

Flow Grouping 을 통한 P2P 트래픽 분석 방법에 관한 연구

김명섭^o, 강훈정, 홍원기

포항공과대학교 컴퓨터공학과

{mount, bluewind, jwkhong}@postech.ac.kr

요 약

고속 네트워크 기반시설에 힘입어 인터넷 사용자가 전 세계적으로 급증하고, 네트워크 기반의 응용 프로그램이 다양하게 개발되어 사용됨에 따라 네트워크 트래픽이 급격히 증가하고 있다. 이에 따라 네트워크의 현재 상태를 파악하고, 확장 계획을 세우는 기능 및 CRM, SLA 지원을 위한 네트워크 트래픽 모니터링의 중요성은 더욱 더 커져가고 있다. 그러나 Web 이나 FTP, TELNET 등 네트워크 응용프로그램들이 한정적인 과거 환경에서와는 달리 peer-to-peer 응용 프로그램이나 streaming media 응용프로그램 등 네트워크 기반 응용프로그램들이 다양하게 사용되고 있는 현재 인터넷 환경에서 과거의 well-known port 기반으로 응용프로그램을 구분하는 단순한 분석방법으로는 현재의 복잡 다양한 트래픽의 정확한 분석이 어렵다. 본 논문에서는 현재 그 수나 그 구조면에서 다양성이 극도로 증가하고 있는 P2P 트래픽의 효율적인 분석 방법을 제시한다. P2P 트래픽의 분석은 우선 P2P 트래픽에 대한 규정과 기존의 트래픽과 다른 특성을 살펴보고, 이들의 효과적 분석을 위한 알고리즘을 제시한다. 제시하는 알고리즘의 핵심은 인터넷 트래픽의 대부분을 차지하고 있는 TCP 트래픽들을 Flow 단위로 구분하고 중요 port number 결정을 내린 후, P2P Application Port Table 과 Flow Relationship Map 을 통한 P2P 응용프로그램 결정의 3 단계 과정으로 이루어진다. 제한된 방법의 검증은 위하여 기존에 개발된 실시간 트래픽 모니터링 시스템인 NG-MON 에 P2P 분석 시스템을 통합하여 대학내에서 발생하는 인터넷 트래픽에서 P2P 트래픽이 차지하는 양과 각 P2P 응용프로그램 별로 차지하는 분포를 측정하였다.

1. 서론

*최근 인터넷 사용자가 전 세계적으로 급격히 증가하고, 네트워크 기반의 응용 프로그램들이 다양하게 개발되어 사용됨에 따라 네트워크 트래픽이 급격히 증가하고 있다. 이를 지원하기 위해 인터넷 서비스 사업자들(Internet Service Provider, ISP)은 네트워크를 계속해서 증설하고 있으며 이 네트워크들은 기존의 수 메가 (Mbps) 급에서 현재 2.5 Gbps - 10 Gbps 급의 백본 네트워크로 대용량화 되고 있고, 앞으로 몇 년 후에는 테라 (Tera bps)급의 네트워크가 구축될 것이다. 그러나 대부분의 ISP 들은 어떠한 트래픽이 누구에 의해서 얼마만큼 발생되는지에 대한 정확한 정보를 파악하지 못하고 있는 실정이다. 이것은 네트워크의 용량이 대용량화 되어 기존의 트래픽 모니터링 방법으로는 고속 트래픽을 실시간으로 모니터링하기 어려워졌고 기존의 텍스트나 이미지 위주의 트래픽이 스트리밍 미디어 (streaming media) 및 피어투피어 (peer-to-peer, P2P) 위주의 트래픽으로 변하고 있어 이 변화하는 특성에 맞는 모니터링 및 분석 시스템이 존재하고 있지 않기 때문

이다. 이러한 네트워크 트래픽에 관한 정보는 ISP 가 네트워크를 구축하고 증설하는데 그리고 서버 구축 및 다양한 서비스 개발에 있어 절실히 필요한 것이다. 따라서 네트워크상에 구성된 자원의 구체적인 사용 현황을 분석할 수 있는 정보를 제공하는 일은 네트워크 관리에 있어 상당히 중요한 일임에 틀림없다. 그러므로 스트리밍 미디어나 P2P 트래픽이 증가하고 있는 현재의 고속 네트워크 환경에 맞는 적절한 모니터링 및 분석 방법이 제시되어야만 한다.

고속 네트워크에서 발생하는 대용량의 데이터를 실시간으로 모니터링할 수 있는 모니터링 시스템은 여러 가지 방법으로 제시되고 있다. 트래픽 모니터링을 위한 DAG Card[1]와 같이 모니터링 전용 하드웨어들도 개발되고 있고, 라우터와 같은 네트워크 장비에서 모니터링을 위한 데이터를 제공하기 위해 NetFlow[2]나 sFlow[3]와 같은 표준들도 제정이 되어 있다. 그리고 실시간 모니터링을 위한 RTFM[4]과 같은 시스템 구조도 제안되어 많은 모니터링 시스템들[5, 6, 7]에서 참고를 하고 있으며, cluster 구조를 이용하여 실시간으로 대용량 트래픽을 모니터링 및 분석하는 NG-MON[8]과 같은 시스템들도 소개되고 있다.

트래픽 분석에 있어서 앞서 기술한 대용량의 데

* 본 연구는 2002년도 POSCO 의 기술개발비 지원으로 수행 되었음

이터를 실시간으로 효과적으로 처리할 수 있는 시스템 구조 제시에 관한 문제와 더불어 해결되어야 할 또 다른 문제는 복잡 다양해진 트래픽을 어떻게 분석할 것인가의 분석 방법에 관한 문제이다. 일반 사용자들에 의해 발생하는 네트워크 트래픽은 과거에는 FTP, HTTP, Telnet 등의 클라이언트/서버 구조의 응용프로그램이 발생시키는 트래픽이 거의 대부분을 차지하였다. 그러나 몇 년 전 국내의 소리바다[9]나, 국외의 넷스터[10], 뉴텔라[11]와 같이 개인 사용자들이 정보를 서로 개인 대 개인(peer-to-peer)으로 공유할 수 있는 네트워크 응용프로그램들이 개발되면서 네트워크 트래픽은 그 흐름의 방향과 종류에 있어 큰 변화를 가져 왔다. 이렇게 네트워크 응용프로그램이 서비스를 제공하는 서버와 서비스를 제공받는 클라이언트로 명확히 구분되지 않고 네트워크에 접속한 누구나 정보를 제공하는 서버, 정보를 제공받는 클라이언트가 될 수 있도록 개발된 네트워크 응용 프로그램을 P2P 응용프로그램이라고 일컫고 이러한 목적으로 다양한 응용프로그램들이 개발되었고 이들이 차지하는 트래픽의 양은 이미 HTTP, FTP의 양을 훨씬 넘어선 상태이다. 이들 P2P 트래픽에 대한 효율적인 분석은 효율적인 네트워크 관리에 있어 현재 트래픽 현황을 파악하는데 중요한 요소가 되었고, P2P를 홈네트워크나 기업간 정보교환 등 다양한 분야에 적용하는데 있어 기초 자료가 될 것이다. 그러나 복잡 다양한 형태를 보이고 있는 P2P 트래픽에 대한 분석은 과거의 분석방법의 틀을 벗어나지 못하고 있어 정확한 P2P 트래픽의 분석결과를 얻기 힘들다. 본 논문에서는 응용 레벨에서 P2P 트래픽의 특성과 종류를 조사하고, 이들의 분석에 있어 현재 네트워크 트래픽 분석 방법의 한계를 제시한다. 그리고 P2P 트래픽 특성에 맞는 효율적인 분석 방법을 제시하고 이를 적용한 분석 결과를 제시한다.

