

BitTorrent Swarming System에서의 가용성 Metric에 대한 분석

정태중, 한진영, 김현철, 권태경, 최양희
서울대학교 컴퓨터 공학부

{tjchung, jyhan, hkim}@mmlab.snu.ac.kr, {tk, yhchoi}@snu.ac.kr

Analysis on Availability Metrics in BitTorrent Swarming System

Taejoong Chung, Jinyoung Han, Hyunchul Kim, Ted "Taekyoung" Kwon, Yahnghee Choi
School of Computer Science and Engineering, Seoul National University

요약

BitTorrent는 파일 공유를 목적으로 하는 Peer-to-Peer 방식의 소프트웨어로, 현재 인터넷 트래픽의 약 27~55%를 차지할 정도로 널리 이용되고 있다. 하지만 flash crowd 상황 이후에 공유집단(이하 swarm)내에서 원하는 파일을 얻지 못하거나 어려워지는 가용성 문제가 심각하다고 밝혀졌다. 이에 따라 BitTorrent에서의 가용성 문제에 대한 연구가 많이 진행되었지만, 각 연구의 관점이나 측정 방법에 따라 각각 가용성을 다르게 정의하고 있다. 본 논문에서는, 대표적으로 사용되는 가용성 정의들이 실제로 swarm의 가용성을 얼마나 잘 반영하는지 분석한다. 분석을 위해 swarm의 모든 peer의 정보까지 측정하는 대규모 measurement 연구를 18일 동안 진행하였고, 약 4만개의 swarm에서 얻어진 210만개의 기록을 분석하였다. 분석 결과 swarm의 전체 peer로부터 수집한 Piece Map 정보를 바탕으로 계산된 가용성(본 논문에서는 Ground Truth로 사용됨)에 비해 대표적으로 사용되는 가용성 정의인 Busy Period Availability와 Seed Availability로 계산된 가용성이 보다 낮게 측정되는 것을 관측하였다. 이것은 각 가용성의 정의가 leecher가 보유한 piece의 유용 가능성을 반영하지 못하였기 때문으로, swarm 전체의 상황을 잘 반영하는 가용성에 대한 정의가 필요함을 나타낸다. 본 연구는 swarm의 가용성 정의들을 평가하는 첫 번째 연구로서의 의의를 갖는다.

I. 서론

BitTorrent[1]는 파일 공유를 목적으로 하는 Peer-to-Peer (이하 P2P) 방식의 소프트웨어이다. IPOQUE에서 최근 발간한 보고서에 따르면, BitTorrent는 현재 인터넷 트래픽의 약 27~55%를 차지할 정도로 널리 이용되고 있다[2]. BitTorrent의 이러한 성공은 사용자들이 폭발적으로 몰리는 flash crowd 상황에서도 효과적으로 파일을 공유할 수 있는 확장성 (scalability)에 기인한다. 하지만 flash crowd 상황 이후에 시간이 점차 지나감에 따라 원하는 파일을 더 이상 얻지 못하거나 얻기 어려워지는 가용성 (availability) 문제가 BitTorrent에서 심각하다고 밝혀졌다[3].

이에 따라 BitTorrent에서의 가용성 문제를 분석하고 해결하기 위해 많은 연구들이 진행되어 왔다[3][4]. 하지만 해당 연구들에서는 각 연구의 관점과 측정 방법에 따라 가용성을 다르게 정의하고 있다. 대표적인 정의 방법으로는 전체 swarm에 seed가 있는지 여부로 가용성을 판단하는 방법이 있다[4]. 또한 데이터를 받기 위해 swarm에 참여한 시간 중 실제로 데이터를 받은 시간 비율을 가용성의 정의로 사용하기도 한다[3].

본 논문에서는, 대표적인 가용성 정의들이 실제로 swarm의 가용성을 얼마나 잘 반영하는지 분석한다. swarm의 가용성은 파일의 공유 가능 정도를 나타내주는 척도로, swarm내에서 공유할 수 있는 piece¹가 얼마나 있는지에 따라 결정된다. 하지만

swarm내의 모든 piece의 분포를 알아내는 것은 쉽지 않기 때문에 일종의 근사치로 위의 가용성 정의 방법들이 사용되는 것이다. 본 논문에서는 swarm의 모든 peer의 정보까지 측정하는 대규모 measurement 연구를 통해 각 가용성 정의가 실제 swarm의 가용성을 얼마나 잘 반영하는지 분석한다. 본 연구는 swarm의 가용성 정의들을 평가하는 첫 번째 연구로서의 의의를 갖는다.

