

## 5. User-generated ubiquitous geoinformation as evidences of tourist dynamics

The approach of the previous case study (Chapter 4) focused on logs generated from individuals' implicit interactions with wireless infrastructures to perform travel surveys on a worldwide scale. This chapter considers other user-generated ubiquitous geoinformation to provide more empirical evidences of travellers' density and flows.

In a case study of Florence and Rome, Italy, we explore the value of explicitly disclosed geographically-referenced photos and implicitly-generated records of mobile phone network usage. We show that the analysis of these spatio-temporal data can supply high-level human behaviour information valuable to social scientists, urban planners and local authorities. Based on the techniques of the replay tool employed in our pervasive game (Chapter 2), we designed "Urban dynamics" a software that performs novel data collection and analysis techniques augmented with visualization and mapping tools. This software illustrates the potential of user-generated electronic trails to remotely reveal the presence and movement of a city's visitors, their spatio-temporal presence, inbound-outbound trajectories, internal flows, and semantic description. For instance, it helped comparing the significance of aggregated cellular network traffic data with georeferenced photos and reveal different presence of tourists in Rome.

We present this case study with the following papers:

Girardin, F., Fiore, F. D., Ratti, C., and Blat, J. (2008). Leveraging explicitly disclosed location information to understand tourist dynamics: A case study. *Journal of Location-Based Services* 2, 1, 41–54.

Girardin, F., Calabrese, F., Dal Fiore, F. , Ratti, C., and Blat, J. (2008). Digital footprinting: Uncovering tourists with user-generated content. *IEEE Pervasive Computing*, 7(4):36–43.

# Leveraging Explicitly Disclosed Location Information to Understand Tourist Dynamics: A Case Study

Fabien Girardin<sup>a</sup>, Filippo Dal Fiore<sup>b</sup>, Carlo Ratti<sup>b</sup> and Josep Blat<sup>a</sup>

<sup>a</sup>Department of Information and Communication Technologies, Universitat Pompeu Fabra, Barcelona, Spain

<sup>b</sup>Department of Urban Studies and Planning, Massachusetts Institute of Technology, Cambridge, USA

**Abstract.** In recent years, the large deployment of mobile devices has led to a massive increase in the volume of records of where people have been and when they were there. The analysis of these spatio-temporal data can supply high-level human behavior information valuable to urban planners, local authorities, and designer of location-based services. In this paper, we describe our approach to collect and analyze the history of physical presence of tourists from the digital footprints they publicly disclose on the web. Our work takes place in the Province of Florence in Italy, where the insights on the visitors' flows and on the nationalities of the tourists who do not sleep in town has been limited to information from survey-based hotel and museums frequentation. In fact, most local authorities in the world must face this dearth of data on tourist dynamics. In this case study, we used a corpus of geographically referenced photos taken in the province by 4280 photographers over a period of 2 years. Based on the disclosure of the location of the photos, we design geovisualizations to reveal the tourist concentration and spatio-temporal flows. Our initial results provide insights on the density of tourists, the points of interests they visit as well as the most common trajectories they follow.

**Keywords:** spatio-temporal data analysis; geovisualization; location-disclosure; location-based services

## 1. Introduction

In the past years, research in location sensing and tracking has been dominated by figuring out where persons and objects are in space (see Hazas et al., 2004 for a review). The wide adoption of the mobile and wireless technologies often referred to as Location-Based Services allows people a new perception of their surrounding physical and social spaces. In parallel, the records of where and when people have been (i.e. spatio-temporal data), produced by these services, have led to improve the understanding of different aspects of mobility and travel. The analysis of these data helps to recognize the modes of mobility (Sohn et al., 2006), define significant places (Ashbrook and Starner, 2002), cluster tourist routes in a city (Asakura and Iryob, 2007), infer travel purposes (Wolf, 2001) or predict a driver's destination as a trip progresses (Krumm and Horwitz, 2006). Ratti et al. (2006) benefit from people's experience of mobile devices to gain a more thorough understanding of urban environments, and this is the field of application undertaken in this paper.

First, let us remark that the large spatio-temporal logs can be seen as personal, as we explain later and their analysis can provide urban planners, local authorities and designers of location-based services with information on how a city gets used by different groups complementing current data sources, which are normally used as the basis of decision and policy making. Knowing who populates different parts of the city at different times can lead to the provision of customized services (or advertising), the rescheduling of monuments opening times or the reallocation of existing service infrastructures. From a quite different perspective, city users themselves could be aware of the current ways in which they populate the city, and adopt different strategies as a result.

Our approach takes advantage of the recent explosion in the use of capture devices (e.g. mobile phones, digital cameras) and collaborative web platforms to share their content (see Torniai et al., (2007) for a review). This people-generated information provides large amounts of digital data linked to the physical world. Recent research showed the potential of the geographically annotated material available on the Web. For instance, Ahern et al., (2007)

and Snavely et al. (2006) developed means of world exploration via photos and maps to foster “virtual tourism”. Using a similar dataset, Rattenbury et al. (2007) showed that the location and time metadata associated with photos and their tags enable the extraction of “place” and “event” semantics. In our study, we focus on the sense of presence extracted from a collection of photographs captured in trips. We consider that uploading, tagging and disclosing the location of a photo can be interpreted as an “I was here” statement indicating the physical presence in space and time. By extension, we propose that people-generated geographically referenced information provide new insights on how people travel and experience the city.

The validation of the relevance of this concept takes body in the Italian Province of Florence. The local authorities of that region aim at better understanding the tourist flows travelling across the cities boundaries. So far, they have been using classical survey-based hotel and museums frequentation data to know where tourists of different nationalities prefer to spend their time, hence money. However, they lack observations of the mobility, nationality and quantity of the “day trippers”, that is the tourists who visit Florence but are “invisible” in the data, as they do not sleep in town. In consequence, we retrieved 81017 photos taken in the region by 4280 photographers over a period of 2 years from the popular photo-sharing web platform Flickr<sup>14</sup>. Based on the time and the disclosed location of the photos, we extracted records of the people presence and movements; performed statistical analysis and designed geovisualizations. This exploratory visual analysis was used as a mean of preliminary investigation. In fact, the results go far beyond the initial expectations of collecting clues on “day trippers” activities, as we shall see, providing new insights, and even a novel paradigm of urban geography.

In the remainder of this paper, we discuss first current work on mobility data collection and constraints to perform travel surveys. These shortcomings suggest that the research community should investigate and evaluate new data and perspectives. Then, we propose a novel approach that takes advantage of spatio-temporal

---

<sup>14</sup> <http://www.flickr.com>



data generated by tourists when publicly sharing their photos on the world-wide web. Next, we describe the types of data that can be collected and their meaning. Afterwards, we present the preliminary results of our analysis supported by geospatial visualizations. They highlight the ability to quantify the concentration of tourists and their movements over time through the major areas of interests. Finally, we conclude with a description of the implications from this case study and discuss the meaning for future work.

## **2. Related work and their limitations**

The recent emergence of location technologies and techniques favoured the development of new approaches to capture and analyze people's mobility (see Wolf, 2004 for a survey). The aim has been to replace traditional travel diaries, paper-and-pencil interview, computer-assisted telephone interviews, and computer-assisted-self-interview, by automatically collecting mobility data. For instance, Wolf et al. (2001) proved the feasibility of using Global Positioning System (GPS) data loggers to improve the quality or completely replace traditional questionnaires. However, this type of mobility survey faces the problem inherent to longitudinal studies such as recruiting a pool of respondents or preventing any fatigue effects. Besides the privacy concerns of continuously and precisely tracking people, Schoenfelder et al. (2002) and Stopher et al. (2003) identify the potential technical drawbacks of a GPS-based approach. They list transmission problems, warm-up times before getting a valid position, and the cost of post-processing the GPS data as issues that impair the quality of the survey. Indeed, Wolf et al. mention that the equipment packages deployed for their pilot study had a lot more problems than anticipated: the units and cabling used were not optimized for durability, resulting in the loss participants due to problems with equipment performance and user operations.

Other mobility studies relied on the mobile phones use of the Global System for Mobile communications (GSM) network to generate mobility data. In a first type of approach, the mobile devices calculate and report their position to a centralized service. The TeleTravel System (TTS) project (Wermuth, 2003) combined a mobile device GSM tracking technology and an electronic travel diary to determine the travel behaviour of the respondents; while

Asakura and Iryob (2005) determined it by interpreting changes in the set of nearby towers and signal strengths of the phones as indicative of position and motion of their owners; Froehlich et al. (2006) augmented the GSM mobility data with context-triggered in-situ survey in which panelists rate the place they are in to study travel routines. A second approach employs the measurement by mobile network operators of traffic intensity and migration of each cell in a GSM network to capture the movement patterns of mobile phone users. Through this technique Ratti et al. (2006) measure the evolution through space and time of the activities in a city to support urban planning. Based on comparable network data, Ahas et al. (2007) retrospectively link the digital track of visitors with visited events and locations retrospectively. While this approach scales because it does not rely on a costly software or hardware deployment, it fails in capturing individual traces. A third approach relies on the deployment of Bluetooth enabled systems. Recent studies have been able to establish the flow of people at strategic locations (O'Neill et al., 2006) as well as to recognize daily user activities and to identify socially significant locations (Eagle and Pentland, 2006).

