

# Peer-to-Peer storage of power-law graphs

Casper Petersen, Noy Rotbart,  
Jakob Grue Simonsen and Christian Wulff-Nilsen

Department of Computer Science, University of Copenhagen  
Universitetsparken 5, 2100 Copenhagen  
{cazz,noyro,simonsen,koolooz}@diku.dk  
phone: (+45) 29611674

## ABSTRACT

The study of web graph compression assumes a centralised data structure which is undesired in peer-to-peer networks. We propose an investigation on web graph dissemination using a theoretically appropriate method, adjacency labeling scheme. These are methods that assigns labels to the vertices of a graph such that adjacency between them can be inferred directly from the assigned labels, without using a centralized data structure. We first provide a lower bound for the family of power-law graphs. This family has been used to model many types of networks, e.g. the Internet AS-level graph. Furthermore, we provide an almost matching upper bound for this family. This upper bound is attained by pre-determining a threshold separating nodes of large degree and small degree that attempts to balance the storage required among the vertices. We then validate the efficiency of our separation by an experimental evaluation using both synthetic data and real-world networks of up to hundreds of thousands of vertices.

## Categories and Subject Descriptors

H.4 [Information Systems Applications]: Miscellaneous;  
D.2.8 [Software Engineering]: Metrics—*complexity measures, performance measures*

## General Terms

Theory

## Keywords

ACM proceedings, L<sup>A</sup>T<sub>E</sub>X, text tagging

## 1. INTRODUCTION

A body of work on web graphs deals with the difficulties of storing them. Due to their size, studying web graphs is difficult to store in main memory. This led first to compression techniques [20, 19], and thereafter to study on the dissemination of such graphs over several machines [37]. Our study sets theoretical limits to the space required for such efforts.

One possible approach to solve these difficulties is to disseminate the structural information of the graph to its vertices. Known as peer to peer, such dissemination allows for an inference of the graph's local topology using only local information stored in each vertex without using costly access to large, global data structures. One way of doing so is via a *labeling scheme*: an algorithm that assigns a bit string—a *label*—to each vertex so that a query between any two vertices can be deduced solely from their respective labels. Labeling schemes are extremely well studied topic [], under the objective of minimizing the *maximum label size*: the maximum number of bits used in a label of any vertex. They were also found practical use in XML search engines [30], mapping services [1] and routing [46].

Similarly to [], we restrict our attention to the atomic operation of *adjacency* queries. Adjacency labeling schemes for general graphs require a label size of  $n/2 + O(1)$  [49, 11]. Trees, planar graphs, and bounded degree graphs, on the other hand, enjoy labels of logarithmic size [12, 33, 3]. To the best of our knowledge, we are the first to study adjacency labeling schemes for classes of graphs whose statistical properties—in particular their *degree distribution*—more closely resemble that of real-world networks. One class of graphs extensively used for modelling real-world networks is *power-law graphs*: roughly,  $n$ -vertex graphs where the number of vertices of degree  $k$  is proportional to  $n/k^\alpha$  for some positive  $\alpha$ . Power-law graphs (also called scale-free graphs in the literature) have been used, e.g., to model the Internet AS-level graph [52, 5], and many other types of network (see, e.g., [48, 29] for overviews). The adequacy of fit of power-law graph models to actual data, as well as the empirical correctness of the conjectured mechanisms giving rise to power-law behaviour, have been subject to criticism (see, e.g., [2, 29]). In spite of such criticism, and because their degree distribution affords a reasonable approximation of the degree distribution of many networks, the class of power-law graphs remains a popular tool in network modelling whose statistical behaviour is well-understood: e.g., for power-law graphs with  $2 < \alpha < 3$ , the range most often seen in the modeling of real-world networks [29], it is known that with high probability the average distance between any two vertices is  $O(\log \log n)$ , the diameter is  $O(\log n)$  and there exists a dense subgraph of  $n^{c/\log \log n}$  vertices [26].

Routing labeling schemes for power-law graphs have been investigated by Brady and Cowen [22], and by Chen et al. [25]. Labeling schemes for other properties than adja-

gency have been investigated for various classes of graphs, e.g., distance [34], and flow [41]. Dynamic labeling schemes were studied by Korman and Peleg [43, 44, 42] and recently by Dahlgaard et. al [31]. Experimental evaluation for some labeling schemes for various properties on general graphs have been performed by Caminiti et. al [24], Fischer [32] and Rotbart et. al [51].

Adjacency labeling schemes are tightly coupled with the graph-theory related concept of induced universal graphs. Given a graph family  $\mathcal{F}$ , the aim is to find smallest  $N$  such that a graph of  $N$  vertices contains all graphs in  $\mathcal{F}$  as induced subgraphs. Kannan, Naor and Rudich [40] showed that an  $f(n)\log n$  adjacency labeling scheme for  $\mathcal{F}$  constructs an induced universal graph for this family of  $2^{f(n)}$  vertices. Some of the adjacency labeling schemes reported earlier contributed a better bound than was known of induced universal graphs (see e.g [18, 12]). In the context of sparse graphs, a body of work on universal graphs<sup>1</sup> for this family was investigated both by Babai et al. [16] and by Alon and Asodi [7].

## 1.1 Our contribution

Our contributions are:

*An  $O(\sqrt[3]{n}(\log n)^{1-1/\alpha})$  adjacency labeling scheme for power-law graphs  $G$ .* The scheme is based on two ideas: (I) a labeling *strategy* that partitions the vertices of  $G$  into high (“fat”) and low degree (“thin”) vertices based on a threshold degree, and (II) a threshold *prediction* that depends only on the coefficient  $\alpha$  of a power-law curve fitted to the degree distribution of  $G$ . Real-world power-law graphs rarely exceed  $10^{10}$  vertices, implying a label size of at most  $10^5$  bits, well within the processing capabilities of current hardware. We claim that our scheme is thus appealing in practice due both to its simplicity and hte small size of its labels. Using the same ideas, we get an asymptotically near-tight  $O(\sqrt[3]{n} \log n)$  adjacency labeling scheme for sparse graphs.

*A lower bound of  $\Omega(\sqrt[3]{n})$  bits on the maximum label size for any adjacency labeling scheme for power-law graphs.* To this end we define a restrictive subclass of power-law graphs and show that it is contained in the bigger class we study for the upper bound; we show that this class requires label size  $\Omega(\sqrt[3]{n})$  for  $n$ -vertex graphs. This lower bound shows that our upper bound above is asymptotically optimal, bar a  $(\log n)^{1-1/\alpha}$  factor. By the connections between adjacency labeling schemes and universal graphs, we also obtain upper and lower bounds for induced universal graphs for power-law graphs.

*An experimental investigation of our labeling scheme.* Using both real-world (23K-325K vertices) and synthetic (300K-1M vertices) data sets, we observe that: (i) Our threshold *prediction* performs close to optimal when using the labeling *strategy* above. (ii) our labeling scheme achieves maximum label size several orders of magnitude smaller than

<sup>1</sup>A graph that contains each graph from the graph family as subgraph, not necessarily induced.

the state-of-the-art labeling schemes for more general graph families.

In addition, our study may contribute to the understanding of the quality of *generative models*—procedures that “grow” random graphs whose degree distributions are with high probability “close” to power-law graphs, such as the Barabasi-Albert model [17] and the Aiello-Chung-Lu model [4]. As a first step, we provide an evidence that the randomized Barabasi-Albert model [17] produces only a small fraction of the power-law graphs possible.

## 2. PRELIMINARIES

Throughout the paper, we consider  $n$ -vertex, undirected, finite graphs. For real  $c > 0$ , a graph is  $c$ -sparse if it has at most  $cn$  edges and *sparse* if it is  $c$ -sparse for some constant  $c$ . For  $0 < c \leq n - 1$ , the set of  $c$ -sparse graphs with  $n$  vertices is denoted by  $\mathcal{S}_{c,n}$ . If  $\mathcal{F}$  is a set of graphs,  $\mathcal{F}_n$  denotes the subset of graphs in  $\mathcal{F}$  having exactly  $n$  vertices. The degree of a vertex  $v$  in a graph is denoted by  $\Delta(v)$ , and for non-negative integers  $k$ , the set of vertices in a graph  $G$  of degree  $k$  is denoted by  $V_k$ . The length of a binary string  $x \in \{0, 1\}^*$  is denoted by  $|x|$ .

