

Challenges of Visually Realistic Augmented Reality

Claus B. Madsen ^a

Computer Graphics Group, Aalborg University, Rendsburggade 14, Aalborg, Denmark

Keywords: Augmented Reality, Illumination, Realism, Shadows.

Abstract: Why is achieving real-time handheld visually realistic Augmented Reality so hard? What are the main challenges? We present an overview of these challenges, and discuss the most important issues involved in developing AR that automatically adapts to changes in the environment, specifically the illumination conditions. We then move on to present how we see a path of research going forward for the immediate future; a path based partly on recent advances in real-time 3D modelling and partly on lessons learned from a decade of Augmented Reality illumination estimation research.

1 INTRODUCTION

Augmented Reality (AR) is gaining new momentum with the advent of easy-to-use APIs such as ARKit, ARCore, and Vuforia. With these APIs it has never been easier to develop robust handheld AR experiences. If the purpose of an AR application is to create the illusion that a virtual object is situated, and realistically visualized in the context of an actual physical scene, then there are 3 main challenges to address: 1) tracking/registering the camera relative to the scene, 2) handling occlusions between real and virtual objects, and 3) render the augmented objects with illumination conditions that are consistent with the real scene; (Azuma et al., 2001). The aforementioned APIs primarily address the tracking/registration problem. The occlusion problem will not really find a solution until AR devices are capable of providing an accurate scene depth value for each pixel in the camera feed, and hence current AR applications basically assume that augmented objects are positioned on an uncluttered flat surface, or floating in the air. The last of the three problems, the illumination consistency, has actually been the subject of extensive research, but so far no really elegant, actually functional solutions exist.

In this paper we will not address the issues of tracking any further. We will also not be spending much energy on the handling of occlusions. The main focus of the paper is to list and discuss the main challenges one faces when trying to develop visually realistic AR applications, attempt to give an overview of

possible research directions that present themselves in the area. In the end we describe an approach we believe to be relevant for research in the immediate future.



Figure 1: A screen shot from the *IKEA Place* app used in an outdoor scenario in broad daylight. The arm chair is an augmentation. The app renders the augmented object with a generic "contact shadow", which is better than not having *any* but in most cases will be inconsistent with the real shadows. In this case the augmented shadow is inconsistent in terms of direction, depth and softness.

For the purposes of the discussion in this paper we will assume that one ultimate goal of AR research is to enable AR applications that are so visually convincing, that the augmentations cannot be distinguished from reality. That is, we ultimately desire real-time interactive AR which is as visually convincing as an A-level Hollywood blockbuster movie. Additionally, we have chosen to limit the paper to only discussing issues related to so-called video see-through AR, as opposed to optical see-through AR. In video see-through AR the real world is experienced as a video of the real world on a display. In optical see-through, the real work is sensed directly through

^a  <https://orcid.org/0000-0003-0762-3713>

a transparent, prism-like optical device which allows augmentations to be optically mixed with the real-world view (as in Microsoft Holodeck headsets, for example).

Sadly, current AR is nowhere near that level. Consider Figure 1. This example is made with the *IKEA Place* app and is a good representative of current state of the art commercial AR. This particular example mainly focuses on the realism of the shadows cast by the augmented object, but obviously all other rendering related effects would be equally relevant, e.g., reflections, refractions, shading, color bleeding, occlusions, depth-of-field, motion blur, etc.

These are all important standard issues for *any* kind of rendering application; but AR is by definition a special kind of rendering application scenario, as part of the final scene is real and part is virtual/rendered, and hence the perceived realism of the virtual part will naturally be held against, and has to be visually consistent with, the real part. In fact, the main challenge in realistic AR can be boiled down to this: how much do we know about the real part of the scene? How good are our models of the geometry, materials and illumination conditions in the real scene?

2 ASPECTS OF AR RENDERING

Rendering requires models of the scene geometry, light/matter models (reflectances) of the materials in the scene, and models of the scene illumination. Let us subsequently address these aspects one by one, focusing on the what is particularly challenging in the case of AR. In this section we will go through these in a little more detail, followed by a brief description of other rendering challenges faced when desiring to achieve realistic AR. A great in-depth literature review on mixed-reality rendering is provided in (Kronander et al., 2015). Understanding the concept of *differential rendering* as presented in (Debevec, 1998) is also a good starting point for understanding the challenges involved. Below we focus on giving a brief overview of aspects relevant for real-time AR.