However, there are several key issues when using technologies to collect travel behaviours. A major concern consists is the privacy and ethical issues related to collecting data without individual's consent (Gutman and Stern, 2007). Indeed, there is an increase in the risk of identifying people or organizations when their data spatial precision improves. Other issues include the length of the survey (e.g. to prevent fatigue effects); the ability to collect individual traces; and the scalability challenge of deploying the tracking system, for instance, when dealing with a variety of cellular network standards and providers. Each of the above-mentioned approaches has to face these aspects. Table 1 shows how GPS, GSM and Bluetooth tracking perform with respect to these important aspects of a travel survey.

Table 1. Mobility data capture techniques with respect to issues in the context of travel surveys. Low=little issue, High=big issue

Mobility data capture techniques / Issues	Scalability	Longevity	Individual traces	Privacy
GPS	High	High	Low	Medium
GSM	High	Medium	Medium	Medium
(device-based)				
GSM	High	Low	High	Medium
(aggregated network-based)				
Bluetooth	High	Low	Low	Medium

Some scholars work in the combination of the approaches to produce more complete spatial data (origin, route, destination) when tracking individual's routes (Kracht, 2004). However, while promising, the approach has not been able to address all the issues. In the remainder of this paper, we argue that a new type of explicitly disclosed spatio-temporal data coming from public web platforms can overcome these constraints and provide an additional insights to understand the dynamic of people in an urban space.

### 3. The raw data and the processing

We study photos uploaded on the popular photo-sharing platform Flickr. People use this service to share and organize photos; an option allows to add them geographical references. Each time a photo is virtually linked to a physical location, the Flickr system assigns a longitude and latitude and retrieves the time of capture from the Exchangeable Image File Format (EXIF) metadata embedded in the photo. The location provided by the user generally indicates where the photo was taken; but sometimes it denotes the photographed object. When the user provides the location of the photo, the zoom level of the map (from 8 for region/city level to the most precise 16 for street level) is recorded as an accuracy attribute, completing the spatio-temporal information.

Figure 1 describes the process of recording and collecting the data. First, tourists take photos during their trips and journeys. Later, they manually associate a position to the photos through a Flickr map

interface or other external map-based services. A minority of tech savvy users has their photos implicitly annotated with a position extracted from data collected by GPS devices embedded or external to the camera. From the publicly available photos in a given area, we retrieve the coordinates, timestamp, accuracy level, and an obfuscated identifier of the owner via the Flickr API. Based on these data, as an example, a chronologically ordered set of photos gives rise to traces that reveal the movements of different individuals in space. In the figure we show as well some analysis tools we discuss later.

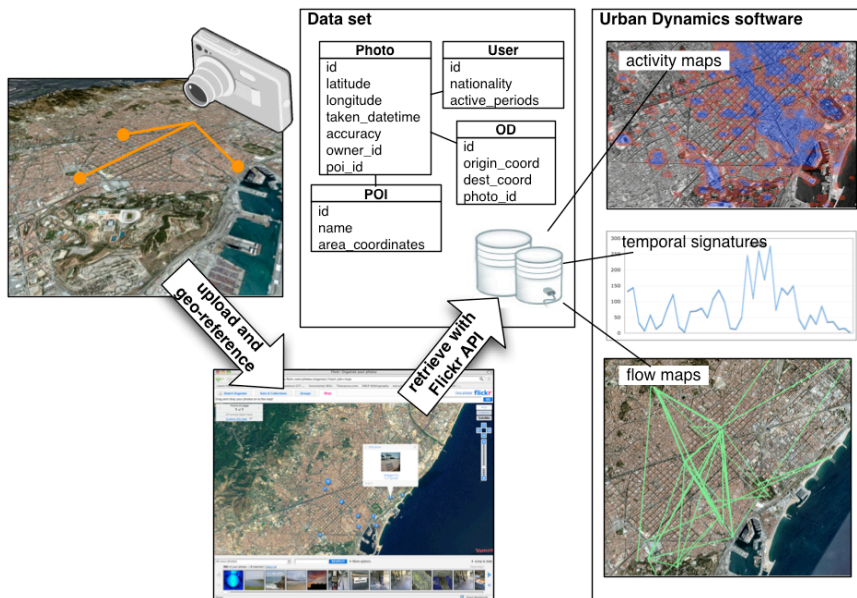


Figure 1. Data flow, from data recording, retrieving, storing to the visualizations. Photographic imagery copyright 2007 DigitalGlobe.

Our data set contains photos from the Province of Florence in Italy over a 2-year period, from November 2005 to November 2007. The timestamps extracted from the camera-generated EXIF metadata do not necessarily match the correct time. Indeed, the capturing devices might not be set with the correct time zone information, and the accuracy of the time embedded in the photo when it is taken is uncertain. However, the temporal analysis we describe later indicates that this seems to have little impact on our analyses. For example, as we consider only photos stamped in 2005, 2006 or

2007, we can assume that the user had set the camera date as its default value is rarely one of these years.

A first step of pre-processing was to separate visitors from residents of the region. To achieve that, we used the presence in the area over time as the discriminating factor. We divided the time in 30-day periods and computed the number of periods each photographer was active in the area. If a photographer took all his/her photos within 30 days, the algorithm considers him/her as a visitor, while if there is an interval greater than 30 days between two photos taken, he/she is categorized as resident. The aim of this strict threshold was to capture the real one-time tourists. Out of a population of 4280 photographers, we identified 3505 one-time visitors.

In addition, we were interested to know more about the nationalities of our photographers. Hence, we took advantage of the Flickr social function that invites (without forcing) users to provide “offline bits” on their city and country of residence. We found out that 65% (2782 out of 4280) of the users actually disclosed this information. While it is hard to predict how much of this data is truthful, the country of residence could be retrieved by automatically parsing the data in most cases. We had to assign manually a country in 11% (306/2782) of the cases, because of spelling errors or idiosyncratic names (e.g. “Big Apple” for New York).

Before going any further, let us remark that collecting and analyzing any kind of mobility data raises serious privacy issues; people are concerned about revealing the history of their whereabouts to un-trusted third party applications. We are concerned as well and our approach addresses this at two different levels. First, the users explicitly disclose the position of their photos on maps and control who gets access to their location data for our data set. Moreover, while obtaining this public information, we applied an obfuscation algorithm to lose the relationship with the identity used in Flickr. Thus, we only analyze anonymized records of digital traces publicly disclosed by individuals.

Our experience of collecting Flickr geo-referenced data provided some insights on the data at hand, which are an indicator of different aspects of their value. As of March 2007, Flickr contained more than 20 million photos linked to a physical location. Cities

such as London or New York contain more than 250'000 photos and 9'000 single photographers each and growing pace of around 400 photos a day for London, and 150 a day for Barcelona. This quantitative richness might push towards an even bigger increase of publicly accessible people-generated location and time of the Flickr data set. At another level, a very specific type of people use Flickr to geographically reference their images. They are generally well travelled and technologically savvy. Therefore, we are dealing with a very specific type of tourists.

In addition, unlike the automatic capture of traces, the manual location disclosure embedded in the act of geotagging of photo provides additional qualities. Positioning photo on a map is not simply adding information on its location, it is also an act of communication containing what people estimate as being relevant for themselves and others. In that sense, a specific richness of this dataset arises in the intentional weight people put in disclosing their photos. We show that they have a tendency to select the highlights of their discovery of the city and discarding the downtimes.

#### **4. Tools used, and initial results**

This section of the paper presents the geovisualizations that are produced by our tools and discusses their meanings. In the research process, we were driven by an overarching research question: what is the spatio-temporal behaviour of tourists from different nationalities, inside the Province of Florence? Flickr seemed an adequate context to extract data from, as it contains the spatial and temporal elements, as well as the nationality. However, we have to take into account the following caveats: we have already discussed the lack of temporal accuracy; second we cannot assume perfect mirroring of the Flickr spatio-temporal patterns and the actual ones, as we do not know where users have actually been between a given photo and the following one. In other words, we can only deal with approximations. With these limitations in mind, we provide an analysis of the data based on map visualizations to reveal the tourist concentration and flow within the Central Italy region and the city of Florence. Each set of maps aims at giving indications on the applicability of the collected data set to better understand the tourist dynamics within a certain area.

Our initial results bring a new perspective on two aspects related to spatial and temporal density and movements of visitors of the Province of Florence and its capital the city of Florence: (a) Characterizing the areas of the city/region where the tourists are concentrated, and (b) Revealing spatio-temporal signatures: activity by day of the week and month of the year, and days of the year. In addition, two types geovisualizations seek to understand better the flow of tourists into, out of and within the Province of Florence: (c) trajectories into and out of the area studied, (d) Patterns of flow within the Province of Florence.

