

Predicting human mobility through the assimilation of social media traces into mobility models

Mariano G. Beiró¹

André Panisson¹

Michele Tizzoni¹

Ciro Cattuto¹

¹ ISI Foundation, Turin, Italy

Abstract

Predicting human mobility flows at different spatial scales is challenged by the heterogeneity of individual trajectories and the multi-scale nature of transportation networks. As vast amounts of digital traces of human behaviour become available, an opportunity arises to improve mobility models by integrating into them proxy data on mobility collected by a variety of digital platforms and location-aware services. Here we propose a hybrid model of human mobility that integrates a large-scale publicly available dataset from a popular photo-sharing system with the classical gravity model, under a stacked regression procedure. We validate the performance and generalizability of our approach using two ground-truth datasets on air travel and daily commuting in the United States: using two different cross-validation schemes we show that the hybrid model affords enhanced mobility prediction at both spatial scales.

1 Introduction

Modelling and understanding human mobility patterns at different spatial scales and aggregation levels - from single individuals to population displacements - is an important research topic because of a vast number of applications, ranging from urban and transportation planning [1, 2] and resource allocation [3, 4] to the prediction of migration flows [5, 6] and epidemic spreading at local, regional or worldwide level [7, 8, 9, 10].

In the last few years, a significant research effort has been made to understand human mobility patterns, both in the laws governing individual human trajectories [11, 12] and collective movements [13, 6, 14]. In the latter case, the most extensively used models are the gravity model [15, 16] and the more recent radiation model [6]. The gravity model assumes that the amount of people travelling between two locations is directly proportional to some power of their populations, and decays as some power of the distance between them. Instead, the radiation model considers human movements as diffusion processes that depend on the population distribution over the space, reproducing Stouffer's theory of intervening opportunities [17]. Both models are static and require some information in order to be adjusted: in the gravity model, parameters are fitted using real mobility data, provided by an independent source; the radiation model, in its original formulation, is parameter-free, but it requires accurate knowledge of the spatial population distribution. Both modelling approaches have been extensively tested, showing advantages and limitations. The gravity model has been successfully used to describe highway flows [18], air-travel [19, 20], commuting [8] and mobile phone calls between cities [21]. However, it has some relevant limitations, as the availability of data for calibration and the lack of a first principle derivation [6, 22]. On the other hand, the radiation model offers very good predictions for commuting patterns between U.S. counties using only population data, but its applicability at different spatial scales has been debated since it does not succeed in capturing commuting inside urban or metropolitan areas [23, 22, 24], and it has never been used to model long distance travel patterns either.

The limitations of these models suggest that the quality of their results can be largely improved if they are supported by additional data [25, 24]. In fact, several works have analyzed records from mobile phone companies to study individual [11] and collective mobility [26, 27, 28], showing that it is possible to infer these flows from human activity. The mobility flows obtained in this way can be successfully used for the prediction of epidemic spreading [29, 30], as a proxy for the real, often inaccessible, mobility data.

In this context, the large volumes of digital traces left by humans over the Internet allows for a better understanding of mobility processes, with immediate benefits. On the one hand, the increase in transport infrastructure

during the last years and the dynamics of change of mobility patterns have brought the requirement of real-time modelling. People travel more, and travel patterns may change very fast, with important consequences for epidemic spreading and planning. A timely modelling of mobility processes might then allow for rapid interventions and for the design of emergency policies. On the other hand, though mobility data is usually available in many developed countries for airline transportation, train trips or commuting, in many underdeveloped countries this information is scarce or does not exist at all. The fact that mobility datasets are aggregated at a particular resolution level also constitutes a limitation for many potential studies.

The scientific community recently recognized that one of the challenges in the modelling of social and epidemic processes is the assimilation of geolocalized data, and the construction of hybrid models combining metapopulation and network models with individual traces [31, 32]. Our approach to human mobility is in line with this perspective, as we analyze the effects of incorporating geolocalized traces from social media into the classical gravity model.

Social media platforms like Flickr [33], Twitter [34] or Foursquare [35] offer the possibility of georeferencing the content shared by users. Thus, they constitute a timely source of disaggregated, high-resolution spatio-temporal data on human mobility. The advantage of social media traces with respect to other sources of digital information is that they can be publicly accessed and at a very low cost. This approach have been taken by recent works in the literature, showing that mobility patterns can be successfully extracted from social media traces. Lenormand *et al.* used traces from Twitter to study highway and roadway transportation networks in Europe [2]; Noulas *et al.* used a Foursquare dataset to analyze the link between user activity and place transitions [36]; Hawelka *et al.* modelled international travel of Twitter users by residence country [14]; Lenormand *et al.* have also used Twitter traces to model commuting from home to work [37]; Grabowicz *et al.* studied the relation between human mobility and interactions using traces from different social networks [38]; Barchiesi et al. [39] used Flickr data from 16,000 individuals in the UK to model the flows between its 20 largest cities, comparing their results with travel data obtained from surveys.

In this work we used a set of 18,000,000 timestamped, georeferenced pictures from Flickr, taken by 40,000 users in the U.S., which are part of the Yahoo Flickr Creative Commons 100M public dataset [40]. We processed the sequences of pictures belonging to each individual user in order to extract user trip paths at different resolution levels. Then, we used these emerging collective flow patterns to feed a learning model based on the gravity law.

Our main contribution is to design a data-driven hybrid model of human mobility, in which social media traces are combined with the classical gravity model under a machine learning approach, by training and cross-validating with real datasets. We evaluate the model for two different human activities and resolution scales: an air-travel network and a daily commuting network. Firstly, we show how individual traces can be adapted to these different resolution scales, by tessellating the space into adequate basins and filtering the correct individual flows. Secondly, we combine these traces with the gravity law and we fit the resulting hybrid model using a subset of the real data. Then, we evaluate the fitted model using the remaining part of the dataset. With a cross-validation procedure, we show that the hybrid model can be fitted using a small portion of the data as training set, to correctly predict the remaining mobility flows. In fact, we observe that the incorporation of Flickr traces into the gravity model improves its performance significantly, measured in terms of the determination coefficient. Our findings show that the Flickr traces are representative of the real human mobility and that they can be assimilated into a more theoretical model as the gravity model. Moreover, this procedure can be applied in other cross-validation contexts in which there is scarce information on mobility, by combining the available data with digital traces from social media.

2 Results

We processed the traces left by 40,000 Flickr users in the U.S. (about 18 million pictures) in order to obtain mobility flow matrices which represent human flows between pairs of geographical nodes, looking for the collective mobility patterns of two types of human activities at different resolution scales: air travel and daily commuting. We also used two real mobility datasets as a ground truth: the RITA dataset of air travel in the U.S. [41], and the commuting data provided by the U.S. Census Bureau [42]. At each resolution scale, flows were aggregated into geographic basins. For air travel, our geographic basins were defined by the presence of an airport, while for the commuting network each basin corresponded to a U.S. county (see the *Methods* for more details on the ground truth mobility datasets).

The construction of the Flickr flow matrices at the airport and county level from the users' traces is described in *Methods: Flickr-based flow matrices*. In short, we added a connection between pairs of basins, i and j , every

time a user took a picture in basin i and the subsequent one in basin j . Each connection is then weighted by the total number of users who travelled between i and j .

The distance between the locations of two consecutive pictures taken by the same user ranges from a few meters to thousands of kilometers, showing a heavy-tailed behavior with an exponential cutoff (see the *Supplementary material* for more details). Thus, it is clear that users' traces can provide information about very different types of human displacements, ranging from a short walk within a city to long distance trips or international travel. In each case, and when modelling a particular type of mobility, it is important to consider only the relevant traces for that type of movements, otherwise, the effects of different activities will be mixed. We analyzed the effect of distance on both the daily commuting and the air travel, and we set a maximum trip distance of 100 km for daily commuting in the U.S. and a minimum trip distance of 500 km for U.S. air travel. The analysis is described in *Methods: Distance thresholds*.

After aggregating the trips by basin and filtering by distance, we obtain the Flickr flow matrices denoted as $\mathbf{F}_r = (f_{ij}^r)$ for the air travel network, and $\mathbf{F}_c = (f_{ij}^c)$ for the commuting network. Here, f_{ij}^c and f_{ij}^r represent the number of Flickr users travelling from basin i to basin j in each of the networks.

Model definition and fitting. A model of human mobility should be able to predict real mobility flows. The aim of our work was to predict mobility flows of the air travel network and the U.S. commuting network, denoted as $\mathbf{Y}_r = (r_{ij})$ and $\mathbf{Y}_c = (c_{ij})$ respectively. Here, r_{ij} represents the number of travellers between two airports i and j , while c_{ij} represents the daily amount of commuters between any two counties i and j . Adopting a machine learning approach, we represent this prediction problem as a regression task in which the model is first trained, and then it is used to estimate the real flows, \mathbf{Y}_r or \mathbf{Y}_c , which are the so called target values.