2.1 Scene Geometry

In the AR literature it is generally accepted that it makes sense to divide the scene geometry into three separate classes, (Debevec, 1998; Debevec, 2002):

1. Augmentations
2. Local scene
3. Distant scene

Augmentations represent the 3D models of the elements to be augmented into the scene. The local scene is the part of the scene that has essential radiometric interactions with augmented objects, e.g., for shadows and occlusions. And lastly, the distant scene is a term referring to all real elements of the scene for which we do not absolutely need a detailed representation, and perhaps simply an image-based representation will have to suffice. Examples of these classes are shown in Figure 2.



Figure 2: Top: the sculpture constitutes the augmentation class; the local scene made up of a model of the cave wall on the left to receive the cast shadow; the rest of the cave is the distant scene and only represented by an image. Rendered with a real-time AR application, (Madsen and Laursen, 2007). Bottom: the pyramid, the teapot, the cylinder and the text are the augmentations; the local scene is the supporting plane and a cylinder modelling the trunk of the tree; the rest of the image is the distant scene. Rendered with Autodesk 3ds Max.

Given this breakdown of scene elements it is obvious that the geometry of the augmentations can be as complex as the application requires, while taking into account that the rendering framework has to be able to handle the complexity. Similarly, if very intricate handling of occlusion between augmented and local scene elements is required, the main challenge is obtain sufficiently detailed geometry models for the local scene. In Figure 2 a small tree trunk is modelled by a simple cylinder to handle the occlusion between the tree and the pyramid. In actual real-time AR appli-

cations the realistic handling of occlusion remains one of the absolute main challenges, and in many cases failure to properly handle this will totally break the illusion.

2.2 Scene Materials

The materials for the augmented geometry in the scene can naturally be as complex as the application requires provided the materials can be rendered on the rendering platform used. Figure 3 shows an example combining refraction and specular reflection mapping. As the distant scene is typically only represented as image information in AR rendering (basically as a backdrop), the main challenge from a materials point of view is the local scene elements. Looking carefully at the teapot in the lower part of Figure 2 it can be seen that the teapot reflects the texture of the ground plane. This has been achieved by projectively mapping the scene image onto the ground plane to get geometrically correct reflections. Hence, the more ambitious one is regarding radiometric interaction between the augmented part of the scene and the local scene, e.g., reflections (augmented objects reflecting local scene elements, and vice versa), color bleeding, etc., the more accurate the models of local scene materials have to be.



Figure 3: Example of floating sculpture augmented into scene with real-time refraction and reflection mapping.

2.3 Scene Illumination

Rendering the augmented objects using illumination conditions that are consistent with the real scene has been proven to be essential for achieving realistic AR. In Figure 2 the top example was rendering with an image-based representation of the illumination conditions, which we will return to in section 3. The bottom part was rendered with completely manually tweaked illumination conditions (a single directional

light source representing the sun, and a hemispherical dome representing the sky). For an outdoor scenario such a simple illumination model can be sufficient, but the problem of establishing a good model of indoor illumination conditions can be really complex if there are many potentially differently shaped artificial light sources, and even contributions of exterior lighting through windows. For realistic rendering the placement, geometry, color and intensity (RGB radiances) are needed.

A challenge in AR is to correctly compute the appearances of shadows cast by augmented objects onto local scene objects, i.e., adjusting pixels in the image of the real scene so as to appear to be in shadow, if they were not in shadow had the augmented object not been there. The correct way to do this for a certain point in the scene is to compute the ratio of irradiance received *with* the augmentation in the scene, to the irradiance received *without* the augmentation present in the scene. These irradiances are computed using the model of the scene illumination.

Related to this it is a challenge in AR to avoid "double shadow", i.e., avoid casting virtual shadow in an area of the real scene that is already in shadow from the same light source. Notice for example in Figure 2 how we have carefully ensured that the vertical cylinder casts its shadow into where a tree outside the field of view is already casting a shadow. Avoiding double shadows is a huge challenge in realistic AR, as it will require tremendously accurate local scene geometry models (potentially even for objects outside the field of view), and similarly accurate illumination models, i.e., extremely precisely located light sources, that might also very likely be outside the field of view. A recent paper, (Wei et al., 2019), presents a really interesting approach to handling this problem based on shadow edges within the field-of-view.