Let  $\mathcal{F}$  be a set of graphs. An *adjacency labeling scheme* (from hereon just *labeling scheme*) for  $\mathcal{G}$  is a pair consisting of an *encoder* and a *decoder*. The encoder is an algorithm that receives  $G \in \mathcal{G}$  as input and outputs a bit string  $\mathcal{L}(v) \in \{0, 1\}^*$  called the *label* of  $v$ . The decoder is an algorithm that receives any two labels  $\mathcal{L}(v), \mathcal{L}(u)$  as input and outputs **true** iff  $u$  and  $v$  are adjacent in  $G$  and **false** otherwise. Note that the graph  $G$  is not an input to the decoder. The *size* of a labeling scheme is the map  $\text{size} : \mathbb{N} \rightarrow \mathbb{N}$  such that  $\text{size}(n)$  is the maximum length of any label assigned by the encoder to any vertex in any graph  $G \in \mathcal{F}_n$ . The *degree distribution* of a graph  $G = (V, E)$  is the mapping  $\text{ddist}_G(k) : \mathbb{N}_0 \rightarrow \mathbb{Q}$  defined by  $\text{ddist}_G(k) := \frac{|V_k|}{n}$ .

We treat the family of *power-law* graphs, which is defined in the literature as the class of  $n$  vertex graphs  $G$  such that  $\text{ddist}_G(k)$  is proportional to  $k^{-\alpha}$  for some real number  $\alpha > 1$ . Ideally, and ignoring rounding, we would like  $\text{ddist}_G(k) = Ck^{-\alpha}$  for all  $k$  for constant  $C$ . As the degree distribution of a graph must be a probability distribution, we have  $\sum_{k=1}^{\infty} Ck^{-\alpha} = C \sum_{k=1}^{\infty} k^{-\alpha} = 1$ , hence  $C = 1/\zeta(\alpha)$  where  $\zeta$  is the Riemann zeta function.

## 3. GRAPH FAMILIES RELATED TO POWER-LAW GRAPHS

In this section we define two families of graphs  $\mathcal{P}_\alpha$  and  $\mathcal{P}'_\alpha$  with  $\mathcal{P}'_\alpha \subseteq \mathcal{P}_\alpha$ . Family  $\mathcal{P}_\alpha$  is rich enough to contain the graphs whose degree distribution is approximately, or perfectly, power-law distributed, and our upper bound on the label size for our labeling scheme holds for any graph in  $\mathcal{P}_\alpha$ . Family  $\mathcal{P}'_\alpha$  is used to show our lower bound. In the following, let  $i_1 = \Theta(\sqrt[3]{n})$  be the smallest integer such that  $\lfloor Cn/i_1^\alpha \rfloor \leq 1$ , and let  $C' \geq (\frac{C}{\alpha-1} + \frac{i_1}{\sqrt[3]{n}} + 5)^\alpha + \frac{C}{\alpha-1}$  be a constant; we shall use  $C'$  in the remainder of the paper.

DEFINITION 1. Let  $\alpha > 1$  be a real number.  $\mathcal{P}_\alpha$  is the

family of graphs  $G$  such that if  $n = |V(G)|$  then for all integers  $k$  between  $\sqrt[n]{n/\log n}$  and  $n-1$ ,  $\sum_{i=k}^{n-1} |V_i| \leq C'(\frac{n}{k^{\alpha-1}})$ .

The class of  $\alpha$ -proper power law graphs contains graphs where the number of vertices of degree  $k$  must be  $C \frac{n}{k^\alpha}$  rounded either up or down and the number of vertices of degree  $k$  is non-increasing with  $k$ . Note that the function  $k \mapsto C \frac{1}{k^\alpha}$  is strictly decreasing.

DEFINITION 2. Let  $\alpha > 1$  be a real number. We say that an  $n$ -vertex graph  $G = (V, E)$  is an  $\alpha$ -proper power-law graph if

1.  $\lfloor Cn \rfloor - i_1 - 1 \leq |V_1| \leq \lceil Cn \rceil$ ,
2.  $\lfloor C \frac{n}{2^\alpha} \rfloor \leq |V_2| \leq \lceil C \frac{n}{2^\alpha} \rceil + 1$ ,
3. for every  $i$  with  $3 \leq i \leq n$ :  $|V_i| \in \{\lfloor C \frac{n}{i^\alpha} \rfloor, \lceil C \frac{n}{i^\alpha} \rceil\}$ , and
4. for every  $i$  with  $2 \leq i \leq n-1$ :  $|V_i| \geq |V_{i+1}|$ .

The family of  $\alpha$ -proper power-law graphs is denoted  $\mathcal{P}'_\alpha$ .

Note that we allow slightly more noise in the sizes of  $V_1$  and  $V_2$  than in the remaining sets; without it, it seems tricky to prove a better lower bound than  $\Omega(n^{\frac{1}{\alpha+1}})$  on label sizes.

Other definitions of power-law graphs are given in the literature. One is the shifted power law distribution [Sank,Eom] which is defined as follows. Let  $c_1 > 0$  be a constant. A graph  $G$  is *power law bounded* for parameters  $1 < \alpha = O(1)$  and  $t \geq 0$  if for every integer  $d \geq 0$ , the number of vertices of  $G$  of degree in  $[2^d, 2^{d+1})$  is at most

$$c_1 n(t+1)^{\alpha-1} \sum_{i=2^d}^{2^{d+1}-1} (i+t)^{-\alpha}.$$

As experimentally verified in [Sank], the value of  $t$  is typically very small. If  $t = O(1)$ , the bound above becomes  $O(n \sum_{i=2^d}^{2^{d+1}-1} i^{-\alpha})$ . In this case, it is easy to see that our family  $\mathcal{P}'_\alpha$  is rich enough to contain these power law bounded graphs (for sufficiently big constant  $C'$ ) and so our upper bound also applies to power law bounded graphs. It is possible to extend our upper bound to super-constant  $t$  where the bound is stronger the smaller  $t$  is; we omit the details. Regarding lower bounds, our family  $\mathcal{P}'_\alpha$  is restrictive enough so that any lower bound for this family also holds for power law bounded graphs when  $t = O(1)$ .

Comment from CWN: papers Sank and Eom should be cited and added to bib file. The papers are, respectively:

Pawel Brach, Marek Cygan, Jakub Lacki, Piotr Sankowski  
Algorithmic Complexity of Power Law Networks arXiv:1507.02426  
[cs.DS]. To appear at SODA'16.

Y.-H. Eom and S. Fortunato. Characterizing and modeling citation dynamics. PLoS ONE, 6(9):e24926, 09 2011

We show the following properties of  $\mathcal{P}'_\alpha$ .

PROPOSITION 1. The maximum degree in an  $n$ -vertex graph in  $\mathcal{P}'_\alpha$  is at most  $\left(\frac{C}{\alpha-1} + 2\right) \sqrt[n]{n} + i_1 + 3 = \Theta(\sqrt[n]{n})$ .