The classical gravity model (described in the *Methods*) estimates the target values as directly proportional to some powers of the population of the origin and destination basins and inversely proportional to some increasing function of the distance between them (typically, a power-law or exponential function). We will indicate this model as $\mathbf{G}(\alpha, \beta, \gamma; \mathcal{P})$, where α, β and γ are the three exponents in the gravity law, and $\mathcal{P} = (p_i)$ is a vector containing the populations of the basins.

Analogously, travel flows estimated from Flickr traces can be described by a model in which we assume that the value y_{ij} , representing the mobility flow between two basins i and j , is proportional to the number of users' trips from i to j . We represent this model as $C_F \cdot \mathbf{F}$, where \mathbf{F} represents the Flickr flow matrix at the corresponding resolution level (i.e., \mathbf{F}_r or \mathbf{F}_c), and C_F is a constant.

Our approach is to combine the flows of human mobility estimated from the two models under a stacked regression procedure [43] in which each model is fitted alone, and then a linear regression determines the weight of each of them. In this way we combine the Flickr traces and the gravity model to improve the prediction of the real mobility patterns.

The incorporation of the traces will be defined by a linear function (we omit the subscript indices to generalize to any of the two resolution levels):

$$\mathbf{H}(\alpha, \beta, \gamma, A, B; \mathcal{P}) = A \cdot \mathbf{G}(\alpha, \beta, \gamma; \mathcal{P}) + B \cdot \mathbf{F},$$

where $\mathbf{H} = (h_{ij})$ represents our hybrid model, \mathbf{G} is the gravity model and \mathbf{F} is the Flickr flows matrix; $\alpha, \beta, \gamma, A, B$ are real-valued fitting parameters. We will fit the model by minimizing the following loss function L :

$$L = \| \mathbf{Y} - \mathbf{H}(\alpha, \beta, \gamma, A, B; \mathcal{P}) \|,$$

where $\| \cdot \|$ denotes the Frobenius norm. Again, we removed the subscript indices, so that this equation represents what is done at each resolution level. The minimization process is made under the stacked regression assumption that each of the two component models can be fitted alone, and then a linear regression determines the weight of each of them. It is important to notice that, as \mathbf{Y}_r and \mathbf{Y}_c are sparse matrices, we only evaluate L for those connections for which there is a positive flow of travellers.

Model evaluation. The performance of a learning model is measured by its capacity to generalize to flows that were not known during the training step. Then, in order to validate our model we used a cross-validation strategy in which we fitted the model using only part of the target flows \mathbf{Y} (training set) and then we tested its performance using the remaining flows (test set). We measured the performance in terms of Pearson correlation coefficient ρ and the determination coefficient r^2 between the target flows and the predicted flows.

Our procedure follows a 10-fold cross-validation scheme [44]: the real dataset is divided into 10 parts or folds, and for each fold we use the 10% of the target values as test set, and the remaining 90% as training set. Thus,

each sample in the dataset was tested once, using a model that was not fitted with that sample. The performance of the different models is shown in Table 1 in terms of the Pearson correlation coefficient ρ and the determination coefficient r^2 between the real flows and the predicted flows. The table shows also the results for the gravity model alone ($A = 1, B = 0$) and the Flickr model alone ($A = 0, B = 1$). To note, all the results were cross-validated, i.e. each individual flow was estimated with a fit that excluded it from the training set.

Notice that the low determination coefficients for the gravity model alone are related to its limitations for capturing the large flows across distant cities, due to the fast decay of the model as a function of distance. This had also been observed in [6]. Indeed, by plotting the prediction ratio of the gravity model for flows above 100 passengers in the red curves of Fig. 1 we clearly observe that the gravity model tends to underestimate large flows. Instead, the hybrid model captures correctly the large distance flows (violet curves of Fig. 1).

Model	Commuting		Air travel	
	ρ	r^2	ρ	r^2
Gravity alone	0.69	0.41	0.68	0.40
Flickr alone	0.69	0.47	0.78	0.62
Hybrid model	0.79	0.62	0.84	0.72

Table 1: **Cross-validated hybrid model performance (10-fold cross-validation).** The table shows the performance of our hybrid model in terms of the Pearson correlation coefficient ρ and the determination coefficient r^2 . We also show the results for the gravity model alone and the Flickr traces on their own. All the values were produced under a 10-fold cross-validation scheme.

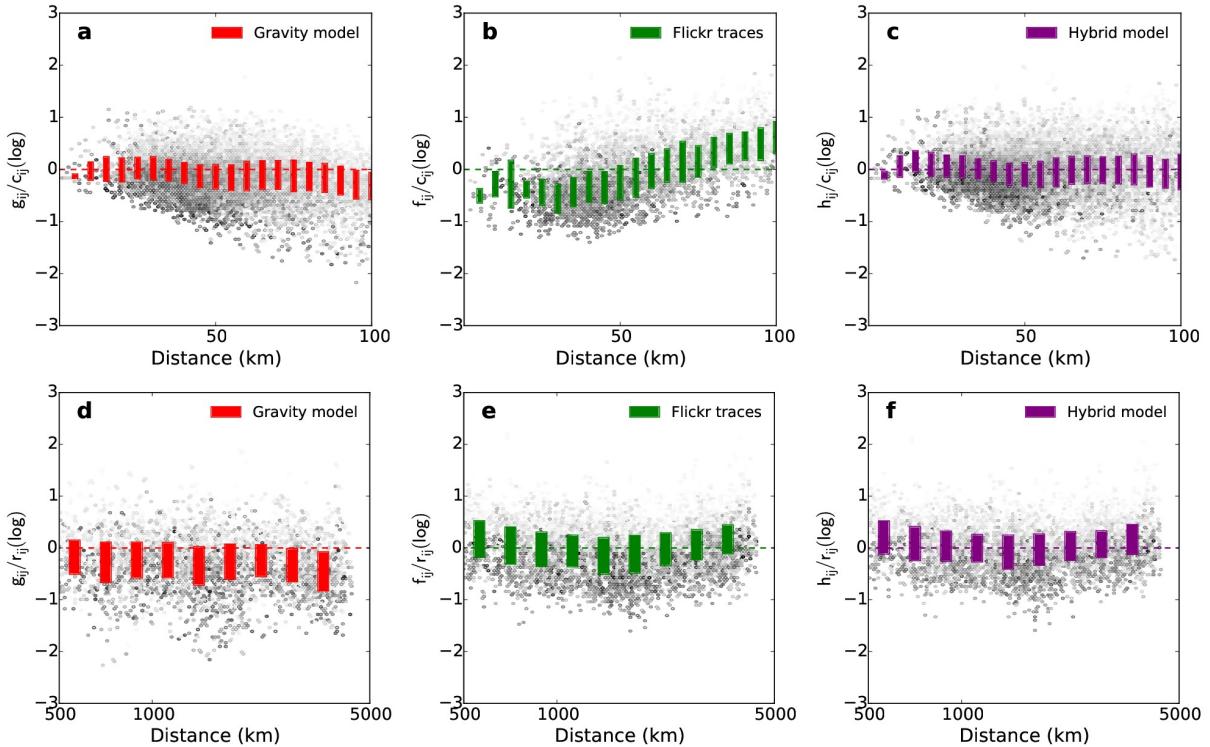


Figure 1: **Prediction ratio for large flows as a function of distance.** The points in these figures represent predicted flows between pairs of nodes under a 10-fold cross-validation scheme. (a, b, c) Commuting network of the U.S. (d, e, f) Air-transportation network of the U.S. (We only show flows above 100 passengers). The color intensity at each point represents the total passenger flows aggregated under a certain distance and prediction ratio.

