Efficient Detection of Points of Interest from Georeferenced Visual Content

Ying Lu^{1*} and Juan A. Colmenares²

University of Southern California, USA ylu720@usc.edu
 Samsung Research America, USA juan.col@samsung.com

Abstract. Many people take photos and videos with smartphones and more recently with 360° cameras at popular places and events, and share them in social media. Such visual content is produced in large volumes in urban areas, and it is a source of information that online users could exploit to learn what has got the interest of the general public on the streets of the cities where they live or plan to visit. A key step to providing users with that information is to identify the most popular k spots in specified areas. In this paper, we propose a clustering and incremental sampling (C&IS) approach that trades off accuracy of top-k results for detection speed. It uses *clustering* to determine areas with high density of visual content, and *incremental sampling*, controlled by stopping criteria, to limit the amount of computational work. It leverages spatial metadata, which represent the scenes in the visual content, to rapidly detect the hotspots, and uses a recently proposed Gaussian probability model to describe the capture intention distribution in the query area. We evaluate the approach with metadata, derived from a non-synthetic, user-generated dataset, for regular mobile and 360° visual content. Our results show that the C&IS approach offers $2.8 \times -19 \times$ reductions in processing time over an optimized baseline, while in most cases correctly identifying 4 out of 5 top locations.

1 Introduction

Many people frequently use their smartphones and more recently their 360° cameras to take photos and videos to capture memorable subjects and situations at places and events (*e.g.*, touristic attractions, concerts, and political rallies). People also upload their visual content (*i.e.*, photos and videos) very often to social media websites, such as Facebook, Flickr, Instagram, and YouTube. Urban areas naturally produce visual content in large volumes. A recent study [1] indicates that over 75% of people in New York City (NYC) own smartphones; *i.e.*, \sim 6.3 million people and \sim 0.3% of roughly 2 billion smartphone users worldwide [2]. With 350+ million photos and videos being uploaded daily to Facebook [3], NYC may produce over 1 million pieces of visual content per day.

Such abundant and continuously generated visual content offers online users the opportunity to learn about subjects, places, and events that have caught the attention of people (physically present) in a given area. For example, users often search the web to know what popular attractions are currently in their city as well as what interesting events have recently happened there. Most web services today answer this type of questions by looking at camera locations and timestamps photos and videos have been tagged with. The results, however, are inherently imprecise because the camera and the subject are usually at different locations and many times far apart (e.g., pictures of the Statue of Liberty are usually taken at a considerable distance from it). In addition, travel and review websites like TripAdvisor and Yelp are often limited to well-known, static landmarks, but interesting events can also happen at ad hoc locations (e.g., an amazing musical performance down the street).

In this paper, we focus on efficiently identifying the top-k most popular *points of interest* (POIs) from photos and videos. POIs are estimated locations of subjects captured in the visual content; the subjects' locations are obtained from analyzing (metadata of) visible scenes. We take into account that POIs are not necessarily static and their appeal may vary over time (*e.g.*, the Barra Olympic Park, in Rio de Janeiro, was the world's focus during the 2016 Summer Olympics, but was abandoned after six months [4]); so, we allow POIs to be identified in specific time intervals.

Background and Baseline Approach Early approaches in the literature (*e.g.*, [5,6]) identify POIs from visual content by extracting and analyzing image features. They are computationally intensive and thus not applicable to large volumes of visual content. To accelerate the process, researchers have proposed other approaches (*e.g.*, [7–13]) that

^{*} Ying Lu did most of this work as an intern at Samsung Research America.

leverage sensor-generated geo-metadata (*e.g.*, GPS locations, timestamps, and compass directions) associated with the visual content. From this group, most of the studies (*e.g.*, [7–10]) detect hotspots based on camera locations. However, using the camera location is insufficient to represent the coverage of a photo or video. On one hand, as mentioned above, the location of the camera and the location of the subject in a photo or video are often not the same and many times are far apart [11–13]. On the other hand, cameramen often move during video recording; thus, a single camera location is inadequate for a video with a trajectory. To avoid this issue, recent studies on POI identification [11–13] represent the visible scene of a photo or individual video frames with the spatial extent of its coverage area at a fine granularity (*i.e.*, geo-tagged at the video frame level). Such spatial extent of a scene is called *field of view* (FoV) [14] and is illustrated in Figure 1.

Among the studies based on the FoV model, the most recent approach [13] – the state-of-the-art method – identifies POIs from georeferenced videos by partitioning the query area into a grid with equally-spaced cells and using a Gaussian probability model to describe the capture intention distribution on the grid. This approach, described in Section 3, has been shown to be able to achieve high accuracy (≤ 1 m), and for that reason we adopt it as our baseline. However, the baseline approach, if implemented naively, takes long time to process big areas with large volumes of visual content as it computes the capture intention contribution of *all* FoVs to *every* cell in the query area (*e.g.*, over 1 hour to detect the top-5 spots in Munich and 19 hours in Los Angeles). Such long processing time would render the approach incapable to offer an interactive user experience.

Challenge The challenge in this work is how to accelerate top-k POI detection without much loss of *accuracy* when compared to the results of the baseline approach. In other words, how to trade off detection efficiency and accuracy. One may attempt to reduce the number of cells to be processed by increasing the cell size. This simple method can reduce the detection time proportionally to the cell count; however, it may significantly deteriorate the result accuracy mainly because larger cells cover the target area with a lower resolution (see details in Section 3.1). Another option is to reduce the number of FoVs simply by using random sampling. But, determining the right sample size that offers good detection speed and accurate results across multiple target areas is not straightforward.

Contributions To overcome the challenge, we propose a series of techniques by exploiting the spatial properties of FoVs. Our contributions are as follows:

- 1. We first introduce two practical optimization techniques (Section 4) that enable significant (up to 2000×) reduction in detection time with no accuracy loss. The first technique seeks to reduce the number of grid cells processed per FoV; thus, it only considers the cells that overlap the *minimum bounding rectangle* (MBR) of each FoV, rather than the entire grid. The second technique makes an adjustment to the probability model to properly handle and efficiently process 360° visual content, which is proliferating as 360° cameras and virtual reality headsets have entered the market. Both techniques are combined in an improved implementation of the adopted baseline approach, referred to as the *optimized baseline*.
- 2. Considering that densely populated areas are expected to contain very large number of FoVs and a fraction of the FoVs may be enough to identify the top-k spots, we propose an approach that combines two well-known techniques: 1) clustering to determine areas with high density of FoVs, where the hotspots are more likely to be, and 2) incremental sampling to limit the number of FoVs to be processed and thus reduce detection time. This clustering and incremental sampling (C&IS) approach, described in Section 5, relies on stopping criteria to make incremental sampling terminate after having some indications that the top-k results have been identified. By combining these techniques, we can flexibly trade off result accuracy and detection speed.
- 3. We conducted extensive experiments with a real-world geo-tagged video dataset [15] recorded with regular mobile phones plus two variants of the same dataset modified assuming the use of both regular and 360° cameras. Experimental results demonstrate that the C&IS approach brings $2.8 \times -19 \times$ improvements in processing time over the optimized baseline, while in most cases correctly identifying 4 out of 5 top locations.

2 Field of View

Visual content can be captured along with metadata representing the scenes in it, particularly their spatial features. This is easily achievable today with smartphones equipped with cameras and all the necessary sensors. A photo and individual frames in a video (e.g., a frame every second in a 30-fps video) can be tagged with a *field of view* (FoV) [14],

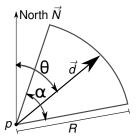


Fig. 1. Two-dimensional field of view (2D FoV).

a piece of metadata that describes the area covered by the captured scene. A two-dimensional FoV, illustrated in Figure 1, is stored as a tuple $\langle p, \theta, R, \alpha \rangle$, where p is the camera's location given by its latitude and longitude coordinates, θ is the camera's orientation or azimuth (*i.e.*, the angle between the north reference line and the camera's shooting direction d measured clockwise), R is the maximum visible distance from p at which an object within the camera's FoV can be recognized, and α is the camera's visible angle. Moreover, FoVs may contain an extra field t with the time at which the scene was captured.

