

25. Log-structured File System

▼ 등장 배경

1. system memory → 계속 커짐
 - a. (page) cache memory도 커짐 ⇒ 더 많은 data가 cached → 매번 disk 접근 X
 - b. disk traffic에 의해 write 성능은 read만큼 좋아지지 않음
↳ load가 필요하(에) disk 접근이 많
2. random I/O와 sequential I/O 사이에 너무 큰 gap
 - seek과 rotational delay costs는 여전히 너무 성능 bad임
 - ⇒ random I/O 성능이 bad
3. 현존하는 file system들이 대부분의 workload에서 성능 거지 똥임
 - ex) file system : one block의 새로운 file을 생성할 때 많은 수의 write를 수행
↳ 22만 블록을 자주 생성하면 성능 좋지 X

▼ Log-structured File system(LFS)

disk에 writing 할 때 in-memory segment에 쓰고자 하는 data 정보 모두 모아 놓음

→ segment가 꽉 차면 disk에서 sequential하게 붙어 있는 sector에 한 번에 작성

(기존에 있는 data overwrite X, 새로운 공간에 계속 작성)

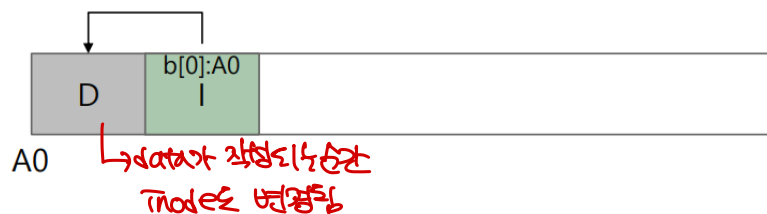
↳ freelist

- 최근 연구 : large I/O → 높은 성능의 flash-based SSD들을 필요로 한다는 것을 보임

↳ 다음 chapter

▼ Writing to disk

- file-system에 어떻게 sequential하게 write를 한 번에 모아서 함?
 - data만 쓰게 되면 inode는 random하게 여기 저기 작성될 수 있음 (정작 random I/O)
 - ⇒ metadata(inode)도 함께 작성! ☆ → inode 위치가 계속해서 변경됨 (고정된 위치에 저장 X)



- 단순히 sequential writing만 하는 것은 성능 개선에 큰 도움이 되지 X

Time T : A라는 주소에 하나의 block을 write 하는데 걸리는 시간

Time T+f : A+1 "

Disk $\Rightarrow (T_{rotation} - f)$ 만큼 기다려야 sector에 data 작성 가능

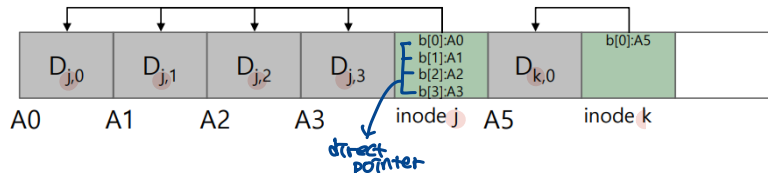
\Rightarrow but, 큰 사이즈의 write를 한다면 성능적인 이득이 있을 것임

- write buffering

- LFS는 되게 큰 memory 영역에다가 request를 차곡차곡 모음

- 충분할만큼 모였을 때 한 번에 disk에 작성

\Rightarrow segment : request가 모인 memory 영역 (크기는 few MB)
 \Rightarrow disk에 한 번에 request 쓰기



\rightarrow segment size는 환경에 따라 다른 값으로 설정됨

▼ Finding inodes \rightarrow Inode Map

LFS : 아예 새로운 block에 작성 \rightarrow sequential한 I/O 작성하기 위함

\rightarrow 특정 inode number를 갖는 inode 위치가 계속 바뀜 (disk 전체에 흩어져 있음)

\Rightarrow update 꾸준히 해줘야 함

			iblock 0				iblock 1				iblock 2				iblock 3				iblock 4			
Super	i-bmap	d-bmap	0	1	2	3	16	17	18	19	32	33	34	35	48	49	50	51	64	65	66	67
			4	5	6	7	20	21	22	23	36	37	38	39	52	53	54	55	68	69	70	71
			8	9	10	11	24	25	26	27	40	41	42	43	56	57	58	59	72	73	74	75
			12	13	14	15	28	29	30	31	44	45	46	47	60	61	62	63	76	77	78	79
0KB	4KB	8KB	12KB	16KB	20KB	24KB	28KB	32KB														

* Inode 어디 찾을?

