

پاییز ۹۳

« به نام راستگوی بی همتا »  
مبانی کامپیوتر و برنامه‌سازی  
پروژه نهایی



دکتر هاشمی و دکتر مرادی

## پایگاه داده

### مقدمه

هر برنامه‌ای پس از اجرا با دریافت و دستکاری مداوم اطلاعات سروکار دارد و یک سری خروجی تولید می‌کند که بعضی از این اطلاعات باید نگهداری شود؛ چرا که نتیجه‌ی محاسبات هستند یا خودشان ورودی بخش‌های دیگری از این برنامه یا برنامه‌های دیگرند. این ذخیره‌سازی اطلاعات باید ویژگی‌هایی داشته باشد تا بتوان به داده‌های ذخیره شده اعتماد و از آنها استفاده کرد. به طور مثال این ذخیره‌سازی باید به نحوی باشد که بازیابی آن به راحتی و با سرعت بالا امکان پذیر شود چرا که برنامه‌های دیگری در فازهای مختلفی از اجرا به این اطلاعات نیاز دارند که باید به محض دریافت درخواست، در زمان مناسب این اطلاعات در اختیارشان قرار گیرد. همچنین پس از مدت طولانی که حجم داده‌ها زیاد می‌شود، باید همچنان سرعت و کارایی قابل قبولی داشته باشند. تصور کنید سرورهای شرکت‌های بزرگ مثل گوگل در کسری از ثانیه مقدار قابل ملاحظه‌ای داده را بازیابی و ذخیره می‌کنند. به عبارت دیگر اضافه، حذف و جست‌وجو کردن داده‌ها باید با سرعت بالا و اطمینان امکان پذیر باشد. چالش دیگری که در نگهداری داده‌ها با آن مواجه هستیم، حفاظت از داده‌ها در مقابل خراب شدن به دلیل اتفاقات غیر قابل پیش‌بینی (مثل قطعی برق، حملات خرابکارانه و...) و ایجاد امکان بازیابی در صورت وقوع است. این خراب شدن ممکن است در هنگام ذخیره‌سازی، تغییر یک داده و یا پس از اینکه داده در حافظه قرار گرفت، اتفاق بیافتد. فرض کنید شما در حال ذخیره یک سری داده هستید که برق دستگاه شما قطع می‌شود. در این حالت اگر به داده‌ها آسیبی وارد شود (به عنوان مثال تنها نصفی از آنها ذخیره شوند)،

دیگر قابل استفاده نیستند. هزینه‌هایی را متصور شوید که شرکت‌ها، بانک‌ها و ... در صورت بروز چنین مشکلی خواهند پرداخت. روش‌های مختلفی برای مدیریت پایداری داده‌ها هست و معمولا در این مواقع داده‌ها به آخرین وضعیت پایدار بازگردانده می‌شوند. که این امر معمولا با نگه‌داشتن پشتیبان و ذخیره کردن سلسله تغییراتی که بر داده‌ها اعمال شده امکان پذیر می‌شود. مشکل دیگر مدیریت دسترسی است. فرض کنید یک سری داده قرار است بین چند نرم افزار که توسط چند برنامه نویس مختلف نوشته شده است به صورت مشترک استفاده شود. برای اینکه همه آن‌ها بتوانند از این داده‌ها استفاده کنند، باید نحوه خواندن و نوشتن اطلاعات با یک استاندارد از قبل تعریف شده باشد تا بی‌نظمی در ساختار داده‌ها به وجود نیاید. تصور کنید اگر هر برنامه ای روشی مخصوص به خود را برای ذخیره سازی داده‌ها داشته باشد، استفاده مشترک از داده‌ها بسیار مشکل خواهد بود. همچنین در این صورت، تمامی مواردی که در بالا ذکر شد (پایداری داده‌ها، امکان پشتیبان‌گیری و ...) باید توسط تک تک نرم افزارهای نوشته شده به صورت مجزا پیاده سازی و مدیریت شود!

