# Bilevel optimization in flow networks – A message-passing approach

Bo Li,[1, 2] David Saad,[1] and Chi Ho Yeung[3]

[1] *Non-linearity and Complexity Research Group, Aston University, Birmingham, B4 7ET, United Kingdom*
[2] *School of Science, Harbin Institute of Technology (Shenzhen), Shenzhen, 518055, China*
[3] *Department of Science and Environmental Studies,*
*The Education University of Hong Kong, 10 Lo Ping Road, Taipo, Hong Kong*

Optimizing embedded systems, where the optimization of one depends on the state of another, is a formidable computational and algorithmic challenge, that is ubiquitous in real world systems. We study flow networks, where bilevel optimization is relevant to traffic planning, network control and design, and where flows are governed by an optimization requirement subject to the network parameters. We employ message-passing algorithms in flow networks with sparsely coupled structures to adapt network parameters that govern the network flows, in order to optimize a global objective. We demonstrate the effectiveness and efficiency of the approach on randomly generated graphs.

Many problems in science and engineering involve hierarchical optimization, whereby some of the variables cannot be freely varied but are governed by another optimization problem [1]. As a motivating example, consider the task of designing a network (e.g., a road or communication network) that maximizes the throughput of commodities or information flow. While the designer controls the network parameters (upper-level optimization), traffic flows are determined by the network users who maximize their own benefit (lower-level optimization) [2]. Therefore, the designer needs to adapt the network intricately, taking into account the reaction of network users. Similarly, many physical systems admit a certain extremization principle for given controllable system parameters, e.g., minimal free energy in thermal equilibrium [3], electric flows in resistor networks that minimize dissipation [4, 5] and, entropy maximization and parameter optimization that are used across disciplines in inference and learning tasks [6, 7]. Adapting system parameters to extremize a given objective requires bilevel optimization, which considers both system parameters and the inherent optimization of the physical variables.

Bilevel optimization is intrinsically difficult to solve [8]. In fact, even the simple instance where both levels are linear programming tasks is NP-hard [9, 10]. Generic methods for bilevel optimization include (i) bilevel programming approach, by expressing the lower-level optimization problem as nonlinear constraints and solving the bilevel problem as global optimization [11, 12]; (ii) gradient-descent method by computing the descent direction of the upper-level objectives while keeping the valid lower-level state variables [13, 14]. The former introduces complicated nonlinear constraints, making the reduced single-level problem difficult in general, while the latter can be challenging in computing the descent direction [8]. Moreover, such generic methods do not utilize existing system structure to simplify the task.

In this Letter, we develop message-passing (MP) algorithms to tackle bilevel optimization in sparse flow networks. The advances presented in this work are three-fold: (i) the derived MP algorithms are intrinsically distributed, scalable, and generally efficient; (ii) they are applicable to bilevel optimization problems with combinatorial constraints, which are difficult for generic bilevel programming approaches; (iii) these algorithms can successfully deal with nonsmooth flow problems, having potential applications for transport based approaches in machine learning [15–17].

*Routing Game.* We focus on a network planning problem in the routing game setting, widely used in modeling route choices of drivers [18]. Users on the road network make their route choices in a selfish and rational manner, where the corresponding Nash equilibrium is generally not the most beneficial for the global utility, measured by the total travel time of all users [2, 19]. The operator's task is to set the appropriate tolls or rewards on network edges to reduce the total travel time while taking into account the reactions of users to the tolls [20–22]. Recently, the idea of reducing traffic congestion by economic incentives to influence drivers' behaviors has regained interest [23–25], partly due to the deployment of smart devices and data availability [26–28]. Here, we focus on the algorithmic aspect of toll optimization.

The road network is represented by a directed graph $G(V, E)$, where $V$ is the set of nodes (junctions) and $E$ the set of directed edges (unidirectional roadways), having one connected component. Users routing from an origin node $i_0$ to a destination node $\mathcal{D}$ would select a path $\mathcal{P} = ((i_0, i_1), (i_1, i_2), ..., (i_{n-2}, i_{n-1}), (i_{n-1}, \mathcal{D}))$ by minimizing their total travel time $\sum_{e \in \mathcal{P}} \ell_e(x_e)$, or alternative cost, where the edge flow $x_e$ represents the number of users choosing edge $e$ and $\ell_e(x_e)$ is the corresponding latency function. It is assumed that $\ell_e$ is monotonically increasing with the edge flow $x_e$. The social cost is defined as the total travel time of all users $H = \sum_{e \in E} x_e \ell_e(x_e)$, which is the overall objective of the bilevel optimization problem.

We consider the limit of a large number of users, where each user controls an infinitesimal fraction of the overall traffic, such that the edge flow $x_e$ is a continuous variable. This is termed the nonatomic game setting [19]. As the equilibrium reached by the selfish decisions of
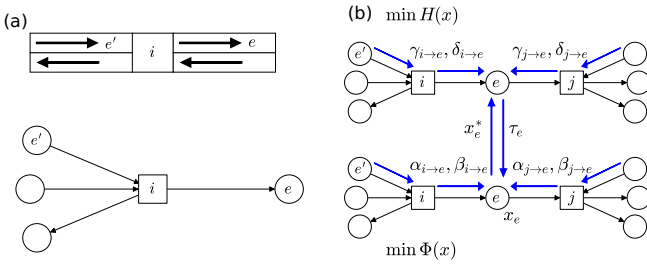
Figure 1. (a) Top: a directed road network section with a junction node $i$. Bottom: the corresponding factor graph representation; node $i$ is a factor node and is marked by a square. (b) Bilevel MP for toll planning. Blue arrows indicate the directions of messages. The equilibrium flow $x_e^*$ is determined in the lower level, while the toll $\tau_e$ is set in the upper level.

users does not generally achieve the lowest social cost, we seek to place tolls $\{\tau_e\}$ on edges to influence users' route choices. Gauging the monetary penalty at the same scale as latency, users will choose a path $\mathcal{P}$ that minimizes the combined total journey cost in latency and tolls $\sum_{e\in\mathcal{P}} \left[\ell_e(x_e) + \tau_e\right]$. If tolls can be placed freely on all edges, marginal cost pricing is known to induce socially optimal flow for nonatomic games [21]. However, it is usually infeasible to set an unbounded toll on every road, which renders marginal cost pricing less applicable. We therefore consider restricted tolls $0 \leq \tau_e \leq \tau_e^{\max}$; an edge $e$ is not chargeable when $\tau_e^{\max} = 0$. For simplicity, we do not consider the income from tolls to contribute to the social cost [29]. In total, $\Lambda_i$ users are traveling from node $i$ to a universal destination $\mathcal{D}$, where the case with multiple destinations is discussed in the supplemental material (SM) [30]. The resulting edge flows satisfy the non-negativity $x_e \geq 0$ and the flow conservation constraints

$$R_i = \Lambda_i + \sum_{e\in\partial_i^{\mathrm{in}}} x_e - \sum_{e\in\partial_i^{\mathrm{out}}} x_e = 0, \qquad (1)$$

where $\partial_i^{\mathrm{in}}$ and $\partial_i^{\mathrm{out}}$ are the sets of incoming and outgoing edges adjacent to node $i$. It has been established that the edge flows in user equilibrium (i.e., the Wardrop's equilibrium [2]) can be obtained by minimizing a potential function $\Phi = \sum_{e\in E} \phi_e(x_e) := \sum_{e\in E} \int_0^{x_e}[\ell_e(y) + \tau_e]\mathrm{d}y$ subject to the constraints of Eq. (1) [31, 32]. We emphasize that the potential function $\Phi(\boldsymbol{x})$ only plays an auxiliary role in defining the equilibrium flows; the values of $\Phi$ do *not* correspond to the routing costs of users.

The lower-level optimization is a nonlinear min-cost flow problem, where edge flows are coupled through the conservation constraints in Eq. (1), represented as factor nodes in Fig. 1(a). We employ the MP approach developed in Ref. [33] to tackle the nonlinear optimization problem. It turns the global optimization of the potential into a local computation of the following message

functions

$$\Phi_{i\to e}(x_e) = \min_{\{x_{e'}\geq 0\}|R_i=0} \sum_{e'\in\partial i\setminus e} \left[\Phi_{e'\to i}(x_{e'}) + \phi_{e'}(x_{e'})\right], \qquad (2)$$

where $\partial i = \partial_i^{\mathrm{in}} \cup \partial_i^{\mathrm{out}}$ and $\Phi_{i\to e}(x_e)$ relates to the optimal potential function contributed by the flows adjacent to node $i$ where the flow on edge $e$ is set to $x_e$, taking into account flow conservation at node $i$. In Eq. (2), denoting $e' = (k, i)$, we can write $\Phi_{e'\to i}(x_{e'}) = \Phi_{k\to e'}(x_{e'})$; therefore only factor-to-variable messages are needed. The message $\Phi_{k\to e'}(x_{e'})$ can be obtained recursively by an expression similar to Eq. (2), but using the incoming messages from its upstream edges $\{l \to k|(l, k) \in \partial k\setminus i\}$. Upon computing the messages iteratively until convergence, we can determine the equilibrium flow $x_e^*$ on edge $e = (i, j)$ by minimizing the edgewise full energy dictated by the nonlinear cost $\phi_e(x_e)$ and messages from both ends of edge $e$, defined as $\Phi_e^{\mathrm{full}}(x_e) = \Phi_{i\to e}(x_e) + \Phi_{j\to e}(x_e) + \phi_e(x_e)$.

This algorithm can be demanding when different values of $x_e$ are needed to determine the profile of the message $\Phi_{i\to e}(x_e)$. To reduce the computational cost, we consider the approximation of the message in the vicinity of some working point $\tilde{x}_{i\to e}$ as

$$\Phi_{i\to e}(\tilde{x}_{i\to e} + \varepsilon_e) \approx \Phi_{i\to e}(\tilde{x}_{i\to e}) + \beta_{i\to e}\varepsilon_e + \frac{1}{2}\alpha_{i\to e}\left(\varepsilon_e\right)^2, \qquad (3)$$

where $\beta_{i\to e}$ and $\alpha_{i\to e}$ are the first and second derivatives of $\Phi_{i\to e}$ evaluated at $\tilde{x}_{i\to e}$, assuming the derivatives exist. For a particular $\tilde{x}_{i\to e}$, the computation of the message function $\Phi_{i\to e}(x_e)$ in Eq. (2) reduces to the optimization of $\beta_{i\to e}$ and $\alpha_{i\to e}$ by using $\{\tilde{x}_{k\to e'}, \beta_{k\to e'}, \alpha_{k\to e'}|e' = (k, i) \in \partial i\setminus e\}$. The working point $\tilde{x}_{i\to e}$ is updated by pushing it towards the minimizer $x_e^*$ of the full energy $\Phi_e^{\mathrm{full}}(x_e)$ gradually [30]. The iterative updates of the coefficients $\{\beta_{i\to e}, \alpha_{i\to e}\}$ and the working points $\{\tilde{x}_{i\to e}\}$ constitute a perturbative version of the original MP algorithm, which only requires to keep track of a few coefficients rather than the full profile of $\Phi_{i\to e}$, making it tractable [33]. It has been shown to work remarkably well in many network flow problems [34], while the algorithm may not converge in problems with nonsmooth characteristics [33]. We discover that the nonnegativity constraints on flows can result in a nonsmooth message function $\Phi_{i\to e}(x_e)$, which makes the approximation of Eq. (3) inadequate. One solution is to approximate $\Phi_{i\to e}(x_e)$ by a continuous and piecewise quadratic function with at most two branches, where each branch $m$ is a quadratic function governed by three coefficients $\{\tilde{x}_{i\to e}, \beta_{i\to e}^{(m)}, \alpha_{i\to e}^{(m)}\}$, as illustrated in Fig. 2(a). Taking into account the nonsmooth structures, MP algorithms converge well even in loopy networks and provide the correct solutions [30]. We demonstrate the case of random regular graphs (RRG) with degree 3 in Fig. 2(b).
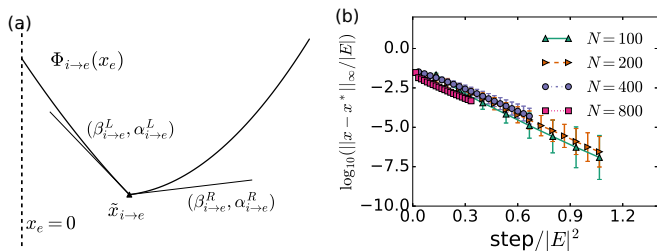
Figure 2. (a) A nonsmooth message function $\Phi_{i\to e}(x_e)$ with one breakpoint. (b) Convergence of the single-level MP algorithm for computing equilibrium flows in routing games in random regular graphs with degree 3 of different sizes $N = |V|$. An affine latency model $\ell_e(x_e) = t_e(1 + sx_e/c_e)$ is considered, where $t_e$ and $c_e$ are the free traveling time and edge capacity, respectively, while $s$ is a sensitivity measure of latency to congestion [19]. Random sequential schedule of MP updates has been used.

For bilevel optimization, we notice that the cost function of the upper layer $H(x)$ has a similar structure as $\Phi(x)$. Therefore, one can apply a similar MP procedure as $H_{i\to e}(x_e) = \min_{\{x_{e'}\}|R_i=0}\sum_{e'\in\partial i\backslash e}[H_{e'\to i}(x_{e'}) + x_{e'}\ell_e(x_{e'})]$. The message $H_{i\to e}(x_e)$ can also be approximated by a piecewise quadratic function with at most one break point, where each branch $m$ has the form $H_{i\to e}^{(m)}(\tilde{x}_{i\to e} + \varepsilon_e) \approx H_{i\to e}^{(m)}(\tilde{x}_{i\to e}) + \gamma_{i\to e}^{(m)}\varepsilon_e + \frac{1}{2}\delta_{i\to e}^{(m)}(\varepsilon_e)^2$. As the equilibrium state is determined in the lower level, the working points $\{\tilde{x}_{i\to e}\}$ in the lower-level MP are also used for the upper level. The landscape of the edgewise full cost $H_e^{\text{full}}(x_e) = H_{i\to e}(x_e) + H_{j\to e}(x_e) + x_e\ell(x_e)$ provides the information for setting the toll. Specifically, the toll is updated by $\min_{\tau_e} H_e^{\text{full}}(x_e^*(\tau_e))$, where the toll-dependent equilibrium flow $x_e^*$ is provided by the lower-level messages. In practice, an approximate $H_e^{\text{full}}$ is sufficiently informative for updating tolls. The basic structure of such bilevel MP is illustrated in Fig. 1(b), while details are provided in the SM [30].

We demonstrate the effectiveness of the proposed bilevel MP algorithm for tasks on RRG in Fig. 3(a), where the set-up is the same as in Fig. 2(b). Experiments on other networks and the cases of multiple destinations are discussed in the SM [30]. Although bilevel message-passing does not generally converge to a set of *unique* optimal tolls due to the non-convex nature of the problem, we found that the social costs are reduced when tolls are updated during MP. The scaling relation in the inset of Fig. 3(a) empirically indicates that the number of updates is $O(|E|^2)$ for achieving a given cost reduction. Moreover, the MP algorithm can be implemented in a fully distributed manner, unlike the generic global optimization approach [30]. Note that we have utilized the special set-up of routing games here, where the social optimum $H_S = \min_x H(x)$ can be obtained *a priori* for this benchmark. Such information may be unavailable in other bilevel-optimization problems. The toll optimiza-
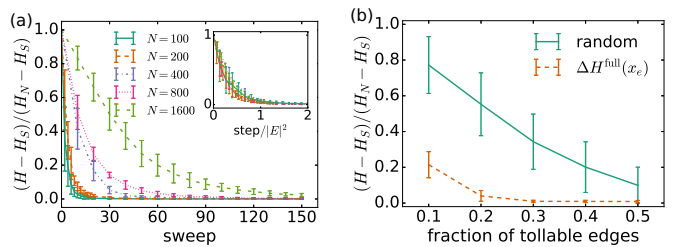


Figure 3. Bilevel MP algorithm for routing games on RRG. (a) Effect of tolls on the fractional social cost reduction $(H(x^*(\tau)) - H_S)/(H_N - H_S)$, where $H_S$ and $H_N$ represent the social costs at the social optimum and the Nash equilibrium without tolls. Tolls $\tau$ are recorded during bilevel MP updates, based on the resulting equilibrium flows $x^*(\tau)$ and social cost $H(x^*(\tau))$. Each data point is the average of 10 different problem realizations. Each sweep consists of $40|E|$ local MP steps and 100 edgewise toll updates in a random sequential schedule. A fixed number of sweeps without toll updates are performed to warm up the system. Inset: panel (a) with $x$-axis as MP steps rescaled by $|E|^2$. (b) Fractional cost reduction as a function of the fraction of tollable edges on an RRG with $N = 200$. A random selection of edges to be charged is compared with selections based on edgewise full cost reduction $H_e^{\text{full}}(x_e^*)$.

tion problem can also be tackled by the bilevel programming approach [11, 12]; however, it requires a treatment with mixed integer programming, which is centralized and generally not scalable, unlike the MP approach [30].

*Combinatorial Problems.* In practice, it may be infeasible to charge for every edge, but desirable to choose a subset of tollable edges for toll-setting [35], which is a difficult combinatorial optimization problem. As the cost landscape is manifested locally by the message functions, we heuristically select the tollable edges according to the largest possible reduction in edgewise full cost $H_e^{\text{full}}(x_e^*)$ due to tolling, which effectively selects the chargeable links as seen in Fig. 3(b). Such combinatorial problems are generally very difficult for traditional bilevel-optimization methods, while MP algorithms can provide approximate solutions in some scenarios.

Another important class of combinatorial problems is the atomic games which consider integer flow variables $\{x_e\}$ [36]. In principle, atomic games can be solved via the same MP procedure as in Eq. (2), where the message $\Phi_{i\to e}(x_e)$ is defined on a one dimensional grid. Using the techniques in Refs. [37–41], the MP approach provides a scalable algorithm to approximately tackle the difficult combinatorial optimization of atomic games in a single level; it can also solve instances of the bilevel toll-optimization problems. However, its performance is sub-optimal in large networks and for cases with heavy loads [30]. Nevertheless, we found some interesting patterns of the optimal tolls in a realistic test case network using this method [30].

*Flow Control.* We consider the problem of tuning

network flows to achieve certain functionality. In this example, resources need to be transported from source nodes to destination along edges in an undirected network $G(V, E)$, where the equilibrium flows $\{x_{ij}^*\}$ minimize the transportation cost $C = \sum_{(i,j)\in E} \frac{1}{2} r_{ij} x_{ij}^2$, subject to flow conservation constraints similar to Eq. (1). The major difference of this model from routing games is that the network is undirected, where edge $(i, j)$ can accommodate either the flow from node $j$ to $i$ or $i$ to $j$. The objective is to control the parameters $\{r_{ij}\}$ to reduce or increase the flows on some edges. The task of reducing edge flows has applications in power grid congestion mitigation in the direct-current (DC) approximation [42], where $r_{ij}$ is related to the reactance of edge $(i, j)$, controllable through devices in a flexible alternating current transmission system (FACTS) [43]. On the other hand, the task of increasing certain edge flows has been used to model the tunability of network functions, which is applicable in mechanical and biological networks [44] as well as learning machines in metamaterials [45].

