

P2P 네트워크 상의 Spanning Tree 기반 그룹 통신

김희정^o, 손영성

한국전자통신연구원, 컴퓨터 소프트웨어 연구소

Spanning Tree-based Group Communication in P2P Networks

HeeJeong Kim^o, YoungSung Son
Computer Software Research Laboratory, ETRI

요 약

PC의 컴퓨팅 능력과 네트워크 성능이 급속도로 향상됨에 따라 이러한 인프라를 바탕으로 떠오르고 있는 Peer-to-Peer (P2P) 시스템은 자원의 복제 및 공유를 통하여 컴퓨팅 시스템 전반에 걸친 신뢰성과 결함감내 능력을 향상시킨다. 그러나 현재까지 연구된 P2P 네트워크는 각 노드가 복제한 데이터의 일관성 유지를 위해 필요한 그룹통신에 대한 고려가 부족하다. 본 논문은 Spanning Tree 방식의 링 구성을 통하여, P2P 시스템에서의 그룹통신을 위해 설계된 메시지 전송 방식을 설명하고 기존 P2P 메시지 전송방식과의 비교를 통하여 성능향상의 근거를 보인다.

1. 서론

하루가 다르게 사용자 PC의 성능이 향상되고, 네트워크 속도와 대역폭이 발전하는 현실에서, 전통적인 서버 집중식 모델은 현재의 인터넷 자산을 충분히 활용하지 못하고 있다. 그러한 현실인식을 바탕으로 떠오르고 있는 Peer-to-Peer (P2P) 기술은 분산컴퓨팅 기술의 성공적인 상업화를 예견하고 있다. P2P는 불특정 다수가 참여하는 분산시스템으로서, 정보검색과 정보전송, 연산에서의 성능향상뿐만 아니라 컴퓨팅 시스템 전반적에 걸친 신뢰성과 결함 감내 능력을 향상시킬 수 있다. 이러한 결과는 분산 시스템에서 결함 감내, 고성능, 고가용성 등의 효과를 목적으로 자주 사용하는 기법이기도 한 자원의 복제 및 공유 특성 때문이다. 이러한 복제 개체들간의 일관성이 필수적인 응용을 개발함에 있어서의 어려움을 경감시키기 위해 도입된 메카이즘이 그룹 통신 시스템이다 [8].

그룹통신시스템이 제공하는 복제 개체의 일관성 보장은 신뢰성 있는 다중전송 프로토콜을 전제로 하지만, 기존의 그룹통신시스템에 적용되었던 신뢰성 있는 다중전송(Reliable Multicast) 기술들은 P2P 시스템에 적합하지 않다. IP 다중전송은 그 제어와 관리 관련 문제로 인하여 여전히 비현실적이며 [9], 응용수준의 다중전송은 이러한 이슈를 극복하지만 QoS 문제가 있다 [10]. 한편, P2P 망을 구성하고 있는 노드들이 상호 협력하여 파일이나 컴퓨팅 능력을 공유하는 Peer-to-Peer 네트워크가 소개되었으나, 그룹통신에 적합하지 않거나 확장성이 없다. [2]

본 논문은 P2P 시스템에서의 그룹통신을 위하여 설계된 확장성 있는 메시지 전송방식을 설명한다.

2. 관련연구

2.1 P2P 기술과 일반적인 P2P 메시지 전송 방식

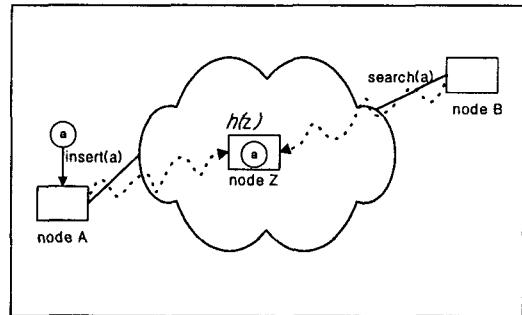
P2P 기술은 불특정 다수가 참여하는 분산시스템

이며 검색의 효율성, 익명성, 저장 비용의 분산 등의 특징을 가진다. P2P 구조에서는 각 클라이언트의 네트워크 상황도 고정적이 아니며 PC의 작업 로드 상황도 안정적이지 않다. 특히 PC의 도메인도 없고 IP가 고정되어 있지 않는 것이 일반적이기 때문에 매 서비스에 새로운 네트워크 상황을 초기화하는 작업 등이 필요하다. 또한, 사용자나 네트워크 상황으로 인한 빈번한 단절상황을 고려해야하는 어려움이 있다. P2P의 파일 공유방식은 Napster [1]와 같이 중앙에 서버에 공유 파일의 인덱스를 두는 방식과 Gnutella [2]와 같이 플러딩(flooding)을 이용해서 파일을 찾는 순수 P2P 방식이 있다. 순수 P2P 방식은 플러딩의 양이 어떤 시점에서부터 감소하기 때문에 실제로 시스템에서 파일을 못 찾을 수도 있다.

