

# Efficient Estimation of Heat Kernel PageRank for Local Clustering

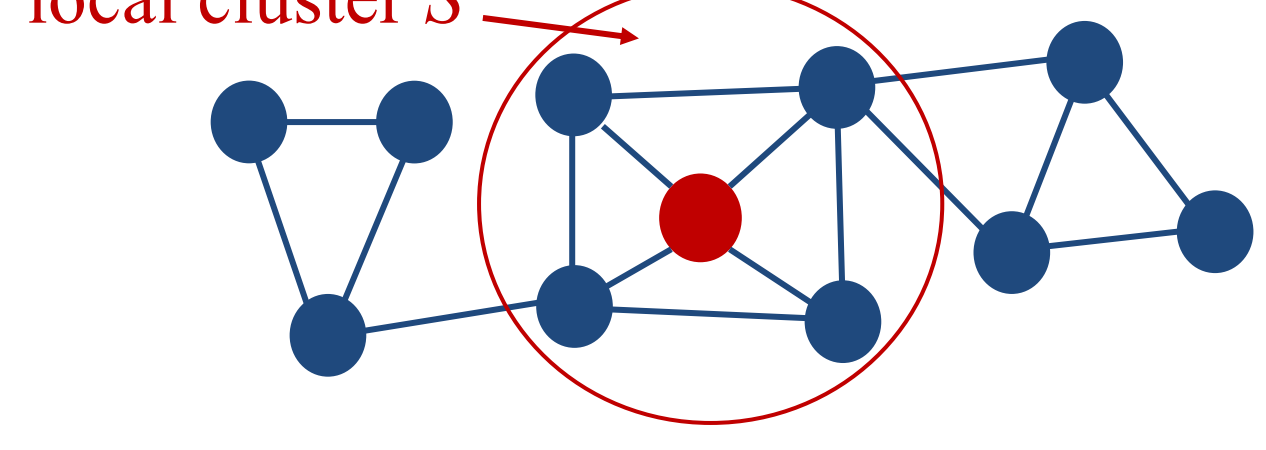
Renchi Yang, Xiaokui Xiao, Zhewei Wei, Sourav Bhowmick, Jun Zhao, Rong-Hua Li



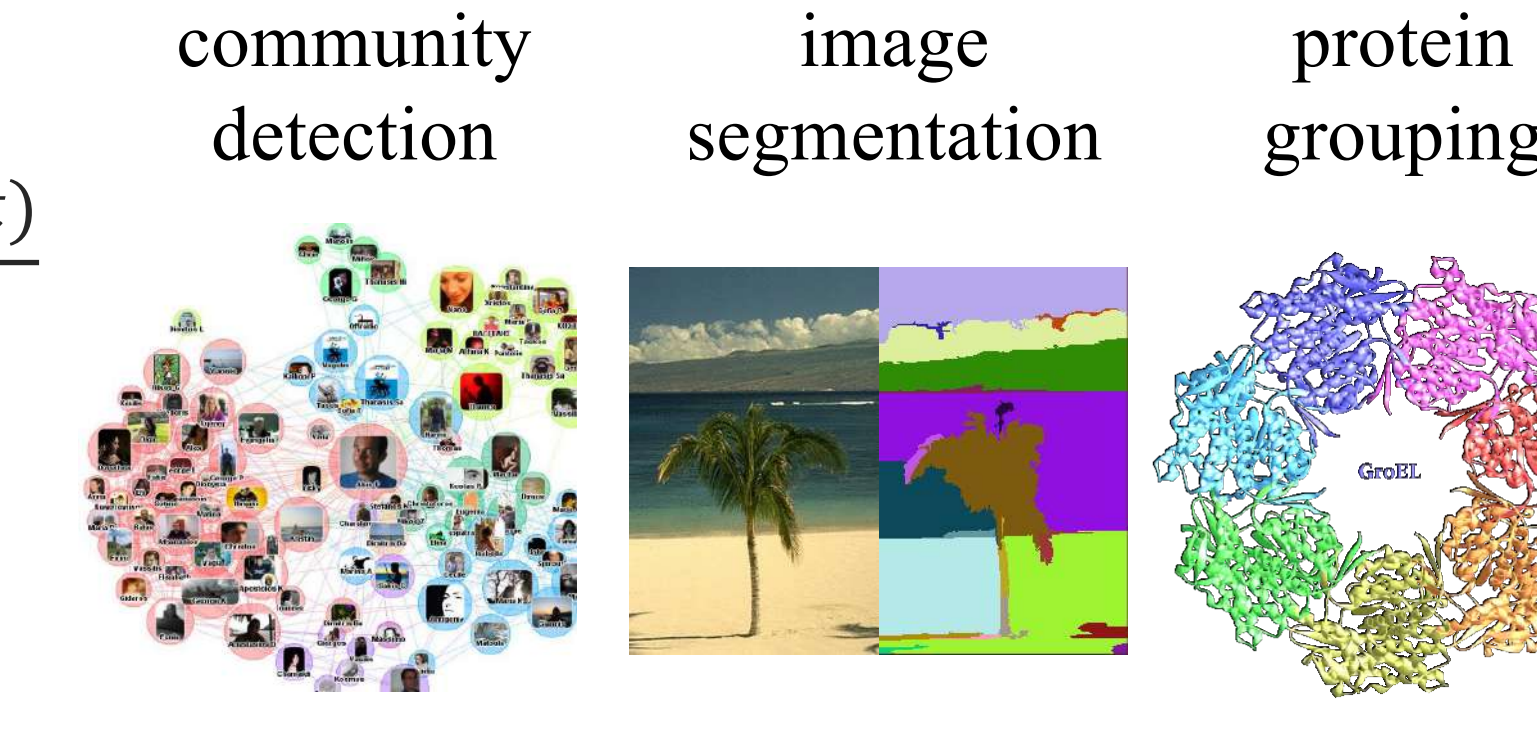
## 1. Heat Kernel-based Local Clustering