PROOF. Let  $n > 0$  be an integer and let  $k' = \lfloor \sqrt[n]{n} \rfloor$ . Furthermore, let  $S_{k'} = \sum_{i=1}^{k'} |V_i|$ , that is  $S_{k'}$  is the number of vertices of degree at most  $k'$ . Let  $S_{k'}^- = (\sum_{i=1}^{k'} \lfloor Cni^{-\alpha} \rfloor) - i_1 - 1$ . Then  $S_{k'} \geq S_{k'}^-$ . We now bound  $S_{k'}^-$  from below. For every  $i$  with  $1 \leq i \leq k'$ ,

$$\begin{aligned} S_{k'}^- + k' &= -i_1 - 1 + \sum_{i=1}^{k'} (\lfloor Cni^{-\alpha} \rfloor + 1) \geq \\ &= -i_1 - 1 + \sum_{i=1}^{k'} Cni^{-\alpha} = -i_1 - 1 + Cn \sum_{i=1}^{k'} i^{-\alpha} \\ &\geq n \left( 1 - C \sum_{i=k'+1}^{\infty} i^{-\alpha} \right) - i_1 - 1 \\ &\geq n \left( 1 - C \int_{k'}^{\infty} x^{-\alpha} dx \right) - i_1 - 1 \\ &= n \left( 1 - C \left[ \frac{1}{\alpha-1} x^{-\alpha+1} \right]_{k'}^{\infty} \right) - i_1 - 1 \\ &= n \left( 1 - \frac{C}{\alpha-1} (\lceil \sqrt[n]{n} \rceil)^{-\alpha+1} \right) - i_1 - 1 \\ &\geq n \left( 1 - \frac{C}{\alpha-1} (\sqrt[n]{n})^{-\alpha+1} \right) - i_1 - 1 \\ &= n - \frac{Cn}{\alpha-1} n^{-1+\frac{1}{\alpha}} - i_1 - 1 \\ &= n - \frac{C}{\alpha-1} \sqrt[n]{n} - i_1 - 1, \end{aligned}$$

giving  $S_{k'} \geq S_{k'}^- \geq n - \frac{C}{\alpha-1} \sqrt[n]{n} - \lceil \sqrt[n]{n} \rceil - i_1 - 1$ . There are thus at most  $\frac{C}{\alpha-1} \sqrt[n]{n} + \lceil \sqrt[n]{n} \rceil + i_1 + 1$  vertices of degree strictly more than  $k' = \lceil \sqrt[n]{n} \rceil$ . Since for every  $1 \leq i \leq n-1$ :  $|V_i| \geq |V_{i+1}|$ , it follows that the maximum degree of any  $\alpha$ -proper power-law graph is at most  $\left(\frac{C}{\alpha-1} + 2\right) \sqrt[n]{n} + i_1 + 3$ .  $\square$

PROPOSITION 2. For  $\alpha > 2$ , all graphs in  $\mathcal{P}'_\alpha$  are sparse.

PROOF. By Proposition 1, the maximum degree of an  $n$ -vertex  $\alpha$ -proper power-law graph is at most  $k' \triangleq \left(\frac{C}{\alpha-1} + 2\right) \sqrt[n]{n} + i_1 + 3$ , whence the total number of edges is at most  $\frac{1}{2} \sum_{k=1}^{k'} k |V_k|$ . By definition,  $|V_k| \leq \lceil \frac{Cn}{k^\alpha} \rceil \leq \frac{Cn}{k^\alpha} + 1$  for  $k \neq 2$  and  $|V_2| \leq \lceil \frac{Cn}{2^\alpha} \rceil + 1$ , and thus

$$\begin{aligned} \frac{1}{2} \sum_{k=1}^{k'} k |V_k| &\leq 1 + \frac{1}{2} \sum_{k=1}^{k'} k \left( \frac{Cn}{k^\alpha} + 1 \right) \\ &\leq 1 + \frac{k'(k'+1)}{4} + Cn \sum_{k=1}^{\infty} k^{-\alpha+1} \\ &= O(n^{2/\alpha}) + Cn\zeta(\alpha-1) = O(n). \end{aligned}$$

$\square$

PROPOSITION 3.  $\mathcal{P}'_\alpha \subseteq \mathcal{P}_\alpha$ .

PROOF. Let  $d = \lfloor (\frac{C}{\alpha-1} + 2) \sqrt[\alpha]{n} + i_1 + 3 \rfloor$ . For any  $\alpha$ -proper power-law graph with  $n$  vertices and for any  $k, |V_k| \leq Ck^{-\alpha}n + 1$  and by Proposition 1,  $|V_k| = 0$  when  $k > d$ .

Let  $k$  be an arbitrary integer between  $\sqrt[\alpha]{n/\log n}$  and  $n-1$ . We need to show that  $\sum_{i=k}^{n-1} |V_i| \leq C'(\frac{n}{k^{\alpha-1}})$ . It suffices to show this for  $k \leq d$ . We have:

$$\begin{aligned} \sum_{i=k}^{n-1} |V_i| &\leq \sum_{i=k}^d (Ci^{-\alpha}n + 1) = d - k + 1 + Cn \sum_{i=k}^d i^{-\alpha} \\ &\leq \left( \frac{C}{\alpha-1} + \frac{i_1}{\sqrt[\alpha]{n}} + 5 \right) \sqrt[\alpha]{n} + Cn \int_k^d x^{-\alpha} dx \\ &\leq \left( \frac{C}{\alpha-1} + \frac{i_1}{\sqrt[\alpha]{n}} + 5 \right) \sqrt[\alpha]{n} + Cn \left[ \frac{1}{\alpha-1} x^{-\alpha+1} \right]_k^\infty \\ &\leq \left( \left( \frac{C}{\alpha-1} + \frac{i_1}{\sqrt[\alpha]{n}} + 5 \right) \left( \frac{\sqrt[\alpha]{n} d^{\alpha-1}}{n} \right) + \frac{C}{\alpha-1} \right) nk^{-\alpha+1} \\ &\leq \left( \frac{C}{\alpha-1} + \frac{i_1}{\sqrt[\alpha]{n}} + 5 \right) \left( \frac{C}{\alpha-1} + \frac{i_1}{\sqrt[\alpha]{n}} + 5 \right)^{\alpha-1} nk^{-\alpha+1} \\ &\quad + \left( \frac{C}{\alpha-1} \right) nk^{-\alpha+1} \\ &\leq C' nk^{-\alpha+1}, \end{aligned}$$

as desired.  $\square$

## 4. THE LABELING SCHEMES

We now construct algorithms for labeling schemes for  $c$ -sparse graphs and for the family  $\mathcal{P}_\alpha$ . Both labeling schemes partition vertices into *thin* vertices which are of low degree and *fat* vertices of high degree. The *degree threshold* for the scheme is the lowest possible degree of a fat vertex. We start with  $c$ -sparse graphs.

**THEOREM 1.** *There is a  $\sqrt{2cn \log n} + 2 \log n + 1$  labeling scheme for  $\mathcal{S}_{c,n}$ .*

PROOF. Let  $G = (V, E)$  be an  $n$ -vertex  $c$ -sparse graph. Let  $f(n)$  be the degree threshold for  $n$ -vertex graphs; we choose  $f(n)$  below. Let  $k$  denote the number of fat vertices of  $G$ , and assign each to each fat vertex a unique identifier between 1 and  $k$ . Each thin vertex is given a unique identifier between  $k+1$  and  $n$ .

For a  $v \in V$ , the first part of the label  $\mathcal{L}(v)$  is a single bit indicating whether  $v$  is thin or fat followed by a string of  $\log n$  bits representing its identifier. If  $v$  is thin, the last part of  $\mathcal{L}(v)$  is the concatenation of the identifiers of the neighbors of  $v$ . If  $v$  is fat, the last part of  $\mathcal{L}(v)$  is a *fat bit string* of length  $k$  where the  $i$ th bit is 1 iff  $v$  is incident to the (fat) vertex with identifier  $i$ .

Decoding a pair  $(\mathcal{L}(u), \mathcal{L}(v))$  is now straightforward: if one of the vertices, say  $u$ , is thin,  $u$  and  $v$  are adjacent iff the identifier of  $v$  is part of the label of  $u$ . If both  $u$  and  $v$  are fat then they are adjacent iff the  $i$ th bit of the fat bit string of  $\mathcal{L}(u)$  is 1 where  $i$  is the identifier of  $v$ .

Since  $|E| \leq cn$ , we have  $k \leq 2cn/f(n)$ . A fat vertex thus has label size  $1 + \log n + k \leq 1 + \log n + 2cn/f(n)$  and a thin vertex has label size at most  $1 + \log n + f(n) \log n$ . To minimize the

maximum possible label size, we solve  $2cn/x = x \log n$ . Solving this gives  $x = \sqrt{2cn/\log n}$  and setting  $f(n) = \lceil x \rceil$  gives a label size of at most  $1 + \log n + (\sqrt{2cn/\log n} + 1) \log n \leq 1 + 2 \log n + \sqrt{2cn \log n}$ .  $\square$

By Proposition 2, graphs in  $\mathcal{P}'_\alpha$  are sparse for  $\alpha > 2$ . This gives a label size of  $O(\sqrt{n \log n})$  with the labeling scheme in Theorem 1. We now show that this label can be significantly improved, by constructing a labeling scheme for  $\mathcal{P}_\alpha$  which contains  $\mathcal{P}'_\alpha$ .