With a smartphone or another sensor-rich camera device, p and θ can be read from the GPS receiver and digital compass, respectively, and α can be obtained based on the properties of the camera and its lens for the current zoom level [16]. In addition, the maximum visible distance can be estimated with the formula R = f.h/y, where f is the camera's focal length, h is the height of the visible object with maximum depth of view, and y is the object's height in the image, which is typically at least $1/20^{th}$ of the image's height [14].

Two-dimensional (2D) FoVs are pie-shaped ($\alpha < 360^{\circ}$) or circular ($\alpha = 360^{\circ}$); they assume that the camera and target are on the same plane, and only consider azimuth rotation (yaw movements). By contrast, FoVs in a three-dimensional space are cone-shaped ($\alpha < 360^{\circ}$) or spherical ($\alpha = 360^{\circ}$), and consider, besides yaw, the other two rotation axis: pitch and roll. In this work, we focus on 2D FoVs.

3 Baseline Approach

An approach to detecting points of interest from georeferenced videos has been recently proposed in [13]. It uses 2D FoVs to represent visible scenes and applies a Gaussian probability model to describe the capture intention distribution in a user-specified area. It has been shown to achieve high accuracy (*i.e.*, a meter or less between detected hotspots and their actual locations). We adopt this state-of-the-art approach as our *baseline*.

The approach first partitions the user-specified area A into a grid with equally-spaced cells. It assumes the centers of the cells are the visual targets (*i.e.*, each cell is represented by its center). Then, it calculates the capture intention of each cell c as follows:

$$\gamma(c,F) = \frac{\sum_{f \in F} ci(c,f)}{\max_{c \in C} \sum_{f \in F} ci(c,f)}$$
(1)

 $\gamma(c,F)$ is the normalized sum of the individual intentions of the FoVs covering the area A to capture the cell c. The set of FoVs overlapping A is denoted as F, and the set of cells forming the grid as C. Finally, the k cells with the highest cumulative capture intention are returned as the top-k points of interest.

The individual contribution of an FoV f to the capture intention of a cell c involves two factors (see Figure 2): the *angular difference* $|\theta - \theta_{pc}|$ between the camera's shooting direction d and the cell's center, and the Euclidean *distance* ||p,c|| between the camera's location and the cell's center. Such contribution is calculated as:

$$ci(c,f) = ci_a(c,f) \times ci_d(c,f)$$
(2)

where ci_a is the probability of capture intention with respect to the angular difference, and ci_d is the probability of capture intention with respect to the distance.

In general, people are more likely to capture the target in the center of the image (along d). Intuitively the capture intention increases as the angular difference decreases. For that reason, the approach adopts a Gaussian distribution to

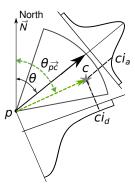
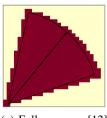
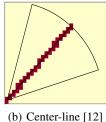
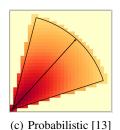


Fig. 2. An FoV's contribution to the capture intention of a cell c involves i) ci_a , the capture intention probability with respect to the angular difference $|\theta - \theta_{pc}|$, and ii) ci_p , the capture intention probability with respect to the distance ||p, c||.







(a) Full-coverage [12] (b)

Fig. 3. Models for estimating the capture intention contributions of an FoV. The darker the color of a cell, the higher its capture intention value.

model the angular capture intention $ci_a(c, f, \sigma_a)$ that an FoV $f \triangleq \langle p, \theta, R, \alpha \rangle$ has for a cell c. The distribution is given by:

$$ci_{a}(c, f, \sigma_{a}) = \begin{cases} \frac{e^{-\frac{(\theta - \theta_{pc})^{2}}{2\sigma_{a}^{2}}}}{\sqrt{2\pi}\sigma_{a}} : & |\theta - \theta_{pc}| \ge \frac{\alpha}{2}, \\ 0 : & \text{otherwise} \end{cases}$$
(3)

where σ_a is the variance.

Similarly, people tend to capture the target close to the camera for better visibility; so, the closer the target the higher the capture intention. In this case, the approach models the distance-based capture intention $ci_d(c, f, \sigma_d)$ of an FoV f on a cell c with the following distribution:

$$ci_d(c, f, \sigma_d) = \begin{cases} \frac{e^{-\frac{\|p, c\|^2}{2\sigma_d^2}}}{\sqrt{2\pi}\sigma_d} : & \|p, c\| \le R, \\ 0 : & \text{otherwise} \end{cases}$$
(4)

where σ_d is the variance.

Algorithm 1 presents the pseudocode of the baseline approach [13]. The CALCCIMATRIX procedure computes the capture intention matrix of FoVs. Specifically, line 12 calculates the individual contribution of an FoV to the capture intention of a cell, using the Equations (2), (3) and (4).

This Gaussian probability model is the key distinguishing feature of the state-of-the-art approach we have adopted as our baseline [13]. Figure 3 illustrates the existing models for calculating the capture intention contributions of an FoV. The full-coverage model [11, 12] gives the same value to every cell (or point) inside an FoV. The center-line model [11, 12] only considers the cells in the FoV along the shooting direction. With the probabilistic model, as shown in Figure 3(c), capture intention values smoothly spread over the space inside the FoV. In particular, the contribution from an FoV to the capture intention of a cell increases as the cell gets closer to the camera's location and to the

Algorithm 1 Naive baseline approach.

```
Input:
      k: Number of top cells to detect.
      A: Area of interest, specified by its minimum and maximum latitude, and minimum and maximum longitude.
      T: Time interval of interest.
      l: Length of each side of the cell forming the grid.
      \sigma_a: Variance of angular capture intention distribution.
      \sigma_d: Variance of distance-based capture intention distribution.
Output:
      K: Set with the top-k cells.
     F \leftarrow \text{GETFoVsInRange}(A \mid T)
                                                                                                                                          > Range query to obtain set of FoVs
     Determine the size of a matrix M: xdim \times ydim from A and l
     M \leftarrow \mathsf{ZEROMATRIX}(xdim, ydim)
                                                                                                                                                             ▶ Initialize matrix
 4: CALCCIMATRIX(F, \sigma_a, \sigma_d, 0, xdim, 0, ydim, M)

    ▷ Calculate the caption intention matrix.

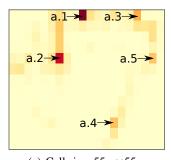
     K \leftarrow \text{GETTOPKCELLS}(M)
                                                                                                     Dobtain top-k cells with the highest capture intentions from the matrix.
 6: return K
      procedure CALCCIMATRIX(F, \sigma_a, \sigma_d, x_{min}, x_{max}, y_{min}, y_{max}, M)
 8:
9:
        for each f in F do
           for y \leftarrow y_{min} to y_{max} do
10:
              for x \leftarrow x_{min} to x_{max} do
11:
                Let center be the center of the cell (x, y)
12:
13:
                 ci \leftarrow CALCCAPTUREINTENTION(f, center, \sigma_a, \sigma_d)
                 if ci > 0 then
14:
                   M(x,y) \leftarrow M(x,y) + ci
15:
                 end if
16:
              end for
17:
           end for
18:
         end for
19:
     end procedure
```

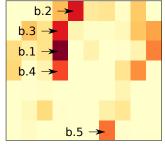
camera's shooting direction. Research indicates that people tend to focus on the center of an image [17]. Additionally, a closer object is likely to be more prominent in an image or video frame. Hence, the probabilistic model can yield better accuracy than the other two models. For example, at Merlion Park (Singapore), with the probabilistic model, the distances between the detected POIs and their ground-truth locations are less than 0.8 meters, while the distance error with the center-line model, which outperforms the full-coverage model [12], is more than 35 meters [13].

