

---

# An Information-theoretic Perspective of Hierarchical Clustering

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A combinatorial cost function for hierarchical clustering was introduced by Das-  
2 gupta [10]. It has received great attention and several new cost functions from sim-  
3 ilar combinatorial perspective have been proposed. In this paper, we investigate  
4 hierarchical clustering from the *information-theoretic* perspective and formulate  
5 a new objective function. We also establish the relationship between these two  
6 perspectives. In algorithmic aspect, we present two algorithms for expander-like  
7 and well-clustered cardinality weighted graphs, respectively, and show that both  
8 of them achieve  $O(1)$ -approximation for our new objective function. For practi-  
9 cal use, we consider non-binary hierarchical clustering problem. We get rid of  
10 the traditional top-down and bottom-up frameworks, and present a new one. Our  
11 new framework stratifies the sparsest level of a cluster tree recursively in guide  
12 with our objective function. Our algorithm called HCSE outputs a  $k$ -level cluster  
13 tree by an interpretable mechanism to choose  $k$  automatically without any hyper-  
14 parameter. Our experimental results on synthetic datasets show that HCSE has  
15 its own superiority in finding the intrinsic number of hierarchies, and the results  
16 on real datasets show that HCSE also achieves competitive costs over the popular  
17 non-binary hierarchical clustering algorithms LOUVAIN and HLP.

## 18 1 Introduction

19 Hierarchical clustering for graphs plays an important role in the structural analysis of a given data  
20 set. Understanding hierarchical structures on the levels of multiple granularities is fundamental in  
21 various disciplines including artificial intelligence, physics, biology, sociology, etc [4, 11, 13, 9].  
22 Hierarchical clustering requires a cluster tree that represents a recursive partitioning of a graph into  
23 smaller clusters as the tree nodes get deeper. A leaf represents a graph node while a non-leaf node  
24 represents a cluster containing its descendant leaves. The root is the largest one containing all leaves.

25 Clustering is usually formulated as an optimization problem with some objective function. For hier-  
26 archical clustering, no cost function with a clear and reasonable combinatorial explanation was de-  
27 veloped until Dasgupta [10] introduced a cost function for cluster trees. In this definition, similarity  
28 or dissimilarity between data points is represented by weighted edges. Taking the similarity-based  
29 metrics as an example, a cluster is a set of nodes with relatively denser intra-links compared with its  
30 inter-links, and in a good cluster tree, heavier edges tend to connect leaves whose lowest common  
31 ancestor (LCA) is as deep as possible. This intuition leads to Dasgupta’s cost function that is a  
32 bilinear combination of edge weights and the sizes of corresponding LCAs.

33 Motivated by Dasgupta’s cost function, Cohen-Addad et al. [8] proposed admissible cost functions.  
34 In their definition, the size of each LCA in Dasgupta’s objective is generalized to be a function of  
35 the sizes of its left and right children. For all similarity-based graphs generated from a minimal  
36 ultrametric, a cluster tree achieves the minimum cost if and only if it is a generating tree that is a

37 “natural” ground truth tree in an axiomatic sense therein. A necessary condition of admissibility of  
38 an objective function is that it achieves the same value for every cluster tree for a uniformly weighted  
39 clique that has no structure in common sense. However, any slight deviation of edge weights would  
40 generally separate the two end-points of a light edge on a high level of its optimal (similarity-based)  
41 cluster tree. Thus, it seems that admissible objective functions, which take Dasgupta’s cost function  
42 as a specific form, ought to be an unchallenged criterion in evaluating cluster trees since they are  
43 formulated by an axiomatic approach.

44 However, an admissible cost function seems imperfect in practice. The arbitrariness of optima of  
45 cluster trees for cliques indicates that the division of each internal node on an optimal cluster tree  
46 totally neglects the *balance* of its two children. Edge weight is the unique factor that decides the  
47 structure of optimal trees. But a balanced tree is commonly considered as an ideal candidate in  
48 hierarchical clustering compared to an unbalanced one. Even clustering for cliques, a balanced  
49 partition should be preferable for each internal node. At least, an optimal cluster tree whose height  
50 is logarithm of graph size  $n$  is intuitively more reasonable than a caterpillar shaped cluster tree  
51 whose height is  $n - 1$ . Moreover, a simple proof would imply that the optimal cluster tree for any  
52 connected graphs is binary. This property is not always useful in practice since a real system usually  
53 has its inherent number of hierarchies and a natural partition for each internal cluster. For instance,  
54 the natural levels of administrative division in a country is usually intrinsic, and it is not suitable to  
55 differentiate hierarchies for parallel cities in the same state. This structure cannot be obtained by  
56 simply minimizing admissible cost functions.

57 In this paper, we investigate the hierarchical clustering from the perspective of information theory.  
58 Our study is based on Li and Pan’s structural information theory [14] whose core concept named  
59 structural entropy measures the complexity of hierarchical networks. We summarize our contribu-  
60 tions as follows.

61 (1) We formulate **a new objective function from the information-theoretic perspective**, which  
62 builds the bridge for combinatorial and information-theoretic perspectives for hierarchical cluster-  
63 ing. For this cost function, the balance of cluster trees will be involved naturally as a factor just  
64 like we design optimal codes, for which the balance of probability over objects is fundamental in  
65 constructing an efficient coding tree. We also define cluster trees with a specific height, which is  
66 coincident with our cognition of natural clustering.

67 (2) For our new objective function, we present **two polynomial-time approximation algorithms**  
68 respectively for two cases of the conductance  $\Phi(G)$  of a cardinality weighted graph  $G$ . Our first  
69 result shows that *any* cluster tree of  $G$  has a approximation factor  $O(\Phi(G)^{-1})$  (Theorem 3.1). So  
70 a "Huffman-merge" heuristic that solely depends on the degrees of vertices achieves this guaran-  
71 tee, and it achieves  $O(1)$ -approximation when  $\Phi(G)$  is a constant. The second result is a  $O(1)$ -  
72 approximation algorithm for  $G$  that can be well clustered into a constant number of expanders (The-  
73 orem 3.2). The main idea of this algorithm is inspired by very recent Manghiuc and Sun’s work [15],  
74 and our approximation factors for our new objective also match their results in these two cases.

75 (3) For practical use, we develop **a new interpretable framework for natural hierarchical clus-**  
76 **tering** that outputs a non-binary cluster tree. The idea of our framework is essentially different from  
77 the traditional recursive division or agglomeration ones. In our framework, the *sparsest level* of the  
78 cluster tree is stratified recursively. This coincide with the intuition that when we differentiate the  
79 hierarchies of a complex system, the clearest level should be stratified first, rather than in a rigid  
80 divisive or agglomerative fashion. Therefore, this framework has much better interpretability than  
81 the traditional ones.

82 (4) We develop **a new non-binary clustering algorithm (HCSE)** under the new clustering frame-  
83 work. To find the sparsest level in each iteration, we formulate two basic operations called *stretch*  
84 and *compress*, respectively. HCSE terminates when a specific criterion that intuitively coincides  
85 with the natural hierarchies is met, and *no* hyperparameter is needed. Our extensive experiments on  
86 both synthetic and real datasets demonstrate that HCSE outperforms the present popular heuristic  
87 algorithms LOUVAIN [3] and HLP [19]. These two algorithms proceed simply by recursively in-  
88 voking flat clustering algorithms based on modularity and label propagation, respectively, and the  
89 hierarchy number is solely determined by the number of rounds when the algorithm terminates. So  
90 their interpretability is quite poor. Our experimental results on synthetic datasets show that HCSE  
91 has a great advantage in finding the intrinsic number of hierarchies, and the results on real datasets  
92 show that HCSE achieves much better costs than HLP and competitive costs to LOUVAIN.

93 **Related work.** The first combinatorial objective function was proposed by Dasgupta [10]. Along  
 94 with this line of study, several alternative objectives have been presented. All of them are bilinear  
 95 functions of edge weights and some function of the corresponding LCAs. For Dasgupta’s cost func-  
 96 tion and for the worst case study, Dasgupta showed that a recursively bipartition applying Arora’s  
 97 seminal algorithm for sparsest cut problem [2] yields  $O(\log^{1.5} n)$ -approximation, and it was im-  
 98 proved by [20] and [5, 8] to  $O(\log n)$  and  $\sqrt{\log n}$ , respectively. It is NP-hard to optimize the cluster  
 99 tree [10] and even a  $O(1)$ -approximation is impossible under the Small Set Expansion hypothesis  
 100 [20, 5]. Beyond the worst case, Cohen-Addad et al. [8] showed that a SVD-based algorithm achieves  
 101 a  $O(1 + o(1))$ -approximation for the stochastic block model with high probability. Manghiuc and  
 102 Sun [15] presented a  $O(1)$ -approximation algorithm for more generalized well-clustered graphs. The  
 103 outline of their method is to utilize a flat clustering algorithm [12] to obtain the underlying clusters  
 104 first, and then some relatively easy heuristics for clustering in and out of these clusters are enough  
 105 for the guarantee. Our proof follows this route also.

106 For other lines of this study, Moseley and Wang [16] studied the dual of Dasgupta’s cost function  
 107 and showed that the average-linkage algorithm achieves a  $(1/3)$ -approximation. This factor has  
 108 been improved by a series of works to 0.336 [6], 0.4246 [7] and 0.585 [1], respectively. Cohen-  
 109 Addad et al. [8] considered maximization of Dasgupta’s cost function for the dissimilarity-based  
 110 metrics. They proved that the average-link and random partitioning algorithms achieve a  $(2/3)$ -  
 111 approximation, which has been improved to 0.667 [6], 0.716 [18] and 0.74 [17], respectively.

112 For non-binary cluster tree construction, the most popular algorithm for practical use is LOUVAIN  
 113 [3]. More recently, a hierarchical label propagation based algorithm HLP has been presented [19].  
 114 Both of these two algorithms construct a non-binary cluster tree with the same framework, that is,  
 115 the hierarchies are formed from bottom to top one by one. In each round, they invoke different flat  
 116 clustering algorithms, Modularity and Label Propagation, respectively.

## 117 2 A cost function from information-theoretic perspective

118 In this section, we introduce Li and Pan’s structural information theory [14] and the combinatorial  
 119 cost functions of Dasgupta [10] and Cohen-Addad et al. [8]. Then we propose a new cost func-  
 120 tion that is developed from structural information theory and establish the relationship between the  
 121 information-theoretic and combinatorial perspectives.

122 **Notations.** Let  $G = (V, E, w)$  be an undirected weighted graph with a set of vertices  $V$ , a set of  
 123 edges  $E$  and a weight function  $w : E \rightarrow \mathbb{R}^+$ , where  $\mathbb{R}^+$  denotes the set of all positive real numbers.  
 124 An unweighted multigraph can be viewed as a cardinality weighted one whose edge weight is the  
 125 number of parallel edges. For each vertex  $u \in V$ , denote by  $d_u = \sum_{(u,v) \in E} w(u, v)$  the weighted  
 126 degree of  $u$ . For a subset of vertices  $S \subseteq V$ , define the volume of  $S$  to be the sum of degrees of  
 127 vertices. We denote it by  $\text{vol}(S) = \sum_{u \in S} d_u$ . We denote by  $G[S]$  the subgraph induced by  $S$ . A  
 128 cluster tree  $T$  for graph  $G$  is a rooted tree with  $|V|$  leaves, each of which is labeled by a distinct  
 129 vertex  $v \in V$ . Each non-leaf node on  $T$  is labeled by a subset  $S$  of  $V$  that consists of all the leaves  
 130 treating  $S$  as an ancestor. For each node  $\alpha$  on  $T$ , denote by  $\alpha^-$  the parent of  $\alpha$ , and by  $|\alpha|$  its size.  
 131 For each pair of leaves  $u$  and  $v$ , denote by  $u \vee v$  the LCA of them on  $T$ .