#### **4.1. The “Urban Dynamics” software**

The extraction of structure, meaning and insight from large, multifaceted, spatio-temporal datasets is a challenging task that requires skills not possessed by many engaged in geovisualization (Dykes et al., 2007). Our approach takes advantage of open and freely-available resources and combines them using de facto standards often based on the extensible Markup Language (XML). We use Google Earth<sup>15</sup> for interactive visual synthesis of encodings generated using a combination of MySQL<sup>16</sup> for data storage and querying to select and aggregate, and a software developed in Java we named “Urban Dynamics” to access, process, transform, aggregate, cluster, sample, filter the raw data stored in the database and to generate outputs. The Keyhole Markup Language<sup>17</sup> (KML) is used to describe visual encodings and define interactions. Google Earth is used as a means of interactively visually analysing and inspecting data, through its spatial and temporal navigation tools, its access to wider contextual data, and its ability to create animations with image overlays. Figure 2 describes the data processing and multiple outputs generated by Urban Dynamics. The data of the photos are stored in a spatial matrix to perform population density analysis and in linked arrays of positions to reveal the traces left by people in their visits. In practical terms, a trace consists in a chronologically ordered set of geographically referenced photos taken by one person over one day. Is this the appropriate place to define the trace? Should it be defined earlier? These data structure

---

<sup>15</sup> <http://earth.google.com/>

<sup>16</sup> <http://www.mysql.com/>

<sup>17</sup> <http://earth.google.com/kml/>

can be divided in temporal periods to ground the base of a spatio-temporal analysis and visualization. The density values stored in the spatial matrix are employed to cluster the main areas of interest in the studied area. After their storage in the database, these areas are mapped in Google Earth with a KML file for identification and labelling. The software also generates from the matrix data heat maps of population density that can be included as an image overlay in KML files visualized in Google Earth. We describe next the spatio-temporal geovisualizations of flows and density generated with Urban Dynamics and the analyses of their results.

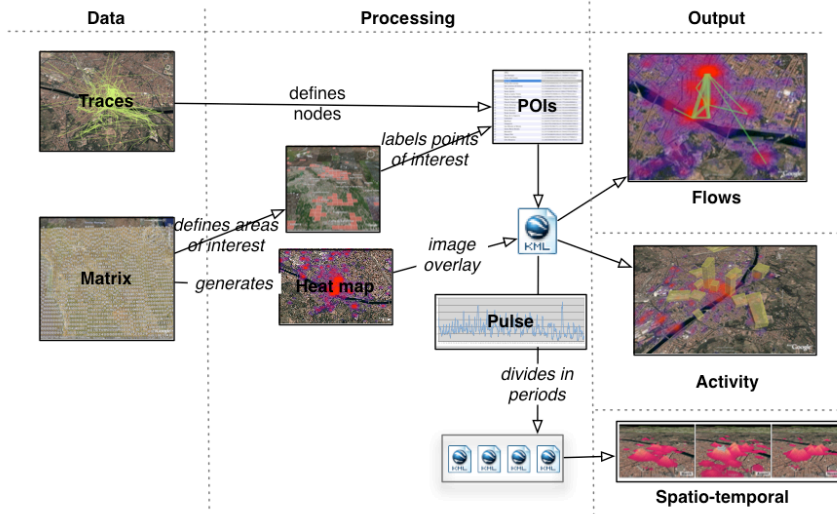


Figure 2. The Urban Dynamics software multiple data, processing and outputs. Photographic imagery copyright 2007 DigitalGlobe.

#### 4.2. Where are the tourists concentrated?

To understand where tourists concentrate, we extracted from our dataset the latitude and longitude of the photos taken by visitors that were geotagged with a granularity more precise than the city resolution level. These data were put in a matrix covering the studied area with each cell of the matrix containing the number of photos and visitors ever present in that zone. To reveal visually the tourist concentrations, we proceeded in two steps. First, we generated from the matrix a Keyhole Markup Language (KML)



document to feed it to the Geocommons<sup>18</sup>, a geospatial data visualization platform. It produced an interactive interface with heat maps of the tourist concentrations on top of the Google Maps system (Figure 3). The results show a zoomable map of the overall tourist activity covering the northern part of Central Italy, the city of Florence and around the famous Basilica di Santa Maria del Fiore.



Figure 3. Interactive zoomable map of the overall tourist activity from the northern part of central Italy to fine-grained information on the Basilica di Santa Maria dal Fiore in Florence. Generated with GeoCommons. Photographic imagery copyright 2007 Cnes/Spot, NASA, DigitalGlobe and TerraMetrics.

These visualizations show that the main activity in the Province of Florence actually takes place in the city of Florence. There are other major points of interest in the region, such as the cities of Pisa and Siena, as well as the Mediterranean coast between La Spezia and Livorno, and the Island of Elba. Within the city, the tourists tend to concentrate around the Ponte Vecchio and the Basilica di Santa Maria del Fiore. While these geovisualizations give a quick and simple overview of the presence of tourists, they do not offer a scale to understand the quantitative meaning of the colours. Therefore, as a second step, we developed an algorithm with the Processing<sup>19</sup> Java Library to generate heat maps for their inclusion as an additional layer in Google Earth. Next, we clustered the presence values in the matrix to define the major areas of tourist concentration. A list of 50 areas each representing a point of interest was computed and manually labelled according to their position on the map. The match of the areas with the matrix values allowed rank the major points of interest of the area studied according to the number of photos taken, the number of tourists and residents, number of photos taken by person and the likeliness to encounter

<sup>18</sup> <http://www.geocommons.com/>

<sup>19</sup> <http://www.processing.org>

residents or tourists (Figure 4). The Basilica di Santa Maria del Fiore is the most visited monument of the city followed by the famous Ponte Vecchio, the Plaza de la Signoria and Palazzo Vecchio. However, the Arco Lorena is the monument that tourists like to photograph the most (18.06 photos per tourist) compared to the Basilica di Santa Maria del Fiore that triggers 7.72 photos per tourist. Fiesole, a town in the suburbs of the city is one of the main point of interest fairly out of the tourist frenzy with a proportion of visitors of 73% compared to the 89% of the Ponte Vecchio.



Figure 4. Presence of tourists in the main areas of interest in downtown Florence. The visualization shows areas of activity (i.e. photos taken) in red and the presence of tourists with yellow polygons. The altitude of the polygons represents the number of individuals present. Photographic imagery copyright 2007 DigitalGlobe.

### 4.3. Temporal activity and spatio-temporal signatures

In our first attempt, we produced maps of tourist concentrations without taking into account the period when data were recorded. In this section, we generated similar maps to compare the activity in time, namely, seasons and weekends versus weekdays, etc. As differential data might be more informative, we took the “tourist pulse” of the area of study, by developing first its spatio-temporal signature. We charted the temporal signature of the number of visitors active in Figure 5, which shows the weekly, monthly and yearly pulse in 2007 with the number. Unsurprisingly, the summer

months are the busiest of the year and tourist are more active on Saturdays and Sundays. A yearly vision of the presence reveals a steady pulse with peaks on the weekends.

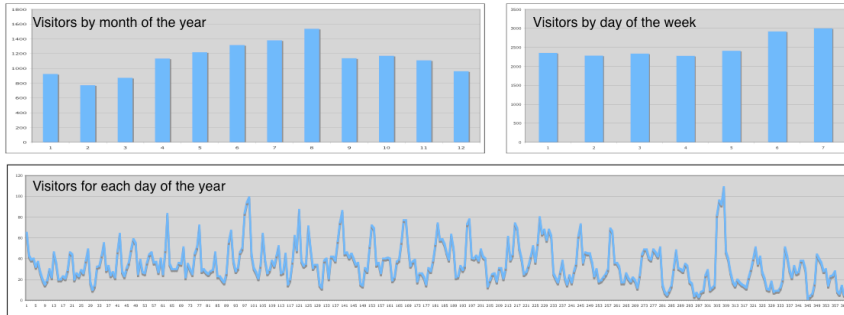


Figure 5. Pulse as the average number of visitors per each month of the year, day of the week and day of the year for 2007.

Based on the 2006 data, we produced animations of monthly activity to understand where and when was taking place the tourist activity (Figure 6). Our geovisualization adds a layer on top of a map of the Province, showing 250 x 250 cells of 20m x 12m whose colors represented the number of photos taken and visitors present. Similarly to Real-time Rome animations based on Erlang values (Reades et al., 2007), we used a 5 persons threshold to reveal the main areas of activity. That is, any cells of the matrix containing less than 5 persons were discarded. In addition we smoothed the visualization of areas of activity with an interpolation algorithm. This type of animation helped revealing the areas and monuments popular at certain seasons. The Boboli gardens, for instance, do not draw many visitors in winter while other points of interests do not lose much attractiveness.

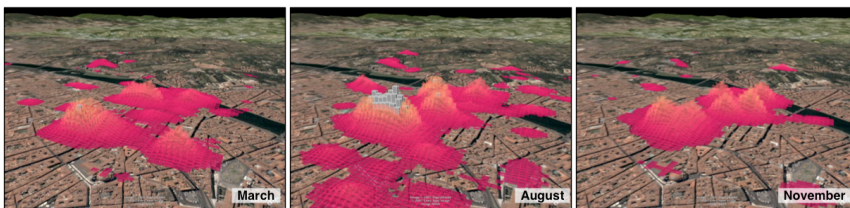


Figure 6. Screenshots of the spatio-temporal animation of the presence of tourists in downtown Florence in 2007. Photographic imagery copyright 2007 DigitalGlobe.

