World Scientific
www.worldscientific.com

# SELECTIVE PRESENTATION OF PERCEPTUALLY IMPORTANT INFORMATION TO AID ORIENTATION AND NAVIGATION IN AN URBAN ENVIRONMENT

MAXIMINO BESSA*, ANTONIO COELHO† and JOSE BULAS CRUZ‡

*Department of Engineering, University of Tras-os-Montes e Alto Douro
Quinta de Prados 5000-911 Vila Real, Portugal*
*maxbessa@utad.pt
†acoelho@utad.pt
‡jcruz@utad.pt

ALAN CHALMERS

*Department of Computer Science, University of Bristol, Bristol, UK*
alan@cs.bris.ac.uk

A map can be a major advantage when exploring unfamiliar environments. With the proliferation of mobile devices, such as PDAs and even mobile phones, the tourist industry is currently exploring the potential of new presentation strategies that will maximize the promotional appeal of tourism in their region. Mobile devices are capable of guiding a tourist when he/she is exploring a city. These mobile devices offer the potential for providing relevant 3D information to enable tourists to locate themselves within the city, rapidly navigate around the unfamiliar environment and explore it interactively. However, the computational resources of current mobile technology prevents the display of full complex 3D content in real-time, and thus selective rendering techniques must be adopted to ensure the viewer is provided with the perceptually most important information at interactive rates.

    This paper presents a series of experiments which help to identify key features of a scene for users to orientate themselves in that environment. Knowledge of these salient key features enable them to be provided to a user at a high quality while the remainder of the scene can be rendered in a much lower quality, saving significant bandwidth and computing power, without the user being perceptually aware of this difference in quality within the image.

*Keywords*: Mobile devices; 3D maps; visual perception; inattentional blindness.

## 1. Introduction

Two-dimensional maps are an ancient media for helping us navigate through unfamiliar environments. With modern GPS assistance it is even possible to (approximately) locate our position on a map. However, a map is an abstract representation of an environment and as such, it simply may not be possible to orientate the map

correctly in the absence of key information, such as road names. The next generation of maps should provide a more realistic representation of our world, a high quality 3D representation, and they should be where we need them most, on our Smart phones, PDAs and Laptops. These maps will be more suitable for tourism purposes.

There are many issues that still need to be addressed if such high quality 3D maps are to be provided at interactive rates on mobile devices. Although the next generation of mobile devices may partially resolve some of the performance issues that exist today, such as bandwidth, storage and the small dimension of displays, this is likely to further increase demand for even more realistic and complex 3D applications. Visual perception is one approach that can be used to alleviate this problem. By knowing exactly what users are looking at on the mobile device it may be possible to render only this part of the image at high quality while the rest of the scene could be rendered at a lower quality for a fraction of the computation cost, without the user being aware of this quality difference.

In this paper we describe visual attention experiments in order to investigate just how visual perception can help to produce perceptually high quality 3D urban models to be displayed on mobile devices.

## 2. Related Work

How we perceive an environment depends on who we are and the task that we are currently performing in that environment.[17] Visual attention is the process by which we humans select a portion of the available visual information for localization, identification and understanding of objects in the environment. It allows our visual system to process visual input preferentially by shifting attention about an image, giving more attention to salient locations and less attention to unimportant regions. Although our eyes are good, they are not perfect, and so when attention is not focused onto items in a scene they can literally go unnoticed. This is known as inattentional blindness.[11]

In recent years, knowledge of the human visual system has been increasingly used to improve the quality of the displayed image, for example, Refs. 4, 5, 12–14. Other research has investigated how complex details in the models can be reduced without any reduction in the viewer's perception of the models, for example, Ref. 9. While for visual navigation systems, Maciel and Shirley used texture mapped primitives to represent clusters of objects to maintain high and approximately constant frame rates.[10] In addition, saliency models have been developed to simulate where people focus their involuntary attention in images[7] and these have been used to reduce overall computation costs.[18]

Task maps, on the other hand, specify what parts of a scene are perceptually important to perform a task. These can be defined either manually[3] or done by predicting nondiffuse objects as important.[6]

Cater *et al.*[3] showed that conspicuous objects in a computer graphics scene that would normally attract the viewer's attention are ignored if they are not relevant to the task at hand. In their experiments, viewers were presented with two animations. One was a full, high-quality rendering, while in the other, only the pixels in visual angle of the fovea (2°) centered around the location of a task within the environment were rendered at high quality. This high quality was blended to a much lower quality in the rest of the image. They showed that when observers were performing the task within an animation, their visual attention was fixed exclusively on the area of the task, and they consistently failed to notice the significant difference in rendering quality between the two animation.

