

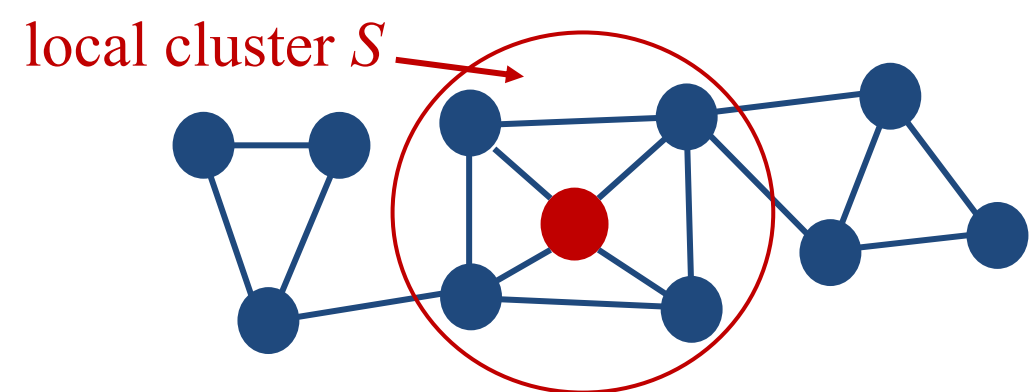
Efficient Estimation of Heat Kernel PageRank for Local Clustering

