

# Fusing Multimedia Data Into Dynamic Virtual Environments

Ruofei Du  
ruofei@cs.umd.edu

Committee: Dr. Varshney, Dr. Zwicker, and Dr. Huang

# Motivation

Popularity of VR and AR devices



# Motivation

Popularity of VR and AR devices



# Motivation

Popularity of VR and AR devices



# Motivation

Popularity of VR and AR devices



# Motivation

Popularity of VR and AR devices



# Motivation

Assorted VR and AR applications



# Motivation

Assorted VR and AR applications





# Motivation

Assorted VR and AR applications



These VR/AR applications have

Huge  
Data

requirements

These VR/AR applications have

# Huge Data

requirements



*Where is the  
3D data going  
to come from?*

# Motivation

Lack of contents



# Motivation

Lack of contents



These VR/AR applications have

# Huge Data

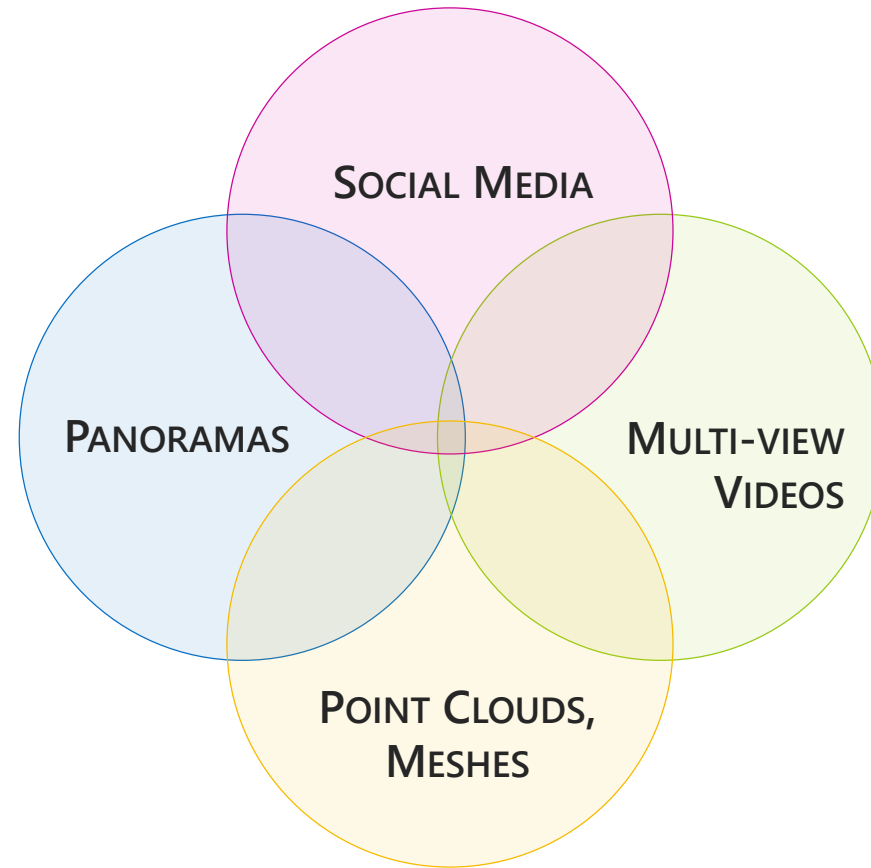
requirements



*Automatically  
fusing multimedia  
data into virtual  
environments*

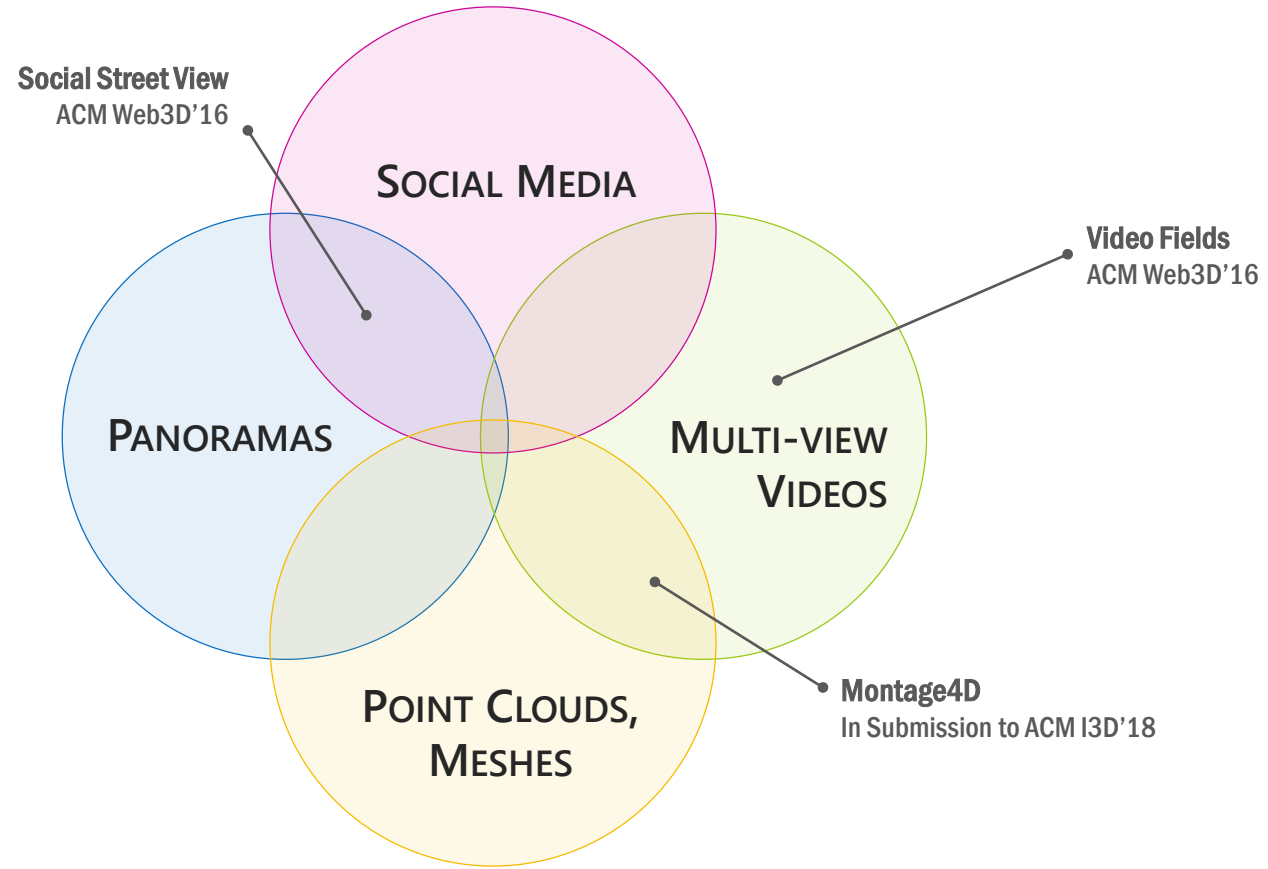
# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



# Proposal

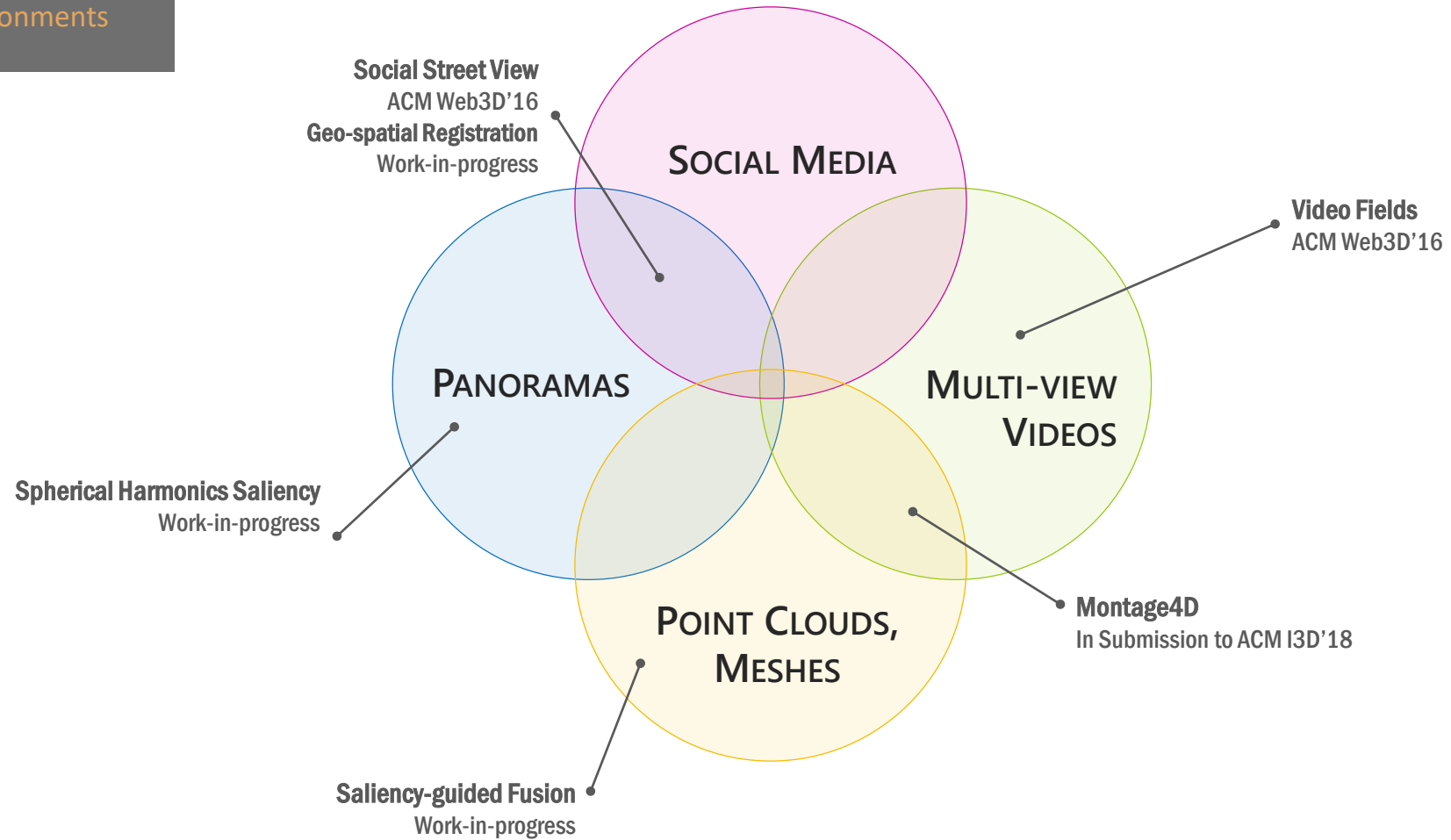
Fusing Multimedia Data Into  
Dynamic Virtual Environments





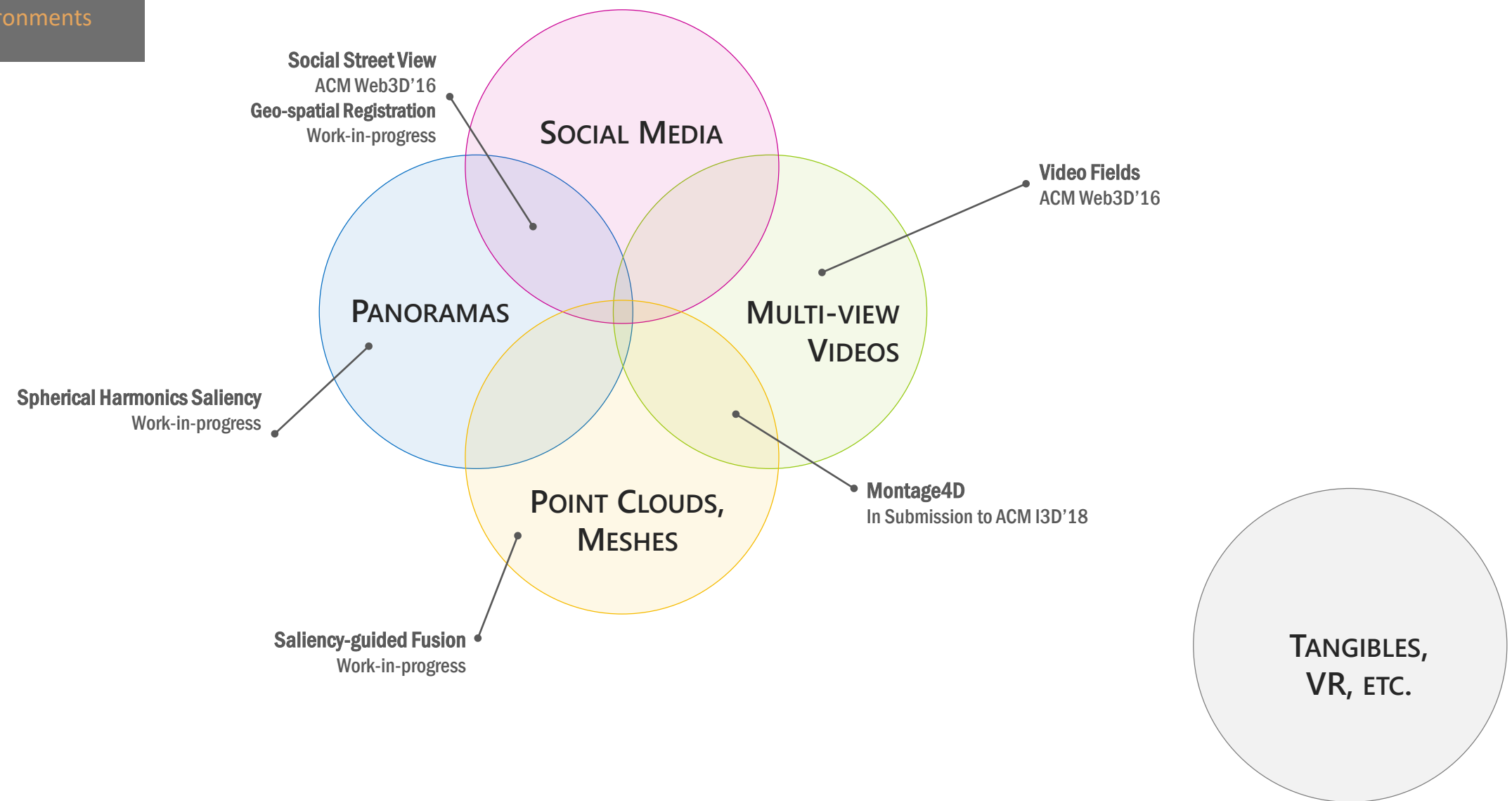
# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



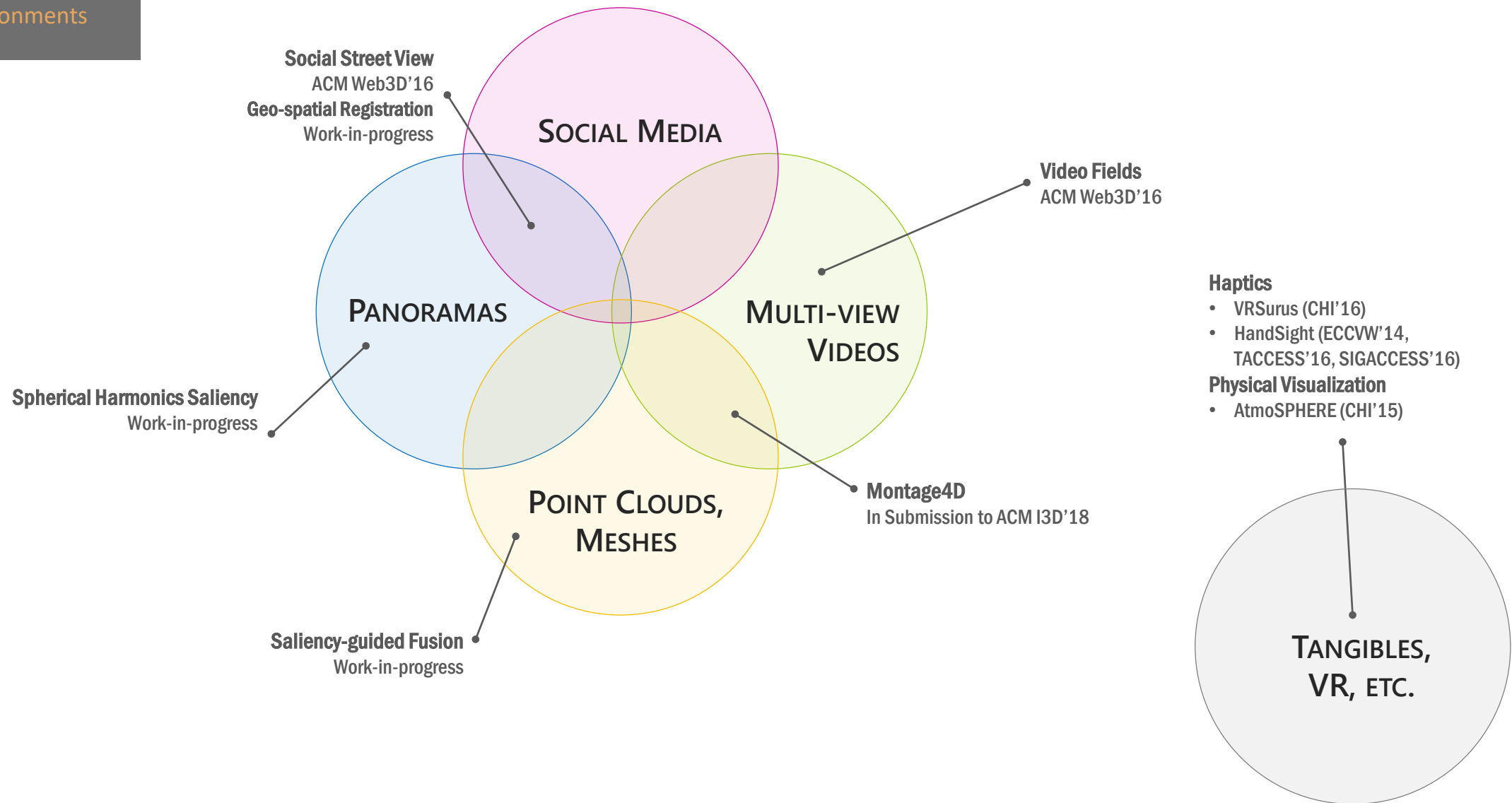
# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



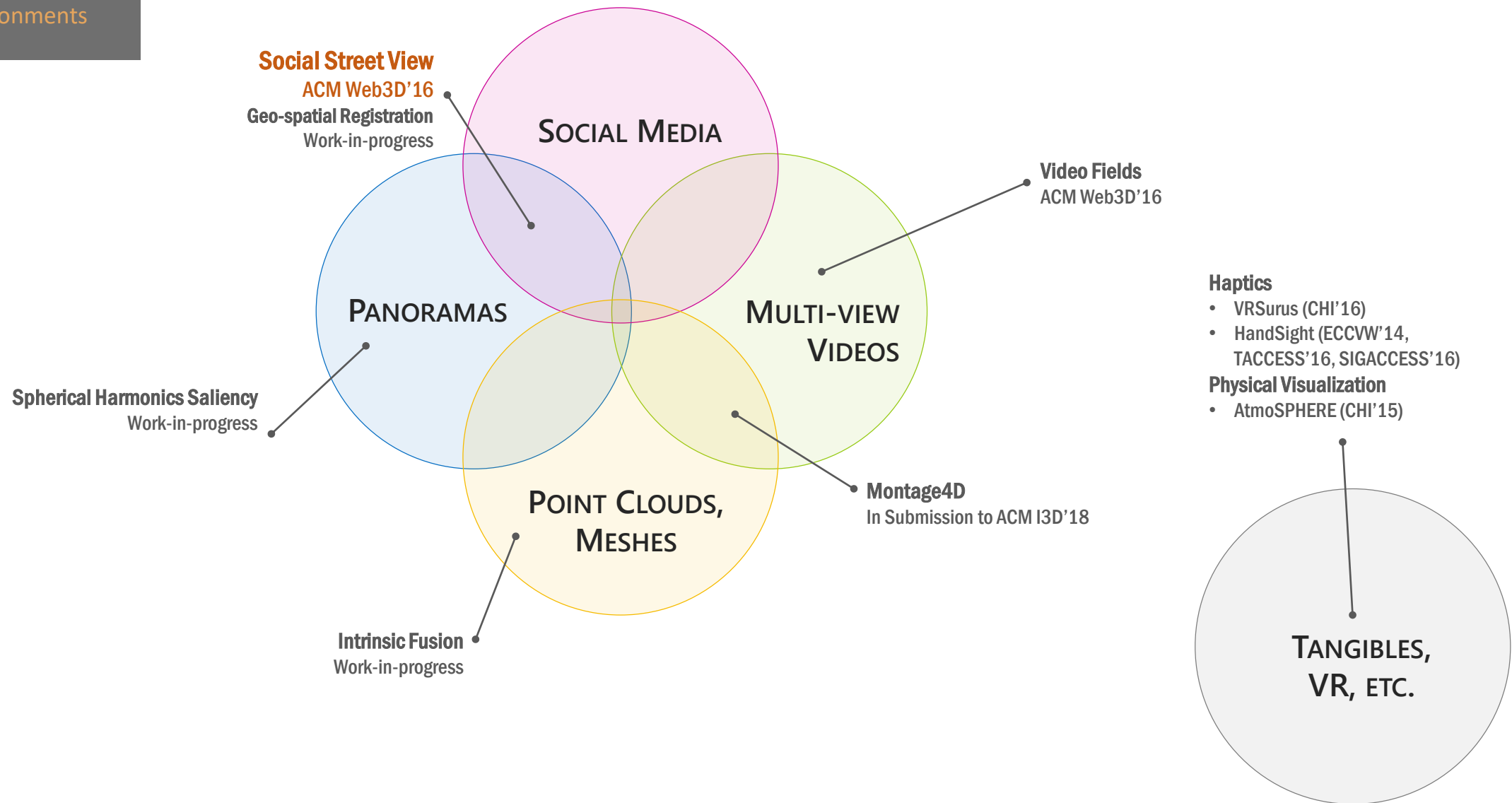
# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



# Social Street View: Blending Immersive Street Views with Geo-tagged Social Media

Ruofei Du and Amitabh Varshney

{ruofei, varshney} @ cs.umd.edu

[www.SocialStreetView.com](http://www.SocialStreetView.com)

In Proceedings of the International Conference of Web3D 2016. (Best Paper Award)

# Introduction

Social Media



But why is social media **so popular**?





Introduction  
Social Media - Versatility



image courtesy: forbes.com



**Metadata** are useful for understanding **spatial relevance**, **relationship amongst users**, **sentiment mining**, and **propagation of influence**.



# Introduction

Meta data - Time of Creation



*image courtesy: Brian Clegg*

# Introduction

Meta data - Location of Creation



image courtesy: squarespace.com

# Introduction

Meta data - Camera Parameters

# Related Work

Linear narrative visualization

The image displays a vertical scroll of social media posts from three different platforms: Weibo (left), Facebook (middle), and Twitter (right). Each platform shows a sequence of posts that tell a story over time. The posts include photos of food, people, and events, along with text updates and interactions like likes and comments. The timeline progresses from top to bottom, showing various moments in a person's life, such as dining, social gatherings, and academic achievements.

image courtesy:  
instagram.com,  
facebook.com,  
twitter.com





# Related Work

Natural Immersive Virtual Reality?



# Related Work

*Karnath et al. and Loomis et al.*

## Related Work

Karnath et al. and Loomis et al.

### Spatial awareness is a function of the temporal not the posterior parietal lobe

Hans-Otto Karnath, Susanne Ferber & Marc Himmelbach

Department of Cognitive Neurology, University of Tübingen, Hoppe-Seyler-Strasse 3, 72076 Tübingen, Germany

Our current understanding of spatial behaviour and parietal lobe function is largely based on the belief that spatial neglect in humans (a lack of awareness of space on the side of the body contralateral to a brain injury) is typically associated with lesions of the posterior parietal lobe. However, in monkeys, this disorder is observed after lesions of the superior temporal cortex<sup>1</sup>, a puzzling discrepancy between the species. Here we show that, contrary to the widely accepted view, the superior temporal cortex is the neural substrate of spatial neglect in humans and monkeys. Unlike the monkey brain, spatial awareness in humans is a function largely confined to the right superior temporal cortex, a location topographically reminiscent of that for language on the left<sup>2</sup>. Hence, the decisive phylogenetic transition from monkey to human brain seems to be a restriction of a formerly bilateral function to the right side, rather than a shift from the temporal to the parietal lobe. One may speculate that this lateralization of spatial awareness parallels the emergence of an elaborate representation for language on the left side.

Spatial neglect is a characteristic failure to explore the side of space contralateral to a brain lesion. Patients with this disorder behave as if one side of the surrounding space had ceased to exist. Since the early post-mortem studies, we have believed that, in humans, lesions located predominantly in the posterior parietal lobe are critical for this disorder. Analyses of computerized tomography scans of right-hemispheric stroke patients with neglect found that superimposed lateral projections of these lesions centred on the inferior parietal lobule (IPL)<sup>3,4</sup> and the temporo-parieto-occipital (TPO) junction<sup>4</sup>. More recent studies have confirmed the validity of this conclusion although evidence for additional pathology leading

## Related Work

Karnath et al. and Loomis et al.

### Spatial awareness is a function of the temporal not the posterior parietal lobe

Hans-Otto Karnath, Susanne Ferber & Marc Himmelbach

Department of Cognitive Neurology, University of Tübingen, Hoppe-Seyler-Strasse 3, 72076 Tübingen, Germany

Our current understanding of spatial behaviour and parietal lobe function is largely based on the belief that spatial neglect in humans (a lack of awareness of space on the side of the body contralateral to a brain injury) is typically associated with lesions of the posterior parietal lobe. However, in monkeys, this disorder is observed after lesions of the superior temporal cortex<sup>1</sup>, a puzzling discrepancy between the species. Here we show that, contrary to the widely accepted view, the superior temporal cortex is the neural substrate of spatial neglect in humans and monkeys. Unlike the monkey brain, spatial awareness in humans is a function largely confined to the right superior temporal cortex, a location topographically reminiscent of that for language on the left<sup>2</sup>. Hence, the decisive phylogenetic transition from monkey to human brain seems to be a restriction of a formerly bilateral function to the right side, rather than a shift from the temporal to the parietal lobe. One may speculate that this lateralization of spatial awareness parallels the emergence of an elaborate representation for language on the left side.

Spatial neglect is a characteristic failure to explore the side space contralateral to a brain lesion. Patients with this disorder behave as if one side of the surrounding space had ceased to exist. Since the early post-mortem studies, we have believed that in humans, lesions located predominantly in the posterior parietal lobe are critical for this disorder. Analyses of computerized topography scans of right-hemispheric stroke patients with neglect that superimposed lateral projections of these lesions centred on the inferior parietal lobule (IPL)<sup>3,4</sup> and the temporo-parieto-occipital (TPO) junction<sup>4</sup>. More recent studies have confirmed the validity of this conclusion although evidence for additional pathology

### Immersive virtual environment technology as a basic research tool in psychology

JACK M. LOOMIS and JAMES J. BLASCOVICH  
University of California, Santa Barbara, California

and

ANDREW C. BEALL  
Massachusetts Institute of Technology, Cambridge, Massachusetts

Immersive virtual environment (IVE) technology has great promise as a tool for basic experimental research in psychology. IVE technology gives participants the experience of being surrounded by the computer-synthesized environment. We begin with a discussion of the various devices needed to implement immersive virtual environments, including object manipulation and social interaction. We review the benefits and drawbacks associated with virtual environment technology, in comparison with more conventional ways of doing basic experimental research. We then consider a variety of examples of research using IVE technology in the areas of perception, spatial cognition, and social interaction.

Human history records a progression of artifacts for representing and recreating aspects of external reality, ranging from language, drawings, and sculpture in earlier times to the more modern artifacts of photographs, movies, television, and audio recordings. Relatively recently, the digital computer and its associated technologies, including three-dimensional (3-D) graphics, have given rise to increasingly realistic artifacts that blur the distinction between reality and its representation (Ellis, 1995).

The ultimate representational system would allow the observer to interact "naturally" with objects and other individuals within a simulated environment or "world," an experience indistinguishable from "normal reality." Although such a representational system might conceivably use direct brain stimulation in the future, it will more likely use digitally controlled displays that stimulate the human sensory organs, the natural conduits to the brain.

Displays of this type, referred to as *virtual displays* (VDs), although far from ideal, exist today. Following the terminology of others (e.g., Durlach & Mavor, 1995; Stanley & Salvendy, 1998), we refer to the corresponding environment represented and stored in the computer and experienced by the user as a *virtual environment* (VE). *Virtual environment technology* (VET) refers inclusively both to VDs and to the VEs so created, including VEs produced by using conventional desktop computer displays.

(*Virtual reality* is widely used as an alternative term, but we prefer VE.) An immersive virtual environment (IVE) is one in which the user is perceptually surrounded by the VE. Ivan Sutherland (1965), one of the originators of 3-D computer graphics, was the first person to conceive and build an immersive VD system. For the history of IVEs, see Ellis (1995), Kalawsky (1993), and Rheingold (1991).

There are two usual implementations of an IVE. The first of these involves placing multiple projection screens and loudspeakers around the user. A popular design is the CAVE (Cruz-Neira, Sandin, & DeFanti, 1993), which involves back-projecting the computer-generated visual imagery onto the translucent walls, floor, and ceiling of a moderately sized cubical room, in which the user is free to move; shutter glasses provide stereoscopic stimulation, so that one sees the VE not as projections on the room surfaces, but as solid 3-D structures within and/or outside of the cube. The second and more common implementation of an IVE involves the use of a head-mounted display (HMD), used in conjunction with a computer and a head tracker (Barfield & Furness, 1995; Biocca & Delaney, 1995; Burdea & Coiffett, 1994; Durlach & Mavor, 1995; Kalawsky, 1993). The head tracker measures the changing position and orientation of the user's head within the physical environment, information that is communicated to the rendering computer, which has stored within it a 3-D representation of the simulated environment (Meyer, Applewhite, & Biocca, 1992). At any given moment, the computer generates and outputs the visual and auditory imagery to the user's HMD from a perspective that is based on the position and orientation of the user's head. The HMD consists of earphones and video displays attached to a support worn on the head; the video display component is based on cathode ray tube (CRT) displays, liquid crystal displays, or laser-based retinal scanners (Barfield, Hendrix, Bjorneseth, Kaczmarek, &

With the support of Grant N00014-95-1-0573 from the Office of Naval Research (to J.M.L.) and National Science Foundation Grants SBR 9872084 (to J.J.B.) and SBR 9873432 (to J.M.L. and J.J.B.), and the authors have developed several immersive virtual displays, use of which has stimulated many of the ideas expressed here. The authors thank Florence Gaunet and Patrick Pèruch for comments on an earlier version of the article. Correspondence concerning this article should be addressed to J. M. Loomis, Department of Psychology, University of California, Santa Barbara, CA 93106 (e-mail: loomis@psych.ucsb.edu).

# Related Work

Karnath et al. and Loomis et al.

## Spatial awareness is a function of the temporal not the posterior parietal lobe

Hans-Otto Karnath, Susanne Ferber & Marc Himmelbach

Department of Cognitive Neurology, University of Tübingen, Hoppe-Seyley-Strasse 3, 72076 Tübingen, Germany

Our current understanding of spatial behaviour and parietal lobe function is largely based on the belief that spatial neglect in humans (a lack of awareness of space on the side of the body contralateral to a brain injury) is typically associated with lesions of the posterior parietal lobe. However, in monkeys, this disorder is observed after lesions of the superior temporal cortex<sup>1</sup>, a puzzling discrepancy that has been noted in other species. Here we show that, contrary to the widely accepted view, the superior temporal cortex is not a system of spatial neglect in humans, but that spatial awareness in humans is a function largely confined to the right inferior parietal cortex, a location topographically reminiscent of that for language on the left<sup>2</sup>. Hence, the decisive phylogenetic transition from monkey to human brain seems to be a restriction of a formerly bilateral function to the right side, rather than a shift from the temporal to the parietal lobe. One may speculate that this lateralization of spatial awareness parallels the emergence of an elaborate representation for language on the left side.

Spatial neglect is a characteristic failure to explore the side space contralateral to a brain lesion. Patients with this disorder behave as if one side of the surrounding space had ceased to exist. Since the early post-mortem studies, we have believed that humans, lesions located predominantly in the posterior parietal lobe are critical for this disorder. Analyses of computerized topography scans of right-hemispheric stroke patients with neglect that superimposed lateral projections of these lesions centred on the inferior parietal lobule (IPL)<sup>3,4</sup> and the temporo-parieto-occipital (TPO) junction<sup>4</sup>. More recent studies have confirmed the validity of this conclusion although evidence for additional pathology

## Immersive virtual environment technology as a basic research tool in psychology

JACK M. LOOMIS and JAMES J. BLASCOVICH  
University of California, Santa Barbara, California

and  
ANDREW C. BEALL  
Massachusetts Institute of Technology, Cambridge, Massachusetts

Immersive virtual environment (IVE) technology has great promise as a tool for basic experimental research in psychology. IVE technology gives participants the experience of being surrounded by the computer-synthesized environment. We begin with a discussion of the various devices needed to implement immersive virtual environments, including object manipulation and social interaction. We review the benefits and drawbacks associated with virtual environment technology, in comparison with more conventional ways of doing basic experimental research. We then consider a variety of examples of research using IVE technology in the areas of perception, spatial cognition, and social interaction.

Human history records a progression of artifacts for representing and recreating aspects of external reality, ranging from language, drawings, and sculpture in earlier times to the more modern artifacts of photographs, movies, television, and audio recordings. Relatively recently, the digital computer and its associated technologies, including three-dimensional (3-D) graphics, have provided a new medium for representing external reality. The ultimate representational system would allow an observer to interact with a virtual world that is an experience indistinguishable from "normal reality." Although such a representational system might conceivably use direct brain stimulation in the future, it will more likely use digitally controlled displays that stimulate the human sensory organs, the natural conduits to the brain.

Displays of this type, referred to as *virtual displays* (VDs), although far from ideal, exist today. Following the terminology of others (e.g., Durlach & Mavor, 1995; Stanney & Salvendy, 1998), we refer to the corresponding environment represented and stored in the computer and experienced by the user as a *virtual environment* (VE). *Virtual environment technology* (VET) refers inclusively both to VDs and to the VEs so created, including VEs produced by using conventional desktop computer displays.

(*Virtual reality* is widely used as an alternative term, but we prefer VE.) An immersive virtual environment (IVE) is one in which the user is perceptually surrounded by the VE. Ivan Sutherland (1965), one of the originators of 3-D computer graphics, was the first person to conceive and build an immersive VD system. For the history of IVEs, see (1995), Kalawsky (1993), and Rheingold (1991). There are two usual implementations of an IVE. The first of these involves the use of a head-mounted display (HMD), which involves back-projecting the computer-generated visual imagery onto the translucent walls, floor, and ceiling of a moderately sized cubical room, in which the user is free to move; shutter glasses provide stereoscopic stimulation, so that one sees the VE not as projections on the room surfaces, but as solid 3-D structures within and/or outside of the cube. The second and more common implementation of an IVE involves the use of a head-mounted display (HMD), used in conjunction with a computer and a head tracker (Barfield & Furness, 1995; Biocca & Delaney, 1995; Burdea & Coiffett, 1994; Durlach & Mavor, 1995; Kalawsky, 1993). The head tracker measures the changing position and orientation of the user's head within the physical environment, information that is communicated to the rendering computer, which has stored within it a 3-D representation of the simulated environment (Meyer, Applewhite, & Biocca, 1992). At any given moment, the computer generates and outputs the visual and auditory imagery to the user's HMD from a perspective that is based on the position and orientation of the user's head. The HMD consists of earphones and video displays attached to a support worn on the head; the video display component is based on cathode ray tube (CRT) displays, liquid crystal displays, or laser-based retinal scanners (Barfield, Hendrix, Bjorneseth, Kaczmarek, &

With the support of Grant N00014-95-1-0573 from the Office of Naval Research (to J.M.L.) and National Science Foundation Grants SBR 9872084 (to J.J.B.) and SBR 9873432 (to J.M.L. and J.J.B.), and the authors have developed several immersive virtual displays, use of Florence Gaunet and Patrick Pèruch for comments on an earlier version of the article. Correspondence concerning this article should be addressed to J. M. Loomis, Department of Psychology, University of California, Santa Barbara, CA 93106 (e-mail: loomis@psych.ucsb.edu).

# Related Work

*Visualization of Geo-tagged Information*

Very little work has been carried out in designing immersive interfaces that interleave **visual navigation of our surroundings** with **social media content**

# Related Work

Visualization of Geo-tagged Information

Jagan Sankaranarayanan  
jagan@cs.umd.edu

Benjamin E. Teitler  
bteitler@cs.umd.edu

Michael D. Lieberman  
codepoet@cs.umd.edu

Hanan Samet  
hjs@cs.umd.edu

Jon Sperling  
Jon.Sperling@hud.gov

## TwitterStand: News in Tweets\*

### ABSTRACT

Twitter is an electronic medium that allows a large user population to communicate with each other simultaneously. Inherent to Twitter is an asymmetrical relationship between friends and followers that provides an interesting social network-like structure among the users of Twitter. Twitter messages, called tweets, are restricted to 140 characters and thus are usually very focused. We investigate the use of Twitter to build a news processing system, called *TwitterStand*, from Twitter tweets. The idea is to capture tweets that correspond to late breaking news. The result is analogous to distributed news wire service. The difference is that the identities of the contributors/reporters are not known in advance and there may be many of them. Furthermore, tweets are not sent according to a schedule: they occur as news is happening, and tend to be noisy while usually arriving at a high throughput rate. Some of the issues addressed include removing the noise, determining tweet clusters of interest bearing in mind that the methods must be online, and determining the relevant locations associated with the tweets.

### Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: Information Storage and Retrieval

### General Terms

Algorithms, Design, Performance

\*This work was supported in part by the National Science Foundation under Grants EIA-08-12377, CCF-08-30618, and HIS-07-13501, as well as NVIDIA Corporation, Microsoft Research, Google, the E.T.S. Walton Visitor Award of the Science Foundation of Ireland, and the National Center for Geocomputation at the National University of Ireland at Maynooth.

†Department of Computer Science, Center for Automation Research, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA.

‡HUD Office of Policy Development & Research (PD&R), 451 7th St. SW, Room 8146, Washington, DC 20410, USA.

### Keywords

Twitter, News, Geotagging, Online clustering

## 1. INTRODUCTION

Twitter<sup>1</sup> is a social networking website that recently has been gaining much attention and following. Twitter is composed of users who send messages (termed *tweets*) to each other, where each tweet contains a maximum of 140 characters. At this time, it is estimated that there are 6 to 7 million users who use Twitter a total of 134 million times a month [4], and this number is increasing at a rapid rate. For example, for the year of 2008, Twitter grew in terms of the number of tweets sent at a rate of 1382% [12] which is a testament to the immense popularity and wide adoption of this service. The popularity of Twitter stems from its availability on a number of different electronic devices (e.g., web, cell phones, etc.), as well as the prevalence of a subculture in Twitter that encourages users to acquire a large friend pool, as well as send tweets on a wide variety of subjects, typically several times a day. The restriction on the lengths of Twitter messages invariably means that the tweets do not necessarily contain well formed ideas, being rather brief, yet complete enough so that users can make sense of the ideas that they convey. Note that tweets also have a mechanism by which the user can link to other objects on the web such as articles, images, videos, etc. (termed *artifacts*) which is typically used to link tweets to related material on the Internet.

The goal of this paper is to demonstrate how to use Twitter to automatically obtain breaking news from the tweets posted by Twitter users, and to provide a map interface for reading this news, since the geographic location of the user as well as the geographic terms comprising the tweets play an important role in *clustering* tweets and establishing clusters' geographic foci. In contrast to news aggregators such as Google News, Bing News, and Yahoo! News, we introduce a system called *TwitterStand* that works exclusively with only the tweets posted by the users of Twitter. The key novelty behind *TwitterStand* is one of mobilizing the millions of users in Twitter to be our eyes and ears in the world, bearing in mind that geographically proximate users often tweet about the same breaking news. In other words, we rely on Twitter users to be either providers of original news content (e.g., the 2008 Southern California earthquake [13] and the 2009 Iranian election [3]), or expressers of opinions on current news topics (i.e., mini blogs), both of which enable *TwitterStand* to automatically identify current news topics and cluster the corresponding tweets into appropriate news stories. We also associate an importance score with each news topic which can

© 2009 Association for Computing Machinery. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or consultant of a firm that has been granted a copyright in this article, or to



# Related Work

TwitterStand: News in Tweets  
Sankaranarayanan et al. SIGGIS 2009



image courtesy: [twitterstand.umiacs.umd.edu](http://twitterstand.umiacs.umd.edu)

# Related Work

## Visualization of Geo-tagged Information

Jagan Sankaranarayanan  
jagan@cs.umd.edu

Benjamin E. Teitler  
bteitler@cs.umd.edu

Michael D. Lieberman  
codepoet@cs.umd.edu

Hanan Samet  
hjs@cs.umd.edu

Jon Sperling  
Jon.Sperling@h

## TwitterStand: News in Tweets\*

### ABSTRACT

Twitter is an electronic medium that allows a large user population to communicate with each other simultaneously. Inherent to Twitter is an asymmetrical relationship between friends and followers that provides an interesting social network-like structure among the users of Twitter. Twitter messages, called tweets, are restricted to 140 characters and thus are usually very focused. We investigate the use of Twitter to build a news processing system, called *TwitterStand*, from Twitter tweets. The idea is to capture tweets that correspond to late breaking news. The result is analogous to a distributed news wire service. The difference is that the identities of the contributors/reporters are not known in advance and there may be many of them. Furthermore, tweets are not sent according to a schedule: they occur as news is happening, and tend to be noisy while usually arriving at a high throughput rate. Some of the issues addressed include removing the noise, determining tweet clusters of interest bearing in mind that the methods must be online, and determining the relevant locations associated with the tweets.

### Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: Information Storage and Retrieval

### General Terms

Algorithms, Design, Performance

\*This work was supported in part by the National Science Foundation under Grants EIA-08-12377, CCF-08-30618, and HIS-07-13501, as well as NVIDIA Corporation, Microsoft Research, Google, the E.T.S. Walton Visitor Award of the Science Foundation of Ireland, and the National Center for Geocomputation at the National University of Ireland at Maynooth.  
†Department of Computer Science, Center for Automation Research, Institute for Advanced Computer Studies, University of Maryland, College Park, MD 20742, USA.  
‡HUD Office of Policy Development & Research (PD&R), 451 7th St. SW, Room 8146, Washington, DC 20410, USA.

### Keywords

Twitter, News, Geotagging, Online

## 1. INTRODUCTION

Twitter<sup>1</sup> is a social network that has been gaining much attention among other, where each tweet contains a message of up to 140 characters. At this time, it is estimated that there are 10 million users who use Twitter, and this number is growing rapidly. For example, for the year of 2008, the number of tweets sent was 1.5 billion, a testament to the immense popularity of this service. The popularity of Twitter is due to its ability to allow users to communicate with each other (via cell phones, etc.), as well as to allow users to send tweets on their mobile devices several times a day. The messages are typically short and contain well formed ideas, and are easy to convey. Note that the user can link to other tweets (e.g., images, videos, etc.) to link tweets to related content. The goal of this paper is to automatically process tweets posted by Twitter users, as well as the geographic information contained in the tweets posted behind Twitter in Twitter to help users to be aware of the same broad topics (i.e., to automatically associate

# Placing Flickr Photos on a Map

Pavel Serdyukov \*†  
Database Group  
University of Twente  
PO Box 217, 7500 AE  
Enschede, The Netherlands  
serdyukovpv@cs.utwente.nl

Vanessa Murdock  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
vmurdock@yahoo-inc.com

Roelof van Zwol  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
roelof@yahoo-inc.com

## ABSTRACT

In this paper we investigate generic methods for placing photos uploaded to Flickr on the World map. As primary input for our methods we use the textual annotations provided by the users to predict the single most probable location where the image was taken. Central to our approach is a language model based entirely on the annotations provided by users. We define extensions to improve over the language model using tag-based smoothing and cell-based smoothing, and leveraging spatial ambiguity. Further we demonstrate how to incorporate GeoNames<sup>1</sup>, a large external database of locations. For varying levels of granularity, we are able to place images on a map with at least twice the precision of the state-of-the-art reported in the literature.

## Categories and Subject Descriptors

H.3.3 [Information Search and Retrieval]

## General Terms

Algorithms, Measurement, Performance, Experimentation

## Keywords

image localisation, language models, Flickr

## 1. INTRODUCTION

Due to the massive production of affordable GPS-enabled cameras and mobile phones [13, 16], location metadata such as *latitude* and *longitude* are automatically associated with the content generated by users. Users have the opportunity to spatially organise and browse their personal media, and photo sharing services are leading the growing enthusiasm for personal location-awareness [22]. Geo-referenced photos

can be organised in a browsable taxonomy of major locations or pin-pointed on a map to identify very small regions. Some of the most popular examples are Flickr Places<sup>2</sup> and Google Panoramio.<sup>3</sup>

While in theory every photo can be anchored to the location it was taken, in practice many photos are location agnostic. Furthermore, the majority of Flickr users do not own location-aware cameras. Thus a large proportion of photos uploaded to Flickr contain no location information even when the photo merits localizing. When uploading photos on Flickr users can still geo-tag their photos by dragging the photos to a particular point on the world map. This process is time-consuming and results in less accurate geo-tagging of photos compared to automatically geo-tagged photos from GPS-enabled cameras. When manually geo-tagging photos, Flickr initially suggests the location of the last uploaded photo or simply displays the world map.

The objective of this paper is to provide a more accurate starting point for geo-tagging photos, uploaded on Flickr, using the textual annotations provided by the user. According to recent literature [2, 21] users spend considerable effort to organise their “memory” geographically by describing photos with *tags* related to locations where they were taken. The location specific tags (such as *Torre Agbar* which is only located in Barcelona), and location related tags (such as *elephants* which are related to locations such as zoos, Africa and Asia) provide essential cues as to where a picture was taken. For photos that are location agnostic (such as *dog*), location information may or may not be provided, but it is normally not relevant to the context of the photo.

The literature related to geo-tagging of photos and its use is extensive. In particular the reverse problem of discovering important landmarks and events, given a geographic co-ordinate has been studied extensively [1, 17, 13]. However the problem of placing images on a map using the textual annotations provided by the user has received less attention. While we focus on Flickr as our primary data source, our approach could be applied to other photo sharing services.

\*Research performed while the author was an intern at Yahoo! Research.  
†Also affiliated with TU Delft, ICT Group  
<sup>1</sup><http://www.geonames.org> visited May 2008

# Related Work

Flickr - World Map  
Serdyukov et al. SIGIR 2009

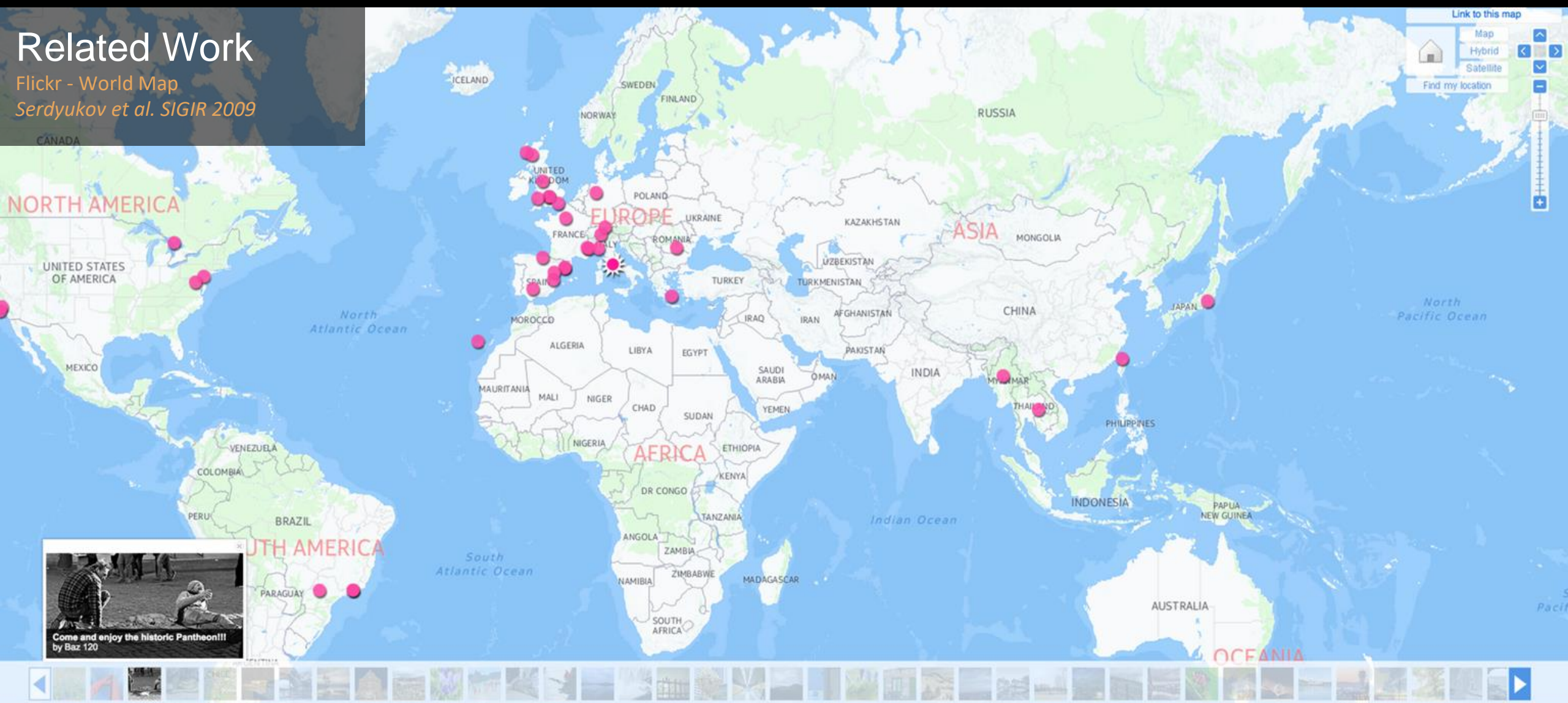


image courtesy: flickr.com

# Related Work

## Visualization of Geo-tagged Information

Jagan Sankaranarayanan  
jagan@cs.umd.edu

Benjamin E. Teitler†  
bteitler@cs.umd.edu

Michael D. Lieberman†  
codepoet@cs.umd.edu

Hanan Samet†  
hjs@cs.umd.edu

Jon Sperling†  
Jon.Sperling@h

### ABSTRACT

Twitter is an electronic medium that allows a large user population to communicate with each other simultaneously. Inherent to Twitter is an asymmetrical relationship between friends and followers that provides an interesting social network-like structure among the users of Twitter. Twitter messages, called tweets, are restricted to 140 characters and thus are usually very focused. We investigate the use of Twitter to build a news processing system, called *TwitterStand*, from Twitter tweets. The idea is to capture tweets that correspond to late breaking news. The result is analogous to

### Keywords

Twitter, News, Geotagging, Online

### 1. INTRODUCTION

Twitter<sup>1</sup> is a social network that has been gaining much attention among a large number of users who send messages to each other, where each tweet contains up to 140 characters. At this time, it is estimated that there are over 10 million users who use Twitter.

CHI 2008 Proceedings · Works in Progress

Category  
H.3 [Information Systems] Storage

General Terms  
Algorithms

\*This work was supported by the National Science Foundation (NSF) under grant IIS-07-11111. The authors would like to thank the anonymous reviewers for their comments.  
†Department of Computer Science, University of Maryland, College Park, MD 20742-4111.  
‡HUD Center for Urban and Environmental Analysis, 7th St.

**Marco Cristani**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
marco.cristani@univr.it

**Alessandro Perina**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
perina@sci.univr.it

**Umberto Castellani**  
Università degli Studi di Verona

**Vittorio Murino**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
Vittorio.murino@univr.it

## Content Visualization and Management of Geo-located Image Databases

### Abstract

In the last years, several algorithms and platforms for photo sharing have been developed. Usually, in order to index huge quantities of images for a fast and intuitive retrieval, additional textual tags attached to the pictures are considered. In this paper, we present a set of solutions for an effective management of geo-located images, i.e. pictures equipped with tags indicating the geographical coordinates of acquisition. This brings towards an intuitive content visualization and management of large geo-located image databases.

### Keywords

Image categorization, geo-located images, interfaces

## Placing Flickr Photos on a Map

Pavel Serdyukov \*†  
Database Group  
University of Twente  
PO Box 217, 7500 AE  
Enschede, The Netherlands  
serdyukovpv@cs.utwente.nl

Vanessa Murdock  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
vmurdock@yahoo-inc.com

Roelof van Zwol  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
roelof@yahoo-inc.com

### ABSTRACT

In this paper we investigate generic methods for placing photos uploaded to Flickr on the World map. As primary input for our methods we use the textual annotations provided by the users to predict the single most probable location where the image was taken. Central to our approach is a language model based entirely on the textual annotations by users. We define a method for placing photos on a map that can be organized in a browsable taxonomy of major locations or pin-pointed on a map to identify very small regions. Some of the most popular examples are Flickr Places<sup>2</sup> and Google Panoramio.<sup>3</sup>

April 5-10, 2008 · Florence, Italy

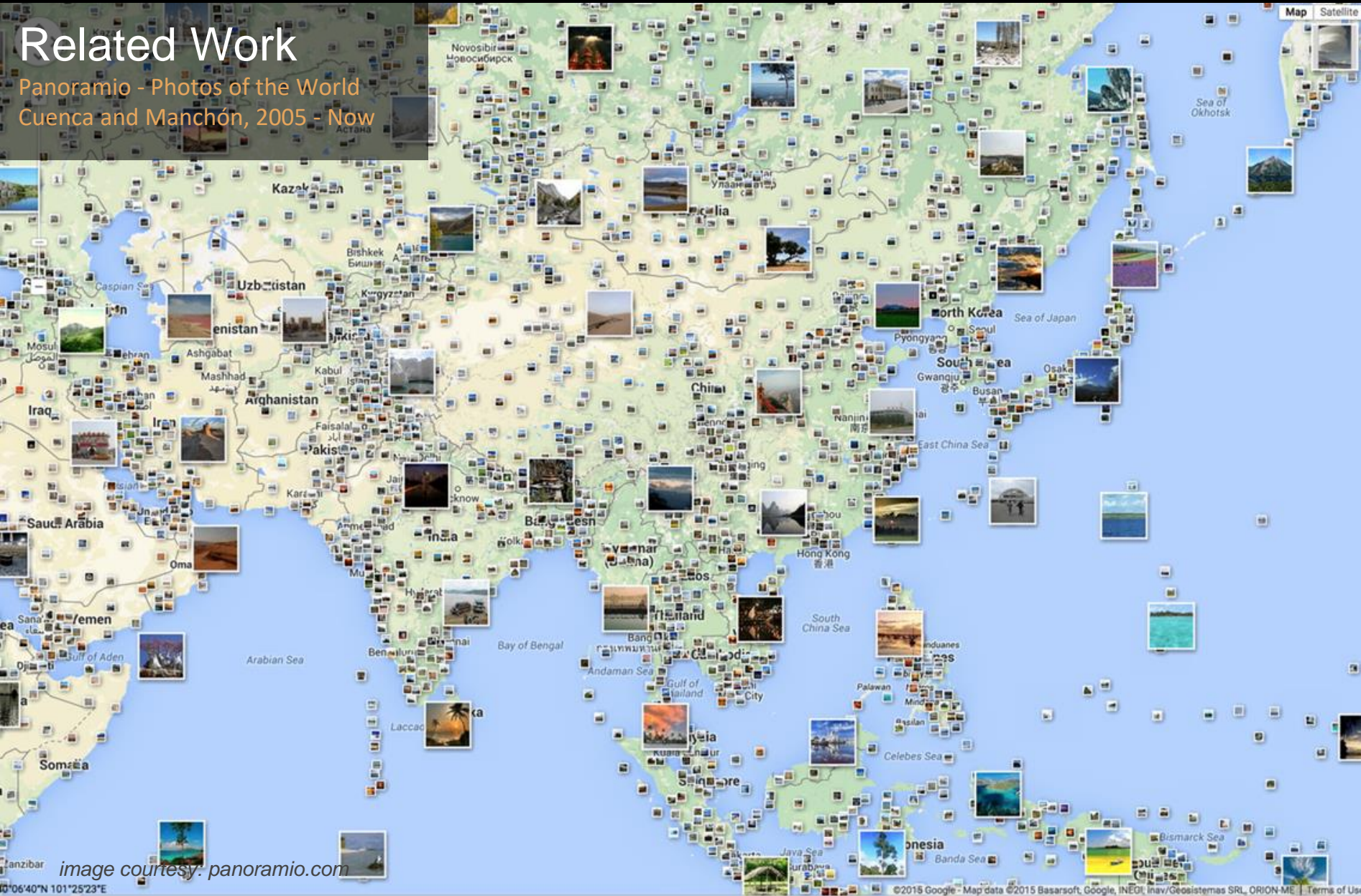
While in theory every photo can be anchored to the location it was taken, in practice many photos are location agnostic. Furthermore, the majority of Flickr users do not own location-aware cameras. Thus a large proportion of photos uploaded to Flickr contain no location information even though the photo merits localizing. When uploading photos to Flickr users can still geo-tag their photos by dragging the photo to a particular point on the world map. This process is time-consuming and results in less accurate geo-tagging of photos compared to automatically geo-tagged photos from location-aware cameras. When manually geo-tagging photos, the user typically suggests the location of the last uploaded photo simply displays the world map.

The objective of this paper is to provide a more accurate method for geo-tagging photos, uploaded on Flickr, using textual annotations provided by the user. According to prior literature [2, 21] users spend considerable effort to geo-tag their “memory” geographically by describing photos with specific tags (such as *Torre Agbar* which is only valid for Barcelona), and location related tags (such as *zoo* which are related to locations such as zoos, Africa) which provide essential cues as to where a picture was taken. Photos that are location agnostic (such as *dog*), where the location may or may not be provided, but it is irrelevant to the context of the photo.

While related to geo-tagging of photos and its use in particular the reverse problem of discovering landmarks and events, given a geographic area, has been studied extensively [1, 17, 13]. However, placing images on a map using the textual annotations by the user has received less attention. Flickr as our primary

# Related Work

Panoramio - Photos of the World  
Cuenca and Manchón, 2005 - Now



Popular Recent Places Indoor

Also show photos not selected for Google Earth

Popular photos in Google Earth

« Previous Next »

image courtesy: panoramio.com

©2015 Google - Map data ©2015 Basarsoft, Google, INEGI, Inav/Geosistemas SRL, ORION-ME | Terms of Use

# Related Work

## Visualization of Geo-tagged Information

### TwitterStand: News in Tweets\*

Jagan Sankaranarayanan  
jagan@cs.umd.edu

Benjamin E. Teitler†  
bteitler@cs.umd.edu

Michael D. Lieberman†  
codepoet@cs.umd.edu

Hanan Samet†  
hjs@cs.umd.edu

Jon Sperling†  
Jon.Sperling@h

#### ABSTRACT

Twitter is an electronic medium that allows a large user population to communicate with each other simultaneously. Inherent to Twitter is an asymmetrical relationship between friends and followers that provides an interesting social network-like structure among the users of Twitter. Twitter messages, called tweets, are restricted to 140 characters and thus are usually very focused. We investigate the use of Twitter to build a news processing system, called *TwitterStand*, from Twitter tweets. The idea is to capture tweets that correspond to late breaking news. The result is analogous to

#### Keywords

Twitter, News, Geotagging, Online

#### 1. INTRODUCTION

Twitter<sup>1</sup> is a social network that has been gaining much attention among millions of users who send messages to each other, where each tweet contains up to 140 characters. At this time, it is estimated that over 10 million users who use Twitter

CHI 2008 Proceedings · Works In Progress

#### Category

H.3 [Information Systems] Storage

#### General Terms

Algorithms

\*This work was supported by the National Science Foundation (NSF) under grant IIS-07-11171. †Present address: Department of Computer Science, University of Massachusetts Lowell, 140 College Ave., Lowell, MA 01854. ‡Present address: Department of Computer Science, University of Maryland, College Park, MD 20742.

## Content Visualization of Geo-located Images

**Marco Cristani**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
marco.cristani@univr.it

**Alessandro Perina**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
perina@sci.univr.it

**Umberto Castellani**  
Università degli Studi di Verona

**Vittorio Murino**  
Università degli Studi di Verona  
Strada le Grazie 15  
Verona, 37134 Italy  
Vittorio.murino@univr.it

**Abstract**  
In this paper we present a photo visualization interface to index and search a large set of geo-located images. This interface is based on a database of images and their associated metadata.

**Key**  
Image

## Placing Flickr Photos on a Map

Pavel Serdyukov \*†  
Database Group  
University of Twente  
PO Box 217, 7500 AE  
Enschede, The Netherlands  
serdyukovpv@cs.utwente.nl

Vanessa Murdock  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
vmurdock@yahoo-inc.com

Roelof van Zwol  
Yahoo! Research  
Diagonal 177  
08018 Barcelona, Spain  
roelof@yahoo-inc.com

#### ABSTRACT

In this paper we investigate how to place photos uploaded to Flickr on a map. We use a language model based on the text of the photos to predict where the image was taken. We use a language model based on the text of the photos to predict where the image was taken.

## PhotoStand: A Map Query Interface for a Database of News Photos

Hanan Samet Marco D. Adelfio Brendan C. Fruin  
Michael D. Lieberman Jagan Sankaranarayanan  
Center for Automation Research, Institute for Advanced Studies,  
Department of Computer Science, University of Maryland  
College Park, MD 20742 USA  
{hjs, marco, brendan, codepoet, jagan}@cs.umd.edu

#### ABSTRACT

PhotoStand enables the use of a map query interface to retrieve news photos associated with news articles that are in turn associated with the principal locations that they mention collected as a result of monitoring the output of over 10,000 RSS news feeds, made available within minutes of publication, and stored in a PostgreSQL database. The news photos are ranked according to their relevance to the clusters of news articles associated with locations at which they are displayed. This work differs from traditional work in this field as the associated locations and topics (by virtue of the cluster with which the articles containing the news photos are associated) are generated automatically without any human intervention such as tagging, and that photos are retrieved by location instead of just by keyword as is the case for many existing systems. In addition, the clusters provide a filtering step for detecting near-duplicate news photos.

#### 1. INTRODUCTION

A demo is presented of PhotoStand (see also the related NewsStand [9, 17, 21, 29], TwitterStand [6, 24], and STEWARD [12] systems) which is an example application of a general framework we are developing for retrieving multimedia data (e.g., text, images, videos) using a map query interface from a database of news articles, photos, and videos (i.e., by location in real-time which differentiates it from other systems).

articles, enabling them to be accessed by spatial queries such as windowing or simple point location; and its *clusterer* [30], which groups articles about the same topic. A key to the NewsStand database system is its pipe server which coordinates its processing modules by assigning batches of articles to them. NewsStand's user interface enables the retrieval of clusters of news articles for display using its map user interface by executing what we term *top-k window queries*. At present, NewsStand handles about 50K articles per day and has a large underlying database of articles currently containing about 300GB of data.

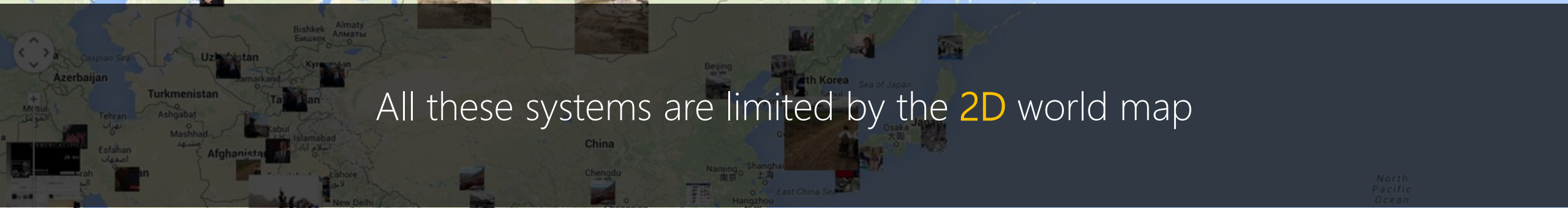
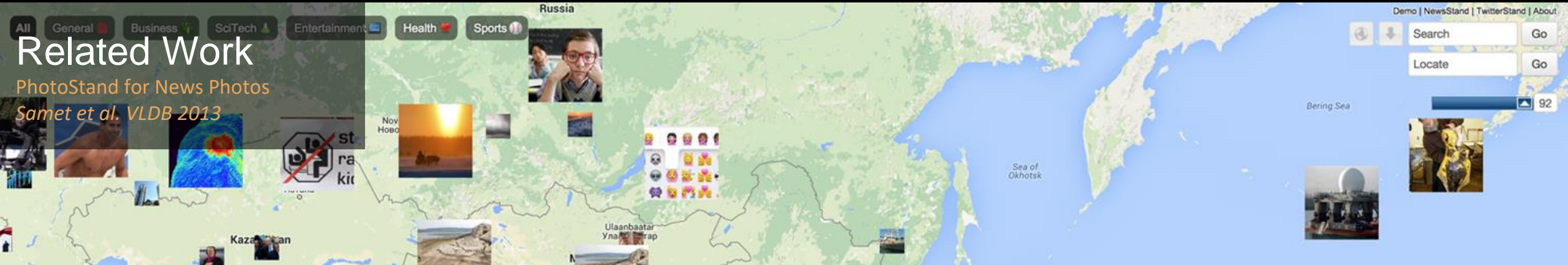
The PhotoStand and TweetPhoto [3] demos are related in the sense that PhotoStand uses photos from news articles in NewsStand, while TweetPhoto uses photos from news tweets in TwitterStand [24]. In addition, the PhotoStand demo demonstrates the database querying capability of NewsStand as well as its capability to do similarity searching for news photos where the first step in the similarity detection process is based on the text associated with the photos, while the second step involves use of the actual image features (e.g., texture, color) to enable detecting near duplicates, thereby avoiding the combinatorial complexity of comparing every photo with every other photo.

The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 indicates how news articles (and consequently news photos) are accessed.

All General Business SciTech Entertainment Health Sports

# Related Work

PhotoStand for News Photos  
Samet et al. VLDB 2013



All these systems are limited by the 2D world map



image courtesy: photostand.umiacs.umd.edu

# Related Work

*Visualization of Geo-tagged Information  
(cont.)*

Previous research also advances server **3D** solutions.



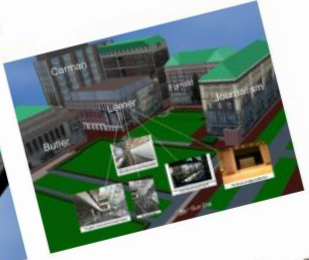
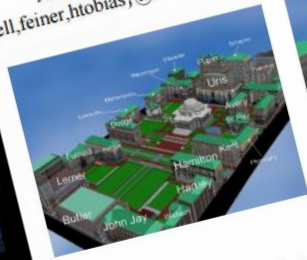
# Related Work

Visualization of Geo-tagged Information  
(cont.)

UIST 2001 (ACM Symp. on User Interface Software and Technology), Orlando, FL, November 11-14, 2001, pp. 101-110

## View Management for Virtual and Augmented Reality

Blaine Bell    Steven Feiner    Tobias Höllerer  
Department of Computer Science  
500 W 120<sup>th</sup> St., 450 CS Building  
Columbia University  
New York, NY 10027  
{bell,feiner,htobias}@cs.columbia.edu



### ABSTRACT

We describe a view-management component for interactive 3D user interfaces. By *view management*, we mean maintaining visual constraints on the projections of objects on the view plane, such as locating related objects near each other, or preventing objects from occluding each other. Our view-management component accomplishes this by modifying selected object properties, including position, size, and transparency, which are tagged to indicate their constraints. For example, some objects may have geometric properties that are determined entirely by a physical simulation and which cannot be modified, while other objects may be annotations whose position and size are flexible.

We introduce algorithms that use upright rectangular extents to represent on the view plane a dynamic and efficient approximation of the occupied space containing projections of visible portions of 3D objects, as well as the space in which objects can be placed to

avoid occlusion. Layout decisions from previous frames are taken into account to reduce visual discontinuities. We present augmented reality and virtual reality examples to which we have applied our approach, including a dynamically labeled and annotated environment.

**CR Categories and Subject Descriptors:** H.5.1 [Information Interfaces—Artificial, augmented, and virtual realities], H.5.2 [Information Interfaces and Presentation] User Interfaces—Graphical User Interfaces, Screen design, I.3.6 [Computer Graphics] Methodology and Techniques—Interaction Techniques, I.3.7 [Computer Graphics] Three-Dimensional Graphics and Realism—Virtual Reality.

**Additional Keywords and Phrases:** view management, environment management, annotation, labeling, wearable computing, augmented reality, virtual environments

### 1. INTRODUCTION

Designing a 3D graphical user interface (UI) requires creating a set of objects and their properties, arranging them in a scene, setting a viewing specification, determining lighting and rendering parameters, and deciding how to light these decisions for each frame. Some of these decisions are fully constrained, for example, a

Registers user-annotated text and images to a particular point in 3D space.

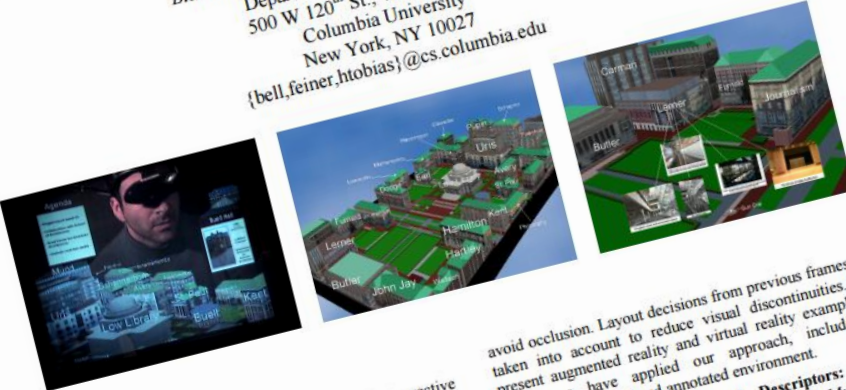
# Related Work

Visualization of Geo-tagged Information  
(cont.)

UIST 2001 (ACM Symp. on User Interface Software and Technology), Orlando, FL, November 11-14, 2001, pp. 101-110

## View Management for Virtual and Augmented Reality

Blaine Bell    Steven Feiner    Tobias Höllerer  
 Department of Computer Science  
 500 W 120<sup>th</sup> St., 450 CS Building  
 Columbia University  
 New York, NY 10027  
 {bell,feiner,htobias}@cs.columbia.edu



**ABSTRACT**  
 We describe a view-management component for interactive 3D user interfaces. By *view management*, we mean maintaining visual constraints on the projections of objects on the view plane, such as locating related objects near each other, or preventing objects from occluding each other. Our view-management component accomplishes this by modifying selected object properties, including position, size, and transparency, which are tagged to indicate their constraints. For example, some objects may have geometric properties that are determined entirely by a physical simulation and which cannot be modified, while other objects may be annotations whose position and size are flexible.

We introduce algorithms that use upright rectangular extents to represent on the view plane a dynamic and efficient approximation of the occupied space containing the projections of visible portions of 3D objects, as well as the space in which objects can be placed to avoid occlusion. Layout decisions from previous frames are taken into account to reduce visual discontinuities. We present augmented reality and virtual reality examples to which we have applied our approach, including a dynamically labeled and annotated environment.

**CR Categories and Subject Descriptors:** H.5.1 [Information Interfaces and Presentation] Multimedia Information Interfaces—Artificial, augmented, and virtual realities; H.5.2 [Information Interfaces and Presentation] User Interfaces—Graphical User Interfaces, Screen design, I.3.6 [Computer Graphics] Methodology and Techniques—Interaction Techniques, I.3.7 [Computer Graphics] Three-Dimensional Graphics and Realism—Virtual Reality.

**Additional Keywords and Phrases:** view management, environment management, annotation, labeling, wearable computing, augmented reality, virtual environments

**1. INTRODUCTION**  
 Designing a 3D graphical user interface (UI) requires creating a set of objects and their properties, arranging them in a scene, setting a viewing specification, determining lighting and rendering parameters, and deciding how to display these decisions for each frame. Some of these decisions may be fully constrained, for example, a building's position and shape must be explicitly defined. In contrast, other decisions may be more flexible.



Both the 3D model and the annotation are predefined for rendering such a scene.

# Related Work

Visualization of Geo-tagged Information  
(cont.)

UIST 2001 (ACM Symp. on User I...)

## View Manager

Bl...



**ABSTRACT**  
We describe a view-man...  
3D user interfaces. By...  
maintaining visual constr...  
on the view plane, such a...  
other, or preventing obje...  
view-management con...  
modifying selected obje...  
size, and transparency, ...  
constraints. For exampl...  
properties that are c...  
simulation and which...  
objects may be anno...  
flexible.  
We introduce algo...  
extents to represent...  
efficient approximat...  
the projections of v...  
the unoccupied spa...

## Photo Tourism: Exploring Photo Collections in 3D

Noah Snavely  
University of Washington

Steven M. Seitz  
University of Washington

Richard Szeliski  
Microsoft Research

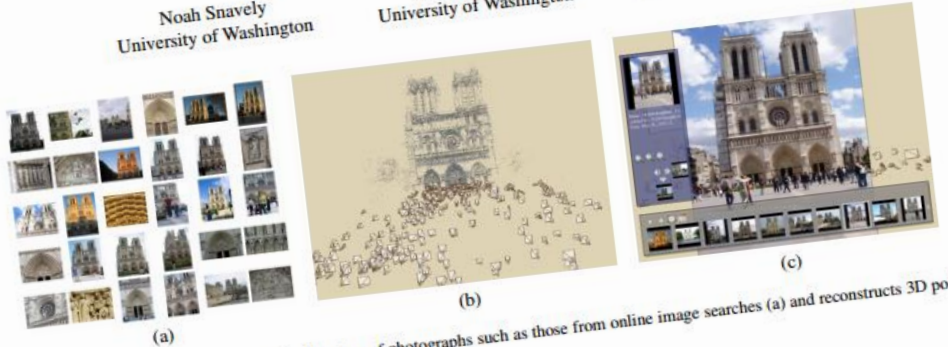


Figure 1: Our system takes unstructured collections of photographs such as those from online image searches (a) and reconstructs 3D points and viewpoints (b) to enable novel ways of browsing the photos (c).