پس درواقع برای نگهداری داده‌ها نیازمند یک ساختار منظم و سامان مند هستیم. پایگاه‌های داده، ساخته شده‌اند تا این نیاز را برآورده کنند، یعنی با حجم بسیار زیادی از داده‌ها کار کرده، قابلیت ذخیره سازی، بازیابی و مدیریت کردن آن‌ها را به همراه تمامی نکاتی که در بالا ذکر شده فراهم نمایند. در اینجا به تعریف پایگاه داده می‌پردازیم.

**پایگاه داده** مجموعه‌ای از رکوردهای ذخیره شده در رایانه با یک روش سیستماتیک است. پایگاه‌های داده از فیلدها، رکوردها و فایل‌ها ساخته می‌شوند. به عنوان مثال یک دفترچه تلفن یک **پایگاه داده** است که لیستی از رکوردها را نگهداری می‌کند که هر رکورد سه فیلد نام، آدرس و شماره تلفن را شامل می‌شوند.

UserID	First Name	Last Name	Email	Phone #
7500848	Stephen	Barrett	sbarrett@mail.com	555-222-3987
7500843	Derek	Clapton	derek@dominos.com	555-735-2406
7500843	John	Didsbury	jdisbury@mail.com	555-769-3987
7500847	Georgia	Grace	gg@mail.com	555-859-9876
7500841	Carly	Rose	crose@mail.com	555-403-1018

Fields

Each row is a separate record

برنامه رایانه‌ای که برای مدیریت و پرسش و پاسخ بین پایگاه‌های داده‌ای و برنامه‌های دیگر استفاده می‌شود را سیستم مدیریت پایگاه داده‌ای یا به اختصار (DBMS<sup>1</sup>) می‌نامیم.

درواقع DBMS درخواست نرم افزارها را مدیریت می‌کند و داده‌هایی که نیاز دارند را از پایگاه داده در اختیار آنها قرار می‌دهد بدون اینکه نرم افزارها بدانند که داده‌ها به صورت فیزیکی در کجا و چگونه نگهداری می‌شوند. همچنین DBMS ها پایداری، امنیت و دیگر ویژگی‌های یک پایگاه داده را تضمین می‌کنند. IBM DB2, MongoDB و Oracle PostgreSQL از نمونه DBMS های موفق هستند.

نرم افزارها باید بتوانند با استاندارد زبانی خاصی با DBMS صحبت کنند. این صحبت کردن برای ایجاد نوعی مجموعه داده یا ایجاد، تغییر و حذف یک رکورد و یا گرفتن داده‌ها از پایگاه داده نیاز است. به این زبان، زبان جست و جو (Query Language) گفته می‌شود. در سال ۱۹۸۶ زبان SQL که نقش هر ۳ نوع زبان بالا را بازی می‌کرد به عنوان یک استاندارد در ANSI قرار گرفت.

پس به طور کلی می‌توان گفت یک پایگاه داده، مجموعه‌ای شامل رکورد‌های مختلف است در صورتی که یک DBMS، سیستمی است که مجموعه‌ای از این پایگاه‌های داده را نگه‌داری و مدیریت می‌کند و امکان برقراری ارتباط با نرم افزارهای دیگر را فراهم می‌کند. همانطور که یک فایل سیستم، فایل‌های مختلف را در یک سیستم عامل مدیریت می‌کند، DBMS هم پایگاه‌های داده‌ای مختلف را مدیریت می‌کند.

<sup>1</sup> Database Management System

گاهی اوقات برای راحتی و آژه دیتابیس بجای DBMS استفاده می شود.

## نگهداری داده ها

این که داده ها بر روی دیسک سخت یا حافظه اصلی کامپیوتر نگه داری شوند، DBMS ها را به ۲ دسته مهم Memory-Resident Databases و On-Disk Databases تقسیم می کند.