132 **Structural entropy of graphs.** Because of the tense space limit, we just give the definition of the  
 133 core concept structural entropy in structural information theory. The idea of this definition is briefly  
 134 introduced in Appendix A. Readers could also refer to [14] for more information on this theory.

135 Given a weighted graph  $G = (V, E, w)$  and a cluster tree  $T$  for  $G$ , the *structural entropy of  $G$  on  $T$*   
 136 is defined as

$$\mathcal{H}^T(G) = - \sum_{\alpha \in T} \frac{g_\alpha}{\text{vol}(V)} \log \frac{\text{vol}(\alpha)}{\text{vol}(\alpha^-)},^1$$

137 where  $\alpha^-$  denotes the parent of tree node  $\alpha$ , and  $g_\alpha$  denotes the sum of weights of edges in  $G$   
 138 with exactly one end-point in the set of vertices corresponding to  $\alpha$ . The *structural entropy of  $G$*  is  
 139 defined as the minimum one among all cluster trees, denoted by  $\mathcal{H}(G) = \min_T \{\mathcal{H}^T(G)\}$ .

<sup>1</sup>For notational convenience, for the root  $\lambda$  of  $T$ , set  $\lambda^- = \lambda$ . So the term for  $\lambda$  in the summation is 0. In this paper, the omitted base of logarithm is always 2.

140 **Combinatorial explanation of structural entropy.** The cost function of a cluster tree  $T$  for graph  
 141  $G = (V, E)$  introduced by Dasgupta [10] is defined to be  $c^T(G) = \sum_{(u,v) \in E} w(u,v)|u \vee v|$ . The  
 142 admissible cost function introduced by Cohen-Addad et al. [8] generalizes the term  $|u \vee v|$  in the  
 143 definition of  $c^T(G)$  to be a general function  $g(|L|, |R|)$ , where  $L$  and  $R$  are the two children of  
 144  $u \vee v$ , respectively. Dasgupta defined  $g(x, y) = x + y$ . For both definitions, the optimal hierarchical  
 145 clustering of  $G$  is in correspondence with a cluster tree of minimum cost in the combinatorial sense  
 146 that heavy edges are cut as far down the tree as possible. The following proposition establishes the  
 147 relationship between structural entropy and this kind of combinatorial form of cost functions.

148 **Proposition 2.1.** *For a weighted graph  $G = (V, E, w)$ , minimizing  $\mathcal{H}^T(G)$  (over  $T$ ) is equivalent*  
 149 *to minimizing the cost function*

$$\text{cost}^T(G) = \sum_{(u,v) \in E} w(u,v) \log \text{vol}(u \vee v). \quad (1)$$

150 We defer the proof of Proposition 2.1 to Appendix B. We call  $\text{cost}(\text{SE})$  the cost function in Propo-  
 151 sition 2.1 from now on. Proposition 2.1 indicates that when we view  $g$  as a function of the LCA  
 152 rather than that of its size and define  $g(u, v) = \log \text{vol}(u \vee v)$ , the “admissible” function becomes  
 153 equivalent to structural entropy in evaluating cluster trees, although it is not admissible any more.

154 So what is the difference between these two cost functions? As stated by Cohen-Addad et al. [8],  
 155 an important axiomatic hypothesis for admissible function, thus also for Dasgupta’s cost function,  
 156 is that the cost for every binary cluster tree of an unweighted clique is identical. So any binary tree  
 157 for clustering on cliques is reasonable, which coincides with the common sense that structureless  
 158 datasets can be organized hierarchically free. However, for structural entropy, the following theorem  
 159 indicates that balanced organization is of importance even though for structureless dataset.

160 **Proposition 2.2.** *For any positive integer  $n$ , let  $K_n$  be the clique of  $n$  vertices with identical weight*  
 161 *on every edge. Then a cluster tree  $T$  of  $K_n$  achieves minimum structural entropy if and only if  $T$  is*  
 162 *a balanced binary tree, that is, the two children clusters of each sub-tree of  $T$  have difference in size*  
 163 *at most 1.*

164 The proof of Proposition 2.2 is a bit technical, and we defer it to Appendix C. The intuition behind  
 165 Proposition 2.2 is that balanced codes are the most efficient encoding scheme for unrelated data. So  
 166 the codewords of the random walk that jumps freely among clusters on each level of a cluster tree  
 167 have the minimum average length if all the clusters on this level are in balance.

168 It is worth noting that the admissible function introduced by Cohen-Addad et al. [8] is defined  
 169 from the viewpoint that a generating tree  $T$  of a similarity-based graph  $G$  that is generated from a  
 170 minimal ultrametric achieves the minimum cost. In this definition, the monotonicity of edge weights  
 171 between clusters on each level from bottom to top on  $T$ , which is given by Cohen-Addad et al. [8] as  
 172 a property of a “natural” ground-truth hierarchical clustering, is the unique factor when evaluating  $T$ .  
 173 However, Proposition 2.2 implies that for  $\text{cost}(\text{SE})$ , besides cluster weights, the balance of cluster  
 174 trees is implicitly involved as another factor. Moreover, for cliques, the minimum cost should be  
 175 achieved on every subtree, which makes an optimal cluster tree balanced everywhere. This optimal  
 176 clustering for cliques is also robust in the sense that a slight perturbation to the minimal ultrametric,  
 177 which can be considered as slight variations to the weights of a batch of edges, will not change the  
 178 optimal cluster tree structure wildly due to the holdback force of balance.

### 179 3 Approximation algorithms for SE-based cost function

180 In this section, we present approximation algorithms for expander-like and well-clustered graphs,  
 181 respectively. These algorithms work for cardinality edge weights (e.g. the multiplicity of edges).

182 **Why cardinality weights?** In general, the term  $\log \text{vol}(u \vee v)$  in Eq. 1 and thus  $\text{cost}(\text{SE})$  may be  
 183 negative when the volume of  $u \vee v$  varies, which may lead to pathosis in approximation analysis.  
 184 The cardinality weight function  $w$  is at least one, which makes  $\text{cost}(\text{SE})$  non-negative. The depen-  
 185 dence of  $\text{cost}(\text{SE})$  on the scale of edge weights violates the scale-invariance principle. However, we  
 186 emphasize that  $\mathcal{H}^T(G)$  is scale-invariant and Proposition 2.1 holds for any scale variation. In this  
 187 paper, we present approximation algorithms for  $\text{cost}(\text{SE})$  in well-defined settings.

188 **Theorem 3.1.** *For any cardinality weighted graph  $G = (V, E, w)$  with conductance  $\Phi(G)$ , it holds*  
 189 *that any cluster tree has a cost  $O(\Phi(G)^{-1}) \cdot \text{OPT}$ , where  $\text{OPT}$  is the minimum  $\text{cost}(\text{SE})$  of  $G$ .*

190 We defer the proof of Theorem 3.1 to Appendix D. When  $\Phi(G)$  is a constant, Theorem 3.1 implies  
 191 that any cluster tree achieves  $O(1)$ -approximation for expanders. Thus, the balance of a cluster  
 192 tree has a significant impact on its cost. Considering balance as an important factor, we present a  
 193 Huffman-merging heuristic (Algorithm 1). It will serve as a subroutine for the algorithm for well-  
 194 clustered graphs.

---

**Algorithm 1:** HuffmanMerge

---

**Input:** a graph  $G = (V, E, w)$

**Output:** A cluster tree  $T$  of  $G$

- 1 Create  $n$  singleton trees;
  - 2 **while** there are at least two trees **do**
  - 3     Select the two trees  $T_1$  and  $T_2$  with the least volumes;
  - 4     Construct a new tree  $T_0$  with  $T_1$  and  $T_2$  as two subtrees of the root;
  - 5 **Return** the resulting binary tree  $T_0$ .
- 

195 Next, we consider well-clustered graphs that are composed by a collection of densely-connected  
 196 components with high inner conductance and weakly interconnections. Our settings for well-  
 197 clustered graphs is the same as those in [15]. We start from the following  $(\Phi_{in}, \Phi_{out})$ -decomposition  
 198 presented by Gharan and Trevisan [12]. Let  $\lambda_k$  be the  $k$ -th smallest eigenvalue of the normalized  
 199 Laplacian matrix of  $G$  and  $\Phi_G(S)$  be the conductance of a vertex set  $S$  in  $G$ .

200 **Lemma 3.1.** ([12], Theorem 1.5) *Let  $G = (V, E, w)$  be a graph such that  $\lambda_k > 0$ , for some  $k \geq 1$ .  
 201 Then, there is a local search algorithm that finds a  $l$ -partition  $\{P_i\}_{i=1}^l$  of  $V$ , for some  $l < k$ , such  
 202 that for every  $1 \leq i \leq l$ ,  $\Phi_G(P_i) = \mathcal{O}(k^6 \sqrt{\lambda_{k-1}})$  and  $\Phi(G[P_i]) = \Omega(\lambda_k^2/k^4)$ .*

203 Lemma 3.1 implies that, when graph  $G$  exhibits a clear clustering structure, there is a partition  
 204  $\{P_i\}_{i=1}^l$  of  $V$  such that for each  $P_i$  both the outer and inner conductance can be bounded. This is  
 205 one of the most crucial insights that we can use  $\{P_i\}_{i=1}^l$  directly to construct a cluster tree.

206 For a high-level description, our algorithm consists of two phases: Partition and Merge. In the  
 207 Partition phase, it invokes the algorithm in Lemma 3.1 to partition  $V$  into sets  $\{P_i\}_{i=1}^l$ . In the Merge  
 208 phase it combines the trees in a "caterpillar style" according to an increasing order of their volumes.  
 209 This algorithm is described as Algorithm 2.

---

**Algorithm 2:** CaterpillarMerge

---

**Input:** A graph  $G = (V, E, w)$ , an integer  $k \geq 2$  such that  $\lambda_k > 0$

**Output:** A cluster tree  $T$  of  $G$

- 1 Apply the partitioning algorithm in Lemma 3.1 on input  $(G, k)$  to obtain  $\{P_i\}_{i=1}^l$  for some  
 $l < k$ ;
  - 2 Sort  $P_1, \dots, P_l$  be such that  $\text{vol}_G(P_i) \leq \text{vol}_G(P_{i+1})$ , for all  $1 \leq i < l$ ;
  - 3 Let  $T_i = \text{HuffmanMerge}(G[P_i])$ ;
  - 4 Initialize  $T = T_1$ ;
  - 5 **for**  $i = 2, \dots, l$  **do**
  - 6     Let  $T$  be the tree with  $T$  and  $T_i$  as its two children;
  - 7 **Return**  $T$ .
- 

210 Note that Algorithm 2 degenerates to Algorithm 1 when  $k = 2$ . For the approximation guarantee,  
 211 we have the following theorem.

