

# 데이터의 가용성을 보장하는 고차원 색인 정보 관리

최현화\*, 이미영\*, 이규철\*\*  
\*한국전자통신연구원 데이터베이스연구팀  
\*\*충남대학교 데이터베이스연구실  
e-mail : [hyunwha@etri.re.kr](mailto:hyunwha@etri.re.kr)

## High Dimensional Index Information Management for Data Availability

Hyun-Hwa Choi\*, Mi-Young Lee\*, Kyu-Chul Lee\*\*  
\*Database Research Team, Electronics and Telecommunications Research Institute  
\*\*Dept. of Computer Engineering, Chung-Nam University

### 요 약

웹 서비스 혹은 클라우드 컴퓨팅 서비스로써 대용량의 멀티미디어 데이터에 대한 내용 기반 검색을 지원하기 위하여, 분산 고차원 색인 구조에 대한 연구가 활발하게 이뤄지고 있다. 이러한 고차원 데이터의 색인 구조에 대한 연구와 달리, 고차원 색인 데이터의 저장 및 관리에 대한 연구는 거의 전무한 것이 현실이다. 지금껏 대부분의 색인 데이터는 빠른 접근을 위하여 트랜잭션 관리 및 데이터의 복구를 제공하지 않은 파일 시스템에서 관리되어 왔다. 그러나, 파일 시스템에 저장된 색인 데이터는 이를 서비스하는 노드의 장애 발생 시에 일부 혹은 전체 데이터에 대한 검색이 이뤄지지 않는 문제점이 있다. 서비스의 가용성 여부가 중요한 요소인 웹 서비스와 클라우드 컴퓨팅 서비스를 위하여, 본 논문에서는 고차원 색인 데이터를 데이터베이스를 통해 관리하여, 안정성 및 가용성을 보장하면서, 고차원 데이터의 색인 및 검색의 성능을 보장하는 방법을 제안하고자 한다.

### 1. 서론

카메라 기술의 발전으로 멀티미디어 데이터가 급증하면서, 인터넷 서비스 및 클라우드 컴퓨팅 서비스로써 내용 기반 검색을 지원하기 위하여 고차원 데이터의 확장성을 지원하는 색인 구조 연구가 활발히 이뤄지고 있다. 여기서, 내용 기반 검색은 멀티미디어 데이터로부터 추출한 색깔, 모양 및 질감 등의 특징을 고차원의 특징 벡터 데이터로 변환하고, 이에 대한 색인을 구축한 후 고차원 특징 벡터 데이터 간의 유사성을 바탕으로 가장 근접한 멀티미디어 데이터를 검색하는 것이다. 고차원 데이터의 색인 구조 연구는 수 십 년간 지속되어온 데이터베이스의 주요한 연구 분야 중 하나이다. 그러나, 최근의 고차원 색인 구조 연구는 특징 벡터 데이터의 차원이 매우 많고, 검색을 위한 고차원 특징 벡터 데이터가 하나의 컴퓨팅 노드에서 수용할 수 없을 만큼 대용량이라는 것에 초점을 두고 있다.

한편, 대부분의 검색 엔진은 색인 데이터를 파일 시스템에 저장하여 관리하고 있다. 이는 데이터베이스가 데이터의 관리 용이성과 안정성이 뛰어난 반면, 파일 시스템과 비교하여 상대적으로 느리기 때문에 검색 엔진의 저장소로 인정받지 못함이였다. 그러나, 색인 데이터의 증가와 잦은 검색 및 변경을 수반하는 색인 정보를 트랜잭션 관리 및 데이터의 복구(회복)를

제공하지 않는 파일 시스템에 저장하면서, 색인 데이터를 관리하는 노드의 장애 발생 시 그 일부 혹은 전체 데이터의 검색이 수행되지 못하는 문제점이 발생하고 있다. 또한, 최악의 경우, 전체 데이터에 대한 색인을 다시 수행하는 일이 발생하기도 한다.

이에 본 논문에서는 대용량의 멀티미디어 데이터의 색인 정보를 분산 저장하는데 있어, 저장소로 데이터베이스를 사용하여 고차원 색인 데이터의 안정성 및 가용성을 제공하고, 고차원 데이터의 색인 및 검색의 성능을 보장하는 방법을 제안하고자 한다.