3.1 Drawbacks

Despite the high accuracy, the naive baseline has a major drawback: The algorithm (i.e., the procedure CALCCI-MATRIX) takes long time to process big areas with large number of FoVs since it computes the capture intention contribution of all FoVs to every cell in the query area. Its time complexity is $\mathcal{O}(N_c \times N_f)$, where N_c is the number of cells in the grid, and N_f is the number of FoVs. For example, according to our evaluation (see Section 6), it takes over 19 hours to find the top-5 cells in Los Angeles area (\sim 1,300 km²) with over 12.6 million cells of \sim 11m×11m and about 52,000 FoVs from the user-generated dataset GeoUGV [15]. Back in 2009, it was estimated that 4,306 photos were uploaded daily to Flickr from Los Angeles [18]. Assuming proportionality to FoV count, the detection time of the top-5 cells in that case would be over 1.5 hours; i.e., the user may need to wait way more than an hour for a visual summary of what happened in Los Angeles in the last 24 hours. Such long processing time is clearly unacceptable for interactive user experience.

Since the approach's complexity is $\mathcal{O}(N_c \times N_f)$, its detection time can be shortened by reducing the number of cells (N_c) and/or the number of FoVs (N_f) that need processing. An easy way to reduce N_c is to increase the cell size. Unfortunately, while it reduces the detection time proportionally to the cell count, the accuracy may deteriorate significantly. For example, Figure 4 shows the capture intention heatmaps at Merlion Park (Singapore) for two cell sizes. Looking at the discrepancies between the two top-5 result sets, Figure 4(b), with larger cells, misses the 3rd and 5th cells (a.3 and a.5) in Figure 4(a) and misidentifies its 3rd and 4th cells (b.3 and b.4). The reason is that the approach works best on a fine grid, with cells much smaller than the FoVs. This way cells can record high-resolution changes in capture intention on the grid. As the cell size increases, the centers of the cells covered by an FoV change, along with the polygon those cells form (see Figure 5). And, since a cell's capture intention depends on the location of its center, the capture intention of a large cell is in general not equal to the sum of the capture intentions of smaller adjacent cells that make up the large cell. In addition, in practice it is difficult to set a single cell size that properly balances processing time and accuracy because FoVs come in different sizes (i.e., their parameters R and α vary).





(a) Cell size: 55m×55m

(b) Cell size: 111m×111m

Fig. 4. Different top-5 results with two cells sizes at Merlion Park.

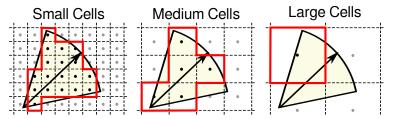


Fig. 5. As the cell sizes increase, the centers of the cells covered by an FoV change, along with the polygon those cells form.

4 Optimized Baseline

In this section, we propose two practical techniques to accelerate the baseline approach with no need to increase the cell size. The optimized procedure implementing these techniques is denoted as CALCCIMATRIX⁺, and it is interchangeable with the naive procedure in Algorithm 1 as they both receive the same parameters.

MBR-based Cell Filtering This technique seeks to reduce the number of cells (N_c) to be processed per FoV. The idea is to only compute the capture intention probability of the cells that fully or partially overlap the *upright minimum* bounding rectangle (MBR) of each FoV (see Figure 6). The reason is that an FoV surely makes no contribution to the capture intention of the cells not covered by its MBR; those cells can then be ignored when processing the FoV.

This simple, yet effective optimization exploits the fact that each FoV often covers a small fraction of the whole area of interest. For example, Los Angeles and Munich metropolitan areas cover 12.5E9 m² and 27.7E9 m², respectively; by contrast, a circular FoV with a maximum visible distance of 100 m only covers 31.4E3 m² (six orders of magnitude less). This technique requires FoVs to be augmented with their upright MBRs. MBRs are aligned with

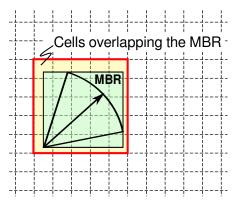


Fig. 6. MBR-based cell filtering. For each FoV, our approach only considers the cells that overlap the upright minimum bounding rectangle (MBR) of the FoV.

the geographic coordinate system and comprise four coordinate values each (maximum and minimum latitude, and maximum and minimum longitude). One option is to store the FoVs along with their MBRs, where MBRs are either provided by the camera devices or calculated by the storage system at ingestion time. Another option is to calculate the FoVs' MBRs on the fly at query time. We consider that storing FoVs with their MBRs is more beneficial given the significant speedups this optimization brings. Our evaluation (Section 6.2) shows that MBR-based cell filtering reduces the processing time up to three orders of magnitude when compared to the naive baseline.

Improving Efficiency for 360° Visual Content The baseline approach can be easily improved even further to process 360° visual content more efficiently. The main observation is that the intention of a circular FoV (with visible angle $\alpha=360^{\circ}$) is to capture its entire surrounding. That is, the angular capture intention of a circular FoV is constant and equal to 1 in every direction ($ci_a \triangleq 1$). Hence, ci_a need not be computed for the cells covered by circular FoVs, saving some time in detecting the top-k cells.

5 Clustering and Incremental Sampling (C&IS)

MBR-based cell filtering reduces the detection time of the baseline approach by limiting the number of cells that need to be processed per FoV. But densely populated areas are expected to contain very large number of FoVs, and a fraction of the FoVs may be sufficient to identify the 5 to 10 top spots. A simple approach to limit the number of FoVs is to first take a sample, uniformly distributed and without repetitions, of the FoV population in the query area, and then call CALCCIMATRIX⁺ to identify the top-k cells using the sample. We call this approach *single sampling*. Clearly, the less FoVs in the sample, the faster the results are produced, but the less accurate those results are likely to be. Single sampling is not flexible in trading off efficiency and accuracy for two reasons. First, it samples the FoV population once with a predetermined sample size. It is difficult to determine the optimal sample size to achieve the best performance (*i.e.*, fast detection time and high accuracy) for a given area. Second, the single sampling approach samples all the FoVs in the entire target area. Different datasets and areas may have different density distributions of FoVs; thus, their optimal sample sizes may be very different. To overcome these drawbacks, we proceed to present our approach that is based on clustering and incremental sampling techniques.

People usually take photos and videos at the same locations (*e.g.*, touristic attractions, stadiums, and concert venues). So there is no surprise that social media visual content is often heavily concentrated around specific spots. We leverage this by focusing on areas with high density of FoVs, where the top-k cells are more likely to be, and ignoring sparse areas. This, in turn, helps reduce the number of cells to be processed and, to a lesser extent, the number of FoVs.

We apply a *clustering* technique to determine the high-density areas. Another reason we use clustering is that we can have different sample sizes for different clusters depending on their FoV distributions. In addition, since each cluster may still contain numerous FoVs, we need to further sample FoVs in each cluster to limit the number of FoVs to be processed. We try to sample the minimum fraction of FoVs that still yields accurate results for the top-k POIs. To that end, we adopt an *incremental sampling* technique. Therefore, we combine both traditional techniques, clustering and incremental sampling (C&IS), into an approach that efficiently detects top-k POIs without significant loss in accuracy.

At the high level, our C&IS approach first identifies a set of clusters of FoVs in the query region, and then incrementally samples FoVs in each cluster. For each cluster, a small fraction of FoVs are sampled at each iteration. FoVs are sampled without repetitions, and each FoV is processed only once independetly of the number of sampling iterations. We know that the more iterations, the more FoVs are considered and the more accurate the results are, but at the expense of longer processing time. Our C&IS approach aims to flexibly trade off detection speed and result accuracy by determining the number of sampling iterations for each cluster with several heuristic stopping criteria. These criteria are based on the convergence of the detected top-k POIs and the capture intention matrix of the culster. With this mechanism, the sample size for a cluster of FoVs gradually increases iteration by iteration approximating to the "optimal" sample size. Further, the sample sizes of different clusters are decided according to their own density distributions of FoVs.

