3,R3,R4		IF	ID	stall	EX1	EX2	MEM

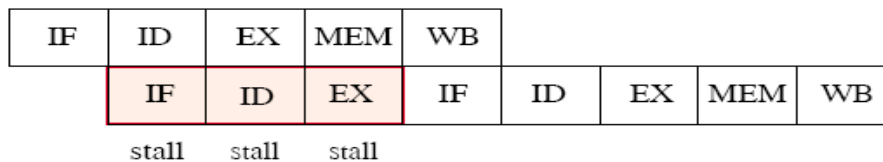
شکل : هازارد داده ای ایجاد کننده مانع

بازهم می توان گفت که زمانبند کامپایلر می تواند مانع این مشکل شود. برای نمونه فرض کنید که عبارتی برای جمع چهار متغیر به صورت  $A+B+C+D$  قرار است اجرا شود. در اینجا بعد از هر عمل جمع به تعداد دو سیکل انتظار برای دریافت نتیجه دستورالعمل قبلی نیاز است. برای رفع مشکل کامپایلر می تواند عمل جمع را به صورت ذیل گروه بندی کند:

$$(A+B) + (C+D)$$

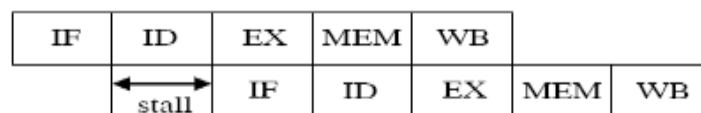
در این صورت جمع دومی می توان در یک سیکل بعد از جمع اولی به اجرا در آید. در مجموع یک سیکل مجموع سریعتر بدست می آید.

مخاطرات کنترلی همانطور که در بالا توضیح داده شد بواسطه دستورالعملهای پرش ایجاد می شوند. در ماشینهای DLX مخاطرات کنترلی می تواند مانعی به طول سه سیکل زمانی ایجاد نماید. پردازنده شروع به آوردن دستورالعمل بعدی از حافظه می کند بدون اطلاع از اینکه پرش انجام می شود. تشخیص عمل پرش بعد از مرحله MEM از دستورالعمل پرش مشخص می شود. در این صورت دستورالعمل بعدی می بایست از درون حافظه بیرون کشیده شود قبل از اینکه هر گونه تغییری در ثباتها و یا حافظه ایجاد شود. برای نمونه به شکل ذیل توجه نمایید.



شکل هازارد ایجاد شده بواسطه دستورالعمل پرش شرطی

از آنجایی که سه سیکل پهنالتی زیادی می باشد، طراحان ماشین سعی زیادی به کاهش این هزینه نموده اند. یک راه، استفاده از سخت افزاری است که در ضمن دیکد دستورالعمل بتواند تشخیص پرش را بدهد. این عمل در صورت آزمون مقایسه با صفر امکان پذیر است. در این صورت stall به یک عدد به صورت زیر تقلیل داده می شود:



شکل تقلیل هازارد ایجاد شده بواسطه دستورالعمل پرش شرطی

## ۲-۴-۴ انواع پیاده سازی های خط لوله:

### ۲-۴-۴-۱ روش سنکرون

برای پیاده سازی این نوع از خط لوله، از یک قطعه ای به نام latch استفاده می کنیم. زمانی که پایه لود مدار فعال باشد، latch ورودی ها را در خروجی قرار می دهد و آن را نگه می دارد. در این سیستم پایه لود (load) و چیزی به اسم clock عمل همزمان سازی را انجام می دهند. پس همان طور که گفتیم زمانی که پایه لود فعال می شود، latch ها هرچه در ورودی خود دارند در خروجی نمایش می دهند.

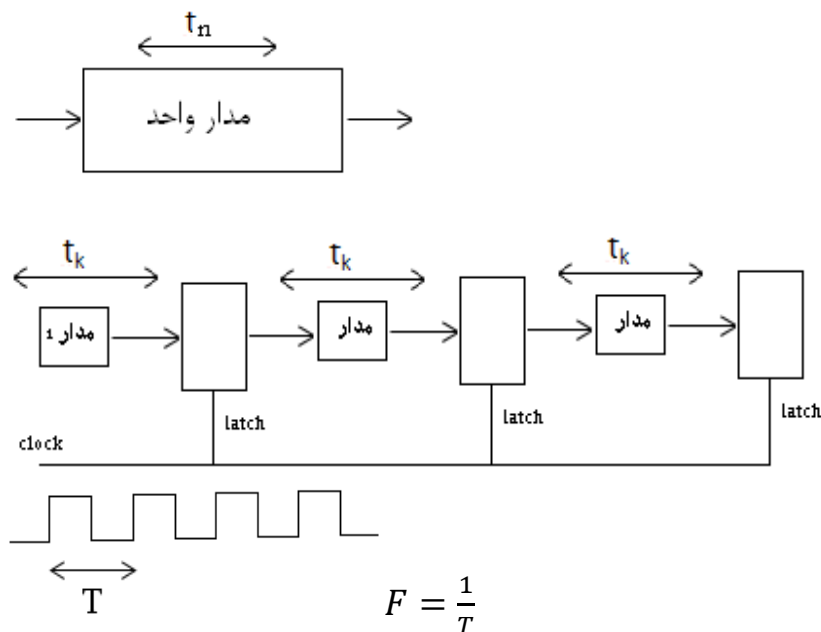
مثال: در مثال شستن لباس ها، قبلاً یک دانه تشت ( ظرفهایی که در آنها لباس می شویند) جواب می داده ولی حالا باید یک ظرف میانی باشد. تا با یک طریقی اینها با هم کار کنند. شاید آدم ها چشم و گوش داشته باشند و بینند که بعدی خالی است و اگر هم داده ها روی هم بیفتد هیچ اشکالی پیش نمی آید چون لباس ها را روی هم بگذاریم خراب نمی شود ولی در مورد داده ها وقتی داده اول وجود دارد و داده ی دوم را می گذاریم روی آن در همان خانه ی حافظه چه اتفاقی می افتد؟ داده اولی پاک می شود و از بین می رود.

اول از همه اینها چشم و گوش ندارند. باید یک فکر اساسی کنیم. در طریقه sync ما یک طبال داریم و کارش این است که هر وقت من طبل زدم شما باید به نفر بعدی سریع تحویل بدهید این ۴ تایی که اینجا نشسته اند هم کور و هم لال هستند هیچی نمی بینند و فقط صدای طبل را می شنوند. هر وقت صدای طبل را شنیدند باید کارشان تمام شود و بدهند به نفر بعدی و وقتی این را تحویل نفر بعدی می دهد از نفر قبلی، نفر ۲ از نفر قبلی تحویل می گیرد و آن نتایج اش را به نفر بعدی یا ۳ می دهد.

سوال: این طبال به کسی باید نگاه کند؟

به کسی که از همه کندتر است باید نگاه کند، زمانی که طبال طبل می زند باید بر اساس مصرف کننده پر تاخیرترین قسمت باشد. هرچقدر آن قسمت بیشتر تاخیر داشته باشد طبال می زند. پس اگر تعداد مراحل هم زیاد باشد بهتر است. چرا؟ چون طبال سریعتر طبل می زند.

مدار سنکرون: مدار واحد را تکه تکه کردیم.



اگر تاخیر قسمت های مختلف یکسان باشد. نتیجه می گیریم کار بهتر انجام می شود. چرا؟ چون با هم همزمان کار انجام می دهند. و به جای طبال اینجا clock را داریم.

در مدل سنکرون بینشان Latch می گذاریم. قبلاً یک مدار واحد داشتیم که زمان تاخیرش  $t_n$  بود و حالا مدار را تکه تکه کردیم و زمان هر مرحله اش  $t_k$  است. فرض می کنیم که اندازه تاخیرهایشان با هم برابر است ولی اگر تاخیر یک قسمت بیشتر باشد بد است. ما باید به آن نگاه کنیم. بهتر است است تاخیر قسمت های مختلف یکسان باشد.

حال ما طبال را در دنیای واقعی داریم ولی اینجا ما یک clock pals (پالس) می گذاریم. Clock یک سری ۰ و ۱ متوالی است که توسط یک دستگاه کریستال ایجاد می شود. فاصله ی بین دو لبه ی بالا رونده را پررود می گویند و با T نشان می دهند و فرکانس هم  $\frac{1}{T}$  است. که T واحدش زمان و فرکانس واحدش HZ است. و از آنجایی که Latch کار قفل را انجام می دهد. و از آنجا که باید فرکانسش فعال شود هر چه در ورودی است در خروجی می گذارد. و در خروجی اش نگه می دارد تا زمانی که پایه فرمانش فعال شود. و در خروجی خویش قفل می کند.

پس این پایه فرمان Latch ها یا پایه قفل Latch ها مثل طبل می ماند. بهتر است که وقتی حساس به لبه بالا رونده باشد و طول پررودمان برابر  $T_k$  باشد و فرکانس  $\frac{1}{T_k}$  باشد.

هرچقدر فرکانسمان بالا رود  $T_k$  کم می شود. پس ما علاقه مند هستیم که فرکانسمان را بالا ببریم چون هر چه فرکانس را بالا ببریم  $T_k$  کم شده و هر چه  $T_k$  کم شود یعنی اینکه تعداد مراحلمان زیاد شده و درباره آن مشکلات خاص و هازارد ها پیش می آید.

قیمت خط لوله و تاخیرش چه می شود؟

فرض کنیم قیمت کل واحدمان  $P_T$ ، قیمت یک Latch را  $P_L$  حساب کنیم، قیمت کل می شود:

$$\text{قیمت کل} = P_T + KP_L$$

چرا K: چون Latch ها یکی نیست چند تاست.

فرض کردیم که قسمت قسمت کردن مدار هزینه ندارد.

تاخیر هر کدام از اینها  $T_k$  است. وقتی که همه را پشت سرهم بگذاریم می شود  $t_n$  و هر یک از این Latch ها یک تاخیر دارد  $t_L$  پس تاخیر کلمان می شود:

$$t_n + t_L = \text{تاخیر کل} \uparrow \frac{t_n}{k} \text{ تاخیر یک مرحله و } t_L \text{ تاخیر یک Latch است.}$$

تعداد K بهینه را می خواهیم حساب کنیم. چون K قیمت ، تاخیر را طوری می خواهیم حساب کنیم که K بهینه باشد. قبلاً گفتیم که K بیشتر باشد بهتر است بعد گفتیم K باشد هزاردار داریم.

$$\frac{\text{قیمت}}{\text{تأخیر}} = \frac{p_T + kp_L}{\frac{t_n}{k} + t_L}$$

وقتی K افزایش یابد تعداد Latch بیشتر می شود و وقتی تعداد Latch ها بیشتر شود هزینه و تاخیر زیاد می شود. بنابراین با افزایش K قیمت کل و تاخیر کل سیستم بیشتر می شود. حال دنبال این هستیم که یک K بدست آوریم که هم تاخیرش کم شود و هم قیمت زیاد افزایش نداشته باشد.

چگونه K بهینه را بدست آوریم؟ از نسبت قیمت به تاخیر نسبت به k مشتق می گیریم و صورت مشتق را مساوی صفر قرار می دهیم.

$$\frac{p_T + kp_L}{\frac{t_n}{k} + t_L} = \frac{p_L \left( \frac{t_n}{k} + t_L \right) - \left( -\frac{t_n}{k^2} \right) (p_T + kp_L)}{\left( \frac{t_n}{k} + t_L \right)^2} = 0$$

$$p_L \left( \frac{t_n}{k} + t_L \right) + \left( \frac{t_n}{k^2} \right) (p_T + kp_L) = 0$$

$$K^* = \sqrt{\frac{p_n t_n}{p_L t_L}}$$

وقتی سنکرون را حساب کنیم همان فرمول را حساب کنیم به این نتیجه می رسیم که اگر قیمت مدار را افزایش داده و قیمت Latch را کاهش دهیم بهتر است تعداد مراحل زیاد شود و از آن طرف اگر قیمت مدار افزایش و تاخیر Latch کاهش یابد بهتر است تعداد مراحل زیاد شود. یعنی وقتی تاخیر و قیمتمان پایین است بهتر است تعداد مراحل کم شود و همین طور بالعکس.

با افزایش قیمت Latch و افزایش تاخیر Latch بهتر است تعداد مراحل کم شود. چرا؟ چون Latch ها تاخیر دارند و از طرفی با کاهش قیمت Latch و کاهش تاخیر Latch بهتر است تعداد مراحل n بیشتر شود.

نکته: ابزارهای طراحی سنکرون در بازار فراوان است و این روش سنکرون خیلی خوب است و اینجا یک اشکالی که بوجود می آید این است که ما باید یک مداراتی اضافه کنیم به اسم Latch که این هم هزینه دارد و این باعث می شود که ما نتوانیم زیاد kمان را افزایش دهیم.

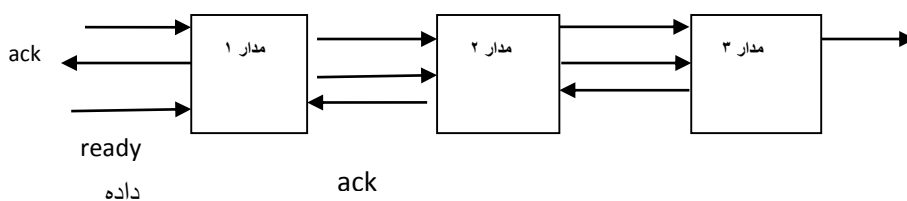
راه حل دیگر:



## ۲-۴-۴-۲ روش آسنکرون:

در این سیستم دیگر نیازی به clock نداریم و از چیزی به اسم hand shaking استفاده می کنیم. در این روش دو قطعه برای یکدیگر سیگنال هایی ارسال می کنند تا از میزان آمادگی یکدیگر با خبر شوند.

