

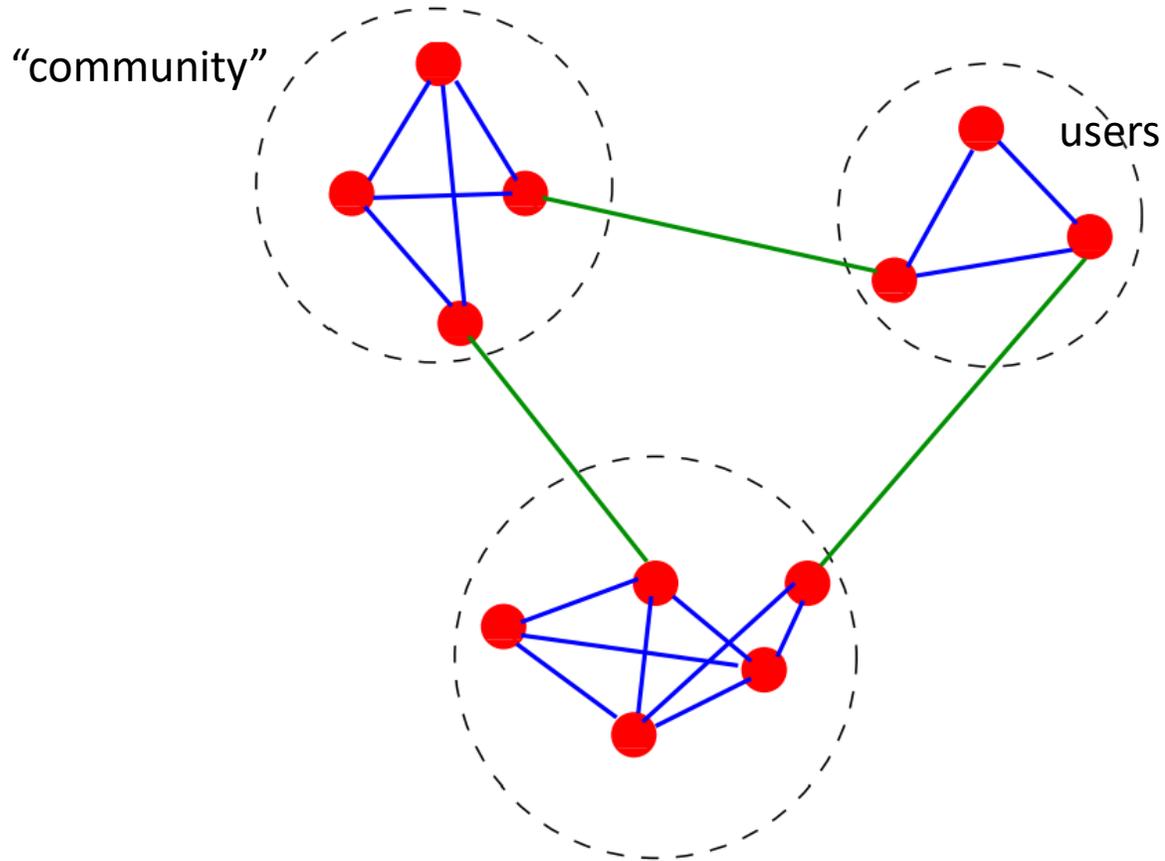
Detecting Communities under Differential Privacy

Hiep H. Nguyen, Abdessamad Imine, and Michael Rusinowitch

2016

Presenter: Zifeng Wang (zifengw2)

Community detection (CD)



Modularity:

