

# A Lightweight Log-based Deferred Update for Linux Kernel Scalability

## Abstract

We propose a novel light weight concurrent update method, LDU, to improve performance scalability for Linux kernel on many-core systems through eliminating lock contentions for update-heavy global data structures during process spawning and optimizing update logs. The proposed LDU is implemented into Linux kernel 4.5 and evaluated using representative benchmark programs. Our evaluation reveals that the Linux kernel with LDU shows performance improvement by ranging from Xx through Xx on a 120 core system.

## 1 Introduction

최근 코어수가 증가하고 있다. 따라서 멀티코어에서 매니코어 시스템으로 바뀌고 있다. 매니코어 시스템에 대한 운영체제 커널의 parallelism은 시스템 전체의 parallelism에서 가장 중요하다. 만약 커널이 scale하지 않으면, 그 위에 동작하는 응용프로그램들도 역시 scale하지 않는다[1]. 이처럼 중요한 운영체제 커널 중 멀티코어 또는 매니코어 환경에서 많이 사용되는 운영체제가 리눅스 커널이다[2]. 하지만 리눅스 커널은 아직 확장성 문제가 있다 [9] [32]. 확장성 문제 중 하나는 락 경쟁 때문에 발생하는 업데이트 직렬화 문제이다 [26] [21]. 그 이유는 업데이트 오퍼레이션은 여러 쓰레드가 동시에 수행되지 못하기 때문이다 [28].

이처럼 업데이트 직렬화 문제를 해결하기 위해 여러 동시적 업데이트 방법들이 연구되고 있다 [6] [26]. 이러한 동시적 업데이트 방법들은 워크로드 특성인 업데이트 비율에 따라 많은 성능 차이를 보인다 [26]. 이 중 높은 업데이트 비율을 가진 자료 구조 때문에 발생하는 확장성 문제를 해결하기 위한 여러 방법이 연구되고 있다. 그 중 하나는 cache communication bottleneck을 줄인 log-based 알고리즘 [35] [25] [7]을 사용하는 것이다. Log-based 알고리즘은 업데이트가 발생하면, data structure의 업데이트 operation을 per-core 또는 atomic하게 log로 저장하

고 read operation을 수행하기 전에 저장된 로그를 수행하는 것이다. 이것은 마치 CoW(Copy On Write)와 유사하다 [27].

S. Boyd-Wickizer et al.는 동기화된 타임스탬프 카운터(synchronized timestamp counters) 기반의 per-core log를 활용하여 update-heavy한 자료구조를 대상으로 동시적 업데이트 문제를 해결함과 동시에 cache communication bottleneck을 줄였다 [7]. 동기화된 타임스탬프 카운터 기반의 per-core log를 활용한 동시적 업데이트 방법은 업데이트 부분만 고려했을 때, per-core에 데이터를 저장함으로써 굉장히 높은 scalability를 가진다[3]. 하지만 per-core 기반의 동기화된 타임스탬프 카운터를 사용한 방법은 결국 timestamp merging and ordering 작업을 야기한다. 만약 코어 수가 늘어 날 경우, 로그를 자료 구조에 적용하는 과정에서 timestamp 때문에 발생하는 추가적인 sequential 프로세싱이 요구된다. 이것은 결국 확장성과 성능을 저해한다.

본 논문은 동기화된 타임스탬프 카운터를 이용하여 따라 생기는 sequential processing 문제를 해결하기 위해 LDU(Lightweight log-based Deferred Update)를 개발하였다. LDU는 타임스탬프 카운터가 필요한 operation log를 업데이트 순간마다 지우고, 매번 로그를 생성하지 않고 재활용하는 방법이다. 이로 인해 synchronized timestamp counter 문제와 cache communication bottleneck 문제를 동시에 해결하였다. 이러한 해결 방법은 분산 시스템에서 사용하는 log기반의 concurrent updates 방식과 shared-memory system의 hardware-based synchronization 기법(compare and swap, test and set, atomic swap)을 조합하여 동시적 업데이트 문제를 해결하였다.

이처럼 동기화된 타임스탬프 카운터를 제거함과 동시에, cache communication bottleneck 줄인 LDU는 log-based 알고리즘들의 장점들을 모두 포함할 뿐만 아니라 추가적인 장점을 가진다. 첫째로, update가 수행하는 시점 즉 로그를 저장하는 순간에는 lock이 필요가 없다. 따라서 lock에 대한 오버헤드 없이 concurrent updates를 수행할 수 있다 둘째로, 저장된 update operation log를 coarse-grained lock과 함께 하나

의 코어에서 수행하기 때문에, cache 효율성이 높아진다 [25]. 셋째로, 기존 여러 자료구조에 쉽게 적용할 수 있는 장점이 있다. 게다가 마지막으로, log를 저장하기 전에 로그를 삭제하므로 보다 빠르게 log의 수를 줄일 수 있다.

우리는 위와 같은 장점을 가지는 LDU를 리눅스 커널에서 high update rate 때문에 scalability 문제를 야기시키는 anonymous reverse mapping과 file reverse mapping에 적용하였다. 또한 우리는 LDU를 Linux 4.5.rc4에 구현하였고, fork-intensive 워크로드인 AIM7 [1], Exim [3] from MOSBENCH [2], Imbench [30]를 대상으로 성능 개선을 보였다. 개선은 stock 리눅스 커널에 비해 120코어에서 각각 x,x,x 배이다.

**Contributions.** This paper makes the following contributions:

- 우리는 리눅스 커널을 위한 새로운 log-based concurrent updates 방법인 LDU를 개발하였다. LDU는 동기화된 타임스탬프 카운터를 이용함에 따라 생기는 시간 정렬과 머징에 의한 추가적인 sequential processing 문제를 해결하였다. 이를 위해 LDU는 로그를 업데이트 순간 지우고 로그를 재활용하는 방법을 개발하였다.
- 우리는 LDU를 practical한 manycore system인 intel xeon 120코어 위에 동작하는 리눅스 커널의 2가지 reverse mapping(anonymous, file)에 적용하여, fork scalability 문제를 해결하였다. Fork 관련 벤치마크 성능은 워크로드 특성에 따라 1.6x부터 2.2x까지 개선되었다.

The rest of this paper is organized as follows. Section 2 describes the background and Linux scalability problem. Section 3 describes the design of the LDU algorithm and Section 4 explains how to apply to Linux kernel. section 5 explains our implementations in Linux and Section 6 shows the results of the experimental evaluation. Finally, section 8 concludes the paper.

## 2 Background and Problem

운영체제 커널의 parallelism은 시스템 전체의 parallelism에서 가장 중요하다. 만약에 커널이 scale하지 않으면, 그 위에 동작하는 응용프로그램들도 역시 scale하지 않는다 [15]. 우리는 이처럼 중요한 부분인 운영체제 커널중 multi-core에 최적화된 리눅스의 scalability를 분석하기 위해, AIM7 multiuser[]를 가지고 scalability를 실험해보았다. AIM7은 최근에도 scalability를 위해 reserach 진영과 리눅스 커널 진영에서도 활발히 사용되고 있는 벤치마크 중 하나이다 [12] [11]. File system scalability를 최소화 하기 위해 temp filesystem [34]을 사용하였다. This workload simultaneously create many processes. Up to 60 core, the stock Linux scale linearly, then they flattens out.

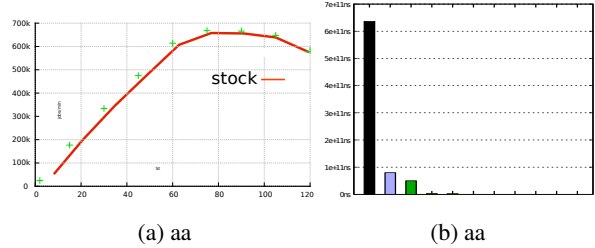


Figure 1: Scalability of AIM7 multiuser. This workload simultaneously create many processes. Up to 60 core, the stock Linux scale linearly, then they flattens out.