مثال: در مثال شستن لباس ها هر کسی با بغل دستی اش در ارتباط است. یعنی دست تکان می دهد به این صورت است که وقتی مدار اصلی را تکه تکه کردیم مدارها خودشان را پشت سر هم می گذاریم. مدار ۱ که آماده شد تحویل مدار ۲ بدهد. و مدار ۲ که آماده شد تحویل مدار ۳ بدهد. در سنکرون مدار نتیجه اش را می گذارد داخل Latch و Latch تحویل مدار ۲ می دهد. ولی در اینجا مدار هر وقت کارش را انجام داد به بعدی می گوید من آماده ام و داده ی جدید را بگیر و بعدی هم خواست می گوید من گرفتم و اگر آماده نبود جوابش را نمی دهد. و صبر می کند کار قبلی اش تمام شود و سپس به بعدی تحویل می دهد. به این کار می گویند دست تکان دادن. در این روش هر کسی با خودش و با بعدی و قبلی خودش ارتباط دارد از طریق سیگنال Ack, ready



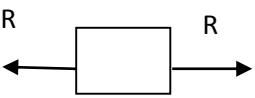
مزیت: قیمت Latch، چون هزینه ای برای Latch ها داشتیم هزینه Latch حذف شده و تأخیر Latch هم نداریم.

بدی: مدارش یک مقدار پیچیده تر می شود چرا؟ باید با طرف راست و چپش هماهنگی کند و این پیچیدگی مداری خواهد داشت. و مداری برای hand shunting (دست تکان دادن) به هر کدام از آنها باید اضافه کنیم. ولی این بهتر از قبلی است.

اشکال بزرگ: چون تحلیل این ها به صورت دستی امکان پذیر نیست و نرم افزار خوبی که این کار را انجام دهد عملاً وجود ندارد پس این روش، روش خوبی نیست.

در سنکرون ابزار زیاد است، چرا؟ چون پیچیدگی کمتر است. چون طبال کار را خیلی آسان می کند clock را که می گذاریم می گوید سر clock کار را انجام بده و تحویل بده ولی در آسنکرون مشخص نیست کجاست؟ در سنکرون همه ی عملیات ها در لبه ی بالا رونده clock انجام می شود و هر گاه لبه ی بالا رونده رخ داد یک سری کار را انجام می دهد و در مراحل بعدی تا لبه ی بالا رونده بعدی خداحافظی می کند. آن لحظه ای که لبه ی بالا رونده اتفاق می افتد آنجا یک کاری اتفاق می افتد، و در آنجا سیستم را می توان راحت طراحی کرد. و سیستم گسسته است.

وقتی آسنکرون است  
آماده است که بفرستد؟



ما نمی دانیم که آیا این هم آماده است که بگیرد و هم

در گذشته انسان ها سیستم را تحلیل می کردند. ولی الان ابزارها تحلیل می کنند. ابزارهایی که تحلیل می کند و وقتی می خواهیم روی سیستم خودمان طراحی کنیم گیر می کند و به مشکل می خورد. و به راحتی تحلیل نمی کند. پس باید نرم افزار یا ابزار تحلیل سیستم آسنکرون تهیه کنیم. ابزارها بیشتر برای sync است و برای Async ابزار کم است.

روال ساخت IC:

وقتی IC می خواهیم بسازیم از خاک رس یا خاک قرمزی که در آن سیلیکتون است استفاده می شود. و این را به کارخانه می برند و تبدیل به IC می شود.

دو زبان HDL داریم: ۱- VHDL ۲- Verilog

چرا در سیستم آسنکرون کسی برایش نرم افزار طراحی نکرد؟

چون گسسته نیست و پیوسته است برای سیستمی که پیوسته است ما همیشه می آییم گسسته سازی می کنیم. وقتی یک سیستم پیوسته را گسسته سازی روی آن انجام می دهیم یک سری جزئیاتی روی آن از بین می رود. و جزئیات که از بین می رود باعث می شود که مدار تغییر کند و ما اصلاً جواب نداریم. با دقت بالا گسسته سازی کنیم مشکل حل می شود. پس عملاً ابزارهایی که این را تبدیل کنیم به IC برای سیستم های sync سخت است. پس عملاً همه ی چیزهایی که داخل com هستند بیشتر sync است و یعنی همه ی آنها طبال را دارند.

## فصل سوم

### چندپردازندگی

#### ۳-۱ طبقه بندی سیستم ها:

یک طبقه بندی بنام Flynn داریم که چندپردازنده ها را به چهار سیستم یا نوع تقسیم می کند:

- 1- SISD: Single Instruction Single Data
- 2- SIMD: Single Instruction Multiple Data
- 3- MISD: Multiple Instruction Single Data
- 4- MIMD: Multiple Instruction Multiple Data

- ۱- پردازنده های معمولی جزء این دسته هستند. یک دستور، یک داده
- ۲- یک دستو به صورت همزمان بر روی چند پردازنده اجرا می شود. کاربردش در پردازش ماتریس می باشد. به طور مثال جمع ماتریس یکبارہ انجام می شود.
- ۳- چندین دستور به طور همزمان بر روی یک داده اجرا می شوند. پیاده سازی این روش انجام نشده است.
- ۴- واقعا چند پردازنده ای است. یعنی چند تا دستور توسط چند پردازنده انجام می شود.

مثال: در SISD (یک دوربین، یک داده) یک دوربین در جاده می گذاریم و هر ماشینی که رد می شود دوربین باید پلاکش را بخواند و پردازش سیگنال یا پردازش دستور را انجام دهد.

در SIMD از ۴ دوربین استفاده می کنیم و ۴ عکس از محیط می گیریم، در همه تصاویر به دنبال پلاک خودرو هستیم. پس یک دستور است و قرار است یک دستورالعمل روی آن اجرا شود. دستورالعمل در همه آنها یکسان است و تصویر را به داده یا پلاک تبدیل می کنند. MISD دو دستورالعمل داریم، یک دستور تشخیص پلاک و دو تشخیص تعداد افراد داخل خودرو.

پس در این صورت یک عکس از خودرو گرفته و به یک پردازنده دستور می دهیم پلاک ماشین را تشخیص بدهد و به پردازنده دعدی دستور می دهیم تعداد افراد داخل ماشین را بدست آورد. پس دو دستور بر روی یک داده انجام می شود. ولی عملا این روش پیاده سازی نشده است. MIMD چند دستور توسط چند پردازنده انجام می شود.

برای بوجود آوردن چندپردازندگی برای Async ها:

در اینجا به جای اینکه واحد F مان یکی F کند ۲ تا F می کند اما حواسمان باشد، ۲ عدد پردازنده داخل سیستم وجود دارد. در قسمت D هر دو را همزمان D می کند پس ما ۲ عدد مدار داریم.

آیا می شود ۲ دستور را همزمان در حافظه انجام داد؟ آن مدار حافظه مان باید یک مقدار سخت تر باشد. باید یک مدار متنوع تری باشد تا بتوانیم ۲ بار از روی آن همزمان بخوانیم یا حافظه مان را Inter leveling می کند. یعنی یک دستوری در حافظه ی اولی می گذاریم و یک دستور را در حافظه ۲ می گذاریم، f می کنند. به این کار Inter leveling کردن می گویند.

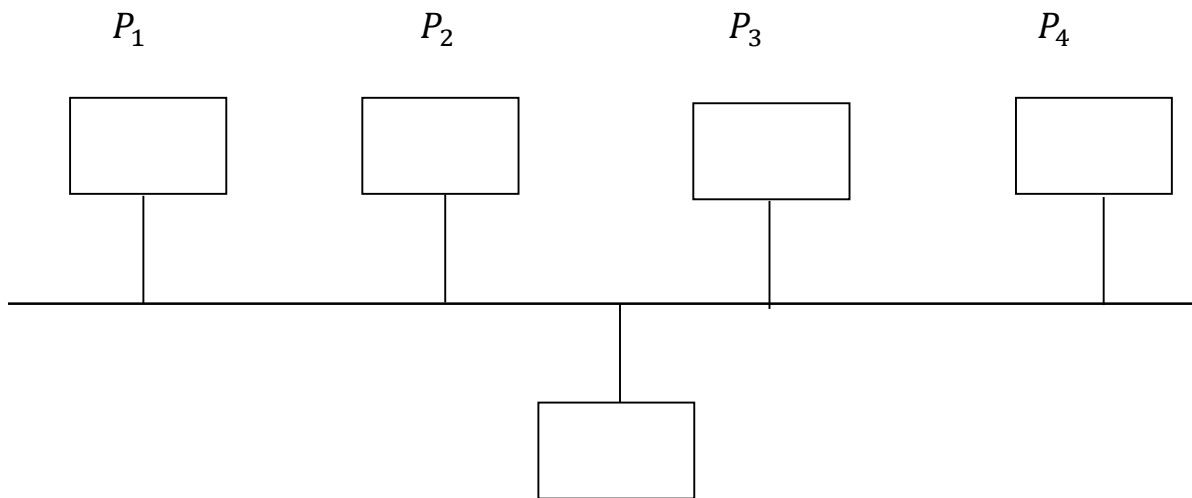
F1	D1	Ex1	Wb1			
F2	D2	Ex2	Wb2			
	F3	D3	Ex3	Wb3		
	F4	D4	Ex4	Wb4		
		F5	D5	Ex5	Wb5	
		F6	D6	Ex6	Wb6	
			F7	D7	Ex7	Wb7
			F8	D8	Ex8	Wb8

درهم ریختن و پشت سر هم قرار ندادن بعضی از این کلال ها به عمل ها کاری نداریم و عملیات را سطح بالاتر نگاه می کنیم و مطمئناً هزاردهایش بیشتر است مشکلاتش بیشتر است و دردسرهای بیشتری دارد و سخت افزار بیشتری دارد و هزینه اش نیز بیشتر است. ما با این ها هیچ کاری نداریم می خواهیم ببینیم اگر این کار را کنیم سرعت چقدر بهتر است همه ی آنها را می دانیم.

**تمرین:** سیستم های grid را با سیستم های چند پردازنده مقایسه کنید؟

در grid سیستم های مختلف با هم ارتباط دارند در صورتی که در چند پردازنده ای cpu ها به طور مستقیم با هم در ارتباط هستند. در grid ارتباط از طریق اینترنت می باشد در صورتی که در چندپردازنده ها ارتباط از طریق bus می باشد. در grid پردازنده ها با هم یک کار را انجام میدهند در صورتی که در چند پردازنده هر پردازند کار خودش را انجام می دهد.

چگونگی اتصال پردازنده ها با یکدیگر:



### ۲-۳ ساختارهای اتصالات متقابل

#### ۱-۲-۳ سیستم چند پردازنده معمولی (ساده):

در این روش همه پردازنده ها از طریق یک BUS بهم متصل شده اند. از لحاظ هزینه، هزینه چندانی ندارد چراکه قبلا گذرگاه بین حافظه و پردازنده وجود داشته است و در این روش تنها چند عدد پردازنده به آن اضافه شده است.

مشکل این روش این است که :

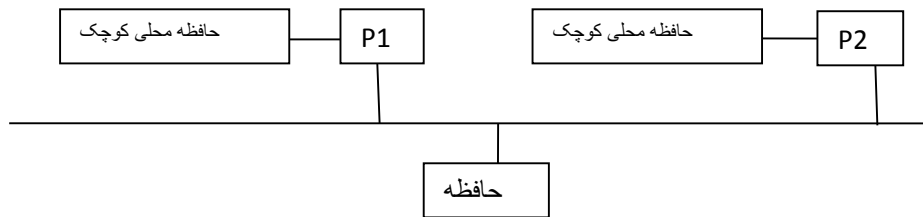
۱- در هر لحظه تنها یک پردازنده می تواند از حافظه استفاده کند. بدلیل ارتباط از طریق یک باس.

۲- اگر تعداد پردازنده ها زیاد شود هزینه بالا رفته ولی کارایی بالا نمی رود.

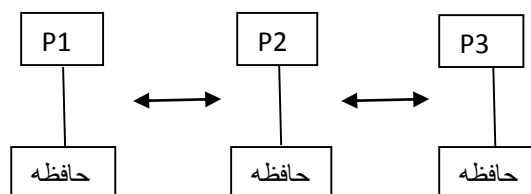
راه حل: یک پردازنده محلی به هر پردازنده اختصاص دهیم. در این صورت هر پردازنده به جای اینکه برای واکنشی تک تک دستورات دفعات زیادی به حافظه مراجعه کند می تواند با یک مرحله مراجعه به حافظه اطلاعات زیادی را واکنشی کرده و در حافظه خودش قرار دهد. در این صورت باس آزاد بوده و سرعت کار بالا می رود. ولی باز مشکلاتی نیز وجود دارد، هدف ما از چند پردازنده ای موازی سازی و ارتباط پردازنده ها باهم می باشد. اما اگر در این جا یک پردازنده به نتیجه کار پردازنده دیگر احتیاج داشت باید صبر کند تا پردازنده دیگر کارش به پایان برسد و سپس نتیجه کار را بر روی حافظه بنویسد و بعد این پردازنده بتواند از آن اطلاعات استفاده کند، و یا می توان گفت اگر دو پردازنده بخواهند بر روی یک داده تغییراتی اعمال کنند

و هر دو داده را واکشی کرده و آن را تغییر دهند، سوالی که در این جا مطرح است این است که کدامیک باید بر روی حافظه نوشته شود؟

یک دفعه به اندازه ی حافظه اش به اندازه ای که می خواهد پردازش کند می رود یکی می گذارد توی حافظه ی محلی اش و همین طور  $p2, p3$  هم همینطور و این وسط پردازش را روی اطلاعات محلی شان انجام می دهد و در این صورت کارایی مان افزایش می یابد.