### Abstract

We present a system for interactively browsing and exploring large unstructured collections of photographs of a scene using a novel 3D interface. Our system consists of an image-based modeling front end that automatically computes the viewpoint of each photograph as well as a sparse 3D model of the scene and image to model correspondences. Our *photo explorer* uses image-based rendering techniques to smoothly transition between photographs, while also enabling full 3D navigation and exploration of the set of images and world geometry, along with auxiliary information such as overhead maps. Our system also makes it easy to construct photo tours of scenic or historic locations, and to annotate images. We demonstrate our system on several large personal photo collections as well as images gathered from Internet photo sharing sites.

**CR Categories:** H.5.1 [Information Interfaces and Presentation]: Multimedia Information Systems—Artificial, augmented, and virtual realities 1.2.10 [Artificial Intelligence]: Vision and Scene Understanding—Modeling and recovery of physical attributes

**Keywords:** image-based rendering, image-based modeling, photo browsing, structure from motion

### 1 Introduction

The goal of image-based rendering is to evoke a visceral sense of a scene. The goal of image-based rendering is to evoke a visceral sense of a scene. The goal of image-based rendering is to evoke a visceral sense of a scene.

is that these approaches will one day allow virtual tourism of the world's interesting and important sites.

During this same time, digital photography, together with the Internet, have combined to enable sharing of photographs on a truly massive scale. For example, a Google image search on "Notre Dame Cathedral" returns over 15,000 photos, capturing the scene from myriad viewpoints, levels of detail, lighting conditions, seasons, decades, and so forth. Unfortunately, the proliferation of shared photographs has outpaced the technology for browsing such collections, as tools like Google ([www.google.com](http://www.google.com)) and Flickr ([www.flickr.com](http://www.flickr.com)) return pages and pages of thumbnails that the user must comb through.

In this paper, we present a system for browsing and organizing large photo collections of popular sites which exploits the common 3D geometry of the underlying scene. Our approach is based on computing, from the images themselves, the photographers' locations and orientations, along with a sparse 3D geometric representation of the scene, using a state-of-the-art image-based modeling system. Our system handles large collections of unorganized photographs taken by different cameras in widely different conditions. We show how the inferred camera and scene information enables the following capabilities:

- **Scene visualization.** Fly around popular world sites in 3D by morphing between photos.
- **Object-based photo browsing.** Show me more images that contain this object or part of the scene.
- **Where was I?** Tell me where I was when I took this picture.
- **What am I looking at?** Tell me about objects visible in this image by transferring annotations from similar images.

# Related Work

Visualization of Geo-tagged Information  
(cont.)

UIST 2001 (ACM Symp. on User Interface Design)

## View Manager

Bl



**ABSTRACT**  
We describe a view-management system for 3D user interfaces. By maintaining visual constraints on the view plane, such as object size, or preventing object overlap, our view-management system can modify selected objects' size, and transparency, and other visual constraints. For example, we can simulate which objects may be annotated in a 3D scene. We introduce algorithms for efficient approximation of the projections of the unoccupied space.

### Photo Tourism: Exploring Photos through Dynamic Views

Noah Snavely  
University of Washington



(a)

Figure 1: Our system takes unstructured collections of photographs and viewpoints (b) to enable novel ways of viewing them.

**Abstract**  
We present a system for interactively browsing unstructured collections of photographs through a 3D interface. Our system consists of a front end that automatically computes correspondences. Our *photo explorer* uses techniques to smoothly transition between enabling full 3D navigation and exploring world geometry, along with auxiliary maps. Our system also makes it easy to explore scenic or historic locations, and to share our system on several large public servers as images gathered from Internet photo-sharing services.

**CR Categories:** H.5.1 [Information Systems and Multimedia] Information Systems and Multimedia—User Interfaces; I.2.10 [Artificial Intelligence]—Modeling and Simulation

**Keywords:** image-based rendering, structure from motion, 3D visualization

## 1 Introduction

# Social Snapshot: A System for Temporally Coupled Social Photography

Robert Patro, Cheuk Yiu Ip, Sujal Bista, and Amitabh Varshney • University of Maryland, College Park

**S**ince the invention of photography, taking pictures of people, places, and activities has become integral to our lives. In the past, only purposeful, precious moments were the primary subjects of photography. But technological advances have brought photography to our everyday lives in the form of compact cameras and even cell phone cameras.

### Social Snapshot's Contributions

Social Snapshot's contributions fit naturally into two categories: technical and social.

The technical contributions are improved algorithms and techniques that enhance our system's novelty and scalability. For example, Social Snapshot produces a textured and colored-mesh reconstruction from a loosely ordered photo collection, rather than the sparse or dense point reconstructions produced by related approaches. In addition, it features locally optimized mesh generation and viewing. Finally, it provides camera network capabilities to support synchronized capture of temporally dynamic data.

The social contributions lead to a new way of thinking about the interplay between data acquisition and social interactions. They also let us define social photography as an active, rather than a passive, endeavor. For example, Social Snapshot encourages collaborative photography as a social endeavor, letting users capture dynamic action by synchronizing their photographs. It leverages social trends such as online media sharing and event organization to spur a novel data acquisition mode.

For a look at some of the previous research on which Social Snapshot is based, see the "Related Work in Scene Visualization and Computer Vision" sidebar on pages 78–79.

The next phase in the photography revolution, 3D photography, can bring users together to socialize and collaboratively take pictures in an entirely new way. However, transforming a photographic scene from 2D to 3D requires introducing multiple images of the same underlying geometry from different viewpoints. The reconstruction of 3D geometry from multiple overlapping images is the classic structure-from-motion (SFM) problem in computer vision. Typically, the instruments used to acquire photographs are tediously calibrated to produce precise measurements.

To simplify 3D photography, our Social Snapshot system performs active acquisition and reconstruction of temporally dynamic data. Using multiple users' cell phone cameras and no preliminary calibration, it achieves approximate but visually convincing renderings of 3D scenes, even though the quality of the input images is low.


# Related Work

Visualization of Geo-tagged Information (cont.)

UIST 2001 (ACM Symp. on User Interface Technology)

## View Manager

Bl



(a)




Figure 1: Our system takes unstructured collections of photographs and viewpoints (b) to enable novel ways of browsing the data.

**Abstract**

We present a system for interactively browsing unstructured collections of photographs in a 3D interface. Our system consists of a front end that automatically computes correspondences. Our *photo explorer* techniques to smoothly transition between enabling full 3D navigation and exploring world geometry, along with auxiliary maps. Our system also makes it easy to explore scenic or historic locations, and to share our system on several large public spaces as images gathered from Internet photo sharing sites.

**CR Categories:** H.5.1 [Information Systems—Multimedia]—Artificial Intelligence; I.2.10 [Artificial Intelligence—Modeling and Simulation]—Modeling and Simulation

**Keywords:** image-based rendering, structure from motion, 3D visualization

### 1 Introduction

Instead of relying on static images embedded in text, suppose you could create an interactive, photorealistic visualization, where, for example, a Wikipedia page is shown next to a detailed 3D model of the described site. When you select an object in the scene via a smooth, photorealistic visualization, when you click on an object in the corresponding descriptive text, the corresponding descriptive text is to create such a visualization.

# Social Snapshot: A System for Temporal Social Photography

Robert Patro, Cheuk Yiu Ip, Sujal Bista, and Amitabh Varshney

**Abstract**

Since the invention of photography, pictures of people, places, and events have become integral to our lives. Photography is not only purposeful, precious moments but also a primary subject of photography. But advances have brought photography into the 21st century in the form of compact camera phones. The next phase of photography revolution, which can bring socialization and connectivity to pictures in an augmented reality (AR) environment. However, traditional photography requires intricate techniques of image capture and processing. The reconstruction of a scene from images is a non-trivial task. Instruments used to acquire images are often calibrated to produce precise results. To simplify 3D photography, our system performs active calibration of temporally dynamic data. It achieves convincing rendering of the scene by calibrating the user's cell phone camera. It achieves this by using a novel technique of image-based rendering, structure from motion, and 3D visualization.

**Keywords:** image-based rendering, structure from motion, 3D visualization

**Links:** [DL](#) [PDF](#)

# 3D Wikipedia: Using online text to automatically label and navigate reconstructed geometry

Bryan C. Russell<sup>1</sup>, Ricardo Martin-Brualla<sup>2</sup>, Daniel J. Butler<sup>2</sup>, Steven M. Seitz<sup>2</sup>, Luke Zettlemoyer<sup>2</sup>

<sup>1</sup>Intel Labs



Figure 1: Given a reference text describing a specific site, for example the Wikipedia article above for the Pantheon, we automatically create a labeled 3D reconstruction, with objects in the model linked to where they are mentioned in the text. The user interface enables coordinated browsing of the text with the visualization (see video).

**Abstract**

We introduce an approach for analyzing Wikipedia and other text, together with online photos, to produce annotated 3D models of famous tourist sites. The approach is completely automated, and leverages online text and photo co-occurrences via Google Image Search. It enables a number of new interactions, which we demonstrate in a new 3D visualization tool. Text can be selected to move the camera to the corresponding objects, 3D bounding boxes provide anchors back to the text describing them, and the overall narrative of the text provides a temporal guide for automatically flying through the scene to visualize the world as you read about it. We show compelling results on several major tourist sites.

**CR Categories:** H.5.1 [Information Systems—Artificial Intelligence]; I.2.7 [Artificial Intelligence—Multimedia]—Artificial Intelligence; I.2.10 [Artificial Intelligence—Modeling and Simulation]—Text analysis; I.2.10 [Artificial Intelligence—Modeling and Simulation]—Modeling and recovery of physical attributes

**Keywords:** image-based modeling and rendering, Wikipedia, natural language processing, 3D visualization

**Links:** [DL](#) [PDF](#)

### 1 Introduction

Tourists have long relied on guidebooks and other reference texts to learn about and navigate sites of interest. While guidebooks are packed with interesting historical facts and descriptions of specific objects and spaces, it can be difficult to fully visualize the scenes they present. The primary cues come from images provided with the text, but coverage is sparse and it can be difficult to understand the spatial relationships between each image viewpoint. For example, the Berlitz and Lonely Planet guides [Berlitz International 2003; Garwood and Hole 2012] for Rome each contain just a single photo of the Pantheon, and have a similar lack of photographic coverage of other sites. Even online sites such as Wikipedia, which do not have space restrictions, have similarly sparse and disconnected visual coverage.

Instead of relying on static images embedded in text, suppose you could create an interactive, photorealistic visualization, where, for example, a Wikipedia page is shown next to a detailed 3D model of the described site. When you select an object in the scene via a smooth, photorealistic visualization, when you click on an object in the corresponding descriptive text, the corresponding descriptive text is to create such a visualization.

# Related Work

Visualization of Geo-tagged Information (cont.)

UIST 2001 (ACM Symp. on User Interface Design)

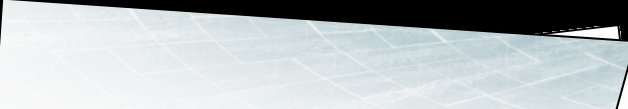
## Photo Tourism: Enabling Novel Views of Images

Noah Snavely  
University of Washington



## Social Snapshot: A System for Temporal Social Photography


Robert Patro, Cheuk Yiu Ip, Sujit Bista, and Anil K. Jain



## 3D Wikipedia: Using online text to automatically label and navigate reconstructed geometry

Bryan C. Russell<sup>1</sup>  
Ricardo Martin-Brualla<sup>2</sup>  
Daniel J. Butler<sup>2</sup>  
Steven M. Seitz<sup>2</sup>  
Luke Zettlemoyer<sup>2</sup>

<sup>1</sup>Intel Labs  
<sup>2</sup>University of Washington



They use offline 3D reconstruction algorithm or existing 3D models for registering the photos onto the meshes. Such methods usually suffer from **hours of processing time** for a single location.

## View Manager

ABSTRACT  
We describe a view-management interface. By maintaining visual constraints on the view plane, such as other, or preventing object view-management constraints. For example, modifying selected object size, and transparency, simulation and which objects may be annotated.



Figure 1: Our system takes unstructured images and viewpoints (a) to enable novel ways of viewing the scene (b).

### Abstract

We present a system for interactively browsing unstructured collections of photographs through a 3D interface. Our system consists of a front end that automatically computes correspondences. Our *photo explorer* enables full 3D navigation and exploration of world geometry, along with auxiliary maps. Our system also makes it easy to explore scenic or historic locations, and to share our system on several large photo collections as images gathered from Internet photo sharing sites.

**CR Categories:** H.5.1 [Information Systems—Multimedia Information Systems]—User Interfaces I.2.10 [Artificial Intelligence—Modeling and Simulation]—Modeling and Simulation

**Keywords:** image-based rendering, structure from motion, 3D visualization

### 1 Introduction

The next phase of photography revolution, can bring socialize and connect pictures in an online environment. However, traditional photographic scene requires intrusions of the geometry from the reconstruction. The reconstruction of images is from motion capture. Instruments used to acquire calibrated to produce precise. To simplify 3D photography system performs active construction of temporally multiple users' cell phone calibration, it achieves convincing rendering of the scene.

**Social Snapshot actively acquires and reconstructs temporally dynamic data. The system enables spatiotemporal 3D photography using commodity devices, assisted by their auxiliary sensors and network functionality. It engages users, making them active rather than passive participants in data acquisition.**

Figure 1: Google Earth interface text describing a specific location in the 3D reconstruction, with the corresponding text in the Wikipedia article above for the Pantheon, we automatically create links to where they are mentioned in the text. The user interface enables coordinated viewing of the scene.

We introduce an approach for analyzing Wikipedia and other text, together with online photos, to produce annotated 3D models of famous tourist sites. The approach is completely automated, and leverages online text and photo co-occurrences via Google Image Search. It enables a number of new interactions, which we demonstrate in a new 3D visualization tool. Text can be selected to move the camera to the corresponding objects, 3D bounding boxes provide anchors back to the text describing them, and the overall narrative of the text provides a temporal guide for automatically flying through the scene to visualize the world as you read about it. We show compelling results on several major tourist sites.

**CR Categories:** H.5.1 [Information Interfaces and Presentation]: Multimedia Information Systems—Artificial, augmented, and virtual realities I.2.7 [Artificial Intelligence]: Natural Language Processing—Text analysis I.2.10 [Artificial Intelligence]: Vision and Scene Understanding—Modeling and recovery of physical attributes

**Keywords:** image-based modeling and rendering, Wikipedia, natural language processing, 3D visualization

**Links:** [DL](#) [PDF](#)

### 1 Introduction

Tourists have long relied on guidebooks and other reference texts to learn about and navigate sites of interest. While guidebooks are packed with interesting historical facts and descriptions of sites, specific objects and spaces, it can be difficult to fully visualize the scenes they present. The primary cues come from images provided with the text, but coverage is sparse and it can be difficult to understand the spatial relationships between each image viewpoint. For example, the Berlitz and Lonely Planet guides [Berlitz International 2003; Garwood and Hole 2012] for Rome each contain just a single photo of the Pantheon, and have a similar lack of photographic coverage of other sites. Even online sites such as Wikipedia, which do not have space restrictions, have similarly sparse and disconnected visual coverage.

Instead of relying exclusively on static images embedded in text, suppose you could create an interactive, photorealistic visualization, where, for example, a Wikipedia page is shown next to a detailed 3D model of the described site. When you select an object location in the scene via a smooth, photorealistic visualization (e.g., "Raphael's tomb") in the text, it flies you to the corresponding descriptive text, and you click on an object in the visualization to create such a visualization. This interactive visualization is to create such a visualization.

So what about **our approach**?

# Social Street View



# Social Street View



*The first immersive social  
media navigation system in  
mixed-reality!*

# Demonstration

The Augmentarium, UMIACS  
6000 x 3000 pixels



# Natural Immersive *Virtual Reality*?

# Algorithm

Adding depth & normal map &  
maximal Poisson-disk sampling



Social Street View enables users to see-through the nearby restaurants.

# Conception, architecting & implementation

## Social Street View

A mixed reality system that can depict geo-tagged social media in immersive 3D web environments

# Blending multiple modalities of

## Street View + Social Media

Depth maps, normal maps, and road orientation  
GPS coordinates and time creation

# 2

# Enhancing visual augmentation

3

## Maximal Poisson-disk sampling

Evaluated by image saliency metrics

# Achieving cross-platform compatibility by

## WebGL + Three.js

smartphones, tablets, desktop, high-resolution large-area wide field of view tiled display walls, as well as head-mounted displays.





# Technical Challenges?





# Architecture

## Social Street View System Flowchart



Street View Panorama and Depth



Geo-tagged Social Media

# Street View Cars - Cameras, LIDARs and GPS

Image courtesy from Google Street View





# Tiles of Panoramic Images

Image courtesy from Google Street View

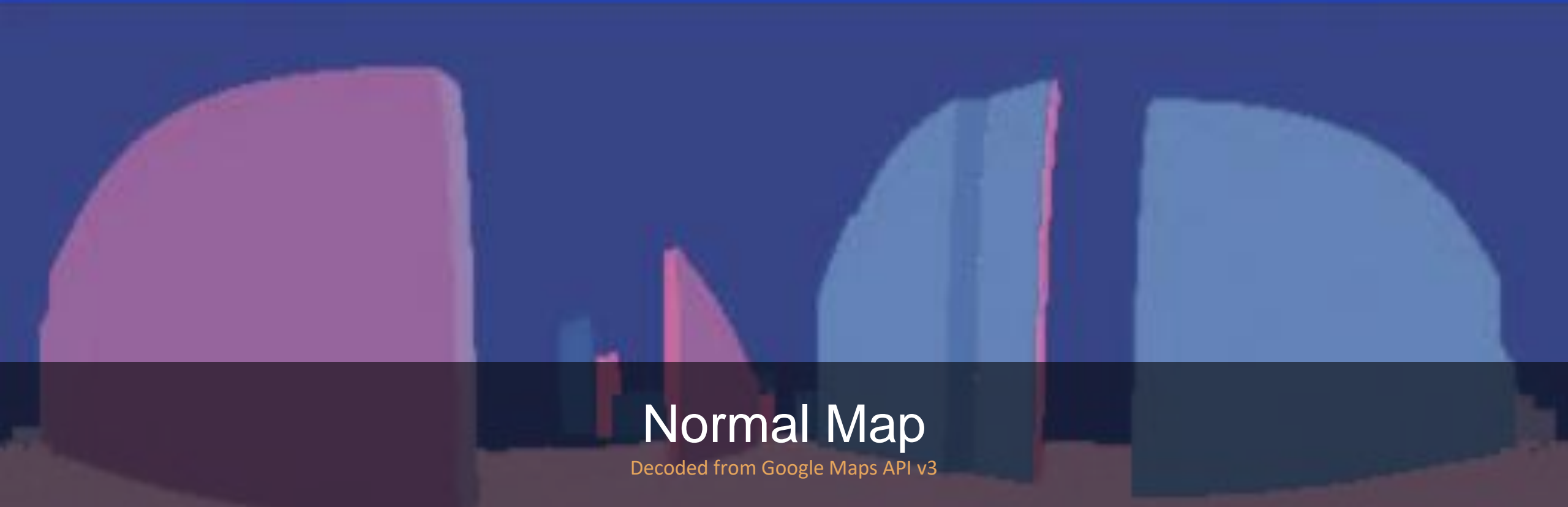


# Panorama



# Depth Map

Decoded from Google Maps API v3 and GSVPanoDepth.js



# Normal Map

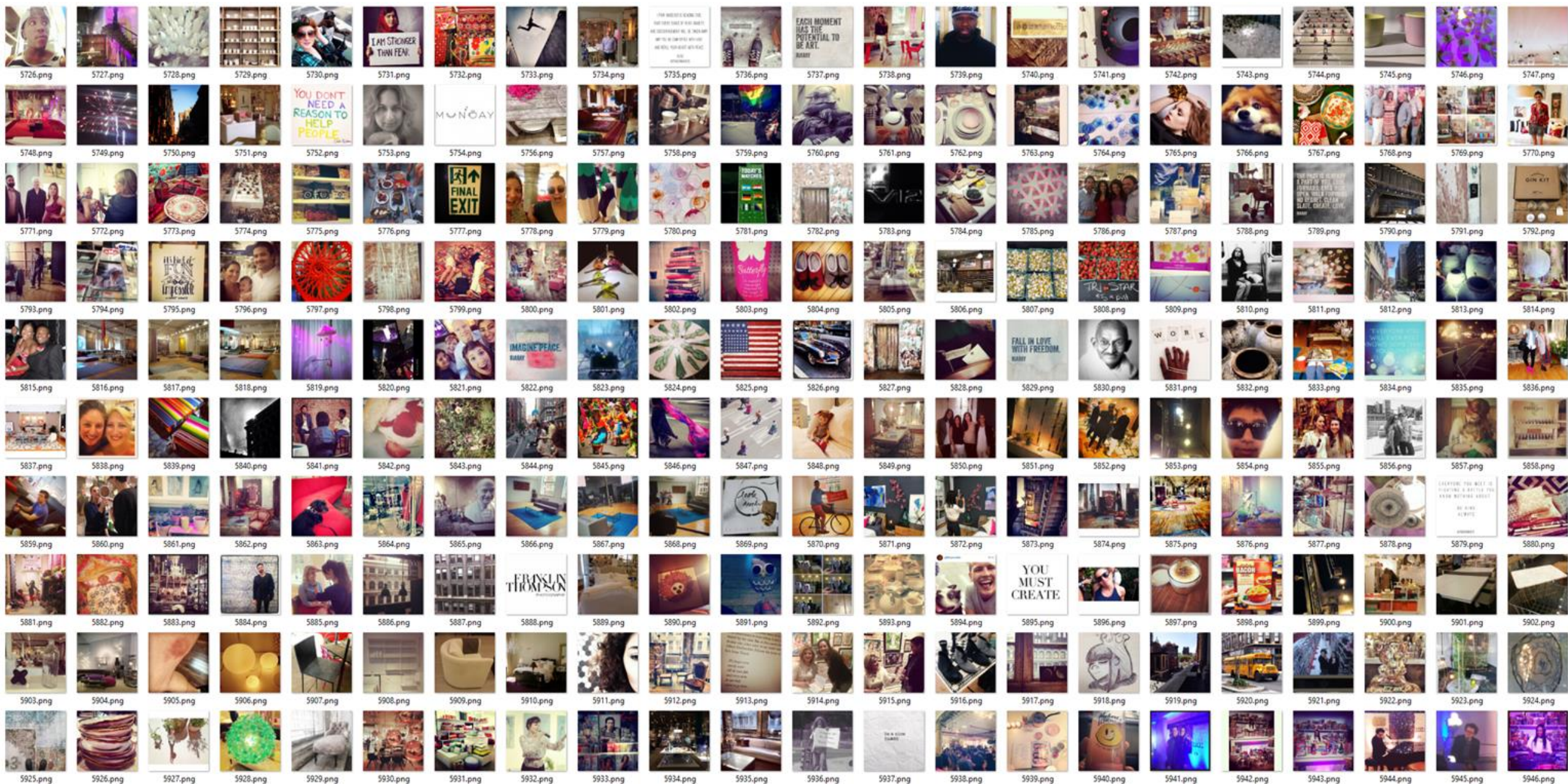
Decoded from Google Maps API v3

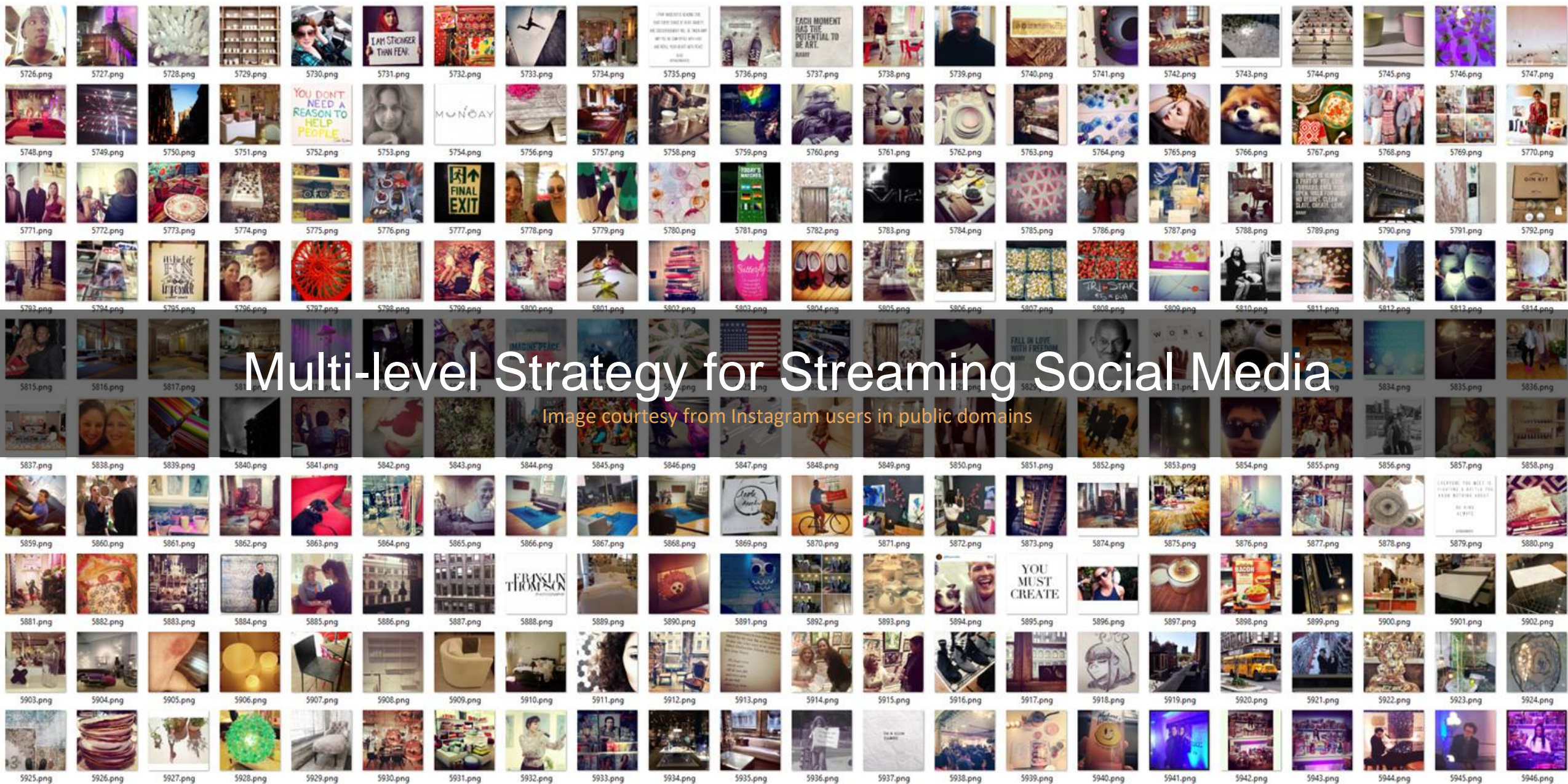




# Road Orientations

Decoded from Google Maps API v3





# Multi-level Strategy for Streaming Social Media

Image courtesy from Instagram users in public domains



8 METERS

16 METERS

.....

1024 METERS





# Haversine Formula

Andrew, 1805

$$\alpha_{ij} = \sin^2\left(\frac{\varphi_i - \varphi_j}{2}\right) + \cos \varphi_i \cdot \varphi_j \cdot \sin^2\left(\frac{\lambda_i - \lambda_j}{2}\right)$$

The distance between social media and street view panorama

**Haversine formula**, which defines the distance on the surface of a sphere.

$$\beta_{ij} = 2 \cdot \operatorname{atan2}\left(\sqrt{\alpha_{ij}}, \sqrt{1 - \alpha_{ij}}\right)$$

$$d_{ij} = R \cdot \beta_{ij}$$

Keywords  Search

Urban Rural Indoor Terrain Misc

NYC London Paris Rome Tokyo

Month  1 - 9

Hour  3 - 20

Distance  0 - 64

2D Depth Side Front Model

Enhance Dark None Enhance Bright

Spring Summer Autumn Winter

All Face-only None face-only

All No-text only With text only Include-video

DirLight  Unlimited

AmbLight  Unlimited

BloomEffect  Unlimited

DarkEffect  Unlimited

LightColor  Unlimited

Number  Unlimited

Radius  Unlimited

Help

Credits



Map Satellite

Welcome to Social Street View

You are standing at: 40.75986105721-73.982681004327  
 Address: 129 W 49th St, New York  
 HashID: ayQskPofZv-friZyq3A Zoom: 5, G: 1



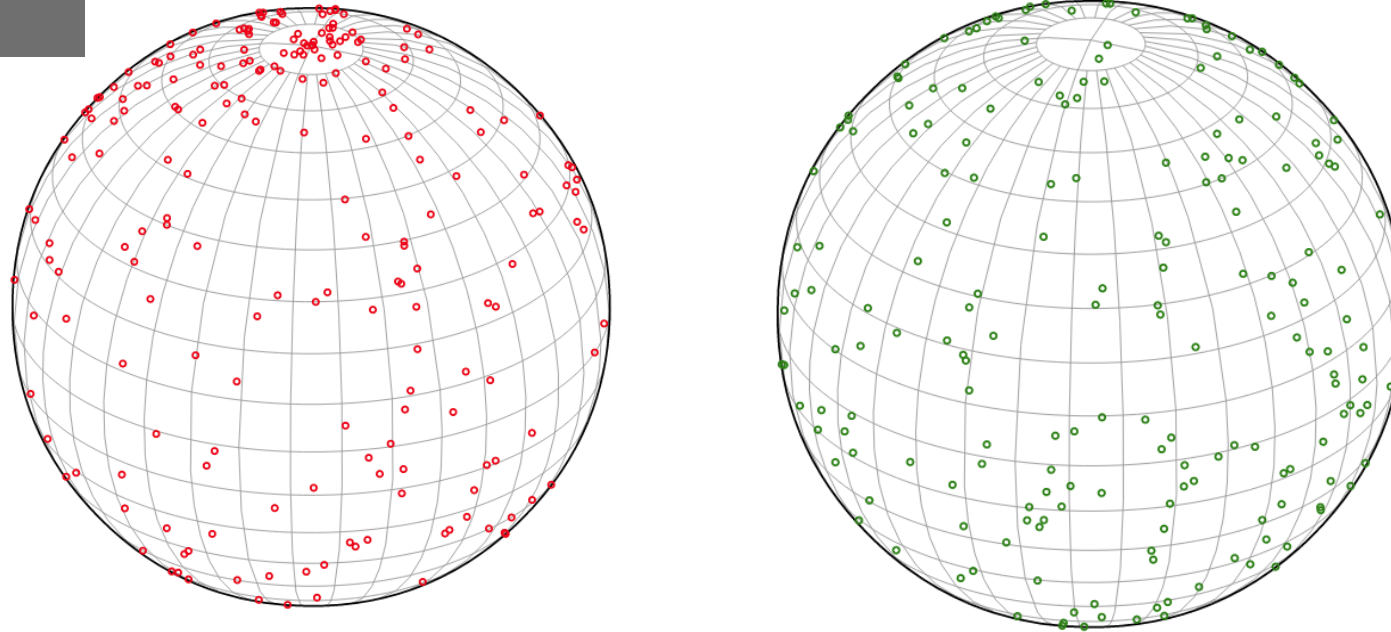




How can we **render and layout** the social media?

# Baseline

Random Uniform Sampling

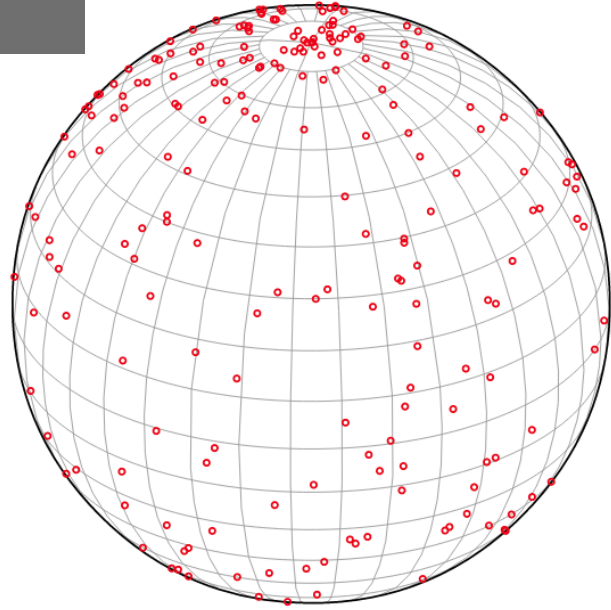


$$\varphi_i = (u_i - \frac{1}{2})\pi, \lambda_i = (2v_i - 1)\pi$$

$$x_i = \cos \varphi_i \cos \lambda_i, y_i = \sin \varphi_i, z_i = \cos \varphi_i \sin \lambda_i$$

# Baseline

Random Uniform Sampling



Without uniform random sampling

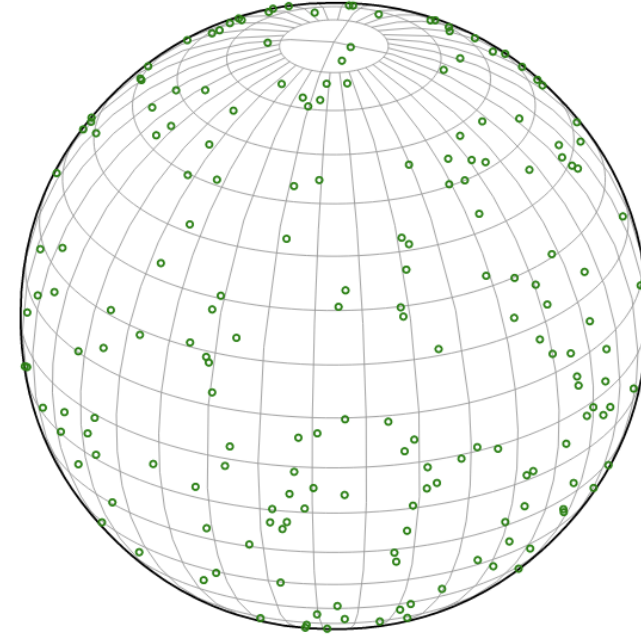
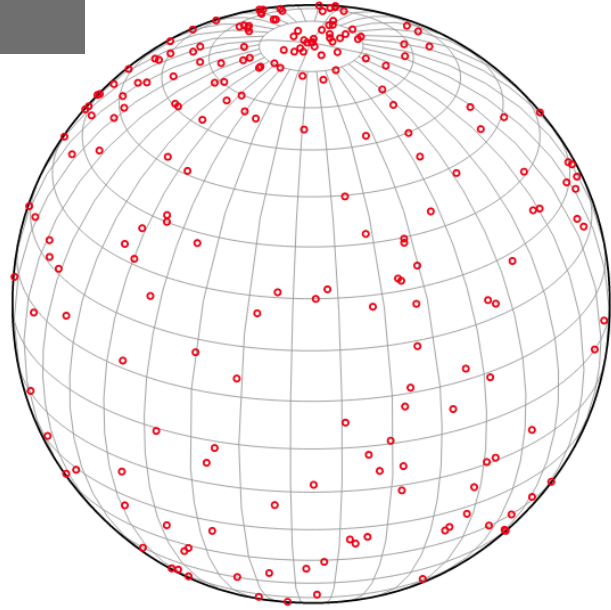
# Without uniform sampling

Accumulation



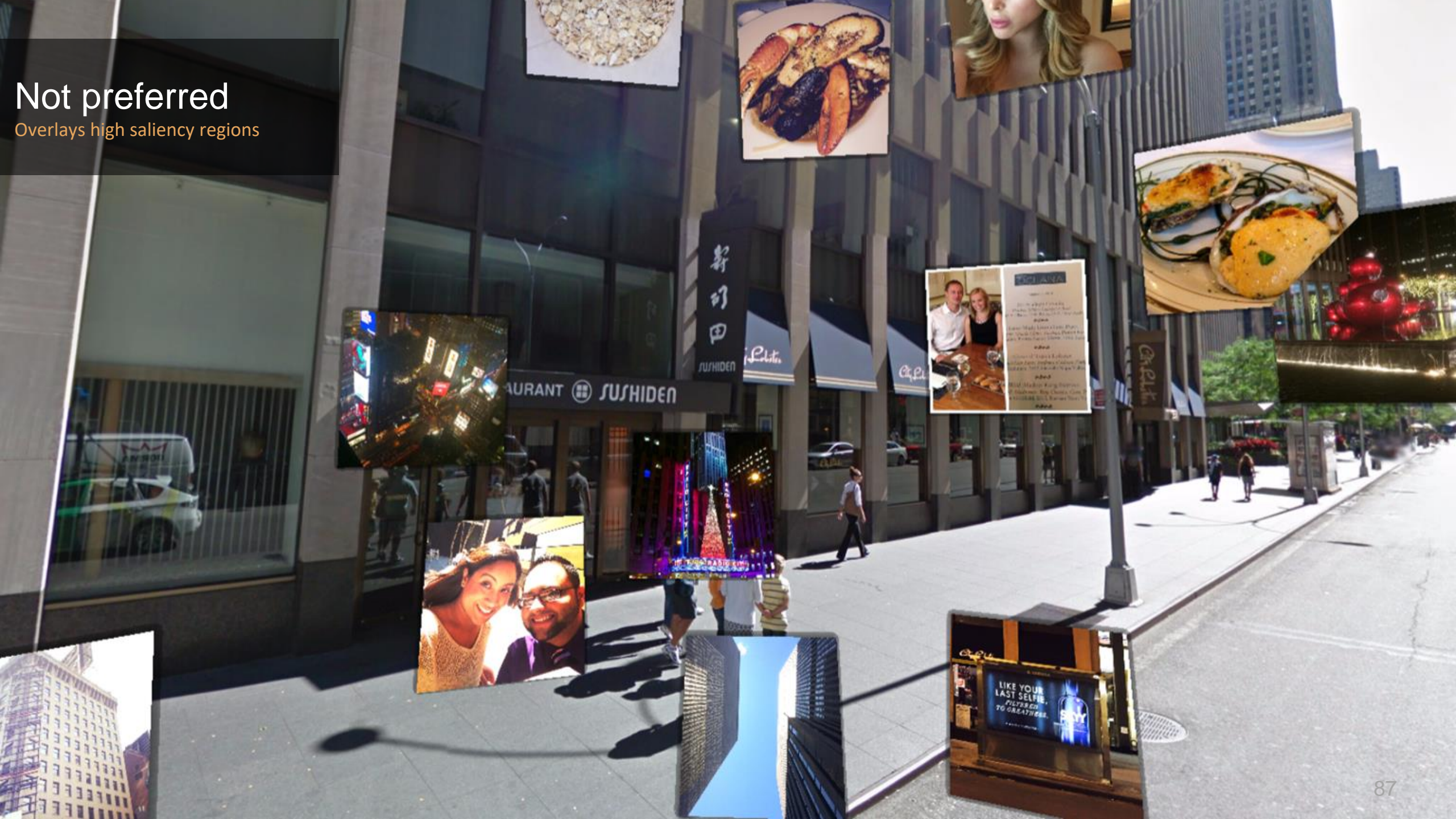
# Baseline

Random Uniform Sampling



With uniform random sampling

Not preferred  
Overlays high saliency regions



# Add depth map

Remove sky and ground (most)



