

Parametrized Power Iteration Clustering for Directed Graphs

Gwendal Debaussart-Joniec, Harry Sevi, Matthieu Jonckheere & Argyris Kalogeratos.

Centre Borelli · École Normale Supérieure Paris-Saclay · CNRS · Université Paris-Saclay | LAAS-CNRS

Abstract

Clustering directed graphs remains difficult because edge directionality breaks the assumptions behind classical spectral clustering. Existing methods often rely on expensive eigen-decomposition, graph symmetrization, or teleportation mechanisms that distort the original dynamics. As the directed nature breaks the reversibility of random walks, adaptation of diffusion-based methods to digraphs is non-trivial.

We propose Parametrized Power-Iteration Clustering **ParPIC**, an *eigen-free framework* that constructs a parametrized reversible random walk operator for any weakly connected digraph, automatically selects the diffusion time via an entropic criterion, and computes low-dimensional embeddings through random projections. On both synthetic and real-world directed graphs, ParPIC matches or outperforms 10 competing methods while scaling significantly better than spectral-based approaches.

Motivation

Classical diffusion geometry and spectral clustering assume:

- Reversible random walks
- Real-valued spectra
- Strong connectivity

These assumptions fail for many directed graphs:

- Random walks are often non-reversible
- Spectra may become complex-valued
- Weak connectivity breaks convergence guarantees

Existing workarounds introduce limitations:

Workarounds	Limitations
Symmetrization	Loose directional information
Teleportation	Distortion of original dynamics, loss of sparsity
Hermitianization	No diffusion interpretation

Goal. Enabling diffusion-based clustering for directed graphs without eigen-decomposition, while preserving edge directionality and scalability.

Notations

$G = (V, E)$: digraph	k : number of clusters
A : adjacency matrix	N : number of vertices
d^+, d^- : out-degree and in-degree	d : embedding dimension
$D_{ii} = d^+(i)$: out-degree matrix	$\nu: V \rightarrow \mathbb{R}, \ \nu\ _1 = 1$: vertex measure
$P = D^{-1}A$: natural random walk	D_ν : diagonal matrix of ν
$A_S = (A + A^T)/2$: symmetrized A	$P_S = (D_S)^{-1}A_S$: random walk of A_S

Method

Parametrized random walk operator for directed graphs:

$$P_\nu = (D_\xi + D_\nu)^{-1}(D_\nu P + P^T D_\nu), \quad \xi = P^T \nu.$$

Key properties. under mild assumptions:

- Reversible with respect to ν
- Converges to a unique stationary distribution π_ν
- Real-valued spectrum
- Continuous over ν

Vertex measure design

$$\nu(i) = \gamma d^+(i) + (1 - \alpha) d^-(i), \quad \gamma \in [0, 1].$$

It is assumed that d^+ and d^- are normalized to sum to 1. Experiments show that $\gamma \approx 0.5$ is a good default choice and that extreme values of γ can lead to poor clustering performance.

Diffusion time selection

We propose an entropy-based criterion to automatically select the diffusion time t that balances local and global structure in the embedding space, avoiding the need for manual tuning.

$$\mathcal{H}(P) = \sum_i \mathcal{H}_i(P), \quad \mathcal{H}_i(P) = - \sum_j P(i, j) \log P(i, j),$$

$$t^* = \underset{t \in \mathbb{N}}{\text{Elbow}} \mathcal{H}(P_\nu^t)$$

Randomized embeddings

To scale to large graphs, we use randomized methods to compute low-dimensional embeddings based on the diffusion operator *without explicit eigen-decomposition*.

$$Z_t = P_\nu Z_{t-1}, \quad Z_0 \sim \frac{1}{\sqrt{N}} \mathcal{N}(0, I), \quad Z_t \in \mathbb{R}^{N \times d},$$

