

## PCT 출원서

(전자적 형태가 원본)

<b>0</b>	수리관청 전용	
<b>0-1</b>	국제출원번호	
<b>0-2</b>	국제출원일자	
<b>0-3</b>	수리관청 명칭 및 "PCT 국제출원"	
<b>0-4</b>	서식 PCT/RO/101 - PCT 출원서	
0-4-1	우측에 기재된 바와 같이 작성되었다.	<b>PCT-SAFE</b> <b>버전 3.51.087.263 MT/FOP 20190401/0.20.5.24</b>
<b>0-5</b>	신청 아래 서명인은 본 국제 출원서가 특허협력조약에 의해 처리될 것을 청구합니다.	
<b>0-6</b>	출원인이 지정한 수리관청	<b>대한민국 특허청 (RO/KR)</b>
<b>0-7</b>	출원인 또는 대리인의 서류참조기호	<b>PCT19-0023</b>
<b>I</b>	발명의 명칭	<b>비휘발성 메모리를 이용한 로그 구조 병합 트리 기반의 데이터 베이스의 데이터 처리 방법</b>
<b>II</b>	출원인	
II-1	이 사람은	<b>오직 출원인 (applicant only)</b>
II-2	우측 지정국에 관한 출원인	<b>모든 지정국 (all designated States)</b>
II-4ko	성명	<b>연세대학교 산학협력단</b>
II-4en	Name:	<b>UNIVERSITY-INDUSTRY FOUNDATION(UIF), YONSEI UNIVERSITY</b>
II-5ko	주소	<b>대한민국 03722</b>
II-5en	Address:	<b>서울시 서대문구 연세로 50 50, Yonsei-ro, Seodaemun-gu, Seoul 03722 Republic of Korea</b>
II-6	국적	<b>대한민국 KR</b>
II-7	거주국	<b>대한민국 KR</b>
II-8	전화번호	<b>82-2-2123-5138</b>
II-9	팩스번호	<b>82-2-2123-8618</b>
II-10	이메일 주소	<b>patent@yonsei.ac.kr</b>
II-11	출원인 코드	<b>2-2005-009509-9</b>

## PCT 출원서

(전자적 형태가 원본)

<b>III-1</b>	출원인 및/또는 발명자	
III-1-1	이 사람은	오직 발명자 (inventor only)
III-1-3	우측 지정국에 관한 발명자	모든 지정국 (all designated States)
III-1-4ko	성명	박상현
III-1-4en	Name (LAST, First):	<b>PARK, Sang Hyun</b>
III-1-5ko	주소	대한민국 <b>08004</b> 서울시 양천구 오목로 300, 204동 3701호
III-1-5en	Address:	<b>204-3701, 300, Omok-ro, Yangcheon-gu, Seoul 08004 Republic of Korea</b>
<b>III-2</b>	출원인 및/또는 발명자	
III-2-1	이 사람은	오직 발명자 (inventor only)
III-2-3	우측 지정국에 관한 발명자	모든 지정국 (all designated States)
III-2-4ko	성명	이지환
III-2-4en	Name (LAST, First):	<b>LEE, Ji Hwan</b>
III-2-5ko	주소	대한민국 <b>03312</b> 서울시 은평구 통일로 1045, 102동 504호
III-2-5en	Address:	<b>102-504, 1045, Tongil-ro, Eunpyeong-gu, Seoul 03312 Republic of Korea</b>
<b>III-3</b>	출원인 및/또는 발명자	
III-3-1	이 사람은	오직 발명자 (inventor only)
III-3-3	우측 지정국에 관한 발명자	모든 지정국 (all designated States)
III-3-4ko	성명	최원기
III-3-4en	Name (LAST, First):	<b>CHOI, Won Gi</b>
III-3-5ko	주소	대한민국 <b>03724</b> 서울시 서대문구 연희로16길 17, 202호
III-3-5en	Address:	<b>#202, 17, Yeonhui-ro 16-gil, Seodaemun-gu, Seoul 03724 Republic of Korea</b>
<b>III-4</b>	출원인 및/또는 발명자	
III-4-1	이 사람은	오직 발명자 (inventor only)
III-4-3	우측 지정국에 관한 발명자	모든 지정국 (all designated States)
III-4-4ko	성명	성한승
III-4-4en	Name (LAST, First):	<b>SUNG, Han Seung</b>
III-4-5ko	주소	대한민국 <b>06116</b> 서울시 강남구 봉은사로21길 45, 3층
III-4-5en	Address:	<b>3Fl., 45, Bongeunsa-ro 21-gil, Gangnam-gu, Seoul 06116 Republic of Korea</b>

## PCT 출원서

(전자적 형태가 원본)

<b>III-5</b>	출원인 및/또는 발명자	
III-5-1	이 사람은	오직 발명자 (inventor only)
III-5-3	우측 지정국에 관한 발명자	모든 지정국 (all designated States)
III-5-4ko	성명	김도영
III-5-4en	Name (LAST, First):	<b>KIM, Do Young</b>
III-5-5ko	주소	대한민국 <b>07229</b> 서울시 영등포구 국회대로55길 28, 704호
III-5-5en	Address:	<b>#704, 28, Gukhoe-daero 55-gil, Yeongdeungpo-gu, Seoul 07229 Republic of Korea</b>
<b>IV-1</b>	대리인 또는 대표자	대리인
	아래에 기재된 자는 관할 국제기관에 대하여 우측에 표시된 자격으로 출원인을 대리하는 것으로 선임되었다.	
IV-1-1ko	성명	특허법인 우인
IV-1-1en	Name:	<b>WOIN PATENT &amp; LAW FIRM</b>
IV-1-2ko	주소	대한민국 <b>06246</b> 서울시 강남구 역삼로 157, 2층 (역삼동, 중평빌딩)
IV-1-2en	Address:	<b>(Yeoksam-dong, Jungpyeong Bldg.) 2Fl., 157 Yeoksamro, Gangnam-gu, Seoul 06246 Republic of Korea</b>
IV-1-3	전화번호	<b>82-2-541-9841</b>
IV-1-4	팩스번호	<b>82-2-541-9842</b>
IV-1-5	이메일 주소	<b>patent@woinlaw.com</b>
IV-1-5(a)	이메일 사용동의 수리관청, 국제조사기관, 국제사무국, 국제예비심사기관이 필요 시 이 이메일 주소를 사용하여 이 국제출원과 관련하여 발행된 통지서를 송부할 것에 동의한다.	서면 통지서에 앞서 선람용 사본 송부
IV-1-6	대리인 코드	<b>9-2006-100082-1</b>
<b>V</b>	지정국	
<b>V-1</b>	본 출원서의 제출로, 규칙 4.9(a)에 따라, 부여될 수 있는 모든 종류의 권리 보호를 위하여, 그리고 해당하는 경우 지역특허 및 국내특허 모두를 위하여 당해 국제출원일에 PCT에 기속되는 모든 계약국이 지정된다.	
<b>V-2</b>	<b>V-2</b> 란은 출원서 제출시 또는 규칙 26의 2.1에 의해 그 이후 출원서 제6기재란에 위 특정 관련 계약국의 국내 선출원에 대한 우선권주장이 포함되어 있을 경우 당해 계약국의 국내법에 의해 해당 국내 선출원의 효력이 상실되는 것을 방지하기 위한 목적으로 당해 계약국의 지정을 제외하는 데에만 사용될 수 있다 (지정 제외시 이의 취소 불가능).	<b>KR</b>
<b>VI-1</b>	선국내출원에 대한 우선권 주장	
VI-1-1	출원일	<b>2019년 11월 01일 (01.11.2019)</b>
VI-1-2	출원번호	<b>10-2019-0138684</b>
VI-1-3	파리협약 당사국명 또는 WTO 회원국명	<b>KR</b>

## PCT 출원서

(전자적 형태가 원본)

<b>VI-2</b>	우선권서류 신청 수리관청에 대하여 위에 명시된 선출원의 인증등본을 준비하여 국제사무국에 송부하여 줄 것을 신청한다.	<b>VI-1</b>	
<b>VI-3</b>	인용에 의한 보완 조약 제11조(1)(iii)(d) 또는 (e)에서 규정하는 국제출원의 요소, 또는 규칙 20.5(a)에서 규정하는 명세서, 청구 범위 또는 도면의 일부가 본 국제출원에는 포함되어 있지 않지만 조약 제11조(1)(iii) 규정의 요소 중 하나 이상이 수리관청에 최초로 접수된 날에 우선권주장의 기초가 된 선출원에 완전히 포함되어 있는 경우, 그 요소 또는 부분은 규칙 20.6 규정에 의한 확인을 조건으로, 규칙 20.6의 규정과 관련하여 본 국제출원에 있어서 인용에 의해 보완된다.		
<b>VII-1</b>	국제조사기관(ISA) 선택	<b>대한민국 특허청 (ISA/KR)</b>	
<b>VIII</b>	선언서	선언서 개수	
<b>VIII-1</b>	발명자의 신원에 관한 선언	-	
<b>VIII-2</b>	국제출원일에 특허출원 및 특허를 받을 수 있는 출원인의 자격에 관한 선언	-	
<b>VIII-3</b>	국제출원일에 선출원의 우선권을 주장할 수 있는 출원인의 자격에 관한 선언	-	
<b>VIII-4</b>	발명자 선언(미국에 대한 지정의 경우에 한함)	-	
<b>VIII-5</b>	신규성을 해치지 아니하는 개시 또는 신규성 상실의 예외에 관한 선언	-	
<b>IX</b>	체크 리스트	용지 수	전자적 파일 첨부
<b>IX-1</b>	출원서(선언서 포함)	<b>5</b>	✓
<b>IX-2</b>	명세서	<b>13</b>	✓
<b>IX-3</b>	청구범위	<b>1</b>	✓
<b>IX-4</b>	요약서	<b>1</b>	✓
<b>IX-5</b>	도면	<b>12</b>	✓
<b>IX-7</b>	용지매수 소개	<b>32</b>	
	첨부 항목	서면 첨부	전자적 파일 첨부
<b>IX-8</b>	수수료 계산 용지	-	✓
<b>IX-9</b>	개별위임장 원본	-	✓
<b>IX-20</b>	요약서에 수반되어야 할 도면 번호	<b>3</b>	
<b>IX-21</b>	국제출원의 출원 언어	<b>한국어</b>	

PCT 출원서

(전자적 형태가 원본)

X-1	출원인, 대리인 또는 대표자의 서명 또는 날인	
X-1-1	성명	특허법인 우인
X-1-2	서명인의 성명	특허법인 우인
X-1-3	권한 (출원서를 통해 서명자의 자격이 명백하지 않은 경우에는 그 자격도 표시)	특허법인 우인

수리관청 전용

10-1	국제출원으로 제출된 서류의 실제 접수일	
10-2	도면	
10-2-1	접수	
10-2-2	미접수	
10-3	국제출원으로 제출된 서류를 완성하는 서류 또는 도면의 추후 기간내 제출에 따른 정정된 실제 접수일	
10-4	PCT 제11조(2)에 따라 제출이 요구된 보완서로서 기간내 제출된 보완서의 접수일	
10-5	국제조사기관(ISA)	ISA/KR
10-6	조사료 납부시까지 지연된 조사용 사본의 송부	

