

Dynamic Adaptive 3D Streaming over HTTP

For the University of Toulouse PhD granted by the INP Toulouse Presented and defended on Friday 29th November, 2019 by Thomas Forgione

> Gilles GESQUIÈRE, president Sidonie CHRISTOPHE, reviewer Gwendal SIMON, reviewer Maarten WIJNANTS, examiner Wei Tsang OOI, examiner Vincent CHARVILLAT, thesis supervisor Axel CARLIER, thesis co-supervisor Géraldine MORIN, thesis co-supervisor

Doctoral school and field: EDMITT: École Doctorale de Mathématiques, Informatiques et Télécommunications de Toulouse Field: Computer science and telecommunication Research unit: IRIT (5505) Thesis supervisors: Vincent CHARVILLAT, Axel CARLIER and Géraldine MORIN Reviewers: Sidonie CHRISTOPHE and Gwendal SIMON

Titre : Transmission Adaptative de Modèles 3D Massifs

Résumé :

Avec les progrès de l'édition de modèles 3D et des techniques de reconstruction 3D, de plus en plus de modèles 3D sont disponibles et leur qualité augmente. De plus, le support de la visualisation 3D sur le web s'est standardisé ces dernières années. Un défi majeur est donc de transmettre des modèles massifs à distance et de permettre aux utilisateurs de visualiser et de naviguer dans ces environnements virtuels. Cette thèse porte sur la transmission et l'interaction de contenus 3D et propose trois contributions majeures.

Tout d'abord, **nous développons une interface de navigation dans une scène 3D avec des signets** – de petits objets virtuels ajoutés à la scène sur lesquels l'utilisateur peut cliquer pour atteindre facilement un emplacement recommandé. Nous décrivons une étude d'utilisateurs où les participants naviguent dans des scènes 3D avec ou sans signets. Nous montrons que les utilisateurs naviguent (et accomplissent une tâche donnée) plus rapidement en utilisant des signets. Cependant, cette navigation plus rapide a un inconvénient sur les performances de la transmission : un utilisateur qui se déplace plus rapidement dans une scène a besoin de capacités de transmission plus élevées afin de bénéficier de la même qualité de service. Cet inconvénient peut être atténué par le fait que les positions des signets sont connues à l'avance : en ordonnant les faces du modèle 3D en fonction de leur visibilité depuis un signet, on optimise la transmission et donc, on diminue la latence lorsque les utilisateurs cliquent sur les signets.

Deuxièmement, **nous proposons une adaptation du standard de transmission DASH (Dynamic Adaptive Streaming over HTTP), très utilisé en vidéo, à la transmission de maillages texturés 3D. Pour ce faire, nous divisons la scène en un arbre k-d où chaque cellule correspond à un adaptation set DASH. Chaque cellule est en outre divisée en segments DASH d'un nombre fixe de faces, regroupant des faces de surfaces comparables. Chaque texture est indexée dans son propre adaptation set à différentes résolutions. Toutes les métadonnées (les cellules de l'arbre k-d, les résolutions des textures, etc.) sont référencées dans un fichier XML utilisé par DASH pour indexer le contenu: le MPD (Media Presentation Description). Ainsi, notre framework hérite de la scalabilité offerte par DASH. Nous proposons ensuite des algorithmes capables d'évaluer l'utilité de chaque segment de données en fonction du point de vue du client, et des politiques de transmission qui décident des segments à télécharger.**

Enfin, **nous étudions la mise en place de la transmission et de la navigation 3D sur les appareils mobiles**. Nous intégrons des signets dans notre version 3D de DASH et proposons une version améliorée de notre client DASH qui bénéficie des signets. Une étude sur les utilisateurs montre qu'avec notre politique de chargement adaptée aux signets, les signets sont plus susceptibles d'être cliqués, ce qui améliore à la fois la qualité de service et la qualité d'expérience des utilisateurs. Title: Dynamic Adaptive 3D Streaming over HTTP

Abstract:

With the advances in 3D models editing and 3D reconstruction techniques, more and more 3D models are available and their quality is increasing. Furthermore, the support of 3D visualization on the web has become standard during the last years. A major challenge is thus to deliver these remote heavy models and to allow users to visualise and navigate in these virtual environments. This thesis focuses on 3D content streaming and interaction, and proposes three major contributions.

First, we develop a 3D scene navigation interface with bookmarks – small virtual objects added to the scene that the user can click on to ease reaching a recommended location. We describe a user study where participants navigate in 3D scenes with and without bookmarks. We show that users navigate (and accomplish a given task) faster when using bookmarks. However, this faster navigation has a drawback on the streaming performance: a user who moves faster in a scene requires higher streaming capabilities in order to enjoy the same quality of service. This drawback can be mitigated using the fact that bookmarks positions are known in advance: by ordering the faces of the 3D model according to their visibility at a bookmark, we optimize the streaming and thus, decrease the latency when users click on bookmarks.

Secondly, we propose an adaptation of Dynamic Adaptive Streaming over HTTP (DASH), the video streaming standard, to 3D textured meshes streaming. To do so, we partition the scene into a k-d tree where each cell corresponds to a DASH adaptation set. Each cell is further divided into DASH segments of a fixed number of faces, grouping together faces of similar areas. Each texture is indexed in its own adaptation set, and multiple DASH representations are available for different resolutions of the textures. All the metadata (the cells of the k-d tree, the resolutions of the textures, etc.) is encoded in the Media Presentation Description (MPD): an XML file that DASH uses to index content. Thus, our framework inherits DASH scalability. We then propose clients capable of evaluating the usefulness of each chunk of data depending on their viewpoint, and streaming policies that decide which chunks to download.

Finally, we investigate the setting of 3D streaming and navigation on mobile devices. We integrate bookmarks in our 3D version of DASH and propose an improved version of our DASH client that benefits from bookmarks. A user study shows that with our dedicated bookmark streaming policy, bookmarks are more likely to be clicked on, enhancing both users quality of service and quality of experience.

Acknowledgments

First of all, I would like to thank my advisors, Vincent CHARVILLAT, Axel CARLIER, and Géraldine MORIN for luring me into doing a PhD (which was a lot of work), for the support, and for the fun (and beers). I took *a little time* to take the decision of starting a PhD, so I also want to thank them for their patience. I also want to thank Wei Tsang OOI, for the ideas, the discussions and for the polish during the deadlines.

Then, I want to thank Sidonie CHRISTOPHE and Gwendal SIMON for reviewing this manuscript: I appreciated the feedback and constructive comments. I also want to thank all the members of the jury, for their attention and the interesting discussions during the defense.

I want to thank all the kids of the lab and elsewhere that contributed to the mood (and beers): Bastien, Vincent, Julien, Sonia, Matthieu, Jean, Damien, Richard, Thibault, Clément, Arthur, Thierry, the other Matthieu, the other Julien, Paul, Maxime, Quentin, Adrian, Salomé. I also want to thank Praveen: working with him was a pleasure.

I would also like to thank the big ones (whom I forgot to thank during the defense, *oopsies*), Sylvie, Jean-Denis, Simone, Yvain, Pierre. They, as well as my advisors, not only helped during my PhD, but they also were my teachers back in engineering school and are a great part of the reason why I enjoyed being in school and being a PhD student.

Then, I also want to thank Sylvie and Muriel, for the administrative parts of the PhD, which can often be painful.

I would also like to thank the colleagues from when I was in engineering school, since they contributed to the skills that I used during this PhD: Alexandre, Killian, David, Maxence, Martin, Maxime, Korantin, Marion, Amandine, Émilie.

Finally, I want to thank my brother, my sister and my parents, for the support and guidance. They have been decisive to my education and I would not be writing this today if it was not for them.