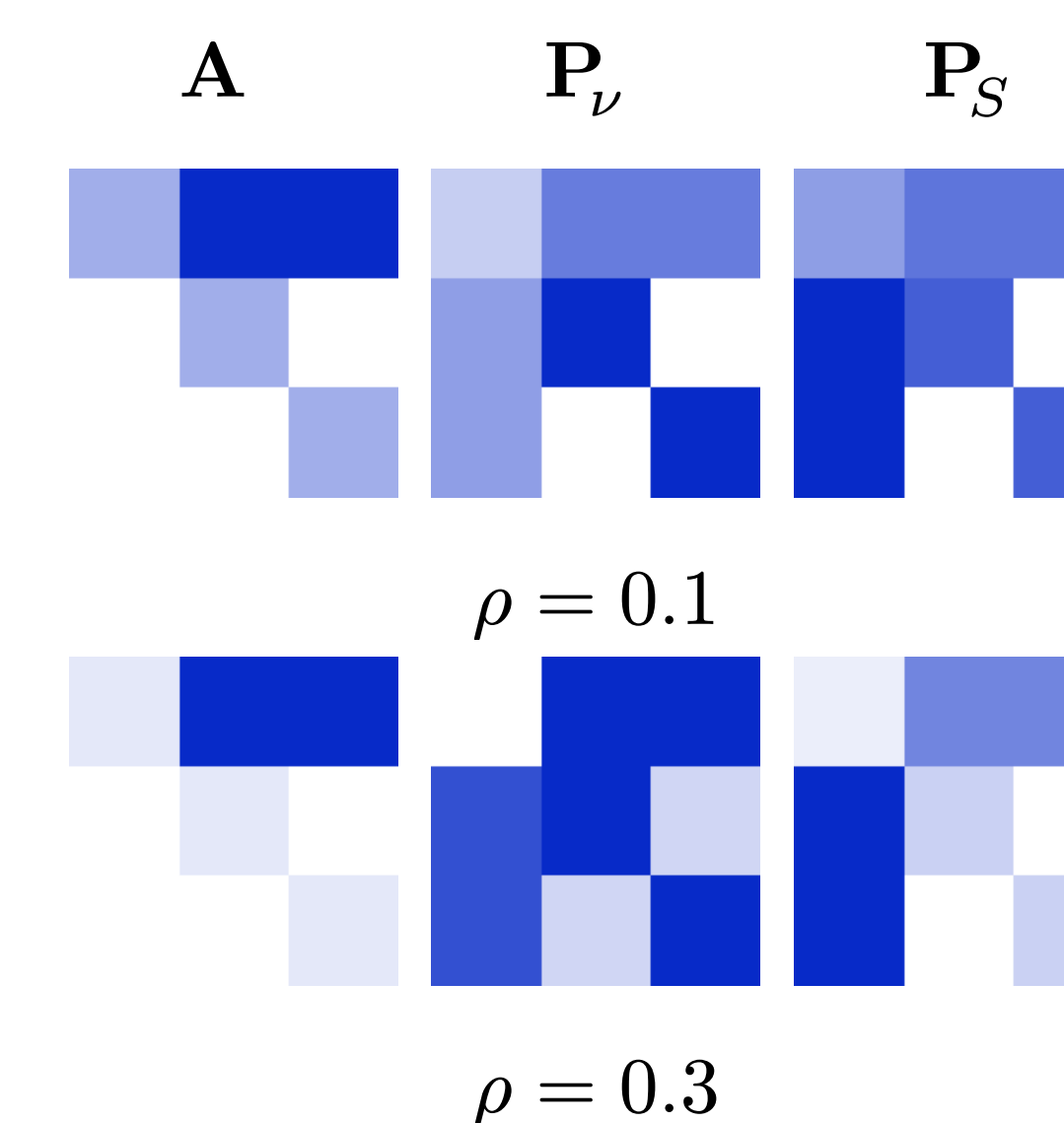
where $d = \sqrt{N}$ is the embedding dimension. For larger graphs, it can be further reduced to $d = c \log(N)$ with minimal loss in clustering performance, with c being a ‘small’ constant.

