# Community Detection in networks by Dynamical Optimal Transport Formulation

**Daniela Leite**[1,*,+]**, Diego Baptista**[1,*]**, Abdullahi Ibrahim**[1]**, Enrico Facca**[2]**, and Caterina De Bacco**[1]

[1]Max Planck Institute for Intelligent Systems, Cyber Valley, 72076 Tübingen, Germany
[2]Univ. Lille, Inria, CNRS, UMR 8524 - Laboratoire Paul Painlevé, F-59000 Lille, France
[+]daniela.leite@tuebingen.mpg.de
[*]These authors contributed equally to this work

## ABSTRACT

Detecting communities in networks is important in various domains of applications. While a variety of methods exists to perform this task, recent efforts propose Optimal Transport (OT) principles combined with the geometric notion of Ollivier-Ricci curvature to classify nodes into groups by rigorously comparing the information encoded into nodes' neighborhoods. We present an OT-based approach that exploits recent advances in OT theory to allow tuning for traffic penalization, which enforces different transportation schemes. As a result, our model can flexibly capture different scenarios and thus increase performance accuracy in recovering communities, compared to standard OT-based formulations. We test the performance of our algorithm in both synthetic and real networks, achieving a comparable or better performance than other OT-based methods in the former case, while finding communities more aligned with node metadata in real data. This pushes further our understanding of geometric approaches in their ability to capture patterns in complex networks.

## 1 Introduction

Complex networks are ubiquitous, hence modeling interactions between pairs of individuals is a relevant problem in many disciplines[1,2]. Among the variety of analysis that can be performed on them, community detection[3–6] is a popular application that involves finding groups (or communities) of nodes that share similar properties. The detected communities may reveal important functional properties of the underlying system. Community detection has been used in diverse areas including, discovering potential friends on social networks[7], evaluating social networks[8], personalized recommendation of item to user[9], detecting potential terrorist activities on social platforms[10], fraud detection in finance[11], study epidemic spreading process[12] and so on.

Several algorithms have been proposed to tackle this problem which utilize different approaches, such as statistical inference[13,14], graph modularity[15], statistical physics[16] and information theory[17]. Here, instead, we adopt a recent approach connecting community detection with geometry, where communities are detected using geometric methods like the Ollivier-Ricci curvature (ORC) and we exploit optimal transport theory to calculate this efficiently.

In Riemannian geometry, a curvature quantifies how geodesic paths converge or diverge, depending on the curvature's sign. In networks, the ORC plays a similar role where edges with negative curvature are traffic bottlenecks, in terms of network flow of the shortest paths. In the opposite case, positively curved edges contribute to transport on the network along with several others, reflecting the fact that they are well connected. Defining communities as robust transport of information along with the network, we could cluster edges based on their curvature: those with positive curvature can be clustered together, while those with negative curvature may be seen as "bridges" connecting different communities. The idea of using Ricci curvature to find communities on networks has been recently proposed in[18,19]. In this work we follow a similar approach, but generalize it for the case of branched[20,21] and congested[22] optimal transport problems, building from recent results[23]. Specifically, our algorithm allows us to efficiently tune the sensitivity to detecting communities in a network, by means of a parameter that controls the flow of information shared between nodes. We perform a comprehensive comparison between the proposed algorithm and existing ones on synthetic and real data. Our algorithm, named ORC-Nextrout, detects communities in synthetic networks with similar or higher accuracy compared to other OT-based methods in the regime where inference is not trivial. This is also observed in a variety of real networks, where the ability to tune between different transportation regimes allows finding at least one result that outperforms other methods, including approaches based on statistical inference and modularity-based community detection.

**Related work.**

The idea of exploring geometrical properties of a graph, and in particular curvature, has been explored in different branches of network science, ranging from biological[24] to communication[25] networks. Intuitively, the Ricci curvature can be seen as the amount of volume through which a geodesic ball in a curved Riemannian manifold deviates to the standard ball in Euclidean spaces[26]. When defined in graphs, it indicates whether edges (those with positive values for the curvature) connect nodes inside a cluster, or if they rather bond different clusters together (those with negative values for the curvature).

Two main discrete graph curvature approaches have been proposed: the Ollivier-Ricci (OR) curvature based on the optimal transport theory introduced by Ollivier,[27,28] and Forman-Ricci curvature introduced by Forman[29]. While the graph Laplacian-based Forman curvature is computationally fast and less geometrical, we focus on OT-based approach due to its more geometric nature. Some applications of the Ollivier-Ricci curvature include network alignment[30] and community detection[18,19,31].

On the other hand, community detection in networks is a fundamental area of network science, with a wide range of approaches proposed for this task[3,4,32]. Our work is inspired by recent OT-based methods[18,19] for community detection. These methods consider the OR curvature to sequentially identify and prune negatively curved edges from a network to identify communities. While our approach also considers OR curvature to prune edges, it controls the flow of information exchanged between nodes by means of parameter, making the edge pruning dynamic. This is detailed in Sec. 2.

## 2 $\beta$-Wasserstein Community Detection Algorithm

In this section, we describe how our approach solves the community detection problem. As previously stated, we rely on optimal transport principles to find the communities. To solve the optimal transport problem applied in our analysis we use the discrete *Dynamic Monge-Kantorovich* model (*DMK*), as proposed by Facca et al.[33,34] to solve transportation problems on networks.

