



链滴







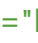


# 计算机中存储体系的设计

作者: [vcjmhg](#)

原文链接: <https://ld246.com/article/1588142182996>

来源网站: [链滴](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

<p>看到这里可能大家会有疑问，三种存储器两两组合方式组中结果不是应该是三种吗？</p>

<p>或者说应该还有一种组合方式即 Cache+ 磁盘存储器？</p>

</blockquote>

<p>坦白来讲此种存储体系理论上讲是可以有的，但是实际应用中，由于 Cache 和虚拟存储器的速度差别太大，强行相互结合，根本不可能发挥出 Cache 存储器的速度优势，实际上还会大大拖累 Cache 的运行。因此此种组合方式在实际应用过程中根本不可能存在。</p>

<p></p>

<h2 id="存储效率">存储效率</h2>

<p>关于速度的评定我们一般通过，<strong>访问周期、存取周期、存储周期、存取时间</strong>等来进行表示。</p>

<p>关于这些指标的计算我们详细可以参照<a href="https://ld246.com/forward?goto=https%3F%2Fbaidu.com%2Fitem%2F%25E5%25AD%2598%25E5%258F%2596%25E5%259125A8%25E6%259C%259F%2F6686295%3Ffr%3Daladdin" target="\_blank" rel="nofollow ugc">百度百科</a></p>

<p>提到了存储器的速度，我们就提另一个指标 <code>存储效率</code>。</p>

<p>首先我们给出命中率的定义：</p>

<blockquote>

<p>在<span class="language-math">M\_1</span>存储器中达到的访问效率，给出以下公式</p>

</blockquote>

<div class="language-math">H=N\_1/\left(N\_1+N\_2\right)</div>

<p>其中：</p>

<p><span class="language-math">N\_1</span>是<span class="language-math">M\_1</span>存储器访问到的效率。</p>

<p><span class="language-math">N\_2</span>是对<span class="language-math">M\_2</span>存储器的访问次数</p>

<p>进而我们引出存储系统的访问效率：</p>

<div class="language-math">e=\frac{T\_1}{T}=\frac{T\_1}{H\cdot T\_1+\left(1-H\right)\cdot T\_2}=\frac{1}{H+\left(1-H\right)\cdot\frac{T\_2}{T\_1}}=f\left(H,\frac{T\_2}{T\_1}\right)</div>

<p>从上边的公式中我们可以得出提高存储系统的两条途径：</p>

<ol>

<li>提高命中率 H</li>

<li>两个存储器的速度差别不应太大</li>

</ol>

<p>在实际操作中有时第二条途径做不多（如虚拟存储器），因此，在实际操作过程中我们通过<strong>提高命中率来提高存取效率。</strong></p>

<blockquote>

<p>那此处可能就又有一个问题：如何提高命中率？</p>

</blockquote>

<p><strong>我们一般采用 <code>预取技术</code> 来提高命中率</strong></p>

<p>措施：在未命中时，把<span class="language-math">M\_2</span>存储器中相邻几个单元成的一个数据块都取出来送入到<span class="language-math">M\_1</span>存储其中。</p>

<p>在引入 <code>预取技术</code> 之后我们可以的出新的命中率公式：</p>

<div class="language-math">H^{\prime}=\frac{H+n-1}{n}</div>

<p>其中：</p>

<p><span class="language-math">H^{\prime}</span>是采用预取技术之后的命中率</p>

<p><span class="language-math">H</span>是原来的命中率</p>

<p><span class="language-math">n</span>是数据块大小与数据重复使用次数的乘积。</p>

<p>为了更加深入的理解命中率我们做下边一个例题</p>

<blockquote>

<p>例：在一个 Cache 存储系统中，当 Cache 的块大小为一个字时，命中率为  $H = 0.8$ ；假设数据重复利用率为 5，Cache 的块大小为 4 个字</p>