본 논문의 구성은 다음과 같다. 제 2 장에서는 P2P 트래픽의 정의 및 종류와 특성 및 현재 트래픽 분석 방법의 P2P 분석에 있어 문제점을 기술한다. 제 3 장에서는 효율적인 P2P 트래픽 분석을 위한 Flow Grouping 기반의 P2P 트래픽 분석 방법을 제시한다. 제 4 장에서는 P2P 트래픽 분석 모듈의 구현과 실시간 고속 트래픽 모니터링 시스템인 NG-MON과의 통합에 관하여 기술한다. 제 5 장에서는 P2P 분석 시스템을 통하여 학내 인터넷 트래픽을 분석한 결과에 대하여 기술한다. 마지막으로 제 6 장에서 결론과 앞으로의 연구방향을 제시함으로써 본 논문을 마무리한다.

2. 관련연구

본 장에서는 우선 P2P 트래픽에 대한 개념을 정립하고, P2P 트래픽의 특성을 살펴본다. 그리고 기존의 응용 레벨에서 이루어지는 트래픽 분석 방

법들이 P2P 트래픽 분석에 적용하기에 적합하지 않은 이유를 설명하고, P2P 트래픽 분석에 관한 여러 관련 연구에 대하여 기술한다.

2.1 P2P traffic 에 대한 정의

P2P 트래픽 분석을 위해서는 우선 P2P 트래픽에 대한 명확한 개념 정립이 필요하다. 본 논문에서는 P2P traffic에 대한 정의를 P2P 응용프로그램들에 의해 발생하는 트래픽이라고 정의한다. 그러면 P2P 응용프로그램들은 어떤 프로그램들인가에 대한 정의가 필요하다. 그림 1은 기존 네트워크 응용프로그램의 대표적인 구조인 클라이언트/서버 구조와 P2P 응용프로그램의 구조의 차이를 나타낸 것이다.

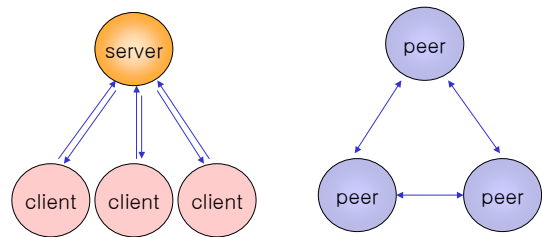


그림 1. client/server 구조 및 P2P 구조

클라이언트/서버 구조에서는 서버와 클라이언트의 명확한 구분이 있고, 서버는 서비스 제공의 역할만 수행하고, 클라이언트는 서비스를 제공받는 역할을 수행하며 클라이언트와 클라이언트 사이에는 직접적인 데이터 교환은 이루어지지 않으며 서버를 통해서만 이루어진다. 데이터의 전달은 클라이언트의 요구와 서버의 응답 방식으로 이루어진다. 그에 비해 P2P 시스템 구조는 peer 들 각각이 서로에게 서비스 제공자이면서 또한 서비스를 이용하는 사용자의 기능을 동시에 수행하는 응용프로그램 구조이다. 즉 peer가 서버의 역할을 수행함과 동시에 클라이언트의 역할을 수행하는 것으로 데이터의 전달은 양방향이다. 데이터의 흐름을 비교해 볼 때 클라이언트/서버 구조는 주로 서버에서 클라이언트로의 데이터 흐름이 대부분을 차지 하였으나 P2P 구조는 서로간에 데이터를 교환하기 때문에 어느 한 peer로 트래픽이 집중되거나 방향성을 띄지 않고 peer들 사이에 자유로운 데이터 교환이 이루어지는 것이다.

2.2 P2P 응용프로그램의 분류 및 특징

현재 P2P 응용 프로그램으로 인식되는 네트워크 응용 프로그램의 형태는 크게 두가지로 분류할 수 있다. 하나는 메신저 응용 프로그램이고 다른 하나는 파일 공유 응용 프로그램이다. MSN 메신저[12], Yahoo 메신저[13]와 같은 메신저 프로그램은 peer들간에 메시지를 주고 받는 기본적인 기능과 1:1 대화, 채팅방을 이용한 그룹 대화 및 파일 전달의 기능에서부터 음성 채팅, 화상 채팅 그리고 응용프로그램 공유에 이르기 까지 다양한 기능을 제공하고 있다.

외국의 웹스터[10]를 잇는 국내의 소리바다[9], 그리고 현재 국내에서 많이 사용되는 당나귀[14]와 같은 파일 공유 프로그램은 기본적으로 사용자가 원하는 파일을 검색하는 기능과 파일 다운로드 기능을 가지고 있으며 추가 기능으로 peer 들 간에 대화할 수 있는 기능을 포함하고 있다. 표 1 은 현재 국내에서 많이 사용하고 있는 P2P 응용 프로그램들의 종류와 기능을 분류한 것이다.

표 1. P2P 응용프로그램 분류

	Instant Messaging Application	File Sharing Application
Function	- Message delivery - 1:1 & multi-chatting - voice & video chatting - application sharing - File transfer	- Searching - File download - Chatting
Application	- MSN Messenger - Yahoo Messenger - Daum Messenger - ICQ, AOL Messenger - ...	- 소리바다 - eDonkey - gnutella - guruturu - ...