**THEOREM 2.** *There is a  $\sqrt[\alpha]{C'n}(\log n)^{1-1/\alpha} + 2 \log n + 1$  labeling scheme for  $\mathcal{P}_\alpha$ .*

PROOF. The proof is very similar to that of Theorem 1. We let  $f(n)$  denote the degree threshold. If we pick  $f(n) \geq \sqrt[\alpha]{n/\log n}$  then by Definition 1 there are at most  $C'n/f(n)^{\alpha-1}$  fat vertices. Defining labels in the same way as in Theorem 1 gives a label size for thin vertices of at most  $1 + \log n + f(n) \log n$  and a label size for fat vertices of at most  $1 + \log n + C'n/f(n)^{\alpha-1}$ . We minimize by solving  $x \log n = C'n/x^{\alpha-1}$ , giving  $x = \sqrt[\alpha]{C'n/\log n}$ . Setting  $f(n) = \lceil x \rceil$  gives a label size of at most  $\sqrt[\alpha]{C'n}(\log n)^{1-1/\alpha} + 2 \log n + 1$ .  $\square$

### 4.1 Labeling scheme for random graphs

Graphs in  $\mathcal{P}_\alpha$  have a fixed degree sequence. However, certain graph generation models for power-law graphs are inherently random; one example is the preferential attachment model. For graphs obtained from such models, their degree sequences are instead probability distributions. In this section, we show that applying our labeling scheme for  $\mathcal{P}_\alpha$  to random graphs with the power law distribution, we get a good expected worst-case label size.

Using the definition of Mitzenmacher, a random variable  $X$  is said to have the *power law* distribution (w.r.t.  $\alpha > 1$ ) if

$$\Pr[X \geq x] \sim cx^{-\alpha+1},$$

for a constant  $c > 0$ , i.e.,  $\lim_{x \rightarrow \infty} \Pr[X \geq x]/cx^{-\alpha+1} = 1$ .

Let  $\epsilon > 0$  be fixed. Consider a graph  $G$  picked from a family  $\mathcal{F}$  of random graphs whose degree sequences have the power law distribution. Order the vertices of  $G$  arbitrarily as  $v_1, \dots, v_n$ . For  $i = 1, \dots, n$ , let indicator variable  $X_i$  be 1 iff  $v_i$  has degree at least  $d = \sqrt[\alpha]{n/\log n}$ . There is a constant  $N_0 \in \mathbb{N}$  (depending on  $\epsilon$ ) such that if  $n \geq N_0$  then for all  $i$ ,

$$E[X_i] = \Pr[X_i = 1] \leq (1 + \epsilon)cd^{-\alpha+1}.$$

With the same labeling scheme as for  $\mathcal{P}_\alpha$  with degree threshold  $d$ , denote by  $E_n$  the expected label size of an  $n$ -vertex

graph from  $\mathcal{F}$ . Then for all  $n \geq N_0$ ,

$$\begin{aligned}
E_n &= \sum_{x=0}^n \Pr \left[ \sum_{i=1}^n X_i = x \right] O((x + d \log n)) \\
&= O \left( d \log n + E \left[ \sum_{i=1}^n X_i \right] \right) \\
&= O \left( d \log n + \sum_{i=1}^n E[X_i] \right) \\
&= O(d \log n + n d^{-\alpha+1}) \\
&= O(\sqrt[n]{n} (\log n)^{1-1/\alpha}).
\end{aligned}$$

**THEOREM 3.** *Let  $\mathcal{F}$  be a family of graphs with degree sequences having the power law distribution w.r.t.  $\alpha > 1$ . Then there is a labeling scheme for  $\mathcal{F}$  such that the expected worst-case label size of any graph  $G \in \mathcal{F}$  is  $O(\sqrt[n]{n} (\log n)^{1-1/\alpha})$  where  $n$  is the number of vertices of  $G$ .*

## 5. LOWER BOUNDS

We now derive lower bounds for the label size of any labeling schemes for both  $\mathcal{S}_{c,n}$  and  $\mathcal{P}_\alpha$ . Our proofs rely on Moon's [49] lower bound of  $\lfloor n/2 \rfloor$  bits for labeling scheme for general graphs. We first show that the upper bound achieved for sparse graphs is close to the best possible. The following proposition is essentially a more precise version of the lower bound suggested by Spinrad [53].

**PROPOSITION 4.** *Any labeling scheme for  $\mathcal{S}_{c,n}$  requires labels of size at least  $\left\lfloor \frac{\sqrt{cn}}{2} \right\rfloor$  bits.*

**PROOF.** Assume for contradiction that there exists a labeling scheme assigning labels of size strictly less than  $\left\lfloor \frac{\sqrt{cn}}{2} \right\rfloor$ . Let  $G$  be an  $n$ -vertex graph. Let  $G'$  be the graph resulting by adding  $\left\lfloor \frac{n^2}{c} \right\rfloor - n$  isolated vertices to  $G$ , and note that now  $G'$  is  $c$ -sparse. The graph  $G$  is an induced subgraph of  $G'$ . It now follows that the vertices of  $G$  have labels of size strictly less than  $\left\lfloor \frac{\sqrt{c \lfloor n^2/c \rfloor}}{2} \right\rfloor \leq n/2$  bits. As  $G$  was arbitrary, we obtain a contradiction.  $\square$

### 5.1 Lower bound for power-law graphs

In the remainder of this section we are assuming that  $\alpha > 2$  and prove the following:

**THEOREM 4.** *For all  $n$ , any labeling scheme for  $n$ -vertex graphs of  $\mathcal{P}_\alpha$  requires label size  $\Omega(\sqrt[n]{n})$ .*

More precisely, we present a lower bound for  $\mathcal{P}'_\alpha$  which is contained in  $\mathcal{P}_\alpha$ . Let  $n \in \mathbb{N}$  be given and let  $H = (V(H), E(H))$  be an arbitrary graph with  $i_1$  vertices where  $i_1 = \Theta(\sqrt[n]{n})$  is defined as in Section 3. We show how to construct an  $\alpha$ -proper power-law graph  $G = (V, E)$  with  $n$  vertices that contains  $H$  as an induced subgraph. Observe that a labeling of  $G$  induces a labeling of  $H$ . As  $H$  was chosen arbitrarily and as any labeling scheme for  $k$ -vertex graphs requires  $\lfloor i_1/2 \rfloor$  label size in the worst case, Theorem 4 follows if we can show the existence of  $G$ .

We construct  $G$  incrementally where initially  $E = \emptyset$ . Partition  $V$  into subsets  $V_1, \dots, V_n$  as follows. The set  $V_1$  has size  $\lfloor Cn \rfloor - i_1$ . For  $i = 2, \dots, i_1 - 1$ ,  $V_i$  has size  $\lfloor Cn/i^\alpha \rfloor$ . Letting  $n' = \sum_{i=1}^{i_1-1} |V_i|$ , we set the size of  $V_i$  to 1 for  $i = i_1, \dots, i_1 + n - n' - 1$  and the size of  $V_i$  to 0 for  $i = i_1 + n - n', \dots, n$ , thereby ensuring that the sum of sizes of all sets is  $n$ . Observe that  $\sum_{i=1}^{i_1-1} \lfloor Cn/i^\alpha \rfloor \leq n$  so that  $n' \leq n - i_1$ , implying that  $n - n' \geq i_1$ . Hence we have at least  $i_1$  size 1 subsets  $V_{i_1}, \dots, V_{i_1+n-n'-1}$  in each of which the vertex degree allowed by Definition 2 is at least  $i_1$ .

Let  $v_1, \dots, v_{i_1}$  be an ordering of  $V(H)$ , form a set  $V_H \subseteq V$  of  $i_1$  arbitrary vertices from the sets  $V_{i_1}, \dots, V_{i_1+n-n'-1}$ , and choose an ordering  $v'_1, \dots, v'_{i_1}$  of  $V_H$ . For all  $i, j \in \{1, \dots, i_1\}$ , add edge  $(v'_i, v'_j)$  to  $E$  iff  $(v_i, v_j) \in E(H)$ . Now,  $H$  is an induced subgraph of  $G$  and since the maximum degree of  $H$  is  $i_1 - 1$ , no vertex of  $V_i$  exceeds the degree bound allowed by Definition 2 for  $i = 1, \dots, n$ .