As an example, we consider the task of flow control such that the relative increments of flows on the targeted edges $\mathcal{T}$ exceed a limit $\theta$ [44], i.e., $\rho_{ij} = \frac{|x_{ij}| - |x_{ij}^0|}{|x_{ij}^0|} - \theta \geq 0, \forall (i, j) \in \mathcal{T}$ (with $x_{ij}^0$ being the flow prior to tuning). It can be achieved by minimizing the hinge loss (upper-level objective) $\mathcal{O} = \sum_{(i,j)\in\mathcal{T}} -\rho_{ij}\Theta(-\rho_{ij}) =: \sum_{(i,j)\in\mathcal{T}} \mathcal{O}_{ij}$, where $\Theta(\cdot)$ is the Heaviside step function. The task of congestion mitigation in power grids can be studied similarly. We adopt the usual MP algorithm to compute the equilibrium flows as

$$C_{i\to j}(x_{ij}) = \min_{\{x_{ki}\}|R_i=0} \left[ \frac{1}{2} r_{ij} x_{ij}^2 + \sum_{k\in\mathcal{N}_i\setminus j} C_{k\to i}(x_{ki}) \right], \tag{4}$$

where $\mathcal{N}_i$ is the set of neighboring nodes adjacent to node $i$. The definition of the message $C_{i\to j}(x_{ij})$ differs from the one of Eq. (2) in that it includes the interaction term on edge $(i, j)$, which yields a more concise update rule here. Similar to Eq. (3), we approximate the message function by a quadratic form $C_{i\to j}(x_{ij}) = \frac{1}{2}\alpha_{i\to j}(x_{ij} - \hat{x}_{i\to j})^2 + \text{const}$, such that the optimization in Eq. (4) reduces to the computation of the real-valued messages $m_{i\to j} \in \{\alpha_{i\to j}, \hat{x}_{i\to j}\}$ by passing the upstream messages $\{m_{k\to i}\}_{k\in\mathcal{N}_i\setminus j}$, as illustrated on the left panel of Fig. 4(a) [30]. Upon convergence, the equilibrium flow $x_{ij}^*$ can be obtained by minimizing the edgewise full cost $C_{ij}^{\text{full}}(x_{ij}) = C_{i\to j}(x_{ij}) + C_{j\to i}(x_{ij}) - \frac{1}{2}r_{ij}x_{ij}^2$.

The variation of the control parameters $\{r_{ij}\}$ will impact on the messages $\{m_{i\to j}\}$, which in turn affects the equilibrium flows $\boldsymbol{x}^*$ and therefore the upper-level objective $\mathcal{O}(\boldsymbol{x}^*)$. Specifically, one considers the effect of the change of $r_{ij}$ on the targeted edge flows $\{x_{pq}^*\}_{(p,q)\in\mathcal{T}}$, derived by computing the gradient $\frac{\partial\mathcal{O}}{\partial m_{i\to j}}$. The targeted edges provide the boundary conditions as $\frac{\partial\mathcal{O}_{pq}}{\partial m_{p\to q}} =$

$\frac{\partial\mathcal{O}_{pq}}{\partial x_{pq}^*}\frac{\partial x_{pq}^*}{\partial m_{p\to q}}, \forall(p, q) \in \mathcal{T}$. As the messages from node $i$ to $j$ are functions of the upstream messages, i.e., $m_{i\to j} = m_{i\to j}(\{m_{k\to i}\}_{k\in\mathcal{N}_i\setminus j})$, the gradients on edge $i \to j$ are passed backward to its upstream edges $\{k \to i\}_{k\in\mathcal{N}_i\setminus j}$ through the chain rule, as illustrated in the right panel of Fig. 4(a). The full gradient on a non-targeted edge $k \to i$ can be obtained by summing the gradients on its downstream edges, computed as

$$\frac{\partial\mathcal{O}}{\partial m_{k\to i}} = \sum_{l\in\mathcal{N}_i\setminus k} \sum_{m_{i\to l}\in\{\alpha_{i\to l}, \hat{x}_{i\to l}\}} \frac{\partial\mathcal{O}}{\partial m_{i\to l}}\frac{\partial m_{i\to l}}{\partial m_{k\to i}}. \tag{5}$$

The gradient messages $\{\frac{\partial\mathcal{O}}{\partial m_{k\to i}}\}$ are passed in a random and asynchronous manner, resulting in a decentralized algorithm.

The gradient with respect to the control parameter on the non-targeted edge $(k, i)$ can be obtained straightforwardly as

$$\frac{\partial\mathcal{O}}{\partial r_{ki}} = \sum_{m\in\{\alpha,\hat{x}\}} \left( \frac{\partial\mathcal{O}}{\partial m_{k\to i}}\frac{\partial m_{k\to i}}{\partial r_{ki}} + \frac{\partial\mathcal{O}}{\partial m_{i\to k}}\frac{\partial m_{i\to k}}{\partial r_{ki}} \right), \tag{6}$$

which serves to update the control parameter in a gradient descent manner $r_{ki} \leftarrow r_{ki} - s\frac{\partial\mathcal{O}}{\partial r_{ki}}$ with certain step size $s$. The gradient for targeted edges can be similarly defined [30]. The control parameters are bounded to be $r_{ij} \in [0.9, 1.1]$, achieved by necessary thresholding after gradient descent updates. In this flow model, the gradient $\frac{\partial\mathcal{O}}{\partial r_{ki}}$ can be calculated exactly, leading to a global gradient descent (GGD) algorithm. However, the GGD approach requires computing the inverse of the Laplacian matrix in every iteration, which can be time-consuming for large networks. On the contrary, the gradients are computed in a local and distributed manner in the MP approach. Similar ideas of gradient propagation of MP have been proposed in Refs [46, 47] in the context of approximate inference, which are usually implemented centrally in the reversed order of MP updates, unlike the decentralized approach presented here.

The gradient computed by the MP algorithm provides an excellent estimation to the exact gradient, as illustrated in Fig. 4(b). For bilevel optimization, we do not wait for the convergence of the gradient-passing, but update the control parameters during the MP iterations to make the algorithm more efficient. It provides approximated gradient information, which is already effective for optimizing the global objective, as shown in Fig. 4(c). The MP approach yields similar success rates in managing the network flows for different thresholds compared to the GGD approach as shown in Fig. 4(d), demonstrating the effectiveness of the MP approach for the bilevel optimization.

In summary, we propose MP algorithms for solving bilevel optimization in flow networks, focusing on applications in the routing game and flow control problems.
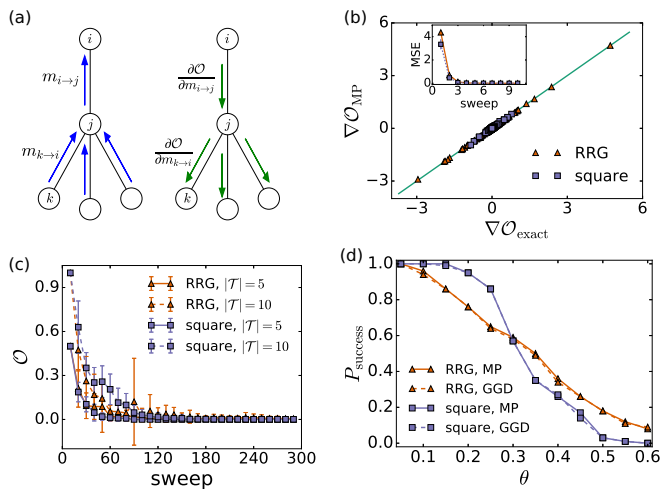
Figure 4. Bilevel optimization for flow control. An RRG ($N = 200$, degree 3) and a square lattice of size $15 \times 15$ are considered. The source and destination nodes, and the targeted edges are randomly selected. (a) Left: MP for solving the lower-level equilibrium flow problem. Right: Computing gradients of the upper-level objective function $\mathcal{O}$. (b) Comparison of the gradients at initial $r$ computed by the MP approach (obtained by fixing $r$ and passing messages $\{m_{i \to j}\}$ and gradients $\{\frac{\partial \mathcal{O}}{\partial m_{i \to j}}\}$) and the GGD approach, with $|\mathcal{T}| = 5, \theta = 0.1$. Inset: mean square error (MSE) of the gradients by the MP approach during iterations, in comparison to the GGD approach. Each sweep consists of $4|E|$ local MP steps. (c) MP for minimizing the upper-level objective function $\mathcal{O}$ with $\theta = 0.1$, where one randomly selected control parameter is updated following the descent direction every $4|E|/10$ steps. (d) Fraction of successfully tuned cases (satisfying $\mathcal{O} = 0$) $P_{\text{success}}$ out of 100 different problem realizations of source/destination nodes, with $|\mathcal{T}| = 5$, as a function of the threshold $\theta$.

In routing games, the objective functions in both levels admit a similar structure, which leads to two sets of similar messages being passed. Updates of the control variables based on localized information appear effective for toll optimization. However, the long-range impact of control-variable changes should be considered in some applications. This is accommodated by a separate distributed gradient-passing process, which is effective and efficient in flow control problems. Leveraging the sparse network structure, the MP approach offers efficient and intrinsically distributed algorithms in contrast to global optimization methods such as nonlinear programming, which is more generic, but is generally not scalable and therefore unsuitable for large-scale systems. The MP approach provides effective algorithms for bilevel optimization problems that are intractable or difficult to solve by global optimization approaches, such as combinatorial problems. We believe that these MP methods provide a valuable tool for solving difficult bilevel optimization problems, especially in systems with sparsely coupled structures.

Source codes of this work can be found in `https://github.com/boli8/bilevelMP_flow`.

[1] Ankur Sinha, Pekka Malo, and Kalyanmoy Deb, "A review on bilevel optimization: From classical to evolutionary approaches and applications," IEEE Transactions on Evolutionary Computation **22**, 276–295 (2018).
[2] J G Wardrop, "Some theoretical aspects of road traffic research." Proceedings of the Institution of Civil Engineers **1**, 325–362 (1952).
[3] Michael Plischke and Birger Bergersen, *Equilibrium Statistical Physics*, 3rd ed. (WORLD SCIENTIFIC, Singapore, 2006).
[4] William Thomson and Peter Guthrie Tait, *Treatise on Natural Philosophy*, 2nd ed., Cambridge Library Collection - Mathematics, Vol. 1 (Cambridge University Press, 2009).
[5] Peter G. Doyle and J. Laurie Snell, *Random Walks and Electric Networks* (Mathematical Association of America, 1984).
[6] E. T. Jaynes, "Information theory and statistical mechanics," Phys. Rev. **106**, 620–630 (1957).
[7] K.P. Murphy, *Machine Learning: A Probabilistic Perspective*, Adaptive Computation and Machine Learning series (MIT Press, Cambridge, Massachusetts, 2012).
[8] Benoît Colson, Patrice Marcotte, and Gilles Savard, "An overview of bilevel optimization," Annals of Operations Research **153**, 235–256 (2007).
[9] Robert G. Jeroslow, "The polynomial hierarchy and a simple model for competitive analysis," Mathematical Programming **32**, 146–164 (1985).
[10] Pierre Hansen, Brigitte Jaumard, and Gilles Savard, "New branch-and-bound rules for linear bilevel programming," SIAM Journal on Scientific and Statistical Computing **13**, 1194–1217 (1992).
[11] Jonathan F. Bard and James E. Falk, "An explicit solution to the multi-level programming problem," Comput-

ers & Operations Research **9**, 77–100 (1982).

[12] Jonathan F. Bard, "Convex two-level optimization," Mathematical Programming **40-40**, 15–27 (1988).

[13] C. D. Kolstad and L. S. Lasdon, "Derivative evaluation and computational experience with large bilevel mathematical programs," Journal of Optimization Theory and Applications **65**, 485–499 (1990).

[14] Gilles Savard and Jacques Gauvin, "The steepest descent direction for the nonlinear bilevel programming problem," Operations Research Letters **15**, 265–272 (1994).

[15] Raif M. Rustamov and James T. Klosowski, "Interpretable graph-based semi-supervised learning via flows," in *Proceedings of the 32nd AAAI Conference on Artificial Intelligence* (New Orleans, Louisiana USA, 2018) pp. 3976–3983.

[16] Montacer Essid and Justin Solomon, "Quadratically regularized optimal transport on graphs," SIAM Journal on Scientific Computing **40**, A1961–A1986 (2018).

[17] Gabriel Peyré and Marco Cuturi, "Computational optimal transport: With applications to data science," Foundations and Trends in Machine Learning **11**, 355–607 (2019).

[18] Michael Patriksson, *The traffic assignment problem : models and methods* (Dover Publications, Inc, New York, 2015).

[19] Tim Roughgarden, *Selfish routing and the price of anarchy* (MIT Press, Cambridge, Massachusetts, 2005).

[20] Martin J. Beckmann, C. B. McGuire, and C. B. Winsten, *Studies in the Economics of Transportation* (Yale University Press, New Haven, 1956).

[21] M.J. Smith, "The marginal cost taxation of a transportation network," Transportation Research Part B: Methodological **13**, 237–242 (1979).

[22] Richard Cole, Yevgeniy Dodis, and Tim Roughgarden, "How much can taxes help selfish routing?" Journal of Computer and System Sciences **72**, 444–467 (2006), network Algorithms 2005.

[23] Serdar Çolak, Antonio Lima, and Marta C. González, "Understanding congested travel in urban areas," Nature Communications **7**, 10793 (2016), article.

[24] "Electronic road pricing," `https://web.archive.org/web/20110605101108/http://www.lta.gov.sg/motoring_matters/index_motoring_erp.htm` (2019).

[25] Naama Barak, "Israel tries battling traffic jams with cash handouts," ISRAEL21c, `https://www.israel21c.org/israel-tries-battling-traffic-jams-with-cash-handouts/` (2019).

[26] Javier Alonso-Mora, Samitha Samaranayake, Alex Wallar, Emilio Frazzoli, and Daniela Rus, "On-demand high-capacity ride-sharing via dynamic trip-vehicle assignment," Proceedings of the National Academy of Sciences **114**, 462–467 (2017).

[27] Sejoon Lim, Hari Balakrishnan, David Gifford, Samuel Madden, and Daniela Rus, "Stochastic motion planning and applications to traffic," The International Journal of Robotics Research **30**, 699–712 (2011).

[28] Bo Li, David Saad, and Andrey Y. Lokhov, "Reducing urban traffic congestion due to localized routing decisions," Phys. Rev. Research **2**, 032059 (2020).

[29] George Karakostas and Stavros G. Kolliopoulos, "The efficiency of optimal taxes," in *Combinatorial and Algorithmic Aspects of Networking*, edited by Alejandro López-Ortiz and Angèle M. Hamel (Springer Berlin Heidelberg, Berlin, Heidelberg, 2005) pp. 3–12.

[30] See Supplemental Material for details, which includes Refs [48–56].

[31] Dov Monderer and Lloyd S. Shapley, "Potential games," Games and Economic Behavior **14**, 124–143 (1996).

[32] Hillel Bar-Gera, "Origin-based algorithm for the traffic assignment problem," Transportation Science **36**, 398–417 (2002).

[33] K. Y. Michael Wong and David Saad, "Inference and optimization of real edges on sparse graphs: A statistical physics perspective," Phys. Rev. E **76**, 011115 (2007).

[34] K. Y. Michael Wong, David Saad, and Chi Ho Yeung, "Distributed optimization in transportation and logistics networks," IEICE Transactions on Communications **E99.B**, 2237–2246 (2016).

[35] Martin Hoefer, Lars Olbrich, and Alexander Skopalik, "Taxing subnetworks," in *Internet and Network Economics*, edited by Christos Papadimitriou and Shuzhong Zhang (Springer Berlin Heidelberg, Berlin, Heidelberg, 2008) pp. 286–294.

[36] R. W. Rosenthal, "The network equilibrium problem in integers," Networks **3**, 53–59 (1973).

[37] Chi Ho Yeung and David Saad, "Competition for shortest paths on sparse graphs," Phys. Rev. Lett. **108**, 208701 (2012).

[38] Chi Ho Yeung, "Efficient algorithm for routing optimization via statistical mechanics," in *2013 IEEE International Conference on Communications Workshops (ICC)* (2013) pp. 1420–1424.

[39] C. De Bacco, S. Franz, D. Saad, and C. H. Yeung, "Shortest node-disjoint paths on random graphs," Journal of Statistical Mechanics: Theory and Experiment **2014**, P07009 (2014).

[40] Chi Ho Yeung, David Saad, and K. Y. Michael Wong, "From the physics of interacting polymers to optimizing routes on the london underground," Proceedings of the National Academy of Sciences **110**, 13717–13722 (2013).

[41] Ho Fai Po, Chi Ho Yeung, and David Saad, "Futility of being selfish in optimized traffic," Phys. Rev. E **103**, 022306 (2021).

[42] A.J. Wood, B.F. Wollenberg, and G.B. Sheblé, *Power Generation, Operation, and Control* (Wiley, Hoboken, New Jersey, 2013).

[43] Xiao-Ping Zhang, Christian Rehtanz, and Bikash Pal, *Flexible AC Transmission Systems: Modelling and Control* (Springer, Berlin, Heidelberg, 2006).

[44] Jason W. Rocks, Henrik Ronellenfitsch, Andrea J. Liu, Sidney R. Nagel, and Eleni Katifori, "Limits of multifunctionality in tunable networks," Proceedings of the National Academy of Sciences **116**, 2506–2511 (2019).

[45] Menachem Stern, Daniel Hexner, Jason W. Rocks, and Andrea J. Liu, "Supervised learning in physical networks: From machine learning to learning machines," Phys. Rev. X **11**, 021045 (2021).

[46] Frederik Eaton and Zoubin Ghahramani, "Choosing a variable to clamp," in *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, Vol. 5, edited by David van Dyk and Max Welling (PMLR, Hilton Clearwater Beach Resort, Clearwater Beach, Florida USA, 2009) pp. 145–152.

[47] Justin Domke, "Learning graphical model parameters with approximate marginal inference," IEEE Transactions on Pattern Analysis and Machine Intelligence **35**, 2454–2467 (2013).

[48] Alessandro Lonardi, Enrico Facca, Mario Putti, and Caterina De Bacco, "Designing optimal networks for multicommodity transport problem," Phys. Rev. Research **3**, 043010 (2021).

[49] Vincenzo Bonifaci, Enrico Facca, Frederic Folz, Andreas Karrenbauer, Pavel Kolev, Kurt Mehlhorn, Giovanna Morigi, Golnoosh Shahkarami, and Quentin Vermande, "Physarum-inspired multi-commodity flow dynamics," Theoretical Computer Science **920**, 1–20 (2022).

[50] Joaquim Dias Garcia, Guilherme Bodin, davide-f, Ian Fiske, Mathieu Besançon, and Nick Laws, "joaquimg/bileveljump.jl: v0.4.1," https://doi.org/10.5281/zenodo.4556393 (2021), 10.5281/zenodo.4556393.

[51] Chi Ho Yeung, K. Y. Michael Wong, and Bo Li, "Coverage versus supply cost in facility location: Physics of frustrated spin systems," Phys. Rev. E **89**, 062805 (2014).

[52] Christopher M. Bishop, *Pattern Recognition and Machine Learning* (Springer-Verlag, Berlin, Heidelberg, 2006).

[53] Chaisak Suwansirikul, Terry L. Friesz, and Roger L. Tobin, "Equilibrium decomposed optimization: A heuristic for the continuous equilibrium network design problem," Transportation Science **21**, 254–263 (1987).

[54] Chi Ho Yeung and K. Y. Michael Wong, "Optimal location of sources in transportation networks," Journal of Statistical Mechanics: Theory and Experiment **2010**, P04017 (2010).

[55] Patrick Rebeschini and Sekhar Tatikonda, "A new approach to laplacian solvers and flow problems," Journal of Machine Learning Research **20**, 1–37 (2019).

[56] G. H. Golub and V. Pereyra, "The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate," SIAM Journal on Numerical Analysis **10**, 413–432 (1973).