of clusters

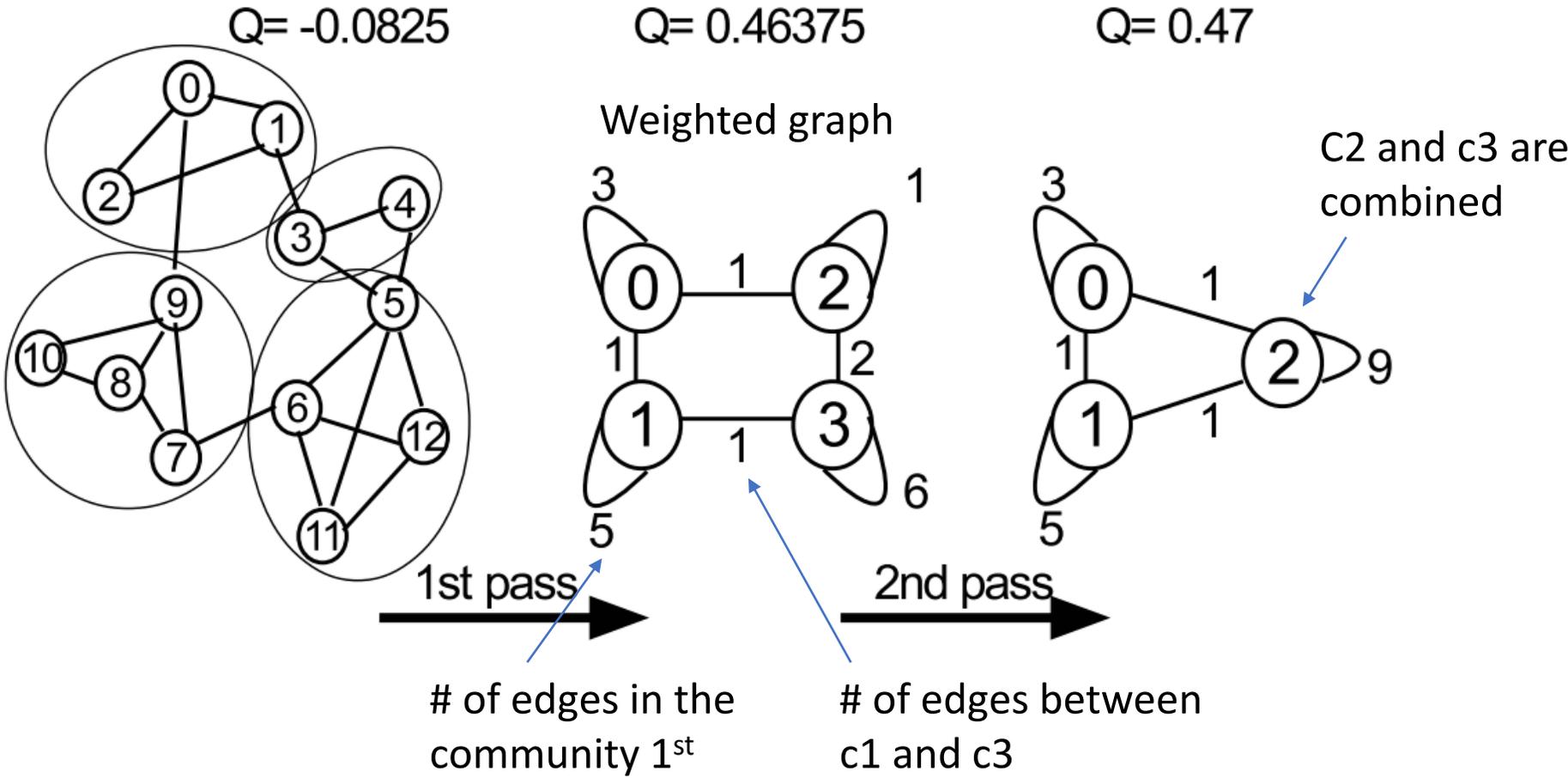
edges joining the community c

$$Q = \sum_{c=1}^{n_c} \left[\frac{l_c}{m} - \left(\frac{d_c}{2m} \right)^2 \right]$$

of edges in the whole graph

Sum of degree of nodes in c

Louvain method for CD



Differential privacy

Definition 2.4 (Differential Privacy). A randomized algorithm \mathcal{M} with domain $\mathbb{N}^{|\mathcal{X}|}$ is (ϵ, δ) -differentially private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all $x, y \in \mathbb{N}^{|\mathcal{X}|}$ such that $\|x - y\|_1 \leq 1$:

$$\Pr[\mathcal{M}(x) \in \mathcal{S}] \leq \exp(\epsilon) \Pr[\mathcal{M}(y) \in \mathcal{S}] + \delta,$$

“privacy budget”



Two categories of DP CD



(a) Input perturbation



(b) Algorithm perturbation

ϵ -DP graph

Definition 3.1: A mechanism \mathcal{A} is ϵ -differentially private if for any two neighboring graphs G_1 and G_2 , and for any output $O \in \text{Range}(\mathcal{A})$,

$$\Pr[\mathcal{A}(G_1) \in O] \leq e^\epsilon \Pr[\mathcal{A}(G_2) \in O]$$

Neighboring graph: $G_1 = (V_1, E_1)$ $V_1 = V_2$
 $G_2 = (V_2, E_2)$ $E_1 \subset E_2$ and $|E_2| = |E_1| + 1$

Input perturbations

- Laplace mechanism

Global sensitivity: $\Delta f = \max_{G_1, G_2} \|f(G_1) - f(G_2)\|_1$

Theorem 3.1: (Laplace mechanism [11]) For any function $f : G \rightarrow \mathbb{R}^d$, the mechanism \mathcal{A}

$$\mathcal{A}(G) = f(G) + \langle \text{Lap}_1\left(\frac{\Delta f}{\epsilon}\right), \dots, \text{Lap}_d\left(\frac{\Delta f}{\epsilon}\right) \rangle \quad (2)$$

satisfies ϵ -differential privacy, where $\text{Lap}_i\left(\frac{\Delta f}{\epsilon}\right)$ are i.i.d Laplace variables with scale parameter $\frac{\Delta f}{\epsilon}$. \square

Laplace distribution: $\text{Lap}(\lambda) : p(x|\lambda) = \frac{1}{2\lambda} e^{-|x|/\lambda}$

Input perturbations

- Geometric mechanism $Geom(\alpha) : Pr(\Delta = \delta | \alpha) = \frac{1-\alpha}{1+\alpha} \alpha^{|\delta|}$

For a function $f : G \rightarrow \mathbb{Z}^d$ The mechanism $\mathcal{A}(G) = f(G) + \{\Delta_i\}_{i=1}^d$
where $\Delta \sim Geom(\alpha)$ and $\alpha = \exp(-\varepsilon)$ satisfies ε -DP.