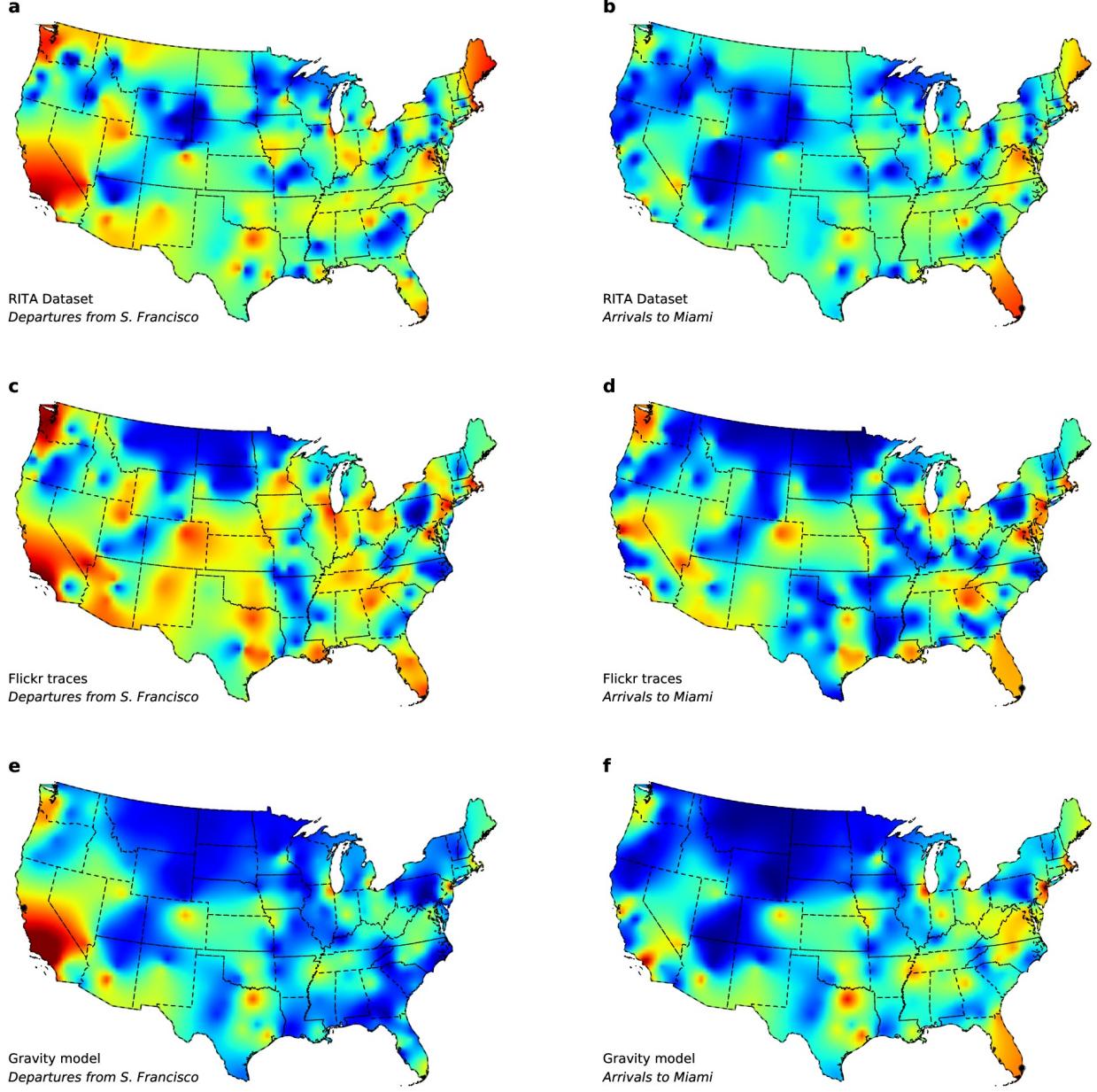


Figure 2: Departures and arrivals in the air transportation network. Heatmaps representing the distribution of predicted trips departing from San Francisco towards any point in the U.S. (left column), or arriving to Miami from any point in the U.S. (right column). **(a, b)** Ground-truth from the RITA air travel dataset. **(c, d)** Flickr traces. **(e, f)** Gravity model.

On the other hand, according to Table 1, the performance of the Flickr model alone is comparable to that of the gravity model for the commuting network, and is significantly higher in the air travel network. This reveals that Flickr users' trips can be a good proxy of collective human mobility at different resolution scales, and can be particularly useful when real data for fitting the gravity model is not available.

A more qualitative evaluation of the geographical spanning of Flickr traces is offered in Fig. 3 and 2, which show the distribution of the origins and destinations of travellers according to the real datasets, the Flickr traces and the gravity model. Fig. 3 shows the commuting patterns in the New York State, where the color intensity

represents the amount of people commuting from one county to New York City (in black). Here, we observe that the gravity law correctly captures the flows from neighbouring counties, which in fact represent more than 95% of the commuting flows, but it underestimates long distance flows. Instead, the traces from Flickr have a slower distance decay, more in accordance with census data. Something similar is observed in the air travel network, as depicted in Fig. 2. Looking at trips departing from San Francisco (left panels), the gravity law correctly predicts that the largest flows are those towards Los Angeles (CA), San Diego (CA), Las Vegas (NV) and Seattle (WA), but those directed to the East Coast are generally underestimated. For the arrivals to Miami (right panels) the gravity shows a good performance. In both cases, the flows from Flickr are a representative sample of all the territory. This large span is also confirmed when we observe the trip distance distribution of Flickr users (see the *Supplementary Information*).

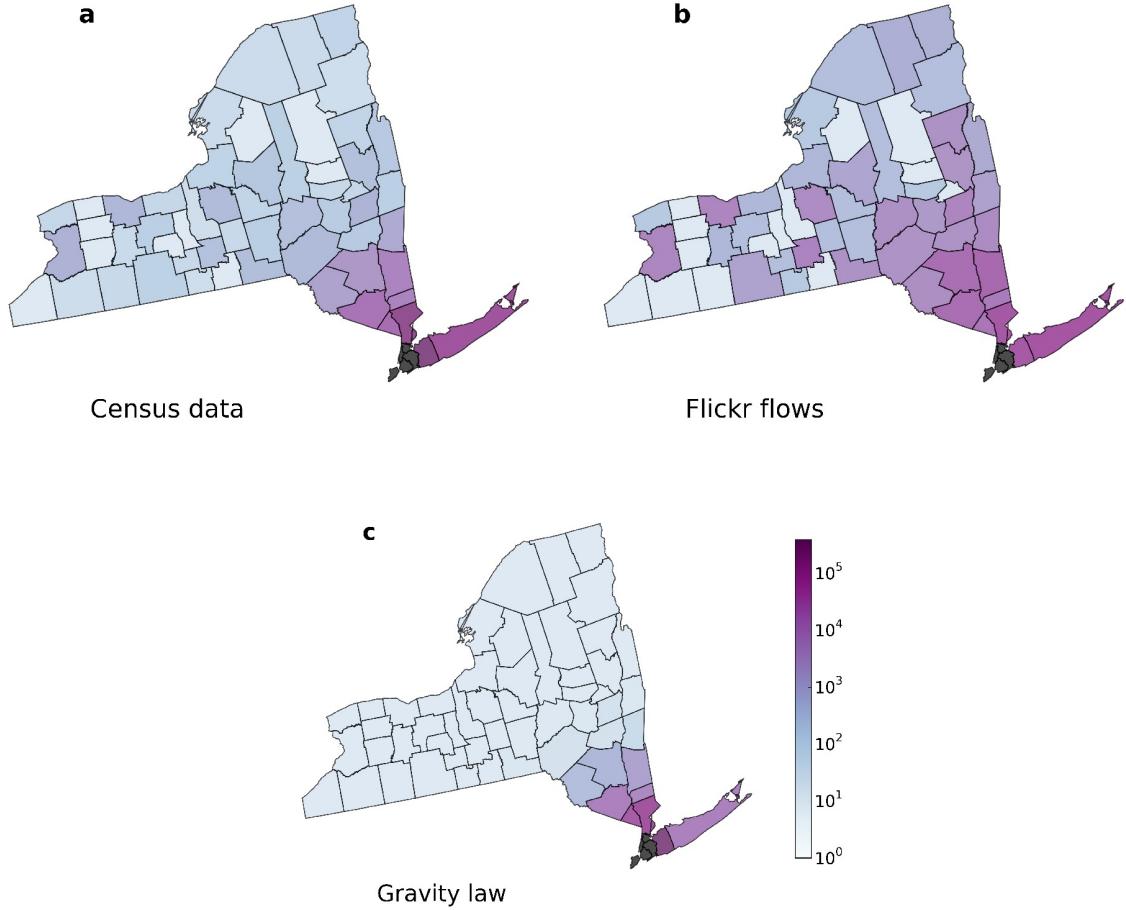


Figure 3: Daily commuting to New York City. Each panel shows the daily amount of commuters arriving to New York City from the different counties of the New York State. Here, we aggregated all the commuters traveling to the 5 boroughs of New York City. (a) U.S. Census Bureau data. (b) Flickr traces. (c) Gravity law. Flickr traces are shown without taking into account any distance threshold.

The incorporation of the Flickr traces into the gravity model produced a significant increase in the predictive performance both for daily commuting as for air travel in the U.S., as shown in the last row of Table 1. We observed a relative increase of $\sim 8 - 20\%$ in the Pearson correlation ρ between the model predicted flows and the real flows. The predictive power has also improved, as verified by a relative increase of $> 50\%$ in the determination coefficient r^2 (i.e., to what extent the model accounts for the real flows). The ratio between predicted flows and real flows as a function of distance has also improved, as shown in the right panels of Fig. 1.