2.4 Post Processing

This section briefly describes a few other important elements of achieving realistic AR,- elements that perhaps not quite as often thought about in current research, (Borg et al., 2014).

When rendering some computer graphics elements into a real image it is important to subject the rendered elements to a level of imaging noise that matches the imaging noise in the real image. Studies have proven this to be important. Especially in low level lighting conditions small digital cameras, such as smart phones, have significant levels of imaging noise, and augmentations should be subject to similar levels of noise to avoid the absence of noise being conspicuous.

The silhouette of the augmented objects also needs to be handled with care so as to not appear fake and "too perfect" compared to the rest of the image. Aliasing artifacts around the silhouette should be handled and some kind of gentle blending of the augmentation with the real image background is necessary. Related to this, it will also be necessary for future realistic AR applications to render augmentation with depth-of-field blurring that is consistent with the scene and the settings of the camera used.

Motion blur is another important element that in principle needs to be addressed.

3 ESTABLISHING THE REQUIRED MODELS

The previous section introduced the most important elements and challenges in achieving realistic AR. This section aims at given a brief overview of current approaches to obtaining the necessary models and dealing with the challenges. Since the problem of easy to use, general purpose, perfectly realistic AR is by no means solved, there is no end-all, be-all approach we can describe. Every single relevant approach that we might describe will be based on a number of assumptions and delimitations.

Before we start, we reiterate that the ultimate goal in the paper is considered to be handheld (or with a head-mounted display) video see-through, real-time Augmented Reality, allowing augmentations to be rendered into *dynamic* (geometrically and illumination-wise) scenes in a visually convincing manner.

3.1 Scene Geometry

As previously described, obtaining the geometry information of the local scene is a major challenge. Most current AR, e.g., *IKEA Place* assumes the local scene only consist of the planar surface detected by APIs such as ARCore or ARKit. In AR research the local scene is sometimes assumed to have been somehow modelled in advance, and hence assumed static, e.g., Figure 2.1 top.

Figure 4 shows an example where a real-time binocular stereo camera has been used to obtain a snapshot of a dynamic local scene. The geometry resolution provided by the stereo camera is sufficient for a realistic 3D model of the main elements of the local scene, but not quite precise enough to correctly model the geometry of the person.

We believe the trend regarding local scene modelling for handheld devices to be moving towards real-



Figure 4: Augmentation rendered into scene with local scene obtained from a real-time stereo camera, (Madsen and Lal, 2013).

time 3D modelling exploiting multiple cameras in the device and combining it with real-time Structure-from-Motion (SfM) techniques, such as the 6D.ai API. These approaches will be able to generate persistent models that extend outside any instantaneous field-of-view by accumulating geometry information over time as the user moves around and points the device in different directions. Future APIs may even very well be able employ cloud services where multiple users feed scanned geometry into a globally persistent model that others users can use and help update.

We believe SfM-based models will be useful for serving as local scene models in AR applications, but we highly doubt that these models will be good enough to convincingly handle occlusions between e.g., humans and augmentations. For this purpose it is much more likely that AI based approaches to segmenting humans from video, like BodyPix 2.0, will prevail for near future AR applications.

An issue regarding scene modelling is that the border between what is local scene and what can be treated distant scene is dependent on the scale of the augmentation. If augmenting a fire hydrant on a side-

walk, then the sidewalk may be the only necessary local scene, and nearby buildings can be considered distant scene. If augmenting a building into a scene, naturally neighbouring buildings suddenly become important local scene elements.

