

효율적인 다중 데이터 접근을 위한 방송 스케줄 생성

(A Broadcast Schedule Generation for Efficient Multiple Data Accesses)

이 상 돈 ^{*}
(Sangdon Lee)

요 약 이동 컴퓨팅 환경에서 클라이언트가 서버에 요구하는 데이터가 데이터 항목의 집합이라면, 클라이언트 측면에서는 모든 데이터 항목을 얻을 때까지의 소요 시간이 중요하다. 본 논문에서는 클라이언트가 여러 개의 데이터를 한번에 접근하는 환경에서 클라이언트가 전체 데이터 항목을 얻는데 소요되는 응답시간을 감소시킬 수 있는 방송 스케줄 생성 기법을 제안한다. 클라이언트에서 요구하는 데이터 항목들이 중복되는 정도에 따라 방송 스케줄에서 데이터 항목들의 상대적인 배치 순서가 중요해지며, 클라이언트에서의 데이터 접근 요구가 편중되는 정도에 따라 방송 스케줄 내에서 각 데이터 항목이 방송되는 빈도가 영향을 받는다. 데이터 집합의 중복 정도와 데이터 접근 편중도를 통합적으로 적용함으로써 효과적으로 다중 데이터 접근을 지원하는 계층적 방송 스케줄을 생성하는 기법을 제시하였다. 실험 결과를 통해, 제안 기법이 여러 환경에서 적용적으로 대응하여 기존의 연구 결과에 대해 성능을 개선함을 보였다.

키워드 : 이동 컴퓨팅, 데이터 방송, 방송 스케줄

Abstract When a client requests a set of data items at once, it is crucial to reduce the total time elapsed till the client receives all data items requested. This paper describes a data broadcasting technique that reduces access time for multiple data accesses from clients. The impact of the relative position of each data item in a broadcast schedule is dependent on the degree of data replication between data requests from clients. The relative broadcast frequencies for data items are affected by the degree of access skewness from clients. This paper proposes a technique for generation of a hierarchical data broadcast schedule, which can support multiple data accesses effectively by considering the data sharing among client requests and data access skewness together. Simulation results show that the proposed technique consistently performs better than the previous research results in various environments.

Key words : Mobile Computing, Data Broadcast, Broadcast Schedule

1. 서 론

데이터 방송 기법은 이동 무선 통신 환경에서 가장 각광받는 데이터 전달 기술이다. 무선 통신 환경의 비대칭성, 즉 대용량의 하향 채널과 제한된 대역폭의 상향 채널의 특성을 가장 효과적으로 활용할 수 있기 때문이

다. 기존의 전통적인 클라이언트-서버 환경과는 달리 서버는 대용량의 하향 채널을 사용하여 데이터를 방송하고, 클라이언트는 방송되는 데이터를 수신하므로써 클라이언트의 수가 증가하더라도 충분한 확장성을 제공할 수 있다는 것이 큰 장점 중의 하나이다. 이러한 데이터 방송 기술은 이동 컴퓨팅 환경의 비대칭성 뿐만 아니라 데이터가 생성되는 분량의 비대칭성[1,7], 즉 서버에서는 많은 데이터를 생성하고, 클라이언트에서는 소량의 데이터(또는 데이터 요구)를 생성하는 응용 특성과 결합되어 일상 생활에 필요한 다양한 응용을 효과적으로 제공할 수 있도록 해준다.

* 이 논문은 1999년도 한국학술진흥재단의 연구비에 의하여 지원되었음 (KRF-99-003-E00295).

^{*} 통신회원 : 목포대학교 정보공학부 교수
sdlee@mokpo.ac.kr
논문접수 : 2002년 2월 5일
심사완료 : 2002년 5월 21일

데이터가 방송되는 이동 컴퓨팅 환경은 새로운 여러 기술적인 이슈들을 제기하였으며[7], 응용의 특성에 따라 요구되는 데이터의 접근이 효과적으로 이루어질 수 있도록 방송 스케줄을 생성하는 것이 가장 중요한 이슈 중의 하나이다. 이동 컴퓨팅 환경에서 데이터 접근 효율성을 평가하는데 중요한 척도는 응답시간(response time)과 적응시간(tuning time)이다[5,8,9,11]. 응답시간은 데이터가 요구되는 시점부터 데이터를 전달 받을 때까지 소요되는 시간을 의미하며, 클라이언트의 휴지(idle) 시간을 결정한다. 적응시간은 클라이언트가 요구한 데이터를 받을 때까지 방송을 들어야만 하는데 소요되는 시간을 의미하며 데이터 전송을 청취하는 동안에 클라이언트가 소비하는 전력을 결정한다. 결국 효율적인 데이터의 접근을 위해서는 응답시간과 적응시간을 최소화시키는 것이 요구된다[1,2,5,8,9,11].

방송(broadcasting) 기법은 방송 채널을 하나의 캐시처럼 간주할 수 있다. 즉 클라이언트의 캐시로부터 서버의 디스크에 이르는 저장 계층 구조의 중간에, 방송되는 데이터들로 구성된 새로운 캐시 계층이 추가되는 것으로 간주할 수 있다. 그러나 이러한 방송 캐시의 중요한 특징 중의 하나는 데이터의 전달이 순차적으로 이루어진다는 것이다. 전통적인 저장구조에서는 캐시 계층에 있는 데이터의 접근 시간은 동일한 것으로 간주된다. 이와 달리 방송 캐시에서의 데이터의 접근 시간은 데이터의 위치에 종속적이다. 그러므로 클라이언트의 응답시간은 방송 캐시 내에서의 상대적인 위치에 의존적이다. 또한 전체 클라이언트의 평균 응답 시간을 최소화하기 위해서는, 빈번히 접근되는 데이터의 방송 빈도를 상대적으로 크게 하여 방송 캐시를 계층화시키는 기법[1,2,11]이 효과적이다.

클라이언트가 어떤 응용을 처리하기 위하여 한번에 하나의 데이터 항목만을 요청한다면, 개별 데이터의 빈도에 기반하여 계층적으로 방송 스케줄을 구성하는 전략은 매우 효과적이다. 그러나 만일, 클라이언트가 한번에 요청하는 데이터가 여러 데이터 항목으로 구성되어 있으며, 클라이언트가 모든 데이터를 받을때까지의 시간이 중요한 요인이라면, 개별 데이터 항목의 요구 빈도에 기반한 방송 스케줄링 기법은 효과가 제한적일 수 있다. 만일 클라이언트에서 동시에 요구하는 데이터 집합이 개별 데이터의 접근 빈도가 매우 큰 것과 매우 작은 것으로 구성되어 있다면 클라이언트가 모든 데이터를 받을 때까지의 소요시간은 결국 접근 빈도가 작은 데이터 항목에 의존적일 것이므로 접근 빈도가 큰 개별 데이터만을 우선적으로 고려하는 기존의 접근 방법은 효율성

이 떨어지게 된다.

이러한 다중 데이터의 접근을 요구하는 응용의 예는 매우 광범위하다. 주식 정보의 방송 환경을 고려하자면 최 상위 10개의 매매 종목과 매수 종목, 그리고 블루칩의 시세를 확인하는 클라이언트의 응용들을 예로 들 수 있다. 이러한 응용들에 해당되는 주식 종목들은 서로 중복될 수 있으며 이동 클라이언트에서의 각 응용은 어느 단일 종목만에 대한 정보를 필요로 하는 것이 아니라 조건에 해당하는 모든 데이터 항목을 전송받는 것을 요구한다.

본 논문에서는 이동 컴퓨팅 환경에서 클라이언트가 여러 개의 데이터 항목에 대해 데이터 접근을 요구하는 환경에서 클라이언트의 응답 시간을 감소시키므로써 다중 데이터 접근을 효율적으로 처리할 수 있는 계층적 방송 스케줄을 생성하는 기법에 대해 기술한다. 이후 2장에서는 이와 관련된 기존의 연구 결과를 기술하며, 3장에서는 다중 데이터의 접근을 위한 주요 개념과 용어를, 그리고 4장에서는 제안하는 방송 스케줄링 방안에 대해 기술한다. 5장에서는 실험을 통해 제안기법을 평가 및 분석하고 6장에서 결론을 맺는다.