We denote a weighted undirected graph as $G = (V, E, W)$, where $V, E, W$ are the set of nodes, edges and weights, respectively. We use the information of a node neighborhood $\mathcal{N}(i) = \{j \in V | (i,j) \in E\}$ to decide whether node $i$ belongs to a given community. We do this by comparing a distribution defined on $\mathcal{N}(i)$ to the ones defined on other nodes close to $i$. This distribution is defined as $m_i^{\alpha}$, where $m_i^{\alpha}(k) := \alpha$ if $k = i$ and $m_i^{\alpha}(k) := (1-\alpha)/|\mathcal{N}(i)|$ if $k \in \mathcal{N}(i)$. Intuitively, the distribution $m$ assigns a unit of mass to $i$ and its connections: $\alpha$ controls how much weight the node $i$ should have, and once this is assigned, its neighbors receive the remaining mass in an even way.
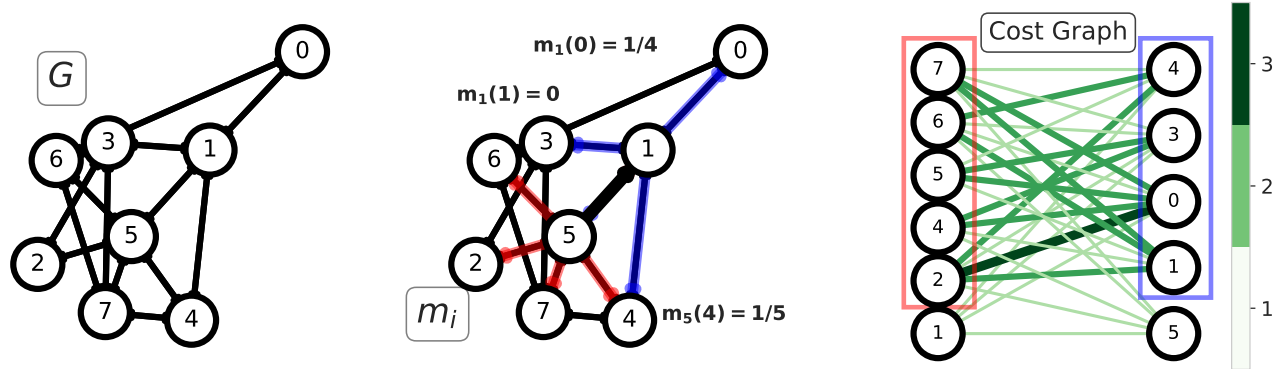
The next step is to compare the distribution $m_i^{\alpha}$ of the node $i$ to that of its neighbors. Consider an edge $(i,j) \in E$ and $m_j$, the distribution defined on the node $j$, neighbor of $i$. We assume that if $i$ and $j$ belong to the same community, then both nodes may have several neighbors in common, and therefore, $m_i$ and $m_j$ should be similar. Notice that this implies an assortativity assumption, where nodes within the same community are more likely to interact than nodes in different communities[2,4]. This has been observed for instance in social or biological networks[13,35]. On the contrary, it may not be appropriate to model disassortative datasets, where nodes tend to connect more often across communities.

To estimate the similarity between $m_i$ and $m_j$ we use OT principles. Specifically, we compute the cost of transforming one distribution into the other. This is related to the cost of moving the mass from one neighborhood to the other, and it is assumed to be the weighted shortest-path distance between nodes belonging to $\mathcal{N}(i)$ and $\mathcal{N}(j)$. A schematic representation of the algorithm can be seen in Fig. 1. The OT problem is solved in an auxiliary graph, the complete bipartite network $B_{ij} = (V_{ij}, E_{ij}, \omega_{ij})$ where $V_{ij} := (V_i, V_j) := (\mathcal{N}(i) \cup \{i\}, \mathcal{N}(j) \cup \{j\})$, $E_{ij}$ is made of all the possible edges between $V_i$ and $V_j$. The weights of the edges are given by the weighted shortest path distance $d$ between two nodes measured on the input network $G$.

The similarity between $m_i$ and $m_j$ is the Wasserstein cost $\mathcal{W}(m_i, m_j, \omega_{ij})$ of the solution of the transportation problem. In its standard version, this number is the inner product between the solution $Q$, a vector of flows defined on edges, and the cost $\omega_{ij}$. In our case, since the DMK model allows to control the flow of information through a hyperparameter $\beta \in (0,2]$, we define the $\beta$-*Wasserstein cost*, $\mathcal{W}_{\beta}(m_i, m_j, \omega_{ij})$, as the inner product of the solution $Q = Q(\beta)$ of the DMK model and the cost $\omega_{ij}$. For $\beta = 1$ we compute the Wasserstein-1 distance between $m_i$ and $m_j$, while for $\beta \neq 1$ the influence of $\beta$ in the solution of the transportation problem can be seen in Fig. 2. When $\beta < 1$, more edges of $B$ tend to be used to transport the mass, thus we observe congested transportation[22]. When $\beta > 1$ fewer edges are used, hence we observe branched transportation, and the $\beta$-*Wasserstein cost* coincides with a branched transport distance[21,36]. The idea of tuning $\beta$ to interpolate between various transportation regimes has been used in several works and engineering applications[23,37–42].

Calculating the Wasserstein cost is necessary to determine our main quantity of interest, the discrete Olliver-Ricci curvature, defined as:

$$\kappa_{\beta}(i,j) := 1 - \frac{\mathcal{W}_{\beta}(m_i, m_j, \omega_{ij})}{d_{ij}} \quad , \tag{1}$$

**Figure 1.** Left): an example graph $G$ where edges have unitary weights. Center): the edge $(1,5)$ (bold black line) is selected to define the OT problem between $m_1, m_5$; neighborhoods of nodes 1 and 5 are highlighted with blue and red edges and are used to build the corresponding distributions $m_1, m_5$. Right): The complete bipartite graph $B_{15}$ where the OT problem is defined. The color intensity of the edges represent the distance between the associated nodes on the graph $G$, as shown by the colorbar. $m_1$ and $m_5$ are both defined for $\alpha = 0$, i.e. no mass is left in 1 and 5.

where $d_{ij}$ is the weighted shortest path distance between $i$ and $j$ as measured in $G$. Intuitively, if $i$ and $j$ are in the same communities, several $k \in V_i$ and $\ell \in V_j$ will be also directly connected. Thus, the Wasserstein distance between $m_i$ and $m_j$ will be shorter than $d_{ij}$, yielding a positive $\kappa_\beta(i,j)$. Instead, when $i$ and $j$ are in different communities, their respective neighbors will be unlikely connected, hence $d(i,j) < \mathcal{W}_\beta(m_i, m_j, \omega_{ij})$, yielding a negative $\kappa_\beta(i,j)$.

The Ricci flow algorithm on a network is defined by iteratively updating the weights of the graph $G$[18,19]. These are updated by combining the curvature and shortest path distance information[27]. We redefine these updates using our proposal for the Ollivier-Ricci curvature:

$$w_{ij}^{(t+1)} := d_{ij}^{(t)} - \kappa_\beta^{(t)}(i,j) \cdot d_{ij}^{(t)}, \tag{2}$$

where $w_{ij}^{(t+1)}$ is the weight of edge $(i,j)$ at time $t$, $w_{ij}^{(0)} = d_{ij}^{(0)}$, and $d_{ij}^{(t)}$ is the shortest path distance between nodes $i$ and $j$ at iteration $t$. At every time step $t$, the weights are normalized by their total sum.

The algorithm ORC-Nextrout dynamically changes the weights of the graph $G$ to isolate communities: intra-community edges will be shortened, while inter-community ones will be enlarged. These changes are reached after different number of iterations, depending on the input data. We choose the one that maximises some predefined quality measure. To find the communities we apply a *network surgery* criterion as proposed by Ni et al.[18] based on the stabilisation of the modularity of the network. Notice that our algorithm does not need prior information about the number of communities: edges will be either enlarged or shortened depending on the optimal transport principles agnostic to community labelings.