- **imap(Inode Map)**

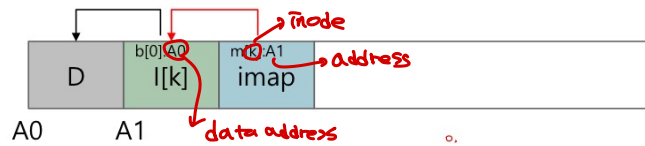
- inode number를 받아서 가장 최신 version의 inode로 disk address에 기록해놓는 table

새로운 위치로!

- inode가 disk에 쓰일 때마다 imap update되어야 함

\rightarrow 고정되어 있지 않음!

- imap은 disk 어디에 존재?



- disk 내에서 고정해놓는다면? → performance 나쁨(여전히 random I/O)

⇒ update된 내용만 imap에 저장해놓기

↳ 모든 새로운 내용이 기록된 곳 바로 옆에 저장
* Imap 위치는 어케 찾을지?

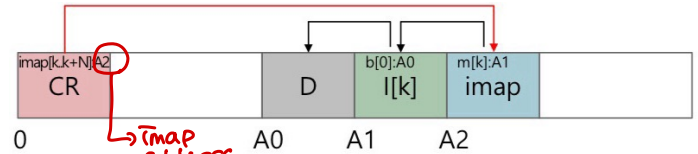
- checkpoint Region(CR)

- imap도 고정해놓는다면? → 여전히 random access

⇒ checkpoint region 사용!

- CR ⇒ 고정된 위치에서 사용

- 부분적으로 저장해놓음
- CR을 매번 update하면 성능이 나빠질 수 있음(random access 가능)



⇒ periodically하게 update(ex. 30초에 한 번) ⇒ 성능에 크게 영향을 주지 않도록!!

but, 이 사이에 crash 발생하면 update 중단되고 정보 손실 가능

▼ Reading a file from disk (w. imap)

1. checkpoint region 읽기
2. CR 이용해서 전체 inode map 읽은 뒤 memory에 넣어둠

⇒ 매번 disk 접근하지 않기에 추가적인 overhead는 미미함

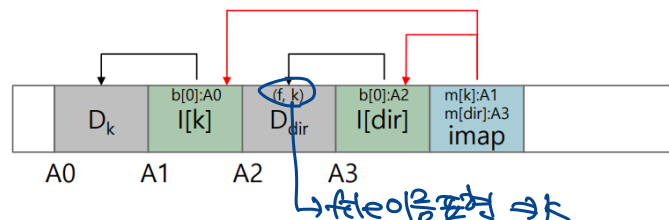
(전체 imap이 cached 되므로) 추가적인 작업은 imap에서 inode address 찾는 일뿐.

▼ Directory 접근

기존 system은 작은 file을 계속해서 생성할 때 좋은 성능을 보여주지 않음

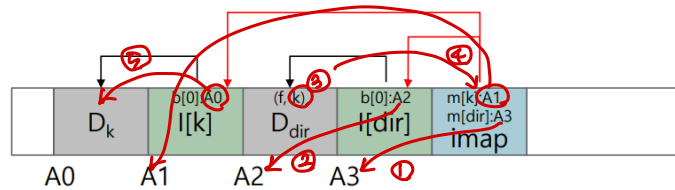
→ LFS ver 확인해보자!

1. Creating a file dir/f ⇒ directory 고려해서 file 생성된다고 생각해보자



- LFS : directory에 해당하는 inode update
 - directory data 영역에 file 영역 이름 표현해야 함

2. accessing file f



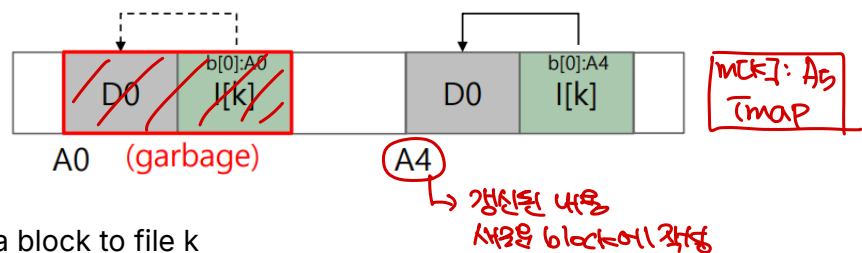
1. memory에 저장되었던 inode map을 통해 dir의 inode 위치 찾을 $\rightarrow A_3$
2. dir inode 읽은 뒤 dir data의 위치 찾을 $\rightarrow b[0]:A_2$
3. data block 읽은 뒤 (f,k)에 매핑되어 있는 inode-number의 name 찾을 $\rightarrow (f, k) \rightarrow \text{file name}$
4. inode number k에 위치한 inode map 찾을 $\rightarrow m[k]:A_1$
5. 그 address에 있는 data block read $\rightarrow b[0]:A_0$