### Local Clustering

- Explodes the local neighbourhood around the seed node only to find  $S$
- $S$  has min conductance  $\phi(S) = \frac{\#(\text{edges cut})}{\sum_{u \in S} d(u)}$



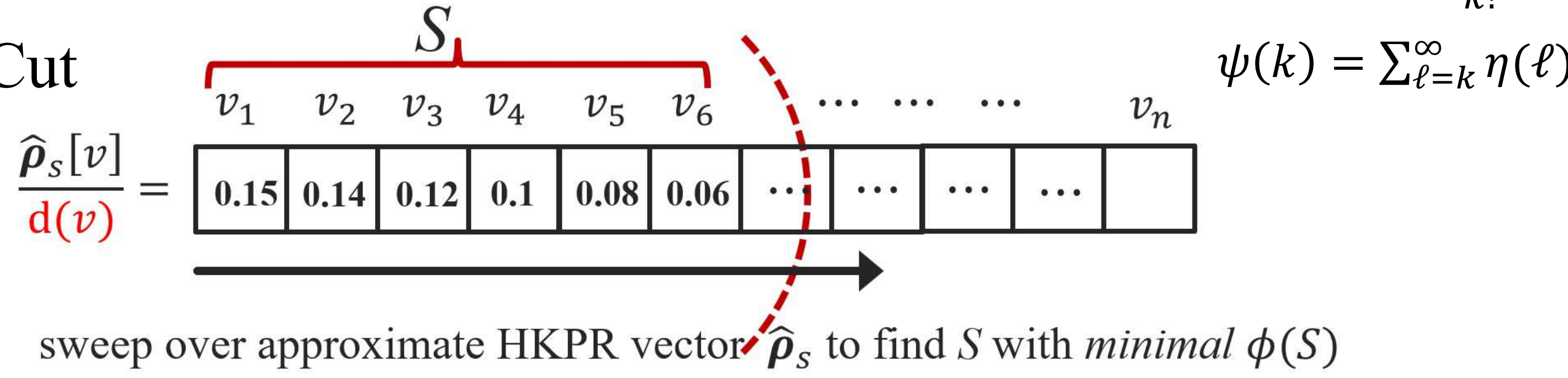
### Applications



### The Heat Kernel PageRank (HKPR) from $s$ to $v$ is

- $\rho_s[v] = \mathbb{P}[\text{Random walk of length-}k \text{ from } s \text{ stops at } v]$
- $k$  follows a Poisson distribution with mean  $t$ ;  $k$ 's probability:  $\eta(k) = \frac{e^{-t} t^k}{k!}$