A pseudo-code of the implementation is shown in Algorithm 1.

## 3 Results on Community Detection problems

### Synthetic Networks

To investigate the accuracy of our model in detecting communities, we consider synthetic networks generated using the *Lancichinetti–Fortunato–Radicchi* (LFR) benchmark[43] and the *Stochastic Block Model* (SBM)[44]. Both models provide community labels used as *ground-truth* information during the classification tasks.

*Lancichinetti–Fortunato–Radicchi benchmark:* this benchmark generates undirected unweighted networks $G$ with disjoint communities. It samples node degrees and community sizes from power law distributions, see Fig. 4 for an example. One of its advantages is that it generates networks with heterogeneous distributions of degrees and community sizes. The main parameters in input are the number of nodes $N$, two exponents $\tau_1$ and $\tau_2$ for the power law distributions of the node degree and community size respectively, the expected degree $d$ of the nodes, the maximum number of communities on the network $K_{max}$ and a fraction $\mu$ of inter-community edges incident to each node. To test the performance of our algorithm, we use the set of LFR networks used and provided by the authors of[18]. We set $\tau_1 = 2$, $\tau_2 = 1$, $d = 20$, $K_{max} = 50$ and $\mu \in [0.05, 0.75]$.

*Stochastic Block Model:* this model probabilistically generates networks with non-overlapping communities. One specifies the number of nodes $N$ and the number of communities $K$, together with the expected degree $d$ of a node and a ratio $r \in [0,1]$.
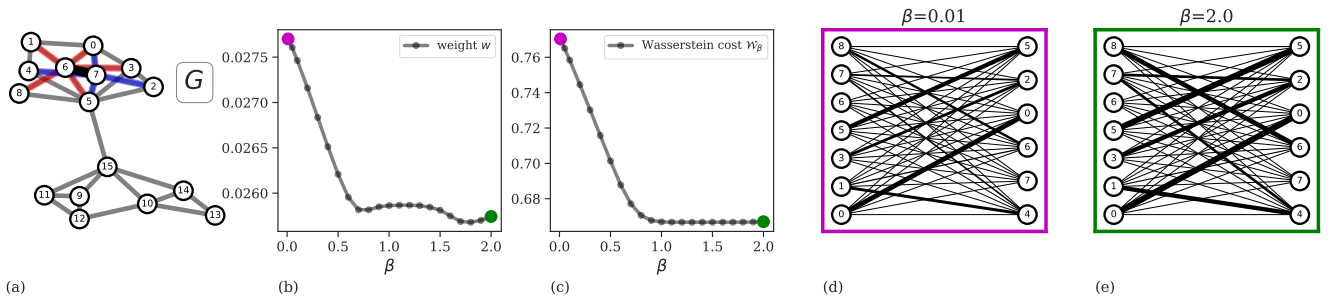
**Algorithm 1** ORC-Nextrout

---

**Input:** $G = (V, E, W)$, traffic rate $\beta$, *MaxIterNum* $\in \mathbb{N}$
**Output:** updated $W$
Initialize: $\mathbf{w}^0 = W$
**for** $t \in range(MaxIterNum)$ **do**
    **for** $e = (i, j) \in E$ **do**
        Calculate $m_i, m_j$
        Build $B_{ij}$
        Get $Q(\beta) \in \mathbb{R}^{|E_{ij}|}$, $Q(\beta) = DMK(B_{ij}, m_i, m_j, \beta)$
        Compute $\kappa_\beta(e)$
        Compute $\mathbf{w}_\beta(e)$
    **end for**
    Update $\mathbf{w}^t = \mathbf{w}_\beta$
**end for**

---



**Figure 2.** Visualization of how $\beta$ impacts an intra-community edge. (a) Example intra-community structure between nodes 6 and 7. (b) The weight of edge $(6,7)$ decreases when $0 < \beta < 0.6$, while for $0.5 < \beta < 2.0$ it reaches a minimum and then slightly increases again. This justifies the better performance in detecting communities obtained for higher values of $\beta$, as shown in Figs. 3a and 3b. (c) A similar decreasing behavior is observed for the $\beta$-Wasserstein cost: for intra-community edges, $\beta > 1$ consolidates traffic in the network as the Wasserstein cost stabilizes. (d-e) Example cost graph $B_{67}$ with fluxes solution of the OT problem (edge thickness is proportional to the amount of flux) in the regimes of small (d) and high (e) values of $\beta$.

Networks are generated by connecting nodes with a probability $r * p_{intra}$ if they belong to different communities; $p_{intra}$ if they are part of the same community, where $p_{intra} = d \times K/N$. Notice that the smaller the ratio $r$ is, the less inter-community connections would exist, which leads to networks with a more distinct community structure.

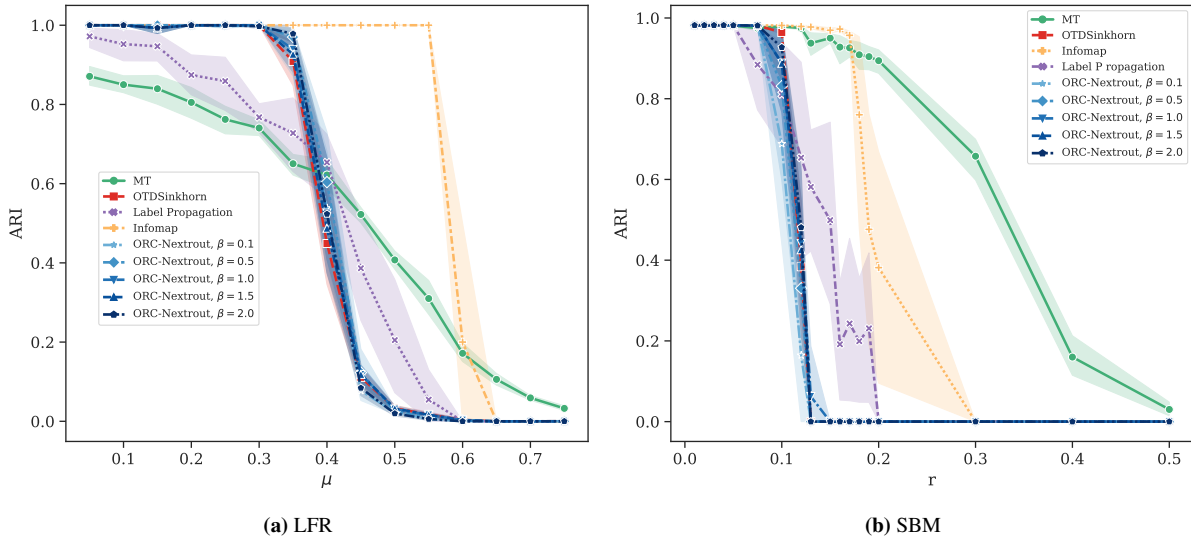We set $N = 500$, $K = 3$, $d = 15$ and $r \in [0.01, 0.5]$ and generate 10 random networks per value of $r$.