Pipeline

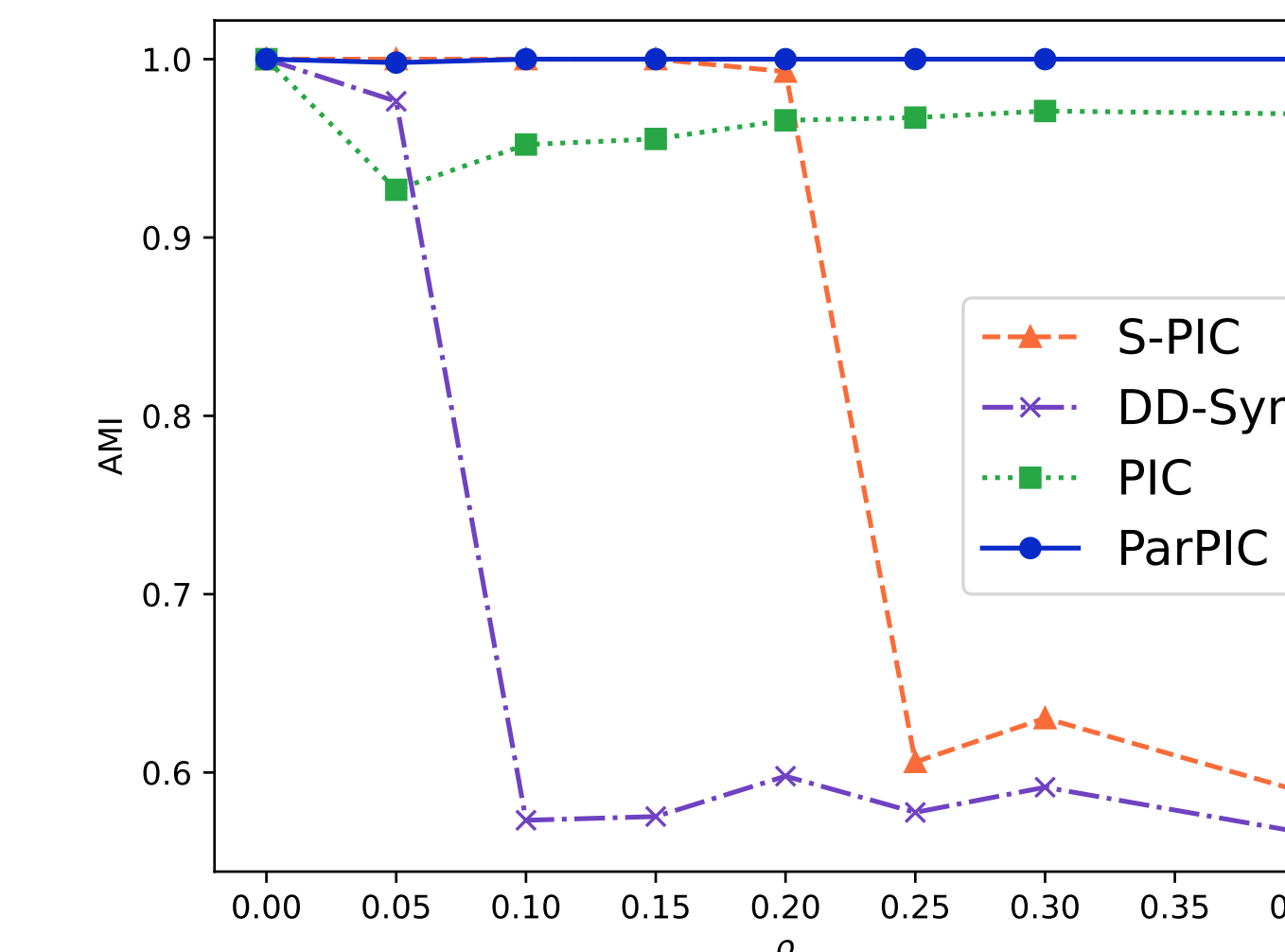
1. Compute P from the directed graph
2. Compute P_ν from P and ν
3. Select diffusion time t using entropy-based criterion
4. Compute low-dimensional embedding
5. Cluster using k-means++