이러한 메신저 프로그램과 파일 공유 프로그램은 그 종류가 엄청나게 많고, 앞으로도 새로운 응용 프로그램들이 수 없이 생겨날 것이다. 각각의 P2P 응용프로그램들에서 사용되는 연결구조 즉 P2P 네트워크 구조 또한 다양한 형태를 띄고 있는데 이들의 대표적인 시스템 구조는 그림 2 에서 보는 바와 같이 Central Arbitrer Type 과 Pure Distributed Type 의 두 가지로 볼 수 있다.

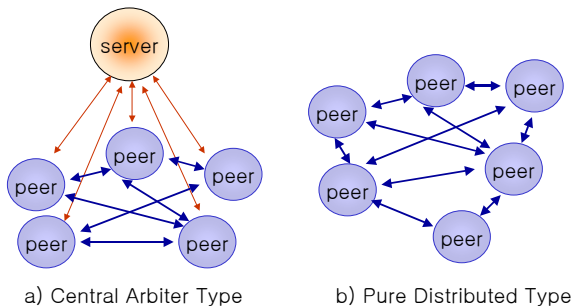


그림 2. P2P 네트워크 구조

Central Arbitrer Type 의 가장 대표적인 것으로 메신저 프로그램들로 이는 중앙의 서버가 peer 들의 등록정보를 관리하고 있으면서 어떤 peer 가 P2P 네트워크에 접속하게 되면 이 peer 와 관련있는 다른 peer 들에게 그 peer 의 접속을 알려줌으로써 peer 들 간에 통신을 가능케 하는 구조이다. Pure Distributed Type 은 Gnutella[11]와 같은 몇몇 파일 공유 프로그램들에서 사용하는 방식으로 중앙의 관리 서버 없이 모든 peer 들이 서로 대등한 관계로 연결되는 구조이다. 이러한 분산 구조에서는 중앙의 관리 서버가 없기 때문에 트래픽이 한 곳에 집중하지 않고 P2P 네트워크가 견고하게 형성되는 장점은 있으나 서로간에 데이터 전송이나 peer 를 검색하는 기능에 있어서는 효율성이 떨어진다. 대부분의 P2P 응용프로그램들은 이 형태를 그대로 따르고 있거나 두 형태를 혼합한 hybrid 형태를 취하고 있다.

현재 이러한 P2P 응용프로그램들이 발생하는 트래픽은 전체 인터넷 트래픽의 상당수를 차지하고 있으며 이는 스트리밍 미디어 트래픽과 더불어 그 양이 계속 증가할 것이다.

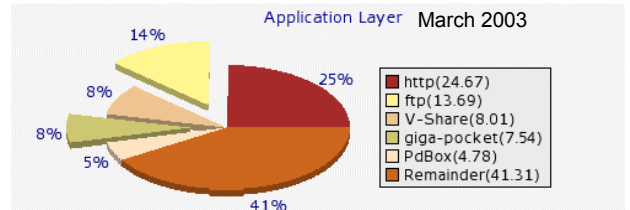


그림 3. 인터넷 트래픽 분포

그림 3 은 대학내의 인터넷 트래픽을 하루 동안 분석한 결과이다. 그림 3 에서 보는 바와 같이 WEB 이나 FTP 가 아직도 많은 분포를 차지하고는 있지만 그 비율은 전체에 비해 40%가 되지 않고, 그 외의 다른 트래픽이 무시할 수 없을 만큼 많이 발생하고 있음을 알 수 있다. 이 나머지 트래픽의 많은 부분을 차지하고 있는 것이 바로 P2P 트래픽이다.

2.3 P2P 트래픽의 분석

이처럼 인터넷 트래픽의 상당수를 차지하는 P2P 트래픽의 분석은 네트워크 트래픽을 분석하는데 있어 상당히 중요한 일이되었다. 그러나 P2P 시스템의 구조에 대한 연구[15]나, 특정 P2P 응용 프로그램의 구조 및 이들이 발생하는 트래픽의 특성에 관한 연구[16]는 간간히 이루어지고 있으나 아직 인터넷 전체 트래픽에서 P2P 트래픽을 분류해 내고 이들 사이의 상관관계나 이들이 어떤 P2P 응용 프로그램에 의해 발생했는지에 대한 분석은 초보 단계를 벗어나지 못하고 있다. 이에 P2P 트래픽을 기존의 트래픽 분석 방법이 그대로 적용하기에는 여러가지 문제를 가지고 있기 때문이다.

그림 4 는 국내에서 많이 사용되고 있는 대표적인 파일 공유 프로그램인 소리바다와 대표적인 메신저 프로그램인 MSN 메신저에서 Peer 간의 통신에 사용되는 전송 계층의 프로토콜 및 port number 들을 나타낸 것이다.

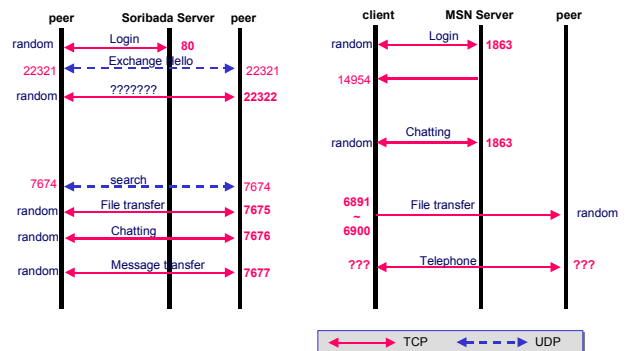


그림 4. 소리바다 및 MSN 메신저의 peer 간 통신

그림 4 에서 보는 바와 같이 P2P 응용프로그램들은 WEB 이나 FTP, TELNET 과같이 한 두개의 고

정된 port number 를 사용하지 않는다. 또한 소리바다와 같은 응용프로그램은 peer 들 간에 hello message 를 보내는 경우와 검색을 수행할 때에는 UDP 를 사용하고, 파일 전송이나 메시지 전송 및 대화의 기능을 수행할 때에는 TCP 를 이용하고 있으며 각 기능들 마다 서로 다른 port number 를 사용한다. 이것은 비단 그림 4 에서 제시하고 있는 응용 프로그램 이외에 많은 다른 P2P 응용 프로그램들 역시 자신들만의 고유한 프로토콜을 사용하고 다수의 well-known port number 를 사용하여 peer 들 간 통신을 수행한다. 프로토콜의 경우는 AOL[18]이나 ICQ[19]에서 사용하는 것처럼 그 포맷이 공개된 경우도 있으나 대부분의 P2P 응용프로그램들은 프로토콜 포맷조차 알려지지 않고, 어떤 경우는 사용하는 port number 조차 무작위로 발생시켜 사용하고 있기 때문에 모니터링 된 트래픽이 어떠한 응용프로그램에 의해 사용되었는지 규명하기란 쉽지 않다.