In addition, research done in peripheral vision showed that it is possible to reduce details in the periphery without disturbing visual processing. Watson[16] evaluated the effectiveness of high detail insets in head-mounted displays. The high detail inset they used was rectangular and it was presented at the fine level of resolution. Their results showed that although observers found their search targets faster and more accurately in a full high resolution environment, this condition was not significantly better than the high-resolution inset displays with either medium or low peripheral resolutions. Loschky[8] used an eye-linked, multiple resolution display that produces high visual resolution only in the region to which the eyes are directed. Their work has shown that the image needs to be updated after an eye saccade within 5 milliseconds of a fixation, otherwise the observer will detect the change in resolution.

## 3. Experiments

On a mobile device, due to limited processing power and bandwidth, we simply cannot render realistic images in real time to produce high quality 3D interactive maps. Nevertheless we can render the more important objects of a scene in higher quality with selective rendering.

In order to accomplish this we have to first find out which are these important objects. To determine what key features are important when someone is trying to orientate himself in an urban environment, a case study was carried out, which is described in the next section. A further study was executed to investigate whether subjects fail to notice a quality difference when all key features used to get a good location and orientation of where a photograph was taken are of high quality and the rest of the scene is rendered of low quality.

### 3.1. *Experiments number 1*

In order to gather information about the key features that would be important in a scene, a case study was carried out[1] to determine which elements of an urban scenario are more important to people when they were trying to orientate themselves. The task chosen was for the subjects to identify the correct spot from where a specific photograph had been taken.

### 3.1.1. *Methodology*

This study had two phases. In each phase eight subjects were tested by being given three photographs in order to identify three different locations, Figs. 1–3. The photographs had a VGA resolution and a dimension of $10\,\mathrm{cm} \times 7\,\mathrm{cm}$, and were all from the city center of Vila Real, Portugal. Note that the subjects used in first phase were different from the second phase. In the first phase the subjects were given the original photograph and asked to identify the exact spot and direction from where the photograph was taken. While completing the task, each subject was also asked to mark, in order of priority, the objects in the urban scenario that helped him/her to locate the spot.

The subjects were driven to the locations where the photographs had been taken from a different direction to that of the photographs. The time it took for the subject to identify the spot was measured and the subject was asked to specify the number of objects that had been used to identify the place where the photograph had been taken. In the second phase the same procedure was followed with one difference;



(a)                                      (b)

Fig. 1.   Location 1 with key features identified by one subject. (a) Unmodified photograph, (b) Modified photograph.



(a)                                      (b)

Fig. 2.   Location two with key features identified by one subject. (a) Unmodified photograph, (b) Modified photograph.

(a)                                          (b)

Fig. 3.  Location three with key features identified by one subject. (a) Unmodified photograph, (b) Modified photograph.

the photograph that was given to the subject had been altered. Some of the key features identified in the first phase were changed or removed in the photograph with an image editor [Fig. 2(b)], or, in the case of the red and blue flag in Fig. 3, the actual item had been removed (in this case by the shop) when phase 2 of the experiment was conducted.

Again a group of eight (different) subjects was given the three modified photographs and asked to identify the exact location from where the photograph had been taken. The details of the task were well explained to the subjects prior to them undertaking the experiment and they were also shown an example picture to ensure that they clearly understood the instructions.

### 3.1.2. *Results*

The time spent, and the number of objects used to identify the correct spot by each subject in the two phases are shown in the next two tables.