우리는 Scalability에 문제가 있는 120코어에서 리눅스의 Lockstat[]를 이용하여 락 경합을 분석하였다. 먼저 멀티 프로세스 기반의 벤치마크인 AIM7을 동작시키고 동시에 120코어 대해서 락 경합을 분석하면 그림 3과 같은 결과를 가진다. AIM7 벤치마크의 경우 상당히 많은 부분이 anonvma에서 쓰기 락 경합이 발생한다. 이는 리눅스 역 매핑(reverse mapping)을 효율적으로 수행하기 위한 자료구인 anonvma를 수많은 fork에 의해 프로세스를 생성하면서 발생하는 락 경합 문제이다. 다음으로 우리는 anonvma의 lock 경합을 줄이기 위해, 임시로 fork에서 anonvma를 호출하는 부분과 read와 관련있는 pageswap이 안되도록 하고, 120코어를 대상으로 다시 Lock 경합을 분석하여 보았다. 이 때 부터 그동안 상대적으로 가려졌던 file reverse mapping에서 많은 락 경합이 발생되었다. 이러한 anonvma reverse page mapping은 리눅스 커뮤니티에서 잘 알려진 락 경합 문제 [5] [36]이고, file mapping에 대한 락 경합 문제는 Silas wikizer가 OpLog 논문을 통해 fork의 scalability 문제의 원인으로 제시한 부분이다. 본 연구의 분석 결과 둘 중 하나가 아니라 두 가지 락 모두 fork의 scalability 문제를 야기 시킨다. 즉 두가지 모두 개선해야지 fork의 scalability가 향상 된다.

두가지 reverse page mapping의 근본적인 문제는 high update operation에 대한 serialization 때문에 발생하는 문제이다. 이러한 reverse page mapping은 page frame reclaiming을 위해 존재하며, 리눅스가 fork(), exit(), and mmap() 시스템콜을 사용할 때 rmap을 update한다. 리눅스 커널은 reader들은 parallel 하게 동작할 수 있는 RW-lock 또는 RCU 같은 락 메커니즘이 있으나, 이러한 락은 결국 high update rate 앞에서는 serialization된다. Update는 exclusive lock을 통해 보호해야하기 때문에, update rate 높은 상황이 발생하면 리눅스 커널은 결국 lock에 의해 serialized되어 scalability가 떨어진다.

이러한 high update rate이 발생하는 상황의 update serialization 문제에 대한 해결 방법들은 존재한다. 근본적인 해결 방법은 update heavy한 상황이 발생하지 않도록 reverse page mapping 알고리즘을 새롭게 디자인하는 방법이 있다. 하지만 이미 오랜 시

간 성능과 안정성이 검증된 알고리즘 citeCorbetL-WNRMAP [18]을 새롭게 디자인 하는 일은 쉽지 않다. 보다 현실적인 방법은 concurrent updates를 위한, non-blocking data structure와 log-based 알고리즘을 사용하는 방법이 있다. Non-blocking algorithms들은 hardware synchronized atomic 연산들을 활용하여 current 하게 update와 read를 수행하게 만든 data structure이다. 하지만 shared memory global value를 multipul CAS로 접근하여 bottlenecks이 생긴다. due to inter-core communication overheads [7]. 최근에는 Deu to the multipul CAS, inter-core communication overheads를 줄인 log-based 방법들이 연구되고 있다. 우리의 LDU도 이러한 log-based 기반 방법 활용하였으며, log-based 방법에 대한 설명은 다음 장에서 자세히 다룬다.

### 3 LDU Design

LDU는 리눅스 커널의 high update rates를 가진 data strcuture의 scalability를 해결하기 위한 log-based 방법 중에 하나이다. 동기화된 타임스탬프 카운터 기반의 per-core log를 활용한 동시적 업데이트방법은 결국 timestamp ordering and merging 작업을 야기한다. 특히 코어 수가 늘어 날 경우, per-core 로그를 자료 구조에 적용하는 과정에서 추가적인 sequential 프로세싱이 요구된다. 이것은 확장성과 성능을 저해한다. 이러한 문제를 해결하기 위해, LDU는 log기반 방식의 concurrent updates 방법과 atomic synchronization 기능을 최소한으로 사용하도록 설계하였다. 따라서, LDU는 synchronized timestamp counter를 사용하는 방식의 timestamp를 제거함과 동시에 cache communication overhead를 최소화 하였다. 예를 들어, high update rates data structure의 scalability 문제를 해결하기 위해, LDU는 update operation 시점에 operation 로그를 atomic한 방법으로 삭제 또는 저장을 한다. 이를 위해 우리는 update-side absorbing log, reuse garbage log라고 불리는 2가지 최적화 방법을 사용하였다. 최적화 이후 남은 update operation에 대한 로그는 워크로드 특성에 따라 global queue 또는 per-core queue에 저장할 수 있다. 본 장에서는 LDU의 Log-based 알고리즘과 어떻게 timestamp를 제거 했는 지에 대한 디자인 측면에 대해서 설명한다.

#### 3.1 Log-based Concurrent updates

Update heavy한 구조 때문에 발생하는 scalability 문제에 대한 해결책 중 하나는 Log-based 알고리즘을 사용하는 것이다. Log-based 알고리즘은 lock을 피하기 위해 update가 발생하면, data structure의 update operation(insert or remove)을 argument와 함께 저장하고, 주기적 또는 read operation을 수행하기 전에 applies the updates in all the logs to the data structure, so reader can read up to date data structure. 이러한 Log-based 방법은 마치 CoW(Copy on Write)와 유사하다. 즉, read

전에 저장된 log가 수행됨으로 read가 간헐적으로 수행되는 data structure에 적합한 방법이다.

Update heavy한 구조를 위한 Log-based 방법은 총 4가지의 장점을 가진다. 첫째로, update가 수행하는 시점 즉 로그를 저장하는 순간에는 lock이 필요가 없다. 따라서 update를 concurrent하게 수행할 수 있을 뿐만 아니라, lock 자체가 가지고 있는 overall coherence traffic is significantly reduced. 둘째로, 저장된 sequantial update operation log를 coarse-grain lock과 함께 하나의 코어에서 수행하기 때문에, cache 효율성이 높아진다. 셋째로, 큰 수정 없이 기존 여러 데이터(tree, queue) structure에 쉽게 적용할 수 있는 장점이 있다. 마지막으로 저장된 log를 실제 수행하지 않고, 여러가지 optimization 방법을 사용하여 적은 operation으로 Log를 줄일 수 있다. LDU도 log-based approach를 따른다. 그러므로 앞에서 설명한 log-based 방법의 장점을 모두 가짐과 동시에 업데이트 순간 time-sensitive log를 지움으로 성능을 향상시킨다.

#### 3.2 Approach

LDU는 update 순간 time-sensitive log를 제거함으로 synchronized timestamp counter 때문에 발생하는 문제를 해결하였다. time-sensitive log란 순서가 바뀌어서는 안되는 operation log이다. Synchronized timestamp가 근본적으로 필요한 이유는 when a process logs an insert operation on one core, migrates to another core, and logs a remove operation, the remove should eventually execute after the insert[]. 본 논문에서는 이러한 time-sensitive log를 설명하기 위해, [15]에서 사용한 심볼방법을 이용하였다. 먼저 ⊕와 ⊖는 각각 insert와 remove update operation을 의미한다. 바로 따라오는 심볼은 object의 명을 의미하고 색깔과 높낮이는 서로 다른 CPU를 의미한다.

⊕<sup>A</sup>, ⊕<sup>B</sup>, ⊕<sup>C</sup>, ⊕<sup>A</sup>, ⊕<sup>C</sup>, ⊕<sup>A</sup>, ⊕<sup>B</sup>, ⊕<sup>C</sup>, ⊕<sup>A</sup>, ⊕<sup>C</sup>

즉 위와 같이 수행되었을 경우 time-sensitive log인 ⊕<sup>A</sup>과 ⊕<sup>C</sup>에 대해서는 항상 timestamp에 맞게 수행되어야 한다. 여기서 중요한 사실은 time-sensitive operation log은 삭제해도 괜찮은 operation log들이다. 여기서 insert-remove operation 또는 remove-insert operation을 가지는 ⊖<sup>B</sup>, ⊕<sup>C</sup> 삭제해도 괜찮으며, 결국 남은 operation log은 같은 object에 대해서 insert 또는 remove operation만 남게 된다. 위의 log는 아래와 같이 non-time-sensitive한 operation만 남게 된다.

⊕<sup>C</sup>, ⊕<sup>A</sup>, ⊕<sup>B</sup>, ⊕<sup>C</sup>, ⊕<sup>A</sup>, ⊕<sup>B</sup>

LDU는 update-side absorbing 기술로 이러한 time-sensitive한 log를 update 순간 바로 지운다. 그 결과 synchronized timestamp의 필요성을 줄였다.