Δ is guaranteed to be integral

Algorithm1: LouvainDP

Basic idea: create a noisy weighted supergraph G_1 from G by grouping nodes with equal size k

Algorithm 1 LouvainDP(G, s)

Input: undirected graph G , group size k , privacy budget ϵ

Output: noisy partition \tilde{C}

Randomly initialize G_1 by node permutation and grouping

- 1: $G_1 \leftarrow \emptyset, n_1 = \lfloor \frac{|V|}{k} \rfloor - 1, V_1 \leftarrow \{0, 1, \dots, n_1\}$
- 2: $\epsilon_2 = 0.1, \epsilon_1 = \epsilon - \epsilon_2, \alpha = \exp(-\epsilon_1)$
- 3: get a random permutation V_p of V
- 4: compute the mapping $M : V_p \rightarrow V_1$

Noisy # of superedges

- 5: compute superedges of $G_1: E_1 = \{e_1(i, j)\}$ where $i, j \in V_1$

- 6: $m_1 = |E_1| + \text{Lap}(1/\epsilon_2), m_0 = \frac{n_1(n_1+1)}{2}$

- 7: $\theta = \lceil \log_{\alpha} \frac{(1+\alpha)m_1}{m_0 - m_1} \rceil$

- 8: $s = (m_0 - m_1) \frac{\alpha^{\theta}}{1+\alpha}$

of all possible edges in G_1

Algorithm1: LouvainDP

Basic idea: create a noisy weighted supergraph G_1 from G by grouping nodes with equal size k

Go through non-zero superedges, add geometric noise, decide if adding edge by threshold θ

Go through zero superedges, draw integral weight and add edge if w is larger than zero

```
7:  $\theta = \lceil \log_{\alpha} \frac{(1+\alpha)m_1}{m_0 - m_1} \rceil$ 
8:  $s = (m_0 - m_1) \frac{\alpha^{\theta}}{1+\alpha}$ 
9: for  $e_1(i, j)$  in  $E_1$  do
10:    $e_1(i, j) = e_1(i, j) + Geom(\alpha)$ 
11:   if  $e_1(i, j) \geq \theta$  then
12:     add  $e_1(i, j)$  to  $G_1$ 
13: for  $s$  edges sampled uniformly at random  $e_1(i, j) \notin E_1$  do
14:   draw  $w$  from the distribution  $Pr[X \leq x] = 1 - \alpha^{x-\theta+1}$ 
15:   if  $w > 0$  then
16:     add edge  $e_1(i, j)$  with weight  $w$  to  $G_1$ 
17: run Louvain method on  $G_1$  to get  $\tilde{C}_1$ 
18: compute  $\tilde{C}$  from  $\tilde{C}_1$  using the mapping  $M$ 
19: return  $\tilde{C}$ 
```

Algorithm1: LouvainDP

Basic idea: create a noisy weighted supergraph G_1 from G by grouping nodes with equal size k

High-pass filtering (Cormode 2012) for private data sampling

```
7:  $\theta = \lceil \log_{\alpha} \frac{(1+\alpha)m_1}{m_0 - m_1} \rceil$ 
8:  $s = (m_0 - m_1) \frac{\alpha^{\theta}}{1+\alpha}$ 
9: for  $e_1(i, j)$  in  $E_1$  do
10:    $e_1(i, j) = e_1(i, j) + \text{Geom}(\alpha)$ 
```

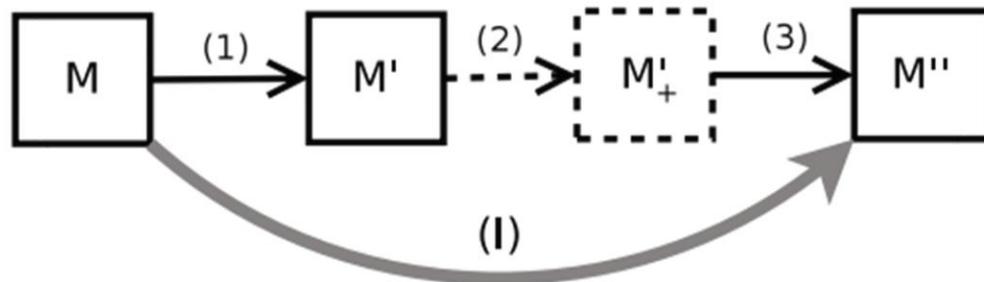


Figure 2: Anonymization process: (1) adding Geometric/Laplace noise, (2) filtering negative entries, (3) sampling, (I) shortcut to generate M'' directly from M .