각 노드는 전체 시스템의 유일한 해쉬함수($h(x)$)를 가진다. 이 해쉬함수는 자신의 대표 id를 만들기 위해서 사용되고 객체 저장, 검색시에도 사용된다. 일반적인 P2P 네트워크에서 객체 저장 방법은 그림 1과 같다. 파일을 저장하기 위해서는 노드A는 저장할 객체(a)의 특성(객체 이름, 크기, 생성날짜, 그외 정보)를 해쉬함수($h(x)$)를 이용해서 저장위치키(destination key)를 결정한다. 이 저장위치키를 이용해서 이 값을 담당하는 노드 Z를 찾아낸다. 그리고 객체(a)를 노드 Z에 저장한다. 노드B에서 해당 객체(a)를 찾기 위해서는 해쉬함수($h(x)$)를 이용해 알아낸 저장위치키를 이용해서 노드Z를 찾아서 객체(a)를 얻는다.

2.2 그룹통신시스템

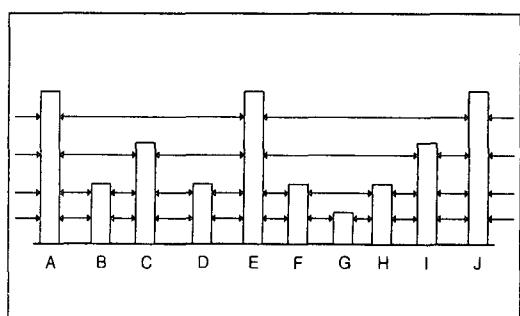
CSCW, 그룹웨어, 네트워크 게임 등 많은 응용에서는 공동의 목적을 위한 공유데이터를 복제함으로써, 시스템의 신뢰성이나 결함감내, 가용성 향상을 도모하지만, 그 복제 데이터의 일관성을 보장하는 것은 상당히 어려운 작업이다. 그룹통신시스템은 현재 일관성을 유지하고 있는 멤버들의 리스트를 관리하면서, 그 멤버쉽이 동적으로 변경될 때마다 멤버들에게 통보하고, 공유 데이터에 대한 업데이트를 모든 멤버들에게 정확히 다중 전송함으로써, 복제 데이터의 일관성을 보장한다 [8].



[그림 1] P2P 파일 저장 방법

3. Magic Square 시스템 모델

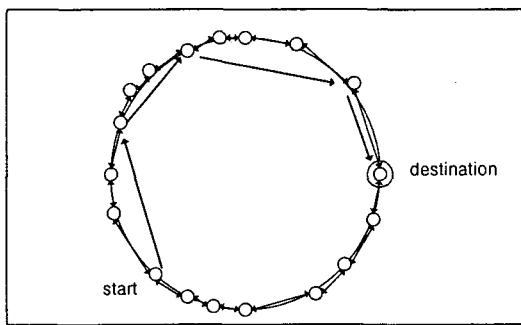
Magic Square 네트워크는 시스템에 참여하는 불특정 다수의 노드들의 접속과 탈락에 영향을 덜 받는 P2P 파일 저장 시스템을 만들기 위해서 고안되었다. 각 노드는 해쉬함수를 이용해서 m비트의 자신의 노드 아이디를 소유한다. 고정 IP를 가진 노드는 IP를 이용해서 노드 아이디를 만들고 가변 IP를 가진 노드는 사용자 지정 정보를 이용해서 노드 아이디를 생성한다. 각 노드는 라우팅 테이블을 이용해 피어를 유지한다. 라우팅 테이블은 각 노드가 시스템에 접속 시에 네트워크에서 자신의 노드 아이디를 중심으로 순차적으로 연결성을 유지하는 테이블이다. 전체 시스템 구성은 그림 2와 같다.