LDU는 time-sensitiv한 log를 제거하기 위해 update-side absorbing이라는 방법을 사용한다. 이 방

법은 만약 같은 object에 대해서 insert와 remove가 발생하였으면, 같은 object에 대해서, insert operation과 remove operation에 대한 log를 update시점에 바로 바로 삭제하는 방법이다. 이처럼 log가 삭제될 수 있는 이유는 operation log가 deferred로 수행하기 문에 가능하다. synchronized timestamp counters 기반의 OpLog도 이러한 log 삭제 방법 수행하여 최적화를 하였으나, operation log가 서로 다른 코어에 존재하는 log 같은 경우에 log를 merge와 검색한 후 삭제를 해야한다. 즉 synchronized timestamp counters 방법은 워크로드에 따라, 최적화를 위해 또 다른 sequential 프로세싱이 요구된다. 하지만 LDU는 individual object를 대상으로 swap atomic operation을 사용하여 shared log 삭제하는 방법을 사용해서 이런 문제가 없다.

업데이트 순간 로그를 지우는 방법은 shared memory system의 swap operation을 사용한다. 이를 위해, LDU는 모든 object에 insert와 remove의 mark 필드를 추가해서 update-side absorbing을 수행 하였다. 예를 들어 만약 A라는 object를 대상으로 insert-remove operation이 수행될 경우 처음 insert operation은 insert mark 필드에 표시하고 queue에 저장한다. 다음 remove operation 부터는 log를 queue에 저장하지 않고 insert에 표시한 mark 필드에 표시한 값만 atomic하게 지워주는 방식으로 진행된다. 다음으로 LDU는 log를 적용할 때, queue안에 log가 존재하더라도, mark 필드가 표시된 log만 실행한다. 이것은 swap이라는 상대적으로 가벼운 연산과 상대적으로 덜 share하는 indivisaul global object의 mark filed를 사용해서 time-sensitive한 operation을 제거할 뿐만 아니라, 동시에 실제 operation을 수행하지 않고 log를 지워주는 효과를 가져주므로 성능이 향상된다.

이처럼 update-side absorbing으로 log를 지울 수 있는 근본 이유는 리눅스 커널(freeBSD 커널 역시) search가 two update operations(one to insert and one to remove an element)이 함수 외부에서 수행하는 구조이기 때문에 가능하다. 즉 리눅스의 data structure의 경우에는 search와 update가 분리되어 있어, 같은 update operation이 연달아 호출되지 않는다. LDU는 이러한 사실을 이용하여 update-side absorbing log를 개발 하였다. 즉 리눅스의 data structure 구조는 search가 update operation 외부에 있기 때문에 insert operation 다음에는 반드시 remove operation이 발생하고 remove operation 다음에는 반드시 insert operation이 발생하는 특징이 있다.

이와 반대로 research 분야에서 많이 연구되는 CSDS(Concurrent search data structures) [20]의 방법들은 search가 update operation 내부에 존재하는 구조로 되어 있다. 따라서 같은 key값의 노드에 대해서 insert-insert 또는 remove-remove 순서로 동작할 수 있다. 이처럼 CSDS 알고리즘일 경우에는 LDU를 적용하지 못하는 limitation이 있다. 하지만 본 연구는 리눅스 커널과 같은 practical한 data structure에 초점을 맞추었기 때문에, research 분야에 많이 연구되고 있는 CSDS 알고리즘은 고려하지 않았다. 아직 CSDS의

non-blocking 알고리즘들은 garbage collector [4], iteration [33], ordering [40] 여러 practical한 이유로 C언어로 구현된 리눅스 커널 등에 적용하기에는 아직 한계가 있다.

LDU의 또 다른 최적화 기법은 update-side absorbing 때문에 취소된 garbage log를 재활용하는 것이다. 이 garbage log는 queue에는 존재하지만 update-side absorbing 때문에 취소되어 garbage가 되어 queue에 보관되어 있는 log이다. 이러한 garbage log일 경우 다음 update operation에 대해서는 해당 로그를 새로 만들어 넣지 않고 기존 로그를 재활용하는 방법이다. 예를 들어 A라는 오브젝트에 대해서 insert-remove-insert 순서로 update가 수행될 경우, 세번째 insert 명령어는 queue에 들어가 있지만 update-side absorbing 때문에 취소되어 insert mark filed가 zero인 경우(garbage object)이다. 이러한 경우, 다음 insert operation에 대해서는 새로운 log를 list에 연산으로 다시 넣지 않고, 해당 object의 mark 필드만 변경하여 log를 재활용하는 방법을 사용한다.

LDU는 queue의 위치(global or per-core)에 대한 의존성을 없애기 위해 log를 항상 non-blocking queue에 저장한다. 따라서 log를 저장하는 시점에 대해서 global lock 또는 per-core lock을 제거 하였다. LDU는 head 포인터에 대한 CAS 연산을 최대한 줄인 multiple producers and single consumer에 기반하는 non-blocking queue를 이용 하였다. 이러한 큐는 이미 Linux 커널에 Lock-less list [39]라는 이름으로 구현되어 있으며, 이미 커널에 많은 부분에 사용되고 있다. 이 큐는 다른 non-blocking list[1][2]와 다르게, insert operation을 항상 처음 노드에 삽입함에 따라, CAS가 발생하는 횟수를 상대적으로 줄일 수 있다. 게다가 single consumer를 고려 했기 때문에, remove를 위한 복잡한 알고리즘이 필요없다. 예를 들어 single consumer는 operation log 전체를 얻기 위해, xchg 명령을 사용하여 head pointer를 NULL로 atomic하게 제거한다. 결론적으로 LDU는 queue의 위치에 따른 의존성을 제거하고, head pointer의 CAS fail 때문에 발생하는 multiple CAS를 최소화 하기 위해, Linux 커널에 있는 Lock-less list를 활용 하였다.

LDU는 log를 저장하기 위해, global 또는 per-core queue 모두 이용할 수 있게 설계하였다. 즉 저장하는 방법이 global queue이면 global non-blocking queue를 이용하고, per-core queue이면 per-core 메모리에 non-blocking queue를 저장하여 사용하였다. 이러한 두가지 queue를 지원하도록 한것은 서로 장단점을 가지고 있기 때문이다. 예를 들어 일반적으로 per-core 방법이 굉장히 높은 scalability를 가진다고 하지만, 본 연구 결과 data structure의 head를 가지는 root object가 굉장히 많이 생성되는 워크로드일 경우, per-core 방법은 너무 많은 메모리를 요구한다. 만약 적은 메모리를 사용할 경우 오히려 성능이 떨어지는 결과를 가진다[section x]. 즉 log-based 방법은 워크로드에 맞는 queue를 사용해야한다.

이와 같이 LDU의 두가지 queue는 서로 장단점



을 가지고 있다. 먼저 global queue의 장점은 굉장히 simple하여 어떠한 data structure간에 쉽게 적용이 가능하다. global queue를 사용한 방법은 LDU의 light-weight한 queue사용과 update-side absorbing log와 reuse garbage log 기법등으로 global head 포인터에 대한 CAS operation 사용을 줄였지만, 여전히 global header에 CAS가 사용되고 있다. 따라서 cache coherence traffic issue가 여전히 남아있어서 scalability 측면으로는 단점을 가지고 있다. Per-core queue일 경우 global head 포인터에 대한 CAS operation을 완전히 제거한 장점을 가진다. 하지만 특정 data structure에는 사용 하지 못하는 data structure의 dependency가 있는 문제가 있다. 즉 LDU를 사용하면 time-sensitvive log가 제거되었다라든 stack과 queue와 같이, 같은 operation에 대한 순서도 중요한 data structure는 사용하지 못한다. 또한 앞에서 설명한 OpLog와 같이 일반적으로 per-core 방법의 근본적인 단점인 메모리를 더 사용하게 된다[].

LDU는 log 때문에 불필요하게 메모리 낭비를 방지하고, To keep the log from growing without end, log-based method periodically apply the time-ordered operations at the head of the log to remove these operation from the log. LDU는 주기적으로 Log를 flush해줌으로써 Log가 쌓여서 발생하는 메모리 낭비를 줄인다.

### 3.3 LDU example

먼저 LDU를 보다 정확하게 설명하기위해, 로그를 global queue 저장한 방법과 per-core queue에 저장한 두가지 방법에 대해서 예를 들어 설명한다. 그림 2 LDU의 하나의 예에 대해서 보여준다. 그림 2은 예로서 queue를 global queue를 사용하였다. 7개의 update operation과 1개의 read operation을 LDU에서 어떻게 처리하는지 보여준다. update operation에 대해서는 앞에서 설명한 심볼을 사용했고 garbage object을 위한 심볼이 추가 되었다. garbage object는 **D**과 같이 rectangle로 표시한다. 이 그림의 실행 순서는 위에서 부터 아래이다. 그림의 왼쪽은 CPU의 operation 순서를 보여주고 오른쪽은 특정 시간에 메모리에 저장된 자료구조의 내용을 보여준다. 초기 data structure의 트리에는 **A** and **B** 두개의 오브젝트가 들어가고, global queue에는 아무것도 없다.