212 **Theorem 3.2.** *Let  $G = (V, E, w)$  be a cardinality weighted graph such that  $\lambda_k > 0$  for some  
 213  $k \geq 1$ . Then Algorithm 2 constructs in polynomial time a cluster tree  $T$  of  $G$  that achieves  
 214  $O\left(\frac{1}{(1-\alpha)^\beta} \log \frac{k}{1-\alpha}\right)$ -approximation for  $\text{cost}^T(G)$ , where  $\alpha = O(k^6 \sqrt{\lambda_{k-1}})$ ,  $\beta = \Omega(\lambda_k^2/k^4)$ . Con-  
 215 sequently, when  $\lambda_k = \Omega(1/\text{poly}(k))$  and  $\lambda_{k-1} = O(1/k^{12})$  such that  $\alpha < 1 - \rho$  for some constant  
 216  $\rho \in (0, 1)$ , Algorithm 2 achieves  $O(\text{poly}(k))$ -approximation. In addition, when  $k$  is a constant,  
 217 Algorithm 2 achieves  $O(1)$ -approximation.*

218 The proof of Theorem 3.2 is given in Appendix E.

219 **4 Practically used non-binary hierarchical clustering algorithm**

220 In this section, we develop a non-binary hierarchical clustering algorithm based on cost(SE) opti-  
 221 mization. At present, all existing algorithms for hierarchical clustering can be categorized into two  
 222 frameworks: top-down division and bottom-up agglomeration [8]. The top-down division approach  
 223 usually yields a binary tree by recursively dividing a cluster into two parts with a cut-related cri-  
 224 terion. But a binary clustering tree is far from a practical one as we introduced in Section 1. For  
 225 practical use, bottom-up agglomeration that is also known as hierarchical agglomerative clustering  
 226 (HAC) is commonly preferable. It constructs a cluster tree from leaves to the root recursively, during  
 227 each round of which the newly generated clusters shrink into single vertices.

228 Our algorithm jumps out of these two frameworks. We establish a new one that stratifies the *sparsest*  
 229 level of a cluster tree recursively rather than in a sequential order. In general, in guide with cost(SE),  
 230 we construct a  $k + 1$ -level cluster tree from the previous  $k$ -level one, during which we find the level  
 231 whose stratification makes the average cost in a local reduced subgraph decrease most, and then  
 232 differentiate it into two levels. The process of stratification consists of two basic operations: *stretch*  
 233 and *compression*. Generally speaking, in stretch steps, given an internal node of a cluster tree, a  
 234 local binary subtree is constructed, while in compression steps, the paths that are overlength from  
 235 the root to leaves on the binary tree is compressed by shrinking tree edges that make the cost reduce  
 236 most. The intuition behind the “stretch-and-compress” scheme is as follows. First, we run a fast and  
 237 simple, but probably rough clustering algorithm to obtain a binary cluster subtree. So intuitively,  
 238 after stretch, we unfold all the potential hierarchies such that the sparsest level is possibly to be  
 239 seen. Second, we compress every overlength path that is supposed to get through each level of this  
 240 subtree, during which, the edge on the sparsest level whose compression makes too many graph  
 241 edges amplify the sizes of their LCAs to a large extent will be retained.

242 We remark that this framework can be collocated with any cost function and any binary cluster tree  
 243 algorithm. For computational efficiency, especially for real networks of large scale more than  $10^4$ ,  
 244 we will adopt in our experiments an HAC construction of binary cluster trees in stretch steps.

245 **Stretch and compress.** Given a cluster tree  $T$  for graph  $G = (V, E)$ , let  $u$  be an internal node  
 246 on  $T$  and  $v_1, v_2, \dots, v_\ell$  be its children. We call this local parent-children structure rooted at  $u$  to  
 247 be a  $u$ -triangle of  $T$ , denoted by  $T_u$ . These two operations are defined on  $u$ -triangles. Note that  
 248 each child  $v_i$  of  $u$  is a cluster in  $G$ . We reduce  $G$  by shrinking each  $v_i$  to be a single vertex  $v'_i$   
 249 while maintaining each inter-link and ignoring each internal edge of  $v_i$ . This reduction captures the  
 250 connections of clusters at this level in the parent cluster  $u$ . The stretch operation constructs a binary  
 251 tree for  $u$ -triangle. We adopt a common HAC construction in this  $u$ -triangle. That is, initially, view  
 252 each  $v'_i$  as a cluster and recursively combine two clusters into a new one for which cost(SE) drops  
 253 most. The sequence of combinations yields a binary subtree  $T'_u$  rooted at  $u$  which has  $v_1, v_2, \dots, v_\ell$   
 254 as leaves. Then the compression operation is proposed to reduce the height of  $T'_u$  to be 2. Let  $\hat{E}(T')$   
 255 be the set of edges on  $T'$ , each of which appears on a path of length more than 2 from the root of  
 256  $T'$  to some leaf. Denote by  $\Delta(e)$  for edge  $e$  be the amount of structural entropy enhanced by the  
 257 shrink of  $e$ . We pick from  $\hat{E}(T'_u)$  the edge  $e$  with least  $\Delta(e)$ . Note that the compression of a tree  
 258 edge makes the grandchildren of some internal node to be children, which must amplify the cost.  
 259 The compression operation picks the least amplification. The processes of stretch and compress are  
 260 illustrated in Figure 3 and stated in Algorithms 5 and 6, respectively (see Appendix G).

261 **Sparsest level.** Let  $U_j$  be the set of  $j$ -level nodes on cluster tree  $T$ , that is,  $U_j$  is the set of  
 262 nodes each of which has distance  $j$  from  $T$ 's root. Suppose that the height of  $T$  is  $k$ , then  
 263  $U_0, U_1, \dots, U_{k-1}$  is a partition for all internal nodes of  $T$ . For each internal node  $u$ , define  
 264  $\mathcal{H}(u) = -\sum_{v:v^--=u} \frac{g_u}{\text{vol}(V)} \log \frac{\text{vol}(v)}{\text{vol}(u)}$ . Note that  $\mathcal{H}(u)$  is the partial sum contributed by  $u$  in  $\mathcal{H}^T(G)$ .  
 265 After a “stretch-and-compress” round on  $u$ -triangle, denote by  $\Delta\mathcal{H}(u)$  the structural entropy by  
 266 which the new cluster tree reduces. Since the reconstruction of  $u$ -triangle stratifies cluster  $u$ ,  $\Delta\mathcal{H}(u)$   
 267 is always non-negative. Define the sparsity of  $u$  to be  $\text{Spar}(u) = \frac{\Delta\mathcal{H}(u)}{\mathcal{H}(u)}$ , which is the relative varia-  
 268 tion of structural entropy in cluster  $u$ . From the information-theoretic perspective, this means that the  
 269 uncertainty of random walk can be measured locally in any internal cluster, which reflects the quality  
 270 of clustering in this local area. At last, we define the *sparsest level* of  $T$  to be the  $j$ -th level such that  
 271 the average sparsity of triangles rooted at nodes in  $U_j$  is maximum, that is  $\arg \max_j \{\overline{\text{Spar}}_j(T)\}$ ,

272 where  $\overline{\text{Spar}}_j(T) = \sum_{u \in U_j} \text{Spar}(u)/|U_j|$ . Then stratification works for the sparsest level of  $T$ . This  
 273 process is illustrated in Figure 4 (see Appendix G).

274 For a given positive integer  $k$ , to construct a cluster tree of height  $k$ , we start from the trivial 1-level  
 275 cluster tree that involves all vertices of  $G$  as leaves. Then we do not stop stratifying at the sparsest  
 276 level recursively until a  $k$ -level cluster tree is obtained. This process is described in Algorithm 3.

---

**Algorithm 3:**  $k$ -Hierarchical clustering based on structural entropy ( $k$ -HCSE)

---

**Input:** a graph  $G = (V, E)$ ,  $k \in \mathbb{Z}^+$   
**Output:** a  $k$ -level cluster tree  $T$

- 1 Initialize  $T$  to be the 1-level cluster tree;
- 2  $h = \text{height}(T)$ ;
- 3 **while**  $h < k$  **do**
- 4      $j' \leftarrow \arg \max_j \{\overline{\text{Spar}}_j(T)\}$ ;   // Find the sparsest level of  $T$  (breaking ties arbitrarily);
- 5     **if**  $\overline{\text{Spar}}_{j'}(T) = 0$  **then**
- 6         **break**;   // No cost will be saved by any further clustering;
- 7     **for**  $u \in U_{j'}$  **do**
- 8          $T_u \leftarrow \text{Stretch}(u\text{-triangle } T_u)$ ;
- 9          $\text{Compress}(T_u)$ ;
- 10     $h \leftarrow h + 1$ ;
- 11    **for**  $j \in [j' + 1, h]$  **do**
- 12         $\text{Update } U_j$ ;
- 13 **return**  $T$

---

277 To determine the height of the cluster tree automatically, we derive the natural clustering from the  
 278 variation of sparsity on each level. Intuitively, a natural hierarchical cluster tree  $T$  should have  
 279 not only sparse boundary on clusters, but also low sparsity for triangles of  $T$ , which means that  
 280 stratification within the reduced subgraphs corresponding to the triangles on the sparsest level makes  
 281 little sense. For this reason, we consider the inflection points of the sequence  $\{\delta_t(\mathcal{H})\}_{t=1,2,\dots}$ ,  
 282 where  $\delta_t(\mathcal{H})$  is the structural entropy by which the  $t$ -th round of stratification reduces. Formally,  
 283 denote  $\Delta_t \mathcal{H} = \delta_{t-1}(\mathcal{H}) - \delta_t(\mathcal{H})$  for each  $t \geq 2$ . We say that  $\Delta_t \mathcal{H}$  is an inflection point if both  
 284  $\Delta_t \mathcal{H} \geq \Delta_{t-1} \mathcal{H}$  and  $\Delta_t \mathcal{H} \geq \Delta_{t+1} \mathcal{H}$  hold. Our algorithm finds the least  $t$  such that  $\Delta_t \mathcal{H}$  is  
 285 an inflection point and fix the height of the cluster tree to be  $t$  (Note that after  $t - 1$  rounds of  
 286 stratification, the number of levels is  $t$ ). This process is described as Algorithm 4.

---

**Algorithm 4:** Hierarchical clustering based on structural entropy (HCSE)

---

**Input:** a graph  $G = (V, E)$   
**Output:** a cluster tree  $T$

- 1  $t \leftarrow 2$ ;
- 2 **while**  $\Delta_t \mathcal{H} < \Delta_{t-1} \mathcal{H}$  **or**  $\Delta_t \mathcal{H} < \Delta_{t+1} \mathcal{H}$  **do**
- 3     **if**  $\max_j \{\overline{\text{Spar}}_j(T)\} = 0$  **then**
- 4         **break**;
- 5      $t \leftarrow t + 1$ ;
- 6 **return**  $t$ -HCSE( $T$ )

---

287 **Time complexity.** The running time of HCSE on graph  $G = (V, E)$  for which  $|V| = n$  and  $|E| = m$   
 288 depends mainly on the iterations of stratification for the sparsest level. For each round of  $t$ -HCSE in  
 289 Algorithm 4, since the change of structure entropy can be calculated incrementally and locally when  
 290 merge siblings, the time complexity for the Stretch process is  $O(mh \log n)$ , where  $h$  is the height  
 291 of the binary tree that Stretch yields. Since at most  $n$  times of shrinking operations on tree edges  
 292 will happen, the time complexity for the Compress process is  $O(hn)$ . Let  $h_{\max}$  be the maximum  
 293 height among the binary trees that appear during all iterations. The time complexity of HCSE (and  
 294 also  $k$ -HCSE) is  $O(kmh_{\max} \log n + kh_{\max}n)$ . In practice,  $k$  is usually very small (we can even  
 295 set  $k = O(1)$  in  $k$ -HCSE). Moreover, the balance property of structural entropy tends to produce a

	$\vec{p}$	HCSE	HLP	LOU
$p_2$	4.5E(-2)	0.89	0.79	<b>0.92</b>
$p_1$	1.5E(-3)	<b>0.93</b>	0.75	0.92
$p_0$	6E(-6)	<b>0.62</b>	0.58	--
$p_2$	5.5E(-2)	0.87	<b>0.89</b>	0.89
$p_1$	1.5E(-3)	<b>0.95</b>	0.87	0.87
$p_0$	4E(-6)	<b>0.72</b>	--	--
$p_2$	6.5E(-2)	0.96	0.95	<b>0.99</b>
$p_1$	4.5E(-3)	0.94	0.81	<b>0.99</b>
$p_0$	2.5E(-6)	<b>0.80</b>	--	--

Table 1: NMI for three algorithms. Each dataset has 2,500 vertices, and the cluster numbers at three levels are 5, 25 and 250, respectively, for which the size of each cluster is accordingly generated at random.  $p_3 = 0.9$  for each graph. “--” means the algorithm does not find this level.