국제 사무국 전용

11-1	국제 사무국의 기록원본 접수일	
------	------------------	--

PCT 위임장

(전자적 형태가 원본)

0-1	PCT 위임장 (특허 협력 조약에 의거하여 제출된 국제 출원) (PCT 규칙 제90.4조)	
0-1-1	우측에 기재된 바와 같이 작성되었다.	<b>PCT-SAFE</b> <b>버전 3.51.087.263 MT/FOP 20190401/0.20.5.24</b>
1	아래에 서명한 출원인	<b>연세대학교 산학협력단</b>
1-1-1	우측에 기재된 사람을 아래의 자격으로 선임한다.	<b>특허법인 우인</b> <b>WOOIN PATENT &amp; LAW FIRM</b>  대한민국 <b>06246</b> 서울시 강남구 역삼로 157, 2층 (역삼동,중평빌딩)  (Yeoksam-dong, Jungpyeong Bldg.) 2Fl., 157 Yeoksamro, Gangnam-gu, Seoul 06246 Republic of Korea
1-2	자격	대리인
1-3	우측 기관에 대하여	모든 관할 국제 기관
1-4	아래의 국제 출원에 관한 서명의 출원인을 대리함	
1-4-1	발명의 명칭	비휘발성 메모리를 이용한 로그 구조 병합 트리 기반의 데이터 베이스의 데이터 처리 방법
1-4-2	출원인 또는 대리인의 서류참조기호	<b>PCT19-0023</b>
1-4-3	국제출원번호(이용 가능한 경우)	
1-4-4	수리관청	<b>대한민국 특허청 (RO/KR)</b>
1-5	그리고 아래 서명인을 대신하여 지불하거나 지불받았다.	

PCT 위임장

(전자적 형태가 원본)

2-1	출원인 서명	
2-1-1	성명	연세대학교 산학협력단
2-1-2	서명인의 성명	연세대학교 산학협력단
2-1-3	권한 (출원서를 통해 서명자의 자격이 명백하지 않은 경우에는 그 자격도 표시)	연세대학교 산학협력단
3	일자	2019년 11월 15일 (15.11.2019)

PCT(부속문서 - 수수료 계산용지)

(전자적 형태가 원본)  
이 페이지는 국제 출원서의 일부가 아니며 페이지수에 포함되지 않는다

0	수리관청 전용				
0-1	국제출원번호				
0-2	수리관청의 우편 소인 일자				
0-4	Form PCT/RO/101 (부속문서) PCT 수수료 계산 용지		PCT-SAFE 버전 3.51.087.263 MT/FOP 20190401/0.20.5.24		
0-4-1	우측에 기재된 바와 같이 작성되었다.				
0-9	출원인 또는 대리인의 서류참조기호		PCT19-0023		
2	출원인		연세대학교 산학협력단		
12	규정 수수료 계산		수수료 금액/개수	총 금액 (CHF)	총 금액 (KRW)
12-1	송달료	T	⇔		45000
12-2-1	조사료	S	⇔		450000
12-2-2	국제조사기관		KR		
12-3	국제 출원 수수료 최초 30장	i1	1330 CHF		
12-4	최초 30장 초과 장수		2		
12-5	최초 30장 초과 1장당 추가 수수료	(X)	15 CHF		
12-6	총 추가금액	i2	30 CHF		
12-7	i1 + i2 =	i	1360 CHF		
12-12	XML 전자출원 감면	R	CHF-300		
12-13	총 국제출원 수수료(i-R)	I	⇔		
12-14	우선권 서류에 대한 수수료 우선권 서류를 요청한 개수		1		
12-15	문서별 수수료	(X)	0 KRW		
12-16	총 우선권 서류 수수료	P	⇔		
12-17	우선권 주장 회복에 대한 수수료 우선권 주장 회복에 대한 요청 개수	RP	0		
	우선권 회복에 대한 수수료 총 금액				
12-19	총 금액 (T+S+I+P+RP)		⇔	1060	495000
12-21	결제 방법		현금		



## PCT

13-2-2	검증 메시지 지정국	녹색입니까? 다음의 국가가 지정되지 않았음을 알립니다: <b>KR</b>
13-2-3	검증 메시지 이름	녹색입니까? 출원인 1: 전자메일 사용승인 항목이 체크되지 않았습니다. 이 전자메일 주소는 전화로 할 수 있는 종류의 연락을 대신하기 위해서만 사용됩니다.
13-2-7	검증 메시지 내용	노란색 처리된(rendering) 출원서 본체에 300dpi 이상의 이미지들이 포함되어 있습니다. 국제공개시 품질저하를 피하기 위해서는 좀 더 낮은 해상도로 이미지들을 저장하시기 바랍니다.
13-2-8	검증 메시지 수수료	녹색입니까? 가장 최근의 수수료표가 사용되었는지 확인하십시오.

## 명세서

### 발명의 명칭: 비휘발성 메모리를 이용한 로그 구조 병합 트리 기반의 데이터 베이스의 데이터 처리 방법

#### 기술분야

- [1] 본 발명이 속하는 기술 분야는 비휘발성 메모리를 이용한 로그 구조 병합 트리 기반의 데이터 베이스 및 그 데이터 처리 방법에 관한 것이다. 본 연구는 2019년도 과학기술정보통신부(정부)의 재원으로 정보통신기술진흥센터의 지원을 받아 수행된 연구사업인 SW스타랩 IoT 환경을 위한 고성능 플래시 메모리 스토리지 기반 인메모리 분산 DBMS 연구개발(No. 1711080997)과 관련된다.

#### 배경기술

- [2] 이 부분에 기술된 내용은 단순히 본 실시예에 대한 배경 정보를 제공할 뿐 종래기술을 구성하는 것은 아니다.
- [3] 키-값 기반의 데이터 베이스는 센서 데이터, 소셜 네트워크 데이터 등과 같이 비정형 데이터를 다루는데 유용하다. 키-값 기반의 데이터 베이스는 로그 구조 병합 트리(Log Structured Merge Tree)를 주로 사용한다.
- [4] 로그 구조 병합 트리(Log Structured Merge Tree, LSM-Tree)는 연속적인 쓰기 연산을 수행하는 워크로드를 위해 설계되었다. LSM-Tree 구조는 하나의 인메모리 데이터 구조와 여러 개의 블록(ex. 디스크 등)에 저장을 위한 이어쓰기(Append) 방식의 데이터 구조로 이루어져 있다.
- [5] LSM-Tree는 키-값 데이터 베이스에서 빈번히 발생하는 삽입 및 수정을 효율적으로 수행한다. 데이터를 우선 로그 형식으로 저장하고, 로그 상의 데이터 정렬, 수정 작업의 처리 등의 병합을 이루는 쓰기 친숙형 구조(Write Friendly Structure)이다. 하지만 나중에 발생하는 병합 동작은 쓰기 증폭을 발생시키며 시스템 성능과 저장장치의 수명에 영향을 준다.
- [6] LSM-Tree는 임의적인 순서로 데이터를 쓰지 않고 순차적으로 데이터를 쓴다. 데이터를 조회할 때 주어진 데이터가 트리 내의 어느 위치에 있는지 알 수 없어서 데이터를 찾기 위해서 상위 레벨부터 순차적으로 검색해야 한다. 디스크에 데이터가 없어도 모든 레벨의 모든 파일을 읽어야 한다.
- [7] (특허문헌1) 미국공개특허공보 US 2017-0344619 (2017.11.30.)
- [8] (특허문헌2) 한국공개특허공보 KR 10-2016-0121819 (2016.10.21.)
- [9] (특허문헌3) 미국공개특허공보 US 2018-0121121 (2018.05.03.)

#### 발명의 상세한 설명

#### 기술적 과제

- [10] 본 발명의 실시예들은 휘발성 메모리 및 비휘발성 메모리를 포함하는 데이터 베이스가 휘발성 메모리의 일정 용량을 초과한 데이터에 관하여 비휘발성

메모리에 저장하고, 비휘발성 메모리의 리스트 구조 및 영속성 버퍼를 통해 플러시 동작 및 컴팩션 동작을 수행함으로써, 데이터 영속성을 유지하면서 쓰기 지연과 읽기 지연을 최소화하는 데 발명의 주된 목적이 있다.

- [11] 본 발명의 명시되지 않은 또 다른 목적들은 하기의 상세한 설명 및 그 효과로부터 용이하게 추론할 수 있는 범위 내에서 추가적으로 고려될 수 있다.

### 과제 해결 수단

- [12] 본 실시예의 일 측면에 의하면, 데이터 베이스의 데이터 처리 방법에 있어서, 상기 데이터 베이스의 휘발성 메모리에 데이터를 저장하는 단계, 및 상기 데이터 베이스의 비휘발성 메모리에 복수의 노드가 연결된 리스트 구조를 생성하고 상기 데이터를 상기 리스트 구조에 저장하는 방식으로 플러시 동작을 수행하는 단계를 포함하는 데이터 베이스의 데이터 처리 방법을 제공한다.
- [13] 상기 데이터 베이스는 키-값 형식으로 데이터를 저장하고, 상기 리스트 구조는 다수의 다음 포인터를 갖는 스킵 리스트일 수 있다.
- [14] 상기 데이터 베이스가 데이터 쓰기를 수행하지 않고, 새로운 리스트 구조를 생성하고 상기 새로운 리스트 구조가 기존 리스트 구조의 노드에 할당된 키-값을 포인팅하는 방식으로 컴팩션 동작을 수행하는 단계를 포함할 수 있다.
- [15] 상기 플러시 동작을 수행하는 단계는, 상기 비휘발성 메모리의 영속성 버퍼(Persistent Buffer)에 키-값을 순차적으로 복사하고, 상기 영속성 버퍼는 상기 노드에 할당된 키-값의 랜덤 접근을 방지하며, 상기 리스트 구조는 상기 리스트 구조에 대응하는 영속성 버퍼의 오프셋을 포인팅할 수 있다.
- [16] 상기 컴팩션 동작을 수행하는 단계는, 상기 비휘발성 메모리에 새로운 리스트 구조를 생성하고 상기 새로운 리스트 구조가 이전 리스트 구조에 대응하는 영속성 버퍼의 오프셋을 포인팅할 수 있다.
- [17] 상기 데이터 베이스의 비휘발성 메모리에 저장된 데이터가 기 설정된 용량 범위를 초과하면, 단계화 정책(Tiering Policy)에 따라 상기 데이터 베이스의 블록 드라이브로 방출(Eviction)할 수 있다.
- [18] 상기 단계화 정책은, (i) 특정 레벨에 있는 데이터를 선택하여 상기 블록 드라이브에 저장하는 제1 단계화 정책, (ii) 데이터 접근이 오래된 데이터를 선택하여 상기 블록 드라이브에 저장하는 제2 단계화 정책, (iii) 모든 데이터를 상기 비휘발성 메모리에 저장하는 제3 단계화 정책, 또는 이들의 조합으로 설정될 수 있다.
- [19] 상기 데이터 베이스에서 시스템 오류가 발생하면, 상기 데이터 베이스의 비휘발성 메모리에 저장된 상기 리스트 구조가 포인팅하는 데이터를 조회한 결과를 통해 데이터를 순차적으로 복구하는 단계를 포함할 수 있다.
- [20] 본 실시예의 다른 측면에 의하면, 프로세서, 휘발성 메모리, 및 비휘발성 메모리를 포함하는 데이터 베이스에 있어서, 상기 휘발성 메모리에 데이터를 저장하고, 상기 비휘발성 메모리에 복수의 노드가 연결된 리스트 구조를

생성하고 상기 데이터를 상기 리스트 구조에 저장하는 방식으로 플러시 동작을 수행하는 것을 특징으로 하는 데이터 베이스를 제공한다.