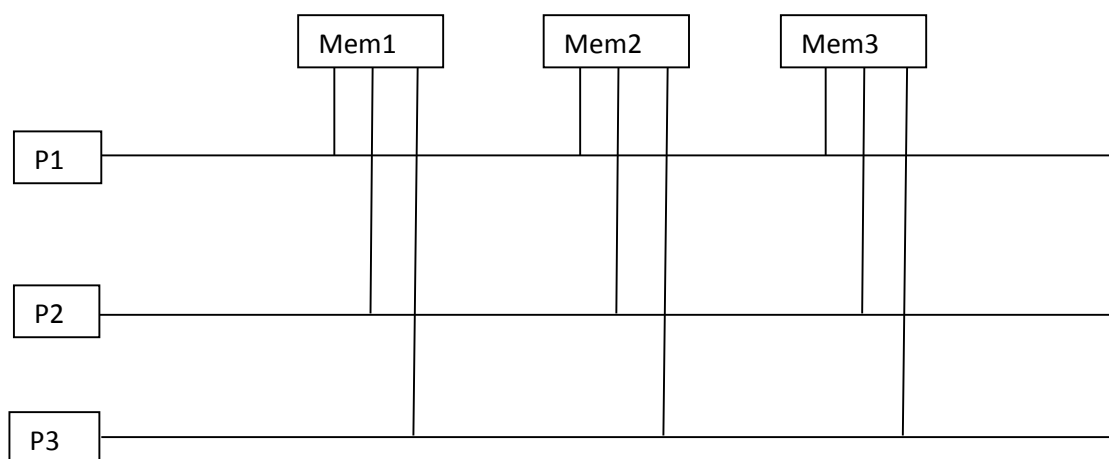


**مشکل:** مثلاً خانه ی حافظه مقدار ۱۰ را دارد و هر دو تا این ۱۰ را برمی دارند و در حافظه محلی خودشان می ریزد و هر کدام از اینها قرار است که یکی \*۲ یکی \*۳ و یکی \*۴ کند و حال یک خانه حافظه را می خواهند تغییر دهند هر کدامشان و هر کدام مقدار داده ای را می آورد و در خانه حافظه خودش و یکی \*۲ می شود و ۲۰ می شود و یکی \*۳ و دیگری ۴۰ و هر سه که می خواهند مقدار را برگردانند و در حافظه اصلی نوشته شود و این مشکل است. اگر اینها خواستند برگردانند و هر کدام تغییری روی داده داشته باشند کدامیک باید بازنویسی شوند؟ و یک پیشنهاد می دهیم: و می گوییم این پردازنده منتظر بماند و هر وقت کارش تمام شد این برود، بخواندش و تغییر دهد و هر موقع این کارش تمام شد دیگری برود ریال بخواند و تغییر دهد. و همین طور پردازنده ها پشت سر هم این کار را انجام دهند و الان چه اشکالی پیش می آید؟ یعنی  $p1$  وقتی تغییر می دهد  $p2$  باید منتظر این باشد که جواب را برگرداند و این روش بدتر می شود و اینکه پردازنده ها هر کدام حافظه خودشان را داشته باشند و بعداً باهم ارتباط داشته باشند و این دیگر چند پردازشی نمی شود و می شود شبکه با شبکه gride آن چند پردازنده ای که ارتباط شبکه ای دارند، البته ما به دنبال این هستیم که پردازنده ها خودشان باهم در ارتباط باشند.



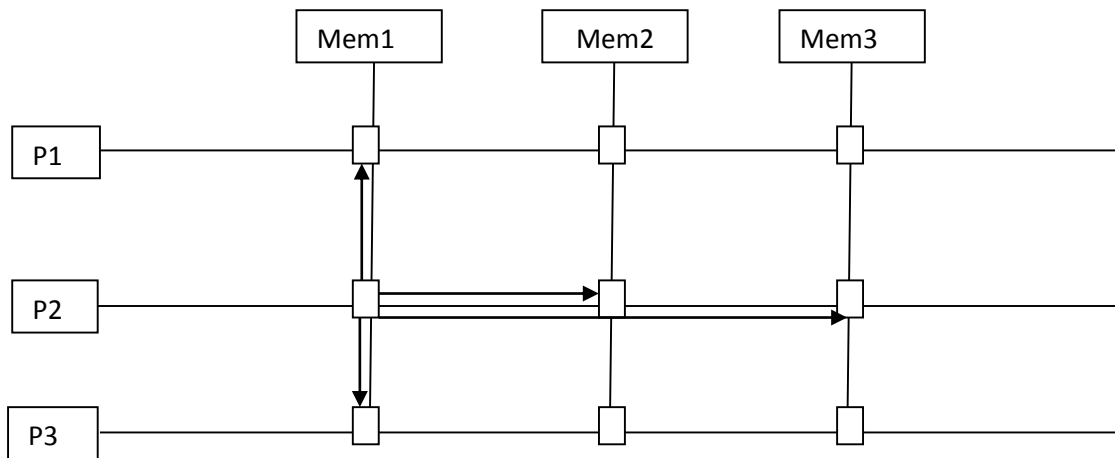
### ۳-۲-۲ حافظه چند پورتی:

حافظه چند پورته : این حافظه ها با حافظه های معمولی یک تفاوت دارند و می توانند همزمان ۲ دستور را باهم پردازش کنند اگر یک دستور به آن بدهیم و ما به حافظه دستور خواندن و نوشتن می دهیم، نباید با دستور Add,Sub اشتباه کنیم در این شکل از یک حافظه چند پورت خارج می شوند که می تواند همزمان چندین cpu به آن وصل باشد و همزمان از چند cpu دستور بگیرد. اگر یک cpu به آن دستور دهد سریعتر عمل می کند و اگر چندتا cpu به آن دستور دهد یک مقدار کندتر عمل می کند. البته بعضی از مدل ها همزمان دارند دستور می دهند و سریعتر هم عمل می کند ولی گران تمام می شود. این حافظه ها از حافظه های معمولی گرانتر است. در اینجا چه فرقی کرده ؟ در اینجا پردازنده ها یمان به چند تا حافظه وصل است باید طوری برنامه هایمان را تغییر دهیم و دسته بندی کنیم که برنامه هایی که مربوط به p1 هستند داخل mem1 باشد و برنامه هایی که مربوط به p2 است داخل mem2 باشد. و برنامه هایی که مربوط به p3 هستند داخل mem3 باشد. پس سوال مطرح می شود که چرا p3 به mem1 وصل است؟ شاید به اطلاعات نیاز پیدا کرد ولی کمتر به اطلاعات نیاز پیدا کرده ولی p3 بیشتر با mem3 کتر می کند. ولی ارتباط هم با mem1 دارد به خاطر بعضی مواقع. البته mem1 توانایی جواب دادن به هر دو یعنی p1, p3 را دارد ولی مایک کار می کنیم که هر کدام از اینها داده های مربوط به خودشان را در حافظه جداگانه داشته باشند و از این طرف هر وقت به حافظه دیگری نیاز پیدا کرد مشکل پیش نیاید.



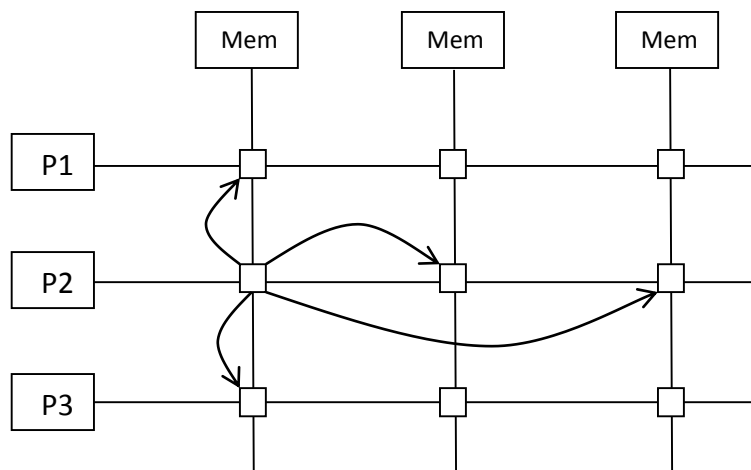
یک مشکل اساسی دارد : هزینه - ما باید برای حافظه چند پورته پول بدهیم و از یک طرف هم مشکل scalability (توسعه پذیری) را داریم و به همین راحتی نمی شود زیاد کرد، وقتی زیاد می شود باز هم به مشکل می خوریم البته تعداد پردازنده های ما با حافظه ها لزوماً برابر نیست و کمتر هم می تواند شود ولی وقتی تعداد پردازنده ها زیاد شود و تعداد حافظه ها کم باشد مثل قبل می شود و اگر نشود حافظه با پورت زیاد گران می شود.

### ۳-۲-۳ حافظه با سوئیچ تقاطعی:



سوئیچ ها نسبتا ساده هستند.

**عملکرد سوئیچ:** به این صورت است که می بیند به این حافظه چیزی وصل است یا نه وقتی  $p1$  می خواهد با حافظه ۱ ارتباط برقرار کند می آید به سوئیچ دستور می دهد و این سوئیچ به حافظه ۱ متصل می شود و در حینی که  $p1$  با  $mem1$  متصل است این  $a$  اگر بخواهد به  $mem1$  متصل شود این سوئیچ اجازه نمی دهد و این سوئیچ به  $a2$  می گوید اگر می خواهید بنویسی بر روی  $mem2$  بنویسید. پس نوشتن به حافظه و خواندن از حافظه از طریق این سوئیچ امکان پذیر می شود.



مزایا :

۱. حافظه تک پورت است.
۲. گسترش پذیری راحت است هر چقدر می خواهیم پردازنده استفاده کنیم و هر چقدر می خواهیم حافظه استفاده کنیم.

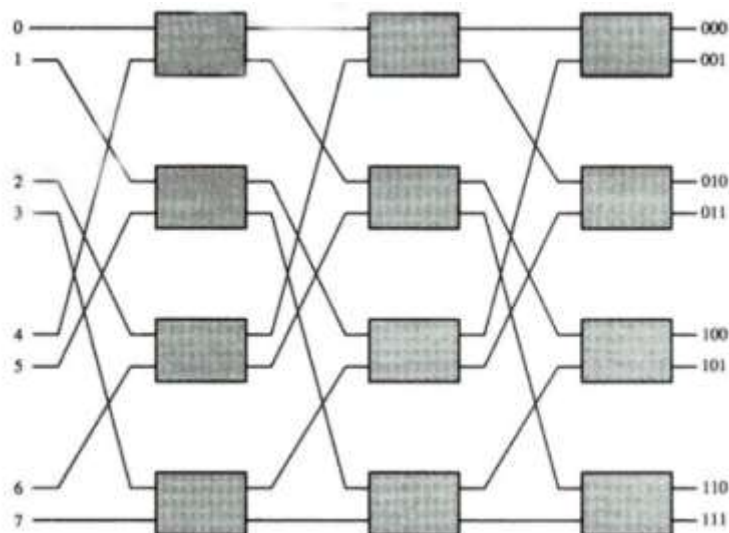
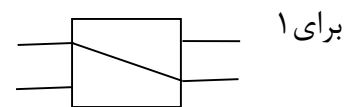
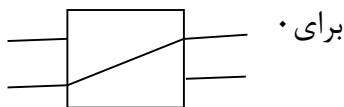


اشکال : ما برای ۳ دستور ۹ تا سوئیچ می کنیم و برای  $n$  دستور  $n^2$  سوئیچ می کنیم. تعداد سوئیچ ها رابطه مستقیمی با تعداد پردازنده هایمان ندارد برای ۴ تایی یا ۵ تایی خوب است اما برای مثلا ۲۵۶ پردازنده که  $256^{256}$  سوئیچ نیاز داریم خوب نیست. به صرفه و معمول نیست. خطایابی سخت و تاخیرمان بالا می رود.

### ۴-۲-۳ استفاده از سوئیچ های بانیان (شبکه سوئیچ های چند طبقه $2 \times 2$ )

مدل **banyan**: در ارتباط پردازنده ها به هم کاربرد دارد و در طراحی سوئیچ و روتر کاربرد دارد. ساختار  $\Omega$  امگا دارد. مدل های دیگری هم دارد. در مثال قبلی به سوئیچ باید نگاه میکرد که دیگران اشغال هستند یا نه ولی در اینجا سوئیچ اگر ورودی اش صفر باشد ورودی به خروجی مستقیم وصل می شود و اگر ورودی اش یک باشد باید به پایین وصل شود. اگر دو تا صفر و دو تا یک باشد، **internal blocking** می شود و همچنین خالی نمی شود پس ما براساس آدرس وصل می شویم ، سمت چپ به پردازنده ها و سمت راست به حافظه ها. از طریق این سوئیچ ها به هم متصل می شویم.

پردازنده فقط شماره حافظه می فرستد. سوئیچ طوری عمل می کند که مستقیم وصل شود.



یک مشکل اساسی دارد:

هزینه: ما باید برای حافظه چندپورته پول بدهیم و از یک طرف هم مشکل توسعه پذیری Scalability را داریم و به همین راحتی نمی شود زیاد کرد و زیاد که می شود باز هم به مشکل برمی خوریم. البته تعداد پردازنده ها با حافظه هالزوما برابر نیست و می تواند کمتر شود. ولی وقتی تعداد پردازنده ها زیاد شود و تعداد حافظه ها کم باشد مثل قبلی می شود و اگر نشود حافظه با پورت زیاد گران می شود.

از سوئیچ استفاده می کنیم:

سوئیچ ها ی مدار نسبتا ساده هستند.

### عملکرد Switch :

Switch به نحوی عمل می کند که می بینید به این حافظه چیزی وصل است یا خیر؟

وقتی P1 می خواهد با حافظه ۱ ارتباط برقرار کند به Switch دستوری دهد و این Switch به حافظه ۱ متصل می شود در این حین که P1 به حافظه ۱ متصل است اگر P2 بخواهد به حافظه ۱ متصل شود این سوئیچ اجازه نمی دهد و به P2 می گوید اگر می خواهی بنویسی بر روی حافظه ۲ بنویس، پس نوشتن به حافظه و خواندن از حافظه از طریق این سوئیچ امکان پذیر می شود

مزایا:

۱- در اینجا از یک جهت مدار، مدار پیچیده ای نیست پس سوئیچ هایمان هزینه بالایی ندارند.