Fig. 4 shows several 2D-histograms for the gravity model alone and for the hybrid model, together with boxplots grouped by the real flow values. The plots compare the model predictions to the real mobility data for commuting and air travel, and each cell represents pairs of nodes (counties in the case of commuting, and airport basins for air travel). Despite the fact that most of the interquartile ranges do not change in average, it is interesting to observe the changes on the right side of the pictures: for the models incorporating Flickr data, the interquartile ranges become shorter for large values of flows, and at the same time they are better aligned with the diagonal line representing the perfect matching between model and real data. These means that the improved model is both more precise and more exact for pairs of nodes connected by large mobility flows. This behavior is clearly seen if we compute the determination coefficient of the model after filtering flows above a certain value. We show this analysis in Fig. 5 (panels a,b), where the x -axis represents the flow threshold and the y -axis is the determination coefficient between the models and the real flows, restricted to pairs of nodes for which the real flow is larger than the threshold. The contribution of the Flickr traces is particularly significant in the air-travel network, where they duplicate the gravity model in a predictive capacity. In the commuting network both models are close, but the hybrid model almost duplicates the gravity performance for flows above 2000 travellers.

In the following two subsections we validate the model under different geographical and random data availability constraints. The *Supplementary Information* also includes a test based on the Sorensen-Dice coefficient, in which we evaluated the model performance for different distances and populations. This test shows that flows towards largely populated basins are the ones most improved by the hybrid model.

Prediction under data availability constraints. We evaluated the power of our regression model as a function of the training set size in order to see the minimum amount of data required to reach the largest improvement in the predictive power of both the gravity law and the Flickr flows. The analysis was performed under a cross-validation scheme with repeated random subsampling (bootstrapping). The advantage of this procedure is that the training set size can be varied in small steps (while k -fold cross-validation only allows for a minimum training set size of 0.5). The results of this test are shown in Fig. 5. For the air transportation network, we observe that with a training size of 1% the gravity law is already close to its best prediction levels, while from that value onwards the assimilation of flows from the Flickr traces starts producing an improvement in the performance. With a training size of 3% this improvement is already quite significant. In the commuting network, training with 1% of the ground-truth flows is also enough for the gravity law, but we need a 10% of Flickr flows in order to observe a significant improvement in the model performance.

Spatial cross-validation: U.S. West Coast vs. U.S. East Coast. To extend the analysis of the generalizability of the model, we performed a geographical 2-fold cross-validation in which we split the contiguous United States into two roughly symmetric parts, taking by reference the meridian -102° . We trained the model in one half of the U.S. and we then used it to predict the mobility flows inside the other half. To make our test independent from the specific choice of the spatial partition, we disregarded crossed flows, that is, the traffic flows connecting two points in different halves of the partition. Table 2 shows the performance in terms of the Pearson correlation ρ and the determination coefficient r^2 . We see that, in this case, the performance of Flickr traces is inferior with respect to the 10-fold cross-validation values of Table 1, denoting the presence of spatial inhomogeneities in the users' activities across the U.S., especially for the commuting network. However, the performance of the hybrid model is still improved by the assimilation of the Flickr flows.

Model	Commuting		Air travel	
	ρ	r^2	ρ	r^2
Gravity alone	0.72	0.41	0.68	0.40
Flickr alone	0.60	0.34	0.72	0.43
Hybrid model	0.78	0.46	0.81	0.49

Table 2: **Geographically cross-validated hybrid model performance (U.S. West Coast vs. U.S. East Coast).** Performance of our hybrid model in terms of the Pearson correlation coefficient ρ and the determination coefficient r^2 . The predicted values were produced by training the model in the West Coast and then validating in the East Coast and viceversa.

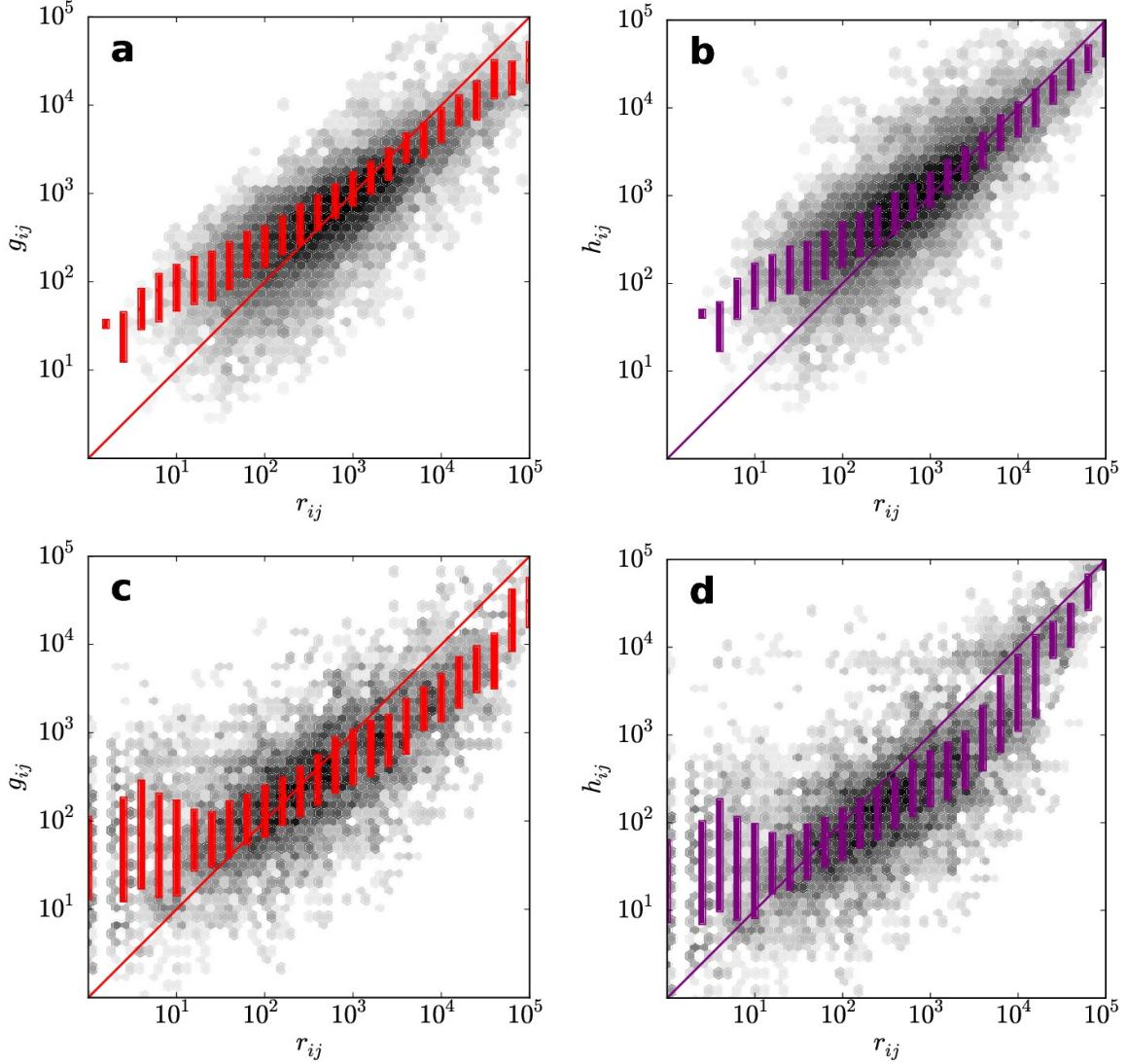


Figure 4: **Mobility models predictions** Each point in the 2D-histograms represents flows with some real/estimated flow value relation. The color of the points in a gray scale represents the frequency values. The boxplots in each panel correspond to the interquartile ranges. **(a, b)** U.S. Commuting network. **(c, d)** U.S. air travel network.

3 Discussion

Theoretical models of mobility offer the possibility of predicting human flows when calibrating data is available, which is the case in most developed countries, but not in all the world. The gravity model is currently the prevailing one, and has been successfully used in different contexts and scales. However, it has some limitations, as the underestimation of large flows and its fast decay with distance. On the other hand, the more recent radiation model requires less calibration data – only the population distribution – but has a limited geographical spanning, being usually applied for commuting at regional scale.

While in some cases first-principled theoretical approaches can be pursued [45], the increasing amount of geolocalized data publicly available through the Internet and social media suggests a different perspective for modelling mobility, which is the incorporation of large volumes of digital traces into theoretical models.