2. 관련 연구

이동 컴퓨팅 환경에서 방송 기법을 통해서 데이터를 전달하는 경우 push/pull 여부와 주기적/비주기적 여부에 따라 크게 4가지로 분류가 가능하며[7] 각 분야별로 많은 연구가 이루어져 왔다. 방송 대역폭을 디스크로 간주하고 데이터의 접근 빈도에 따라 그룹으로 분류한 후 자주 접근되는 데이터들을 작고 빠른 디스크에 배치하므로써 전체적인 응답시간을 감소시키는 방송 디스크(broadcast disk) 기법[1,2]이 연구되었으며, 데이터 접근 빈도가 편중된 분포를 보일수록 계층적인 방송 디스크 구성 기법이 유리함을 보였다. 이 연구에서의 가정과 달리 접근 패턴이 동적이고 상향 링크가 존재하는 경우에는 방송 대상이 되는 데이터 집합을 효과적으로 선정할 필요가 있으며, 접근 패턴을 미리 알고 있는 경우 정적으로 방송 데이터를 선정하거나[9] 또는 시스템의 작업 부하에 따라 동적으로 구성하는 것[10]이 효과적이다.

시스템이 낮은 작업 부하를 가지는 경우에는 pull에 기반한 데이터 전달 방법이 효과적이며, 작업 부하가 클 경우 push에 기반한 데이터 전달 방법이 효과적이다. 그러므로 두가지 방법을 효과적으로 결합하여 다양한 작업부하에 대응할 수 있는 방안들[2]이 연구되었다.

대부분 데이터 방송 기법들이 방송되는 데이터의 길이가 균일한 데이터를 다루고 있는 반면에 다른 길이를

갖는 이질적인 데이터의 방송에 대한 연구가 이루어졌다[3,11]. 가변길이 데이터 항목이 방송 스케줄 내에서 다시 방송되는 시간 간격인 “방송간격(spacing)” 개념에 기반하여 평균 응답시간을 최소화시키도록 데이터의 방송 빈도를 결정하는 이상적인 기준과 현실적으로 적용 가능한 알고리즘들이 제시[11]되었다. 또한 이질적인 데이터 요구에 대한 주문형(On-Demand) 데이터 방송 환경을 대상으로 평균 응답시간이 아닌 스트레치(stretch)¹⁾라는 기준을 사용한 방송 데이터의 스케줄링 기법에 대한 연구[3]가 이루어졌다. 그러나 이러한 연구들에서는 방송되는 데이터 항목들은 서로 겹치지 않는 것으로 가정하고 있다.

클라이언트들이 서로 겹칠 수 있는 데이터 집합을 요구하는 경우 이들 데이터 집합 간의 상호 거리를 최소화시킬 수 있도록 다중 데이터 사이의 관계를 표현하는 연구[4] 결과가 제시되었다. 이 연구에서는 Gray 코딩[6]에 기반하여 데이터 집합들을 방송 스케줄 내에서 집중화시키므로써 클라이언트의 응답시간을 줄일 수 있음을 보였다.

다중 데이터의 접근 문제는 부분적으로 서로 중복될 수 있는 가변길이의 데이터를 효율적으로 방송하는 문제와 매우 유사하다. 데이터 집합은 동일한 길이를 갖는 가변 개수의 항목의 구성으로 간주될 수 있기 때문이다. 그러므로 본 논문에서 추구하는 다중 데이터 접근을 위한 계층적 방송 스케줄 구성과 가장 밀접한 연관성이 있는 기존 연구는 [11]과 [4]이다. 만일 클라이언트들에서 요구되는 데이터 집합 사이에 부분적인 중복이 없다면 각 데이터 집합은 가변 길이를 갖는 하나의 데이터 항목으로 간주할 수 있으며 이 경우 [11]에서 가정한 환경과 동일하므로 [11]의 접근 방법을 적용할 수 있다. 그러나 다중 데이터 요구 사이에는 요구되는 데이터 항목이 중복되는 것이 일반적이다. 이 경우 [4]의 연구와 같이 방송 스케줄 내에서 다중 데이터 요구 사이의 중복 정보를 사용한 데이터 집중(clustering)을 적용하는 것이 필요하다. 그러나 [4]는 방송 스케줄 내에서 각 데이터 항목이 한번씩만 포함되는 평면(flat) 방송 스케줄만을 고려하고 있으므로 클라이언트의 데이터 요구가 편중되는 경우 효과적이지 못하다.

1) 데이터 요구가 하나만 있을 경우 이를 처리하는데 요구되는 시간을 서비스 시간이라 할 때 스트레치(stretch)는 서비스 시간에 대한 응답시간의 비율로서 정의[3]된다. 이 기준은 길이가 다른 작업들 사이의 공평성(fairness)이 중요한 경우 적용될 수 있다.

3. 다중 데이터 접근을 위한 방송 스케줄

3.1 시스템 환경 및 성능 척도

시스템 구성 요소는 방송 서버와 이동 클라이언트로 구성된다. 방송 서버는 방송 채널을 통해 대상 데이터 집합을 주기적으로 방송한다. 클라이언트는 일련의 데이터를 요구하고²⁾ 방송 채널을 감시한다. 방송 채널에 클라이언트가 요청했던 데이터가 나타나면 클라이언트는 이 데이터를 자신의 캐시에 보관한다. 요청하였던 데이터 집합에 포함된 모든 데이터를 수신하면 클라이언트는 이 데이터를 사용하여 클라이언트에서 필요한 작업을 수행한다. 클라이언트가 데이터 집합을 요청한 시점부터 요청한 데이터 집합을 모두 수신할 때까지의 시간을 “다중 데이터 접근 시간(Multiple Data Access Time: 이하 MDAT을 사용함)”이라 정의한다. 각 클라이언트 입장에서 보면 자신의 MDAT를 최대한 줄이는 것이 중요하다. 그러나 시스템 전체 입장에서는 모든 사용자의 평균 MDAT를 최소화시키는 것이 중요하다. 그러므로 본 논문에서 목표로 하는 성능 척도로서 평균 MDAT를 사용한다.

3.2 용어 및 기호 정의

이 절에서는 본 논문에서 사용되는 용어와 기호에 대해서 기술한다.

방송 서버에는 총 m 개의 유일한 데이터 항목의 집합으로 구성되는 데이터베이스(D)가 존재하며 데이터 항목들은 동일한 크기를 갖는다. 즉 $D = \{d_1, d_2, \dots, d_m\}$ 이다. 이동 클라이언트는 여러 데이터 항목에 대한 접근을 요청하며 이를 질의(Q_i)라 한다. Q_i는 클라이언트가 요구하는 K_i개의 유일한 데이터 항목의 집합으로 구성된다. 즉 $Q_i = \{d_{ij}\}, d_{ij} \in D, (1 \leq j \leq K_i)$ 이다. 또한 질의 Q_i의 크기(|Q_i|)는 Q_i를 구성하는 데이터 항목의 개수이다. 즉 |Q_i| = K_i이다. Q_i가 요청되는 횟수는 질의의 접근 빈도(F_i)라 하며, 이때 Q_i에 대한 클라이언트로부터의 접근 확률(P_i)은 $P_i = (F_i / \sum F_j), 1 \leq j \leq n$ 이고 $\sum P_i = 1$ 이다.

각 질의가 요청한 데이터 집합을 방송한 후 다시 이 데이터 집합을 방송할 때까지의 소요 시간을 방송 간격

2) 본 논문에서는 [1,4,11]에서의 연구와 같이 push를 통해 데이터가 주기적으로 방송되는 환경을 가정한다. 방송 서버는 주기적인 데이터의 방송 순서를 결정하는 시점에서 클라이언트들에서 요구되는 데이터 항목들과 접근 특성들을 미리 알고 있다고 가정한다. 이러한 push를 기반으로 하는 방송 환경에서 클라이언트로부터의 데이터 요청은 상향 채널을 통해서 서버에 물리적으로 데이터를 요청하는 것이 아니라 해당 데이터 항목이 방송 채널에 존재하는지 여부를 감시하는 논리적인 데이터 요청을 의미한다.