<ol>

- <li>计算 Cache 存储系统的命中率是多少? </li>
- <li>假设  $T_2=5T_1$ , 分别计算访问效率。 </li>

</ol>

</blockquote>

<p></p>

## <p>存储体系具体操作过程中可能会遇到**容量不足**或者**速度不够**的问题, 下边我们分成两个专题分别对其进行讨论和解决。</p> <p>第一种方案: **增加主存容量**</p> <p>解决容量不足的问题, 我们最容易想到的便是**增加主存容量**。正如缺啥补, 缺少容量我们增加主存容量必然可以在一定程度上解决容量不足的问题, 但是提到任何一种方案, 考虑经济因素便是耍流氓。我们要知道主存的价格并不是非常便宜, 随着主存容量的增加, 存储器价总量必然增加。从而导致单位容量存储器价格增加。显然这不是最终的解决的方案。</p> <p>第二种方案: 采用两级存储器</p> <p>既然仅仅增加主存容量我们不能根本上解决该问题, 因此我们也提出了第二种方案来进行补充--主存和辅存结合在一起采用两级存储结构, 利用**低价格**的辅存来扩充存储容量。</p> <p>当然能够用辅存扩充容量且能达到主存的速度, 其背后是有如下依据的: </p> <p>首先计算机中所有的信息可以分为如下三类: </p> - <ol> - <li>活跃的信息, 即当前正在使用的信息</li> - <li>待命的信息, 将要使用的信息</li> - <li>静止的信息, 已被使用而不再处理的信息</li> </ol> <p>因此我们按照不同的使用频率, 可以将不同的信息放置到不同的存储器中。可将活跃的信息和部待命的信息放在主存中, 其余部分放到辅存, 以减少主存容量的要求, 从而实现扩充容量的同时又不显著的增加成本。</p> <p>由于信息本身的使用频率不是固定的, 或者说同一个信息在不同的时间段可能被分为不同的类别因此在使用两级存储之后, 主辅存之间的信息调出和调入可能会有一些问题。</p> <p>程序的调入和调出由程序员考虑和安排, 一定程度上增加了程序员的负担。因此使用辅助机构自定位, 从而引出了**虚拟存储器**。</p> <p>可以将高速辅存(如磁盘)伪装成主存访问, 信息在主存、辅存之间的调进和调出完全由辅助机自动完成。</p> <p>存储器的速度往往是整个计算机速度的第一个瓶颈。</p> <p>方法一, 同样也是从主存入手, 提高主存的速度。此法也在采用, 但此法随存储器速度的提高成也会显著增加。</p> <p>方法二也是引入二级存储结构, 只不过引入的是 `Cache+主存` 此体系的结构如下: </p> <p></p> <p>让 CPU 直接与它速度匹配的 Cache 访问, 此法也需要 CPU 和 cache 之间利用辅助机构来完成ache 与主存之间信息的调入与调出。</p> <blockquote> <p>上边解决两个问题都需要用到辅助机构, 那么我们到此处可能会有疑问了, 那辅助机构到底有什么功能那? </p> </blockquote> <p>首先, 第一个是**地址映像功能**: </p>

解决  $M_2$  中的信息采用何种规则调入到  $M_1$  中（即调入规则的问题）。



第二个功能叫做地址变换功能：

根据映像规则，将包括  $M_2$  在内的大空间的地址转成 CPU 能够直接访问的  $M_1$  中的地址（即地址变换问题）。

第三个功能是进行页面替换：

在  $M_1$  中装满信息的条件下，采用何种算法，算出用  $M_1$  的部分信息，使  $M_2$  中的部分功能调到  $M_1$  中（即替换算法问题）。

