

Recap: Beyond Physical Memory

Swap space (part of disk) can be utilized to extend physical memory.

Implementation:

- Swap space is divided into blocks of the same size as pages; pages can be stored in swap space.
- Page table is extended to indicate whether a page is in physical memory or in swap space.
- To access a page not in physical memory causes page-fault, which is handled by a handler in OS kernel; the detailed procedure has been discussed.

Examples in xv6-riscv: get arguments of system call

syscall.c:

```
// Retrieve an argument as a pointer.  
// Doesn't check for legality, since  
// copyin/copyout will do that.  
void argaddr(int n, uint64 *ip)  
{  
    *ip = argraw(n);  
}
```

//note: the obtain address is VA (virtual address)

```
static uint64 argraw(int n)  
//get the n-th arg of current process'  
//syscall  
{  
    struct proc *p = myproc();  
    switch (n) {  
        case 0: return p->trapframe->a0;  
        case 1: return p->trapframe->a1;  
        case 2: return p->trapframe->a2;  
        case 3: return p->trapframe->a3;  
        case 4: return p->trapframe->a4;  
        case 5: return p->trapframe->a5;  
    }  
    return -1;  
}
```

Examples in xv6-riscv

Riscv.h:

```
#define PGSIZE 4096 // bytes per page
```

```
#define PGSHIFT 12 // bits of offset within a page
```

```
//get the VA of the starting of the page containing VA a
```

```
#define PGROUNDDOWN(a) (((a)) & ~(PGSIZE-1))
```

```
//get the VA of the starting of the next page of the page containing VA a
```

```
#define PGROUNDUP(sz) (((sz)+PGSIZE-1) & ~(PGSIZE-1))
```

Examples in xv6-riscv: copy from kernel to user space

vm.c:

```
Int copyout(pagetable_t pagetable, uint64
dstva, char *src, uint64 len)
{
    uint64 n, va0, pa0;
    while(len > 0){
        va0 = PGROUNDDOWN(dstva);
        //find PA of VA va0
        pa0 = walkaddr(pagetable, va0);
        if(pa0 == 0)
            return -1;
        n = PGSIZE - (dstva - va0);
```

```
        if(n > len)
            n = len;
        memmove((void*)(pa0 + (dstva - va0)),
src, n);
        len -= n;
        src += n;
        dstva = va0 + PGSIZE;
    }
    return 0;
}
```

Examples in xv6-riscv: address translation

vm.c:

```
pte_t *
walk(pagetable_t pagetable, uint64 va, ...)
{
    if(va >= MAXVA)
        panic("walk");
    for(int level = 2; level > 0; level--) {
        pte_t *pte = &pagetable[PX(level, va)];
        if(*pte & PTE_V) {
            pagetable = (pagetable_t)PTE2PA(*pte);
        } else {
            ...
        }
    }
    return &pagetable[PX(0, va)];
}
```

Note:

```
// Return the address of the PTE in page table pagetable
// that corresponds to virtual address va. If alloc!=0,
// create any required page-table pages.
//
// The risc-v Sv39 scheme has three levels of page-table
// pages. A page-table page contains 512 64-bit PTEs.
// A 64-bit virtual address is split into five fields:
// 39..63 -- must be zero.
// 30..38 -- 9 bits of level-2 index.
// 21..29 -- 9 bits of level-1 index.
// 12..20 -- 9 bits of level-0 index.
// 0..11 -- 12 bits of byte offset within the page.
```

Examples in xv6-riscv: address translation

vm.c:

```
Uint64 walkaddr(pagetable_t pagetable, uint64 va)
{
    pte_t *pte;
    uint64 pa;
    if(va >= MAXVA)
        return 0;
    pte = walk(pagetable, va, 0);
    if(pte == 0) return 0;
    if((*pte & PTE_V) == 0) return 0; //valid?
    if((*pte & PTE_U) == 0) return 0; //user allowed to access?
    pa = PTE2PA(*pte);
    return pa;
}
```

Note:

// Look up a virtual address,
return the physical address,

// or 0 if not mapped.