Experiments

Synthetic graphs



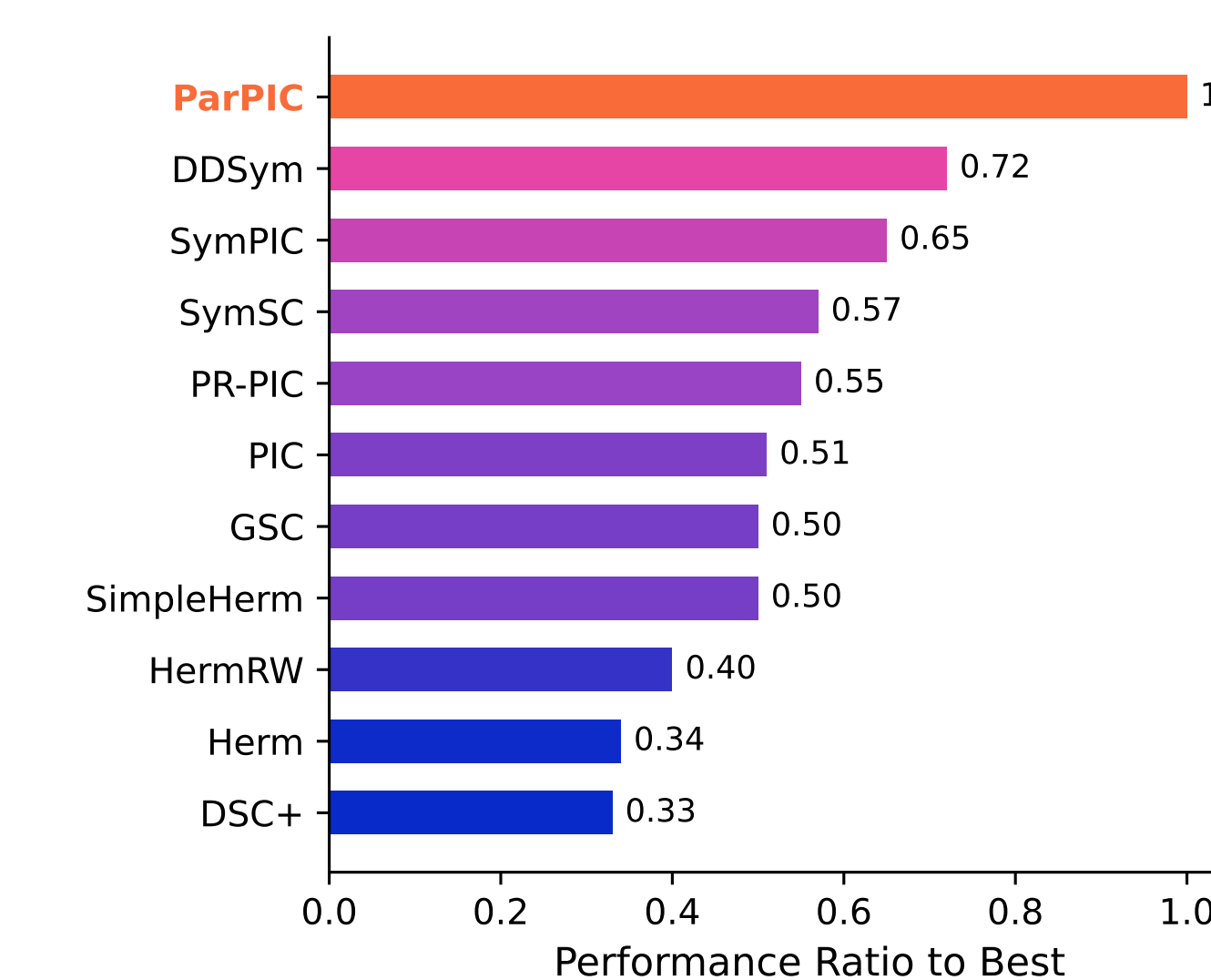
Expected matrices for a DiSBM.