#### 4.4. Inbound and outbound trajectories

As we discussed earlier, the classic survey-based hotel and museums frequentation data used by tourist authorities to know where tourists of different nationalities come from and go do not capture “day trippers”, the tourists who visit Florence but do not sleep in town. Our data brought a new perspective on the issue to understand where the visitors had been prior to entering the area and where they were heading after leaving it. We retrieved the images the 3505 visitors took prior to entering, and after leaving the area. As a result, we collected more than 1.8 million geographically referenced photos taken worldwide. Then, we recovered the inbound and outbound trajectories through the rules described in Table 2. A specific problem was to find the proper threshold in time to infer a precise origin of an inbound movement and the destination of an outbound one. The test of different thresholds of 24, 48 and 72 hours revealed that more than 83% of the inbound and outbound movements took place within 24 hours. The movements exceeding this interval were discarded, as they might not have reflected the proper origin and destination of a journey via Florence.

Table 2. Description of the algorithm to detect inbound and outbound traces

Type of trajectories	Algorithm
Inbound	Detect the position and time of the last photo before entering the area, and the position and time of the first photo after entering the area.
Outbound	Detect the position and time of the first photo after leaving the area, and the position and time of the last photo before leaving the area.

With the similar process and tools as the ones described in section 4.1. we plotted two types of geospatial visualizations to provide information on where tourists are prior to entering the city of Florence and after leaving it. The first result is an interactive zoomable heat maps that reveals density of inbound and outbound destination of a day trip in Florence. With the manual annotation of the heat maps we were able to define the 28 main areas and cities of origin and destination. Each area had a bounding box and assigned with the number of incoming and outgoing visitors computed. With this quantitative understanding of the trajectories, our software generated flow maps to offer another perspective on the movement

by rendering the trajectories of tourist. They are generated from KML files of origin and destination traces loaded in the Google Earth software. For flows leading out of Florence (Figure 7), each green trace links the position of the last photo taken in Florence with the position of the first photo taken outside the city within 24 hours. Figure 7b displays through an analogous strategy, the aggregates outbound movement from the city of Florence. The comparison of inbound and outbound trajectories shows striking similarities. The most connected cities, Rome, Pisa and Venice generated a relatively equal amount of incoming and outgoing movements. Florence acts as connecting city between the north and the more southern part of Italy with Rome as the main destination.

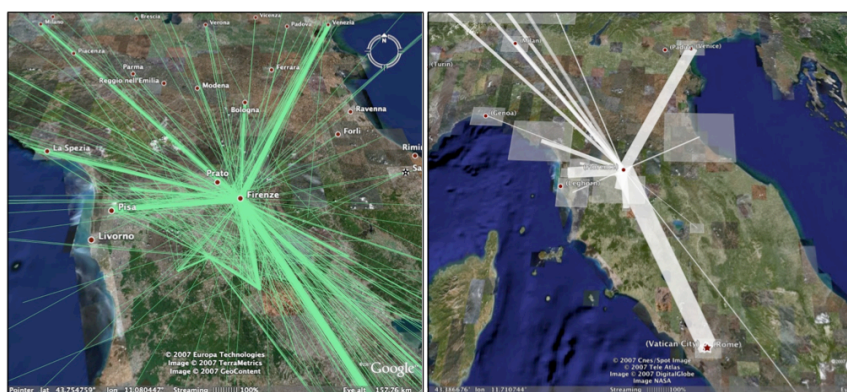


Figure 7. Outbound trajectories from the city of Florence. The map displays the individual traces in green and their aggregation in white. Photographic imagery copyright 2007 Cnes/Spot, NASA, DigitalGlobe and TerraMetrics.

#### 4.5. Patterns of flow

In addition to revealing the origins and destinations of the tourists, the analysis of the dataset allows to gather insights on the flows within the boundaries of the northern part of Central Italy and the city of Florence. We could trace the Flickr users from the digital footprints they leave along their path. Aggregating these personal footprints and formatting them in KML reveals the travel behaviours of specific types of visitors. For instance, Figure 8 shows that tourists from the USA follow a specific graph constituted by the nodes of Florence, Siena, Pisa, Genova and Perugia. On the other hand, Italians are more adventurous in their



exploration of the area. For instance it shows the popularity of the Island of Elba, an attraction not visited by American tourists.

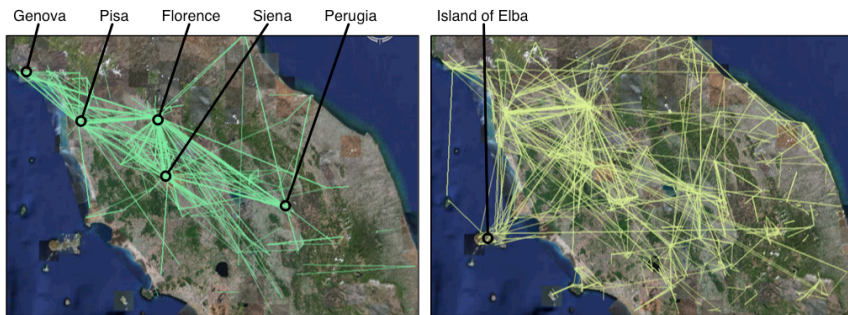


Figure 8: Movements of Americans in green compared to Italians in yellow in northern part of central Italy. The comparison of the visualization shows the distinct was to apprehend the space. Photographic imagery copyright 2007 Cnes/Spot, NASA, DigitalGlobe and TerraMetrics.

In practice, our basic flow maps can easily become too cluttered and prevent the detection of flow patterns. As a solution, we use the main areas of interests of the Province of Florence and the city. Our visualization software matches these areas of interests with the traces left by visitors to gain a quantitative measurement of the movements between the main attractions of the area studies. Figure 9 shows the result of this process for downtown Florence. Our software generates the traces left by each visitor and aggregates them in correspondence to the main points of interests they connect. A KML file is generated with the heat map as an image overlay object and the flows between each point of attraction plotted as green lines with a weight proportional to the density of the flow. This combined vision of the density and flow of tourists shows that they consume a very limited space of the city. The tour to the unique relatively distant area of interest, the viewpoint of the Piazzale Michelangelo, seems to take place from and the most popular monument of the city, the Basilica di Santa Maria del Fiore.

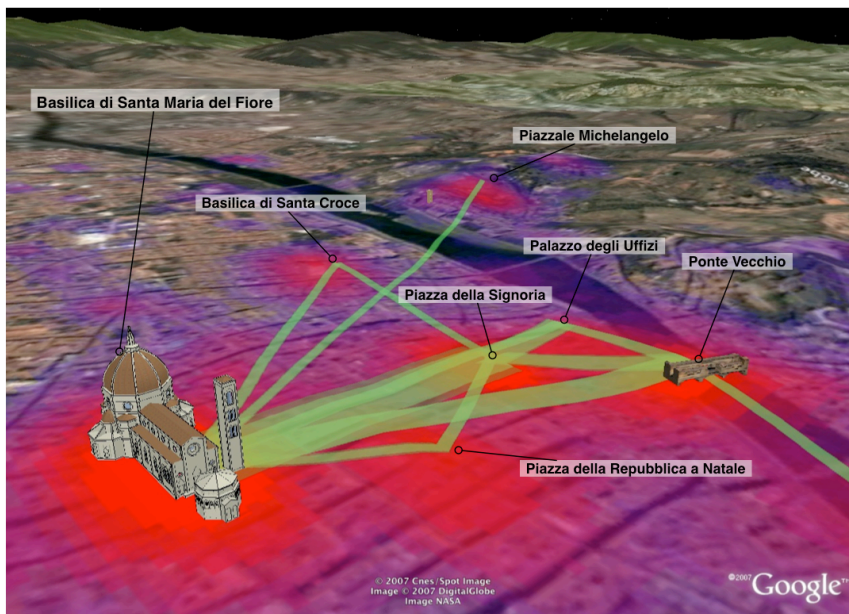


Figure 9: Downtown Florence, main areas of tourist activity (in red) and flows of visitors between them (in green). The width of the green line represents the amount of visitors who moved from one attraction to another. Photographic imagery copyright 2007 Cnes/Spot, NASA, and DigitalGlobe.

## 5. Conclusion and future work

The explosion in the use of captures devices (e.g. mobile phone, digital cameras) and the emergence of content sharing platforms is leading to the emergence of a wealth of publicly available user-generated geospatial data. Our case study specifically featured the value of Flickr and its geographically reference photos in the Italian Province of Florence with the goal of performing urban and mobility analysis. Our approach takes advantage of individual traces made publicly available to analyze behaviours in a urban space. Their exploratory analysis enables us to quantify by the amount of photos taken and the presence of individuals the attractiveness over time of the major points of interests of an urban space. In addition, the mapping of these data shows the flows entering, leaving and visiting a city like Florence. Not only the photos allow us to quantify the movements between attractions, they also help detecting the distinct patterns of mobility among groups of tourists of different nationalities. This type of insight is

limited in most travel surveys and urban sensing infrastructure by privacy-sensitive or aggregated information.

The careful processing of the user-generated data, and the use of appropriate geovisualizations are necessary ingredients of the analyses performed. We have indicated where some improvements of the processing and the visualizations might come from. For instance, we need to find out more on the representativeness of the user of Flickr and their veracity in disclosing the photos they share. User-generated data require special attention to self-selection in the data sample. Indeed, those people who proactively upload digital information on open platforms such as Flickr are very likely to represent a tech-savvy sub-section of any given population. By collecting further data on the socio-demographics of these users, it might be possible to understand better how much the behaviour findings from research based on a self-selected sample can be generalized to a given population.