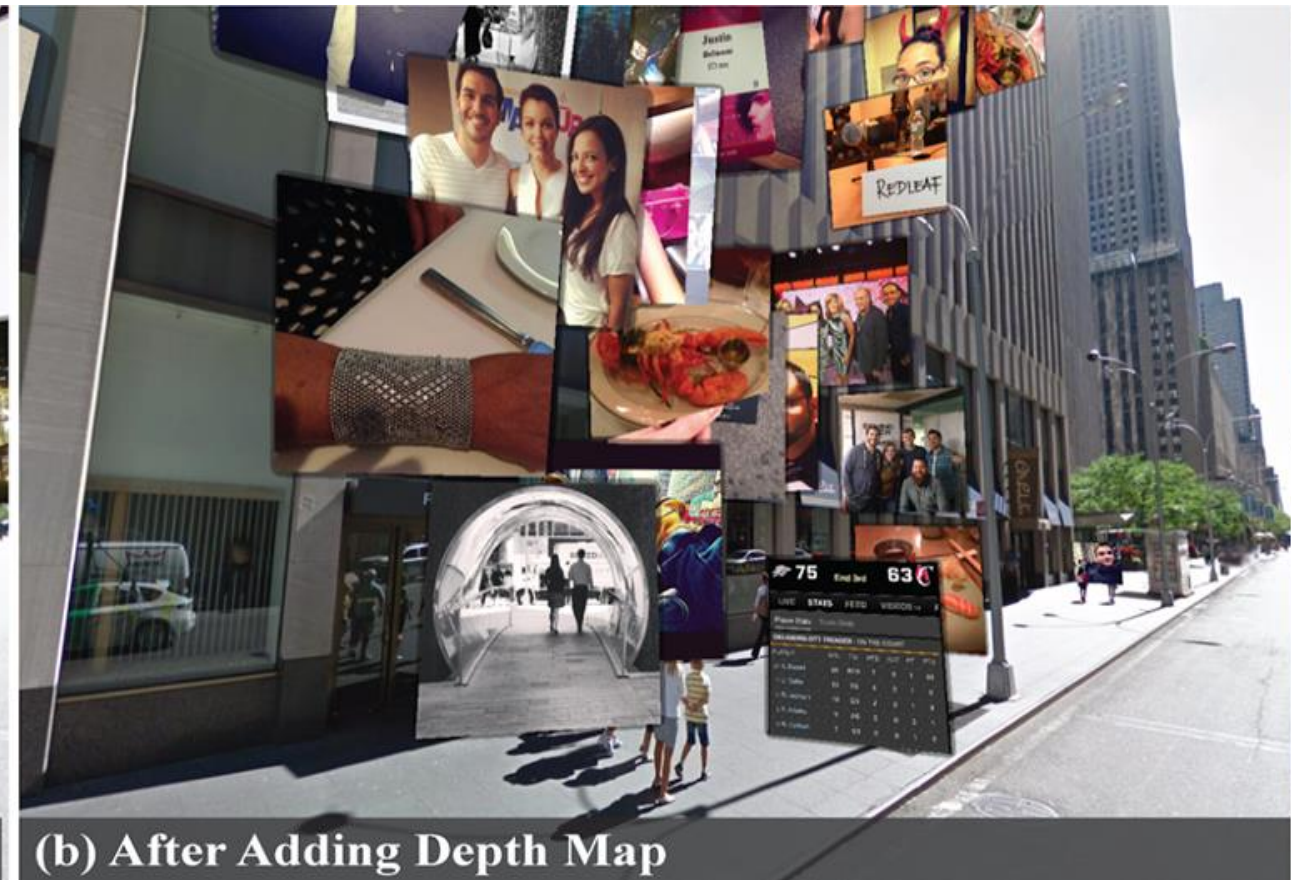
Sky:  $\Omega_s = \{q_i \mid \forall q_i \in \Omega \wedge d_i = \infty\}$

Distance Limit:  $\forall \tilde{p}_i \in T, D_{min} < \tilde{d}_i < D_{max}$



# Add depth map

Remove sky and ground (most)



Can we ensure each image **aligns with** the building geometries?

# Add normal map

Remove all ground, align images



Define the ground region:

$$\Omega_g = \{q_i \mid \forall q_i \in \Omega \wedge \|\mathbf{n}_i - \mathbf{n}_g\| < \delta\}$$



Can we further reduce **visual clutter** and **occlusion**?

# Maximal Poisson-disk Sampling

*Gamito et al. Remove visual clutter and occlusion*

$$\forall \tilde{p}_i \in \tilde{P}, \tilde{P} \subseteq T, \forall S \subseteq \Omega : Pr(q_i \in S) = \int_S di \quad (10)$$

The Poisson-disk distribution is **uniform**.

$$\forall p_i, p_j \in \tilde{T}, p_i \neq p_j : \|p_i - p_j\| \geq r \quad (11)$$

Minimum distance between each pair of social media is **greater than R**.

$$S(X) = \{ \tilde{p}_j \in T : \| \tilde{p}_i - \tilde{p}_j \| \geq r, \tilde{p}_i \in \tilde{P} \} : S(X) = \emptyset \quad (12)$$

Terminates the sampling when it reaches the **maximal coverage**.

# Dart-throwing Algorithm

*PixelPie by Ip et al. using vertex and fragment shaders*



# Pixel-Pie Algorithm

Remove when occlusion occurs

---

**Algorithm 1** Maximal Poisson-disk sampling by dart-throwing

---

**Input:** The minimum distance  $r$  between sampled points

**Output:** A set  $\tilde{P}$  of points which satisfy equation (10)-(12)

1: Set  $\tilde{P} \leftarrow \emptyset$ , empty region  $\tilde{R} \leftarrow T$

2: **repeat**

Pixel-Pie by Ip. et al. uses GPU depth-testing feature for efficient sampling.

3: Generate some random points  $P' \subseteq R$  by rasterizing them as circular disks into a depth map in vertex shader.

4: Remove any point  $\tilde{p} \in \tilde{P}'$  whose corresponding point  $\tilde{q}$  violates  $\tilde{q} \in \Omega_g \vee D_{min} < \tilde{d} < D_{max}$

5: Identify and remove the occluded disks from  $\tilde{P}$  by reading the depth map in the shader.

6:  $\tilde{P} \leftarrow \tilde{P} \cup \tilde{P}'$

7: Update the empty region  $R$  in the fragment shader.

8: **until**  $R \leftarrow \emptyset$

---



# Project Social Media Pictures

By Maximal Poisson-disc Sampling

---

## ALGORITHM 1: Social Media Layout using Poisson-disk Samples

---

**Input:**  $N$  sorted social media images  $\hat{S} = \{s_i \mid i = 1 \dots N\}$ , acquired from SSV servers.

**Output:** A set of image planes to display social media:  $I = I_1 \dots I_M, M \leq N$ .

Generate the set of candidate sample points  $\tilde{P}$  by the PixelPie algorithm;

**Sampling** in regions which are not *Sky* nor *Ground*, and with a limitation of depth values; sort points in  $\tilde{P}$  in descending order according to their corresponding values in the depth map  $D$  so that the closest sample point is laid out first;

Set  $I \leftarrow \emptyset$ ;  
Place social media as **billboards upon** the building geometries;

**Cast soft shadows** as if lighting were 45 degree to the normal vectors.  
perspective visual effects;

Rotate  $I_i$  so that it is perpendicular to the normal vector  $\mathbf{n}_i \leftarrow \mathbf{N}(u_i, v_i)$ ;

Add  $I_i$  to the result set:  $I \leftarrow I \cup I_i$ ;

**end**

---

# Sampling Comparison

*Remove visual clutter and occlusion*



# Algorithm

Adding depth & normal map &  
maximal Poisson-disk sampling



In addition, our system allows users to walk around and explore live social media streams.

This algorithm works well in **dense urban areas** such as the Manhattan District,

What if there are **no buildings** in the scene?

# Scenic Landscapes

Using orientation of the road

---

**ALGORITHM 2:** Social Media Layout using Road Orientations

---

**Input:**  $|O|$  road orientations with  $o_i \in [0, 2\pi]$ .  $K$  social media to be placed for each orientation. Typically,  $|O| = 2$  for a road with two orientations.

**Output:** A set of image planes to display social media:

$$I = I_1 \dots I_M, M \geq K \cdot |O|.$$

Set  $I \leftarrow \emptyset$ ;

for  $i \leftarrow 1 \dots |O|$  do

Set  $\tilde{\mathbf{q}} \leftarrow (R \cos(o_i - \frac{\pi}{2}), h, R \sin(o_i - \frac{\pi}{2}))$  with radius  $R$ ;

(Optional based on user's preference) Add a frontal image plane to  $I$  at  $\mathbf{q}_i$ ;

Set the translation  $\mathbf{t} \leftarrow (T \cos(o_i + \frac{\pi}{2}), 0, T \sin(o_i + \frac{\pi}{2}))$  with constant  $T$ ;

for  $k \leftarrow 1 \dots K$  do

Set  $\tilde{\mathbf{q}} \leftarrow (kR \cos o_i, h, kR \sin o_i)$ ;

Add a left side image plane to  $I$  at position  $\mathbf{q}' \leftarrow \tilde{\mathbf{q}} + \mathbf{t}$ ;

Add a right side image plane to  $I$  at position  $\mathbf{q}' \leftarrow \tilde{\mathbf{q}} - \mathbf{t}$ ;

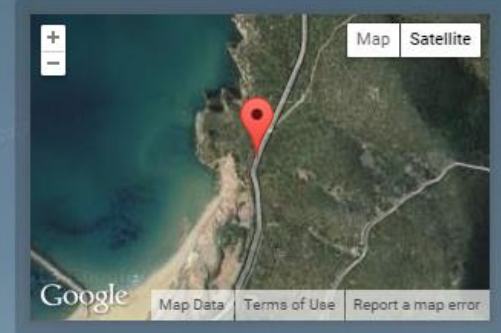
end

end

---

# Scenic Landscapes

*Using orientation of the road*



# Scenic Landscapes

*Using orientation of the road*



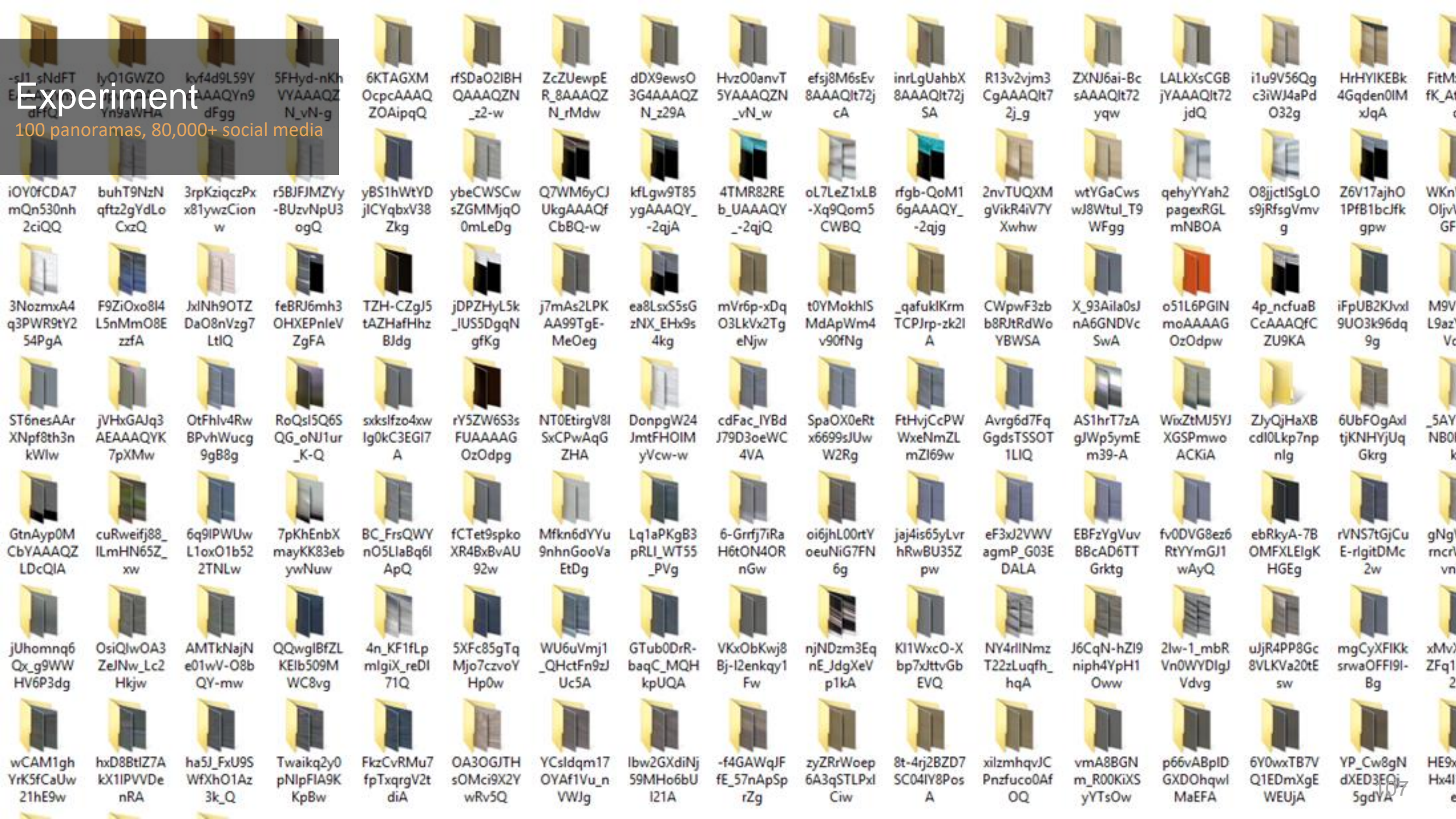




What if ...

Temporal information is used for filtering and rendering?





Experiment  
100 panoramas, 80,000+ social media

-sJ1\_sNdFT  
E dHFC

IOY0fCDA7  
mQn530nh  
2ciQQ

3NozmxA4  
q3PWR9tY2  
54PgA

ST6nesAAr  
XNpf8th3n  
kWlw

GtnAyp0M  
CbYAAAQZ  
LDcQIA

jUhomnq6  
Qx\_g9WW  
HV6P3dg

wCAM1gh  
YrK5fCaUw  
21hE9w

6KTAGXM  
OpcpAAAQ  
ZOaipqQ

yBS1hWtYD  
jICYqbxV38  
Zkg

TZH-CZgJ5  
tAZHafHhz  
BJdg

sxkslfzo4xw  
lg0kC3EGi7  
A

BC\_FrsQWY  
nO5LlAbq6l  
ApQ

4n\_KF1fLp  
mlgiX\_reDl  
71Q

FkzCvRMu7  
fpTxqrgV2t  
diA

r5DAoO2IBH  
QAAAQZN  
\_z2-w

ybeCWSCw  
sZGMMjqO  
0mLeDg

jDPZHyL5k  
\_IUS5DgqN  
gfKg

rY5ZW6S3s  
FUAAAAG  
OzOdpq

fCTet9spko  
XR4BxBvAU  
92w

5XFc85gTq  
Mjo7czvoY  
Hp0w

OA3OGJTH  
sOMci9X2Y  
wRv5Q

ZcZUewpE  
R\_8AAAQZ  
N\_rMdw

Q7WM6yCJ  
UkgAAAQf  
CbBQ-w

j7mAs2LPK  
AA99TgE-  
MeOeg

NT0EtirgV8l  
SxCPwAqG  
ZHA

Mfkn6dYYu  
9nhnGooVa  
EtDg

WU6uVmjl  
\_QHctFn9zj  
Uc5A

YCslqdm17  
OYAf1Vu\_n  
VWJg

FitM

WKn

M9V

\_5AY

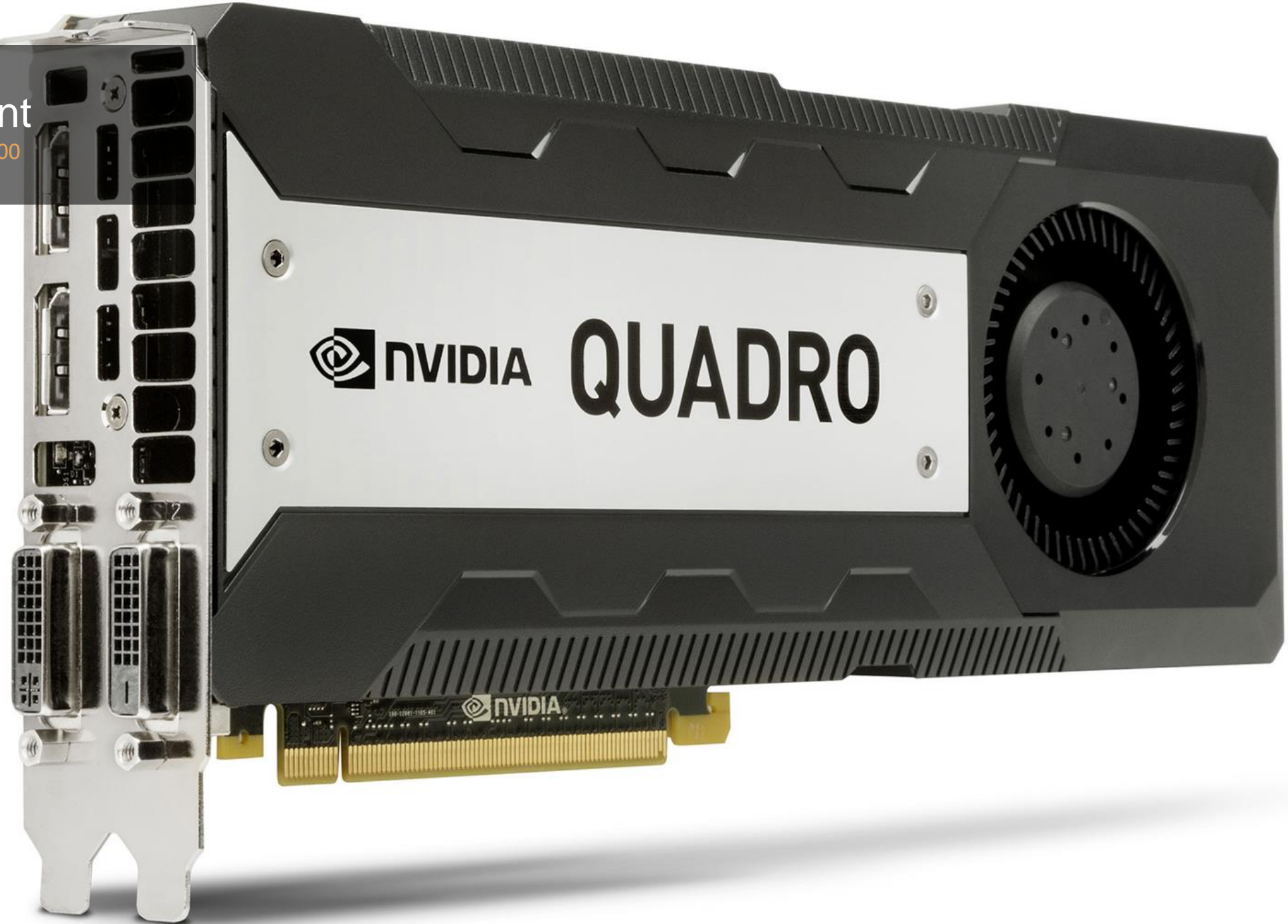
gNg

xMv

HE9x

# Experiment

Nvidia Quadro K6000



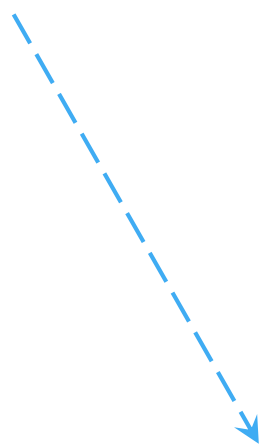
# Experiment

Conditions of panoramic images

Immersive high-resolution screens



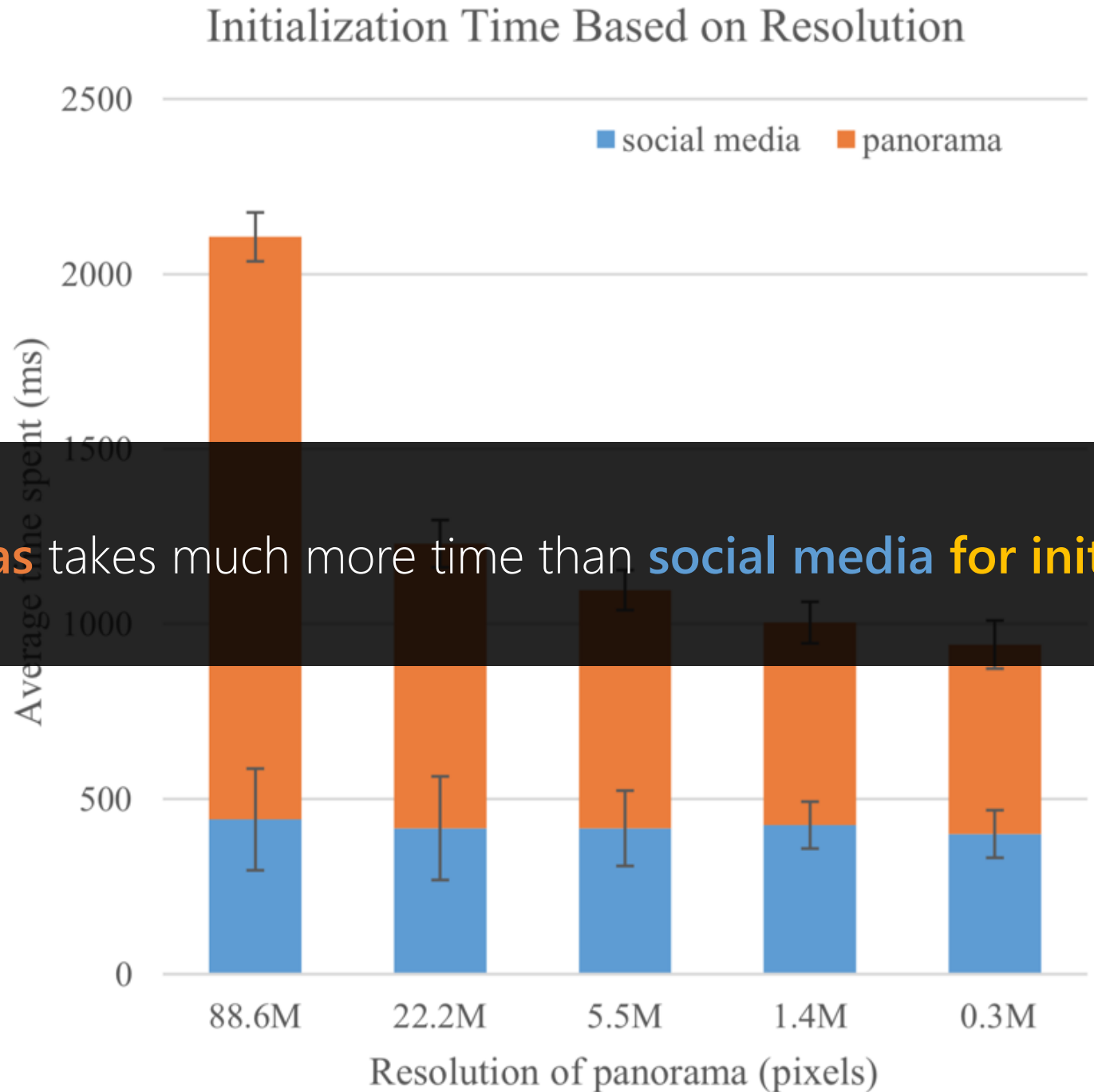
Pixels	Resolution	Number of tiles	File size
88.6M	13312 × 6656	26 × 13	~ 5M
22.2M	6656 × 3328	13 × 7	~ 2M
5.5M	3328 × 1664	7 × 4	~ 800K
1.4M	1664 × 832	4 × 2	~ 300K
0.3M	832 × 416	2 × 1	~ 90K



Common Consumer-level Displays

# Initialization Time

Panorama takes a while to load

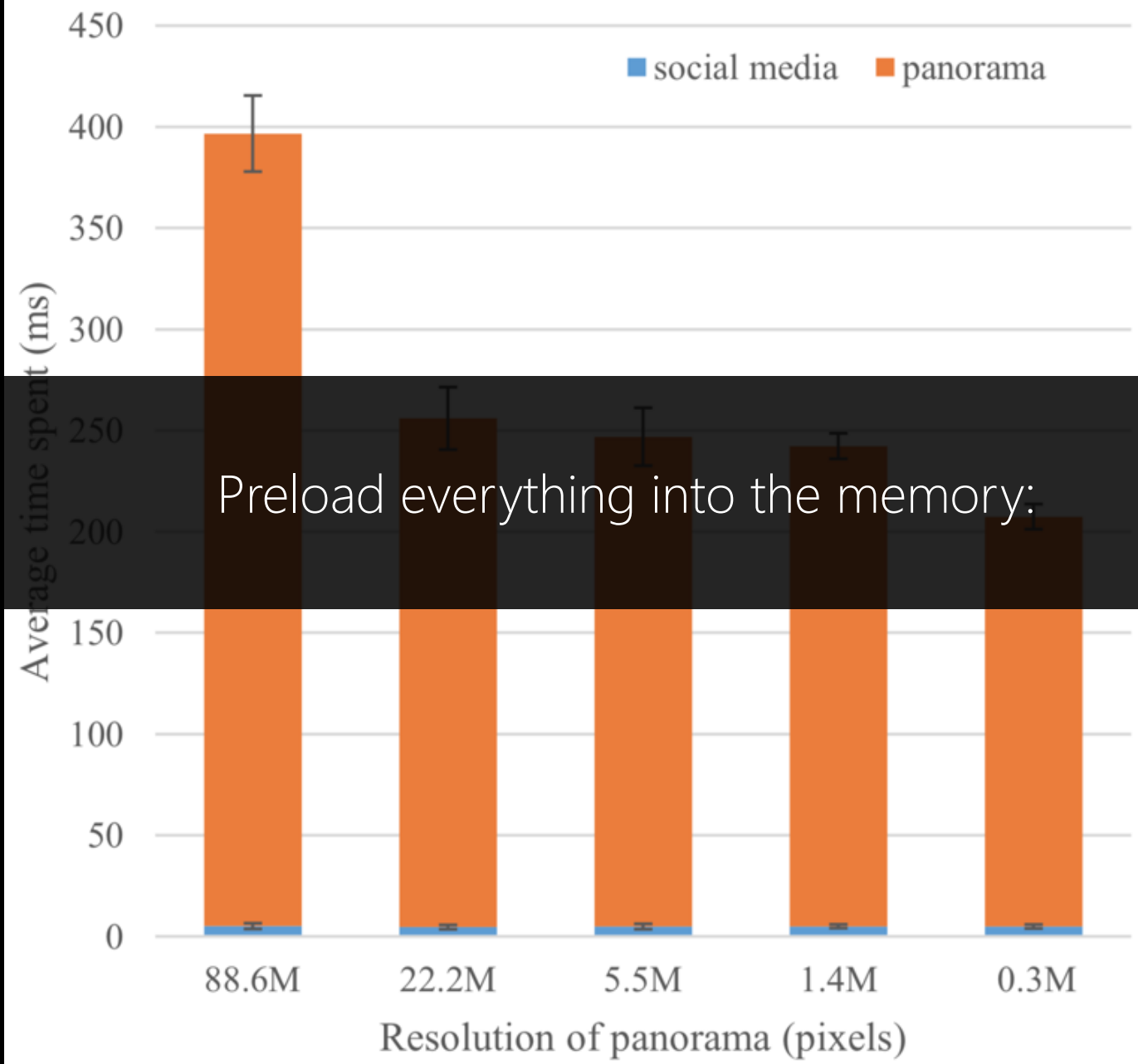


**Panoramas** takes much more time than **social media** for initialization.

# After Prefetching

¾ - ⅘ time reduced

## Initialization Time After Prefetching

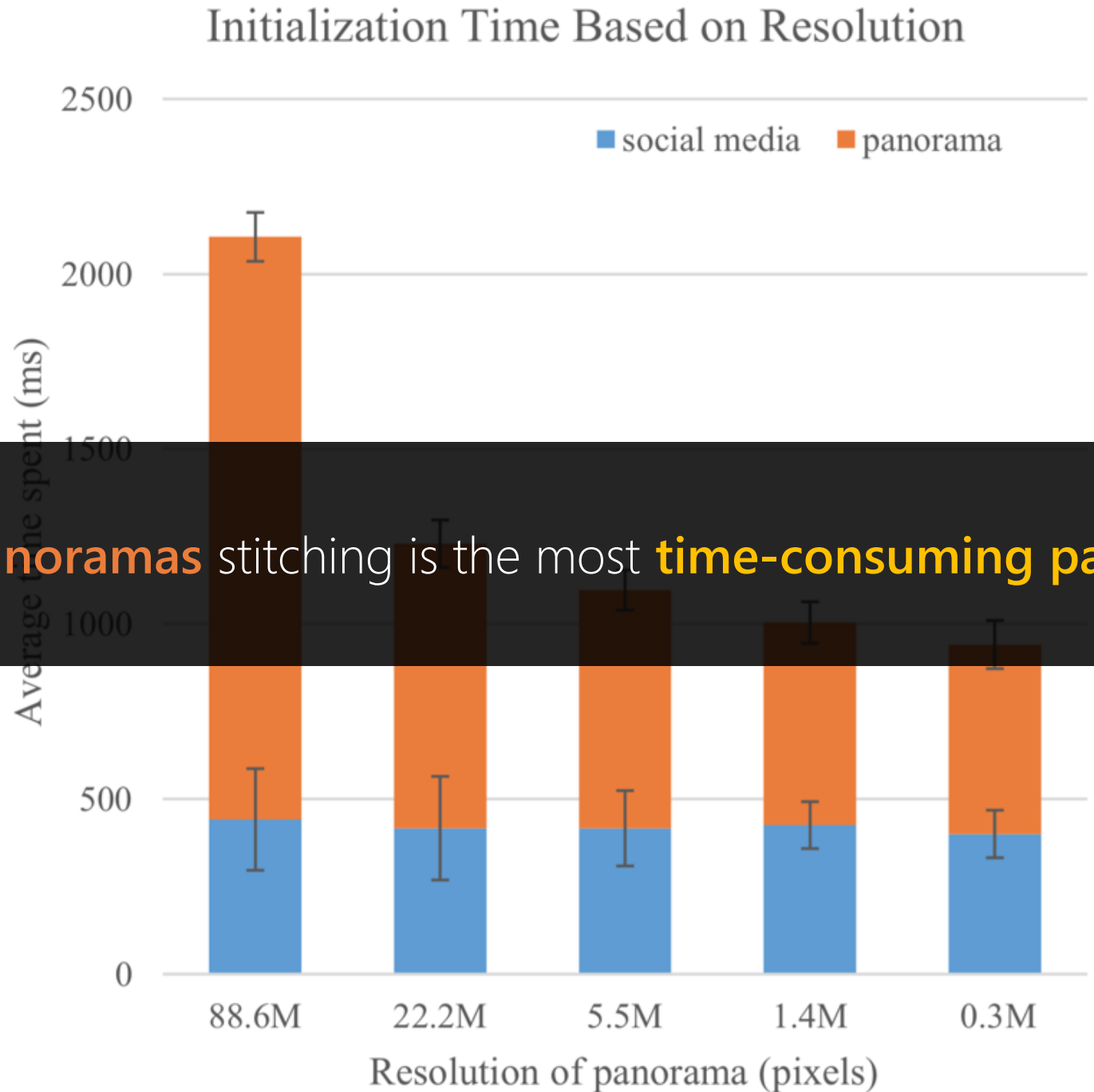


Preload everything into the memory:

~250ms

# Initialization Time

Panorama takes a while to load

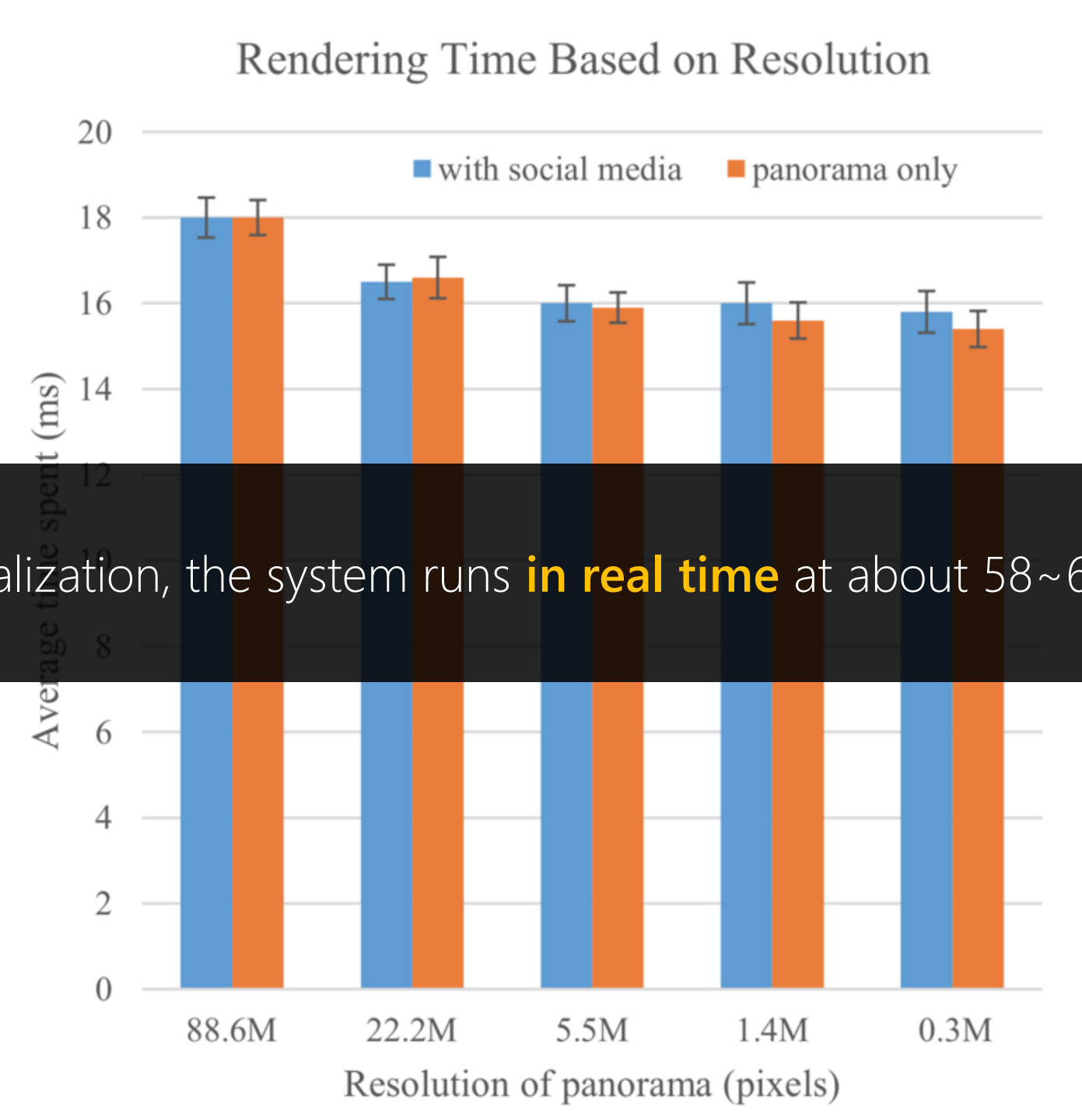


**Panoramas** stitching is the most **time-consuming part**.



# Rendering Time

Almost 58~60 FPS



After initialization, the system runs **in real time** at about 58~60 FPS.