- [21] 본 실시예의 또 다른 측면에 의하면, 프로세서에 의해 실행 가능한 컴퓨터 프로그램 명령어들을 포함하는 비일시적(Non-Transitory) 컴퓨터 판독 가능한 매체에 기록되어 데이터 처리를 위한 컴퓨터 프로그램으로서, 상기 컴퓨터 프로그램 명령어들이 데이터 베이스의 적어도 하나의 프로세서에 의해 실행되는 경우에, 상기 데이터 베이스의 휘발성 메모리에 데이터를 저장하는 단계, 및 상기 데이터 베이스의 비휘발성 메모리에 복수의 노드가 연결된 리스트 구조를 생성하고 상기 데이터를 상기 리스트 구조에 저장하는 방식으로 플러시 동작을 수행하는 단계를 포함한 동작들을 수행하는 컴퓨터 프로그램을 제공한다.

### 발명의 효과

- [22] 이상에서 설명한 바와 같이 본 발명의 실시예들에 의하면, 휘발성 메모리 및 비휘발성 메모리를 포함하는 데이터 베이스가 휘발성 메모리의 일정 용량을 초과한 데이터에 관하여 비휘발성 메모리에 저장하고, 비휘발성 메모리의 리스트 구조 및 영속성 버퍼를 통해 플러시 동작 및 컴팩션 동작을 수행함으로써, 데이터 영속성을 유지하면서 쓰기 지연과 읽기 지연을 최소화하는 효과가 있다.
- [23] 여기에서 명시적으로 언급되지 않은 효과라 하더라도, 본 발명의 기술적 특징에 의해 기대되는 이하의 명세서에서 기재된 효과 및 그 잠정적인 효과는 본 발명의 명세서에 기재된 것과 같이 취급된다.

### 도면의 간단한 설명

- [24] 도 1은 기존의 로그 구조 병합 트리 기반의 데이터 베이스를 예시한 도면이다.
- [25] 도 2는 기존의 로그 구조 병합 트리 기반의 데이터 베이스가 데이터 컴팩션 동작을 수행하는 것을 예시한 도면이다.
- [26] 도 3은 본 발명의 일 실시예에 따른 데이터 베이스를 예시한 블록도이다.
- [27] 도 4는 본 발명의 일 실시예에 따른 데이터 베이스의 내부 데이터 구조를 예시한 도면이다.
- [28] 도 5는 본 발명의 다른 실시예에 따른 데이터 베이스의 데이터 처리 방법을 예시한 흐름도이다.
- [29] 도 6은 본 발명의 다른 실시예에 따른 데이터 베이스가 비휘발성 메모리에 생성한 스kip 리스트를 예시한 도면이다.
- [30] 도 7은 본 발명의 다른 실시예에 따른 데이터 베이스가 스kip 리스트에 대해 컴팩션을 수행한 것을 예시한 도면이다.
- [31] 도 8은 본 발명의 다른 실시예에 따른 데이터 베이스가 영속성 버퍼를 통해 순차적 복사를 수행한 것을 예시한 도면이다.
- [32] 도 9는 본 발명의 다른 실시예에 따른 데이터 베이스가 영속성 버퍼를 통해

바이트 어드레싱 컴팩션을 수행한 것을 예시한 도면이다.

- [33] 도 10 내지 도 12는 본 발명의 실시예들에 따라 수행된 모의실험 결과를 도시한 것이다.

### 발명의 실시를 위한 형태

- [34] 이하, 본 발명을 설명함에 있어서 관련된 공지기능에 대하여 이 분야의 기술자에게 자명한 사항으로서 본 발명의 요지를 불필요하게 흐릴 수 있다고 판단되는 경우에는 그 상세한 설명을 생략하고, 본 발명의 일부 실시예들을 예시적인 도면을 통해 상세하게 설명한다.
- [35] 도 1은 기존의 로그 구조 병합 트리(LSM-Tree) 기반의 데이터 베이스를 예시한 도면이고, 도 2는 기존의 로그 구조 병합 트리 기반의 데이터 베이스가 데이터 컴팩션 동작을 수행하는 것을 예시한 도면이다.
- [36] LSM-Tree를 이용한 대표적인 데이터 베이스로는 LevelDB와 RocksDB가 있다.
- [37] LSM-Tree는 삽입 연산이 수행되면 먼저 메모리 영역에 데이터를 저장한다. 메모리의 일정 용량까지 데이터가 쌓이면 메모리의 내용을 디스크로 플러시(Flush)를 수행한다. 플러시되는 데이터는 디스크에 저장되어 있던 기존 데이터와 병합 정렬을 하여 기록된다. 디스크 영역의 각 레벨이 임계치를 넘으면 병합 정렬을 실행하여 하위 레벨을 생성한다.
- [38] LSM-Tree 기반의 데이터 베이스는 키-값 형태로 데이터를 저장한다. LSM-Tree 기반의 데이터 베이스에 데이터의 삽입 연산 요청이 들어오면 데이터를 메모리에 기록하기 전에 우선적으로 로그 파일에 로그를 기록한다. 로그를 기록한 다음 메모리 영역에 있는 맴테이블(Memtable)에 데이터를 저장한다. 쓰기 요청이 계속되어 맴테이블(Memtable)에 데이터가 일정 용량까지 기록되면, 맴테이블(Memtable)은 변경이 불가능한 불변 맴테이블(Immutable Memtable, Read-Only Memtable)로 변경된다. 불변 맴테이블이 가득 차게 되면 블록(디스크) 영역으로 플러시가 발생한다.
- [39] 플러시 동작을 수행하면, 맴테이블의 파일은 키 순서에 따라 정렬되어 SST(Stored String Table) 파일로 변경된다. SST 파일은 복수의 블록을 갖는다. 블록의 예시로는 데이터를 저장하는 데이터 블록(Data Block), 데이터 블록의 위치를 인덱싱하는 인덱스 블록(Index Block), 인덱스 블록의 위치를 처리하는 푸터 블록(Footer Block) 등이 있다.
- [40] SST 파일은 디스크 영역에서 컴팩션(Compaction)을 통해 업데이트된다. 한 번 생성된 SST 파일은 사라지지 않을 수 있다. 하위 레벨에 상주하는 SST 파일일수록 상위 레벨의 SST 파일보다 오래된 데이터가 위치할 수 있다.
- [41] 트랜잭션 수행 도중에 시스템 오류 또는 전원 차단 등과 같은 문제가 발생하면, 아직 디스크에 반영되지 않고 버퍼에 남아있는 데이터는 유실된다. 시스템이 재부팅된 후 데이터 베이스가 복구를 수행할 때, 트랜잭션이 어떤 갱신 연산을 수행했는지 기록하는 로그를 사용한다. 로그 기록 방식으로는

WAL(Write-Ahead-Logging) 규칙이 있다. WAL은 트랜잭션으로 인해 변경된 데이터가 디스크에 기록되기 전에 관련된 로그를 로그 파일에 기록하는 규칙이다.

- [42] LSM-Tree 기반의 데이터 베이스는 두 개의 명령어를 수행한다. 하나는 메모리에서 디스크로 넘어가는 플러시 명령어이고, 다른 하나는 디스크의 레벨들을 조정하는 컴팩션 명령어이다.
- [43] 플러시 명령어를 수행할 때, 불변 메모이블은 단일 SST 파일로 변경된다. 대량의 데이터가 한꺼번에 입력되면, 플러시 속도의 균형을 맞추고 SST 파일의 레벨의 용량 한계치를 유지하기 위해서 의도적으로 플러시 속도를 조절한다. 이러한 의도된 지연을 'Write Stall'이라고 한다. 표 1에 누적된 Write Stall이 예시되어 있다.
- [44] [표1]

Data size	Accumulated write stalls on SSD(us)
2 GB	61,902,809
4 GB	310,872,755
6 GB	569,200,925
8 GB	901,254,333

- [45] LSM-Tree의 각각의 레벨은 특성이 구분된다. 컴팩션 비용과 디스크 쓰기는 특정 레벨에 집중된다. 계층적 저장 구조는 상위 레벨에서 하위 레벨로 점진적인 데이터 누적을 야기한다. SST 파일의 수명은 해당하는 레벨에서 존재하는 동안 컴팩션을 수행하는 횟수를 의미한다. 컴팩션을 수행하는 동안 SST 파일이 특정 레벨에서 삭제되지 않으면, 해당 SST 파일의 수명은 높게 나타난다. 표 2에 SST 파일의 수명, 컴팩션 파일의 개수, 컴팩션 파일의 비율, 및 컴팩션 동안 쓰기량이 예시되어 있다.

[46] [표2]

Level	Average lifetime of SST files	The number of compacted files	The ratio of compacted files	The amount of write during compaction (MB)
L0	4.498	35434	99.989%	113230
L1	10.679	118038	99.986%	196361
L2	51.163	379120	99.961%	511559
L3	403.795	481843	99.728%	650031
L4	4340.665	339725	96.629%	599837

- [47] 컴팩션 명령어를 수행한 각 레벨의 결과를 살펴보면, 데이터 컴팩션을 수행하는 짧은 시간 동안에 상위 레벨(ex.  $L_0$  to  $L_3$ )의 SST 파일들은 생성되고 삭제됨을 나타낸다. 컴팩션 파일의 개수와 컴팩션 파일의 비율을 보면, LSM-Tree의 계층적 구조로 인하여 컴팩션 파일의 크기가 작지 않음을 나타낸다.
- [48] 컴팩션 동안 예비 자원을 수반하는 상위 레벨에서의 쓰기량은 중요하지 않아 보이지만, 디스크 I/O를 피할 수 없다. 상당히 많은 데이터를 입력했음에도, 하위 레벨  $L_4$ 의 컴팩션 파일의 개수와 쓰기량은  $L_3$ 보다 작은 값을 나타낸다.
- [49] 본 실시예에 따른 데이터 베이스는  $L_4$ 와 같은 하위 레벨에서 높은 값을 갖는 SST 파일의 수명에 집중해서, NVM을 통해 상위 레벨에서 낭비되는 디스크 I/O를 감소시킨다. NVM은 바이트 어드레싱을 이용하여 쓰기 증폭(Write Amplification)을 해결하고 연속적 컴팩션을 수행하게 한다.
- [50] 도 3은 본 발명의 일 실시예에 따른 데이터 베이스를 예시한 블록도이고, 도 4는 본 발명의 일 실시예에 따른 데이터 베이스의 내부 데이터 구조를 예시한 도면이다.
- [51] 도 3에 도시한 바와 같이, 데이터 베이스(10)는 프로세서(100), 휘발성 메모리(200), 및 비휘발성 메모리(300)를 포함한다. 데이터 베이스(10)는 도 3에서 예시적으로 도시한 다양한 구성요소들 중에서 일부 구성요소를 생략하거나 다른 구성요소를 추가로 포함할 수 있다. 예컨대, 데이터 베이스(10)는 단계화 정책에 따라 블록 디바이스(400)를 추가로 포함할 수 있다.
- [52] 데이터 베이스(10)는 데이터를 가공 및 저장하는 장치이다. 데이터 베이스(10)는 키-값 형식으로 데이터를 저장하고 읽을 수 있다. 키-값의 처리 명령은 (SET K, V), (DEL K, V) 등으로 정의될 수 있다.
- [53] 프로세서(100)는 휘발성 메모리(200), 비휘발성 메모리(300), 및 블록 디바이스(400)에 기 정의된 명령어를 전송하여, 각종 신호 및 데이터 흐름을 제어한다.
- [54] 휘발성 메모리(200)는 저장된 정보를 계속 유지하기 위하여 전원 공급이 필요한 메모리이다. 예컨대, 휘발성 메모리(300)로는 DRAM(Dynamic Random Access Memory) 등이 있다.
- [55] 비휘발성 메모리(300)는 전원이 공급되지 않아도 저장된 정보를 계속 유지하는