关于这三部分功能下边我们会分成三个小专题分别来进行讲述。

## 存储体系中的页式管理

因为存储器本身的物理结构并不适合信息的存储于组织，因此我们在对存储器使用之前需要对其象出一层逻辑结构来进行信息组织，因此此处我们引入 `页式管理` 的概念。

## 页的定义

首先我们要明白 `页式管理` 中“页”是何含义。

`页式管理` 中将 **虚拟存储空间** 和 **实际存储空间** 等分成 **固定大小** 的页，使虚拟页可装入主存中不同的实际页面位置。

## 页式管理中的地址表示

`页式管理` 将物理的存储结构进行一层抽象，引入了“页”来进行管理，从而得对存储地址的访问也需要某种新的地址表示。

- 虚地址（逻辑地址，程序地址）：包括  $M_2$  在内的空间地址。



其中，

$N_v$ ：虚页号

$N_r$ ：页内地址

- 实地址（物理地址）：为 CPU 能直接访问的  $M_1$  中地址。



其中：

$n_v$ ：实页号

$n_r$ ：页内地址

当然引用两种地址表示不是目的，为了解决 `虚地址` 和 `实地址`，因此下边我们引入了 `页表` 的概念。

页表是一种特殊的数据结构，放在系统空间的页表区，存放逻辑页与物理页帧的对应关系。

`页表` 具有以下特点：

- `页表` 所需行数和虚页号数相等，虚页号与页表行号对应，因此无需虚页号字



。</li>

<li><code>页表</code> 中每行内容分为两个字段：实页号<span class="language-math"> $n_v$ </span>以及装入位(1 位，其中 0 表示该页已装入，1 表示该页未被装入)</li>

</ol>

<p>针对上边的两个特点我们可以设计出结构如下的页表：</p>

<p></p>

<p>此时可能会有同学问了，<strong>虚页号到哪里去了？？</strong></p>

<p>这个很容易解释，根据第一个特点，虚页号和页表行号一致的，而页表行号实际上是从零开始逐“+1”的。也即我们看到的上边那个两行的页表实际上等同于下边这种三行结构的页表。</p>

<p></p>

<p>## 页式管理的地址变换（关键）</p>

<p>① 根据<span class="language-math"> $N_v$ </span>去查页表中的某一行  $m$ 。<br>

② 查该行的装入位。<br>

③ 装入位=1 时，命中。表示该虚页已装入。</p>

<ul>

<li>从该行中送出<span class="language-math"> $n_v$ </span>（实页号）。</li>

<li>再将<span class="language-math"> $N_r$ </span>直送<span class="language-math"> $n_r$ </span>,</li>

<li>即完成<span class="language-math"> $N_v N_r$ </span>  $\rightarrow$  <span class="language-math"> $n_v n_r$ </span>。</li>

</ul>

<p>④ 装入位=0 时，失效，表示该虚页未装入<span class="language-math"> $M_1$ </span>中。</p>

<p></p>

<p>为了加深理解我们做下边一个例题，</p>

<blockquote>

<p>某虚拟存储器共 8 个页面，每页为 1024 个字，实际主存为 4096 个字，采用页表法进行地址映。映像表的内容如下表所示：</p>

<p></p>

<p>1) 列出会发生页面失效的全部虚页号；<br>

2) 按以下虚地址计算主存实地址：<br>

0、3728、1023、1024、2055、7800、4096、6800</p>

</blockquote>

<p></p>

<h2 id="地址映像及其变换">地址映像及其变换</h2>

<p>在上一小节，我们简单介绍了虚地址和实地址之间的转换。但实际上如果存储系统使用的是不同址映像方式的话，可能地址变换方式可能也会有很大的不同。一般来说我们常用的地址映像方式有以四种：</p>