# Bilevel Optimization in Flow Networks – A Message-passing Approach – Supplemental Material

Bo Li,[1,2] David Saad,[1] and Chi Ho Yeung[3]

[1] *Non-linearity and Complexity Research Group, Aston University, Birmingham, B4 7ET, United Kingdom*
[2] *School of Science, Harbin Institute of Technology (Shenzhen), Shenzhen, 518055, China*
[3] *Department of Science and Environmental Studies,*
*The Education University of Hong Kong, 10 Lo Ping Road, Taipo, Hong Kong*

## I. MESSAGE-PASSING ALGORITHMS FOR NON-ATOMIC ROUTING GAMES

In this section, we provide details of the message-passing (MP) algorithm for non-atomic routing games in road networks, modeled by a directed graph $G(V, E)$. We denote $\mathbb{R}, \mathbb{Z}, \mathbb{N}$ as the sets of real numbers, integers and non-negative integers, respectively.

### A. Problem Setting and Notation for Non-atomic Games

A directed edge $e$ in the directed graph is represented by an ordered tuple $e = (i, j)$, where node $i$ is the head and node $j$ is the tail of edge $e$, i.e., $i = h(e), j = t(e)$. Note that there can be at most two directed edges connecting node $i$ and node $j$, i.e., $e = (i, j), e' = (j, i)$.

We write the set of incoming edges to node $i$ as $\partial_i^{\mathrm{in}} = \{e | e \in E, t(e) = i\}$, the set of outgoing edges from node $i$ as $\partial_i^{\mathrm{out}} = \{e | e \in E, h(e) = i\}$ and the set of edges adjacent to node $i$ as $\partial i = \partial_i^{\mathrm{in}} \cup \partial_i^{\mathrm{out}}$. For convenience, we define the incident operator $B : E \to V$, with matrix elements

$$B_{i,e} = \begin{cases} 1, & \text{if } e \in \partial_i^{\mathrm{in}} \\ -1, & \text{if } e \in \partial_i^{\mathrm{out}} \\ 0, & \text{otherwise.} \end{cases} \tag{S1}$$

Consider the scenario where all users travel to a universal destination $\mathcal{D}$. The edge flows $\{x_e \in \mathbb{R}\}$ resulting from users' path choices satisfy the flow conservation constraints,

$$R_i := \Lambda_i + \sum_{e \in \partial_i^{\mathrm{in}}} x_e - \sum_{e \in \partial_i^{\mathrm{out}}} x_e$$
$$= \Lambda_i + \sum_{e \in \partial i} B_{i,e} x_e = 0, \quad \forall i \neq \mathcal{D}, \tag{S2}$$

and the non-negativity constraints

$$x_e \geq 0, \quad \forall e. \tag{S3}$$

Due to the flow conservation constraint, any resource on a leaf node $i$ with only one outgoing edge (i.e., $|\partial_i^{\mathrm{out}}| = 1, |\partial_i^{\mathrm{in}}| = 0$) must be transmitted to its only neighboring node $j$. Similarly, if a leaf node $i$ with only one incoming edge (i.e., $|\partial_i^{\mathrm{in}}| = 1, |\partial_i^{\mathrm{out}}| = 0$) is the destination node, then traffic must first arrive at its only neighboring node $j$, and then go through the edge $(j, i)$ to the destination. In the former case, one can remove the leaf node $i$ and add $\Lambda_i$ resources to its neighboring node $j$. In the latter, one can simply set node $j$ as the destination. By preprocessing the network using the above reduction, we can reduce the network to have no leaf nodes.

Denoting $\ell_e(x_e)$ as the latency function on edge $e$ (assumed to be a non-decreasing function of $x_e$) and $\tau_e$ to be the corresponding toll, the Wardrop equilibrium can be obtained by minimizing the following potential function

$$\Phi(\boldsymbol{x}) = \sum_{e \in E} \int_0^{x_e} \big[\ell_e(y) + \tau_e\big] \mathrm{d}y =: \sum_{e \in E} \phi_e(x_e), \tag{S4}$$

subject to the flow conservation Eq. (S2) and non-negativity constraints Eq. (S3). We have assumed the same gauge between latency and toll can be used for all users (more precisely the edge cost for a user is $\ell_e(x_e) + \chi \tau_e$ with $\chi$ being a coefficient converting money to time which is set to one in Eq. (S4) in some appropriate unit).
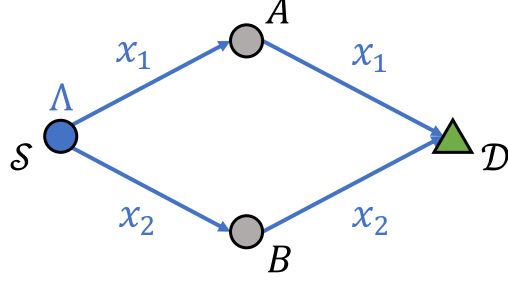
Figure S1. A small network with 4 nodes. There are $\Lambda > 0$ (where $\Lambda \in \mathbb{R}$) users originating from the source node $\mathcal{S}$ to the destination node $\mathcal{D}$, where $x_1$ users choose the path $\mathcal{P}_1 = ((\mathcal{S}, A), (A, \mathcal{D}))$ and $x_2$ users choose the path $\mathcal{P}_2 = ((\mathcal{S}, B), (B, \mathcal{D}))$.

The social cost of the routing game is defined as

$$H(\boldsymbol{x}) = \sum_{e \in E} x_e \ell_e(x_e) =: \sum_{e \in E} \sigma_e(x_e), \tag{S5}$$

where the corresponding minimizer is the social optimum. Tolls are not assumed to contribute to the social cost $H(\boldsymbol{x})$.

The network planners only need to know the aggregated network flow on each edge and the corresponding latency in order to determine the social cost and set the tolls, while the specific paths where users choose are not that relevant. For this reason, we do not address the problem of finding individual routes for each user. The individual routing problem can be tackled by a multi-commodity formalism [1, 2], or by using some physics-inspired algorithms [3].

## B. Intuition of the potential function $\Phi(\boldsymbol{x})$

To gain some intuition of the role of the potential function $\Phi(\boldsymbol{x})$ in finding the Nash equilibrium, we consider the following simple scenario as shown in Fig. S1, where there are $\Lambda > 0$ amount of users originating from the source node $\mathcal{S}$ to the destination node $\mathcal{D}$, and there are only two non-overlapping paths from node $\mathcal{S}$ to node $\mathcal{D}$, which are $\mathcal{P}_1 = ((\mathcal{S}, A), (A, \mathcal{D}))$ and $\mathcal{P}_2 = ((\mathcal{S}, B), (B, \mathcal{D}))$. Denote $x_1$ and $x_2$ as the flow on path $\mathcal{P}_1$ and path $\mathcal{P}_2$ respectively. The flow conservation constraint asserts that $\Lambda = x_1 + x_2$.

Users choosing $\mathcal{P}_1$ experience a cost $C_1(x_1) = \ell_{\mathcal{S}A}(x_1) + \tau_{\mathcal{S}A} + \ell_{A\mathcal{D}}(x_1) + \tau_{A\mathcal{D}}$, while users choosing $\mathcal{P}_2$ experience a cost $C_2(x_2) = \ell_{\mathcal{S}B}(x_2) + \tau_{\mathcal{S}B} + \ell_{B\mathcal{D}}(x_2) + \tau_{B\mathcal{D}}$. There are 3 possible scenarios of the Wardrop (Nash) equilibrium $x_1^*, x_2^*$, depending on the network parameters:

$$
\begin{array}{llll}
\text{Case I:} & C_1(x_1^*) < C_2(x_2^*), & x_1^* > 0, x_2^* = 0, \\
\text{Case II:} & C_1(x_1^*) > C_2(x_2^*), & x_1^* = 0, x_2^* > 0, \\
\text{Case III:} & C_1(x_1^*) = C_2(x_2^*), & x_1^* > 0, x_2^* > 0.
\end{array} \tag{S6}
$$

In Case I, all users choose $\mathcal{P}_1$ and there is no incentive for them to move to $\mathcal{P}_2$ as the corresponding cost $C_2(x_2 = 0)$ is higher. Similar analysis applies to Case II. In Case III, a user (controlling $\mathrm{d}x$ amount of traffic) choosing $\mathcal{P}_1$ also has no incentive to switch to $\mathcal{P}_2$; if she did so, her cost will become $C_2(x_2 + \mathrm{d}x) \geq C_2(x_2) = C_1(x_1)$, which is unfavorable.

Now we turn to the optimization problem as stated above

$$\min_{x_1, x_2} \Phi(x_1, x_2) = \int_0^{x_1} \left[ \ell_{\mathcal{S}A}(y) + \tau_{\mathcal{S}A} \right] \mathrm{d}y + \int_0^{x_1} \left[ \ell_{A\mathcal{D}}(y) + \tau_{A\mathcal{D}} \right] \mathrm{d}y \tag{S7}$$

$$+ \int_0^{x_2} \left[ \ell_{\mathcal{S}B}(y) + \tau_{\mathcal{S}B} \right] \mathrm{d}y + \int_0^{x_2} \left[ \ell_{B\mathcal{D}}(y) + \tau_{B\mathcal{D}} \right] \mathrm{d}y$$

$$\text{s. t. } x_1 \geq 0, x_x \geq 0,$$

$$x_1 + x_2 = \Lambda,$$

which can be solved by extremizing the Lagrangian function $\mathcal{L} = \Phi(x_1, x_2) - \mu(x_1 + x_2 - \Lambda) - \nu_1 x_1 - \nu_2 x_2$, together with the Karush-Kuhn-Tucker (KKT) conditions $\nu_1, \nu_2 \geq 0, \nu_1 x_1 = 0, \nu_2 x_2 = 0$. The extremum $\frac{\mathrm{d}\mathcal{L}}{\mathrm{d}x_1} = \ell_{\mathcal{S}A}(x_1) +

$\tau_{SA} + \ell_{AD}(x_1) + \tau_{AD} - \mu - \nu_1 = 0$ yields $C_1(x_1) = \mu + \nu_1$. Similarly, $C_2(x_2) = \mu + \nu_2$. The optimal dual parameters $(\mu^*, \nu_1^*, \nu_2)$ has 3 possible scenarios under the KKT condition

$$
\begin{array}{lll}
\text{Case I:} & \nu_1^* = 0, \nu_2^* > 0, & \\
\text{Case II:} & \nu_1^* > 0, \nu_2^* = 0, & \\
\text{Case III:} & \nu_1^* = 0, \nu_2^* = 0, &
\end{array}
\tag{S8}
$$

which exactly corresponds to the Wardrop equilibrium stated in Eq. (S6). Therefore, we have established that in this simple example, the Wardrop equilibrium can be identified by minimizing $\Phi(\boldsymbol{x})$.

We would also like to emphasize that the potential function $\Phi(\boldsymbol{x})$ does not carry any meaningful information about the routing costs (of either individuals of their aggregation), which can be seen by its definition. Therefore, the magnitude of $\Phi(\boldsymbol{x})$ is *not* a performance measure of the traffic network. It plays an auxiliary role in defining the equilibrium flow as $\boldsymbol{x}^* = \arg\min_{\boldsymbol{x}} \Phi(\boldsymbol{x})$ s.t. Eq. (S2) and Eq. (S3). Roughly speaking, in the lower-level optimization problem, only $\arg\min_{\boldsymbol{x}} \Phi(\boldsymbol{x})$ carries a correspondence to the routing problem, but not $\min_{\boldsymbol{x}} \Phi(\boldsymbol{x})$.

On the other hand, the social cost $H(\boldsymbol{x})$ corresponds to the aggregated travel latency of all users, which is a useful measure of the routing condition on a road network.

## C. MP Equations for Smooth Message Functions

The MP equation for minimizing the potential $\Phi(\boldsymbol{x})$ reads

$$
\Phi_{i \to e}(x_e) = \min_{\{x_{e'} \geq 0\}|R_i = 0} \sum_{e' \in \partial i \backslash e} \left[ \Phi_{e' \to i}(x_{e'}) + \phi_{e'}(x_{e'}) \right],
\tag{S9}
$$

where the message function $\Phi_{i \to e}(x_e)$ is called the cavity energy in the jargon of statistical physics. Denoting $e = (i,j)$ and $e' = (k,i)$, we can write $\Phi_{e' \to i}(x_{e'}) = \Phi_{k \to e'}(x_{e'})$.

In this framework, one needs to keep track of the profile of the message functions $\Phi_{i \to e}(x_e)$, which is only practical if they are restricted to a certain family of functions and are easy to optimize. One can approximate the message function $\Phi_{i \to e}(x_e)$ by its series expansion around the working point $\tilde{x}_{i \to e}$ [4]

$$
\begin{aligned}
\Phi_{i \to e}(x_e) &= \Phi_{i \to e}(\tilde{x}_{i \to e} + \varepsilon_e) \\
&\approx \Phi_{i \to e}(\tilde{x}_{i \to e}) + \beta_{i \to e}|_{\tilde{x}_{i \to e}} \cdot \varepsilon_e + \frac{1}{2} \alpha_{i \to e}|_{\tilde{x}_{i \to e}} \cdot \left(\varepsilon_e\right)^2,
\end{aligned}
\tag{S10}
$$

where $\beta_{i \to e}$ and $\alpha_{i \to e}$ are the first and second derivatives of $\Phi_{i \to e}$ evaluated at the working point $\tilde{x}_{i \to e}$, assuming the message function $\Phi_{i \to e}(x_e)$ is smooth in the vicinity of $\tilde{x}_{i \to e}$. The MP equations have been derived in [4] for undirected flow networks. Here, we extend it to directed graph with non-negativity flow constraints.

Similarly, the interaction term $\phi_{e'}(x_{e'})$ is also approximated as $\phi_{e'}(x_{e'}) \approx \phi_{e'}(\tilde{x}_{i \to e}) + \phi_{e'}'(\tilde{x}_{k \to e'})\varepsilon_e + \frac{1}{2}\phi_{e'}''(\tilde{x}_{k \to e'})\left(\varepsilon_e\right)^2$. To solve the local optimization problem in Eq. (S9) over the variables on edges $\{k \to e'|e' \in \partial i \backslash e\}$, we introduce the Lagrangian

$$
\begin{aligned}
L_{i \to e} &= \sum_{e' \in \partial i \backslash e} \left[ \frac{1}{2}\alpha_{k \to e'}\left(\varepsilon_{e'}\right)^2 + \beta_{k \to e'}\varepsilon_{e'} + \frac{1}{2}\phi_{e'}''(\tilde{x}_{k \to e'})\left(\varepsilon_{e'}\right)^2 + \phi_{e'}'(\tilde{x}_{k \to e'})\varepsilon_{e'} \right] \\
&\quad + \mu_{i \to e}R_i + \sum_{e' \in \partial i \backslash e} \lambda_{e'}(\tilde{x}_{k \to e'} + \varepsilon_{e'}),
\end{aligned}
\tag{S11}
$$

where $\mu_{i \to e}$ and $\lambda_{e'}$ are the Lagrange multipliers for the flow conservation constraint $R_i = 0$ and flow non-negativity constraint $x_{e'} \geq 0$, respectively. Solving the extremum equation $\frac{\partial L_{i \to e}}{\partial \varepsilon_{e'}} = 0$ gives

$$
\varepsilon_{e'}^*(\mu_{i \to e}) = \max\left( \frac{-1}{\alpha_{k \to e'} + \phi_{e'}''} \left(\mu_{i \to e}B_{i,e'} + \phi_{e'}' + \beta_{k \to e'}\right), -\tilde{x}_{k \to e'} \right),
\tag{S12}
$$

and the corresponding optimal cavity flow is

$$
x_{k \to e'}^*(\mu_{i \to e}) = \tilde{x}_{k \to e'} + \varepsilon_{e'}^*(\mu_{i \to e}) = \max\left( \tilde{x}_{k \to e'} - \frac{\mu_{i \to e}B_{i,e'} + \phi_{e'}' + \beta_{k \to e'}}{\alpha_{k \to e'} + \phi_{e'}''}, 0 \right).
\tag{S13}
$$

The Lagrange multiplier (or the dual variable) $\mu_{i \to e}$ needs to satisfy

$$R_{i \to e}(\mu_{i \to e}; x_e) := \sum_{e' \in \partial i \setminus e} B_{i,e'} x^*_{k \to e'}(\mu_{i \to e}) + B_{i,e} x_e + \Lambda_i = 0. \tag{S14}$$

The function $R_{i \to e}(\mu; x_e)$ is a non-increasing piece-wise linear function of $\mu$. To determine the value of $\mu^*_{i \to e}$ at the optimum, we need to find the root of $R_{i \to e}(\mu; x_e)$, which can be done in finite steps by following the breakpoints of the piece-wise linear function $R_{i \to e}(\mu; x_e)$. Upon obtaining the optimal dual variable $\mu^*_{i \to e}$, the messages $\beta_{i \to e}$ and $\alpha_{i \to e}$ are calculated by

$$\beta_{i \to e} = \frac{\partial \Phi^*_{i \to e}(x_e)}{\partial x_e} = \frac{\partial L^*_{i \to e}}{\partial x_e} = B_{i,e} \mu^*_{i \to e}, \tag{S15}$$

$$\begin{aligned}
\alpha_{i \to e} &= \frac{\partial^2 \Phi^*_{i \to e}(x_e)}{\partial x_e^2} = B_{i,e} \frac{\partial \mu^*_{i \to e}}{\partial x_e} = B_{i,e} \left[ \left. \frac{\partial x_e}{\partial \mu} \right|_{\mu = \mu^*_{i \to e}} \right]^{-1} \\
&= -\left[ \left. \frac{\partial}{\partial \mu} \sum_{e' \in \partial i \setminus e} B_{i,e'} x^*_{k \to e'}(\mu) \right|_{\mu = \mu^*_{i \to e}} \right]^{-1} \\
&= \left[ \sum_{e' \in \partial i \setminus e} \frac{1}{\alpha_{k \to e'} + \phi''_{e'}} \Theta \left( (\alpha_{k \to e'} + \phi''_{e'}) \tilde{x}_{k \to e'} - (\mu^*_{i \to e} B_{i,e'} + \phi'_{e'} + \beta_{k \to e'}) \right) \right]^{-1}, \tag{S16}
\end{aligned}$$

where $\Theta(\cdot)$ is the Heaviside step function. The shadow price interpretation of Lagrangian multiplier has been used in Eq. (S15) and the inverse function theorem has been used in Eq. (S16). In the implementation of the algorithm, we take $x_e = \tilde{x}_{i \to e}$ in solving Eq. (S14).

### 1. Destination node $\mathcal{D}$

There are two ways to treat the destination node $\mathcal{D}$:

- Method I: Since the destination node $\mathcal{D}$ has no constraint, it will absorb all incoming flows (like a grounded node in an electric circuit). So it has no preference for network flows of the incident edges, such that $\Phi_{\mathcal{D} \to e}(x_e) = 0$ and

$$\alpha_{\mathcal{D} \to e} = 0, \quad \beta_{\mathcal{D} \to e} = 0. \tag{S17}$$

- Method II: Alternatively, one can set an explicit constraint on the flows to the destination node $\mathcal{D}$

$$R_{\mathcal{D}} := \Lambda_{\mathcal{D}} + \sum_{e \in \partial \mathcal{D}} B_{\mathcal{D},e} x_e = 0, \tag{S18}$$

where $\Lambda_{\mathcal{D}} = -\sum_{i \neq \mathcal{D}} \Lambda_i$. Then the messages from the destination node $\mathcal{D}$ are calculated in the same way as other nodes given by Eqs. (S15) and (S16).

Method I is used for the experiments of routing games in the main text.