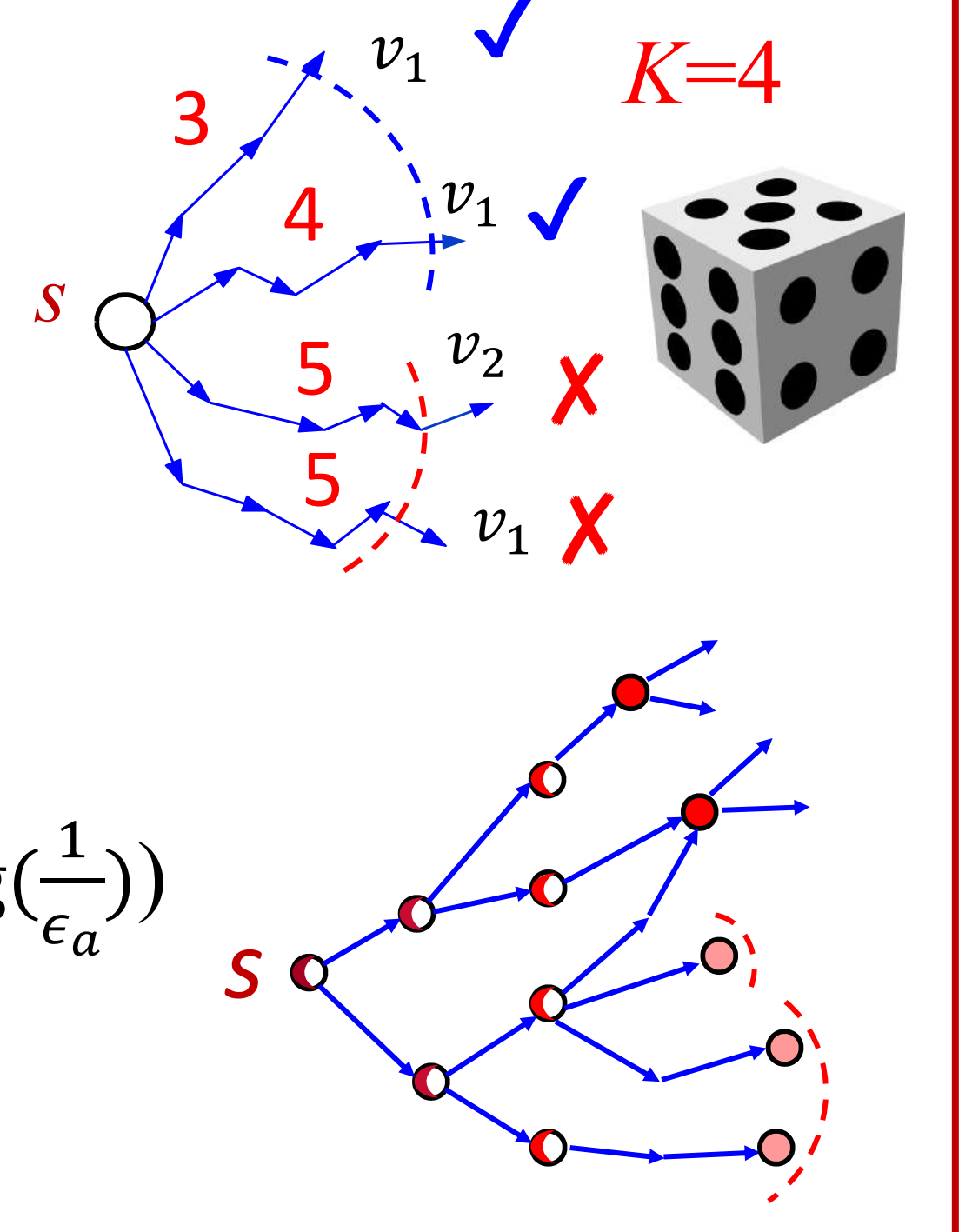
### Sweep Cut



## 2. Existing Approximate Solutions

### ClusterHKPR

- Sets max random walk length  $K = O(\frac{\log(1/\epsilon)}{\log \log(1/\epsilon)})$
- $16 \log(n)/\epsilon^3$  truncated random walks from  $s$
- $\hat{\rho}_s[v]$  = Fraction of random walks stopping at  $v$



### HK-Relax

- Sets initial residual  $r_s[s, 0] = e^{-t}$
- At  $k$ -th hop from  $s$ ,  $r_s[v, k] \rightarrow$  reserve ( $k \leq 2t \log(\frac{1}{\epsilon_a})$ ); distributes  $\frac{t}{k+1} \times r_s[v, k]$  to neighbors evenly
- $\hat{\rho}_s[v]$  = Sum of reserves at  $v$

Algorithm	Accuracy Guarantee	Complexity
ClusterHKPR	$\mathbb{P}\left\{\begin{aligned}  \hat{\rho}_s[v] - \rho_s[v]  &\leq \epsilon \cdot \rho_s[v], & \text{if } \rho_s[v] > \epsilon \\  \hat{\rho}_s[v] - \rho_s[v]  &\leq \epsilon, & \text{otherwise} \end{aligned} \right\} \geq 1 - \epsilon$	$O(\frac{t \log(n)}{\epsilon^3})$
HK-Relax	$\left  \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right  \leq \epsilon_a$	$O(\frac{te^t \log(1/\epsilon_a)}{\epsilon_a})$
Our solutions	$\mathbb{P}\left\{\begin{aligned} \left  \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right  &\leq \epsilon_r \cdot \frac{\rho_s[v]}{d(v)}, & \text{if } \frac{\rho_s[v]}{d(v)} > \delta \\ \left  \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right  &\leq \epsilon_r \cdot \delta, & \text{otherwise} \end{aligned} \right\} \geq 1 - p_f$	$O(\frac{t \log(n/p_f)}{\epsilon_r^2 \cdot \delta})$

## 3. The Basic Ideas

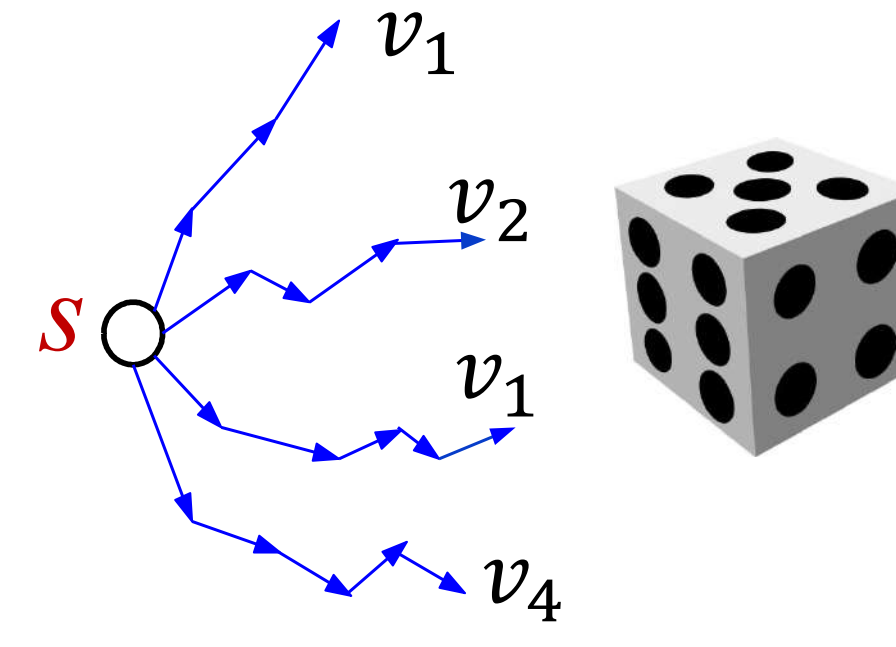
### $(d, \epsilon_r, \delta)$ -approximate HKPR

