

类型	缓存什么	被缓存在何处	延迟(周期数)	由谁管理
CPU寄存器	4字节或8字节字	芯片上的CPU寄存器	0	编译器
TLB	地址翻译	芯片上的TLB	0	硬件MMU
L1高速缓存	64字节块	芯片上的L1高速缓存	4	硬件
L2高速缓存	64字节块	芯片上的L2高速缓存	10	硬件
L3高速缓存	64字节块	芯片上的L3高速缓存	50	硬件
虚拟内存	4KB页	主存	200	硬件+OS
缓冲区缓存	部分文件	主存	200	OS
磁盘缓存	磁盘扇区	磁盘控制器	100 000	控制器固件
网络缓存	部分文件	本地磁盘	10 000 000	NFS客户
浏览器缓存	Web页	本地磁盘	10 000 000	Web浏览器
Web缓存	Web页	远程服务器磁盘	1 000 000 000	Web代理服务器

图 6-23 缓存在现代计算机系统中无处不在。TLB: 翻译后各缓冲器(Translation Lookaside Buffer); MMU: 内存管理单元(Memory Management Unit); OS: 操作系统(Operating System); AFS: 安德鲁文件系统(Andrew File System); NFS: 网络文件系统(Network File System)

6.4 高速缓存存储器

早期计算机系统的存储器层次结构只有三层: CPU 寄存器、DRAM 主存储器和磁盘存储。不过, 由于 CPU 和主存之间逐渐增大的差距, 系统设计者被迫在 CPU 寄存器文件和主存之间插入了一个小的 SRAM 高速缓存存储器, 称为 L1 高速缓存(一级缓存), 如图 6-24 所示。L1 高速缓存的访问速度几乎和寄存器一样快, 典型地是大约 4 个时钟周期。

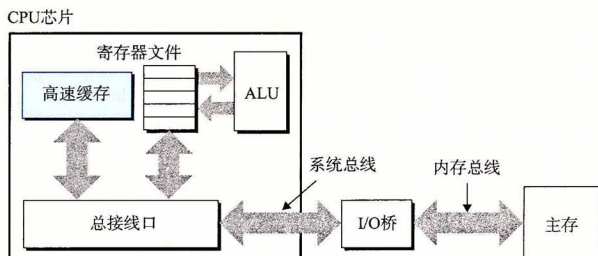


图 6-24 高速缓存存储器的典型总线结构

随着 CPU 和主存之间的性能差距不断增大, 系统设计者在 L1 高速缓存和主存之间又插入了一个更大的高速缓存, 称为 L2 高速缓存, 可以在大约 10 个时钟周期内访问到它。有些现代系统还包括有一个更大的高速缓存, 称为 L3 高速缓存, 在存储器层次结构中, 它位于 L2 高速缓存和主存之间, 可以在大约 50 个周期内访问到它。虽然安排上有相当多的变化, 但是通用原则是一样的。对于下一节中的讨论, 我们会假设一个简单的存储器层次结构, CPU 和主存之间只有一个 L1 高速缓存。

6.4.1 通用的高速缓存存储器组织结构

考虑一个计算机系统, 其中每个存储器地址有 m 位, 形成 $M=2^m$ 个不同的地址。如图 6-25a 所示, 这样一个机器的高速缓存被组织成一个有 $S=2^s$ 个高速缓存组(cache set)的