// Can only be used to look up
user pages.

// va must be the VA of the
starting of a page

L12: Beyond Physical Memory: Policies

(based on Ch. 22)

Replacement Policy

We have seen the issue of **replacement policy** in two contexts: TBL cache and swap

Indeed, it is a common and challenging problem

- Miss rates must be low, every miss is costly
- Data structure must be fast, potentially updated every memory access
- Would like to take advantage of hardware support for anything that happens on every memory access

Our first attempt will be FIFO, which despite its simplicity is not a good choice, it experiences a paradoxical corner-case behavior known as **Bélády's anomaly**

How to decide which page to evict?

Framing the Problem

Main memory can only hold a subset of all pages in the system, the rest are pushed to the disk (swap space)

In a sense, main memory is a "cache" of the systems "popular pages"

We will discuss replacement policies from the perspective of a cache, but keep in mind the following interpretation for swap

"cache" hit = page is in main memory

"cache" miss = **page fault** (page is in swap space)

Reference Strings (as Input)

We will examine the performance of policies for specific memory access patterns called **page reference strings**

Suppose the memory accesses on a system look like the following (page size = 128B)

Memory Address	112	58	114	200	220	200	290	58	224
Page Number	0	0	0	1	1	1	2	0	1

Every time there is a page fault the page must be brought to main memory

Consecutive repeated page accesses can never cause additional misses, for this reason they are not interested in them when evaluating policy performance and we remove them from the reference string

Example Reference String (with consecutive repeated pages removed): 0, 1, 2, 0, 1

Comparison to Optimal

To evaluate policies, we need a point of comparison, what are the fewest possible misses for a given cache size

Optimal (OPT) – assume we can predict the future and always chose to evict the page that will be used furthest out (or any page that will never be used again)

Page	0	1	2	0	1	3	0	3	1	2	1
Hit/Miss	M	M	M	H	H	M	H	H	H	M	H
Evict						2				3	
Cache[0]	-	-	0	0	0	0	0	0	0	0	0
Cache[1]	-	0	1	1	1	1	1	1	1	1	1
Cache[2]	0	1	2	2	2	3	3	3	3	2	2

For reference string 0,1,2,0,1,3,0,3,1,2,1 and a cache size of 3, no policy will every have fewer than 5 misses

FIFO

FIFO – evict the page that has been in memory the longest

Advantage: easy to implement and fast ($O(1)$ operations), can be implemented in a fixed size cache using a circular buffer

Page	0	1	2	0	1	3	0	3	1	2	1
Hit/Miss	M	M	M	H	H	M	M	H	M	M	H
Evict						0	1		2	3	
Queue Head	-	-	0	0	0	1	2	2	3	0	0
	-	0	1	1	1	2	3	3	0	1	1
Queue Tail	0	1	2	2	2	3	0	0	1	2	2

Belady's Anomaly

Intuitively we would expect larger cache to perform better (or at least as well) as smaller cache

FIFO has a strange corner case known as **Belady's Anomaly**, where sometimes **a larger cache does worse**

Cache size 3 (9 misses)

Page	1	2	3	4	1	2	5	1	2	3	4	5
Hit/Miss	M	M	M	M	M	M	M	H	H	M	M	H
Evict				1	2	3	4			1	2	
Queue Head	-	-	1	2	3	4	1	1	1	2	5	5
	-	1	2	3	4	1	2	2	2	5	3	3
Queue Tail	1	2	3	4	1	2	5	5	5	3	4	4

Cache size 4 (10 misses)

Page	1	2	3	4	1	2	5	1	2	3	4	5
Hit/Miss	M	M	M	M	H	H	M	M	M	M	M	M
Evict							1	2	3	4	5	1
Queue Head	-	-	-	1	1	1	2	3	4	5	1	2
	-	-	1	2	2	2	3	4	5	1	2	3
	-	1	2	3	3	3	4	5	1	2	3	4
Queue Tail	1	2	3	4	4	4	5	1	2	3	4	5