### 2. Working Points

We also need a scheme to update the the working points $\{\tilde{x}_{i \to e}\}$ at which the messages $\{\alpha_{i \to e}, \beta_{i \to e}\}$ are defined. Here, we suggest to update the working point $\tilde{x}_{i \to e}$ such that it gets closer to the equilibrium flow $x_e^*$ [4]

$$\begin{aligned}
x_e^* &= \arg \min_{x_e \geq 0} \left[ \Phi_{i \to e}(x_e) + \Phi_{j \to e}(x_e) + \phi_e(x_e) \right] \\
&= \arg \min_{x_e \geq 0} \left[ \frac{1}{2} \left( \alpha_{i \to e} + \frac{1}{2} \phi''_e(\tilde{x}_{i \to e}) \right) (x_e - \tilde{x}_{i \to e})^2 + \left( \beta_{i \to e} + \frac{1}{2} \phi'_e(\tilde{x}_{i \to e}) \right) (x_e - \tilde{x}_{i \to e}) \right. \\
&\quad \left. + \frac{1}{2} \left( \alpha_{j \to e} + \frac{1}{2} \phi''_e(\tilde{x}_{j \to e}) \right) (x_e - \tilde{x}_{j \to e})^2 + \left( \beta_{j \to e} + \frac{1}{2} \phi'_e(\tilde{x}_{j \to e}) \right) (x_e - \tilde{x}_{j \to e}) \right] \\
&= \max \left( \frac{(\alpha_{i \to e} + \frac{1}{2} \phi''_{i \to e}) \tilde{x}_{i \to e} + (\alpha_{j \to e} + \frac{1}{2} \phi''_{j \to e}) \tilde{x}_{j \to e} - (\beta_{i \to e} + \beta_{j \to e} + \frac{1}{2} \phi'_{i \to e} + \frac{1}{2} \phi'_{j \to e})}{\alpha_{i \to e} + \alpha_{j \to e} + \frac{1}{2} \phi''_{i \to e} + \frac{1}{2} \phi''_{j \to e}}, 0 \right). \tag{S19}
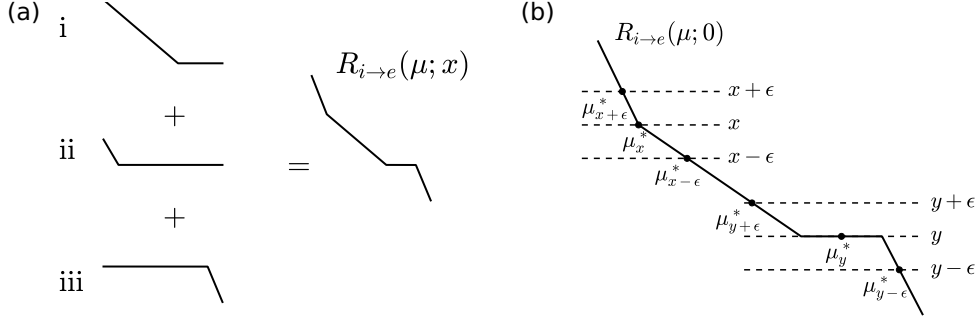\end{aligned}$$

Figure S2. (a) The net resource $R_{i \to e}(\mu; x)$ defined in Eq. (S14) is a non-increasing piecewise-linear function of the Lagrange multiplier $\mu$. Cases (i) and (ii) correspond to edges $e'$ incoming to node $i$ ($B_{i,e'} = 1$), while case (iii) corresponds to edge $e'$ outgoing of node $i$ ($B_{i,e'} = -1$). (b) The roots of $R_{i \to e}(\mu; x_e)$ in the vicinity of $x_e = x$ and $x_e = y$. It is assumed that edge $e$ is an outgoing edge of node $i$ (with $B_{i,e} = -1$) such that finding the root of $R_{i \to e}(\mu; x_e)$ is equivalent to solving $R_{i \to e}(\mu; 0) = x_e$. For infinitesimal $\epsilon$, if the flow $x_e$ changes from $x + \epsilon$ to $x - \epsilon$, the solution of the Lagrange multiplier changes continuously from $\mu^*_{x+\epsilon}$ to $\mu^*_{x-\epsilon}$. On the other hand, there is a plateau at $R_{i \to e}(\mu; 0) = y$, so that when the flow $x_e$ changes from $y + \epsilon$ to $y - \epsilon$, the solution of the Lagrange multiplier changes discontinuously from $\mu^*_{y+\epsilon}$ to $\mu^*_{y-\epsilon}$.

Furthermore, a learning rate $s$ is applied to update the working point

$$\tilde{x}^{\text{new}}_{i \to e} \leftarrow s x^*_e + (1 - s)\tilde{x}^{\text{old}}_{i \to e}, \tag{S20}$$

such that $\tilde{x}_{i \to e}$ does not jump too drastically; otherwise the messages $\alpha_{i \to e}$ and $\beta_{i \to e}$ will approximate the curvature and slope of the message function $\Phi_{i \to e}(x_e)$ less precisely.

## D. Non-Smooth Message Functions

### 1. Qualitative Picture

The MP algorithms in Sec. I C work well if the smoothness assumption of the message function $\Phi_{i \to e}(x_e)$ holds. However, it is not always the case in the routing game problem, where the non-smoothness is induced by the non-negativity constraints of Eq. (S3). Direct implementation of the MP algorithms in Sec. I C leads to oscillations of the messages when the traffic patterns are sparse. In fact, similar non-convergence phenomena have been noticed in the system with a non-smooth energy function [4].

To better understand this phenomenon, we examine $R_{i \to e}(\mu; x_e)$ as a function of the Lagrange multiplier $\mu$ in Eq. (S14), of which the root $\mu^*$ (satisfying $R_{i \to e}(\mu^*, x_e) = 0$) will determine $\beta_{i \to e}$ and $\alpha_{i \to e}$ in Eqs. (S15) and (S16). The function $R_{i \to e}(\mu; x_e)$ is non-increasing piecewise-linear function as illustrated in Fig. S2(a). Assuming edge $e$ is an outgoing edge of node $i$ (with $B_{i,e} = -1$), finding the root of $R_{i \to e}(\mu; x_e)$ is equivalent to solving $R_{i \to e}(\mu; 0) = x_e$. Consider the configuration in Fig. S2(b), where the solution of $R_{i \to e}(\mu; 0) = y$ occurs at a plateau, such that the solution (denoted as $\mu^*_y$) is degenerate; when the flow $x_e$ changes infinitesimally from $y + \epsilon$ to $y - \epsilon$, the solution of the Lagrange multiplier changes discontinuously from $\mu^*_{y+\epsilon}$ to $\mu^*_{y-\epsilon}$. In this case, the slope $\beta_{i \to e}$ of the cavity energy $\Phi_{i \to e}(x_e)$ changes discontinuously from $x_e = y + \epsilon$ to $x_e = y - \epsilon$, while the curvature $\alpha_{i \to e}$ is ill-defined at $x_e = y$. The profiles of the message function $\Phi_{i \to e}(x_e)$ in the smooth and non-smooth cases are illustrated in Fig. S3. If the normal messages $\{\beta_{i \to e}, \alpha_{i \to e}\}$ are used when the message function $\Phi_{i \to e}(x_e)$ is non-smooth, the solution will be jumping between the two branches, resulting in non-convergence behaviors of the MP algorithms as observed in Ref. [4].

### 2. Criteria for Non-smooth Message Function

As mentioned above, the message function $\Phi_{i \to e}(x_e)$ is non-smooth if the solution of $\mu$ in Eq. (S14) is degenerate. This occurs if the optimal flow $x^*_{k \to e'}(\mu_{i \to e})$ of all descendant edges $e' \in \partial i \backslash e$ are inactive, i.e., lying in the zero branch of the function in Eq. (S13); when $B_{i,e}x_e + \Lambda_i = 0$, the flow conservation equation $R_{i \to e}(\mu; x_e) = \sum_{e' \in \partial i \backslash e} B_{i,e'} x^*_{k \to e'}(\mu) = 0$ has degenerate solutions. In this case, all the resources $\Lambda_i$ are transmitted along edge $e$, while the flows on all other
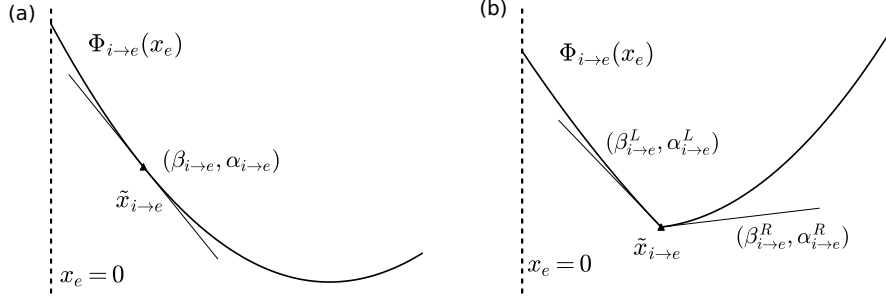
Figure S3. (a) Smooth message function $\Phi_{i\to e}(x_e)$, corresponding to $x_e = x$ in Fig. S2(b). (b) Non-smooth message function $\Phi_{i\to e}(x_e)$ with one breakpoint, corresponding to $x_e = y$ in Fig. S2(b), where the first and second derivatives of $\Phi_{i\to e}(x_e)$ are discontinuous near $x_e = y$.
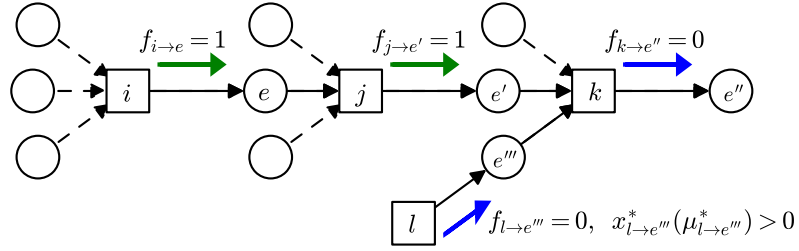


Figure S4. Illustration of the effective leaf edges. Arrows with dashed lines correspond to edges with zero optimal cavity flow $x^*(\mu^*)$ in the MP calculation (expression given in Eq. (S13)) of one of its downstream edges. Edge $i \to e$ is a primary effective leaf (assuming $\Lambda_i > 0$), as all the upstream optimal flows are zero. Edge $j \to e'$ is a general effective leaf, as its upstream edges are either effective leaves or attain zero optimal cavity flows. Edge $k \to e''$ is a non-effective leaf, as its upstream edge $l \to e'''$ is a non-effective leaf and it has a non-zero optimal flow $x^*_{l\to e'''}(\mu^*_{l\to e'''}) > 0$.

edges $\partial i \backslash e$ adjacent to node $i$ are idle. When $\Lambda_i > 0$, edge $i \to e$ is a leaf in the subgraph with edges holding non-zero flows, therefore we call edge $i \to e$ a primary effective leaf in such cases. Since $\Lambda_i \geq 0$ and $x_e \geq 0$, only the out-going edge $i \to e$ from node $i$ (with $B_{i,e} = -1$) can be a primary effective leaf.

The leaf state can also propagate from primary effective leaves to downstream edges. We define edge $i \to e$ to be a general effective leaf if and only if $\forall e' \in \partial i \backslash e$, either (i) the optimal flow $x^*_{k\to e'}(\mu^*_{i\to e}) = 0$ in Eq. (S13) or (ii) edge $k \to e'$ is a general effective leaf. A primary leaf is by default a general effective leaf. If an edge $i \to e$ is an effective leaf, we denote $f_{i\to e} = 1$, otherwise $f_{i\to e} = 0$. An example of effective leaf configurations is shown in Fig S4.

It can be proved by contradiction that only out-going edges $i \to e$ with $B_{i,e} = -1$ can be general effective leaves under the condition $\Lambda_i \geq 0$. The set of effective leaves in the upstream of edge $i \to e$ is $EL_{i\to e} = \{e'|e' \in \partial i \backslash e, f_{k\to e'} = 1\}$, while the set of non-effective leaves is $NEL_{i\to e} = \{e'|e' \in \partial i \backslash e, f_{k\to e'} = 0\}$.

Since there is at most one plateau in the function $R_{i\to e}(\mu; x)$, the cavity message $\Phi_{i\to e}(x_e)$ has at most one breakpoint. For an effective leaf edge $i \to e$, we always use the breakpoint of the message function (denoted as $\tilde{x}^b_{i\to e}$) as the working point, such that $\Phi_{i\to e}(x_e)$ has the following expression

$$\Phi_{i\to e}(x_e) = \begin{cases} \frac{1}{2}\alpha^L_{i\to e}(x_e - \tilde{x}^b_{i\to e})^2 + \beta^L_{i\to e}(x_e - \tilde{x}^b_{i\to e}) + E_{i\to e}(\tilde{x}^b_{i\to e}) & x < \tilde{x}^b_{i\to e}, \\ \frac{1}{2}\alpha^R_{i\to e}(x_e - \tilde{x}^b_{i\to e})^2 + \beta^R_{i\to e}(x_e - \tilde{x}^b_{i\to e}) + E_{i\to e}(\tilde{x}^b_{i\to e}) & x > \tilde{x}^b_{i\to e}. \end{cases} \tag{S21}$$

For a primary effective leaf edge $i \to e$, the breakpoint is $\tilde{x}^b_{i\to e} = \Lambda_i$. For a general effective leaf edge $i \to e$, the breakpoint is most likely (but not always) located at the value of effective resource defined as

$$\Lambda^{\text{eff}}_{i\to e} := \Lambda_i + \sum_{e' \in EL_{i\to e}} B_{i,e'}\tilde{x}^b_{i\to e}. \tag{S22}$$

## E.   MP Equations for Non-smooth Message Functions

The MP equations for non-smooth message functions can be obtained with the information on effective leaf status of upstream edges, where one replaces the quadratic expansion $\Phi_{i\to e}(\tilde{x}_{i\to e} + \varepsilon_e) \approx \Phi_{i\to e}(\tilde{x}_{i\to e}) + \beta_{i\to e}\varepsilon_e + \frac{1}{2}\alpha_{i\to e}(\varepsilon_e)^2$ by the piecewise quadratic counterpart in Eq. (S21) when edge $i \to e$ is determined to be an effective leaf, and the double-sided message parameters $\{\alpha_{i\to e}^L, \beta_{i\to e}^L, \alpha_{i\to e}^R, \beta_{i\to e}^R\}$ are maintained and passed to its downstream edges.

For updating the messages, the computation of $\min_{\{x_{e'}\geq 0\}|R_i=0}\sum_{e'\in\partial i\backslash e}\left[\Phi_{k\to e'}(x_{e'}) + \phi_{e'}(x_{e'})\right]$ can be tedious if there are multiple effective leaf edges in $\{k \to e'|e' \in \partial i\backslash e\}$, where one needs to solve for a quadratic optimization of every case, where one branch of each non-smooth message function is selected each time (there are $2^{|EL_{i\to e}|}$ such cases in total). To simplify this process, we propose to firstly fix the flow $x_{e'}$ of effective leaves $EL_{i\to e}$ to be their breakpoints $\tilde{x}_{k\to e'}^b$ and then optimize non-effective leaf edges $NEL_{i\to e}$

$$\min_{\{\varepsilon_{e'}|e'\in NEL_{i\to e}\}} \sum_{e'\in NEL_{i\to e}} \left[\Phi_{k\to e'}(\tilde{x}_{k\to e'} + \varepsilon_{e'}) + \phi(\tilde{x}_{k\to e'} + \varepsilon_{e'})\right], \tag{S23}$$

$$\text{s. t. } 0 = \sum_{e'\in NEL_{i\to e}} B_{i,e'}\left[\tilde{x}_{k\to e'} + \varepsilon_{e'}\right] + \sum_{e'\in EL_{i\to e}} B_{i,e'}\tilde{x}_{k\to e'}^b + B_{i,e}x_e + \Lambda_i$$

$$= \sum_{e'\in NEL_{i\to e}} B_{i,e'}\left[\tilde{x}_{k\to e'} + \varepsilon_{e'}\right] + \Lambda_{i\to e}^{\text{eff}} + B_{i,e}x_e, \tag{S24}$$

$$0 \leq \tilde{x}_{k\to e'} + \varepsilon_{e'}. \tag{S25}$$

We then perturb the optimal solution by perturbing some of the flows $x_{e'}$ of upstream edges $e' \in \partial i\backslash e$ by an infinitesimal amount $dx$ as $x_{e'} = x_{k\to e'}^* + \eta_{e'}dx$ for non-effective leaves and $x_{e'} = \tilde{x}_{k\to e'}^b + \eta_{e'}dx$ for effective leaves with $\eta_{e'} = 0, \pm 1$. For non-smooth message function, $\eta_{e'} = -1$ and $\eta_{e'} = 1$, corresponding to the left and the right branch of $\Phi_{k\to e'}(x_{e'})$, respectively. To obey the flow conservation constraint $R_i = 0$, the perturbation coefficient $\eta_{e'}$ must satisfy $\sum_{e'\in\partial i\backslash e} B_{i,e'}\eta_{e'} = 0$.

If the perturbation configuration $\{\eta_{e'}^*\}$ leading to the lowest energy of $\sum_{e'\in\partial i\backslash e}\left[\Phi_{e'\to i}(x_{e'}) + \phi_{e'}(x_{e'})\right]$ reduces the outcome of Eq. (S23), we need to consider adding the effective leafs $k \to e'$ with $\eta_{e'}^* \neq 0$ as active optimization variables in addition to the non-effective leafs. Specifically, we define $\eta^{\text{active}} = \{e'|e' \in EL_{i\to e}, \eta_{e'}^* \neq 0\}$, and proceed to solve

$$\min_{\{\varepsilon_{e'}|e'\in NEL_{i\to e}\cup\eta^{\text{active}}\}} \sum_{e'\in NEL_{i\to e}\cup\eta^{\text{active}}} \left[\Phi_{k\to e'}(\tilde{x}_{k\to e'} + \varepsilon_{e'}) + \phi(\tilde{x}_{k\to e'} + \varepsilon_{e'})\right], \tag{S26}$$

$$\text{s. t. } 0 = \sum_{e'\in NEL_{i\to e}\cup\eta^{\text{active}}} B_{i,e'}\left[\tilde{x}_{k\to e'} + \varepsilon_{e'}\right]$$

$$+ \sum_{e'\in EL_{i\to e}\backslash\eta^{\text{active}}} B_{i,e'}\tilde{x}_{k\to e'}^b + B_{i,e}x_e + \Lambda_i, \tag{S27}$$

$$0 \leq \tilde{x}_{k\to e'} + \varepsilon_{e'}, \tag{S28}$$

where we use $\Phi_{k\to e'} = \frac{1}{2}\alpha_{k\to e'}^L\varepsilon_{e'}^2 + \beta_{k\to e'}^L\varepsilon_{e'}$ if $\eta_{e'}^* = -1$ and use $\Phi_{k\to e'} = \frac{1}{2}\alpha_{k\to e'}^R\varepsilon_{e'}^2 + \beta_{k\to e'}^R\varepsilon_{e'}$ if $\eta_{e'}^* = 1$.

