

基于动态自适应性的语义对等网

刘 晔, 马 慧

许昌学院计算机科学与技术学院, 河南许昌 461000

摘要 针对 P2P 网络的聚类性以及网络的逻辑拓扑与物理拓扑不匹配所造成的通信延迟增加、网络开销增大等问题, 根据节点的物理位置及节点所包含信息的语义相似度, 提出了一种 3 层结构的语义对等网模型。该模型基于节点物理位置与属性特征动态调整网络结构, 自适应地改善自身的资源搜索性能。仿真结果显示, 该 P2P 覆盖网能够以较低的查找时延和代价获得较高的查全率。

关键词 对等网络; 语义; 物理拓扑; 自适应性

中图分类号 TP393.02

文献标识码 A

doi 10.3981/j.issn.1000-7857.2011.32.011

Semantic P2P Networks Based on Dynamic Self-adaptability

LIU Ye, MA Hui

College of Computer Science and Technology, Xuchang University, Xuchang 461000, Henan Province, China

Abstract In view of the clustering of P2P networks and the issues of increased communication latency and network costs caused by mismatch between the logical topology and the physical topology, this paper proposes a mixed three-layer P2P network model. The model adjusts the network structure dynamically based on the position and the property features of nodes, so that the performance of searching resources can be improved adaptively. The simulation result shows that the P2P networks can achieve a high recall ratio with low search delay and cost.

Keywords P2P networks; semantic; physical topology; self-adaptability

0 引言

如何提高资源搜索性能一直是非结构化 P2P 系统研究领域的一个热点, 而网络模型是影响资源搜索性能的重要因素。非结构化 P2P 网络^[1-2]允许节点随机接入, 结构简单, 然而其采用的洪泛搜索机制使带宽资源浪费严重、网络负载加重。混合型 P2P 网络 Gnutella 0.6^[3]兼具了集中式与分布式 P2P 网络的优点, 但其模型中超级节点既要负责节点管理, 又要进行消息路由, 负载过重。文献[4]提出, P2P 网络逻辑拓扑与物理拓扑的严重不匹配延长了节点间的通信延迟, 增加了网络的开销, 从而降低了资源搜索效率。文献[5]表明, P2P 网络的拓扑图不是完全随机的, 而是具有明显的聚类性。

本文根据实际网络特征把所有节点按物理位置划分成若干个域, 在每个域内根据节点的语义相关度划分若干个群, 建立了 3 层结构的语义对等网模型。该模型既考虑了网络逻辑拓扑与物理拓扑的不匹配, 同时缩小了资源查找的范

围。在 Peersim 模拟器上搭建 P2P 模拟环境, 对资源查找算法进行仿真分析, 并与 Gnutella 0.6 进行对比。

1 属性相似度

在 P2P 网络中, 节点之间的资源查找与现实世界中的社会活动具有一定的相似性。现实世界中的资源查找一般在一个特定的群体内进行, 该群体的人有某种共同的爱好与兴趣, 这样查找成功的概率非常大, 同样在 P2P 网络中, 一些节点拥有的资源具有一定的语义相关性, 如果资源查找在具有语义相关性的节点间进行, 不仅能提高资源查找的效率, 而且可以减少网络流量。

定义 1 (语义相似度) 设资源空间 R 与查询空间 S 都采用 VSM 模型^[6-7]表示, 即 $\forall r \in R, \forall s \in S, V_r = [W_{r_1}, \dots, W_{r_M}], V_s = [W_{s_1}, \dots, W_{s_N}]$, 其中 W_r, W_s 分别表示 i 项在 R 与 S 中的权重, 其值采用 TF-idf^[8]计算, 则节点 r_1 与 r_2 的语义相似度计算公

收稿日期: 2011-08-11; 修回日期: 2011-11-04

基金项目: 河南省教育厅科技攻关项目(2009B520026)

作者简介: 刘晔, 讲师, 研究方向为计算机网络和数字媒体, 电子信箱: eduliuye@126.com

式可表示为

$$\text{similarity}(r_1, r_2) = \frac{\sqrt{\sum_{k=1}^N (w_{1k} - w_{2k})^2}}{\sqrt{\sum_{k=1}^N w_{1k}^2} \cdot \sqrt{\sum_{k=1}^N w_{2k}^2}} \quad (1)$$