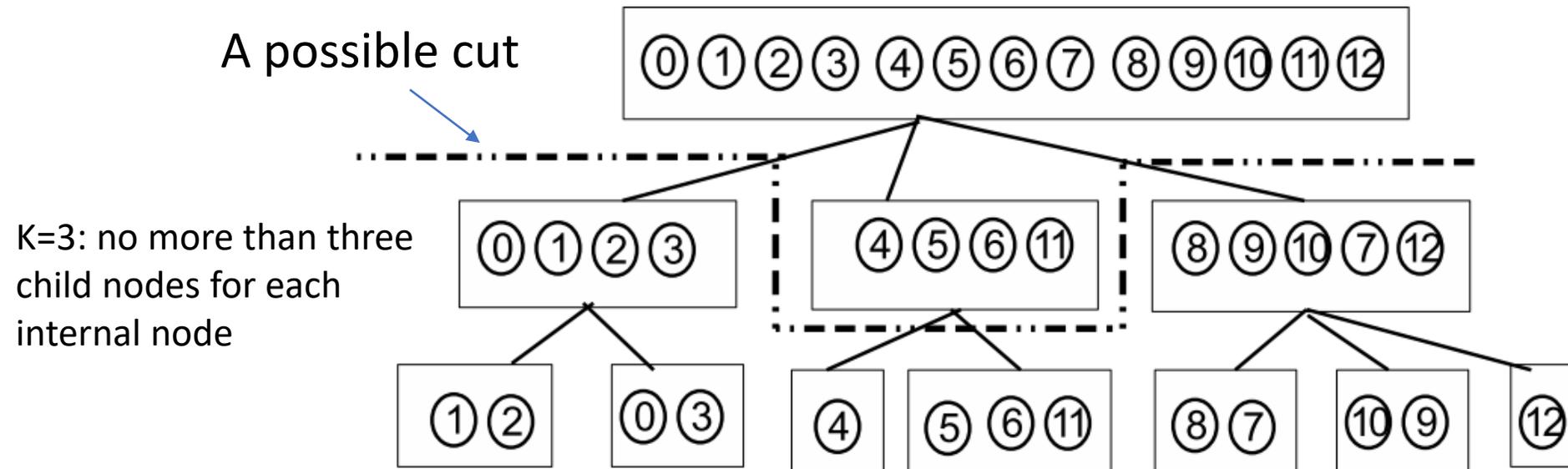
1
uniformly at random $e_1(i, j) \notin E_1$ **do**
distribution $Pr[X \leq x] = 1 - \alpha^{x-\theta+1}$

with weight w to G_1

G_1 to get \tilde{C}_1
ing the mapping M

Algorithm perturbations

Basic idea: heuristically detect cohesive groups of nodes privately



Challenges:

- Efficiently find good split of nodes with high modularity under DP restrictions
- Merge small groups

Preliminary: Exponential mechanism

Theorem 3.2: (Exponential mechanism [21]) Given a score function $u : (G \times \mathcal{O}) \rightarrow \mathbb{R}$ for a graph G , the mechanism \mathcal{A} that samples an output O with probability proportional to $\exp\left(\frac{\epsilon \cdot u(G, O)}{2\Delta u}\right)$ satisfies ϵ -differential privacy. \square

Global sensitive $\Delta u = \max_{O, G_1, G_2} |u(G_1, O) - u(G_2, O)|$

Sampling output $O_i \sim \exp\left(\frac{\epsilon \cdot u(G, O_i)}{2\Delta u}\right) / \sum_j \exp\left(\frac{\epsilon \cdot u(G, O_j)}{2\Delta u}\right)$