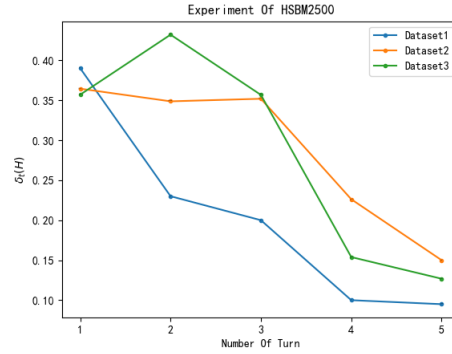


Figure 1:  $\delta_t(\mathcal{H})$  variations for HCSE. It can be observed easily that the inflection points for all the three datasets appear on  $t = 4$ , which is also the ground-truth number of hierarchies.

296 balanced binary tree after stretch, which makes  $h_{\max} = O(\log n)$ . Therefore, in this case, the time  
 297 complexity is merely  $O(m \log^2 n)$ .

## 298 5 Experiments

299 We evaluate experimentally our practically used non-binary clustering algorithm both on synthetic  
 300 networks generated from the Hierarchical Stochastic Block Model (HSBM) and on real datasets,  
 301 respectively. We compare HCSE with the popular practical algorithms LOUVAIN [3] and HLP  
 302 [19]. To avoid over-fitting to higher levels, which possibly results in under-fitting to lower levels,  
 303 LOUVAIN admits a sequential input of vertices. Usually, to avert the worst-case trap, the vertices  
 304 come randomly, and the resulting cluster tree depends on their order. HLP invokes the common  
 305 LP algorithm recursively, and so it cannot be guaranteed to avoid under-fitting in each round. This  
 306 can be seen in our experiments on synthetic datasets, for which these two algorithms usually miss  
 307 ground-truth levels. For real datasets, as far as we know, no public real datasets have clear ground  
 308 truth for hierarchical clustering. We do the comparative experiments on real networks. Some of  
 309 them have (overlapping, possibly hierarchical) ground truth, e.g., Amazon, while others do not  
 310 have. We evaluate the resulting cluster trees for the Amazon network by Jaccard index, and show  
 311 the results in Appendix F, For other networks without ground truth, we evaluate results by both  
 312 cost(SE) and Dasgupta’s cost function cost(Das). All the source code can be downloaded from  
 313 <https://github.com/samwu-learn/HCSE>.

314 **Synthetic datasets generated from HSBM.** Our experiments on synthetic datasets utilize 4-level  
 315 HSBM. For simplicity, let  $\vec{p} = (p_0, p_1, p_2, p_3)$  be the probability vector for which  $p_i$  is the proba-  
 316 bility of generating edges for vertex pairs whose LCA on the ground-truth cluster tree has depth  $i$ .  
 317 Note that the 0-depth node is the root. We compare the Normalized Mutual Information (NMI) at  
 318 each level of the ground-truth cluster tree to those of three algorithms. Note that the randomness in  
 319 LOUVAIN, and breaking-ties rule as well as convergence of HLP make different results, we choose  
 320 the most effective strategy and pick the best results in five runs for both of them. Compared to their  
 321 uncertainty, our algorithm HCSE yields stable results.

322 Table 1 demonstrates the results in three groups of probabilities, for which the hierarchical structures  
 323 get clearer one by one. Each dataset has 2,500 vertices, and the cluster numbers at three levels are  
 324 5, 25 and 250, respectively, for which the size of each cluster is accordingly generated at random.  
 325  $p_3 = 0.9$  for each graph. Our algorithm HCSE is always able to find the right number of levels,  
 326 while LOUVAIN always misses the top level, and HLP misses the top level in two groups. The  
 327 inflection points for choosing the intrinsic hierarchy number  $t = 4$  of hierarchies are demonstrated  
 328 in Figure 1.



Networks	HCSE	HLP	LOUVAIN
CSphd	1.30E4 / <b>5.19E4</b> / 5	1.54E4 / 5.58E4 / 4	<b>1.28E4</b> / 7.61E4 / 5
fb-pages-government	2.48E6 / <b>1.18E8</b> / 4	2.53E6 / 1.76E8 / 3	<b>2.43E6</b> / 1.33E8 / 4
email-univ	1.16E5 / <b>2.20E6</b> / 3	1.46E5 / 6.14E6 / 3	<b>1.14E5</b> / 2.20E6 / 4
fb-messages	1.58E5 / <b>4.50E6</b> / 4	1.76E5 / 8.12E6 / 3	<b>1.52E5</b> / 4.96E6 / 4
G22	<b>5.56E5</b> / <b>2.68E7</b> / 4	6.11E5 / 4.00E7 / 3	5.63E5 / 2.80E7 / 5
As20000102	2.64E5 / <b>2.36E7</b> / 4	3.62E5 / 7.63E7 / 3	<b>2.42E5</b> / 2.42E7 / 5
bibd-13-6	<b>7.41E5</b> / <b>2.56E7</b> / 3	8.05E5 / 4.41E7 / 2	7.50E5 / 2.75E7 / 4
delaunay-n10	4.65E4 / <b>3.39E5</b> / 4	4.87E4 / 3.55E5 / 4	<b>4.24E4</b> / 4.25E5 / 5
p2p-Gnutella05	9.00E5 / <b>1.48E8</b> / 3	1.01E6 / 2.78E8 / 3	<b>8.05E5</b> / 1.49E8 / 5
p2p-Gnutella08	5.59E5 / <b>5.51E7</b> / 4	6.36E5 / 1.28E8 / 4	<b>4.88E5</b> / 6.03E7 / 5

Table 2: “cost(SE) / cost(Das) /  $k$ ” for three algorithms, where  $k$  is the hierarchy number that the algorithm finds.

329 **Real datasets.** We do our experiments on a series of real networks <sup>2</sup> without ground truth. We  
330 compare cost(SE) and cost(Das), respectively. Since the different level numbers given by the three  
331 algorithms influence the costs seriously, that is, lower costs are obtained just due to greater heights,  
332 we only list in Table 2 the networks for which the three algorithms yield similar level numbers that  
333 differ by at most 1 or 2. It can be observed that HLP does not achieve optima for any network,  
334 while HCSE performs best w.r.t. cost(Das) for all networks, but does not outperform LOUVAIN  
335 for most networks. This is mainly due to the fact that LOUVAIN always finds no less number  
336 of hierarchies than HCSE, and the better cost benefits from its depth. Moreover, we emphasize that  
337 there is no evidence to indicate that the lower cost(SE) or cost(Das) is, the better a non-binary cluster  
338 tree is. Our experiments on these real datasets are just demonstrations of the effectiveness for our  
339 interpretable mechanism in hierarchical clustering.

## 340 6 Conclusions and future discussions

341 In this paper, we investigate the hierarchical clustering problem from an information-theoretic per-  
342 spective and propose a new objective function that relates to the combinatorial cost functions raised  
343 by Dasgupta [10]. For optimization of this function, we present two  $O(1)$ -approximation algorithms  
344 for expander-like and well-clustered cardinality weighted graphs, respectively. For practical use, we  
345 propose a new interpretable non-binary hierarchical clustering framework that stratifies the sparsest  
346 level of the cluster tree recursively, which can be collocated with any cost function. We also present  
347 an interpretable strategy to find the intrinsic number of levels without any hyper-parameter. The ex-  
348 perimental results on  $k$ -level HSBM demonstrate that our algorithm HCSE has a great advantage in  
349 finding  $k$  compared to the popular but strongly heuristic algorithms LOUVAIN and HLP. Our results  
350 on real datasets show that HCSE also achieves competitive costs compared to these two algorithms.

351 There are several directions that are worth further study. The first problem is about the relationship  
352 between the concavity of  $g$  of the cost function and the balance of the optimal cluster tree. It can be  
353 checked that for cliques, being concave is not a sufficient condition for total balance. Whether is it a  
354 necessary condition? Moreover, is there any explicit necessary and sufficient condition for total bal-  
355 ance of the optimal cluster tree for cliques? The second problem is about approximation algorithms  
356 for both structural entropy and cost(SE) in the worst case. Due to the non-linear and volume-related  
357 function  $g$ , many previous proof techniques for approximation algorithms seems unavailable. The  
358 third one is about more precise characterizations for “natural” hierarchical clustering whose depth is  
359 limited. Since any reasonable choice of  $g$  makes the cost function achieve optimum on some binary  
360 tree, a blind pursuit of minimization of cost functions seems not to be a rational approach. More  
361 criteria in this scenario need to be studied.

