

# 语义覆盖网最佳规模的数学分析

刘 震 周浩浩 邓 苏

(国防科技大学 C<sup>4</sup>ISR 技术国防科技重点实验室 长沙 410073)

**摘 要** 现在大量研究者通过语义覆盖网构建来提高 P2P 网络资源查询效率,但在语义覆盖网最佳规模大小上缺乏研究。考虑运用数学方法对语义覆盖网络进行数据建模,对路由算法的路由性能指标的求解方法进行研究,并分析语义覆盖网规模与路由性能指标之间的关系。通过模型的分析 and 求解,得出了社区的最佳规模大小,为语义覆盖网构建与研究提供了有力的支撑。

**关键词** 对等网,语义覆盖网,社区规模,数学分析

**中图法分类号** TP311 **文献标识码** A

## Mathematic Analysis of Semantic Overlay Network's Optimal Scale

LIU Zhen ZHOU Hao-hao DENG Su

(C<sup>4</sup>ISR Key Laboratory, National University of Defense Technology, Changsha 410073, China)

**Abstract** At present, a lot of researchers have focused on semantic overlay network's construction for improving query efficiency in P2P networks. But the research on optimal scale of semantic overlay network is lacked. This paper tried to use mathematic method to model the semantic overlay network. And then model solution method of properties which impacts the routing arithmetic was researched. At the same time, the relationship between semantic overlay network's scale and the capability properties of routing arithmetic was analyzed. In the base of this, the optimal scale of the community would appear, which supports the construction and research on semantic overlay network effectively.

**Keywords** Peer-to-peer, Semantic overlay network, Community scale, Mathematic analysis

## 1 引言

P2P 网络相对传统的 C/S 架构的网络具有巨大的优势,拥有广阔的市场前景。在第一代非结构化 P2P 系统,如 Gnutella<sup>[1]</sup>,Napster 等和第二代结构化的基于 DHT 的 P2P 系统,如 PAST<sup>[2]</sup>,CFS<sup>[3]</sup>等的基础上,相关研究者对 P2P 的拓扑模型和查询路由由算法进行了大量的改进。

Arturo Crespo 和 Hector Garcia-Molina 最早提出了语义覆盖网络—SON(Semantic Overlay Network)的模型,即将占有资源相似的节点聚集在一起组成群(Cluster)。在语义覆盖网络建立之后,请求查询报文将只被转发到符合该请求内容主题的群中,这就降低了查询代价,同时提高了查询结果的准确度。于婧、张建辉和顾小卓在文献[4]中提出了一种基于主题和物理位置相近原则的层次化对等语义覆盖网络结构 TP-PH,它充分结合了结构化 P2P 网络高效的定位和非结构化 P2P 网络的复杂检索功能,采用分布式哈希表机制将相同主题节点组织成主题区域,在同一主题区域内通过物理位置相近原则进行群的划分,从而在物理网络拓扑基础上建立语义 P2P 覆盖网络结构,提高了查全率并缩短了平均查询时延,是一种支持复杂查询、高性能的语义覆盖网络结构。Elena Paganì, Gian Paolo Rossi 和 Enrico Bertoso 在文献[5]中提出一

种基于本体的路由查询覆盖网络 ORION,它是完全分布式的自组织的覆盖网络,而且反映了节点之间的语义关系。该系统用较低的查询消耗获得了较高的查询效率。陈汉华、金海、宁小敏等人在文献[6]中构建了一种基于语义相似度的 P2P 覆盖网络 SemreX,该系统是一种 P2P 网络环境下的文献检索系统,针对该系统,本文提出一种基于语义相似度的 P2P 拓扑管理和查询路由由算法,它能够有效地提高系统的搜索效率。Paraskevi Raftopoulou 和 Euripides G. M. Petrakis 在文献[7]中从提高语义覆盖网络的群(社区)聚集程度出发,提出了增强聚集度的有效途径。

以上研究者通过语义覆盖网络结构、语义覆盖网构建、语义路由由算法等各个方面的研究,来改进和提高 P2P 网络系统的查询效率。但这些研究没有考虑语义覆盖网络的规模问题,语义覆盖网络的规模大小直接影响着路由算法的效率:若社区规模较小,则查全率较低,若社区规模较大,则每次查询需要处理的消息数较多。本文试图运用数学方法,在对 P2P 网络进行数据建模的基础上,分析对路由算法影响较大的一些拓朴性质,并据此得出适用于路由算法的拓朴结构和语义社区的最佳规模大小,从而达到改进 P2P 网络系统整体性能的目的。