۲- در روش قبل برای ۸ پردازنده ۶۴ سوئیچ می خواستیم ولی در اینجا ۱۲ سوئیچ می خواهیم.

برای  $n$  پردازنده  $\frac{n}{2}$  سطر داریم یعنی برای ۸ پردازنده ۴ سطر داریم.

تعداد ستون ها هم برابر تعداد بیت های ما باشد پس برابر  $\log n$  می باشد که ستون اول به بیت اول، ستون دوم به بیت دوم، ستون سوم به بیت سوم .

$$\text{تعداد سوئیچ ها} = \frac{n}{2}$$

$$\text{تعداد بیت} = \text{تعداد ستون} = \log n$$

برتری ها:

تعداد کمتری سوئیچ نیاز دارد. بنابراین هم مدار سوئیچ ساده است و هم قیمت کمتری دارد.

## مشکلات:

## ۱- Internal blocking

فرض کنید ۱۰۰ می خواهد از پردازنده P0 وارد سوئیچ شود. ۱ وارد سوئیچ اول شده و از پورت دوم سوئیچ خارج می شود. سپس وارد پورت اول سوئیچ دوم می شود و اینجا به مشکل بر می خورد چون این سوئیچ قبلاً توسط ۱۰۱ اشغال شده است. با وجود اینکه حافظه ۱۰۰ اشغال نیست ولی امکان اتصال وجود ندارد. در اینجا Internal blocking رخ داده است.

**راه حل:** اگر ورودی ها را مرتب کنیم Internal blocking رخ نمی دهد. اما در اینجا درخواست ها مرتب نبودند. در شبکه Banyan برای مرتب کردن درخواست ها از یک شبکه Batcher قبل از شبکه Banyan استفاده می کنیم. ( منظور از شبکه اتصال آنهاست نه شبکه کامپیوتری) یعنی اینکه ابتدا شماره درخواست هایمان از حافظه را به Batcher می دهیم تا آنها را به ترتیب ۱ و ۴ و ۵ مرتب کند و مشکل Internal blocking حل می شود.

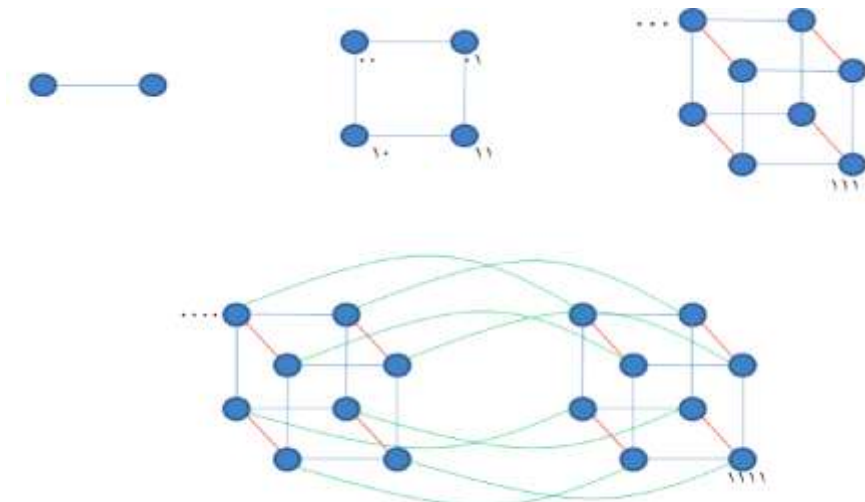
## ۲- output contention

پردازنده های نمی توانند تشخیص دهند کدام حافظه اشغال و کدام حافظه آزاد است. درخواست می دهند ولی امکان دارد حافظه اشغال باشد و در خروجی برخورد بوجود آید. یعنی دو پردازنده نمی توانند همزمان به یک حافظه وصل شوند. که به این مشکل output contention گویند.

## ۳-۲-۵: Hyper cube: روش

در این روش پردازنده هایی که شماره باینری آنها فقط و فقط در یک بیت تفاوت دارند به هم وصل می شوند.

اگر  $n$  پردازنده داشته باشیم ، هر پردازنده به  $\log_2^n$  پردازنده دیگر متصل می شود.



به بحث پردازش موازی (Parallel) مرتبط است.

البته قبل از آن یک مفهوم جدید را مطرح می کنیم و دنبال این هستیم که موازی حل شود. (مسئله کوله پشتی، مسئله کوتاهترین مسیر و.....)

شبکه ارتباطات را با گراف مقایسه می کنیم. پردازنده ها با همدیگر با گراف مدل می شوند. مباحث در هر دو بکار می رود. یعنی در چند پردازنده ای و شبکه کاربرد دارد. مثلا Banyan در ارتباط بین پردازنده ها و حافظه ها بود. عینا در سوئیچ ها و روترها هم استفاده می شود.

این بحث هم در ارتباط پردازنده ها و هم در ارتباط شبکه ها مطرح می شود.

NOC: Network On Chip

SOC: System On Chip

SOC: تمام چیزهایی که داخل یک سیستم بود، داخل یک Chip به کار می رود.

Chip همان تراشه یا IC است. مثلا در Computer ما پردازنده داخل یک Chip است. ولی یک Chip هم برای کارت گرافیک می گذاریم و یک Chip هم برای LAN و.... و می گوئیم کل یک سیستم داخل یک Chip پیاده سازی می شود. همه چیز، حافظه، ارتباط با دستگاه ها، همه چیز داخل یک Chip قرار می گیرد که به آن SOC می گوئیم و یک بحث قدیمی است.

NOC: داخل یک IC ما چند پردازنده داشته باشیم، چند حافظه داخل یک Chip داشته باشیم و این ها با همدیگر ارتباط داشته باشند. یعنی شبکه ای بینشان موجود باشد که خود بحث مفصلی دارد. اینکه شبکه بین آنها چطور باشد. Switchها چطور باشد، چون در شبکه واقعی Switch داریم.

تمرین: ۱- سوئیچ بچربانیاں Batcher-Banyan را با رسم مدار توضیح دهید.

**سوئیچ بانیاں (Banyan switch)**

یک سوئیچ چند مرحله ای است که در آن میکروسوئیچهای هر مرحله، اطلاعات را با استفاده از دنباله باینری آدرس، به سمت خروجی سوئیچ می کنند. اگر ورودی سوئیچ صفر باشد، ورودی مستقیما به خروجی وصل می شود. اگر ورودی سوئیچ یک باشد، ورودی به پایین وصل می شود.

اگر ورودی ها ۲ تا صفر یا ۲ تا یک باشد، مشکل بلوکه شدن داخلی اتفاق می افتد یکی از راه هایی که برای حل این مشکل ارائه شده، استفاده از مرتب کننده قبل از سوئیچ می باشد .

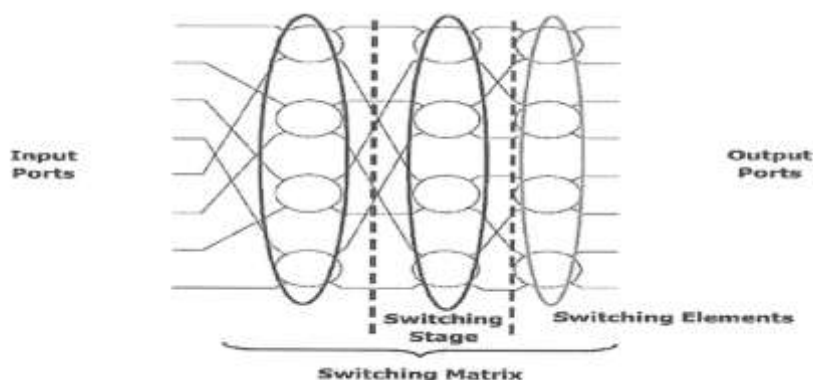
از مزایای سوئیچ بانیان میتوان سادگی مدار و همچنین کم هزینه بودن به دلیل تعداد کمتر سوئیچ را نام برد. برای  $N$  ورودی و  $N$  خروجی به  $\log N$  در مبنای ۲ مرحله نیاز داریم در هر مرحله به  $N/2$  میکروسوئیچ نیاز است، در شکل ساختار این سوئیچ دیده می شود.

به عبارت دیگر مثلا برای ۸ پردازنده به ۱۲ سوئیچ نیازمندیم:

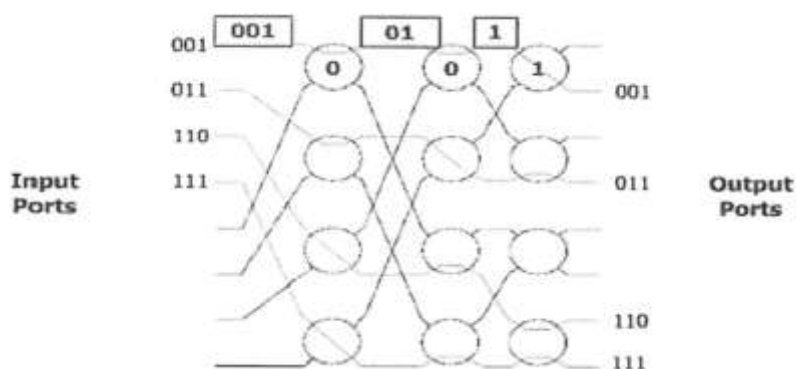
$$\text{تعداد سوئیچ ها} = \frac{n}{2} \cdot \log n$$

$$\text{تعداد سطرها} = \frac{n}{2}$$

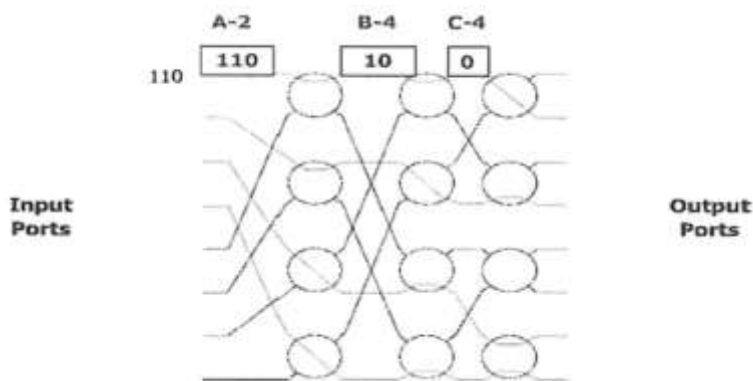
تعداد ستون ها همان تعداد بیت هایمان می باشد یعنی  $\log n$



مرحله اول، بسته اطلاعاتی را بر اساس با ارزشترین بیت آدرس، مسیریابی و هدایت می کند، مرحله دوم بر اساس بیت بعدی، الی آخر.



فرض کنید، بسته ای وارد باب ۱ شده است و آدرس بسته ۱۱۰ یا باب ۶ است، اولین میکروسوئیچ (یا A-2) بر اساس اولین بیت که یک است آن را مسیریابی می کند. دومین میکروسوئیچ (B-4) آن را بر اساس دومین بیت که باز هم یک است. مسیریابی می کند و سومین میکروسوئیچ (یا C-4) بر اساس سومین بیت (صفر) آن را مسیریابی می نماید.



و لذا، بسته به باب خروجی (۶) مطابق آدرسش می رسد .

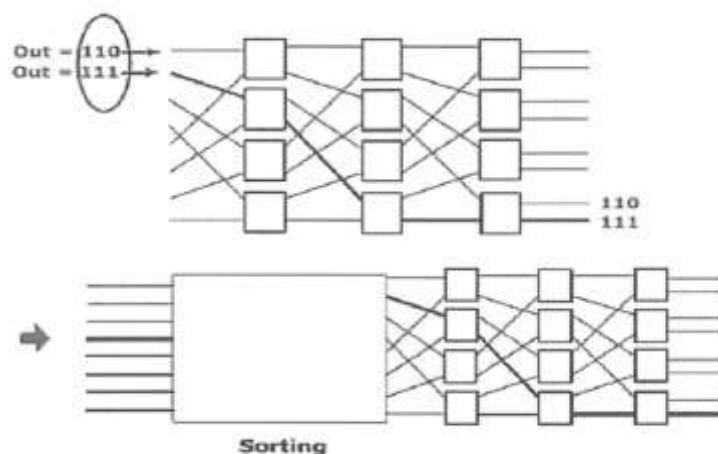
### سوئیچ بچر – بانیان (Batcher – Banyan switch)

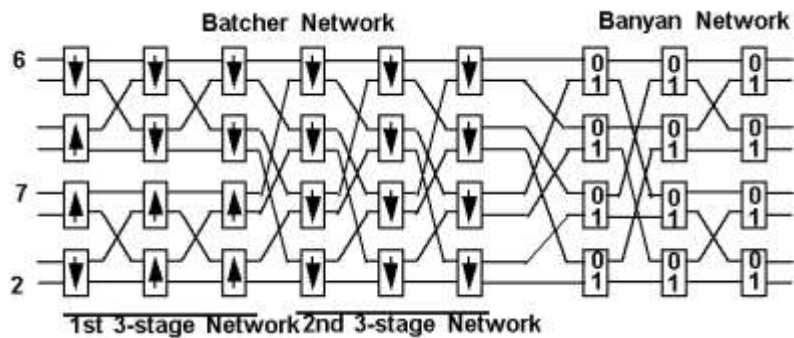
مشکلی که در سوئیچ Banyan وجود دارد که حتی اگر دو بسته آدرسهای خروجی یکسان نداشته باشند ( به باب خروجی یکسانی فرستاده نشود ) ممکن است به علت یک برخورد داخلی نتوانند مسیر یابی شوند .

این مشکل با ردیف کردن (sort کردن) بسته ها با توجه به باب مقصد ، توسط Batcher حل شده است .