▼ Garbage Collection \rightarrow LFS : 매번 data를 새로운 위치에 씀 \Rightarrow 기존에 존재하던 data?

회수가 되지 않은 영역을 모아서 중간에 update 해줌

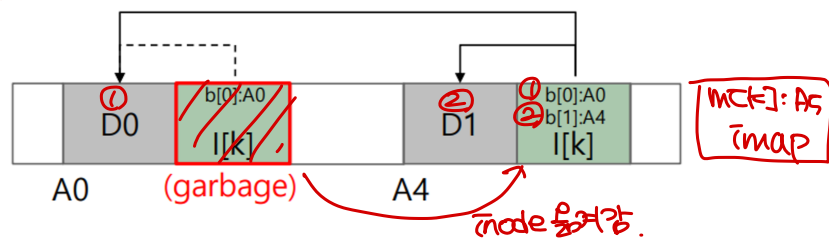
\rightarrow data, inode, 다른 자료구조들의 오래된 version 찾을 뒤 초기화

- example
 - update data block



- appending a block to file k

\rightarrow data 추가



\Rightarrow 많은 양의 공간을 clean up하게 되면 연속적인 writing 가능

- ① M개의 old segment clear
- ② 어떤 block이 사용되지... reverse한 block 옮겨감
- ③ N개의 새로운 segment compact \rightarrow 새로운 위치 disk에 작성 (N < M)
- ④ M개의 오래된 segment free \rightarrow 연속적인 write를 위한 공간으로 사용

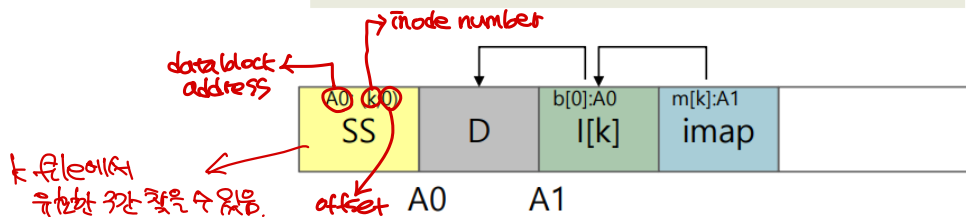
- 살아있는 block인지, 죽어있는 block인지 어떻게 check?

\Rightarrow segment summary block

```

(N, T) = SegmentSummary[A];
inode = Read(imap[N]);
if (inode[T] == A)
    // block D is alive
else
    // block D is garbage

```



- 어떤 block을 초기화 + 언제 초기화?
 - 언제 초기화?
 - idle time ⇒ cpu가 여유있을 때
 - disk가 꽉 찼을 때
 - 어떤 block을 clean?
 - hot segments : 자주 overwritten
 - cold segments : 몇 가지 dead block을 가지고 있지만 아직 살아있는 놈 몇 개 있는 segment
- ⇒ cold 초기화 한 뒤 hot 초기화

▼ Crash Recovery

LFS가 disk에 쓰고 있는 동안 crash 일어난다면? → journaling 호출

1. CR에 쓰고 있을 때 발생

⇒ LFS는 두 가지 CR을 가지고 있으며 번갈아 사용함 (disk의 양끝에 CR을 두고 사용)

- header에 write (with timestamp)
- CR의 body에 작성
- last block에 write(with timestamp)

→ 여가 더 느리면 crash가 일어났다고 판단

→ 항상 일관적인 timestamp를 가진 가장 최근의 CR을 골라서 사용함

2. segment에 쓰고 있을 때 발생

⇒ LFS는 30초마다 update 하기에 file system이 이 중간에 침입 가능

→ roll forward

⇒ CR에 기록된 (로그의 끝을 찾아서 다음 segment를 읽고 만약 valid한 update가 보인다면 사용함.

⇒ LFS는 file system을 update하므로 last checkpoint 이전에 쓰여진 모든 data + metadata를 recovery 가능