이라 하며, S_{ij} 는 Q_i 를 j 번째 방송한 후 다시 Q_i 를 방송할 때까지의 방송 간격을 말한다. 또한 **데이터 갭(G_{ijk})**는 S_{ij} 내에서 Q_i 를 구성하는 데이터 항목들 사이에 방송되는 $Q_j(j \neq i)$ 에 속하는 연속된 데이터 항목의 개수로서 G_{ijk} 는 S_{ij} 의 k 번째 데이터 갭을 의미한다. 또한 $|G_{ijk}|$ 는 갭 G_{ijk} 의 크기를 나타낸다. **데이터 집합의 최소거리(D_i)**는 Q_i 의 방송 간격에서 최대 데이터 갭을 뺀 값. 즉 $D_i = S_{ij} - \max(|G_{ijk}|)$, ($1 \leq k \leq S_{ij}$)을 의미한다. 그림 1은 이상에서 기술한 방송 간격 및 데이터 집합의 최소거리와 갭의 예를 보여준다.



그림 1 방송 간격 및 데이터 집합의 최소거리와 갭의 예

데이터의 방송은 주기적으로 반복되며 한 방송 주기 내에서 각 데이터 항목들이 방송되는 순서를 **방송 스케줄(σ)**이라 한다. 한 질의 Q_i 가 데이터 항목들에 대한 접근을 요청한 시점부터 요청한 데이터를 모두 전송받을 때까지의 소요시간의 평균을 질의 Q_i 의 **평균 다중 데이터 접근시간(T_i)** 또는 **평균 MDAT**라 한다. Q_i 에 갭이 M 개 있을 때 T_i 는 다음의 수식으로 유도된다[3].

$$T_i = S_{ij} - \frac{1}{S_{ij}} \times \sum_{j=1}^M \sum_{k=0}^{G_{ijk}} k \quad (\text{수식 1})$$

또한 한 클라이언트에 있는 모든 질의의 **평균 다중 데이터 접근 시간(T)** 또는 **평균 MDAT**는 수식 2와 같이 나타내어진다.

$$T = \sum_{i=1}^M T_i \times P_i = \sum_{i=1}^M \left(S_{ij} - \frac{1}{S_{ij}} \times \sum_{j=1}^M \sum_{k=0}^{G_{ijk}} k \right) \times P_i \quad (\text{수식 2})$$

만일 $S_{ij} = S_i (1 \leq j \leq F_i)$ 라면, G_{ijk} 는 G_{ik} 가 되며 수식 2는 아래와 같이 변환된다.

$$T = \sum_{i=1}^M T_i \times P_i = \sum_{i=1}^M \left(S_i - \frac{1}{S_i} \times \sum_{j=1}^M \sum_{k=0}^{G_{ijk}} k \right) \times P_i \quad (\text{수식 3})$$

이상의 기본 개념을 바탕으로 본 논문에서 해결하고자 하는 문제를 다시 정의하면 이동 클라이언트에서 질의 $Q_i (1 \leq i \leq M)$ 와 데이터 접근 확률 $P_i (1 \leq i \leq M)$ 이 주어졌을 때 전체 질의의 평균 다중 데이터 접근 시간

(T)을 최소화시키는 계층적 방송 스케줄(σ)를 구하는 문제이다. 즉 다음의 예제 1에 대한 답을 구하는 것이 본 논문의 목표이다.

[예제 1] 세 질의 Q_1, Q_2, Q_3 가 각각 $\{d_1, d_2, d_3, d_8\}, \{d_2, d_3, d_4, d_5\}, \{d_1, d_5, d_6, d_7\}$ 의 데이터 집합을 비율 2:1:1로 접근 요구하는 경우 효과적인 계층적 방송 스케줄은 무엇인가?

4. 문제 분석 및 제안 기법

4.1 문제의 심층 분석 및 관찰

각 질의의 방송 간격이 동일할 때 방송 스케줄이 최소의 응답시간을 갖는 것으로 알려져 있다[11]. 물론 각 질의의 방송 간격을 항상 동일하게 할 수는 없지만 이상적인 최적의 방송 빈도를 유도해 내기 위한 기초로 사용될 수 있다. 그러므로 본 논문에서도 방송 간격이 동일하다는 가정을 사용한다.

한 질의의 평균 MDAT T_i 는 갭의 크기에 의존적이다. 또한 방송 간격에도 의존적이다. 한편 T_i 는 최대 갭의 크기가 클수록 감소하는 것을 [4]에서 검증한 바 있다. T_i 를 최소화하기 위해서는 방송 간격(S_i)이 최소이면서 최대 갭을 갖도록 해야 하며, 최대 갭을 갖는다는 것은 데이터 집합의 최대거리(D_i)를 최소로 해야함을 의미한다.

방송 간격의 하한(lower bound)은 한 질의의 데이터 집합 중 방송을 위해 선정된 데이터 항목의 개수로 결정된다. 또한 D_i 를 최소로 하기 위해서는 한 질의의 모든 데이터 항목을 연속적으로 방송 스케줄에 배치해야 함을 의미한다. 만일 질의 간에 데이터 공유가 존재하는 경우, D_i 를 최소로 하기 위해서는 중복된 데이터를 전부 방송해야 하며, 이 경우 방송 간격이 길어진다. 방송 간격을 최소로 하기 위해 중복된 데이터를 방송하지 않는 경우에는 D_i 가 커진다. 물론 방송 간격과 D_i 를 동시에 최소화시킬 수 있는 경우가 존재하기는 하나 매우 특수한 경우에 해당하므로 일반적으로 방송 간격과 D_i 간에는 이해득실이 존재한다. 만일 질의 간의 데이터 공유가 전혀 없다면, 다른 질의에 영향 받지 않고 한 질의의 모든 데이터 항목을 연속적으로 배치함으로써 D_i 를 최소로 조정하는 것이 가능하므로 스케줄링 문제는 결국 적절한 방송 간격을 결정하는 문제로 귀납된다. 이 경우에는 질의의 데이터 집합을 가변길이를 갖는 한 항목처럼 간주할 수 있으므로 가변길이 데이터 항목에 대한 방송 간격, 즉 방송 빈도를 결정하는 [11]의 접근 방법이 사용될 수 있다.

3 평면 스케줄을 대상으로 [4]에서 MDAT와 유사한 개념인 "평균 접근시간"을 유도하였으며 비슷한 과정으로 MDAT를 유도하는 것이 가능하므로 본 논문에서는 유도과정을 생략한다.

4.2 접근 방법

전체 질의의 평균 MDAT를 최소화하기 위해 고려해야 하는 주요 사항은 다음과 같다.

첫째, 질의 사이에 데이터가 공유되는 경우 방송 스케줄 내에서 질의를 구성하는 데이터 간의 연관도에 따라 상대적인 위치를 결정하는 것이다. 질의에서 요구하는 데이터가 단일 데이터 항목인 경우나 또는 중복 정도에 상관없이 질의가 요구하는 데이터 집합을 항상 전부 방송하는 경우에는 방송 스케줄 내에서 데이터 항목 간의 상대적 위치를 고려할 필요 없이 상대적인 빈도만을 고려하는 것으로 충분하다. 그러나 중복되는 데이터의 일부분을 방송하기로 결정하는 경우에는 데이터의 중복 정도와 상대적인 중요도(즉 데이터 접근 빈도)에 따라 적절한 상대적 순서를 따르는 것이 유리하다.

둘째, 중복되는 데이터의 일부만 방송하는 경우 각 질의의 데이터 집합 중 방송에 실제로 포함될 데이터 집합의 선정이 필요하다. 전송한 바와 같이 방송될 데이터 항목들은 방송 간격(S_i)에 영향을 미치고 또한 D_i 에도 영향을 미치기 때문이다.

셋째, 방송 스케줄 내에서 각 질의를 위한 데이터 집합의 방송 빈도를 결정하는 것이다. 계층적 스케줄의 경우 적절한 방송 빈도의 결정은 성능에 직접적인 영향을 미치는 중요한 요인 중의 하나이다.

다음 절에서는 이상에서 설명한 주요 고려사항의 결정에 대해 기술한다.