As we can see, the average time spent trying to reach the right spot is very similar in phases one and two. This result was not as expected, as we thought that by removing the key features identified in the first phase, the subjects would require significantly more time to find the right spot. However, with the key features missing, the subjects simply used other features within the environment to help with their orientation. Of note, however, is that in the first phase, all the subjects

Table 1.  Results of the first phase, the average time spent in seconds, the standard deviation, and the average number of items that subjects used to identify each photograph.

| First Phase | Time Average (in seconds) | Standard Deviation | Average No. of Items Used |
|---|---|---|---|
| 1st photo | 70.75 | 26.18 | 2.88 |
| 2nd photo | 70.63 | 25.12 | 3.13 |
| 3rd photo | 62.13 | 10.83 | 3.13 |

Table 2.   Results of the second phase, the average time
spent in seconds, the standard deviation, and the average
of items that subjects used to identify each photograph.

| Second Phase | Time Average (in seconds) | Standard Deviation | Average No. of Items Used |
|---|---|---|---|
| 1st photo | 70.13 | 24.39 | 2.75 |
| 2nd photo | 60.75 | 22.39 | 3.25 |
| 3rd photo | 84.25 | 18.57 | 3.75 |

reached the right spot with an error of about a 1 m, while in the second phase, this
error increased to approximately 3 m.

The elements in the photographs were classified in one of the following
categories:

- Urban furniture — Lamps, seats, bins, etc.
- Buildings — All buildings[1] characteristics, such as doors, windows, balconies,
  geometry, volume and even the whole building.
- Publicity — All types of advertisements.
- Other — Cars, trees, temporary elements.

The elements used by the subjects to find the right spot are shown in
Figs. 4 and 5.

As the results show, when most of the key features identified in the first phase
were changed or removed, the subjects identified other key features. For example, in
the second photograph, Fig. 2, all the subjects identified a total of nine elements of



Fig. 4.   Items identified by the subjects in the first phase of the experiment.

Fig. 5.   Items identified by the subjects in the second phase of the experiment.

urban furniture, Fig. 4. When the urban furniture was removed, Fig. 2, the number
of urban furniture items identified dropped to 1 and at the same time the number of
building elements increased from 12 to 20. We can also see that the total number of
elements identified to find the right spot is about the same (23 elements in original
photograph and 22 in the changed photograph). This was similar for the other two
photographs.

Publicity seems to be an important key feature as subjects always noticed it
when it was present. (This is perhaps not surprising as it is precisely what was
hoped for.) We also noticed that the geometry of the features, for example, the size
of the buildings, was always used by the subjects to gain an approximate orientation,
and subsequently, when the subjects wanted to have a more precise idea of the spot
where the photograph was taken, the key features, in particular urban furniture
and publicity, were used. The minimum number of elements used to identify any
spot was two and the maximum eight.

Another useful result from this study is that the subjects always chose the
closest elements to identify the correct spot. Only if they did not recognize any
key feature near them did they look further afield to try to identify another key
element.

Figure 6 show the saliency maps for photograph 1 before and after it has been
modified using the Itti and Koch's model.[7] Highly salient objects are in white
and less salient ones in increasing dark gray. In the unmodified photograph, the
fountains were a key feature used in the orientation and yet they do not appear as
a high saliency object in the saliency map, Fig. 6 (left). Similarly in Figs. 7 and 8,
although some of the features identified by the users do appear as salient in the Itti
and Koch model, a significant number do not.

Fig. 6.   Saliency maps for photograph 1 (left) unmodified, (right) modified.



Fig. 7.   Saliency maps for photograph 2 (left) unmodified, (right) modified.



Fig. 8.   Saliency maps for photograph 3 (left) unmodified, (right) modified.

## 3.2.  *Experiments number 2*

This experiment, built on the work done in the first experiment, was to investigate whether, if all key features used to reach a good location and orientation of where a photograph was taken are high quality and the rest of the scene is rendered low quality, the subjects failed to notice this quality difference.[2]

### 3.2.1.  *Methodology*

The mobile device chosen for the experiment was the PALM T5 which has a display resolution of $320 \times 480$ pixels. We chose the same location for our experiment as used in first experiment so we could directly compare the results.

Fig. 9.    High quality image of second location.



Fig. 10.    Low quality image of second location.

Three images were prepared for experiment:

- High quality image (HQ): This image was obtained from the original photograph which was resized for the PDA resolution, Fig. 9.
- Low quality image (LQ): This image had only 25% of the quality of the HQ image, Fig. 10. A pilot study had been run to show that for an LQ image below 50% of the quality of the HQ image, when simply viewing the images, sub-jects could easily determine the difference between the HQ and LQ images on the PDA.

Fig. 11.   Selective quality image of second location.