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

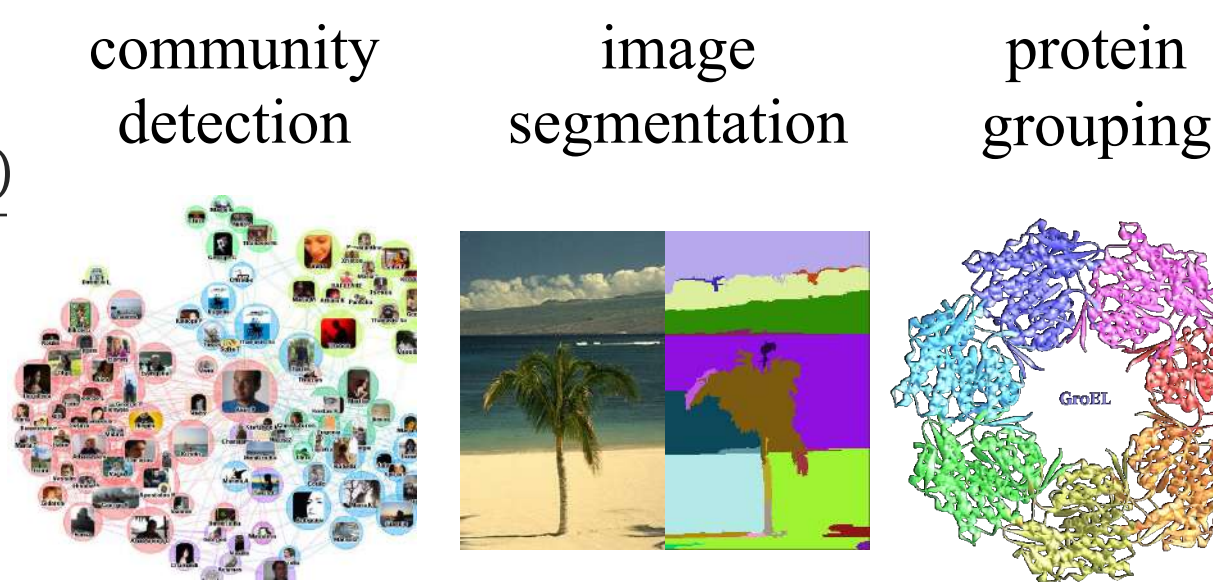
1. Heat Kernel-based Local Clustering

Local Clustering

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



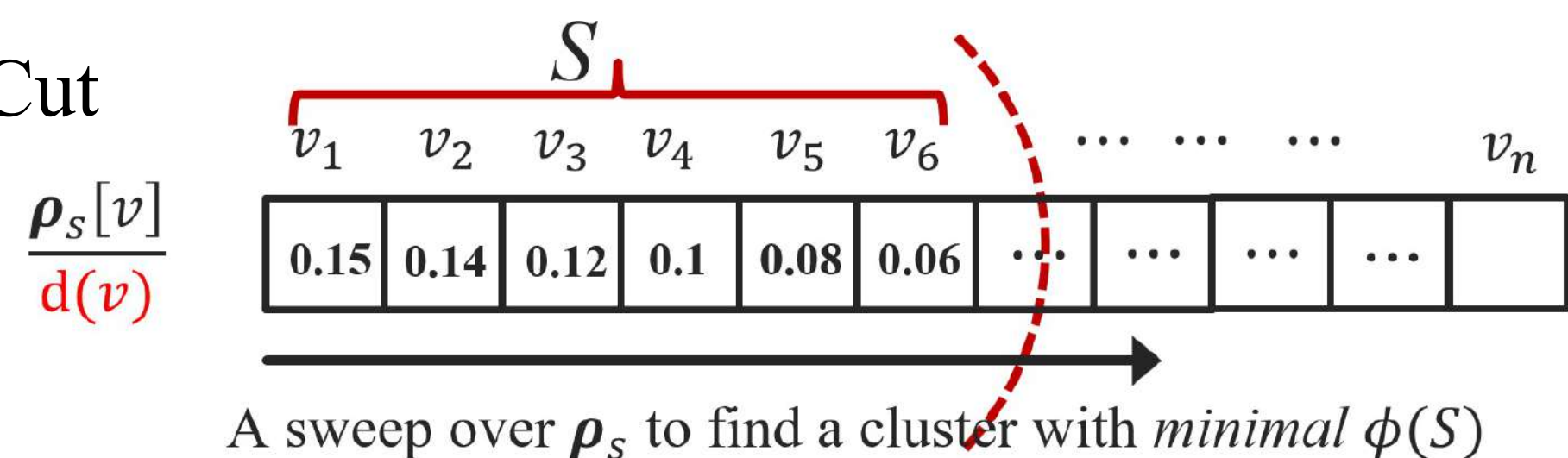
Applications



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



3. The Basic Ideas

(d, ϵ_r, δ) -approximate HKPR

- $\forall v \in G$ s.t. $\rho_s[v]/d(v) > \delta$,
- $\forall v \in G$ s.t. $\rho_s[v]/d(v) \leq \delta$,

$$\left| \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right| \leq \epsilon_r \cdot \frac{\rho_s[v]}{d(v)};$$

$$\left| \frac{\hat{\rho}_s[v]}{d(v)} - \frac{\rho_s[v]}{d(v)} \right| \leq \epsilon_r \cdot \delta.$$