Contents

Acknowledgments	7
Introduction	11
1 Open problems	13
2 Thesis outline	
1 Foreword	15
1.1 What is a 3D model?	
1.1.1 Rendering a 3D model	17
1.2 Similarities and differences between video and 3D	
1.2.1 Chunks of data	
1.2.2 Data persistence	
1.2.3 Multiple representations	19
1.2.4 Media types	19
1.2.5 Interaction	19
1.2.6 Relationship between interface, interaction and streaming	21
1.3 Implementation details	21
1.3.1 JavaScript	22
1.3.2 Rust	23
2 Related work	27
2.1 Video	28
2.1.1 DASH: the standard for video streaming	28
2.1.2 DASH-SRD	30
2.2 3D streaming	
2.2.1 Compression and structuring	
2.2.2 Viewpoint dependency	33
2.2.3 Texture streaming	

2.2.4 Geometry and textures	
2.2.5 Streaming in game engines	
2.2.6 NVE streaming frameworks	
2.3 3D bookmarks and navigation aids	
3 Bookmarks, navigation and streaming	
3.1 Introduction	
3.2 Impact of 3D bookmarks on navigation	
3.2.1 Our NVE	
3.2.2 3D bookmarks	
3.2.3 User study	
3.2.4 Experimental results	
3.3 Impact of 3D bookmarks on streaming	
3.3.1 3D model streaming	
3.3.2 3D bookmarks	50
3.3.3 Comparing streaming policies	
4 DASH-3D	55
4.1 Introduction	
4.2 Content preparation	
4.2.1 The MPD File	
4.2.2 Adaptation sets	
4.2.3 Representations	59
4.2.4 Segments	
4.3 Client	
4.3.1 Segment utility	
4.3.2 Offline parameters	
4.3.3 Online parameters	
4.3.4 Utility for geometry segments	
4.3.5 Utility for texture segments	
4.3.6 DASH adaptation logic	
4.3.7 JavaScript client	
4.3.8 Rust client	65
4.4 Evaluation	
4.4.1 Experimental setup	
4.4.2 Experimental results	69
4.5 Conclusion	
5 Bookmarks for DASH-3D on mobile devices	
5.1 Introduction	
5.2 Desktop and mobile interactions	

5.2.1 Desktop interaction	75
5.2.2 Mobile interaction	76
5.3 Adding bookmarks into DASH NVE framework	77
5.3.1 Bookmark interaction and visual aspect	77
5.3.2 Segments utility at bookmarked viewpoint	78
5.3.3 MPD modification	79
5.3.4 Loader modifications	80
5.4 Conclusion	80
Bibliography	81

Introduction

During the last years, 3D acquisition and modeling techniques have made tremendous progress. Recent software uses 2D images from cameras to reconstruct 3D data, e.g. Meshroom is a free and open source software which got almost 200.000 downloads on fosshub, which use *structure-from-motion* and *multi-view-stereo* to infer a 3D model. More and more devices are specifically built to harvest 3D data: for example, LIDAR (Light Detection And Ranging) can compute 3D distances by measuring time of flight of light. The recent research interest for autonomous vehicles allowed more companies to develop cheaper LIDARs, which increase the potential for new 3D content creation. Thanks to these techniques, more and more 3D data become available. These models have potential for multiple purposes, for example, they can be printed, which can reduce the production cost of some pieces of hardware or enable the creation of new objects, but most uses are based on visualization. For example, they can be used for augmented reality, to provide user with feedback that can be useful to help worker with complex tasks, but also for fashion (for example, Fittingbox is a company that develops software to virtually try glasses, as in Figure 1).



Figure 1: My face with augmented glasses

3D acquisition and visualization is also useful to preserve cultural heritage, and software such as Google Heritage or 3DHop are such examples, or to allow users navigating in a city (as in Google Earth or Google Maps in 3D). Sketchfab (see Figure 2) is an example of a website allowing users to share their 3D models and visualize models from other users.



🔗 music instrument guitar gibson sg michaelsito1995 mikebit michaelsito gibsonsg electric

Figure 2: Sketchfab interface

In most 3D visualization systems, the 3D data are stored on a server and need to be transmitted to a terminal before the user can visualize them. The improvements in the acquisition setups we described lead to an increasing quality of the 3D models, thus an increasing size in bytes as well. Simply downloading 3D content and waiting until it is fully downloaded to let the user visualize it is no longer a satisfactory solution, so adaptive streaming is needed. In this thesis, we propose a full framework for navigation and streaming of large 3D scenes, such as districts or whole cities.

1 Open problems

The objective of our work is to design a system which allows a user to access remote 3D content. A 3D streaming client has lots of tasks to accomplish:

- Decide what part of the content to download next,
- Download the next part,
- Parse the downloaded content,
- Add the parsed result to the scene,
- Render the scene,
- Manage the interaction with the user.

This opens multiple problems which need to be considered and will be studied in this thesis.

Content preparation

Before streaming content, it needs to be prepared. The segmentation of the content into chunks is particularly important for streaming since it allows transmitting only a portion of the data to the client. The downloaded chunks can be rendered while more chunks are being downloaded. Content preparation also includes compression. One of the questions this thesis has to answer is: *what is the best way to prepare 3D content so that a streaming client can progressively download and render the 3D model?*

Streaming policies

Once our content is prepared and split in chunks, a client needs to determine which chunks should be downloaded first. A chunk that contains data in the field of view of the user is more relevant than a chunk that is not inside; a chunk that is close to the camera is more relevant than a chunk far away from the camera. This should also include other contextual parameters, such as the size of a chunk, the bandwidth and the user's behaviour. In order to propose efficient streaming policies, we need to know *how to estimate a chunk utility, and how to determine which chunks need to be downloaded depending the user's interactions*?

Evaluation

In such systems, two commonly used criteria for evaluation are quality of service, and quality of experience. The quality of service is a network-centric metric, which considers values such as throughput and measures how well the content is served to the client. The quality of experience is a user-centric metric: it relies on user perception and can only be measured by asking how users feel about a system. To be able to know which streaming policies are best, one needs to know *how to compare streaming policies and evaluate the impact of their parameters on the quality of service of the streaming system and on the quality of experience of the final user?*

Implementation

The objective of our work is to setup a client-server architecture that answers the above problems: content preparation, chunk utility, streaming policies. In this regard, we have to find out *how do we build this architecture that keeps a low computational load on the server so it scales up and on the client so that it has enough resources to perform the tasks described above?*

2 Thesis outline

First, in Chapter 1, we give some preliminary information required to understand the types of objects we are manipulating in this thesis. We then proceed to compare 3D and video content: video and 3D share many features, and analyzing video setting gives inspiration for building a 3D streaming system.

In Chapter 2, we present a review of the state of the art in multimedia interaction and streaming. This chapter starts with an analysis of the video streaming standards. Then it reviews the different 3D streaming approaches. The last section of this chapter focuses on 3D interaction.

Then, in Chapter 3, we present our first contribution: an in-depth analysis of the impact of the UI on navigation and streaming in a 3D scene. We first develop a basic interface for navigating in 3D and then, we introduce 3D objects called *bookmarks* that help users navigating in the scene. We then present a user study that we conducted on 51 people which shows that bookmarks ease user navigation: they improve performance at tasks such as finding objects. We analyze how the presence of bookmarks impacts the streaming: we propose and evaluate streaming policies based on precomputations relying on bookmarks and that measurably increase the quality of experience.

In Chapter X, we present the most important contribution of this thesis: DASH-3D. DASH-3D is an adaptation of DASH (Dynamic Adaptive Streaming over HTTP): the video streaming standard, to 3D streaming. We first describe how we adapt the concepts of DASH to 3D content, including the segmentation of content. We then define utility metrics that rate each chunk depending on the user's position. Then, we present a client and various streaming policies based on our utilities which can benefit from DASH format. We finally evaluate the different parameters of our client.

In Chapter X, we present our last contribution: the integration of the interaction ideas that we developed in Chapter X into DASH-3D. We first develop an interface that allows desktop as well as mobile devices to navigate in streamed 3D scenes, and that introduces a new style of bookmarks. We then explain why simply applying the ideas developed in Chapter X is not sufficient and we propose more efficient precomputations that enhance the streaming. Finally, we present a user study that provides us with traces on which we evaluate the impact of our extension of DASH-3D on the quality of service and on the quality of experience.

Chapter 1

Foreword

A 3D streaming system is a system that progressively collects 3D data. The previous chapter voluntarily remained vague about what *3D data* actually are. This chapter presents in detail the 3D data we consider and how they are rendered. We also give insights about interaction and streaming by comparing the 3D setting to the video one.

1.1 What is a 3D model?

The 3D models we are interested in are sets of textured meshes, which can potentially be arranged in a scene graph. Such models can typically contain the following:

- Vertices, which are 3D points,
- Faces, which are polygons defined from vertices (most of the time, they are triangles),
- Textures, which are images that can be used to paint faces in order to add visual richness,
- **Texture coordinates**, which are information added to a face, describing how the texture should be painted over it,
- Normals, which are 3D vectors that can give information about light behaviour on a face.

The Wavefront OBJ is a format that describes all these elements in text format. A 3D model encoded in the OBJ format typically consists in two files: the material file (.mtl) and the object file (.obj).

The material file declares all the materials that the object file will reference. A material consists in name, and other photometric properties such as ambient, diffuse and specular colors, as well as texture maps, which are images that are painted on faces. Each face corresponds to a material. A simple material file is visible on Listing 3.