- Selective quality image (SQ): This image was created from the HQ and LQ images with all the key features in HQ and the rest of the image in LQ. The percentage of high quality in the SQ images never exceeded 30 % of the image, Fig. 11.

As in the first experiment, the task we chose for the subjects to perform was to correctly identify the spot and direction from where a photograph was taken. In our experiment, the subjects were told that they had to perform the experiment in the shortest time possible. The person who identified the correct spot and direction in the fastest time won a DVD player.

A total of 40 subjects participated in the experiment. These 40 subjects were divided in two groups: a main group composed of 24 subjects, and a control group of 16 subjects. The main group performed the experiment using the SQ image and the time to do the experiment and the error in their results were measured. After performing the task, each subject was asked to choose between two images as to which one they thought they had seen when performing the experiment. The two images, HQ and LQ, were shown at the same time using two PDAs.

The same procedure was used with control group, but this time half of the group performed the task with the LQ image and the other half with the HQ image. After they had performed the task, each control group subject was asked to choose which image they thought they had seen. Again both the LQ and HQ images were shown at the same time using two PDAs.

The subjects were always driven to the locations where the photographs had been taken, but started the experiment from a different direction to that of the photographs.

### 3.2.2. *Results experiment number 2*

The following tables show the results for the main and control groups. From Tables 3–5, we can see that the performance of all subjects was very close despite the fact they performed the experiment with different quality images. The average time to perform the experiment, when compared with the results presented from the first experience, shows a significant reduction. For example, the main group took on average 52 sec for the first location and 31.91 sec for the second location compared with 70.75 sec and 70.63 sec shown in the first experiment. It could well be that the prize we offered significantly improved the committment of subjects when performing our task.

The results from the control group show that the LQ and HQ images were easily distinguishable even when performing the experiment, with the HQ image being correctly identified 100% of the time.

From Table 5, we can see in the first location 14 subjects, 58.3%, chose the HQ image as the one they thought they saw and 10, 31.6%, LQ. For the second image, 10 of the 23, 43.47%, chose HQ, while 13.56.52%, selected LQ.

We analyzed the results using the chi-square test. For them to have significance, the value of $p$ should be 0.05 (less than 5% chance of random occurrence). For a $df = 1$ (*df* is related to the number of subjects) we obtained a value of $p = 0.41$ for

Table 3.   Results of the main group that performed the task using the SQ image.

|  | No. of Subjects | Time Average (in seconds) | % of HQ | % of LQ |
|---|---|---|---|---|
| Location 1 | 24 | 52.63 | 58.33 | 41.67 |
| Location 2 | 23 | 31.91 | 43.48 | 56.52 |

Table 4.   Results of the control group that performed the task using the HQ image.

|  | No. of Subjects | Time Average (in seconds) | % of HQ | % of LQ |
|---|---|---|---|---|
| Location 1 | 8 | 54.00 | 100 | 0 |
| Location 2 | 8 | 36.38 | 100 | 0 |

Table 5.   Results of the control group that performed the task using the LQ image.

|  | No. of Subjects | Time Average (in seconds) | % of HQ | % of LQ |
|---|---|---|---|---|
| Location 1 | 8 | 37.75 | 0 | 100 |
| Location 2 | 8 | 25.88 | 25 | 75 |

Fig. 12.   Saliency maps, from top SQ/LQ/HQ.

the first location and a $p = 0.53$ for the second location. These values thus show there was no significant difference in our results.

Saliency maps were created from each image (SQ/LQ/HQ) to see if image degradation influences the salient parts of an image. As we can see in Fig. 12, the most salient parts (in white) remain constant, despite the resolution of the images.

One other key problem which we identified during our experiments was the poor clarity of the picture on the PDA display, especially when the experiments were conducted in sunshine.

## 4.  Discussion

The first experiment showed that it may be possible to use saliency maps to predict points of interest in a photograph when someone is trying to orientate himself, since some key features used by subjects are presented on the saliency map. The absent key features on the saliency maps can be completed maybe with the use of task maps. With an ordered list of the key features, from the most to the less important, it should be possible to produce these task maps automatically.