到稿日期:2010-02-05 返修日期:2010-05-01

刘 震 博士,讲师,主要研究方向为信息管理、P2P、智能决策支持技术,E-mail:hero12251976@163.com;周浩浩 硕士生,主要研究方向为 P2P、语义覆盖网;邓 苏 教授,主要研究方向为信息管理、智能决策支持技术。

## 2 模型构建

语义覆盖网可以理解为一个语义社区,后面用社区代指语义覆盖网,在语义覆盖网络构建之初,假设已有  $M(M \gg 2)$  个社区存在于网络中。机制如下:

### 1. 初始化

每个社区有  $m_0$  个节点(初始建立的时候规模较小),社区之间有  $\frac{M(M-1)}{2}$  条边相连,即每个社区之间有一条直接相连的边存在。

### 2. 节点的偏好加入

每个单位时间内系统加入一个节点,对于节点  $i$  加入社区  $j$  的概率与该社区的节点个数有关,且满足偏好连接:

$$\Pi(v_j) = \frac{v_j}{\sum_k v_k}$$

式中,  $\sum_k v_k = N = Mm_0 + t$ ;  $v_j$  表示社区  $j$  的节点个数。

在社区内该节点与  $m(1 \leq m \leq m_0)$  个节点相连,并以概率  $\alpha$  与其他  $M-1$  个社区内的  $n_0(1 \leq n_0 \leq m)$  节点相连,总共连接  $n_0$  个。注:若不对网络进行划分,则社区个数为 1,即整个网络为一个社区。

### 3. 边的择优连接

#### A. 社区内择优连接

新加入的节点  $j$  与社区内节点  $i$  相连的概率与节点  $i$  在社区内的度有关,即

$$\Pi(s_{ij}) = \frac{s_{ij}}{\sum_k s_{kj}}$$

式中,  $s_{ij}$  为节点  $i$  在社区  $j$  内的度。

#### B. 社区间择优连接

新节点与社区  $k(k \neq j)$  内的节点  $i$  相连的概率与  $i$  的社区间的度有关,即

$$\Pi(l_{ik}) = \frac{l_{ik}}{\sum_{m,n,n \neq j} l_{im}}$$

式中,  $l_{ik}$  为社区  $k$  内的节点  $i$  的社区间的度。

所以在  $t$  个单位时间之后,网络中有  $N = Mm_0 + t$  个节点,边的数量为  $\frac{Mm_0(m_0-1) + M(M-1)}{2} + mt + \alpha n_0 t$ 。

## 3 模型的理论分析与求解

### 3.1 社区规模的求解

对于节点在社区间的数量分布,就社区  $j$ ,有

$$\frac{\partial v_j}{\partial t} = \frac{v_j}{\sum_k v_k} = \frac{v_j}{Mm_0 + t}$$

假设  $t_j$  时刻社区加入第一个新节点,此时社区的节点个数应当为  $m_0 + 1$  个,故可得方程的解为:

$$v_j(t) = \frac{(m_0 + 1)(Mm_0 + t)}{t_j}$$

假设社区的节点个数是连续的且是均匀的,则节点个数小于  $k$  的概率可以表示为:

$$P(v_j(t) < k) = P(t_j > \frac{(m_0 + 1)(Mm_0 + t)}{k})$$

$$\text{且有, } P_j(t_j) = \frac{1}{Mm_0 + t}$$

故可得

$$\begin{aligned} P(t_j > \frac{(m_0 + 1)(Mm_0 + t)}{k}) \\ &= 1 - P(t_j \leq \frac{(m_0 + 1)(Mm_0 + t)}{k}) \\ &= 1 - \frac{(m_0 + 1)(Mm_0 + t)}{k(Mm_0 + t)} = 1 - \frac{(m_0 + 1)}{k} \end{aligned}$$

所以节点在各个社区内的分布为:

$$P(k) = \frac{\partial P(v_j(t) < k)}{\partial k} = (m_0 + 1)k^{-2}$$

即,节点在各社区的分布满足幂律分布。从实际角度看,可以得出的结论是社区的规模有大有小。

### 3.2 网络度分布的求解

对于社区  $j$  内的节点  $i$ ,有

$$\frac{\partial s_{ij}}{\partial t} = m \frac{v_j}{\sum_k v_k} \frac{s_{ij}}{\sum_k s_{kj}}$$