II. Swarm의 가용성 정의

2.1 Seed Availability

대부분의 가용성 관련 연구들은 swarm안에 데이터의 모든 piece를 소유하고 있는 seed의 존재여부로 swarm의 가용성을 측정한다[4]. 즉, seed가 하나라도 존재를 한다면 해당 데이터를 유용할 수 있기 때문에 swarm내에 seed가 존재하는지 여부로 가용성을 판단하는 것이다. 우리는 이를 Seed Availability라고 부르고, swarm내 모든 peer 집합 X 에 대하여 (1)과 같이 표현 가능하다.

$$\text{Seed Availability} = \overline{f(x_i)} \quad \dots (1)$$

$$\left(\text{단, } x_i \in X \mid f(x_i) = \begin{cases} 1 & \text{if there is a seed} \\ 0 & \text{else} \end{cases} \right)$$

2.2 Busy Period Availability

[3]에서는, 유저가 데이터를 받기 위해 swarm에 참가한 시간을 실제로 다운을 받는데 소모한 시간인 Busy Period와 가용할 수 있는 piece가 존재하지

본 연구는 기초기술연구회의 NAP 과제 지원으로 수행되었음. 이 연구를 위해 연구장비를 지원하고 공간을 제공한 서울대학교 컴퓨터연구소에 감사 드립니다.

¹ BitTorrent에서는 file을 여러 개의 piece로 나누어서 공유한다.

않음으로 다운받지 못하는 유휴 상태인 Idle Period 로 나눈다. 그리고 이에 따라 swarm 의 가용성을 전체 swarm 에 참가한 시간 대 Busy Period 의 비율로 나타낼 수 있다. 우리는 이를 Busy Period Availability 로 부르고, 이는 (2)와 같이 표현할 수 있다.

$$Busy\ Period\ Availability = \frac{BusyPeriod}{BusyPeriod + IdlePeriod} \dots (2)$$

2.3 Piece Map Availability

BitTorrent 의 실제 공유 단위는 전체 파일이 아니라 파일을 여러 조각으로 나눈 piece 이다. 그리고 유저는 자신이 가지고 있는 piece 종류를 나타내는 Piece Map 을 관리한다. Swarm 의 모든 유저의 Piece Map 정보들을 조합하면 swarm 에서 유용 가능한 piece 분포를 알 수 있게 된다. 우리는 이를 Piece Map Availability 라 부른다. 본 가용성은 swarm 의 모든 peer 의 Piece Map 을 받아와야 하기 때문에 측정하는 것이 쉽지 않지만, swarm 의 piece 분포를 직접적으로 반영하기 때문에 이 가용성 정보를 Ground Truth 로 사용하였다. 본 가용성은 swarm 내 모든 peer 집합 X 에 대하여 (3)과 같이 표현 가능하다.

$$Piece\ Map\ Availability = \frac{1}{M} f(x_i) \dots (3)$$

(단, $x_i \in X \mid f(x_i) = \text{piece의 개수}, \max[f(x_i)] = M$)

III. 측정 결과 및 분석

각 가용성의 정의가 실제 swarm 의 가용성을 얼마나 잘 반영하는지 분석하기 위해, swarm 의 모든 peer 의 정보까지 측정하는 대규모 measurement 연구를 진행하였다. Measurement 를 위해 본 연구진에서는 현재 널리 사용되는 BitTorrent 클라이언트인 Azureus [5]을 수정하여 사용하였다. 이때, swarm 내의 각종 정보를 tracker 에서 뿐만 아니라 Peer Exchange² 를 통해 찾은 peer 들에게서도 받아오도록 하였다. <표 1>에 measurement 환경이 요약되어 있다.