This case study could open a new perspective in urban cartography and lead to a new urban paradigm based on the analysis of publicly available people generated geo-referenced data. The explicit act of sharing this information with the ability to control and cloak the data greatly reduces privacy concerns inherent to current travel surveys techniques, as we have shown in this paper. In addition, unlike the classic automatic capture of traces, the act of geotagging a photo can be interpreted as an act of communication because people only give information they estimate as being relevant for themselves and others. In that sense, a richness of people-generated geographically referenced data relies on the specific effort people put in to disclosing the location information. Our geovisualizations show that they have a tendency to select the highlights of the their discovery of the city and discard the downsides.

In the future, similar types of spatially anchored user-generated content might surface to become relevant for travel, mobility and urban studies. Sources might range from implicit data from the usage of radio-frequency networks (GPS, GSM, WiFi, Bluetooth) to more explicit information contained in geospatial web applications (e.g. geo-referencing in Flickr), and the emerging social applications based on the disclosure of presence and location information. Researchers, urban planners, local authorities, and



others might consider the aggregate analysis of these different channels.

As part of further work, we consider several extensions to the promising results presented in this paper. First we aim at correlating our dataset with other spatio-temporal data such as GSM network usage, hotel and attractions surveys. Second, we believe that the recording of such data could be used to inform the design and deployment of location-based services to enhance the tourist experience. So far location-aware applications have tended to concentrate on using a mobile carrier's immediate geographic location. We believe that using a position history to tailor results from requests for information should enhance them (Mountain and Raper, 2001). This work shows that there might be potential in taking advantage of the digital traces people constantly leave behind them to, for instance, reveal the temporal character of a space, its attractiveness among a certain group of people, and its level of connectivity with other spaces.

## Acknowledgements

The results reported in this paper are part of a broader research undertaken in the Wireless City, Florence and Beyond project. We would like to thank the authorities of the Province of Florence for their support. Also, we are indebted to many people at the Massachusetts Institute of Technology and the Universitat Pompeu Fabra for providing extremely stimulating research environments and for their generous feedback. In particular, thanks to Assaf Biderman, Francisco Calabrese, Nicolas Nova, Bernd Resch and Toni Navarrete for letting us pick their brains. Of course, any shortcomings are our sole responsibility.

## References

- Ahas, R., Aasa, A., Silm, S., Tiru, M. 2007. Mobile positioning data in tourism studies and monitoring: case study in Tartu, Estonia. In: Sigala, M., Mich, L., Murphy, J. (Eds.), *Springer Computer Science: Information and Communication Technologies in Tourism*, ISBN: 978-3-211-69564-7, pp. 119-128.
- Ahern, S., Naaman, M., Nair, R., and Yang, J. H. 2007. World explorer: visualizing aggregate data from unstructured text in geo-referenced

- collections. In *JCDL '07: Proceedings of the 2007 conference on Digital libraries*, pages 1–10, New York, NY, USA. ACM Press.
- Asakura, Y. and Iryob, T. 2007. Analysis of tourist behaviour based on the tracking data collected using a mobile communication instrument. *Transportation Research Part A: Policy and Practice*, 41(7):684–690.
- Ashbrook, D. and Starner, T. 2002. Learning significant locations and predicting user movement with gps. In *ISWC '02: Proceedings of the 6th IEEE International Symposium on Wearable Computers*, page 101, Washington, DC, USA. IEEE Computer Society.
- Dykes, J., Slingsby, A., and Clarke, K. 2007. Interactive visual exploration of a large spatio-temporal dataset: Reflections on a geovisualization mashup. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1176–1183.
- Eagle, N. and Pentland, A. S. 2006. Reality mining: sensing complex social systems. *Personal Ubiquitous Comput.*, 10(4):255–268.
- Froehlich, J., Chen, M. Y., Smith, I. E., and Potter, F. (2006). Voting with your feet: An investigative study of the relationship between place visit behavior and preference. In *Ubicomp*, pages 333–350.
- Gutman, M. and Stern, P., editors 2007. Putting People on the Map: Protecting Confidentiality with Linked Social-Spatial Data. *National Academies Press*.
- Hazas, M., Scott, J., and Krumm, J. 2004. Location-aware computing comes of age. *IEEE Computer*, 37(2):95–97.
- Kracht, M. 2004. Tracking and interviewing individuals with gps and gsm technology on mobile electronic devices. In *Seventh International Conference on Travel Survey Methods*.
- Krumm, J. and Horvitz, E. 2006. Predestination: Inferring destinations from partial trajectories. In *Ubicomp*, pages 243–260.
- Mountain, D. and Raper, J. 2001. Modelling human spatio-temporal behaviour: a challenge for location-based services. In *6th Internat. Conference on GeoComputation*. University of Queensland, Brisbane, Australia.
- O'Neill, E., Kostakos, V., Kindberg, T., gen. Schieck, A. F., Penn, A., Fraser, D. S., and Jones, T. 2006. Instrumenting the city: Developing methods for observing and understanding the digital cityscape. In *Ubicomp*, pages 315–332.

- Ratti, C., Pulselli, R. M., Williams, S., and Frenchman, D. 2006. Mobile landscapes: Using location data from cell-phones for urban analysis. *Environment and Planning B: Planning and Design*, 33(5):727 – 748.
- Reades, J., Calabrese, F., Sevtsuk, A., and Ratti, C. 2007. Cellular census: Explorations in urban data collection. *IEEE Pervasive Computing*, 6(3):30–38.
- Rattenbury, T., Good, N., and Naaman, M. 2007. Towards automatic extraction of event and place semantics from flickr tags. In *SIGIR '07: Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 103–110, New York, NY, USA. ACM Press.
- Schoenfelder, S., Axhausen, K., Antille, N., and Bierlaire, M. 2002. Exploring the potentials of automatically collected gps data for travel behaviour analysis - a swedish data source. In Möltgen, J. and Wytzisk, A., editors, *GI-Technologien für Verkehr und Logistik*, number 13 in IfGIprints, pages 155–179. Institut für Geoinformatik, Universität Münster.
- Snavely, N., Seitz, S. M., and Szeliski, R. 2006. Photo tourism: exploring photo collections in 3d. *ACM Trans. Graph.*, 25(3):835–846.
- Sohn, T., Varshavsky, A., LaMarca, A., Chen, M. Y., Choudhury, T., Smith, I., Consolvo, S., Hightower, J., Griswold, W. G., and de Lara, E. 2006. Mobility detection using everyday gsm traces. In *UbiComp*, pages 212–224.
- Stopher, P., Bullock, P., and Jiang, Q. 2003. Visualising trips and travel characteristics from gps data. *Road & Transport Research*, 12(2):3–14.
- Torniai, C., Battle, S., and Cayzer, S. 2007. Sharing, discovering and browsing geotagged pictures on the web. Technical report, HP Labs.
- Wermuth, M., Sommer, C., and Kreitz, M. 2003. *Impact of New Technologies in Travel Surveys*, chapter 27, pages 455–482. Pergamon, Oxford.
- Wolf, J. 2004. Applications of new technologies in travel surveys. In *7th International Conference on Travel Survey Methods, Costa Rica*.
- Wolf, J., Guensler, R., and Bachman, W. 2001. Elimination of the travel diary: experiment to derive trip purpose from global positioning system travel data. *Transportation Research Record*, (1768):125–134.

# Digital footprinting: uncovering the presence and movements of tourists from user-generated content

Fabien Girardin<sup>a,b</sup>, Francesco Calabrese<sup>b</sup>, Filippo Dal Fiore<sup>b</sup>,  
Carlo Ratti<sup>b</sup> and Josep Blat<sup>a</sup>

<sup>a</sup>Department of Information and Communication Technologies, Universitat  
Pompeu Fabra, Barcelona, Spain

<sup>b</sup>Department of Urban Studies and Planning, Massachusetts Institute of  
Technology, Cambridge, USA

**Abstract.** In recent years, the large deployment of mobile devices has led to a massive increase in the volume of records of where people have been and when they were there. The analysis of these spatio-temporal data can supply high-level human behavior information valuable to social scientists, urban planners and local authorities. This paper explores this hypothesis by reporting on new information revealed by this pervasive user-generated content. We present novel techniques, methods and tools we have been developing to explore the significance of these new types of data. In a case study of Rome, Italy, we showcase the ability to uncover the presence and movements of tourists from geo-referenced photos they explicitly make public, as well as from network data implicitly generated by users of mobile phones.

## 1. Introduction

Nowadays every click of every move of every user who interacts with any software may be gathered in a database and submitted to a second-degree data-mining operation. Along with the growing of ubiquity of mobile technologies, the logs produced have helped to create and define new methods of observing, recording, and analyzing the city and its human dynamics [1]. In effect, these personal devices create a vast, geographically-aware sensor web [2], which uses the accumulation of tracks to reveal both individual and social behaviors with unprecedented detail [3]. The low cost and high availability of these ‘digital footprints’ will challenge the social sciences, which have never before had access to the volumes of data used in the natural sciences [4], but the benefits to fields where an in-depth understanding of large group behavior could be equally great.