4.2.1 데이터 집합간의 연관성 표현

Gray 코드[6]는 다차원 공간을 일차원으로 매핑하는 기법의 하나로서, 숫자를 이진 비트열로 표현하고 이 비트열들이 정확히 1비트씩 차이가 나도록 연속적으로 이진 비트열들을 배치한다. 여러 가지 표현 방법들이 있으나 일반적으로 가장 많이 사용되는 것은 “이진 사상(binary-reflected)” Gray 코드이다. 이진 비트열의 Gray 코드 값은 Gray 코드로 표현된 비트 열에서의 상대적인 위치를 의미한다.

본 논문에서는 데이터 집합간의 연관성 표현을 위해 [4]에서 제안된 기법을 사용한다. [4]에서는 각 질의가 요구하는 데이터 항목을 비트로 표현하고 질의의 데이터 집합에 대한 Gray 코드를 생성한 후 Gray 코드값에 따라 데이터 항목들을 순서대로 배치하는 평면 방송 스케줄 구성을 제안하고 데이터 집합간의 연관성을 표현하지 않은 방송 스케줄보다 우수한 성능을 보이는 것을 보였다. 그림 2는 [4]에서 제안한 기법에 따라 $P_1 > P_2 > P_3$ 일 때 예제 1의 질의들에 대해서 각 질의의 데이터 항목들을 비트열로 표현하고 Gray 코드 값에 따라

정렬한 결과를 보여준다. 결국 방송 스케줄이 Gray 코드 값의 오름차순인 $\langle d_7, d_6, d_5, d_4, d_3, d_2, d_1, d_8 \rangle$ 이거나 또는 내림차순인 $\langle d_8, d_1, d_2, d_3, d_4, d_5, d_6, d_7 \rangle$ 으로 구성된다.

데이터항목	Gray 코드
	000
d6,d7	001
d5	011
d4	010
d2,d3	110
	111
d1	101
d8	100

그림 2 3-비트 이진 사상 Gray 코드 값에 기반한 예제 1의 데이터 집중 예

4.2.2 질의를 위한 방송 대상 데이터 집합의 선정

방송 대상 데이터 집합($BQ_i, 1 \leq i \leq M$)은 각 질의 Q_i 내의 데이터 항목 중 실제로 방송 스케줄에 추가될 항목의 집합을 의미한다. Q_i 와 BQ_i 사이에는 다음의 관계가 성립한다. 즉 방송 대상 데이터 집합은 각 질의의 부분집합이며, 방송 대상 데이터 집합은 질의 내의 모든 데이터 항목들을 포함해야 한다.

$$BQ_i \subseteq Q_i, \text{ 그리고 } \bigcup_{i=1}^M BQ_i = \bigcup_{i=1}^M Q_i$$

또한 BQ_i 의 크기는 해당 질의의 방송 간격 뿐 아니라 다른 질의의 방송 간격에도 영향을 미친다. BQ_i 를 결정하는 방법에는 중복된 데이터의 방송 여부에 따라 다음과 같이 구분 가능하다.

가. 중복 허용

각 질의에서 중복된 데이터를 방송 대상에 포함시킨다. 그러므로 방송 데이터 집합은 각 질의와 동일하다. 이 경우 BQ_i 는 다음과 같다.

$$BQ_i = Q_i$$

예제 1의 경우를 적용하면 $BQ_1 = Q_1 = \{d_8, d_1, d_2, d_3\}$, $BQ_2 = Q_2 = \{d_2, d_3, d_4, d_5\}$ 그리고 $BQ_3 = Q_3 = \{d_5, d_6, d_7, d_1\}$ 이 된다.

나. 완전 중복 제거

각 질의에 중복된 데이터를 완전히 제거한다. 접근 빈도가 큰 질의의 D_i 를 최소화시킬 수 있도록 BQ_i 는 다음과 같이 결정된다. 여기서 $i < j$ 이면 $P_i > P_j$ 가 항상 성립한다. 즉 질의들이 접근 빈도의 내림차순으로 정렬되어 있음을 의미한다. 이것은 접근 빈도가 큰 질의, Q_i 에 대

해서 우선적으로 D_i 를 감소시키기 위함이다.

$$BQ_i = Q_i - \bigcup_{j=0}^{i-1} Q_j, \quad (Q_0 = \{\}, \quad i \geq 1)$$

예제 1의 경우를 적용하면 $BQ_1 = \{d_6, d_1, d_2, d_3\}$, $BQ_2 = \{d_4, d_5\}$ 그리고 $BQ_3 = \{d_6, d_7\}$ 이 된다.

다. 부분 중복 제거

각 질의에 중복된 데이터의 일부를 경우에 따라 방송 대상 집합에 포함 여부를 결정하는 것이다. 본 논문에서는 부분 중복 제거의 경우는 고려 대상에서 배제하였다.

4.2.3 데이터 집합의 방송 빈도 결정

만일 데이터 중복이 존재하지 않는 경우에는 최적의 방송 빈도가 $\sqrt{P_i/|Q_i|}$ 에 비례하는 것을 [11]에서 유도한 바 있다. 여기서 $|Q_i|$ 는 질의의 크기이다. 그러나 이 경우 질의간 인접도(또는 데이터 중복의 정도)를 고려하지 않고 있다. 한편 예제 1과 같은 유형의 문제는 최적 선형 배열 문제(OLAP: Optimal Linear Arrangement Problem)로 구분되는 NP-complete 문제임을 [4]에서 보인바 있다. 그러므로 본 논문에서는 [11]의 관찰을 기반으로 하여 단순한 질의의 길이 대신 “가중치 상대적 코드 거리”를 사용하는 방안을 제안한다.

먼저 “상대적 코드 거리”의 개념을 설명한 후 “가중치 상대적 코드 거리”에 대해서 기술한다.

가. 상대적 코드 거리

“상대적 코드 거리” RD_{ij} 는 Gray 코드 값의 감소 순으로 정렬된 두 순서 데이터 집합 $R_i = \{d_{i1}, \dots, d_{i\alpha}\}$, $R_j = \{d_{j1}, \dots, d_{j\beta}\}$ 이 주어질 때 R_i 의 마지막 데이터 항목($d_{i\alpha}$)으로부터 R_j 의 마지막 데이터 항목($d_{j\beta}$)까지의 Gray 코드로 표현된 비트열 내에서의 상대적인 거리를 의미한다. RD_{ij} 는 다음과 같이 정의된다⁴⁾. 함수 $Index(d_{ij})$ 는 Gray 코드값의 감소 순으로 정렬된 배열에서 데이터 항목 d_{ij} 의 인덱스를 반환하는 함수이다. 또한 N 은 최대 인덱스, 즉 상이한 Gray 코드 값을 갖는 항목의 총수이다. 그리고 같은 코드값을 갖는 데이터 항목은 같은 인덱스를 갖는 것으로 가정한다.

$$RD_{ij} = \begin{cases} Index(d_{jr}) - Index(d_{ia}) \\ \text{(if } Index(d_{jr}) - Index(d_{ia}) \geq 0) \\ N - \{Index(d_{jr}) - Index(d_{ia})\} \\ \text{(otherwise)} \end{cases}$$

예를 들어 예제 1로부터 Gray 코드 값에 기반한 데이터 집중 과정을 거친 후 완전 중복 제거를 통해 얻어진 각 질의를 위한 방송 데이터 집합 $BQ_1 = \{d_6, d_1, d_2, d_3\}$, $BQ_2 = \{d_4, d_5\}$ 그리고 $BQ_3 = \{d_6, d_7\}$ 에 대해서 상대

적 코드 거리를 구하는 예를 들어보자. 그림 3의 A로부터 $RD_{12} = Index(d_5) - Index(d_3) = 5 - 3 = 2$ 가 된다. 유사하게 $RD_{13} = Index(d_7) - Index(d_3) = 6 - 3 = 3$ 이 얻어짐을 알 수 있다. 일련의 과정을 거쳐 R_{ij} 를 구한 것을 행렬로 표현하면 그림 3의 (B)와 같이 얻어진다. RD_{ij} 는 Gray 코드를 사용하여 얻어진 데이터 집합간의 연관성 정도를 R_i 와 R_j 간의 상대적 위치 차이로 표현한 것이다. 연관성이 큰(작은) 데이터 집합은 상대적으로 작은(큰) 상대적 코드 거리를 가지게 된다.