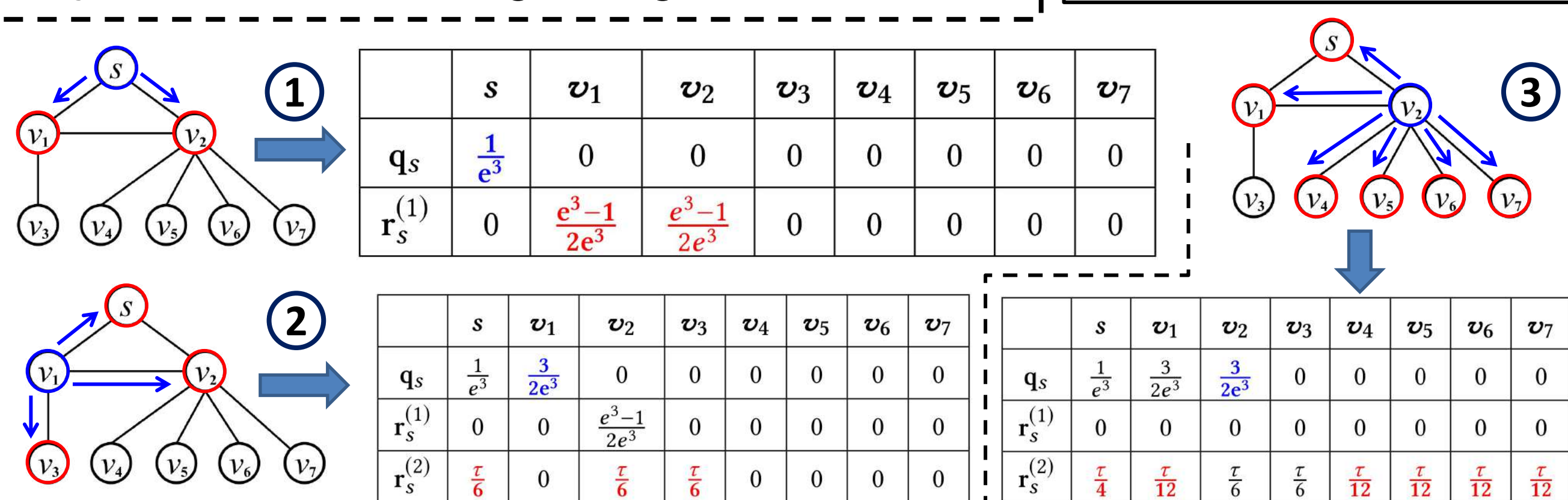
Monte-Carlo Random Walks

- Starting from s , at k -th step, with probability $\frac{\eta(k)}{\psi(k)}$, the random walk stops; OTRW, jumps to a random neighbor
- Performing $\omega = \frac{2(1+\epsilon_r/3)\log(n/p_f)}{\epsilon_r^2 \delta}$ random walks from s produces a (d, ϵ_r, δ) -approximate HKPR vector

HK-Push

- $\rho_s[v] = \text{estimation } q_s[v] + \text{error! } \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v]$
- $h_u^{(k)}[v] = \sum_{\ell=0}^{\infty} \frac{\eta(k+\ell)}{\varphi(k)} \cdot \mathbf{P}^\ell[u, v]$
- A traversal of G from s , $\frac{\eta(k)}{\psi(k)}$ portion of residue $r_s^{(k)}[v]$ to $q_s[v]$, remaining to neighbors' residues

- Each node v : a *reserve* $q_s[v]$, and *residue* $r_s^{(k)}[v]$ at k -th step
- It first set $r_s^{(0)}[s] = 1$ and stops pushing until $\forall u \in G, r_s^{(k)}[u]/d(u) \leq rmax$



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

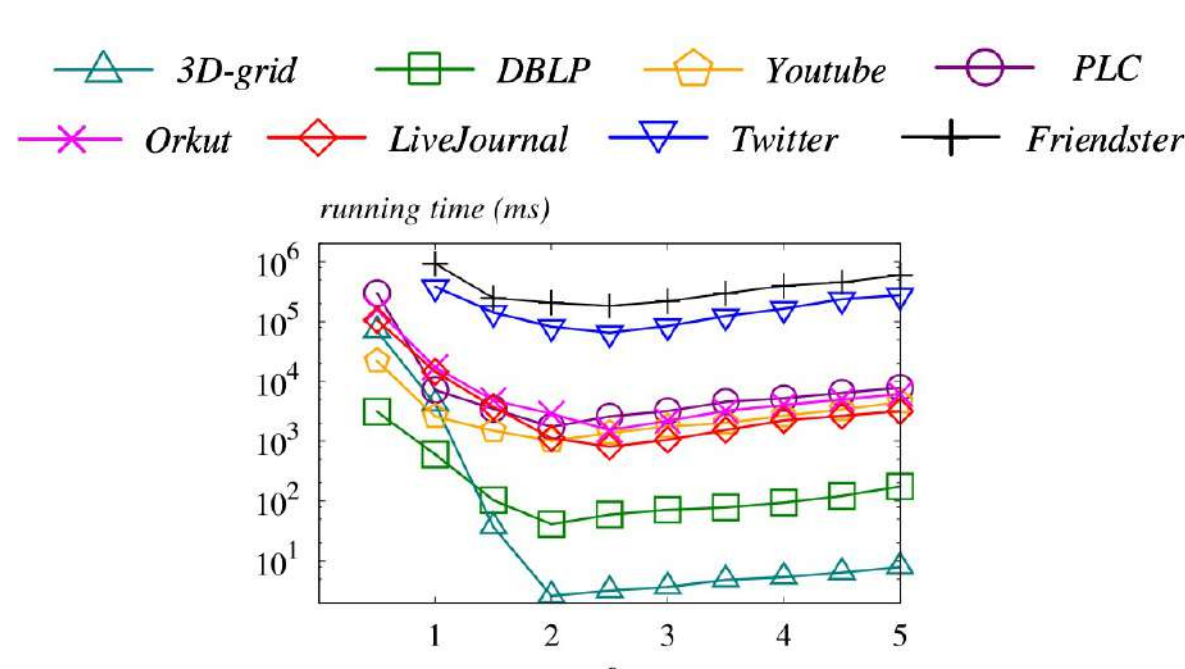


Figure 1: Running time of TEA+ vs c .