Batcher سوئیچی طراحی نمود که قبل از سوئیچ Batcher می آید و بسته ها را متناظر به آدرس نهایی آنها ردیف می نماید که به مجموعه این دو سوئیچ Batcher – Banyan سوئیچ گفته می شود .

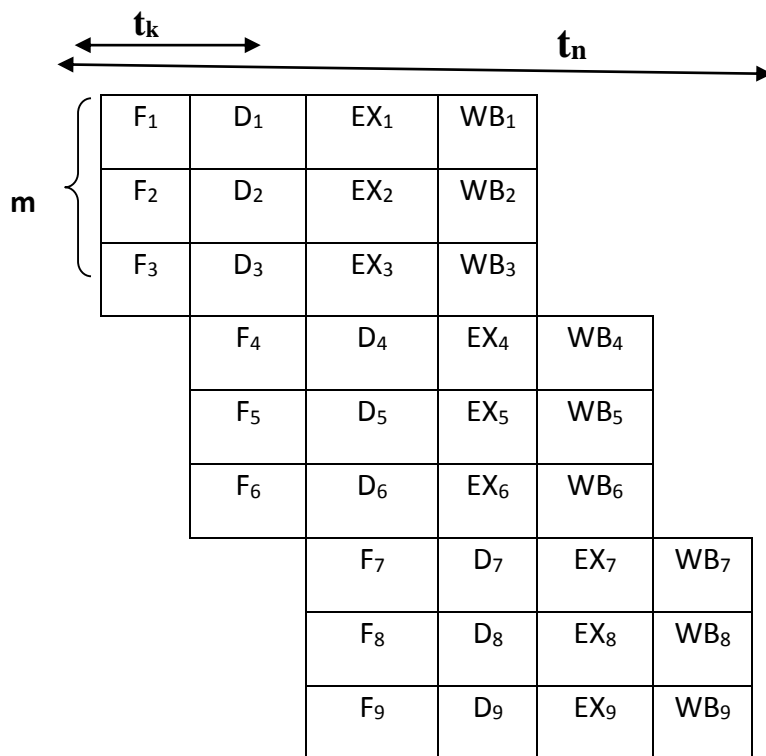
یکی از مشکلاتی که سوئیچ بانیان داشت این بود که احتمال برخورد داخلی حتی برای دو بسته ای که به پورت خروجی یکسانی ارسال شده اند وجود دارد. برای حل این مشکل می توان بسته های ورودی را بر اساس پورت مقصد (خروجی) مرتب کرد. برای این کار از سوئیچ مرتب کننده Batcher قبل از سوئیچ Banyan استفاده می شود. در این سوئیچ از یک قطعه سخت افزاری به نام تله Trap استفاده می کنند. این قطعه از عبور همزمان بسته های تکراری به سوئیچ بانیان جلوگیری می کند. در هر زمان برای هر مقصد فقط یک بسته می تواند ارسال شود، اگر تعداد این بسته ها بیش از یکی باشد باید منتظر بماند.





## ۲-۴ Super scalar

در super scalar همزمان چند دستور با هم اجرا می شوند، پس مراحل اجرا به صورت زیر است. مثلاً در این شکل تعداد دستوراتی که همزمان اجرا می شوند برابر  $m$  یعنی ۳ تا است. یعنی ۳ دستور با هم fetch می شوند و سپس با هم decode می شوند و سپس اجرا و back write می شوند.



$$t_n = kt_k$$

$$m = 3$$

زمان اجرای  $n$  دستور  $= t_n + t_k + \dots$

$m$  دستور اول

$$\frac{n-m}{m} \text{ ها } t_k \text{ تعداد}$$

شرح فرمول بالا: زمان اجرای  $m$  دستور اول  $t_n$  و  $\frac{n-m}{m}$  دستور بعدی هر کدام یک  $t_k$  را به زمان اجرا اضافه می کنند. یعنی به تعداد  $\frac{n-m}{m}$  در این فرمول  $t_k$  داریم.

$$\text{super scaler با } n \text{ دستور اول با زمان اجرای } = kt_k + \frac{n-m}{m} t_k = t_k \left( k + \frac{n-m}{m} \right) = t_k \left( k + \frac{n}{m} - 1 \right)$$

$$\text{pipe line با } n \text{ دستور با زمان اجرای } = t_k (n + k - 1)$$

$$\text{دستور بدون خط لوله } = nkt_k$$

محاسبه کارایی super scaler نسبت به super pipeline:

$$\frac{\text{کارایی Super scaler}}{\text{کارایی pipeline}} = \frac{\text{زمان اجرای با خط لوله}}{\text{زمان اجرای با super scaler}}$$

$$= \lim_{n \rightarrow \infty} \frac{t_k (n + k - 1)}{t_k \left( k + \frac{n}{m} - 1 \right)} = m$$

✓ نتیجه: کارایی super scaler نسبت به خط لوله  $m$  برابر بیشتر است.

$m$ : تعداد مراحل که موازی اجرا می شوند.

### ۳-۵ Super pipeline

در خط لوله مشابه شکل زیر یک دستور در ۴ مرحله اجرا می شود.

F <sub>1</sub>	D <sub>1</sub>	EX <sub>1</sub>	WB <sub>1</sub>		
	F <sub>2</sub>	D <sub>2</sub>	EX <sub>2</sub>	WB <sub>2</sub>	
		F <sub>3</sub>	D <sub>3</sub>	EX <sub>3</sub>	WB <sub>3</sub>

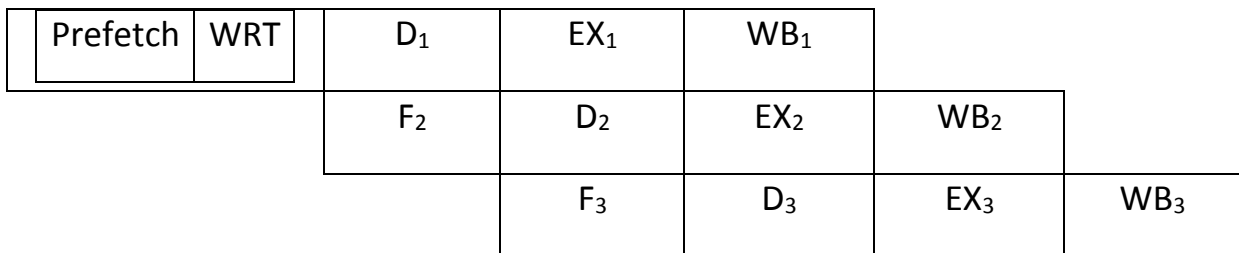
pipeline مزایای زیادی دارد پس برای بهبود کارایی در خط لوله می خواهیم تعداد قطعات را زیاد کنیم اما به ۲ دلیل نمی توان چون اگر تعداد قطعات بیشتر شود:



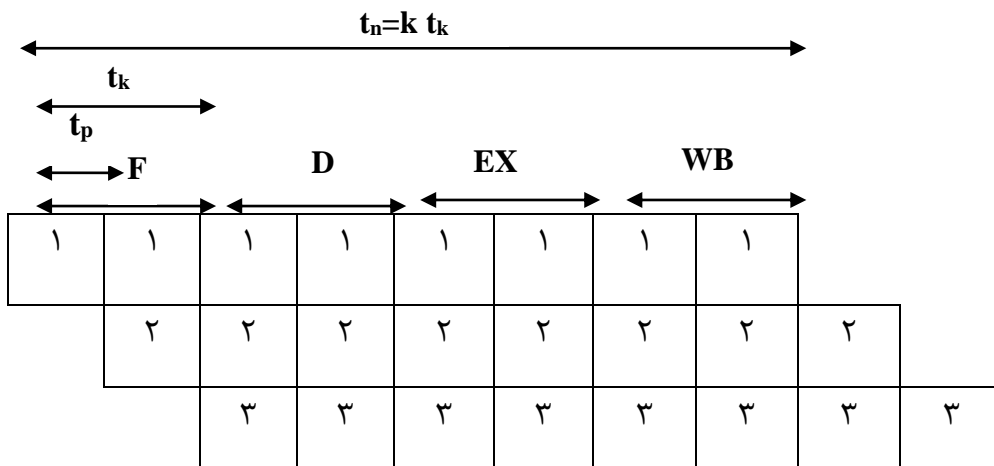
۱- مخاطرات بیشتر می شود.

۲- هزینه ها بالاتر می رود به دلیل اینکه در خط لوله از روش سنکرون استفاده می کنیم و نیاز به latch داریم.

به دلیل این مشکلات می گوییم به جای اینکه fetch را در یک مرحله انجام دهی fetch را به صورت چند مرحله ای انجام بده یعنی در چند مرحله کوچکتر انجام شود و بقیه قطعات نیز به همین ترتیب به مراحل و قطعات کوچکتری تقسیم می شوند. مثلاً fetch به دو مرحله prefetch و WRT=write to register تقسیم می شود.



مثلاً می توان هر قطعه را به دو قطعه کوچکتر تقسیم نمود مانند شکل زیر:



$P=2$

$t_k = p t_p$

با دستور  $n$  زمان اجرای  $t_n + t_p + \dots + t_p = k p t_p + (n - 1) t_p = (k p + n - 1) t_p$

super pipeline

یک  $t_n$  برای دستور اول و تعداد  $(n-1)$  تا  $t_p$  برای  $n-1$  دستور بعدی. چون تعداد دستورات را  $n$  در نظر گرفتیم با اجرای دستور اول  $n-1$  دستور باقی می ماند.

$$\frac{\text{کارایی } Super\ pipeline}{\text{کارایی } pipeline} = \frac{\text{زمان اجرای باخط لوله}}{\text{زمان اجرای با } super\ pipeline}$$

$$= \lim_{n \rightarrow \infty} \frac{pt_p(n+k-1)}{t_p(kp+n-1)} = p$$

در واقع در superpipeline تعداد مراحل pipeline را با روش دیگر وهزینه کمتر، افزایش می دهیم. مثلا هر مرحله pipeline را ۲ مرحله ای می کنیم.

✓ نتیجه: Super pipeline نسبت به خط لوله معمولی p برابر سریعتر اجرا می شود.

#### ۴-۵ Super pipeline-scalar

معماری superpipeline ارتقا روش pipeline است که در آن منابع دو برابر می باشد. تا در نتیجه هم پوشانی بیشتری بین مراحل ایجاد گردد و در نهایت تعداد بیشتری دستور در یک کلاک انجام گیرد. در این تکنیک هر مرحله به دو زیر مرحله تقسیم میشود. در این روش Fetch دستور اول می تواند همزمان با Fetch دستور دوم اجرا شود. کارایی و بهره وری این روش بالاتر است. زیرا دو عمل مشابه می توانند همزمان اجرا شوند. و نیز زمان انتظار و زمان پاسخ دهی کاهش می یابد.

اشکال:

این روش یک اشکال عمده دارد، آن هم نیاز به سخت افزار اضافه است که باعث می شود هزینه بالاتر رود. و نیز افزایش سخت افزار حدی امکان پذیر است. زیرا وابستگی را بالا برده و نیز باعث بالا رفتن سربار Hazard میشود.

فرمول میزان بالا رفتن سرعت در معماری superpipeline:

$$S = ((1 - n) + k^2) / (((1 - n) * k) * Z)$$

k : تعداد مراحل pipeline

n : تعداد دستور

معماری superscaler علاوه بر ویژگی های معماری superpipeline ، توانایی اجرای همزمان چندین دستور العمل را در یک مرحله از خط لوله را داراست.

پردازنده چند شاخه که بصورت پویا زمانبندی و برنامه ریزی میشوند را پردازنده superscaler می گویند.

در ساده ترین نوع این پردازنده ها، دستورات بصورت ترتیبی اجرا می شوند. و خود پردازنده تصمیم می گیرد که دو سیکل ساعت مورد نظر، یک یا چند دستور العمل را روی خط لوله بفرستد.

دستیابی به کارایی خوب superscaler به تلاش هوشمندانه کامپایلر در راستای مهمانپذیری و برنامه ریزی دستورات و از میان برداشتن وابستگی بین آنها نیاز دارد، تا به این صورت نرخ ارسال همزمان دستورات بهبود یابد.

فرمول میزان بالا رفتن سرعت در معماری superscaler (در صورتی که مرحله خط لوله به دو مرحله تقسیم شود):

$$S = (2 / (1 - n) + k) / ((1 - n) + k)$$

**تفاوت کارایی روش superscaler-superpipeline:**

روش Super Scaler کارایی بالاتری نسبت به superpipeline دارد، ولی دو برابر شدن سخت افزار باعث دو برابر شدن کارایی نمی شود که به دلیل Hazard است.

## فصل چهارم

### ارزیابی کارایی پردازنده ها

#### مقدمه:

این فصل توضیح خواهد داد که چگونه کارایی سرعت یک کامپیوتر را اندازه بگیریم و آن را گزارش کنیم. در این فصل سعی خواهیم نمود که به سؤالات زیر جواب دهیم:

- چرا کارایی مهم است؟
- چگونه ما می توانیم کارایی را به دقت تعریف کنیم؟
- چگونه طراحی سخت افزار کارایی نرم افزار را تحت تأثیر قرار میدهد؟
- چگونه در دنیای واقعی کارایی اندازه گرفته شود؟
- چرا بعضی از سخت افزارها بهتر از بقیه عمل می کنند
- کدام یک از فاکتورهای کارایی به سخت افزار وابسته اند؟
- ما برای اجرا شدن سریعتر یک برنامه به یک ماشین جدید نیاز داریم یا نیازمند یک سیستم - عامل جدید هستیم؟
- مجموعه دستورات یک ماشین چگونه کارایی را تحت تأثیر قرار میدهند؟

در این فصل عامل های تعیین کننده در کارایی کامپیوترها مورد بحث قرار خواهد گرفت. بررسی کردن کارایی به این دلیل مهم است که کارایی سخت افزار اغلب کلیدی ترین اثر را در کل سیستم متشکل از سخت افزار و نرم افزار برعهده دارد. ما معمولاً در این فصل کارایی و سرعت را به جای هم استفاده می کنیم. ولی در بعضی از موارد این دو عبارت کاملاً معادل هم نیستند. به عنوان مثال اگر هدف ما استفاده از ماشینی باشد که تعداد کار بیشتری را در واحد زمان انجام دهد، ماشینی که سریعتر است بعضی مواقع نمی تواند هدف ما را برآورده کند و ما مجبور هستیم در این مواقع از ماشینی استفاده کنیم که قابلیت های بیشتری داشته باشد (کارا تر باشد). (ما در این کتاب اغلب از اصطلاح کارایی استفاده خواهیم کرد.