Accordingly, this paper illustrates the potential of user-generated electronic trails to reveal – remotely – the presence and movement of visitors in a city. We anticipate that the validation of these trails with respect to existing surveys may lead to an improved understanding of several aspects of urban mobility and travel. We therefore present several novel data collection techniques, analytical methods, and visualization tools that we have been developing to uncover these dynamics in the city. While the nature of digital footprints renders the information derived both more credible and reliable, we must further consider how to validate this pervasively user-generated content.

In previous work, we showed that explicitly disclosed spatio-temporal data from open platforms provide novel insights on the dynamics of visitors in an urban space [5]. Understanding population dynamics by type, by neighborhood, or by region would enable the provision of customized services (or advertising) as well as the accurate timing of urban service provision, such as the scheduling of monument opening times based on the presence of tourists, based on daily, weekly, or monthly demand. In general, more synchronous management of service infrastructures clearly could play a particularly important in tourism management.

## **2. Working with Digital Footprints**

Visitors to a city have many ways of leaving voluntary or involuntary electronic trails: prior to their visits tourists generate server log entries when they consult digital maps [6] or travel web sites; during their visit they leave traces on wireless networks whenever they use their mobile phones; and after their visit they may add online reviews and photos. Broadly speaking then, there are two types of footprint: active and passive. Passive tracks are left through interaction with infrastructure, such as the mobile phone network, that produces entries in locational logs, while active prints come from the users themselves when they expose locational data in photos, messages, and sensor measurements.

In this paper, we consider two types of digital traces from the city of Rome, Italy: geo-referenced photos made publicly available on through the photo-sharing web site Flickr<sup>20</sup>, and aggregate records of wireless network events generated by mobile phone users making calls and sending text messages on the Telecom Italia Mobile (TIM) system.

## **2.1. Explicit footprints: georeferenced photos**

People using the Flickr service to share and organize photos also have the option to add geographical attributes. Each time a photo is anchored to a physical location, Flickr assigns longitude and latitude values together with an accuracy attribute derived from the zoom level of the map in use to position the photos. So photos positioned on a map when the user is zoomed in at the street level receive a higher accuracy estimate than ones positioned when the user had pulled back in the online map view. The system also adds metadata embedded by the camera into the image using the Exchangeable Image File Format (EXIF) information, completing the spatio-temporal information.

Flickr also provides a public Application Programming Interface (API) that enables anyone to query their public data store for photos. We elected to analyze three years of data – from November 2004 to November 2007 – for the city of Rome since it is a very popular and highly photographed tourist destination. For the 3-year period of interest, we were able to extract 144,501 geo-referenced photos that had been uploaded by 6,019 different users. For each of these publicly available photos, we retrieved the geographical coordinates, timestamp, accuracy level, and an obfuscated identifier of the owner.

Since we were particularly interested in the behavior of tourists in Rome, we separated the photographers into two groups based on their presence in the city over time. Discriminating between locals and visitors required dividing the study period into 30-day blocks: if a photographer took all his or her photos within a period of thirty days, the algorithm considered them to be a visitor, but if they uploaded photographs at intervals of more than 30 days then they

---

<sup>20</sup> <http://www.flickr.com>

would be categorized as a resident. From our population of 6,019 photographers, 4,719 were classed as one-time visitors.

To find out more about the nature of our photographers, we took advantage of a social function in Flickr that invites users to voluntarily provide additional information about themselves such as their city and country of residence. In some cases, because of spelling errors or user idiosyncrasies – such as using “The Big Apple” to mean New York City – we were forced to manually process the city or country information. However, after cleaning we found that 59% of the users had disclosed meaningful origin information. They break into several main populations: 991 Italians, 1171 other Europeans, 807 North Americans, 104 South Americans, 71 Asians, and 70 from Australia and New-Zealand.

## **2.2. Implicit footprints: wireless network usage**

Previous research has shown the wide diffusion of mobile phones and the widespread coverage of mobile phone wireless networks in urban areas make these technologies very interesting as means to identify and track both groups and individuals [7, 8]. Our collaboration with TIM took advantage of new a system called LocHNESSs (Localizing & Handling Network Event Systems), which is a software platform that localizes and stores user-generated events as they occur on the mobile network. Calls in progress, SMS transmissions, and call handovers are all captured through external probes that localize and collate incoming messages before transmitting the results to LocHNESSs. The messages are then aggregated to produce raster-format maps of the distribution of users. A detailed introduction to the platform can be found in Calabrese et al [9].

TIM installed the LocHNESSs platform and related probes on a set of Base Station Controllers (BSC) located in the northeast quadrant of the city, covering an area of approximately 100Km<sup>2</sup>. The system permits users to be reliably localized to within an area of 250 by 250 m<sup>2</sup> and then assigned to the corresponding grid reference. LocHNESSs divided the users into two groups – Italians and foreigners – based on the country code information embedded in their International Mobile Subscriber Identity (IMSI) number. Over a period of three months, timed to coincide with the Venice

Biennale from September to November 2006, the system calculated these attributes every five minutes and transmitted the results to servers at the Massachusetts Institute of Technology (MIT).

### **2.3. Processing and Visualization**

Dykes et al. [10] suggest that the large volumes of data coming from these types of sources can only be interpreted through geovisualization, which is to say that after collection, mapping and visualization is critical to interpreting and explaining user behaviors. We elected to use Google Earth to support visual synthesis and our preliminary investigation of digital traces. Accordingly, data collected by the LoCHNESs platform and from the Flickr service was stored on a MySQL server, enabling us to flexibly query and aggregate the data further as required. Using software developed in-house, we then exported the aggregate results in a format compatible with Google Earth for interactive visual exploration. Precise digital satellite imagery from Telespazio was added as image overlay. The application of these techniques and tools to process digital footprints allows us to uncover the presence of crowds, the patterns of movement over time, and to perform a comparison of user behaviors to generate new hypotheses.

## **3. Analysis of Digital Footprints**

### **3.1. Spatial presence**

To map the spatial distribution of users, data is stored in a matrix covering the entire study area. Each cell in the matrix includes data about the number of photos taken, the number of photographers present, and the number of phone calls made by foreigners over a given period of time. The geovisualization shown in Figure reveals the main areas of tourist activity in part of central Rome over the 3-month period of September to November 2006.

The presence of photographers is pictured in left and areas of heavy mobile phone usage by foreigners are depicted on the right. The union between visiting photographers and foreign mobile phone customers quickly uncovers the area's major visitor-attractions such as the Coliseum and the main train station next to Piazza della Repubblica. Intriguingly, it appears that the Coliseum attracts the photographers in their sightseeing activity while foreign mobile



phone users, typically on the move, tend to be active around the train station.

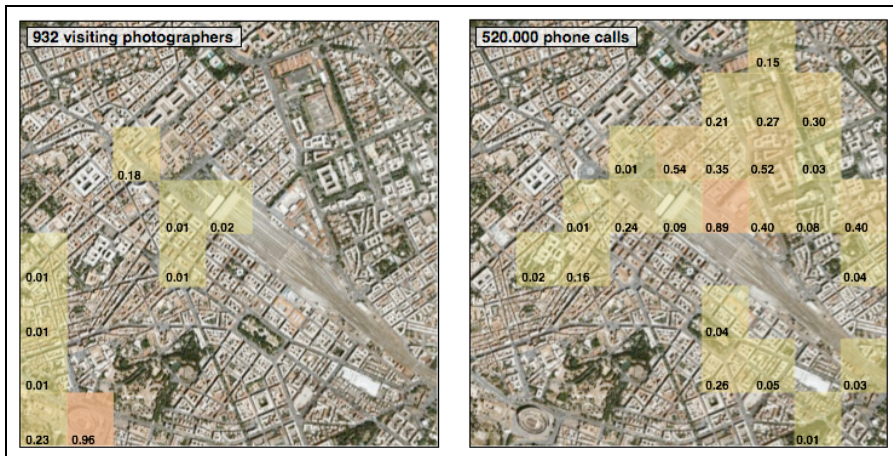


Figure 1. Geovisualizations of the presence of 932 tourist photographers (left) and 520 000 phone calls from foreign mobile phones (right) in the Coliseum-Piazza della Repubblica area from September to November 2006. Both type of data cover the train station area in the proximity of the Piazza della Repubblica. The values in each cell are normalized.

### 3.2. Temporal presence

Turning to the temporal patterns obtained from the digital traces, we compared the number of photographers and the volume of phone activity for each day of the week over the 3-month study period. Figure shows the difference between the average weekly distribution of phone calls made by visitors and the presence of visiting photographers in the areas around the Coliseum and Piazza della Repubblica. The histograms show the normalized variation between the average number of calls and the average number of photographs for each day of the week, and the average amount for the whole week.

The resulting temporal signatures for the Coliseum area show related trends for both data sets, with higher activity over the weekend than on weekdays. However, the Piazza della Repubblica area reveals a markedly different pattern: photographers, though fewer in number than at the Coliseum, also tend to be active on the weekend, whereas the foreign mobile phone users are much more active during the weekdays.