The primal and dual variables in the optimum satisfy

$$x_{k\to e'}^*(\mu) = \tilde{x}_{k\to e'} + \varepsilon_{e'}^*(\mu) = \max\left(\tilde{x}_{k\to e'} - \frac{\mu B_{i,e'} + \phi_{e'}' + \beta_{k\to e'}}{\alpha_{k\to e'} + \phi_{e'}''}, 0\right), \tag{S29}$$

$$R_{i\to e}(\mu; x_e) = \sum_{e'\in NEL_{i\to e}\cup\eta^{\text{active}}} B_{i,e'}x_{k\to e'}^*(\mu)$$

$$+ \sum_{e'\in EL_{i\to e}\backslash\eta^{\text{active}}} B_{i,e'}\tilde{x}_{k\to e'}^b + B_{i,e}x_e + \Lambda_i = 0. \tag{S30}$$

If the solution $\mu^*$ in Eq. (S30) is non-degenerate, we have

$$\beta_{i\to e} = B_{i,e}\mu^*, \tag{S31}$$

$$\alpha_{i\rightarrow e} = \left[ \sum_{e' \in NEL_{i\rightarrow e} \cup \eta^{\text{active}}} \frac{1}{\alpha_{k\rightarrow e'} + \phi''_{e'}} \right.$$

$$\left. \times \Theta\left( (\alpha_{k\rightarrow e'} + \phi''_{e'})\tilde{x}_{k\rightarrow e'} - (\mu^* B_{i,e'} + \phi'_{e'} + \beta_{k\rightarrow e'}) \right) \right]^{-1}, \tag{S32}$$

in which case the edge $i \rightarrow e$ is not an effective leaf with $f_{i\rightarrow e} = 0$.

On the other hand, if the solution $\mu^*$ is degenerate, we need to consider $x_e = \tilde{x}_{i\rightarrow e} - \mathrm{d}x$ to solve for $\beta^L_{i\rightarrow e}, \alpha^L_{i\rightarrow e}$, and consider $x_e = \tilde{x}_{i\rightarrow e} + \mathrm{d}x$ to solve for $\beta^R_{i\rightarrow e}, \alpha^R_{i\rightarrow e}$, and identify edge $i \rightarrow e$ as an effective leaf with $f_{i\rightarrow e} = 1$.

It can also be shown that $\beta^R_{i\rightarrow e} > \beta^L_{i\rightarrow e}$ and the non-smooth message function $\Phi_{i\rightarrow e}(x_e)$ is convex.

### 1. Update of the Working Points

If the message function $\Phi_{i\rightarrow e}(x_e)$ is non-smooth, we would like to bring the working point $\tilde{x}_{i\rightarrow e}$ to the vicinity of the breakpoint of the two branches. To determine whether an edge $i \rightarrow e$ is an effective leaf, we perform the following procedure: We check the two following criteria: (i) each edge $k \rightarrow e'$ in the upstream edge set $\partial i \backslash e$ satisfies either $f_{k\rightarrow e'} = 1$ or $\tilde{x}_{k\rightarrow e'} = 0$; (ii) the difference between the current working point and the effective resource $|\tilde{x}_{i\rightarrow e} - \Lambda^{\text{eff}}_{i\rightarrow e}|$ is smaller than some threshold ($\Lambda^{\text{eff}}_{i\rightarrow e}$ is defined in Eq. (S22)). If both criteria (i) and (ii) are met, then we use the effective resource as the working point $\tilde{x}_{i\rightarrow e} = \Lambda^{\text{eff}}_{i\rightarrow e}$, and perform the optimization $\min_{\{x_{e'} \geq 0\}|R_i=0} \sum_{e' \in \partial i \backslash e} \left[ \Phi_{k\rightarrow e'}(x_{e'}) + \phi_{e'}(x_{e'}) \right]$; if it results in degenerate solutions of the Lagrangian multiplier $\mu^*$ for the flow conservation constraint, then edge $i \rightarrow e$ is determined as an effective leaf and the double-sided messages $\{\beta^L_{i\rightarrow e}, \alpha^L_{i\rightarrow e}, \beta^R_{i\rightarrow e}, \alpha^R_{i\rightarrow e}\}$ are computed. Otherwise, edge $i \rightarrow e$ is a non-effective leaf and the normal messages $\{\beta_{i\rightarrow e}, \alpha_{i\rightarrow e}\}$ are recorded.

If criteria (i) and (ii) are not met, we use the current value of the working point $\tilde{x}_{i\rightarrow e}$ to solve for the messages. Similarly, if the optimization leads to degenerate solutions of $\mu^*$, then edge $i \rightarrow e$ is determined as an effective leaf. Otherwise, edge $i \rightarrow e$ is a non-effective leaf.

Similar to the case of smooth message functions in Sec. I C, the working point is updated as

$$\tilde{x}^{\text{new}}_{i\rightarrow e} \leftarrow sx^*_e + (1-s)\tilde{x}^{\text{old}}_{i\rightarrow e}, \tag{S33}$$

where $x^*_e = \arg\min_{x_e \geq 0} \left[ \Phi_{i\rightarrow e}(x_e) + \Phi_{j\rightarrow e}(x_e) + \phi_e(x_e) \right]$.

### F. Results of MP Algorithm for Routing Game

The MP algorithm for solving the (single-level) equilibrium flow problem is summarized in Algorithm. 1.

---

**Algorithm 1:** Message-passing algorithm for equilibrium flows in routing games (single-level, single destination, Method I to treat the destination node $\mathcal{D}$ )

---

**Input:** Road network $G(V, E)$ (pre-processed to remove dangling nodes), node parameters $\{\Lambda_i\}$ defining the resources, edge parameters defining the latency function $\ell_e(x_e)$ and the edge-wise potential $\phi_e(x_e)$ (defined in Eq. (S4)), maximal number of iterations $T$.

Initialize the messages $\{\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}\}$ and $\{f_{i \to e}\}$ (effective leaf states) randomly.

Initialize the working points $\{\tilde{x}_{i \to e}\}$ randomly.

**for** $t$ in $1 : T$ **do**

    Randomly select a node $i$ and one of its adjacent edge $e \in \partial i$.

    **Begin Subroutine (a) (update the messages $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, f_{i \to e}$ using $\phi_e(\cdot)$):**

    Compute $\Lambda_{i \to e}^{\text{eff}}$ using Eq. (S22).

    **if** $i = \mathcal{D}$ **then**

        | Set $\alpha_{i \to e} = 0, \beta_{i \to e} = 0, f_{i \to e} = 0$.

    **else if** $\exists\, e' \in \{\partial i \backslash e\}, \tilde{x}_{k \to e'} \neq 0, f_{k \to e'} = 0$ and $|\tilde{x}_{i \to e} - \Lambda_{i \to e}^{eff}| >$ threshold **then**

        | Update the messages $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, f_{i \to e}$ with the working point evaluated at the up-to-date value of $\tilde{x}_{i \to e}$,
        | following the procedures outlined in Sec. I E.

    **else**

        | /* edge $i \to e$ is a potential effective leaf. */

        | Replace $\tilde{x}_{i \to e}$ by $\Lambda_{i \to e}^{\text{eff}}$ in Sec. I E to compute the messages $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, f_{i \to e}$.

        **if** $f_{i \to e} = 1$ **then**

            | /* This is an optional step, corresponding to a more conservative strategy of identifying the
            |    effective leaf state. */

            | Using the newly computed $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}$, check whether $\Lambda_{i \to e}^{\text{eff}}$ is the minimum of $\Phi^{\text{full}}(x_e)$.
            | If yes, adopt the messages $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, f_{i \to e}$.
            | If no, reset $f_{i \to e} = 0$.

        **end**

        **if** $f_{i \to e} = 0$ **then**

            | /* Do not consider edge $i \to e$ as an effective leaf. */

            | Use the up-to-date value of $\tilde{x}_{i \to e}$ to recompute the messages $\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}$.

        **end**

    **end**

    **End Subroutine (a)**

    **if** $f_{i \to e} = 1$ **then**

        | Set $\tilde{x}_{i \to e} = \Lambda_{i \to e}^{\text{eff}}$.

    **end**

    Using the up-to-date messages $\{\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, f_{i \to e}\}$ and $\{\tilde{x}_{i \to e}\}$ to determine the approximated forms of $\Phi_{i \to e}(x_e)$
    and $\Phi_{j \to e}(x_e)$.

    Compute $x_e^* = \arg\min_{x_e \geq 0} [\Phi_{i \to e}(x_e) + \Phi_{j \to e}(x_e) + \phi_e(x_e)]$.

    Update the working point as $\tilde{x}_{i \to e}^{\text{new}} \leftarrow s x_e^* + (1 - s)\tilde{x}_{i \to e}^{\text{old}}$ with a learning rate $s$.

    **if** messages converge **then**

        | Exit the for loop.

    **end**

**end**

For each edge $e$, compute the equilibrium flow as $x_e^* = \arg\min_{x_e \geq 0} [\Phi_{i \to e}(x_e) + \Phi_{j \to e}(x_e) + \phi_e(x_e)]$.

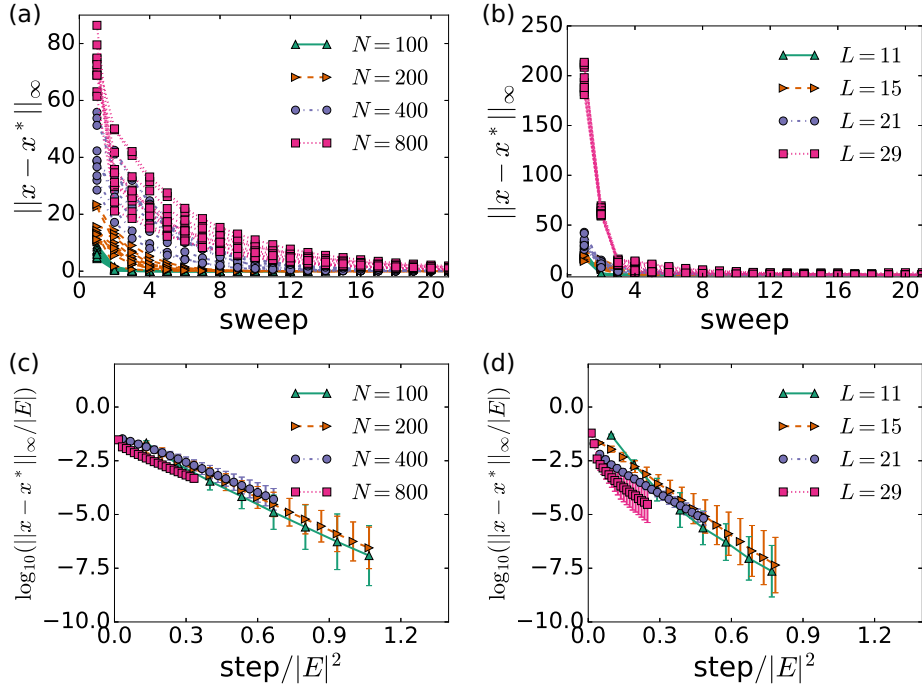**Output:** Convergence status and equilibrium flows $\{x_e^*\}$

---

Figure S5. Convergence of the MP algorithm for routing games in networks to the equilibrium flows $\boldsymbol{x}^*$. Random regular graphs of degree 3 are considered in (a)(c), while small-world networks obtained by rewiring square lattices with randomly chosen shortcut edges (rewiring probability $p_{rw} = 0.05$) are considered in (b)(d), respectively. The flows adjacent to the destination node $\mathcal{D}$ are unconstrained (Method I in Sec. I C 1), reminiscent of a grounded node in electric circuits. In (a)(b), the MP process for specific problem realizations is illustrated. Each sweep comprises $40|E|$ local MP updates. Panels (c)(d) are obtained by averaging over 10 problem realizations, and rescaling the MP runtime by $|E|^2$.

We further report results of the MP algorithms described above. Taking into account the possible non-smooth structure of message functions $\Phi_{i \to e}(x_e)$, the MP algorithm converges well for various types of graphs and resource distributions. We demonstrate the effectiveness of the algorithm in Figures S5 and S6, where random regular graphs and small-world networks are considered. The small-world networks are obtained by rewiring square lattices with randomly chosen shortcut edges [5]. In Fig. S5, the flows adjacent to the destination node $\mathcal{D}$ are unconstrained (Method I in Sec. I C 1). The MP algorithms converge to the correct equilibrium flows $\boldsymbol{x}^*$, and the empirical complexity for computing the equilibrium flows up to a certain error $|\boldsymbol{x}^{MP} - \boldsymbol{x}^*|$ is roughly $O(|E|^2)$.

In Fig. S6, we use Method II in Sec. I C 1, i.e., we put an explicit constraint to the flows adjacent to the destination node $\mathcal{D}$ as

$$R_{\mathcal{D}} = \Lambda_{\mathcal{D}} + \sum_{e \in \partial \mathcal{D}} B_{\mathcal{D},e} x_e = 0, \tag{S34}$$

where $\Lambda_{\mathcal{D}} = -\sum_{i \neq \mathcal{D}} \Lambda_i$. In this approach, the MP algorithms converge much faster; the empirical complexity for computing the equilibrium flows up to a certain error $|\boldsymbol{x}^{MP} - \boldsymbol{x}^*|$ is roughly $O(|E|)$. However, there exists some networks where MP with Method II does not converge, while MP with Method I converges successfully. For the experiments in the main text, we use Method I to treat the destination node for its better convergence properties.

## G. Extension to The Case of Multiple Destination

The case of multiple destinations can be studied similarly. The traffic flows can be classified into different classes according to their destinations. Let $N_d$ denotes the number of destinations, and $x_e^a$ denote the flow on edge $e$ targeted at the $a$-th destination (or the $a$-th class), the lower-level optimization problem (for solving equilibrium flows) is
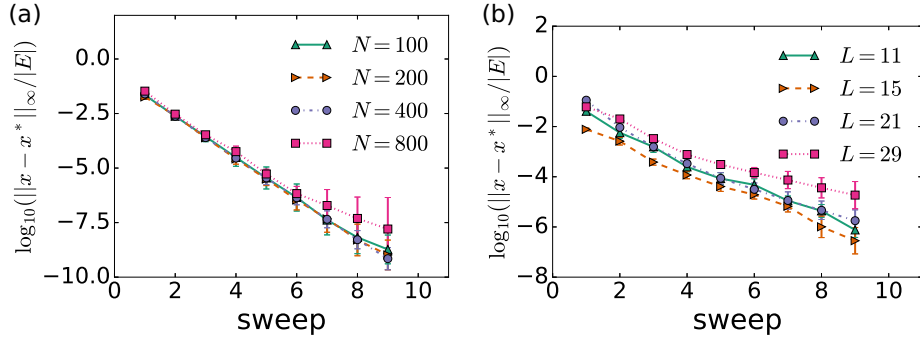
Figure S6. Same setting as in Fig. S5, except that the flows adjacent to the destination node $\mathcal{D}$ explicitly obey the constraint $R_\mathcal{D} = \Lambda_\mathcal{D} + \sum_{e \in \partial \mathcal{D}} B_{\mathcal{D},e} x_e = 0$, where $\Lambda_\mathcal{D} = -\sum_{i \neq \mathcal{D}} \Lambda_i$ (Method II in Sec. I C 1). Each sweep comprises of $40|E|$ local MP updates.
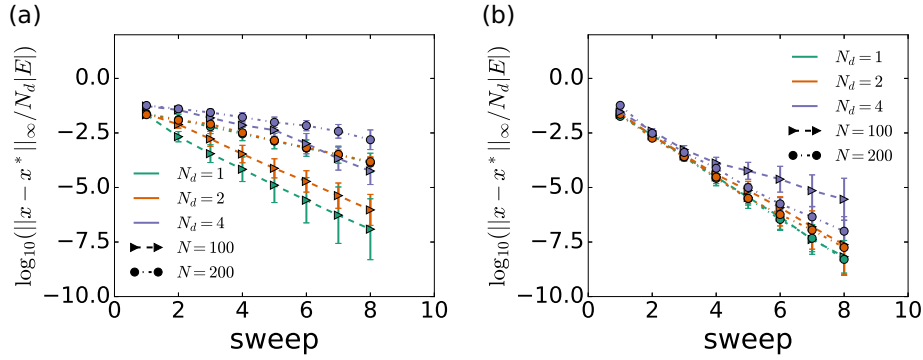


Figure S7. Message-passing algorithm for routing games with $N_d$ destinations converges to the equilibrium flows $\boldsymbol{x}^*$. (a) Method I in Sec. I C 1) is used to treat the destination nodes. (b) Method II in Sec. I C 1) is used to treat the destination nodes. Each sweep comprises of $40|E|$ local MP updates.

defined as

$$\min_{\boldsymbol{x}} \; \Phi(\boldsymbol{x}) = \sum_{e \in E} \int_0^{\sum_a x_e^a} \ell_e(y) \mathrm{d}y = \sum_{e \in E} \phi_e \Big( \sum_{a=1}^{N_d} x_e^a \Big), \tag{S35}$$

$$\text{s.t. } R_i^a := \Lambda_i^a + \sum_{e \in \partial i} B_{i,e} x_e^a = 0, \quad \forall i, a, \tag{S36}$$

$$x_e \geq 0, \quad \forall e, a. \tag{S37}$$

To accommodate the nonlinear interactions of flows of different classes $\{x_e^a\}$ in $\phi_e(\sum_a x_e^a)$, we adopt a coordinate-descent like approach in the MP algorithm as follows. In the treatment of the flows of class $a$, we fix the flows of other classes $\{x_e^b\}_{b \neq a}$ to their working points and compute the messages of class $a$ as

$$\Phi_{i \to e}^a(x_e^a) = \min_{\{x_{e'}^a \geq 0\} | R_i^a = 0} \sum_{e' \in \partial i \setminus e} \left[ \Phi_{k \to e'}^a(x_{e'}^a) + \phi_{e'} \Big( \sum_{b \neq a} \tilde{x}_{k \to e'}^b + x_{e'}^a \Big) \right]$$

$$= \min_{\{x_{e'}^a \geq 0\} | R_i^a = 0} \sum_{e' \in \partial i \setminus e} \left[ \Phi_{k \to e'}^a(\tilde{x}_{k \to e'}^a + \varepsilon_{e'}^a) + \phi_{e'} \Big( \sum_b \tilde{x}_{k \to e'}^b + \varepsilon_{e'}^a \Big) \right]. \tag{S38}$$

We further adopt the approximations $\Phi_{k \to e'}^a(\tilde{x}_{k \to e'}^a + \varepsilon_{e'}^a) \approx \Phi_{k \to e'}^a(\tilde{x}_{k \to e'}^a) + \beta_{k \to e'}^a \varepsilon_{e'}^a + \frac{1}{2} \alpha_{k \to e'}^a (\varepsilon_{e'}^a)^2$ (augmented by a piecewise quadratic function if it is non-smooth) and $\phi_{e'}(\sum_b \tilde{x}_{k \to e'}^b + \varepsilon_{e'}^a) \approx \phi_{e'}(\sum_b \tilde{x}_{k \to e'}^b) + \phi_{e'}'(\sum_b \tilde{x}_{k \to e'}^b)\varepsilon_{e'}^a + \frac{1}{2}\phi_{e'}''(\sum_b \tilde{x}_{k \to e'}^b)(\varepsilon_{e'}^a)^2$, and solve for the coefficients $\{\alpha_{i \to e}^a, \beta_{i \to e}^a\}$ as in the single-destination scenario. The resulting MP algorithm has the same structure as the one of single class described above, where its efficacy is shown in Fig. S7.

## H. Bilevel Optimization in Routing Games

The toll optimization problem (for single destination) of the upper-level planner is defined as

$$\min_{\boldsymbol{\tau}} H(\boldsymbol{x}^*(\boldsymbol{\tau})) = \sum_{e \in E} x_e^*(\boldsymbol{\tau}) \ell_e\big(x_e^*(\boldsymbol{\tau})\big), \tag{S39}$$

$$\text{s. t. constraints of } \boldsymbol{\tau} \text{ and:} \tag{S40}$$

$$\boldsymbol{x}^*(\boldsymbol{\tau}) = \arg\min_{\boldsymbol{x}} \Phi(\boldsymbol{x}; \boldsymbol{\tau}) = \arg\min_{\boldsymbol{x}} \sum_{e \in E} \int_0^{x_e} \big[\ell_e(y) + \tau_e\big] \mathrm{d}y, \tag{S41}$$

$$\text{s.t. } x_e \geq 0, R_i = 0, \forall e, i. \tag{S42}$$

Here, the social cost $H(\boldsymbol{x})$ is the objective function of the upper-level problem, i.e., the overall target of the central planner is to reduce this social cost. The potential function $\Phi(\boldsymbol{x}; \boldsymbol{\tau})$ is the objective function of the lower-level problem, which governs the equilibrium flow $\boldsymbol{x}^*(\boldsymbol{\tau})$ for a given toll configuration $\boldsymbol{\tau}$.