<ol>

<li>全相联</li>

<li>直接相联</li>

<li>组相联</li>

<li>段相联</li>

</ol>

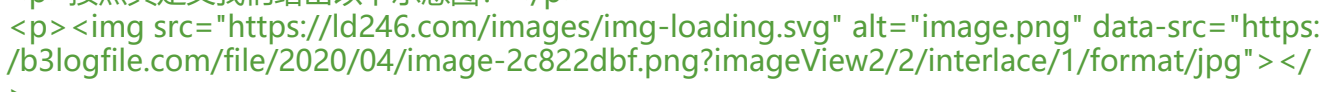
下边我们分四个小专题分别对其进行介绍。

## 全相联映象及其变换

首先我们给出`全相联映象的含义：

- 对辅存中的任何一个页面都可以放到主存中任何一个页面的映像规则，称为全相联映射。

按照其定义我们给出以下示意图：

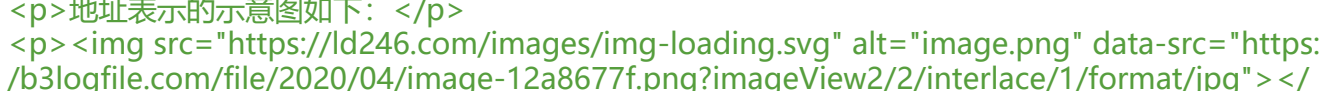


我们不难 `全相联映射` 具有如下特点：

辅存中的任何一个页面都可能映射到主存中。

### 地址变换

地址表示的示意图如下：



当然上边具体变换方法不只一种，除了上一章节介绍的 `页表法` 外还有一种方叫做 `目录表法`

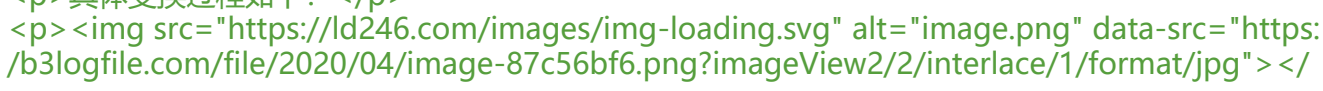
### 全相联目录表法

该方法要求用相联存储器作为目录表（相联存储器是一种可按照内容的特征字段来访问的一种存储器，是一个小容量高速存储器）。

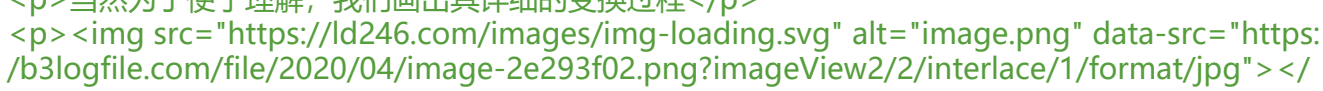
目录表的**结构**具有如下特点：

- 目录表的行数与主存页面数相等
- 目录表中每行的内容：
  - $N_v$ 为相联比较字段；
  - $n_v$ 为主存页号（非相联比较字段）。

具体变换过程如下：



当然为了便于理解，我们画出其详细的变换过程



具体过程如下：

- 将虚地址中的  $N_v$  送目录表中去进行相联比较（一个  $t_m$ ）。
- 当有某个比较器比较相等(比较的时候是比较的是  $N_v$  相关联的字段)时，将该行  $n_v$  送出，同时  $N_r \rightarrow n_r$ ，实现了  $N_v N_r \rightarrow n_v n_r$  的变换命中)。
- 若没有相等的，不命中，等待调入。

通过以上过程我们不难分析出全相连目录表法具有如下**特点**：

- 产生页面冲突的可能性很小（因为页表的行数和主存行数相同）；

<li>与页表放入主存中相比，查表的速度快（因为目录表是放到相联存储器中的，输入高速存储器）

</li>

</ol>

<p>当然该方法也有以下<strong>缺点：</strong></p>

<ol>

<li>可扩展性差（因为相联存储器容量在一定程度上限制了主存的扩展，不能太大超过相联存储器的数）</li>

<li>主存容量增加时，目录表的造价高，速度降低（成本原因，在工程上很常见）。</li>

</ol>

<h2 id="直接映象及其变换">直接映象及其变换</h2>

<p>该方法现将辅存<em><strong>按照主存大小分为若干块</strong></em>，在辅存的每一块都有与主存相同的页面数，辅存每块内的页面<em><strong>只能调入</strong></em>到与主存相同的页面上。</p>