**Results.** To evaluate the performance of our method in recovering the communities, we use the *Adjusted Rand Index* (ARI)[45]. ARI compares the obtained community partition with the *ground truth* clustering. It takes values ranging from 0 to 1, where ARI=0 is equivalent to random community assignment, and ARI=1 denotes perfect matching with the ground truth communities, hence the higher this value the better the recovery of communities.

We test our algorithm for different types of information spreading in our OT-based model, as controlled by the parameter $\beta$, using the software developed in[46][1]. We used $\beta = 1$, i.e., standard Wasserstein distance; $\beta \in \{0.1, 0.5\}$ for congested transportation, enforcing broad spreading across the neighbors; and $\beta \in \{1.5, 2\}$ to favor branching schemes, where fewer edges are used to decide which community a node should belong to. For the OT-based algorithms, we run 15 iterations and choose the one with the best ARI scores. In some cases, high scores are reached in fewer iterations.
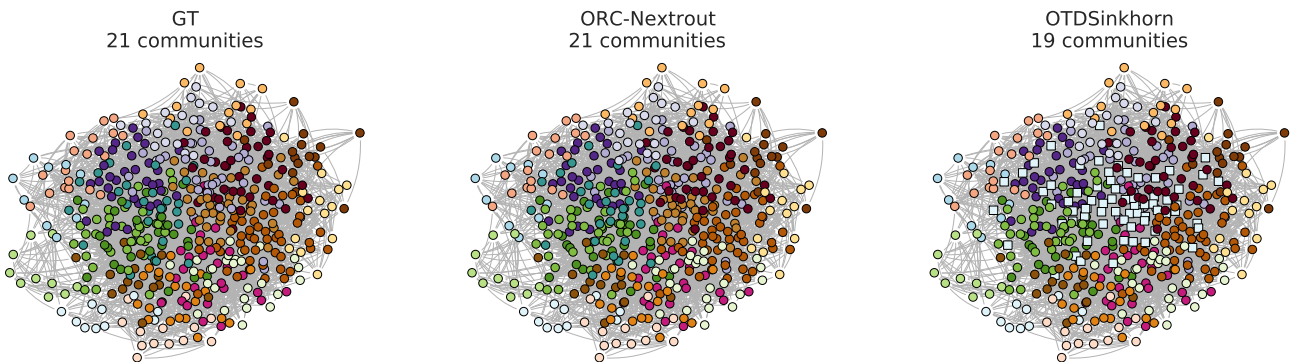
The results in Fig. 3 show the performance in both LFR and SBM benchmarks with OT-based methods, our method for various $\beta$ and one based on the Sinkhorn algorithms (OTDSinkhorn)[47,48]. Our main goal is to assess the impact of tuning between different transportation regimes (as done by $\beta$) in terms of community detection via OT principles. Nevertheless, to better contextualize the performance of OT-based algorithms in the wide spectrum of community detection methods, we also include comparisons with algorithms that are not OT-based. Namely, we consider a probabilistic model with latent variables (MT)[13], and with two modularity-based algorithms, Label Propagation[14] and Infomap[17]. Our algorithm outperforms

---

[1]Source code at https://gitlab.com/enrico_facca/dmk_solver

OTDSinkhorn for various values of $\beta$ in an intermediate regime where OT-based inference is not trivial. This happens in both benchmark LFR and SBM, as shown in Fig. 3. For lower and higher values of the parameters, performance is similar and close to the two extremes of ARI = 0 and 1. OT-based methods have a similar sharp decay in performance from the regime where inference is easy to the more difficult one, as also observed in[18]. The other community detection methods have smoother decay, but with lower performance in the regime where OT-based approaches strive, except for Label Propagation and MT, which are more robust in this sense. In the intermediate regime where inference is not trivial (i.e. along the sharp decay of OT-based methods), we observe that different values of $\beta$ give higher performance than OTDSinkhorn. For SBM the highest performance is achieved consistently for high $\beta = 2$, while for LFR the best $\beta$ varies with $\mu$. A qualitative example where ORC-Nextrout is performing better than OTDSinkhorn, in an instance of LFR of this intermediate regime, is shown in Fig. 4.



**(a)** LFR

**(b)** SBM

**Figure 3.** Results on LFR and SBM synthetic data. Performance in detecting ground-truth communities is measured by the ARI score. Markers and shadows are the averages and standard deviations over 10 network realisations with the same value of the parameter used in generation. Markers' shape denote different algorithms. a) LFR graph with $N = 500$ nodes and different values of $K$ ranging from $(17,22)$. b) SBM with $N = 500$ nodes, $K = 3$ communities and average degree $d = 15$. The parameter $r$ is the ratio of inter-community with intra-community edges.
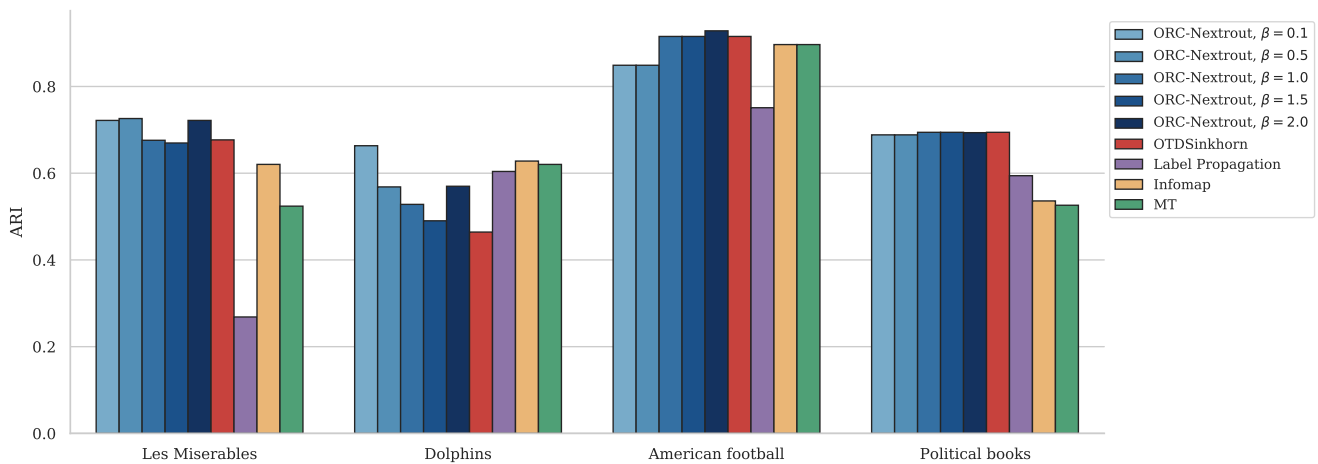