Index	데이터항목	Gray 코드
1	d8	100
2	d1	101
3	d2,d3	110
4	d4	010
5	d5	011
6	d6,d7	001

(A)

	1	2	3
1	0	2	3
2	4	0	1
3	3	5	0

(B)

그림 3 A) 유효한 Gray 코드값, B) 상대적 코드 거리 행렬의 예

나. 가중치 상대 코드 거리

이상에서 기술한 “상대적 코드 거리”를 방송 스케줄의 결정에 사용하는 경우 데이터가 중복된 질의의 상대적인 인접도가 스케줄 내에서의 배치 순서에 영향을 미친다. 그러나 데이터가 중복되지 않은 경우에는 인접도와 무관하게 스케줄이 결정되어야 한다. 그러므로 이 경우 상대적 코드거리의 차이는 의미가 없게 된다. 단순히 코드 거리만을 스케줄의 생성 과정에서 고려하는 경우에는 이러한 독립성을 표현하지 못하는 이상 현상을 가진다. 예를 들어 $Q_4 = \{d_1, d_2\}$, $Q_5 = \{d_3, d_4\}$, $Q_6 = \{d_5, d_6\}$ 인 경우 중복되는 데이터 요구가 이들간에 존재하지 않으므로 Q_4, Q_5, Q_6 이 서로 독립적이나, Gray 코드의 표현시 상대적인 코드 거리가 코드 내에서의 위치에 따라 달라지므로 방송 순서의 결정에 불필요하게 큰 영향을 미치게 된다.

직관적으로 관찰해보면 방송 스케줄 내에서의 질의의 방송 순서가 “상대 코드 거리”에 의존하는 정도는 데이터의 공유도가 커짐에 따라 커지고, 또한 질의의 접근 확률이 편중될수록(즉 θ 가 클수록) 작아지는 것이 바람직하다. 그러므로 본 논문에서는 “상대 코드 거리”에 데이터의 공유 정도와 θ 의 가중치를 적용한 “가중치 상대 코드 거리(WRD_{ij})”를 사용한다. 질의간에 데이터가 중복되는 정도는 “평균 데이터 공유도”로 표현한다.

4) RD_{ij} 는 R_i 와 R_j 의 전후관계에 독립적이므로 첫 번째 항은 R_i 보다 R_j 가 앞에 있는 경우, 두 번째 항은 순서가 바뀐 경우를 대상으로 한다.

질의 Q_i 의 데이터 공유도(SF_i)는 Q_i 내의 데이터 항목 중 다른 질의와 중복되는 데이터 항목 수의 비율로서 다음과 같이 표현된다.

$$SF_i = \frac{|Q_i - \bigcup_{j \neq i} Q_j|}{|Q_i|}$$

또한 평균 데이터 공유도(SF)는 위의 수식으로부터 다음과 같이 구한다.

$$SF = \sum_{i=1}^N SF_i / N \quad (N: \text{총질의수})$$

데이터 편중도의 영향은 $(1.0 - \alpha\theta)$ 로 표현하는데 여기서 α 는 θ 를 0과 1 사이의 수로 표현하기 위한 정규화 상수이다.

평균 데이터 공유도와 데이터 편중도의 영향을 고려한 Q_i 와 Q_j 의 가중치 상대 거리(WRD_{ij})는 다음의 수식을 사용하여 구한다.

$$WRD_{ij} = (RD_{ij})^{SF * (1.0 - \alpha * \theta)}$$

위 수식의 영향을 살펴보면 WRD_{ij} 는 데이터 공유도가 증가할수록 RD_{ij} 의 영향이 증가하고 데이터 편중도가 커질수록 그 영향이 줄어드는 것을 관찰할 수 있다. 또한 만일 데이터 공유가 존재하지 않는 경우($SF=0$), 또는 데이터 접근 편중도가 매우 큰 경우에는 방송 스케줄링은 RD_{ij} 에 독립적으로 결정되며, 데이터 접근이 균일하다면($\theta=0.0$), WRD_{ij} 는 데이터 공유의 정도(SF)에 전적으로 의존적이 됨을 유추할 수 있다.

4.3 제안 기법

이상의 개념들을 바탕으로 본 논문에서 제안하는 다중 데이터 접근을 위한 계층적 방송 스케줄을 생성하는 절차에 대해서 기술한다.

가. 질의를 구성하는 데이터의 연관 관계

[4]에서의 접근 방법과 같이 Gray 코드를 사용한 데이터 집중을 적용한다. 우선 접근 확률(P_i)의 내림차순으로 N개의 전체 질의를 정렬한 후, 질의 내 모든 데이터 항목에 대하여 이를 접근하는 질의를 식별하도록 길이 N의 비트 열로 표현한다. 이후 각 데이터 항목에 대하여 Gray 값에 따라 정렬한다.

나. 방송 대상 데이터 집합의 선정

4.2.2절에서 기술한 중복 허용 또는 완전 중복 제거 방법을 적용한다

다. 방송 빈도의 결정

동일 방송 간격 개념에 기반하여 [11]에서 제안된 온라인 알고리즘을 적용한다. 여기서 각 기호의 의미는 다음과 같다. C는 현재 시간을 의미한다. 바로 이전에 방송 데이터 집합 BQ_i 가 방송되었다고 가정한다. $R(j)$ 는 방송 대상 데이터 집합 $BQ_j(1 \leq j \leq M)$ 가 방송된 최근

의 시간을 나타낸다. 만일 BQ_j 가 전혀 방송되지 않은 경우 $R(j) = -1$ 로 초기화된다. BQ_j 가 방송될 때마다 $R(j)$ 의 값은 갱신된다. 함수 $G(j)$ 는 다음과 같이 정의된다. 여기서 $C - R(j)$ 는 현재 시간과 BQ_j 가 이전에 방송되었던 시간 사이의 방송 간격을 의미한다.

$$G(j) = (C - R(j))^2 P_j / WRD_{ij}, \quad 1 \leq j \leq M$$

1) 모든 BQ_j 에 대해서 최대 $G(j)$ 를 구하고 이를 G_{max} 라 한다.

2) $G(j) = G_{max}$ 인 BQ_j 를 선택한다. 이때 동일한 값이 하나 이상 있는 경우 임의의 하나를 선택한다.

3) BQ_j 를 모두 방송한다.

4) $R(j) = C$ 로 갱신한다.

위의 수식에서 BQ_j 의 방송 시점을 결정하는 함수 $G(i)$ 에서 $C - R(j)$ 는 방송 간격이다. 만일 BQ_j 가 오랜 동안 방송되지 않은 경우 $C - R(j)$ 가 계속 증가할 것이므로 이 데이터는 WRD_{ij} 에 상관없이 $G(i)$ 의 값이 최대가 되어 결국은 방송 대상에 포함된다. 그러므로 결과적으로 모든 질의가 방송 대상에 포함되는 것을 보장한다. WRD_{ij} 는 질의간에 데이터가 중복되는 정도와 질의의 접근 빈도에 따라 결정된다. 위의 수식에서 WRD_{ij} 가 증가할수록 $G(j)$ 가 감소하므로 BQ_j 의 방송 빈도는 감소하는 것을 유추할 수 있다.

기본적으로는 [11]에서 제안된 온라인 알고리즘과 유사하지만, 이 알고리즘이 개별 데이터 항목에 대해서 방송 여부를 결정하는 반면에 제안된 알고리즘에서는 방송 대상 데이터 집합인 BQ_j 를 동일 시점에서 방송하도록 스케줄링 결정과정에서 고려한다는 점과 방송 빈도를 결정하는데 있어 실제 방송되는 데이터 집합의 크기인 $|BQ_j|$ 대신에 4.2.3절에서 기술한 가중치 상대거리(WRD_{ij})를 사용한다는 점이 다른 점이다.