[그림 2] 스kip 리스트로 연결된 피어 집합

Magic Square 네트워크에서 피어의 라우팅 테이블을 이용해서 양방향 스kip 리스트 (bi-directed skip list) 형태로 구성한다. 스kip 리스트는 검색 효율을 위해서 자동 조정 (self-balanced) 기능을 가진 자료 구조이다 [7]. 각 피어는 자신의 네트워크 요구 처리 능

력에 따라서 라우팅 테이블의 크기를 정한다. 그림 3은 피어의 라우팅 테이블을 구성한 경우에 스kip 리스트 형태로 구성됨을 나타낸다. Magic Square 네트워크에서 큰 라우팅 테이블을 가지는 노드 A, E, J는 노드 건너뛰기진행 방법에 사용되어 효과적인 메시지 전송에 사용된다. 노드 A에서 노드 J로의 메시지 전송은 노드 E를 통해서 두번 건너뛰고 노드 B에서 노드 G로의 메시지 전송은 노드 C, E, F를 통해서 네번 건너뛰게 된다.

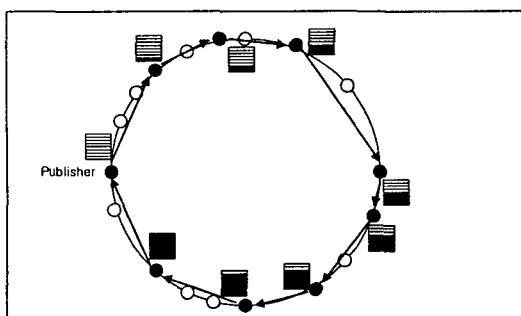


[그림 3] 메시지 전송 방법

4. 그룹통신 메시지 전송 방식

4.1 메시지 전송 방법

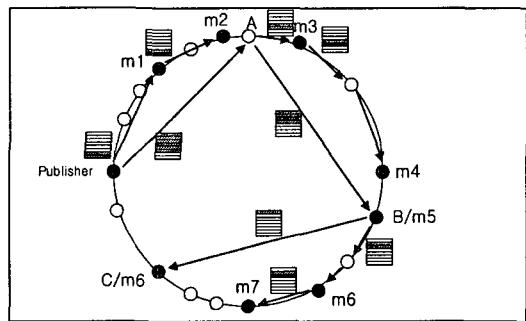
Magic Square 네트워크에서 메시지 전송은 시작 지점(start)에서 목적지점(destination)까지 Peer를 건너뛰며 진행된다. 우선 그림 4와 같이 시작지점에서 라우팅 테이블을 검색하여 목적지점과 가장 가까운 Peer를 정해서 메시지를 전송한다. 메시지를 받은 Peer는 자신의 아이디와 비교한 뒤 라우팅 테이블을 검색해서 위의 작업을 반복해서 수행한다.



[그림 4] 순차 전송 방식

4.2 순차 전송 방식 (Basic Scheme)

그룹 통신 서비스에서 일반적인 메시지 전송은 한 전송 멤버에서 다수의 수신 멤버에게 메시지를 전송하고 모든 멤버가 그 메시지를 받았을 경우에 끝난다(commit). 기본적인 IP multicast 를 이용할 경우는 순차적으로 1대1 통신을 수행하기 때문에 전송 멤버의 블록킹 현상이 발생한다. 이를 해결하기 위해서 Magic Square 네트워크의 특성을 이용한 그림 4의 순차 전송 방식을 설명한다. 메시지 송신 멤버는 메시지의 헤더에 수신 멤버를 모두 기록해서 메시지를 구성한다. 송신 멤버는 메시지 헤더에 있는 멤버를 순차적으로 찾아서 메시지를 시계방향(clock-wise)으로 전송한다. 각 멤버는 메시지를 받아 수신여부를 메시지 헤더에 기록한 뒤 다음 멤버로 재전송한다.



[그림 5] Spanning Tree 방식

4.3 Spanning Tree 전송방식 (Spanning Tree Scheme)

순차 전송 방식은 모든 멤버에 순차적으로 메시지를 전송하기 때문에 Magic Square 네트워크의 효과적인 Peer 건너뛰기 진행방식으로 메시지 전송 방법 사용할 수 없다. 효과적인 전송을 위해서는 큰 라우팅 테이블을 가진 노드만을 건너뛰어야 하나 순차 전송 방식은 작은 라우팅 테이블을 가진 노드(멤버)를 많이 거치게 되어 메시지 재전송이 발생하게 된다. 이를 개선하기 위해서 spanning tree 전송방식을 고안하였다. Spanning tree 전송 방식은 메시지 전송시에 메시지 전송 효율을 높이기 위해서 큰 라우팅 테이블을 가진 노드(A,B,C)에서 그림 5 와 같이 메시지를 나누어서 병렬 전송을 한다. 전체 Magic Square 네트워크는 노드 A, B, C 를 기준으로 크게 영역이 나