注意到

$$\begin{cases} v_j(t) = \frac{(m_0 + 1)(Mm_0 + t)}{t_j} \\ \sum_k v_k = Mm_0 + t \\ \sum_k s_{kj} = 2mt \frac{1}{M} + m_0(m_0 - 1) \\ \approx 2mt \frac{1}{M} \quad (t \text{ 较大时后一项可以忽略}) \end{cases}$$

故

$$\frac{\partial s_{ij}}{\partial t} \approx \frac{M(m_0 + 1)s_{ij}}{2t_j t}$$

设在  $t_i$  时刻节点  $i$  加入社区  $j$  中,初始化为  $s_{ij}(t_i) = m$ ,故得

$$s_{ij}(t) \approx \frac{Mm(m_0 + 1)}{t_j} \left(\frac{t}{t_i}\right)^{0.5}$$

基于上式,可得到度小于  $k$  的概率为:

$$P(s_{ij}(t) < k) = P(t_i > \frac{\theta^2 t}{k^2})$$

$$\text{式中, } \theta = \frac{Mm(m_0 + 1)}{t_j}$$

假设节点加入的时间是均匀的,则

$$P_i(t_i) = \frac{1}{Mm_0 + t}$$

$$P_j(t_j) = \frac{1}{Mm_0 + t}$$

由上面两式可得

$$P(t_i > \frac{\theta^2 t}{k^2}) = 1 - P(t_i \leq \frac{\theta^2 t}{k^2}) = 1 - \frac{M^2 m^2 (m_0 + 1)^2 t}{k^2 (Mm_0 + t)^3}$$

$$\text{即 } P(k) = \frac{\partial P(s_{ij}(t) < k)}{\partial k} = \frac{2M^2 m^2 (m_0 + 1)^2 t}{(Mm_0 + t)^3} k^{-3}$$

同理,对于社区间的连接,有

$$\frac{\partial l_{ik}}{\partial t} = \frac{M-1}{M} \alpha n_0 \frac{l_{ik}}{\sum_{m,n,n \neq j} l_{im}}$$

其中,

$$\sum_{m,n,n \neq j} l_{im} = 2 \frac{M-1}{M} \alpha n_0 t + [M(M-1) - (M-1)]$$

$$\text{当 } \alpha \neq 0 \text{ 时, } l_{ik}(t) = \alpha n_0 \left(\frac{t + \beta}{t_i + \beta}\right)^{0.5}$$

$$\text{式中, } \beta = \frac{M(M-1) - (M-1)}{2\alpha n_0 (M-1)}$$

$$\text{当 } 2 \frac{M-1}{M} \alpha n_0 t \gg M(M-1) - (M-1) \text{ 时,有 } \frac{\partial l_{ik}}{\partial t} \approx \frac{l_{ik}}{2t};$$

又,已知初始条件:在  $t_j$  时  $l_{ik}(t_j) = \alpha n_0$

故得  $l_{ik}(t) \approx \alpha n_0 (\frac{t}{t_j})^{0.5}$

即  $P(k) = \frac{2(\alpha n_0)^2 t}{M m_0 + t} k^{-3}$

以上求出的即为社区内的度分布和社区间节点相连的度分布情况。

由以上结果可知,社区内的度分布和社区间节点相连的度分布均服从幂律分布,与 BA 网络模型得出的分布情况类似,符合实际情况。

### 3.3 平均路径长度以及节点平均度数的求解

关于节点的平均度数,我们主要求解社区内的节点平均度数。

首先,介绍节点度的概率分布函数:

$$\Phi(x) = \sum_{k=1}^{\infty} p_k x^k$$

式中,  $p_k$  表示社区中任选一点度数为  $k$  的概率,并且有,  $\Phi(1) = 1$ 。

使得整个概率中间的总和为 1。则在某一社区  $j$  中,节点的平均度数可表示为:

$$\begin{aligned} \bar{k} &= \sum_{k=1}^{\infty} k p_k = \sum_{k=1}^{\infty} \frac{2M^2 m^2 (m_0 + 1)^2 t}{(M m_0 + t)^3} k^{-2} \\ &= \frac{2M^2 m^2 (m_0 + 1)^2 t}{(M m_0 + t)^3} \sum_{k=1}^{\infty} k^{-2} = \Phi' \end{aligned} \quad (1)$$