<sup>2</sup><http://networkrepository.com/index.php>

## 362 References

- 363 [1] Noga Alon, Yossi Azar, and Danny Vainstein. Hierarchical clustering: A 0.585 revenue ap-  
364 proximation. In Jacob D. Abernethy and Shivani Agarwal, editors, *Conference on Learning*  
365 *Theory, COLT 2020, 9-12 July 2020, Virtual Event [Graz, Austria]*, volume 125 of *Proceedings*  
366 *of Machine Learning Research*, pages 153–162. PMLR, 2020.
- 367 [2] Sanjeev Arora, Satish Rao, and Umesh V. Vazirani. Expander flows, geometric embeddings  
368 and graph partitioning. *J. ACM*, 56(2):5:1–5:37, 2009.
- 369 [3] Vincent D. Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast  
370 unfolding of communities in large networks. *Journal of statistical mechanics: theory and*  
371 *experiment*, 2008(10):P10008, 2008.
- 372 [4] Peter F Brown, Vincent J Della Pietra, Peter V Desouza, Jennifer C Lai, and Robert L Mercer.  
373 Class-based n-gram models of natural language. *Computational linguistics*, 18(4):467–480,  
374 1992.
- 375 [5] Moses Charikar and Vaggos Chatziafratis. Approximate hierarchical clustering via sparsest  
376 cut and spreading metrics. In Philip N. Klein, editor, *Proceedings of the Twenty-Eighth Annual*  
377 *ACM-SIAM Symposium on Discrete Algorithms, SODA 2017, Barcelona, Spain, Hotel Porta*  
378 *Fira, January 16-19*, pages 841–854. SIAM, 2017.
- 379 [6] Moses Charikar, Vaggos Chatziafratis, and Rad Niazadeh. Hierarchical clustering better than  
380 average-linkage. In Timothy M. Chan, editor, *Proceedings of the Thirtieth Annual ACM-SIAM*  
381 *Symposium on Discrete Algorithms, SODA 2019, San Diego, California, USA, January 6-9,*  
382 *2019*, pages 2291–2304. SIAM, 2019.
- 383 [7] Vaggos Chatziafratis, Grigory Yaroslavtsev, Euiwoong Lee, Konstantin Makarychev, Sara Ah-  
384 madian, Alessandro Epasto, and Mohammad Mahdian. Bisect and conquer: Hierarchical clus-  
385 tering via max-uncut bisection. In Silvia Chiappa and Roberto Calandra, editors, *The 23rd*  
386 *International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 Au-*  
387 *gust 2020, Online [Palermo, Sicily, Italy]*, volume 108 of *Proceedings of Machine Learning*  
388 *Research*, pages 3121–3132. PMLR, 2020.
- 389 [8] Vincent Cohen-Addad, Varun Kanade, Frederik Mallmann-Trenn, and Claire Mathieu. Hierar-  
390 chical clustering: Objective functions and algorithms. *Journal of the ACM (JACM)*, 66(4):1–42,  
391 2019.
- 392 [9] Aron Culotta, Pallika Kanani, Robert Hall, Michael Wick, and Andrew McCallum. Author  
393 disambiguation using error-driven machine learning with a ranking loss function. In *Sixth*  
394 *International Workshop on Information Integration on the Web (IIWeb-07), Vancouver, Canada,*  
395 *2007*.
- 396 [10] Sanjoy Dasgupta. A cost function for similarity-based hierarchical clustering. In Daniel Wicks  
397 and Yishay Mansour, editors, *Proceedings of the 48th Annual ACM SIGACT Symposium on*  
398 *Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, pages 118–127.  
399 ACM, 2016.
- 400 [11] Michael B Eisen, Paul T Spellman, Patrick O Brown, and David Botstein. Cluster analysis  
401 and display of genome-wide expression patterns. *Proceedings of the National Academy of*  
402 *Sciences*, 95(25):14863–14868, 1998.
- 403 [12] Shayan Oveis Gharan and Luca Trevisan. Partitioning into expanders. In Chandra Chekuri,  
404 editor, *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms,*  
405 *SODA 2014, Portland, Oregon, USA, January 5-7, 2014*, pages 1256–1266. SIAM, 2014.
- 406 [13] Alexander N Gorban, Balázs Kégl, Donald C Wunsch, Andrei Y Zinovyev, et al. *Principal*  
407 *manifolds for data visualization and dimension reduction*, volume 58. Springer, 2008.
- 408 [14] Angsheng Li and Yicheng Pan. Structural information and dynamical complexity of networks.  
409 *IEEE Trans. Inf. Theory*, 62(6):3290–3339, 2016.

- 410 [15] Bogdan-Adrian Manghiuc and He Sun. Hierarchical clustering:  $O(1)$ -approximation for well-  
411 clustered graphs. In Marc’ Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang,  
412 and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*  
413 *34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, De-*  
414 *cember 6-14, 2021, virtual*, pages 9278–9289, 2021.
- 415 [16] Benjamin Moseley and Joshua R. Wang. Approximation bounds for hierarchical clustering:  
416 Average linkage, bisecting k-means, and local search. In Isabelle Guyon, Ulrike von Luxburg,  
417 Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett,  
418 editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*  
419 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages  
420 3094–3103, 2017.
- 421 [17] Stanislav Naumov, Grigory Yaroslavtsev, and Dmitrii Avdiukhin. Objective-based hierarchical  
422 clustering of deep embedding vectors. In *Thirty-Fifth AAAI Conference on Artificial Intel-*  
423 *ligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelli-*  
424 *gence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence,*  
425 *EAAI 2021, Virtual Event, February 2-9, 2021*, pages 9055–9063. AAAI Press, 2021.
- 426 [18] Mirmahdi Rahgoshay and Mohammad R. Salavatipour. Hierarchical clustering: New bounds  
427 and objective. *CoRR*, abs/2111.06863, 2021.
- 428 [19] Ryan A Rossi, Nesreen K Ahmed, Eunyee Koh, and Sungchul Kim. Fast hierarchical graph  
429 clustering in linear-time. In *Companion Proceedings of the Web Conference 2020*, pages 10–  
430 12, 2020.
- 431 [20] Aurko Roy and Sebastian Pokutta. Hierarchical clustering via spreading metrics. *J. Mach.*  
432 *Learn. Res.*, 18:88:1–88:35, 2017.

## 433 Checklist

- 434 1. For all authors...
- 435 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
436 contributions and scope? [\[Yes\]](#) See Abstract and Section 1
- 437 (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 3 and 6
- 438 (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
- 439 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
440 them? [\[Yes\]](#)
- 441 2. If you are including theoretical results...
- 442 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 3
- 443 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Appendices
- 444 3. If you ran experiments...
- 445 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
446 imental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See Section  
447 5
- 448 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
449 were chosen)? [\[N/A\]](#)
- 450 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
451 ments multiple times)? [\[N/A\]](#)
- 452 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
453 of GPUs, internal cluster, or cloud provider)? [\[No\]](#) We do not compare the speed of  
454 computing, and so due to the space limit we have omitted resource introduction.
- 455 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 456 (a) If your work uses existing assets, did you cite the creators? [\[N/A\]](#) All codes are written  
457 by the authors, and all datasets we use is public. We have provided the URLs that link  
458 to the datasets we use.

- 459 (b) Did you mention the license of the assets? [N/A]  
 460 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
 461  
 462 (d) Did you discuss whether and how consent was obtained from people whose data  
 463 you're using/curating? [No] All datasets we used are public.  
 464 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 465 information or offensive content? [No] No such content is included.  
 466 5. If you used crowdsourcing or conducted research with human subjects...
- 467 (a) Did you include the full text of instructions given to participants and screenshots, if  
 468 applicable? [N/A]  
 469 (b) Did you describe any potential participant risks, with links to Institutional Review  
 470 Board (IRB) approvals, if applicable? [N/A]  
 471 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 472 spent on participant compensation? [N/A]

## 473 A A brief introduction to structural information

The idea of structural information is to encode a random walk with a certain rule by using a high-dimensional encoding system for a graph  $G$ . It is well known that a random walk, for which a neighbor is randomly chosen with probability proportional to edge weights, has a stationary distribution on vertices that is proportional to vertex degree.<sup>3</sup> So to position a random walk under its stationary distribution, the amount of information needed is typically the Shannon's entropy, denoted by

$$\mathcal{H}^{(1)}(G) = - \sum_{v \in V} \frac{d_v}{\text{vol}(V)} \log \frac{d_v}{\text{vol}(V)}.$$

474 By Shannon's noiseless coding theorem,  $\mathcal{H}^{(1)}(G)$  is the limit of average code length generated from  
 475 the *memoryless* source for one step of the random walk. However, dependence of locations may  
 476 shorten the code length. For each level on cluster trees, the uncertainty of locations is measured by  
 477 the entropy of the stationary distribution on the clusters of this level. Consider an encoding for every  
 478 cluster, including the leaves. Each non-root node  $\alpha$  is labeled by its order among the children of  
 479 its parent  $\alpha^-$ . So the amount of self-information of  $\alpha$  within this local parent-children substructure  
 480 is  $-\log(\text{vol}(\alpha)/\text{vol}(\alpha^-))$ , which is also roughly the length of Shannon code for  $\alpha$  and its siblings.  
 481 The codeword of  $\alpha$  consists of the sequential labels of nodes along the unique path from the root  
 482 (excluded) to itself (included). The key idea is as follows. For one step of the random walk from  $u$  to  
 483  $v$  in  $G$ , to indicate  $v$ , we omit from  $v$ 's codeword the longest common prefix of  $u$  and  $v$  that is exactly  
 484 the codeword of  $u \vee v$ . This means that the random walk takes this step in the cluster  $u \vee v$  (and  
 485 also in  $u \vee v$ 's ancestors) and the uncertainty at this level may not be involved. Therefore, intuitively,  
 486 a quality similarity-based cluster tree would trap the random walk with high frequency in the deep  
 487 clusters that are far from the root, and long codeword of  $u \vee v$  would be omitted. This shortens the  
 488 average code length of the random walk. Note that we ignore the uniqueness of decoding since a  
 489 practical design of codewords is not our purpose. We utilize this scheme to evaluate and differentiate  
 490 hierarchical structures.

491 Then we formulate the above scheme and measure the average code length as follows. Given a  
 492 weighted graph  $G = (V, E, w)$  and a cluster tree  $T$  for  $G$ , note that under the stationary distribution,  
 493 the random walk takes one step out of a cluster  $\alpha$  on  $T$  with probability  $g_\alpha/\text{vol}(V)$ . Therefore, the  
 494 aforementioned uncertainty measured by the average code length is

$$\mathcal{H}^T(G) = - \sum_{\alpha \in T} \frac{g_\alpha}{\text{vol}(V)} \log \frac{\text{vol}(\alpha)}{\text{vol}(\alpha^-)},$$

495 which is defined as the structural entropy of  $G$  on  $T$ . To minimize this uncertainty, the structural  
 496 entropy  $\mathcal{H}(G)$  of  $G$  is defined as the minimum one among all cluster trees. Note that the structural  
 497 entropy of  $G$  on the trivial 1-level cluster tree is consistent with the previously defined  $\mathcal{H}^{(1)}(G)$ . It  
 498 doesn't have any non-trivial cluster.

<sup>3</sup>For connected graphs, this stationary distribution is unique, but not for disconnected ones. Here, we consider this canonical one for all graphs.

499 **B Proof of Proposition 2.1**

500 *Proof.* For each internal node  $\alpha$  on  $T$ , denote by  $\partial(\alpha)$  the sets of edges in  $G$  with exactly one  
 501 end-point in the set of vertices corresponding to  $\alpha$ . So  $g_\alpha = \sum_{e \in \partial(\alpha)} w(e)$ . Note that

$$\begin{aligned} \mathcal{H}^T(G) &= - \sum_{\alpha \in T} \frac{g_\alpha}{\text{vol}(V)} \log \frac{\text{vol}(\alpha)}{\text{vol}(\alpha^-)} \\ &= - \sum_{\alpha \in T} \sum_{(u,v) \in \partial(\alpha)} \frac{w(u,v)}{\text{vol}(V)} \log \frac{\text{vol}(\alpha)}{\text{vol}(\alpha^-)} \\ &= - \sum_{(u,v) \in E} \left( \frac{w(u,v)}{\text{vol}(V)} \sum_{\alpha: (u,v) \in g_\alpha} \log \frac{\text{vol}(\alpha)}{\text{vol}(\alpha^-)} \right). \end{aligned}$$

502 For a single edge  $(u, v) \in E$ , all the terms  $\log(\text{vol}(\alpha)/\text{vol}(\alpha^-))$  for leaf  $u$  satisfying  $(u, v) \in g_\alpha$   
 503 sum (over  $\alpha$ ) up to  $\log(d_u/\text{vol}(u \vee v))$  along the unique path from  $u$  to  $u \vee v$ . It is symmetric for  $v$ .  
 504 Therefore, considering ordered pair  $(u, v) \in E$ ,

$$\begin{aligned} \mathcal{H}^T(G) &= - \sum_{\text{ordered } (u,v) \in E} \frac{w(u,v)}{\text{vol}(V)} \log \frac{d_u}{\text{vol}(u \vee v)} \\ &= \frac{1}{\text{vol}(V)} \left( - \sum_{u \in V} d_u \log d_u + \sum_{\text{ordered } (u,v) \in E} w(u,v) \log \text{vol}(u \vee v) \right) \\ &= \frac{1}{\text{vol}(V)} \left( - \sum_{u \in V} d_u \log d_u + 2 \cdot \sum_{(u,v) \in E} w(u,v) \log \text{vol}(u \vee v) \right). \end{aligned}$$

505 The second equality follows from the fact  $\sum_{u \in V} d_u = \sum_{\text{ordered } (u,v) \in E} w(u,v) = \text{vol}(V)$  and the  
 506 last equality from the symmetry of  $(u, v)$ . Since the first summation is independent of  $T$ , Proposition  
 507 2.1 follows.  $\square$

508 **C Proof of Proposition 2.2**

509 We restate Proposition 2.2 as follows.

510 **Theorem 2.2.** *For any positive integer  $n$ , let  $K_n$  be the clique of  $n$  vertices with identical weight on*  
 511 *every edge. Then a cluster tree  $T$  of  $K_n$  achieves minimum structural entropy if and only if  $T$  is a*  
 512 *balanced binary tree, that is, the two children clusters of each sub-tree of  $T$  have difference in size*  
 513 *at most 1.*

514 Note that a balanced binary tree (BBT for abbreviation) means the tree is balanced on every internal  
 515 node. Formally, for an internal node of cluster size  $k$ , its two sub-trees are of cluster sizes  $\lfloor k/2 \rfloor$  and  
 516  $\lceil k/2 \rceil$ , respectively.