We next add additional edges to  $G$  in three phases to ensure that it is an  $\alpha$ -proper power law graph while maintaining the property that  $H$  is an induced subgraph of  $G$ . For  $i = 1, \dots, n$ , during the construction of  $G$  we say that a vertex  $v \in V_i$  is *unprocessed* if its degree in the current graph  $G$  is strictly less than  $i$ . If the degree of  $v$  is exactly  $i$ ,  $v$  is *processed*.

**Phase 1:** Let  $V' = V \setminus (V_1 \cup V_H)$ . Phase 1 is as follows: while there exists a pair of unprocessed vertices  $(u, v) \in V' \times V_H$ , add  $(u, v)$  to  $E$ .

When Phase 1 terminates,  $H$  is clearly still an induced subgraph of  $G$ . Furthermore, all vertices of  $V_H$  are processed. To see this, note that the sum of degrees of vertices of  $V_H$  when they are all processed is  $O(i_1^2) = O(n^{2/\alpha})$  which is  $o(n)$  since  $\alpha > 2$ . Furthermore, prior to Phase 1, each of the  $\Theta(n)$  vertices of  $V'$  have degree 0 and can thus have their degrees increased by at least 1 before being processed.

**Phase 2:** Phase 2 is as follows: while there exists a pair of unprocessed vertices  $(u, v) \in V' \times V'$ , add  $(u, v)$  to  $E$ . At termination, at most one vertex of  $V'$  remains unprocessed. If such a vertex exists we process it by connecting it to  $O(\sqrt[n]{n})$  vertices of  $V_1$ ; as  $|V_1| = \Theta(n)$  there are enough vertices of  $V_1$  to accomodate this. Furthermore, prior to adding these edges, all vertices of  $V_1$  have degree 0, and hence the bound allowed for vertices of this set is not exceeded.

**Phase 3:** In Phase 3, we add edges between pairs of unprocessed vertices of  $V_1$  until no such pair exists. If no unprocessed vertices remain we have the desired  $\alpha$ -proper power law graph  $G$ . Otherwise, let  $w \in V_1$  be the unprocessed vertex of degree 0. We add a single edge from  $w$  to another vertex  $w'$  of  $V_1$ , thereby processing  $w$  and moving  $w'$  from  $V_1$  to  $V_2$ . Note that the sizes of  $V_1$  and  $V_2$  are kept in their allowed ranges due to the first two conditions in Definition 2. This proves Theorem 4.

## 6. DISTANCE LABELING SCHEME

In this section we propose a distance labeling scheme for power law graphs.

The *distance* between two nodes in an undirected graph is the length of the shortest path connecting the two nodes, if it exists, and  $\infty$  if no such path exists.

Let  $f : \mathbb{N} \rightarrow \mathbb{N}$  be a map such that  $f(n) \leq -1$  for all  $n$ . An  $f(n)$ -distance labelling scheme is a labelling scheme such that, for any graph  $G$ , its decoder given labels  $\mathcal{L}(u)$  and  $\mathcal{L}(v)$  of two nodes  $u$  and  $v$  will output the distance between  $u$  and  $v$  if the distance is at most  $f(|V(G)|)$ , and output “no” if the distance is strictly greater than  $f(|G|)$ . If  $f(n) = n-1$ , an  $f(n)$ -distance labelling scheme is simply called a *distance labelling scheme*.

For sparse graphs, Alstrup et al. [10] obtain a distance labelling scheme with maximum label size  $O(\frac{n}{D} \log^2 D)$  where  $D = (\log n)/(\log \frac{m+n}{n})$  and  $m$  is the number of edges in the graph. Using similar methods, Gawrychowski et al. obtain an upper bound of [35]  $O(\frac{n}{D} \log D)$  with sublinear decoding time. Few general results on lower bounds exist. The lower bound of  $\Omega(\sqrt{n})$  for adjacency given in the present paper is trivially also a lower bound for distance; for total label size, the best known lower bound remains  $\Omega(n^{3/2})$  as proved by Gavaille et al. citeGavaille2001.

We now devise an  $f(n)$ -distance labelling scheme for  $c$ -sparse graphs, works particular well (i.e., has shorter labels than any known labelling schemes) for small distances. As all power-law graphs will in general be sparse, and power-law graphs in general have very small expected distances, the labelling scheme should work well for practical purposes in power-law graphs.

LEMMA 1. *Let  $c > 0$ . For any computable  $f : \mathbb{N} \rightarrow \mathbb{N}$   $f(n) \leq -1$  for all  $n$ , there is an  $f(n)$ -distance labelling scheme for the family of  $c$ -sparse graphs that assigns labels of length at most  $O(n^{f(n)/(f(n)+1)} \log f(n))$ .*

PROOF. As for adjacency labelling, the scheme is based on *thin* and *fat* nodes. Let  $G$  be a  $c$ -sparse graph. Call a node of  $G$  *fat* if it has degree at least  $n^{1/(f(n)+1)}$  and *thin* otherwise. The label of each node  $v$  now contains (i) a table of distances to all fat nodes (if the distance is more than  $f(n)$ , it is simply ignored), and (ii) a table of distances to all thin nodes  $w$  that are at most distance  $f(n)$  away from  $v$  where the shortest path between  $v$  and  $w$  does not pass through any fat node. Clearly, as  $f(n)$  is computable and distances in  $G$  are computable, there is a computable encoder assigning labels. Furthermore, as all nodes of  $G$  are either thin or fat, it is clearly possible for an encoder to compute all distances less than or equal to  $f(n)$  between any pair of nodes. Note that as all distances we care for are bounded above by  $f(n)$ , each such distance can be stored using at most  $\log f(n)$  bits.

As the sum of degrees of fat nodes are at most  $2cn$  in a  $c$ -sparse graph, there can be at most

$$\frac{2cn}{n^{\frac{1}{f(n)+1}}} = 2cn^{1-\frac{1}{f(n)+1}} = 2cn^{\frac{f(n)}{f(n)+1}}$$

fat nodes in  $G$ . Thus, a table of distances to all fat nodes takes up at most  $O(n^{\frac{f(n)}{f(n)+1}} \log f(n))$  bits.

Similarly, for each node  $v$  there are at most  $(n^{1/(f(n)+1)})^{f(n)} = n^{f(n)/(f(n)+1)}$  nodes at distance at most  $f(n)$  away from  $v$  where the shortest path consists only of thin nodes. Hence, the associated table of distances takes up at most  $O(n^{f(n)/(f(n)+1)} \log n)$  bits.

In total, each label thus has size at most  $O(n^{f(n)/(f(n)+1)} \log n)$  bits.  $\square$

For  $f(n) = \log n$ , Lemma 1 yields a labelling scheme having label size at most  $O(n^{(\log n)/(1+\log n)} \log \log n)$ . Unsurprisingly, as we are only considering distances up to  $f(n)$ , this label size is asymptotically smaller than for the labelling schemes working for all distances in sparse graphs, e.g. the largest label sizes of [35] for sparse graphs is  $O((n/\log n) \log \log n)$ . For power law random graphs with *expected JGS: FIXME!* [Make sure this ties in with Christian's section](#), Chung and Lu show in [26] that, subject to mild conditions, the diameter of power law graphs with  $\alpha > 2$  is almost surely  $\Theta(\log n)$ . We thus expect our labelling scheme to have superior performance for such graphs.

## 7. SCALE FREE GRAPHS FROM GENERATIVE MODELS

The Barabási-Albert (BA) model is a well-known generative model for power-law graphs that, roughly, grows a graph in a sequence of time steps by inserting a single vertex at each step and attaching it to  $m$  existing vertices with probability weighted by the degree of each existing vertex [17]. The BA model generates graphs that asymptotically have a power-law degree distribution ( $\alpha = 3$ ) for low-degree nodes [21]. Graphs created by the BA model have low arboricity<sup>2</sup> [36] we use that fact to prove the following highly efficient labelling scheme.