Algorithm 2 presents the pseudocode of the C&IS approach. We first find the FoVs that overlap with the target region (line 1) via a range query supported by an R-tree index [14, 19]. Then, we identify c clusters (line 3) from a

Algorithm 2 Clustering and incremental sampling approach.

Input:

- V: Video database.
- k: Number of top cells to detect.
- A: Area of interest.
- T: Time interval of interest.
- l: Length of each side of the cell forming the grid.
- σ_a : Parameter for Gaussian distribution of angular capture intention.
- σ_d : Parameter for Gaussian distribution of distance capture intention.
- c: Number of clusters to identify.
- f_c : Fraction of the FoV population used in cluster identification.
- f_i : Fraction of the FoV population used in each iteration.

Output:

K: Set with the top-k cells.

```
1: F \leftarrow \text{GETFoVsInRange}(V, A, T)
                                                                                             2: S \leftarrow \text{GETRANDOMSAMPLENOREPETITIONS}(F, f_c)
 3: C \leftarrow \text{IDENTIFYCLUSTERS}(S, c)
 4: K \leftarrow \text{HEAP}(k)
                                                                                    ▶ Initialize the global heap for top-k cells.
 5: for each cluster in C do
       radius \leftarrow CALCRMSRADIUS(cluster)
                                                                                            ⊳ Root Mean Square (RMS) radius
 6:
       A_c \leftarrow \text{GETBOUNDINGRECT}(cluster.center, radius)
 7:
       [x_{min}, x_{max}, y_{min}, y_{max}] \leftarrow \text{CALCCELLRANGE}(A, A_c, l)
 8:
 9:
       F_c \leftarrow \text{GETFoVsInRange}(F, A_c, T)
       M_c \leftarrow \text{ZeroMatrix}(xdim, ydim)
10:
11:
       iter \leftarrow 0
12:
       do
          S' \leftarrow \text{GETRANDOMSAMPLENOREPETITIONS}(F_c, f_i)
13:
14:
          M' \leftarrow \text{ZEROMATRIX}(xdim,ydim)
15:
          CALCCIMATRIX<sup>+</sup>(S', \sigma_a, \sigma_d, x_{min}, x_{max}, y_{min}, y_{max}, M')
16:
          M_c \leftarrow M_c + M'
                                                                                                ▶ Update matrix for the cluster.
          K_c \leftarrow \text{GETTOPKCELLS}(M_c)
                                                                                            ▷ Obtain top-k cells for the cluster.
17:
          Exclude S' from F_c
                                                                      > Same FoVs will not be reused at different iterations.
18:
19:
          iter \leftarrow iter + 1
       while SATISFYSTOPCRITERIA(K_c, M_c, iter)
20:
21:
       p \leftarrow min(iter \times f_i, 1)
                                                               > The percentage of FoVs in the cluster that are considered.
       K.\mathsf{UPDATE}(K_c,p)
                                                                                                ▶ Update the global top-k cells.
22:
23: end for
24: return K
```

uniformly random sample (with no repetitions) containing a fraction f_c of the FoV population (line 2). Our current implementation uses k-means as the clustering algorithm. In our experience, $c \in [k, 2 \times k]$ and $f_c \in [0.2, 0.5]$ work reasonably well. Next, we calculate top-k POIs for each cluster and use them to update the final top-k POI results (lines 5–23). For each cluster, we obtain all the FoVs F_c that belong to it (line 9). FoVs in F_c are then incrementally sampled to update the capture intention matrix of the cluster (lines 12-20), Once the *stopping criteria* are satisfied (line 20), we obtain from this matrix the top-k POIs identified in the cluster so far. The top-k POIs K_c of the cluster are used to update the global top-k POIs K (line 22). Since different clusters may have different sample sizes, for fair comparison among the top-k cells of different clusters, the capture intention value ci of a cell in a cluster is estimated as ci/p, where p is the total fraction of the cluster's FoV population that was considered (i.e., $iter \times f_i$, line 21), assuming the clusters have the same sample size.

Stopping Criteria The stopping criteria are responsible to tell the algorithm to cut the iterations short and stop processing a cluster, after having some *indication* that the top-k cells for the cluster have already been identified. They tend to reduce the detection time, but often at the expense of some loss in accuracy.

First of all, the iteration number is constrained to not exceed the maximum iteration number (i.e., $iter \leq 1/f_i$). Besides that, we consider two heuristic stopping criteria.

- The first criterion monitors changes in capture intention. It evaluates whether the difference between the maximum capture intentions in M_c from one iteration to the next is less than a threshold, in which case the algorithm stops processing the current cluster. We refer to this criterion as the difference in maximum capture intention.
- The second criterion monitors changes in the locations of the top-k cells from one iteration to the next. It calculates the sum of the distances between the closest pairs of cells in the two top-k result sets from subsequent iterations. Cells are taken in pairs, one from each result set, in a closest-pair-first manner, each cell is considered only once, and the distance between their centers is accumulated. The criterion then checks whether such sum is less than a threshold, in which case it tells the algorithm to stop processing the current cluster. We refer to this criterion as the sum of minimum distances between top-k results.

Section 6.3 examines interesting aspects of the C&IS approach, particularly the trade-offs between detection time and accuracy.

6 Evaluation

In this section, we first show the performance gains of the optimized baseline over the naive baseline. Then, we evaluate the clustering and incremental sampling (C&IS) approach and study how its parameters influence detection time and the accuracy of top-k results.

6.1 Experimental Setup

Test Datasets We use the datasets in Table 1. They were generated from the GeoUGV dataset³ with a Python tool we developed. GeoUGV [15] is a real-world dataset that includes 2,397 videos and 208,976 FoVs, collected by \sim 300 users in more than 20 cities across the globe between 2007 and 2016. We modify GeoUGV by varying the maximum visible distance (R) and visible angle (α) of the FoVs. We also augment it with the FoVs' upright MBRs for the optimized baseline, featuring MBR-based cell filtering, to use. Note that camera's location (p) and orientation (θ) remain the same.

Table 1. Test Datasets.

Name	Description
DS _(100%,60°)	100% of FoVs with visible angle $\alpha = 60^{\circ}$.
	30% of FoVs with $\alpha=160^\circ$ and 70% of FoVs with $\alpha=360^\circ$.
DS _(70%,160°)	70% of FoVs with $\alpha=160^\circ$ and 30% of FoVs with $\alpha=360^\circ$.

³ Available at http://mediag.usc.edu/dataset/

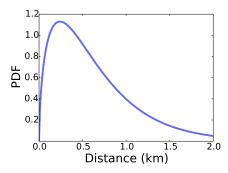


Fig. 7. Probability density function of Gamma distribution (shape=1.6, scale=0.4) used to generate R values.

Table 2. Target regions.

Locations	Area (km²)	N_c	N_{f}
Merlion Park	0.78	6,446	13,483
Munich	309.81	940,608	16,554
Singapore	709.96	1,441,396	64,296
Los Angeles	1,297.05	12,641,118	51,977

Table 3. Average number of cells coveredby MBRs of FoVs.

Locations	Data Sets				
Locations	(100%, 60°)	$(30\%, 160^{\circ})$	$(70\%, 160^{\circ})$		
Merlion Park	361	1,240	993		
Munich	517	2,310	1,756		
Singapore	519	2,310	1,750		
Los Angeles	519	2,314	1,771		

Instead of having a fixed R for all the FoVs (like in GeoUGV), we allow R values to vary to make our datasets slightly more realistic. We follow the intuition that people more often than not take videos and pictures of subjects that are close rather than far away, especially in social media [12]. To validate our intuition, we analyzed a sample of videos from the GeoUGV dataset. We first select all the well-known points of interest (e.g., Chinesischer Turm in Munich, Marina Bay Sands Skypark in Singapore, the White House in Washington, D.C.) in the dataset. Then we search all the videos within a large range (say 2km) around each point of interest. From the sample, we only considered the FoVs of video frames that captured one of those spots.