- $\forall v \in G$  s.t.  $\rho_s[v]/d(v) > \delta$ ,  $\left| \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right| \leq \epsilon_r \cdot \frac{\rho_s[v]}{d(v)}$
- $\forall v \in G$  s.t.  $\rho_s[v]/d(v) \leq \delta$ ,  $\left| \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right| \leq \epsilon_r \cdot \delta$

### Monte-Carlo Random Walks

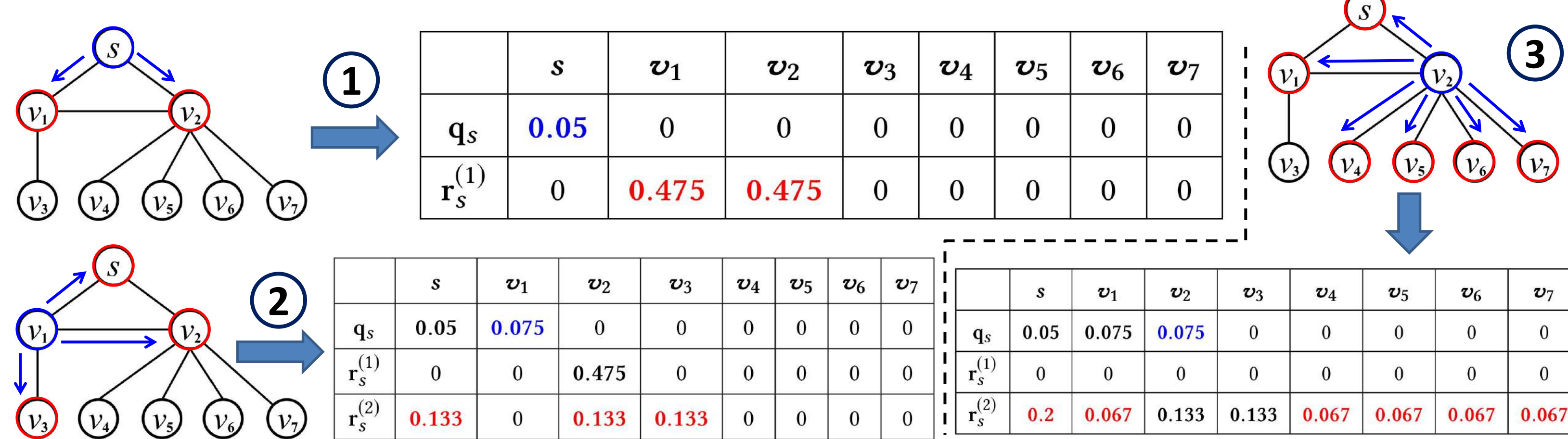
- At  $k$ -th hop, stops with probability  $\frac{\eta(k)}{\psi(k)}$ ; otherwise, jumps to a random neighbor.  $\omega = \frac{2(1+\epsilon_r/3)\log(n/p_f)}{\epsilon_r^2 \delta}$  random walks.

- $\hat{\rho}_s[v]$  = Fraction of random walks stopping at  $v$ .  $\hat{\rho}_s[v]$  is a  $(d, \epsilon_r, \delta)$ -approximate HKPR vector.