Random

Random – evict a page at random

Advantages

- easy to implement and fast
- no corner cases such as Belady's anomaly possible because algorithm is not influenced by order of pages (its random)

Page	0	1	2	0	1	3	0	3	1	2	1
Hit/Miss	M	M	M	H	H	M	M	H	M	H	H
Evict						0	1		3		
Cache[0]	-	-	0	0	0	1	2	2	2	2	2
Cache[1]	-	1	1	1	1	2	3	3	0	0	0
Cache[2]	0	0	2	2	2	3	0	0	1	1	1

LRU

LRU – evict the page used the least recently

Advantages

- Better performance than FIFO
- No Belady's anomaly possible

Disadvantage – slow, accounting to do on every page access

Page	0	1	2	0	1	3	0	3	1	2	1
Hit/Miss	M	M	M	H	H	M	H	H	H	M	H
Evict						2				0	
Least Recent	-	-	0	1	2	0	1	1	0	2	3
	-	0	1	2	0	1	3	0	3	3	2
Most Recent	0	1	2	0	1	3	0	3	1	1	1

Clock Algorithm

LRU performs the closest in optimal for typical workloads, but it is really costly on every memory access

Clock Algorithm - is a way to approximate LRU cheaply

Add extra **use bit** to page table entry, on every memory access MMU set use bit of page to 1

When looking for page to evict

- visit pages in round-robin order
- if page use bit is 0, chose that page to evict
- else set use bit to 0 and continue search

Clock Algorithm: Example

Frames after accessing: 0 1 2 0 1

Frame #	0	1	2
Page	0	1	2
Use bit	1	1	1

No eviction; all pages are marked as used recently

To access page: 3

Frame #	0	1	2
Page	0 -> 3	1	2
Use bit	1 -> 0 -> 1	1 -> 0	1 -> 0

Page 0 (not exactly the LRU) is evicted; then page 3 (new) is marked as used recently

To access pages: 1 3

Frame #	0	1	2
Page	3	1	2
Use bit	1	0 -> 1	0

No eviction; pages 1 and 3 are marked as used recently

To access pages: 0

Frame #	0	1	2
Page	3	1	2 -> 0
Use bit	1	1 -> 0	0 -> 1

Page 2 (the LRU) is evicted; page 0 (new) is marked as used recently

Dirty Bit Optimization

If a page has only been read from (never written to) there is no reason to write it back to the swap space (if the swap page is large enough to contain a copy of it already) when it is evicted

Add a **dirty bit** to the page table, initialize to 0, on every write to the page table the MMU (hardware) sets the bit to 1