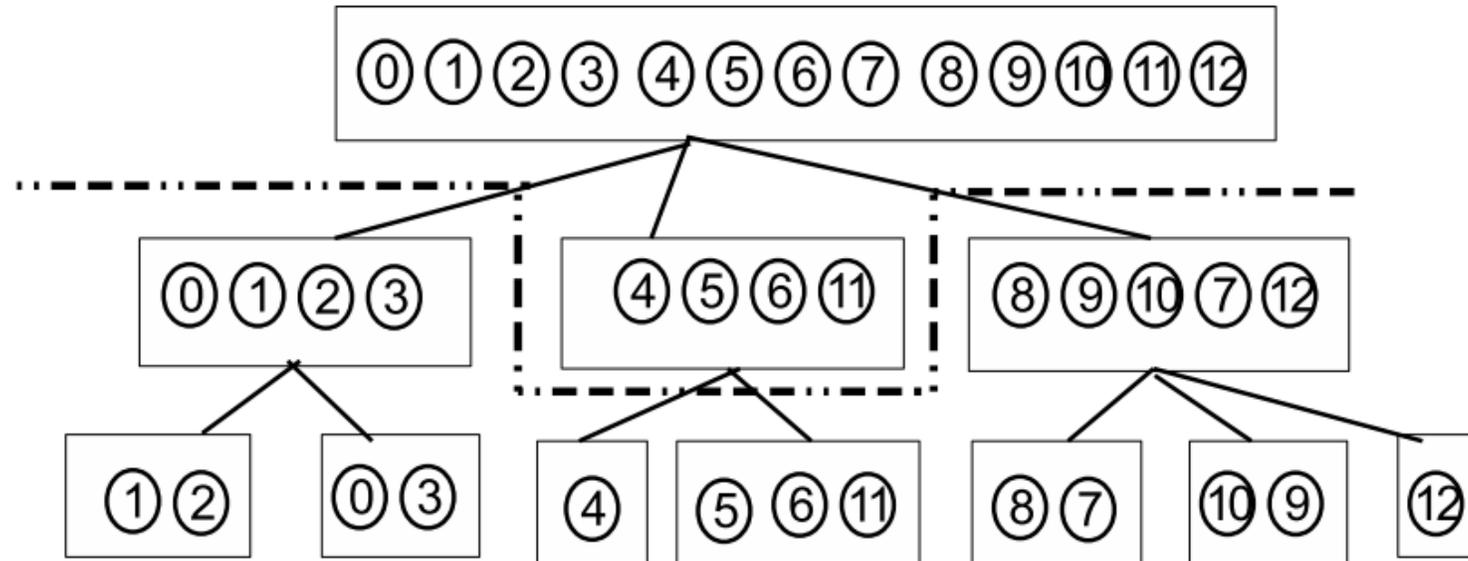
Preliminary: Composability of DP

Theorem 3.3: (Sequential and parallel compositions [22])
Let each A_i provide ϵ_i -differential privacy. A sequence of $A_i(D)$ over the dataset D provides $\sum_{i=1}^n \epsilon_i$ -differential privacy.
Let each A_i provide ϵ_i -differential privacy. Let D_i be arbitrary disjoint subsets of the dataset D . The sequence of $A_i(D_i)$ provides $\max_{i=1}^n \epsilon_i$ -differential privacy. \square

Sequential: same D , different A

Parallel: different D , different A

Algorithm: ModDivisive



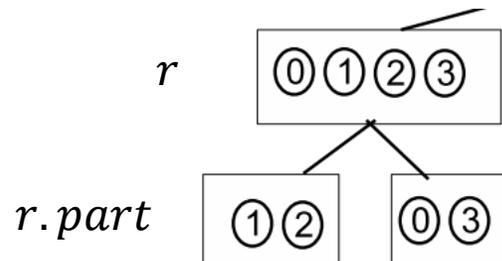
Stage 1st : DP sampling a tree of depth $maxL$, use budget ϵ_1

Stage 2nd : find best cut for all levels, use budget $\epsilon_m * maxL$

Algorithm: ModDivisive

Allocate the privacy budget for all levels in the tree: $\sum_i eA[i] = \epsilon_1$ cuz the sequential composability of DP

Go across each node in Q , try to make a node partition by **ModMCMC** algorithm. Then attach each partition as the child nodes



Algorithm 3 ModDivisive

Input: graph G , group size k , privacy budget ϵ , max level $maxL$, ratio λ , BestCut privacy at each level ϵ_m

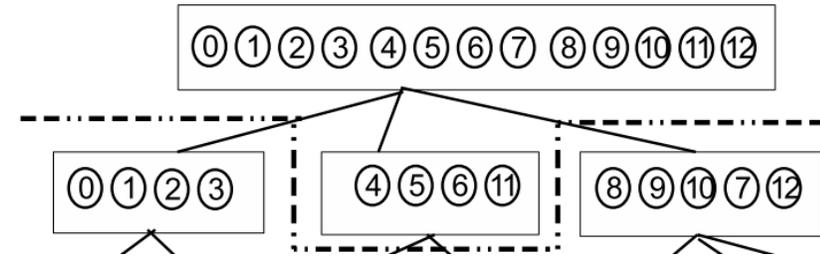
Output: noisy partition \tilde{C}

```

1: compute the array  $eA[0..maxL-1]$  s.t.  $\sum_i eA[i] = \epsilon_1$ ,  $eA[i] = eA[i+1] * \lambda$  where  $\epsilon_1 = \epsilon - maxL * \epsilon_m$ 
2: initialize the root node with nodeset  $V$ 
3:  $root = NodeSet(G, V, k)$ 
4:  $root.level = 0$ 
5: queue  $Q \leftarrow root$ 
6: while  $Q$  is not empty do
7:    $r \leftarrow Q.dequeue()$ 
8:   if  $r.level < maxL$  then
9:      $r.part = ModMCMC(G, r.S, k, eA[r.level])$ 
10:    for subset  $S_i$  in  $r.part$  do
11:       $P_i = NodeSet(G, S_i, k)$ 
12:       $P_i.level = r.level + 1$ 
13:       $r.child_i \leftarrow P_i$ 
14:       $Q.enqueue(P_i)$ 
15:  $\tilde{C} \leftarrow BestCut(root, \epsilon_m)$ 
16: return  $\tilde{C}$ 

```

Algorithm: ModMCMC



```

8:   if  $r.level < maxL$  then
9:      $r.part = \text{ModMCMC}(G, r.S, k, eA[r.level])$ 
10:    for subset  $S_i$  in  $r.part$  do
11:       $P_i = \text{NodeSet}(G, S_i, k)$ 
12:       $P_i.level = r.level + 1$ 
13:       $r.child_i \leftarrow P_i$ 
14:       $Q.enqueue(P_i)$ 

```

Transform partition P_{i-1} to P_i with the specified probability by exponential mechanism

Algorithm 4 ModMCMC

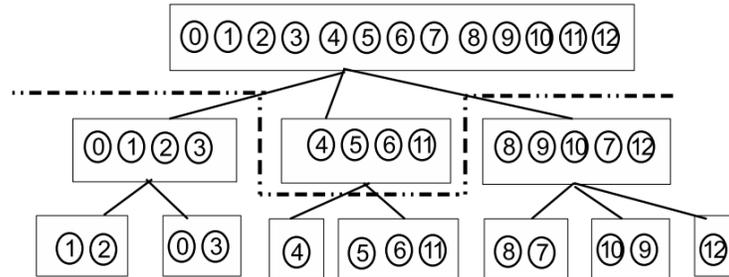
Input: graph G , nodeset $r.S$, group size k , privacy budget ϵ_p