**Figure 4.** Example of community structure on a synthetic LFR network. The rightmost panel shows the ground-truth community structures to be predicted in an LFR network generated using $\mu = 0.35$. Square-shaped markers denote nodes that are assigned to communities different than those in ground-truth. In middle and last panels, ORC-Nextrout with $\beta = 2$ perfectly retrieves the 21 communities, while OTDSinkhorn predicts only 19 communities with an ARI score of 0.73, wrongly assigning ground-truth dark green and light brown (square-shaped) nodes to the light blue community.

## Analysis of real networks

Next, we evaluate our model on various real datasets[49] containing node metadata that can be used to assess the recovery of communities. While failing to recover communities that align well with node metadata should not by automatically interpreted as a model's failure[50] (e.g. the inferred communities and the chosen node metadata may capture different aspect of the data), having a reference community structure to compare against allows to inspect quantitatively difference between models. These real networks differ on structural features like number of nodes, average degree, number of communities and other standard network properties as detailed in Table 1. Specifically, we consider i) a network of co-appearances of characters in the novel *Les Misérables*[51] (Les Miserables). Edges are built between characters that encounter each other. ii) A network of 62 bottlenose dolphins in a community living off Doubtful Sound, in New Zealand[52] (Dolphins). Nodes represent dolphins, and edges indicate frequent associations between them. This network is clustered into four groups, conjectured as clustered from one population and three sub-populations based on the interactions between dolphins of different sex and ages[53]. The dolphins were observed between 1994 and 2001. iii) A network of Division I matches of American Football during a regular season in the fall of 2000[54] (American football). Nodes represent teams and edges are games between teams. Teams can be clustered according to their football college conference memberships. iv) A network of books on US politics published around the 2004 presidential election and sold by an online bookseller[55] (Political books). Nodes represent the books and the edges between books are frequent co-purchasing of books by the same buyers. Books are clustered based on their political spectrum as neural, liberal or conservative.

**Table 1. Real networks description.** We report statistics for the real networks used in our experiments. $N$ and $E$ denote the number of nodes and edges, respectively. $K$ is the number of communities in the ground truth data. AvgDeg, AvgBtw and AvgClust are the average degree, betweenness centrality and average clustering coefficient, respectively.

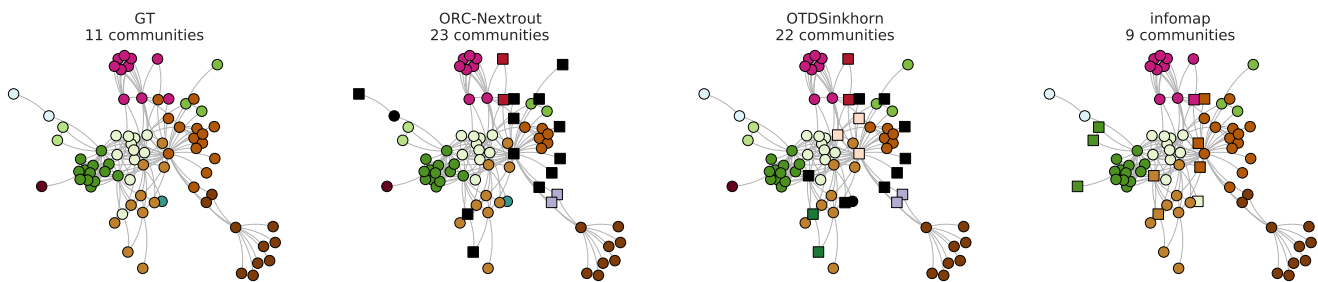| Dataset | $N$ | $E$ | $K$ | AvgDeg | AvgBtw | AvgClust |
|---|---|---|---|---|---|---|
| Les Miserables | 77 | 254 | 11 | 6.6 | 0.0219 | 0.5731 |
| Dolphins | 62 | 159 | 4 | 5.1 | 0.0393 | 0.2590 |
| American football | 115 | 613 | 12 | 10.7 | 0.0133 | 0.4032 |
| Political books | 105 | 441 | 3 | 8.4 | 0.0202 | 0.4875 |



**Figure 5.** Results on real data. Performance in terms of recovering communities using metadata information is calculated in terms of the ARI score. ORC-Nextrout shows competing results against all methods with different optimal $\beta$ across datasets.

OT-based algorithms outperform other community detection algorithms in detecting communities aligned with the node metadata, as shown in Fig. 5. In particular, ORC-Nextrout has the highest accuracy performance considering the best performing $\beta$. The impact of tuning this parameter is noticeable from these plots, as the best performing value varies across datasets. In the Les Miserables and Dolphins networks, $\beta < 1$ has better performance, while in American Football the best performing value is for $\beta > 1$. Performance is similar across OT-based methods in the Political Books network. In Fig. 6 we show the communities
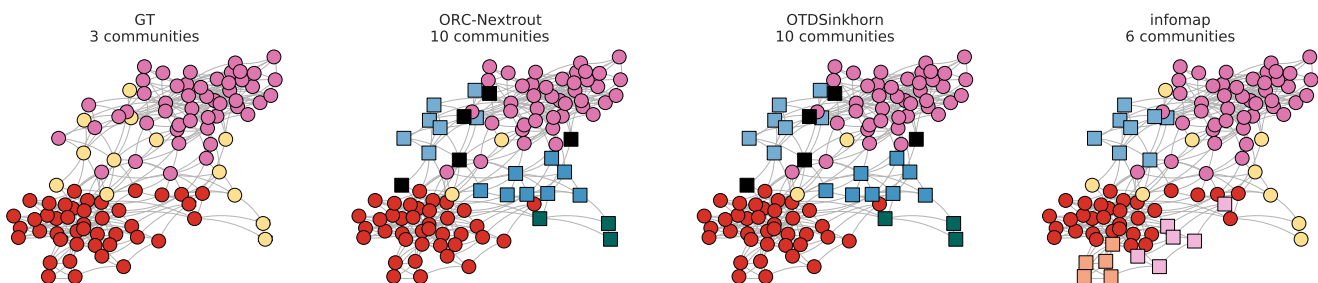
detected by the best performing ORC-Nextrout version together with OTDSinkhorn and Infomap in Les Miserables and Political books. Focusing on Les Miserables, we see how ORC-Nextrout successfully detects three characters in the green communities, in particular a highly connected node in the center of the figure (in dark green). Notice that these are placed in the same (pink or black) community by OTDSinkhorn. Thus ORC-Nextrout achieves a higher ARI than OTDSinkhorn. Both OT-based approaches retrieve well communities exhibiting clustering patterns, with many connections within community. Instead, they both divide the communities with a hub and spokes structure due to the lack of common connections within the group.