查询请求与节点间的语义相似度计算方法同式(1)。

2 节点性能评价

在本文中节点的物理性能(CPU,存储能力与带宽)与稳定值作为评价节点性能的主要参数。

定义 2 (物理性能) 设语义群内所有节点的 CPU 频率总和为 C , 内存总和为 M , 带宽总和为 B , 语义群内某节点 P_i 的 CPU 频率为 c_i , 内存为 m_i , 带宽为 b_i , 为这 3 个参数分配的权重系数分别为 $\lambda_1, \lambda_2, \lambda_3$, 则节点 P_i 的物理性能为

$$\text{Capability}(P_i) = \sqrt{\left(\lambda_1 \frac{c_i}{C}\right)^2 + \left(\lambda_2 \frac{m_i}{M}\right)^2 + \left(\lambda_3 \frac{b_i}{B}\right)^2} \quad (2)$$

$\text{Capability}(P_i)$ 反映了节点 P_i 的物理性能, 该值越大, 节点 P_i 的物理性能就越高。

定义 3 (稳定值) 设语义群内某节点 P_i 某段时间的登陆次数为 $\text{LandTotalNum}_{P_i}$, 某次登陆在线时间为 LandTime_{P_i} , 在线时间总和为 TotalOnTime_{P_i} , 则节点 P_i 的稳定值为

$$\text{Stability}(P_i) = \frac{\text{在线时间大于或等于 } \text{AverOnTime}_{P_i} \text{ 的登陆次数}}{\text{LandTotalNum}_{P_i}} \quad (3)$$

$$\text{AverOnTime}_{P_i} = \text{TotalOnTime}_{P_i} / \text{LandTotalNum}_{P_i} \quad (4)$$

定义 3 (综合性能) 假设某节点 P_i 的物理性能为 $\text{Capability}(P_i)$, 稳定性为 $\text{Stability}(P_i)$, 节点性能的权重系数为 λ_p , 稳定性权重系数为 λ_s , λ 依实际情况动态调整, 且 $\sum \lambda = 1$ 。则节点 P_i 的综合性能为

$$W(P_i) = \text{Capability}(P_i) \times \lambda_p + \text{Stability}(P_i) \times \lambda_s \quad (5)$$

式中, λ_p, λ_s 表示节点性能与稳定性 2 项指标对节点综合性能的影响程度, 且系数越大, 影响越大。

3 基于语义与物理拓扑的对等网

基于语义与物理拓扑的对等网由物理位置邻近节点所构成的若干个域组成, 在每个域内, 具有相似语义特征的节点构成群, 网络内节点数为 n 的对等网的形式化描述为

$$\text{P2P networks} = \{ \langle \text{CN}, \text{DCN}, \text{CND}, \delta, \text{RN} \rangle \} \quad (6)$$

式中, CN (Character Nodes) = $\{m_i, i=1, 2, \dots, n\}$; DCN (Domain Character Nodes) 表示经过节点间距离计算, 满足节点间邻近阈值 δ 的构成域 (Domain) 的属性节点集; CND (Character Node Distance) 表示节点间距离计算函数, 用于计算属性节点间的距离; δ 表示检验属性节点是否属于同一域的邻近阈值; RN (Routing Node) 表示域中的路由节点, 负责查找请求在域间的路由。其中

$$\text{Domain} = \{ \langle \text{DCN}, \text{GCN}, \text{CSF}, \eta, \text{IN} \rangle \} \quad (7)$$

式中, GCN (Group Character Nodes) 表示经过节点语义相似度计算后, 满足语义相似度阈值 η 的构成群 (Group) 的属性节点集; CSF (Character Similarity Function) = $\langle F: x, y \rightarrow z \rangle (x, y \in \text{GCN}, z \in (0, 1))$, 表示语义相似度计算函数, 用于计算属性节点的语义相似度; $\eta \in \{0, 1\}$ 表示检验属性节点是否属于同一群的语义相似度阈值; IN (Intelligent Node) 表示群中的智能节点, 负责管理群中所有终端节点 (Terminal Node, TN) 以及群内的资源查找。基于语义与物理拓扑的对等网结构如图 1 所示。