메모리이다. 비휘발성 메모리(300)는 리스트 구조(310)를 포함할 수 있다. 리스트 구조(310)는 헤드(Head)와 리어(Rear)를 갖고, 키의 주소와 값의 주소를 각각 갖는다. 리스트 구조(310)는 다수의 다음 포인터를 갖는 스킵 리스트로 구현될 수 있다. 스킵 리스트는 각 노드마다 키 길이, 키, 값 길이, 및 값을 갖는다. 비휘발성 메모리(300)는 영속성 버퍼(320)를 포함할 수 있다.

[56] 블록 디바이스(400)는 블록 단위로 임의 접근이 가능한 저장매체이다. 예컨대, 블록 디바이스(400)로는 HDD(Hard Disk Drive), SSD(Solid State Drive) 등이 있다.

[57] 데이터 베이스(10)는 휘발성 메모리(200)에 1차적으로 데이터를 저장한다. 저장된 데이터가 기 설정된 용량을 초과하면, 일부 데이터를 비휘발성 메모리(300)에 2차적으로 저장한다. 데이터 베이스(10)는 비휘발성 메모리(300)에 복수의 노드가 연결된 리스트 구조(310)를 생성하고 데이터를 리스트 구조(310)에 저장하는 방식으로 플러시 동작을 수행할 수 있다.

[58] 도 5는 본 발명의 다른 실시예에 따른 데이터 베이스의 데이터 처리 방법을 예시한 흐름도이다.

[59] 단계 S210에서 데이터 베이스는 휘발성 메모리에 데이터를 저장한다.

[60] 단계 S220에서 데이터 베이스는 비휘발성 메모리에 플러시(flush) 동작을 수행한다. 플러시는 제1 저장소에서 제2 저장소로 복사하는 동작이다. 예컨대, 휘발성 메모리에서 비휘발성 메모리로 데이터를 복사한다.

[61] 단계 S230에서 데이터 베이스는 컴팩션(Compaction) 동작을 수행한다. 컴팩션은 병합 과정으로 특정 레벨의 임계치까지 데이터가 차면 해당 레벨의 데이터를 하위 레벨로 내려주는 동작이다.

[62] 단계 S240에서 데이터 베이스는 정책에 따라 블록 드라이브로 데이터를 방출(Eviction)하는 동작을 수행한다.

[63] 단계 S250에서 데이터 베이스는 오류가 발생하면 데이터를 복구(Recovery)하는 동작을 수행한다.

[64] 도 6은 본 발명의 다른 실시예에 따른 데이터 베이스가 비휘발성 메모리에 생성한 스킵 리스트를 예시한 도면이다.

[65] 비휘발성 메모리(Non-Volatile Memory, NVM)는 중간 유연 레벨에 해당하며, 스킵 리스트는 기존의 블록 파일 형식(ex. SST 파일)을 대체할 수 있다.

[66] 비휘발성 메모리(Non-Volatile Memory, NVM)는 바이트 어드레싱(Byte Addressability)이 가능하다. 바이트 어드레싱이 가능한 비휘발성 메모리의 예시로는 STT-MRAM(Spin-Transfer Torque Magnetic Random Access Memory), PCM(Phase-Change Memory) 등이 있다.

[67] 데이터 베이스 시스템에서 NVM이 일관적으로 동작하려면, 데이터 베이스 시스템은 'cflush' 및 'mfence' 명령어를 사용해야 한다. 'cflush' 명령어는 메모리에 캐시 라인을 플러시하는 명령어이며, NVM에 데이터를 완전하게 저장하는 것을 보장한다. 'mfence' 명령어는 명령어들의 재순서 배정으로부터 프로세서를 보호하는 메모리 장벽 명령어이다.



- [68] PMDK(Persistent Memory Development Kit) API는 NVM를 이용하여 키-값 기반의 데이터 베이스를 제공한다.
- [69] 도 7은 본 발명의 다른 실시예에 따른 데이터 베이스가 스kip 리스트에 대해 컴팩션을 수행한 것을 예시한 도면이다.
- [70] 데이터 베이스가 데이터 쓰기를 수행하지 않고, 새로운 리스트 구조를 생성하고 새로운 리스트 구조가 기존 리스트 구조의 노드에 할당된 키-값을 포인팅하는 방식으로 컴팩션 동작을 수행한다.
- [71] 도 8은 본 발명의 다른 실시예에 따른 데이터 베이스가 영속성 버퍼를 통해 순차적 복사를 수행한 것을 예시한 도면이다.
- [72] 데이터 베이스의 플러시 동작에 관한 알고리즘은 표 3와 같다.
- [73] [표3]

---

**Input:**

memtableIter : key-value iterator from immutable memtable

**Output:**

fileMetadata : metadata including new output number, smallest/largest key, filesize

*/\* Step 1. Iteration about all data in memtable \*/*

1: buf ← [] *// append all KV-pair, then memcpy at a time*

2: pBuf ← create new persistent buffer

3: startOffset ← get new start offset from persistent buffer

4: pos ← 0 *// relative position from start\_offset*

5: pSkiplist ← create new persistent skip list with new ID

*/\* Step 2. Iteration about all data in memtable \*/*

6: **for** KV-pair **in** memtableIter **do**

7:   (key, value) ← **Encode**(KV-pair)

8:   entryLength ← key.length + value.length

9:   **Append**(buf ← key, value)

*/\* Insertion into skip list node through the relative position of pBuf \*/*

10:   pSkiplist.current ← **Node**(startOffset + pos, pSkiplist.prev)

11:   **Persist**(pSkiplist.current)

12:   pSkiplist.current.next ← pSkiplist.prev.next

13:   **Persist**(pSkiplist.current.next)

14:   pSkiplist.prev.next ← pSkiplist.current

15:   **Persist**(pSkiplist.prev.next)

16:   pSkiplist.current ← pSkiplist.current.next

17:   **Persist**(pSkiplist.current)

18:   **Update**(pos ← pos + entryLength)

19: **end**

*/\* Step 3. Finally, add null node for recovery \*/*

20: pSkiplist.current ← **NullNode**()

21: **Persist**(pSkiplist.current)

*/\* Step 4. For persistent buffer, bulk copy \*/*

22: **Copy**(pBuf ← buf **until** buf.size)

23: **Persist**(pBuf)

24: **Update**(fileMetadata)

---

- [74] 플러시 동작은, 비휘발성 메모리의 연속성 버퍼에 키-값을 순차적으로 복사하고, 연속성 버퍼는 노드에 할당된 키-값의 랜덤 접근을 방지한다. 리스트 구조는 리스트 구조에 대응하는 연속성 버퍼의 오프셋을 포인팅한다.
- [75] 도 9는 본 발명의 다른 실시예에 따른 데이터 베이스가 연속성 버퍼를 통해 바이트 어드레싱 컴팩션을 수행한 것을 예시한 도면이다.
- [76] 데이터 베이스의 바이트 어드레싱 컴팩션 동작에 관한 알고리즘은 표 4와 같다.
- [77] [표4]

---

**Input:**

**mergeIter** : merge and sorted iterater from skip list IDs to be compacted

**Output:**

**fileMetadataList**: list of metadata including new output number, smallest/largest key, filesize

```

1: for KV-pair in mergeIter do
    /* Step 1. Get key-value entry by buffer pointer */
2:   if not KV-pair.isBufferPointerValid() then
3:     | break // end of iteration or invalid status
4:   else
5:     | bufPtr ← KV-pair.GetBufferPointer()
    /* Step 2. Check whether the persistent skip list is valid */
6:   if pSkiplist is Null then
7:     | pSkiplist ← CreatePersistentSkiplist(new ID)
8:     | Insert(pSkiplist.metadata ← inputIDs) // for recovery
    /* Step 3. Failure-atomic pointing operation with persistent buffer */
9:   pSkiplist.current ← Node(bufPtr, pSkiplist.prev)
10:  Persist(pSkiplist.current)
11:  pSkiplist.current.next ← pSkiplist.prev.next
12:  Persist(pSkiplist.current.next)
13:  pSkiplist.prev.next ← pSkiplist.current
14:  Persist(pSkiplist.prev.next)
15:  pSkiplist.current ← pSkiplist.current.next
16:  Persist(pSkiplist.current)
    /* Step 4. Check whether persistent skiplist is full */
17:  if pSkiplist.isFull() then
18:    | pSkiplist.current ← NullNode()
19:    | Persist(pSkiplist.current)
20:    | Insert(fileMetadataList ← new metadata)
21:    | pSkiplist ← Null
22: end

```

---

- [78] 컴팩션 동작은, 비휘발성 메모리에 새로운 리스트 구조를 생성하고 새로운 리스트 구조가 이전 리스트 구조에 대응하는 연속성 버퍼의 오프셋을 포인팅한다. 컴팩션 동작을 수행할 때, 연속성 버퍼는 키 길이, 키, 값 길이, 및

값에 대해 쓰기를 수행하지 않는다. NVM에 대해서 'cflush' 및 'mfence' 명령어를 사용하여 영속성을 확보한다. 영속성 버퍼는 SST 파일과 달리 인덱스 블록과 푸터 블록을 포함하지 않는다. 영속성 버퍼의 오프셋을 통해 데이터 검색 성능을 향상시킨다. 컴팩션을 수행하기 전에 스킵 리스트의 최하위 레벨에 대한 반복자들을 생성하고, 바이트 어드레싱 컴팩션 과정에서 반복자를 병합한다.

[79] 데이터 베이스는 저장소 단계화(Storage Tiering)를 수행한다. 데이터 베이스는 블록 드라이브를 포함한다. 데이터 베이스는 비휘발성 메모리에 저장된 데이터가 기 설정된 용량 범위를 초과하면, 단계화 정책(Tiering Policy)에 따라 데이터 베이스의 블록 드라이브로 방출(Eviction)한다.