In this work we followed this approach and we showed that geolocalized traces collected from social media as

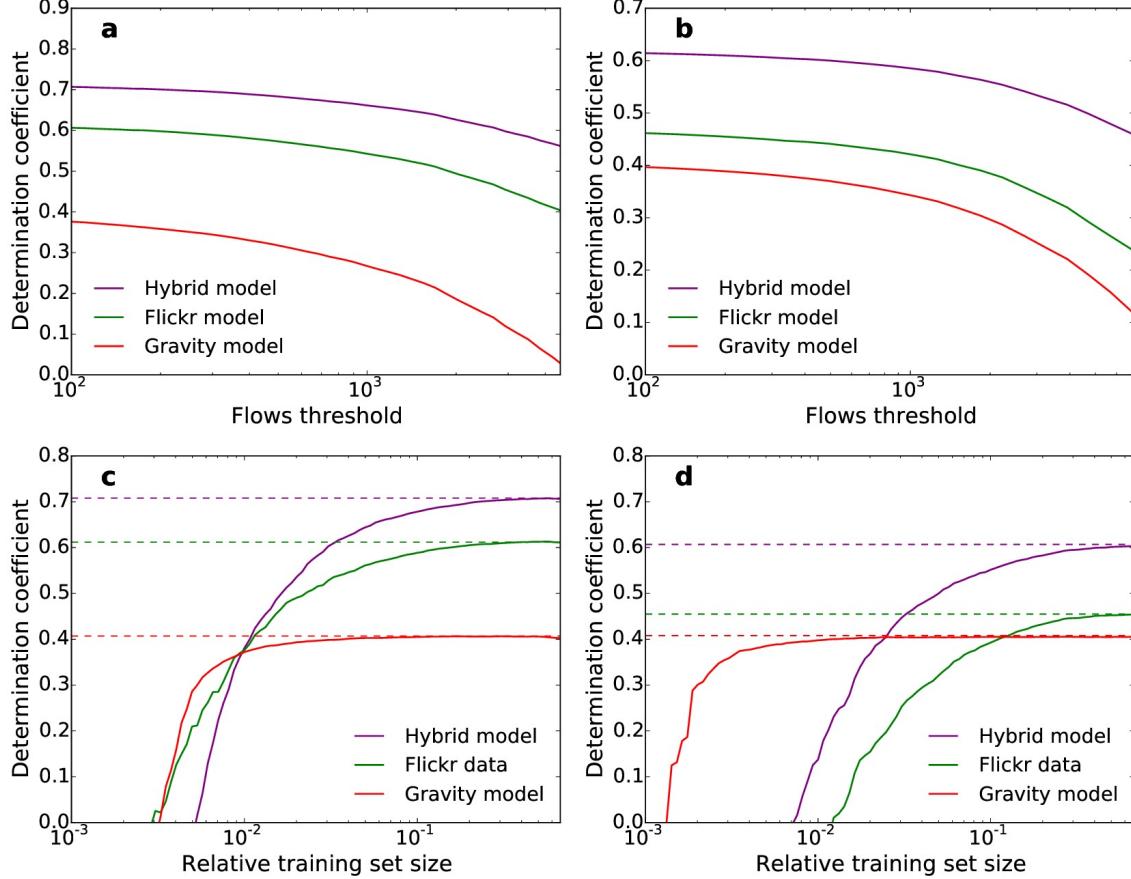


Figure 5: (a, b) **Determination coefficient for the prediction of large flows, as a function of the flow threshold.** Performance of the models measured by the determination coefficient r^2 when restricted to pairs of nodes for which the real flows are above a given threshold value. (a) Air transportation data. (b) Commuting data. (c, d) **Prediction error rate for different training set sizes.** Prediction error in terms of the determination coefficient, as a function of the training set relative size. The training set sizes were varied in a logarithmic scale. Dotted lines correspond to the maximum attained values. (c) Air transportation data. (d) Commuting data.

those from Flickr can successfully inform a predictive model of collective human mobility. Even though individual mobility trajectories can be noisy, their aggregation is representative of collective mobility patterns. Our results showed both qualitatively and quantitatively the high performance of Flickr traces in the predictions, at different distances and for different resolution levels.

The use of traces from Flickr as a proxy for human mobility has also been explored by Barchiesi et al. [39], who successfully modelled the flows between the 20 largest cities in the UK using traces from 16,000 individuals, and considering a threshold value to fix the cities' diameters. They showed the agreement of their predictions with travel data obtained from surveys. Our results strengthen the capabilities of digital traces for predicting human mobility, as we show that they can be assimilated into the gravity to produce a data-driven model. Our model was tested at different resolution levels, for flows among over 200 nodes for the U.S. air travel network, and over 3,000 counties for the commuting network. Our predicted flows span over cities and counties of variable area and population, and at quite different resolution scales.

The assimilation of Flickr traces to produce a data-driven model addressed the limitations of the gravity model, as the availability of data and the underestimation of some large distance flows [6, 22]. Moreover, considering the potential data availability limitations in many situations, our cross-validation with bootstrapping approach showed that a small sample of data ($\sim 1\%$) is enough for its calibration. Thus, by assimilating the traces from Flickr users into the classical gravity model, we improved its predictions in several ways: we handled the underestimation issues,

we accounted for missing calibration data and we reinforced the model by providing it with individual traces.

Finally, we note that we used an open dataset as is the Flickr 100M dataset [40], which can be publicly accessed by anyone. However, Flickr is far from being the only appropriate platform for modelling mobility, and our research suggests the possibility of assimilating data from several sources simultaneously, in order to capture other types of activities and users. A future line of research also involves exploring other assimilation techniques different from the simple stacked regression.

In conclusion, our results expose two new directions in the modelling of mobility flows: first, that the assimilation of digital traces can reinforce classical mobility models as the gravity improving its predictions significantly. Second, our cross-correlation study proved that with a small number of real flows our model can be adjusted to give predictions close to its optima. These contributions are in line with current challenges in the study of spreading processes and social systems, as the assimilation of geolocalized data into network models and the construction of models at appropriate resolution scales.

4 Methods

4.1 Dataset description

We use two large real datasets as ground-truth for human mobility at quite different resolution levels:

- **Air travel in the U.S.** The RITA dataset from the *Airline Origin and Destination Survey* [41] collected by the *U.S. Bureau of Transport Statistics*, contains a 10% sample of all the domestic itineraries in the U.S. We used the subset correspondent to 2014, which comprises ~ 14 million air tickets between 466 airports.
- **Commuting in the U.S.** This dataset obtained from the *U.S. Census Bureau* [42] contains data on ~ 36 million commuters from the period 2009 – 2013 at the county level. It specifies a “*home county*” and a “*working county*” for each commuter.

4.2 Ground-truth flow matrices

The ground-truth flow network for the air transportation in the U.S. was built using the RITA dataset [41]. We define an airport basin as the area covered by an airport (i.e., the set of points for which that airport is the closest one). Thus, the partition of the U.S. territory into basins is built as the Voronoi tessellation given by the airport coordinates. We note that when two airports are at less than 30 km of distance, we consider that they represent a single airport basin (because they serve the same metropolitan area) and we replace them by a single airport before computing the Voronoi tessellation. We also generalized this idea to connected components of airports at less than 30 km.

Each ticket in the RITA dataset contains an itinerary formed by several (*coupons*). Each coupon contains information about the origin and destination airport of the trip and the number of passengers, and it also points out whether the destination was a stopover or either the passengers remained there. Those destinations in which the passengers did other than just a stopover are marked as *trip breaks* (TB); the first and last airports from an air ticket are always TB’s. By removing the stopovers from the itinerary we manage to clean the flow network from the presence of airport hubs, which do not represent the real passengers destination. For each itinerary, we obtain a list of destinations (d_1, d_2, \dots, d_n) (d_1 is the departing city), and we use it to build the flow network between the airport basins.

For the commuting dataset [42], each sample contains the “*home county*” and “*working county*” of one worker. We consider the list of U.S. counties as nodes for the commuting flow network, and for each pair of counties (c_i, c_j) we put a weighted link counting the number of workers that live in one of them and work in the other.

We use the notation $\mathbf{Y}_r = (r_{ij})$ and $\mathbf{Y}_c = (c_{ij})$ for the adjacency matrices which describe the airport and commuting ground-truth flows, respectively. We put the diagonal elements of both matrices to zero, because we shall not consider users that take a flight inside the same airport basin (we only observed 5 cases) or commuters which do not change county from home to work.

The air travel matrix \mathbf{Y}_r contains 204 airport basins with 30,472 links, with an amount of $\sum_{ij} r_{ij} = 40$ million flows. The commuting matrix \mathbf{Y}_c covers the 3,144 U.S. counties and has 55,578 links, describing the pattern of about 74 million commuters.

4.3 Flickr-based flow matrices

The traces left by Flickr users when they take and upload pictures are given by the coordinates of their geotagged pictures, ordered by the time in which they were taken. We only consider timestamped, geolocalized pictures taken in the U.S. For each user we obtain an array of pictures (p_1, p_2, \dots, p_n) sorted by timestamp. At a particular resolution level (airport or county) we will consider that the user makes a trip when two consecutive pictures have coordinates belonging to different basins. Then, we aggregate all the users' trips into a Flickr flow matrix.