در پایگاه های داده ای که به صورت On-Disk کار می کنند، داده ها در دیسک سخت نگهداری می شوند. سرعت این دیسک ها بسیار پایین تر از حافظه اصلی کامپیوتر است. برای بالا بردن سرعت، معمولاً از تکنیک نگهداری موقت داده ها در حافظه اصلی (Caching)، بافر کردن و شاخص گذاری<sup>۲</sup> استفاده می شود. دسترسی به دیسک ۱۰۰۰۰۰ بار کندتر از دسترسی به حافظه اصلی است و CPU های معمولی امروزه می توانند در زمان یک دسترسی به دیسک، تعدادی بسیار بیش تر از ۱ میلیون پردازش را انجام دهند. با توجه با این معلوم می شود که گلوگاه (bottleneck)<sup>۳</sup> در دیتابیس ها دسترسی به دیسک است و به دلیل این زیاد بودن هزینه ی دسترسی به دیسک سخت در مقابل دسترسی به حافظه ی اصلی، (و انجام محاسبات توسط CPU)؛ طراحان سعی می کنند با روش های مختلف، از CPU cycle و حافظه ی اصلی هزینه بدهند تا میزان دسترسی به دیسک را کاهش پیدا کند و کارایی مطلوب حاصل شود.

پایگاه های داده In-Memory داده ها را در حافظه اصلی کامپیوتر (RAM) نگه داری می کنند. این موضوع باعث می شود خواندن و نوشتن با سرعت بسیار بالاتری انجام شود. این نوع از پایگاه های داده در چند سال اخیر با توجه به ارزان شدن حافظه های اصلی و اهمیت سرعت بالا، بسیار مورد توجه قرار گرفته و رقابت میان شرکت های بزرگ صنعت پایگاه داده مانند SAP، IBM، oracle و ... در این زمینه مدتی است که آغاز شده

---

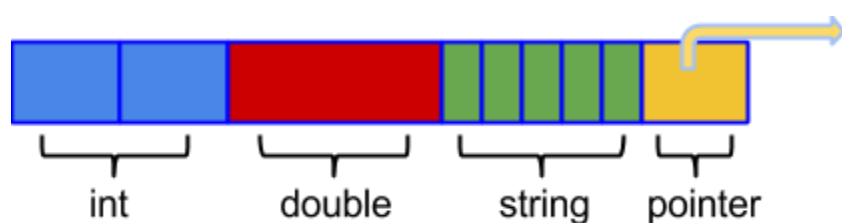
<sup>2</sup> indexing

بخشی از یک برنامه که بیشترین تاثیر را روی سرعت اجرا و کارایی دارد<sup>3</sup>

است. و مسائل مربوط نیز از جمله عناوین داغ تحقیقاتی در مباحث پایگاه داده هستند. در این نوع از پایگاه های داده، ۲ نکته قابل توجه است:

۱. چون همه داده ها در حافظه اصلی نگهداری می شوند، نیاز به حافظه زیاد است و همانطور که می دانید، تهیه حافظه اصلی با حجم بالا گران قیمت تر از تهیه همان حجم با دیسک سخت است.

۲. یکی از ویژگی هایی که برای نگه داری داده ها گفته شد، پایداری آن ها در مقابل اتفاقات است. در این نوع از DBMS ها اگر برق قطع شود، تمام داده ها از دست می روند. مدیریت این موضوع یکی از نکاتی است که باید در طراحی این نوع از DBMS ها مدنظر باشد. ( لازم به ذکر است که اخیرا چنین حافظه هایی با قابلیت حفظ اطلاعات حتی در صورت نبود برق نیز معرفی شده اند<sup>۴</sup>)



داده ها در این نوع از DBMS ها به صورت بخشی از حافظه ی اصلی بوده و می توان آن ها را به صورت لیست های