The object file declares the 3D content of the objects. It declares vertices, texture coordinates and normals from coordinates (e.g. v 1.0 2.0 3.0 for a vertex, vt 1.0 2.0 for a texture coordinate, vn 1.0 2.0 3.0 for a normal). These elements are numbered starting from 1. Faces are declared by using the indices of these elements. A face is a polygon with an arbitrary number of vertices and can be declared in multiple manners:

- f 1 2 3 defines a triangle face that joins the first, the second and the third declared vertices;
- f 1/1 2/3 3/4 defines a similar triangle but with texture coordinates, the first texture coordinate is associated to the first vertex, the third texture coordinate is associated to the second vertex, and the fourth texture coordinate is associated with the third vertex;
- f 1//1 2//3 3//4 defines a similar triangle but referencing normals instead of texture coordinates;
- f 1/1/1 2/3/3 3/4/4 defines a triangle with both texture coordinates and normals.

An object file can include materials from a material file (mtllib path.mtl) and apply the materials that it declares to faces. A material is applied by using the usemtl keyword, followed by the name of the material to use. The faces declared after a usemtl are painted using the material in question. An example of object file is visible on Listing 1.

1	mtllib cube.mtl
2	
3	usemtl cubemtl
4	
5	v -0.5 -0.5 -0.5
6	v -0.5 -0.5 0.5
7	v -0.5 0.5 -0.5
8	v -0.5 0.5 0.5
9	v 0.5 -0.5 -0.5
10	v 0.5 -0.5 0.5
11	v 0.5 0.5 -0.5
12	v 0.5 0.5 0.5
13	
14	vt 0.0 0.0
15	vt 0.0 1.0
16	vt 1.0 0.0
17	vt 1.0 1.0
18	
19	f 1/1 2/3 4/4 3/2
20	f 2/1 6/3 8/4 4/2
21	f 6/1 5/3 7/4 8/2
22	f 5/1 1/3 3/4 7/2
23	f 4/1 8/3 7/4 3/2
24	f 2/1 1/3 5/4 6/2

1	newmtl	cubemtl
2	Ka 1.0	1.0 1.0
3	Kd 1.0	1.0 1.0
4	Ks 1.0	1.0 1.0
5	map_Kd	cube.png

Listing 3: A material file describing a material



Figure 3: A rendering of the cube

Listing 2: An object file describing a cube



1.1.1 Rendering a 3D model

A typical 3D renderer follows Algorithm X. The first task the renderer needs to perform is sending the data to the GPU: this is done in the loading loop during an initialization step. This step can be slow, but it is generally acceptable since it only occurs once at the beginning of the program. Then, the renderer starts the rendering loop: at each frame, it renders the whole scene: for each object, it binds the corresponding material to the GPU and then renders the object. During the rendering loop, there are two things to consider regarding performances:

- the more faces in a geometry, the slower the draw call;
- the more objects in the scene, the more overhead caused by the CPU/GPU communication at each step of the loop.

The way the loop works forces objects with different materials to be rendered separately. An efficient renderer keeps the number of objects in a scene low to avoid introducing overhead. However, an important feature of 3D engines regarding performance is frustum culling. The frustum is the viewing

volume of the camera. Frustum culling consists in skipping the objects that are outside the viewing volume of the camera in the rendering loop. Algorithm X is a variation of Algorithm Y with frustum culling.

A renderer that uses a single object avoids the overhead, but fails to benefit from frustum culling. An optimized renderer needs to find a compromise between a too fine partition of the scene, which introduces overhead, and a too coarse partition, which introduces useless rendering.

- ensures to have objects that do not spread across the whole scene, since that would lead to a useless frustum culling, and many objects to avoid rendering the whole scene at each frame;
- but not too many objects to avoid suffering from the overhead.

1.2 Similarities and differences between video and 3D

The video streaming setting and the 3D streaming setting share many similarities: at a higher level of abstraction, both systems allow a user to access remote content without having to wait until everything is loaded. Analyzing similarities and differences between the video and the 3D scenarios as well as having knowledge about video streaming literature are the key to developing an efficient 3D streaming system.

1.2.1 Chunks of data

In order to be able to perform streaming, data need to be segmented so that a client can request chunks of data and display it to the user while requesting another chunk. In video streaming, data chunks typically consist in a few seconds of video. In mesh streaming, some progressive mesh approaches encode a base mesh that contains low resolution geometry and textures and different chunks that increase the resolution of the base mesh. Otherwise, a mesh can also be segmented by separating geometry and textures, creating chunks that contain some faces of the model, or some other chunks containing textures.

1.2.2 Data persistence

One of the main differences between video and 3D streaming is data persistence. In video streaming, only one chunk of video is required at a time. Of course, most video streaming services prefetch some future chunks, and keep in cache some previous ones, but a minimal system could work without latency and keep in memory only two chunks: the current one and the next one.

Already a few problems appear here regarding 3D streaming:

- depending on the user's field of view, many chunks may be required to perform a single rendering;
- chunks do not become obsolete the way they do in video, a user navigating in a 3D scene may come back to a same spot after some time, or see the same objects but from elsewhere in the scene.

1.2.3 Multiple representations

All major video streaming platforms support multi-resolution streaming. This means that a client can choose the quality at which it requests the content. It can be chosen directly by the user or automatically determined by analyzing the available resources (size of the screen, downloading bandwidth, device performances)



Figure 4: The different qualities available for a Youtube video

Similarly, recent work in 3D streaming have proposed different ways to progressively stream 3D models, displaying a low quality version of the model to the user without latency, and supporting interaction with the model while details are being downloaded. Such strategies are reviewed in Section X.

1.2.4 Media types

Just like a video, a 3D scene is composed of different media types. In video, those media are mostly images, sounds, and subtitles, whereas in 3D, those media are geometry or textures. In both cases, an algorithm for content streaming has to acknowledge those different media types and manage them correctly.

In video streaming, most of the data (in terms of bytes) are used for images. Thus, the most important thing a video streaming system should do is to optimize images streaming. That is why, on a video on Youtube for example, there may be 6 available qualities for images (144p, 240p, 320p, 480p, 720p and 1080p) but only 2 qualities for sound. This is one of the main differences between video and 3D streaming: in a 3D setting, the ratio between geometry and texture varies from one scene to another, and leveraging between those two types of content is a key problem.

1.2.5 Interaction

The ways of interacting with content is another important difference between video and 3D. In a video interface, there is only one degree of freedom: time. The only things a user can do is letting the video

play, pausing, resuming, or jumping to another time in the video. There are also controls for other options that are described on this help page.

All the keyboard shortcuts are summed up in Figure X. Those interactions are different if the user is using a mobile device.

When it comes to 3D, there are many approaches to manage user interaction. Some interfaces mimic the video scenario, where the only variable is the time and the camera follows a predetermined path on which the user has no control. These interfaces are not interactive, and can be frustrating to the user who might feel constrained.

Some other interfaces add 2 degrees of freedom to the timeline: the user does not control the camera's position but can control the angle. This mimics the 360 video scenario. This is typically the case of the video game *nolimits 2: roller coaster simulator* which works with VR devices (oculus rift, HTC vive, etc.) where the only interaction available to the user is turning the head.

Finally, most of the other interfaces give at least 5 degrees of freedom to the user: 3 being the coordinates of the camera's position, and 2 being the angles (assuming the up vector is unchangeable, some interfaces might allow that, giving a sixth degree of freedom). The most common controls are the trackball controls where the user rotate the object like a ball (live example here) and the orbit controls, which behave like the trackball controls but preserving the up vector (live example here). These types of controls are notably used on the popular mesh editor MeshLab and SketchFab, the YouTube for 3D models.



Figure 5: Screenshot of MeshLab

Another popular way of controlling a free camera in a virtual environment is the first person controls (live example here). These controls are typically used in shooting video games, the mouse rotates the camera and the keyboard translates it.

1.2.6 Relationship between interface, interaction and streaming

In both video and 3D systems, streaming affects interaction. For example, in a video streaming scenario, if a user sees that the video is fully loaded, they might start moving around on the timeline, but if they see that the streaming is just enough to not stall, they might prefer not interacting and just watch the video. If the streaming stalls for too long, the user might seek somewhere else hoping for the video to resume, or get frustrated and leave the video. The same types of behaviour occur in 3D streaming: if a user is somewhere in a scene, and sees more data appearing, they might wait until enough data have arrived, but if they see nothing happens, they might leave to look for data somewhere else.