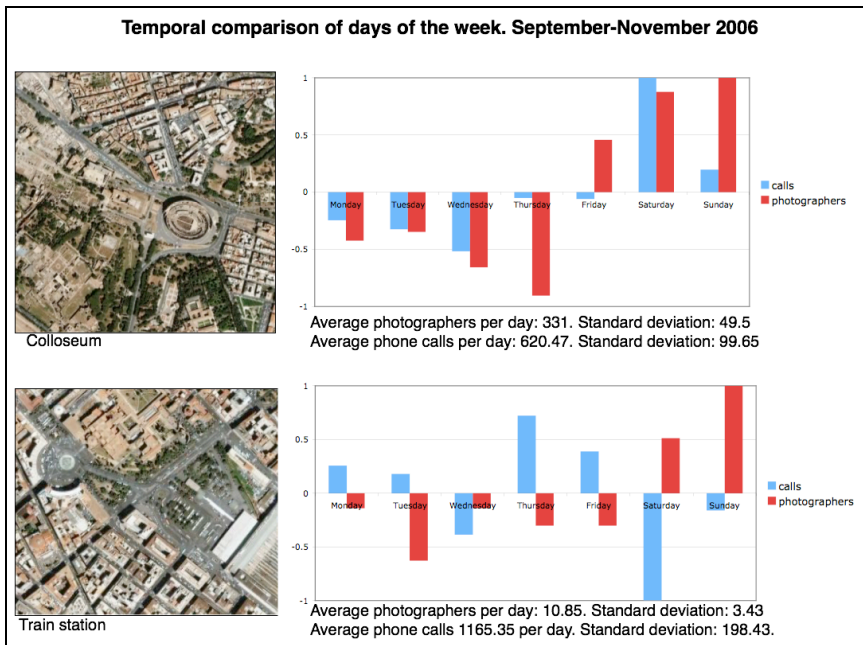


Figure 2. Comparison of the temporal signature of foreigners phone activity and number of tourist photographers. It reveals similar patterns of low below average activity during the week days and a rise of presence over the weekend at the Coliseum. In opposition, the temporal signature of the train station shows a higher presence of foreigners calling from their mobile phones during the week, while photographers indicate a reverse pattern and increased presence over the weekend.

These temporal signatures provide further evidences to the different types of presence that occur at the tourist points of interest. It can be further hypothesized that the Coliseum attracts sightseeing (i.e. photographers) activities over the weekend and the neighborhood of the train station provides facilities for visitors on the move (e.g. people on business trips) during the weekdays.

### 3.3. Desire lines from digital traces

The study of digital footprints also enables us to uncover the digital ‘desire lines’ embodied in people’s paths through the city. Based on the time stamp and location of photos, our software organizes the images chronologically in order to reconstruct the movement of the photographers. More precisely, we start by revealing the most

active areas obtained by spatial clustering of the data<sup>21</sup>. Next, we aggregate these individual paths to generate desire lines that capture the sequential preferences of visitors. The location of each user activity (i.e. photo) is checked to see if it is contained in a cluster, and in the case of a match, the point is added to the trace generated by the owner of the photo. This process produces multiple directed graphs that support better quantitative analysis, enabling us to obtain the number of sites visited by season, the most visited and photographed points of interests, as well as where do photographers start and end their journeys.

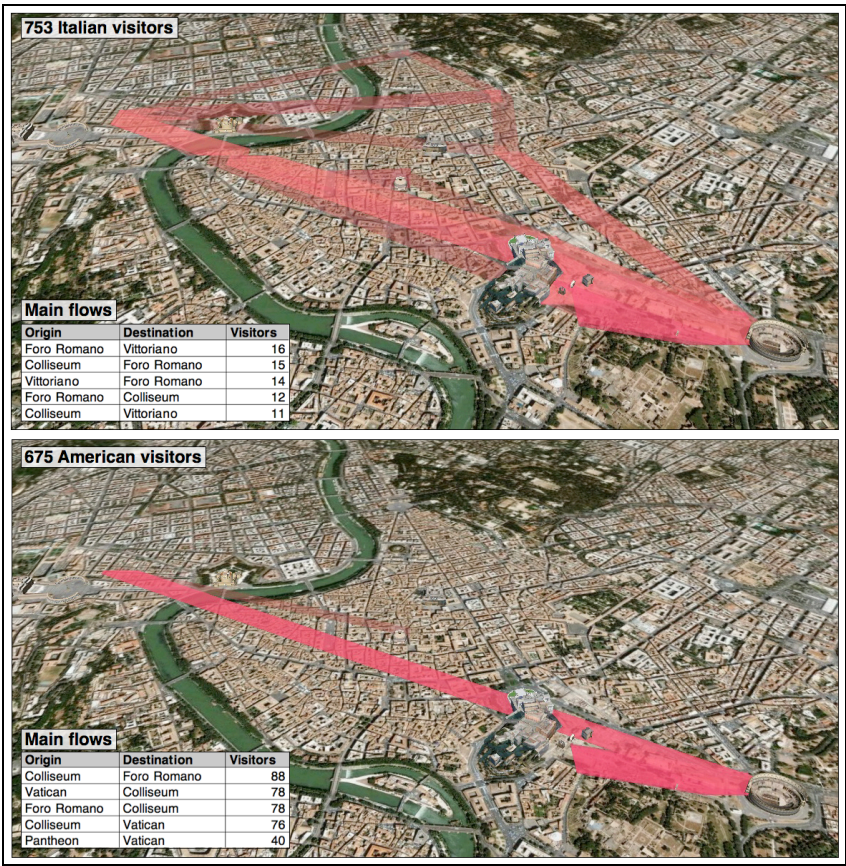


Figure 3. Geovisualisation of the main paths taken by photographers between points of interests of the city. Significantly, the 753 visiting Italian photographers (top) are active across many areas of the city, while the 675 American visitors

<sup>21</sup> Direct aggregation of the traces was not providing results easy to visualise and analyze.

stay on a narrow desire line between Vatican, the Forum, and the Coliseum. Note that different scales apply for each geovisualization.

Formatting this data according to the open Keyhole Markup Language<sup>22</sup> standard enables us to import the data into Google Earth to explore the traveling behaviors of specific types of visitors. The resulting visualization in Figure suggests the main points of interest in the city as a whole. Building asymmetric matrices of the number of photographers who moved from one point of interest ‘x’ to another point of interest ‘y’ reveals the predominant sequence of site visits. In addition, queries can be based on the nationality of the users, the number of days of activity in the city, the number of photos taken, and areas visited during a trip.

### **3.4. Semantic description**

Previous work has demonstrated that spatially- and temporally-annotated material available on the Web can be used to extract “place” and “event” related semantic information [11]. In a similar vein, we analyzed the tags associated with the user-originating photos to reveal clues of people’s perception of their environment and the semantics of their perspective of urban space. For instance, the word “ruins” is one of the most-used tags to describe photos in Rome. Mapping the distribution of this tag for 2,866 photos uncovers the most ancient and ‘decayed’ part of the city: the Coliseum and the Forum (Figure).

---

<sup>22</sup> <http://www.opengeospatial.org/standards/kml/>





Figure 4. Geovisualization of the areas defined by the position of the 2886 photos with the tag "ruins" uploaded by 260 photographers. It reveals the Coliseum and Forum areas known for their multitude ancient ruins.

#### 4. The significance of user-generated data

These aggregate spatio-temporal records seem to lead to an improved understanding of different aspects of mobility and travel. Although the results are still fairly coarse, we have clearly shown the potential for geographically-referenced digital footprints to reveal patterns of mobility and preference amongst different visitor groups. However, in the context of our study, traditional methods such hotel occupancy and museum surveys to observe and quantify the presence and movements of visitors would help us to better define the usefulness of pervasive user-generated content. Fortunately, the Rome Tourism office supplied us with monthly ticket receipts for the Coliseum in 2006.

Figure compares sales figures with the mobile usage and photographic activity. Ticket receipts show that there are slightly more Coliseum visitors in October than September, with a major drop in attendance during November. This pattern matches the activity of foreign-registered mobile phones in the area, but does not coincide with the activity of photographers. We hypothesize that these discrepancies arise from the fact that the datasets are capturing the activity of different sets of visitors: one set of data is generated

by visitors paying to see an attraction, another by mobile service customers wealthy enough to pay international roaming charges, and the third by a technology-savvy community of people who are familiar with digital photography, and mapping and social networking software. Due to the large difference in the nature of the activity producing the data that we compare, it might be that correlating with user-generated content does not reinforce existing tourism and travel knowledge, but reveals new dimensions of user behavior.