式中,  $\sum_{k=1}^{\infty} k^{-2} = \frac{\pi^2}{6}$

即节点平均度数等于函数  $\Phi(x)$  在  $x$  等于 0 时的一阶导数。

关于平均路径长度的计算,进行复杂网络研究的研究人员还没有找出一个实用的方法,有的得出了计算公式,但其因过于复杂而未得到广泛的应用。

为了对平均路径长度进行求解,假设每个节点的度均为节点的平均度数。平均路径长度是指网络中某节点到其他节点的最短路径长度的平均。故对于某个节点来说,其与网络中其他节点的连接可看作是一棵模型树,树的第一层次为该节点第二层次有  $\bar{k}$  个节点,到了第三层次,应该有  $\bar{k}^2$  个节点,第  $n$  层次应当有  $\bar{k}^{n-1}$  个节点,而对于这棵模型树其总的节点的数量应为网络或者社区的节点总数。即有,  $N = \bar{k} + \bar{k}^2 + \bar{k}^3 + \dots + \bar{k}^{n-1} = \bar{k} \frac{\bar{k}^{n-1} - 1}{\bar{k} - 1}$ 。

由上可知  $n$  为网络中其他节点与该节点的最大距离。

故,平均路径长度可表示为

$$\begin{aligned} \bar{l} &= \frac{\bar{k} + 2\bar{k}^2 + 3\bar{k}^3 + \dots + (n-1)\bar{k}^{n-1}}{\bar{k} + \bar{k}^2 + \bar{k}^3 + \dots + \bar{k}^{n-1}} \\ &= \frac{\bar{k} + 2\bar{k}^2 + 3\bar{k}^3 + \dots + (n-1)\bar{k}^{n-1}}{\bar{k} \frac{\bar{k}^{n-1} - 1}{\bar{k} - 1}} \end{aligned}$$

对于  $\bar{k} + 2\bar{k}^2 + 3\bar{k}^3 + \dots + (n-1)\bar{k}^{n-1}$  的求解如下:

由于  $1 + x + x^2 + \dots + x^n = \frac{1-x^{n+1}}{1-x}$ , 对其求导可得

$0 + 1 + 2x + \dots + nx^{n-1} = (\frac{1-x^{n+1}}{1-x})'$ , 两边各乘以  $x$  可得

$$\begin{aligned} x + 2x^2 + \dots + nx^n &= x(\frac{1-x^{n+1}}{1-x})' \\ &= \frac{x[1-x^{n+1} - nx^{n-1}(1-x)]}{(1-x)^2} \end{aligned}$$

故

$$\begin{aligned} &\bar{k} + 2\bar{k}^2 + 3\bar{k}^3 + \dots + (n-1)\bar{k}^{n-1} \\ &= \bar{k} \frac{[1-\bar{k}^{n-1} - (n-1)\bar{k}^{n-2}(1-\bar{k})]}{(1-\bar{k})^2} \end{aligned}$$

由此得平均路径长度为:

$$\begin{aligned} \bar{l} &= \frac{\bar{k} + 2\bar{k}^2 + 3\bar{k}^3 + \dots + (n-1)\bar{k}^{n-1}}{\bar{k} \frac{\bar{k}^{n-1} - 1}{\bar{k} - 1}} \\ &= \frac{\bar{k} \frac{[1-\bar{k}^{n-1} - (n-1)\bar{k}^{n-2}(1-\bar{k})]}{(1-\bar{k})^2}}{\bar{k} \frac{\bar{k}^{n-1} - 1}{\bar{k} - 1}} \\ &= \frac{[1-\bar{k}^{n-1} - (n-1)\bar{k}^{n-2}(1-\bar{k})]}{(\bar{k}-1)(\bar{k}^{n-1}-1)} \end{aligned}$$

式中,  $n$  为网络中其他节点与该节点的最大距离。

### 3.4 模型的数据求解示例

假设在构建之初,已有  $M=10$  个社区存在于网络中,且每个社区有  $m_0=10$  个全耦合的节点。当节点加入社区中时,与社区内的  $m=5$  个节点相连,并以  $\alpha=0.2$  概率与其他社区内的  $n_0=3$  个节点建立连接。