Those examples show how streaming can affect interaction, but interaction also affects streaming. In a video streaming scenario, if a user is watching peacefully without interacting, the system just has to request the next chunks of video and display them. However, if a user starts seeking at a different time of the streaming, the streaming would most likely stall until the system is able to gather the data it needs to resume the video. Just like in the video setup, the way a user navigates in a networked virtual environment affects the streaming. Moving slowly allows the system to collect and display data to the user, whereas moving frenetically puts more pressure on the streaming: the data that the system requested may be obsolete when the response arrives.

Moreover, the interface and the way elements are displayed to the user also impacts his behaviour. A streaming system can use this effect to enhancing the quality of experience by providing feedback on the streaming to the user via the interface. For example, on Youtube, the buffered portion of the video is displayed in light grey on the timeline, whereas the portion that remains to be downloaded is displayed in dark grey. A user is more likely to click on the light grey part of the timeline than on the dark grey part, preventing the streaming from stalling.

1.3 Implementation details

During this thesis, a lot of software has been developed, and for this software to be successful and efficient, we chose appropriate languages. When it comes to 3D streaming systems, we need two kind of software.

- Interactive applications which can run on as many devices as possible so we can easily conduct user studies. For this context, we chose the **JavaScript** language, since it can run on many devices and it has great support for WebGL.
- Native applications which can run fast on desktop devices, in order to prepare data, run simulations and evaluate our ideas. For this context, we chose the **Rust** language, which is a somewhat recent language that provides both the efficiency of C and C++ and the safety of functional languages.

1.3.1 JavaScript

THREE.js

On the web browser, it is now possible to perform 3D rendering by using WebGL. However, WebGL is very low level and it can be painful to write code, even to render a simple triangle. For example, this tutorial's code contains 121 lines of JavaScript, 46 being code (not comments or empty lines) to render a simple, non-textured triangle. For this reason, it seems unreasonable to build a system like the one we are describing in raw WebGL. There are many libraires that wrap WebGL code and that help people building 3D interfaces, and THREE.js is a very popular one (56617 stars on github, making it the 35th most starred repository on GitHub as of November 26th, 2019) THREE.js acts as a 3D engine built on WebGL. It provides classes to deal with everything we need:

- the **Renderer** class contains all the WebGL code needed to render a scene on the web page;
- the **Object** class contains all the boilerplate needed to manage the tree structure of the content, it contains a transform (translation and rotation) and it can have children that are other objects;
- the **Scene** class is the root object, it contains all of the objects we want to render and it is passed as argument to the render function;
- the **Geometry** and **BufferGeometry** classes are the classes that hold the vertex buffers, we will discuss it more in the next paragraph;
- the **Material** class is the class that holds the properties used to render geometry (the most important information being the texture), there are many classes derived from Material, and the developer can choose what material they want for their objects;
- the **Mesh** class is the class that links the geometry and the material, it derives the Object class and can thus be added to a scene and rendered.

A snippet of the basic usage of these classes is given in Listing 4.

```
1
     // Computes the aspect ratio of the window.
2
     let aspectRatio = window.innerWidth / window.innerHeight;
 3
 4
     // Creates a camera and sets its parameters and position.
     let camera = new THREE.PerspectiveCamera(70, aspectRatio, 0.01, 10);
 5
6
     camera.position.z = 1;
7
8
     // Creates the scene that contains our objects.
9
     let scene = new THREE.Scene();
10
11
     // Creates a geometry (vertices and faces) corresponding to a cube.
     let geometry = new THREE.BoxGeometry(0.2, 0.2, 0.2);
12
13
     // Creates a material that paints the faces depending on their normal.
14
15
     let material = new THREE.MeshNormalMaterial();
16
17
     // Creates a mesh that associates the geometry with the material.
     let mesh = new THREE.Mesh(geometry, material);
18
19
20
     // Adds the mesh to the scene.
21
     scene.add(mesh);
22
23
     // Creates the renderer and append its canvas to the DOM.
     renderer = new THREE.WebGLRenderer({ antialias: true });
24
     renderer.setSize(window.innerWidth, window.innerHeight);
25
     document.body.appendChild(renderer.domElement);
26
27
28
     // Renders the scene with the camera.
29
     renderer.render(scene, camera);
```

Listing 4: A THREE.js hello world

Geometries

Geometries are the classes that hold the vertices, texture coordinates, normals and faces. THREE.js proposes two classes for handling geometries:

- the **Geometry** class, which is made to be developer friendly and allows easy editing but can suffer from performance issues;
- the **BufferGeometry** class, which is harder to use for a developer, but allows better performance since the developer controls how data is transmitted to the GPU.

1.3.2 Rust

In this section, we explain the specificities of Rust and why it is an adequate language for writing efficient native software safely.

Borrow checker

Rust is a system programming language focused on safety. It is made to be efficient (and effectively has performances comparable to C but with some extra features. C++ users might see it as a language like C++ but that forbids undefined behaviours. The most powerful concept from Rust is *ownership*. Basically, every value has a variable that we call its *owner*. To be able to use a value, you must either be its owner or borrow it. There are two types of borrow, the immutable borrow and the mutable borrow (roughly equivalent to references in C++). The compiler comes with the *borrow checker* which makes sure you only use variables that you are allowed to use. For example, the owner can only use the value if it is not being borrowed, and it is only possible to either mutably borrow a value once, or immutably borrow a value many times. At first, the borrow checker seems particularly efficient to detect bugs in concurrent software, but in fact, it is also decisive in non concurrent code. Consider the piece of C++ code in Listing 5 and Listing 6.

```
1 auto vec = std::vector<int> {1, 2, 3};
2 for (auto value: vec)
3 vec.push back(value);
```

Listing 5: Undefined behaviour with for each syntax

```
1 auto vec = std::vector<int> {1, 2, 3};
2 for (auto it = std::begin(vec); it < std::end(vec); it++)
3 vec.push_back(*it);
```

Listing 6: Undefined behaviour with iterator syntax

This loop should go endlessly because the vector grows in size as we add elements in the loop. But the most important thing here is that since we add elements to the vector, it will eventually need to be reallocated, and that reallocation will invalidate the iterator, meaning that the following iteration will provoke an undefined behaviour. The equivalent code in Rust is in Listing 7 and Listing 8.

```
1 let mut vec = vec![1, 2, 3];
2 for value in &vec {
3 vec.push(value);
4 }
```

Listing 7: Rust version of Listing 5

```
1
    let mut vec = vec![1, 2, 3];
2
    let iter = vec.iter();
3
    loop {
4
        match iter.next() {
5
             Some(x) => vec.push(x),
6
             None => break,
7
        }
8
    }
```

Listing 8: Rust version of Listing 6

What happens is that the iterator needs to borrow the vector. Because it is borrowed, it can no longer be borrowed as mutable since mutating it could invalidate the other borrowers. And effectively, the borrow checker will crash the compiler with the error in Listing 9.

```
error[E0502]: cannot borrow `vec` as mutable because it is also borrowed as
 1
     immutable
      --> src/main.rs:4:9
2
3
4
     3 |
             for value in &vec {
5
                          - - - -
6
       L
                          T
                          immutable borrow occurs here
7
       L
                          immutable borrow later used here
8
       L
9
                 vec.push(*value);
     4
                 mutable borrow occurs here
10
       T
```

Listing 9: Error given by the compiler on Listing 7

This example is one of the many examples of how powerful the borrow checker is: in Rust code, there can be no dangling reference, and all the segmentation faults coming from them are detected by the compiler. The borrow checker may seem like an enemy to newcomers because it often rejects code that seem correct, but once they get used to it, they understand what is the problem with their code and either fix the problem easily, or realize that the whole architecture is wrong and understand why.

It is probably for those reasons that Rust is the *most loved programming language* according to the Stack Overflow Developer Survey

Tooling

Moreover, Rust comes with many programs that help developers.

- **rustc** is the Rust compiler. It is comfortable due to the clarity and precise explanations of its error messages.
- **cargo** is the official Rust's project and package manager. It manages compilation, dependencies, documentation, tests, etc.
- **racer**, **rls** (Rust Language Server) and **rust-analyzer** are software that manage automatic compilation to display errors in code editors as well as providing semantic code completion.
- **rustfmt** auto formats code.
- **clippy** is a linter that detects unidiomatic code and suggests modifications.