2. Existing Approximate Solutions

ClusterHKPR

- Samples a random walk length $k \leq O(\frac{\log(1/\epsilon)}{\log \log(1/\epsilon)})$
- Starts the k -step truncated random walk from s
- Repeats the above process for $16 \log n / \epsilon^3$ times

HK-Relax

- Injects initial residual ($r_s[s, 0] = e^{-t}$) to the seed node s
- At k -th step from s ($k \leq K = 2t \log(1/\epsilon_a)$), converts $r_s[v, k] \rightarrow \rho_s[v]$; distributes $\frac{t}{k+1} \cdot r_s[v, k]$ to neighbors evenly
- Stops until all $\frac{r_s[v, k]}{d(v)} \leq \frac{e^t \epsilon_a}{2K} / (\sum_{i=0}^{K-k} \frac{k!}{(k+i)!} t^i)$

Table 1: Theoretical guarantee of our solution against that of the state-of-the-art solutions.

Algorithm	Accuracy Guarantee	Time Complexity
ClusterHKPR [11]	with probability at least $1 - \epsilon$, $\begin{cases} \hat{\rho}_s[v] - \rho_s[v] \leq \epsilon \cdot \rho_s[v], & \text{if } \rho_s[v] > \epsilon \\ \hat{\rho}_s[v] - \rho_s[v] \leq \epsilon, & \text{otherwise,} \end{cases}$	$O\left(\frac{t \log(n)}{\epsilon^3}\right)$
HK-Relax [17]	$\frac{1}{d(v)} \hat{\rho}_s[v] - \rho_s[v] < \epsilon_a$	$O\left(\frac{t e^t \log(1/\epsilon_a)}{\epsilon_a}\right)$
Our solutions	with probability at least $1 - p_f$, $\begin{cases} \frac{1}{d(v)} \hat{\rho}_s[v] - \rho_s[v] \leq \epsilon_r \cdot \frac{\rho_s[v]}{d(v)}, & \text{if } \frac{\rho_s[v]}{d(v)} > \delta \\ \frac{1}{d(v)} \hat{\rho}_s[v] - \rho_s[v] \leq \epsilon_r \cdot \delta, & \text{otherwise,} \end{cases}$	$O\left(\frac{t \log(n/p_f)}{\epsilon_r^2 \cdot \delta}\right)$

4. The TEA+ Algorithm

A Combination of HK-Push and Random Walks

- $\rho_s[v] = q_s[v] + \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v] \rightarrow \mathbb{E}[X]$
- $\alpha = \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u]$
- Rolling a biased dice, with probability $r_s^{(k)}[u]/\alpha$, it shows v_5
- Starting from v_5 , at ℓ -th step, with probability $\frac{\eta(k+\ell)}{\psi(k+\ell)}$, the random walk stops, otherwise, jumps to a random neighbor
- If the random walk ends at v , $X = 1$, otherwise 0
- needs $\alpha \omega$ random walks $\rightarrow O(\alpha \omega t)$ time

Optimizations

- Balancing HK-Push and random walks via
 - Maximum #hops from the seed node s , $K = c \cdot \frac{\log(1/\delta/\epsilon_r)}{\log(d)}$
 - Maximum #push-operations $n_p = \omega \cdot t/2$
- Pruning random walks
 - $rb_s^{(k)}[u] = \min\{r_s^{(k)}[u], \beta_k \cdot \epsilon_r \delta \cdot d(u)\}$, $r_s^{(k)}[u] = r_s^{(k)}[u] - rb_s^{(k)}[u]$
 - $\rho_s[v] = q_s[v] + \sum_{u \in G} \sum_{k=0}^K rb_s^{(k)}[u] \cdot h_u^{(k)}[v] + \sum_{u \in G} \sum_{k=0}^K r_s^{(k)}[u] \cdot h_u^{(k)}[v]$