We calculated the distance between the point of interest and the camera's location (p) in over 600 FoVs. The results confirmed our intuition: videos were taken more often at a closer distance than at a longer distance. We also observed that the distance follows a distribution skewed to the right. For that reason, we generate R values from a Gamma distribution with shape and scale parameters equal to 1.6 and 0.4, respectively, ensuring that R < 2 km (see Figure 7). A single R value is obtained per video and assigned to all its FoVs.

We also vary α to simulate video content recorded with 360° video cameras. Samsung Gear 360^{4} and other cameras alike can operate in two modes: a dual-lens mode and a single-lens mode. We experimented with a Samsung Gear 360° and determined that $\alpha = 360^{\circ}$ with both lenses and $\alpha \approx 160^{\circ}$ with one lens. We hence use both angular values in different proportions in our datasets $DS_{(30\%,160^{\circ})}$ and $DS_{(70\%,160^{\circ})}$.

Cell Size, Target Regions, and Other Parameters We use in the experiments cells of 0.0001° of latitude by 0.0001° of longitude (*i.e.*, $\sim 11 \text{m} \times 11 \text{m}$). Moreover, our evaluation targets the regions in Table 2; they have relatively large numbers of FoVs in the GeoUGV dataset, while differing in area size. Note, however, that FoV counts are modest compared to what it is expected in reality (*e.g.*, >200 thousand photos and videos uploaded daily to Instagram in NYC area). Moreover, Table 3 lists the average number of cells covered per FoV in the different regions for our datasets.

⁴ http://www.samsung.com/global/galaxy/gear-360/

Table 4. Detection of top-5 cells on $DS_{(100\%,60^{\circ})}$.

Locations	Naive Baseline	Optimized Baselin	e
Locations	Proc. time	Proc. time (Speedup)	Diff.
Merlion Park	11.71s	2.50s (4.7×)	0.0
Munich	5,043.92s (1.4h)	8.63s (584.5×)	0.0
Singapore	30,886.66s (8.6h)	30.95s (998.0×)	0.0
Los Angeles	69,595.80s (19.3h)	26.19s (2657.3×)	0.0

Table 5. Detection of top-10 cells on $DS_{(100\%,60^{\circ})}$.

Locations	Naive Baseline	Optimized Baseline		
Locations	Proc. time	Proc. time (Speedup)	Diff.	
Merlion Park	11.72s	2.53 (4.6×)	0.03m	
Munich	5,044.13s (1.4h)	8.66 (582.5×)	0.0	
Singapore	30,887.22s (8.6h)	30.99 (996.7×)	0.0	
Los Angeles	69,601.41s (19.3h)	26.29 (2647.4×)	0.0	

Table 6. Detection of top-5 cells on $DS_{(30\%,160^{\circ})}$.

Locations	Naive Baseline	Optimized Baseli	ne
Locations	Proc. time	Proc. time (Speedup)	Diff.
Merlion Park	11.71s	3.90s (3.00×)	0.01m
Munich	4,975.64s (1.4 h)	25.48s (195.3×)	0.0
Singapore	31,009.70s (8.6 h)	89.23s (347.5×)	0.0
Los Angeles	61,129.10s (17.0 h)	75.94s (805.0×)	0.0

Table 7. Detection of top-5 cells on $DS_{(70\%,160^{\circ})}$.

Locations	Naive Baseline	Optimized Baseline		
Locations	Proc. time	Proc. time (Speedup)	Diff.	
Merlion Park	13.97s	5.42s (2.6×)	0.0	
Munich	4,995.01s (1.4 h)	31.05s (160.9×)	0.0	
Singapore	31,031.10s (8.6 h)	113.44s (273.5×)	0.0	
Los Angeles	56,627.50s (15.7 h)	94.70s (598.0×)	0.0	

We retrieve from the database all the FoVs covering the target regions in the time interval from 2010-03-18 to 2016-06-28. We verified across the experiments that the query time is just a small fraction of the total processing time. Moreover, the variance parameters for the angular and distance-based capture intention distributions are $\sigma_a=15^\circ$ and $\sigma_d=25$ m, respectively, in all the experiments. These values are suggested in [13] to effectively identify the points of interest.

Test Platform and Implementation We conduct our experiments on a MacBook Pro laptop running OS X 10.9.5 and equipped with a 2.6GHz dual-core Intel Core i5-4288U processor, 8GB of RAM (1600MHz DDR3), and a 512GB SSD. We use MySQL Community Server (GPL) v.5.7.15 (with MyISAM engine) to store the data. The table schema is omitted due to space limitations.

The approaches described in the paper have been implemented in C++11, and compiled using gcc with -03 optimization option. We use Boost uBLAS library for matrix operations, and libmysqlclient library to access the database.

6.2 Optimized Baseline

In this experiment, we evaluate the optimized baseline (Section 4) vs. the naive baseline (Section 3). The performance metrics for comparison are: 1) *total processing time*, which includes the time taken by both the detection procedure

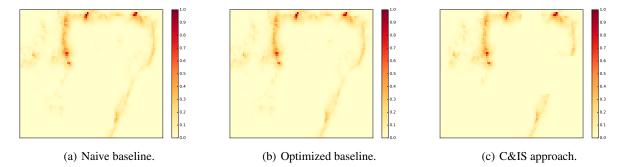


Fig. 8. Heatmaps at Merlion Park (Singapore) for $DS_{(100\%,60^\circ)}$. Although virtually indistinguishable, the heatmaps have some differences, especially that obtained with the clustering and incremental sampling (C&IS) approach. and the range query that retrieves the FoVs from the database; and 2) *difference between top-k results*, measured by the sum of the distances between the closest pairs of cells in two top-k result sets, as described in 5. The latter allows us to measure the accuracy of the optimized baseline compared to the naive baseline. Note that our results report average values across 30 runs, except when individual runs take longer than an hour to complete.

Tables 4 and 5 report the average processing time for the detection of the top-5 and top-10 cells on the $DS_{(100\%,60^\circ)}$ dataset. We observe that the optimized baseline, exploiting the MBR-based cell filtering technique, offers important speedups $(4.6\times$ to more than $2500\times$) over the naive baseline with essentially no difference in the top-k results. The last point is illustrated in Figures 8(a) and 8(b) that show identical heatmaps produced by the baseline and the optimized approaches at Merlion Park. Similarly, Tables 6 and 7 show that, on the datasets with 360° visual content, the optimized baseline also brings significant reductions $(2.6\times-805\times)$ in processing time while producing the same top-k results.

Our results show that the naive baseline takes very long time to detect top spots in large areas (*i.e.*, Munich, Singapore, and Los Angeles) with modest FoV counts, compared to those expected in reality. The reason is that the naive baseline, as explained in Section 3, uses a double nested loop to iterate over the FoVs computing the contribution that each FoV makes to the capture intention of *every* cell in the query area. By contrast, using MBR-based cell filtering the optimized baseline only processes the cells that overlap the MBR of each FoV. From Tables 2 and 3, we can see that this represents an important reduction in the number of cells to be considered per FoV in our test datasets (*e.g.*, from 6,446 cells in Merlion Park to an average of 1,240 or less, and from ~12.6 million cells in Los Angeles to an average of 2,310 or less). Therefore, MBR-based cell filtering, besides being simple and practical, proves to be very effective in reducing the processing time.