The communities detected in the Political books datasets highlight the tendency of OT-based methods to extract a larger number of communities (10) than those observed from node metadata (3). Among these extra communities, 3 are made of a few nodes, while 5 of them are made of one isolated node each. This is related to the fact that OT-based methods perform particularly well for networks with internally densely-connected community structures, but may be weaker for community structures that are sparsely connected[19]. One could potentially assign these nodes to larger communities, for instance by preferential attachment as done in[19], thus in practice reducing the number of communities. Devising a principled method or criterion to do this automatically is an interesting topic for future work. This tendency is further corroborated by the fact that OT-based algorithms recover robustly the two communities that are mostly assortative (red and pink in the figure), while they struggle to recover the disassortative community depicted in the centre (yellow). This is community has several connections with nodes in the other two communities and has been separated into smaller groups by OT-based approaches, as described above. This also highlights the need for methods that are robust against situations where a combination of assortative and disassortative communities coexist in a network.

**(a)** Les Miserables



**(b)** Political books



**Figure 6.** Communities in real networks. We show the communities inferred by ORC-Nextrout ($\beta = 0.5, 1.5$ for top and bottom rows respectively), OTDSinkhorn and Infomap and compare against those extracted using node attributes (GT). The visualization layout is given by the *Fruchterman-Reingold force-directed* algorithm[56], therefore, groups of well-connected nodes are located close to each other. Dark nodes represent individual nodes who are assigned to isolated communities by OT-based methods. Square-shaped markers denote nodes assigned to communities different than those obtained from node metadata.

## Conclusion

Community detection on networks is a relevant and challenging open area of research. Several methods have been proposed to tackle this issue, with no "best algorithm" that fits well every type of data. We focused here on a recent line of work that exploits principles from Optimal Transport theory combined with the geometric concept of Ollivier-Ricci curvature applied to discrete graphs. Our method is flexible in that it tunes between different transportation regimes to extract the information

necessary to compute the OR curvature on edges. On synthetic data, our model is able to identify communities more robustly than other OT-based methods based on the standard Wasserstein distance in the regime where inference is not trivial. On real data, our model shows either better or comparable performance in recovering community structure aligned with node metadata compared to other approaches, thanks to the ability to tune the parameter $\beta$.

A relevant advantage of OT-based methods is that the number of communities is automatically learned from data, contrarily to other approaches that need this as an input parameter. In this respect, our model has the tendency of overestimating this number, similarly to other OT-based methods. Understanding how to properly incorporate small-size communities into larger ones in a principled and automatic way is an interesting topic for future work. Similarly, it would be interesting to quantify the extent to which various $\beta$ capture different network topologies. To address this, one could for instance use methods to calculate the structural distance between networks[57] and correlate this against the values of the best performing $\beta$.

There are a number of directions in which this model could be extended. Nodes can be connected in more than one way, as in multilayer networks. Our model could be extended by considering a different $\beta$ for each edge type, as done in[42]. Similarly, real networks are often rich in additional information, e.g. attributes on nodes. It would be interesting to incorporate a priori additional information to inform community detection[58,59]. This information can potentially be used to mitigate the problem of overestimation of the number of communities, as explained above.

## Methods

### Optimal Transport Formulation

Consider the proability distributions $q$ that take pairs of vertices and also satisfy the constraints $\sum_i q_{ij} = m_j, \sum_j q_{ij} = m_i$. In other words, these are the joint distributions whose marginals are $m_i$ and $m_j$. We call these distributions *transport plans* between $m_i$ and $m_j$. The Optimal Transport problem we are interested in is that of finding a transport plan $q^*$ that minimises the quantity $\sum_{i \sim j} q_{ij} d_{ij}$, where $i \sim j$ means that nodes $i$ and $j$ are neighbors and $d_{ij}$ is the cost of transporting mass from $i$ to $j$, e.g. the distance between these two nodes. The quantity $\mathscr{W}_\beta(m_i, m_j, d) := \sum_{i,j} q_{ij}^* d_{ij}$, defined for this optimal $q^*$, is the *Wasserstein distance* between $m_i$ and $m_j$.

### The Dynamical Monge-Kantorovich model

It was recently proved[33,34] that solutions of the optimal transport problem previously stated can be found by turning that problem into a system of differential equations. This section is dedicated to describe this dynamical formulation.

Let $G = (V, E, W)$ be a weighted graph, with $N$ the number of nodes and $E$ the number of edges in $G$. Let **B** be the *signed incidence* matrix of $G$. Let $f^+$ and $f^-$ be two $N$-dimensional discrete distributions such that $\sum_{i \in V} f_i = 0$ for $f = f^+ - f^-$; let $\mu(t) \in \mathbb{R}^E$ and $u(t) \in \mathbb{R}^N$ be two time-dependent functions defined on edges and nodes, respectively. The discrete *Dynamical Monge-Kantorovich model* can be written as:

$$f_i = \sum_e B_{ie} \frac{\mu_e(t)}{w_e} \sum_j B_{ej} u_j(t), \tag{3}$$

$$\mu'_e(t) = \left[ \frac{\mu_e(t)}{w_e} \left| \sum_j B_{ej} u_j(t) \right| \right]^\beta - \mu_e(t), \tag{4}$$

$$\mu_e(0) > 0, \tag{5}$$

where $|\cdot|$ is the absolute value element-wise. Equation (3) corresponds to Kirchhoff's law, Eq. (4) is the discrete dynamics with $\beta$ a traffic rate controlling the different routing optimization mechanisms; Eq. (5) is the initial distribution for the edge conductivities.