From the results of the first experiment we can only see a few key features are needed to accomplish a good orientation. We expected, all key features were of high quality if people had the same experience as if all the image features were of high quality. But, as the second experiment shows, users performing a orientation task cannot distinguish between SQ and HQ/LQ images (Fig. 13). The SQ image presented to subjects had a very low percentage of the image rendered in high quality (less than 30%). If the amount of high quality increases we believe that subjects will tend more to select the high quality image when asked which image they saw when performing the experiment. This high quality should, of course, be

Fig. 13.   Number of subjects that, when asked, choose HQ or LQ image for locations 1 and 2.

increased around the key feature objects. Further experiments are needed to confirm this hypothesis, and if it is true, then the results can be used to save significant computational costs when rendering 3D virtual environments. To accomplish this, we will simply need an ordered list of key features that can be used by subjects to perform an orientation task.

In addition, our perceptual approach can also be used in conjunction with other techniques, such as impostors,[15] which lower geometrical detail of certain features in a scene. Using our perceptual criteria to determine the key objects of a large scene, such as a city, only the geometry of these need to be rendered in detail, the less perceptually important geometry can be shown using imposters, with a significantly lower computational cost. Furthermore, time will be saved in the actual modeling process as modelers will only need to actually include details for objects which may be perceptually important.

Where the 3D map application relies on a distributed architecture then streaming could be a better option. Streaming will allow the key features to be delivered first and then the other features later in lower resolution.

## 5. Conclusions

The results from the first experiment have shown that a user does not need all the features of an environment to obtain a good orientation. In fact, only the overall geometry and a few key features are required, and these features are not necessarily those most salient to the human visual system. Other features in the image are all but ignored. Furthermore, the time taken to orientate oneself does not increase when key features are absent, but the accuracy of the orientation does.

The selection of features used in the orientation task is mostly driven by proximity and size. Certainly the nearest or biggest features are the most important. For example, a big factory chimney or a nearby traffic light will probably be chosen as key features by the user. There is much more work that needs to be done

before these results can be incorporated into a perceptually driven selective renderer for navigating on mobile devices. The creation of an ordered list of features for rendering is the key issue for this approach.

Results from the second experiment show that subjects do not perform any better whilst orientating themselves in an urban environment when using a high quality image or a selective quality image. The absence of significance, from the results of the main group experiments, indicates that they failed to recognize the difference between the SQ and the HQ/LQ images, whereas the results, from the control group, show that subjects can constantly distinguish between the LQ and HQ images. Our hypothesis was that if all the key features in the SQ image, needed for a correct and rapid identification of the spot in the orientation of the photograph, were of high quality, a significant number of people would choose the HQ image when shown HQ and LQ. We assumed this would be because when they were performing the task they would focus their attention only on these high quality key features. As we can see in Fig. 11, although the key feature, for example, the lamp post, is of high quality, when a user looks at this feature, the foveal angle also includes some of the low quality background.

Future work will consider making the entire foveal angle around key features of high quality with the rest of the image of low quality. In addition, instead of only photographs we will consider a detailed virtual model of the locations, with level of detailed techniques providing the high and low quality differences.

Other issues that will also be considered are: how different lighting and weather conditions may affect the perception of the scene and whether the method of orientation is different for men and women. The empirical data which will be gathered will then be used to modify the Itti and Koch model [7] to develop an urban navigation saliency map based on the GPS position of the specific user and the prevalent weather to ensure the user can orientate himself/herself rapidly to a high precision and from there navigate efficiently through the unfamiliar urban environment.

## Acknowledgments

## References

1. M. Bessa, A. Coelho and A. Chalmers, Alternate feature location for rapid navigation using a 3d map on a mobile device, in *MUM '04: Proc. 3rd Int. Conf. Mobile and Ubiquitous Multimedia*, New York, NY, USA (ACM Press, 2004), pp. 5–9.
2. M. Bessa, A. Coelho and A. Chalmers, Selective rendering quality for an efficient navigational aid in virtual urban environments on mobile platforms, in *MUM '05: Proc. 4th Int. Conf. Mobile and Ubiquitous Multimedia* (2005).

3. K. Cater, A. Chalmers and P. Ledda, Selective quality rendering by exploiting human inattentional blindness: looking but not seeing, in *VRST '02: Proc. ACM Symp. Virtual Reality Software and Technology*, New York, NY, USA (ACM Press, 2004), pp. 17–24.