Glium

When we need to perform rendering for 3D content analysis or for evaluation, we use the glium library. Glium has many advantages over using raw OpenGL calls. Its objectives are:

• to be easy to use: it exposes functions that are higher level than raw OpenGL calls, but still low enough level to let the developer free;

- to be safe: debugging OpenGL code can be a nightmare, and glium does its best to use the borrow checker to its advantage to avoid OpenGL bugs;
- to be fast: the binary produced use optimized OpenGL functions calls;
- to be compatible: glium seeks to support the latest versions of OpenGL functions and falls back to older functions if the most recent ones are not supported on the device.

Conclusion

In our work, many tasks will consist in 3D content analysis, reorganization, rendering and evaluation. Many of these tasks require long computations, lasting from hours to entire days. To perform them, we need a programming language that has good performances. In addition, the extra features that Rust provides ease tremendously development, and this is why we use Rust for all tasks that do not require having a web interface.

Chapter 2

Related work

In this chapter, we review the part of the state of the art on multimedia streaming and interaction that is relevant for this thesis. As discussed in the previous chapter, video and 3D share many similarities and since there is already a very important body of work on video streaming, we start this chapter with a review of this domain with a particular focus on the DASH standard. Then, we proceed with presenting topics related to 3D streaming, including compression and streaming, geometry and texture compromise, and viewpoint dependent streaming. Finally, we end this chapter by reviewing the related work regarding 3D navigation and interfaces.

2.1 Video

Accessing a remote video through the web has been a widely studied problem since the 1990s. The Real-time Transport Protocol (RTP, (Schulzrinne et al. 1996)) has been an early attempt to formalize audio and video streaming. The protocol allowed data to be transferred unilaterally from a server to a client, and required the server to handle a separate session for each client.

In the following years, HTTP servers have become ubiquitous, and many industrial actors (Apple, Microsoft, Adobe, etc.) developed HTTP streaming systems to deliver multimedia content over the network. In an effort to bring interoperability between all different actors, the MPEG group launched an initiative, which eventually became a standard known as DASH, Dynamic Adaptive Streaming over HTTP. Using HTTP for multimedia streaming has many advantages over RTP. While RTP is stateful (that is to say, it requires keeping track of every user along the streaming session), HTTP is stateless, and thus more efficient. Furthermore, an HTTP server can easily be replicated at different geographical locations, allowing users to fetch data from the closest server. This type of network architecture is called CDN (Content Delivery Network) and increases the speed of HTTP requests, making HTTP based multimedia streaming more efficient.

2.1.1 DASH: the standard for video streaming

Dynamic Adaptive Streaming over HTTP (DASH), or MPEG-DASH (Stockhammer 2011), (Sodagar 2011)) is now a widely deployed standard for adaptively streaming video on the web (DASH 2014), made to be simple, scalable and inter-operable. DASH describes guidelines to prepare and structure video content, in order to allow a great adaptability of the streaming without requiring any server side computation. The client should be able to make good decisions on what part of the content to download, only based on an estimation of the network constraints and on the information provided in a descriptive file: the MPD.

DASH structure

All the content structure is described in a Media Presentation Description (MPD) file, written in the XML format. This file has 4 layers: the periods, the adaptation sets, the representations, and the segments. An MPD has a hierarchical structure, meaning it has multiple periods, and each period can have multiple adaptation sets, each adaptation set can have multiple representation, and each representation can have multiple segments.

Periods

Periods are used to delimit content depending on time. It can be used to delimit chapters, or to add advertisements that occur at the beginning, during or at the end of a video.

Adaptation sets

Adaptation sets are used to delimit content according to the format. Each adaptation set has a mimetype, and all the representations and segments that it contains share this mime-type. In videos, most of the time, each period has at least one adaptation set containing the images, and one adaptation set containing the sound. It may also have an adaptation set for subtitles.

Representations

The representation level is the level DASH uses to offer the same content at different levels of quality. For example, an adaptation set containing images has a representation for each available quality (it might be 480p, 720p, 1080p, etc.). This allows a user to choose its representation and change it during the video, but most importantly, since the software is able to estimate its downloading speed based on the time it took to download data in the past, it is able to find the optimal representation, being the highest quality that the client can request without stalling.

Segments

Until this level in the MPD, content has been divided but it is still far from being sufficiently divided to be streamed efficiently. A representation of the images of a chapter of a movie is still a long video, and keeping such a big file is not possible since heavy files prevent streaming adaptability: if the user requests to change the quality of a video, the system would either have to wait until the file is totally downloaded, or cancel the request, making all the progress done unusable.

Segments are used to prevent this issue. They typically encode files that contain two to ten seconds of video, and give the software a greater ability to dynamically adapt to the system. If a user wants to seek somewhere else in the video, only one segment of data is potentially lost, and only one segment of data needs to be downloaded for the playback to resume. The impact of the segment duration has been investigated in many work, including (Sideris et al. 2015)", "stohr2017sweet. For example, (Stohr et al. 2017) discuss how the segment duration affects the streaming: short segments lower the initial delay and provide the best stalling quality of experience, but make the total downloading time of the video longer because of overhead.

Content preparation and server

Encoding a video in DASH format consists in partitioning the content into periods, adaptation sets, representations and segments as explained above, and generating a Media Presentation Description file (MPD) which describes this organization. Once the data are prepared, they can simply be hosted on a static HTTP server which does no computation other than serving files when it receives requests. All the intelligence and the decision making is moved to the client side. This is one of the DASH strengths:

no powerful server is required, and since static HTTP server are mature and efficient, all DASH clients can benefit from it.

Client side adaptation

A client typically starts by downloading the MPD file, and then proceeds on downloading segments from the different adaptation sets. While the standard describes well how to structure content on the server side, the client may be freely implemented to take into account the specificities of a given application. The most important part of any implementation of a DASH client is called the adaptation logic. This component takes into account a set of parameters, such as network conditions (bandwidth, throughput, for example), buffer states or segments size to derive a decision on which segments should be downloaded next. Most of the industrial actors have their own adaptation logic, and many more have been proposed in the literature. A thorough review is beyond the scope of this state-of-the-art, but examples include (Chiariotti et al. 2016) who formulate the problem using queuing theory, or (Huang et al. 2019) who use a formulation derived from the knapsack problem.

2.1.2 DASH-SRD

Being now widely adopted in the context of video streaming, DASH has been adapted to various other contexts. DASH-SRD (Spatial Relationship Description, (Niamut et al. 2016)) is a feature that extends the DASH standard to allow streaming only a spatial subpart of a video to a device. It works by encoding a video at multiple resolutions, and tiling the highest resolutions as shown in Figure ref{sota:srd-png}. That way, a client can choose to download either the low resolution of the whole video or higher resolutions of a subpart of the video.



Figure 6: DASH-SRD (Niamut, Thomas, D'Acunto, Concolato, Denoual, and Lim 2016)

For each tile of the video, an adaptation set is declared in the MPD, and a supplemental property is defined in order to give the client information about the tile. This supplemental property contains

many elements, but the most important ones are the position (x and y) and the size (width and height) describing the position of the tile in relation to the full video. An example of such a property is given in Listing 10.

1	<period></period>
2	<adaptationset></adaptationset>
3	<supplementalproperty <="" schemeiduri="urn:mpeg:dash:s rd:2014" th=""></supplementalproperty>
	value="0,0,0,5760,3240,5760,3240"/>
4	<role schemeiduri="urn:mpeg:dash:role:2011" value="main"></role>
5	<representation height="2160" id="1" width="3840"></representation>
6	<baseurl>full.mp4</baseurl>
7	
8	
9	<adaptationset></adaptationset>
10	<supplementalproperty <="" schemeiduri="urn:mpeg:dash:s rd:2014" th=""></supplementalproperty>
	value="0,1920,1080,1920,1080,5760,3240"/>
11	<role schemeiduri="urn:mpeg:dash:role:2011" value="supplementary"></role>
12	<representation height="1080" id="2" width="1920"></representation>
13	<baseurl>part.mp4</baseurl>
14	
15	
16	

Listing 10: MPD of a video encoded using DASH-SRD

Essentially, this feature is a way of achieving view-dependent streaming, since the client only displays a part of the video and can avoid downloading content that will not be displayed. While Figure ref{sota:srd-png} illustrates how DASH-SRD can be used in the context of zoomable video streaming, the ideas developed in DASH-SRD have proven to be particularly useful in the context of 360 video streaming (see for example (Ozcinar, De Abreu, and Smolic 2017)). This is especially interesting in the context of 3D streaming since we have this same pattern of a user viewing only a part of a content.

2.2 3D streaming

In this thesis, we focus on the objective of delivering large, massive 3D scenes over the network. While 3D streaming is not the most popular research field, there has been a special attention around 3D content compression, in particular progressive compression which can be considered a premise for 3D streaming. In the next sections, we review the 3D streaming related work, from 3D compression and structuring to 3D interaction.