기존의 대부분의 응용 레벨에서의 트래픽 분석은 port number 를 기준으로 한 분석이었다. 클라이언트/서버구조의 대부분의 트래픽에서는 source port 나 destination port 들 중 하나는 1024 이하의 well-known port 를 사용하고 있고, 이들이 사용하는 port number 의 의미가 IANA[24]에 등록된 용도를 잘 따르고 있기 때문에 well-known port number 와 IANA 등록 정보를 이용하면 쉽게 트래픽이 어떤 응용프로그램에 의해 발생하였는지 쉽게 구분할 수 있었다.

그러나 P2P 트래픽의 경우는 1024 번 이상의 port number 를 사용하고, 그 프로토콜 정보가 IANA 에 등록되지 않은 경우가 대부분이기 때문에 기존의 port number 기반의 분석은 효과적이지 못하다. 또한 하나의 P2P 응용 프로그램들이 그 기능에 따라 다양한 port number 를 사용하는 것은 P2P 트래픽을 구분해 내는데 큰 어려움있다.

P2P 트래픽 분석과 함께 트래픽 분석에 있어 기존의 port 기반의 분석이 어려운 것이 바로 스트리밍 미디어 트래픽의 분석이다.[22] 스트리밍 미디어 트래픽은 제어 트래픽과 데이터 트래픽의 두 부분으로 나누어 볼 수 있는데 제어 트래픽은 WEB 의 80 번 port 와 같이 정해진 port number 를 이용하여 전송되지만 데이터 트래픽은 무작위로 발생시킨 port number 를 이용하여 전송한다. 이 경우는 현재 일반적으로 많이 사용되는 스트리밍 미디어 프로토콜이 MMS[19], RTSP[20]등 몇 가지로 한정되어 있고, 그 포맷이 어느 정도 공개되어 있어 mmdump[21]나 smmon[22]에서와 같이 제어 프로토콜 분석을 통하여 이들이 사용하는 데이터 트래픽의 port number 를 알아냄으로써 데이터를 전송하는 스트리밍 트래픽을 결정할 수 있다. 그러나 그 종류에 있어서나 통신 프로토콜에 있어서 다양성이 극대화 되고 있는 P2P 트래픽의 분석은 이와 같은 방법을 사용하기에는 무리가 있다.

본 논문에서는 이러한 인터넷 트래픽의 상당부

분을 차지하고 있는 복잡 다양한 P2P 트래픽의 효율적인 분석 방법에 관하여 논하고자 한다.

3. P2P 트래픽 분석 알고리즘

본 논문에서는 P2P 트래픽 분석에 있어서 기존의 port number 기반의 분석 방법이 가지는 문제점을 해결할 수 있는 분석 알고리즘을 제시한다. 이 방법은 현재의 P2P 트래픽 분석의 정확성을 높이는 데 기여할 수 있다.

P2P 트래픽의 한 특성이 다양한 port number 들을 사용한다는 것이다. 그런데 이 port number 들이 1024 이하의 well-known port 가 아니라 시스템에서 자동으로 생성되는 random port 와 같은 범위라는 것이다.

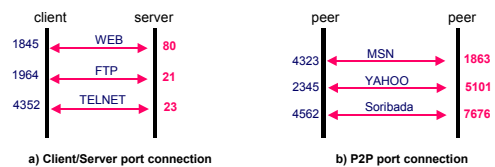


그림 5. P2P 트래픽의 port number 사용상 특성

그림 5 에서 보는 바와 같이 WEB 이나 FTP, TELNET 같은 경우는 각각 80, 21/22, 23 번의 1024 이하의 well-known port 를 사용하고 클라이언트 port 는 1024 이상의 port number 를 사용한다. 따라서 기존에는 이러한 트래픽이 수집되면 1024 이하로 사용되는 port number 만을 조사하여 응용프로그램을 판별할 수 있었다. 그러나 P2P 트래픽의 경우는 클라이언트 port 와 서버 port 모두 1024 이상의 port number 를 사용하고 있기 때문에 이들 중 어떤 port 를 기준으로 응용프로그램의 이름을 결정할 것인가의 문제를 안고 있다. 이를 해결하기 위한 방법으로 본 논문에서는 P2P 트래픽 분석에 있어서 1024 이상의 두 source port 와 destination port 중에서 응용프로그램 결정에 영향을 미치는 중요 port 를 구분하는 방법을 제시한다. MSN 메신저의 경우 접속을 시도하는 peer 는 random port 를 생성하지만 접속을 기다리는 peer 는 port number 1863 를 공통으로 사용하고 있기 때문에 캡처된 패킷이 MSN 패킷이라면 두 port 들 중에서 1863 port 가 중요 port 로 결정되고, 이를 바탕으로 이 패킷이 MSN 트래픽임을 확인할 수 있다. 소리바다 파일전송 트래픽의 경우 7676 번이 중요 port 로 결정되는 것이다.

본 논문에서 제시하는 P2P 트래픽 분석 알고리즘은 이러한 중요 port number 를 패킷에서 결정함으로써 분석이 시작된다. 그림 6 은 P2P 트래픽 분석 알고리즘을 플로우 차트로 나타낸 것이다. 우선 수집된 패킷의 헤더로부터 분석에 필요한 값들을 추출하여 패킷 헤더 정보를 구성한다. 이 패킷 헤더 정보의 내용은 NG-MON[8]에서 규정한 내용을 그대로 따른다. 그리고 이 패킷 헤더 정보들을 바탕으로

Flow 정보를 구성한다. Flow의 정의는 NG-MON에서 정의한 바와 같이 5-tuple 정보(source IP address, destination IP address, source port number, destination port number, protocol number)가 같은 패킷들의 모임이다. 이렇게 구성된 Flow 정보에서 중요 port number를 결정한다. 이 중요 port number가 결정된 flow들을 이들 사이의 시간적, 공간적 상관관계에 따라 Flow Relation Map을 구성한다. 그리고 이 Flow Relation Map의 정보를 바탕으로 Flow는 application 별로 Grouping이 된다. 또한 실제 P2P 응용프로그램에 대한 광범위한 조사를 통하여 각 P2P 응용 프로그램이 사용하는 중요 port number들을 찾아내어 P2P Application Port Table (P2P-APT)를 구성한다. Grouping된 Flow들은 P2P-APT를 바탕으로 해당 P2P 응용 프로그램의 이름을 결정한다.