4. J. Ferwerda and P. *et al.*, A model of visual adaptation for realistic image synthesis, in *Proc. SIGGRAPH 1996* (ACM, 1996), pp. 249–258.

5. D. Greenberg, K. Torrance, P. Shirley, J. Arvo, J. Ferwerda, S. Pattanaik, A. Lafortune, B. Walter, S. Foo and B. Trumbore, A framework for realistic image synthesis, in *Proc. SIGGRAPH 1997 (Special Session)* (ACM, 1997), pp. 477–494.

6. J. Haber, K. Myszkowski, H. Yamauchi and H.-P. Seidel, Perceptually guided corrective splatting, in *Computer Graphics Forum*, eds. A. Chalmers and T.-M. Rhyne, Vol. 20, Manchester, UK (September 2001), (Eurographics, Blackwell), pp. C142–C152.

7. L. Itti and C. Koch, A saliency-based search mechanism for overt and covert shifts of visual attention, *Vis. Res.* **40**(10–12) (2000) 1489–1506.

8. L. C. Loschky, G. W. McConkie, J. Yang and M. E. Miller, Perceptual effects of a gaze-contingent multi-resolution display based on a model of visual sensitivity, in *Advanced Displays and Interactive Displays Fifth Ann. Symp. 2001*, pp. 53–58.

9. D. Lubeke and B. Hallen, Perceptually driven simplification for interactive rendering, in *Proc. 12th Eurographics Workshop on Rendering*, (Eurographics, 2001), pp. 221–223.

10. P. Maciel and P. Shirley, Visual navigation of large environments using textured clusters, in *Proc Symp. Interactive 3D Graphics* (1995), pp. 95–102.

11. A. Mack and I. Rock, *Inattentional Blindness* (MIT Press, 1998).

12. A. Mcnamara, A. Chalmers, T. Troscianko and I. Gilchrist, Comparing real and synthetic scenes using human judgements of lightness, in *12th Eurographics Workshop on Rendering* (2000), pp. 207–219.

13. K. Myszkowski, T. Tawara, H. Akamine and H. Seidel, Perception-guided global illumination solution for animation rendering, in *Proc. SIGGRAPH 2001* (ACM, 2001), pp. 221–230.

14. M. Ramasubramanian, S. Pattanaik and D. Greenberg, A perceptually based physical error metric for realistic image synthesis, in *Proc. SIGGRAPH 1999* (ACM, 1999), pp. 73–82.

15. F. Sillion, G. Drettakis and B. Bodelet, Efficient impostor manipulation for real-time visualization of urban scenery, in *Computer Graphics Forum (Proc. Eurographics '97)*, eds. D. Fellner and L. Szirmay-Kalos, Vol. 16 (September 1997), pp. 207–218.

16. B. Watson, N. Walker, L. F. Hodges and M. Reddy, Satiotemporal sensitivity and visual attention for efficient rendering of dynamic environments, *ACM Trans. Comput. Graph.* **20**(1) (2001) 39–65.

17. A. Yarbus, Eye movements during perception of complex objects, *Eye Movements and Vision*, Chap. VII:171–196 (1967).

18. H. Yee, S. Pattanaik and D. Greenberg, Spatiotemporal sensitivity and visual attention for efficient rendering of dynamic environments, *ACM Trans. Comput. Graph.* **20**(1) (2001) 39–65.

**Maximino Bessa** received a degree in electrical engineering from the University of Tras-os-Montes e Alto Douro, Portugal in 2003. Currently he is a Ph.D. student at the University of Tras-os-Montes e Alto Douro, Portugal.

**Jose Bulas-Cruz** received the bachelor degree in electric engineering from the University of Porto (Portugal), Ph.D. in electrical engineering from the University of Bristol (UK). Currently is full Professor at the University of Tras-os-Montes e Alto Douro, Portugal, where he is also vice-chancellor for technology and innovation.

**Antonio Coelho** is Assistant Lecturer at the University of Tras-os-Montes e Alto Douro, Portugal, since 1996. He is currently conducting a doctoral research at the University of Porto, Portugal, on the topic of "Expeditious modelling of three-dimensional scenes for virtual reality systems".

**Alan Chalmers** received a M.Sc. from Rhodes University, South Africa in 1984 and a Ph.D. from the University of Bristol in 1991. Currently he is a Professor of computer graphics at the University of Bristol, UK.