When comparing the results of the optimized baseline across datasets, we notice that the speedups for $DS_{(30\%,160^\circ)}$ and $DS_{(70\%,160^\circ)}$ (*i.e.*, the datasets with 360° visual content) are smaller than those for $DS_{(100\%,60^\circ)}$. This is explained again by the differences in the number of cells per FoV in Table 3 – up to 519 in average for $DS_{(100\%,60^\circ)}$, but between \sim 1,000 and \sim 2,300 for $DS_{(30\%,160^\circ)}$ and $DS_{(70\%,160^\circ)}$. Another interesting observation is that the speedup is larger when there are more circular FoVs with $\alpha=360^\circ$ than with $\alpha=160^\circ$ (*e.g.*, $805\times$ vs. $598\times$ for Los Angeles), even though circular FoVs cover more cells. This suggests that computing the angular capture intention is a more dominant factor than the number of cells covered by the FoVs' MBRs, making the case for the optimization of Section 4. Finally, we have observed that the naive and optimized baselines are both insensitive to k (the number of top cells being detected). The reason is that maintaining a heap with the top-k cells throughout the detection process is not the dominant factor.

6.3 Clustering and Incremental Sampling

In this section, we evaluate the clustering and incremental sampling (C&IS) approach, presented in Section 5. We use 50% of the FoV population for cluster identification ($f_c=0.5$), and 5% of the FoV population in each cluster as the incremental sample per iteration ($f_i=0.05$), for a maximum of 20 iterations per cluster. Other parameters, such as cell size, σ_a and σ_d , remain the same. As before, we report average values across 30 runs, unless stated otherwise.

We first study how the number of clusters and iterations influence the processing time and accuracy of the C&IS approach on the target regions for the different test datasets. We use various cluster counts (from 1 to 10), and stop processing each cluster at different numbers of iterations (5, 10, 15, and 20).

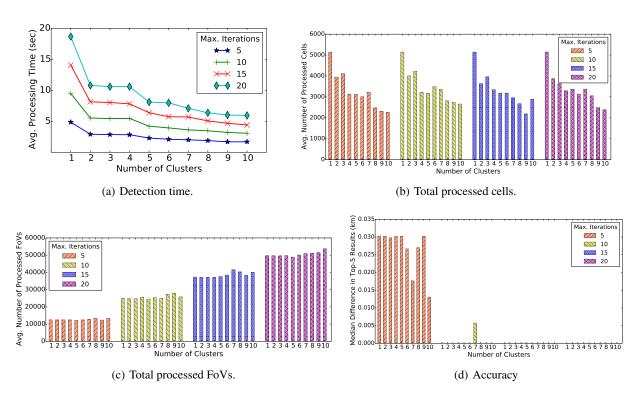
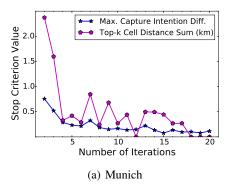


Fig. 9. Effects of varying cluster counts and maximum iterations in the detection of top-5 cells in Los Angeles for $DS_{(100\%,60^\circ)}$ dataset.



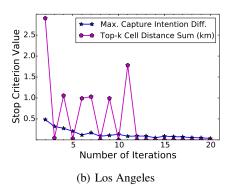
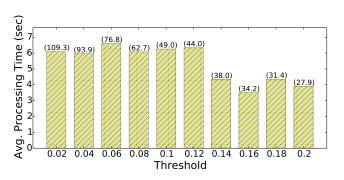
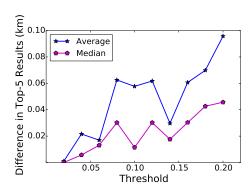


Fig. 10. Sample behavior of stopping criteria as iterations progress while processing FoVs in a cluster for $DS_{(100\%,60^{\circ})}$.





- (a) Detection time with average total number of iterations in parenthetical values
- (b) Loss in accuracy.

Fig. 11. Effect of varying the threshold of the difference in maximum capture intention for Los Angeles area and $DS_{(100\%,60^{\circ})}$.

We observe that the more clusters, the less processing time; this is exemplified in Figure 9(a) for Los Angeles area and $DS_{(100\%,60^\circ)}$. The time reduction is mainly because the number of processed cells decreases as the number of clusters increases, as shown in Figure 9(b). By contrast, the number of processed FoVs does not vary greatly with the cluster count (see Figure 9(c)). Our results then suggest that clustering filters out cells in sparse regions of the query area and helps focus the detection effort in the most popular regions. We also notice that beyond a certain number of clusters (*e.g.*, 6 in Figure 9(a)) there is no significant additional gain in detection time. In addition, as expected we observe that the more iterations, the longer the detection time, which is also illustrated in Figure 9(a). The reason is that the number of processed FoVs increases with the iterations performed per cluster (see Figure 9(b)).

More importantly, we notice that in many cases the C&IS approach does not need to reach the maximum number of iterations and process all the FoVs in the clusters in order to correctly identify the top-k cells. For example, Figure 9(d) shows that it can obtain the right top-5 cells with just 10 iterations (*i.e.*, 50% of the FoVs in each cluster). Consequently, there is opportunity for stopping criteria to cut the iterations short and identify the top-k cells without much loss in accuracy.

Next, we consider the stopping criteria introduced in Section 5; they are: 1) the difference in maximum capture intention, and 2) the sum of minimum distances between successively identified top-k result sets. We evaluate both criteria across the target areas and test datasets, with a fixed cluster count equal to 6 since it yields reasonably good processing times. As exemplified in Figure 10, our results indicate that the former is more stable, usually decreases and tends to converge as iterations proceed. Therefore, the difference in maximum capture intention is the only stop criterion we use in the rest of this section. By further investigating the criterion, we observe that as expected, the higher

Table 8. Detection of top-5 cells on $DS_{(100\%,60^{\circ})}$.

Locations	C&IS		Single Sampling		
Locations	Time (Speedup)	Diff.	Sample	Time (Speedup)	Diff.
Merlion Park	$0.89s(2.8\times)$	251.9m	58%	1.53s (1.6×)	121.1m
Munich	$1.51s(5.7\times)$	16.3m	30%	$2.77s(3.1\times)$	11.5m
Singapore	$5.88s(5.3\times)$	743.5m	20%	$6.93s~(4.5\times)$	1,536.8m
Los Angeles	$4.59s(5.7\times)$	30.0m	23%	$6.41s(4.1\times)$	62.9m

Table 9. Detection of top-5 cells on $DS_{(30\%,160^{\circ})}$.

Locations	C&IS			Single Sampling		
Locations	Time (Speedup)	Diff.	Sample	Time (Speedup)	Diff.	
Merlion Park	$1.17s(3.3\times)$	1,259.1m	66%	$2.57s(1.5\times)$	939.4m	
Munich	1.34s (19.0×)	8.4m	21%	$3.33s(7.7\times)$	4.8m	
Singapore	6.84s (13.0×)	0.5m	17%	$10.88s~(8.2\times)$	0.5m	
Los Angeles	4.62s (16.4×)	8.4m	18%	7.90s (9.6×)	5.3m	

Table 10. Detection of top-5 cells on $DS_{(70\%,160^{\circ})}$.

Locations	C&IS	C&IS		Single Sampling		
Locations	Time (Speedup)	Diff.	Sample	Time (Speedup)	Diff.	
Merlion Park	1.37s (4.0×)	1,036.6m	76%	3.40s (1.6×)	209.0m	
Munich	1.90s (16.3×)	16.4m	24%	$4.18s(7.4\times)$	9.5m	
Singapore	10.61s (10.7×)	3.6m	22%	$14.23s~(8.0\times)$	54.2m	
Los Angeles	6.09s (15.6×)	22.4m	19%	$10.07s (9.4 \times)$	13.0m	

the threshold, the shorter the detection time, but the lower the accuracy of the top-k results. The reason is that the higher the threshold, the lesser iterations are needed to satisfy the criterion. Figure 11 is illustrative of this observation.

We conclude this section with the evaluation of C&IS approach's overall performance for the detection of the top-5 cells in the target regions across the test datasets. As before, we use 6 clusters, and we adopt 0.1 as the threshold for the stop criterion based on the difference in maximum capture intention. We choose these parameters because they offer just reasonably good detection time across the datasets, as opposed to highly tuned parameters that yield the best possible performance for particular cases. The reason for doing this is that we want to assess how effective and robust the C&IS approach really is in practice.