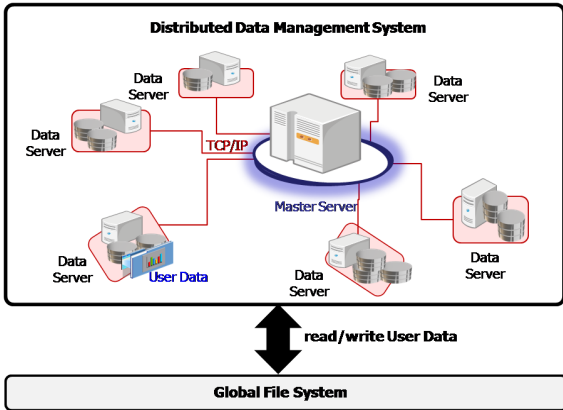
본 논문의 구성은 다음과 같다. 먼저, 2 장에서는 본 논문에서 대상으로 하고 있는 분산 데이터 관리 시스템 및 분산 고차원 색인 구조에 대해서 설명한다. 3 장에서는 데이터베이스를 통해 고차원의 색인 데이터를 관리하면서, 고차원 색인 데이터의 검색을 지원하는 방법에 대해서 설명한다. 끝으로 4 장에서는 결론을 맺는다.

### 2. 관련 연구

#### 2.1 분산 데이터 관리 시스템

웹 검색 혹은 클라우드 컴퓨팅 서비스를 위한 클러스터 환경 기반 분산 데이터 관리 시스템의 일반적인 구조는 그림 1 과 같다. 분산 데이터 관리 시스템은 그림 1 에서 보는 것과 데이터 서버 관리 및 데이터

분배를 담당하는 마스터 서버와 마스터 서버로부터 할당받은 데이터에 대한 검색, 삭제 및 삽입을 지원하는 N 개의 데이터 서버로 구성된다. 마스터 서버는 데이터 서버를 주기적으로 모니터링하여, 데이터 서버의 장애를 감지한다. 특정 데이터 서버의 장애 발생 시, 마스터 서버는 서비스가 중지된 데이터를 다른 데이터 서버에 재할당함으로써 데이터 서비스의 가용성을 지원한다. 일반적으로 장애 발생에 따른 데이터 재할당 시에 데이터의 복구 과정이 수행된다.



(그림 1) 분산 데이터 관리 시스템 구조

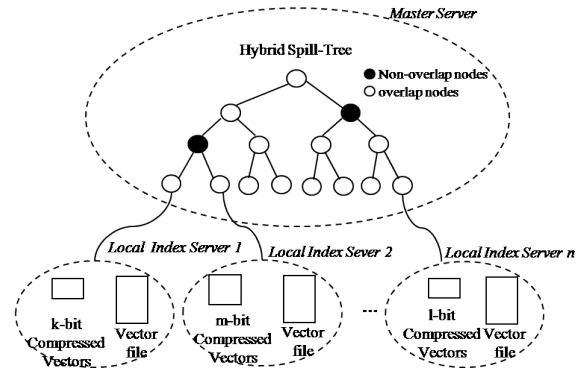
가장 대표적인 분산 데이터 관리 시스템에는 GLORY-DB[1], Bigtable[2], HBase[3] 가 있다. 그러나, 이러한 분산 데이터 관리 시스템들은 멀티미디어 데이터에 대한 내용 기반 검색을 현재 지원하지 않고 있다. 한편 데이터 관리 시스템은 글로벌 파일 시스템(global file system, distributed file system)을 바탕으로 데이터를 관리하는데, 대표적인 글로벌 파일 시스템에는 GFS[4], HDFS[5]가 있다.

### 2.2 분산 고차원 색인

최근 대용량의 고차원 데이터 검색을 위한 분산 고차원 색인 연구가 활발히 이뤄지고 있으나, 대부분 P2P 환경을 대상으로 하고 있다. 웹 서비스와 클라우드 컴퓨팅 서비스를 지원하기 위해서는 메시지 전송 및 네트워크 부하가 적은 클러스터 환경 하에서의 분산 고차원 색인 기술이 필요하다. 클러스터 환경에서의 분산 고차원 색인에는 트리를 이용한 데이터 분할 기법과 특징 벡터의 압축 버전을 이용한 필터링 기법을 병합한 색인[1][6]과 Hybrid spill-tree 를 분산 관리하는 색인[7]이 있다. 본 논문에서는 필터링 기법을 조합한 색인을 대상으로 색인 데이터의 저장 및 관리 방법에 대해서 설명한다.