اظهار نظر کردن در مورد کارایی یک سیستم واقعاً کار مشکلی است. سیستم های بزرگ نرم افزاری با جزئیات زیاد به همراه بازه وسیعی از تکنیک های افزایش سرعت که توسط طراحان سخت افزار به کار گرفته می شوند، از جمله مواردی هستند که صحبت کردن درباره کارایی را مشکل می نمایند. تقریباً کار غیر ممکن است که برگه راهنمایی از مجموعه دستورات یک کامپیوتر را بردارید و با یک برنامه مشخص، تعیین نمایید که آن کامپیوتر برنامه شما را چقدر

سریعتر اجرا خواهد کرد. در واقع برای برنامه های کاربردی مختلف، ممکن است مترهای مختلفی مناسب باشد و قسمت های مختلف یک کامپیوتر ممکن است نقش متفاوتی در تعیین سرعت کل سیستم داشته باشند.

البته در انتخاب بین کامپیوترهای مختلف، کارآیی تقریباً همیشه یک مشخصه مهم به حساب می آید. اندازه گیری و مقایسه ماشین های مختلف برای خریداران و در نتیجه برای طراحان یک امر حیاتی محسوب میشود که این مطلب را فروشندگان کامپیوترها به خوبی درک می کنند. اغلب، فروشندگان دوست دارند که شما کامپیوترهای آنها را خیلی ساده نگاه کنید و این نگاه کردن ممکن است بدون توجه به نیاز خریدار، که برنامه های کاربردی خودش را بر روی آن اجرا خواهد کرد انجام بگیرد. در بعضی مواقع ادعاهایی در مورد کامپیوترها می شود که دید مفیدی برای بعضی از برنامه های کاربردی در اختیار قرار نمی دهند. بنابراین فهمیدن چگونگی اندازه گیری کارآیی یک کامپیوتر و محدودیت های اندازه گیری در انتخاب یک ماشین، بسیار مهم می باشند. علاقمندی ما نسبت به کارآیی کامپیوتر ماورای این است که فقط کارآیی کامپیوتر را از بیرون آن اندازه بگیریم. ما نیاز به این داریم که بفهمیم چه چیزهایی کارآیی یک کامپیوتر را تحت تأثیر قرار می دهند.

فهمیدن اینکه چرا بعضی قسمت های نرم افزار همان طور که ما انتظار داریم کار نمی کنند، چرا یک مجموعه از دستورات می تواند به گونه ای پیاده سازی شود که بهتر عمل کند، و اینکه چگونه بعضی از قابلیت های سخت افزار کارآیی را تحت تأثیر قرار می دهند، همگی از مسائلی است که ما علاقمند به پیدا کردن جواب مناسب برای آنها هستیم.

### تعریف کارآیی

وقتی که گفته می شود کارآیی این کامپیوتر بالاتر از آن یکی است، منظور چیست؟ هر چند این سؤال ساده به نظر می رسد ولی یک مقایسه با هواپیماهای مسافربری نشان می دهد که این سؤال چقدر قابل تأمل است.

۱ تعدادی از هواپیماهای مسافربری را به همراه سرعت متوسط آنها، برد مسافتی و ظرفیتشان - شکل ۲ نشان می دهد. اگر بخواهیم بینیم که کدام یک از هواپیماهای موجود در جدول کارآیی بالاتری دارد، اول باید کارآیی را تعریف کنیم. به طور مثال، با در نظر گرفتن اندازه گیریهای مختلف، می بینیم که و هواپیمایی که بیشترین ظرفیت را دارد 747، concordie هواپیمایی که بالاترین سرعت متوسط را دارد می باشد.

۱. مقایسه کارآیی هواپیماهای مختلف - شکل ۲

بنابراین اگر ما بخواهیم کارآیی این دو هواپیما را مقایسه کنیم. به دو صورت می توانیم این کار را انجام دهیم. می توانیم به این صورت تعریف کنیم که هواپیمای کارآتر هواپیمایی است که یک مسافر را از نقطه ای به نقطه دیگر در کمترین زمان جابجا میکند. پس اگر شما علاقمند باشید که زودتر به مقصد کارآتر می باشد. روش دوم مقایسه، برای این مثال تعداد مسافر جابجا concordie برسد واضح است که

شده در واحد زمان است. بر اساس این معیار همان طور که در ستون آخر جدول نیز نشان داده شده است 747 کارا تر می باشد.

Plane	DC to Paris	Speed	Passengers	Throughput (pmp)
747	6.5 hours	610 mph	470	286,700
Concorde	3 hours	1350 mph	132	178,200

به طور مشابه ما می توانیم کارایی کامپیوترها را به طرق مختلف تعریف نمائیم. اگر شما یک برنامه را بر روی دو کامپیوتر مختلف اجرا کنید، خواهید گفت که کامپیوتری سریعتر است که آن برنامه را سریعتر اجرا می کند. ولی اگر برنامه خود را بر روی یک ترمینال در مرکز کامپیوتری اجرا کنید که دارای دو کامپیوتر است که به طور مشترک توسط چندین نفر استفاده می شوند در این صورت خواهید گفت که کامپیوتری سریعتر است که تعداد زیادتری برنامه را در روز اجرا کند. همان طور که یک کاربر که مدت زمان بین (execution time) یا همان زمان اجرا (response time) دوست دارد که زمان پاسخ شروع و خاتمه یک کار است، بر روی یک کامپیوتر برای برنامه او کمتر شود، در یک مرکز کامپیوتر، افزایش پیدا کند. (throughput) مدیر مرکز اغلب دوست دارد که تعداد کار انجام گرفته در واحد زمان —

**مثال:** برنامه دلخواه ما بر روی کامپیوتر A که فرکانس کلاک آن 400 مگاهرتز است در 10 ثانیه اجرا می شود. ما میخواهیم به یک طراح کامپیوتر کمک کنیم که یک ماشین به نام B بسازد که برنامه ما را در 6 ثانیه انجام دهد. طراح کامپیوتر به ما گفته است که افزایش فرکانس کلاک امکان پذیر است ولی این افزایش فرکانس طراحی بخش های دیگر CPU را تحت تأثیر قرار خواهد داد و باعث خواهد شد که ماشین B برای اجرای این برنامه تعداد کلاک هایی در حدود 1/2 برابر ماشین A نیاز داشته باشد. ما چه فرکانس کلاکی را برای رسیدن به این هدف از طراح بخواهیم؟

$$F_A = 4 \text{ GHz}$$

$$A \text{ زمان اجرا} = 10 \text{ s}$$

$$B \text{ زمان اجرا} = 6 \text{ s}$$

$$1.2 \times \text{تعداد چرخه ساعت} = \text{تعداد چرخه ساعت } B$$

$$\text{آهنگ پالس ساعت} \times \text{زمان اجرا} = \text{تعداد چرخه ساعت}$$

$$A \text{ تعداد چرخه ساعت} = 10 \times 4 \times 10^9$$

$$B \text{ تعداد چرخه ساعت} = 1.2 \times 4 \times 10^9 = 48$$

$$F_B = \frac{\text{تعداد چرخه ساعت}}{\text{زمان اجرا}} = \frac{48 \times 10^9}{6} = 8 \times 10^9 \text{ GHz}$$

### ۱-۴ (Clock per Instructor) CPI

تعداد پالس (کلاک) متوسط لازم برای اجرای هر دستور. (( این مقدار، در هر برنامه با برنامه دیگر متفاوت است.))

تعداد کلاک متوسط لازم برای اجرای هر دستور (CPI) × تعداد دستورات برنامه = تعداد کلاک های لازم برای اجرای یک برنامه

عبارت تعداد کلاک های متوسط مورد نیاز برای هر دستور CPI نامیده می شود. چون دستورات مختلف بسته به اینکه چه کاری انجام می دهند، ممکن است زمان اجرای متفاوتی داشته باشند CPI متوسط زمان اجرای ((با واحد کلاک)) همه دستوراتی است که در داخل یک برنامه قرار دارند. روشی برای مقایسه دو پیاده سازی مختلف از یک مجموعه دستورات را نیز در اختیار قرار می دهد، چون در این دو پیاده سازی تعداد دستورات مورد نیاز برنامه مطمئناً یکی خواهد بود.

**مثال:** فرض کنید ما دو پیاده سازی از یک مجموعه دستورات واحد در اختیار داشته باشیم. ماشین A برای یک برنامه دارای پریود کلاک ۱ نانوثانیه بوده، و CPI آن ۲ می باشد و ماشین B دارای پریود کلاک ۲ نانوثانیه بوده، و CPI آن ۱.۲ می باشد. برای این برنامه کدام ماشین سریعتر است و چقدر؟

ما این تعداد دستوراتی که هر کدام از ماشین ها برای این برنامه اجرا می کنند یکی است. ما این تعداد دستورات را I می نامیم. در ابتدا تعداد کلاک های پردازنده را برای هر کدام از ماشین ها حساب می کنیم:

$$A = I \times 2 = \text{تعداد کلاک های لازم برای اجرای برنامه بر روی ماشین A}$$

$$B = I \times 1.2 = \text{تعداد کلاک های لازم برای اجرای برنامه بر روی ماشین B}$$

حال زمان اجرا برنامه را برای هر کدام از ماشین ها حساب می کنیم:

$$A = \text{پریود کلاک} \times \text{متوسط تعداد کلاکهای لازم برای اجرای برنامه (CPI)} = \text{زمان اجرای برنامه بر روی ماشین A}$$

$$= I \times 2 \times 1 \text{ ns} = 2 \times I \text{ ns}$$

همین طور برای ماشین B داریم:

$$B = \text{پریود کلاک} \times \text{متوسط تعداد کلاکهای لازم برای اجرای برنامه (CPI)} = \text{زمان اجرای برنامه بر روی ماشین B}$$

$$= I \times 1.2 \times 2 \text{ ns} = 2.4 \times I \text{ ns}$$

واضح است که ماشین A سریعتر است چون زمان اجرای پایین تری دارد. برای اینکه بدانیم A چقدر سریعتر از B است، به صورت زیر عمل می کنیم:

$$\frac{\text{کارایی ماشین A}}{\text{کارایی ماشین B}} = \frac{\text{زمان اجرای برنامه بر روی ماشین B}}{\text{زمان اجرای برنامه بر روی ماشین A}} = \frac{2.4 \times \text{Ins}}{2 \times \text{Ins}}$$

بنابراین ماشین A برای این برنامه، ۲۰٪ برابر سریعتر از ماشین B می باشد.

CPI متوسط:

$$\text{CPI متوسط} = \frac{\sum_{i=1}^n \text{CPI}_i}{n}$$

**مثال:** یک طراح کامپایلر می خواهد بین دو قطعه کد برای یک ماشین مشخص یکی را انتخاب نماید، بر اساس پیاده سازی سخت افزار، سه کلاس مختلف از دستورات وجود دارد: کلاس A، کلاس B، و کلاس C، که به ترتیب برای اجرا شدن نیاز به ۱، ۲، و ۳ کلاک دارند، قطعه کد اول دارای ۵ دستور است: ۲ مورد از کلاس A، ۱ مورد از کلاس B، و مورد از کلاس C. و قطعه کد دوم دارای ۶ دستور است: ۴ مورد از کلاس A، ۱ مورد از کلاس B، ۱ مورد از کلاس C:

الف) کدام قطعه کد سریعتر اجرا میشود و چقدر؟

ب) CPI را برای هر قطعه کد حساب کنید؟

الف) برای بدست آوردن تعداد کلاک از فرمول زیر استفاده می کنیم:

$$\text{تعداد کلاک های لازم} = \sum_{i=1}^n (\text{CPI}_i \times C_i)$$

که در آن  $\text{CPI}_i$  تعداد کلاک های لازم برای اجرای دستورات از نوع کلاس  $i$  می باشد، و  $C_i$  تعداد دستورات کلاس  $i$  می باشد.

کلاک  $10 = 2 \times 3 = 1 \times 2 = 2 \times 1 =$  تعداد کلاک های لازم برای قطعه کد ۱

کلاک  $9 = 1 \times 3 = 1 \times 2 = 4 \times 1 =$  تعداد کلاک های لازم برای قطعه کد ۲

چون هر دو قطعه کد بر روی یک ماشین اجرا می شوند و در نتیجه پریود کلاک برای هر دو یکی است، بنابراین قطعه کدی سریعتر است که تعداد کلاک کمتری لازم داشته باشد. بنابراین در این مثال با اینکه قطعه کد ۲ تعداد دستورات بیشتری دارد ولی سریعتر اجرا می شود.