517 For cliques, since the weights of each edge are identical, we assume it safely to be 1. By Theorem  
 518 2.1, minimizing the structural entropy is equivalent to minimizing the cost function (over  $T$ )

$$\begin{aligned} \text{cost}^T(G) &= \sum_{(u,v) \in E} \log \text{vol}(u \vee v) \\ &= \sum_{(u,v) \in E} \log((n-1)|u \vee v|) \\ &= \sum_{(u,v) \in E} \log(n-1) + \sum_{(u,v) \in E} \log |u \vee v| \end{aligned}$$

519 Since the first term in the last equation is independent of  $T$ , the optimization turns to minimizing  
 520 the last term, which we denote by  $\Gamma(T)$ . Grouping all edges in  $E$  by LCA of two end-points, the  
 521 cost  $\Gamma(T)$  can be written as the sum of the cost  $\gamma$  at every internal node  $N$  of  $T$ . Formally, for every

522 internal node  $N$ , let  $A, B \subseteq V$  be the leaves of the sub-trees rooted at the left and right child of  $N$ ,  
 523 respectively. We have

$$\begin{aligned}\Gamma(T) &= \sum_N \gamma(N) \\ \gamma(N) &= \left( \sum_{x \in A, y \in B} 1 \right) \cdot \log(|A| + |B|) \\ &= |A| \cdot |B| \cdot \log(|A| + |B|)\end{aligned}$$

524 Now we only have to show the following lemma.

525 **Lemma C.1.** *For any positive integer  $n$ , a cluster tree  $T$  of  $K_n$  achieves minimum cost  $\Gamma(T)$  if and*  
 526 *only if  $T$  is a BBT.*

527 *Proof.* Lemma C.1 is proved by induction on  $|V|$ . The key technique of tree swapping we use here  
 528 is inspired by Cohen-Addad et al [4]. The basis step holds since for  $|V| = 2$  or  $3$ , the cluster tree is  
 529 balanced and unique. It certainly achieves the minimum cost exclusively.

530 Now, consider a clique  $G = (V, E)$  with  $n = |V| \geq 4$ . Let  $T_1$  be an arbitrary unbalanced cluster  
 531 tree and  $\lambda$  be its root. We need to prove that the cost  $\Gamma(T_1)$  does not achieve the minimum. Without  
 532 loss of generality, we can safely assume the root node is unbalanced, since otherwise, we set  $T_1$  to  
 533 be the sub-tree that is rooted at an unbalanced node. Let  $T_2$  be a tree with root  $\lambda$  whose left and  
 534 right sub-trees are BBTs such that they have the same sizes with the left and right sub-trees of  $T_1$ ,  
 535 respectively. Let  $V_{ll}, V_{lr}, V_{rl}$  and  $V_{rr}$  be the sets of nodes on the four sub-trees at the second level of  
 536  $T_2$  and  $n_{ll}, n_{lr}, n_{rl}$  and  $n_{rr}$  denote their sizes, respectively. Our proof is also available when some of  
 537 them are empty. We always assume  $n_{ll} \leq n_{lr}$  and  $n_{rl} \geq n_{rr}$ . Next, we construct  $T_3$  by swapping  
 538 (transplanting)  $V_{lr}$  and  $V_{rl}$  with each other. Finally, let  $T_4$  be a tree with root  $\lambda$  whose left and right  
 539 sub-trees are BBTs after balancing the left and right sub-trees of  $T_3$ . So  $T_4$  is a BBT. Then we only  
 540 have to prove that  $\Gamma(T_1) > \Gamma(T_4)$ . Note that the strict “ $>$ ” is necessary since we need to negate all  
 541 unbalanced cluster trees.

542 Then we show that the transformation process that consists of the above three steps makes the cost  
 543 decrease step by step. Formally,

- 544 (a)  $T_1$  to  $T_2$ . The sub-trees of  $T_1$  become BBTs in  $T_2$ . Since the number of edges whose  
 545 end-points treat the root as LCA is the same, by induction we have  $\Gamma(T_1) \geq \Gamma(T_2)$ .
- 546 (b)  $T_2$  to  $T_3$ . We will show that  $\Gamma(T_2) > \Gamma(T_3)$  in Lemma C.2.
- 547 (c)  $T_3$  to  $T_4$ . The sub-trees of  $T_3$  become BBTs in  $T_4$ . For the same reason as (a), we have  
 548  $\Gamma(T_3) \geq \Gamma(T_4)$ .

549 Putting them together, we get  $\Gamma(T_1) > \Gamma(T_4)$  and Lemma C.1 follows.

550

□

551 **Lemma C.2.** *After swapping  $V_{lr}$  and  $V_{rl}$ , we obtain  $T_3$  from  $T_2$ , for which  $\Gamma(T_2) > \Gamma(T_3)$ .*

552 *Proof.* We only need to consider the changes in cost of three nodes: root and its left and right  
 553 children, since the cost contributed by each of the remaining nodes does not change after swapping.  
 554 Ignoring the unchanged costs, define

$$\begin{aligned}\text{cost}(T_2) &= n_l n_r \log n + n_{ll} n_{lr} \log n_l + n_{rl} n_{rr} \log n_r \\ &= n_l n_r \log n + \left\lfloor \frac{n_l}{2} \right\rfloor \left\lceil \frac{n_l}{2} \right\rceil \log n_l + \left\lceil \frac{n_r}{2} \right\rceil \left\lfloor \frac{n_r}{2} \right\rfloor \log n_r,\end{aligned}$$

555 where  $n_l = n_{ll} + n_{lr}$ ,  $n_r = n_{rl} + n_{rr}$ . Both of them are at least 1. Similarly, define

$$\begin{aligned}\text{cost}(T_3) &= (n_{ll} + n_{rl})(n_{lr} + n_{rr}) \log n + n_{ll} n_{rl} \log (n_{ll} + n_{rl}) + n_{lr} n_{rr} \log (n_{lr} + n_{rr}) \\ &= \left\lfloor \frac{n}{2} \right\rfloor \left\lceil \frac{n}{2} \right\rceil \log n + \left\lfloor \frac{n_l}{2} \right\rfloor \left\lceil \frac{n_r}{2} \right\rceil \log \left( \left\lfloor \frac{n_l}{2} \right\rfloor + \left\lceil \frac{n_r}{2} \right\rceil \right) + \left\lceil \frac{n_l}{2} \right\rceil \left\lfloor \frac{n_r}{2} \right\rfloor \log \left( \left\lceil \frac{n_l}{2} \right\rceil + \left\lfloor \frac{n_r}{2} \right\rfloor \right)\end{aligned}$$

556 Denote

$$\begin{aligned}
\Delta &= \Gamma(T_2) - \Gamma(T_3) \\
&= \text{cost}(T_2) - \text{cost}(T_3) \\
&= \left\lfloor \frac{n_l}{2} \right\rfloor \left\lceil \frac{n_l}{2} \right\rceil \log \left( \frac{n_l}{n} \right) + \left\lceil \frac{n_r}{2} \right\rceil \left\lfloor \frac{n_r}{2} \right\rfloor \log \left( \frac{n_r}{n} \right) \\
&\quad - \left\lfloor \frac{n_l}{2} \right\rfloor \left\lceil \frac{n_r}{2} \right\rceil \log \left( \frac{\left\lfloor \frac{n_l}{2} \right\rfloor + \left\lceil \frac{n_r}{2} \right\rceil}{n} \right) - \left\lceil \frac{n_l}{2} \right\rceil \left\lfloor \frac{n_r}{2} \right\rfloor \log \left( \frac{\left\lceil \frac{n_l}{2} \right\rceil + \left\lfloor \frac{n_r}{2} \right\rfloor}{n} \right) \quad (2)
\end{aligned}$$

557 So we only have to show that  $\Delta > 0$ . We consider the following three cases according to the odevity  
558 of  $n_l$  and  $n_r$ .

559 **Case 1:**  $n_l$  and  $n_r$  are even.

560 **Case 2:**  $n_l$  and  $n_r$  are odd.

561 **Case 3:**  $n_l$  is odd while  $n_r$  is even.

562 The case that  $n_l$  is even while  $n_r$  is odd is symmetric to **Case 3**.

563 For **Case 1**, if both  $n_l$  and  $n_r$  are even, then notations of rounding in Eq. (2) can be removed and  $\Delta$   
564 can be simplified as

$$\Delta = \frac{n_l^2}{4} \log \left( \frac{n_l}{n} \right) + \frac{n_r^2}{4} \log \left( \frac{n_r}{n} \right) + \frac{n_l n_r}{2}.$$

565 Let  $p = n_l/n, q = n_r/n$ , and so  $p + q = 1$ . Recall that  $T_1$  is unbalanced on the root  $\lambda$ , so is  $T_2$ .  
566 Thus  $p \neq q$ . Multiplying by  $\frac{4}{n^2}$  on both sides, we only have to prove that

$$p^2 \log p + q^2 \log q + 2pq > 0.$$

That is,

$$\frac{p}{q} \log p + \frac{q}{p} \log q + 2 > 0.$$

567 Let  $g(x) = \frac{x}{1-x} \log x$ . Then we only need to show that  $g(p) + g(q) + 2 > 0$  when  $p \neq q$ . Since

$$\begin{aligned}
g'(x) &= \frac{(1-x) + \ln x}{\ln 2 \cdot (1-x)^2}, \\
g''(x) &= -\frac{x^2 - 2x \ln x - 1}{\ln 2 \cdot x(1-x)^3}.
\end{aligned}$$

It is easy to check that  $g''(x) > 0$  when  $0 < x < 1$ . So  $g(x)$  is strictly convex in the interval  $(0, 1)$ .  
Since  $p \neq q$ ,

$$g(p) + g(q) > 2g\left(\frac{p+q}{2}\right) = -2.$$

568 Thus  $\Delta > 0$  holds.

569 For **Case 2**, if both  $n_l$  and  $n_r$  are odd, then  $\Delta$  can be split into two parts  $\Delta = \Delta_1 + \Delta_2$ , in which

$$\begin{aligned}
\Delta_1 &= \frac{n_l^2}{4} \log \left( \frac{n_l}{n} \right) + \frac{n_r^2}{4} \log \left( \frac{n_r}{n} \right) + \frac{n_l n_r}{2} \\
\Delta_2 &= -\frac{1}{4} \log \left( \frac{n_l}{n} \right) - \frac{1}{4} \log \left( \frac{n_r}{n} \right) - \frac{1}{2}
\end{aligned}$$

570 Since we have shown that  $\Delta_1 > 0$ , if we can prove  $\Delta_2 \geq 0$ , then the lemma will hold for **Case 2**.