[80] 단계화 정책은, (i) 특정 레벨에 있는 데이터를 선택하여 블록 드라이브에 저장하는 제1 단계화 정책(Leveled Tiering), (ii) 데이터 접근이 오래된 데이터를 선택하여 블록 드라이브에 저장하는 제2 단계화 정책(LRU Tiering), (iii) 모든 데이터를 비휘발성 메모리에 저장하는 제3 단계화 정책(No Tiering), 또는 이들의 조합으로 설정될 수 있다. 특정 레벨은 구현되는 설계에 따라 통계적인 방식으로 산출되어 설정될 수 있다.

[81] 데이터 베이스를 복구하는 동작에 관한 알고리즘은 표 5와 같다.

[82] [표5]

---

**Input:** None  
**Output:**  
fileMetadataList: list of metadata including new output number, smallest/largest key, filesize

```

/* Step 1. Get persistent pointer from NVM pool */
1: pools ← GetPoolsInDirectory()
2: persistentPtrs ← GetPersistentPointers(pools)
3: logFiles ← GetLogFilesInDirectory()
/* Step 2. Restore metadata of persistent skip list from persistent pointer */
4: for pptr in persistentPtrs do
5:   restoredSkiplist ← pptr.GetSkiplist()
6:   last ← restoredSkiplist.SeekToLastNode()
   /* Detect system failure occurs */
7:   if last.isNullNode() then
   /* Failure when flush data into skip list at LO */
8:     if restoredSkiplist.ID = logFiles.ID then
9:       Reset(restoredSkiplist)
10:      memtableIter ← logFiles.createIteratorAboutID()
11:      Flush(memtableIter)
   /* Failure when byte-addressable compaction at lower levels */
12:    else
13:      IDs ← restoredSkiplist.metadata.IDs // Get IDs to be compacted
14:      mergeIter ← Iterator(IDs)
15:      while mergeIter.key ≤ restoredSkiplist.current.key do
16:        | mergeIter ← mergeIter.next // pass already inserted data
17:      end
18:      ByteAddressableCompaction(mergeIter) // resume for remaining mergeliter
19:      metadata ← restoredSkiplist.generateMetadata()
20:      Insert(fileMetadata ← metadata)
21: end

```

---

- [83] 데이터 베이스에서 시스템 오류가 발생하면, 데이터 베이스의 비휘발성 메모리에 저장된 리스트 구조가 포인팅하는 데이터를 조회한 결과를 통해 데이터를 순차적으로 복구한다. NVM으로부터 영속성 포인터를 획득하고, 영속성 포인터를 이용하여 스킵 리스트를 복구한다. 스킵 리스트에서 노드를 찾고, 메타데이터(ID), ID에 대한 반복자를 획득하고, 이미 입력된 데이터는 생략하고, 남은 바이트 어드레싱 컴팩션을 고려한다. 스킵 리스트의 메타 데이터를 복구하고, 메타 데이터를 리스트에 삽입한다.
- [84] 도 10 내지 도 12는 본 발명의 실시예들에 따라 수행된 모의실험 결과를 도시한 것이다.
- [85] 도 10에 도시된 바와 같이, 본 실시예에 따른 데이터 베이스는 쓰기 지연과 읽기 지연 측면에서 성능이 향상됨을 알 수 있다.
- [86] 도 11을 참조하면, 제3 단계화 정책(No Tiering), 제2 단계화 정책(LRU Tiering),

제1 단계화 정책(Leveled Tiering) 순으로 쓰기 지연과 읽기 지연 측면에서 성능이 향상됨을 알 수 있다.

- [87] 도 12를 참조하면, 본 실시예에 따른 데이터 베이스(TLSM)는 쓰기 지연(Write Stall)과 컴팩션의 쓰기량 측면에서 성능이 향상됨을 알 수 있다.
- [88] 데이터 베이스에 포함된 구성요소들이 도 3에서는 분리되어 도시되어 있으나, 복수의 구성요소들은 상호 결합되어 적어도 하나의 모듈로 구현될 수 있다. 구성요소들은 장치 내부의 소프트웨어적인 모듈 또는 하드웨어적인 모듈을 연결하는 통신 경로에 연결되어 상호 간에 유기적으로 동작한다. 이러한 구성요소들은 하나 이상의 통신 버스 또는 신호선을 이용하여 통신한다.
- [89] 데이터 베이스는 하드웨어, 펌웨어, 소프트웨어 또는 이들의 조합에 의해 로직회로 내에서 구현될 수 있고, 범용 또는 특정 목적 컴퓨터를 이용하여 구현될 수도 있다. 장치는 고정배선형(Hardwired) 기기, 필드 프로그램 가능한 게이트 어레이(Field Programmable Gate Array, FPGA), 주문형 반도체(Application Specific Integrated Circuit, ASIC) 등을 이용하여 구현될 수 있다. 또한, 장치는 하나 이상의 프로세서 및 컨트롤러를 포함한 시스템온칩(System on Chip, SoC)으로 구현될 수 있다.
- [90] 데이터 베이스는 하드웨어적 요소가 마련된 컴퓨팅 디바이스에 소프트웨어, 하드웨어, 또는 이들의 조합하는 형태로 탑재될 수 있다. 컴퓨팅 디바이스는 각종 기기 또는 유무선 통신망과 통신을 수행하기 위한 통신 모듈 등의 통신장치, 프로그램을 실행하기 위한 데이터를 저장하는 메모리, 프로그램을 실행하여 연산 및 명령하기 위한 마이크로프로세서 등을 전부 또는 일부 포함한 다양한 장치를 의미할 수 있다.
- [91] 도 5에서는 각각의 과정을 순차적으로 실행하는 것으로 기재하고 있으나 이는 예시적으로 설명한 것에 불과하고, 이 분야의 기술자라면 본 발명의 실시예의 본질적인 특성에서 벗어나지 않는 범위에서 도 5에 기재된 순서를 변경하여 실행하거나 또는 하나 이상의 과정을 병렬적으로 실행하거나 다른 과정을 추가하는 것으로 다양하게 수정 및 변형하여 적용 가능할 것이다.
- [92] 본 실시예들에 따른 동작은 다양한 컴퓨터 수단을 통하여 수행될 수 있는 프로그램 명령 형태로 구현되어 컴퓨터 판독 가능한 매체에 기록될 수 있다. 컴퓨터 판독 가능한 매체는 실행을 위해 프로세서에 명령어를 제공하는 데 참여한 임의의 매체를 나타낸다. 컴퓨터 판독 가능한 매체는 프로그램 명령, 데이터 파일, 데이터 구조 또는 이들의 조합을 포함할 수 있다. 예를 들면, 자기 매체, 광기록 매체, 메모리 등이 있을 수 있다. 컴퓨터 프로그램은 네트워크로 연결된 컴퓨터 시스템 상에 분산되어 분산 방식으로 컴퓨터가 읽을 수 있는 코드가 저장되고 실행될 수도 있다. 본 실시예를 구현하기 위한 기능적인(Functional) 프로그램, 코드, 및 코드 세그먼트들은 본 실시예가 속하는 기술분야의 프로그래머들에 의해 용이하게 추론될 수 있을 것이다.
- [93] 본 실시예들은 본 실시예의 기술 사상을 설명하기 위한 것이고, 이러한

실시예에 의하여 본 실시예의 기술 사상의 범위가 한정되는 것은 아니다. 본 실시예의 보호 범위는 아래의 청구범위에 의하여 해석되어야 하며, 그와 동등한 범위 내에 있는 모든 기술 사상은 본 실시예의 권리범위에 포함되는 것으로 해석되어야 할 것이다.

## 청구범위

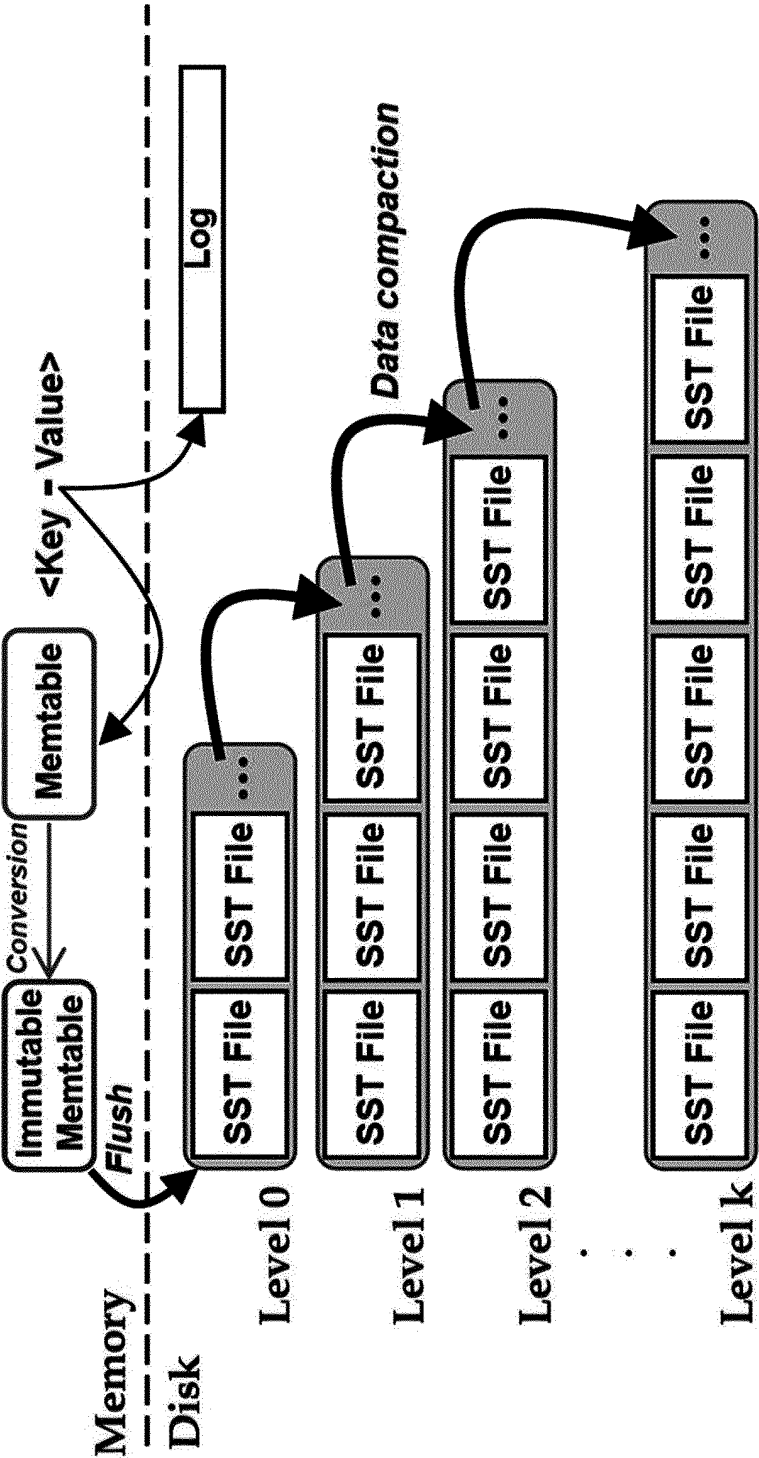
- [청구항 1] 프로세서, 휘발성 메모리, 및 비휘발성 메모리를 포함하는 데이터 베이스에 있어서,  
 상기 휘발성 메모리에 데이터를 저장하고,  
 상기 비휘발성 메모리에 복수의 노드가 연결된 리스트 구조를 생성하고  
 상기 데이터를 상기 리스트 구조에 저장하는 방식으로 플러시 동작을 수행하는 것을 특징으로 하는 데이터 베이스.
- [청구항 2] 제1항에 있어서,  
 상기 데이터 베이스는 키-값 형식으로 데이터를 저장하고,  
 상기 리스트 구조는 다수의 다음 포인터를 갖는 스킵 리스트이며,  
 상기 데이터 베이스가 데이터 쓰기를 수행하지 않고, 새로운 리스트 구조를 생성하고 상기 새로운 리스트 구조가 기존 리스트 구조의 노드에 할당된 키-값을 포인팅하는 방식으로 컴팩션 동작을 수행하는 것을 특징으로 하는 데이터 베이스.
- [청구항 3] 제1항에 있어서,  
 상기 플러시 동작은,  
 상기 비휘발성 메모리의 연속성 버퍼에 키-값을 순차적으로 복사하고,  
 상기 연속성 버퍼는 상기 노드에 할당된 키-값의 랜덤 접근을 방지하며,  
 상기 리스트 구조는 상기 리스트 구조에 대응하는 연속성 버퍼의 오프셋을 포인팅하며,  
 상기 컴팩션 동작은,  
 상기 비휘발성 메모리에 새로운 리스트 구조를 생성하고 상기 새로운 리스트 구조가 이전 리스트 구조에 대응하는 연속성 버퍼의 오프셋을 포인팅하는 것을 특징으로 하는 데이터 베이스.
- [청구항 4] 제1항에 있어서,  
 상기 데이터 베이스는 블록 드라이브를 포함하며,  
 상기 데이터 베이스는 상기 비휘발성 메모리에 저장된 데이터가 기 설정된 용량 범위를 초과하면, 단계화 정책(Tiering Policy)에 따라 상기 데이터 베이스의 블록 드라이브로 방출(Eviction)하며,  
 상기 단계화 정책은,  
 (i) 특정 레벨에 있는 데이터를 선택하여 상기 블록 드라이브에 저장하는 제1 단계화 정책, (ii) 데이터 접근이 오래된 데이터를 선택하여 상기 블록 드라이브에 저장하는 제2 단계화 정책, (iii) 모든 데이터를 상기 비휘발성 메모리에 저장하는 제3 단계화 정책, 또는 이들의 조합으로 설정되는 것을 특징으로 하는 데이터 베이스.