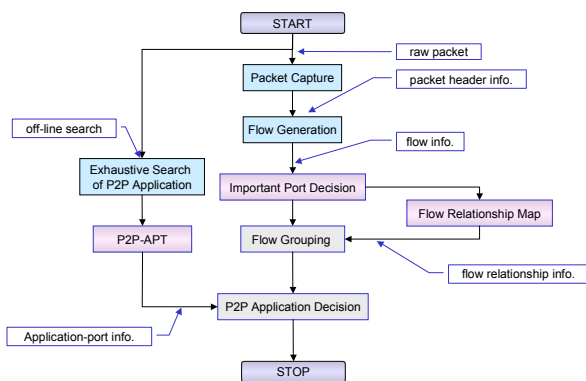


그림 6. 트래픽 분석 알고리즘

그리고 분석의 효율성을 위하여 각 P2P 응용 프로그램들이 사용하는 port number들 중에서 주요 port number를 그 P2P 응용프로그램의 대표 port number로 규정하고 특정 P2P 응용프로그램에 속하는 flow 정보에 이 대표 port number의 tag를 붙여 분석 시스템에서 쉽게 응용 프로그램을 결정할 수 있게 한다.

3.1 중요 port number 결정 방법

이 방법은 대부분의 인터넷 트래픽이 UDP가 아닌 TCP를 이용하고 있고, P2P 응용프로그램 또한 신뢰성 있는 데이터의 전송을 위하여 대부분 TCP를 이용한 데이터 전송을 하고 있다는 사실을 이용한 것이다. TCP 통신은 그림 7에서 보는 바와 같이 데이터의 전송을 맺기 전에 먼저 3-way handshaking을 통하여 두 머신간에 통신 연결을 설립하는 작업을 하게 된다.

이때 서버는 특정 port로 socket를 구성하고, listen mode로 들어가 클라이언트로부터의 연결을 기다리게 된다. 클라이언트는 주로 무작위로 port number를 하나 생성하여 지정된 서버의 port로 SYN 패킷을 보냄으로써 연결 요청이 시작되고, 서버는 SYN-ACK를 보내어 연결 응답을 하고 다시 클라이언트는 ACK 패킷을 보내어 연결을 완료한다.

이때 서버가 클라이언트로부터 연결을 기다리는 listening port가 TCP 통신에 있어 중요한 port number이고 random하게 생성되는 클라이언트 port는 트래픽의 분석에 있어 무의미한 port인 것이다. 따라서 TCP 통신에 있어 SYN 패킷의 경우 destination port number가 그리고 SYN-ACK 패킷의 경우 source port number가 중요한 port number라는 것을 알 수 있다.

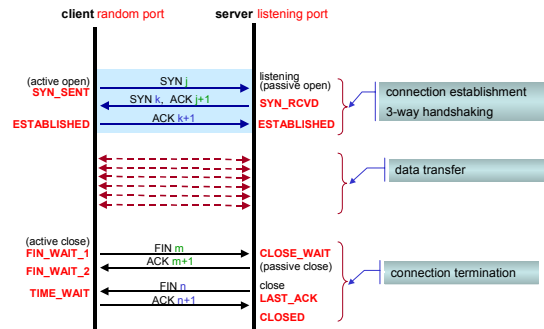


그림 7. TCP 통신 순서

실제 인터넷 환경에 있어 SYN 패킷을 port scan이나 다른 용도로 사용할 수 있기 때문에 SYN 패킷 다음에 반드시 SYN-ACK 패킷이 전달되는 것을 확인함으로써 불순한 트래픽이 P2P 트래픽으로 오인되는 경우등을 제거하는 작업은 실제 구현에 있어 고려되어야 하는 문제이다.

UDP 패킷의 경우는 TCP와 같이 listening port의 개념이 없기 때문에 source port, destination port에서 중요 port를 찾기 힘들다. 이런 경우는 한 host로부터 발생하는 또는 관련된 여러 host로부터 발생하는 UDP flow들간의 상관관계를 바탕으로 중요 port를 결정할 수 있다. 조사 결과 UDP flow의 형태는 그 개수는 많지만 패턴은 단순하게 나타나는 경우가 많아 중요 port를 결정하는 것은 어려운 일은 아닐 것이다. 현재 이에 관한 연구는 진행되고 있다.

3.2 P2P-APT을 통한 flow의 Grouping

TCP/UDP flow에서 중요 port를 결정한 후 이 중요 port를 바탕으로 이 flow가 P2P 트래픽인지의 결정은 P2P Application Port Table (P2P-APT)을 바탕으로 이루어진다. P2P-APT는 각 P2P 응용 프로그램별로 광범위한 조사를 통하여 이들이 사용하는 Port number들을 찾아내고 이들 중 가장 대표적으로 사용되는 port number를 각 P2P 응용프로그램의 대표 port로 결정한 Table이다. 표 2는 현재 국내에서 많이 사용하고 있는 P2P 응용 프로그램들을 Ethereal[23] 패킷 분석 도구를 이용하여 이들이 사용하는 port number들을 조사하여 작성한 P2P-APT이다.

표 2에서 나타난 바와 같이 대부분의 P2P 응용 프로그램들은 1024 이상의 port들을 사용하고, TCP 통신이 대부분을 차지하고 있다. 또한 사용하는 port number의 숫자도 1개에서 많게는 10개 이상을

사용하는 경우도 있었다. 이렇게 광범위한 조사를 통하여 작성된 P2P-APT 를 바탕으로 P2P 응용 프로그램 별로 port number 들의 Group 을 형성한다. 그리고 각 flow 의 중요 port number 를 통하여 이 flow 가 속한 P2P 응용 프로그램을 결정하는 것이다. 이를 통하여 각 flow 들을 P2P 응용 프로그램 별로 grouping 을 하여 P2P 트래픽을 효율적으로 구분 및 분석을 할 수 있다.

표 2 P2P-APT 의 한 예

Application Name	TCP		UDP	
	representative port	well-known ports	representative port	well-known ports
MSN Messenger	1863	1863, 6981-6990, 14594		
Yahoo Messenger	5101	5101, 5050		
Genie Messenger	10000	10000, 10003, 10004		
AIM/ICQ	5190	5190		
Soribada	22322	22322, 7675, 7676, 7677	22321	22321, 7674
eDonkey	4661	4661, 4662, 6667		
Guruguru	9292	9292, 9999, 31200, 22000, 22400, 21700		
V-share	8404	8403, 8404, 1212, 8903, 8908, 8909, 15561		
Shareshare	6399	6399	6777	6398, 6733, 6777

모든 P2P 응용 프로그램들이 서로 다른 port 들을 사용하지 않을 수도 있기 때문에 P2P 응용프로그램의 결정하는 시점에서 어떤 port number 가 서로 다른 두 P2P-APT Group 에 속하게 되는 경우가 발생할 수 있다. 이 경우는 그 flow 와 관련된 두 host 가 발생시킨 다른 flow 들의 정보를 참조하여 해당 flow 의 소속을 결정 짓는다. 이를 위하여 실시간으로 발생하는 flow 정보를 바탕으로 Flow Relationship Map 을 구성한다. 또한 어떤 P2P 응용 프로그램들은 P2P-APT 를 구성할 수 없을 정도로 random 하게 port 들을 생성하여 peer 들간 통신을 하는데 이 경우도 그 flow 를 발생시킨 host 의 다른 flow 들을 참조하여 소속 P2P 응용 프로그램을 결정할 수 있다. flow 와 flow 사이의 연관관계 즉 Flow Relationship 에 관한 연구는 현재 진행 중이다.