**How effective** is the Maximal Poisson-disk Sampling for reducing the visual clutter?

# Saliency Map

Hou et al. TPAMI 2011

**How much** social media **covers** the salient regions of the panorama?



# Saliency Map

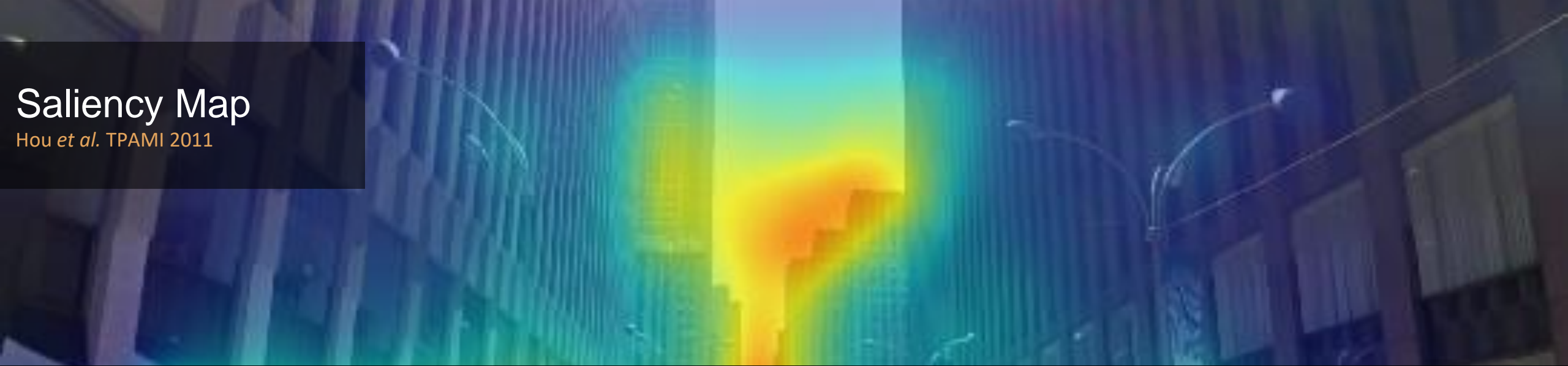
Hou et al. TPAMI 2011

**Saliency maps** use color, intensity, and orientation contrasts to estimate **where humans will look at**.



# Saliency Map

Hou et al. TPAMI 2011

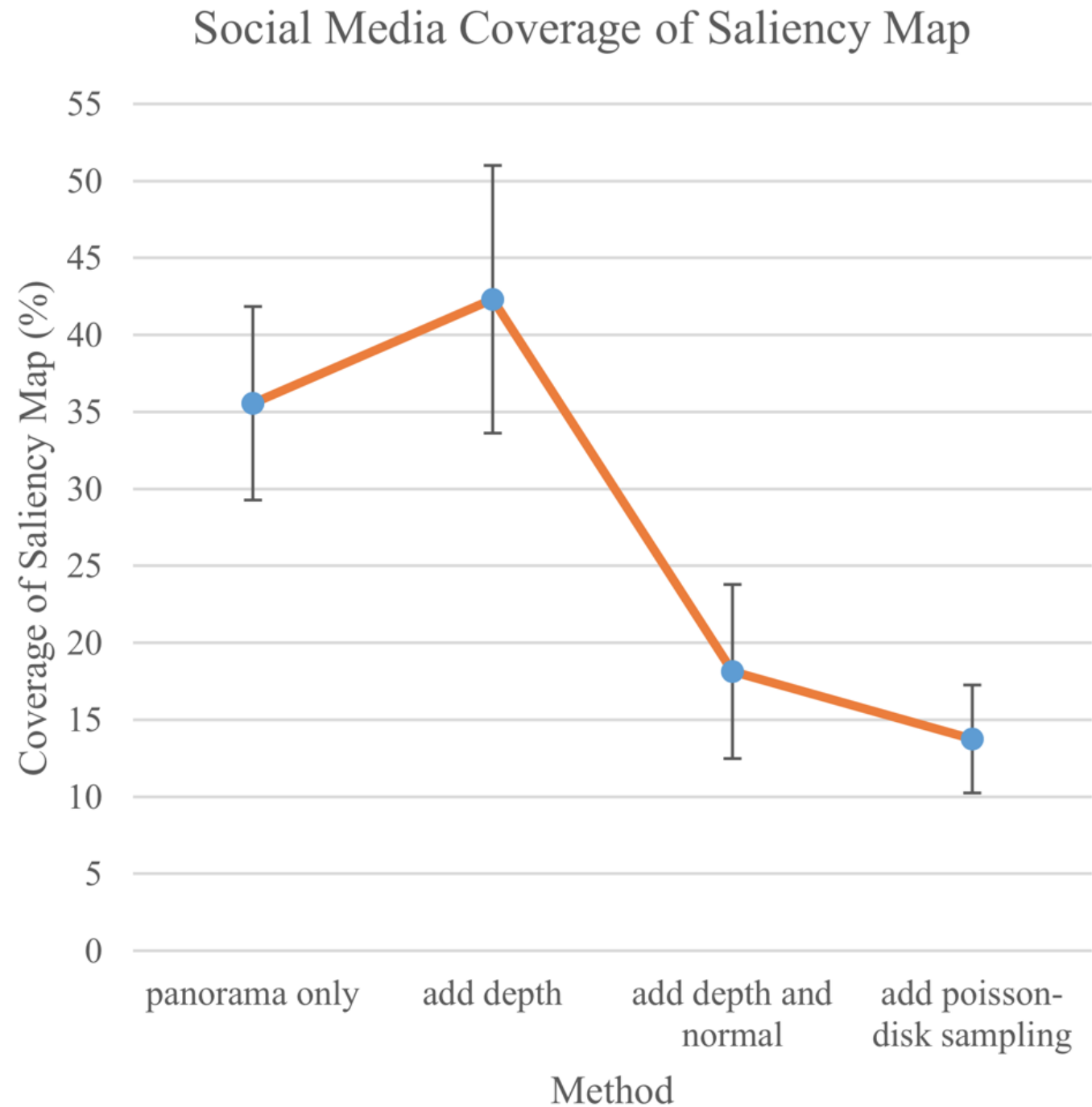


We prefer the social media to cover **less high-saliency regions**.



# Social Media Coverage

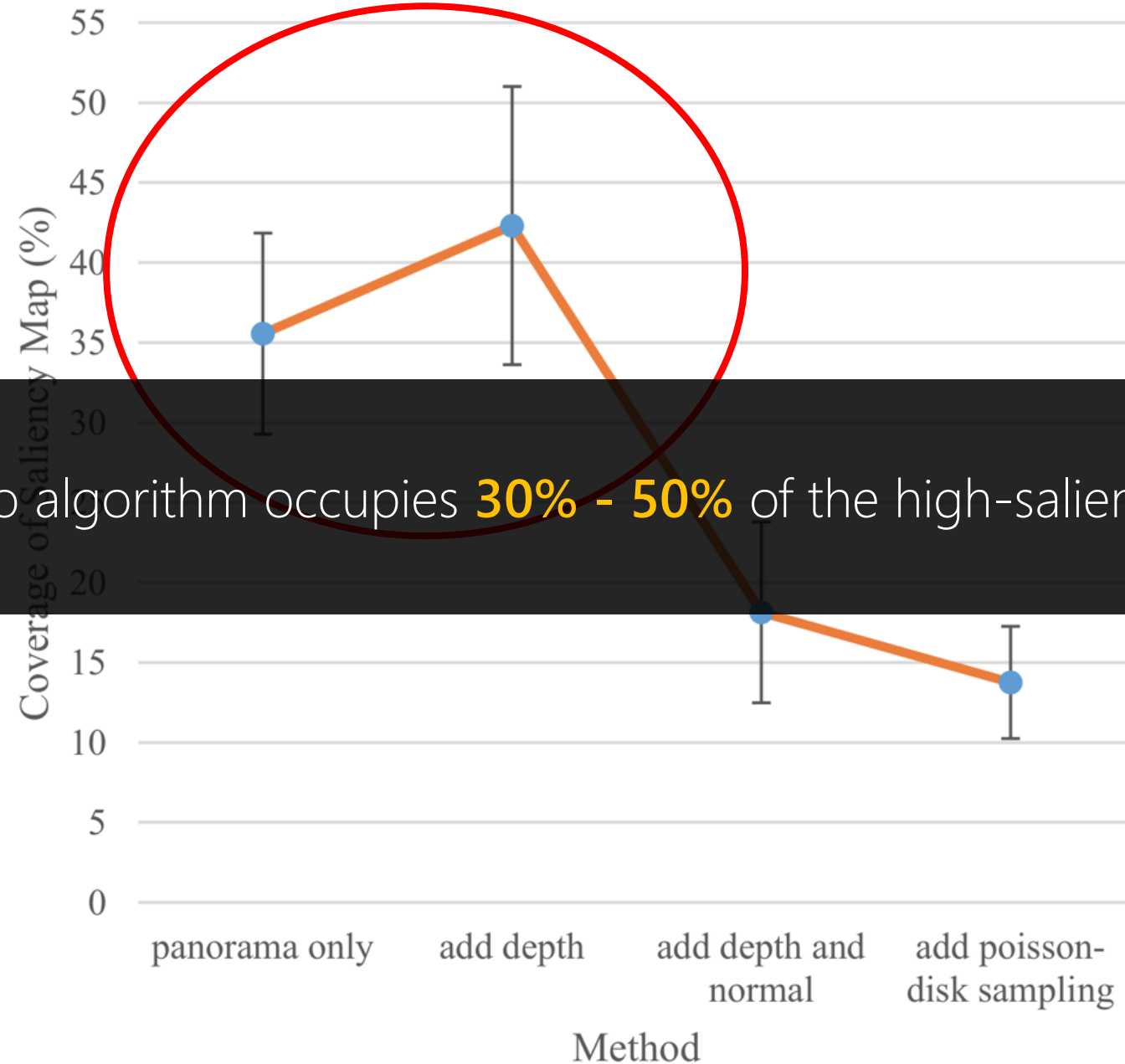
100 panoramas for each algorithm



# Social Media Coverage

100 panoramas for each algorithm

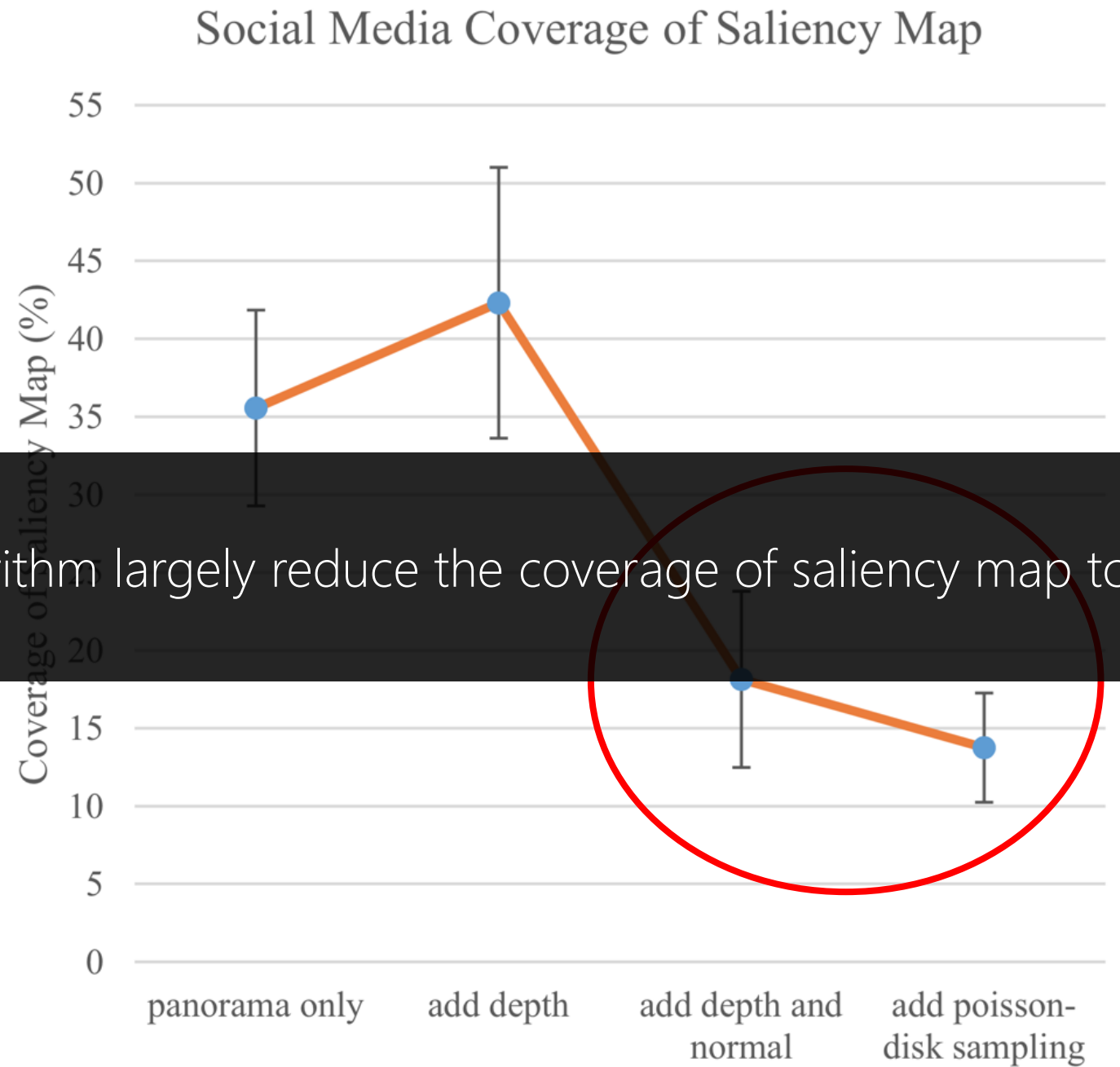
## Social Media Coverage of Saliency Map



The first two algorithm occupipes 30% - 50% of the high-saliency regions.

# Social Media Coverage

100 panoramas for each algorithm



The last two algorithm largely reduce the coverage of saliency map to **less than 20%**.

In the last algorithm, all the images are separated with each other and uniformly distributed on the building surfaces.



What could be the **potential applications** of Social Street View?



“

Stuck in traffic on our way to  
Cabo with this awesome view

#roadtrip #cabo #view #mexico

”

Daniela on *Instagram*  
July 12, 2014



# Application

Immersive story telling



However, we can hardly enjoy the view given the small posted image.

# Business Advertising

Museum, restaurant, real-estate ...



“

... dinner started off with  
*amazing oysters* paired with my  
favorite Ruinart blanc de blancs  
champagne

”

By frankiextah on *Instagram*

# Business Advertising

Museum, restaurant, real-estate ...

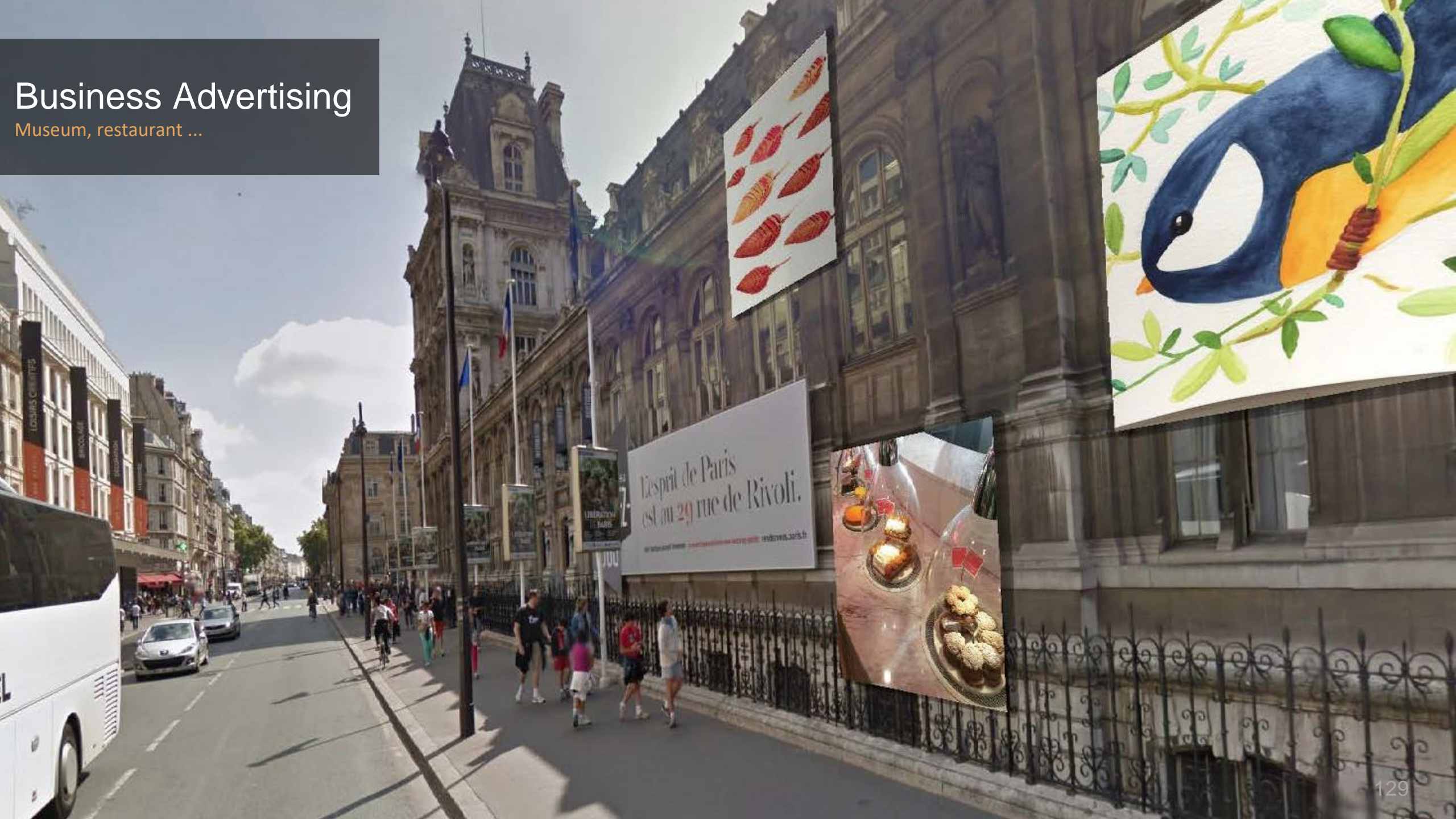


Social Street View enhances future consumers' **visual memories** and makes it easier for them to seek the "amazing oysters" around this place.



# Business Advertising

Museum, restaurant ...



L'esprit de Paris  
est au 29 rue de Rivoli.

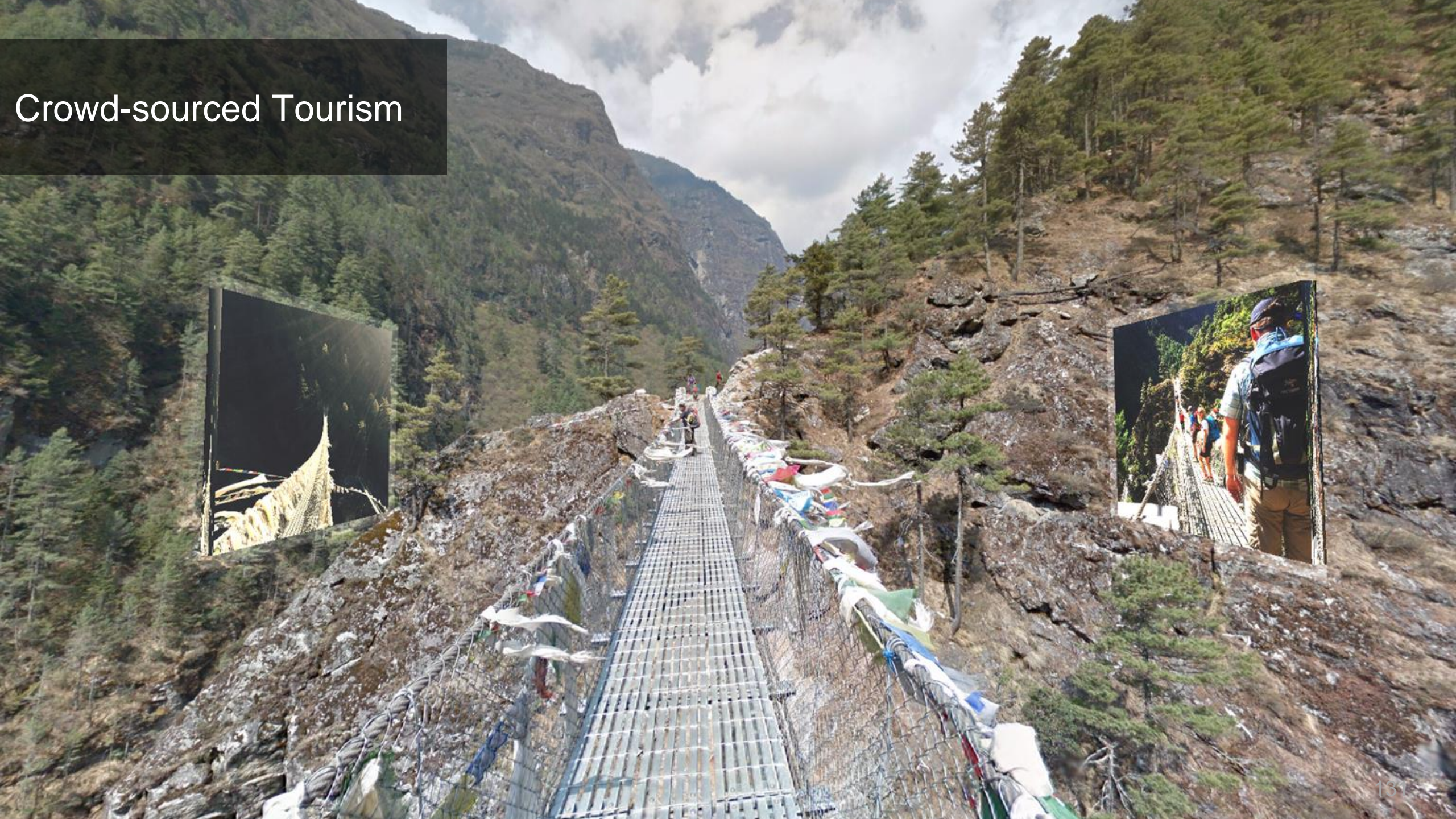


# Learning Culture

Taierzhuang, Chinese Spring Festivals



# Crowd-sourced Tourism



## Crowd-sourced Tourism



Social Street View allows users to explore **novel views** (e.g. from top of buildings), where you cannot see with only the panoramas.



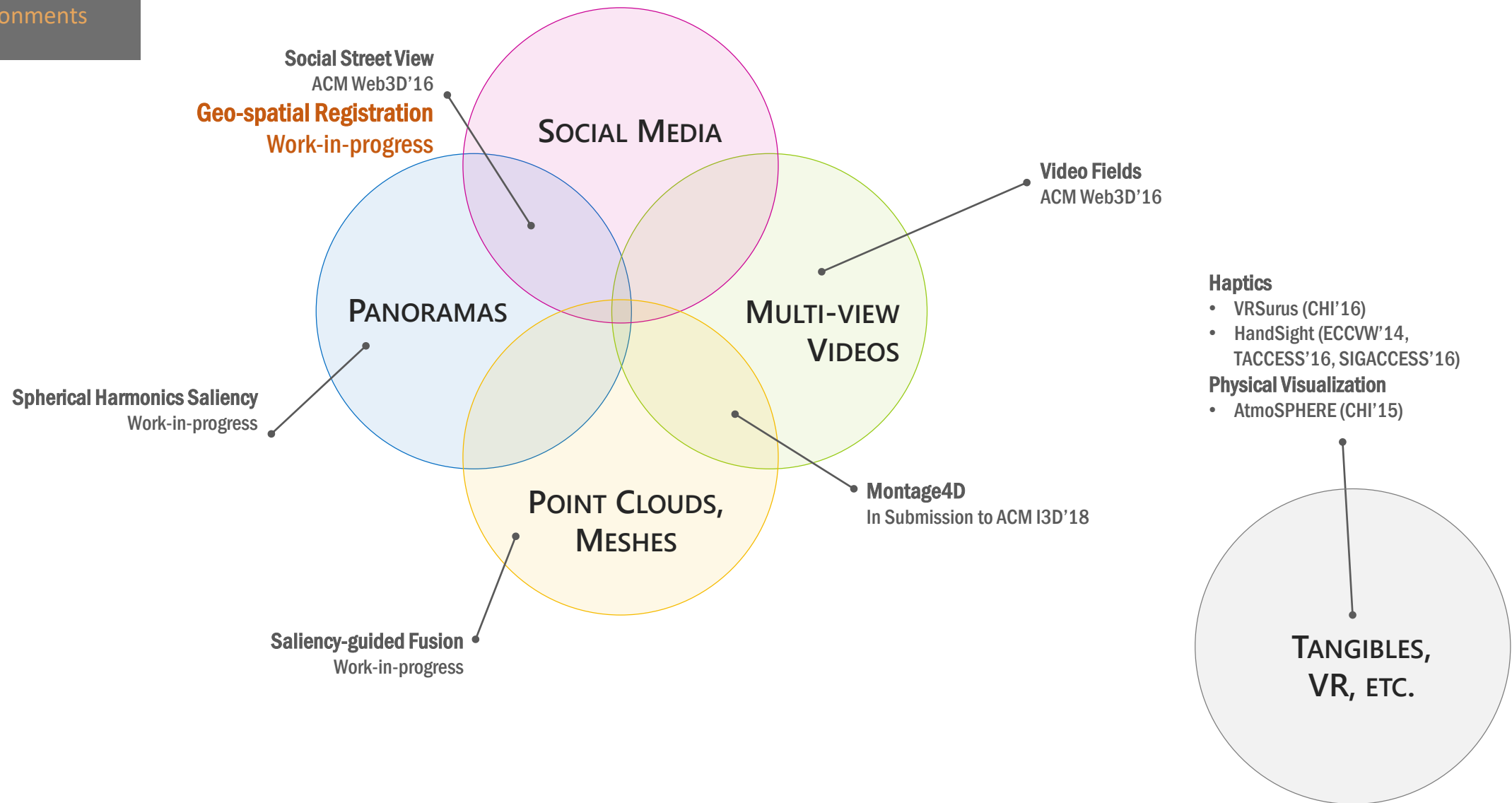
In conclusion, **Social Street View** provides a novel solution for automatically fusing geo-tagged social media in an immersive 3D environments.



We envision **social media** and **panoramic videos** are significant parts of the **big data** for VR and AR applications

# Proposal

Fusing Multimedia Data Into  
Dynamic Virtual Environments



# Challenges

Registration of billboards



The social media layout is **re-generated**, every time, when users walk around the street views.

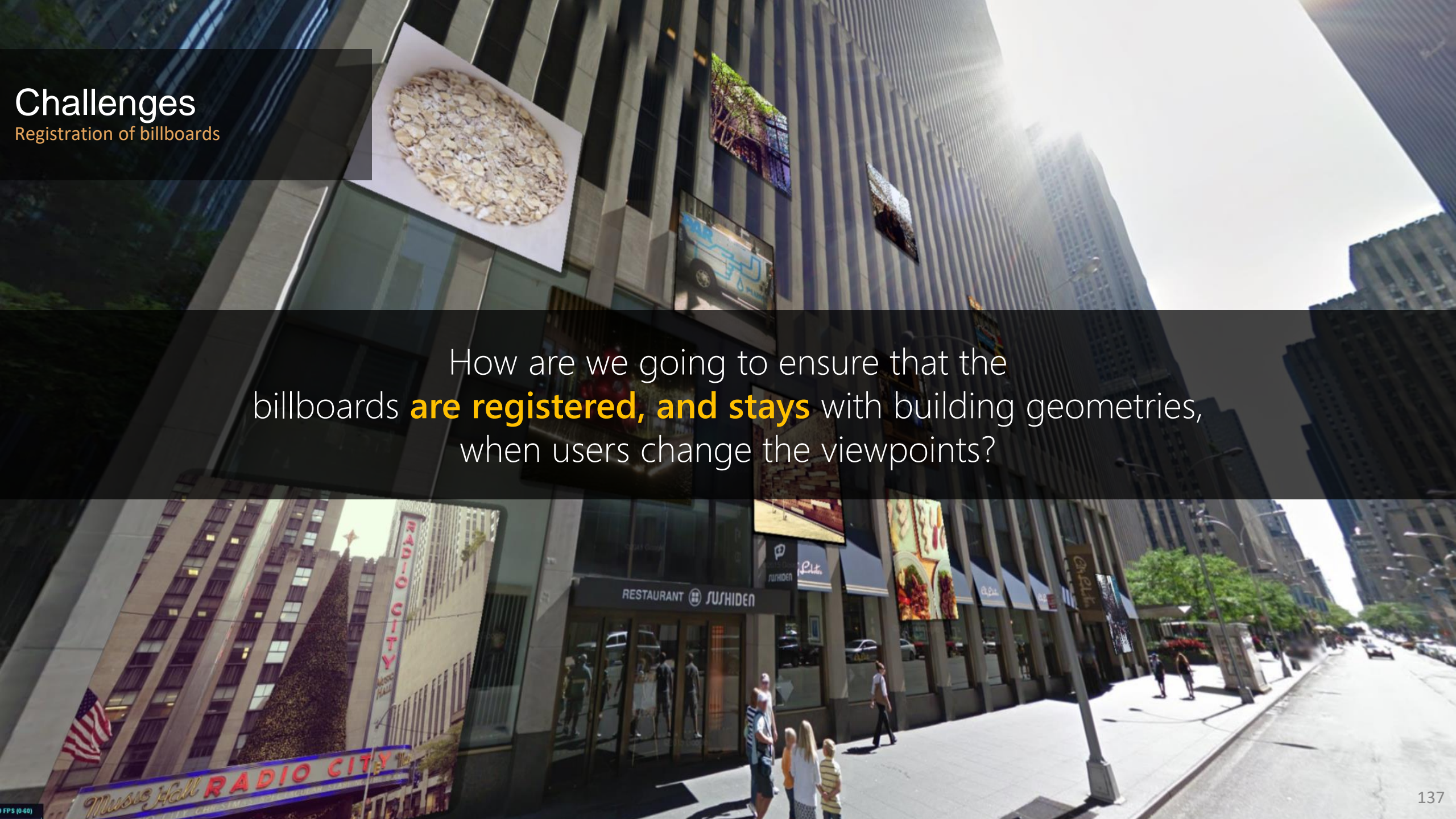




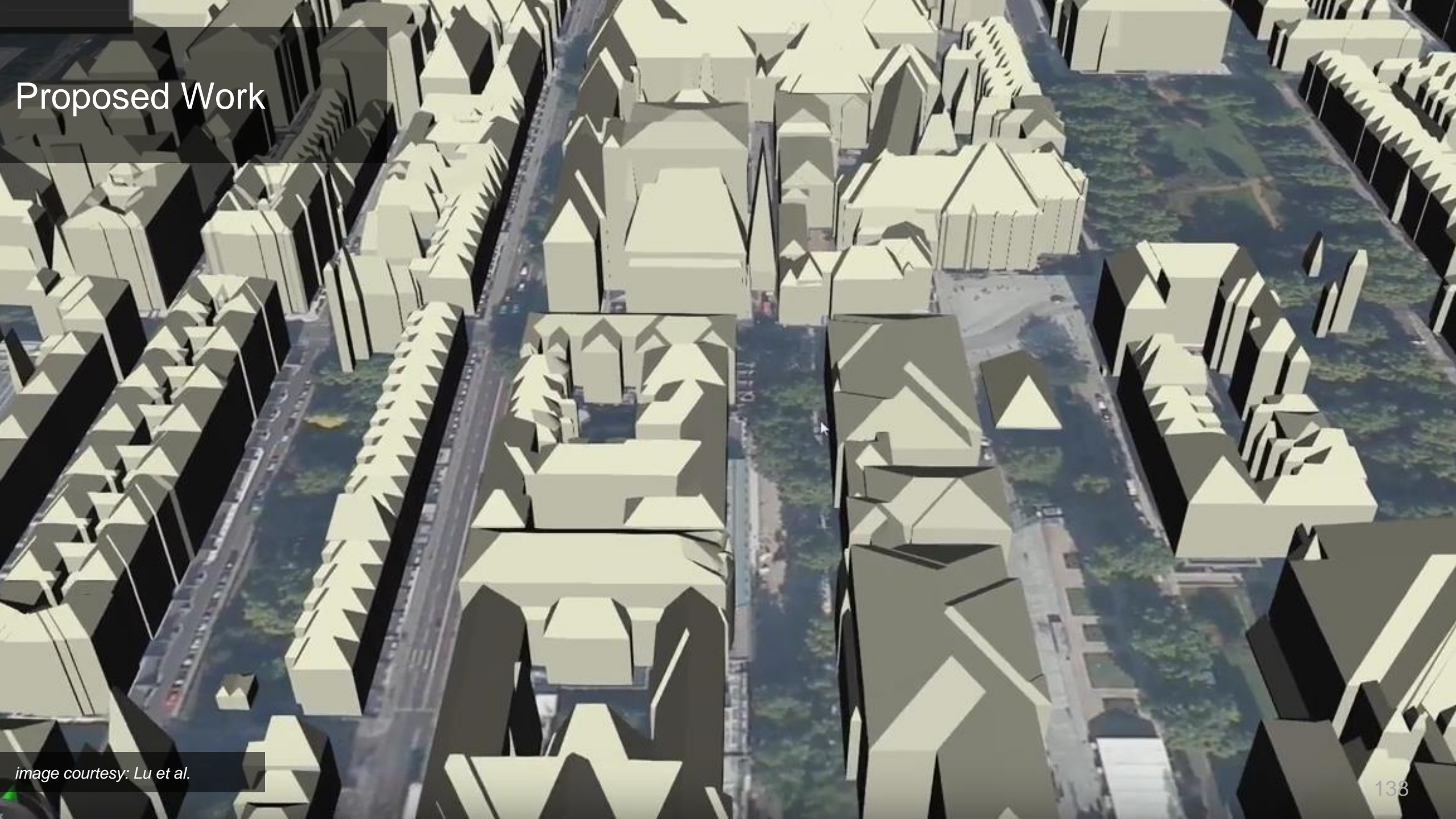
# Challenges

Registration of billboards

How are we going to ensure that the billboards **are registered, and stays** with building geometries, when users change the viewpoints?



# Proposed Work



An aerial view of a city with 3D building models overlaid on a satellite map. The buildings are rendered in a light beige color, and the streets and surrounding areas are visible in the background. The models are semi-transparent, allowing the underlying satellite imagery to be seen through them.

## Proposed Work

Estimate simple **3D building geometries** from the depth and normal maps.

# Visual Registration

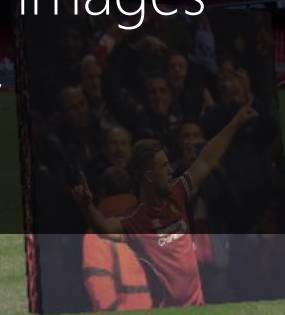
Image Feature



# Visual Registration

Image Feature

How to register some of the social media images  
**with the immersive panoramas,**  
to create a novel layout?