## 요약서

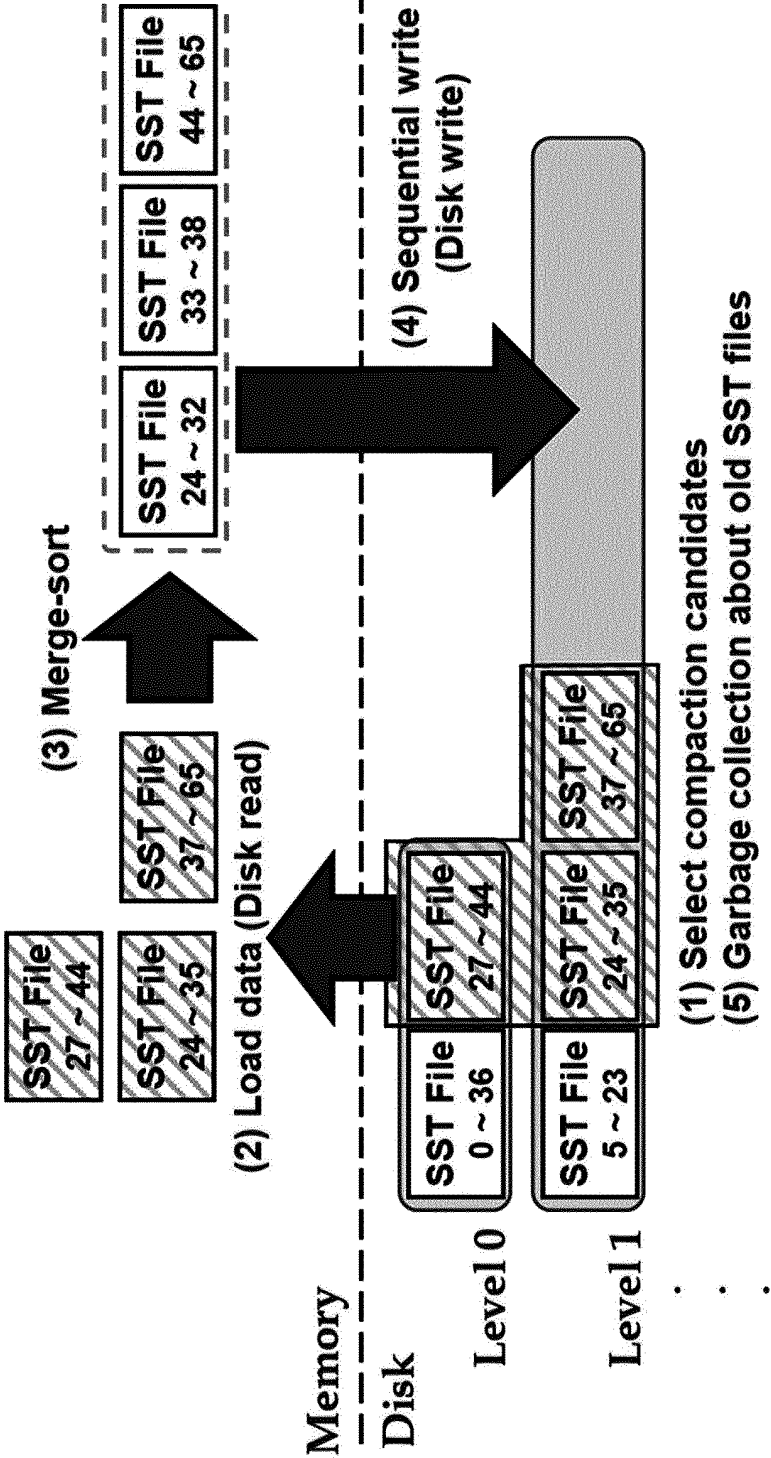
본 실시예들은 휘발성 메모리의 일정 용량을 초과한 데이터에 관하여 비휘발성 메모리에 저장하고, 비휘발성 메모리의 리스트 구조 및 영속성 버퍼를 통해 플러시 동작 및 컴팩션 동작을 수행함으로써, 데이터 영속성을 유지하면서 쓰기 지연과 읽기 지연을 최소화할 수 있는 데이터 베이스를 제공한다.



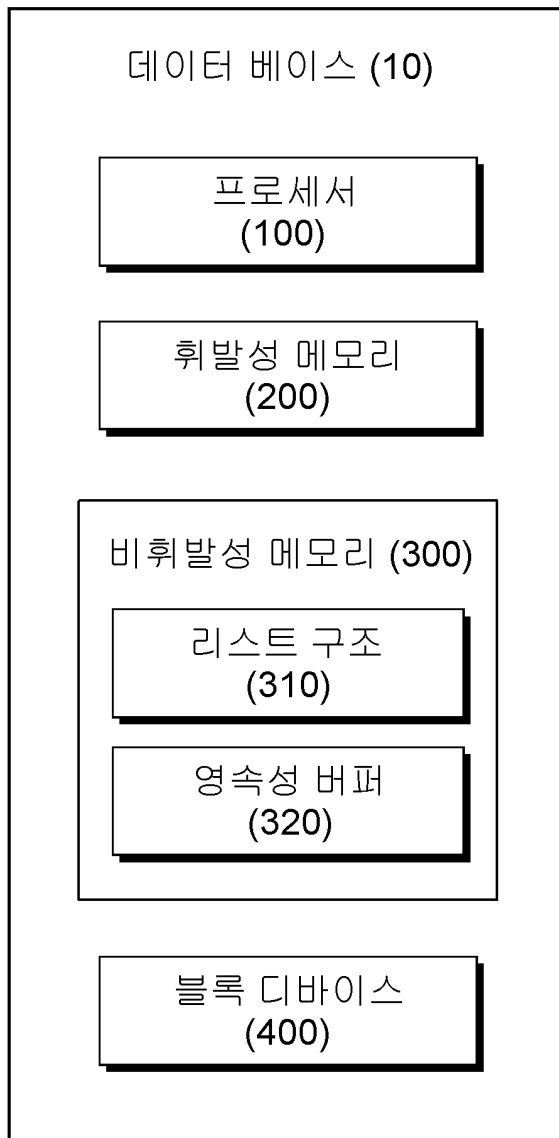
[도1]



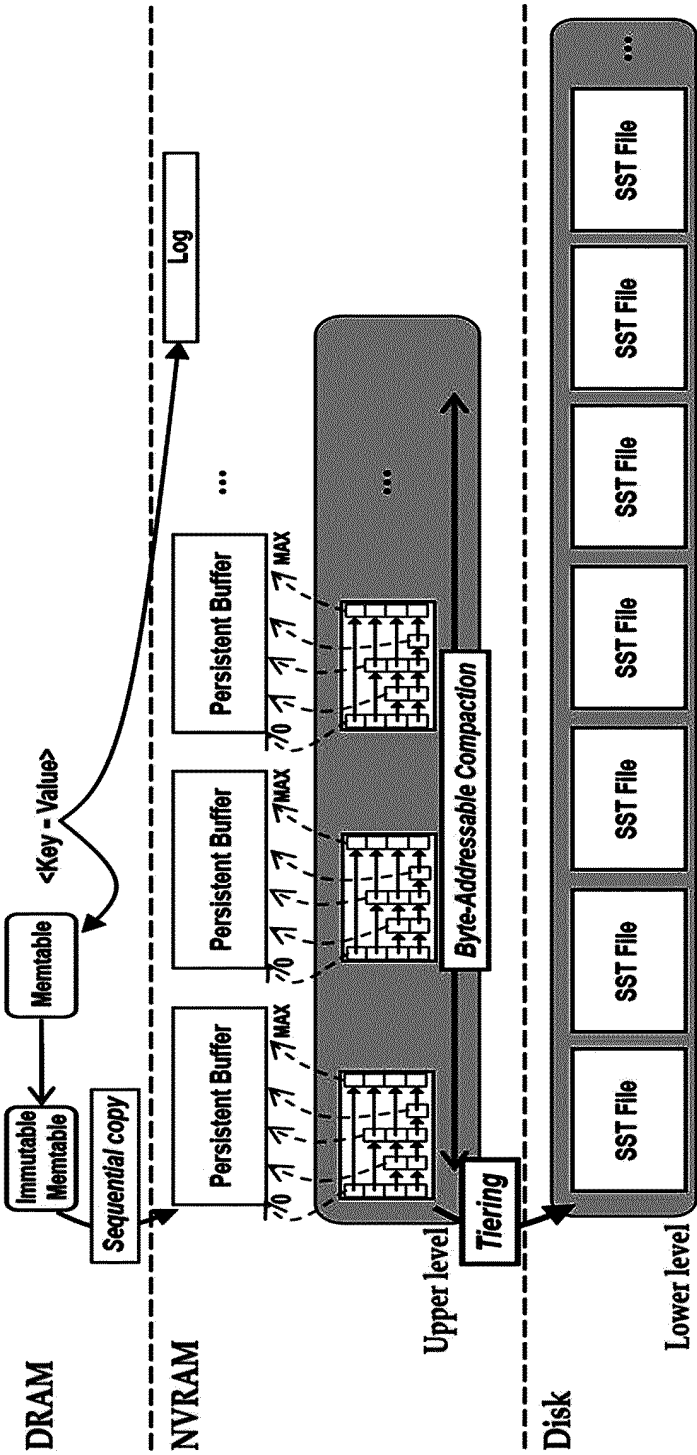
[도2]



