

Data Assignment Method Based on Probability Analysis in Peer-to-Peer Computing System

Ying Yang^{1, a}, Geng Bao^{1, b}

¹College of Computer & electronic Information, Guangxi University, P.R.China

^ayingy2004@126.com, ^bgengbao@126.com,

Keywords: Data Assignment, Probability analysis, P2P System

Abstract. With greatly improved the efficiency of computing performance, Peer-to-Peer systems (P2P) have become hot research in recent years. In this paper, we analyze the probability of a peer to become a seed consisting of all sharing data and resource. Based on the modeled network, a distributed method for data processing is proposed to more efficiently assign the available resources, the data delivery traffic is more evenly distributed and fault tolerance property is well achieved.

Introduction

Peer-to-Peer computing systems have become hot research in recent years. As decentralized architecture, it avoids the drawbacks of conventional client-server model, such as the computation bottleneck on servers, and therefore greatly improved the efficiency of system computing performance [1][2]. On receiving various requests in a P2P network, a peer faces the challenge as how to assign its data and resources [3][4]. For instance, which peer would get service and which fragment of the shared content would be delivered. An effective data assignment method can lead to the resource availability increased, and much more fault tolerance when some peers depart.

In this paper, we analyze the probability of a peer to become a seed which is consisting of all share data and resource. The probability is formalized by the duration time of a peer in a P2P network modeled by uniform, exponential and normal random variables respectively. Based on the probability analysis, an innovative data assignment method is proposed to achieve more available data and resource management, the data delivery traffic is more evenly distributed and fault tolerance property is well achieved.

System Model

Definition1: A P2P network is denoted as an undirected graph $G = (V, E)$, where each peer such as u or v , is denoted as a node, $u \in V, v \in V$. An edge $(u, v) \in E$ represents a bidirectional communication link between u and v . f_s is the size of a shared content.

Definition2: Each node's current state is denoted as a 6-tuple $(B_{uu}, B_{ul}, B_{du}, B_{dl}, f_l, t_s)$, where B_{uu} describes the bandwidth used for node u to upload data in F while B_{ul} represents the left bandwidth, B_{du} shows the occupied bandwidth to download F while B_{dl} means the left bandwidth for content downloading, f_l is the size of the left content, which is denoted by the percentage of f_s , and t_s is the time for a session to content sharing of F when peer u connected in the system.

Node u can not only download file F from other peers but also upload pieces of F concurrently in a distributed P2P system. The node u with $f_l = 0$ is denoted as a seed due to contain the whole file F . The seed node only sends data for a particular file F to non-seed nodes. For those nodes with $f_l \neq 0$, they are attempting to derive data fragment to minimize f_l and as a source to provide data to others. A P2P system is shown in Figure 1 for the sharing of a file F among 4 peers. The network topology represents the request connections in a file sharing session. Along with each node, an information table provides its current state. For node A, its $f_l = 0$ reveals that A can be treated as a seed. Nodes B and C are in the data transmission process and part of the content is attained. Node

D is at the beginning of file downloading because f_i is the same as f_s .

Suppose that the uploading bandwidth of A is free of occupation ($B_{uu} = 0$). With the request from node B and C of sharing of file F, node A can clearly view their tables as in Figure 1. Thus, node A needs to make a decision how to assign its limited bandwidth for package delivery to peers B and C. Moreover, peers, such as B and C, face the same problem to handle requests from others. Figure 2 describes the case for n peers (from b_1 to b_n) send their requests as content sharing to peer a in a P2P distributed system. Peer a is a seed that consists of the entire sharing content. Although a P2P setting is loosely organized and fully decentralized, a participant can be aware of the configuration of those peers connecting to it. Hence, peer a in Figure 2 is transparent after message exchange.