# Storytelling

Baja California Sur, Mexico



# Storytelling

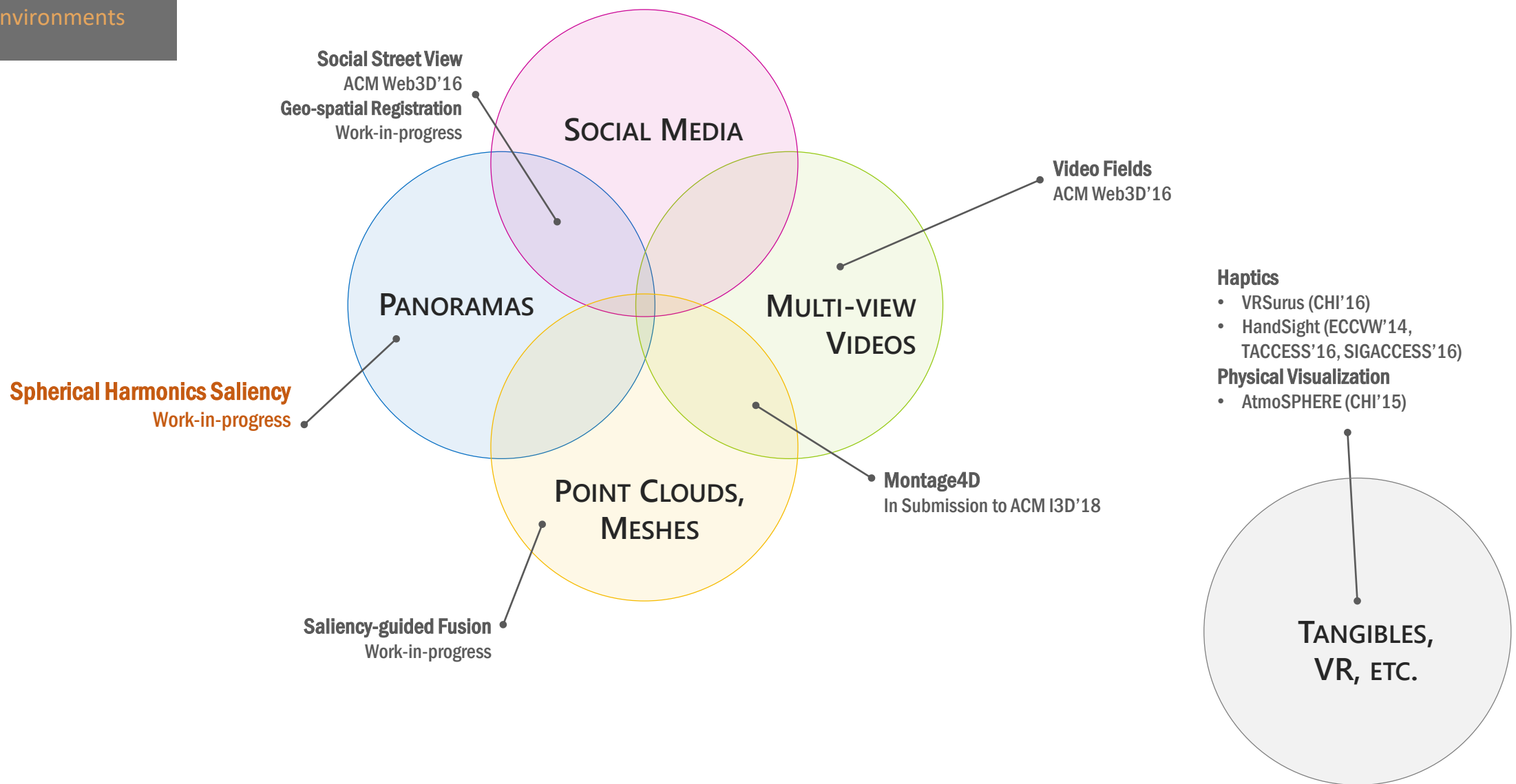
Baja California Sur, Mexico

Develop fast algorithms for  
**reconstructing simple geometries,**  
**feature matching,**  
and **navigating the social street views without resampling.**



# Proposal

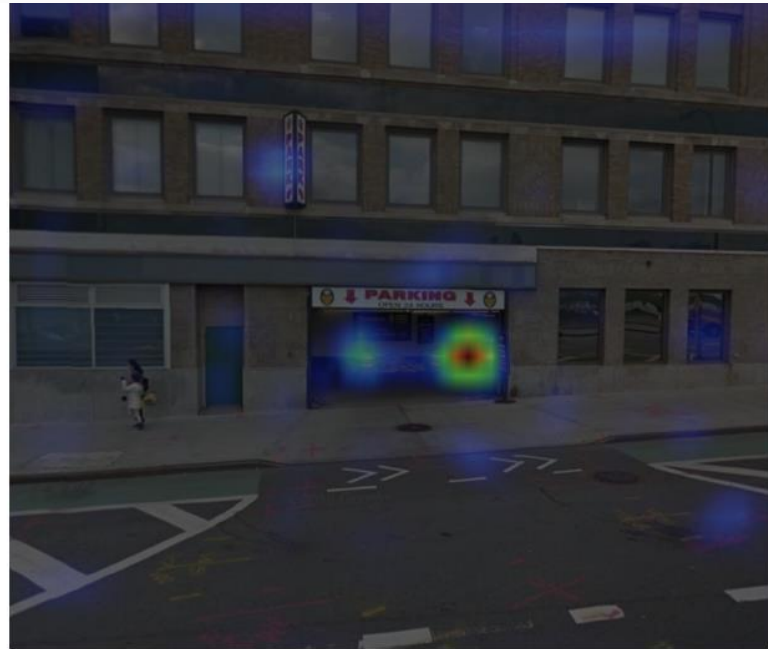
Fusing Multimedia Data Into Dynamic Virtual Environments





# Motivation

## Social Media Overlay



How to efficiently generate the **saliency map for panoramas?**

# Motivation

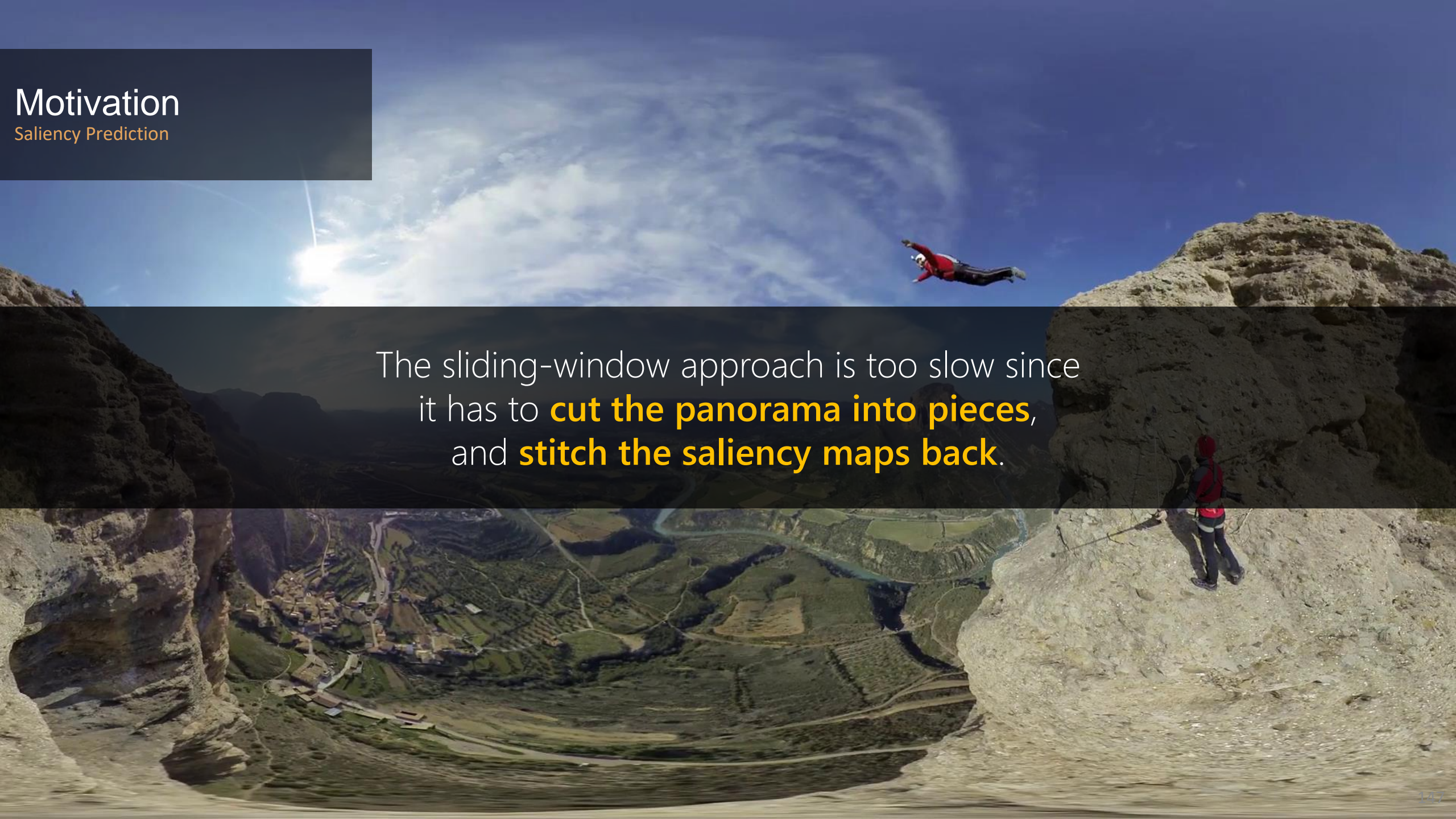
Saliency Prediction

A wide-angle, panoramic photograph of a valley. In the foreground, a person in a red jacket and black pants is skydiving, floating horizontally against a blue sky with wispy clouds. Below, a winding river flows through a lush green valley. A small village with several buildings is visible in the distance. The valley is framed by steep, rocky cliffs on both sides. The overall scene is bright and clear, with a strong sense of depth and perspective.

Traditional saliency methods can hardly deal with **spherical rotations** and **horizontal clipping**

# Motivation

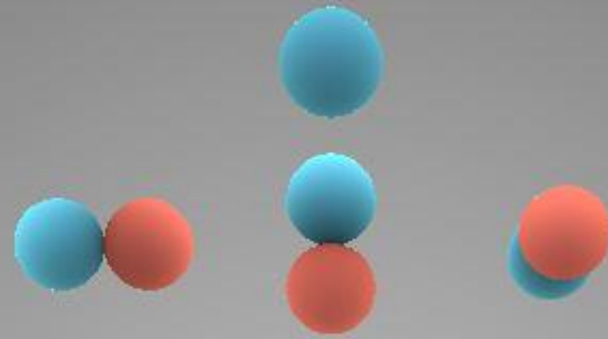
Saliency Prediction



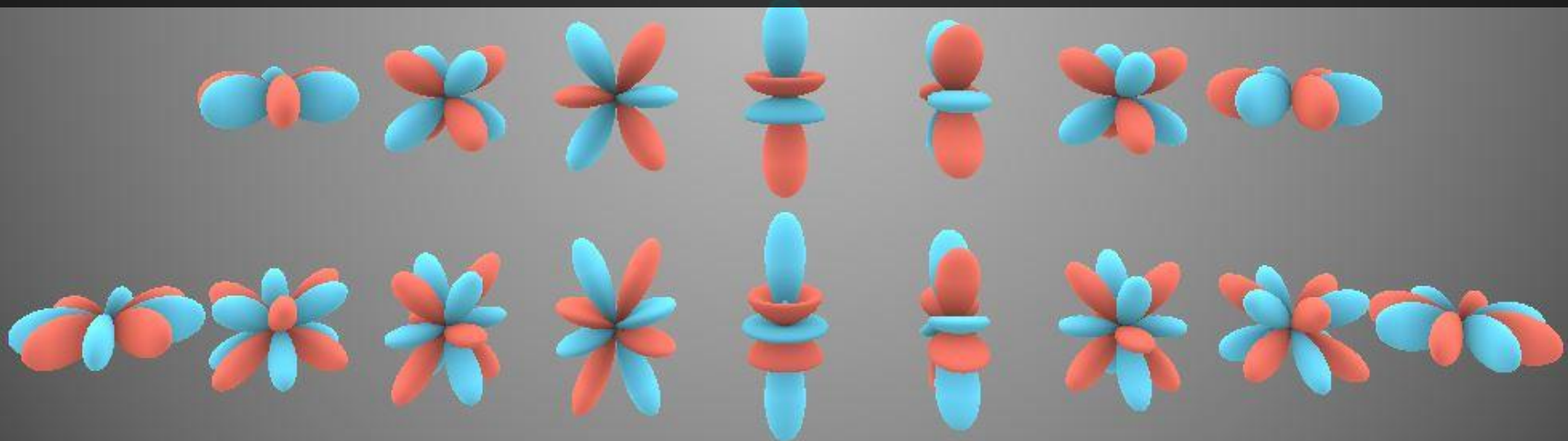
The sliding-window approach is too slow since it has to **cut the panorama into pieces**, and **stitch the saliency maps back**.

# Spherical Harmonics

Orthogonal function on the sphere  
Du. [www.shadertoy.com/starea](http://www.shadertoy.com/starea)

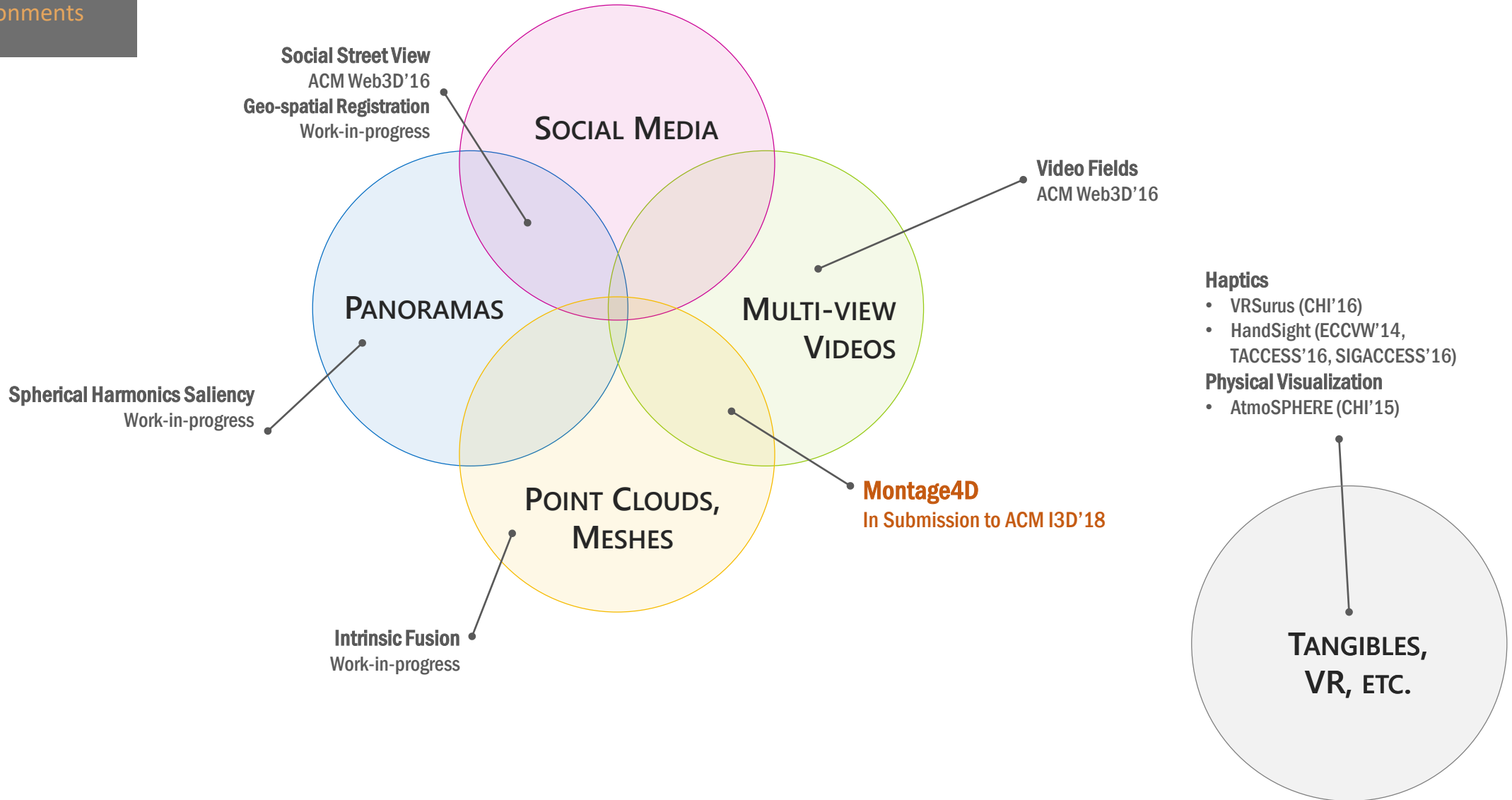


Spherical Harmonics, a complete set of orthogonal functions on the sphere, can be used for estimating the **panoramic** saliency maps **in one pass**.



# Proposal

Fusing Multimedia Data Into Dynamic Virtual Environments



# Montage4D: Interactive Seamless Fusion of Multiview Textures



Ruofei Du, Ming Chuang, Wayne Chang, Hugues Hoppe, and Amitabh Varshney

University of Maryland, College Park

Microsoft Research, Redmond



THE AUGMENTARIUM  
VIRTUAL AND AUGMENTED REALITY LABORATORY  
AT THE UNIVERSITY OF MARYLAND

Microsoft

Research

UMIACS



COMPUTER SCIENCE  
UNIVERSITY OF MARYLAND

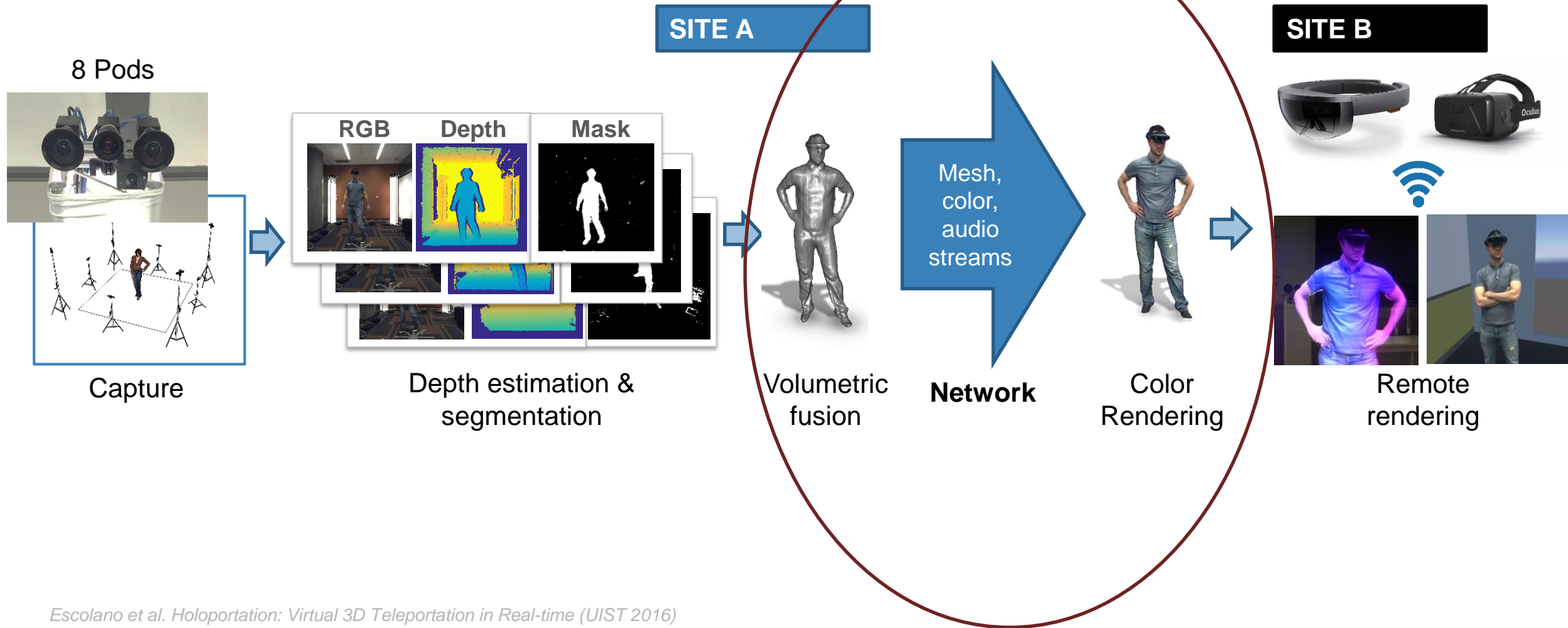


UNIVERSITY OF  
MARYLAND



# Introduction

Fusion4D and Holoportation





Fusing multiview video textures onto dynamic meshes  
with real-time constraint is **a challenging task**

# 30%

of the participants does not believe the 3D reconstructed person looks real

# Motivation

Visual Quality Matters



# Motivation

Visual Quality Matters



# Related Work

3D Texture Montage

## Color Map Optimization for 3D Reconstruction with Consumer Depth Cameras

Qian-Yi Zhou\*

Vladlen Koltun†

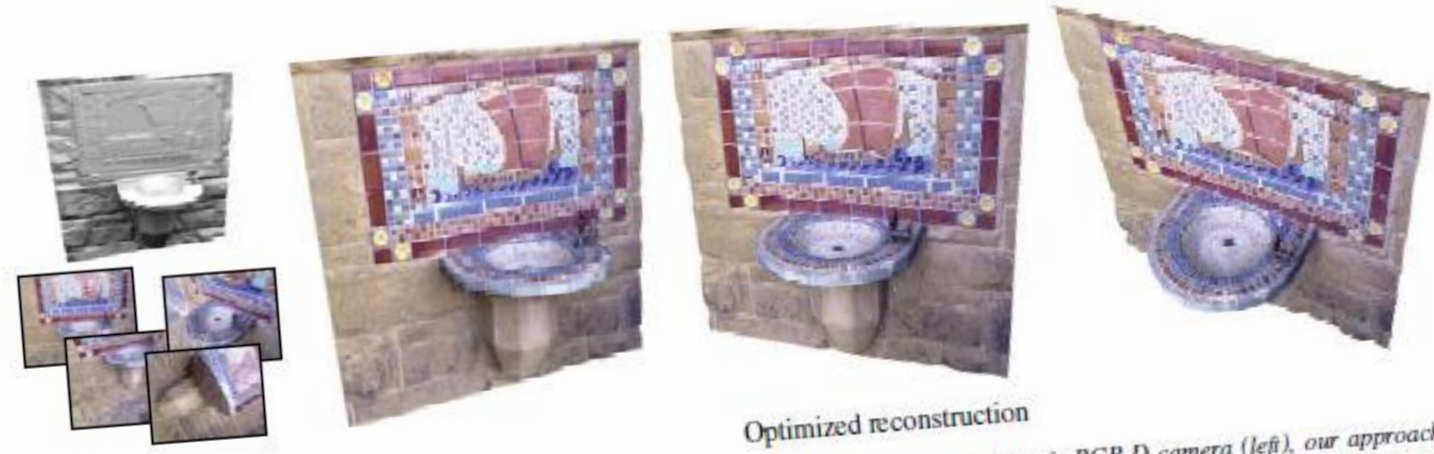


Figure 1: Given a geometric model and corresponding color images produced by a consumer-grade RGB-D camera (left), our approach optimizes a photometrically consistent mapping of the images to the model.

### Abstract

We present a global optimization approach for mapping color images onto geometric reconstructions. Range and color videos produced by consumer-grade RGB-D cameras suffer from noise and optical distortions, which impede accurate mapping of the acquired color data to the reconstructed geometry. Our approach addresses these sources of error by optimizing camera poses in tandem with non-rigid correction functions for all images. All parameters are optimized jointly to maximize the photometric consistency of the reconstruction. This optimization can be performed on a GPU.

### 1 Introduction

Consumer depth cameras are now widely available. Tens of millions of such cameras have been shipped and miniaturized versions are being developed for integration into laptops, tablets, and smartphones. As a result, millions of people can now create high-fidelity geometric models of real-world objects and scenes [Newcombe et al. 2011; Chen et al. 2013; Zhou and Koltun 2013; Zhou et al. 2013].

However, capturing an object's geometry is not sufficient for reproducing its appearance. A visually faithful reconstruction must also incorporate the apparent color of every point on the object. In this paper, we present a global optimization approach for color map optimization for 3D reconstruction systems.

# Related Work

3D Texture Montage

4.02801v1 [cs.CV] 11 Apr 2016

Volume 0 (1981), Number 0 pp. 1-7

## Seamless Montage for Texturing Models

Ran Gal<sup>1</sup> Yonatan Wexler<sup>2</sup> Eyal Ofek<sup>2</sup> Hugues Hoppe<sup>2</sup> Daniel Cohen-Or<sup>1</sup>

<sup>1</sup>Tel-Aviv University, Tel-Aviv, Israel  
<sup>2</sup>Microsoft, Redmond, USA

### Abstract

We present an automatic method to recover high-resolution texture over an object by mapping detailed photographs onto its surface. Such high-resolution detail often reveals inaccuracies in geometry and registration, as well as lighting variations and surface reflections. Simple image projection results in visible seams on the surface. We minimize such seams using a global optimization that assigns compatible texture to adjacent triangles. The key idea is to search not only combinatorially over the source images, but also over a set of local image transformations that compensate for geometric misalignment. This broad search space is traversed using a discrete labeling algorithm, aided by a coarse-to-fine strategy. Our approach significantly improves resilience to acquisition errors, thereby allowing simple and easy creation of textured models for use in computer graphics.

Categories and Subject Descriptors (according to ACM CCS): I.3.7 [Computer Graphics]: Three-Dimensional Graphics and Realism—Color, shading, shadowing, and texture

### 1. Introduction

Our goal is to generate a seamless texture over a surface from a set of photographs taken with an ordinary handheld camera from different views around the object.

Our approach favors smoothness in the texture assignment and penalizes the introduction of sharp seams. Figure 1(c) shows results of their approach, which indeed alleviates many artifacts. However, as shown in Figure 1(c), some seams often remain due to misregistration and imperfect geometry. The errors in Figure 1(b,c) suggest that it is unlikely that such techniques alone (e.g. graph-cut, MRF) can solve the problem of stitching a set of photographs into a seamless texture.

# Related Work

3D Texture Montage

4.02801v1 [cs.CV] 11 Apr 2016

Volume 0 (1981), Number 0 pp. 1-7

## Seamless Montage for Texturing Models

Ran Gal<sup>1</sup> Yonatan Wexler<sup>2</sup> Eyal Ofek<sup>2</sup> Hugues Hoppe<sup>2</sup> Daniel Cohen-Or<sup>1</sup>

<sup>1</sup>Tel-Aviv University, Tel-Aviv, Israel  
<sup>2</sup>Microsoft, Redmond, USA

Dataset name	Input photographs	Mesh triangles	Processing time (minutes)
Doll	9	12,000	15
Striped cat	8	13,500	15
Monkey	12	15,000	15
Brown cat	15	18,500	30
Yellow car	16	17,000	25
Tall cat	10	11,000	7

### 1. Introduction

Our goal is to generate a seamless texture over a surface from a set of photographs taken with an ordinary handheld camera from a set of views around the object.

mization that favors smoothness in the texture assignment and penalizes the introduction of sharp seams. Figure 1(c) shows results of their approach, which indeed alleviates many artifacts. However, as shown in Figure 1(c), some seams often remain due to misregistration and imperfect geometry. The errors in Figure 1(b,c) suggest that it is unlikely that such techniques alone (e.g. graph-cut, MRF) can produce a seamless texture from a set of photographs.

# Related Work

3D Texture Montage

4.02801v1 [cs.CV] 11 Apr 2016

Volume 0 (1981), Number 0 pp. 1-7

## Seamless Montage for Texturing Models

Ran Gal<sup>1</sup> Yonatan Wexler<sup>2</sup> Eyal Ofek<sup>2</sup> Hugues Hoppe<sup>2</sup> Daniel Cohen-Or<sup>1</sup>

<sup>1</sup>Tel-Aviv University, Tel-Aviv, Israel  
<sup>2</sup>Microsoft, Redmond, USA

Dataset name	Input photographs	Mesh triangles	Processing time (minutes)
Doll	9	12,000	15
Striped cat	8	13,500	15
Monkey	12	15,000	15
Brown cat	15	18,500	30
Yellow car	16	17,000	25
Tall cat	10	11,000	7

### 1. Introduction

Our goal is to generate a seamless texture over a surface from photographs taken with an ordinary handheld camera from a set of views around

mization that favors smoothness in the texture assignment and penalizes the introduction of sharp seams. Figure 1(c) shows results of their approach, which indeed alleviates many artifacts. However, as shown in Figure 1(c), some seams often remain due to misregistration and imperfect geometry. The errors in Figure 1(b,c) suggest that it is unlikely that techniques alone (e.g. graph-cut,)



# Related Work

3D Texture Montage

## High-Quality Streamable Free-Viewpoint Video

Alvaro Collet Ming Chuang Pat Sweeney Don Gillett Dennis Evseev David Calabrese  
Hugues Hoppe Adam Kirk Steve Sullivan  
Microsoft Corporation



Figure 1: Examples of reconstructed free-viewpoint video acquired by our system.

### Abstract

We present the first end-to-end solution to create high-quality free-viewpoint video encoded as a compact data stream. Our system records performances using a dense set of RGB and IR video cameras, generates dynamic textured surfaces, and compresses these to a streamable 3D video format. Four technical advances contribute to high fidelity and robustness: multimodal multi-view stereo fusing RGB, IR, and silhouette information; adaptive meshing guided by automatic detection of perceptually salient areas; mesh tracking to create temporally coherent subsequences; and encoding of tracked textured meshes as an MPEG video stream. Quantitative metrics include accuracy, texture fidelity, and

Our goal is to transform free-viewpoint video from research prototype into a rich and accessible form of media. Several system components must work together to achieve this goal: capture rigs must be easy to reconfigure and support professional production workflows; reconstruction must be automatic and scalable to high processing throughput; and results must be compressible to a data rate close to common media formats. Visual quality from any angle must be on par with traditional video, and the format must be renderable in real-time on a wide range of consumer devices.

In this paper, we discuss how we address these challenges to create an end-to-end system for realistic, streamable free-viewpoint video at significantly higher quality than the state of the art. Our approach does not require prior knowledge of the scene content. It handles detailed hand

Volume 0 (1981), Number 0 pp. 1-7

Search

Database name

Doll

Striped

Monkey

Brown

Yellow

Tall c

1. Introduction

Our goal is to genera

4.02801v1 [cs.CV] 11 Apr 2016

# Related Work

3D Texture Montage

## High-Quality Streamable Free-Viewpoint Video

Alvaro Collet Ming Chuang Pat Sweeney Don Gillett Dennis Evseev David Calabrese  
 Hugues Hoppe Adam Kirk Steve Sullivan  
 Microsoft Corporation



Figure 1: Examples of reconstructed free-viewpoint video.

Scene	Frames	Avg points	Avg time/frame (s)	Avg frames/keyframe	Output Mbps
D.Duo	886	2.25M	29.1	24.6	12.1
Dress	1157	1.44M	27.7	27.5	8.6
Kendo	740	2.01M	26.0	35.2	8.3
Haka	173	2.70M	28.2	57.7	12.0
Lincoln	508	1.60M	27.5	72.6	7.9

**Abstract**  
 We present a system for capturing free-viewpoint video records per scene. The system generates a streamable free-viewpoint video to high fidelity using RGB, IR, and depth data by automatically tracking and creating textured 3D models.

om research pro-  
 Several system  
 goal: capture rig  
 ional production  
 scalable to high  
 resible to a data  
 lity from any an-  
 e format must be  
 er devices.

allenges to create  
 free-viewpoint video  
 Our approach  
 does not require prior knowledge of the scene content. It handles  
 detailed hand

4.02801v1 [cs.CV] 11 Apr 2016

Volume 0 (1981), Number 0 pp. 1-7

Search

Rate

Dataset name

- Doll
- Striped
- Monkey
- Brown
- Yellow
- Tall c

### 1. Introduction

Our goal is to generate

# Related Work

## 3D Texture Montage



Screen-space optical flow could fix some misregistration issues, but heavily relies on **RGB features**, and fails when **changing viewpoints**.

# Related Work

3D Texture Montage

Up to now, few systems but *Holoportation* could fuse dynamic meshes **with multiple cameras** in real time. This recent SIGGRAPH 2017 paper produces excellent dynamic reconstruction results, but uses **a single RGBD camera**.

## Real-time Geometry, Albedo and Motion Reconstruction Using a Single RGBD Camera

KAIWEN GUO and FENG XU, Tsinghua University  
TAO YU, Beihang University and Tsinghua University  
XIAOYANG LIU, QIONGHAI DAI, and YEBIN LIU, Tsinghua University

This paper proposes a real-time method that uses a single-view RGBD input to simultaneously reconstruct a casual scene with a detailed geometry model, surface albedo, per-frame non-rigid motion and per-frame low-frequency lighting, without requiring any template or motion priors. The key observation is that accurate scene motion can be used to integrate temporal information to recover the precise appearance, whereas the intrinsic appearance can help to establish true correspondence in the temporal domain to recover motion. Based on this observation, we first propose a shading-based scheme to leverage appearance information for motion estimation. Then, using the reconstructed motion, a volumetric albedo fusing scheme is proposed to complete and refine the intrinsic appearance of the scene by incorporating information from multiple frames. Since the two schemes are iteratively applied during recording, the reconstructed appearance and motion become increasingly more accurate. In addition to the reconstruction results, our experiments also show that additional applications can be achieved, such as relighting, albedo editing and free-viewpoint rendering of a dynamic scene, since geometry, appearance and motion are all reconstructed by our technique.

CCS Concepts: • Computing methodologies → Reconstruction; Motion capture;

Additional Key Words and Phrases: single-view, surface reconstruction, real-time, non-rigid, albedo

ACM Reference format:  
Kaiwen Guo, Feng Xu, Tao Yu, Xiaoyang Liu, Qionghai Dai, and Yebin Liu. 2017. Real-time Geometry, Albedo and Motion Reconstruction Using a Single RGBD Camera. *ACM Trans. Graph.* XX, X, Article XX (March 2017), 13 pages. DOI: XXXXXXXX.YYYYYYY

### 1 INTRODUCTION

Dynamic scene reconstruction involves capturing and reproducing various aspects of the real visual world, including static geometry, detailed motion, and intrinsic or observed appearance. Simultaneously reconstructing all of these aspects or even part of them enables

This work was supported by the National key foundation for exploring scientific instrument No. 2013YQ140517 and NSFC (No. 61671268, 61522111 and 61531014). Authors' addresses: K. Guo, X. Liu; email: {gkw11, liu-xy14}@mails.tsinghua.edu.cn; Q. Dai, Y. Liu (corresponding author); emails: {qhldai, liuyebin}@tsinghua.edu.cn; Department of Automation and TNLi, Tsinghua University, Beijing, 100084, China; F. Xu (corresponding author); email: fengxu@sem.tsinghua.edu.cn; Software and TNLi, Tsinghua University, Beijing, 100084, China; Y. Liu (corresponding author); email: yebliu@sem.tsinghua.edu.cn



Fig. 1. Our system can capture fast and natural motions, geometry, and face albedo and simultaneously render them in new lighting environments in real time.

important applications in computer vision and graphics. For example, reconstructed geometry and surface motion as well as observed appearance can be used for free-viewpoint video. Reconstructed kinematic motion can be transferred to new objects or used to generate new photo-realistic animations. The intrinsic appearance of a dynamic scene/object can be used in applications such as appearance editing and relighting. The real-time reconstruction of geometry, motion and appearance enables more realistic rendering of virtual reality scenarios, for example, Holoportation [7].

Although considerable efforts have been devoted to dynamic scene reconstruction, the problem remains challenging because of the extraordinarily large solution space, necessitating a carefully designed capture environment [5, 34], high-quality lighting equipment [3, 9] and many video cameras [15]. Several recent works have successfully eliminated various constraints on acquisition by using convenient capture equipment, such as a single Kinect [32] or binocular camera [37]. However, they require many cameras to constrain the problem space.

# Related Work

3D Texture Montage

## Holoportation: Virtual 3D Teleportation in Real-time

S. Orts Escolano	C. Rhemann	S.R. Fanello	D. Kim	A. Kowdle	W. Chang
Y. Degtyarev	P. Davidson	S. Khamis	M. Dou	V. Tankovich	C. Loop
Q. Cai	P. Chou	S. Mennicken	J. Valentin	P. Kohli	V. Pradeep
	S. Wang	Y. Lutchyn	C. Keskin	S. Izadi	

Microsoft Research (contact:shahrami@microsoft.com)



Figure 1. Holoportation is a new immersive telepresence system that combines the ability to capture high quality 3D models of people, objects and environments in real-time, with the ability to transmit these and allow remote participants wearing virtual or augmented reality displays to see, hear and interact almost as if they were co-present.

**ABSTRACT**  
 We present an end-to-end system for augmented and virtual reality telepresence, called Holoportation. Our system demonstrates high-quality, real-time 3D reconstructions of an entire space, including people, furniture and objects, using a set of new depth cameras. These 3D models can also be transmitted in real-time to remote users. This allows users wearing virtual or augmented reality displays to see, hear and interact with remote participants in 3D, almost as if they were present in the same physical space. From an audio-visual perspective, communicating and interacting with remote users edges closer to face-to-face communication. This paper describes the Holoportation technical system in full, its key interactive capabilities, the application scenarios it enables, and an initial qualitative study of using this new communication medium.

**Author Keywords** Telepresence; Non-rigid

from delivering an experience close to *physical co-presence*. For example, despite the myriad of telecommunication technologies, we still spend over \$1 trillion/year globally on business travel, with the US making up \$300 billion/year and over 482 million flights/year<sup>1</sup>. And this is without counting the cost on the environment. Indeed telepresence has been cited as key in battling carbon emissions in the future<sup>2</sup>.

However, despite the promise of telepresence, clearly we are still spending a great deal of time, money, and CO<sub>2</sub> getting on planes to meet face-to-face. Somehow much of the subtleties of face-to-face co-located communication — eye contact, body language, physical presence — are still lost in even high-end audio and video conferencing. There is still a clear gap between even the highest fidelity telecommunication tools and physically being there.