그림 2는 시간 순서대로 각단계에 대해서 LDU의 수행 흐름을 보여준다. 먼저 그림에서 concurrent updates 단계에서는 Core0, Core1 and Core2가 각각 **C**, **D** and **E**에 해당하는 3개의 updates operation을 log로 저장한다. Log를 저장하는 queue를 non-blocking queue를 사용하기 때문에, 이 단계에서의 update는 lock이 필요없다. 즉 모든 thread는 lock에 대한 contention없이 병렬적으로 실행된다. T1 시점이 되면 메모리의 트리는 **A** and **B**의 노드를 가지게 되고, queue에는 **C**, **D** and **E**가 insert mark 필드가 true로 된 상태로 저장되어 있다. update-side absorbing 단계에서는 **D** 명령이 수행되면, 새로운

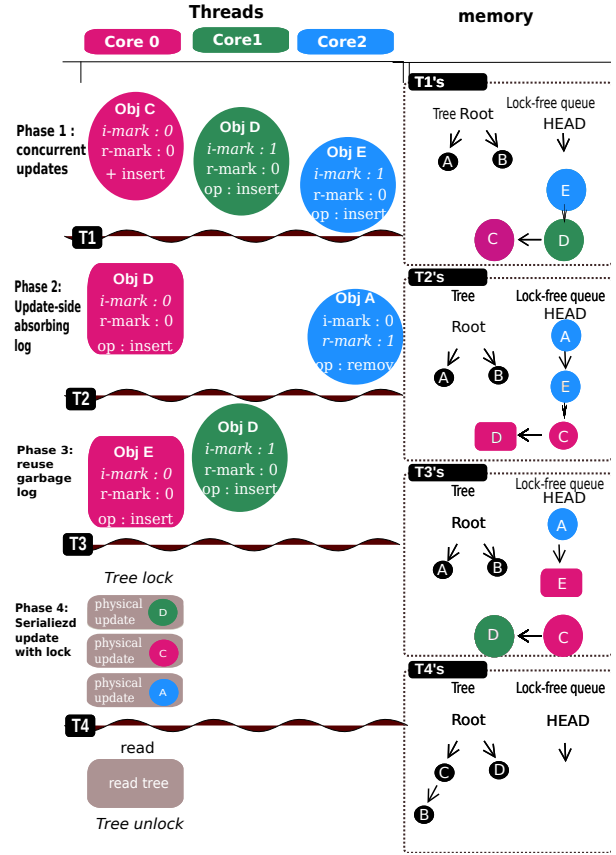


Figure 2: LDU example showing six update operations and one read operation. The execution flows from top to bottom. Memory represents original data structure and logging queue at T1, T2 and T3, respectively.

object를 queue에 넣지 않고, mark 필드만 atomic하게 true로 수정한다. 만약 처음 수행되는 **A** operation이면 queue에 넣는다. T2 시점이 되면 메모리의 트리는 변함없고, queue에는 **C**, **D**, **E** and **A** 저장된다. 여기서 **D**는 mark필드가 false로 되어 garbage가 된 object이다. reuse garbage log 단계에서는 **D**와 같이 queue에는 들어가 있지만, garbage object인 경우 새롭게 queue에 operation log를 넣지 않고, mark 필드 변경으로 재사용하는 것이다. 따라서 **D**는 **D** 상태로 변경된다. T3 시점이 되면 queue에는 **C**, **D**, **E** and **A** 저장된다. 마지막 단계인 저장된 log를 실행하는 단계에서는 tree에서 사용한 방법으로 lock을 건후 하나의 코어에서 전체 로그를 update한다. 여기서 garbage log인 **D**를 제외하고 queue에 저장된 순서대로 **C**, **E** and **A**의 operation을 수행한다. T4 시점이 되면 tree에는 모든 operation이 수행되어, **B**, **C** and **D** 값이 존재하게 된다. 결국 로그를 모두 수행하면 마지막 reader는 같은 tree의 값을 읽게 된다.

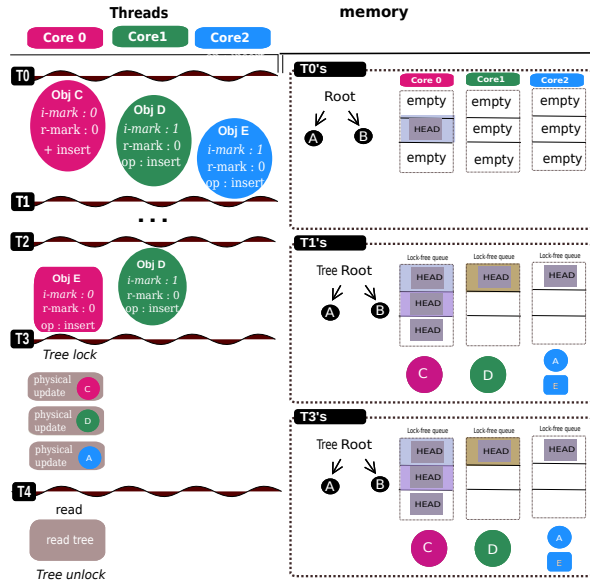


Figure 3: PLDU The execution flows from top to bottom. Memory represents original data structure and logging queue at T1, T2 and T3, respectively.

그림 3는 per-core queue에 log를 저장되는 흐름을 설명한다. global queue와의 차이점만 설명하면,

### 3.4 The LDU Algorithm

#### 3.4.1 logical update: inserting logs

그림 4는 concurrent updates를 수행하는 코드와 앞에서 설명한 update-side absorbing logs와 reuse garbage logs에 대한 코드이다. `logical_insert`과 `logical_remove`의 코드에서 가장 먼저하는 것은 이미 update operation이 수행되어 log가 queue에 들어 있는지에 대해서 체크를 한다. 이것은 update 함수가 수행하는 도중 어느때나 log를 flush하는 `synchronize` 함수가 수행 될 수 있기 때문에, `xchg` 함수를 통해 atomic하게 체크를 한다. 또한 `xchg`를 통해 old mark 필드가 true라면 이미 queue에 log가 있다는 뜻이므로, atomic하게 false로 수정을 한다. 이것은 앞에서 설명한 update-side absorbing log 기능이며, 1개의 atomic operation으로 update를 수행할 수 있다. 만약 mark 필드가 체크되지 않았다면, 두번째로 queue에 존재하는 로그인지 체크를 한 후, 존재하는 로그라면 이미 앞에서 mark 필드를 설정하여 이미 로그를 재활용했으므로 queue에 로그를 넣지 않는다. 만약 처음 사용된 log라면 non-blocking queue에 넣는다. 이 코드는 항상 옳게 동작한다. 그 이유는 리눅스 update operation이 항상 insert-remove 순서나 remove-insert 순서로 실행되기 때문이다.

```

1 bool ldu_logical_insert(struct object_struct *obj,
2 void *head)
3 {
4     ...
5     // Phase 1 : update-side absorbing logs
6     // atomic swap due to synchronize update's logs
7     if(!xchg(&remove->mark, 0)){
8         BUG_ON(insert->mark);
9         // Phase 2 : reuse garbage log
10        insert->mark = 1;
11        if(!test_and_set_bit(LDU_INSERT,
12                             &obj->ldu.used)){
13            // Phase 3 : insert log to queue
14            //...save argument and operation
15            ldu_insert_queue(root, insert);
16        }
17    }
18    ...
19 }
20
21 bool ldu_logical_remove(struct object_struct *obj,
22 void *head)
23 {
24     ...
25
26     // Phase 1 : update-side absorbing logs
27     // atomic swap due to synchronize update's logs
28     if(!xchg(&insert->mark, 0)){
29         BUG_ON(remove->mark);
30         // Phase 2 : reuse garbage log
31         remove->mark = 1;
32         if(!test_and_set_bit(LDU_REMOVE,
33                              &obj->ldu.used)){
34            // Phase 3 : insert log to queue
35            //...save argument and operation
36        }
37    }
38    ...
39 }

```

Figure 4: LDU logical update algorithm. `logical_insert` represents non-blocking insert function. It may be called by original insert position without locks. The fastpath is that when their object was removed by `logical_remove`, `logical_insert` just changes node's marking field.