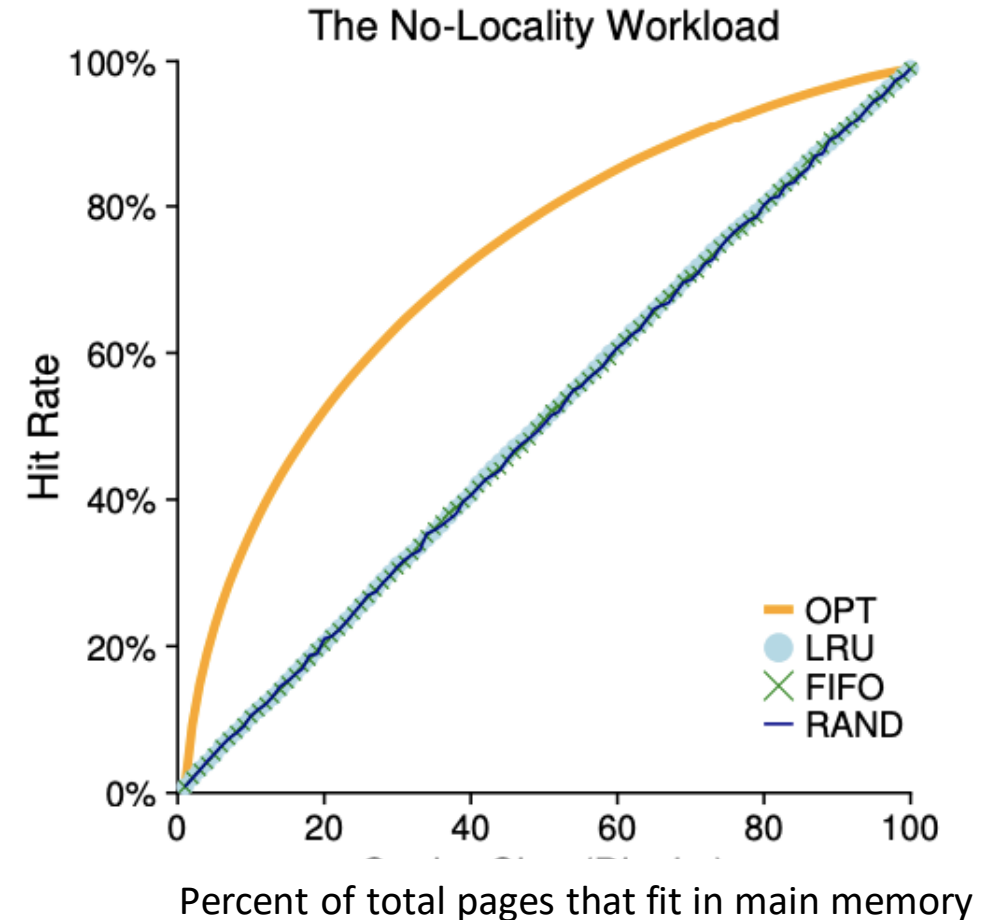
Can also make clock algorithm more efficient

- First try to find a page with dirty bit and use bit both 0, because it is lower cost to replace

What if All Accesses Were Random?

Without locality, none of the policies (except the cheating OPT) provide cost-effective benefit

If main memory is 50% of total pages, only 50% hit rate



80-20 Workload

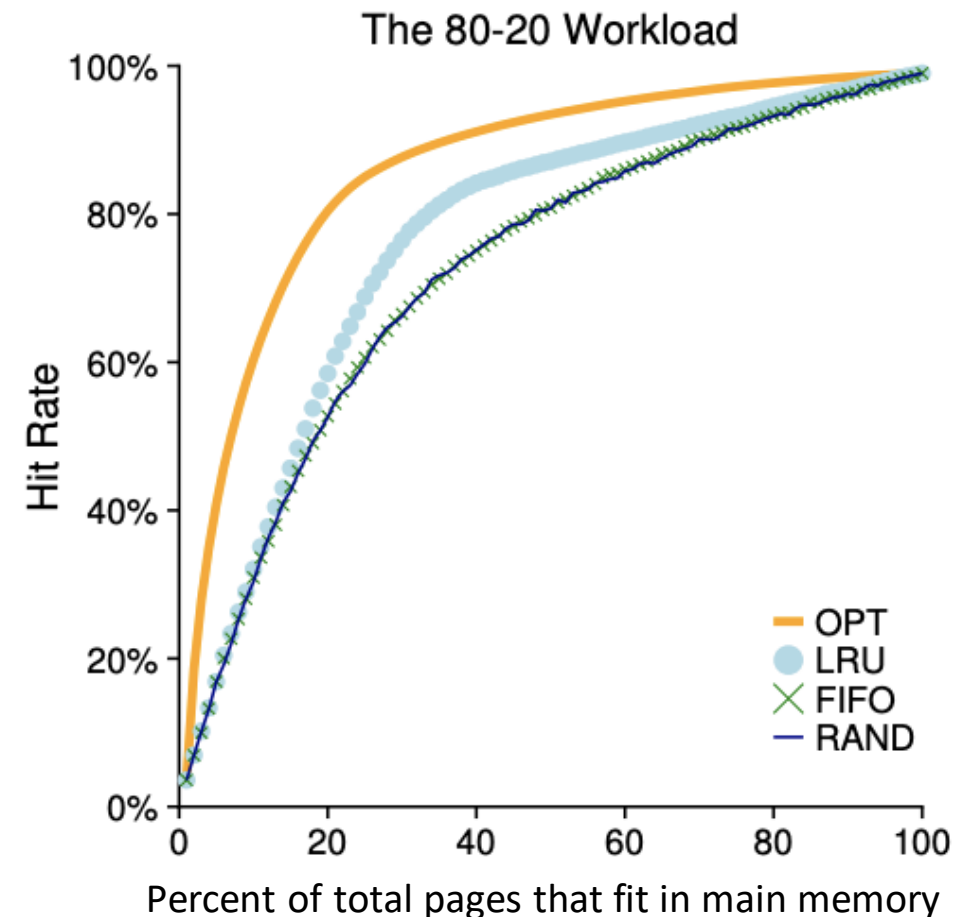
More typical, the **80-20 workload** (80% of accesses are to 20% of pages)