<p>该方法与全相联映射相比我们发现，直接映象辅存中的某块具体某一页能够调度的位置是固定的或者说是固定的几个位置）。</p>

<p>下面我们画出直接映象的规则示意图：</p>

<p></p>

<p>其中：</p>

<ul>

<li><span class="language-math"> $N_d$ </span>：块号</li>

<li><span class="language-math"> $N_v$ </span>：块内页号</li>

</ul>

<p>从图中我们可以看到，该方法先将辅存按主存大小分为若干块，在辅存的每块内都有与主存相同页面数，辅存每块内的页面只能调入到与主存相同的页面上（可以简单理解为一个萝卜一个坑 :smile::mile:)</p>

<h3 id="地址变换-">地址变换</h3>

<p>在展开说明直接映象的地址变换之前我们先了解下直接相联请情况下，<strong>实地址和虚地是如何表示的。</strong></p>

<p>由于直接映象引入了\*\*“块”<strong>，因此我们在给定虚地址时肯定要给出</strong>块号\*，然后通过块号（<span class="language-math"> $N_d$ </span>）、块内页号(<span class="language-math"> $N_v$ </span>)以及页内地址(<span class="language-math"> $N_r$ </span>)三部来最终定位到具体的数据地址，此处我们给出示意图便于理解。</p>

<p></p>

<p>实地址和硬件是高度相关的，因此基本不会发生变化，和前边一样由两部分组成实页号（<span class="language-math"> $n_v$ </span>）页内地址（<span class="language-math"> $n_r$ </span>），示意图如下：</p>