- **County-level Flickr flow matrix:** We assign each picture to a county by considering the county borders as defined in the MAF/TIGER geographic database of the U.S. Census Bureau. We count one flow between counties (i, j) when two consecutive pictures (p_i, p_{i+1}) are taken in counties i and j respectively. We do not consider successive pictures in which the user does not change county.
- **Airport-level Flickr flow matrix:** We choose the airport basin closest to each picture. We count one flow between two airport basins (i, j) whenever two successive pictures are taken in basins i and j respectively. We shall not consider successive pictures in which the user does not change airport basin.

We note the adjacency matrices of these networks as $\mathbf{F}_c = (f_{ij}^c)$ and $\mathbf{F}_r = (f_{ij}^r)$. Both of them have zero diagonals. In total, we observed $\sum_{ij} f_{ij}^r \approx 350,000$ trips between airport basins and $\sum_{ij} f_{ij}^c \approx 520,000$ trips between counties. The flow networks contain $\approx 26,000$ and $\approx 150,000$ nonzero elements, respectively.

4.4 Distance thresholds

As the activity of Flickr users involves different modalities of mobility, it has to be correctly filtered before comparing it against the ground-truth flows. We observed that an important variable for this task is the distance between nodes. In Fig. 6 we compare the ground-truth flows for air travel and commuting against the Flickr flows at the county level and the airport level respectively. We observe that the Flickr flows have good agreement with the ground-truth when we consider distances above 500 km for air travel, and below 100 km for commuting. In fact, if we remove those flows from the data, we observe that the Flickr users trip distance distributions are consistent with the ground-truth trip distributions. The threshold distances are also revealed as the value for which the correlation between the real flows and the Flickr predicted flows is maximum.

4.5 Gravity model

The gravity model considers that the flow between two nodes (i, j) is directly proportional to some power of their populations and inversely proportional to an increasing function of the distance between them:

$$g_{ij} = K \cdot \frac{P_i^\alpha \cdot P_j^\gamma}{d(i, j)^\beta}.$$

We adjusted the gravity model using a linear regression in the logarithmic scale and following the approach of Balcan et al. [8]: we chose a power law of the distance $f(d) = d^{-\beta}$ for the air travel network and an exponential decay $f(d) = e^{-\beta \cdot d}$ for the commuting network, which provided the best results. The population information for the fitting was extracted from the public GeoNames database [46], and the population of a basin was computed as the sum of the populations of all cities inside that basin.

Acknowledgements

This work has been partially funded by the EC FET-Proactive Project MULTIPLEX (Grant No. 317532) to M.T. and C.C. The authors also acknowledge support from the “Lagrange Project” of the ISI Foundation funded by the Fondazione CRT and from the “S3 Project” funded by the Compagnia di San Paolo.

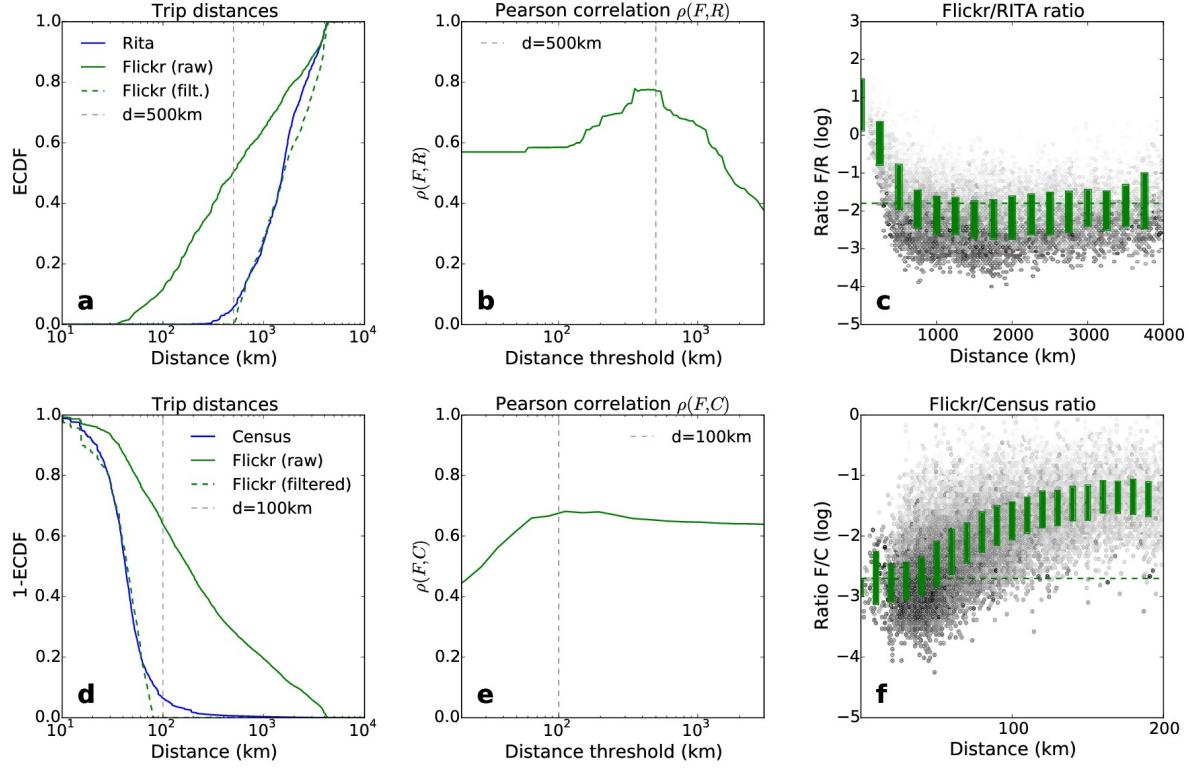


Figure 6: **Distance calibration for the U.S. air transportation network and the U.S. commuting network.** A correct fitting of the human mobility networks based on geolocalized traces requires aggregating the latter at the appropriate resolution level (basins or counties, respectively) and filtering those ones associated with the correspondent mobility type. Distance is an important calibrator for this task, as can be seen from: **(a,d)** Cumulative distribution of the trip distances in Flickr compared with the ground-truth (the RITA dataset in the case of air travel, and the census information for commuting); **(b,e)** Pearson correlation between Flickr flows and the ground-truth flows when considering Flickr trips above a distance threshold; **(c,f)** Flow ratio between Flickr and the ground-truth flows for a pair of (source, destination) nodes –basins or counties– as a function of the distance between them. (The green lines represent the linear regression coefficient between Flickr and the ground-truth flows. **(a,b,c)** Air transportation network; **(d,e,f)** Commuting network.

Supplementary Information

Distribution of Flickr users' trip distances

Figure 7 plots the trip distances of the Flickr users in the U.S., showing a power-law distribution with exponential cut-off.

Large flows between airports predicted by Flickr

Figure 8 shows that the Flickr trips are highly correlated with the air travel dataset when observing the pairs of airport basins with large flows of passengers. The green flows represent the connections predicted to have more than 10,000 passengers by Flickr, after a simple linear re-scaling with the real data.

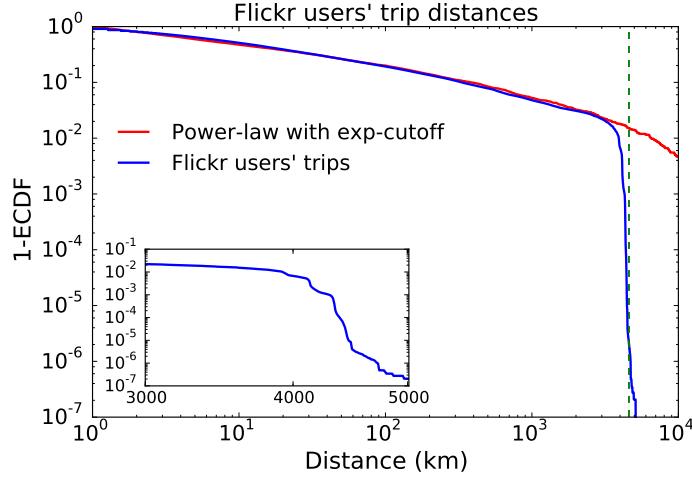


Figure 7: **Trip distances.** The trip distances of Flickr users follow a power-law distribution with exponential cutoff, bounded by the extension of the continental U.S., whose diameter (maximum distance between two points) is about 4,000 kilometers. This distribution is in accordance with previous research on individual human mobility [11].

