

# Computing the Heat Kernel of a Graph for a Local Clustering Algorithm

Olivia Simpson and Fan Chung

We present an efficient local clustering algorithm that finds cuts in large graphs by performing a sweep over a heat kernel vector. We show that for a subset  $S$  of Cheeger ratio  $\phi$ , many vertices in  $S$  may serve as seeds for a heat kernel random walk which will find a cut of conductance  $O(\sqrt{\phi})$ . Further, the random walk process is performed in time sublinear in the size of the graph.

## Local Clustering in Graphs

The goal of a local clustering algorithm is to compute a cluster in a graph near a specified vertex with size and volume constraints. For example, in **social networks** local clustering algorithms are used to identify a small and dense community around a particular member. In **protein networks** local clustering is used to isolate a group of interacting proteins to analyze a component of a biological system.

A good cluster,  $S$ , will be a subset of nodes in the graph with small *Cheeger ratio*:

$$\partial(S) = \{u \sim v : u \in S, v \notin S\}$$

$$\text{vol}(S) = \sum_{v \in S} d_v$$

$$\Phi(S) = \frac{|\partial(S)|}{\min(\text{vol}(S), \text{vol}(V \setminus S))}$$

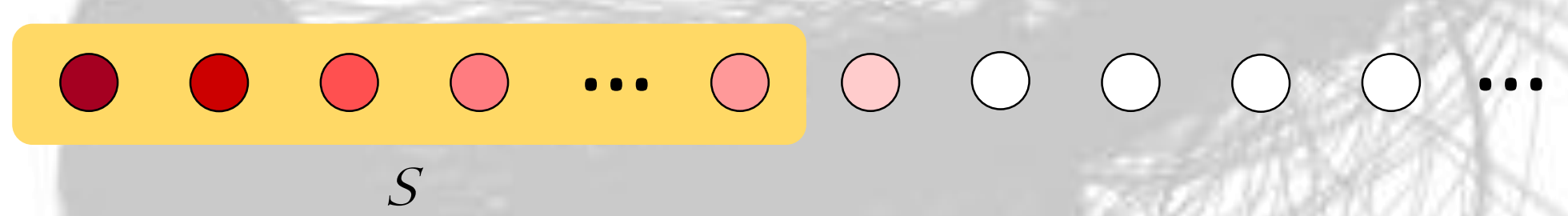
### Single Sweep Algorithms

Consider a probabilistic function  $f: V \rightarrow \mathbb{R}$  and order the vertices by decreasing probability-per-degree

$$\frac{f(v_1)}{d_{v_1}} \geq \frac{f(v_2)}{d_{v_2}} \geq \dots \geq \frac{f(v_n)}{d_{v_n}}$$

For  $i = 1$  to  $n$ :

- Let  $S$  be the set of the first  $i$  nodes in the ordering
- If  $\Phi(S), \text{vol}(S)$  are within desired constraints, output  $S$