A few pages get most accesses

Most pages get few accesses

Results from locality (either spatial or temporal)

LRU performs closest to optimal, FIFO and RAND perform the same

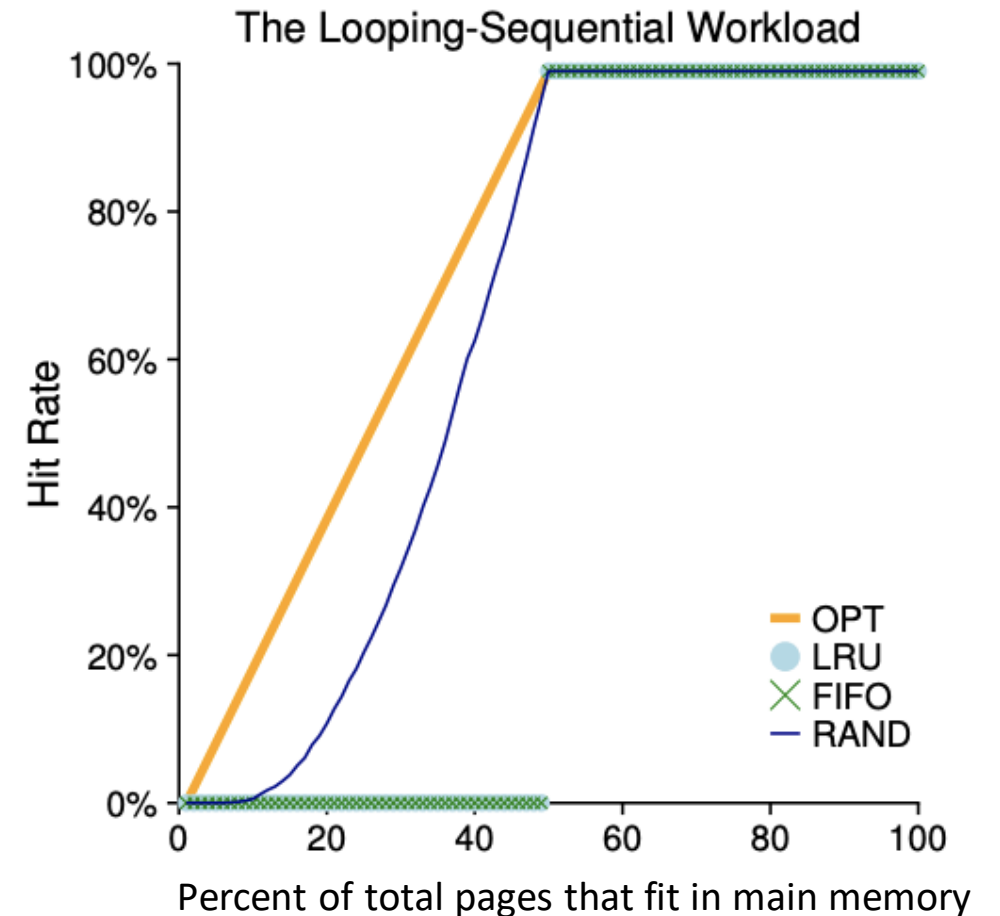


Looping-Sequential Workload

Assume loop repeatedly reads pages 1 to 49 in increasing order

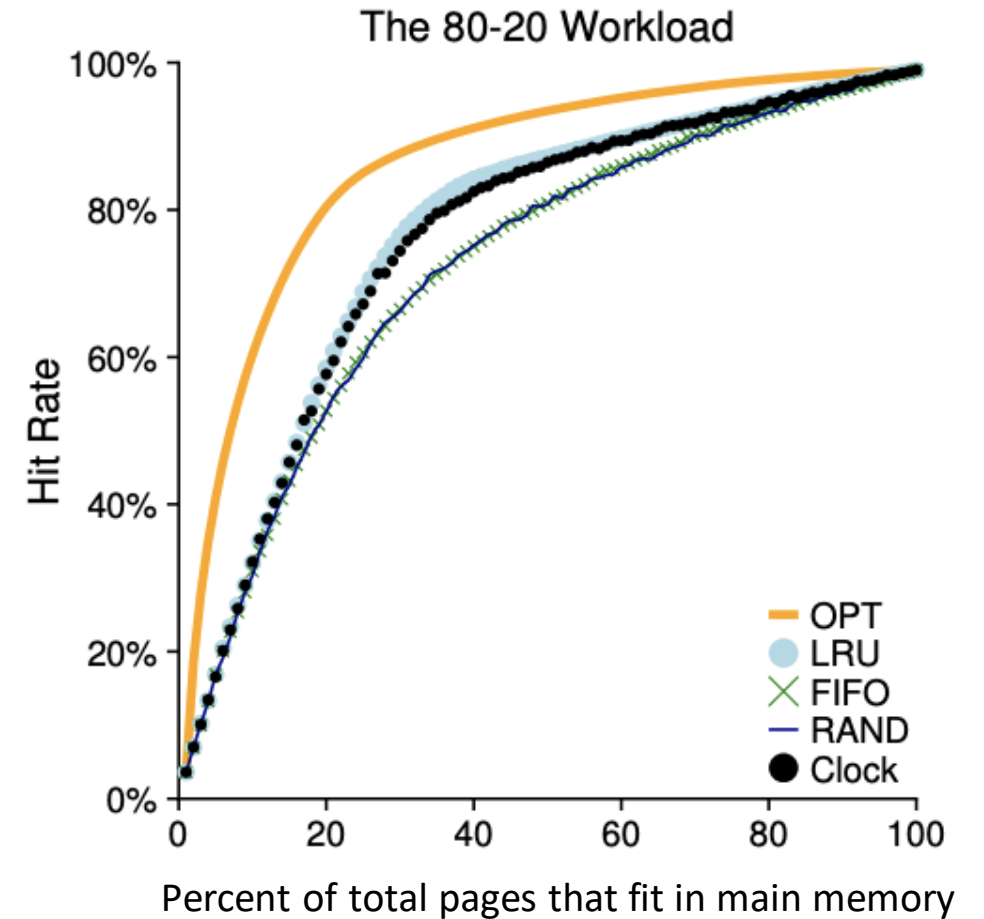
When cache size is below 50, LRU and FIFO have same corner case that causes every access to be a miss

Random avoids corner cases



Clock Policy with 80-20 Workload

Clock is a very close approximation of LRU



Thrashing

Example: System has two processes that both sequentially read from N pages in a continuous loop. The System only has enough main memory to store N pages.

Thrashing is when process is spending more time on handling page-fault than executing

- Starts at one processes and snowballs into several processes thrashing
- Multiprogramming and multitasking make it worse, not better
- Sudden and extreme drop in system performance

Need to reduce the amount of multiprogramming and multitasking to avoid thrashing