پیوندی نگه داری کرد. شکل فوق، به عنوان مثال، یک رکورد را نشان می دهد که دارای فیلدهای مختلفی است. در آخرین قسمت مربوط به این رکورد، اشاره گری به رکورد بعدی قرار می گیرد. حال فرض کنید می خواهیم رکوردی که در قسمت string آن "ICSP" نوشته شده است را پیدا کنیم، با توجه به خاصیتی که لیست های پیوندی دارند، باید از اولین گره شروع کنیم و به جلو حرکت کنیم. فرض کنید در بین ۱ میلیارد رکورد دنبال ۱ رکورد خاص هستیم. با توجه به نکاتی که در بالا گفته شد، سرعت بالا در زمان ذخیره سازی و بازیابی یکی از نکات بسیار مهم است. استفاده از ایندکس و Hash Table که در فاز قبلی پیاده سازی کردید اینجا بسیار اهمیت پیدا می کند. فرض کنید نرم افزاری در بانک اطلاعاتی و در میان مجموعه داده هایی که تمامی رکورد

<sup>4</sup> non-volatile random access memory

های آن ساختار بالا را دارند، دنبال همان کلمه "ICSP" می گردند. به میسر کردن این دسترسی سریع به کمک ساختار داده‌های کمکی را ایندکسینگ<sup>۵</sup> (فهرست‌بندی یا شاخص‌گذاری) گفته می شود. روش های مختلفی برای ایندکسینگ وجود دارد. یکی از انواع ایندکس گذاری، Hash indices است. در پایگاه های داده با توجه به استفاده ای که هر فایل دارد، بروی بعضی از فیلدهای آن ایندکسینگ را فعال می کنند. فرض کنید در مثال بالا، می دانیم جست‌وجو و بازیابی اطلاعات با شرط اسم خاص بسیار پر کاربرد است. اگر یک Hash Table برای این فیلد و برای این مجموعه رکورد خاص ساخته شود، در زمان بازیابی فقط کافیست مقدار Hash متناظر کلمه ای که دنبال آن هستیم را پیدای کنیم و مستقیماً به آن آدرس مراجعه کنیم.

در DBMS هایی که بر پایه حافظه اصلی هستند، با توجه به اینکه تمامی داده ها در حافظه نگهداری می شوند، لازم است که به نحوی از آنها در مقابل خرابی و قطعی برق حفاظت کنیم. برای اینکار در زمان‌هایی مشخص یک تصویر<sup>۶</sup> (کپی) از پایگاه داده را در دیسک ذخیره می کنند ( این فرآیندها بسیار زمان‌بر هستند و معمولاً به صورت مستقل و به صورت موازی با DBMS انجام می شوند؛ پیچیدگی‌های مربوط به این مساله از جمله مسائلی است که همچنان مورد تحقیق است) و تا زمان ذخیره‌سازی بعدی query هایی که داده‌ها را تغییر داده‌اند ذخیره می کنند<sup>۷</sup>. به این طریق در صورت از دست رفتن داده‌ها می توان با بارگذاری مجدد آخرین snapshot و اجرای دوباره‌ی این تغییرات بخش اعظم داده‌ها را بازیابی کرد.

آنچه که شما در نهایت در پروژه‌ی نهایی خود پیاده‌سازی میکنید درواقع یک کتابخانه برای کمک به مدیریت داده‌ها در برنامه‌های به زبان C است که از جهات بسیاری یادآور یک پیاده‌سازی بسیار ساده از سیستم پایگاه داده ساکن در حافظه‌ی اصلی است. البته با تفاوت‌های اساسی از جمله اینکه این سیستم، یک فرآیند مستقل

---

<sup>5</sup> indexing

<sup>6</sup> snapshot

<sup>7</sup> logging

نبوده و به عنوان بخشی از برنامه‌ی کاربر اجرا می‌شود (یک کتابخانه است) و کاربر با اجرای توابع کتابخانه‌ی شما بر روی ساختارهای مربوط انبوه داده‌های خود را مدیریت می‌کند.