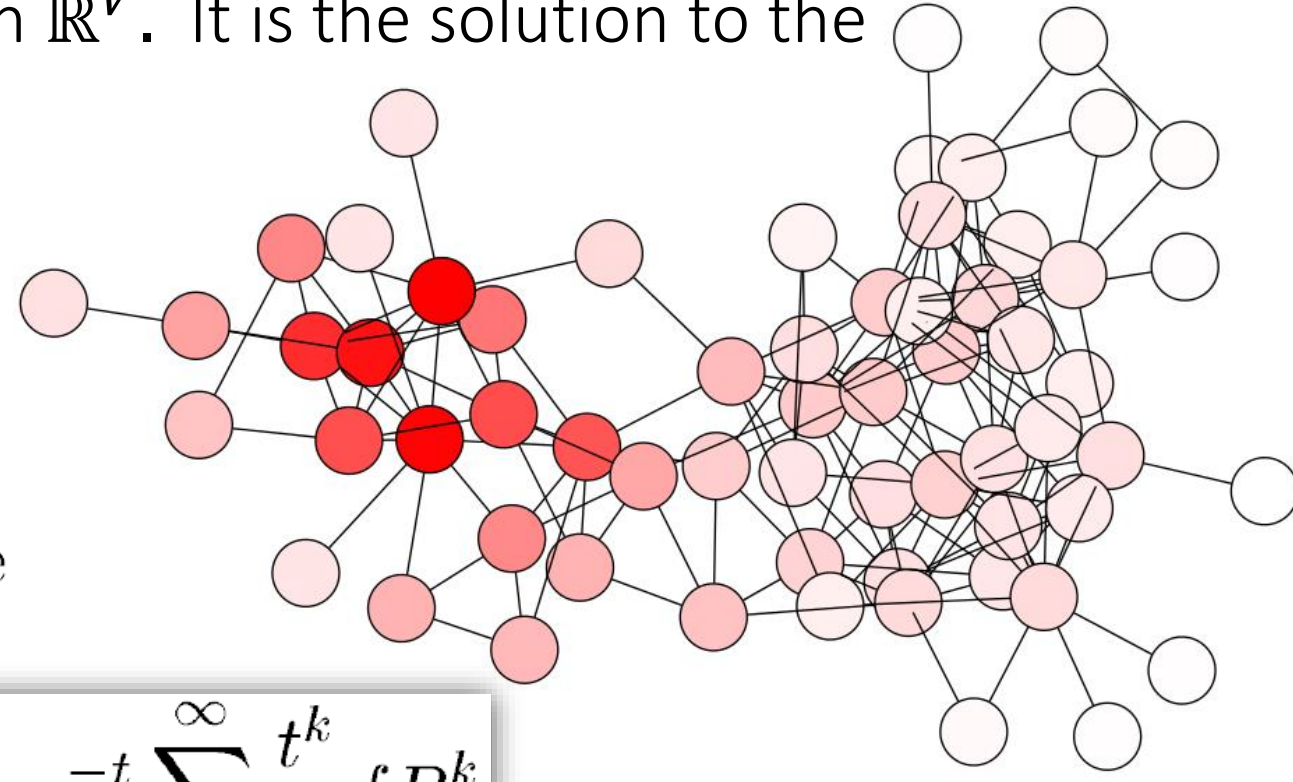
## The Heat Kernel of a Graph

The heat kernel is a vector in  $\mathbb{R}^V$ . It is the solution to the heat equation

$$\frac{\partial}{\partial t} \rho_{t,f} = -\rho_{t,f}(I - P)$$

$$(P)_{uv} = \begin{cases} 1/d_u & \text{if } u \sim v \\ 0 & \text{otherwise} \end{cases}$$

$$\rho_{t,f} = e^{-t} \sum_{k=0}^{\infty} \frac{t^k}{k!} f P^k$$



When viewed as a function, a single sweep of the heat kernel finds a cluster with Cheeger ratio  $O(\sqrt{\phi})$  where  $\phi$  is a constraint.

### Approximating the Heat Kernel in Sublinear Time

- Let  $f: V \rightarrow \mathbb{R}$  be a probabilistic function with all probability on a single vertex,  $u$
- Let  $X$  be the random variable that takes on the distribution after  $k$  random walk steps starting from  $u$  with probability  $p_k = e^{-t} \frac{t^k}{k!}$
- Then the expected value of  $X$  is exactly  $\rho_{t,f}$

#### ApproxHKPRseed( $G, t, u, \epsilon$ )

input: a graph  $G$ ,  $t \in \mathbb{R}^+$ , seed vertex  $u \in V$ , error parameter  $0 < \epsilon < 1$ .  
output:  $\rho$ , an  $\epsilon$ -approximation of  $\rho_{t,u}$ .

initialize a 0-vector  $\rho$  of dimension  $n$ , where  $n = |V|$   
 $r \leftarrow \frac{16}{\epsilon^3} \log n$   
 $K \leftarrow c \cdot \frac{\log(\epsilon^{-1})}{\log \log(\epsilon^{-1})}$  for some choice of constant  $c$

for  $r$  iterations do

**Start**

simulate a  $P$  random walk from vertex  $u$  where  $k$  steps are taken with probability  $e^{-t} \frac{t^k}{k!}$  and  $k \leq K$

let  $v$  be the last vertex visited in the walk

$\rho[v] \leftarrow \rho[v] + 1$

**End**

end for

return  $1/r \cdot \rho$

The algorithm returns an  $\epsilon$ -approximate vector of  $\rho_{t,f}$  in time

$$O\left(\frac{\log(\epsilon^{-1}) \log n}{\epsilon^3 \log \log(\epsilon^{-1})}\right)$$

**Definition 1.** Let  $G$  be a graph on  $n$  vertices, and let  $f: V \rightarrow \mathbb{R}$  be a vector over the vertices of  $G$ . Let  $\rho_{t,f}$  be the heat kernel pagerank vector according to  $f$  and  $t$ . Then we say that  $v \in \mathbb{R}^n$  is an  $\epsilon$ -approximate vector of  $\rho_{t,f}$  if