! NO random walks for this directly estimated as $0.5 \cdot \epsilon_r \delta \cdot d(v)$ with an error of $0.5 \cdot \epsilon_r \delta \cdot d(v)$ estimated by $\alpha \cdot \omega$ random walks with a much smaller α

Analysis

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

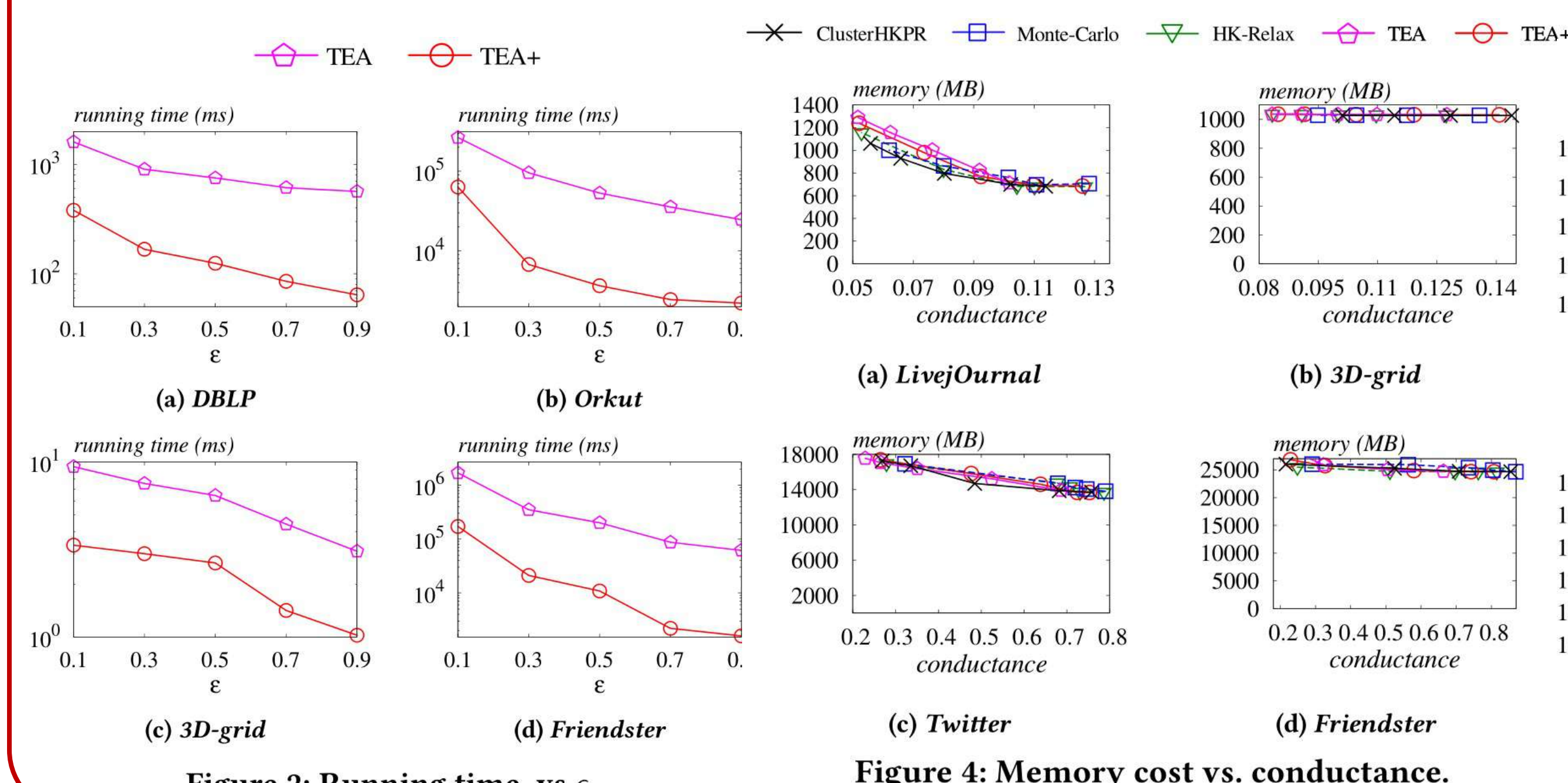


Figure 2: Running time vs ϵ_r .