Output: sampled partition $r.part$

- 1: initialize $r.part$ with a random partition P_0 of k groups
 - 2: **for** each step i in the Markov chain **do**
 - 3: pick a neighboring partition P' of P_{i-1} by randomly selecting node $u \in r.S$ and moving u to another group.
 - 4: accept the transition and set $P_i = P'$ with probability

$$\min\left(1, \frac{\exp(\frac{\epsilon_p}{2\Delta Q} Q(P', G))}{\exp(\frac{\epsilon_p}{2\Delta Q} Q(P_{i-1}, G))}\right)$$
 - 5: // until equilibrium is reached
 - 6: **return** a sampled partition $r.part = P_i$
-

Algorithm: ModDivisive



Algorithm 5 BestCut

Input: undirected graph G , root node $root$, privacy budget at each level ϵ_m

Output: best cut C

```

1: stack  $S \leftarrow \emptyset$ , queue  $Q \leftarrow root$ 
2: while  $Q$  is not empty do
3:    $r \leftarrow Q.dequeue()$ 
4:    $S.push(r)$ 
5:   for child node  $r_i$  in  $r.children$  do
6:      $Q.enqueue(r_i)$ 
7: dictionary  $sol \leftarrow \emptyset$ 
8: while  $S$  is not empty do
9:    $r \leftarrow S.pop()$ ,  $r.mod_n = r.mod + Laplace(\Delta Q/\epsilon_m)$ 
10:  if  $r$  is a leaf node then
11:     $sol.put(r.id, (val = r.mod_n, self = True))$ 
12:  else
13:     $s_m = \sum_{r_i \in r.children} sol[r_i.id].mod_n$ 
14:    if  $r.mod_n < s_m$  then
15:       $sol.put(r.id, (val = s_m, self = False))$ 
16:    else
17:       $sol.put(r.id, (val = r.mod_n, self = True))$ 
18: list  $C \leftarrow \emptyset$ , queue  $Q \leftarrow root$ 
19: while  $Q$  is not empty do
20:    $r \leftarrow Q.dequeue()$ 
21:   if  $sol[r.id].self == True$  then
22:      $C = C \cup \{r\}$ 
23:   else
24:     for child node  $r_i$  in  $r.children$  do
25:        $Q.enqueue(r_i)$ 
return  $C$ 

```

Get a cut on the final partition to reach the best modularity

Algorithm 3 ModDivisive

Input: graph G , group size k , privacy budget ϵ , max level $maxL$, ratio λ , BestCut privacy at each level ϵ_m

Output: noisy partition \tilde{C}

```

1: compute the array  $eA[0..maxL-1]$  s.t.  $\sum_i eA[i] = \epsilon_1$ ,  $eA[i] = eA[i+1] * \lambda$  where  $\epsilon_1 = \epsilon - maxL * \epsilon_m$ 
2: initialize the root node with nodeset  $V$ 
3:  $root = NodeSet(G, V, k)$ 
4:  $root.level = 0$ 
5: queue  $Q \leftarrow root$ 
6: while  $Q$  is not empty do
7:    $r \leftarrow Q.dequeue()$ 
8:   if  $r.level < maxL$  then
9:      $r.part = ModMCMC(G, r.S, k, eA[r.level])$ 
10:    for subset  $S_i$  in  $r.part$  do
11:       $P_i = NodeSet(G, S_i, k)$ 
12:       $P_i.level = r.level + 1$ 
13:       $r.child_i \leftarrow P_i$ 
14:       $Q.enqueue(P_i)$ 
15:  $\tilde{C} \leftarrow BestCut(root, \epsilon_m)$ 
16: return  $\tilde{C}$ 

```

Experiments

- Aims: verify if the DP CD gets good clustering quality in good efficiency

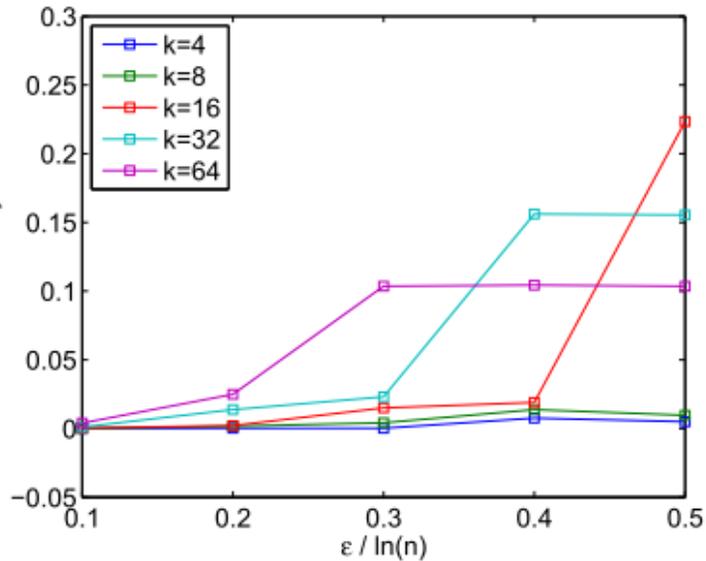
TABLE II: Characteristics of the test graphs