Tables 8, 9, and 10 present the average processing time obtained with the C&IS approach. The (extra) speedups in the tables are with respect to the optimized baseline's processing times, reported in Tables 4, 6 and 7. As in the previous section, we also report the difference between the top-k results obtained with the C&IS approach and the optimized baseline. It is calculated as the sum of minimum distances between the two top-k result sets, which is described in Section 5 and denoted here as $\sum d_{min}$.

Our results show that C&IS brings important $(2.8 \times -19 \times)$ reductions in processing time over the optimized baseline, but at the expense of accuracy. Even though the top-k results can be reasonably accurate (i.e., $\sum d_{min} \leq 30 \ m$), like for Munich and Los Angeles, we see large $\sum d_{min}$ values for Merlion Park with the three test datasets and Singapore with DS_(100\%,60°).

As a point of comparison, Tables 8, 9, and 10 also include results with the *single sampling* approach, described at the beginning of Section 5. To be fair, in each case we configure single sampling to operate on a FoV sample whose size is equal to the average number of FoVs processed by the C&IS approach – *i.e.*, both approaches use roughly similar fractions of the FoV population. We observe that in average C&IS is faster than single sampling, and judging from the $\sum d_{min}$ values, both approaches offer somewhat similar accuracy.

We further investigate the accuracy of C&IS and single sampling approaches by examining the percentage of top-5 cells they correctly identify, compared to those obtained with the baseline. We consider that a cell has been "correctly" identified if it is 20 m or less from the closest actual top cell (*i.e.*, it is adjacent to the actual top cell). Here, cells also are taken in pairs, one from each top-k result set, in a closest-pair-first manner, and each cell is considered only once.

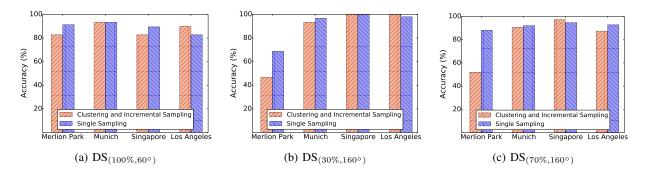


Fig. 12. Percentage of top-5 cells correctly identified by the clustering and incremental sampling approach and the single sampling approach.

Figure 12 indicates that in most cases both C&IS and single sampling detect at least 80% of the top-k cells (*i.e.*, 4 out of 5). For example, Figure 8(c) shows a heatmap produced by C&IS, on Merlion Park for $DS_{(100\%,60^\circ)}$, that is virtually indistinguishable to the naked eye from the heatmaps in Figures 8(a) and 8(b) by the baseline and optimized approaches; although there is some difference. Nevertheless, it is precisely at Merlion Park, but with $DS_{(30\%,160^\circ)}$ and $DS_{(70\%,160^\circ)}$ (*i.e.*, 360° visual content), where the C&IS approach performs rather poorly with the current parameters.

7 Discussion

We have shown that the optimized baseline offers significant reductions in detection time over the naive baseline, with virtually no loss in accuracy. It incorporates practical optimizations that have proven very effective. Particularly, MBR-based cell filtering makes a strong case for storing the FoVs along with their upright minimum bounding rectangles.

We have also shown that the C&IS approach (Section 5) is able to reduce detection time even further, but by sacrificing the accuracy of the results. C&IS was effective in a number cases (e.g., Munich and Los Angeles). But in other cases (e.g., Merlion Park), the accuracy loss was noticeable in terms of distance. This is because the popular subjects at different locations in Merlion Park have comparable capture intentions, as calculated from the GeoUGV dataset (see Figure 8(a)). Due to the use of random sampling, the C&IS approach may miss popular POIs in areas where other POIs exist with similarly high capture intentions. This observation suggests that C&IS is more suitable for detecting POIs whose capture intentions are more diverse.

According to our evaluation, the main reason for accuracy loss is incremental sampling, which has two main components: the FoV sampling method, and the stopping criterion. First, we draw samples of the FoV population uniformly at random (without repetitions). Besides its simplicity and being a popular choice, this method was chosen because it sets a reference point for more sophisticated methods. We expect that more advanced sampling methods will better guide the incremental selection of FoVs and help cap the accuracy loss. Second, the adopted stopping criterion based on the difference of maximum capture intention is rather simple and in some cases, clearly insufficient to ensure good accuracy. Other more elaborate stopping criteria may be more effective.

In Merlion Park, we also observed that the accuracy of the results improved by reducing the fraction of circular FoVs ($\alpha=360^{\circ}$) in the dataset. That suggests that the angular capture intention $ci_a \triangleq 1$ of circular FoVs, which makes cells less differentiable, has some negative effect when, rather than all, only a sample of FoVs is considered.

Note that using the optimized baseline on small regions (like Merlion Park) and the C&IS approach on large ones is also a practical alternative.

8 Related Work

Generally speaking, previous work focuses on two problems: 1) detection of points or regions of interests, and 2) identification and retrieval of the top-k of those spots of interest. Below we summarize research efforts in those areas most relevant to our work.

8.1 Detection of Points and Regions of Interest

Some approaches *e.g.*, [5,6]) identify POIs from visual content by extracting and analyzing image features. For example, Duygulu et al. use content-based techniques to extract image features which are then matched to keywords taken from bag-of-words vocabularies. They are computationally intensive and thus not applicable to large volumes of visual content. By contrast, instead of analyzing the visual content, other approaches (*e.g.*, [7,10,12,13]) accelerate the detection of interesting locations or objects by leveraging sensor-generated metadata (*e.g.*, GPS locations) or keyword tags associated with the visual content. The approaches presented in this paper belongs to this group. For example, Liu et al. [10] propose a filter-refinement framework to discover hot topics based on the spatio-temporal distributions of geo-tagged videos from YouTube. Zheng et al. [7] built a landmark recognition engine which models and identify the landmarks automatically from geotagged photos at the world scale. Unlike our work that considers FoVs, these frameworks only use the camera locations to describe the visual contents whereas the locations that are of interest to people may be far away.

The two recent studies [12, 13] are the most closely related to our work. Hao et al. [11, 12] represent each video frame as a camera view (*i.e.*, a vector pointing along the camera shooting direction) and propose two methods to detect POIs: 1) a cluster-based method and 2) a grid-based method. The cluster-based method computes the intersection points of all the camera views and from these intersections, infers clouds of points as POIs. The grid-based method, on the other hand, divides the space into grid cells, generates a heatmap based on how often a cell appears in different camera views, and then identifies the popular the places. A sector-based cell filtering technique is applied to accelerate the detection. However, this study processes all the FoVs, which is still not efficient for large-scale FoVs. Further authors in [11, 12] assume that people's intention is to only capture targets at the center of the scene (*i.e.*, aligned with the camera's shooting direction). However, targets may be located in different places withing the visible area. To overcome this limitation, Zhang et al. [13] propose an FoV-based approach that applies a probabilistic model to describe people's capture intention (see Section 3). They experimentally show that their approach offers much higher accuracy than the approaches in [11, 12]. For that reason, we adopted this approach as our baseline, and have shown that our improvements offer significant speedup.

Other efforts try to detect points or regions of interest from other data sources. For example, Vu et al. [20] present a framework for estimating social point-of-interest boundaries from spatio-temporal information in geo-tagged tweets. Ye et al. [21] and Yuan et al. [22] provide POI recommendation approaches based on users' check-in behaviors. Gao et al. [23] build a content-aware POI recommendation system by relating the content information on location-based social networks (i.e., POI properties, user interests and sentiment indications) to check-in actions.

Note that most of the works mentioned above focus on identifying all the points of interests in an area. Our work, however, focuses on top-k detection.