### HK-Push

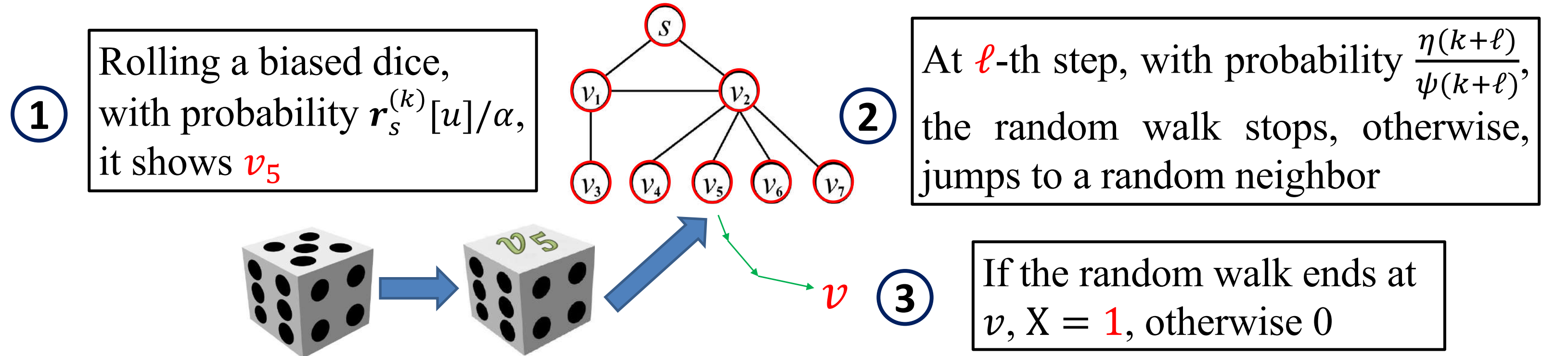
- $\rho_s[v] = \underbrace{q_s[v]}_{\text{estimation}} + \underbrace{\sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v]}_{\text{error!}}$
- Each node  $v$ : a reserve  $q_s[v]$  and a  $k$ -hop residue  $r_s^{(k)}[v]$
- Initially, sets  $r_s^{(0)}[s] = 1$ ; at  $k$ -th hop, converts  $\frac{\eta(k)}{\psi(k)} \times r_s^{(k)}[v] \rightarrow q_s[v]$ , pushes the remaining residue to neighbors evenly
- $\hat{\rho}_s[v]$  = Sum of reserves at  $v$



### A combination of HK-Push & Random Walks

$$\rho_s[v] = q_s[v] + \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v] \rightarrow \frac{\mathbb{E}[X]}{\#(\text{all random walks})}$$

$$\alpha = \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \quad \text{Total portion of random walks that not stopped yet}$$



- needs  $\alpha \omega$  random walks  $\rightarrow O(\alpha \omega t)$  time

### Optimization 1: Balancing HK-Push and random walks

- Max #hops  $K: \frac{1}{m/n}^{K/c} = \epsilon_r \cdot \delta$
- If  $\text{cost}(\text{HK-Push}) > 0.5 * \text{cost}(\text{random walks})$ , switch to random walks

### Optimization 2: Pruning random walks

- If  $r_s^{(k)}[u]$  becomes 0,  $\rho_s[v] = q_s[v] + \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v]$  meaning it's accurate and no need for random walks
- Otherwise, pruning  $\beta_k \cdot \epsilon_r \cdot \delta \cdot d(u)$  portion of random walks

$$\text{Time: } O(n_p) + O(\alpha \omega t) \rightarrow O(\frac{t \log(n/p_f)}{\epsilon_r^2 \delta}), \text{ Space: } O(m + n + \frac{t \log(n/p_f)}{\epsilon_r^2 \delta})$$

## 5. Experimental Results

Table 7: Statistics of graph datasets.

Dataset	$n$	$m$	$d$
DBLP	317,080	1,049,866	6.62
Youtube	1,134,890	2,987,624	5.27
PLC	2,000,000	9,999,961	9.99
Orkut	3,072,441	117,185,083	76.28
LiveJournal	3,997,962	34,681,189	17.35
3D-grid	9,938,375	29,676,450	5.97
Twitter	41,652,231	1,202,513,046	57.74
Friendster	65,608,366	1,806,067,135	55.06