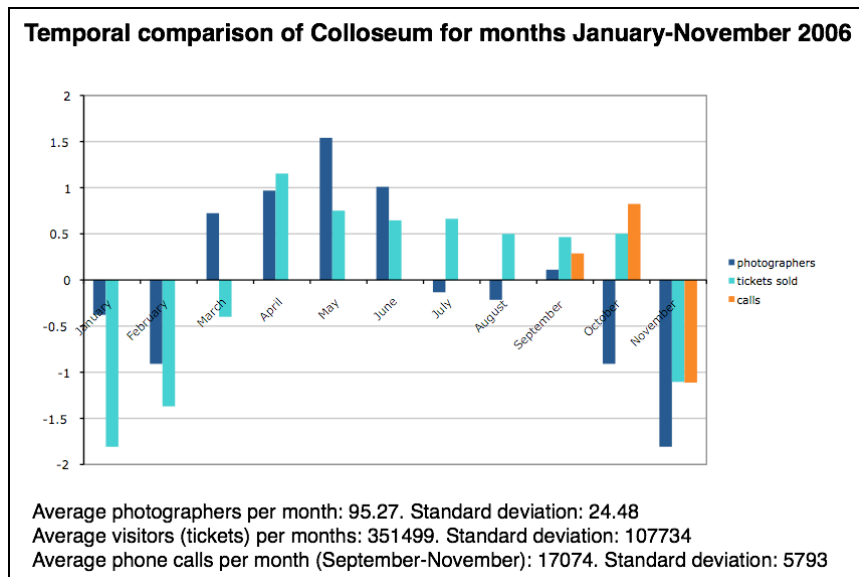


Figure 5. Comparison of the presence of visitors to the Coliseum area between January and November of 2006 using the number of tickets sold, number of calls made, and number of photographers active in the zone. The values represent the variations from the monthly average, scaled by the standard deviation.

## 5. Challenges in user-generated datasets

The analysis of user-generated content does not come without challenges, particularly when dealing with data quality and privacy issues.

### 5.1. Data quality issues

Our data processing techniques have tried account for the fluctuating quality of user-generated data, which can substantially impede our ability to generate accurate information. For instance,

the timestamps extracted from the camera-generated EXIF metadata do not necessarily match the real time at which a photo was taken – the user must not only have set the clock on their camera to the local time, but they must take the time to set the clock in the first place. User-generated data points can also be apparently idiosyncratic and, for instance, indicate not the point where photo was taken but the location of the photographed object.

The inclusion of mobile phone data introduces challenging scale issues since the resolution of the phone and photo datasets vary substantially. Correlation with ticket sales from the Coliseum also fails to account for the fact that users can easily photograph the arena or make a call from the vicinity of the monument without bothering to pay the entry fee. Challenges also arise from the fact that only phone activity handled by a subset of BSCs in Rome were monitored, leading to the risk of “border effects”, where calls near the border of a monitored area may be handled by other BSCs and thus not counted by the LoCHNESs platform. This last consideration applies in particular to areas to the south west of the Coliseum (see also Figure).

## **5.2. Privacy issues**

The use of photograph and mobile data can also be expected to raise privacy and ethical concerns related to collecting data without the individual’s consent. However, our approach addresses these concerns on two levels: first, our photography dataset includes only information that users explicitly disclose on an open platform; and second, all data is aggregated in a way that removes all traces of the individual. On the Flickr service users have direct control over who can access their locational data, but we supplemented this by applying an obfuscation algorithm to erase the relationship with the web identity of the individual and their digital trails. Thus, we could only analyze anonymous records of information already publicly disclosed by individuals.

Collection and analysis of aggregate network usage data fully-complied with the 2002 Directive of the European Parliament and Council on privacy. Data was only reported to us in aggregate, and so we received no data about an individual’s identity or trajectory. In effect, we could only count the total number of people – either

Italian or foreign – that used a mobile phone at a given point in the city and at a given moment. Individual users could not in any way be identified based on the data that we collected and analyzed, and consequently we avoided the significant privacy issues that have been raised by other methodologies [8].

## **6. Discussion**

The explosion in the use of capture and communication devices (e.g. mobile phones, digital cameras) and the introduction of content-sharing platforms has led to the emergence of a wealth of georeferenced-data. This user-generated content provides new opportunities for urban studies and the social sciences to understand the behavior of visitors and residents in an urban context. From a methodological perspective, the data we have analyzed in this paper has a clear advantage over more traditional location data obtained rather through controlled studies where subjects carried sensors and were thus aware of being tracked. Although we could not determine the sample used, our mobile phone data covers the usage habits of more than one million people and thus represents a step-change in the scale of localizable data collection efforts.

These collection methods also contain several important potential advantages over other pervasive tracking systems. Solutions that require people to carry a separate GPS-enabled device not only remind users that their movements are being followed – which might encourage them to pursue ‘high-brow’ activities during their visit – but also the tracking solutions that persons must carry generate fatigue effects and do not always function well in urban areas because of signal multipath and urban canyon obstructions. The alternative of a distributed, but fixed web of sensors entails onerous maintenance and data transmission costs. These issues strongly suggest that the research community should investigate and evaluate the use of these new data types as well as considering approaches that do not rely on the deployment of ad-hoc and costly infrastructures.

This paper therefore seeks to illustrate the value of explicitly-disclosed geographically-referenced photos and implicitly-generated records of mobile phone network usage. We used user-originated digital footprints to uncover some new aspects of the



presence and movement of tourists during their visit to Rome. And we introduced several novel tools and techniques for this analysis, although these results demonstrate that further development is required in order to validate our observations and to lead to new insights into factors such as the temporal usage-signature of a space, its attractiveness to different groups of people, and the degree of similarity to usage of other spaces.

The explicit character of photo geo-tagging and manual disclosure to the world also provides additional dimensions of interest: positioning a photo on a map is not simply adding information about its location; it is an act of communication which embodies locations, times, and experiences that individuals consider to be relevant to themselves and others. There is a very real richness to the ‘intentional weight’ that people attach to disclosing their photos, and the results clearly show that Flickr users have a tendency to point out the highlights of their visit to the city while skipping over the lowlights of their trip.

However, our analysis and visualization are meant to complement, not replace, traditional surveys and other means of data collection. In the pre-digital age tourism officials could know how many visitors spent a night in a hotel, but now we can also use feedback mechanisms on public web sites to estimate how much they enjoyed their stay. Similarly, we could know how many tourists visited a given attraction; but now we can also know infer their experience of it through act of uploading a photography and the semantics of their description of it. Direct observation enables us to know the number of tourists in an area; but through the mobile phone network we can know their nationalities.

The shortcomings of single-site ticket sales as a correlating dataset requires us to pursue alternate strategies for relating our mobile and photographic data to real-world activity with traditional surveys. An additional research avenue is the understanding of the circumstances under which users tag their content with a street address or when they are tagged to a larger region. An initial analysis of our Flickr dataset suggests that the 123 German users tended to provide more accurate locational information than their 175 Spanish counterparts.

The results of further analysis may reveal distinct profiles of geo-referencing and geo-tagging photos. These profiles might be based on culture or nationality, the type of tourist in terms of their length of stay or familiarity with the city, their level of technical expertise or spatial orientation ability, and the type of task or type of environment visited. Other questions that should be considered relate to the types of situations during which users are more or less likely to use their mobile devices for data generation. Answers to these types of questions should allow us to define better the meaning of the data and to explore further their potential usage in social sciences and urban studies.

## Acknowledgments

We would like to thank Barcelona Media, Telecom Italia for their support and Telespazio for their satellite digital imagery. Also, we are indebted to many people at the Massachusetts Institute of Technology and the Universitat Pompeu Fabra for providing extremely stimulating research environments and for their generous feedback. In particular, thanks to Assaf Biderman, Liang Liu, Nicolas Nova, Jon Reades and Andrea Vaccari for letting us pick their brains. Of course, any shortcomings are our sole responsibility.

## References

1. O'Neill, E., Kostakos, V., Kindberg, T., gen. Schieck, A. F., Penn, A., Fraser, D. S., and Jones, T. Instrumenting the city: Developing methods for observing and understanding the digital cityscape. In *Ubicomp* (2006), pp. 315–332.
2. Goodchild, M. F. (2007). Citizens as voluntary sensors: Spatial data infrastructure in the world of web 2.0. *International Journal of Spatial Data Infrastructures Research*, 2:24–32.
3. Eagle, N. and Pentland, A. S. (2006). Reality mining: sensing complex social systems. *Personal Ubiquitous Comput.*, 10(4):255–268.
4. Latour, B. (2007). Beware, your imagination leaves digital traces, Column for Times Higher Education Supplement, 6th of April 2007. <http://www.bruno-latour.fr/poparticles/poparticle/P-129-THES-GB.doc>
5. Girardin, F., Dal Fiore, F., Blat, J., and Ratti, C. (2007). Understanding of tourist dynamics from explicitly disclosed location information. In The 4th International Symposium on LBS &

TeleCartography.

6. Fisher, D. Hotmap: Looking at geographic attention. *IEEE Trans. Vis. Comput. Graph.* 13, 6 (2007), 1184–1191.
7. Ratti, C., Pulselli, R. M., Williams, S., and Frenchman, D. (2006). Mobile landscapes: Using location data from cell-phones for urban analysis. *Environment and Planning B: Planning and Design*, 33(5):727 – 748.
8. González, M., C. Hidalgo, and A Barabási (2008), “Understanding individual human mobility patterns”, *Nature*, Vol.453, pp.779-782
9. Calabrese, F., Ratti, C., Real Time Rome. *Networks and Communication Studies*, vol. 20, nos. 3 & 4, (2006), pp. 247–258.
10. Dykes, J., Slingsby, A., and Clarke, K. (2007). Interactive visual exploration of a large spatio-temporal dataset: Reflections on a geovisualization mashup. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1176–1183.
11. Rattenbury, T., Good, N., and Naaman, M. (2007). Towards automatic extraction of event and place semantics from flickr tags. In SIGIR '07: Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval, pages 103–110, New York, NY, USA. ACM Press.