For $\beta = 1$ the dynamical system described by Eqs. (3)-(5) is known to reach a steady state, i.e., the updates of $\mu_e$ and $u_e$ converge to stationary functions $\mu^*$ and $u^*$ as $t$ inscreases. The flux function $q$ defined as $q_e^* := \mu_e^* |u_i^* - u_j^*| / w_e$ is the solution of the optimal transport problem presented in the previous section. Notice that $\mu$ and $u$ depend on the chosen traffic rate $\beta$, and thus, so does $q = q(\beta)$. Therefore we can introduce a generalized version of the distance $\mathscr{W}$:

$$\mathscr{W}_\beta(m_i, m_j, w) := \sum_{i,j} q_{ij}^*(\beta) w_{ij}.$$

We then redefine the proposed Ollivier-Ricci curvature as:

$$\kappa_\beta(i, j) := 1 - \frac{\mathscr{W}_\beta(m_i, m_j, w)}{d_{ij}}.$$

### Probability distributions on neighborhoods

ORC-Nextrout takes in input a graph and a forcing term. While the graph encapsulates the neighborhood information provided by the nodes $i$ and $j$, the forcing function is related to the distributions one needs to transport. Analogously to what proposed by[18], we define this graph to be the weighted *complete bipartite* $B_{ij} = (V_{ij}, E_{ij}, \omega_{ij})$. The weights in $\omega_{ij}$ change iteratively based on the curvature. Notice that a bipartite graph must satisfy $\mathcal{N}(i) \cap \mathcal{N}(j) = \varnothing$, which does not hold true if $i$ and $j$ have common neighbors (this is always the case since $i \in \mathcal{N}(j)$). Nonetheless, this condition does not have great repercussions in the solution of the optimal transport problem since the weights corresponding to these edges (of the form $(i,i)$) are equal to 0. As for the forcing function, we define it to be $f := f^+ - f^- = m_i - m_j$.

### Other methods

To evaluate the performance of ORC-Nextrout, we compare with some of the well-established community detection algorithms including: Infomap[17], MULTITENSOR[13] (MT), discrete Ricci flow[18] (OTDSinkhorn), and Label propagation[14]. We briefly describe each of these algorithms as follows;

- The *Discrete Ricci flow* (here addressed as OTDSinkhorn)[18] is an iterative node clustering algorithm that deforms edge weights as time progresses, by shrinking sparsely traveled links and stretching heavily traveled edges. These edge weights are iteratively updated based on neighborhood transportation Wasserstein costs, in a similar way to what proposed in this manuscript. After a predefined number of iterations, heavily traveled links are removed from the graph. Communities are then obtained as the connected components of this modified network.

- MULTITENSOR (MT)[13] is an algorithm to find communities in multilayer networks. It is a probabilistic model with latent variables regulating community structure and runs with a complexity of $O(EK)$ with assortative structure (as we consider here), where $K$ is the number of communities. This model assumes that the nodes inside the communities can belong to multiple groups (mixed-membership). In this implementation we use their validity for single layer networks (a particular case of a multilayer network).

- *Infomap*[17] employs information theoretic approach for community detection. This method uses the map equation to attend patterns of flow on a network. This flow is simulated using random walkers' traversed paths. Based on the theoretic description of these paths, nodes with quick information flow are then clustered into the same groups. The algorithm runs in $O(E)$.

- *Label propagation*[14] assigns each node to same community as majority of its neighbors. Its working principle start by initializing each node with a distinct label and converges when every node has same label as majority of its neighboring node. The algorithm has a complexity scaling as $O(E)$.

## Data Availability

The real data can be obtained from network data[49] and the synthetic one from the corresponding author upon request.

## References

1. Huang, X., Chen, D., Ren, T. & Wang, D. A survey of community detection methods in multilayer networks. *Data Min. Knowl. Discov.* **35**, 1–45 (2021).

2. Newman, M. *Networks* (Oxford university press, 2018).

3. Fortunato, S. Community detection in graphs. *Phys. reports* **486**, 75–174 (2010).

4. Fortunato, S. & Hric, D. Community detection in networks: A user guide. *Phys. reports* **659**, 1–44 (2016).

5. Weber, M., Jost, J. & Saucan, E. Detecting the coarse geometry of networks. *In NeurIPS 2018 Work.* (2018).

6. Samal, A. *et al.* Comparative analysis of two discretizations of ricci curvature for complex networks. *Sci. reports* **8**, 1–16 (2018).

7. Zhu, J., Wang, B., Wu, B. & Zhang, W. Emotional community detection in social network. *IEICE Transactions on Inf. Syst.* **100**, 2515–2525 (2017).

8. Wang, D., Li, J., Xu, K. & Wu, Y. Sentiment community detection: exploring sentiments and relationships in social networks. *Electron. Commer. Res.* **17**, 103–132 (2017).

9. Li, C. & Zhang, Y. A personalized recommendation algorithm based on large-scale real micro-blog data. *Neural Comput. Appl.* **32**, 11245–11252 (2020).

10. Waskiewicz, T. Friend of a friend influence in terrorist social networks. In *Proceedings on the international conference on artificial intelligence (ICAI)*, 1 (The Steering Committee of The World Congress in Computer Science, Computer . . . , 2012).

11. Pinheiro, C. A. R. Community detection to identify fraud events in telecommunications networks. *SAS SUGI proceedings: customer intelligence* (2012).

12. Chen, J., Zhang, H., Guan, Z.-H. & Li, T. Epidemic spreading on networks with overlapping community structure. *Phys. A: Stat. Mech. its Appl.* **391**, 1848–1854 (2012).

13. De Bacco, C., Power, E. A., Larremore, D. B. & Moore, C. Community detection, link prediction, and layer interdependence in multilayer networks. *Phys. Rev. E* **95**, 042317 (2017).

14. Raghavan, U. N., Albert, R. & Kumara, S. Near linear time algorithm to detect community structures in large-scale networks. *Phys. review E* **76**, 036106 (2007).

15. Clauset, A., Newman, M. E. & Moore, C. Finding community structure in very large networks. *Phys. review E* **70**, 066111 (2004).

16. Reichardt, J. & Bornholdt, S. Statistical mechanics of community detection. *Phys. review E* **74**, 016110 (2006).

17. Rosvall, M. & Bergstrom, C. T. Maps of random walks on complex networks reveal community structure. *Proc. national academy sciences* **105**, 1118–1123 (2008).

18. Ni, C.-C., Lin, Y.-Y., Luo, F. & Gao, J. Community detection on networks with ricci flow. *Sci. reports* **9**, 1–12 (2019).

19. Sia, J., Jonckheere, E. & Bogdan, P. Ollivier-ricci curvature-based method to community detection in complex networks. *Sci. reports* **9**, 1–12 (2019).

20. Santambrogio, F. Optimal channel networks, landscape function and branched transport. *Interfaces Free. Boundaries* **9**, 149–169 (2007).