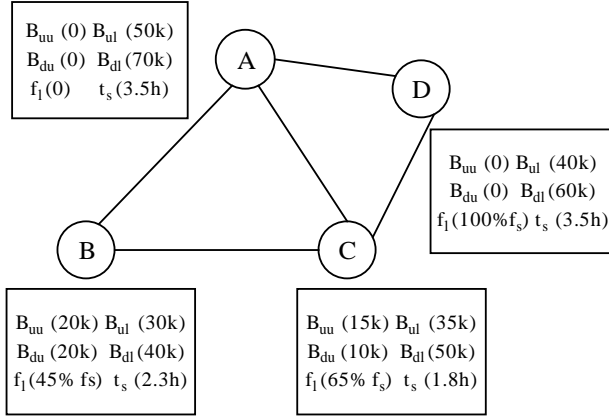


Fig. 1 The state of a peer in P2P

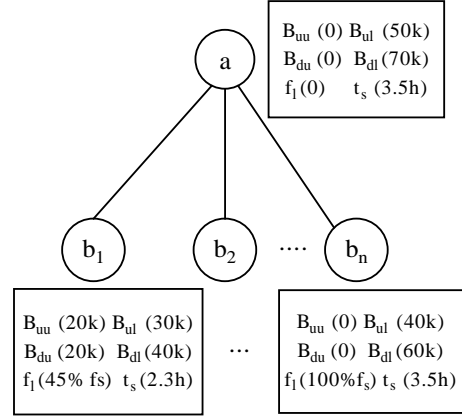


Fig. 2 n peers send their requests to a

Probability Analysis

Let B_i ($i=1,2, \dots, n$) be the event that peer b_i becomes a seed. The value of B_i is either 1 or 0 where 1 represents the defined event to be true, t_i is the time interval from its request to peer a to the time that B_i is true. Thus, we have $t_i = \frac{f_i}{\min(B_{ul_a}, B_{dl_i}) + B_{du_i}}$. The content obtained by peer b_i

can either from other peers, through a channel by B_{du_i} bandwidth, or from the seed a. The second path with the source node a has the maximum bandwidth supported by $\min(B_{ul_a}, B_{dl_i})$ for packet delivery. f_i denotes the content contained by a but not in peer b_i . Let T_a, T_{b_i} be the event that peer a, b_i which will last t_i time longer respectively. Then $B_i = T_a * T_{b_i}$. In a P2P system, peer a and peer b_i can be assumed to be two independent hosts. It follows that T_a and T_{b_i} are two independent events. Therefore, the probability of event B_i to be true can be depicted as

$$P(B_i) = P(T_a) * P(T_{b_i}) \quad (1)$$

Let X_a be the random variable that denotes the time for peer a online. Because peer a has been in the session for t_{s_a} time period, the probability of event T_a to be true is under the conditional probability case. Hence, we have

$$P(T_a) = \frac{P(X_a > t_{s_a} + t_i)}{P(X_a > t_{s_a})} \quad (2)$$

Let X_{b_i} be the random variable that denotes the time for peer b_i online. From Equation 1 and 2, we can rewrite $P(B_i)$ as

$$P(B_i) = \frac{P(X_a > t_{s_a} + t_i)}{P(X_a > t_{s_a})} * \frac{P(X_{b_i} > t_{s_{b_i}} + t_i)}{P(X_{b_i} > t_{s_{b_i}})} \quad (3)$$

Suppose that there are m pieces together for a shared file. They comprise a group $G_i \{g_1, g_2, \dots, g_n\}$. In a session whenever a peer receives any content sharing requests from others, Its

decision for the shared pieces assignment among all applicants depends on its information table as peer a in Figure 2. The peer receives n content sharing requests. Those n peers form a group B_i $\{b_1, b_2, \dots, b_n\}$. For a peer b_i , part of the resource file is stored in its local storage area. Thus, each peer b_i may have different content request upon the source peer a. Although the topology of a P2P network is constantly changed, after an initial setting up process for content sharing of a particular file unit, the network can be treated as a steady state. In other words, the number of peers arrived equalizes those departed and the network keeps in a balanced situation. Let X be a uniform random variable to denote the connecting time for a peer on the interval (α, β) . Its probability density

function is given by $f(x) = \frac{1}{\beta - \alpha}$ For a given $t \in (\alpha, \beta)$, $P(X > t) = 1 - P(X \leq t) = 1 - \frac{t - \alpha}{\beta - \alpha}$. Thus

$$P(B_i) = \frac{\beta - t_{s_a} - t_i}{\beta - t_{s_a}} * \frac{\beta - t_{s_{b_i}} - t_i}{\beta - t_{s_{b_i}}} \quad (4)$$