5. 실험 평가 및 분석

5.1 비교 대상 기법

비교 대상 기법은 데이터 연관성의 표현 여부, 방송 데이터의 선정 방법 및 방송 빈도를 결정하는 요소에 따라 표 1에 나타난 바와 같이 4가지이다. 본 논문에서 제안하는 기법은 VM-NR과 GM-NR이다. 표 1에서 GF-NR은 [4]에서 제안한 기법이며, VM-R과 VM-NR은 [11]에서의 제안 기법을 적용한 것이다. 첫 번째 제안 기법인 VM-R은 [11]의 접근 방법을 그대로 따르는 반면에 VM-NR은 [11]을 변형하여 적용한 것으로서, 일단 질의내의 데이터 항목들을 Gray 코드 값에 따라

표 1 실험 대상 기법의 비교

비교대상기법	데이터 집중 적용	데이터 연관성	방송 데이터 선정	방송 빈도 결정요인
GF_NR[4]	YES	YES	완전 중복 제거	-
VM_R[11]	NO	NO	중복 허용	데이터 길이
VM_NR(제한기법 1)	YES	NO	완전 중복 제거	데이터 길이
GM_NR(제한기법 2)	YES	YES	완전 중복 제거	가중치 상대 코드거리

정렬한 후 방송 데이터 집합을 선정할 때는 질의간의 데이터 중복을 제거한 점이 다른점이다. 그러므로 한 방송 데이터 집합 내에서는 데이터 항목들이 Gray 코드 값에 따라 정렬되나 방송 데이터의 방송 순서를 결정할 때는 방송 데이터 집합간의 연관성은 전혀 고려하지 않고 다만 방송 데이터의 접근 빈도만을 고려한다. 본 논문에서 제안하는 두 번째 기법은 GM-NR으로서 질의의 데이터 항목들을 Gray 코드 값에 따라 정렬하고 데이터 항목의 중복을 제거할 뿐만 아니라 방송 순서를 결정할 때 이들 간의 연관성 역시 고려한다.

5.2 실험 환경 및 고려사항

다양한 실험을 위하여 CSIM을 사용한 시뮬레이터를 구성하였다. 시뮬레이터는 크게 서버와 클라이언트 모듈로 구성된다. 서버 모듈은 각 기법에 따라 방송 스케줄을 생성하고 주기적으로 데이터 항목을 방송한다. 클라이언트 모듈은 질의를 생성하고 질의 내의 모든 데이터 항목을 수신할 때까지 방송 채널을 청취한다. 별도로 언급된 경우를 제외하고 총 질의의 수는 20이며 데이터베이스 크기는 500이다. 질의의 발생 빈도는 Zipf 분포를 따르는 것으로 가정하였다⁵⁾. 정규화 상수 α 는 1/1.6으로 고정하였다⁶⁾. 질의의 크기(즉 질의에서 요구하는 데이터 항목의 개수)는 크게 세가지 유형으로 구분하였다. 첫 번째는 모든 질의가 동일한 개수의 데이터 항목을 갖는 균등길이(U: Uniform), 그리고 $P_i \geq P_j$ 일 때 $|Q_i| \leq |Q_j|$ 인 점진적 증가 유형(I: Increasing), $P_i \geq P_j$ 일 때 $|Q_i| \geq |Q_j|$ 인 점진적 감소 유형(D: Decreasing)의 세가지 경우를 고려한다. $|Q_i|$ 는 평균 길이(LM: Length Mean)와 편차(LS: Length's Standard Deviation)를 갖는

균등 분포 함수를 적용하여 결정된다. 실험에서 질의의 크기 편차는 1.0으로 하였다⁷⁾.

5.3 실험 결과

이 절에서는 실험 결과를 분석 및 비교한다. 비교 대상 접근 방법은 5.1에서 구분한 바와 같이 4가지로서, 실험 결과에서 각 기법은 <적용기법, 질의의 크기 유형, 길이 평균>으로 표기한다. 예를 들어 VM-NR-I4는 VM-NR 기법을 점진적 증가 유형(I)의 길이 분포와 평균 길이 4를 적용한 실험결과임을 나타낸다.

가. 데이터 공유도의 영향

그림 4, 그림 5 및 그림 6은 접근 편중도(θ)가 0.0일 때, 즉 질의의 접근빈도가 균등 분포를 따를 때 SF의 변화에 따른 각 기법의 평균 MDAT를 나타낸다. 질의의 길이 유형에 무관하게 SF가 매우 낮을 때(0.0 ~ 0.3)는 모든 기법은 비슷한 성능을 보인다. 그러나 SF가 커짐에 따라 GF-NR과 GM-NR이 가장 좋은 성능을 보인다. 이 두 가지 기법들은 SF가 증가함에 따라서 질의 간의 연관성을 스케줄링에 효과적으로 반영하기 때문이다. 한편 GF-NR과 GM-NR은 θ 가 0.0인 경우 거의 동일한 성능을 보인다. 반면에 VM-NR과 VM-R은 질의간의 연관성을 전혀 고려하지 않고 있으므로 SF가 커짐에 따라 비효율적으로 스케줄링이 이루어지게 된다.

특히 SF가 커짐에 따라 VM-R 보다 VM-NR이 상대적으로 더 나은 성능을 보이는데 이는 VM-NR의 경우 중복되는 데이터를 제거함으로써 질의의 방송 간격을 줄이는 효과가 있기 때문이다. 그리고 중복을 제거하는 효과는 SF가 커질수록 더욱 증가함을 알 수 있다. 또한 비록 VM-NR에서 한 방송 데이터 집합 내에 포함되는 데이터 항목들에 대해서 순서를 조정하기는 하지만 전체적인 방송 순서를 결정하는데 있어 데이터 집합간의 연관성을 고려하지 않으므로 GM-NR보다 성능이 낮아지는 것을 관찰할 수 있다.

7) 다양한 편차를 사용하여 실험한 결과 전반적인 경향은 매우 유사하였다.

5) Zipf 분포는 균일하지 않은 데이터 접근을 모델링하기 위해 보편적으로 사용되며 접근 편중도 θ 가 증가함에 따라 접근 패턴은 점차적으로 편중되는 양상을 보인다. 데이터 i 의 접근 확률은 $(1/i)^\theta$ 에 비례한다.
6) 데이터 공유의 정도와 데이터 접근 편중도가 방송 빈도에 미치는 상대적인 영향은 정규화 상수 α 에 의해 결정된다. 본 실험에서는 최대 θ 를 α 로 사용하였으며 추후 연구에서 언급하였듯이 보다 정확한 비율의 결정에 대한 추가적인 연구가 필요하다.

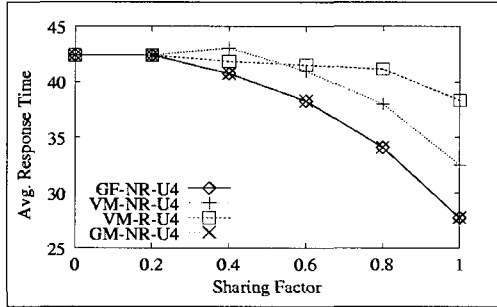


그림 4 데이터 공유도에 따른 성능 비교(균등길이, $\theta=0.0$, $LS=1.0$)

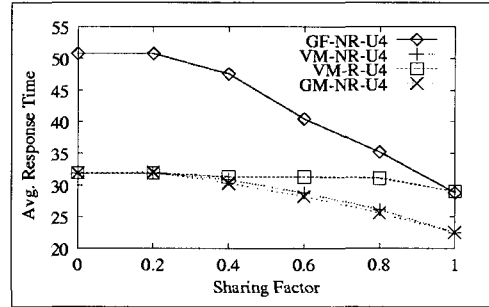


그림 7 데이터 공유도에 따른 성능 비교(균등길이, $\theta=0.8$, $LS=1.0$)

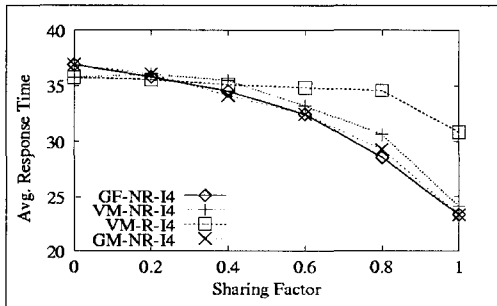


그림 5 데이터 공유도에 따른 성능 비교(점진적 증가, $\theta=0.0$, $LS=1.0$)

