

۱.

الف) حالتی که در آن پردازنده بی کار (idle) است و دستوری انجام (execute) نمی دهد و تنها بخش مربوط به تشخیص وقفه ها و دیکود از الگوریتم فون نایمن انجام می شود)

ب) با اجرای دستور مربوطه در ISA توسط سیستم عامل، چون این دستور در kernel mode اجرا می شود. با وقفه خارجی (external interrupt) نیز از آن خارج می شود که میتواند توسط سخت افزار خارجی تولید شود.

ج) به طور کلی وقتی scheduler سیستم عامل، پروسسی برای انجام نداشته باشد به این حالت می رود. مانند وقتی که پردازنده منتظر I/O است و کار دیگری برای انجام ندارد، برای کاهش توان مصرفی و دمای پردازنده به این حالت میرود.

الف)

استفاده از interrupt-chaining به این صورت که برای هر چند I/O یک generic interrupt routine در نظر بگیریم (مشابه برداری) که در صورتی که یکی از آن ها درخواست وقفه داد، تک تک آن ها را بررسی کند. (مشابه سر کشی)

ب)

سرکشی: وقتی دستگاه I/O در فاصله زمانی بسیار کم آماده باشد (حالت busy کوتاه یا به طور مکرر درخواست داشته باشد) یا رسیدگی به درخواست I/O ها مهم (critical) باشد؛ مانند از دست رفتن داده (پرسیدن بافر) یا صدمه به دستگاه (مشکل در برق یا دما).

وقفه: وقتی فاصله زمانی میان مواقعی که I/O آماده است زیاد باشد، همچنین I/O ها حیاتی نباشند.

روش بیان شده: وقتی تعداد I/O ها زیاد و تأثیر سربار برداری (که فضای اشغالی است) زیاد می شود.

۲.

سخت افزاری: ایجاد مشکل در منبع تغذیه، ارور خواندن از حافظه (ارور parity به طور مثال)، درخواست I/O (کار با کیبورد و موس)

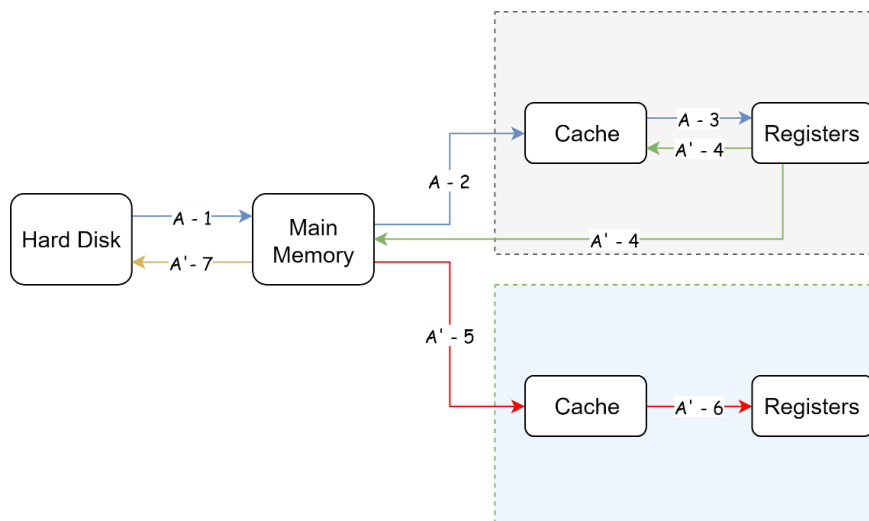
نرم افزاری: فراخوانی سیستمی (system call) برای درخواست انجام کاری توسط سیستم عامل از طرف برنامه کاربر، تولید ارور (exception یا trap) در برنامه (مانند تقسیم بر صفر یا ورودی غیر منتظره)

۳.

الف) ایجاد مکانیسمی برای اولویت دهی به وقفه ها (interrupt priority leveling) در سیستم عامل. این سیستم میتواند این گونه باشد که رسیدگی به درخواست کاربر ۱ یا ۲ حین رسیدگی به درخواست کاربر ۳ انجام شود و همچنین رسیدگی به درخواست کاربر ۱ حین رسیدگی به درخواست کاربر ۲ نیز انجام شود. (nested interrupt)

ب) در صورتی که کاربر با اولویت بالاتر مدام قبل از اتمام رسیدگی کامل به کاربر با اولویت پایین تر درخواست دهد. در این صورت درخواست کاربر اولویت پایین تر هیچ وقت انجام نمی شود.