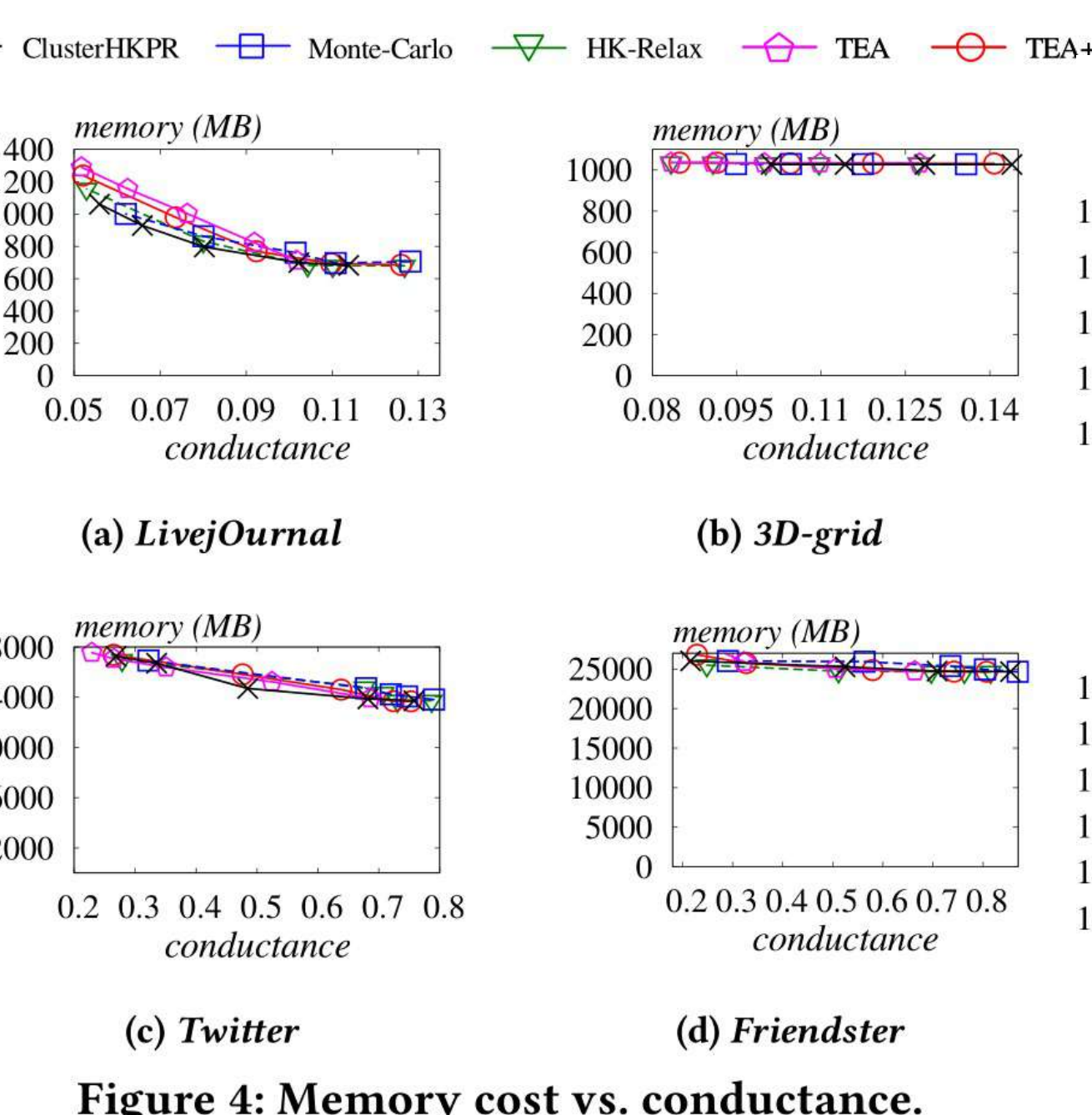


Figure 4: Memory cost vs. conductance.