ب) برای بدست آوردن CPI از فرمول زیر استفاده می کنیم:



$$\text{CPI} = \frac{\text{تعداد کلاک ها}}{\text{تعداد دستورات}}$$

$$\text{CPI} = \frac{10}{5} = 2 \text{ قطعه کد اول}$$

$$\text{CPI} = \frac{9}{6} = 1.5 \text{ قطعه کد دوم}$$

مثال:

$$15 = \text{زمان اجرای A}$$

$$\text{تعداد دستورالعمل B} = 0.6 \times A$$

$$\text{CPI}_B = 1.1 \text{ CPI}_A$$

$$\text{اندازه پالس ساعت} \times \text{CPI}_A \times \text{تعداد دستورالعمل A} = \text{زمان اجرای A}$$

$$\text{اندازه پالس ساعت} \times \text{CPI}_B \times \text{تعداد دستورالعمل B} = \text{زمان اجرای B}$$

$$0.6 \times A \times 1.1 \text{ CPI}_A \times \text{اندازه پالس ساعت} = 0.6 \times 1.1 \times 15 = 9.9$$

## ۲-۴ معیار MFLOP (Million Floating Point Operation Per Second):

این معیار تعداد عمل ممیز شناور (برحسب میلیون عمل) در یک ثانیه را نشان داده، و برای کارهای علمی و تحقیقاتی ((که محاسبه دارند)) معیار مهمی است.

این علامت یک مقیاس قدرت محاسباتی است، که معمولاً برای به کامپیوترهای بسیار بزرگ بکار می رود.

### ۱-۲-۴ MIPS (million instructions per second):

این واحد تعداد دستور ((برحسب میلیون دستور)) در یک ثانیه را نمایش داده، که واحدی برای سنجش کارایی یک پردازنده است. این معیار جایگزین معیار زمان است.

یک نوع معماری پردازنده از نوع معماری ریسک بوده، که توسط شرکت میپس تکنولوژی طراحی شده است. نسخه‌های اولیه آن ۳۲ بیتی بوده، و بعدها نسخه ۶۴ بیتی آن هم عرضه شد. چندین نسخه از این معماری وجود دارد که عبارتند از: MIPS 1، MIPS 2، MIPS 3، MIPS 4، MIPS 5، MIPS 32 و MIPS 64. نسخه‌های فعلی MIPS 32 و MIPS 64 هستند. در MIPS 32 و MIPS 64 علاوه بر مجموعه

دستورالعمل‌ها، یک مجموعه ثابت کنترلی هم وجود دارد. پردازنده‌های MIPS عمدتاً در سامانه‌های توکار از جمله دستگاه‌های ویندوز سی‌ای، مسیریاب‌ها، کنسول‌های بازی‌های ویدیویی همانند پلی‌استیشن، پلی‌استیشن ۲ و پلی‌استیشن پرتابل استفاده شده‌اند.

$$\text{MIPS} = \frac{\text{تعداد دستور}}{10^6 \times \text{زمان اجرا}}$$

MIPS در یک کامپیوتر خاص برای برنامه‌های متفاوت، مقادیر متفاوتی خواهد داشت.

$$\text{MIPS} = \frac{\text{تعداد دستور}}{\frac{\text{CPI} \times \text{تعداد دستور}}{\text{نرخ کلاک}} \times 10^6} = \frac{\text{نرخ کلاک}}{\text{CPI} \times 10^6}$$

### ۳-۴ معیار SPEC (System Performance Evaluation Cooperative) :

این معیار مجموعه‌ای شامل ۱۲ "آزمون صحیح" و ۱۷ "آزمون ممیز شناور"، و متمرکز بر کارایی پردازنده بوده، که طی ۵ نسل تکامل یافته است.

برای ساده‌سازی استفاده از نتیجه این آزمون‌ها، نتایج این ۱۲ آزمون صحیح، باهم تلفیق شده، و در قالب یک عدد واحد گزارش می‌شود.

عملاً به عنوان معیار استفاده شده، که یکسری نمودار به ما می‌دهد که با هم مقایسه می‌شوند. این معیار، بهترین معیار برای ارزیابی کارایی یک سیستم واقعی است.

### ۴-۴ قانون امدال (Amdahl's law):

قانون امدال در معماری کامپیوتر قانونی است که می‌توان به کمک آن ارزش کارهایی که قرار است انجام شوند را از قبل پیش‌بینی نمود. این کارها عمدتاً مربوط به بهبودهایی است که در ساختار پردازنده‌ها داده می‌شود. در این قانون، مقدار بهبودهایی که انجام می‌شود به میزان مؤثر بودن آنها محدود می‌شود.

این قانون نشان می‌دهد که مقدار تسریع در یک سامانه پردازش موازی دارای یک "حد حداکثر" است. قانون امدال برای سامانه‌های موازی، با بار محاسباتی ثابت به کار می‌رود. امدال نشان داد که اگر بار محاسباتی ثابت باشد، و در زمانی که واحدهای موازی به بی‌نهایت میل کند، مقدار تسریع دارای یک حد حداکثر است و از آن مقدار بیشتر نمی‌شود.

برای توضیح این قانون، سیستم واحدی دانشگاه را در نظر می‌گیریم. در این سیستم بعضی دروس یک

واحدی، بعضی دو واحدی، بعضی سه واحدی و بعضی نیز چهار واحدی اند. همانطور که می دانیم در محاسبه معدل ترم، دروسی بیشتر تأثیر دارند که تعداد واحد آنها بالاتر است. بنابراین در این سیستم ارزش دروس چهار واحدی بیشتر از سه واحدی، ارزش دروس سه واحدی بیشتر از دو واحدی، ارزش دروس دو واحدی نیز بیشتر از یک واحدی است. بنابراین ما بهتر است بیشترین وقت و سرمایه گذاری را بر روی دروسی انجام دهیم که بیشترین تأثیر را در معدل دارند.

در معماری کامپیوتر قانون امدال می گوید: قبل از اینکه بهبود و یا کاری را انجام دهی باید به میزان مؤثر بودنش نیز توجه کنی، و آن کار را وقتی انجام دهی که تأثیرش در مقابل کاری که انجام می دهی قابل توجه باشد. در قانون امدال، اگر جایی را بهبود دادید، فقط آنجا بهبود می یابد.

در قالب فرمول، قانون امدال به صورت زیر بیان می شود:

$$\text{زمان اجرای بهبود نیافته} + \frac{\text{زمان اجرای بهبود یافته}}{\text{میزان بهبود}} = \text{زمان اجرای برنامه بعد از بهبود}$$

**مثال:** فرض کنید که تصمیم گرفته باشیم سخت افزار پردازنده را تغییر دهیم تا سرعت اجرای دستورات Floating Point دو برابر شود. دو برابر کردن سرعت یک ایده خوب است ولی باید بینیم در قبال کاری که انجام می دهیم چقدر بهبود در کل زمان اجرای یک برنامه به وجود می آید. به عبارتی باید بینیم دستورات Floating Point چه میزان در برنامه ما نقش دارند. فرض کنید برای یک برنامه فرضی فقط ۱۰ درصد زمان اجرای برنامه (T) شامل دستورات Floating Point باشد. برای این برنامه تأثیر بهبود سخت افزار را بررسی کنید.

با استفاده از قانون امدال می توان نوشت:

$$\text{زمان اجرای برنامه بعد از بهبود} = \frac{0.10 T}{2} + 0.90 T = 0.95 T$$

یعنی زمان اجرای کل برنامه فقط ۵ درصد بهبود خواهد داشت. این میزان بهبود خیلی بالا نیست بنابراین بر طبق قانون امدال دو برابر کردن سرعت اجرای دستورات Floating Point به امید بهبود زمان اجرای برنامه در این مثال ارزش فنی ندارد.

نتیجه قانون امدال: قسمت هایی از برنامه را بهبود دهیم که بیشتر استفاده می شوند چون بیشترین تأثیر را

در کل زمان اجرای برنامه دارند. قانون امدال بیان می کند که موارد پر استفاده را بهبود دهیم و یا به عبارتی می گوید: Make the common case fast.

$$\frac{80}{5} + 20 = 36$$

مثال:

	A	B	C
CPI	۱	۲	۳

	A	B	C	
کامپایلر ۱	۲	۱	۲	$5 \times 10^9$
کامپایلر ۲	۴	۱	۱	$6 \times 10^9$

$$F = 4 \text{ GHz}$$

$$۱ \text{ چرخه ساعت کد } = (2 \times ۱ + 1 \times 2 + 2 \times ۳) = 10 \times 10^9$$

$$۲ \text{ تعداد چرخه ساعت کد } = (4 \times ۱ + 1 \times 2 + ۱ \times ۳) = 9 \times 10^9$$

کد ۱ سریع تر است.

$$۱ \text{ زمان اجرای } = \frac{10 \times 10^9}{5 \times 10^9} = 2$$

$$۲ \text{ زمان اجرای } = \frac{9 \times 10^9}{6 \times 10^9} = 1.5$$

$$\text{MIPS } 1 = \frac{5 \times 10^9}{2 \times 10^6} = 2500$$

$$\text{MIPS } 2 = \frac{6 \times 10^9}{1.5 \times 10^6} = 4000$$

MIPS دومی بیشتر است.

تمرین : تفاوت معماری دو معماری RISC و CISC :

۱- بدون آرگومان (پشته ای)

۲- تک آرگومان (انباره ای)

۳- دو آرگومان (CISC)

۴- سه آرگومان (RISC)

تفاوت های CISC و RISC :

با ظهور کامپیوترهای براساس پردازنده PowerPC به عنوان رقیبی جدی برای کامپیوترهای مبتنی بر پردازنده های اینتل این بحث داغ شده است که پردازنده های RISC بهترند یا پردازنده های CISC.

RISC علامت اختصاری کلماتی است که در فارسی به صورت «پردازنده های کم دستورالعمل» ترجمه کرده ایم و بدین معناست که تعداد دستورالعملهایی که چنین پردازنده هایی می توانند بفهمند کمتر از تعداد دستورالعملهای پردازنده های معمول یا پردازنده های مشهور به CISC است. در مقابل، CISC علامت اختصاری کلماتی است که در فارسی به صورت «پردازنده های پر دستورالعمل» ترجمه کرده ایم. سازندگان پردازنده های RISC ادعا می کنند و به طور کلی از CISC ها قویترند. اما اینتل و سایر طرفداران CISC سرسختانه عکس این مطلب می اندیشند.

آنچه در این مباحثات نادیده گرفته شده است اختلافات اصلی بین تراشه های RISC و CISC است. آیا تعداد دستورالعمل اختلاف اصلی بین این دو گروه از پردازنده هاست؟ اگر چنین باشد، آیا بدین معنی نیست که پردازنده RISC نرم افزار را آهسته تر اجرا می کند، چون باید تعداد بیشتری دستورالعمل را برای انجام کارهای پیچیده اجرا کند؟ مزیت تعداد دستورالعمل کم، اگر مزیتی باشد، چیست؟ اصلاً مجموعه دستورالعمل چیست؟

تحقیق زیر سعی می کند به این پرسشها پاسخ دهد و مرزهای بین پردازنده های RISC و CISC را مشخص سازد. اگر تا به حال پاسخ به این سوالات را نمی دانسته اید حقایقی که خواهد آمد ممکن است شما را متعجب کند.

طرحان کامپیوترهای کم دستور (RISC) مدعی اند که به دو تا سه برابر کارایی سیستم کامپیوترهای پیچیده دستور (CISC) با همان میزان پالس ساعت دست یافته اند. این افزایش کارایی از طریق لوله گذاری

(Pipelining) که یک شکل شناخته شده موازی بودن و بهینه سازی مترجم هاست و معماری خاص RISC حاصل شده است. دستوره های RISC طوری طراحی شده که طول پایدار و با ثبات و زمان های اجرای باثباتی دارند. معماری RISC، ثابت بودن طول دستورالعمل و ثابت بودن زمان اجرا باعث می شوند تا بهینه سازی کامپایلر، که عامل اساسی در تعیین سرعت یک سیستم کامپیوتری است، ساده تر شود.

طرفداران CISC، نه! طراحان CISC با دقت در حال ارزیابی خویشند تا دستوراتی را که بیشتر از همه اجرا می شوند را مشخص کنند و با بهینه سازی میکروکد یا پیاده سازی سخت افزاری (hardwiring) دستورالعمل ها، سرعت اجرای آنها را بالا ببرند. مدافعین CISC برای این باورند که می توانند به میانگین سیکل ساعت هر دستور، مانند مشابه آن در طراحی، RISC، دست یابند. این به معنای دستیابی به کارایی، بدون قربانی کردن سازگاری با تعداد وسیع نرم افزارهای کامپیوترهای شخصی، است. (قابل ذکر است که هنوز تعداد نرم افزارهای همه منظوره برای ماشین های RISC بسیار کمتر از انواع CISC است).

افزون بر این بزرگ بودن مجموعه رجیسترهای ماشین های RISC باعث بروز اختلال در یک محیط تک کاربردی چند کاره می شود، چرا که محتوای یک مجموعه رجیستر، هنگامی که کاربر بین وظایف سوئیچ می کند، باید ذخیره شود. در CISC پیچیدگی سخت افزار و در RISC پیچیدگی نرم افزاری وجود دارد.

برای عمل ضرب در سیستم های با معماری ریسک باید حتما از یک واحد سخت افزاری برای ضرب استفاده کرد. در سیستم های RISC برای انجام عملیات ضرب سیستم کامپایلر به این گونه رفتار می کند که یک حلقه FOR ایجاد کرده و ۲ را دو بار در خود ضرب می کند که این امر باعث سادگی سخت افزاری و سرعت بالا می شود. در تلفن های هوشمند از تکنولوژی RISC و در سیستم های دسکتاپ از تکنولوژی CISC استفاده شده است. تکنولوژی RISC باعث کم مصرف شدن باطری و پردازش سریع است و تکنولوژی CISC دارای قدرت پردازش بالا و مصرف انرژی بیشتری می شود. در کامپیوترهای قدیمی مک از پردازنده PPC استفاده می کردند که RISC بود. پردازنده های SPARC هم که در سرورهای Oracle (یا Sun قدیم) ارائه می شدند RISC هستند.