2.2.1 Compression and structuring

According to (Maglo et al. 2015), mesh compression can be divided into four categories:

• single-rate mesh compression, seeking to reduce the size of a mesh;

- progressive mesh compression, encoding meshes in many levels of resolution that can be down-loaded and rendered one after the other;
- random accessible mesh compression, where different parts of the models can be decoded in an arbitrary order;
- mesh sequence compression, compressing mesh animations.

Since our objective is to stream 3D static scenes, single-rate mesh and mesh sequence compressions are less interesting for us. This section thus focuses on progressive meshes and random accessible mesh compression.

Progressive meshes were introduced in (Hoppe 1996) and allow a progressive transmission of a mesh by sending a low resolution mesh first, called *base mesh*, and then transmitting detail information that a client can use to increase the resolution. To do so, an algorithm, called *decimation algorithm*, starts from the original full resolution mesh and iteratively removes vertices and faces by merging vertices through the so-called *edge collapse* operation (Figure X).

Every time two vertices are merged, a vertex and two faces are removed from the original mesh, decreasing the model resolution. At the end of this content preparation phase, the mesh has been reorganized into a base mesh and a sequence of partially ordered edge split operations. Thus, a client can start by downloading the base mesh, display it to the user, and keep downloading refinement operations (vertex splits) and display details as time goes by. This process reduces the time a user has to wait before seeing a downloaded 3D object, thus increases the quality of experience.



Figure 7: Four levels of resolution of a mesh

(Lavoué, Chevalier, and Dupont 2013) develop a dedicated progressive compression algorithm based on iterative decimation, for efficient decoding, in order to be usable on web clients. With the same objective, (Limper et al. 2013) proposes pop buffer, a progressive compression method based on quantization that allows efficient decoding.

Following these, many approaches use multi triangulation, which creates mesh fragments at different levels of resolution and encodes the dependencies between fragments in a directed acyclic graph. In (Cignoni et al. 2005), the authors propose Nexus: a GPU optimized version of multi triangulation that pushes its performances to make real time rendering possible. It is notably used in 3DHOP (3D Heritage Online Presenter, (Potenziani et al. 2015)), a framework to easily build web interfaces to present 3D objects to users in the context of cultural heritage. Each of these approaches define its own compression and coding for a single mesh. However, users are often interested in scenes that contain multiple meshes, and the need to structure content emerged.

To answer those issues, the Khronos group proposed a generic format called glTF (GL Transmission Format, (Robinet and Cozzi 2013)) to handle all types of 3D content representations: point clouds, meshes, animated models, etc. glTF is based on a JSON file, which encodes the structure of a scene of 3D objects. It contains a scene graph with cameras, meshes, buffers, materials, textures and animations. Although relevant for compression, transmission and in particular streaming, this standard does not yet consider view-dependent streaming, which is required for large scene remote visualization and which we address in our work.

2.2.2 Viewpoint dependency

3D streaming means that content is downloaded while the user is interacting with the 3D object. In terms of quality of experience, it is desirable that the downloaded content falls into the user's field of view. This means that the progressive compression must encode a spatial information in order to allow the decoder to determine content adapted to its viewpoint. This is typically called *random accessible mesh compression*. (Maglo, Grimstead, and Hudelot 2013) is such an example of random accessible progressive mesh compression. (Cheng and Ooi 2008) proposes a receiver driven way of achieving viewpoint dependency with progressive mesh: the client starts by downloading the base mesh, and from then is able to estimate the importance of the different vertex splits, in order to choose which ones to download. Doing so drastically reduces the server computational load, since it only has to send data, and improves the scalability of this framework.

In the case of streaming a large 3D scene, view-dependent streaming is fundamental: a user will only be seeing one small portion of the scene at each time, and a system that does not adapt its streaming to the user's point of view is bound to induce a low quality of experience.

A simple way to implement viewpoint dependency is to request the content that is spatially close to the user's camera. This approach, implemented in Second Life and several other NVEs (e.g., (Liang, Zimmermann, and Ooi 2011)), only depends on the location of the avatar, not on its viewing direction. It exploits spatial coherence and works well for any continuous movement of the user, including turning. Once the set of objects that are likely to be accessed by the user is determined, the next question is in what order should these objects be retrieved. A simple approach is to retrieve the objects based on distance: the spatial distance from the user's virtual location and rotational distance from the user's view.

More recently, Google integrated Google Earth 3D module into Google Maps (Figure 8). Users are now able to go to Google Maps, and click the 3D button which shifts the camera from the aerial view. Even though there are no associated publications to support this assertion, it seems clear that the streaming is view-dependent: low resolution from the center of the point of view gets downloaded first, and higher resolution data gets downloaded for closer objects than for distant ones.



Figure 8: Screeshot of the 3D interface of Google Maps

Other approaches use level of details. Level of details have been initially used for efficient 3D rendering (Lindstrom et al. 1996). When the change from one level of detail to another is direct, it can create visual discomfort to the user. This is called the *popping effect* and level of details have the advantage of enabling techniques, such as geomorphing (Hoppe 1998), to transition smoothly from one level of detail to another. Level of details have then been used for 3D streaming. For example, (Guthe and Klein 2004) propose an out-of-core viewer for remote model visualization based by adapting hierarchical level of details (Erikson, Manocha, and Baxter III 2001) to the context of 3D streaming. Level of details can also be used to perform viewpoint dependant streaming, such as (Meng and Zha 2003).

2.2.3 Texture streaming

In order to increase the texture rendering speed, a common technique is the *mipmapping* technique. It consists in generating progressively lower resolutions of an initial texture. Lower resolutions of the textures are used for polygons which are far away from the camera, and higher resolutions for polygons closer to the camera. Not only this reduces the time needed to render the polygons, but it can also reduce the aliasing effect. Using these lower resolutions can be especially interesting for streaming. (Marvie and Bouatouch 2003) proposes the PTM format which encode the mipmap levels of a texture that can be downloaded progressively, so that a lower resolution can be shown to the user while the higher resolutions are being downloaded.

Since 3D data can contain many textures, (Simon et al. 2019) propose a way to stream a set of textures by encoding them into a video. Each texture is segmented into tiles of a fixed size. Those tiles are then ordered to minimize dissimilarities between consecutive tiles, and encoded as a video. By benefiting from the video compression techniques, the authors are able to reach a better rate-distortion ratio than webp, which is the new standard for texture transmission, and jpeg.

2.2.4 Geometry and textures

As discussed in Chapter 1.1, most 3D scenes consist in two main types of data: geometry and textures. When addressing 3D streaming, one must handle the concurrency between geometry and textures, and the system needs to address this compromise.

Balancing between streaming of geometry and texture data is addressed by (Tian and AlRegib 2008)"), (Guo et al. 2017), and (Yang, Lee, and Kuo 2004). Their approaches combine the distortion caused by having lower resolution meshes and textures into a single view independent metric. (Portaneri et al. 2019) also deals with the geometry / texture compromise. This work designs a cost driven framework for 3D data compression, both in terms of geometry and textures. The authors generate an atlas for textures that enables efficient compression and multi-resolution scheme. All four works considered a single mesh, and have constraints on the types of meshes that they are able to compress. Since the 3D scenes we are interested in in our work consist in soups of textured polygons, those constraints are not satisfied and we cannot use those techniques.

2.2.5 Streaming in game engines

In traditional video games, including online games, there is no requirement for 3D data streaming. Video games either come with a physical support (CD, DVD, Blu-Ray) or they require the downloading of the game itself, which includes the 3D data, before letting the user play. However, transferring data from the disk to the memory is already a form of streaming. This is why optimized engines for video games use techniques that are reused for streaming such as level of details, to reduce the details of objects far away for the point of view and save the resources to enhance the level of detail of closer objects.

Some other online games, such as Second Life, rely on user generated data, and thus are forced to send data from users to others. In such scenarios, 3D streaming is appropriate and this is why the idea of streaming 3D content for video games has been investigated. For example, (Li et al. 2011) proposes an online game engine based on geometry streaming, that addresses the challenge of streaming 3D content at the same time as synchronization of the different players.

2.2.6 NVE streaming frameworks

An example of NVE streaming framework is 3D Tiles (Schilling, Bolling, and Nagel 2016), which is a specification for visualizing massive 3D geospatial data developed by Cesium and built on top of glTF. Their main goal is to display 3D objects on top of regular maps, and their visualization consists in a top-down view, whereas we seek to let users freely navigate in our scenes, whether it be flying over the scene or moving along the roads.