PROPOSITION 5. *The family of graphs generated by the BA model has an  $O(m \log n)$  adjacency labeling scheme.*

PROOF. Let  $G = (V, E)$  be an  $n$ -vertex graph resulting by the construction by the BA model with some parameter  $m$  (starting from some graph  $G_0 = (V_0, E_0)$  with  $|V_0| \ll n$ ). While it is not known how to compute the arboricity of a graph efficiently, it is possible in near-linear time to compute a partition of  $G$  with at most twice<sup>3</sup> the number of forests in comparison to the optimal [13]. We can thus decompose the graph to  $2m$  forests in near linear time and label each forest using the recent  $\log n + O(1)$  labeling scheme for trees [9], and achieve a  $2m(\log n + O(1))$  labeling scheme for  $G$ .  $\square$

Note that if the encoder operates at the same time as the creation of the graph, Proposition 5 can be strengthened to

<sup>2</sup>the arboricity of a graph is the minimum number of spanning forests needed to cover its edges.

<sup>3</sup>More precisely, for any  $\epsilon \in (0, 1)$  there exist an  $O(|E(G)|/\epsilon)$  algorithm [45] that computes such partition using at most  $(1 + \epsilon)$  times more forests than the optimal.

yield an  $m \log n$  labeling scheme: simply store the identifiers of the  $m$  vertices attached with every vertex insertion. Theorem 4 and Proposition 5 strongly suggest that, for each sufficiently large  $n$ , the number of power-law graphs with  $n$  vertices is vastly larger than the number of graphs with  $n$  vertices created by the BA model. In contrast, other generative models such as Waxman [54], N-level Hierarchical [23], and Chung’s [27] (Chapter 3) do not seem to have an obvious smaller label size than the one in Proposition 2.

## 8. EXPERIMENTAL STUDY

We now perform an experimental evaluation of our labeling scheme on a number of large networks. The source code for our experiments can be found at:

[www.diku.dk/~simonsen/suppmat/podc15/powerlaw.zip](http://www.diku.dk/~simonsen/suppmat/podc15/powerlaw.zip)

### 8.1 Experimental Framework

**Performance Indicators.** Recall that our labeling scheme consists of two ideas: separation of the nodes according to some threshold, and selecting a threshold depending on the power-law parameter  $\alpha$ . In our labeling scheme, the threshold is  $\lceil \sqrt[\alpha]{Cn/(\alpha-1)} \rceil$ . We call this the *predicted* threshold; it is an approximation to the theoretically optimal threshold choice when degree distributions follow the power-law curve  $k \mapsto Cn/k^\alpha$  perfectly. The approximation uses integration similar to what is done in, e.g., the proof of Proposition 3. For a concrete graph  $G$ , it is conceivable that some other threshold  $n_0$ , different from the predicted threshold, would yield a labeling scheme with smaller size. Let  $\max_t(n_0)$  and  $\max_f(n_0)$  be the maximum label sizes of thin, resp. fat vertices in  $G$  when the threshold is set at  $1 \leq n_0 \leq n-1$ . Clearly the maximum label size with the threshold  $n_0$  is  $\max\{\max_t(n_0), \max_f(n_0)\}$ . Observe further that  $\max_t(n_0)$  and  $\max_f(n_0)$  are monotonically increasing, resp. decreasing functions of  $n_0$ . Hence, the  $n_0$  for which  $\max\{\max_t(n_0), \max_f(n_0)\}$  is minimal is where the curves of  $\max_t(n_0)$  and  $\max_f(n_0)$  intersect. We call this  $n_0$  the *empirical* threshold. We set up the following performance indicators to gauge (1) the difference in label size with predicted and empirical threshold, and (2) the label size obtained by our labeling scheme on several data sets, compared to other labeling schemes.

**Performance Indicator 1:** We measure the label sizes for the labeling schemes with the predicted and empirical thresholds. We interpret a small relative difference between these label sizes means that the predicted threshold can achieve small label sizes without examining the global properties of the network other than the power-law parameter  $\alpha$ .

**Performance Indicator 2:** We measure the label sizes attained by our labeling schemes to other labeling schemes, namely state-of-the art labeling schemes for the classes of bounded-degree, sparse and general graphs using the labeling schemes suggested in [3], Theorem 1 and [11]. We interpret small label sizes for our scheme, especially in comparison with “small” classes like the class of bounded-degree graphs, as a sign that our labeling scheme efficiently utilizes the extra information about the graphs: namely that their degree distribution is reasonably well-approximated by a power-law.

**Test Sets.** We employ both real-world and synthetic data sets.

The six *synthetic* data sets are created by first generating a power-law degree sequence using the method of Clauset et al. [29, App. D], subsequently constructing a corresponding graph for the sequence using the Havel-Hakimi method [39]. We use the range  $2 < \alpha < 3$  as suggested in [29] as this range of  $\alpha$  occurs most commonly in modeling of real-world networks. We generate graphs of 300,000 and 1M vertices denoted  $s300^{\alpha=x}$  and  $s1M^{\alpha=x}$  respectively, for  $x \in \{2.2, 2.4, 2.6, 2.8\}$ .

The three *real-world* data sets originate from articles that found the data to be well-approximated by a power-law. The WWW data set [6] contains information on links between webpages within the nd.edu domain. The ENRON data set [47] contains email communication between Enron employees (vertices are email addresses; there is a link between two addresses if a mail has been sent between them). The INTERNET data set [50] provides a snapshot the Internet structure at the level of autonomous systems, reconstructed from BGP tables. For all of these sets, we consider the underlying simple, undirected graphs. For each set, standard maximum likelihood methods were used to compute the parameter  $\alpha$  of the best-fitting power-law curve [29]. Additional information on the data sets can be found in Table 1.

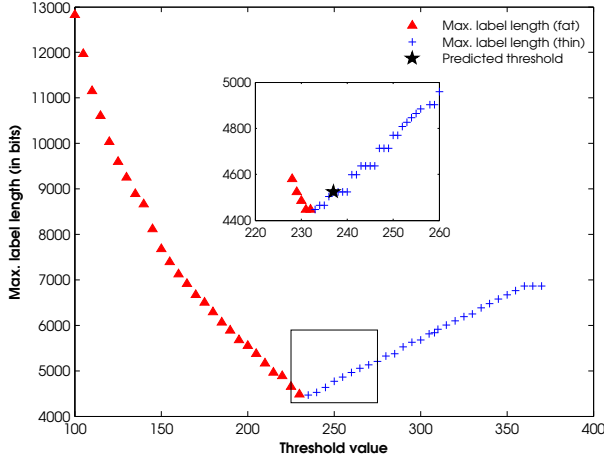
Real-Life					
Data set	$ V $	$ E $	$\alpha$	$\Delta_{\max}$	Source
WWW	325,729	1,117,563	2.16	10,721	[6]
ENRON	36,692	183,830	1.97	1,383	[47]
INTERNET	22,963	48,436	2.09	2,390	[50]
Synthetic					
$s1M^{\alpha=2.4}$	1,000,000	1,127,797	2.4	42,683	–
$s1M^{\alpha=2.6}$	1,000,000	878,472	2.6	12,169	–
$s1M^{\alpha=2.8}$	1,000,000	751,784	2.8	1,692	–
$s300^{\alpha=2.2}$	300,000	491,926	2.2	10,906	–
$s300^{\alpha=2.4}$	300,000	327,631	2.4	3,265	–
$s300^{\alpha=2.6}$	300,000	261,949	2.6	1,410	–
$s300^{\alpha=2.8}$	300,000	227,247	2.8	1,842	–

Table 1: Data sets and their properties. All graphs are undirected and simple.  $\Delta_{\max}$  is the maximum degree of any vertex in the data set.

### 8.2 Findings

Figure 1 shows the distribution of maximum label sizes for one synthetic and one real-world data set. The maximum label size for the predicted and empirical thresholds as well as upper bounds on the label sizes from different label schemes in the literature can be seen in Table 2 for two synthetic data sets and all three real-world data sets. Plots for the remaining data sets can be found in Appendix ??.