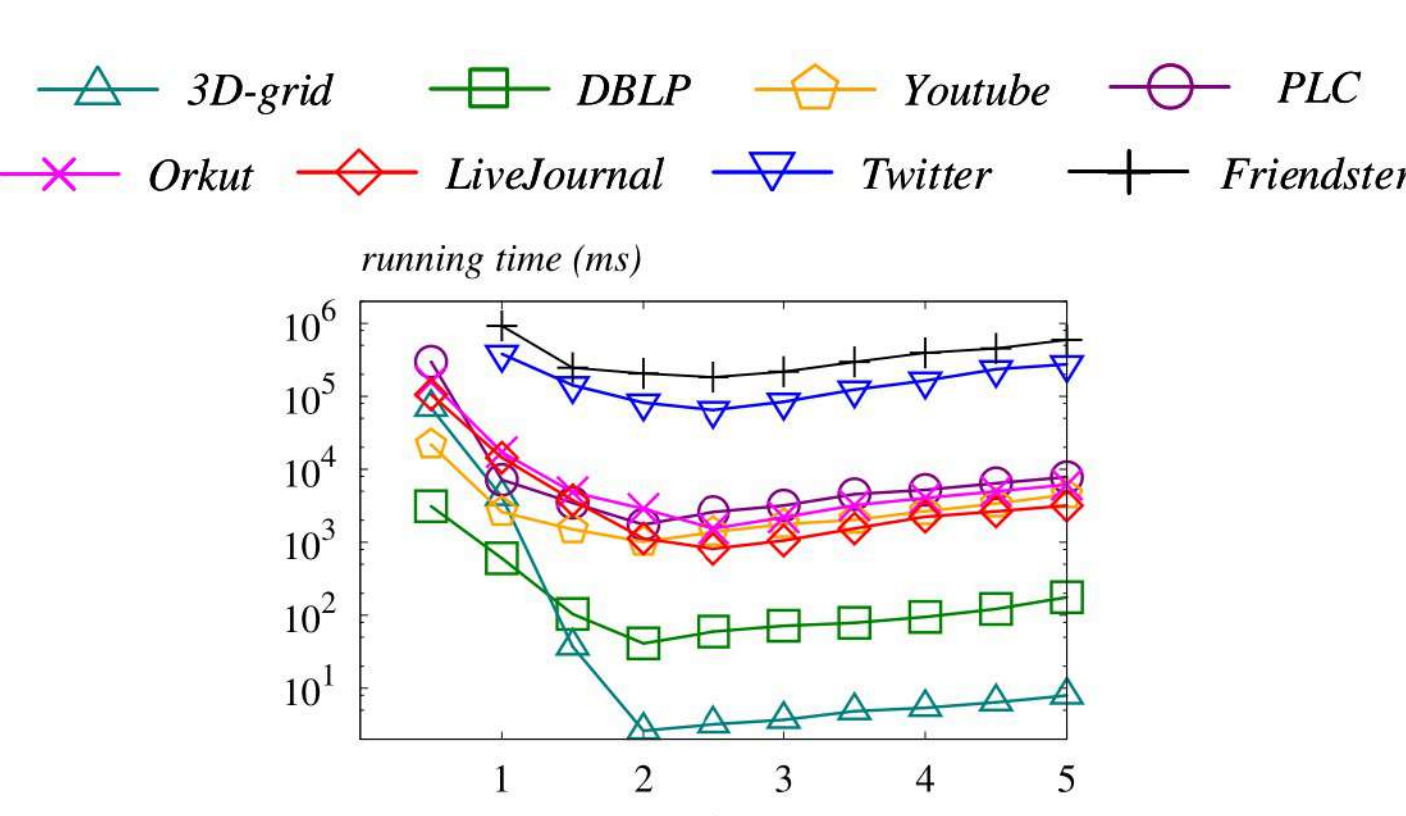


Figure 1: Running time of TEA+ vs c.

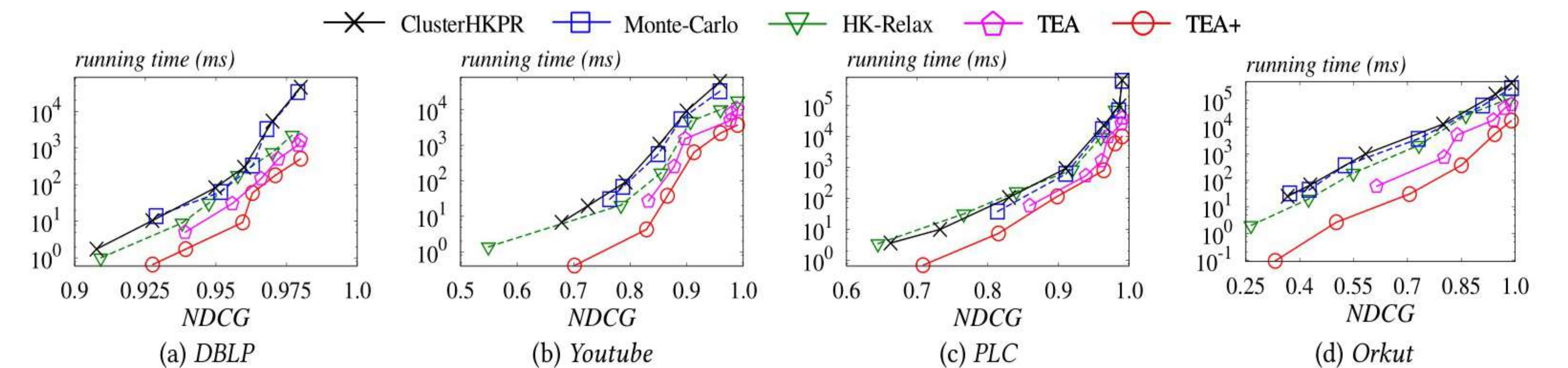


Figure 6: Running time vs. NDCG for computing normalized HKPR (best viewed in color).