The computation of the social optimum $H(x)$ has a similar form for computing the equilibrium flows, we therefore use a parallel MP procedure for the social cost

$$H_{i \to e}(x_e) = \min_{\{x_{e'} \geq 0\}|R_i} \sum_{e' \in \partial i \backslash e} \left\{ H_{k \to e'}(x_{e'}) + \sigma(x_{e'}) \right\},$$

$$= \min_{\{x_{e'} \geq 0\}|R_i} \sum_{e' \in \partial i \backslash e} \left\{ \frac{1}{2} \gamma_{k \to e'} \big(\varepsilon_{e'}\big)^2 + \delta_{k \to e'} \varepsilon_{e'} + \frac{1}{2} \sigma_{e'}''(\tilde{x}_{k \to e'}) \big(\varepsilon_{e'}\big)^2 + \sigma_{e'}'(\tilde{x}_{k \to e'}) \varepsilon_{e'} \right\}, \tag{S43}$$

where $H_{k \to e'}(x_{e'})$ assumes a quadratic approximation and needs to be augmented by a piecewise quadratic function if it is non-smooth. It results in an upper-level MP algorithm having the same structure as the one for computing the equilibrium flows. The difference is that, since the flow is not directly driven by the central planner, the working point $\tilde{x}_{i \to e}$ is not updated at the upper level MP. Instead, the central planner updates the toll $\tau_e$ such that selfish users are attracted to the solution with a lower social cost.

When the toll $\tau_e$ on edge $e$ is adapted, the marginal Nash-equilibrium flow $x_e^N$ changes accordingly

$$x_e^N(\tau_e) = \arg\min_{x_e \geq 0} \big[\Phi_{i \to e}(x_e) + \Phi_{j \to e}(x_e) + \phi_e(x_e) + \tau_e x_e\big]. \tag{S44}$$

For smooth message functions $\Phi_{i \to e}(x_e)$ and $\Phi_{j \to e}(x_e)$

$$x_e^N(\tau_e) = \max\left( \frac{\big(\alpha_{i \to e} + \frac{1}{2}\phi_{i \to e}''\big)\tilde{x}_{i \to e} + \big(\alpha_{j \to e} + \frac{1}{2}\phi_{j \to e}''\big)\tilde{x}_{j \to e} - \big(\beta_{i \to e} + \beta_{j \to e} + \tau_e + \frac{1}{2}\phi_{i \to e}' + \frac{1}{2}\phi_{j \to e}'\big)}{\alpha_{i \to e} + \alpha_{j \to e} + \frac{1}{2}\phi_{i \to e}'' + \frac{1}{2}\phi_{j \to e}''}, 0 \right), \tag{S45}$$

which is a piecewise linear function of $\tau_e$ with two branches. For non-smooth cavity functions, $x_e^N(\tau_e)$ can also be obtained straightforwardly, which is a piecewise linear function of $\tau_e$ with multiple branches.

The goal of toll-adaptation of $\tau_e$ is to decrease the social cost $H(\boldsymbol{x})$, which amounts to decrease the full social cost on edge $e$

$$\tau_e^* = \arg\min_{\tau_e} H_e^{\text{full}}(x_e^N(\tau_e)), \tag{S46}$$

$$H_e^{\text{full}}(x_e) := H_{i \to e}(x_e) + H_{j \to e}(x_e) + \sigma_e(x_e), \tag{S47}$$

where the optimization in Eq. (S46) needs to obey necessary constraints on tolls (e.g., the restriction $0 \leq \tau_e \leq \tau_e^{\max}$ is considered in the main text).

As $H_e^{\text{full}}(x_e^N(\tau_e))$ is a convex function of $x_e^N$, it is sufficient to adapt $\tau_e$ such that $x_e^N(\tau_e)$ gets as close to the marginal socially optimal flow $x_e^G$ as possible, where $x_e^G$ is given by

$$x_e^G = \arg\min_{x_e \geq 0} \big[H_{i \to e}(x_e) + H_{j \to e}(x_e) + \sigma_e(x_e)\big]. \tag{S48}$$

The search for the optimal toll $\tau_e^*$ can be done efficiently by utilizing the property that $x_e^N(\tau_e)$ is a piecewise linear function of $\tau_e$.
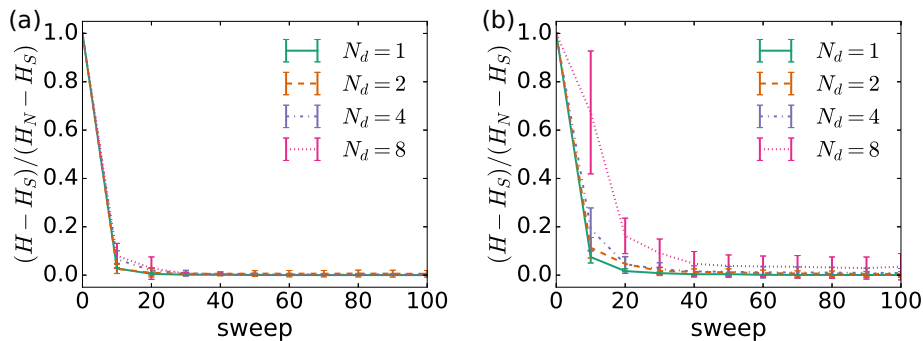
Figure S8. The effect of tolls on the reduction in fractional social cost in routing games on random regular graphs with $N_d$ destinations, where tolls are restricted as $0 \leq \tau_e \leq 1$. Each sweep comprises $40N_d|E|$ local MP updates and 100 edgewise toll updates in a random sequential schedule. (a) $N = 100$. (b) $N = 200$.

The resulting bilevel MP algorithm is described in the main text, where the lower-level messages $\{\alpha_{i \to e}^{(m)}, \beta_{i \to e}^{(m)}, \tilde{x}_{i \to e}\}$ ($m \in \{L, R\}$) and upper-level messages $\{\gamma_{i \to e}^{(m)}, \delta_{i \to e}^{(m)}\}$ are passed along edges to compute the equilibrium flows $x_e^N$ and related quantities. These messages facilitate the computation of $H_e^{\text{full}}(x_e)$ in the upper level, which is used to update the toll variables $\tau_e$. In practice, the update of tolls is less frequent then the update of other messages. In the experiments shown in the main text, for every $\frac{2}{5}N_d|E|$ MP iterations, we randomly select an edge $e$ and update its toll.

The resulting bilevel MP algorithm is summarized in Algorithm 2.

### 1. Extension to Multiple Destinations

The toll optimization problems with multiple destinations can be tackled by the proposed bilevel MP algorithm using the approximations in Sec. I G. The results are shown in Fig. S8, which demonstrate the effectiveness of the algorithm in reducing the social cost by adapting tolls.

## I. The Bilevel Programming Approach

Here we demonstrate the results of the bilevel programming approach to the toll optimization problem. It is achieved by expressing the solution of Eq. (S41) as constraints imposed by the Karush–Kuhn–Tucker (KKT) condition, and solve the bilevel optimization as a global nonlinear programming problem [6]. Such an approach is intrinsically difficult as (i) the constraints by the KKT conditions can be nonlinear and non-convex; (ii) the complementary slackness conditions are combinatorial, which requires a treatment with mixed integer programming (e.g., through branch and bound). Therefore, the bilevel programming approach offers a centralized algorithm that is generally non-scalable.

It is difficult to directly compare the MP algorithms to the bilevel programming approach, as the convergence rate for either approach is difficult to establish. Besides, the bilevel programming approach is a centralized optimization method, which has a different space complexity per iteration. Nevertheless, we present the results of CPU run time (CPU in used: i5-3317U) of bilevel programming (package in used: bileveljump.jl [7] with the IPOPT solver [8]) on the toll optimization problem in Fig. S9. There exist cases where bilevel programming fails to find the solution in a single trial, and the run times vary significantly among different problem realizations.

The bilevel programming approach is more generic and flexible than the MP approach, but it does not offer a decentralized algorithm as the MP approach. Besides, the MP algorithms can be extended to the scenarios with discrete variables, which is very difficult for the global optimization approach. It is also difficult to treat the toll selection problem in the main text with the bilevel programming method, especially when the socially optimum is not known a priori in some variants of toll-setting or network-design problems [9].

---

**Algorithm 2:** Message-passing algorithm for toll optimization in routing games (bilevel, single destination, Method I to treat the destination node $\mathcal{D}$ )

---

**Input:** Road network $G(V, E)$ (pre-processed to remove dangling nodes), node parameters $\{\Lambda_i\}$ defining the resources, edge parameters defining $\ell_e(x_e)$, $\phi_e(x_e)$ and $\sigma_e(x_e)$, maximal number of iterations $T$, time interval $t_{\text{update\_intv}}$ for updating tolls, time interval $t_{\text{dump\_intv}}$ for saving intermediate tolls.

Initialize the messages $\{\alpha_{i\to e}^{(m)}, \beta_{i\to e}^{(m)}\}$ and $\{f_{i\to e}^{\text{lw}}\}$ (effective leaf states of the lower layer) randomly.

Initialize the messages $\{\gamma_{i\to e}^{(m)}, \delta_{i\to e}^{(m)}\}$ and $\{f_{i\to e}^{\text{up}}\}$ (effective leaf states of the upper layer) randomly.

Initialize the working points $\{\tilde{x}_{i\to e}\}$ and tolls $\{\tau_e\}$.

**for** $t$ in $1:T$ **do**

    Randomly select a node $i$ and one of its adjacent edge $e \in \partial i$.

    `/* update the lower-level messages: */`

    Run Subroutine (a) in Algorithm 1 to update the messages $\alpha_{i\to e}^{(m)}, \beta_{i\to e}^{(m)}, f_{i\to e}^{\text{lw}}$ using $\phi_e(\cdot)$.

    **if** $f_{i\to e} = 1$ **then**
        | Set $\tilde{x}_{i\to e} = \Lambda_{i\to e}^{\text{eff}}$.
    **end**

    `/* update the upper-level messages: */`

    **if** $i = \mathcal{D}$ **then**
        | Set $\gamma_{i\to e} = 0, \delta_{i\to e} = 0, f_{i\to e}^{\text{up}} = 0$.
    **else**

        Replace $\{\alpha_{i\to e}^{(m)}, \beta_{i\to e}^{(m)}, f_{i\to e}^{\text{lw}}$ using $\phi_e(\cdot)\}$ by $\{\gamma_{i\to e}^{(m)}, \delta_{i\to e}^{(m)}, f_{i\to e}^{\text{up}}\}$ in Sec. I E.

        Follow the same procedures therein to update the messages $\gamma_{i\to e}^{(m)}, \delta_{i\to e}^{(m)}, f_{i\to e}^{\text{up}}$ with the working point evaluated at the up-to-date value of $\tilde{x}_{i\to e}$.

    **end**

    `/* update the working points: */`

    Using the up-to-date messages $\{\alpha_{i\to e}^{(m)}, \beta_{i\to e}^{(m)}, f_{i\to e}^{\text{lw}}\}$, $\{\tilde{x}_{i\to e}\}$ and $\{\tau_e\}$ to determine the approximated forms of $\Phi_{i\to e}(x_e)$ and $\Phi_{j\to e}(x_e)$.

    Compute $x_e^* = \arg\min_{x_e \ge 0} [\Phi_{i\to e}(x_e) + \Phi_{j\to e}(x_e) + \phi_e(x_e)]$.

    Update the working point as $\tilde{x}_{i\to e}^{\text{new}} \leftarrow s x_e^* + (1-s)\tilde{x}_{i\to e}^{\text{old}}$ with a learning rate $s$.

    `/* update the tolls: */`

    **if** $t$ mod $t_{\text{update\_intv}} = 0$ **then**

        Using the up-to-date messages $\{\gamma_{i\to e}^{(m)}, \delta_{i\to e}^{(m)}, f_{i\to e}^{\text{up}}\}$ and $\{\tilde{x}_{i\to e}\}$ to determine the approximated forms of $H_{i\to e}(x_e)$ and $H_{j\to e}(x_e)$.

        Compute $x_e^G = \arg\min_{x_e \ge 0} [H_{i\to e}(x_e) + H_{j\to e}(x_e) + \sigma_e(x_e)]$.

        Determine the (estimated) equilibrium flow as a function of toll:

        $x_e^N(\tau_e) = \arg\min_{x_e \ge 0} [\Phi_{i\to e}(x_e) + \Phi_{j\to e}(x_e) + \phi_e(x_e | \tau_e)]$, which is a piecewise linear function.

        Find $\tau_e \in [0, \tau_e^{\max}]$ such that $x_e^*(\tau_e)$ is as close to $x_e^G$ as possible, and update $\tau_e$ accordingly.

    **end**

    **if** $t$ mod $t_{\text{save\_intv}} = 0$ **then**
        | Save the up-to-date tolls.
    **end**

    **if** messages converge **then**
        | Exit the for loop.
    **end**

**end**

For each edge $e$, compute the equilibrium flow as $x_e^* = \arg\min_{x_e \ge 0} [\Phi_{i\to e}(x_e) + \Phi_{j\to e}(x_e) + \phi_e(x_e)]$.

**Output:** Convergence status, equilibrium flows $\{x_e^*\}$, the optimized tolls $\{\tau_e\}$

---

## II. MESSAGE-PASSING ALGORITHMS FOR ATOMIC ROUTING GAMES

Message-passing algorithms can be extended to include atomic routing games, where each player controls one unit of traffic. The atomic games differ from the above-studied non-atomic games in that the flow variable of edge $e$ is an integer $x_e \in \mathbb{Z}$. The integer constraints make even the single-level optimization a difficult integer programming problem. An existing mixed integer programming approach can be used to solve the equilibrium flow problem of
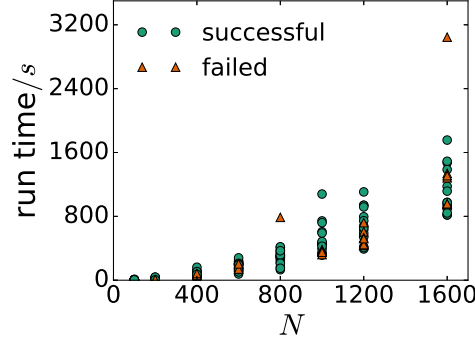
Figure S9. Run time of bilevel programming on the toll optimization problem in random regular graphs. For each network size, 20 different problem realizations are considered; red triangles represent the cases where bilevel programming fails to find the solution in single trials, green dots represent successful trials.

atomic routing games for moderate network sizes. However, the bilevel optimization (such as the toll-setting problem) appears to be much more difficult to solve.

Message-passing algorithms have also been successfully applied to solve network flow and routing problems (as a single-level optimization problem) with integer flow variables [10–13]. Here, we generalize the MP algorithms for integer flows to solve the Wardrop equilibrium problem in directed networks. The problem of atomic games is defined similarly as the non-atomic games introduced in Sec. I A, except that there are additional integer constraints $x_e \in \mathbb{Z}$, and that the potential function $\Phi(\boldsymbol{x})$ (to be minimized) in Eq. (S4) is replaced by

$$\Phi(\boldsymbol{x}) = \sum_{e \in E} \sum_{y=1}^{x_e} \left[ \ell_e(y) + \tau_e \right] =: \sum_{e \in E} \phi_e(x_e). \tag{S49}$$

The MP equation for minimizing the potential $\Phi(\boldsymbol{x})$ can be written similarly to the non-atomic games as

$$\Phi_{i \to e}(x_e) = \min_{\{x_{e'} \in \mathbb{N}\} | R_i = 0} \sum_{e' \in \partial i \setminus e} \left[ \Phi_{e' \to i}(x_{e'}) + \phi_{e'}(x_{e'}) \right], \tag{S50}$$

where the message function $\Phi_{i \to e}(x_e)$ is defined on a 1-dimensional grid $x_e \in \{0, 1, 2, ...\}$ to be computed recursively.

To simplify the above message-passing calculations, we can adopt the perturbation approach similar to the case of non-atomic game, by considering only a few flow values $x_e$ near a certain working point $\tilde{x}_{i \to e} \in \mathbb{N}$ [11]

$$\Phi_{i \to e}(x_e) = \Phi_{i \to e}(\tilde{x}_{i \to e} + m) =: h_{i \to e}^m, \qquad m \in \mathbb{Z}, -M \leq m \leq M, \tag{S51}$$

where $M$ is a small integer parameter determining the scope to look around $\tilde{x}_{i \to e}$. The MP algorithm proceeds by iteratively solving Eq. (S50) for different directed edges $\{i \to e | i \in V, e \in E\}$. The optimization problem for a particular edge $i \to e$ is approximately achieved by searching the grid $\{x_{e'} \in \mathbb{N}\}_{\forall e' \in \partial i \setminus e}$ in the vicinity of $\{\tilde{x}_{e' \to i}\}_{\forall e' \in \partial i \setminus e}$. The complexity of each iteration is proportional to $(2M + 1)^{|\partial i| - 1}$. In our experiments, we use $M = 1$ following the suggestion of [11]. To some extent, the implementation is more straightforward than that in non-atomic games, as there is no need to treat the non-smooth cavity energy function.

Similarly to the non-atomic game case, the working point $\tilde{x}_{i \to e}$ is also gradually pushed towards the minimizer of the edgewise full energy

$$x_e^* = \arg \min_{x_e \in \mathbb{N}} \Phi_e^{\text{full}}(x_e) = \arg \min_{x_e \in \mathbb{N}} \left[ \Phi_{i \to e}(x_e) + \Phi_{j \to e}(x_e) + \phi_e(x_e) \right], \tag{S52}$$

$$\tilde{x}_{i \to e}^{\text{new}} \leftarrow \tilde{x}_{i \to e}^{\text{old}} + \text{sign}(x_e^* - \tilde{x}_{i \to e}^{\text{old}}), \tag{S53}$$

in which case the working point $\tilde{x}_{i \to e}$ is expected to get closer and closer to the equilibrium point $x_e^*$ and the perturbation scheme will become more accurate.

In Eq. (S52), we can only access the flow values in the set $A_e := \{\tilde{x}_{i \to e} + m | -M \leq m \leq M\} \cap \{\tilde{x}_{j \to e} + m | -M \leq m \leq M\}$. If $A_e = \varnothing$, instead of applying Eq. (S53), we push $\tilde{x}_{i \to e}$ towards $\tilde{x}_{j \to e}$ incrementally to increase the chance of overlap in future steps

$$\tilde{x}_{i \to e}^{\text{new}} \leftarrow \tilde{x}_{i \to e}^{\text{old}} + \text{sign}(\tilde{x}_{j \to e}^{\text{old}} - \tilde{x}_{i \to e}^{\text{old}}). \tag{S54}$$

The message-passing algorithm in our study is known as the min-sum algorithm, where caution is needed when there are degenerate energy minima such that taking the minimum is ambiguous (see Sec. 8.4.5 of [14] for example). This issue does not impact on MP algorithms for non-atomic games due to the convex nature of the single-level problem, but complicates the use of MP algorithms for atomic games. In [10, 11], a small random bias field $\xi_e$ is added to non-linear cost $\phi_e(x_e)$ to break the degeneracy

$$\phi'_e(x_e) \leftarrow \phi_e(x_e) + \xi_e|x_e|. \tag{S55}$$

Alternatively, we can also select one of the degenerate solution based on its consistency with the constraints. For example, in Eq. (S52), suppose there are two flow values $\{x_e^{*(1)}, x_e^{*(2)}\}$ corresponding to $\min \Phi_e^{\text{full}}(x_e)$. We would choose the solution that is more consistent with the flow conservation constraint

$$R_{i \to e}^{(n)} = \Lambda_i + \sum_{e'=(k,i)\in\partial i\setminus e} B_{i,e'}\tilde{x}_{k\to e'} + B_{i,e}x_e^{*(n)}, \quad n = 1, 2, \tag{S56}$$

that is, we would assign the solution $x_e^{*(n)}$ with minimal value of $|R_{i\to e}^{(n)}|$ to $x_e^*$ in Eq. (S52).