그림 2 는 다중 길이 시그니처 기반 분산 색인을 나타낸다. 상위 트리는 대용량의 데이터에 대한 클러스터링 정보를 통해, 검색 시 접근해야하는 데이터의 양을 감소시킨다. 한편, 상위 트리의 말단 노드에 매핑된 분산 컴퓨팅 노드는 VA-file 기반으로 유사 검색을 수행한다. 여기서, 각 말단 노드에서 관리되는 VA-file 은 서로 다른 bit 수를 사용하여 특징 벡터를 압축

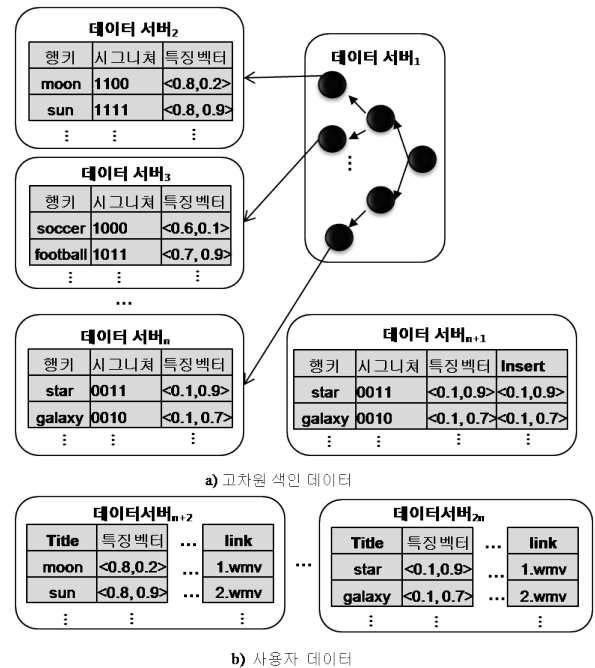
할 수 있다. 다른 길이의 bit 수에 의한 VA-file 을 기반으로 한 검색의 병렬화는 전체 검색 성능을 향상시키는 주요 요인이다.



(그림 2) 다중 길이 시그니처 기반 분산 색인 구조

### 3. 고차원 색인 데이터의 분산 관리

본 장에서는 대용량의 멀티미디어 데이터의 내용 기반 검색을 지원하기 위하여, 분산 데이터 관리 시스템에서 고차원 색인 데이터를 관리하는 방법에 대해서 자세히 설명한다.



(그림 3) 내용 기반 검색을 위해 관리되는 데이터

분산 데이터 관리 시스템에서 관리하는 데이터는 크게 4 종류이다. 먼저, 사용자 정의에 맞춰 데이터를 저장 및 관리하는 일반 사용자 데이터, 고차원 색인 구조에서 상위 트리에 해당하는 데이터의 클러스터링에 따른 분산 정보, 상위 트리의 말단 노드별로 분산된 고차원 색인 데이터, 마지막으로 대용량의 고차원 색인 구축 후 발생하는 데이터의 변경을 저장하는 변

경 고차원 색인 데이터가 그것이다. 그림 3 은 앞에서 설명한 데이터를 관리하기 위한 테이블 구조의 예를 포함한다.

분산 데이터 관리 시스템의 각 데이터 서버는 마스터 서버로부터 할당받은 데이터가 어떤 종류냐에 따라 메모리에 유지하는 정보가 달라진다. 고차원 색인 구조에서 상위 트리 정보를 할당 받은 데이터 서버는 데이터를 모두 읽어 메모리 상에 트리를 구성한다. 여기서 트리 정보는 하나의 데이터 파티션에 저장하여 오직 한 데이터 서버에서만 상위 트리가 구축되도록 한다. 이는 트리 구조를 여러 노드에 걸쳐 분산시키는 방법은 물론 다중 컴퓨팅 노드들에 걸쳐 구축된 트리의 병렬 탐색이 어렵기 때문이다.