## پیوست ۱

روش CRT - I (روش چینی تسریع شده):

$$X \equiv_{(M_1, M_2, \dots, M_n)} (r_1, r_2, \dots, r_n)$$

X از رابطه زیر بدست می آید:

$$X = r_1 + m_1 < k_1(r_2 - r_1) + k_2 m_2(r_3 - r_2) + \dots + k_{n-1} m_2 m_3 \dots m_{n-1}(r_n - r_{n-1}) >_{m_2 m_3 \dots m_n}$$

که  $k_i$  به صورت زیر با محاسبه معکوس ضربی بدست می آید:

$$< k_1 m_1 >_{m_2 m_3 \dots m_n} \equiv 1$$

$$< k_2 m_1 m_2 >_{m_3 m_4 \dots m_n} \equiv 1$$

$$< k_{n-1} m_1 m_2 \dots m_{n-1} >_{m_n} \equiv 1$$

مثال:

$$X \equiv_{(2,3,5,7)} (1,0,4,0) \Rightarrow |k_1 m_1|_{m_2 m_3 m_4} \equiv 1 \Rightarrow |k_1 \times 2|_{105} \equiv 1 \Rightarrow k_1 = 53$$

$$|k_2 m_1 m_2|_{m_3 m_4} \equiv 1 \Rightarrow |k_2 \times 2 \times 3|_{35} \equiv 1 |6k_2|_{35} \equiv 1 \Rightarrow k_2 = 6$$

$$|k_3 m_1 m_2 m_3|_{m_4} \equiv 1 \Rightarrow |k_3 \times 2 \times 3 \times 5|_7 \equiv 1 \Rightarrow |30k_3|_7 \equiv 1 \Rightarrow k_3 = 4$$

$$X = r_1 + m_1 < k_1(r_2 - r_1) + k_2 m_2(r_3 - r_2) + k_3 m_3 m_2(r_4 - r_3) >_{m_2}$$

$$X = 1 + 2 < 53 \times (0 - 1) + 6 \times 3 \times (4 - 0) + 4 \times 5 \times 3 \times (0 - 4) >_{105}$$

$$X = 1 + 2 \times (1 - 221)_{105} = 1 + 2 \times 94 = 189 \Rightarrow X = 189$$

روش Mixed - Radix conversion:

$$* X \equiv_{(M_1, M_2, \dots, M_n)} (r_1, r_2, \dots, r_n)$$

فرض کنیم  $(M_1, M_2, \dots, M_n)$  پیمانه‌های ما هستند که نسبت به هم اول هستند و  $(r_1, r_2, \dots, r_n)$ مانده‌های ما در پیمان مذکور باشند بطوری که:  $r_i \equiv^{m_i} |X|$  هدف یافتن عدد  $X$  است.فرض کنیم  $X$  نمایشی بصورت

$$** X = z_1 + z_2 m_1 + z_3 m_2 m_1 + \dots + z_n m_{n-1} m_{n-2} \dots m_1$$

داشته باشد که در آن  $0 \leq z_i \leq m_i$  می باشد.

هدف یافتن ضرایب مجهول  $Z_i$  بصورتی است که معادله \* برقرار باشد.

$$\text{اگر از دو طرف رابطه ** به پیمانه } m_1 \text{ هم‌نهنستی بگیریم: } X \equiv^{M_1} Z_1 \quad \text{پس: } z_1 = r_1$$

سپس با محاسبه  $Z_1$  رابطه \*\* را بصورت

$$** X = z_1 + z_2 m_1 + \dots + z_n m_{n-1} m_{n-2} \dots m_1$$

$m_2$  هم‌نهنستی می گیریم که بدست می آید

$$\dots + 0 + 0 + \dots \equiv^{M_2} X - Z_1 \quad \text{چون اینها ضرایب } m_2 \text{ دارند.}$$

$$\text{پس: } X - Z_1 \equiv^{M_2} z_2 m_1$$

با ضرب رابطه فوق در معکوس ضربی  $m_1$  در پیمانه  $m_2$ :

$$|m_1^{-1}|_{m_2} (X - Z_1) \equiv^{m_2} z_2$$

چون می دانیم  $X \equiv^{M_2} r_2$  پس رابطه با  $r$  بصورت:  $Z_2 \equiv^{m_2} |m_1^{-1}|_{m_2} (r_2 - r_1)$  به همین ترتیب  $Z_3$ :

$$Z_3 \equiv^{m_3} |(m_2 m_1)^{-1}|_{m_3} \times (r_3 - (z_2 m_1 + z_1))_{m_3}$$

و به همین ترتیب تا  $Z_n$ :

$$z_n \equiv^{m_n} |(m_1 m_2 \dots m_{n-1})^{-1}|_{m_n} \times (r_n - (z_{n-1} m_{n-2} m_{n-3} \dots m_1 + z_{n-2} m_{n-3} m_{n-4} \dots m_1 + \dots + z_2 m_1 + z_1))_{m_n}$$

با استفاده از رابطه  $|mi mj|_{mk} = ||mi^{-1}|_{mk} |mj|_{mk}|_{mk}$  ی

$$z_1 = r_1 \quad \text{و} \quad z_2 = ||mi^{-1}|_{m_2} (r_2 - z_1)_{m_2}$$

$$* z_3 = ||m_2^{-1}|_{m_3} |m_1^{-1}|_{m_2} (r_3 - (z_2 m_1 + z_1))_{m_3}$$

$$* z_n = ||m_{n-1}^{-1}|_{m_n} (|m_{n-2}^{-1}|_{m_n} (|m_{n-3}^{-1}|_{m_n} (\dots (|m_2^{-1}|_{m_n} (|m_1^{-1}|_{m_n} (r_n - z_1) - z_2 \dots) - z_{n-1})_{m_n}$$

مثال برای روش Mixed-Radix:

$$X \equiv^{(2,3,5,7)} (1,0,4,0)$$

$$R = (1,0,4,0)$$

$$M = (2,3,5,7)$$



$$X = (x_1, x_2, x_3, x_4) = ?$$

$$z_1 = r_1 = 1 \Rightarrow z_1 = 1$$

$$z_2 = ||m_i^{-1}|_{m_2}(r_2 - z_1)|_{m_2} \Rightarrow z_2 = ||2^{-1}|_3(0 - 1)|_3 = |2x - 1|_3 = 1 \Rightarrow z_2 = 1$$

$$* z_3 = ||m_2^{-1}|_{m_3}|m_1^{-1}|_{m_3}|m_3(r_3 - z_1) - z_2|_{m_3} = ||3^{-1}|_5|2^{-1}|_5(4 - 1) - 1|_5|(2) \times ((3) \times 3 - 1)|_5 = 1 \Rightarrow z_3 = 1$$

$$z_4 = ||m_3^{-1}|_{m_4}(|m_2^{-1}|_{m_4}(|m_1^{-1}|_{m_4}(r_4 - z_1) - z_2)z_3)|_{m_4} = ||5^{-1}|_7(|3^{-1}|_7(|2^{-1}|_7(0 - 1) - 1) - 1)|_7 = |3 \times (5 \times (4 \times (-1) - 1) - 1)|_7 = |-78|_7 = 6 \quad z_4 = 6$$

$$X = Z_4 m_3 + m_2 + m_1 + Z_3 m_2 + m_1 + Z_2 m_1 + Z_1 = 6 \times 5 \times 3 \times 2 + 1 \times 3 \times 2 + 1 \times 2 + 1 = 189 \quad X = 189$$

روش معکوس با استفاده از مجموعه پیمانه های  $\{2n + 1, 2n, 2n - 1\}$

این روش پیمانه ها بصورت  $\{2n + 1, 2n, 2n - 1\}$  هستند حالت خاصی از آن پیمانه ها  $\{2^m + 1, 2^m, 2^m - 1\}$  است در سخت افزار و پیاده سازی مدار است کاربرد زیادی دارد.

$$X \equiv_{(m_1, m_2, m_3)} (r_1, r_2, r_3)$$

در این حالت جواب بصورت  $X = |\sum_{i=1}^3 w_i r_i|_{m_1 m_2 m_3}$  خواهد بود که اگر فرض کنیم

$$\begin{cases} M_1 = 2N + 1 \\ M_2 = 2N \\ M_3 = 2N - 1 \end{cases}$$

آنگاه  $w_i$  بصورت زیر محاسبه می شوند.

$$w_1 = \frac{(m_1 + 1)(m_2 m_3)}{2}$$

$$w_2 = m_1(m_2 - 1)m_3$$

$$w_3 = \frac{m_1 m_2 \times (m_3 + 1)}{2}$$

در نهایت خواهیم داشت:

$$X = |\sum_{i=1}^3 w_i r_i|_{m_1 m_2 m_3}$$

$$X \equiv^{(9,8,7)} (1,2,3)$$

مثال:

$$\text{چون} \begin{cases} 7 = 2 \times 4 - 1 \\ 8 = 2 \times 4 \\ 9 = 2 \times 4 + 1 \end{cases} \xrightarrow{\{2n+1, 2n, 2n-1\}} \begin{cases} m_3 = 2n - 1 = 7 \\ m_2 = 2n = 8 \\ m_1 = 2n + 1 = 9 \end{cases}$$

$$\text{طبق فرمولهای قبل} \begin{cases} w_1 = \frac{(7+1)(8 \times 9)}{2} = \frac{560}{2} = 280 \\ w_2 = 7 \times (8-1) \times 9 = 441 \\ w_3 = \frac{7 \times 8 + (10)}{2} = 288 \end{cases}$$

$$X = |w_1 k_1 + w_2 k_2 + w_3 k_3|_{9 \times 8 \times 7}$$

$$X = |280 \times 1 + 441 \times 2 + 288 \times 3|_{504}$$

$$X = |2026|_{504} = 15 \Rightarrow X = 15$$

تمرین: چگونه با RNS تقسیم انجام دهیم؟ ابتدا مفهومی معکوس ضربی را تعریف می کنیم.

تعریف: معکوس ضربی  $X^{-1}$  را معکوس ضربی  $X$  در پیمانه  $M$  می نامیم اگر  $XX^{-1} \equiv^M 1$  مثلا معکوس ضربی  $7$  در پیمانه  $11$  می شود  $8$  چون:  $7 \times 8 \equiv^{11} 1$  برای تقسیم  $X$  در پیمانه  $m$  بر  $y$  بصورت زیر عمل می کنیم:

$$\frac{x}{y} \equiv^m q \quad \text{فرض کنیم}$$

$$\Rightarrow x \equiv^m y \times q$$

اگر دو طرف رابطه بالا را در معکوس ضربی  $y$  ضرب کنیم.

$$\Rightarrow x \times y^{-1} \equiv^m y \times y^{-1} \times q \xrightarrow{y \times y^{-1} \equiv^m 1} x \times y^{-1} \equiv^m q$$

البته باید پیمانه طوری باشد که با  $y$  مقسوم علیه مشترک نداشته باشند چون در اینصورت هنگامی که  $y$  را به پیمانه انتخاب شده ببریم صفر می شود و تقسیم بر صفر تعریف نشده است.

پس از این رابطه می توانیم بجای  $\frac{x}{y} \equiv^m q$  استفاده کنیم حال فرض کنیم پیمانه سه تایی  $(m_2, m_1,$

$m_0)$  را داریم و می خواهیم نشان دهیم RNS  $\frac{x}{y}$  است پس بجای  $\frac{x}{y}$  می توانیم  $x \times y^{-1}$  را به پیمانه  $(m_2,$

$m_1, m_0)$  قرار دهیم:

$$\frac{x}{y^{-1}} \rightarrow x \begin{cases} x = m_2 k_2 + x_2 \\ x = m_1 k_1 + x_1 \\ x = m_0 k_0 + x_0 \end{cases}$$

$$y^{-1} \begin{cases} y^{-1} = m_2 k_2 + y_2^{-1} \\ y^{-1} = m_1 k_1 + y_1^{-1} \\ y^{-1} = m_0 k_0 + y_0^{-1} \end{cases}$$

$$x \times y^{-1} = m_2^2 k_2 k_2 + x_2 \times y_2^{-1} + m_2 k_2 y_2^{-1} + m_2 k_2 x_2 = m_2 (m_2 k_2 k_2 + k_2 y_2^{-1} + k_2 x_2) + x = m_2 k_2 + x_2 \times y_2^{-1}$$

$$x \times y^{-1} = m_1^2 k_1 k_1 + m_1 k_1 y_1^{-1} + m_1 k_1 x_1 + x_1 y_1^{-1} = m_1 (m_1 k_1 k_1 + k_1 y_1^{-1} + k_1 x_1) + x = m_1 k_1 + x_1 \times y_1^{-1}$$

$$x \times y^{-1} = m_0^2 k_0 k_0 + m_0 k_0 y_0^{-1} + m_0 k_0 x_0 + x_0 y_0^{-1} = m_0 (m_0 k_0 k_0 + k_0 y_0^{-1} + k_0 x_0) + x = m_0 k_0 + x_0 \times y_0^{-1}$$

$$X \times y^{-1} \equiv_{(m_2, m_1, m_0)} (x_2 \times y_2^{-1}, x_1 \times y_1^{-1}, x_0 \times y_0^{-1})$$

RNS؟ مثال: تقسیم ۲۰ بر ۵ به پیمانه (۷و۳و۲) با روش

ابتدا معکوس ضربی ۵ را به پیمانه ۷ و ۳ و ۲ بدست می آوریم:

$$5 \times 3 \equiv^7 1$$

پس معکوس ضربی ۵ به پیمانه ۷ برابر ۳ می شود

$$5 \times 2 \equiv^3 1$$

پس معکوس ضربی ۵ به پیمانه ۳ برابر ۲ می شود

$$5 \times 1 \equiv^2 1$$

پس معکوس ضربی ۵ به پیمانه ۲ برابر ۱ میشود

پس معکوس ضربی ۵ به پیمانه (۷و۳و۲) برابر (۳و۲و۱) می شود حال حاصلضرب ۲۰ در معکوس ضربی ۵ را بدست می آوریم

۲۰ به پیمانه (۷و۳و۲) برابر (۶و۲و۰) می شود پس :

حاصلضرب به پیمانه (۷و۳و۲) می شود: (۴و۱و۰)