Figure 9: Screenshot of 3D Tiles interface

3D Tiles, as its name suggests, is based on a spacial partitionning of the scene. It started with a regular octree, but has then been improved to a k-d tree (see Figure ref{sote:3d-tiles-partition}).



Figure 10: With regular octree (depth 4)



Figure 11: With k-d tree (depth 4)

In citeyear{3d-tiles-10x}, 3D Tiles streaming system was improved by preloading the data at the camera's next position when known in advance (with ideas that are similar to those we discuss and implement in Chapter ref{bi}, published in citeyear{bookmarks-impact}) and by ordering tile requests depending on the user's position (with ideas that are similar to those we discuss and implement in Chapter ref{d3}, published in citeyear{dash-3d}).
(Zampoglou et al. 2018) is another example of a streaming framework: it is the first paper that proposes to use DASH to stream 3D content. In their work, the authors describe a system that allows users to access 3D content at multiple resolutions. They organize the content, following DASH terminology, into periods, adaptation sets, representations and segments. Their first adaptation set codes the tree structure of the scene graph. Each further adaptation set contains both geometry and texture information and is available at different resolutions defined in a corresponding representation. To avoid requests that would take too long and thus introduce latency, the representations are split into segments. The authors discuss the optimal number of polygons that should be stored in a single segment. On the one hand, using segments containing very few faces will induce many HTTP requests from the client, and will lead to poor streaming efficiency. On the other hand, if segments contain too many faces, the time to load the segment is long and the system loses adaptability. Their approach works well for several objects, but does not handle view-dependent streaming, which is desirable in the use case of large NVEs.

2.3 3D bookmarks and navigation aids

One of the uses for 3D streaming is to allow users interacting with the content while it is being downloaded. However, devising an ergonomic technique for browsing 3D environments through a 2D interface is difficult. Controlling the viewpoint in 3D (6 DOFs) with 2D devices is not only inherently challenging but also strongly task-dependent. In their review, (Jankowski and Hachet 2015) distinguish between several types of camera movements: general movements for exploration (e.g., navigation with no explicit target), targeted movements (e.g., searching and/or examining a model in detail), specified trajectory (e.g., a cinematographic camera path). For each type of movement, specialized 3D interaction techniques can be designed. In most cases, rotating, panning, and zooming movements are required, and users are consequently forced to switch back and forth among several navigation modes, leading to interactions that are too complicated overall for a layperson. Navigation aids and smart widgets are required and subject to research efforts both in 3D companies (see sketchfab.com, cl3ver.com among others) and in academia, as reported below.

Translating and rotating the camera can be simply specified by a *lookat* point. This is often known as point-of-interest (POI) movement (or *go-to*, *fly-to* interactions) (Mackinlay, Card, and Robertson 1990). Given such a point, the camera automatically moves from its current position to a new position that looks at the POI. One key issue of these techniques is to correctly orient the camera at destination. In Unicam (Zeleznik, Forsberg, and Strauss 1997), the so-called click-to-focus strategy automatically chooses the destination viewpoint depending on 3D orientations around the contact point. The more recent Drag'n Go interaction (Moerman, Marchal, and Grisoni 2012) also hits a destination point while offering control on speed and position along the camera path. This 3D interaction is designed in the screen space (it is typically a mouse-based camera control), where cursor's movements are mapped to camera movements following the same direction as the on-screen optical-flow.



Figure 12: Screenshot of the drag'n go interface (Moerman, Marchal, and Grisoni 2012) (the percentage widget is for illustration)

Some 3D browsers provide a viewpoint menu offering a choice of viewpoints (Todd 2004), (Burtnyk et al. 2006). Authors of 3D scenes can place several viewpoints (typically for each POI) in order to allow easy navigation for users, who can then easily navigate from viewpoint to viewpoint just by selecting a menu item. Such viewpoints can be either static, or dynamically adapted: (Jankowski and Decker 2012) report that users clearly prefer navigating in 3D using a menu with animated viewpoints than with static ones.



Figure 13: Screenshot of an interface with menu for navigation (Burtnyk, Khan, Fitzmaurice, and Kurtenbach 2006)

Early 3D VRML environments (Rezzonico and Thalmann 1996) offer 3D bookmarks with animated transitions between bookmarked views. These transitions prevent disorientation since users see how they got there. Hyperlinks can also ease rapid movements between distant viewpoints and naturally support non-linear and non-continuous access to 3D content. Navigating with 3D hyperlinks is faster due to the instant motion, but can cause disorientation, as shown by the work of (Ruddle et al. 2000). (Eno, Gauch, and Thompson 2010) examine explicit landmark links as well as implicit avatar-chosen links in Second Life. These authors point out that linking is appreciated by users and that easing linking would likely result in a richer user experience. (Jankowski and Decker 2012) developed the Dual-Mode

User Interface (DMUI) that coordinates and links hypertext to 3D graphics in order to access information in a 3D space.



Figure 14: The two modes of DMUI (Jankowski and Decker 2012)

The use of in-scene 3D navigation widgets can also facilitate 3D navigation tasks. (Chittaro and Venkataraman 2006) propose and evaluate 2D and 3D maps as navigation aids for complex virtual buildings and find that the 2D navigation aid outperforms the 3D one for searching tasks. The View-Cube widget (Khan et al. 2008) serves as a proxy for the 3D scene and offers viewpoint switching between 26 views while clearly indicating associated 3D orientations. Interactive 3D arrows that point to objects of interest have also been proposed as navigation aids in (Chittaro and Burigat 2004), (Burigat and Chittaro 2007)): when clicked, the arrows transfer the viewpoint to the destination through a simulated walk or a faster flight.

Chapter 3

Bookmarks, navigation and streaming



Figure 16: Viewport display of a bookmark



Figure 18: A preview is shown when the mouse hovers a bookmark



Figure 17: Arrow display of a bookmark



Figure 19: A coin is hidden behind the curtain

Figure 15: 3D bookmarks propose to move to a new viewpoint; when the user clicks on the bookmark, his viewpoint moves to the indicated viewpoint.

In this chapter, we present our first contribution: an analysis of the impact of bookmarks on navigation and streaming.

We implement a 3D navigation interface where the keyboard translates the camera and the mouse rotates it. We augment this interface with 3D bookmarks. When the user's cursor hovers a bookmark, a preview of the point of view is displayed to the user, and when the user clicks on a bookmark, the camera smoothly moves from its current position to the bookmarked point of view. We conduct a within-subject user-study on 51 participants, where each user starts with a tutorial to get used to the 3D navigation controls, and then tries successively to perform a task with and without bookmarks. We show that not only the presence of bookmarks causes a faster task completion, but also that it allows users to see a larger part of the scene during the same time span.

However, in a streaming scenario, this phenomenon leads to higher network requirements to maintain the same quality of service. In the last part of this chapter, we simulate a streaming setup and we show that knowing the positions of the bookmarks beforehand allows us to precompute information that we can reuse during streaming to compensate for the harm caused by the faster navigation with bookmarks.

3.1 Introduction

Navigating in NVE with a large virtual space (most times through a 2D interface) is sometimes cumbersome. In particular, a user may have difficulties reaching the right place to find information. The content provider of the NVE may want to highlight certain interesting features for the users to view and experience, such as a vantage point in a city, an excavation at an archaeological site, or an exhibit in a museum. To allow users to easily find these interesting locations within the NVE, *3D bookmarks* or *bookmarks* for short, can be provided. A bookmark is simply a 3D virtual camera (with position and camera parameters) predefined by the content provider, and can be presented to users in different ways, including as a text link (URL), a thumbnail image, or a 3D object embedded within the NVE itself.

When users click on a bookmark, NVEs commonly provide a "fly-to" animation to transit the camera from the current viewpoint to the destination (Mackinlay, Card, and Robertson 1990; Rezzonico and Thalmann 1996) to help orient the users within the 3D space. Clicking on a bookmark to fly to another viewpoint leads to reduced data locality. The 3D content at the bookmarked destination viewpoint may overlap less with the current viewpoint. In the worst case, the 3D objects corresponding to the current and destination viewpoints can be completely disjoint. Such movement to a bookmark may lead to a *discovery latency* (Varvello et al. 2011), in which users have to wait for the 3D content for the new viewpoint to be loaded and displayed. An analogy for this situation, in the context of video streaming, is seeking into a segment of video that has not been prefetched yet.