Both methods facilitate the convergence of the MP algorithms to a local minima in atomic routing games. We refer readers to Ref. [11] for more details on MP algorithms in network flow problems with integer constraints.

## A. Bilevel MP Algorithms for Atomic Games

The method to treat bilevel optimization (in particular, the toll-setting problem) in non-atomic games can also be applied to atomic games. This is achieved by considering a parallel message-passing process for minimizing the social cost (in the upper-level)

$$H_{i\to e}(x_e) = \min_{\{x_{e'}\in\mathbb{N}\}|R_i=0} \sum_{e'\in\partial i\setminus e} \left[ H_{e'\to i}(x_{e'}) + \sigma_{e'}(x_{e'}) \right], \tag{S57}$$

where $\sigma_{e'}(x_{e'}) = x_{e'}\ell_e(x_{e'})$. The minimizer of the edgewise full cost $x_e^S = \arg\min_{x_e} H_e^{\text{full}}(x_e) = \arg\min_{x_e} \left[ H_{i\to e}(x_e) + H_{j\to e}(x_e) + \sigma_e(x_e) \right]$ is informative of the min social-cost flow during the upper-level MP updates. We then update the toll $\tau_e$ incrementally such that the toll-dependent equilibrium flow $x_e^*(\tau)$ (given by Eq. (S52)) gets closer to $x_e^S$.

We demonstrate the effectiveness of the bilevel MP algorithm for toll-setting problems in atomic games in Fig. S10. We observe that the algorithm is effective for small networks with light and moderate loads, while it becomes less effective for heavy loads. The performance also deteriorate in large networks (not shown). We conjecture that a more non-local toll update method is needed to improve the performance in such cases, which remains to be explored in future studies. Nevertheless, the bilevel optimization problems in atomic games are intrinsically difficult combinatorial optimization problems, where the complex energy landscape is prohibitive for most optimization algorithms.

## B. Optimal Toll Configurations for Some Instances

In this section, we demonstrate the optimal tolls found by the bilevel MP algorithms for some cases in the Sioux Falls road network [15], which is a popular benchmark network for transportation research, consisting of 24 nodes and 76 directed edges.

In this example, upon obtaining the optimal tolls $\{\tau_e^{\text{opt}}\}$ in the bilevel MP algorithm, we keep only the tolls with relatively large values $\tau_e^{\text{new}} \leftarrow \tau_e^{\text{opt}}\Theta(\tau_e^{\text{opt}} \geq \epsilon)$ (where $\epsilon$ is a small threshold). If this does not incur an increment in social cost, we keep $\{\tau_e^{\text{new}}\}$, otherwise, we keep $\{\tau_e^{\text{opt}}\}$.

The results are shown in Fig. S11. In Case I, we observed that the shortest paths from the source nodes to destination are very congested, and the optimal tolls are placed on these paths such that some users are re-routed to alleviate congestion. In Case II, interestingly, tolls are also placed on some non-shortest paths in the optimal configuration found by the bilevel MP algorithm.

## III. MESSAGE-PASSING ALGORITHMS FOR FLOW CONTROL IN UNDIRECTED NETWORKS

In this section, we provide the details of the MP algorithm for flow control in undirected networks. In a simple undirected graph $G(V, E)$, nodes $i$ and $j$ can be connected by at most one edge $(i, j)$, where the order of node $i$ and $j$
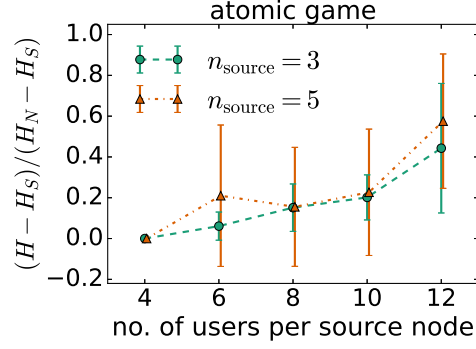
Figure S10. The effect of tolls on the fractional reduction in social cost. Random regular graphs of degree 3 and size $N = 50$ are considered. Each data point is the average of 10 different problem realizations. For each problem realization, we consider 5 different trials of bileve MP processes (with different random starting points), and recorded a few of the tolls during each process; we then kept the toll configuration corresponding to the smallest social cost. In each experiment, three or five nodes (dashed lines or dash-dotted lines) are randomly selected as the sources (origins), on which the users aim to route to a universal destination.
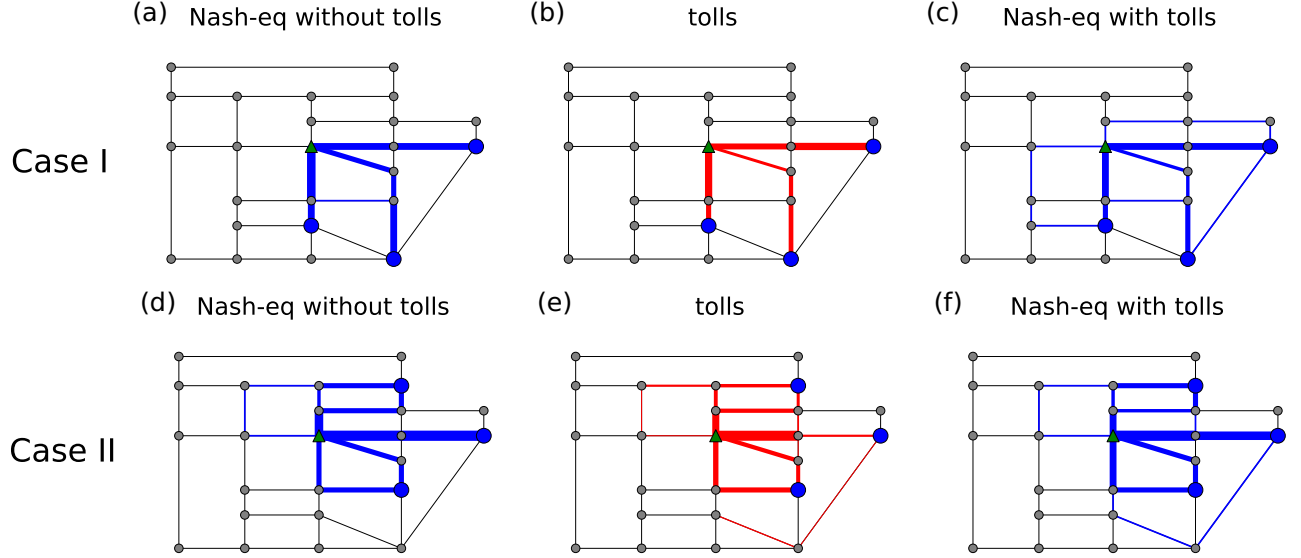


Figure S11. Optimal tolls and network flows in the Sioux Falls network found by the bilevel MP algorithm. Note that the underlying network is a directed graph (each undirected edge in the figures comprises two directed edges in both directions). We set the central node as the destination node (marked as a green triangle), and randomly select 3 nodes as the source nodes (marked as blue circles) in either case. In Case I (top), each source node has 4 users, while in Case II (bottom), each source has 6 users. In (a)(c) and (d)(f), the network flow patterns in Nash (Wardrop) equilibrium are shown, where a non-zero flow is marked as a blue edge (the edge width is proportional to the flow magnitude). In (b) and (e), the optimal tolls are shown, where a non-zero toll is marked as a red edge (the edge width is proportional to the toll magnitude).

does not matter (in contrast, edges $(i, j)$ and $(j, i)$ are two different edges in a directed graph). In this case, edge $(i, j)$ can either transmit resources from node $j$ to node $i$ or from node $i$ to node $j$. Denoting $x_{ij}(=-x_{ji})$ as the flow from node $j$ to node $i$; if $x_{ij} < 0$ (or $x_{ji} = -x_{ij} > 0$), the resources are being transmitted from node $i$ to node $j$. We also assume that the underlying graph does not have any leaf nodes, by recursively trimming leaf nodes and absorbing their resources into neighboring nodes.

## A. MP For Lower-level Optimization

The equilibrium flow (in the lower-level optimization problem) is the minimizer of the problem

$$\min_{\boldsymbol{x}} C(\boldsymbol{x}) = \sum_{(i,j)} \frac{1}{2} r_{ij} x_{ij}^2, \tag{S58}$$

$$\text{s.t. } R_i = \Lambda_i + \sum_{j \in \mathcal{N}_i} x_{ij} = 0, \quad \forall i \neq \mathcal{D}, \tag{S59}$$

where the reference node $\mathcal{D}$ can be arbitrarily chosen.

The above optimization problem can be mapped onto its dual problem as

$$\min_{\boldsymbol{\mu}} C^{\text{dual}}(\boldsymbol{\mu}) = \sum_{(i,j)} \frac{1}{2r_{ij}} (\mu_j - \mu_i)^2 - \sum_i \Lambda_i \mu_i, \tag{S60}$$

$$=: \frac{1}{2} \boldsymbol{\mu}^\top L \boldsymbol{\mu} - \boldsymbol{\Lambda}^\top \boldsymbol{\mu}, \tag{S61}$$

where $\mu_i$ is the Lagrange multiplier (or dual variable) associated with the flow conservation constraint $R_i = 0$, and $L$ is the Laplacian matrix with matrix element

$$L_{ij} := \left( \sum_{k \in \mathcal{N}_i} \frac{1}{r_{ik}} \right) \delta_{ij} - \frac{1}{r_{ij}}. \tag{S62}$$

The solution of the dual problem can be obtained by solving the system of linear equations $L\boldsymbol{\mu}^* = \boldsymbol{\Lambda}$, and the equilibrium flow $x_{ij}^*$ is related to the optimal Lagrange multiplier $\boldsymbol{\mu}^*$ through $x_{ij}^* = \frac{\mu_j^* - \mu_i^*}{r_{ij}}$. The drawback of such an approach is that (i) solving the systems of linear equations usually needs a centralized solver; (ii) to compute the response of the equilibrium flow to changes of the control parameters $\{r_{ij}\}$ in bilevel optimization, one needs to evaluate the pseudo-inverse of the Laplacian matrix per iteration, which can be very computationally demanding for large networks. Instead, we proposed to use MP for computing the equilibrium flows and tackle the related bilevel optimization problem, which is a scalable and efficient decentralized algorithm.

For the lower-level equilibrium flow problem in Eq. (S58), the MP algorithm amounts to computing the message functions

$$C_{i \to j}(x_{ij}) = \min_{\{x_{ki}\}|R_i=0} \left[ \frac{1}{2} r_{ij} x_{ij}^2 + \sum_{k \in \mathcal{N}_i \setminus j} C_{k \to i}(x_{ki}) \right], \tag{S63}$$

where the definition of the message function $C_{i \to j}(x_{ij})$ differs from the one in Eq. (S9) in that it includes the interaction term on edge $(i, j)$, which yields a more concise update rule in this problem. Similar to routing games, we approximate the message function by a quadratic form $C_{i \to j}(x_{ij}) = \frac{1}{2} \alpha_{i \to j} (x_{ij} - \hat{x}_{i \to j})^2 + \text{const}$, such that the local optimization in Eq. (S63) reduces to the computation of the real-number messages $m_{i \to j} \in \{\alpha_{i \to j}, \hat{x}_{i \to j}\}$ by passing the upstream messages $\{m_{k \to i}\}_{k \in \mathcal{N}_i \setminus j}$. Here the message function $C_{i \to j}(x_{ij})$ is always smooth. The messages $\alpha_{i \to j}, \hat{x}_{i \to j}$ are computed as [4, 16, 17]

$$\alpha_{i \to j} = \frac{1}{\sum_{k \in \mathcal{N}_i \setminus j} \alpha_{k \to i}^{-1}} + r_{ij}, \tag{S64}$$

$$\hat{x}_{i \to j} = \frac{\Lambda_i + \sum_{k \in \mathcal{N}_i \setminus j} \hat{x}_{k \to i}}{1 + r_{ij} \sum_{k \in \mathcal{N}_i \setminus j} \alpha_{k \to i}^{-1}}. \tag{S65}$$

### 1. The Reference Node $\mathcal{D}$

Similar to the MP algorithm in routing games, there are two methods to deal with possible boundary conditions of the reference node $\mathcal{D}$:

- Method I: Since node $\mathcal{D}$ has no constraints on its adjacent flows, it will absorb all incoming flows resulting in $C_{\mathcal{D} \to j}(x_{\mathcal{D}j}) = 0$ and

$$\alpha_{\mathcal{D} \to j} = 0, \quad \hat{x}_{\mathcal{D} \to j} = 0. \tag{S66}$$

In this treatment, the Lagrange multiplier of node $\mathcal{D}$ can be set to $\mu_{\mathcal{D}} = 0$ in the dual problem (as node $\mathcal{D}$ is unconstrained), which corresponds to a grounded node in the electric network interpretation of the problem.

- Method II: Alternatively, one can set an explicit constraint on the flows $\{x_{\mathcal{D}j}\}$ to the reference node $\mathcal{D}$

$$R_{\mathcal{D}} := \Lambda_{\mathcal{D}} + \sum_{j \in \mathcal{N}_{\mathcal{D}}} x_{\mathcal{D}j} = 0, \tag{S67}$$

where $\Lambda_{\mathcal{D}} = -\sum_{i \neq \mathcal{D}} \Lambda_i$. Then the messages from the reference node $\mathcal{D}$ are calculated in the same way as for other nodes.

Similar to the routing games, Method II results in an MP algorithm with a faster convergence rate, but it may fail to converge for some graphs while Method I can still provide valid solutions.

Method II is used in the experiments for undirected flow networks in the main text.

### 2. Computation of the Equilibrium Flows from Messages

Upon convergence of the messages, the equilibrium flow $x_{ij}^*$ can be obtained by minimizing the edgewise full cost $C_{ij}^{\text{full}}(x_{ij}) = C_{i \to j}(x_{ij}) + C_{j \to i}(x_{ij}) - \frac{1}{2} r_{ij} x_{ij}^2$, giving rise to

$$x_{ij}^* = \frac{\alpha_{j \to i} \hat{x}_{j \to i} - \alpha_{i \to j} \hat{x}_{i \to j}}{\alpha_{i \to j} + \alpha_{j \to i} - r_{ij}}. \tag{S68}$$

### 3. Results of the MP Algorithm on Undirected Flow Networks

In Fig. S12, we demonstrate the performance of the MP algorithms in undirected flow networks. The MP algorithms converge in different networks, including square lattices with many short loops. The iterations needed to obtained a given precision seems to depend on the topologies of the networks, e.g., square lattices appear to converge slower than random regular graphs. The method used to treat the boundary reference node $\mathcal{D}$ also impacts on the number of iterations needed, where it is observed that in general Method II makes MP converge faster than Method I. We conjecture that the influence of single-node boundary conditions in Eq. (S66) takes more iteration steps to diffuse messages to the bulk of the network.

## B. MP For Bilevel Optimization

The bilevel optimization problem on undirected networks aims to tune the flows of targeted edges $\mathcal{T}$ such that they exceed or drop below certain limits, depending on the application. We consider the former case, where the goal is to control the flows on $\mathcal{T}$ such that $\rho_{ij}(x_{ij}) = \frac{|x_{ij}| - |x_{ij}^0|}{|x_{ij}^0|} - \theta \geq 0, \forall (i,j) \in \mathcal{T}$ (with $x_{ij}^0$ being the flow before tuning). The task in the upper-level is to minimize the objective

$$\min_{\boldsymbol{r}} \mathcal{O}(\boldsymbol{x}^*(\boldsymbol{r})) = \sum_{(i,j) \in \mathcal{T}} \mathcal{O}_{ij}(x_{ij}^*(\boldsymbol{r})) := \sum_{(i,j) \in \mathcal{T}} -\rho_{ij}(x_{ij}^*(r)) \Theta\big(-\rho_{ij}(x_{ij}^*(r))\big). \tag{S69}$$

As mentioned in the main text, the impact of the variation of the control parameters $r_{ij}$ on the upper-level objective $\mathcal{O}$ is mediated through the messages $m_{i \to j} \in \{\alpha_{i \to j}, \hat{x}_{i \to j}\}$ along the pathways from the targeted edges to edge $(i,j)$. For a targeted edge $(p,q) \in \mathcal{T}$, the boundary conditions of the gradient with respect to the messages is given by $\frac{\partial \mathcal{O}_{pq}}{\partial m_{p \to q}} = \frac{\partial \mathcal{O}_{pq}}{\partial x_{pq}^*} \frac{\partial x_{pq}^*}{\partial m_{p \to q}}$, where the components admit the following expressions

$$\frac{\partial \mathcal{O}_{pq}}{\partial x_{pq}^*} = -\Theta\big(-\rho_{pq}(x_{pq}^*)\big) \frac{\text{sgn}(x_{pq}^*)}{|x_{pq}^0|}, \tag{S70}$$

$$\frac{\partial x_{pq}^*}{\partial \alpha_{p \to q}} = \frac{-\hat{x}_{p \to q}}{\alpha_{p \to q} + \alpha_{q \to p} - r_{pq}} - \frac{-x_{pq}^*}{\alpha_{p \to q} + \alpha_{q \to p} - r_{pq}}, \tag{S71}$$

$$\frac{\partial x_{pq}^*}{\partial \hat{x}_{p \to q}} = \frac{-\alpha_{p \to q}}{\alpha_{p \to q} + \alpha_{q \to p} - r_{pq}}. \tag{S72}$$
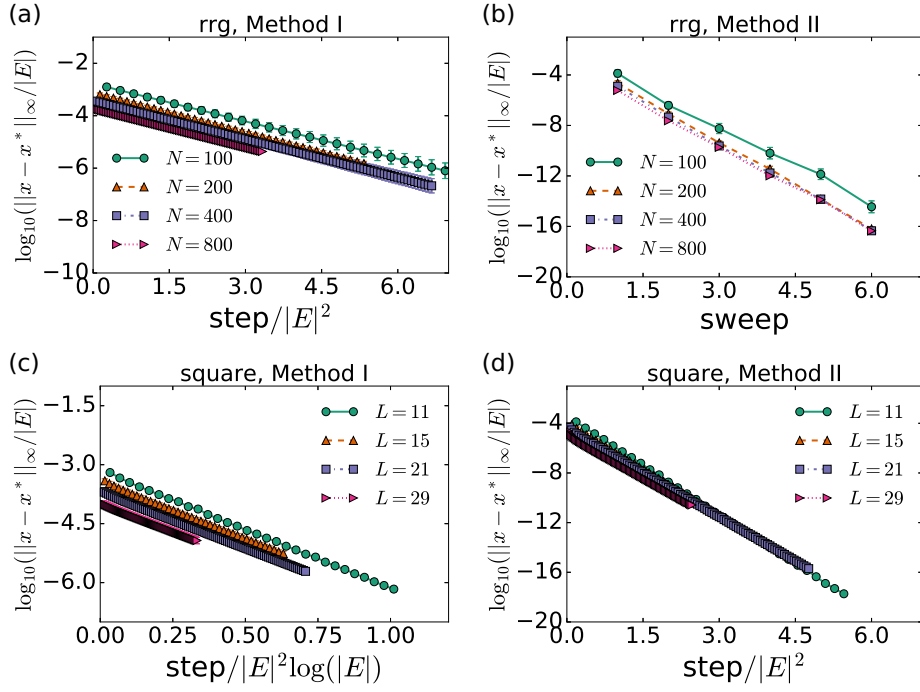
Figure S12. MP algorithm in undirected flow networks converges to the equilibrium flows $\boldsymbol{x}^*$. Random regular graphs (RRG) of degree 3 are considered in (a)(b), while square lattices (size $L \times L$) are considered in (c)(d). Method I in Sec. III A 1 is used in (a)(c) to treat the reference node $\mathcal{D}$, while Method II is used in (b)(d). In (b), each sweep comprises $40|E|$ local MP updates.