4. P2P 트래픽 분석 시스템 설계 및 구현

본 논문에서 제시하는 P2P 트래픽 분석 시스템은 크게 주어진 flow 를 바탕으로 P2P 응용프로그램을 결정하는데 있어 3 개의 중요한 모듈로 구성되며 내용은 그림 8 의 P2P-APT, Port Relationship Mapper, Important Port decider 의 3 개이다.

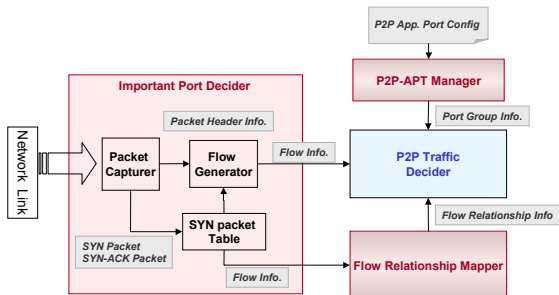


그림 8. P2P 트래픽 분석 시스템 모듈 구성도

Packet Capture 모듈과 Flow Generator 모듈은 Important Port Decider 모듈에 포함되며 이 모듈들은 네트워크의 상태에 따라서 하나의 시스템에 구성될 수도 있고 여러 시스템으로 분산되어 설치될 수도 있다. 예를 들어 라우터나 스위치의 mirroring function 을 이용하는 경우에는 하나의 시스템에 구성될 수 있으나 splitter 를 이용하여 여러 link 를 캡처해야하는 경우에는 Packet Capturer 나 Flow Generator 의 성능에 따라 다수의 시스템으로 구성해야만 한다. Important Port Decider module 은 3 장에서 제시된 알고리즘을 적용하기 위하여 SYN Packet Table 을 구성하여 각 flow 의 방향성을 저장하고 이를 바탕으로 각 flow 별로 중요 port number 를 결정한다. 이렇게 결정된 flow 정보들은 Flow Relationship Mapper 에 전달되고, Flow Relationship Mapper 는 각 flow 별로 상관관계를 규명하고 그 정보를 저장하게 된다. UDP flow 의 경우에 중요 port number 를 결정함에 있어 이 Flow Relation Mapper 의 정보를 이용하게된다.

P2P-APT Module 은 각 P2P 응용프로그램별로 조사된 Port 들을 Grouping 하여 저장하고 있는 곳이다. P2P-APT 에 대한 정보는 그림 9 와 같이 XML file 을 이용하여 저장한다. 여기에서는 각 P2P 응용 프로그램별로 사용하는 port 정보들과 이들의 대표 port 정보를 기록한다.

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<p2p-port-convert>
  <!-- MSN Messenger login & chatting -->
  <p2p protocol="tcp" port="1863" targetPort="1863" /> <!-- msn login -->
  <!-- MSN Messenger file transfer -->
  <p2p protocol="tcp" port="6981" targetPort="1863" /> <!-- msn file transfer -->
  <p2p protocol="tcp" port="6982" targetPort="1863" /> <!-- msn file transfer -->
  <!-- MSN end -->

  <!-- Yahoo Messenger start -->
  <p2p protocol="tcp" port="5101" targetPort="5101" /> <!-- yahoo chatting -->
  <p2p protocol="tcp" port="5050" targetPort="5101" /> <!-- msg transfer -->
  <!-- Yahoo Messenger end -->