571 Due to the convexity of logarithmic function, this holds clearly since

$$2 \log \left( \frac{n}{2} \right) \geq \log n_l + \log n_r.$$

572 For **Case 3**, if  $n_l$  is odd while  $n_r$  is even,

$$\Delta = \frac{n_l^2 - 1}{4} \log \left( \frac{n_l}{n} \right) + \frac{n_r^2}{4} \log \left( \frac{n_r}{n} \right) - \left[ \frac{(n_l - 1)n_r}{4} \log \left( \frac{n - 1}{2n} \right) + \frac{(n_l + 1)n_r}{4} \log \left( \frac{n + 1}{2n} \right) \right].$$

573 Multiplying the above equation by  $4 \ln 2$ , without changing its sign, yields

$$(4 \ln 2)\Delta = (n_l^2 - 1) \ln \left( \frac{n_l}{n} \right) + n_r^2 \ln \left( \frac{n_r}{n} \right) - \left[ (n_l - 1)n_r \ln \left( \frac{n-1}{2n} \right) + (n_l + 1)n_r \ln \left( \frac{n+1}{2n} \right) \right]$$

574 Splitting the right hand side into two parts,

$$\begin{aligned} A &= n_l^2 \ln \left( \frac{n_l}{n} \right) + n_r^2 \ln \left( \frac{n_r}{n} \right) + 2n_l n_r \ln 2 \\ B &= -\ln \left( \frac{n_l}{n} \right) - (n_l + 1)n_r \ln \left( 1 + \frac{1}{n} \right) - (n_l - 1)n_r \ln \left( 1 - \frac{1}{n} \right) \end{aligned}$$

575 Since  $n$  is odd and the root  $\lambda$  of  $T_2$  is unbalanced, we only need to consider the case that  $n_l =$   
576  $(n - i)/2$ ,  $n_r = (n + i)/2$  (Note that  $n_l$  and  $n_r$  are symmetric. So if  $(n - i)/2$  is even, exchange  
577  $n_l$  and  $n_r$ ), where both  $n$  and  $i$  are odd satisfying  $n > i \geq 3$ . Next we show that in this case,  
578  $A \geq \ln(1/5) + 4^2 \ln(4/5) + 2 \cdot 4 \ln 2$  and  $B > \ln 2 - 3/4 - (2/3) \cdot (1/5^2)$ . By calculation,  
579  $\Delta = A + B > 0$  for **Case 3**.

580 **Claim C.1.**  $A \geq \ln(1/5) + 4^2 \ln(4/5) + 2 \cdot 4 \ln 2$  for odd integers  $n > i \geq 3$ .

581 *Proof.* Substituting  $n_l = (n - i)/2$ ,  $n_r = (n + i)/2$  into the  $A$  yields

$$A = C(n, i) \triangleq \left( \frac{n-i}{2} \right)^2 \ln \left( \frac{n-i}{2n} \right) + \left( \frac{n+i}{2} \right)^2 \ln \left( \frac{n+i}{2n} \right) + 2 \cdot \frac{n-i}{2} \cdot \frac{n+i}{2} \ln 2.$$

582 Treat  $n$  as a continuous variable, we have

$$\frac{\partial C(n, i)}{\partial n} = \frac{1}{2} \left[ (n+i) \ln \left( 1 + \frac{i}{n} \right) + (n-i) \ln \left( 1 - \frac{i}{n} \right) - \frac{i^2}{n} \right]$$

583 Multiplying the above equation by  $2/n$  and setting  $x = i/n$  yields

$$\begin{aligned} f(x) &\triangleq (1+x) \ln(1+x) + (1-x) \ln(1-x) - x^2, \\ f'(x) &= \ln(1+x) - \ln(1-x) - 2x, \\ f''(x) &= \frac{2x^2}{1-x^2}. \end{aligned}$$

584 It is easy to check that  $f(0) = 0$  and  $f'(0) = 0$ . When  $0 < x < 1$ ,  $f''(x) > 0$ . Thus  $f'(x) > 0$  and  
585  $f(x) > 0$ . This means that  $\partial C(n, i)/\partial n > 0$  for all  $n > 0$ . So  $C(n, i) \geq C(i+2, i)$  for  $n \geq i+2$   
586 (When  $i$  is fixed, the minimum value of  $n$  can be taken to  $i+2$ , which makes  $n_l = (n - i)/2$  and  
587  $n_r = (n + i)/2$  integral). The curves of  $C(n, i)$  for varying  $i$  are plotted in Figure 2.

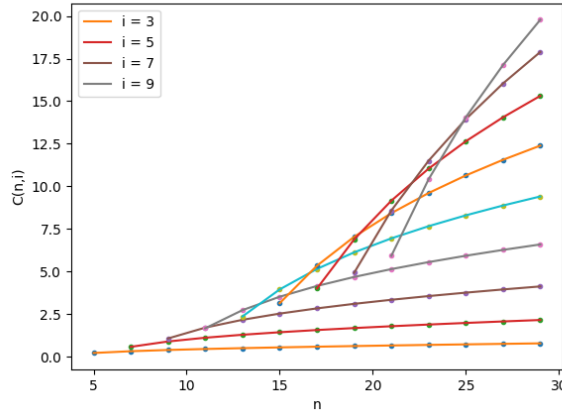


Figure 2: Functions  $C(n, i)$



588 When  $n = i + 2$ , we get  $n_l = (n - i)/2 = 1$  and  $n_r = (n + i)/2 = n - 1$ . Substituting them into  
 589  $A$  yields

$$\begin{aligned} D(n) &\triangleq \ln\left(\frac{1}{n}\right) + (n-1)^2 \ln\left(1 - \frac{1}{n}\right) + 2(n-1) \ln 2, \\ \frac{dD}{dn} &= 1 - \frac{2}{n} + 2 \ln 2 + 2(n-1) \ln\left(1 - \frac{1}{n}\right). \end{aligned}$$

590 When  $n > 2$ , it is easy to check that  $dD/dn > 0$ . So the minimum value of  $d(n)$ , which is also the  
 591 minimum value of  $C(i + 2, i)$ , is achieved at  $n = i + 2 = 5$ . So  $A = C(n, i) \geq C(i + 2, i) \geq$   
 592  $C(5, 3) = \ln(1/5) + 4^2 \ln(4/5) + 2 \cdot 4 \ln 2$ .  $\square$

593 **Claim C.2.**  $B > \ln 2 - 3/4 - (2/3) \cdot (1/5^2)$ .

594 *Proof.* Due to the facts that

$$\begin{aligned} \ln\left(1 + \frac{1}{n}\right) &< \frac{1}{n} - \frac{1}{2n^2} + \frac{1}{3n^3}, \\ \ln\left(1 - \frac{1}{n}\right) &< -\frac{1}{n} - \frac{1}{2n^2} - \frac{1}{3n^3}, \end{aligned}$$

595 we have

$$\begin{aligned} B &= -\ln\left(\frac{n_l}{n}\right) - (n_l + 1)n_r \ln\left(1 + \frac{1}{n}\right) - (n_l - 1)n_r \ln\left(1 - \frac{1}{n}\right) \\ &> -\ln\left(\frac{n_l}{n}\right) + \frac{n_l n_r}{n^2} - \frac{2n_r}{n} - \frac{2n_r}{3n^3} \\ &> -\ln\left(\frac{n_l}{n}\right) + \frac{n_l n_r}{n^2} - \frac{2n_r}{n} - \frac{2}{3n^2}. \end{aligned}$$

596 Let  $\alpha = n_l/n$ , then

$$\begin{aligned} B &> -\ln \alpha + \alpha(1 - \alpha) - 2(1 - \alpha) - \frac{2}{3n^2} \\ &\geq \ln 2 - \frac{3}{4} - \frac{2}{3n^2}. \end{aligned}$$

597 When  $n \geq 5$ ,  $B > \ln 2 - 3/4 - (2/3) \cdot (1/5^2)$ .  $\square$

598 Combining Claims C.1 and C.2, Lemma C.2 follows.  $\square$

599 This completes the proof of Proposition 2.2.

## 600 D Proof of Theorem 3.1

*Proof.* Note that  $\text{cost}^T(G)$  for any cluster tree  $T$  has a trivial upper bound. That is,

$$\text{cost}^T(G) = \sum_{e \in E} \text{cost}^T(e) \leq \sum_{e \in E} w_e \cdot \log(\text{vol}(G)) \leq \frac{\text{vol}(G) \cdot \log(\text{vol}(G))}{2},$$

601 where  $\text{cost}^T(e) = w_e \log \text{vol}(\text{LCA}_T(e))$ . Let  $T^*$  be the optimal cluster tree that achieves the mini-  
 602 mum cost, we present here a lower bound for  $\text{cost}^{T^*}(G)$ . Referring to the dense branch technique  
 603 [10, 15], we start with the root node  $A_0$  and walk along  $T^*$  recursively as follows: at every internal  
 604 node  $A_i$ , walk down to the node  $A_{i+1}$  of higher volume between its two children. This process stops  
 605 when we reach node  $A_k$  such that  $\text{vol}_G(A_k) \leq \frac{2\text{vol}(G)}{3}$ . Denote  $A \triangleq A_k$  as well as  $B \triangleq V \setminus A_k$ . By  
 606 construction, it holds that  $\text{vol}_G(A) > \frac{\text{vol}(G)}{3}$  and  $\text{vol}_G(B) \geq \frac{\text{vol}(G)}{3}$ . Moreover,  $\text{vol}_G(A_i) > \frac{2\text{vol}(G)}{3}$

607 for every  $0 \leq i < k$ . The basic idea behind the dense branch is that the  $cut(A, B)$  has a significant  
 608 contribution to  $cost^{T^*}(G)$ .

$$\begin{aligned}
 cost^{T^*}(G) &= \sum_{e=\{u,v\}} w_e \cdot \log(\text{vol}_G(u \vee v)) \\
 &\geq \sum_{\substack{e=\{u,v\} \\ e \in E\{A,B\}}} w_e \cdot \log(\text{vol}_G(u \vee v)) \\
 &\geq w(A, B) \cdot \log\left(\frac{2\text{vol}(G)}{3}\right). \\
 &\geq \Phi(G) \cdot \frac{\text{vol}(G)}{3} \cdot \log\left(\frac{2\text{vol}(G)}{3}\right).
 \end{aligned}$$

Let  $T$  be an arbitrary cluster tree, and  $T^*$  be an optimal tree. We have

$$\frac{cost^T(G)}{cost^{T^*}(G)} \leq \frac{3}{2\Phi(G)} \cdot \frac{\log(\text{vol}(G))}{\log\left(\frac{2\text{vol}(G)}{3}\right)} = O(\Phi(G)^{-1}).$$

609

□

## 610 E Proof of Theorem 3.2

611 *Proof.* To prove Theorem 3.2, we only have to prove the following lemma. Then the theorem follows  
 612 from a simplification of the approximation factor.