则可求出节点在各个社区内的分布函数为:

$$P(v) = 11v^{-2}$$

其分布图如图 1 所示。

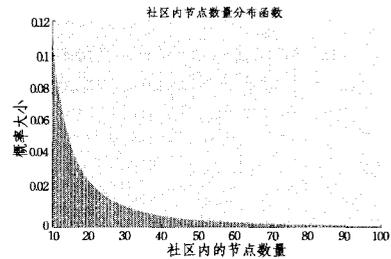


图 1 社区内的节点数量分布

从图 1 中可以看出社区内节点的数量越多,其概率就越小,表明按照此机制构建的语义覆盖网络,其社区的规模呈幂律分布。

由上亦可求出社区内节点的度分布,由计算可知,度分布函数为:

$$p(k) = \frac{0.72t}{100+t} k^{-3}$$

其分布图如图 2 所示。

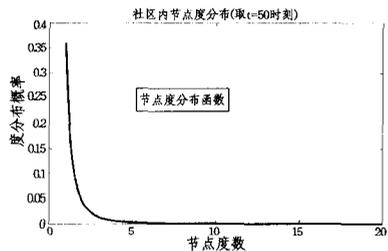


图 2 社区内节点度分布

以下是节点平均度数和平均路径长度的求解。

由公式

$$\bar{k} = \frac{2M^2 m^2 (m_0 + 1)^2 t}{(M m_0 + t)^3} \sum_{k=1}^{\infty} k^{-2} = \frac{M^2 m^2 (m_0 + 1)^2 t \pi^2}{3(M m_0 + t)^3}$$
 可知:

当  $t=200$  时,节点平均度数为  $\bar{k}=3.1846$ 。

由公式  $\bar{l} = \frac{[1-\bar{k}^{n-1} - (n-1)\bar{k}^{n-2}(1-\bar{k})]}{(\bar{k}-1)(\bar{k}^{n-1}-1)}$  可知:

当  $t=200$  时,社区内平均路径长度为  $\bar{l}=4.7330$ 。

## 4 社区最佳规模大小分析

社区规模的大小直接影响路由算法的性能,从理论上讲,如果规模过大,每次查询所要处理的消息数就较多;规模较小,则查全率就会较低。若要求查全率和每次查询处理的消息数满足一定的条件,则社区规模的大小就应当在一个范围内变化。考虑实际情况,所要做的工作是在保证查全率达到某个数值的情况下,调整社区的规模大小,使得每次查询所要处理的消息数最少。而且,对于不同的查询策略,其社区规模的最佳大小也应当相应的不同。

社区规模、查全率、消息数以及平均路径长度之间的关系如图 3 所示。

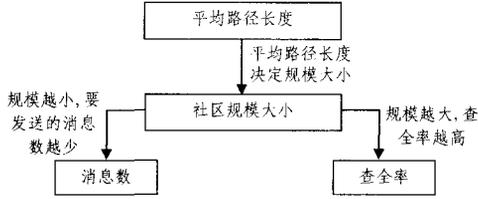


图 3 社区规模、查全率、消息数以及平均路径长度之间的关系

社区的最佳规模大小主要针对有社区网络存在的基于洪泛策略的研究。在保证查全率达到某个数值的情况下,调整社区的规模大小,使得每次查询所要处理的消息数最少。每次查询需要处理的消息数和查全率的数学表示为:

$$Q = \begin{cases} k \frac{\bar{k}^\gamma - 1}{\bar{k} - 1}, & \text{当 } \bar{l} > \gamma \text{ 时;} \\ k \frac{\bar{k}^{\bar{l}} - 1}{\bar{k} - 1}, & \text{当 } \bar{l} \leq \gamma \text{ 时} \end{cases}$$

$$\eta = \frac{\bar{k}^{\bar{l}+1} - \bar{k}^2}{(\bar{k} - 1)(m_0 + 1) \sum_{k=1}^{\infty} k^{-1}}$$