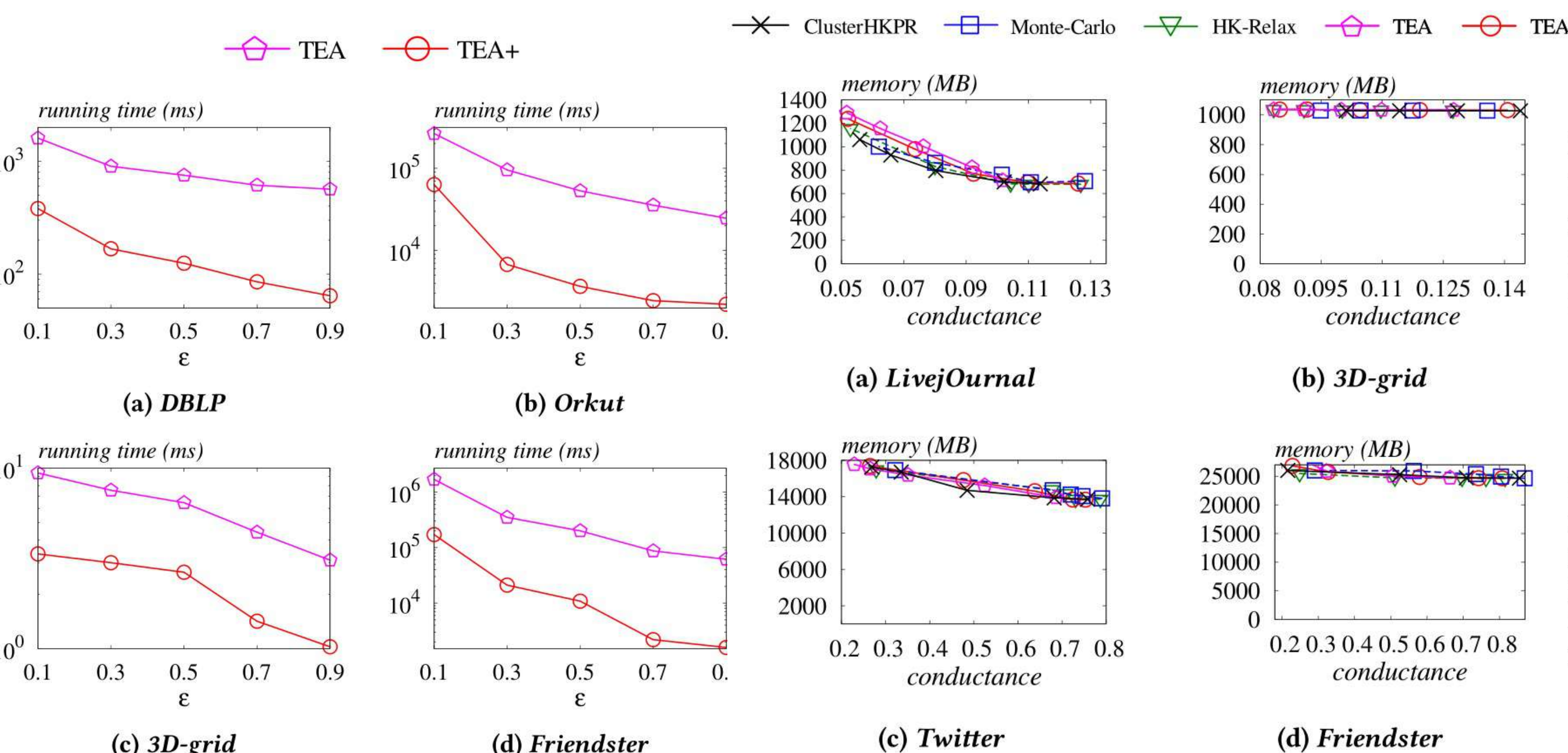


Figure 2: Running time vs  $\epsilon$ .