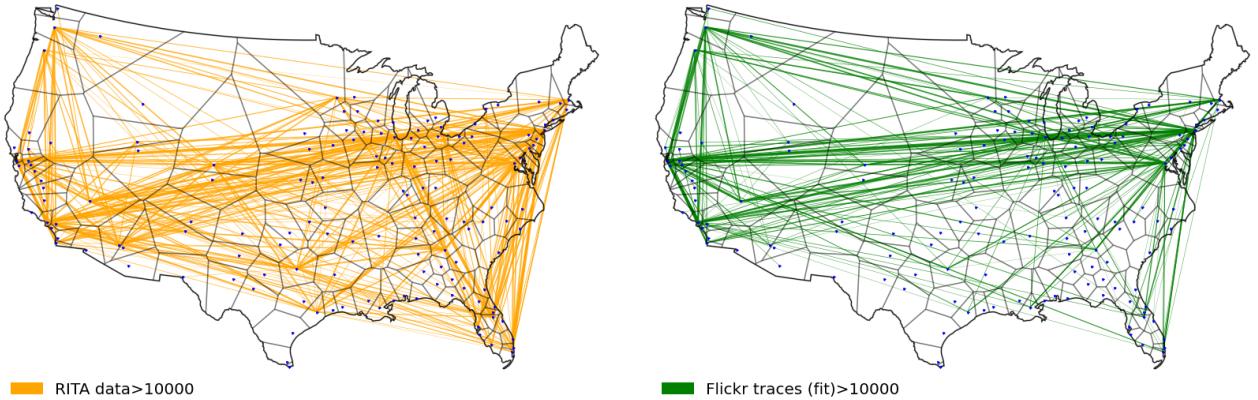


Figure 8: **Large air transportation flows captured by Flickr.** **a** Air flight connections with more than 10,000 according to the RITA dataset. **b** Flickr flows with more than 10,000 trips (after linear re-scaling).

Large commuting flows predicted by Flickr

Figure 9 shows that the Flickr trips are highly correlated with the commuting census data when observing the pairs of counties with large flows of commuters. The green flows represent the connections predicted to have more than 1,000 passengers by Flickr, after a linear re-scaling with the real data.

A Sorenson-Dice coefficient based test

In the main text we thoroughly explored the performance of the hybrid model under geographical and random sampling (*bootstrapping*) constraints. Here we will analyze the performance for different flows subsets organized by distance and destination population, using a test based on the Sorenson-Dice similarity coefficient. The Sorenson-Dice similarity between two sets A and B is defined as

$$s(A, B) = \frac{2 \cdot |A \cap B|}{|A| + |B|} .$$

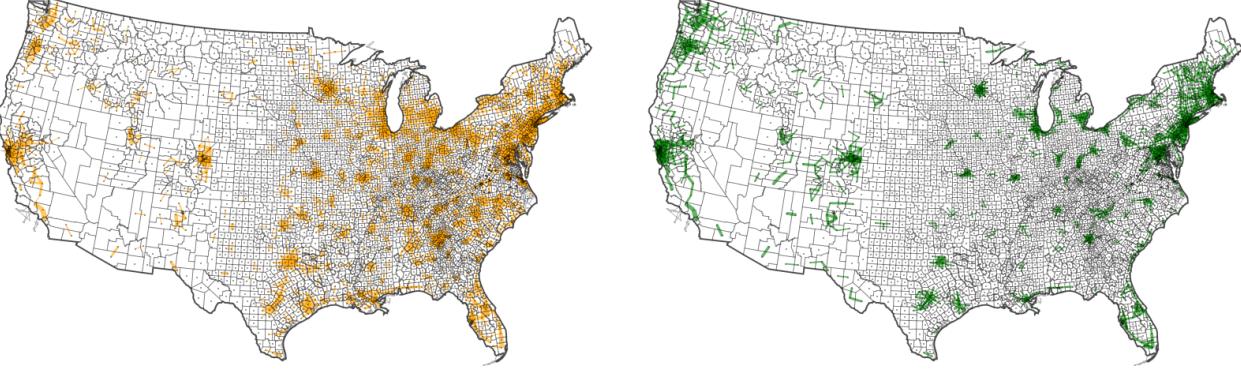


Figure 9: **Large commuting flows captured by Flickr.** **a** Commuting flows with more than 1,000 according to the U.S. Census dataset. **b** Flickr flows with more than 1,000 trips (after linear re-scaling).

Here we follow the modification introduced in [47] (named as “*Common part of commuters*” or *CPC*), and also used by [22, 24, 48], for comparing the predicted flows $\mathbf{H} = (h_{ij})$ against the real flows $\mathbf{Y} = (y_{ij})$:

$$CPC(\mathbf{H}, \mathbf{Y}) = \frac{2 \cdot \sum_{ij} \min(h_{ij}, y_{ij})}{\sum_{ij} h_{ij} + \sum_{ij} y_{ij}}.$$

Fig. 10 shows the goodness of fit in a grid of (distance, population) ranges, where each cell represents a subset of pairs of origin-destination basin filtered by distance and by destination population. In the left and central pictures, darker green colors represent a higher similarity between the predicted and real flows; the right pictures show in red those cells that were improved by the assimilation of Flickr traces, while those in blue are worse fitted when using Flickr data. We observe that in the air travel network (lower pictures) the improvement of the hybrid model is almost constant, but it is more evident for large population basins. In the commuting network, the gravity model outperforms the hybrid for small cities (up to $\approx 10,000$ inhabitants), but for larger cities the hybrid model gives better predictions. However, by improving the largest flows of people (which are usually associated to highly populated cities) the hybrid model can give a better estimation of the total amount of human flow at both resolution levels.

References

- [1] C. Roth, S.M. Kang, M. Batty, and M. Barthélemy. Structure of urban movements: polycentric activity and entangled hierarchical flows. *PLoS ONE*, 6(1):e15923, 2011.
- [2] M. Lenormand, A. Tugores, P. Colet, and J.J. Ramasco. Tweets on the road. *PLoS ONE*, 9(8):e105407, 08 2014.
- [3] L. Song, D. Kotz, R. Jain, and X. He. Evaluating next-cell predictors with extensive wi-fi mobility data. *IEEE Trans. Mob. Comput.*, 5:1633–1649, 2006.
- [4] H. Abou-zeid, H.S. Hassanein, and S. Valentin. Optimal predictive resource allocation: Exploiting mobility patterns and radio maps. In *Global Communications Conference (GLOBECOM), 2013 IEEE*, pages 4877–4882. IEEE, 2013.
- [5] E.G. Ravenstein. The laws of migration. *J. Stat. Soc. Lond.*, pages 167–235, 1885.
- [6] F. Simini, M.C. González, A. Maritan, and A-L. Barabási. A universal model for mobility and migration patterns. *Nature*, 484(7392):96–100, 2012.

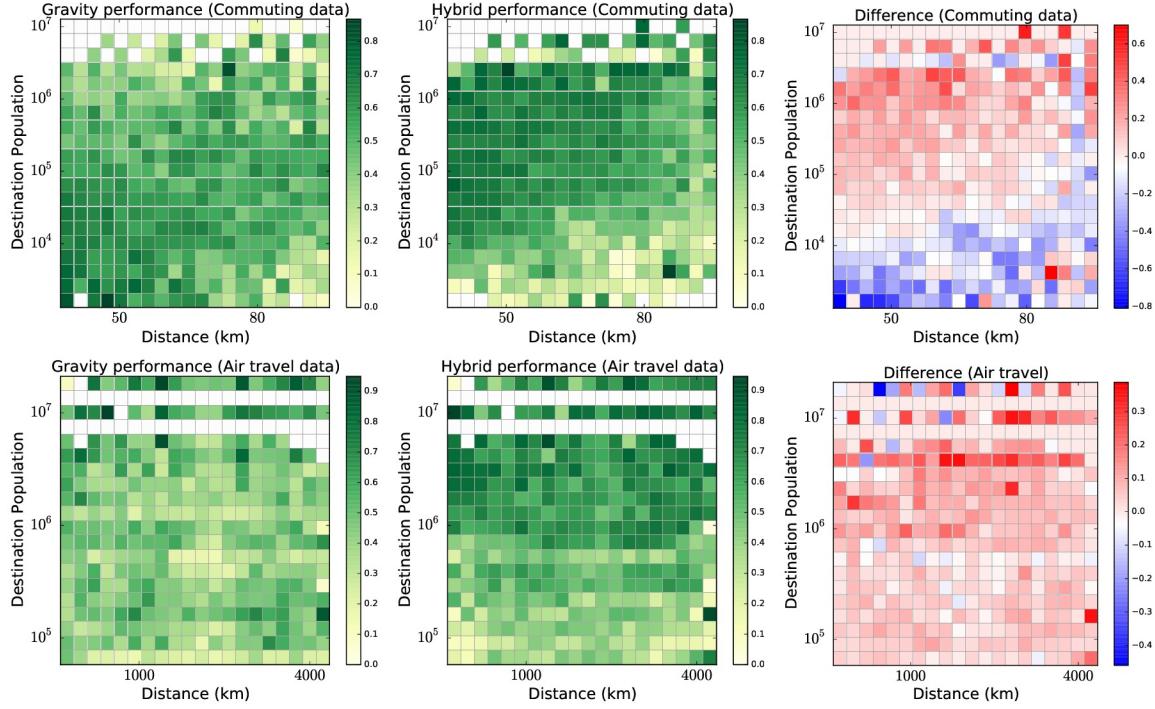


Figure 10: **Sorense-Dice coefficient grids for the gravity model and the hybrid model.** Error of flow estimates for the gravity model (left) and the hybrid model (center), and the difference in performance between them (right) for each distance-population cell, compared to the ground-truths of commuting (top grids) and air travel (bottom grids). The Sorense-Dice coefficients were computed following Eq. 6 in [22].