	Nodes	Edges	Com	Mod
as20graph	6,474	12,572	30	0.623
ca-AstroPh	17,903	196,972	37	0.624
amazon	334,863	925,872	257	0.926
dblp	317,080	1,049,866	375	0.818
youtube	1,134,890	2,987,624	13,485	0.710

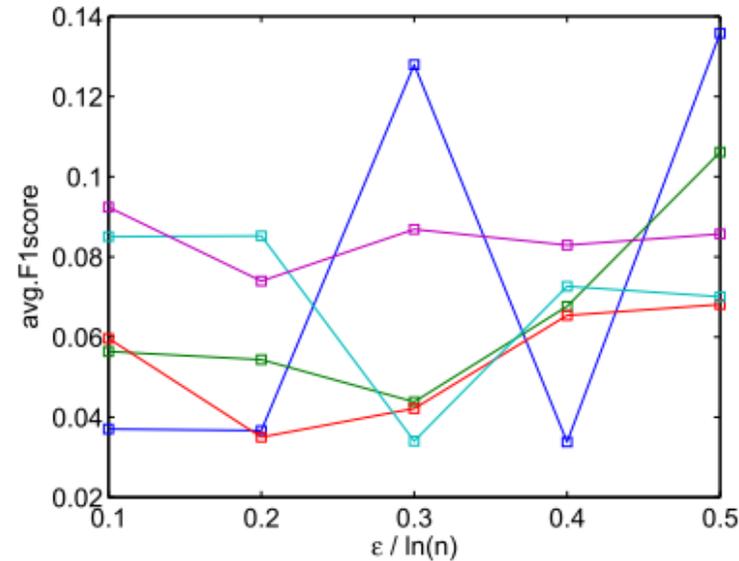
Obtained by Louvain method

Performance of LouvainDP

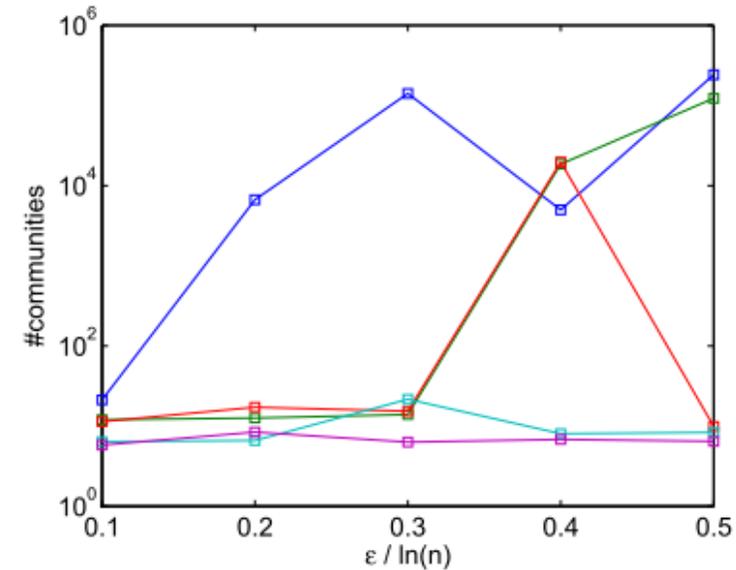
Group size k



(a)



(b)



(c)