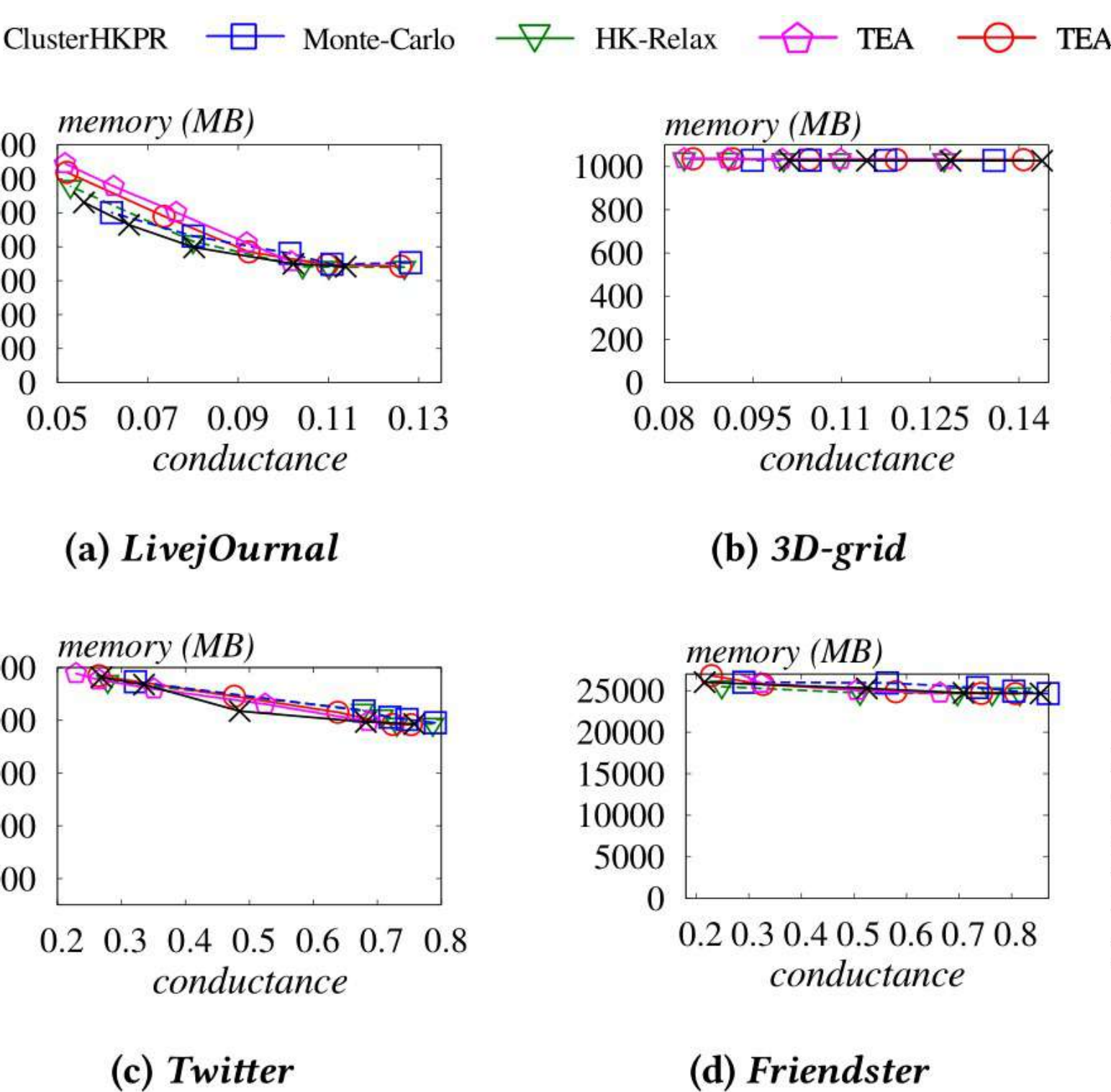


Figure 4: Memory cost vs. conductance.

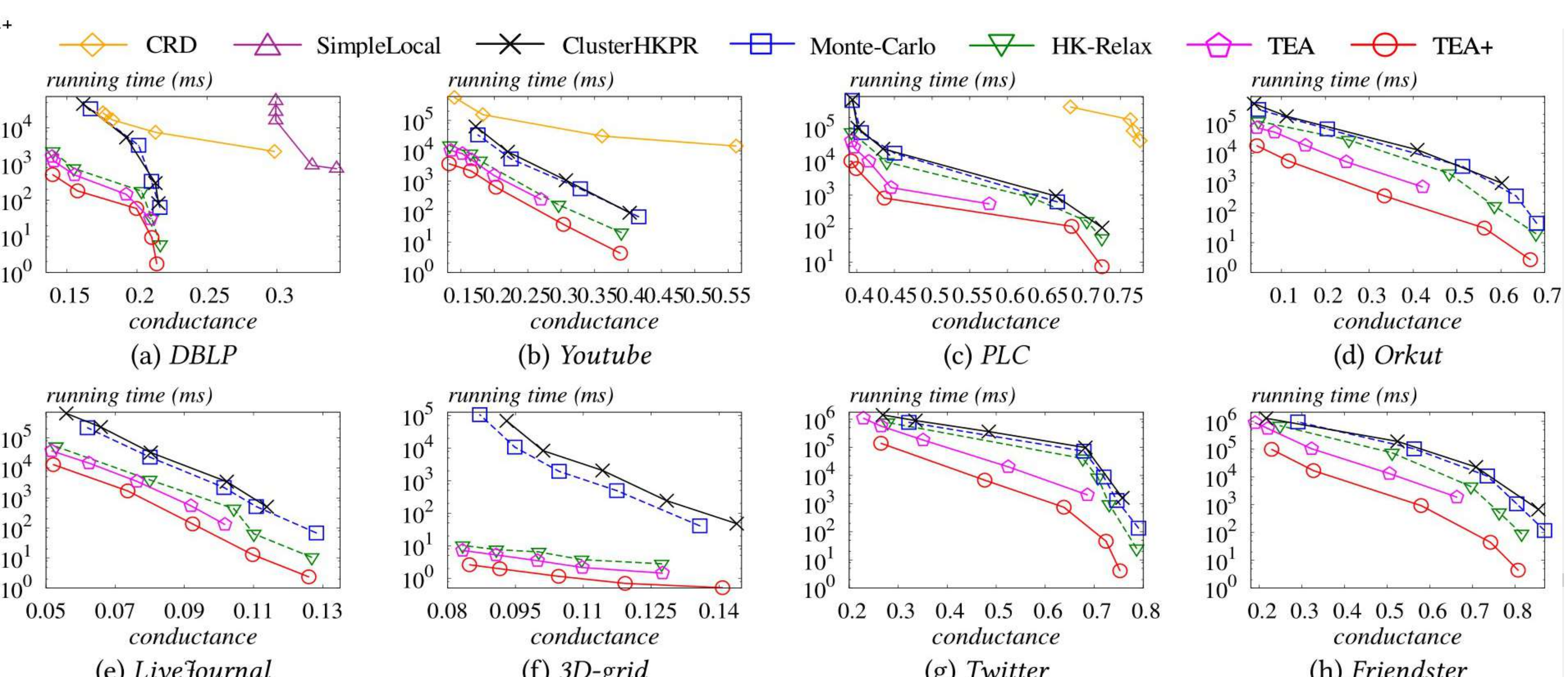


Figure 3: Running time vs conductance for local clustering queries (best viewed in color).