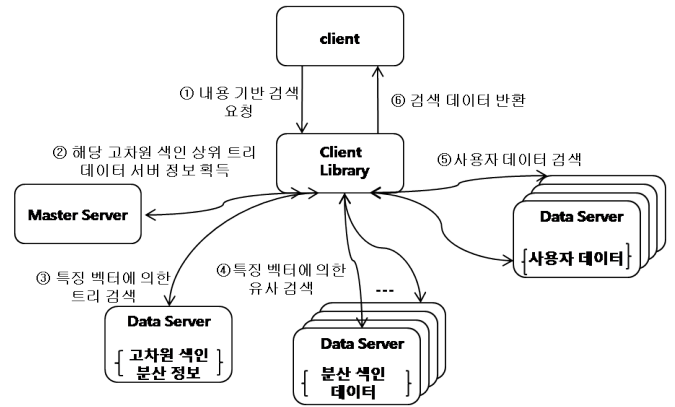
한편, 구축된 상위 트리의 말단 노드에 해당하는 고차원 색인 데이터를 할당 받은 데이터 서버는 특징 벡터의 압축 버전인 시그니처 전체를 메모리에 유지토록 한다. 이는 특징 벡터 데이터가 32bit 의 실수로 이뤄진 반면, 시그니처는 4~8bit 로 구성되어 데이터의 크기가 1/8 까지 축소되기 때문에 많은 양의 고차원 특징 벡터의 시그니처를 메모리에 유지시킬 수 있다. 메모리 내 시그니처 기반 유사 검색은 빠른 시간 내에 많은 후보 특징 벡터 데이터를 필터링 한다. 한편, 데이터 서버의 메모리 내 시그니처 관리로, 새로운 특징 벡터 데이터의 삽입이 많은 경우 메모리 부족이 발생할 수 있다. 이 경우 상위 트리의 해당 노드의 분할 및 현 데이터 서버에서 관리하고 있는 특징 벡터 데이터와 시그니처를 분할하여야만 새로운 특징 벡터 데이터를 삽입할 수 있다. 뿐만 아니라, 많은 고차원 특징 벡터 데이터의 삭제로 메모리 내 시그니처 개수가 특정 임계치보다 적은 경우, 상위 트리 내의 해당 노드와 이웃 노드를 합병할 수 있다.

특징 벡터 데이터의 변경을 저장한 변경 고차원 색인 데이터를 할당받은 데이터 서버 또한 시그니처를 메모리에 유지한다. 이는 고차원 데이터의 유사 검색 요청 시에 새로 삽입된 특징 벡터 데이터도 검색 대상이 되도록 하기 위함이다. 변경 고차원 색인 데이터의 유사 검색은, 상위 트리의 탐색 결과인 한 개 이상의 분산 고차원 색인 데이터에 대하여 유사 검색과 병렬로 수행되기 때문에 검색 성능 저하는 거의 없으면서 검색의 정확도를 높여준다.

분산된 고차원 색인 데이터에 대한 유사 검색은 시그니처 기반 유사 검색에 의해 필터링된 후 남겨진 적은 양의 데이터에 대하여 데이터베이스 검색을 통해 최종 결과를 도출한다. 이때 데이터베이스 검색은 적은 양의 데이터에 한하여 수행되고, 시그니처와 특징 벡터 데이터가 동일한 데이터 서버를 통해 관리되기 때문에 검색 시간이 오래 걸리지 않는다. 게다가 분산 데이터 관리 시스템의 데이터 할당 알고리즘은 대부분 데이터 위치를 고려되어 있다.

대용량의 고차원 데이터에 대한 색인이 구축된 후부터 발생하는 특징 벡터 데이터의 삽입 및 삭제는 색인에 즉시 반영하지 않고 따로 관리토록 한다. 이는 고차원 색인 데이터 변경을 실시간 반영하는 경우, 특정 데이터 서버에서 관리하는 고차원 색인 데이터

의 분할 혹은 병합이 발생하여 서비스의 중지를 초래할 수 있기 때문이다. 그리하여, 색인 데이터의 변경에 대한 반영은 주기적 혹은 특정 조건을 만족하는 시점에 하도록 한다. 여기서 특정 조건의 예로, 특정 데이터 서버에서 관리하기 어려운 만큼의 특징 벡터 데이터 변경이 발생하는 경우를 들 수 있다.



(그림 4) 내용 기반 검색 처리 절차

앞에서 설명한 분산 데이터 관리 시스템에서의 내용 기반 검색의 처리 절차는 그림 4 와 같다. 먼저, 우리는 내용 기반 검색에 따른 사용자 데이터 검색을 통합하여 결과를 반환하는 것을 가정하였다.

사용자는 멀티미디어 데이터로부터 추출한 특징 벡터 데이터와 이로부터 얻고자 하는 정보를 명시한 질의를 전달한다. 먼저, 마스터 서버로부터 해당 테이블의 고차원 색인의 상위 트리를 서비스하는 데이터 서버 정보를 얻는다. 획득한 데이터 서버에 질의 특징 벡터 데이터를 전달하여 트리 탐색을 통한 유사 검색을 요청한다. 그러면, 한 개 이상의 데이터 서버가 반환된다. 반환된 데이터 서버에 질의 특징 벡터 데이터를 전달하여 유사 검색을 요청한다. 이때 변경 고차원 색인 데이터를 서비스하는 데이터 서버에게도 유사 검색을 요청한다. 반환되는 유사 검색 결과를 병합하여, 최종 결과를 산출한다. 산출된 유사 검색 결과는 해당 사용자 테이블의 행 키로, 이를 이용하여 사용자가 요청한 정보를 검색하여 사용자에게 반환한다.