  <!-- Genie Messenger start -->
  <p2p protocol="tcp" port="10000" targetPort="10000" /> <!-- Genie login/chatting -->
  <p2p protocol="tcp" port="10003" targetPort="10000" /> <!-- Genie login/chatting -->
  <p2p protocol="tcp" port="10004" targetPort="10000" /> <!-- Genie login/chatting -->
  <!-- Genie Messenger end -->
</p2p-port-convert>
```

그림 9. P2P Application Port Table 설정 file 의 내용

P2P Traffic decider 는 각 Important Port Decider 로 부터 전달된 flow 를 가지고, P2P-APT 와 Flow Relationship Mapper 의 정보를 바탕으로 그 flow 가 어떤 P2P 응용프로그램에 속하는지 결정한다.

4.1 NG-MON 실시간 트래픽 모니터링 시스템

P2P 트래픽 분석 시스템의 구현은 트래픽 모니터링 시스템인 NG-MON[8]의 한 모듈로 구성한다. NG-MON 은 본 연구팀에서 구현한 고속 네트워크 상에서 실시간 트래픽 모니터링 시스템 구조로 그림 10 과 같이 모니터링 시스템의 기능을 5 단계로 구분하여 각 단계별로 모니터링 시스템의 성능과 네트워크 환경에 따라 다수의 시스템들을 클러스터

로 구성하고 각 단계 사이는 Pipeline 개념을 적용하여 고속 네트워크 상황에 맞게 유연하게 구현할 수 있게 설계되었다.

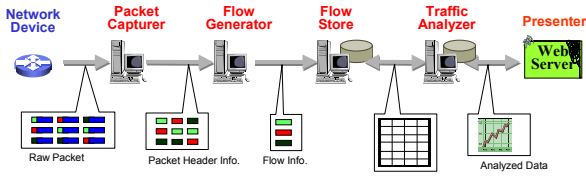


그림 10. NG-MON 구조

5 가지 단계는 Packet Capture Phase, Flow Generation Phase, Flow Store Phase, Traffic Analysis Phase, Presentation Phase 로 구성되며, 앞의 3 단계는 트래픽 분석 목적에 관계없이 동일하며 뒤의 2 단계는 분석 목적에 따라 다양하게 구성될 수 있다. 현재의 NG-MON 의 구현은 각 host 별, 사용자별, subnet 별, 시간별, protocol 별로 트래픽 사용량에 대한 추이를 분석하여 WEB 을 통하여 결과를 나타내는 Throughput Analysis 기능을 수행하고 있다. 그리고 이 시스템은 현재 학내 인터넷 접속부에 설치되어 학내에서 발생하는 인터넷 트래픽을 실시간으로 모니터링하고 분석하는 기능을 담당하고 있다.

4.2 NG-MON 과 P2P 분석 모듈의 통합

P2P 트래픽 분석 시스템은 NG-MON 시스템의 한 모듈로 통합되어 Throughput 분석에 있어 Application 별 트래픽 분포를 파악하는데 도움을 주고 있다. NG-MON 과 P2P 트래픽 분석 모듈의 통합은 그림 11 에서 보는 바와 같이 Packet Capture 와 Flow Generator 에 Important Port Decider 모듈들이 들어가게 되고, P2P-APT 와 Flow Relationship Mapper 는 Flow Store 에 통합되어 Traffic Analyzer 가 Application level 분석이 이루어지기 전에 P2P 트래픽에 대한 Flow Grouping 을 완료한다.

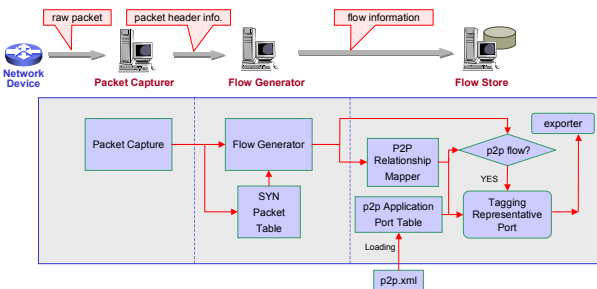


그림 11. P2P Analysis Module Integration into NG-MON

그리고 Traffic Analyzer 는 P2P-APT 정보를 참조하여 Grouping 된 flow 들이 소속된 P2P 응용 프로그램의 이름을 결정하게 된다. 통합에 있어 기존 NG-MON 의 기능을 최대한 활용하는 방향으로 필요한 추가 기능만을 구현하는 것으로 이루어 졌다.

그림 12 는 NG-MON Throughput Analyzer 를 통하여 분석된 Application Level 에서 P2P 트래픽 분석

결과를 나타내는 화면이다.

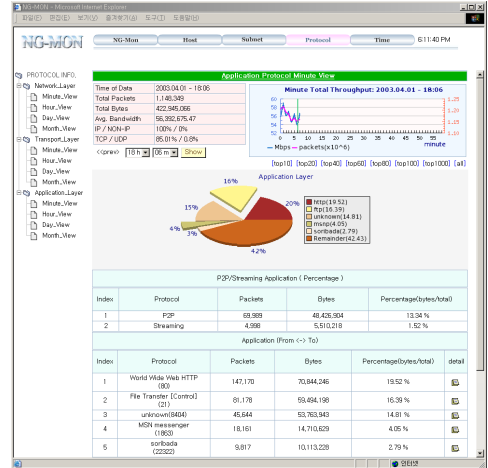


그림 12 P2P 트래픽 분석 결과 화면

5. P2P 트래픽 분석 결과

본 장에서는 본 논문에서 제시된 P2P 트래픽 분석 시스템을 통하여 학내 인터넷 트래픽을 응용 프로그램 레벨에서 분석한 결과를 설명한다. 분석은 학내 인터넷 연결부에서 하루 동안 학내에서 발생한 인터넷 트래픽을 분석한 것이다.

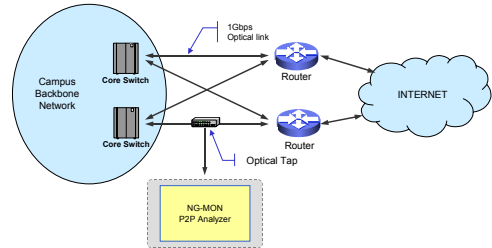


그림 13 학내 인터넷 연결 구조

학내의 인터넷 연결구조는 그림 13 에서 보는 바와 같이 2 개의 100Mbps 메트로 이더넷으로 구성되어 있으며 두대의 Core 기가스위치를 거쳐 두대의 라우터가 매쉬구조로 연결되어 있다. 코어 스위치와 라우터 사이는 Gigabit Ethernet 으로 연결되어 있으며 데이터의 수집은 Optical Tap 을 이용하여 수집하였다. 실험에서 수집된 데이터는 네개의 Gbps Ethernet Link 중에서 하나를 선택하여 하루 동안 수집한 데이터를 분석한 값이다. 네개의 Gbps Ethernet Link 를 모두 수집할 수 있는 방안을 현재 구축 중이며 이것이 완료되면 좀 더 정확한 분석 자료를 제공할 수 있을 것이다.

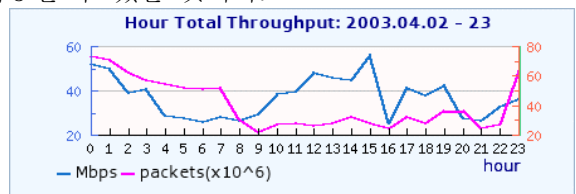


그림 14. 하루동안 트래픽 변화 추이

하루동안 수집된 총 패킷의 수는 960,710,117 이고, 총 bytes 수는 418,434,615,465 였다. 그리고 평균 bandwidth 는 38,743,945.88 였으며 TCP 와 UDP 의 비율은 82.3% : 9.9%였다. 그림 14 는 하루 동안 수집된 트래픽의 시간적인 추이를 보여준다. P2P 트래픽 분석을 위하여 우리는 현재 국내에서 많이 사용하고 있는 P2P 응용프로그램 18 가지의 동작형태를 조사하여 P2P-APT 을 구성하였다. 이 18 가지 P2P 응용 프로그램에는 MSN 메신저 프로그램을 비롯하여 소리바다와 같은 국내외 유명한 메신저 프로그램과 파일공유 프로그램들이 포함된다.

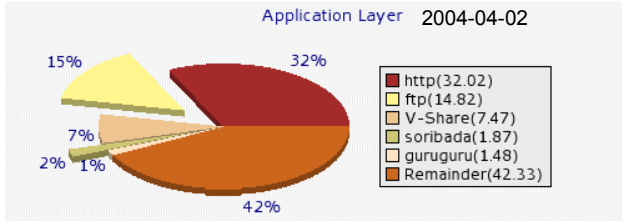


그림 15 응용 레벨 트래픽 분포

그 결과 그림 15 에서 보는 바와 같이 아직도 HTTP 와 FTP 가 가장 많은 분포를 차지 하고 있지만 3 위는 8404 번 port 를 사용하는 V-share[25] 트래픽이 차지 하였다. 4 위는 소리바다, 5 위는 구루구루가 차지 하였다. 그리고 나머지 부분을 차지하는 42.33%의 트래픽들 중에 상당수도 역시 P2P 가 차지할 것이라고 여겨진다. 그리고 전체 트래픽 중에서 조사된 18 가지 P2P 응용프로그램에 해당하는 트래픽은 전체의 15.82%를 차지하였다.

이 분석결과는 하루동안의 분석 데이터이기 때문에 전체 인터넷 트래픽에서 P2P 가 차지하는 분포라고 말할 수는 없지만 앞으로 P2P 트래픽의 증가 추세를 가늠할 수 있는 지표로서 역할은 충분히 할 것이다. 좀 더 정확한 분석을 위해 현재 1 주 일이나 1 달 정도의 장기간 데이터를 바탕으로 P2P 트래픽의 상세한 분포를 조사할 예정이다.