$$\epsilon = \{0.1 \ln n, 0.2 \ln n, \dots, 0.5 \ln n\}$$

of nodes in the graph

Performance of ModDivisive

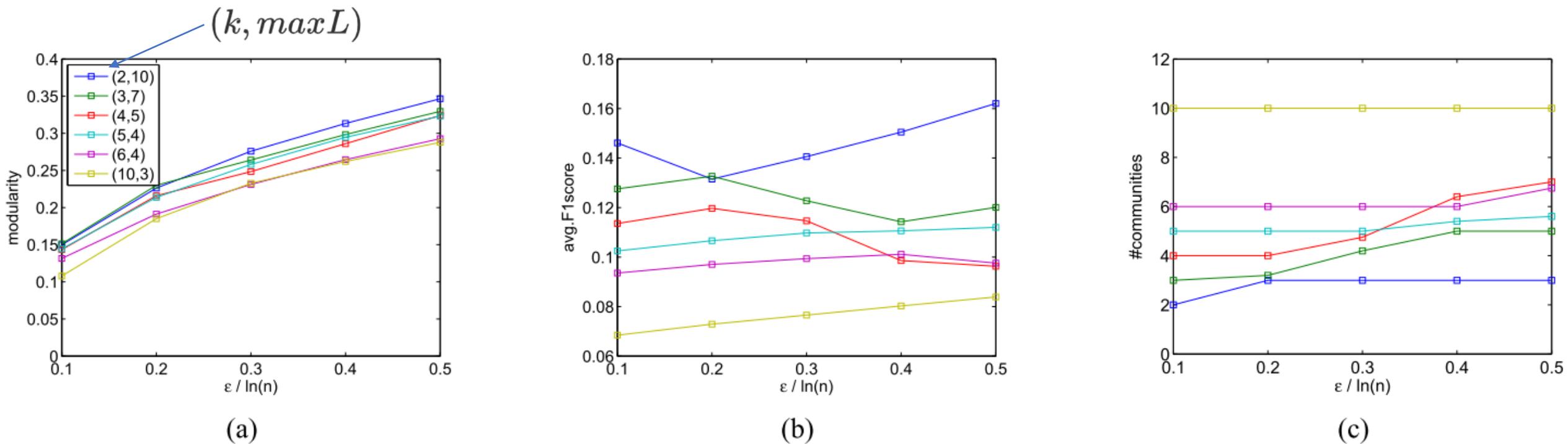
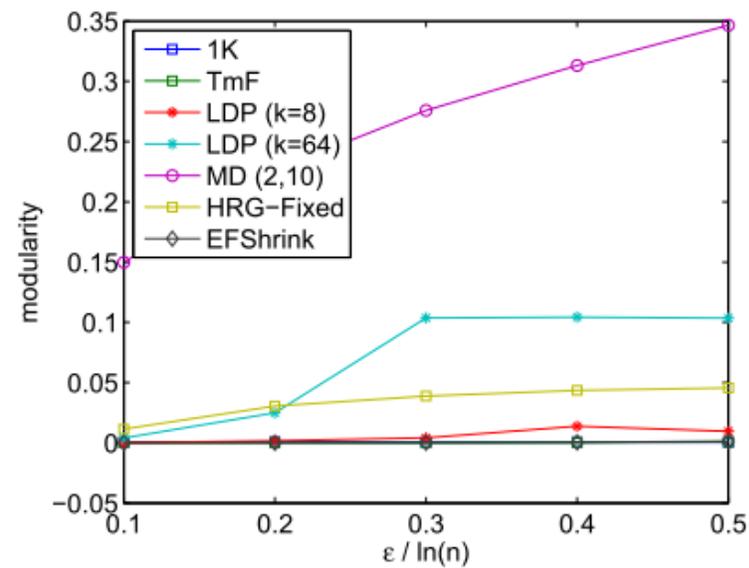
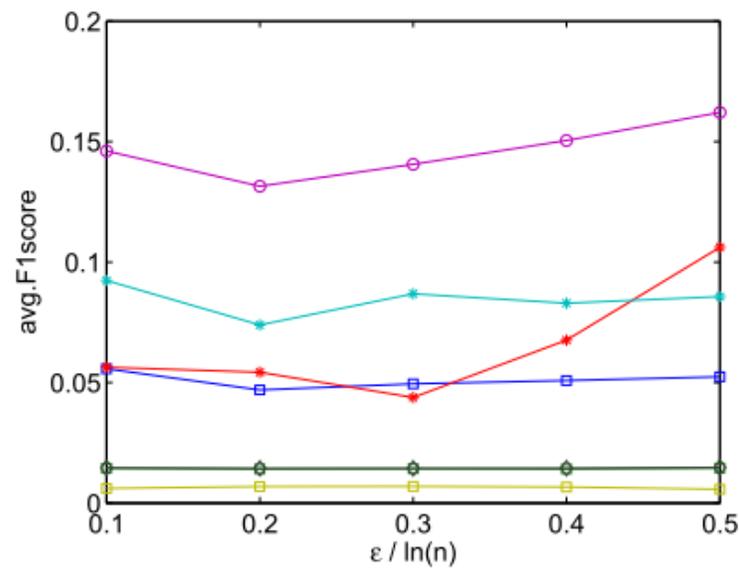


Fig. 6: ModDivisive on youtube with $\lambda = 2.0$, $K = 50$ ($0.5 \ln n = 7.0$)

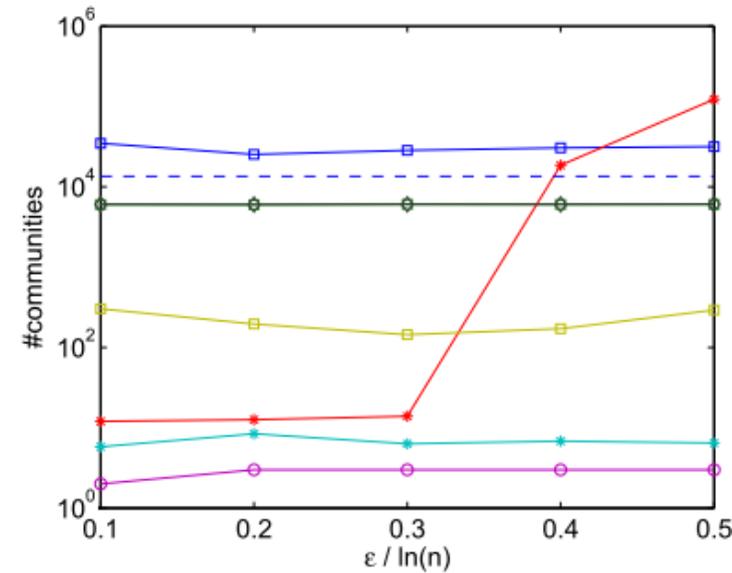
Performance comparison



(a)



(b)



(c)

Fig. 13: Quality metrics and the number of communities (youtube) ($0.5 \ln n = 7.0$)