21. Facca, E., Cardin, F. & Putti, M. Branching structures emerging from a continuous optimal transport model. *J. Comput. Phys.* **447**, 110700 (2021).

22. Brasco, L., Carlier, G. & Santambrogio, F. Congested traffic dynamics, weak flows and very degenerate elliptic equations. *J. de mathématiques pures et appliquées* **93**, 652–671 (2010).

23. Baptista, D., Leite, D., Facca, E., Putti, M. & De Bacco, C. Network extraction by routing optimization. *Sci. reports* **10**, 1–13 (2020).

24. Sandhu, R. *et al.* Graph curvature for differentiating cancer networks. *Sci. reports* **5**, 1–13 (2015).

25. Wang, C., Jonckheere, E. & Banirazi, R. Interference constrained network control based on curvature. In *2016 American Control Conference (ACC)*, 6036–6041 (IEEE, 2016).

26. Ni, C.-C., Lin, Y.-Y., Gao, J., Gu, X. D. & Saucan, E. Ricci curvature of the internet topology. In *2015 IEEE conference on computer communications (INFOCOM)*, 2758–2766 (IEEE, 2015).

27. Ollivier, Y. Ricci curvature of markov chains on metric spaces. *J. Funct. Analysis* **256**, 810–864 (2009).

28. Ollivier, Y. A survey of ricci curvature for metric spaces and markov chains. In *Probabilistic approach to geometry*, 343–381 (Mathematical Society of Japan, 2010).

29. Forman, R. Bochner's method for cell complexes and combinatorial ricci curvature. *Discret. Comput. Geom.* **29**, 323–374 (2003).

30. Ni, C.-C., Lin, Y.-Y., Gao, J. & Gu, X. Network alignment by discrete ollivier-ricci flow. In *International Symposium on Graph Drawing and Network Visualization*, 447–462 (Springer, 2018).

31. Ye, Z., Liu, K. S., Ma, T., Gao, J. & Chen, C. Curvature graph network. In *International Conference on Learning Representations* (2019).

32. Lancichinetti, A. & Fortunato, S. Community detection algorithms: a comparative analysis. *Phys. review E* **80**, 056117 (2009).

33. Facca, E., Cardin, F. & Putti, M. Towards a stationary monge–kantorovich dynamics: The physarum polycephalum experience. *SIAM J. on Appl. Math.* **78**, 651–676 (2018).

34. Facca, E., Daneri, S., Cardin, F. & Putti, M. Numerical solution of monge–kantorovich equations via a dynamic formulation. *J. Sci. Comput.* **82**, 1–26 (2020).

35. Asikainen, A., Iñiguez, G., Ureña-Carrión, J., Kaski, K. & Kivelä, M. Cumulative effects of triadic closure and homophily in social networks. *Sci. Adv.* **6**, eaax7310 (2020).

36. Xia, Q. Optimal paths related to transport problems. *commcont* **5**, 251–279 (2003).

37. Baptista, D. & De Bacco, C. Principled network extraction from images. *Royal Soc. Open Sci.* **8**, 210025 (2021).

38. Baptista, D. & De Bacco, C. Convergence properties of optimal transport-based temporal networks. In *International Conference on Complex Networks and Their Applications*, 578–592 (Springer, 2021).

39. Lonardi, A., Facca, E., Putti, M. & De Bacco, C. Designing optimal networks for multicommodity transport problem. *Phys. Rev. Res.* **3**, 043010 (2021).

40. Lonardi, A., Putti, M. & De Bacco, C. Multicommodity routing optimization for engineering networks. *Sci. Reports* **12**, 1–11 (2022).

41. Lonardi, A., Facca, E., Putti, M. & De Bacco, C. Infrastructure adaptation and emergence of loops in network routing with time-dependent loads. *arXiv preprint arXiv:2112.10620* (2021).

42. Ibrahim, A. A., Lonardi, A. & Bacco, C. D. Optimal transport in multilayer networks for traffic flow optimization. *Algorithms* **14**, 189 (2021).

43. Lancichinetti, A., Fortunato, S. & Radicchi, F. Benchmark graphs for testing community detection algorithms. *Phys. review E* **78**, 046110 (2008).

44. Holland, P. W., Laskey, K. B. & Leinhardt, S. Stochastic blockmodels: First steps. *Soc. networks* **5**, 109–137 (1983).

45. Hubert, L. & Arabie, P. Comparing partitions. *J. classification* **2**, 193–218 (1985).

46. Facca, E. & Benzi, M. Fast iterative solution of the optimal transport problem on graphs. *SIAM J. on Sci. Comput.* **43**, A2295–A2319 (2021).

47. Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Adv. neural information processing systems* **26** (2013).

48. Flamary, R. *et al.* Pot: Python optimal transport. *J. Mach. Learn. Res.* **22**, 1–8 (2021).

49. Network data, http://www-personal.umich.edu/ mejn/netdata/.

50. Peel, L., Larremore, D. B. & Clauset, A. The ground truth about metadata and community detection in networks. *Sci. advances* **3**, e1602548 (2017).

51. Knuth, D. E. *The Stanford GraphBase: a platform for combinatorial computing*, vol. 1 (AcM Press New York, 1993).

52. Lusseau, D. *et al.* The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations. *Behav. Ecol. Sociobiol.* **54**, 396–405 (2003).

53. Lusseau, D. & Newman, M. E. Identifying the role that animals play in their social networks. *Proc. Royal Soc. London. Ser. B: Biol. Sci.* **271**, S477–S481 (2004).

54. Girvan, M. & Newman, M. E. Community structure in social and biological networks. *Proc. national academy sciences* **99**, 7821–7826 (2002).

55. Books about us politics dataset, http://www.orgnet.com/.

56. Fruchterman, T. M. & Reingold, E. M. Graph drawing by force-directed placement. *Software: Pract. experience* **21**, 1129–1164 (1991).

57. Xiao, X., Chen, H. & Bogdan, P. Deciphering the generating rules and functionalities of complex networks. *Sci. reports* **11**, 1–15 (2021).

58. Contisciani, M., Power, E. A. & De Bacco, C. Community detection with node attributes in multilayer networks. *Sci. reports* **10**, 1–16 (2020).

59. Newman, M. E. & Clauset, A. Structure and inference in annotated networks. *Nat. communications* **7**, 11863 (2016).

## Acknowledgements

## Author contributions statement

All authors contributed to developing the models, conceived the experiments, analyzing the results and reviewing the manuscript. DL, DB and AI conducted the experiments.

## Additional information

**Accession codes**: open source codes and executables are available at https://github.com/Danielaleite/ORC-Nextrout.
**Competing interests**. The authors declare no competing interests.