- for every vertex  $v \in V$  in the support of  $v$ ,

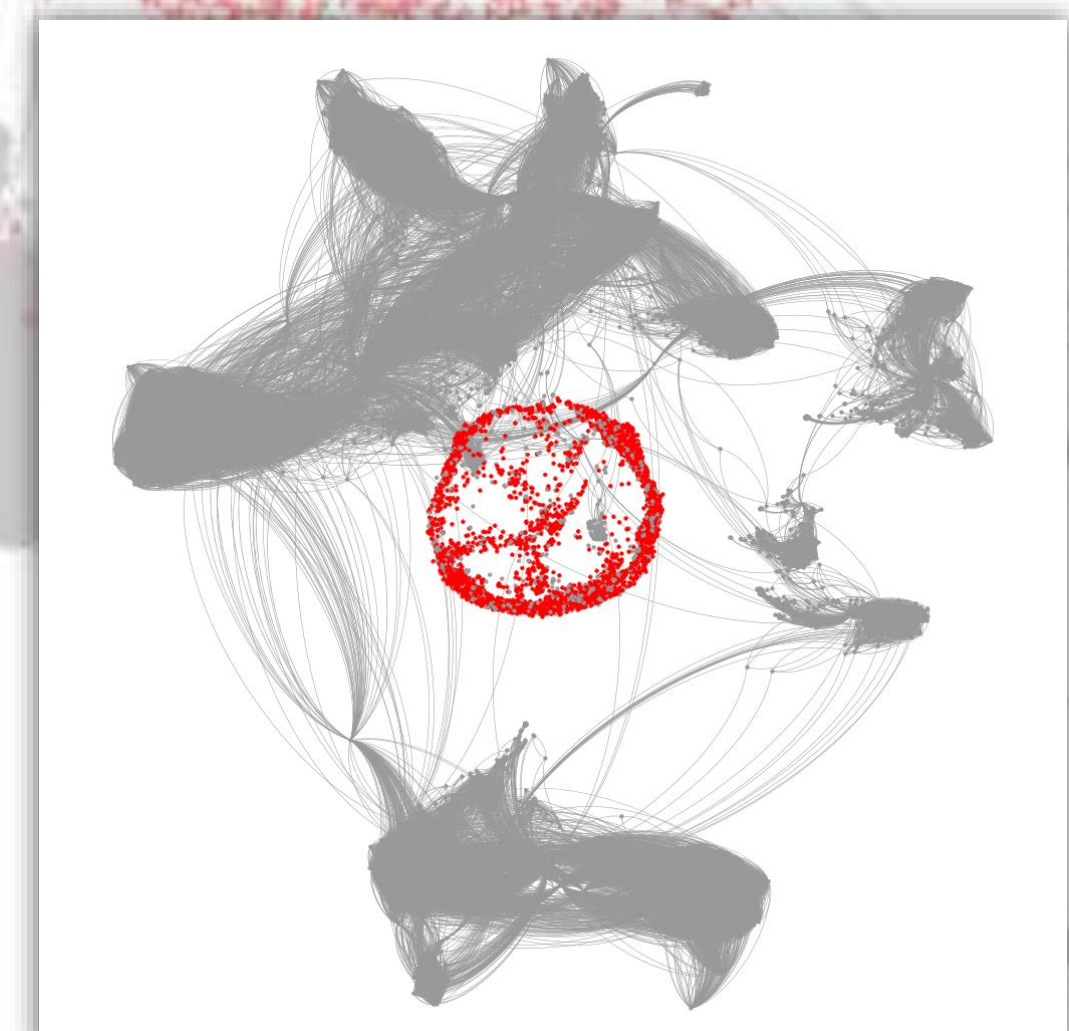
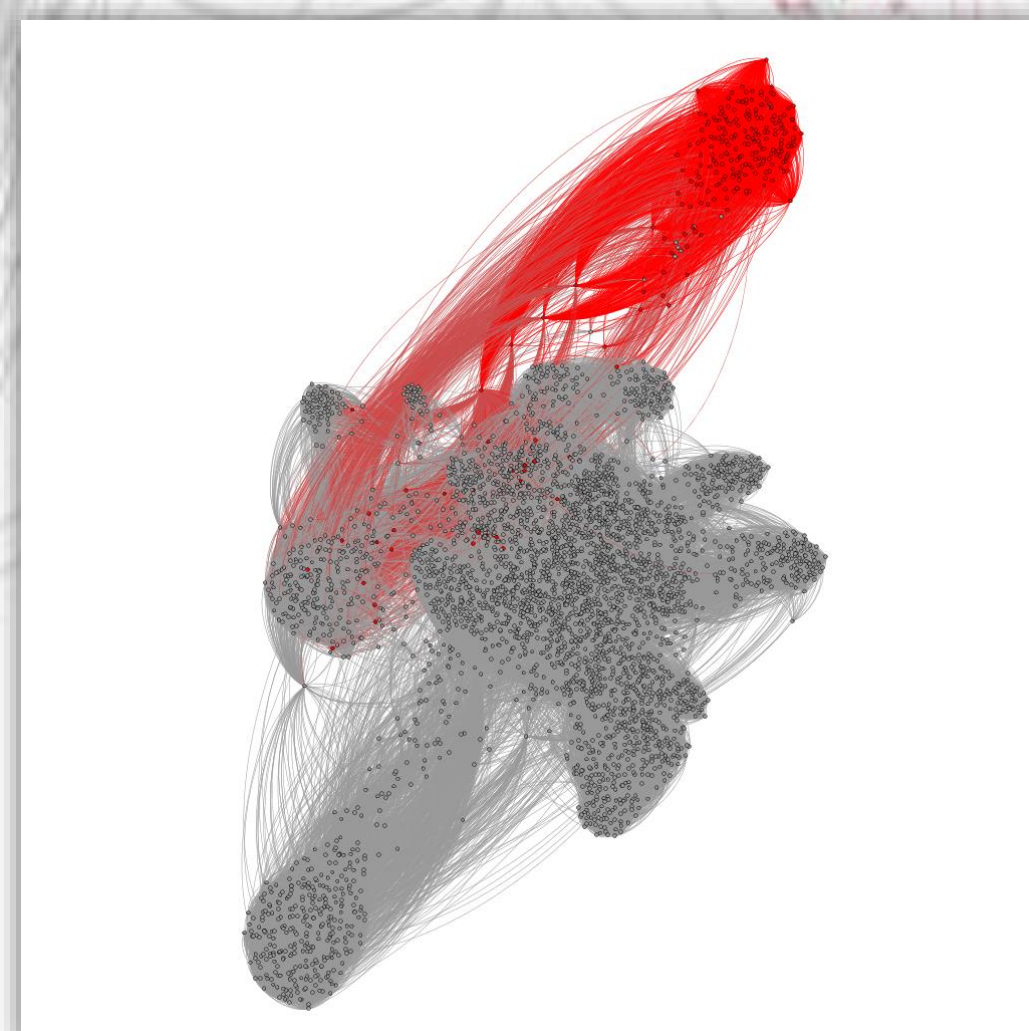
$$(1 - \epsilon)\rho_{t,f}(v) - \epsilon \leq v(v) \leq (1 + \epsilon)\rho_{t,f}(v),$$