#### 3.4.2 physical update: applying logs

그림 5는 concurrent updates를 통해 저장된 log들을 적용하는 deferred update 함수에 대한 코드이다. `synchronize_ldu`는 마치 CoW 처럼 read 전에 호출되기도 하지만, 로그가 쌓이는것을 방지하기 위해 주기적으로 호출되어 로그를 비워준다. 이 함수는 항상 object의 lock이 걸린 상황에서 수행된다. LDU는 single consumer가 독점적으로 log를 처리한다. 따라서 제일 먼저 하는일은 queue에 저장된 log를 얻기 위해 atomic `xchg` operation으로 log의 head 포인터를 얻는다. 그리고 이 순간 부터는 single consumer가 독점적으로 log를 처리하게 된다. LDU는 주기적으로

```

1 void synchronize_ldu(struct obj_root *root)
2 {
3     ...
4
5     //atomic remove first, lock-less list
6     entry = xchg(&head->first, NULL);
7
8     //iterate all logs
9     llist_for_each_entry(dnode, entry, ll_node) {
10         //get log's arguments
11         ...
12         //atomic swap due to update-side removing
13         if (xchg(&dnode->mark, 0))
14             ldu_physical_update(dnode->op_num, arg,
15                                 ACCESS_ONCE(dnode->root));
16         clear_bit(dnode->op_num, &vma->dnode.used);
17         // one more check due to reuse garbage log
18         if (xchg(&dnode->mark, 0))
19             ldu_physical_update(dnode->op_num, arg,
20                                 ACCESS_ONCE(dnode->root));
21     }
22 }

```

Figure 5: LDU physical update algorithm. `synchronize_ldu` may be called by reader and converts update log to original data structure traversing the lock-less list.

log를 비워줘야 하므로, update operation이 수행되는 도중에도 `synchronize_ldu`가 호출될 수 있다. 따라서 항상 original data structure에 적용하기 전에 mark 필드를 초기화해야 한다. 또한 log를 사용했으면, used 비트를 초기화 하여 해당 log가 queue 없음을 표시한다. 하지만 mark 필드 초기화 과정과 used 비트를 초기화 하는 과정에서, reuse garbage logs 기능 때문에, 언제든지 다시 concurrent updates operation을 통해 mark 필드가 설정될 수 있다. 따라서 한번 더 mark 필드를 체크하여 log를 비워준다.

### 3.4.3 LDU queues

LDU의 로그를 저장하기 위한 2가지 queue에 대한 코드는 그림 6와 같다. 앞에서 설명했듯이 LDU는 어떤 queue를 사용해도 문제 없도록 설계하였다. 따라서 queue 대한 의존성이 없다. global queue를 사용할 경우 LDU는 항상 head의 first 포인터에 CAS연산으로 성공할 때까지 반복 수행 후 데이터를 넣는다. per-core queue를 사용할 경우는 per-core의 hash 테이블을 가지고 온 후 hash 테이블이 넣으려는 log의 data structure와 같은지 체크하고, 다르면 해당 코어에 저장된 log를 lock과 함께 update를 수행한다. 전체 코어의 log를 flush할 필요없이 해당 코어에 저장된 log만 flush하면 되는 이유는 이미 update-side absorbing에 의해 time-sensitive한 operation이 삭제되었기 때문이다.

```

1 bool insert_log_global_queue(struct obj_root *root,
2                             struct ldu_node *dnode)
3 {
4     ...
5     //do {
6     // new_last->next = first = head->first;
7     //} while (cmpxchg(&head->first, first, new_first) != first);
8     llist_add(&dnode->ll_node, &p->list);
9     ...
10 }
11
12 bool insert_log_per_core_queue(struct obj_root *root,
13                               struct ldu_node *dnode)
14 {
15     ...
16
17     slot = &get_cpu_var(obj_root_slot);
18     p = &slot->obj[hash_ptr(root, HASH_ORDER)];
19     empty = READ_ONCE(p->list.first);
20     // is empty list?
21     if (!empty) {
22         ldu = llist_entry(first, struct ldu, ll_node);
23         // is hash complicit?
24         if (ldu->root != dnode->root) {
25             ...
26             lock = READ_ONCE(ldu->lock);
27             entry = xchg(&p->list->head->first, NULL);
28             llist_add(&dnode->ll_node, &p->list);
29             // flush log as a direct mapped cache
30             object_lock(&lock);
31             synchronize_ldu(entry);
32             object_unlock(&lock);
33             goto out;
34         }
35     }
36
37     llist_add(&dnode->ll_node, &p->list);
38 out:
39     ...
40 }

```

Figure 6: LDU physical update algorithm. `synchronize_ldu` may be called by reader and converts update log to original data structure traversing the lock-less list.

## 4 Concurrent updates for Linux kernel

### 4.1 Case study:reverse mapping

리눅스 커널의 프로세스간 공유자원 중 하나인 reverse page mapping(rmap)은 fork, exit, mmap이 수행될 때 update가 많이 발생하는 data structure이다. Linux's reverse mapping reads walked through each process's pages mappings and selected pages to unmap when it swaps a physical page out to disk, migrates other cpu, or truncates a file. After all a page's mappings were removed could it be selected for pageout [?]. Rmap의 anonymous page와 file page의 관리는 최근 interval trees로 되어 있으며, 이것은 reverse page 성능 향상을 위해 지속적으로 최적화가 [17] [18] 이루어지

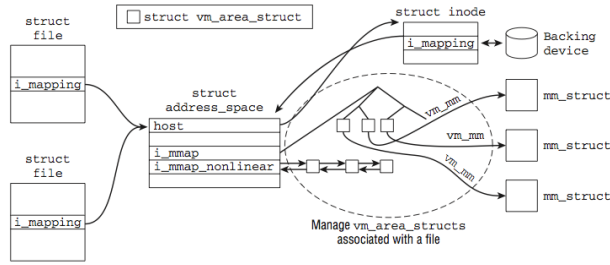


Figure 4-7: Tracking the virtual address spaces into which a given interval of a file is mapped with the help of a priority tree.

Figure 7: An example of applying the LDU to file reverse mapping.

고 있다.

리눅스는 rmap의 interval tree를 보호하기 위해 rw semaphore를 사용한다. 따라서, 프로세스들이 fork, exit, mmap를 simultaneously 수행하면 interval tree를 보호하는 read-write semaphore 때문에 scalability가 떨어진다. anonymous page를 위한 rmap과 file mapped page를 위한 rmap 모두 문제가 있다. 이번 장은 이러한 문제를 해결하기 위해, 어떻게 LDU를 리눅스의 rmap에 적용했는지에 대한 내용을 설명하며, 보다 practical한 내용을 다룬다.

## 4.2 anonymous mapping

Anonymous rmap의 공유데이터는 서로 상당히 복잡하게 연결되어 있다. 그림 x-x는 이렇게 복잡하게 연결된 공유데이터를 보여준다.

Anonymous rmap을 위한 GLDU는 anon\_vma의 root structure에서 log를 저장하도록 하였다. anon\_vma는 anon\_vma의 root의 lock을 사용함에 따라, anon\_vma의 root의 lock이 걸리면 자식들은 모두 root의 lock 때문에 많은 contention이 생긴다. 따라서 GLDU는 anon\_vma는 항상 root의 lock를 사용하기 때문에, 로그를 anon\_vma의 root에 저장하도록 하여 불필요한 iteration을 제거하였다.

Anonymous rmap에는 per-core 방식인 PLDU를 적용하기 힘들다. 그 이유는 log를 저장하는 anon\_vma root의 object가 너무이 생겨서 global per-core hash table의 충돌이 너무 많이 발생하기 때문이다. 이것은 또 다른 lock이 필요하므로 오히려 성능이 떨어진다. 만약 per-core 방식을 적용하려면 하나의 per-core memory에 저장한 후 차후 log를 object별 분리 하는 방법이 있으나, 이 경우는 mix되어 있는 로그를 구별 하면서 해당 data structure만 보호하는 lock을 사용하기가 어려우므로, global lock을 사용해야한다. 이 방법도 역시 global lock에 의한 또 다른 scalability 문제를 야기한다.