3.2 Scene Materials

Work on *inverse global illumination*, e.g., (Yu et al., 1999), demonstrates how complex it is to perform proper non-diffuse-only material model estimation in mixed reality scenes based, as it requires highly detailed 3D models of surfaces and multiple views of each point on those surfaces. As of yet, this type of approach is not relevant for AR application attempting to achieve real-time performance. In most (all?) AR research the material of local scene elements is considered to be perfectly diffuse and the appearances are lifted directly from the image/video of the scene. In some research the appearances are then used to estimate the diffuse albedos of the surfaces by taking into account knowledge of the illumination in the scene, (Madsen and Laursen, 2007). For future research in realistic AR we believe more research needs to go into estimating surface reflectances from multiple views as the user is moving the device around and looking in different directions. With good tracking it will be possible to integrate these different appearance measurements into an estimated reflectance model, perhaps by having a small set of categories of material, e.g., diffuse such as concrete and brick wall, glossy such as lacquered surfaces, and highly specular surfaces such as windows and puddles.

3.3 Scene Illumination

Some AR research aiming at achieving visual realism has adopted the concept of light probes originally developed for movie productions, see Figure 5, (Debevec, 2002; Kanbara and Yokoya, 2004; Jacobs and Loscos, 2004). The idea being that the illumination conditions at the location where something has to be augmented are captured in a omni-directional image. This image then has to be in High Dynamic Range format, as it requires many orders of magnitude dynamic range to capture illumination conditions where for example the sun is 5 orders of magnitude more bright than the sky, (Dutre et al., 2002).

Clearly it is not possible to use a static light probe for AR in a dynamic scene. The ARCore API is now able to accumulate a light probe image of the scene as the users moves around, but the image is not in a HDR format. ARCore also offers functionality for estimating the direction a dominant directional light source



Figure 5: Latitude-longitude mapping of omni-directional HDR light probe image applicable for image-based lighting.

in scenes, (ARCore, 2020), although it remains to be evaluated how well this novel functionality works in various scenarios.

Related research is aimed at developing Machine Learning approaches to "guesstimate" illumination conditions from a video frame, simply by training on video sequences and associated videos of how real diffuse and specular spheres look like in that scene, (LeGendre et al., 2019; Hold-Geoffroy et al., 2016). We are currently looking into how to classify weather conditions from video images so as to be able to use various outdoor daylight models for outdoor AR.

Other research has applied a more model based approach and tried to estimate the illumination conditions from automatically detected shadows in real scenes, (Madsen and Lal, 2013). Figure 4 was an example of this. In a recent paper it has been demonstrated how that approach could be adopted to run in real-time on a handheld device for very simple scenes, Figure 6, (Bertolini and Madsen, 2020). Other related work in this area, demonstrating how the various sensors on mobile devices can be utilized is (Barreira et al., 2018).

We believe the current trend for illumination estimation for AR to be moving in the direction of using machine learning to estimate simplified illumination parameters for augmented reality simply by employing massive amounts of training data. We also believe that this approach might very well make it possible to have AR objects shaded in a manner which is generally consistent with the real scene. Nevertheless, based on more than a decade of experience with working in this area we still maintain that visual elements such as directions, depths and softness of shadows, color bleeding etc. will require elements of more model-based approaches.

4 WHAT WILL WE DO?

As a concrete direction for immediate future research in the area of realistic AR we propose to focus on handheld AR for outdoor urban environments, ad-



Figure 6: Realtime tracking of areas in real sun and real shadow, and using these areas for estimating the radiances of the sun and the sky, for rendering a virtual apple into a scene on an iPhone.

addressing the challenges of continuously, and in real-time, model the dynamic illumination conditions.

4.1 Geometry

In terms of geometric models for the local and distant scene elements we propose to investigate real-time local 3D models acquired by something similar to the aforementioned 6Dai technique, see Figure 7, and how to fuse them with larger scale city models similar to what can be found via, e.g., Google Earth, i.e., previously scanned and textured models of streets and buildings. We believe there will be an advantage in utilizing a mix of these two sources of geometric and appearance information for the scene.

Such a fused model will have high detail near the device, and be up to date with objects such as cars etc., and more coarse detail further away from the device. The very fact that this will enable geometric information about the part of the scene which is outside the current field of view, a part that might never have been with the field of view, will be very important for robust illumination estimation.



Figure 7: Realtime generated 3D mesh of local scene captured from a smartphone.