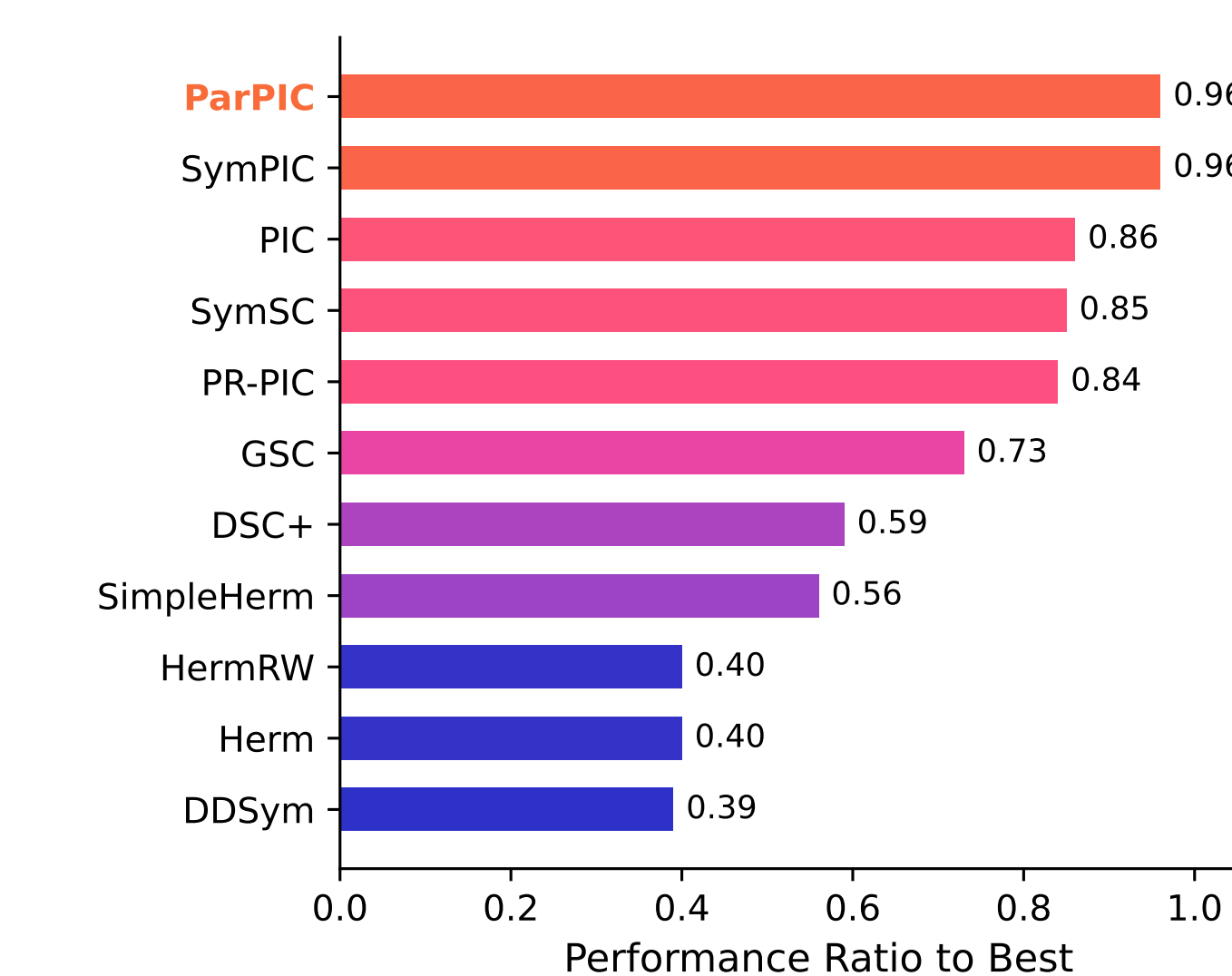
Sensitivity to out-degree heterogeneity.

The probability of an edge between the first cluster and the two other clusters is denoted ρ .