이때, 고차원 색인 데이터를 서비스하는 특정 데이터 서버의 장애 시, 이를 마스터 서버가 모니터링을 통해 감지한다. 그러면, 마스터 서버는 일반 데이터와 마찬가지로 고차원 색인 데이터를 다른 데이터 서버가 서비스 할 수 있도록 재할당한다. 이때 데이터베이스에서 제공하는 트랜잭션 및 데이터 복구(회복) 메커니즘에 의해 고차원 색인 데이터의 일관성을 유지한 채 서비스가 재개된다. 뿐만 아니라, 분산 데이터 관리 시스템은 마스터 서버에 대한 복구 메커니즘이 존재하기 때문에, 전체적으로 클러스터 내 컴퓨팅 노드의 장애 발생과 상관없이 데이터 서비스의 가용성이 보장된다. 결론적으로, 고차원 색인 데이터의 안정성 및 가용성이 보장되어, 고차원 색인 데이터의 파일 시스템 저장에 따른 서비스의 장시간 중지 혹은

색인 재 구축의 문제점이 해결될 수 있다.

#### 4. 결론

본 논문에서는 기존의 고차원 색인 데이터의 파일 시스템 저장 시 데이터 손실에 따른 서비스 중지 혹은 색인 재 구축 문제를 해결하기 위하여, 고차원 색인 데이터를 데이터베이스로 관리하는 방법을 제안하였다.

본 논문에서는 고차원 색인 데이터를 데이터베이스를 이용하여 관리함에 따른 검색 성능 저하를 줄이기 위하여, 다수의 컴퓨팅 노드에서 관리되는 고차원 색인 데이터, 즉 상위 트리 및 각 말단 노드에 해당하는 특징 벡터들의 시그니처를 메모리에서 관리한다. 그리하여, 유사 검색 수행 시에 트리와 시그니처를 기반으로 필터링된 적은 양의 후보 특징 벡터만이 데이터베이스 검색을 통해 정제되도록 제안하였다. 또한, 고차원 색인 구축 후 발생하는 데이터 변경에 따른 고차원 색인 데이터 변경을 따로 관리하여 주기적으로 혹은 특정 시점에 반영하는 방법을 제안하였다. 이는 고차원 색인 데이터의 실시간 반영에 따른 오버헤드로 고차원 색인 데이터의 검색이 중지되지 않도록 하기 위함이다. 대용량의 고차원 색인 데이터를 메모리에 올려 서비스함으로써 데이터베이스 검색을 포함한 전체 유사 검색의 성능은 크게 저하되지 않는다. 한편, 고차원 색인 데이터는 데이터베이스에서 제공하는 트랜잭션 및 복구 메커니즘을 통해 안정성 및 가용성이 제공된다.

그리하여, 본 논문에서 제안하는 고차원 색인 데이터의 저장 및 관리 방법은 서비스의 가용성을 보장해야 하는 웹 서비스 혹은 클라우드 컴퓨팅 서비스와 같은 응용에 효과적인 고차원 색인 데이터 관리 방법이라 하겠다.

#### 참고문헌

- [1] 최현화, 박경현, 이훈순, 이미영, “GLORY-DB: 대용량 인터넷 서비스를 위한 분산 데이터 관리 시스템”, KDBC, 2008
- [2] Fay Chang, et al., “Bigtable: A Distributed Storage System for Structured Data”, USENIX Symposium on Operating System Design and Implementation, 2006
- [3] Hbase, <http://hadoop.apache.org/>
- [4] Sanjay Ghemawat, et. al., “The Google File System”, In Proceeding of the 19<sup>th</sup> ACM Symposium on Operating Systems Principles 2003
- [5] HDFS, <http://hadoop.apache.org/>
- [6] 최현화, 이미영, 이규철, “대용량 멀티미디어 데이터의 내용 기반 검색을 위한 고확장 지원 색인 기법”, 한국 콘텐츠 학회, Vol. 7, No. 1, pp. 726-730, 2009
- [7] Ting Liu, Charles Rosenberg, Henry A. Rowley, “Clustering Billions of Images with Large Scale Nearest Neighbor Search”, Proc. IEEE WACV, 2007

-----  
본 연구는 정보통신부 및 정보통신연구진흥원의 IT 신성장동력핵심 기술개발사업의 일환으로 수행하였음. [2007-S-016-1, 저비용 대규모 글로벌 인터넷 서비스 솔루션]