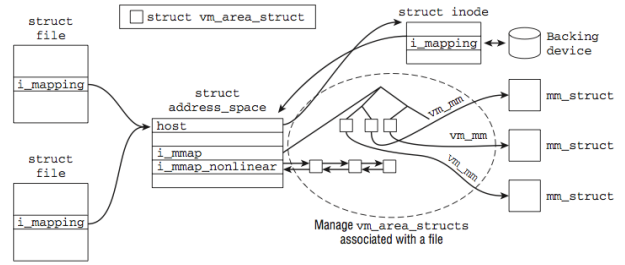


Figure 4-7: Tracking the virtual address spaces into which a given interval of a file is mapped with the help of a priority tree.

Figure 8: An example of applying the LDU to file reverse mapping.

## 4.3 file mapping

file rmap의 공유데이터는 anonymous rmap보다는 덜 복잡하게 연결되어 있다. 그림 x는 file rmap의 공유 데이터를 보여준다.

File mapping의 GLDU는 address\_space object이 interval tree의 head를 가지고 있으므로, 그 곳에 operation log저장한다. File mapping의 PLDU는 per-core hash 테이블로 구현하였고, PLDU를 적용하는데 문제가 없다. 그 이유는 상대적으로 log의 header를 가지고 있는 address\_space가 상대적으로 적게 생성이 되어 hash 충돌이 적게 나타나기 때문에 scalability가 떨어지지 않는다.

## 5 Implementation

We implemented the new deferred update algorithm in Linux 3.19.rc4 kernel, and our modified Linux is available as open source. LDU's scheme is based on deferred processing, so it needs a garbage collector for delayed free.

## 6 Evaluation

This section answers the following questions experimentally:

- Does LDU's design matter for applications?
- Why does LDU's scheme scale well?
- What about LDU's read-write ratio?

### 6.1 Experimental setup

To evaluate the performance of LDU, we use well-known three benchmarks: AIM7 Linux scalability benchmark, Exim email server in MOSBENCH and Imbench. We selected these three benchmarks because they are



fork-intensive workloads and exhibit high reverse mapping lock contentions. Moreover, AIM7 benchmark has widely been used in practical area not only for testing the Linux but also for improving the scalability. To evaluate LDU for real world applications, we use Exim which is the most popular email server. A micro benchmark, Lmbench, has been selected to focus on Linux fork operation-intensive fine grained evaluations. Finally, we wanted to focus on Linux fork performance and scalability; therefore, we selected lmbench, a micro benchmark.

In order to evaluate Linux scalability, we used four different experiment settings. First, we used the stock Linux as the baseline reference. Second, we used ordered Harris lock-free list while we apply unordered Harris lock-free list for the third setting (see section 5). Finally, we used combination of unordered Harris lock-free list for anonymous mapping and our LDU for file mapping. Since we cannot obtain detailed implementation of Oprog, we could not include comparison between LDU and Oprog in this paper.

We compare our LDU implementation to a concurrent non-blocking Harris linked list [24]; therefore, we implement the Harris. The code refers from sysnchrobench [23] and ASCYLIB [20], and we convert their linked list to Linux kernel style. Because both synchrobench and ASCYLIB leak memory, we implement additional garbage collector for the Linux kernel using Linux’s work queues and lock-less list.

In order to further improve performance, we move their ordered list to unordered list. A feature of the Harris linked list is all the nodes are ordered by their key. Zhang [40] implements a lock-free unordered list algorithm, whose list is each insert and remove operation appends an intermediate node at the head of the list; these approach is practically hard to implement. Indeed, Linux does not require contains operation because the Linux data structures such as list, tree and hash table not depended on search key; they depend on their unique object. This feature can eliminates the ordered list in Harris linked list. Therefore, we perform each insert operation appends an intermediate node at the first node of the list; on the other hand, each remove operation searches from head to their node.

We ran the three benchmarks on Linux 3.19.rc4 with stock Linux with the automatic NUMA balancing feature disabled because the Harris linked list has the iteration issue [33]. All experiments were performed on a 120 core machine with 8-socket, 15-core Intel E7-8870 chips equipped with 792 GB DDR3 DRAM.

## 6.2 AIM7

AIM7 forks many processes, each of which concurrently runs. We used AIM7-multiuser, which is one of work-

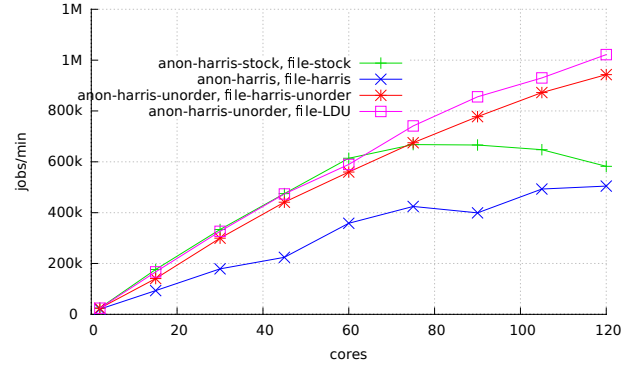


Figure 9: Scalability of AIM7 multiuser for different method. The combination LDU with unordered harris list scale well; in contrast, up to 60 core, the stock Linux scale linearly, then it flattens out.

load in AIM7. The multiuser workload is composed of various workloads such as disk-file operations, process creation, virtual memory operations, pipe I/O, and arithmetic operation. To minimize IO bottlenecks, the workload was executed with tmpfs filesystems, each of which is 10 GB. To increase the number of users during our experiment and show the results at the peak user numbers, we used the crossover.

The results for AIM7-multiuser are shown in Figure 9, and the results show the throughput of AIM7-multiuser with four different settings. Up to 60 core, the stock Linux scales linearly while serialized updates in Linux kernel become bottlenecks. However, up to 120core, unordered harris list and our LDU scale well because these workloads can run concurrently updates and can reduce the locking overheads due to reader-writer semaphores(anon\_vma, file). The combination of LDU with unordered harris list has best performance and scalability outperforming stock Linux by 1.7x and unordered harris list by 1.1x. While the unordered harris list has 19% idle time(see Table 1), stock Linux has 51% idle time waiting to acquire both anon\_vma’s rwsem and file’s immmap\_rwsem. We can notice that although LDU has 23% idle time, the throughput is higher than unordered harris list. In this benchmark, the ordered harris list has the lowest performance and scalability because their CAS fails frequently.

## 6.3 Exim

To measure the performance of Exim, shown in Figure 10, we used default value of MOSBENCH to use tmpfs for spool files, log files, and user mail files. Clients run on the same machine and each client sends to a different user to prevent contention on user mail file. The

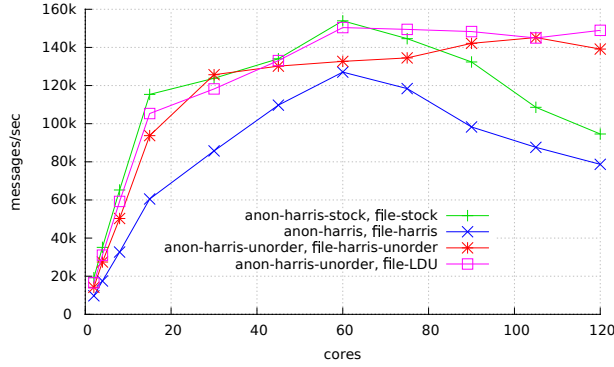


Figure 10: Scalability of Exim. The stock Linux collapses after 60 core; in contrast, both unordered harris list and our LDU flatten out.

AIM7	user	sys	idle
Stock(anon, file)	2487 s	1993 s	4647 s(51%)
H(anon, file)	1123 s	3631 s	2186 s(31%)
H-unorder(anon, flie)	3630 s	2511 s	1466 s(19%)
H-unorder(anon), L(file)	3630 s	1903 s	1662 s(23%)
EXIM	user	sys	idle
Stock(anon, file)	41 s	499 s	1260 s(70%)
H(anon, file)	47 s	628 s	1124 s(62%)
H-unorder(anon, file)	112 s	1128 s	559 s(31%)
H-unorder(anon), L(file)	87 s	1055 s	657 s(37%)
lmbench	user	sys	idle
Stock(anon, file)	11 s	208 s	2158 s(91%)
H(anon, file)	11 s	312 s	367 s(53%)
H-unorder(anon, file)	11 s	292 s	315 s(51%)
H-unorder(anon), L(file)	12 s	347 s	349 s(49%)

Table 1: Comparison of user, system and idle time at 120 cores.

Exim was bottlenecked by per-directory locks protecting file creation in the spool directories and by forks performed on different cores [9]. Therefore, although we eliminate the fork problem, the Exim may suffer from contention on spool directories.