بهتر است حافظه اصلی write-back و کش write-through باشد (با فرض UMA و اینکه کش بین پردازنده ها مشترک نیست)؛ یعنی همزمان با اعمال تغییر در کش پردازنده اول،  $A'$  در حافظه اصلی هم نوشته شود تا پردازنده دوم بتواند از مقدار آن استفاده کند. چون در دیسک سخت نوشته نشد، سرعت چندان کاهش نمیابد (نوشتن در دیسک سخت بعد از اتمام برنامه یا نیاز به جایگزین شدن داده آپدیت شده در رم انجام می شود). لازم به توضیح است، پردازنده اول بعد از خواندن  $A$  باید بیت dirty آن را در رم set کند تا پردازنده دوم بداند داده تغییر کرده ولی آپدیت نشده و در صورت لزوم منتظر بماند که نسخه آپدیت شده (یعنی  $A'$ ) در حافظه اصلی نوشته شود و همچنین در پایان داده در هارد آپدیت شود.

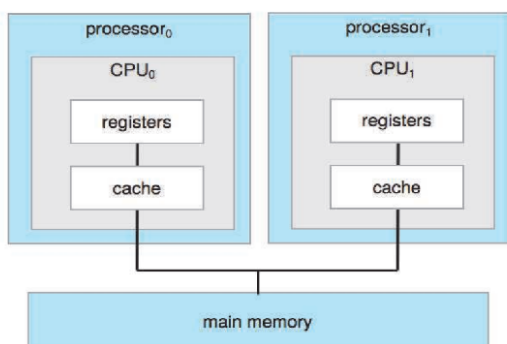


(لیبل های روی فلش ها نشان دهنده داده و زمان جابجایی هستند)

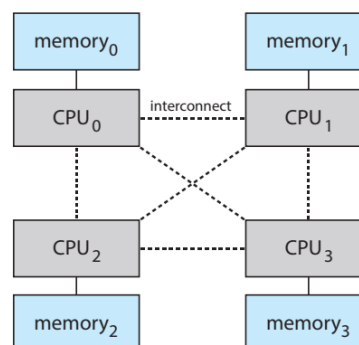
در مواقع انجام اعمال I/O با داده های حجیم بین حافظه اصلی و جانبی، به جای استفاده از وقفه برای جابجایی هر بلوک داده از حافظه جانبی به پردازنده و از آن جا به رم، با مشخص کردن بافر و نشانگر (pointer) و شمارنده برای I/O، کنترل کننده دستگاه می تواند داده ها را مستقیماً و بدون دخالت پردازنده به رم منتقل کند (البته چون رابط رم مشترک با پردازنده است باید برای جابجایی هر بسته چند کلاک باس را قرض بگیرد (cycle stealing) که در این مدت اگر پردازنده رم را بخواهد باید منتظر بماند) و در پایان نیز برا مطلع کردن پردازنده از اتمام فرآید، یک وقفه تولید کند. چون همزمان با این جابجایی داده، پردازنده می تواند به اجرای دستورات دیگر پردازد و با داده های کش خود کار کند، استفاده از DMA باعث افزایش همروندی می شود.

برای افزایش بهره وری در سیستم های چندپردازنده ای از NUMA استفاده می شود. در UMA تمام پردازنده ها دسترسی یکنواخت به حافظه اصلی دارند. بدلیل وجود باس مشترک برای دسترسی به این حافظه، پردازنده ها باید منتظر هم بمانند. پس با افزایش تعداد پردازنده ها کارایی پایین می آید.

در NUMA، برای رفع این مشکل و افزایش بهره‌وری در این سیستم‌ها، از حافظه‌های جداگانه (محلی) برای هر پردازنده استفاده شده و پردازنده‌ها به هم متصل می‌شوند. سرعت دسترسی هر پردازنده به حافظه محلی خود بیشتر و بدون انتظار است. مشکل NUMA آن است که سرعت دسترسی یک پردازنده به حافظه محلی پردازنده دیگر کم است و باعث کاهش بهره‌وری کل سیستم می‌شود، که البته با مدیریت دستورات اجرا شده در هر پردازنده توسط سیستم عامل تأثیر این مشکل را می‌توان کم کرد.



UMA multiprocessing architecture.



NUMA multiprocessing architecture.

۷.

(با فرض واحد های SI)

$$\frac{10 \text{ MB/s}}{4 \text{ KB}} = 2.5 \text{ K/s} \text{ تعداد بلوک در یک ثانیه}$$

$$2.5 \text{ K/s} \times 5000 = 12.5 \text{ M/s} \text{ در ثانیه DMA تنظیمات انجام کلاک}$$

$$\frac{12.5 \text{ M/s}}{500 \text{ MHz}} = 2.5\% \text{ درصد زمان سپری شده توسط پردازنده برای انتقال}$$

$$\frac{2.5 \times 2^{10} \times 5000}{500 \times 10^6} = 2.56\% \text{ (اگر واحد های داده ها بر حسب کیبی بایت و مبی باید باشد، پاسخ نهایی 2.56\% خواهد شد)}$$