Clustering performance



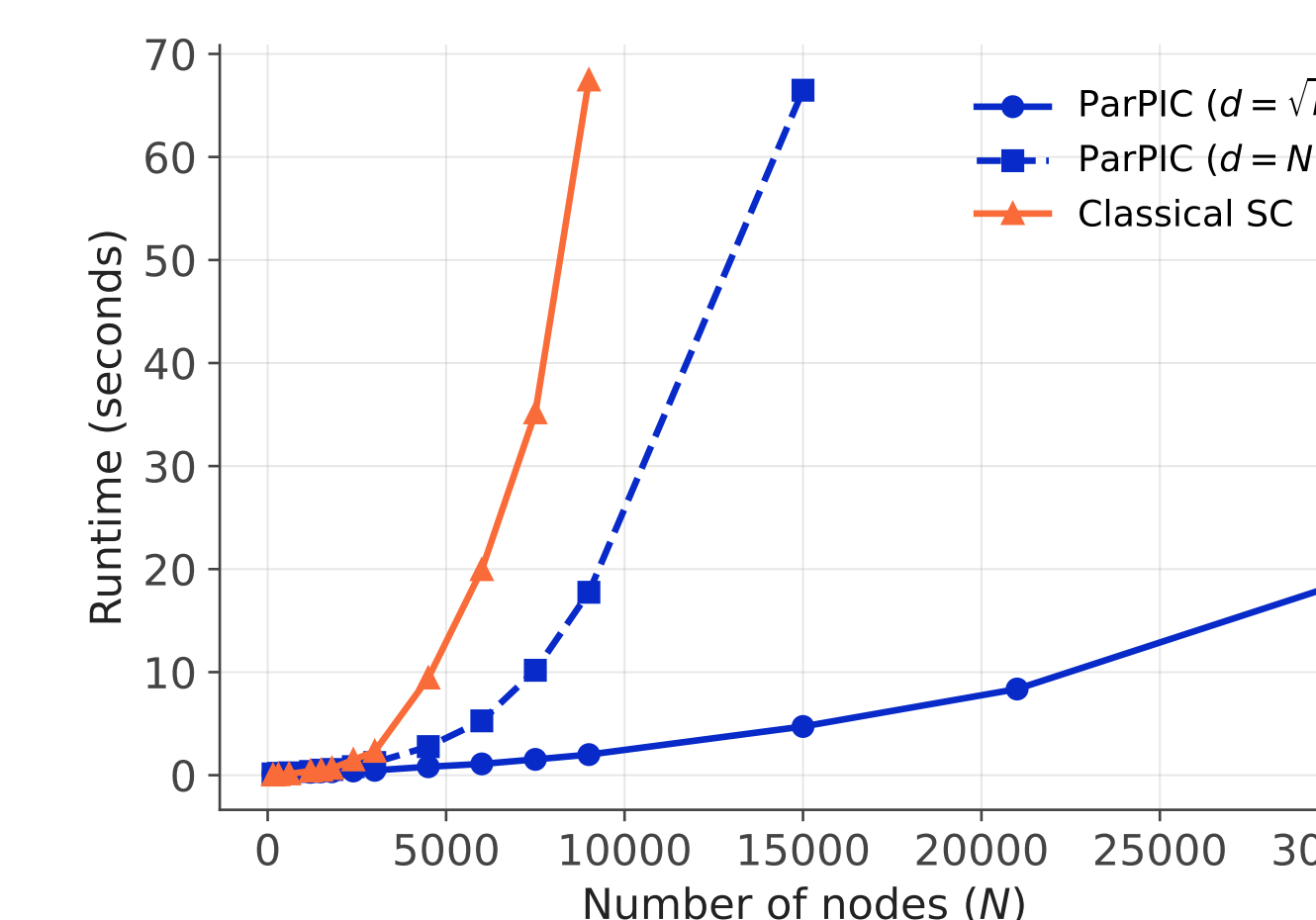
Directed graphs



k -NN directed graphs

Method families — Hermitian spectral: Herm, HermRW, SimpleHerm | RW-spectral: DDSym, DSC+, GSC, SymSC | Power-Iteration: PIC, PR-PIC, S-PIC, ParPIC. Performance Ratio to Best is defined as the ratio of the performance of each method to the best performance among all methods. The mean is shown and computed with respect to AMI.

Runtime compared to spectral clustering



Runtimes on synthetic graphs.

Take-aways. ParPIC provides a scalable and performant, diffusion-based clustering framework by extending PIC to directed graphs.



Interested ?
Check out our paper!