Results shown in Figure 10 show that Exim scales well for all methods up to 60 core but not for higher core counts. The stock Linux shows performance degradation for more than 60 core. Both unordered harris list and our LDU do not suffer from performance loss because they do not acquire the `anon_vma` semaphore and `i_mmap` semaphore in fork. LDU performs better due to the fact that it uses both update-side absorbing and lock-less list, outperforming stock Linux by 1.6x and unordered harris list by 1.1x. Even though we applied scalable solution,

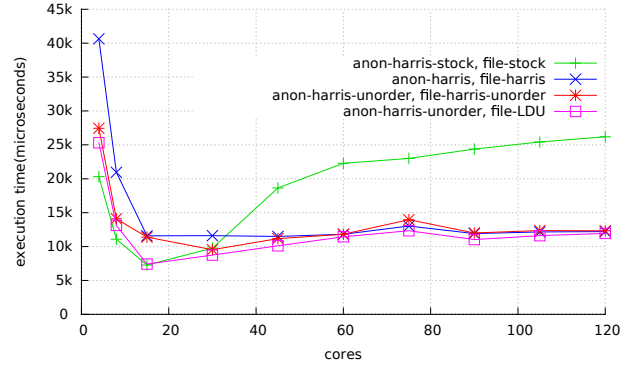


Figure 11: Execution time of lmbench's fork micro benchmark. The fork micro benchmark drops down for all methods up to 15 core but either flattens out or goes up slightly after that. At 15 core, the stock Linux goes up; the others flatten out

Exim shows limitation on scalability improvement since the main bottleneck is per-directory lock contention on spool directories. The unordered harris list has 31% idle time, whereas LDU has 37% idle time due to their efficient concurrent updates.

## 6.4 Lmbench

Lmbench has various workloads including process creation workload(fork, exec, sh -c, exit). This workload is used to measure the basic process primitives such as creating a new process, running a different program, and context switching. We configured process create workload to enable the parallelism option which specifies the number of benchmark processes to run in parallel [30]; we used 100 processes.

The results for lmbench are shown in Figure 11, and the results show the execution times of the fork microbenchmark in lmbench with four different methods. Three methods outperform stock Linux by 2.2x at 120 cores; however, before 30 core, two harris list have lower performance due to their execution overheads. While stock Linux has 90% idle time, other methods have approximately 50% idle time since stock Linux waits to acquire reverse mapping locks such as `anon_vma`'s `rwsem` and mapping's `i_mmap_rwsem`.

## 6.5 Read-write ratio

LDU는 update-heavy한 데이터 구조를 위한 방법이다. 즉 간헐적으로 read가 발생할 때 굉장히 높은 scalability를 보여준다. 하지만 To understand read의 비율을 보기 위해 우리는 아래와 같은 실험을 하였다.

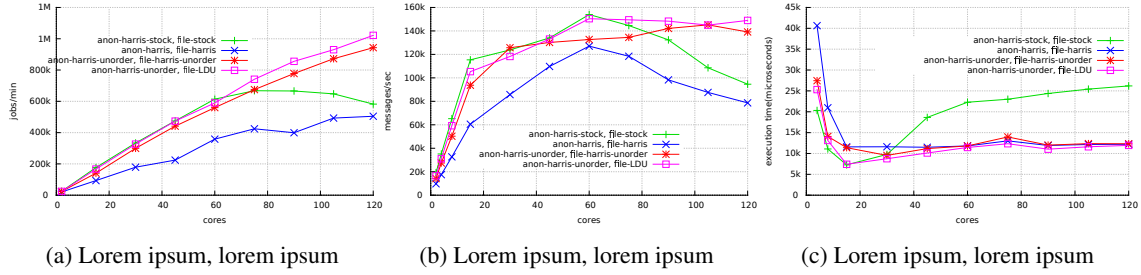


Figure 12: Caption place holder3

## 7 Discussion and future work

우리는 high update rate data structure만 고려 하였지만, LDU는 보다 범용적인 상황에 적용 수 있게 만들어져야 한다. 그 이유는 실제 high update rate를 가진 상황은 자주 발생하는 상황이 아니다. 따라서 data structure의 operation의 update rate이 적을 경우에는 불필요하게 log를 저장하므로 오히려 성능이 떨어지는 역효과가 생긴다. 이것은 임베디드 시스템 등 다양한 분야에 사용되는 리눅스 커널이 때문에 문제가 있다. 따라서 항상 log-based로 처리하지 않고 update rate이 증가하면, 판단하여 log-based로 처리하도록 해야 한다. 즉 high update rate과 상관없이 좋은 성능을 보이는 범용적인 알고리즘으로 발전 시켜야 한다.

LDU의 두번 째 타입은 Per-core에 log를 저장하는 PLDU이다. PLDU의 원리는 기본적으로 log를 per-core 메모리에 저장하되 timestamp가 필요한 log만 shared memory system의 atomic 연산으로 처리를 하는 것이다. When a process logs an insert operation on one core, migrates to another core, and logs a remove operation, the remove should eventually execute after the insert, so this condition needs the timestamp[1].

즉 timestamp 없이 per-core로 merge된 로그가 아래와 같은 상황이 발생하면 문제가 발생한다.

A B C A C

여기서 흥미로운 사실은 같은 object에 대해서 insert-remove operation을 update 할 때 수행해 주면 결국 per-core에 남은 operation은 insert 또는 remove operation만 남게 된다. 즉 per-core log는 timestamp가 없어도 되는 operation만 남게 된다. PLDU는 로그를 update-side absorbing 기술로 로그를 바로 바로 지워 줌으로 timestamp를 없이도 per-core log를 사용할 수 있도록 하였다.

현재 LDU는 data structure를 보호하기 위해 exclusive lock를 사용하도록 구현하였다. 하지만 concurrent read를 위해 concurrent update와 concurrent read를 고려해서 구현해야 한다.

## 8 Related work

**Operating system scalability.** [15]In order to improve Linux scalability, researchers have been optimized memory management in Linux by finding and fixing scalability bottlenecks [8] [10]. Shared address spaces in multithreaded applications easily become scalability bottlenecks since kernel operations including mmap and munmap system calls and page faults handling require per-process locks for synchronization. Multithreaded application, for example, can become bottleneck by kernel operations on their shared address space, whose operations are the mmap and munmap system calls and page faults. These operations are synchronized by a single per-process lock. BonsaiVM [13] solved this address space problem by using the RCU; RadixVM [14] created a new VM using refcache and radix tree, which enable munmap, mmap, and page fault on non-overlapping memory regions to scale perfectly. Alternatively, to avoid contention caused by shared address space locking, system programmers change their multithreaded applications to use processes [9].

**Scalable data structure.** [21]Though sufficient level of performance scalability has been achieved for reader intensive operations through RCU and Hazard pointer, solutions to scalability for update-heavy operations has not been satisfiable. A recent paper by Arbel and Attiya [6] shows a new design of concurrent search tree called the Citrus tree. The Citrus tree combines RCU and fine-grained locks, and it supports concurrent write operations that traverse the search tree by using RCU concurrently. When increasing the update rate, Citrus tree still suffers from bottlenecks. RLU [26] presents a new synchronization mechanism that allows unsynchronized sequences of reads to execute concurrently with updates. In high update rate, Oplog can achieve substantially multi-core scaling for update-heavy data structures. Our work focus on update-heavy data structures and uses non-blocking method to store the operation log instead of per-core processing.

**Scalable lock.** [38] [12] [11]One method for the concurrent update is using the non-blocking algo-

gorithms [24] [22] [37], which are based on CAS.

MCS [31], a scalable exclusive lock, is used in the Linux kernel [16]. to avoid unfairness at high contention levels, so this scalable exclusive lock can be used for fine grained locking in Linux. However, in MCS, since only one thread may hold the lock at a time, it can cause low scalability in case of long critical regions. Reader-writer lock [19] allows either number of readers to execute concurrently or single writer to execute. Thus, readers-writer locks allow better scalability in case of read-mostly objects.

In read-mostly data structures, RCU [29] can be quite useful since it allows read operations to proceed without read locks, and delays freeing of data structures to avoid races. One drawback is that as the update rate increases, their performance and scalability decrease due to a single writer and their synchronization function. Consequently, scalable exclusive lock, reader-writer lock and RCU require serialization for updates and thus show significant limitation on scalability.

## 9 Conclusion

We propose a concurrent update algorithm, LDU, for update-heavy data structures scalable for many-core systems. To achieve the scalability during process spawning, we applied deferred log processing with global log queue and update-side absorbing to Linux reverse mapping. Evaluation results using the AIM7, Exim and Imbench reveal that LDU shows better performance up to 2.2 times compared to existing solutions. LDU is implemented on to Linux kernel 3.19 and available as open-source from <https://github.com/KMU-embedded/scalablelinux>.