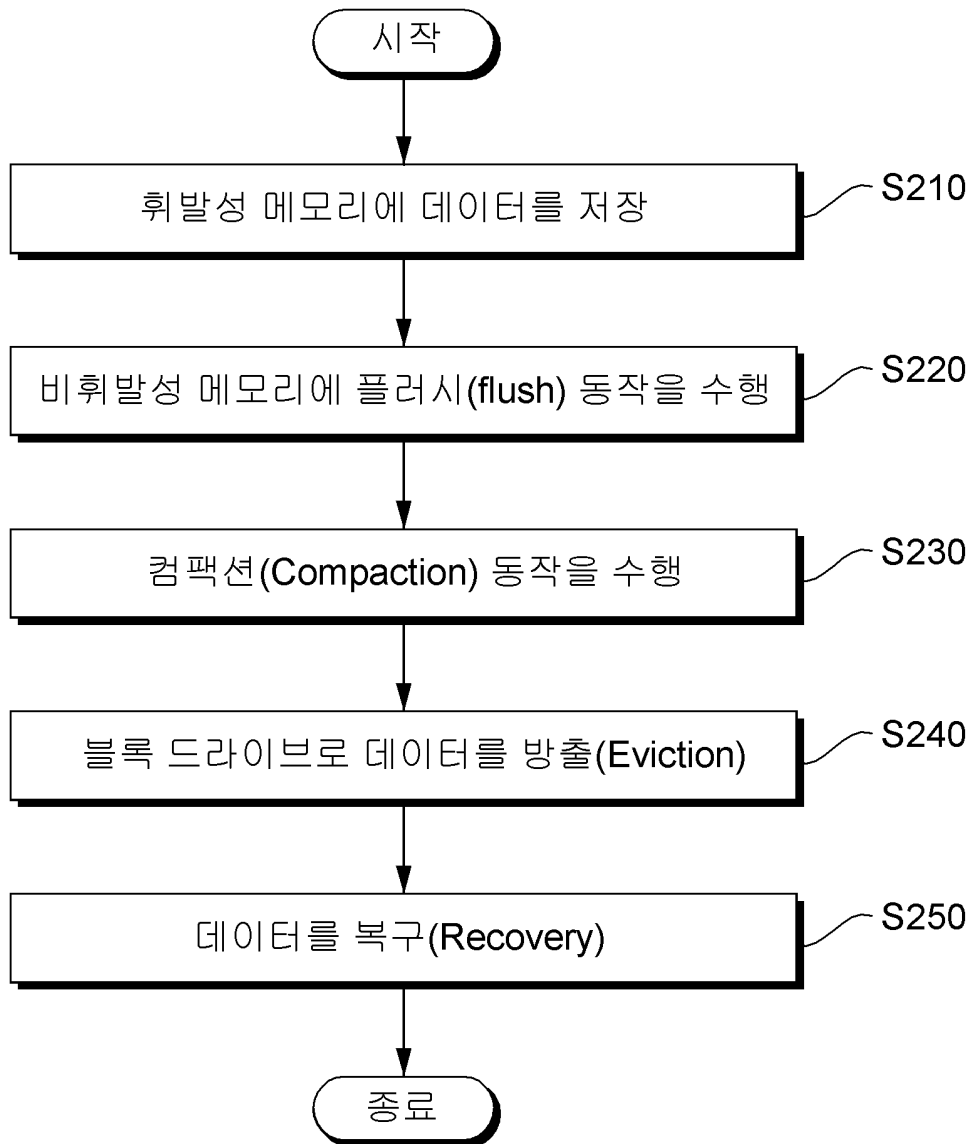
[도3]



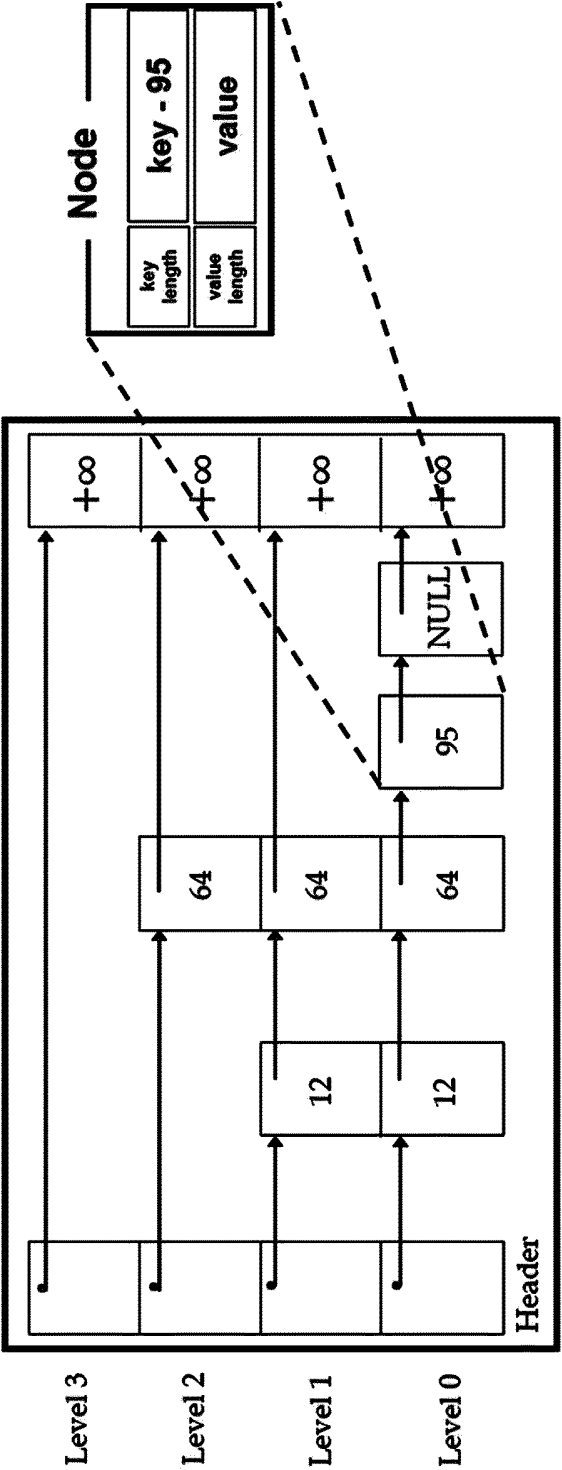
[도4]



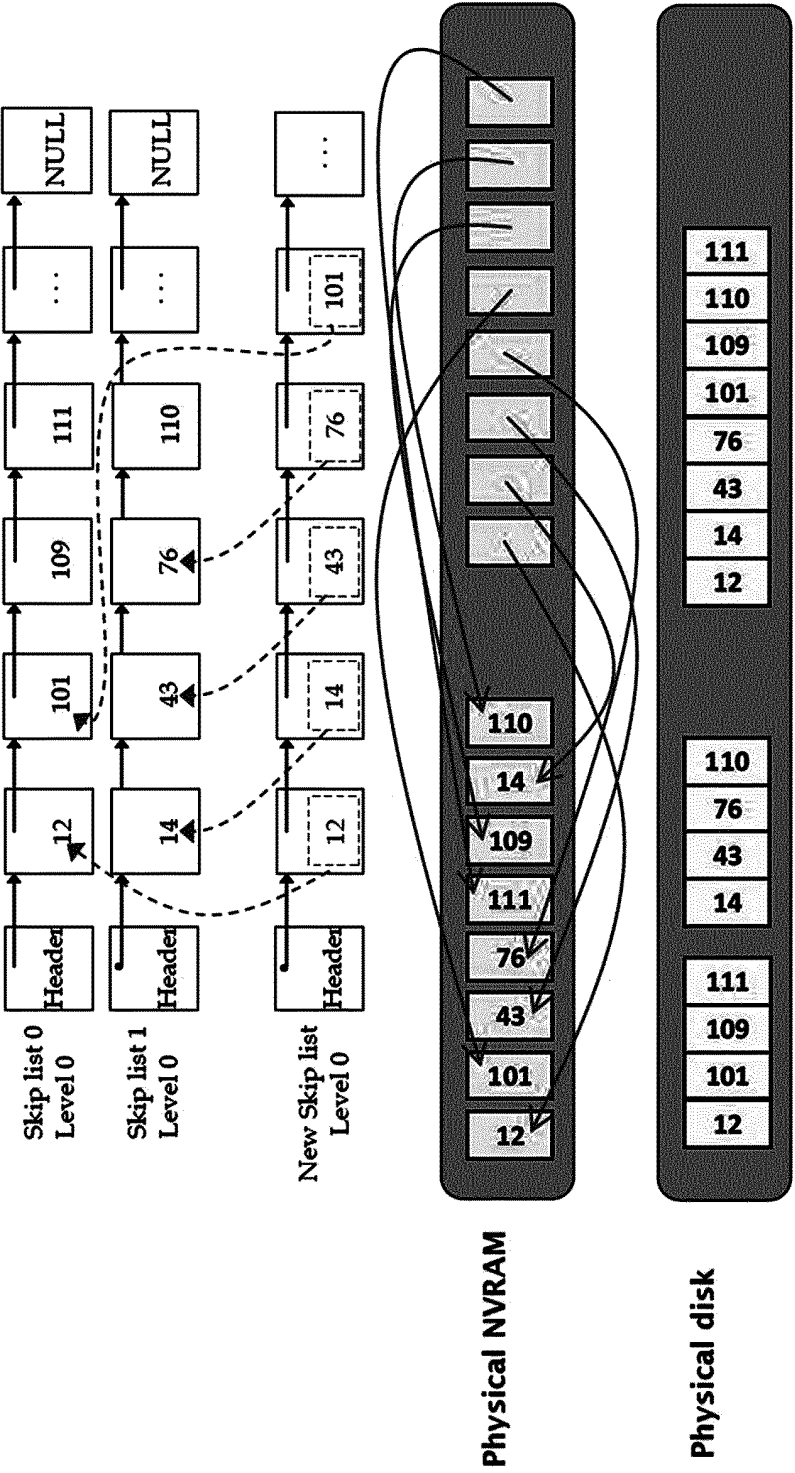
[도5]



[도6]