8.2 Retreival of Top-k Points of Interests and Objects

Peng et al. [24] propose a probabilistic field of view (pFoV) model for smartphone photos to capture the uncertainty in camera sensor data. Given a database of POIs, a set of geotagged photos represented in the pFOV model, and a query photo, authors identify the most prominent POI captured in the query photo. Skovsgaard et al. [25,26] focus on retrieving the top-k points of interest from spatial-keyword data (*e.g.*, geo-tagged twitter data). Toyama et al. [27] introduce an image database that supports image indexing and search based on the image camera locations and recording time.

The research work mentioned above assume that the points of interest are given, whereas in this paper we focus on detecting those spots without prior knowledge.

9 Conclusions

In this paper, we have presented an efficient approach to detecting top-k points of interest from geotagged visual content in a user-specified area. Based on clustering and incremental sampling, it trades off accuracy of top-k results for detection speed. We provided a thorough evaluation of the speedups as well as accuracy losses of the proposed approach. Our results show that the C&IS approach offers $2.8 \times -19 \times$ reductions in processing time over an optimized baseline, while in most cases correctly identifying 4 out of 5 top locations.

We also introduced two simple, yet effective optimization techniques that enable significant reduction in detection time with no accuracy loss. In particular, the MBR-based cell filtering technique makes a strong case for storing FoVs along with their MBRs.

In the future, we plan to study advanced sampling methods and more sophisticated stopping criteria that could offer accuracy improvements. Moreover, despite the obtained speedups, our prototype implementation could be optimized even further. For example, we could leverage multicore processors and process the clusters in parallel. Caching techniques could also be used to reduce response time since multiple users are likely to make the same or similar queries.

Acknowledgments

We thank our colleagues in the Smart Systems team at Samsung Research America for their feedback. Ying Lu's work was partially funded by the NSF grants IIS-1320149, CNS-1461963, and the USC Integrated Media Systems Center. Except for funding, neither sponsor contributed to or influenced any part of this paper. Nothing herein represents the views and opinions of the sponsors.

References

- [1] de Blasio, B., Menin, J.: New York City Mobile Services Study, Research Brief. New York City Department of Consumer Affairs. (November 2015) http://wwwl.nyc.gov/assets/dca/MobileServicesStudy/Research-Brief.pdf.
- [2] Kissonergis, P.: Smartphone ownership, usage and penetration by country. http://thehub.smsglobal.com/smartphone-ownership-usage-and-penetration(October 2015)
- [3] Lowe, L.: 125 amazing social media statistics you should know in 2016. https://socialpilot.co/blog/125-amazing-social-media-statistics-know-2016/(September 2016)
- [4] Sport Illustrated: Watch: What Rio's abandoned olympic venues look like today. http://www.si.com/olympics/2017/02/04/rio-de-janeiro-abandoned-olympic-venues-photos-videos (March 2017)
- [5] Duygulu, P., Barnard, K., Freitas, J.F.G.d., Forsyth, D.A.: Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In: Proc. of the 7th European Conference on Computer Vision. (2002) 97–112
- [6] Sivic, J., C. Russel, B., A. Efros, A., Zisserman, A., T. Freeman, W.: Discovering objects and their location in images. In: Proc. of the 10th IEEE Int'l Conference on Computer Vision. (2005)
- [7] Zheng, Y., Zhao, M., Song, Y., Adam, H., Buddemeier, U., Bissacco, A., Brucher, F., Chua, T., Neven, H.: Tour the world: Building a web-scale landmark recognition engine. In: 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009). (June 2009) 1085–1092
- [8] Ji, R., Xie, X., Yao, H., Ma, W.Y.: Mining city landmarks from blogs by graph modeling. In: Proc. of the 17th ACM Int'l Conference on Multimedia. (2009) 105–114
- [9] Kennedy, L., Naaman, M., Ahern, S., Nair, R., Rattenbury, T.: How flickr helps us make sense of the world: Context and content in community-contributed media collections. In: Proc. of the 15th ACM Int'l Conference on Multimedia. (2007) 631–640
- [10] Liu, K., Xu, J., Zhang, L., Ding, Z., Li, M.: Discovering hot topics from geo-tagged video. Neurocomput. 105 (April 2013) 90–99
- [11] Hao, J., Wang, G., Seo, B., Zimmermann, R.: Keyframe presentation for browsing of user-generated videos on map interfaces. In: Proc. of the 19th Int'l Conference on Multimedia (MM'11). (November 2011) 1013–1016
- [12] Hao, J., Wang, G., Seo, B., Zimmermann, R.: Point of interest detection and visual distance estimation for sensor-rich video. IEEE Trans. Multimedia **16**(7) (2014) 1929–1941
- [13] Zhang, Y., Zimmermann, R.: Efficient summarization from multiple georeferenced user-generated videos. IEEE Transactions on Multimedia 18(3) (March 2016) 418–431
- [14] Ay, S.A., Zimmermann, R., Kim, S.H.: Viewable scene modeling for geospatial video search. In: Proc. of the 16th ACM Int'l Conference on Multimedia (MM'08). (2008) 309–318
- [15] Lu, Y., To, H., Alfarrarjeh, A., Kim, S.H., Yin, Y., Zimmermann, R., Shahabi, C.: GeoUGV: User-generated mobile video dataset with fine granularity spatial metadata. In: Proc. of the 7th Int'l Conference on Multimedia Systems (MMSys'16). (2016)
- [16] Graham, C.H., Bartlett, N.R., Brown, J.L., Mueller, C.G., Hsia, Y., Riggs, L.A.: Vision and Visual Perception. John Wiley & Sons Inc. (1965)

- [17] Judd, T., Ehinger, K.A., Durand, F., Torralba, A.: Learning to predict where humans look. In: IEEE 12th International Conference on Computer Vision, ICCV. (2009) 2106–2113
- [18] Crandall, D.J., Backstrom, L., Huttenlocher, D., Kleinberg, J.: Mapping the world's photos. In: Proc. of the 18th Int'l Conference on World Wide Web. (2009) 761–770
- [19] Lu, Y., Shahabi, C., Kim, S.H.: Efficient indexing and retrieval of large-scale geo-tagged video databases. GeoInformatica **20**(4) (2016) 829–857
- [20] Vu, D.D., To, H., Shin, W., Shahabi, C.: GeoSocialBound: an efficient framework for estimating social POI boundaries using spatio-textual information. In: Proc. of the 3rd Int'l ACM SIGMOD Workshop on Managing and Mining Enriched Geo-Spatial Data (GeoRich@SIGMOD 2016). (June 2016)
- [21] Ye, M., Yin, P., Lee, W.C., Lee, D.L.: Exploiting geographical influence for collaborative point-of-interest recommendation. In: Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '11, New York, NY, USA, ACM (2011) 325–334
- [22] Yuan, Q., Cong, G., Ma, Z., Sun, A., Thalmann, N.M.: Time-aware point-of-interest recommendation. In: Proceedings of the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval. SIGIR '13, New York, NY, USA, ACM (2013) 363–372
- [23] Gao, H., Tang, J., Hu, X., Liu, H.: Content-aware point of interest recommendation on location-based social networks. In: AAAI. (2015) 1721–1727
- [24] Peng, P., Shou, L., Chen, K., Chen, G., Wu, S.: The knowing camera: Recognizing places-of-interest in smartphone photos. In: Proc. of the 36th Int'l ACM SIGIR Conference on Research and Development in Information Retrieval. (2013) 969–972
- [25] Bøgh, K.S., Skovsgaard, A., Jensen, C.S.: Groupfinder: A new approach to top-k point-of-interest group retrieval. PVLDB 6(12) (2013) 1226–1229
- [26] Skovsgaard, A., Jensen, C.S.: Top-k point of interest retrieval using standard indexes. In: Proc. of the 22nd ACM SIGSPATIAL Int'l Conference on Advances in Geographic Information Systems. (2014) 173–182
- [27] Toyama, K., Logan, R., Roseway, A.: Geographic location tags on digital images. In: Proc. of the 11th ACM Int'l Conference on Multimedia. (2003) 156–166