For detailed occlusion handling, for example partial occlusion of augmentations behind people, trees, etc., we believe this will need to be handled as a mixture between depth-based approaches and 2D image-based approaches. The depth information will come from essentially having RGBd data (per pixel depth) available on the device from multi-ocular stereoscopy combined with Structure-from-Motion, and the really detailed occlusion handling, e.g., silhouettes of people, will be assisted by image-based segmentation. In fact, the whole issue of building 3D models of the environment for AR purposes will most likely in the near future benefit highly from recent advances in deep learning approaches to scene segmentation, recognition, and scene understanding.

4.2 Materials

As mentioned previously the immediate future in AR will probably treat local scene materials as perfectly diffuse. The estimation of illumination parameters will be based on this assumption, and the rendering of augmentations into the video stream will be based on this assumption. Nevertheless, we believe that there is a huge potential in doing more research into estimating surface material properties from multiple observations, assuming the AR application user is moving his/her device around and filming areas from multiple viewpoints. This will allow for rough classification into a small set of surface material types, for example diffuse/glossy/specular. Such estimation can obviously also benefit from deep learning based material classification, so that the AR application would be able to identify what elements in the scene are e.g., windows, puddles, polished tiles, etc. and use this to create realism enhancing effects such as augmented objects reflecting in glossy surfaces.

4.3 Illumination

In terms of estimating and modelling the illumination conditions there still are two main competing

approaches: 1) a machine learning based approach where scene consistent illumination is estimated and applied to augmentations, as championed for example in (LeGendre et al., 2019), or 2) model-based approaches where some sort of parametric illumination model is tuned to the scene based on extraction of various properties from the image stream, e.g., (Bertolini and Madsen, 2020).

For outdoor daytime AR we propose to employ an adaptive daylight model; a daylight model that uses the time, data, compass reading and geo-location to compute the direction vector to the sun. The local scene model would be used to classify which areas of the scene should be in shadow (if the sun is even actually shining). The adaptive part of the daylight models should be that it is tweaked to adapt to the actual conditions in terms of the weather; e.g., whether it is a clear blue sky day, partly overcast, completely overcast, or rainy. Or even if there is snow. This adaptation we believe is possible through machine learning approaches based on monitoring the video feed on the device, and there are already examples of work in this area, e.g., (Lu et al., 2017).

The actual estimation of the sun and sky radiances we would in the direction of fusing existing shadow based approaches with an inverse rendering inspired approach comparing the current appearance of surfaces in the local scene model with their appearance as stored in the cloud model (Google Earth). We might want to look into doing laser range finder based capture of huge point clouds for the streets and buildings, while simultaneously capturing the corresponding illumination conditions with omnidirectional HDR cameras. This would enable computation of surface reflectances. These stored models and reflectances could then, at run-time on the handheld device make it realistic to estimate the illumination conditions at that particular time. The viability of such an approach was tentatively demonstrated in for example (Jensen et al., 2006).

We believe this combination of 1) streamed, previously acquired, static models, and 2) run-time acquired additional geometry and illumination estimation, offers a realistic promise of easy to use, real-time handheld AR which can run on off-the-shelf current smartphones. For outdoor scenarios, that is. Indoor scenarios are still much more complicated from an illumination point of view. The only comfort we have is that initial perceptual experiments are indicating that human tolerance to imperfections in illumination correctness is higher for indoor scenarios. Probably because it is more difficult to judge what actually looks correct, as long as the rendered augmentations are largely consistent with the real scene.

5 CONCLUSIONS

In this paper we have attempted to give an overview of the primary challenges involved in developing realistic AR on handheld devices, which can dynamically adapt the changing illumination conditions.

We fundamentally believe a lot more work is required on perceptual studies into how tolerant humans are to various aspects of imperfections in visual quality of AR. That said, we have proposed what we believe to be the best path for future research. A path that involves mixing geometry capture on the device using Structure-from-Motion techniques with streamed, pre-captured gross models of the environment. Dynamically adaptive illumination estimation would then be based on inverse rendering techniques by comparing real-time scene appearance with stored scene reflectances combined with a parametric daylight model.

One day in the future it will be possible to hunt visually convincing augmented dinosaurs in the streets, -that's the dream!

ACKNOWLEDGEMENTS

This work was partially funded by the LER project no. EUDP 2015-I under the Danish national EUDP programme, and partially by the DARWIN project under the Innovation Fund Denmark, case number: 6151-00020B. This funding is gratefully acknowledged. The author would also like to take this opportunity to thank colleagues and students, past and present, for inspiration.