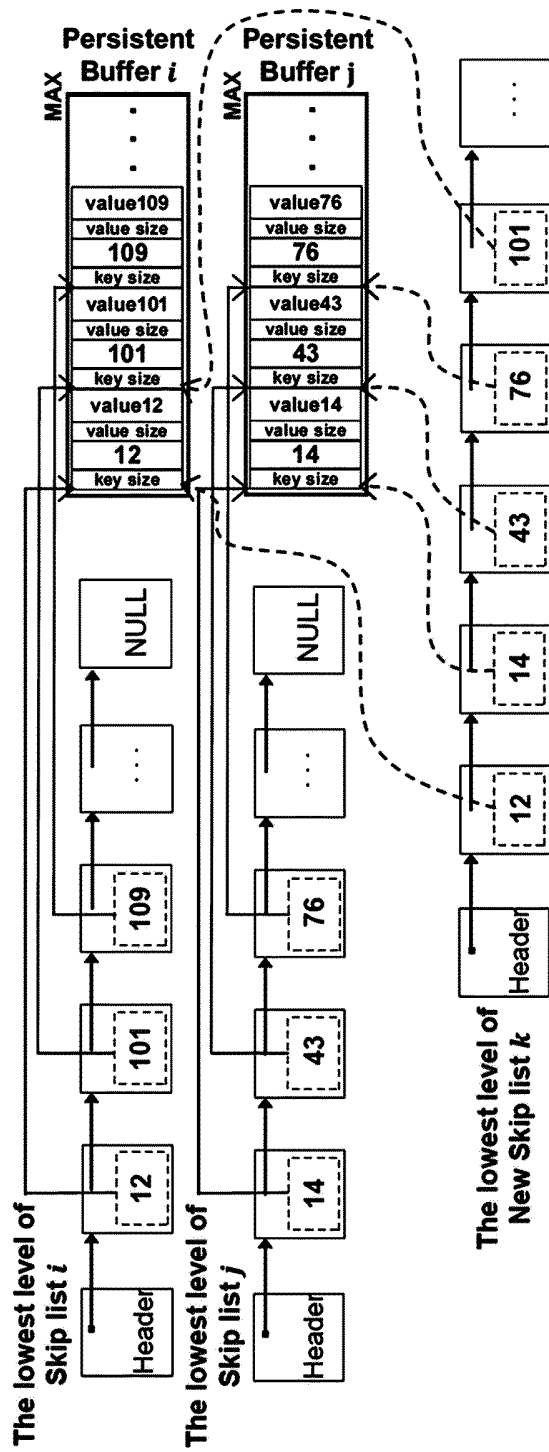
[도7]



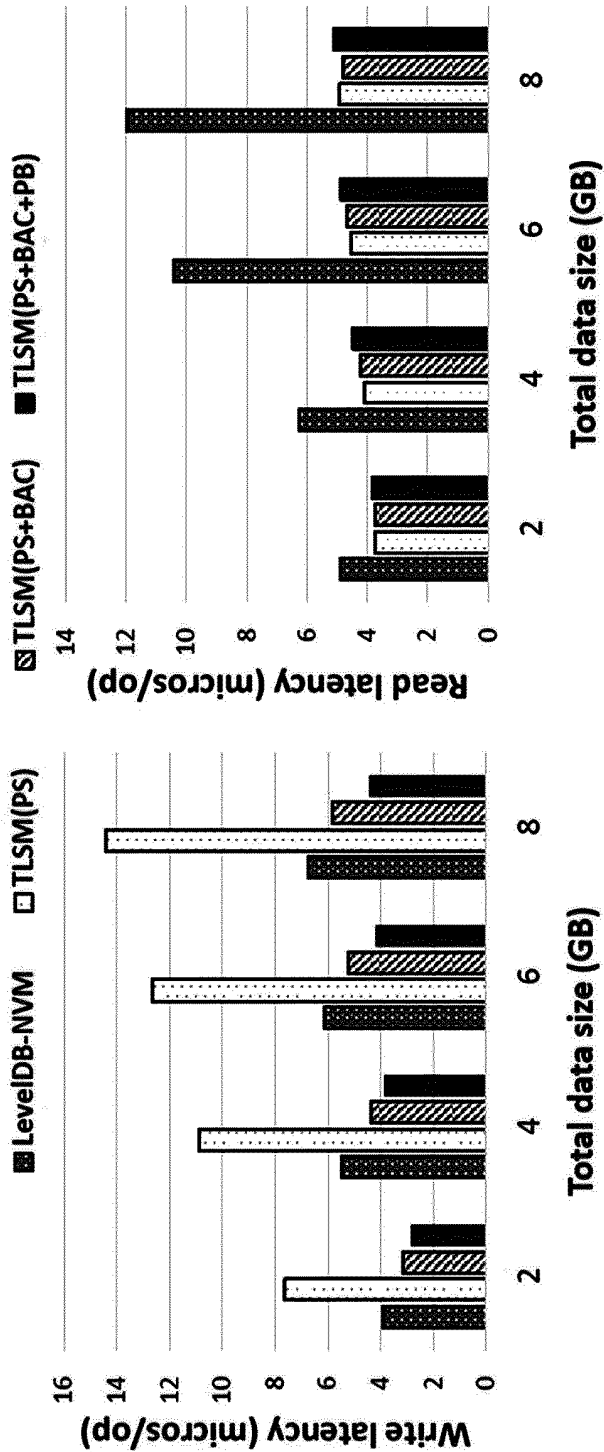




[도9]



[도10]

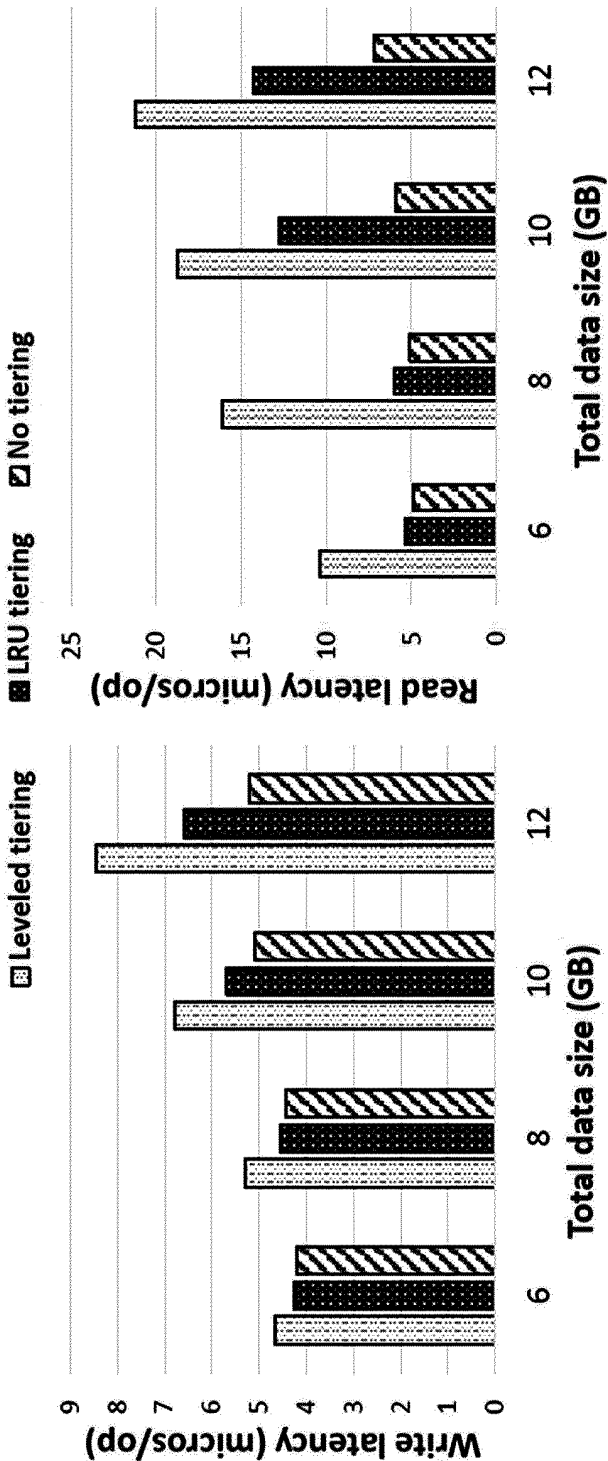


write latency on several models      read latency on several models

(a)

(b)

[51]



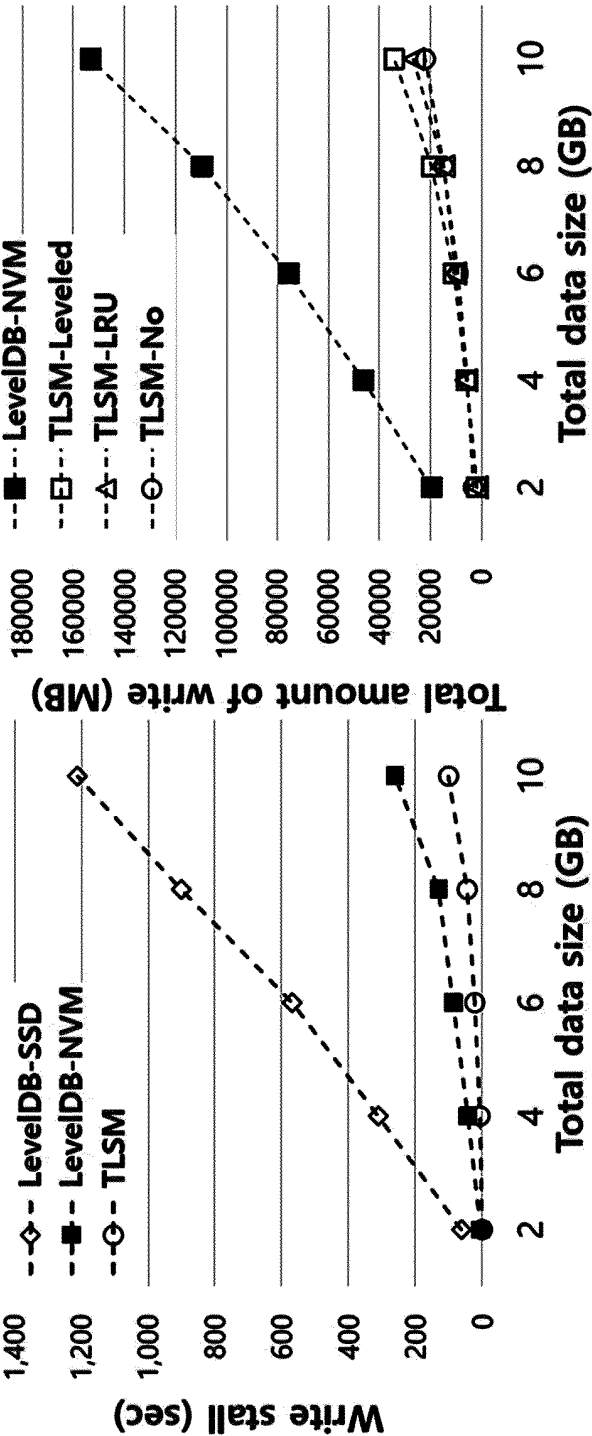
write latency on tiering options

read latency on tiering options

(a)

(b)

[도12]



Accumulation of write stall

(a)

Amount of write from compaction

(b)