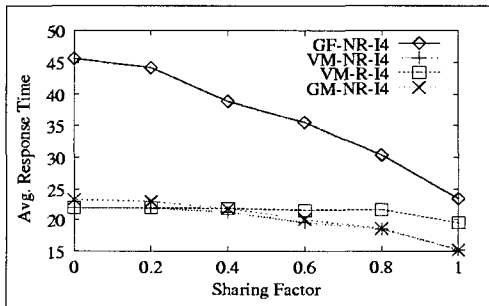


그림 8 데이터 공유도에 따른 성능 비교(점진적 증가, $\theta=0.8$, $LS=1.0$)

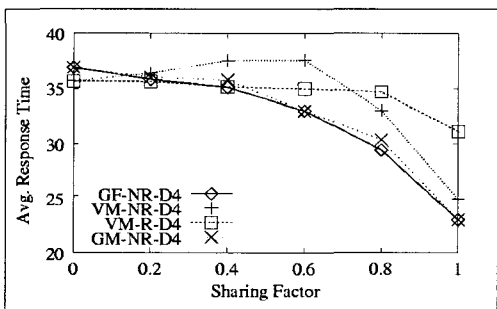


그림 6 데이터 공유도에 따른 성능 비교(점진적 감소, $\theta=0.0$, $LS=1.0$)

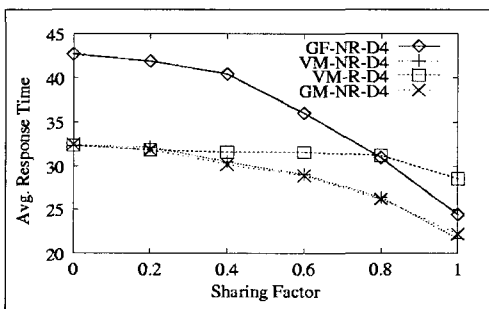


그림 9 데이터 공유도에 따른 성능 비교(점진적 감소, $\theta=0.8$, $LS=1.0$)

그림 7과 그림 8, 그림 9는 θ 가 0.8일 때 SF에 따른 각 기법의 평균 MDAT를 보여준다. 이 결과에서는 이전과는 달리 가장 좋은 성능을 보이는 기법은 VM-NR과 GM-NR 기법임을 알 수 있다. 이는 계층적 스케줄의 구성을 통해 빈번한 접근이 요구되는 데이터 집합의 방송 빈도를 효과적으로 반영하기 때문이다. VM-R 기

법의 경우 중복된 데이터를 방송에 포함시킴으로서 전체적으로 방송 간격이 증가하여 상대적으로 낮은 성능을 나타낸다. GF-NR은 낮은 SF의 영역에서 θ 가 증가함에 따라서 가장 좋지 않은 성능을 보인다. 이는 비록 GF-NR이 데이터의 집중을 통해서 질의간 관련성을 효과적으로 표현하고 있으나 데이터 접근 편중도가 증가함

에 따라 상대적으로 자주 접근되는 질의를 효과적으로 지원하지 못하는 평면 스케줄의 일반적인 한계를 극복하지 못하고 있음을 알 수 있다. 그러나 SF가 증가하면 GF-NR과 다른 기법과의 상대적 성능 격차는 줄어든다. 이상의 결과로부터 다중 데이터의 접근문제는 결국 데이터 집합 간에 얼마나 중복되는 데이터가 존재하는가 하는 데이터 공유의 정도를 효과적으로 다루어야 하는 문제이며, 또한 중복된 데이터를 통제 절차 없이 중복 방송하는 것이 성능 저하의 요인이 되는 것을 단적으로 보여주고 있다. 또한 데이터의 공유도와 데이터 접근 편중도에 따라 기존의 기법이나 이의 변형은 상대적 우위가 달라지는 것을 관찰 할 수 있다. 그러므로 데이터의 공유 정도나 데이터의 편중도와 같은 방송 환경에 적응적으로 대응할 수 있는 새로운 기법의 개발이 필요하며 본 논문에서 제안하고 있는 GM-NR 기법이 변화하는 외부 요인에 효과적으로 적응하는 것을 관찰 할 수 있다.

나. 접근 편중도의 영향

데이터 공유도를 고정시키고 θ 를 변화시키면서 θ 가 미치는 영향을 분석하였다. 그림 10과 그림 11은 데이터

공유도(SF)가 0.0일 때, 질의의 평균 MDAT를 보여준다. 점진적 감소유형의 경우 다른 경우와 유사하므로 생략하였다. GF-NR 기법에 비해 다른 세 가지 기법이 우수한 결과를 보임을 알 수 있다. SF가 0.0이므로 공유되는 데이터가 없어서 VM-R과 VM-NR은 완전히 일치하는 기법이다. 이들 기법에서는 방송 데이터의 길이(BQil)를 방송 빈도의 결정에 사용한다. GM-NR의 경우 SF=0.0이므로 결국 가중치 상대거리의 적용 결과 질의의 길이를 무시하는 효과가 있다. 이러한 이유로 그림 12의 점진적 증가의 경우 약간의(그러나 무시할만한) 성능 차이를 보인 것이다. GF-NR의 경우 SF=0.0인 경우 질의간의 연관성이 전혀 없으므로 단순한 평면 스케줄과 동일하다. 그러므로 θ 가 증가할수록 성능이 저하하게 된다.

반면에 그림 12부터 그림 14는 SF=0.8 인 경우 접근 편중도에 따른 평균 접근 시간의 변화를 보인다. 그림 12는 균등 길이 분포의 경우로서 낮은 θ 에서는 GF-NR 이 VM-R과 VM-NR 기법보다 나은 성능을 보이다가 θ 가 증가함에 따라서 이러한 현상이 역전되는

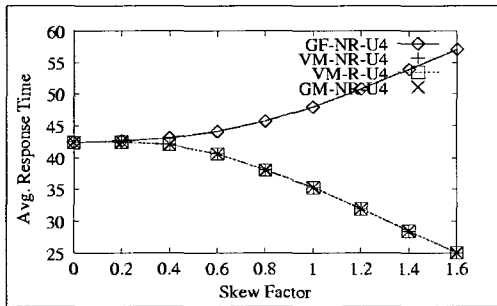


그림 10 데이터 접근 편중도에 따른 성능 비교(균등 길이, SF=0.0, LS=1.0)

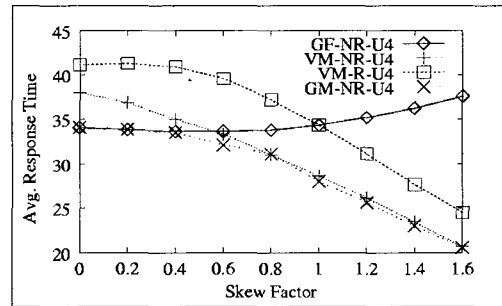


그림 12 데이터 접근 편중도에 따른 성능 비교(균등 길이, SF=0.8, LS=1.0)

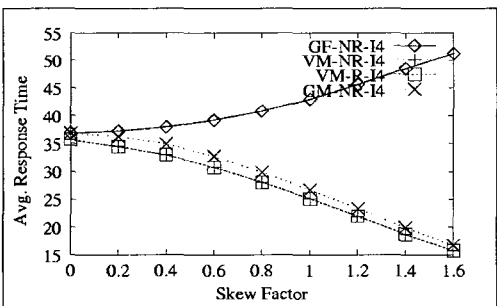


그림 11 데이터 접근 편중도에 따른 성능 비교(점진적 증가, SF=0.0, LS=1.0)

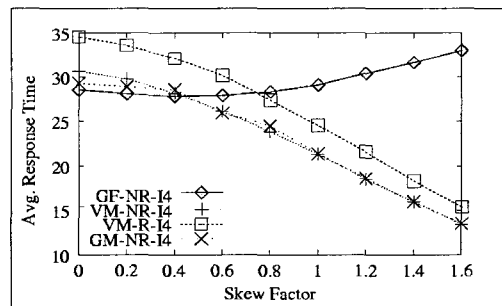


그림 13 데이터 접근 편중도에 따른 성능 비교(점진적 증가, SF=0.8, LS=1.0)