REFERENCES

- ARCore (2020). Using arcore to light models in a scene. <https://developers.google.com/ar/develop/unity/light-estimation>. Accessed: January 7th, 2020.
- Azuma, R. T., Baillot, Y., Behringer, R., Feiner, S., Julier, S., and MacIntyre, B. (2001). Recent advances in augmented reality. *IEEE Transactions on Computer Graphics and Applications*, 21(6):34 – 47.
- Barreira, J., Bessa, M., Barbosa, L., and Magalhaes, L. (2018). A context-aware method for authentically simulating outdoors shadows for mobile augmented reality. *IEEE Transactions on Visualization and Computer Graphics*, 24(3):1223–1231.
- Bertolini, F. and Madsen, C. B. (2020). Real time outdoor light estimation for mobile augmented reality. In *Proceedings: International Conference on Graphics Theory and Applications*. Accepted.

- Borg, M., Paprocki, M., and Madsen, C. (2014). Perceptual evaluation of photo-realism in real-time 3d augmented reality. In *Proceedings of GRAPP 2014: International Conference on Computer Graphics Theory and Applications*, pages 377–386. Institute for Systems and Technologies of Information, Control and Communication.
- Debevec, P. (1998). Rendering synthetic objects into real scenes: Bridging traditional and image-based graphics with global illumination and high dynamic range photography. In *Proceedings: SIGGRAPH 1998, Orlando, Florida, USA*.
- Debevec, P. (2002). Tutorial: Image-based lighting. *IEEE Computer Graphics and Applications*, pages 26 – 34.
- Dutre, P., Bala, K., and Bekaert, P. (2002). *Advanced Global Illumination*. A. K. Peters, Ltd., Natick, MA, USA.
- Hold-Geoffroy, Y., Sunkavalli, K., Hadap, S., Gambaretto, E., and Lalonde, J.-F. (2016). Deep outdoor illumination estimation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2373–2382.
- Jacobs, K. and Loscos, C. (2004). State of the art report on classification of illumination methods for mixed reality. In *EUROGRAPHICS*, Grenoble, France.
- Jensen, T., Andersen, M., and Madsen, C. B. (2006). Estimation of dynamic light changes in outdoor scenes without the use of calibration objects. In *Proceedings: International Conference on Pattern Recognition, Hong Kong*, page (4 pages).
- Kanbara, M. and Yokoya, N. (2004). Real-time estimation of light source environment for photorealistic augmented reality. In *Proceedings of the 17th ICPR, Cambridge, United Kingdom*, pages 911–914.
- Kronander, J., Banterle, F., Gardner, A., Miandji, E., and Unger, J. (2015). Photorealistic rendering of mixed reality scenes. *Computer Graphics Forum*, 34(2):643–665.
- LeGendre, C., Ma, W.-C., Fyffe, G., Flynn, J., Charbonnel, L., Busch, J., and Debevec, P. (2019). Deeplight: Learning illumination for unconstrained mobile mixed reality. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lu, C., Lin, D., Jia, J., and Tang, C. (2017). Two-class weather classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2510–2524.
- Madsen, C. and Lal, B. (2013). Estimating outdoor illumination conditions based on detection of dynamic shadows. In Csurka, G., Kraus, M., Mestetskiy, L., Richard, P., and Braz, J., editors, *Computer Vision, Imaging and Computer Graphics*, pages 33–52. Springer Publishing Company, United States.
- Madsen, C. B. and Laursen, R. (2007). A scalable gpu-based approach to shading and shadowing for photo-realistic real-time augmented reality. In *Proceedings: International Conference on Graphics Theory and Applications, Barcelona, Spain*, pages 252 – 261.
- Wei, H., Liu, Y., Xing, G., Zhang, Y., and Huang, W. (2019). Simulating shadow interactions for outdoor augmented reality with rgbd data. *IEEE Access*, 7:75292–75304.
- Yu, Y., Debevec, P., Malik, J., and Hawkins, T. (1999). Inverse global illumination: Recovering reflectance models of real scenes from photographs. In *Proceedings of the 26th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '99*, page 215–224, USA. ACM Press/Addison-Wesley Publishing Co.