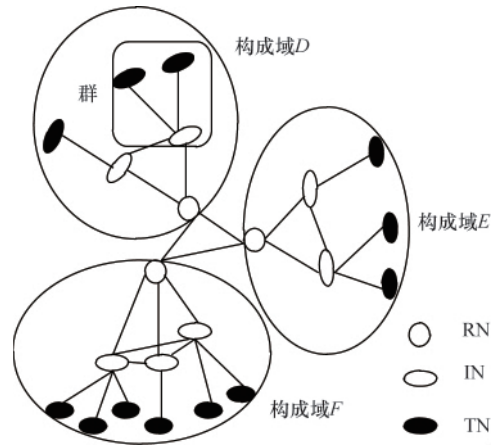


图 1 基于语义与物理拓扑的对等网模型

Fig. 1 P2P network model based on semantic and physical topology

3.1 对等网的构建

在传统混合 P2P 网络中, 某节点加入网络时只需任意选择 1 个域并在该域的超级节点注册, 当要加入的域节点数超过某个阈值时, 就选择另一个域加入。在该 P2P 网络中, 节点加入时不仅要考虑节点的物理位置, 还要考虑节点的语义特征, 不能随意加入某个域内的某个语义群, 要选择与物理位置最近且语义相似的语义群的智能节点相连。

对等网生成算法如下所示。

输入: 节点 N_t 的加入请求

输出: 对等网

Procedure from_network (Node N_t)

Let D represent the set of domain, D_i represent the domain in D

While(true) do {

 Listening to adding request of nodes;

 If(listened adding request of N_t) {

 If($\text{CND}(N_t, \text{RN}_i) \leq \text{CND}(N_t, \text{RN}_k)$) $\text{RN}_k \in D$

 Getconn(N_t, RN_i);

 If($\text{CSF}(N_t, \text{group}_{ij}) \geq \text{CSF}(N_t, \text{group}_{ik})$) and $\text{CSF}(N_t, \text{group}_{ij}) \geq \eta$ and ($\text{Capacity_IP}_{ij} + 1 \leq \text{MAX}$) $\text{group}_{ik} \in D_i$

 /* group_{ik} 表示域 D_i 中的群 k */

 }

比率。查找跳数指单位时间内所有搜索请求的查询跳数的平均值。查找时延与查全率主要衡量该网络构建方法在资源搜索方面的性能,而查找跳数主要衡量资源搜索的代价。

本文在 Peersim 模拟器上搭建 P2P 模拟环境,通过仿真对该对等网的综合性能进行分析,并与 Gnutella 0.6 进行对比。网络初始化结点为 2000 个,分为 10 个域,每个域根据节点语义相似性分为 4 个群。整个网络包含 5000 个资源,设置查询时的语义相似度的阈值为 0.8,节点间按相似度阈值 0.6 生成群。

在仿真中只要某节点拥有被查找的资源就认定搜索成功。仿真结果如图 2—图 4 所示。

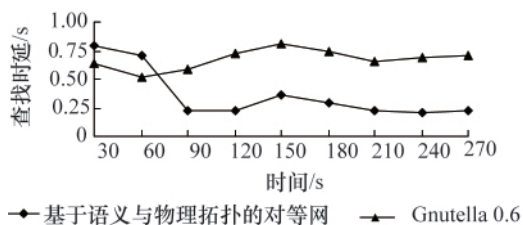


图 2 基于语义与物理拓扑的对等网与 Gnutella 0.6 查找时延的对比

Fig. 2 Search delay comparisons between P2P network based on semantic and physical topology and Gnutella 0.6

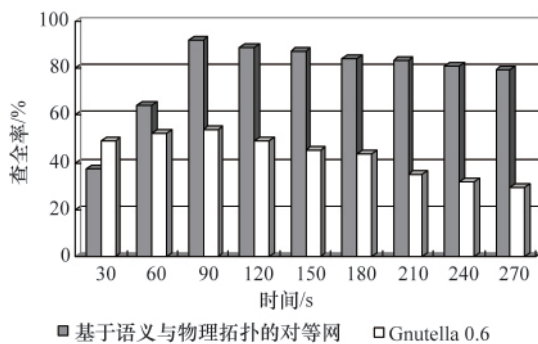


图 3 基于语义与物理拓扑的对等网与 Gnutella 0.6 查全率的对比

Fig. 3 Recall ratio comparisons between P2P network based on semantic and physical topology and Gnutella 0.6