In this paper, we describe Holoportation, a system that attempts to close this gap. Holoportation is a new communication tool that leverages consumer augmented reality (AR) and

## Normal Weighted Blending

$$w_i = \underbrace{V}_{\text{Visibility test}} \cdot \max(0, \underbrace{\hat{n}}_{\text{Normal vector}} \cdot \underbrace{\hat{v}_i}_{\text{Texture camera view direction}})^\alpha$$

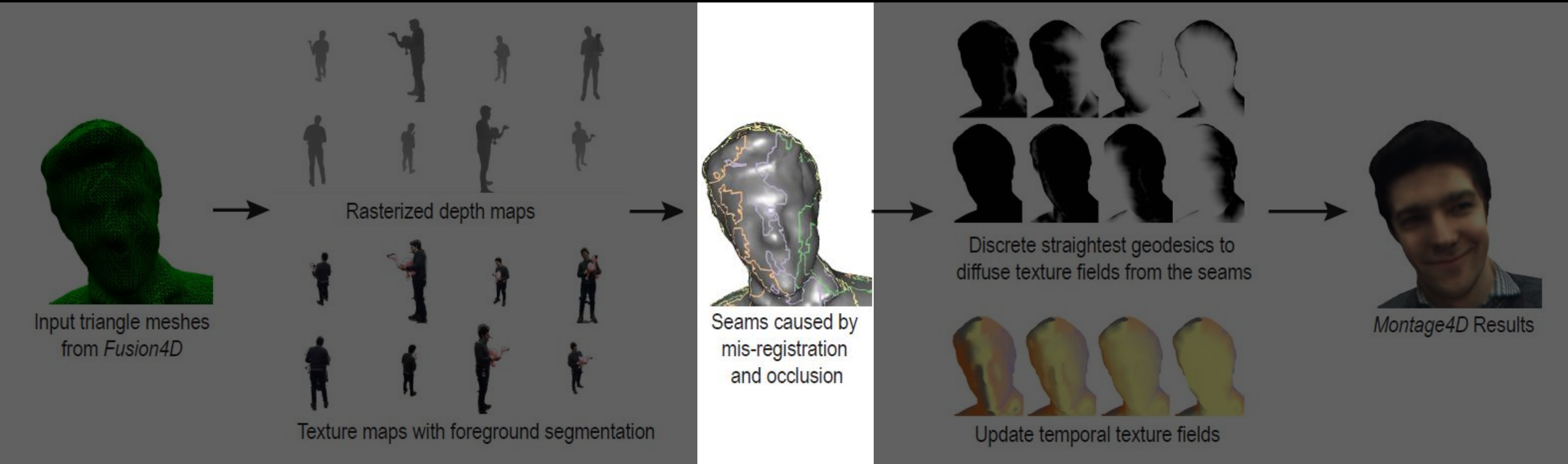
## Majority Voting for color correction

For each vertex, and for each texture, test **if the projected color agrees with more than half of the other textures**, if not, set the texture weight field to 0.

What is our approach for real-time **seamless**  
texture fusion?

# Workflow

Identify and diffuse the seams



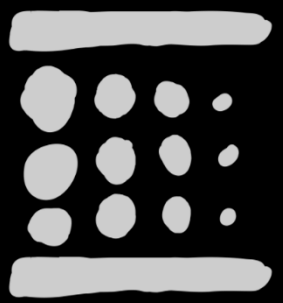
**Figure 2:** *The workflow of the Montage4D rendering pipeline.*

What are the **causes** for the *blurring* and *seams*?

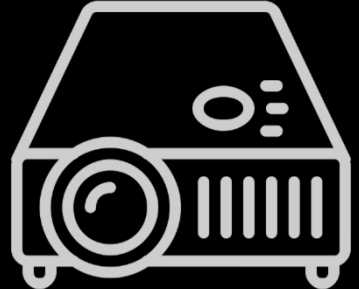


# Motivation

Causes for blurring



Blurring



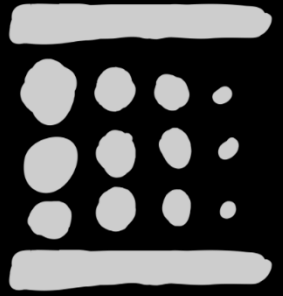
Texture projection errors  
**Inaccurate camera calibration**



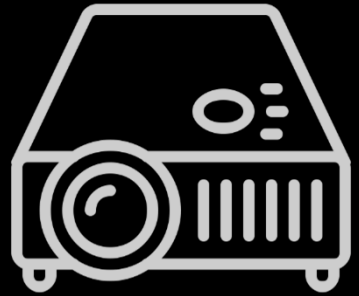
Normal-weighted blending

# Motivation

Causes for blurring



Blurring



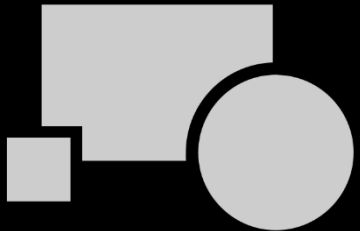
Texture projection errors



~~Normal-weighted blending~~  
**View-dependent rendering**

# Motivation

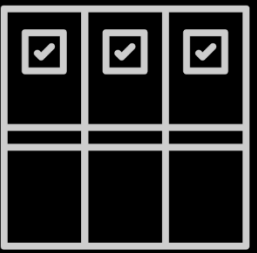
Causes for Seams



*Self-occlusion*



*Seams*



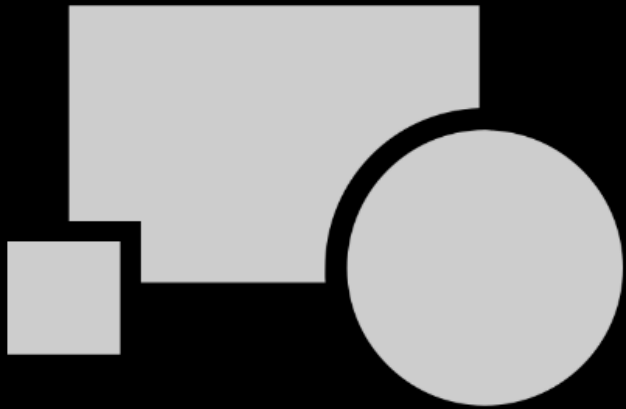
*Majority-voting*



*Field-of-View*

# Seams

Causes



## Self-occlusion

One or two vertices of the triangle are occluded in the depth map while the others are not.

# Seams

Causes



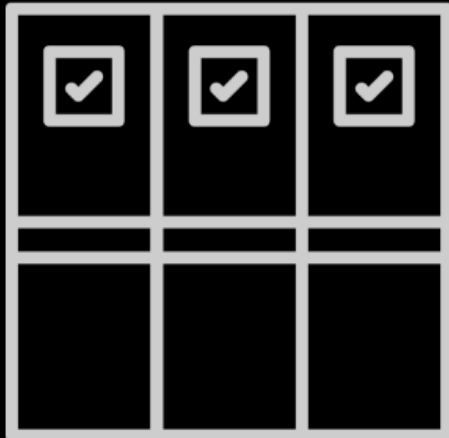
Raw projection mapping results



Seams after occlusion test

# Seams

## Causes



## Majority Voting

The triangle vertices have different results in the majority voting process, which may be caused by either mis-registration or self-occlusion.

# Seams

Causes



Raw projection mapping results



Seams after occlusion test



Seams after majority voting test

# Seams

Causes



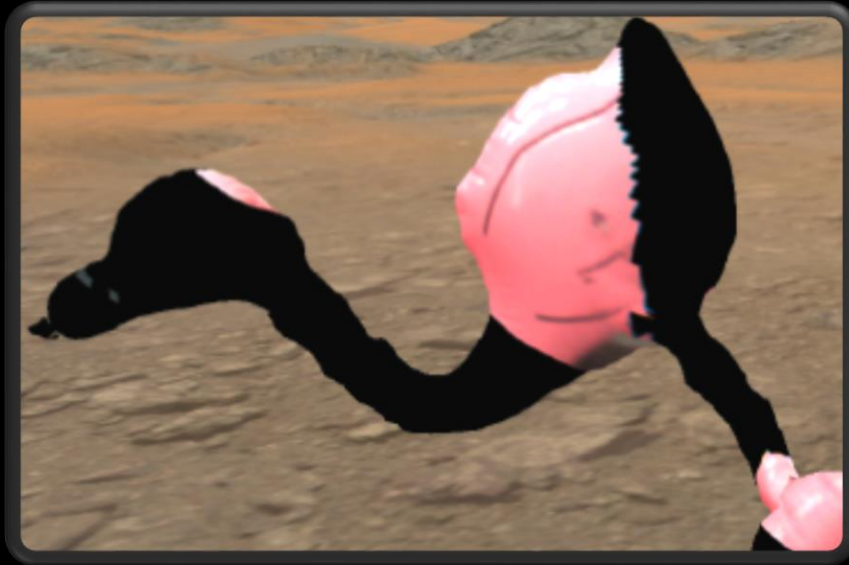
## Field of View

One or two triangle vertices lie outside the camera's field of view or in the subtracted background region while the rest are not.

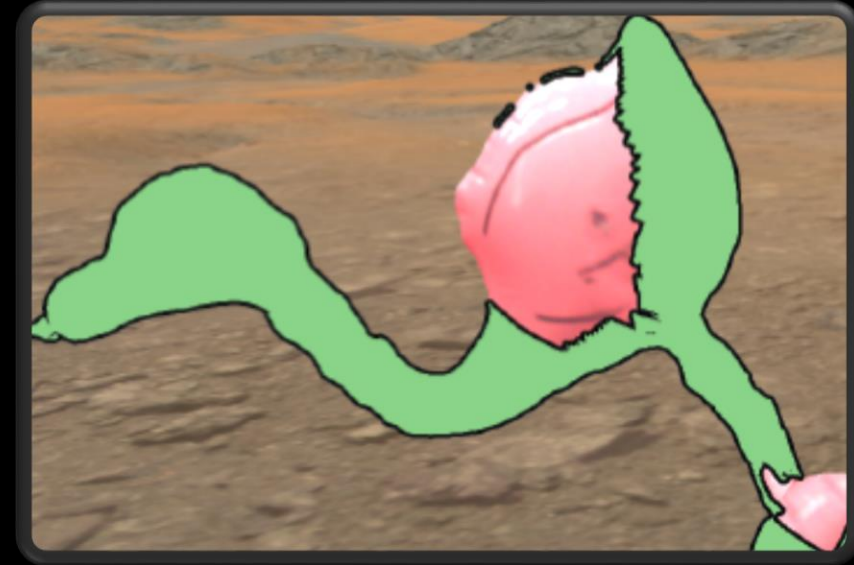


# Seams

Causes



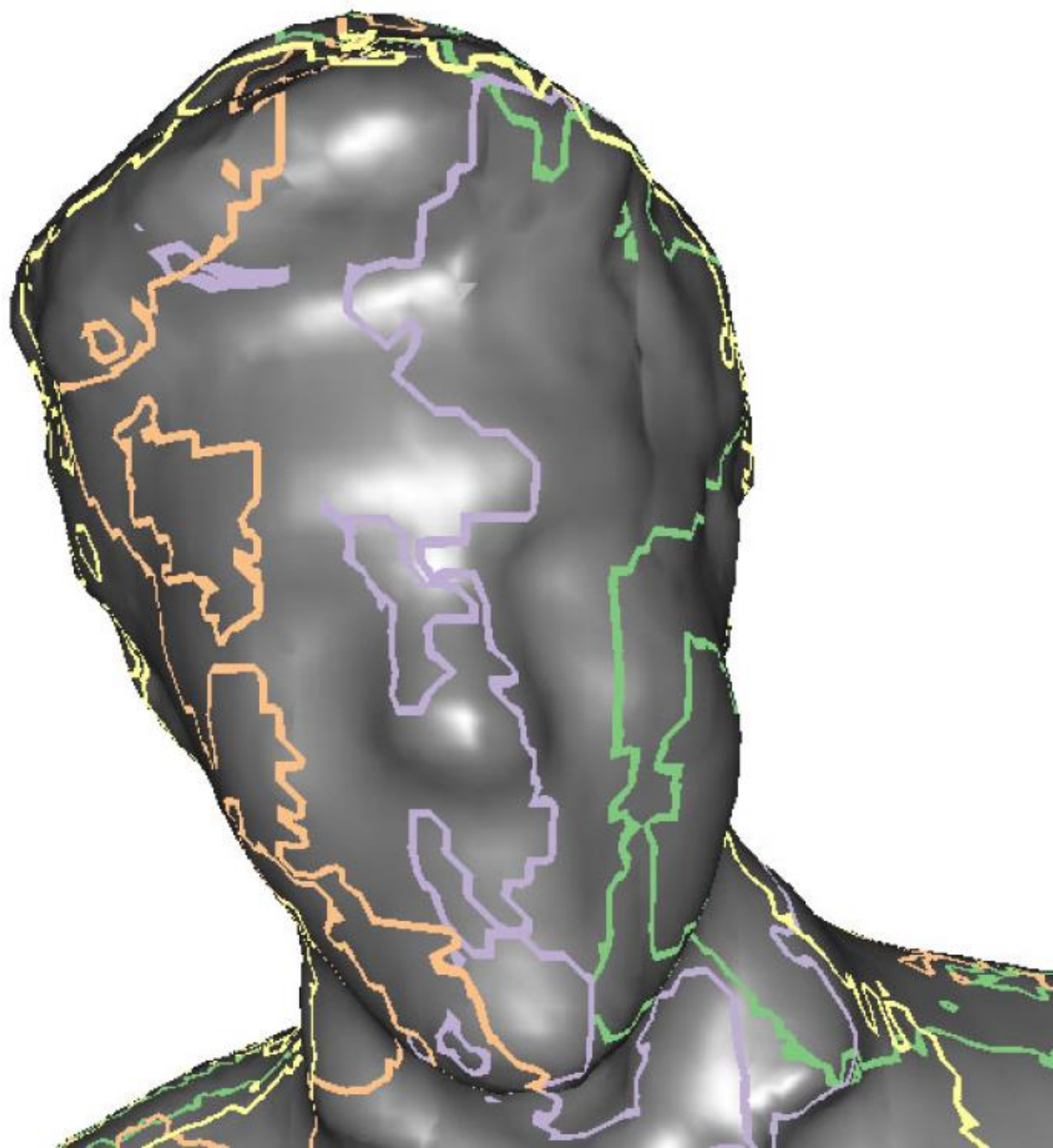
Raw projection mapping results



Seams after field-of-view test

# Seams

Causes

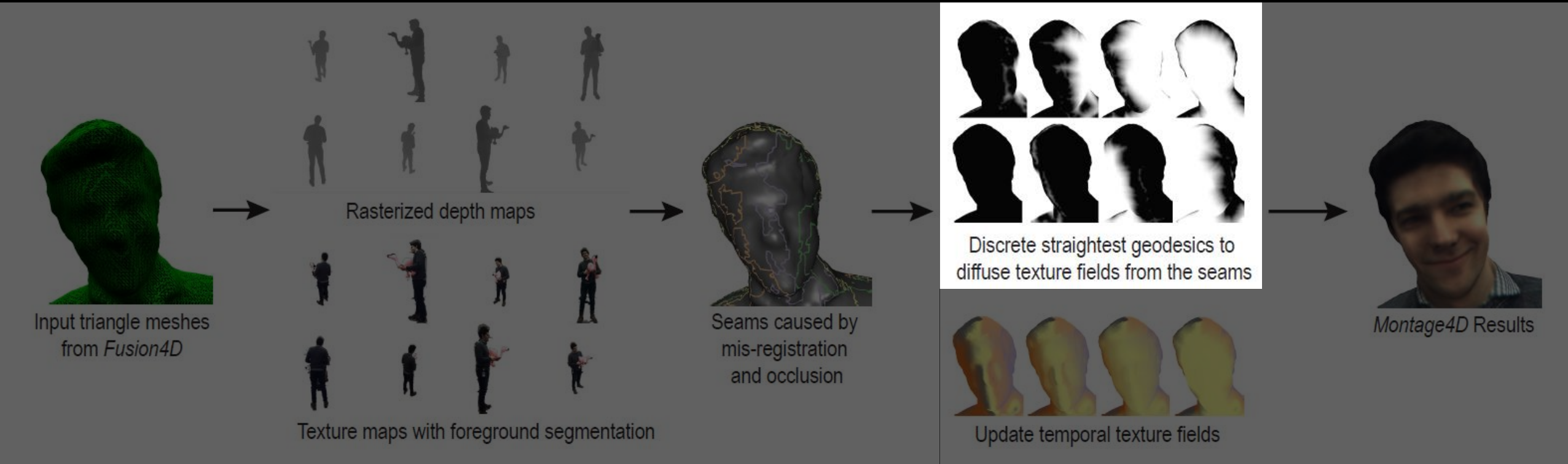


For a static frame, how can we **get rid of** the annoying seams at interactive frame rate?

How can we spatially smooth the **texture (weight) field near the seams** so that we cannot see visible seams in the results?

# Workflow

Identify and diffuse the seams



**Figure 2:** *The workflow of the Montage4D rendering pipeline.*

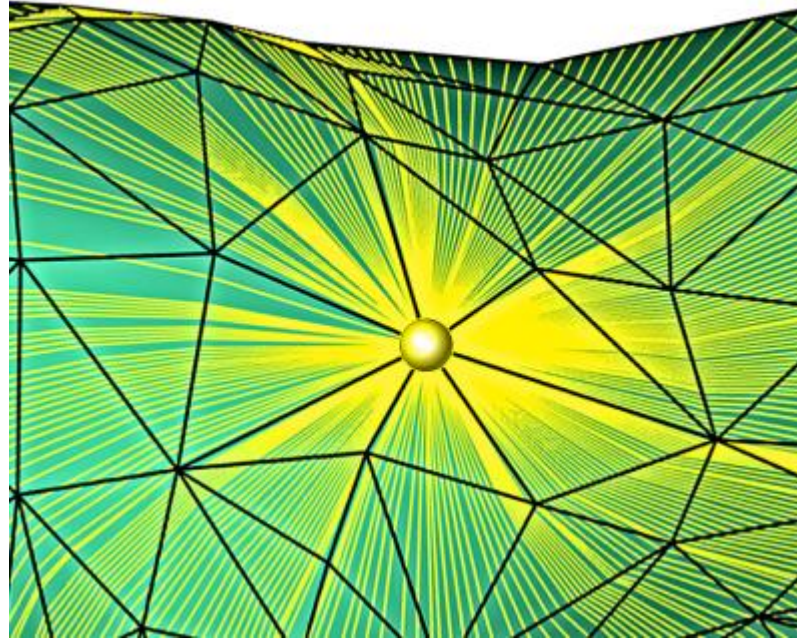
# Geodesics

For diffusing the seams

Geodesic is the **shortest** route between two points on the surface.

# Geodesics

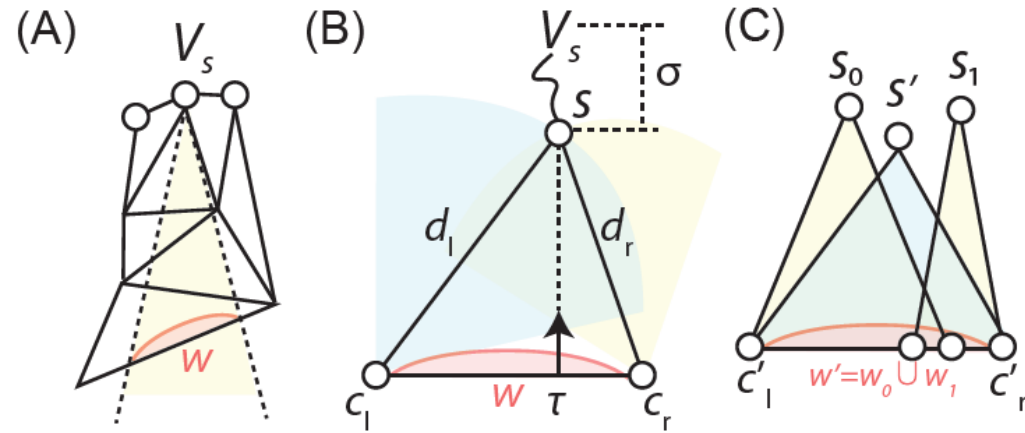
For diffusing the seams



On triangle meshes, this is challenging because of the computation of **tangent directions**.  
And shortest paths are defined on **edges** instead of the vertices.

# Geodesics

For diffusing the seams



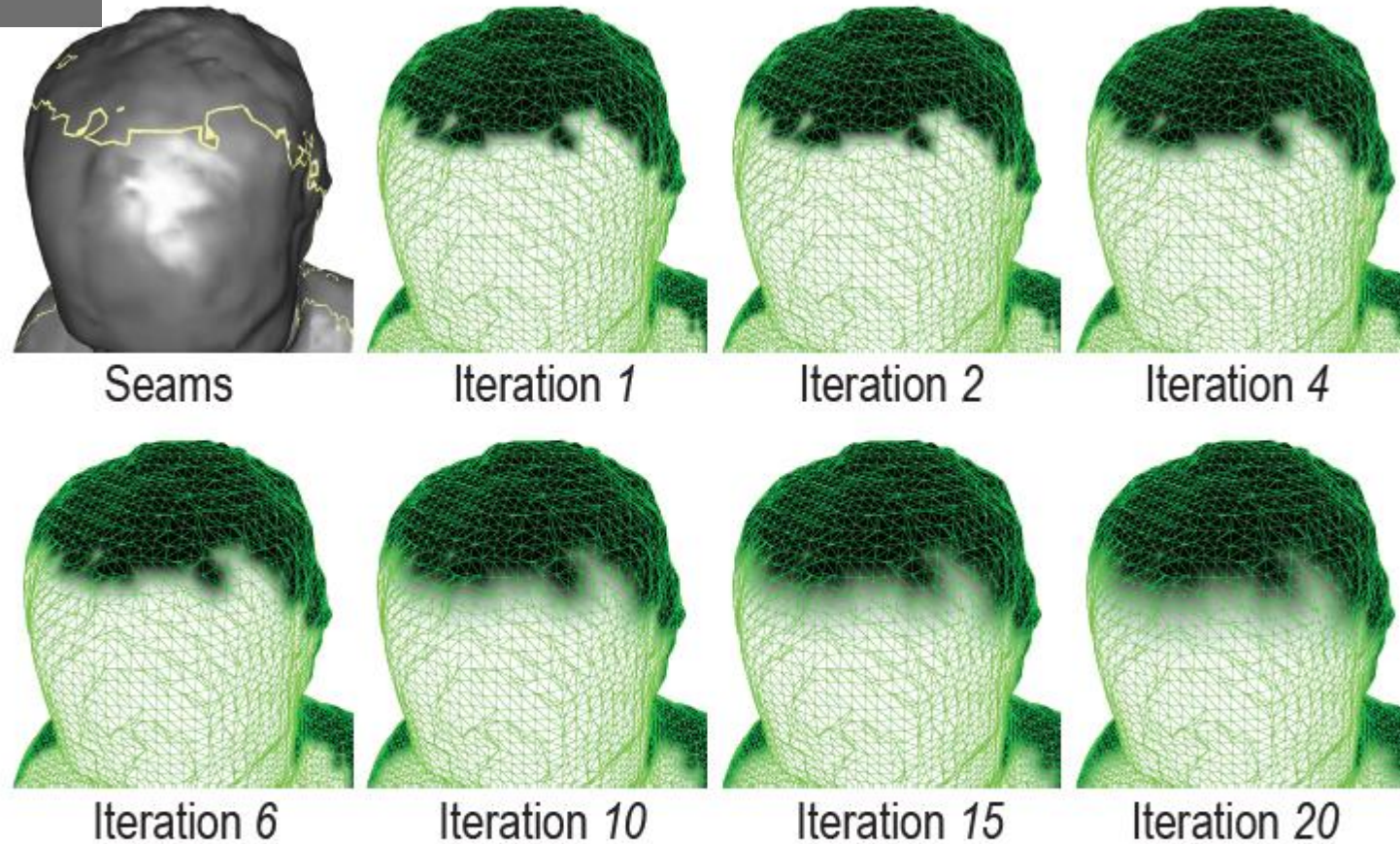
**Figure 5:** Illustration of computing the approximate geodesics. (A) shows the concept of the geodesic window from a single source vertex. (B) shows the components within a window. (C) shows the merging process of two overlapping windows for approximation.

We use the algorithm by *Surazhsky* and *Hoppe* for computing **the approximate geodesics**.  
The idea is to maintain **only 2~3 shortest paths** along each edge to reduce the computational cost.



# Approximate Geodesics

For diffusing the seams



**Figure 6:** *Examples of the initial seam triangles and the propagation process for updating the geodesics.*

# View-dependent Rendering

Spatially highlight the close views

$$\mathcal{I}_{\mathbf{v}}^i = \underbrace{\mathcal{V}_{\mathbf{v}}}_{\text{Visibility test}} \cdot \underbrace{g^i}_{\text{Geodesics}} \cdot \underbrace{\gamma_{\mathbf{v}}^i}_{\text{Global weight Visibility\% of view } i} \cdot \max(0, \underbrace{\hat{\mathbf{v}} \cdot \hat{\mathbf{v}}_i}_{\text{Texture camera view direction}})^{\alpha}$$

Texture camera view direction

Global weight  
Visibility% of view i

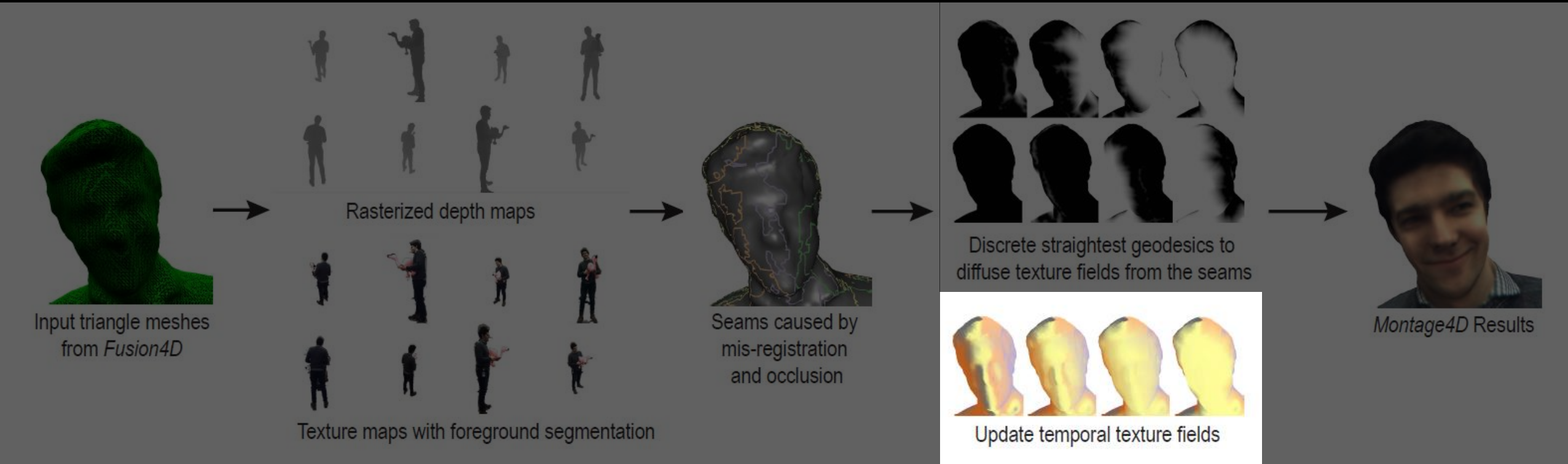
Visibility test

Geodesics

User camera's view direction

# Workflow

Identify and diffuse the seams



**Figure 2:** *The workflow of the Montage4D rendering pipeline.*

# Temporal Texture Field

Temporally smooth the texture fields

Temporal smoothing factor of 0.02

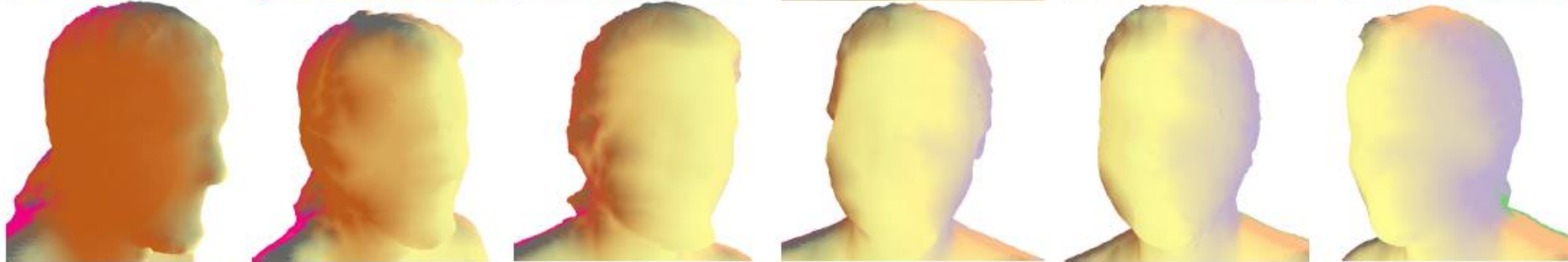
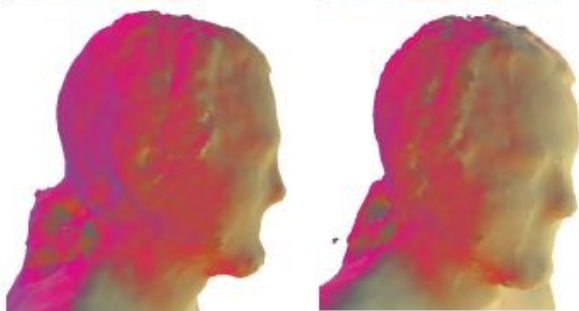
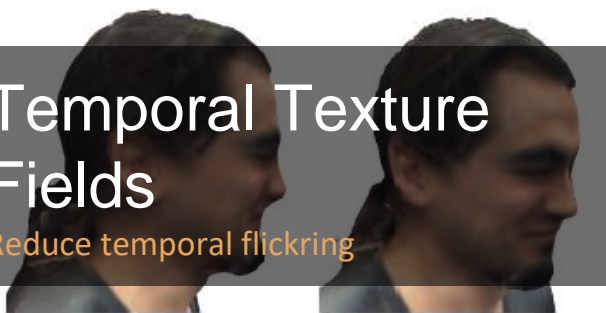
$$\mathcal{I}_{\mathbf{v}}^i(t) = \mathcal{I}_{\mathbf{v}}^i(t-1) + \lambda \nabla \mathcal{I}_{\mathbf{v}}^i(t)$$

Texture field of the previous frame

The gradient between the ideal texture field of the *current* frame, and the value of the *previous* frame

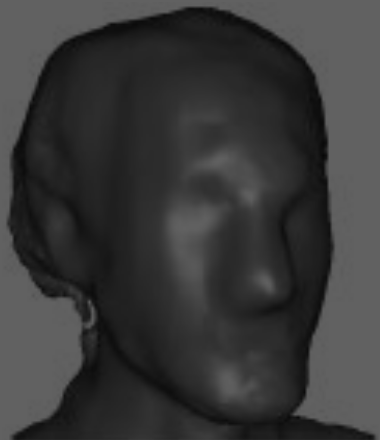
# Temporal Texture Fields

Reduce temporal flickering



Color Scheme for the Texture Fields

*Fusion4D Inputs*  
*Dou et al.*



Representative  
Projection #1



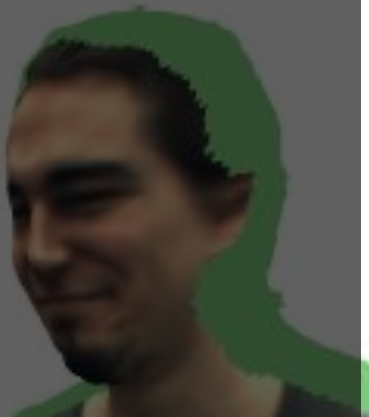
Representative  
Projection #2



*Holoportation*  
*Orts-Escolano et al.*



*Montage4D*  
Results



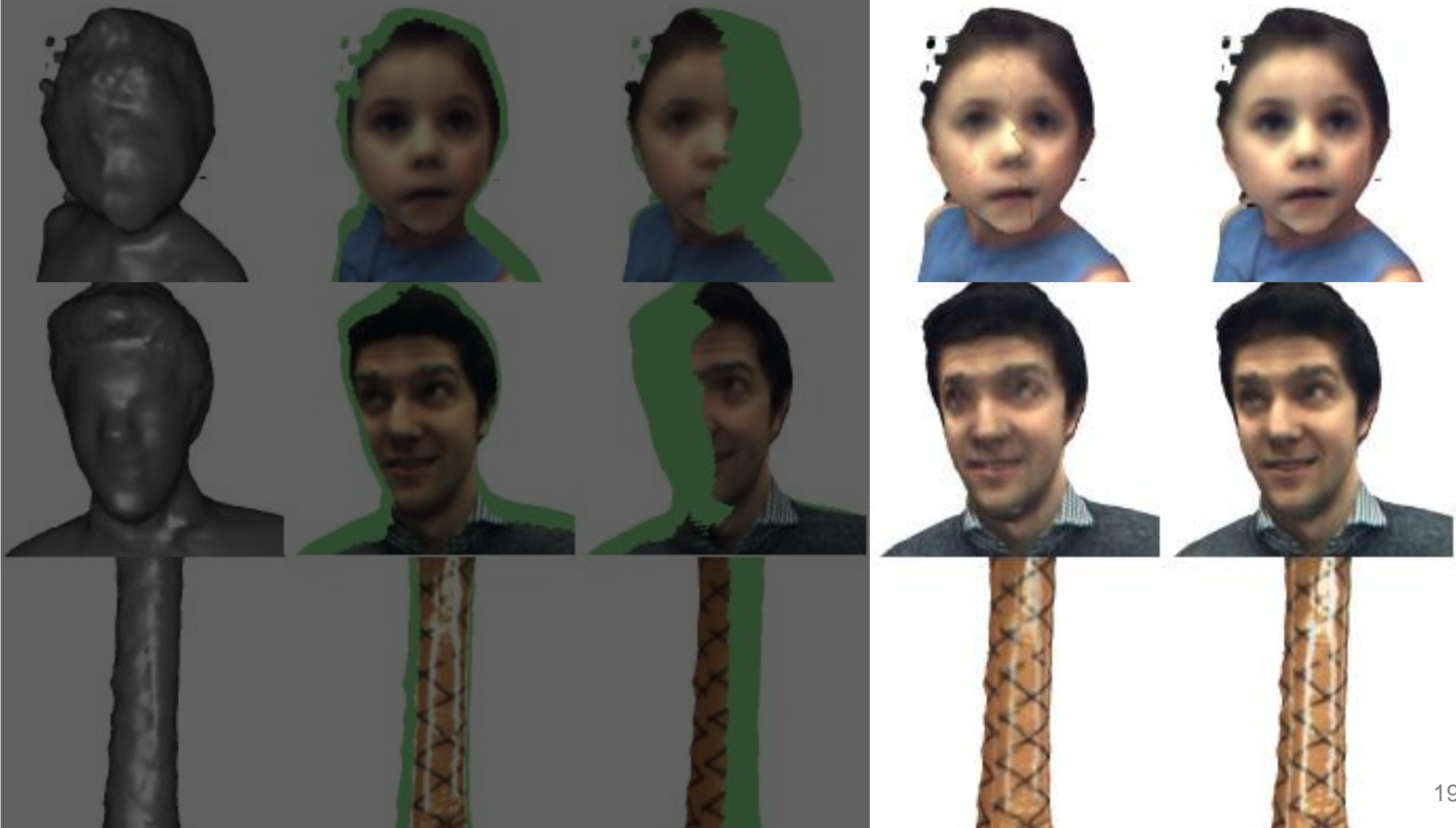
*Fusion4D Inputs*  
*Dou et al.*

Representative  
Projection #1

Representative  
Projection #2

*Holoportation*  
*Orts-Escolano et al.*

*Montage4D*  
Results



*Plenoptic Sampling*  
Chai et al.\*



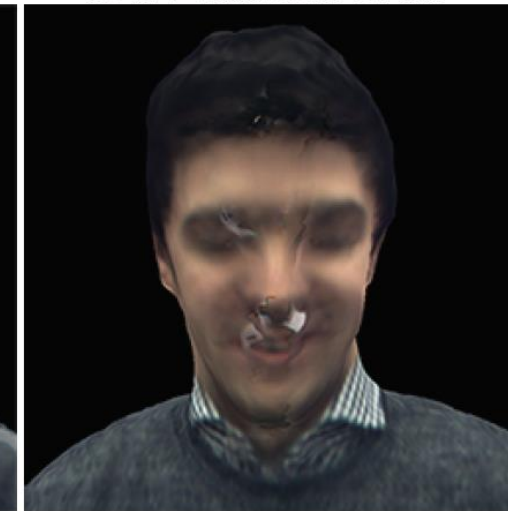
*Filtered Blending*  
Eisemann et al.



*Floating Textures*  
Eisemann et al.



*Holoportation*  
Orts-Escalano et al.



*Montage4D*  
Results





*Plenoptic Sampling*  
Chai et al.\*

*Filtered Blending*  
Eisemann et al.

*Floating Textures*  
Eisemann et al.

*Holoportation*  
Orts-Escalano et al.

*Montage4D*  
Results



*Plenoptic Sampling*  
Chai et al.\*

*Filtered Blending*  
Eisemann et al.

*Floating Textures*  
Eisemann et al.

*Holoportation*  
Orts-Escalano et al.

*Montage4D*  
Results



With additional computation for seams,  
geodesics, and temporal texture fields, is  
**our approach** still in real time?

# Experiment

RMSE = Root mean squared error

FPS = frames per second

**Table 1:** Comparison of root-mean-square error (RMSE) and frame rates between Holoporation and Montage4D methods

Dataset	Frames	#vertices / frame	#triangles / frame	Holoporation		Montage4D	
				RMSE	FPS	RMSE	FPS
Timo	837	131K	251K	5.63%	227.2	3.27%	135.0
Yury	803	132K	312K	5.44%	222.8	3.01%	130.5
Sergio	837	215K	404K	7.74%	186.8	4.21%	114.3
Girl	1192	173K	367K	7.16%	212.56	3.73%	119.4
Julien	526	157K	339K	12.63%	215.18	6.71%	120.6

# Experiment

RMSE = Root mean squared error  
FPS = frames per second

**Table 1:** Comparison of root-mean-square error (RMSE) and frame rates between Holoporation and Montage4D methods

Dataset	Frames	#vertices / frame	#triangles / frame	Holoporation		Montage4D	
				RMSE	FPS	RMSE	FPS
Timo	837	131K	251K	5.63%	227.2	3.27%	135.0
Yury	803	132K	312K	5.44%	222.8	3.01%	130.5
Sergio	837	215K	404K	7.74%	186.8	4.21%	114.3
Girl	1192	173K	367K	7.16%	212.56	3.73%	119.4
Julien	526	157K	339K	12.63%	215.18	6.71%	120.6

The root-mean-squared-error of the rendering results have been great **reduced**.