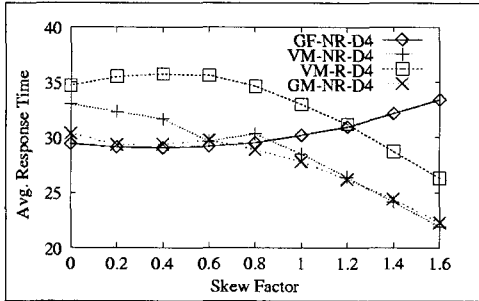


그림 14 데이터 접근 편중도에 따른 성능 비교(접진적 감소, SF=0.8, LS=1.0)

것을 관찰 할 수 있다. 또한 VM-NR이 VM-R 보다 일관성 있게 나은 성능을 보인다. 질의의 길이 유형이 상이한 그림 13, 그림 14의 결과도 비슷한 현상을 보이며, 이 두 기법의 성능 차이는 SF가 증가함에 따라 더욱 커지게 된다. 그러나 흥미로운 사실은 GM-NR 기법의 경우 작은 θ 영역에서는 GF-NR과 유사한 성능을 보이다가 θ 가 증가함에 따라 VM-NR과 유사한 성능을 보이는 점이다. 이러한 현상은 그림 13와 그림 14에서 보듯이 질의의 길이 유형이 달라지더라도 일관성있게 관찰되고 있다. 이러한 현상의 원인은 GM-NR에서의 가중치 상대 코드거리가 질의간에 데이터를 공유하는 정도와 질의가 접근되는 빈도를 모두 고려하고 있으며 또한 환경에 유연하게 적응하는 것을 보여주는 것이다.

다. 질의의 수에 따른 응답시간의 변화

그림 15는 질의 수 변화에 따른 각 기법의 평균 MDT를 보여주는 그래프이다. VM-NR기법은 GM-NR과 거의 유사하므로 비교 대상에서 제외하였다. 여기서 SF=0.5이고 $\theta=1.0$ 이며 상이한 질의 크기에 대해서 점진적 질의 크기 증가의 경우에 대한 결과이다. 각 질의 수가 증가함에 따라서 각 기법의 응답시간 역시

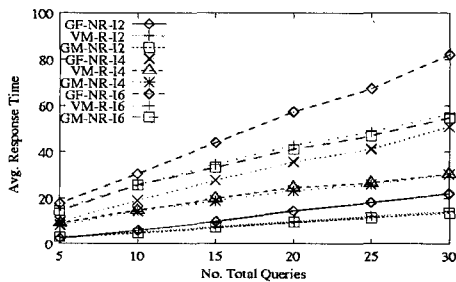


그림 15 질의 수 변화에 따른 응답시간

점진적으로 증가함을 알 수 있다. 또한 질의의 크기가 클수록 응답시간 역시 증가하는 것을 관찰할 수 있다. 각 경우에서 GM-NR과 VM-R 기법이 유사한 성능을 보이는 반면에 GF-NR 기법이 가장 낮은 성능을 보이는 것을 알 수 있다.

6. 결론 및 향후 연구

본 논문에서는 이동 컴퓨팅 환경에서 클라이언트가 여러 데이터 항목을 요구하는 다중 데이터 접근 시간을 감소시키는 계층적 스케줄링 방안에 대하여 기술하였다. 여러 클라이언트에서의 다중 데이터 접근 요구 사이에 데이터가 공유되는 경우 방송 스케줄 내에서 각 데이터 항목의 상대적인 위치가 중요한 영향을 미치며, 다중 데이터 접근 요구가 균등하지 않은 경우 데이터의 방송 빈도를 달리하는 계층적 방송 스케줄이 효과적이다. 이를 위하여 Gray 코드를 기반으로 각 데이터 항목의 상대적 위치를 결정하고, 다중 데이터 요구 사이에 데이터가 공유되는 정도와 데이터 편중도의 영향을 모두 고려하여 방송 빈도를 결정하는 방송 스케줄 구성 방안을 제안하였다.

실험을 통하여 제안 기법의 성능을 평가하였으며 특히 데이터의 공유가 작고 데이터 편중도가 큰 경우에는 기존의 계층적 방송 스케줄에 대응하는 성능을 보이며, 데이터 공유가 크고 데이터 편중도가 작은 경우에는 데이터 요구 사이의 연관성을 고려한 평면 스케줄에 대응하는 성능을 보이는 적응성을 지원함으로써, 제안 기법이 다양한 데이터 공유와 데이터 편중이 일어나는 환경에서 기존의 연구 결과들보다 우수한 성능을 보이는 것으로 관찰되었다. 결국 이동 컴퓨팅 환경에서 효율적인 다중 데이터의 접근을 지원하기 위해서는 요구되는 데이터 집합 사이의 상호 연관성과 데이터 접근의 편중도를 동시에 반영하는 계층적 방송 스케줄의 구성이 매우 효과적임을 실험을 통하여 검증하였다.

추후 방송 데이터의 선정시 질의 간의 중복을 부분적으로 제거하는 방안의 효율성에 대한 연구가 요구되며, 질의들 간의 데이터 공유 정도의 편차가 큰 경우 평균 데이터 공유도가 갖는 오차를 줄이는 방안에 대한 연구도 필요할 것이다. 또한 데이터 공유의 정도와 데이터 접근 편중도가 성능에 미치는 비율을 방송 빈도의 결정 과정에 보다 정확히 반영하는 연구도 필요할 것으로 전망된다.

참고 문헌

- [1] S. Acharya et al., "Broadcast Disks : Data Management for Asymmetric Communication Environments," Proc. ACM SIGMOD Conf. on Management of Data, pp. 199-210, 1995.
- [2] S. Acharya et al., "Balancing Push and Pull for Data Broadcast," Proc. ACM SIGMOD Conf. Proc. on Management of Data, pp. 183-194, 1997.
- [3] A. Acharya, S. Muthukrishnan, "Scheduling On-demand Broadcasts: New Metrics and Algorithms," proc. MOBICOM, 1998
- [4] Yon Don Chung, Myoung Ho Kim, "A Wireless Data Clustering Methods for Multipoint Queries," Decision Support Systems(30), pp. 469-482, 2001.
- [5] Anindya Datta, Aslihan Celik, Jeong G. Kim, Debra E. VanderMeer, Vijay Kumar, "Adaptive Broadcast Protocols to Support Power Conservant Retrieval by Mobile Users," IEEE Int'l Conference on Data Engineering," April 1997, pp. 124-133.
- [6] C. Faloutsos, "Gray codes for partial match and range queries," IEEE Transactions on Software Engineering, Vol. 14, No. 10, pp. 1381-1393, 1988.
- [7] Michael Franklin, and Stan Zdonik. "Data in Your Face: Push Technology in Perspective," ACM SIGMOD Intl. Conference on Management of Data (SIGMOD 98), Seattle, WA, June, 1998.
- [8] T. Imielinski, S. Viswanathan, B. Badrinath, "Energy Efficient Indexing on Air," Proc. ACM SIGMOD Conf. 1994.
- [9] Tomasz Imielinski, S. Viswanathan, "Adaptive Wireless Information Systems" Proc. of SIGDBS Conference, Tokyo, Japan, October 1994.
- [10] K. Stathatos, N. Roussopoulos, J.S. Baras, "Adaptive Data Broadcasting Using Air-Cache," 1st International Workshop on Satellite-based Information Services, Rye, New York, Nov. 1996, pp. 30-37
- [11] N. H. Vaidya and S. Hameed, "Scheduling data broadcast in asymmetric communication environments," Tech. Rep. 96-022, Computer Science Department, Texas A&M University, College Station, November 1996.



이 상 돈

1984년 2월 서울대학교 전자계산기 공학과 학사. 1986년 2월 서울대학교 전자계산기공학과 석사. 1996년 2월 서울대학교 컴퓨터공학과 박사. 1987년 5월 ~ 1997년 8월 한국통신 멀티미디어 연구소 선임연구원. 1997년 9월 ~ 현재 국립목포대학교 정보공학부 조교수. 2001년 1월 ~ 2002년 8월 미국 브라운대학교 객원 교수. 관심분야는 이동 데이터베이스, 스트림 데이터베이스, 웹 데이터 관리