The duration time for a peer stays connecting to a network can be deployed by the exponential random variable with parameter λ . The arrival and departure of peers comply with a Poisson process. Each peer is independently and exponentially distributed online for content sharing. The probability density function for an exponential random variable X is given by $f(x) = \lambda e^{-\lambda x}$. The property of the exponential random variable makes that the probability of a peer using another t time longer does not involve any previous performance. Thus

$$P(B_i) = P(X_a > t_i) * P(X_{b_i} > t_i) = e^{-\lambda t_i} * e^{-\lambda t_i} = e^{-2\lambda t_i} \quad (5)$$

Suppose that the connection time for a peer in a P2P system approximately satisfies the normal distribution with parameters μ and σ^2 . Let X be the modified normal random variable to represent the duration time of a peer and the probability density function can be represented as

following $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/\sigma^2}$ Thus, we have

$$P(X > t) = 1 - P(X \leq t) = 1 - \frac{1}{\sqrt{2\pi}\sigma} \int_0^t e^{-(x-\mu)^2/\sigma^2} dx \quad (6)$$

Data Assignment Method

A data assignment method based on probability in P2P network is proposed as follow,

1. Let the pieces of content requested by peer b_i constitute the group $P_i \{p_1, p_2, \dots, p_n\}$
2. $f(B_i)$ = Extract node (B)
3. $f(G_i)$ = Extract block (B)
4. Block g_i is sent to peer b_i from the source peer a with bandwidth $B = \min(Bu_{a_i}, Bdl_i)$
5. The selected peer has the maximum value of B_i among the group, which is defined in Equation 3.
6. The value of $P(B_i)$ can be calculated by Equation 4, 5 or 6 according to the employed network modeled by uniform, exponential and normal random variables respectively for the duration time of a peer
7. After the decision of peer b_i to be the one served, the block among G_{b_i} requested in peer a is removed and assigned to G_i .
8. Compared with other blocks in the group G_{b_i} , G_i has the maximum number of peers that would like to copy it.

Whenever a peer has part of content available, our method can be performed inside the peer. It ends when all its outgoing channel bandwidth consumed. Two functions are conducted to select the right peer as the receiver and the block to be sent. In the method, the peer that has the high probability to become a potential seed gets service firstly, which is implemented by the function

of $f(B_i)$. The maximum bandwidth for data transfer is $\min(Bul_a, Bdl_i)$. The block of content to be copied firstly is the one most popular requested. If there is still some abundant Bul_a left, another loop of peer and block selection is executed.

Conclusions

When more peers containing the whole file, the P2P system is more stable because the data delivery traffic is evenly distributed and fault tolerance property is well achieved. In this paper, a novel method solving the peer and the piece of content selection problem is proposed in order to maintain more available seeds. The peer with a higher probability to become a potential seed has the priority to be firstly served. The piece of content that has been mostly requested will be delivered first. The ongoing work is to use this model and the method in resource management of grid environment.

Acknowledgements

The paper was supported by the Fund of Guangxi Natural Science (2013GXNSFAA019344), Gui Financial Education Number [2013] 19. It was also supported by the Fund for the fourth batch of distinguished experts in Nanning City.

References

- [1] Amidala, Himaja, Shankar, Karthik; El Taeib, Tarik, An efficient peer-to-peer platform for large scale data processing in network applications, Applications and Technology Conference, LISAT 2016, 816-821
- [2] Gang Chen, Tianlei Hu, Dawei Jiang, BestPeer++: A Peer-to-Peer Based Large-Scale Data Processing Platform, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, 2014 (26) 1316-1331,
- [3] Hoda Mashayekhi, Jafar Habibi, L-overlay: A layered data management scheme for peer-to-peer computing, PEER-TO-PEER NETWORKING AND APPLICATIONS, 2014 (7) 199-212
- [4] Wei Xiang Goh, Kian-Lee Tan, Generalized data processing on peer-to-peer overlays, Proceedings of the IEEE International Conference on Cloud Engineering, IC2E 2013, 318-327,
- [5] C Lucchese, C Mastroianni, S Orlando, toward a public-resource computing framework for distributed data mining, CONCURRENCY AND COMPUTATION-PRACTICE & EXPERIENCE 2010 (22) 658-682
- [6] Muyong Cai, Xiangming Wen, Wei Zheng, Different-Strategy Management of Malicious Nodes in the Peer-to-Peer Network, International Conference on Environmental Science and Information Application Technology, 2009, 575-578