## References

- [1] Aim benchmarks. <http://sourceforge.net/projects/aimbench>.
- [2] Mosbench. 2010. <https://pdos.csail.mit.edu/mosbench/mosbench.git>.
- [3] Exim internet mailer. 2015. <http://www.exim.org/>.
- [4] Samy Al Bahra. Nonblocking algorithms and scalable multicore programming. *Commun. ACM*, 56(7):50–61, July 2013.
- [5] Tim Chen Andi Kleen. Scaling problems in fork. In *Linux Plumbers Conference, September*, 2011.
- [6] Maya Arbel and Hagit Attiya. Concurrent updates with rcu: Search tree as an example. In *Proceedings of the 2014 ACM Symposium on Principles of Distributed Computing*, PODC '14, pages 196–205, New York, NY, USA, 2014.
- [7] Silas Boyd-Wickizer. Optimizing communications bottlenecks in multiprocessor operating systems kernels. In *PhD thesis, Massachusetts Institute of Technology*, 2013.
- [8] Silas Boyd-Wickizer, Haibo Chen, Rong Chen, Yandong Mao, Frans Kaashoek, Robert Morris, Aleksey Pesterev, Lex Stein, Ming Wu, Yuehua Dai, Yang Zhang, and Zheng Zhang. Corey: An operating system for many cores. In *Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation*, OSDI'08, pages 43–57, Berkeley, CA, USA, 2008. USENIX Association.
- [9] Silas Boyd-Wickizer, Austin T. Clements, Yandong Mao, Aleksey Pesterev, M. Frans Kaashoek, Robert Morris, and Nickolai Zeldovich. An analysis of Linux scalability to many cores. In *9th USENIX Symposium on Operating System Design and Implementation*, pages 1–16, Vancouver, BC, Canada, October 2010.
- [10] Silas Boyd-Wickizer, M. Frans Kaashoek, Robert Morris, and Nickolai Zeldovich. Non-scalable locks are dangerous. In *In Proceedings of the Linux Symposium*, 2012.
- [11] D. Bueso and S. Norto. An overview of kernel lock improvements. 2014. <http://events.linuxfoundation.org/sites/events/files/slides/linuxcon-2014-locking-final.pdf>.
- [12] Davidlohr Bueso. Scalability techniques for practical synchronization primitives. *Commun. ACM*, 58(1):66–74, December 2014.
- [13] Austin T. Clements, M. Frans Kaashoek, and Nickolai Zeldovich. Scalable address spaces using RCU balanced trees. In *Architectural Support for Programming Languages and Operating Systems (ASPLoS 2012)*, pages 199–210, London, UK, February 2012.
- [14] Austin T. Clements, M. Frans Kaashoek, and Nickolai Zeldovich. Radixvm: Scalable address spaces for multithreaded applications. In *Proceedings of the 8th ACM European Conference on Computer Systems*, EuroSys '13, pages 211–224, New York, NY, USA, 2013.
- [15] Austin T. Clements, M. Frans Kaashoek, Nickolai Zeldovich, Robert T. Morris, and Eddie Kohler.



- The scalable commutativity rule: Designing scalable software for multicore processors. *ACM Trans. Comput. Syst.*, 32(4):10:1–10:47, January 2015.
- [16] J. Corbet. *MCS locks and qspinlocks*. 2014.
  - [17] Jonathan Corbet. The object-based reverse-mapping vm. 2003. <https://lwn.net/Articles/23732/>.
  - [18] Jonathan Corbet. The case of the overly anonymous anon\_vma. 2010. <https://lwn.net/Articles/383162/>.
  - [19] P. J. Courtois, F. Heymans, and D. L. Parnas. Concurrent control with “readers” and “writers”. *Communications of the ACM*, 14(10):667–668, October 1971.
  - [20] Tudor David, Rachid Guerraoui, and Vasileios Trigonakis. Asynchronized concurrency: The secret to scaling concurrent search data structures. In *Proceedings of the Twentieth International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS ’15*, pages 631–644, New York, NY, USA, 2015.
  - [21] Mike Dodds, Andreas Haas, and Christoph M. Kirsch. A scalable, correct time-stamped stack. In *Proceedings of the 42nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL ’15*, pages 233–246, New York, NY, USA, 2015. ACM.
  - [22] Mikhail Fomitchev and Eric Ruppert. Lock-free linked lists and skip lists. In *Proceedings of the Twenty-third Annual ACM Symposium on Principles of Distributed Computing, PODC ’04*, pages 50–59, 2004.
  - [23] Vincent Gramoli. More than you ever wanted to know about synchronization: Synchrobench, measuring the impact of the synchronization on concurrent algorithms. In *Proceedings of the 20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, PPOPP 2015*, pages 1–10, New York, NY, USA, 2015.
  - [24] Timothy L. Harris. A pragmatic implementation of non-blocking linked-lists. In *Proceedings of the 15th International Conference on Distributed Computing, DISC ’01*, pages 300–314, London, UK, UK, 2001.
  - [25] Danny Hendler, Itai Incze, Nir Shavit, and Moran Tzafrir. Flat combining and the synchronization-parallelism tradeoff. In *Proceedings of the Twenty-second Annual ACM Symposium on Parallelism in Algorithms and Architectures, SPAA ’10*, pages 355–364, New York, NY, USA, 2010. ACM.
  - [26] Alexander Matveev, Nir Shavit, Pascal Felber, and Patrick Marlier. Read-log-update: A lightweight synchronization mechanism for concurrent programming. In *Proceedings of the 25th Symposium on Operating Systems Principles, SOSP ’15*, pages 168–183, New York, NY, USA, 2015.
  - [27] Paul McKenney. Msome more details on read-log-update. 2016. <https://lwn.net/Articles/667720/>.
  - [28] Paul E McKenney. Is parallel programming hard, and, if so, what can you do about it? 2011.
  - [29] Paul E. McKenney and John D. Slingwine. Read-copy update: Using execution history to solve concurrency problems. In *Parallel and Distributed Computing and Systems*, pages 509–518, Las Vegas, NV, October 1998.
  - [30] Larry W McVoy, Carl Staelin, et al. Imbench: Portable tools for performance analysis. In *USENIX annual technical conference*, pages 279–294. San Diego, CA, USA, 1996.
  - [31] John M. Mellor-Crummey and Michael L. Scott. Scalable reader-writer synchronization for shared-memory multiprocessors. In *Proceedings of the Third PPOPP*, pages 106–113, Williamsburg, VA, April 1991.
  - [32] Changwoo Min, Sanidhya Kashyap, Steffen Maass, and Taesoo Kim. Understanding manycore scalability of file systems. In *2016 USENIX Annual Technical Conference (USENIX ATC 16)*, pages 71–85, Denver, CO, June 2016. USENIX Association.
  - [33] Erez Petrank and Shahar Timnat. Lock-free data-structure iterators. In *DISC ’13 Proceedings of the 27th International Conference on Distributed Computing, Jerusalem, Israel*. 2013.
  - [34] C. Rohland. Tmpfs is a file system which keeps all files in virtual memory. 2001. [git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/Documentation](http://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git/tree/Documentation).
  - [35] Ori Shalev and Nir Shavit. Predictive log-synchronization. *SIGOPS Oper. Syst. Rev.*, 40(4):305–315, April 2006.
  - [36] Dave Hansen Tim Chen, Andi Kleen. Linux scalability issues. In *Linux Plumbers Conference, September*, 2013.
  - [37] Shahar Timnat, Anastasia Braginsky, Alex Kogan, and Erez Petrank. Wait-free linked-lists. In *Proceedings of the 17th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*,

PPoPP '12, pages 309–310, New York, NY, USA, 2012.

- [38] Tianzheng Wang, Milind Chabbi, and Hideaki Kimura. Be my guest: Mcs lock now welcomes guests. In *Proceedings of the 21st ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, PPoPP '16, pages 21:1–21:12, New York, NY, USA, 2016. ACM.
- [39] Huang Ying. Lock-less list. 2011. <https://lwn.net/Articles/423366/>.
- [40] Kunlong Zhang, Yujiao Zhao, Yajun Yang, Yujie Liu, and Michael Spear. Practical non-blocking unordered lists. In *DISC '13 Proceedings of the 27th International Conference on Distributed Computing, Jerusalem, Israel*. 2013.