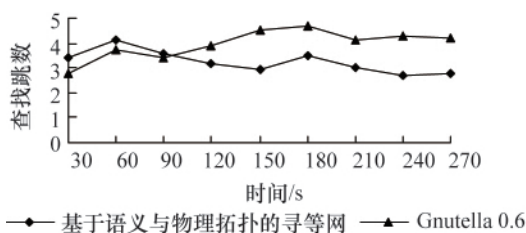


图 4 基于语义与物理拓扑的对等网与 Gnutella 0.6 查找跳数的对比

Fig. 4 Hop comparisons between P2P network based on semantic and physical topology and Gnutella 0.6

仿真结果显示,在试验初始阶段,2 个对等网的查找时延、查全率与查找跳数都有较大幅度的波动;随着时间的推移,基于语义与物理拓扑的对等网的查找时延比 Gnutella 0.6 缩短了约 60%,查全率较 Gnutella 0.6 提高了约 50%,查找跳数比 Gnutella 0.6 减少了约 40%,这是由于本文提出的基于语义与物理拓扑的对等网中节点根据所处物理位置以及语义相似性进行自聚,网络的物理结构与逻辑结构更加匹配,节点在群内或域内就能查找到自己所需资源,从而缩短了查找时延,提高了查全率,大幅提高了网络性能。

5 结论

本文根据 P2P 网络的聚类性,同时考虑 P2P 网络的物理拓扑构造了基于动态自适应的三层语义对等网,详细描述了节点间语义相似度的计算与节点性能评价方法、P2P 对等网的形成以及对等网的动态自适应性。仿真结果显示,基于语义与物理拓扑的对等网能够有效地提高资源查找效率,缩短查找时延,提高资源查全率,具有良好的实用性。

参考文献 (References)

- [1] Clip2 Distributed Search Services. The Gnutella Protocol Specification v0.4 [EB/OL]. [2007-12-20]. <http://wenku.baidu.com/view/1635bad97f1922791688e82a.html>.
- [2] Zhuang L, Pan C J, Guo Y Q, et al. Connection management based on gnutella network[J]. *Journal of Software*, 2005, 16(1): 158-164.
- [3] Ralf S, Klaus W. P2P 系统及其应用[M]. 王玲芳, 陈焱, 译. 北京: 机械工业出版社, 2008.
Ralf S, Klaus W. Peer-to-peer systems and applications[M]. Wang Lingfang, Chen Yan, trans. Beijing: China Machine Press, 2008.
- [4] Xiao L, Liu Y H, Ni L M. Improving unstructured peer-to-peer system by adaptive connection establishment [J]. *IEEE Transactions on Computers (TC)*, 2005, 54(9): 1091-1103.
- [5] Fast A, Jensen D, Levine B N. Creating social networks to improve peer-to-peer networking [C]. Proc of 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining. Chicago: ACM, 2005: 568-573.
- [6] Schutze H, Silverstein C. A comparison of projections for efficient document clustering [C]. Proc of ACM SIGIR. New York: ACM, 1997: 74-81.
- [7] Wikipedia. Vector space model [EB/OL]. [2008-06-20]. http://en.wikipedia.org/wiki/vector_space_model.
- [8] Salton G, Buckley C. Term weighting approaches in automatic text retrieval[R]. New York: Cornell University, 1987.

(责任编辑 孙秀云,代丽)

《科技导报》“卷首语”栏目征稿

“卷首语”栏目每期邀请一位中国科学院院士和中国工程院院士就重大科技现象、事件,以及学科发展趋势、科学研究热点和前沿问题等,撰文发表个人的见解、意见和评论。本栏目欢迎院士投稿,每篇文章约 2000 字,同时请提供作者学术简历、工作照和签名电子文档。投稿邮箱: kjjdbjb@cast.org.cn。