式中,  $\bar{k}$  为节点的平均度数,  $\gamma$  为 TTL 值,  $\bar{l}$  为社区内节点的平均路径长度。

由于社区规模的大小与平均路径长度有很大关联,求出平均路径长度即可求出社区规模的大小,即社区内节点的总数。平均路径长度与社区规模大小的关系见第 3 节。求解思路为:在确定查全率的最低值后,可以得出平均路径长度的范围,从而在消息数最小的情况下可求出社区的最佳规模。

求解思路如图 4 所示。



图 4 社区最佳规模的求解思路

当社区规模最佳时,查全率与每次查询所要处理的消息数应当达到一定的条件,查全率以及消息数是通过平均路径长度与社区的规模建立联系的。在上述的假设条件下,我们得到查全率和每次查询所要处理的消息数与平均路径长度的关系,如图 5,图 6 所示。

由图可知,查全率是随着平均路径长度的增大而基本呈增大趋势的,从客观上讲是因为社区规模较大时查全率较高。亦可知,每次查询需要处理的消息量是随着平均路径长度的减小而减小的。

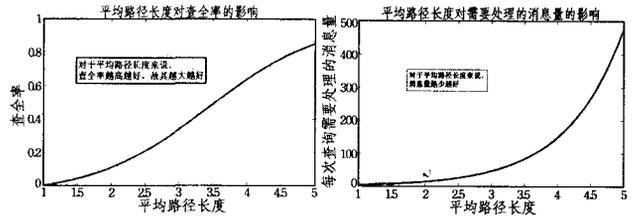


图 5 平均路径长度对查全率的影响 图 6 平均路径长度对每次查询需要处理的消息量的影响

在实际求解中,假设查全率必须达到 60%,如图 7 所示,平均路径长度就需大于 3.8,由此可以得出在此条件下每次查询需要处理的最少的消息数,如图 8 所示,最少 120。可由此时的平均路径长度求出社区的最佳规模大小  $N=182.1317$ 。

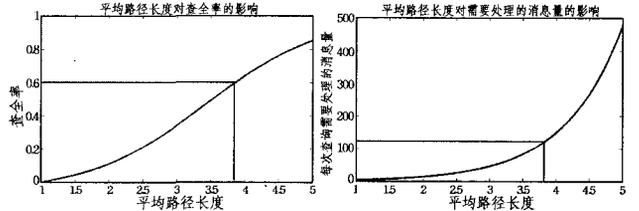


图 7 通过查全率获得平均路径长度(查全率 > 60%) 图 8 通过平均路径长度获得需要处理的消息量(平均路径长度 > 3.8)

**结束语** 本文在分析 P2P、语义覆盖网研究现状的基础上,运用数学方法对 P2P 语义覆盖网络进行建模,并对模型加以分析和求解,求出社区规模分布、网络的节点度分布以及平均路径长度、节点平均度数等。同时在分析路由算法主要性能指标查全率和每次查询所要转发的消息数的求解方法的基础上,建立了路由算法性能指标与模型求出的统计量之间的关系,并根据模型的求解对社区的最佳规模大小进行了分析。

本文的研究为语义覆盖网的构建、优化和分割提供了有力的理论和技术支撑。

## 参考文献

- [1] Ripeanu M, Foster I, Iamnitchi A. Mapping the gnutella network; properties of large-scale peer-to-peer systems and implications for system design[J]. IEEE Internet Computing, 2002, 6(1): 50-57
- [2] Rowstron A, Druschel P. Storage management and caching in PAST, a large-scale, persistent peer-to-peer storage utility [C]// SP'01, Banff; ACM Press, 2001; 188-201
- [3] Dabek F, Kaashoek M, Karger D. Widearea cooperative storage with CFS[C] // SOS P'01. New York, ACM Press, 2001; 202-215
- [4] 于婧, 张建辉, 顾小卓, 等. 基于主题和物理位置相近原则的层次化对等语义覆盖网络结构[J]. 电子与信息学报, 2008, 30(8): 521-524
- [5] Pagani E, Rossi G P, Pertoso E. ORION: Ontology-based query routing in Overlay Networks[J]. Journal of Parallel and Distributed Computing, 2009, 69: 28-38
- [6] 陈汉华, 金海, 宁小敏, 等. SemreX: 一种基于语义相似度的 P2P 覆盖网络[J]. 软件学报, 2006, 17(5): 1170-1181
- [7] Ogston, Elth. An analysis of interest-community facilitated peer-to-peer search[J]. Lecture Notes in Computer Science, 2008; 98-110