The gradient with respect to the control parameter $r_{pq}$ on the targeted edge $(p, q)$ is computed as

$$\frac{\partial \mathcal{O}_{pq}}{\partial r_{pq}} = \frac{\partial \mathcal{O}_{pq}}{\partial x_{pq}^*}\left[\frac{\partial x_{pq}^*}{\partial r_{pq}} + \sum_{m \in \{\alpha, \hat{x}\}} \left(\frac{\partial x_{pq}^*}{\partial m_{p \to q}}\frac{\partial m_{p \to q}}{\partial r_{pq}} + \frac{\partial x_{pq}^*}{\partial m_{q \to p}}\frac{\partial m_{q \to p}}{\partial r_{pq}}\right)\right]. \tag{S73}$$

For a non-targeted edge $(k, i) \notin \mathcal{T}$, we need to first evaluate $\frac{\partial \mathcal{O}}{\partial m_{k \to i}}$, which can be obtained by summing the gradients on its downstream edges $\{i \to l | l \in \mathcal{N}_i \backslash k\}$, computed as

$$\frac{\partial \mathcal{O}}{\partial m_{k \to i}} = \sum_{l \in \mathcal{N}_i \backslash k} \sum_{m_{i \to l} \in \{\alpha_{i \to l}, \hat{x}_{i \to l}\}} \frac{\partial \mathcal{O}}{\partial m_{i \to l}}\frac{\partial m_{i \to l}}{\partial m_{k \to i}}$$
$$= \sum_{(p,q) \in \mathcal{T}} \sum_{l \in \mathcal{N}_i \backslash k} \sum_{m_{i \to l} \in \{\alpha_{i \to l}, \hat{x}_{i \to l}\}} \frac{\partial \mathcal{O}_{pq}}{\partial m_{i \to l}}\frac{\partial m_{i \to l}}{\partial m_{k \to i}}. \tag{S74}$$

where the gradient propagation of messages $\frac{\partial m_{i \to l}}{\partial m_{k \to i}}$ admits the following forms

$$\frac{\partial \alpha_{i \to l}}{\partial \alpha_{k \to i}} = \frac{\alpha_{k \to i}^{-2}}{\left(\sum_{n \in \mathcal{N}_i \backslash l} \alpha_{n \to i}^{-1}\right)^2}, \tag{S75}$$

$$\frac{\partial \alpha_{i \to l}}{\partial \hat{x}_{k \to i}} = 0, \tag{S76}$$

$$\frac{\partial \hat{x}_{i \to l}}{\partial \alpha_{k \to i}} = \frac{\hat{x}_{i \to l} r_{il} \alpha_{k \to i}^{-2}}{1 + r_{il} \sum_{n \in \mathcal{N}_i \backslash l} \alpha_{n \to i}^{-1}}, \tag{S77}$$

$$\frac{\partial \hat{x}_{i \to l}}{\partial \hat{x}_{k \to i}} = \frac{1}{1 + r_{il} \sum_{n \in \mathcal{N}_i \backslash l} \alpha_{n \to i}^{-1}}. \tag{S78}$$

The resulting gradient w.r.t. to the control parameter $r_{ki}$ (note that edge $(k,i) \notin \mathcal{T}$) can be computed as

$$\frac{\partial \mathcal{O}}{\partial r_{ki}} = \sum_{(p,q) \in \mathcal{T}} \frac{\partial \mathcal{O}_{pq}}{\partial r_{ki}} = \sum_{(p,q) \in \mathcal{T}} \left[ \frac{\partial \mathcal{O}_{pq}}{\partial r_{k \to i}} + \frac{\partial \mathcal{O}_{pq}}{\partial r_{i \to k}} \right] \tag{S79}$$

$$\frac{\partial \mathcal{O}_{pq}}{\partial r_{k \to i}} = \sum_{m \in \{\alpha, \hat{x}\}} \frac{\partial \mathcal{O}_{pq}}{\partial m_{k \to i}} \frac{\partial m_{k \to i}}{\partial r_{ki}}, \tag{S80}$$

where we have broken down the gradient as $\frac{\partial \mathcal{O}_{pq}}{\partial r_{ki}} = \frac{\partial \mathcal{O}_{pq}}{\partial r_{k \to i}} + \frac{\partial \mathcal{O}_{pq}}{\partial r_{i \to k}}$ for programming convenience.

In Eqs. (S73) and (S80), the terms $\frac{\partial m_{k \to i}}{\partial r_{ki}}$ are computed as

$$\frac{\partial \alpha_{k \to i}}{\partial r_{ki}} = 1, \tag{S81}$$

$$\frac{\partial \hat{x}_{k \to i}}{\partial r_{ki}} = \hat{x}_{k \to i} \frac{-\sum_{n \in \mathcal{N}_k \setminus i} \alpha_{n \to k}^{-1}}{1 + r_{ki} \sum_{n \in \mathcal{N}_k \setminus i} \alpha_{n \to k}^{-1}}, \tag{S82}$$

which closes the equations for the gradient computations.

In this formalism, we keep track of the influence of the control parameter $r_{ij}$ on a specific targeted edge $(p,q) \in \mathcal{T}$ in the form of $\frac{\partial \mathcal{O}_{pq}}{\partial r_{ij}}$. The algorithmic complexity scales linearly with the number of targeted edges $|\mathcal{T}|$. This is acceptable as long as $|\mathcal{T}|$ is small. Further reducing the complexity (e.g., by considering the influence of the control parameter $r_{ij}$ on the overall objective function) is an interesting future research direction, which requires a more careful treatment of the boundary condition of the gradients.

We also consider the constraints on the control parameters, being in the range $r_{ij} \in [0.9, 1.1]$. More complex constraints can also be considered, e.g., by introducing a penalty function to the global objective, which will yield additional terms in the message-passing equations.

The resulting bilevel MP algorithms for flow control is summarized in Algorithm 3.

### C. The Global Optimization Approach

In this section, we provide details of the global optimization approach on undirected flow networks used in the main text. As mentioned before, e.g., in Sec. III A, we have $x_{ij}^* = \frac{\mu_j^* - \mu_i^*}{r_{ij}}$, where $\boldsymbol{\mu}^* = L(\boldsymbol{r})^\dagger \boldsymbol{\Lambda}$ and $L(\boldsymbol{r})^\dagger$ is the pseudo-inverse of the Laplacian matrix $L(\boldsymbol{r})$ defined in Eq. (S62). To compute the gradient $\frac{\partial \mathcal{O}}{\partial r_{ij}}$, we need to evaluate the response of $\boldsymbol{\mu}^*$ to the variation of the control parameters $\boldsymbol{r}$, i.e.,

$$\frac{\partial \mathcal{O}}{\partial r_{ij}} = \sum_{(p,q) \in \mathcal{T}} \frac{\partial \mathcal{O}}{\partial x_{pq}^*} \left[ \frac{1}{r_{pq}} \left( \frac{\partial \mu_q^*}{\partial r_{ij}} - \frac{\partial \mu_p^*}{\partial r_{ij}} \right) - \frac{1}{r_{pq}^2} (\mu_q^* - \mu_p^*) \delta_{(i,j),(p,q)} \right]. \tag{S83}$$

Furthermore, it requires to compute $\frac{\partial L(\boldsymbol{r})^\dagger}{\partial r_{ij}}$ for $\frac{\partial \boldsymbol{\mu}^*}{\partial r_{ij}}$. Assuming the underlying graph will not fragment into multiple components when adapting $\boldsymbol{r}$, $L(\boldsymbol{r})$ has a constant rank, then we have [18]

$$\frac{\partial L(\boldsymbol{r})^\dagger}{\partial r_{ij}} = -L^\dagger \frac{\partial L}{\partial r_{ij}} L^\dagger + L^\dagger (L^\dagger)^\top \left( \frac{\partial L^\top}{\partial r_{ij}} \right) (I - LL^\dagger)$$
$$+ (I - L^\dagger L) \left( \frac{\partial L^\top}{\partial r_{ij}} \right) (L^\dagger)^\top L^\dagger. \tag{S84}$$

Using the property of the pseudo-inverse of the Laplacian $L^\dagger L = LL^\dagger = I - \frac{1}{N} \mathbf{1} \mathbf{1}^\top$ and $\frac{\partial L}{\partial r_{ij}} \cdot \mathbf{1} = 0$ (with $\mathbf{1}$ as the all-one vector), we have

$$\frac{\partial L(\boldsymbol{r})^\dagger}{\partial r_{ij}} = -L^\dagger \frac{\partial L}{\partial r_{ij}} L^\dagger, \tag{S85}$$

which closes the equations for calculating the gradients.

---

**Algorithm 3:** Message-passing algorithm for (undirected) flow control (bilevel, Method II to treat the reference node $\mathcal{D}$ )

---

**Input:** Undirected network $G(V, E)$ (pre-processed to remove dangling nodes), node parameters $\{\Lambda_i\}$ defining the resources, maximal number of iterations $T$, time interval $t_{\text{update\_intv}}$ for performing gradient descent.

Initialize the messages $\{\alpha_{i\rightarrow j}, \hat{x}_{i\rightarrow j}\}$ randomly.

Initialize the gradients $\{\partial\mathcal{O}_{pq}/\partial\alpha_{i\rightarrow j}, \partial\mathcal{O}_{pq}/\partial\hat{x}_{i\rightarrow j}, \partial\mathcal{O}_{pq}/\partial r_{i\rightarrow j}, \mathcal{O}_{pq}/\partial r_{pq}\}$ randomly.

Initialize the edge control parameters $\{r_{ij}\}$.

**for** $t$ in $1 : T$ **do**

    Randomly select a node $i$ and one of its adjacent node $j \in \mathcal{N}_i$.

    Update the messages $m_{i\rightarrow j} \in \{\alpha_{i\rightarrow j}, \hat{x}_{i\rightarrow j}\}$ as:

$$\alpha_{i\rightarrow j} \leftarrow \frac{1}{\sum_{k\in\mathcal{N}_i\setminus j} \alpha_{k\rightarrow i}^{-1}} + r_{ij}, \qquad \hat{x}_{i\rightarrow j} \leftarrow \frac{\Lambda_i + \sum_{k\in\mathcal{N}_i\setminus j} \hat{x}_{k\rightarrow i}}{1 + r_{ij}\sum_{k\in\mathcal{N}_i\setminus j} \alpha_{k\rightarrow i}^{-1}}.$$

    **for** $(p,q) \in \mathcal{T}$ **do**

        **if** $(i,j) = (p,q)$ or $(j,i) = (p,q)$ **then**

            /* $(i,j)$ is a targeted edge */

            Update the gradients $\partial\mathcal{O}_{pq}/\partial m_{p\rightarrow q}, \partial\mathcal{O}_{pq}/\partial m_{q\rightarrow p}, \partial\mathcal{O}_{pq}/\partial r_{pq}$ as

$$\frac{\partial\mathcal{O}_{pq}}{\partial m_{p\rightarrow q}} \leftarrow \frac{\partial\mathcal{O}_{pq}}{\partial x_{pq}^*} \frac{\partial x_{pq}^*}{\partial m_{p\rightarrow q}}$$

$$\frac{\partial\mathcal{O}_{pq}}{\partial m_{q\rightarrow p}} \leftarrow \frac{\partial\mathcal{O}_{pq}}{\partial x_{pq}^*} \frac{\partial x_{pq}^*}{\partial m_{q\rightarrow p}}$$

$$\frac{\partial\mathcal{O}_{pq}}{\partial r_{pq}} \leftarrow \frac{\partial\mathcal{O}_{pq}}{\partial x_{pq}^*} \left[\frac{\partial x_{pq}^*}{\partial r_{pq}} + \sum_{m\in\{\alpha,\hat{x}\}} \left(\frac{\partial x_{pq}^*}{\partial m_{p\rightarrow q}} \frac{\partial m_{p\rightarrow q}}{\partial r_{pq}} + \frac{\partial x_{pq}^*}{\partial m_{q\rightarrow p}} \frac{\partial m_{q\rightarrow p}}{\partial r_{pq}}\right)\right]$$

        **else**

            /* $(i,j)$ is a non-targeted edge */

            Update the gradients $\partial\mathcal{O}_{pq}/\partial m_{i\rightarrow j}, \partial\mathcal{O}_{pq}/\partial r_{ij}, \forall (p,q) \in \mathcal{T}$ as

$$\frac{\partial\mathcal{O}_{pq}}{\partial m_{i\rightarrow j}} \leftarrow \sum_{l\in\mathcal{N}_j\setminus i} \sum_{m_{j\rightarrow l}\in\{\alpha_{j\rightarrow l}, \hat{x}_{j\rightarrow l}\}} \frac{\partial\mathcal{O}_{pq}}{\partial m_{j\rightarrow l}} \frac{\partial m_{j\rightarrow l}}{\partial m_{i\rightarrow j}},$$

$$\frac{\partial\mathcal{O}_{pq}}{\partial r_{i\rightarrow j}} \leftarrow \sum_{m_{i\rightarrow j}\in\{\alpha_{i\rightarrow j}, \hat{x}_{i\rightarrow j}\}} \frac{\partial\mathcal{O}_{pq}}{\partial m_{i\rightarrow j}} \frac{\partial m_{i\rightarrow j}}{\partial r_{ij}}.$$

        **end**

    **end**

    **if** $t \bmod t_{\text{update\_intv}} = 0$ **then**

        $\frac{\partial\mathcal{O}}{\partial r_{ij}} \leftarrow 0.$     /* To compute the full gradient: */

        **for** $(p,q) \in \mathcal{T}$ **do**

            **if** $(i,j) = (p,q)$ or $(j,i) = (p,q)$ **then**

                $\frac{\partial\mathcal{O}}{\partial r_{ij}} \leftarrow \frac{\partial\mathcal{O}}{\partial r_{ij}} + \frac{\partial\mathcal{O}_{pq}}{\partial r_{pq}}.$

            **else**

                $\frac{\partial\mathcal{O}}{\partial r_{ij}} \leftarrow \frac{\partial\mathcal{O}}{\partial r_{ij}} + \frac{\partial\mathcal{O}_{pq}}{\partial r_{i\rightarrow j}} + \frac{\partial\mathcal{O}_{pq}}{\partial r_{j\rightarrow i}}.$

            **end**

        **end**

        $r_{ij} \leftarrow r_{ij} - s\frac{\partial\mathcal{O}}{\partial r_{ij}}$     /* Update the control parameter with step size $s$. */

        $r_{ij} \leftarrow \min(\max(r_{ij}, r_{\min}), r_{\max})$     /* To respect the box constraints of the control parameters. */

    **end**

    **if** messages converge **then**

        Exit the for loop.

    **end**

**end**

For each edge $e$, compute the equilibrium flow as $x_{ij}^* = \frac{\alpha_{j\rightarrow i}\hat{x}_{j\rightarrow i} - \alpha_{i\rightarrow j}\hat{x}_{i\rightarrow j}}{\alpha_{i\rightarrow j} + \alpha_{j\rightarrow i} - r_{ij}}$.

**Output:** Convergence status, equilibrium flows $\{x_{ij}^*\}$ and optimal control parameters $\{r_{ij}\}$

# REFERENCES

[1] Alessandro Lonardi, Enrico Facca, Mario Putti, and Caterina De Bacco, "Designing optimal networks for multicommodity transport problem," Phys. Rev. Research **3**, 043010 (2021).

[2] Vincenzo Bonifaci, Enrico Facca, Frederic Folz, Andreas Karrenbauer, Pavel Kolev, Kurt Mehlhorn, Giovanna Morigi, Golnoosh Shahkarami, and Quentin Vermande, "Physarum-inspired multi-commodity flow dynamics," Theoretical Computer Science **920**, 1–20 (2022).

[3] Chi Ho Yeung, David Saad, and K. Y. Michael Wong, "From the physics of interacting polymers to optimizing routes on the london underground," Proceedings of the National Academy of Sciences **110**, 13717–13722 (2013).

[4] K. Y. Michael Wong and David Saad, "Inference and optimization of real edges on sparse graphs: A statistical physics perspective," Phys. Rev. E **76**, 011115 (2007).

[5] Bo Li, David Saad, and Andrey Y. Lokhov, "Reducing urban traffic congestion due to localized routing decisions," Phys. Rev. Research **2**, 032059 (2020).

[6] Benoît Colson, Patrice Marcotte, and Gilles Savard, "An overview of bilevel optimization," Annals of Operations Research **153**, 235–256 (2007).

[7] Joaquim Dias Garcia, Guilherme Bodin, davide-f, Ian Fiske, Mathieu Besançon, and Nick Laws, "joaquimg/bileveljump.jl: v0.4.1," https://doi.org/10.5281/zenodo.4556393 (2021), 10.5281/zenodo.4556393.

[8] Andreas Wächter and Lorenz T. Biegler, "On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming," Mathematical Programming **106**, 25–57 (2006).

[9] Athanasios Migdalas, "Bilevel programming in traffic planning: Models, methods and challenge," Journal of Global Optimization **7**, 381–405 (1995).

[10] Chi Ho Yeung and David Saad, "Competition for shortest paths on sparse graphs," Phys. Rev. Lett. **108**, 208701 (2012).

[11] Chi Ho Yeung, "Efficient algorithm for routing optimization via statistical mechanics," in *2013 IEEE International Conference on Communications Workshops (ICC)* (2013) pp. 1420–1424.

[12] C. De Bacco, S. Franz, D. Saad, and C. H. Yeung, "Shortest node-disjoint paths on random graphs," Journal of Statistical Mechanics: Theory and Experiment **2014**, P07009 (2014).

[13] Chi Ho Yeung, K. Y. Michael Wong, and Bo Li, "Coverage versus supply cost in facility location: Physics of frustrated spin systems," Phys. Rev. E **89**, 062805 (2014).

[14] Christopher M. Bishop, *Pattern Recognition and Machine Learning* (Springer-Verlag, Berlin, Heidelberg, 2006).

[15] Chaisak Suwansirikul, Terry L. Friesz, and Roger L. Tobin, "Equilibrium decomposed optimization: A heuristic for the continuous equilibrium network design problem," Transportation Science **21**, 254–263 (1987).

[16] Chi Ho Yeung and K. Y. Michael Wong, "Optimal location of sources in transportation networks," Journal of Statistical Mechanics: Theory and Experiment **2010**, P04017 (2010).

[17] Patrick Rebeschini and Sekhar Tatikonda, "A new approach to laplacian solvers and flow problems," Journal of Machine Learning Research **20**, 1–37 (2019).

[18] G. H. Golub and V. Pereyra, "The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate," SIAM Journal on Numerical Analysis **10**, 413–432 (1973).