- for every vertex with  $v(v) = 0$ , it must be that  $\rho_{t,f}(v) \leq \epsilon$ .

## Local Clusters with Heat Kernel Sweeps

Network	Size	Constraints	PR	HKPR	eHKPR
dolphins	$ V  = 62,$ $ E  = 159$	$\phi = 0.08$ $s = 20$ $\text{vol} = 100$	$\phi = 0.226$ $s = 23$ $\text{vol} = 106$	$\phi = 0.164$ $s = 24$ $\text{vol} = 110$	$\phi = 0.083$ $s = 20$ $\text{vol} = 96$
polbooks	$ V  = 105,$ $ E  = 441$	$\phi = 0.05$ $s = 30$ $\text{vol} = 270$	$\phi = 0.08$ $s = 48$ $\text{vol} = 415$	$\phi = 0.246$ $s = 49$ $\text{vol} = 403$	$\phi = 0.052$ $s = 50$ $\text{vol} = 422$
power	$ V  = 4941,$ $ E  = 6594$	$\phi = 0.05$ $s = 200$ $\text{vol} = 600$	$\phi = 0.375$ $s = 6$ $\text{vol} = 16$	$\phi = 0.003$ $s = 1564$ $\text{vol} = 4342$	$\phi = 0.347$ $s = 85$ $\text{vol} = 300$
facebook	$ V  = 4039,$ $ E  = 88234$	$\phi = 0.05$ $s = 200$ $\text{vol} = 2800$	$\phi = 0.419$ $s = 3063$ $\text{vol} = 88140$	$\phi = 0.001$ $s = 1094$ $\text{vol} = 67326$	$\phi = 0.057$ $s = 258$ $\text{vol} = 35266$
enron	$ V  = 36692,$ $ E  = 183831$	$\phi = 0.05$ $s = 100$ $\text{vol} = 1000$	$\phi = 0.488$ $\text{vol} = 183612$	-	$\phi = 0.037$ $\text{vol} = 3579$

The outputs of a single sweep algorithm using three different vectors: a Personalized PageRank vector (PR), an exact heat kernel vector (HKPR) and an  $\epsilon$ -approximate heat kernel vector (eHKPR).



Local clusters detected in the Facebook ego network (from the SNAP dataset). The image on the left is a local cluster detected with a sweep of an  $\epsilon$ -approximate heat kernel vector, while the image on the right is a local cluster detected with a sweep of a Personalized PageRank vector. Both clusters are colored red.