<p></p>

<p>同样的，由于引入了“块”的概念因此在进行地址变换的时候我们需要一个新的表--<code>块</code>。</p>

<p><code>块表</code> 具有如下<strong>特点</strong>：</p>

<ol>

<li>块表长度与主存页面数相等。</li>

<li>块表行中的内容：块号<span class="language-math"> $N_d$ </span>。</li>

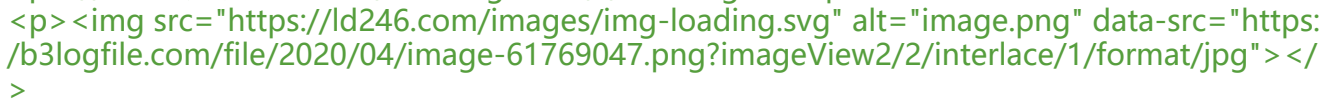
</ol>

<p></p>

<p>有了前边两点的铺垫之后，我们下边可以详细讲述一下地址变换的过程：</p>



给出地址变换的一个简单示意图：



<ol>

<li>根据块内页号  $N_v$  去查块表中的  $N_v$  行，获取到对应的块号

<li>将虚地址中的块号  $N_d$  与所选块表中的  $N_d$  比较。

<li>相同时则命中，直接将  $N_v \rightarrow n_v$ ， $N_r \rightarrow n_r$ 。

<li>不同时，则表示未命中

</ol>

从上边的分析中我们可以总结出直接相连映象具有如下特点：

<ol>

<li>可将查表与访存同时进行，有利于访问速度的提高（命中时）。

<li>产生页面冲突的可能性极大（因无灵活的存放余地）。

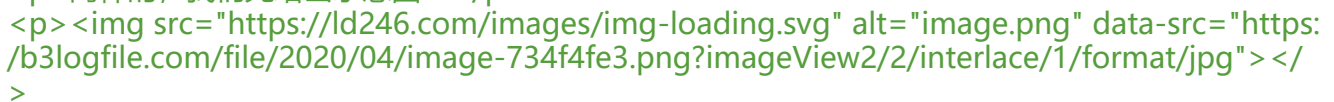
</ol>

组相联映象及地址变换

组相连映像 则是先**将主存**分为页面数相的若干组，再将辅存按主存划分为若干区，**组内**采用全相联映象，组间采用直接映。

这个很明显通过分组融合了全相连映象和直接映象，说起来还有点小期待:happy::happy:。

同样的，我们先给出示意图：



其中：

<ul>

<li>  $N_d$ :区号

<li>  $s$ :组号

<li>  $q$ :组内页号——辅存

<li>  $s'$ :组号

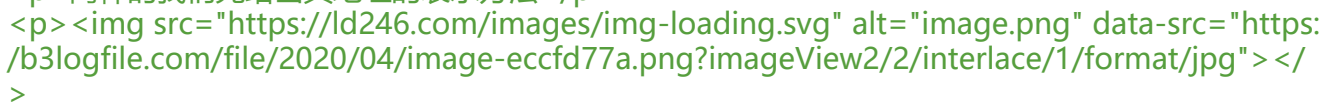
<li>  $q'$ :组内页号——主存

</ul>

从图中我们可以明显看到：组间是直接相联的，而组内则是全相联。

地址变换

同样的我们先给出其地址的表示方法



由于映象方式的改变，我们在进行地址变换的时候，需要引入一个**新的表**--  
随机存储器表。

该表具有如下特点：

<ol>

<li>表的行数与组数**相等**（本例 2 组，即 2 行）。

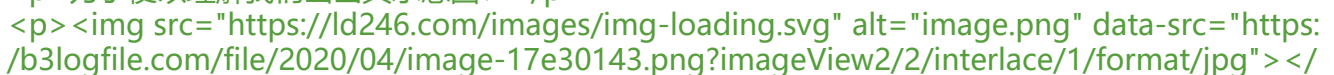
<li>每行大字段数与组内页面数**相等**（本例 2 个）。

<li>每个大字段又分为三个小字段。

<li>每个大字段还有一个比较器。

</ol>

为了便以理解我们画出其示意图：



>

<p>具体的变换过程如下： </p>

<p></p>

>

<ol>

<li>根据虚地址中的  $s$  去查 RAM 表中的某一行。 </li>

<li>将虚地址中的  $N_d$ 、 $q$  同时送各比较器与所选行中的  $N_d$ 、 $q$  进行比较。 </li>

<li>当有一个比较器相等时：

<ol>

<li>将  $s \rightarrow s'$ （组间直接）。 </li>

<li>将相等大字段中  $q'$  送出作  $q'$ 。 </li>

<li>再将  $N_r \rightarrow n_r$ ，即实现了将虚址( $N_d, s, q, N_r$ )  $\rightarrow (s', q', n_r)$ 命中时的地址变换 </li>

>

</ol>

</li>

</ol>

<p>最后我们总结出组相联映象的特点： </p>

<p>既有直接映象中对号入座部分(组间直接)，可减少查表范围，缩短查表时间，又有全相联中的灵活存放规则(组内全相联)，从而可降低页面冲突。 </p>

<p><strong>缺点： </strong> 控制机构复杂 </p>

## 段相联映象

<p>所谓 <code>段相连映象</code>，其实跟组相联特别相似--对主辅的划分与组相联映象相同，为区分两种不同的映象规则，将组相联中的组改为段，段间采用全相联，段内采用直接映象。 </p>

<p>其示意图如下： </p>

<p></p>

>

<p></p>

>

## 四种映象规则的关系

<ol>

<li>

<p>在组相联映象中，当每组只有一页时，此时的组相联就是直接映象。 </p>

<p>当把主存只分为一个组时，此时的组相联也就是全相联映象，即直接映象和全相联映象是组相联映象的两个特例。 </p>

</li>

<li>

<p>在段相联映象中，当每段只有一页时，此时的段相联映象就是全相联映象。 </p>

<p>当把主存只分为一个段时，此时的段相联也就是直接映象。 </p>

</li>

</ol>

## 总结

<p>本文主要从存储体系由来、分类以及原理角度来讲解一个存储体系的设计，努力做到全面。当然于个人水平有限，文章难免可能会有错误，如若发现，恳请指出，不胜感激。 </p>