613 **Lemma E.1.** Let  $\alpha = \max_i \{\Phi_G(P_i)\}$  and  $\beta = \min_i \{\Phi(G[P_i])\}$ . Algorithm 2 achieves  
 614  $\left(\left(\log\left(\frac{1}{1-\alpha}\right) + 1\right) + \frac{2\alpha}{1-\alpha} \left(1 + \log\frac{k}{1-\alpha}\right)\right) \cdot \frac{3}{2\beta \log\left(\frac{4}{3}\right)}$ -approximation.

*Proof.* We group the edges of  $G$  into two categories: let  $E_1$  be the set of edges in the induced  
 subgraphs  $G[P_i]$  for all  $1 \leq i \leq l$ , i.e.,

$$E_1 \triangleq \cup_{i=1}^l E[G[P_i]],$$

and  $E_2$  be the remaining crossing edges. Then we have

$$cost^T(G) = \sum_{e \in E_1} cost^T(e) + \sum_{e \in E_2} cost^T(e).$$

We denote by  $\text{vol}(G[P_i])$  the volume of the induced graph  $G[P_i]$ , by  $\text{vol}_G(P_i)$  the volume of  $P_i$  in  
 $G$ , and by  $\text{parent}^T(P_i)$  the parent of  $P_i$  on  $T$ . Then it holds for every  $P_i$  that

$$\text{vol}_G(\text{parent}^T(P_i)) \leq k \cdot \text{vol}_G(P_i).$$

By the construction of  $T$  we have that

$$\text{vol}_G(\text{parent}^T(P_i)) = \sum_{j=1}^i \text{vol}_G(P_j) \leq i \cdot \text{vol}_G(P_i) \leq k \cdot \text{vol}_G(P_i).$$

Note that

$$\begin{aligned}
 w(P_i, V \setminus P_i) &= \text{vol}_G(P_i) - \text{vol}(G[P_i]) \leq \alpha \cdot \text{vol}_G(P_i), \\
 (1 - \alpha) \cdot \text{vol}_G(P_i) &\leq \text{vol}(G[P_i]),
 \end{aligned}$$

and thus

$$\text{vol}_G(\text{parent}^T(P_i)) \leq k \cdot \text{vol}_G(P_i) \leq \frac{k}{1-\alpha} \text{vol}(G[P_i]).$$

615 Combining the above, we have that

$$\begin{aligned}
\sum_{e \in E_1} \text{cost}^T(e) &\leq \sum_{e \in E_1} w_e \cdot \log(\text{vol}_G(P_i)) \\
&\leq \sum_{e \in E_1} w_e \cdot \log\left(\frac{1}{1-\alpha} \text{vol}(G[P_i])\right) \\
&= \sum_{e \in E_1} \left( w_e \cdot \log \frac{1}{1-\alpha} + w_e \cdot \log(\text{vol}(G[P_i])) \right) \\
&\leq \left( \log \frac{1}{1-\alpha} + 1 \right) \cdot \sum_{j=1}^k \frac{\text{vol}(G[P_i]) \cdot \log(\text{vol}(G[P_i]))}{2},
\end{aligned}$$

616 and

$$\begin{aligned}
\sum_{e \in E_2} \text{cost}^T(e) &\leq \sum_{j=1}^k w(P_i, V \setminus P_i) \cdot \log(\text{vol}_G(\text{parent}^T(P_i))) \\
&\leq \sum_{j=1}^k \frac{\alpha}{1-\alpha} \text{vol}(G[P_i]) \log\left(\frac{k}{1-\alpha} \text{vol}(G[P_i])\right) \\
&\leq \sum_{j=1}^k \frac{\alpha}{1-\alpha} \left( 1 + \log \frac{k}{1-\alpha} \right) \text{vol}(G[P_i]) \log(\text{vol}(G[P_i])) \\
&= \frac{2\alpha}{1-\alpha} \left( 1 + \log \frac{k}{1-\alpha} \right) \cdot \sum_{j=1}^k \frac{\text{vol}(G[P_i]) \cdot \log(\text{vol}(G[P_i]))}{2}.
\end{aligned}$$

Let  $T^*$  be the optimal cluster tree of  $G$ , and  $OPT_G$  be the optimal value. We have

$$OPT_G = \text{cost}_G(T^*) \geq \sum_{i=1}^l \sum_{e \in E(G[P_i])} \text{cost}_{T^*}(e) \geq \sum_{i=1}^l OPT_{G[P_i]}.$$

Denote by

$$h(\alpha, k) = \left( \left( \log \left( \frac{1}{1-\alpha} \right) + 1 \right) + \frac{2\alpha}{1-\alpha} \left( 1 + \log \frac{k}{1-\alpha} \right) \right).$$

617 We have

$$\begin{aligned}
\text{cost}^T(G) &= \sum_{e \in E_1} \text{cost}^T(e) + \sum_{e \in E_2} \text{cost}^T(e) \\
&\leq h(\alpha, k) \cdot \sum_{j=1}^k \frac{\text{vol}(G[P_i]) \cdot \log(\text{vol}(G[P_i]))}{2} \\
&\leq h(\alpha, k) \cdot \sum_{j=1}^k \frac{\text{vol}(G[P_i]) \cdot \log(\text{vol}(G[P_i]))}{2\Phi_{G[P_i]} \cdot \frac{1}{3} \text{vol}(G[P_i]) \cdot \log(\frac{2}{3} \text{vol}(G[P_i]))} OPT_{G[P_i]} \\
&\leq h(\alpha, k) \cdot \max_i \frac{3 \log(\text{vol}(G[P_i]))}{2\Phi_{G[P_i]} \cdot \log(\frac{2}{3} \text{vol}(G[P_i]))} \sum_{j=1}^k OPT_{G[P_i]} \\
&\leq h(\alpha, k) \cdot \max_i \frac{3 \log(\text{vol}(G[P_i]))}{2\Phi_{G[P_i]} \cdot \log(\frac{2}{3} \text{vol}(G[P_i]))} OPT_G \\
&\leq h(\alpha, k) \cdot \frac{3}{2\beta \log(\frac{4}{3})} OPT_G
\end{aligned}$$

618 Lemma E.1 follows. □

619 Note that  $h(\alpha, k) = O\left(\frac{1}{(1-\alpha)} \log \frac{k}{1-\alpha}\right)$ , Theorem 3.2 follows. □

620 **F Experimental results on Amazon network**

621 We do our experiments on Amazon network <sup>4</sup> for which the set of ground-truth clusters has been  
 622 given. For two sets  $A, B$ , the *Jaccard Index* of them is defined as  $J(A, B) = |A \cap B| / |A \cup B|$ . We  
 623 pick the largest cluster which is a subgraph with 58283 vertices and 133178 edges. We run HCSE  
 624 algorithm on it. For each ground-truth cluster  $c$  that appears in this subgraph, we find from the  
 625 resulting cluster tree an internal node that has maximum Jaccard index with  $c$ . Then we calculate  
 626 the average Jaccard index  $\bar{J}$  over all such  $c$ . We also calculate  $\text{cost}(\text{SE})$  and  $\text{cost}(\text{Das})$ . The results  
 627 are demonstrated in Table 3. HCSE performs better for  $\bar{J}$  and  $\text{cost}(\text{SE})$ , while LOUVAIN performs  
 628 better for  $\text{cost}(\text{Das})$ . Because of unbalance in over-fitting and under-fitting traps, HLP outperforms  
 629 none of the other two algorithms for all criteria.

index	HCSE	HLP	LOUVAIN
$\bar{J}$	<b>0.20</b>	0.16	0.17
$\text{cost}(\text{SE})$	<b>1.85E6</b>	2.05E6	1.89E6
$\text{cost}(\text{Das})$	5.57E8	3.99E8	<b>3.08E8</b>

Table 3: Comparisons of the average Jaccard index ( $\bar{J}$ ), cost function based on structural entropy ( $\text{cost}(\text{SE})$ ) and Dasgupta’s cost function ( $\text{cost}(\text{Das})$ ).

630 **G Some figures and pseudocodes**

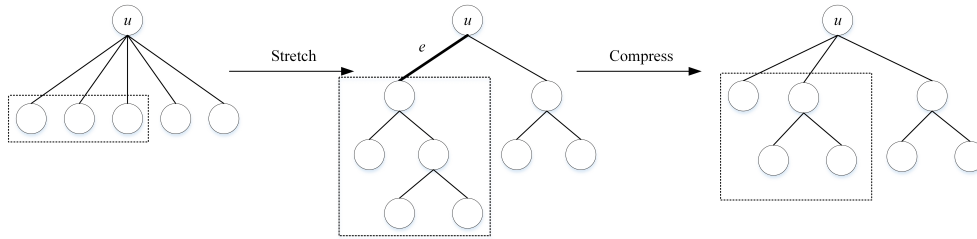


Figure 3: Illustrations of stretch and compress for a  $u$ -triangle. A binary cluster tree is constructed first by stretch, and then edge  $e$  is compressed, which yields a non-binary tree.

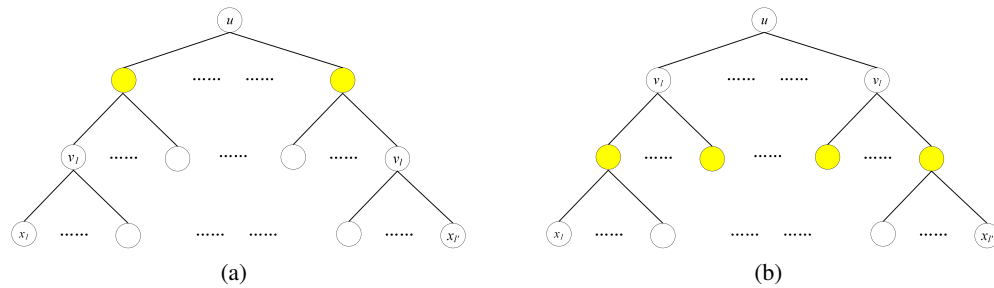


Figure 4: Illustration of stratification for a 2-level cluster tree. The preference of (a) and (b) depends on the average sparsity of triangles at each level.

<sup>4</sup><http://snap.stanford.edu/data/>

---

**Algorithm 5:** Stretch

---

**Input:** a  $u$ -triangle  $T_u$

**Output:** a binary tree rooted at  $u$

```
1 Let  $\{v_1, v_2, \dots, v_\ell\}$  be the set of leaves of  $T_u$ ;  
2 Compute  $\eta(a, b)$  which is the cost reduced by merging siblings  $a, b$  into a single cluster;  
3 for  $t \in [\ell - 1]$  do  
4    $(\alpha, \beta) \leftarrow \arg \max_{(a,b) \text{ are siblings}} \{\eta(a, b)\}$ ;  
5   Add a new node  $\gamma$ ;  
6    $\gamma.\text{parent} \leftarrow \alpha.\text{parent}$ ;  
7    $\alpha.\text{parent} = \gamma$ ;  
8    $\beta.\text{parent} = \gamma$ ;  
9 return  $T_u$ 
```

---

---

**Algorithm 6:** Compress

---

**Input:** a binary tree  $T$

```
1 while  $T$ 's height is more than 2 do  
2    $e \leftarrow \arg \min_{e' \in \hat{E}(T)} \{\Delta(e')\}$ ;  
3   Denote  $e = (u, v)$  where  $u$  is the parent of  $v$ ;  
4   for  $w \in v.\text{children}$  do  
5      $w.\text{parent} \leftarrow u$ ;  
6   Delete  $v$  from  $T$ ;
```

---