Table 2 shows the maximum label sizes achieved using different labeling schemes on our data sets. “Predicted” shows the experimental maximum label size obtained by running our scheme on the graphs, “Empirical” is the label size attained by using the empirical threshold. The remaining columns show non-experimental upper bounds for different label schemes: “Bound” is the upper bound guaranteed in Theorem 2, “C-sparse” is the labeling scheme for sparse graphs defined in Theorem 1, “BD” is the  $\lceil \frac{\Delta}{2} \rceil \lceil \log n \rceil$



(a)  $\text{syn300}^{\alpha=2.2}$

Figure 1: Maximum label sizes of different threshold values for the  $\text{syn300}^{\alpha=2.2}$  and ENRON data sets. The triangles and crosses represent that for the tested threshold the largest label belong to fat, resp. thin node. The star indicate the position of the predicted threshold.

Data set	Pre.	Emp.	Bound	$C$ -sp.	BD [3]	AKTZ
$\text{s1M}^{\alpha=2.4}$	4,841	4,821	25,012	30,079	426,820	500,006
$\text{s1M}^{\alpha=2.6}$	3,361	3,201	15,282	26,551	121,680	500,006
$\text{s1M}^{\alpha=2.8}$	2,101	2,061	10,081	24,566	16,920	500,006
$\text{s300}^{\alpha=2.2}$	4,523	4,447	24,878	18,885	103,607	150,006
$\text{s300}^{\alpha=2.4}$	2,775	2,680	14,404	15,420	31,008	150,006
$\text{s300}^{\alpha=2.6}$	1,958	1,920	9,151	13,792	13,395	150,006
$\text{s300}^{\alpha=2.8}$	1,350	1,312	6,244	12,849	17,499	150,006
WWW	5,245	3,060	29,225	28,445	101,840	162,870
ENRON	2,609	2,577	15,835	9,735	11,056	18,352
INTERNET	1,426	1,156	8,181	4,700	17,925	11,487

Table 2: Label size in bits of labeling schemes. The two leftmost columns are experimental results; the remaining are upper bounds on label sizes computed from the characteristics of the data sets. Pre. stands for predicted, Emp. for empirical, AKTZ for the labeling scheme in [11].

bounded degree graph labeling of [3], and AKTZ is the  $\lceil n/2 \rceil + 6$  general graph labeling of [11]. Both “Empirical” and “Bound” using simple concatenation of labels to represent the fat bit string<sup>4</sup>.

Our findings are as follows. For Performance Indicator (i), our labeling scheme obtains maximum label size at most 3% larger than what would have been obtained by using the empirical threshold for all synthetic data sets. This is expected—the synthetic data sets are graphs generated specifically to have power-law distributed degree distribution. For the real-world data sets, the labeling scheme obtains maximum label size at most 23% larger than by using the empirical threshold; this larger deviation is likely due to

<sup>4</sup>Our labeling schemes introduced in this paper all make use of a succinctly represented “fat bit string”; for our experiments, we use simple concatenation of labels instead of a bit string; this incurs a  $(\log n)/\alpha$  factor on the label size, but simplifies the implementation.

degree distributions of the data sets being close to, but not quite, power-law distributions due to natural phenomena or noise. E.g., for the ENRON data set there is sudden drop in frequency between nodes of degree  $< 158$  and  $\geq 158$ .

For Performance Indicator (ii), both our experimental results and theoretical upper bounds for our labeling scheme are several orders of magnitudes lower than for labeling schemes aimed at more general classes of graphs, as expected. Of the more general classes of graphs, it is most interesting to compare the upper bound of bounded degree graphs—the most restrictive class of graphs that both contains the class of power-law graphs and has an efficient labeling scheme described in the literature [3]. As seen in Table 2, the upper bound on our labeling schemes for both power-law graphs and sparse graphs have better upper bounds on label sizes, but only marginally so for data sets with low maximum degree and low values of the power-law parameter  $\alpha$ , e.g. ENRON ( $\alpha = 1.97$ ). It is interesting to note that the actual label sizes obtained in the experiments (the two leftmost columns of Table 2) are substantially lower than the upper bounds, that is, the labeling scheme performs much better in practice than suggested by theory (down to less than a kilobyte per vertex for all data sets). This phenomenon may be due to the degree distribution of the graphs of the data sets having only minor deviation from a power-law for small vertex degrees; our upper bounds on the label size are derived by using the very rich family  $\mathcal{P}_\alpha$  that allows very large deviation from a power-law for degrees between 1 and  $\sqrt[3]{n/\log n} - 1$ .

Finally, note that our labeling scheme supports adjacency for *directed* graphs by using one more bit per edge in each label to store the edge orientation. For data sets whose natural interpretation is as a directed graph (e.g., the WWW set where edges are outgoing and incoming links), the results of Table 2 thus carry over with just one more bit added to the numbers in the two leftmost columns.

## 9. CONCLUSION AND FUTURE WORK

We have devised adjacency labeling schemes for sparse graphs and graphs whose degree distribution approximately follows a power-law distribution. We have proven lower bounds for the class of power-law graphs showing that our labeling scheme is almost asymptotically optimal. Furthermore, we have shown experimentally that the labeling scheme for power-law graphs obtain results in practice requiring very little space (labels smaller than a kilobyte per vertex for real-world graphs with several hundreds of thousands of vertices).

### 9.1 Future work

It would be of interest to test the performance of the labeling scheme on more real-world data sets, and in particular investigating *dynamic* labeling schemes on such sets: if vertices can enter and exit the network, labels need to be recomputed efficiently. As our labeling scheme can be extended to handle directed graphs by using a single bit more per label, it would be interesting to investigate the overhead incurred by distributing the storage of the graph topology to the labels (as per our labeling scheme) compared to the substantial body of work on storing directed power-law graphs directly in main memory (so-called “web-graph compression”) [38,



14, 15, 28]. The label sizes attained in Sec. 8.1 can be reduced by using the succinctly represented “fat bit string” as well as an additional rule that prevents storing an edge in two labels; doing so would yield a small multiplicative reduction in label size, making our labeling scheme even more practical. Labeling schemes for other properties than adjacency may be investigated for power-law graphs, e.g. for distance as has been done for other classes of graphs [8] and briefly considered for power-law graphs in the context of routing algorithms [25]. Finally, labeling schemes for power law graphs can likely be devised for the realistic case where the scheme only has incomplete knowledge of the graph, for example when the expected frequency of vertices of each degree is known, but not the exact frequency of each vertex.