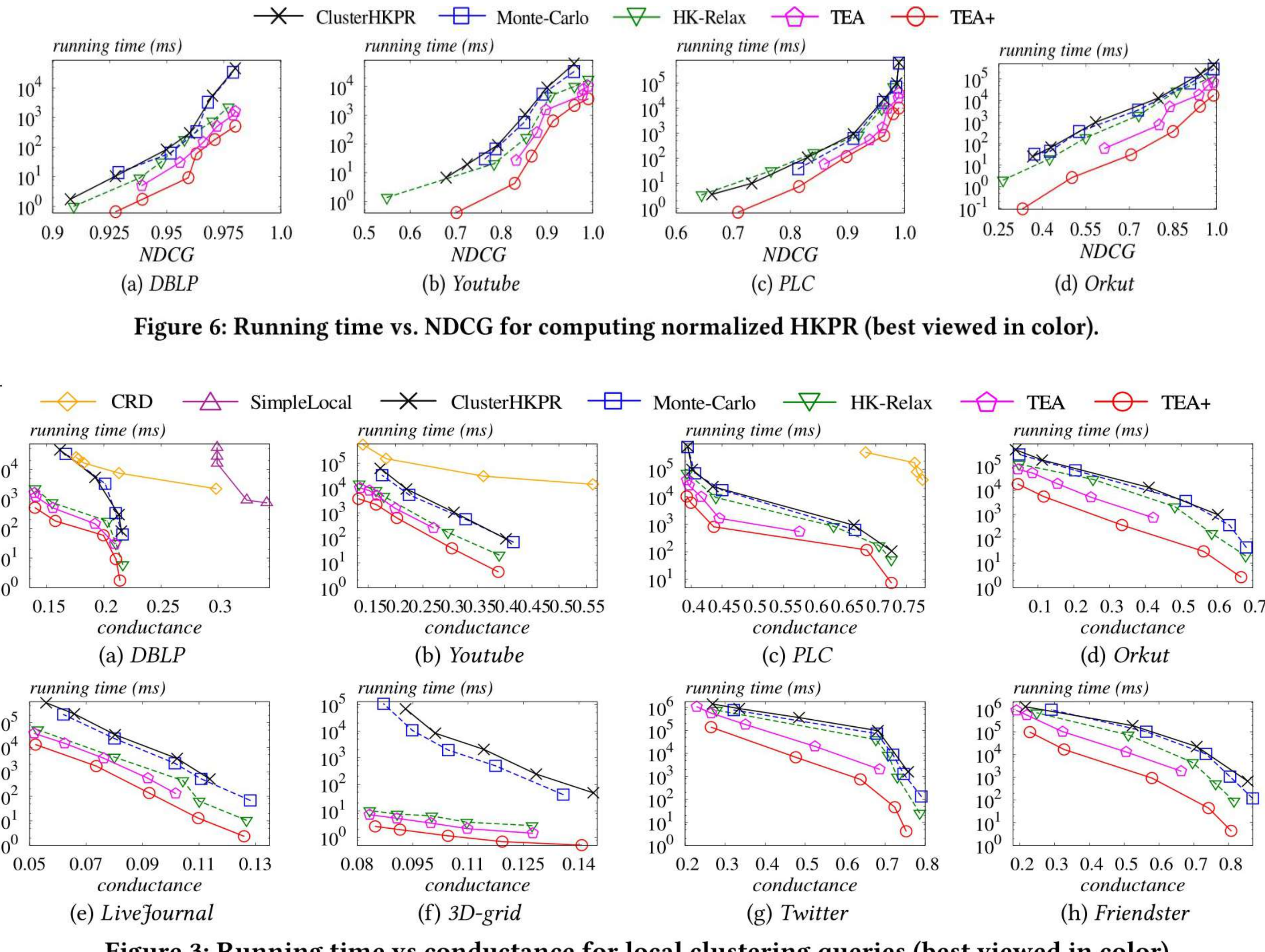


Figure 3: Running time vs conductance for local clustering queries (best viewed in color).

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