<표 1> 측정 기록

항목	값
측정 기간	2010-04-30 ~ 2010-05-16
Swarm 개수	약 4 만개
기록 개수	약 210 만개
평균 seed 개수	0.68 개
평균 leecher 개수	0.81 개

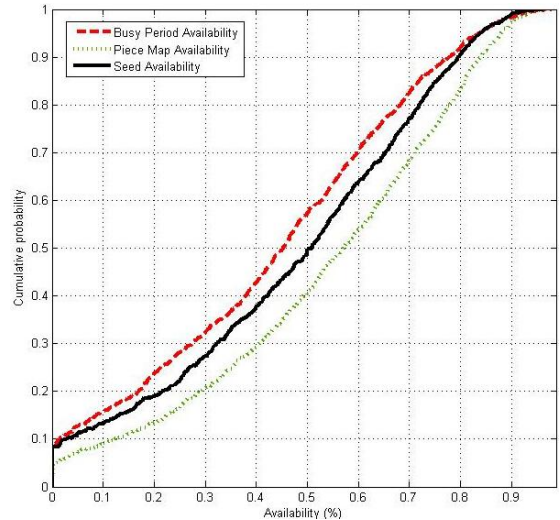
관찰한 swarm 의 가용성을 세가지 방법 (Busy Period, Seed Availability, Piece Map Availability)으로 구해서 그 결과를 <그림 1>에 나타내었다. <그림 1>에서 볼 수 있듯이 가용성이 50% 미만인 swarm 의 비율이, Busy Period Availability 로 구하였을 때는 58%, Seed Availability 로 구하였을 때는 50%로 나타났다. Piece Map Availability 의 경우에는 40% 정도가 가용성이 50% 미만이었는데, 이는 Busy Period Availability 나 Seed Availability 가 실제 swarm 의 가용성보다 낮게 평가되고 있다는 것을 나타낸다. 즉, 예를 들어 Busy Period Availability 의 경우 가용성이 50% 미만인

² 연결된 peer 끼리 각각 알고 있는 다른 peer 들의 위치를 공유하는 것을 말한다.

swarm 이 대부분 (약 60%)를 차지하기 때문에 가용성 문제가 심각하다고 주장할 수 있는데, 실제로 swarm 의 가용성은 약 40%에 해당되기 때문에 심각한 정도가 다르게 판단될 수 있다.

Busy Period Availability 와 Seed Availability 가 실제 swarm 의 가용성과 차이가 있는 이유로는, BitTorrent 에서 데이터를 공유할 때 seed 뿐만 아니라 leecher 에게서도 데이터를 받을 수 있는데, leecher 가 갖고 있는 piece 의 유용 가능성이 반영되어 있지 않기 때문이다.

<그림 1> Swarm 가용성의 누적 분포 함수



IV. 결론 및 향후 과제

본 논문에서는 대규모 measurement 연구를 통해 얻어진 결과를 바탕으로 BitTorrent 의 대표적인 가용성 정의가 실제 가용성을 얼마나 잘 반영하는지 비교 분석하였다. Swarm 내의 모든 peer 의 정보를 얻는 것이 쉽지 않기 때문에 일종의 근사치로 사용되는 Busy Period Availability 와 Seed Availability 방법은 leecher 들의 piece 보유 여부를 반영하지 않고 swarm 전체 상황을 반영하지 못하기 때문에 실제 swarm 의 가용성과는 차이가 있음을 보였다. 따라서 swarm 의 가용성을 보다 정확하게 측정하는 방법이 필요하며, 본 연구팀에서는 swarm 에서 얻어진 지역적 정보를 바탕으로 가용성을 가장 정확히 측정할 수 있는 기법을 후속 연구로 진행하고 있다.

참고 문헌

[1] B. Cohen, " Incentives build robustness in bittorrent." 1st Workshop on Economics of p2p Systems, 2003.

[2] " The impact of p2p file sharing, voice over ip, instant messaging, oneclick hosting and media streaming on the internet." <http://www.ipoque.com/resources/internet-studies/internet-study-2008> 2009

[3] Daniel S. Menasche, Antonio A. A. Rocha, Bin Li, Don Towsley, Arun Venkataramani, " Content availability and bundling in swarming systems", *CoNEXT' 09*,

[4] Ashwin R. Bharambe Cormac Herley Venkata N. Padmanabhan, " Analyzing and Improving a BitTorrent Network' s Performance Mechanisms", *INFOCOM 2006*

[5] Azureus - now called vuze - open source bittorrent client. <http://azureus.sourceforge.net/>.