6. 결론

네트워크가 고속화되고 이를 이용하는 응용프로그램의 수가 다양해지고 이를 이용하는 응용프로그램의 형태도 클라이언트/서버 구조에서 모든 peer 가 직접 서비스를 제공하고, 제공 받을 수 있는 P2P 구조로 변화하고 있다. 이러한 P2P 응용프로그램들이 발생시키는 트래픽은 WEB 이나 FTP 와는 다른 형태를 띄고 있어 기존의 단일 port number 만을 기준으로 응용 프로그램을 판별하는 방법으로 그 분석이 어렵다. 본 논문에서는 P2P 트래픽이 기존의 클라이언트 서버 구조의 트래픽과는 다른 특성을 정리하였고, 이들을 분석해 내기 위해 고려해야 할 고려사항들을 언급하였다. 그리고 이 광범위한 P2P 응용 프로그램들의 조사를 통한 Flow Grouping 과 flow 에서 중요 port number 결정 방법, 그리고 flow

간의 연관관계 구성등을 통한 P2P 트래픽의 효율적인 분석 방법을 제시하였고 이 분석 방법을 바탕으로 현재 학내에서 발생하는 인터넷 트래픽의 분포를 조사 함으로써 제시된 분석 방법의 타당성을 검증하였다.

본 논문에서 제시하고 있는 P2P 트래픽 분석 방법은 아직 다양한 P2P 응용프로그램에 대하여 효과적으로 적용하기에는 미흡한 점이 많이 있다. P2P 네트워크에 접속하면서 random port 를 생성하여 연결하는 경우는 Flow Relationship Map 을 통하여 Flow 에 대한 Grouping 을 하여야만 하는데 이에 대한 연구는 아직 초보 단계를 벗어나지 못하고 있다. 그리고 본 논문에서 제시한 분석 방법은 패킷의 payload 는 전혀 고려하지 않기 때문에 속도나 처리해야 하는 데이터 면에서 효과적이라고 할 수 있지만 100% 정확한 트래픽 분류작업을 위해서는 payload 를 이용하는 방법도 고려되어야 할 것이다. 본 논문을 위해서 구현된 시스템은 prototype 수준으로 실제 인터넷 트래픽의 장시간 분석을 위한 안정된 시스템의 구현과 지속적인 트래픽 분석을 통하여 좀 더 효율적인 P2P 트래픽 분석 방법에 대한 연구등을 앞으로 지속적으로 해 나가야 할 것이다.

[참고 문헌]

- [1] Ian D Graham and John G Cleary, "Cell level measurements of ATM traffic," Proceedings of the Australian Telecommunications Networks and Applications Conference, pp. 495-500, December 1996.
- [2] Cisco, White Papers, "NetFlow Services and Applications," http://www.cisco.com/warp/public/cc/pd/iosw/ioft/neflct/tech/napps_wp.htm.
- [3] P. Phaal, S. Panchen, N. McKee, "InMon Corporation's sFlow: A Method for Monitoring Traffic in Switched and Routed Networks", IETF RFC 3176, September 2001.
- [4] N. Brownlee, C. Mills, G. Ruth, "Traffic Flow Measurement: Architecture", IETF RFC 2722, October 1999.
- [5] N. Brownlee, "Traffic Flow Measurement: Experiences with NeTraMet", IETF RFC2123, March 1997.
- [6] Ken Keys, David Moore, Ryan Koga, Edouard Lagache, Michael Tesch, and k claffy, "The Architecture of CoralReef: An Internet Traffic Monitoring Software Suite," PAM Workshop 2001, April 23-24, 2001.
- [7] Argus, <http://www.qosient.com/argus/>.
- [8] Se-Hee Han, Myung-Sup Kim, Hong-Taek Ju and James W. Hong, "The Architecture of NG-MON: A Passive Network Monitoring System", LNCS 2506, DSOM 2002, October, 2002, pp16-27.
- [9] 소리바다, <http://www.soribada.com/>.
- [10] Napster, <http://www.napster.com>.
- [11] Gnutella, <http://gnutella.wego.com>.
- [12] MSN Messenger, <http://messenger.msn.co.kr/>,

Microsoft.

- [13] Yahoo Messenger, <http://kr.messenger.yahoo.com/>.
Yahoo.
- [14] eDonkey2000, <http://www.edonkey2000.com>.
- [15] Matei Ripeanu, "Peer-to-Peer Architecture Case Study: Gnutella Network", techreports TR-2001-26, University of Chicago, July, 2001.
- [16] Subhabrata Sen, Jia Wang, "Analyzing Peer-to-Peer Traffic Across Large Networks", IMW2002 Workshop, 2002.
- [17] AOL, <http://www.aol.com/>.
- [18] ICQ, <http://web.icq.com/>.
- [19] Microsoft, Windows Media Technology,
<http://www.microsoft.com/windows/windowsmedia/default.asp>.
- [20] H. Schulzrinne, A. Rao, and R. Lanphier, "Real Time Streaming Protocol (RTSP)," RFC 2336, April 1998.
- [21] Jacobus van der Merwe, Ramon Caceres, Yang-hua Chu, and Cormac Sreenan "mmdump- A Tool for Monitoring Internet Multimedia Traffic," ACM Computer Communication Review, 2000.
- [22] Hun-Jeong Kang, Hong-Taek Ju, Myung-Sup Kim and James W. Hong, "Towards Streaming Media Traffic Monitoring and Analysis", APNOMS2002, 2002, pp 503-504.
- [23] Ethereal, <http://www.ethereal.com/>.
- [24] IANA, <http://www.iana.org>.
- [25] V-share, <http://www.v-tv.co.kr/>.