## 10. REFERENCES

- [1] I. Abraham, D. Delling, A. V. Goldberg, and R. F. Werneck. A hub-based labeling algorithm for shortest paths in road networks. In *Experimental Algorithms*, pages 230–241. Springer, 2011.
- [2] D. Achlioptas, A. Clauset, D. Kempe, and C. Moore. On the bias of traceroute sampling: Or, power-law degree distributions in regular graphs. *J. ACM*, 56(4), 2009.
- [3] D. Adjiashvili and N. Rotbart. Labeling schemes for bounded degree graphs. In *Automata, Languages, and Programming*, pages 375–386. Springer, 2014.
- [4] W. Aiello, F. Chung, and L. Lu. A random graph model for power law graphs. *Experimental Mathematics*, 10(1):53–66, 2001.
- [5] A. Akella, S. Chawla, A. Kannan, and S. Seshan. Scaling properties of the internet graph. In *Proceedings of the Twenty-Second ACM Symposium on Principles of Distributed Computing, PODC 2003*, pages 337–346, 2003.
- [6] R. Albert, H. Jeong, and A.-L. Barabási. Diameter of the world-wide web. *Nature*, 401(6749):130–131, 1999.
- [7] N. Alon and V. Asodi. Sparse universal graphs. *J. Comput. Appl. Math.*, 142(1):1–11, May 2002.
- [8] S. Alstrup, P. Bille, and T. Rauhe. Labeling schemes for small distances in trees. *SIAM J. Disc. Math.*, 19(2):448–462, 2005.
- [9] S. Alstrup, S. Dahlgaard, and M. B. T. Knudsen. Optimal induced universal graphs and adjacency labeling for trees. In *Proceedings of the 58th Symposium on Foundations of Computer Science, FOCS '15*, Washington, DC, USA, 2015. IEEE Computer Society.
- [10] S. Alstrup, S. Dahlgaard, M. B. T. Knudsen, and E. Porat. Sublinear distance labeling for sparse graphs. *CoRR*, abs/1507.02618, 2015.
- [11] S. Alstrup, H. Kaplan, M. Thorup, and U. Zwick. Adjacency labeling schemes and induced-universal graphs. *To appear in the 47th symposium on Theory of computing (STOC)*, 2015.
- [12] S. Alstrup and T. Rauhe. Small induced-universal graphs and compact implicit graph representations. In *Proceedings of the 43rd Symposium on Foundations of Computer Science, FOCS '02*, pages 53–62, Washington, DC, USA, 2002. IEEE Computer Society.
- [13] S. R. Arikati, A. Maheshwari, and C. D. Zaroliagis. Efficient computation of implicit representations of sparse graphs. *Discrete Applied Mathematics*, 78(1):1–16, 1997.
- [14] Y. Asano, T. Ito, H. Imai, M. Toyoda, and M. Kitsuregawa. Compact encoding of the web graph exploiting various power laws. In *Advances in Web-Age Information Management*, pages 37–46. Springer, 2003.
- [15] Y. Asano, Y. Miyawaki, and T. Nishizeki. Efficient compression of web graphs. In *Computing and Combinatorics*, pages 1–11. Springer, 2008.
- [16] L. Babai, F. R. Chung, P. Erdős, R. L. Graham, and J. Spencer. On graphs which contain all sparse graphs. *Ann. Discrete Math*, 12:21–26, 1982.
- [17] A.-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- [18] S. Bhatt, F. R. Graham Chung, T. Leighton, and A. Rosenberg. Universal Graphs for Bounded-Degree Trees and Planar Graphs. *SIAM Journal on Discrete Mathematics*, 2(2):145–155, 1989.
- [19] P. Boldi, M. Rosa, M. Santini, and S. Vigna. Layered label propagation: A multiresolution coordinate-free ordering for compressing social networks. In *Proceedings of the 20th international conference on World Wide Web*, pages 587–596. ACM, 2011.
- [20] P. Boldi and S. Vigna. The webgraph framework i: compression techniques. In *Proceedings of the 13th international conference on World Wide Web*, pages 595–602. ACM, 2004.
- [21] B. Bollobás, O. Riordan, J. Spencer, and G. E. Tusnády. The degree sequence of a scale-free random graph process. *Random Struct. Algorithms*, 18(3):279–290, 2001.
- [22] A. Brady and L. J. Cowen. Compact routing on power law graphs with additive stretch. In *ALLENEX*, volume 6, pages 119–128. SIAM, 2006.
- [23] K. L. Calvert, M. B. Doar, and E. W. Zegura. Modeling internet topology. *Communications Magazine, IEEE*, 35(6):160–163, 1997.
- [24] S. Caminiti, I. Finocchi, and R. Petreschi. Engineering tree labeling schemes: A case study on least common ancestors. In *Algorithms-ESA 2008*, pages 234–245. Springer, 2008.
- [25] W. Chen, C. Sommer, S.-H. Teng, and Y. Wang. A compact routing scheme and approximate distance oracle for power-law graphs. *ACM Transactions on Algorithms*, 9(1):4, 2012.
- [26] F. Chung and L. Lu. The average distance in a random graph with given expected degrees. *Internet Mathematics*, 1(1):91–113, 2004.
- [27] F. R. Chung and L. Lu. *Complex Graphs and Networks*, volume 107. American mathematical society Providence, 2006.
- [28] F. Claude and G. Navarro. Fast and compact web graph representations. *ACM Transactions on the Web (TWEB)*, 4(4):16, 2010.
- [29] A. Clauset, C. R. Shalizi, and M. E. Newman. Power-law distributions in empirical data. *SIAM review*, 51(4):661–703, 2009.
- [30] E. Cohen, H. Kaplan, and T. Milo. Labeling dynamic xml trees. *SIAM Journal on Computing*, 39(5):2048–2074, 2010.

- [31] S. Dahlgaard, M. B. T. Knudsen, and N. Rotbart. Dynamic and multi-functional labeling schemes. In *Algorithms and Computation*, pages 141–153. Springer, 2014.
- [32] J. Fischer. Short labels for lowest common ancestors in trees. In *Algorithms-ESA 2009*, pages 752–763. Springer, 2009.
- [33] C. Gavaille and A. Labourel. Shorter implicit representation for planar graphs and bounded treewidth graphs. In *Algorithms-ESA 2007*, pages 582–593. Springer, 2007.
- [34] C. Gavaille, D. Peleg, S. Pérennec, and R. Razb. Distance labeling in graphs. *Journal of Algorithms*, 53:85–112, 2004.
- [35] P. Gawrychowski, A. Kosowski, and P. Uznanski. Even simpler distance labeling for (sparse) graphs. *CoRR*, abs/1507.06240, 2015.
- [36] G. Goel and J. Gustedt. Bounded arboricity to determine the local structure of sparse graphs. In *Graph-Theoretic Concepts in Computer Science*, pages 159–167. Springer, 2006.
- [37] J. E. Gonzalez, Y. Low, H. Gu, D. Bickson, and C. Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *OSDI*, volume 12, page 2, 2012.
- [38] J.-L. Guillaume, M. Latapy, and L. Viennot. Efficient and simple encodings for the web graph. In *Advances in Web-Age Information Management*, pages 328–337. Springer, 2002.
- [39] S. L. Hakimi. On realizability of a set of integers as degrees of the vertices of a linear graph. i. *Journal of the Society for Industrial & Applied Mathematics*, 10(3):496–506, 1962.
- [40] S. Kannan, M. Naor, and S. Rudich. Implicit representation of graphs. In *SIAM Journal On Discrete Mathematics*, pages 334–343, 1992.
- [41] M. Katz, N. A. Katz, A. Korman, and D. Peleg. Labeling schemes for flow and connectivity. *SIAM Journal on Computing*, 34(1):23–40, 2004.
- [42] A. Korman. General compact labeling schemes for dynamic trees. *Distributed Computing*, 20(3):179–193, 2007.
- [43] A. Korman and D. Peleg. Compact separator decompositions in dynamic trees and applications to labeling schemes. In *Distributed Computing*, pages 313–327. Springer, 2007.
- [44] A. Korman and D. Peleg. Labeling schemes for weighted dynamic trees. *Inf. Comput.*, 205(12):1721–1740, Dec. 2007.
- [45] L. Kowalik. Approximation scheme for lowest outdegree orientation and graph density measures. In *Algorithms and computation*, pages 557–566. Springer, 2006.
- [46] D. Krioukov, K. Fall, and X. Yang. Compact routing on internet-like graphs. In *INFOCOM 2004. Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies*, volume 1. IEEE, 2004.
- [47] J. Leskovec, K. J. Lang, A. Dasgupta, and M. W. Mahoney. Community structure in large networks: Natural cluster sizes and the absence of large well-defined clusters. *Internet Mathematics*, 6(1):29–123, 2009.
- [48] M. Mitzenmacher. A brief history of generative models for power law and lognormal distributions. *Internet Mathematics*, 1(2):226–251, 2004.
- [49] J. Moon. On minimal n-universal graphs. In *Proceedings of the Glasgow Mathematical Association*, volume 7, pages 32–33. Cambridge University Press, 1965.
- [50] M. Newman. Network data. <http://www-personal.umich.edu/~mejn/netdata/>, 2013. [Online; accessed 02-Jan-2015].
- [51] N. Rotbart, M. V. Salles, and I. Zotos. An evaluation of dynamic labeling schemes for tree networks. In *Experimental Algorithms*, pages 199–210. Springer, 2014.
- [52] G. Siganos, M. Faloutsos, P. Faloutsos, and C. Faloutsos. Power laws and the as-level internet topology. *IEEE/ACM Trans. Netw.*, 11(4):514–524, 2003.
- [53] J. P. Spinrad. *Efficient graph representations*. American mathematical society, 2003.
- [54] B. M. Waxman. Routing of multipoint connections. *Selected Areas in Communications, IEEE Journal on*, 6(9):1617–1622, 1988.