뇌진다. Publisher 가 멤버들(m1~m6)에게 메시지를 전송할 경우에 노드 A가 메시지를 받으면 멤버 m3,m4 를 향한 메시지와 노드 B 를 향한 메시지를 나눠서 동시에 전송한다. 노드 B와 노드 C 에서도 동일한 작업을 수행한다.

5. 성능평가

이 장에서는 본 논문에서 소개한 메시지 전송 방식의 성능을 해석해본다. P2P 네트워크에서는 메시지 전송 시에 목적지점에 도착하는데 거친 노드의 개수를 흙(hop) 수로 정하고 이를 평가 기준으로 삼았다. 또한, 동시에 발생하는 메시지의 개수도 중요한 평가 척도가 된다.

N 개의 노드로 구성된 Magic Square 네트워크에서 메시지의 전송시의 평균 흉수는 $O(\log N)$ 이다.

M 개의 멤버로 구성된 그룹에서의 메시지 전송시의 평균 흉수는 다음과 같다. 순차 전송 방식은 항상 하나의 메시지가 전송되나 모든 멤버에게 메시지를 전송하기 위해서는 M 번의 재전송이 발생한다. 따라서 순차 전송 방식에서는 모든 멤버가 메시지를 받아 전송을 종료(commit)하려면 $M * \log N$ 의 흉수가 필요하다.

Spanning Tree 전송 방식에서는 전송중에 메시지가 나눠져서 최대 $O(\text{Height})$ 개의 메시지가 전송된다. 하지만, 모든 메시지가 최종 전달되는 시간은 $\log N$ 이 된다.

[Table 1] The Comparison between two schemes

	Basic	Spanning tree
메시지 전송 경로	Clockwise skiplist sequence	Spanning tree sequence
메시지 재전송 횟수	$O(M)$	Avg: $O(\log N)$ Worst: $O(N)$
동시 메시지 개수	1	$O(\text{Height})$
평균 흉수	$O(\log N)$	$O(\log N/M)$
최종전송흉수	$O(M * \log N)$	$O(\log N)$

6. 결론

본 논문에서는 그룹통신시스템의 기반 기술인

신뢰성 있는 다중전송을 P2P 네트워크에서 구성하는 방법에 대해서 소개하였다. P2P 네트워크에 참여하는 모든 노드를 일정한 방식으로 정렬하고 노드들간의 상호 연결성을 라우팅 테이블로 구성하고 메시지를 재전송하여 신뢰성있는 멀티캐스트 전송방식을 대체 할 수 있는 순차 전송 방식과 Spanning Tree 전송 방식을 소개하였으며, 기존 P2P 메시지 전송방식과의 비교를 통하여 성능 향상을 보였다.

[참고문헌]

- [1] Napster. <http://www.napster.com>
- [2] Gnutella. <http://gnutella.wego.com/>
- [3] I. Clark, et al., "Freenet: A distributed anonymous information storage and retrieval system in designing privacy enhancing technologies," In Proc. International Workshop on Design Issues in Anonymity and Unobservability, LNCS 2009, 2001.
- [4] A. Rowstron, et al., "Pastry: Scalable, decentralized object location and routing for large-scale peer-to-peer systems," 18th IFIP/ACM International Conference on Distributed Systems Platforms, 2001.
- [5] I. Stocia, et al., "Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications," SIGCOMM, 2001.
- [6] S. Ratnasamy, et al., "A Scalable Content-Addressable Network," SIGCOMM, 2001.
- [7] William Pugh, "Skip Lists: A Probabilistic Alternative to Balanced Trees," Communication of ACM, June 1990.
- [8] Gregory V. Chockler, et al., "Group Communication Specifications: A Comprehensive Study," ACM Computing Surveys, Dec. 2001.
- [9] S. Deering, "Host Extensions for IP multicasting," Internet RFC 1112, Available at <http://www.ietf.org/rfc/rfc1112.txt>, 1989
- [10] J. Jannotti, et al., "Overcast: Reliable Multicasting with an Overlay Network," In Proceedings of OSDI, Oct. 2000.