- [7] A. Wesolowski, B. Eagle, A.J. Tatem, D.L. Smith, A.M. Noor, R.W. Snow, and C.O. Buckee. Quantifying the impact of human mobility on malaria. *Science*, 338(6104):267–270, 2012.
- [8] D. Balcan, V. Colizza, B. Gonçalves, H. Hu, J.J. Ramasco, and A. Vespignani. Multiscale mobility networks and the spatial spreading of infectious diseases. *PNAS*, 106(51):21484–21489, 2009.
- [9] S. Merler and M. Ajelli. The role of population heterogeneity and human mobility in the spread of pandemic influenza. *Proc. Roy. Soc. Lond. B*, 277(1681):557–565, 2010.
- [10] R. Margutti and A. Parisi. Impact of human mobility on the periodicities and mechanisms underlying measles dynamics. *J. Roy. Soc. Int.*, 12(104), 2015.
- [11] M.C. Gonzalez, C.A. Hidalgo, and A-L. Barabasi. Understanding individual human mobility patterns. *Nature*, 453(7196):779–782, 2008.
- [12] I. Rhee, M. Shin, S. Hong, K. Lee, S.J. Kim, and S. Chong. On the levy-walk nature of human mobility. *IEEE/ACM Trans. Netw.*, 19(3):630–643, 2011.
- [13] A. Noulas, S. Scellato, R. Lambiotte, M. Pontil, and C. Mascolo. A tale of many cities: universal patterns in human urban mobility. *PLoS ONE*, 7(5):e37027, 2012.
- [14] B. Hawelka, I. Sitko, E. Beinat, S. Sobolevsky, P. Kazakopoulos, and C. Ratti. Geo-located twitter as proxy for global mobility patterns. *Cartography and Geographic Information Science*, 41(3):260–271, 2014.
- [15] G.K. Zipf. The p1 p2/d hypothesis: on the intercity movement of persons. *Am. Soc. Rev.*, 11(6):677–686, 1946.
- [16] W. Alonso. *A Theory of Movements: I, Introduction*. Working paper No. 266. Institute of Urban & Regional Development, University of California, Berkeley, 1976.

- [17] S.A. Stouffer. Intervening opportunities: A theory relating mobility and distance. *Am. Soc. Rev.*, 5(6):845–867, 1940.
- [18] W. Jung, F. Wang, and H.E. Stanley. Gravity model in the korean highway. *Europhys. Lett.*, 81(4):48005, 2008.
- [19] T. Grosche, F. Rothlauf, and A. Heinzl. Gravity models for airline passenger volume estimation. *J. of Air Trans. Manag.*, 13(4):175–183, 2007.
- [20] Y. Liu, Z. Sui, C. Kang, and Y. Gao. Uncovering patterns of inter-urban trip and spatial interaction from social media check-in data. *PLoS ONE*, 9(1):e86026, 2014.
- [21] G. Krings, F. Calabrese, C. Ratti, and V.D. Blondel. Urban gravity: a model for inter-city telecommunication flows. *J. Stat. Mech. Theory Exp.*, 2009(07):L07003, 2009.
- [22] A.P. Masucci, J. Serras, A. Johansson, and M. Batty. Gravity versus radiation models: On the importance of scale and heterogeneity in commuting flows. *Phys. Rev. E*, 88:022812, 2013.
- [23] X. Liang, J. Zhao, L. Dong, and K. Xu. Unraveling the origin of exponential law in intra-urban human mobility. *Sci. Rep.*, 3, 2013.
- [24] Y. Yang, C. Herrera, N. Eagle, and M.C. González. Limits of predictability in commuting flows in the absence of data for calibration. *Sci. Rep.*, 4, 2014.
- [25] N.M. Truscott, J. andAND Ferguson. Evaluating the adequacy of gravity models as a description of human mobility for epidemic modelling. *PLoS Comput. Biol.*, 8(10):e1002699, 2012.
- [26] F. Calabrese, G. Di Lorenzo, L. Liu, and C. Ratti. Estimating origin-destination flows using mobile phone location data. *IEEE Perv. Comp.*, 10(4):36–44, 2011.
- [27] V. Palchykov, M. Mitrović, H. Jo, J. Saramäki, and R.K. Pan. Inferring human mobility using communication patterns. *Sci. Rep.*, 4, 2014.
- [28] L. Alexander, S. Jiang, M. Murga, and M.C. González. Origin–destination trips by purpose and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging Technologies*, 2015.
- [29] M. Tizzoni, P. Bajardi, A. Decuyper, G. Kon Kam King, C.M. Schneider, V. Blondel, Z. Smoreda, M.C. González, and V. Colizza. On the use of human mobility proxies for modeling epidemics. *PLoS Comput Biol*, 10(7):e1003716, 2014.
- [30] L. Bengtsson, J. Gaudart, X. Lu, S. Moore, E. Wetter, K. Sallah, S. Rebaudet, and R. Piarroux. Using mobile phone data to predict the spatial spread of cholera. *Sci. Rep.*, 5, 2015.
- [31] S. Riley, K. Eames, V. Isham, D. Mollison, and P. Trapman. Five challenges for spatial epidemic models. *Epidemics*, 10:68–71, 2015.
- [32] B. Gonçalves and N. Perra. *Social Phenomena: From Data Analysis to Models*. Springer, 2015.
- [33] Flickr. <https://www.flickr.com/>.
- [34] Twitter. <https://twitter.com/>.
- [35] Foursquare. <https://es.foursquare.com/>.
- [36] A. Noulas, S. Scellato, C. Mascolo, and M. Pontil. An empirical study of geographic user activity patterns in foursquare. In *Proc. of the 5th Int'l AAAI Conference on Weblogs and Social Media*, pages 570–573, 2011.
- [37] M. Lenormand, M. Picornell, O.G. Cantú-Ros, A. Tugores, T. Louail, R. Herranz, M. Barthelemy, E. Frías-Martínez, and J.J. Ramasco. Cross-checking different sources of mobility information. *PLoS ONE*, 9(8):e105184, 08 2014.

- [38] P.A. Grabowicz, J.J. Ramasco, B. Gonçalves, and V.M. Eguíluz. Entangling mobility and interactions in social media. *PLoS ONE*, 9(3):e92196, 03 2014.
- [39] D. Barchiesi, T. Preis, S. Bishop, and H.S. Moat. Modelling human mobility patterns using photographic data shared online. *Royal Society Open Science*, 2(8), 2015.
- [40] B. Thomee, D.A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L. Li. The new data and new challenges in multimedia research. *arXiv preprint arXiv:1503.01817*, 2015.
- [41] Bureau of Transportation Statistics. Commuting (journey to work). Available at: http://www.rita.dot.gov/bts/data_and_statistics/index.html.
- [42] US Census Bureau. Airline origin and destination survey. Available at: <http://www.census.gov/hhes/commuting/>.
- [43] L. Breiman. Stacked regressions. *Machine learning*, 24(1):49–64, 1996.
- [44] T. Hastie, R. Tibshirani, and J. Friedman. The elements of statistical learning, 2009.
- [45] Yihui Ren, Mária Ercsey-Ravasz, Pu Wang, Marta C González, and Zoltán Toroczkai. Predicting commuter flows in spatial networks using a radiation model based on temporal ranges. *Nature communications*, 5, 2014.
- [46] Geonames. <http://geonames.org/>.
- [47] M. Lenormand, S. Huet, F. Gargiulo, and G. Deffuant. A universal model of commuting networks. *PLoS ONE*, 7(10):e45985, 2012.
- [48] A.. Wesolowski, W.P. O'Meara, N. Eagle, A.J. Tatem, and C.O. Buckee. Evaluating spatial interaction models for regional mobility in sub-saharan africa. *PLoS Comput Biol*, 11(7):e1004267, 07 2015.