# Experiment

RMSE = Root mean squared error

FPS = frames per second

**Table 1:** Comparison of root-mean-square error (RMSE) and frame rates between Holoporation and Montage4D methods

Dataset	Frames	#vertices / frame	#triangles / frame	Holoporation		Montage4D	
				RMSE	FPS	RMSE	FPS
Timo	837	131K	251K	5.63%	227.2	3.27%	135.0
Yury	803	132K	312K	5.44%	222.8	3.01%	130.5
Sergio	837	215K	404K	7.74%	186.8	4.21%	114.3
Girl	1192	173K	367K	7.16%	212.56	3.73%	119.4
Julien	526	157K	339K	12.63%	215.18	6.71%	120.6

Our frame rate is slower, but still capable at **over 90 FPS**, for dynamic VR rendering.

# Experiment

Break-down of a typical frame

**Table 2:** *Timing comparison between Holoportation and Montage4D for a typical frame*

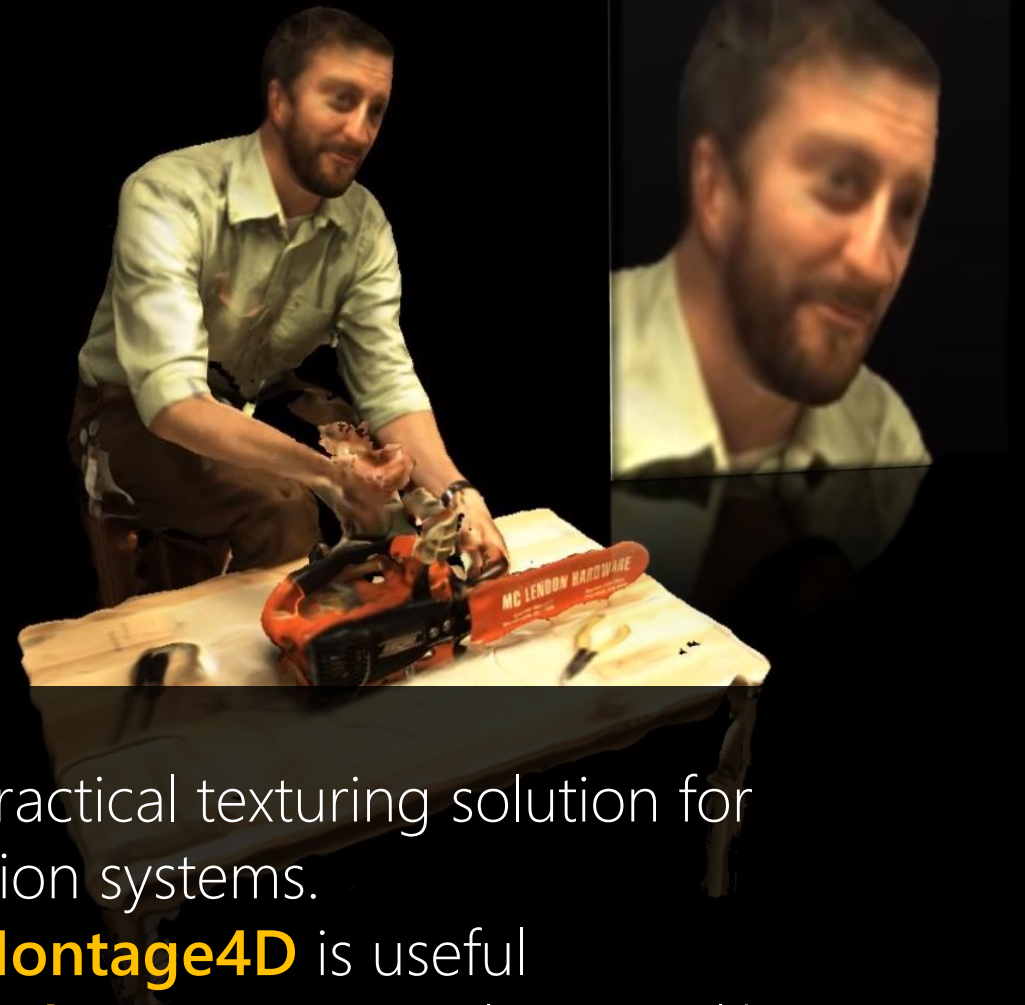
Procedure	Timing (ms)	
	Holoportation	Montage4D
Communication between CPU and GPU	4.83	9.49
Rendering and Texture Sampling	0.11	0.30
Rasterized Depth Maps calculation	0.14	0.13
Seams Identification	N/A	0.01
Approximate Geodesics estimation	N/A	0.31
Other events	0.12	0.18
Total	5.11	10.40

Most of the time is used in communication between CPU and GPU

**Before**



**After**



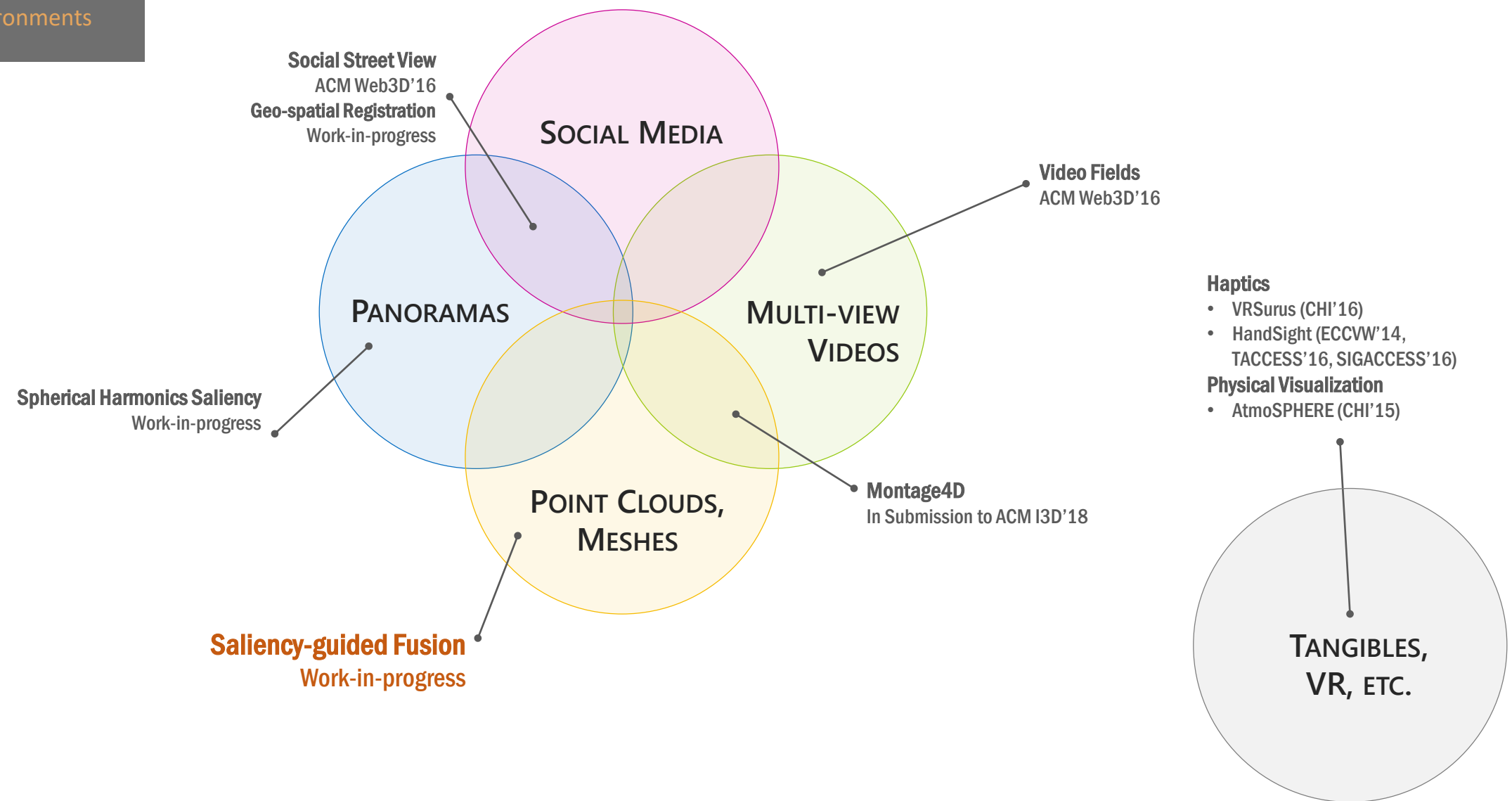
In conclusion, **Montage4D** provides a practical texturing solution for real-time 3D reconstruction systems.

In the future, we envision that **Montage4D** is useful for **fusing the massive multi-view video data** into VR applications like *remote business meeting, immersive education, and family gathering.*



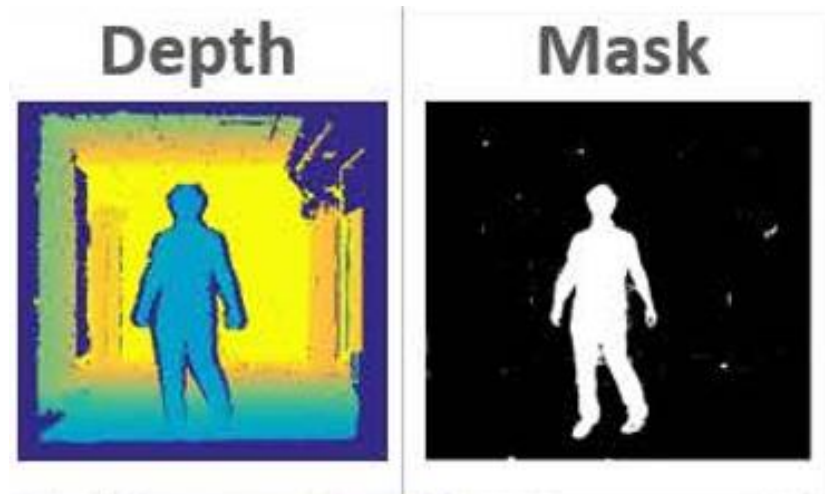
# Proposal

Fusing Multimedia Data Into Dynamic Virtual Environments



# Saliency-guided Fusion

Image Feature



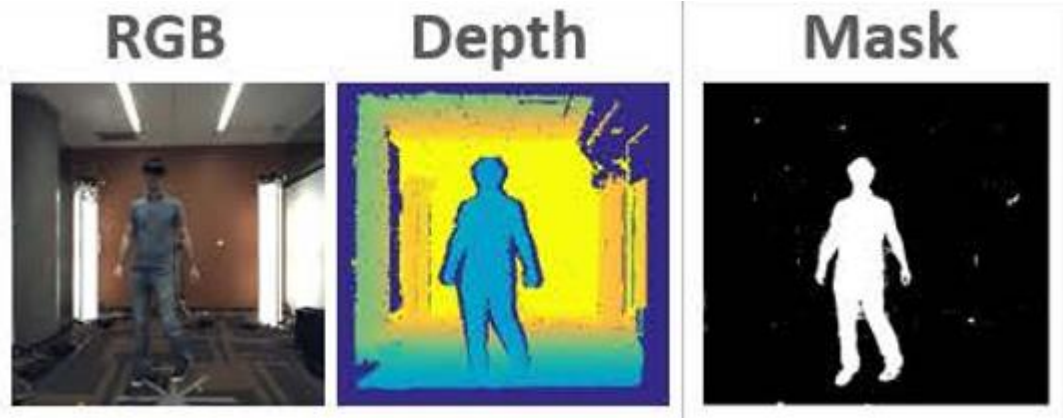
Depth Estimation and Segmentation



Volumetric Fusion

# Saliency-guided Fusion

Image Feature



Full RGBD information and segmentation



Colored Volumetric Fusion

# Colors Terms

## Optimizing Energy Functions

$$E(G) = \lambda_{depth}E_{depth}(G) + \lambda_{color}E_{color}(G) + \lambda_{corr}E_{corr}(G) + \lambda_{rot}E_{rot}(G) + \lambda_{smooth}E_{smooth}(G)$$

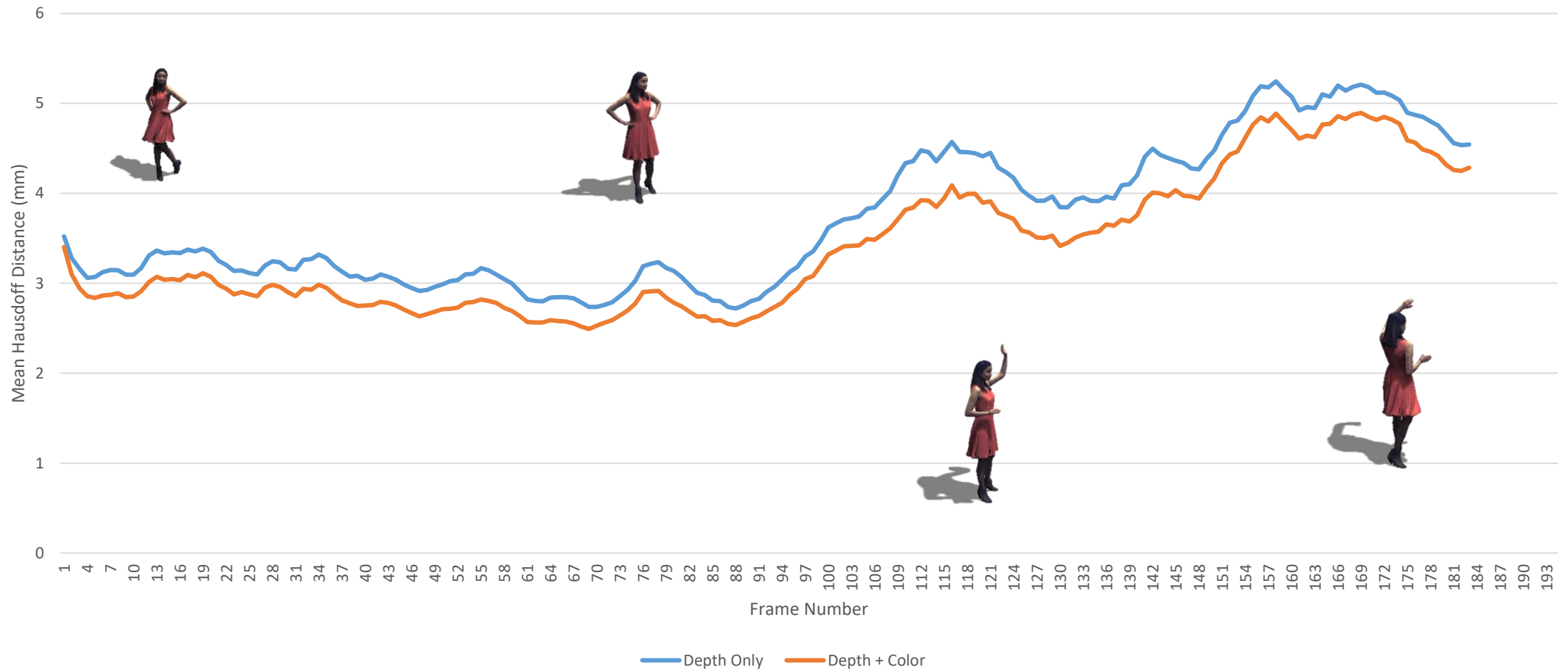
$$E_{depth}(G) = \sum_{n=1}^N \sum_{m \in V_n(G)} \min_{x \in P(\text{Depth}_n)} \|\text{Depth}(v_m; G) - x\|^2$$

$$E_{color}(G) = \sum_{n=1}^N \sum_{m \in V_n(G)} \min_{c \in P(\text{Color}_n)} \|\text{Color}(v_m; G)_{LAB} - c_{LAB}\|^2$$

# Saliency-guided Fusion

Image Feature

Quantitative Comparison Between *Fusion4D (Depth Only)* and *Intrinsic4D (Depth + Color)*

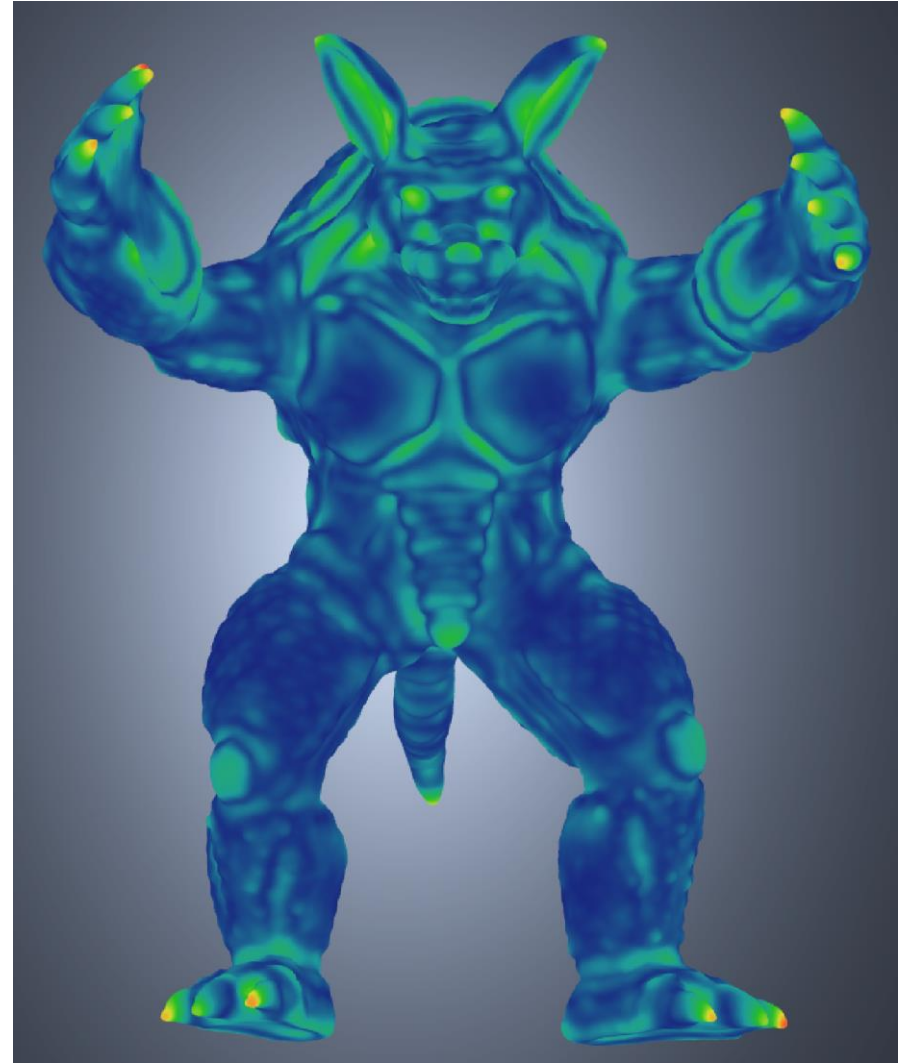


The frame rate **drops** from 30 FPS to 10 FPS with the color term...

How to **speed up** the dynamic reconstruction procedure,  
with the performance improvement?

# Mesh Saliency

Lee, Varshney, Jacobs. SIGGRAPH 2005



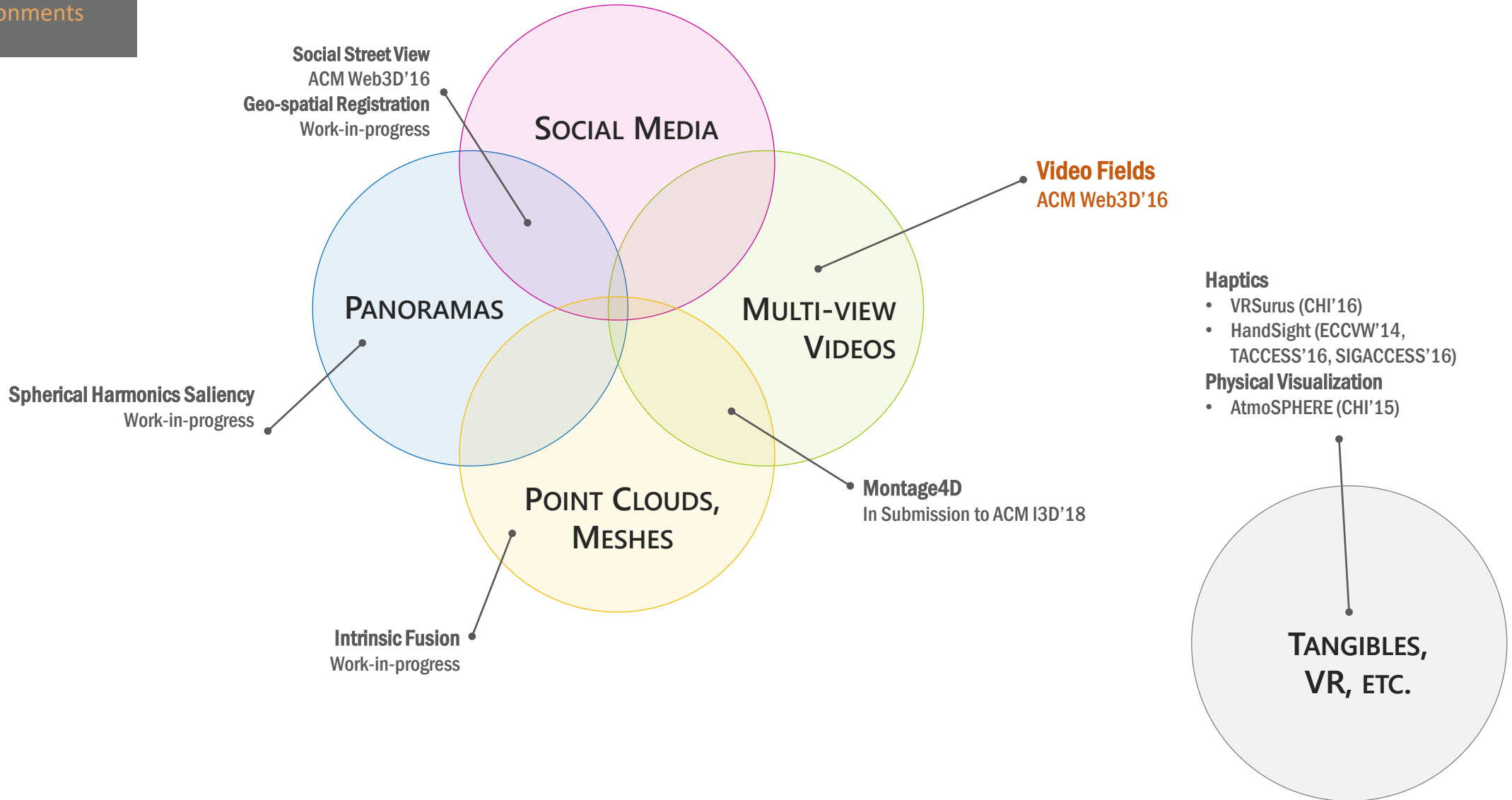


Use color-term optimization and more voxels for reconstructing **high-saliency** regions

**Saliency** is the vital key to balance the tradeoff between quality and speed.

# Proposal

Fusing Multimedia Data Into Dynamic Virtual Environments



# Video Fields: Fusing Multiple Surveillance Videos Into a Dynamic Virtual Environment

**Ruofei Du, Sujal Bista, and Amitabh Varshney**

[www.Video-Fields.com](http://www.Video-Fields.com)  
[www.Augmentarium.com](http://www.Augmentarium.com)

Augmentarium | Department of Computer Science | UMIACS  
University of Maryland, College Park

In Proceedings of the 21st Annual ACM SIGGRAPH Web3D Conference, 2016



**THE AUGMENTARIUM**  
VIRTUAL AND AUGMENTED REALITY LABORATORY  
AT THE UNIVERSITY OF MARYLAND



**GVIL**

**UMIACS**



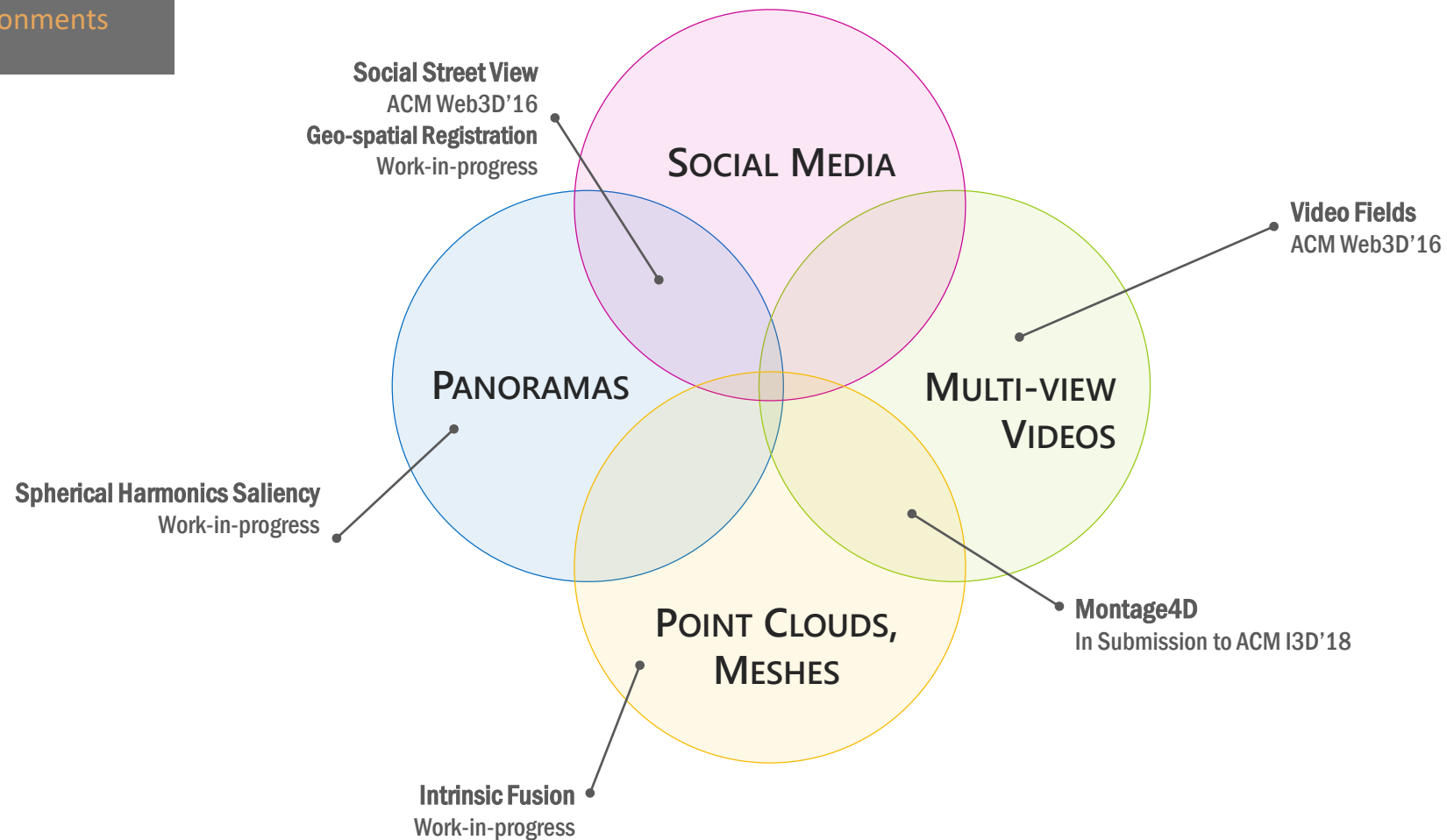
**COMPUTER SCIENCE**  
UNIVERSITY OF MARYLAND



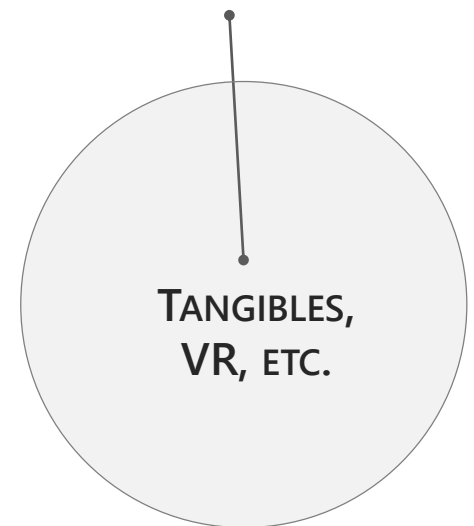
**UNIVERSITY OF  
MARYLAND**


# Proposal

Fusing Multimedia Data Into Dynamic Virtual Environments



- Haptics**
- VRSurus (CHI'16)
  - HandSight (ECCV'14, TACCESS'16, SIGACCESS'16)
- Physical Visualization**
- AtmoSPHERE (CHI'15)





# VRSurus: Enhancing Interactivity and Tangibility of Puppets in Virtual Reality (CHI 2016)

Ruofei Du and Liang He  
University of Maryland, College Park



THE AUGMENTARIUM  
VIRTUAL AND AUGMENTED REALITY LABORATORY  
AT THE UNIVERSITY OF MARYLAND



GVIL

UMIACS



COMPUTER SCIENCE  
UNIVERSITY OF MARYLAND



UNIVERSITY OF  
MARYLAND

# Demo

VRSurus



# Timeline

Paper deadline goals

Date	Project
Oct. – Nov. 2017	Spherical Harmonics Videos ACM I3D 2018
Nov. – Feb. 2018	Geo-spatial Registration with Social Street View ACM ToG / IEEE TVCG
Mar. – Aug. 2018	Saliency-guided Real-time 3D Reconstruction ACM ToG / IEEE TVCG



# Paper List

## Preliminary Work

1. **Du, R.**, Varshney, A. Social Street View: Blending Immersive Street Views with Geo-tagged Social Media. In proceedings of the 21st International Conference on Web3D, 2016. pp. 77-85. ACM. **(Best Paper Award)**
2. **Du, R.**, Bista, S., Varshney, A. Video Fields: Fusing Multiple Surveillance Videos into a Dynamic Virtual Environment. In proceedings of the 21st International Conference on Web3D, 2016. pp. 165-172. ACM.
3. **Du, R.**, Ming, C., Chang, W., Hoppe, H., Varshney, A. Montage4D: Interactive Seamless Fusing Multiview Video Textures. In submission to ACM SIGGRAPH Symposium on Interactive Graphics (ACM I3D) 2018.
4. **Du, R.**, He, L. VRSurus: Enhancing Interactivity and Tangibility of Puppets in Virtual Reality. In Proceeding of the of Human Factors in Computing Systems. (CHI '16 EA) pp. 2454-2461. ACM
5. **Du, R.**, Wills, K., Potasznik, M, Froehlich, J.E. AtmoSPHERE: Representing Space and Movement Using Sand Traces in an Interactive Zen Garden. In Proceeding of the of Human Factors in Computing Systems (CHI '16 EA). pp. 1627-1632. ACM.
6. Stearns, L., **Du, R.**, Oh, U., Catherine, Z., Findlater, L., David, R., Froehlich, J.E. Evaluating Haptic and Auditory Directional Guidance to Assist Blind Persons in Reading Printed Text Using Finger-Mounted Cameras. In ACM Transactions on Accessible Computing, 8(5), pp. 1-39. 2016.
7. Stearns, L., **Du, R.**, Oh, U., Wang, Y., Findlater, L., Chellappa, R., Froehlich, J.E. The Design and Preliminary Evaluation of a Finger-Mounted Camera and Feedback System to Enable Reading of Printed Text for the Blind. In Proceeding of the European Conference on Computer Vision (ECCV '14 Workshop). pp. 615–631. 2014.
8. Findlater, L., Stearns, L., **Du, R.**, Oh, U., Wang, Y., Chellappa, R., Froehlich, J.E. Supporting Everyday Activities for Persons With Visual Impairments Through Computer Vision-Augmented Touch In Proceedings of the 17th International ACM SIGACCESS Conference on Computers & Accessibility, pp. 383-384, 2015.

# Acknowledgement

Collaborators, advisors, labmates,  
and committee members

Committee    Augmentarium  
MSR    Collaborators



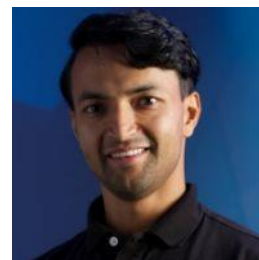
Amitabh Varshney  
varshney@cs.umd.edu



Matthias Zwicker  
zwicker@cs.umd.edu



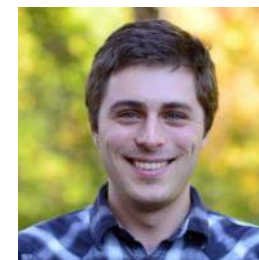
Furong Huang  
furongh@cs.umd.edu



Sujal Bista  
Sujal@cs.umd.edu



Hsueh-Chien Cheng  
hccheng@cs.umd.edu



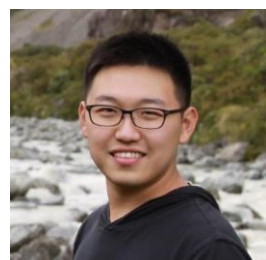
Eric Krokos  
Sujal@cs.umd.edu



Xuetong Sun  
xtsun@cs.umd.edu



Xiaoxu Meng  
xmeng525@umiacs.umd.edu



Sida Li  
sidali@umiacs.umd.edu



Eric Lee  
ericlee@umiacs.umd.edu



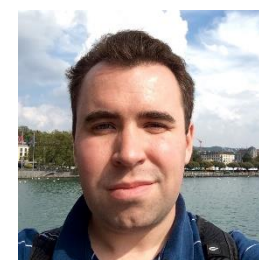
Jon Froehlich  
jonf@cs.washington.edu



Leah Findlater  
leahkf@uw.edu



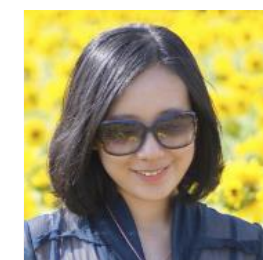
Rama Chellappa  
rama@umiacs.umd.edu



Lee Stearns  
lee@leestearns.com



Liang He  
edigahe@gmail.com



Sai Yuan  
syuan@umd.edu



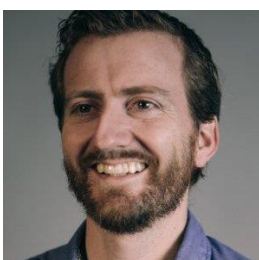
Wayne Chang  
wechang@microsoft.com



Marek Kowalski  
marek.kol4@gmail.com



Zoey Chen  
qiuyuchen14@gmail.com



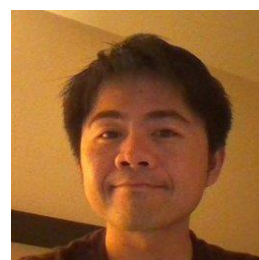
Spencer Fowers  
sfowers@microsoft.com



Jeff Kramer  
jekramer@Microsoft.com



Ben Cutler  
bcutler@microsoft.com



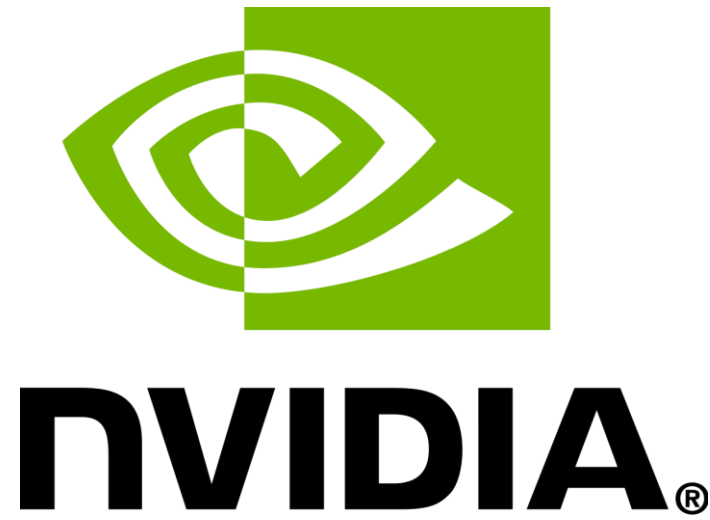
Ming Chuang  
mingchuang82@gmail.com



Hugues Hoppe  
hhoppe@gmail.com

# Acknowledgement

NSF | Nvidia | MPower | UMIACS



**UMIACS**  
University of Maryland  
Institute for Advanced  
Computer Studies

Microsoft®  
**Research**





# Thank you

Ruofei Du  
ruofei@cs.umd.edu

Committee: Dr. Varshney, Dr. Zwicker, and Dr. Huang



**THE AUGMENTARIUM**  
VIRTUAL AND AUGMENTED REALITY LABORATORY  
AT THE UNIVERSITY OF MARYLAND



GVIL

UMIACS



COMPUTER SCIENCE  
UNIVERSITY OF MARYLAND



UNIVERSITY OF  
MARYLAND

# Fusing Multimedia Data Into Dynamic Virtual Environments

Ruofei Du  
ruofei@cs.umd.edu

Committee: Dr. Varshney, Dr. Zwicker, and Dr. Huang