In this chapter, we explore the impact of bookmarks on NVE navigation and streaming, and make several contributions. First, we conducted a crowdsourcing experiment where 51 participants navigated in 3 virtual scenes to complete a task. This experiment serves two purposes: (i) it validates our intuition that bookmarks significantly reduce the number of interactions and navigation time (in average the time needed to complete the task for users with bookmarks is half the time for users without bookmarks); (ii) it produces a set of user interaction traces that we use for subsequent simulation experiments. Second, we quantified the effect of bookmarking on prefetching and visual quality in our experiments. We showed that, without prefetching, the number of correctly rendered pixels right after clicking on bookmarks can drop up to 10% on average. If we prefetch the 3D content from the bookmarks according to the probability of access, we do not limit this drop by more than 5%. Finally, we proposed a method to improve the visual quality after clicking on bookmarks, by exploiting the fact that the visible faces at the bookmark can be precomputed, and by fetching the visible faces only after a bookmark is clicked. We showed that, if the fetching is done during the 1 or 2 seconds of the "flyto" camera movement from the current viewpoint to the bookmarked viewpoint, it suffices to increase the number of correctly rendered pixels to more than 20%, without wasting bandwidth on prefetching. Our key message is that, in addition to easing navigation, bookmarking allows precomputation of visible faces and can significantly reduce interaction latency, without resorting to prefetching, which may waste bandwidth by prefetching 3D data that will not be needed.

The rest of the chapter consists of the following sections. Section X describes the 3D bookmarks that we use in our work, along with our experiments to validate the usefulness of bookmarking. Section X describes the streaming and prefetching mechanisms that we used to simulate our experiments as well as our main findings. Finally, we conclude in Section X.

3.2 Impact of 3D bookmarks on navigation

We now describe an experiment that we conducted on 51 participants, with two goals in mind. First, we want to measure the impact of 3D bookmarks on navigation within an NVE@. Second, we want to collect traces from the users so that we can replay them for reproducible experiments for comparing streaming strategies in Section ref{bi:system}.

3.2.1 Our NVE

To ease the deployment of our experiments to users in distributed locations on a crowdsourcing platform, we implement a simple web-based NVE client using THREE.js The NVE server is implemented with node.js. The NVE server streams a 3D scene to the client; the client renders the scene as the 3D content is received.

The user can navigate within the NVE in the following way; the camera can be translated using the arrow keys along four directions: forward, backward, to the left, and to the right. Alternatively, the keys W, A, S and D can also be used for the same actions. This choice was inspired by 3D video games, which often use these keys in conjunction with the mouse to move an avatar. The virtual camera can rotate in four different directions using the keys I, K, J and L. The user can also rotate the camera by dragging the mouse in the desired direction. Finally, following the UI of popular 3D games, we also give users the possibility to lock their pointer and use their mouse as a virtual camera. The mouse movement controls the camera rotation. The user can always choose to lock the pointer, or unlock it using the escape key. The interface also includes a button to reset the camera back to the starting position in the scene.

3.2.2 3D bookmarks

Our NVE supports 3D bookmarks. A 3D bookmark, or bookmark for short, is simply a fixed camera location (in 3D space), a view direction, and a focal. Bookmarks visible from the user's current view-point are shown as 3D objects in the scene. Figure X depicts some bookmarks from our NVE.

The user can click on a bookmark object to automatically move and align its viewpoint to that of the bookmark. The movement follows a Hermite curve joining the current viewpoint to the viewpoint of the bookmark. The tangent of the curve is the view direction. The user can hover the mouse pointer over a bookmark object to see a thumbnail view of the 3D scene as seen from the bookmark. (Figure X, bottom left).

In our work, we consider two different possibilities for displaying bookmarks: viewports (Figure X top left) and arrows (Figure X top right). A viewport is displayed as a pyramid where the top corre-

sponds to the optical center of its viewpoint and the base corresponds to its image plane. The arrows are view dependent. The bottom of the arrow turns towards the current position, to better visualize the relative position of the bookmark.

Bookmarks allow the user to achieve a large movement within the 3D environment using a single action (a mouse click). Since bookmarks are part of the scene, they are visible only when not hidden by other objects from the scene. We chose size and colors that are salient enough to be easily seen, but not too large to limit the occlusion of regions within the scene. When reaching the bookmark, the corresponding arrow or viewport is not visible anymore, and subsequently will appear in a different color, to indicate that it has been clicked (similar to web links).

3.2.3 User study

We now describe in details our experimental setup and the user study that we conducted on 3D navigation.

Models

We use four 3D scenes (one for the tutorial and three for the actual experiments) which represent recreated scenes from a famous video game. Those models are light (a few thousand of triangles per model) and are sent before the experiment starts. We keep the models small so that users can perform the task with acceptable latency from any country using a decent internet connection. Our NVE does not stream the 3D content for these experiments, in order to avoid unreliable conditions caused by the network bandwidth variation, which might affect how the users interact.

Task design

Since we are interested in studying how efficiently users navigate in the 3D scene, we ask our participants to complete a task which forces them to visit, at least partially, various regions in the scene. To this end, we hide a set of 8 coins on the scene: participants are asked to collect the coins by clicking on them. In order to avoid any bias due to the coins position, we predefined 50 possible coin locations per scene, and randomly select 8 out of these 50 positions each time a new participant starts the experiment.

Experiment

Participants are first presented with an initial screen to collect some preliminary information: age, gender, the last time they played 3D video games, and self-rated 3D gaming skills. We ask those questions because we believe that someone who is used to playing 3D video games should browse the scene more easily, and thus, may not need to use our bookmarks.

Then, the participants go through a tutorial to learn how the UI works, and how to complete the task. The different interactions (keyboard navigation, mouse navigation, bookmarks interaction) are progressively introduced to participants, and the tutorial ends once the participant completes an easy version of the task. The tutorial is always performed on the same scene.

Then, each participant has to complete the task three times. Each task is performed on a different scene, with a different interface. Three interfaces are used. A **NoReco** interface lets the participant navigates without any bookmarks. The other two interfaces allow a participant to move using bookmarks displayed as viewports (denoted as **Viewports**) and arrows (denoted as **Arrows**) respectively.

The coins are chosen randomly, based on the coin configurations that were used by previous participants: if another participant has done an experiment with a certain set of coins, on a certain scene, with a certain type of bookmarks, the current participant will do the experiment with the same set of coins, on the same scene, but with a different type of bookmarks. This policy allows us to limit the bias that could be caused by coin locations.

Once a participant has found all coins, a button is shown on the interface to let the participant move to the next step. Alternatively, this button may appear one minute after the sixth coin was found. This means that a user is authorized to move on without completing the task, in order to avoid potential frustration caused by not finding the remaining two coins.

After completing the three tasks, the participants have to answer a set of questions about their experience with the bookmarks (we refer to the bookmarks as *recommendations* in the experiments). Table 1 shows the list of questions.

	Questions	Answers
1	What was the difficulty level WITHOUT recommendation?	3.04 / 5 ± 0.31
2	What was the difficulty level WITH recommendation?	2.15 / 5, ±0.30
3	Did the recommendations help you to find the coins?	42 Yes, 5 No
4	Did the recommendations help you to browse the scene?	49 Yes, 2 No
5	Do you think recommendations can be helpful?	49 Yes, 2 No
6	Which recommendation style do you prefer and why?	32 Arrows, 7 Viewports
7	Did you enjoy this?	36 Yes, 3 No

Table 1: List of questions in the questionnaire and summary of answers. Questions 1 and 2 have a 99%confidence interval.

Participants

The participants were recruited on microworkers.com, a crowdsourcing website. There were 51 participants (36 men and 15 women), who are in average 30.44 years old.

3.2.4 Experimental results

We now present the results from our user study, focusing on whether bookmarks help users navigating the 3D scene.

Questionnaire

We had 51 responses to the questionnaire. The answers are summarized in Table ref{bi:questions}. Note that not all questions were answered by all participants.

The participants seem to find the task to be of average difficulty (3.04/5) when they have no bookmarks to help their navigation. They judge the task to be easier in average (2.15/5) with bookmarks, which indicates that bookmarks ease the completion of the task.

Almost all users (49 out of 51) think the bookmarks are useful for browsing the scene, and most users (42 out of 51) think bookmarks are also useful to complete the given task. This is slightly in contradiction with our setup; even if coins may appear in some bookmarked viewpoints (which is normal since the viewpoints have been chosen to get the most complete coverage of the scene), most of the time no coin is visible in a given bookmark, and there are always coins that are invisible from all bookmarks.

The strongest result is that almost all users (49 out of 51) find bookmarks to be helpful. In addition, users seem to have a preference for **Arrows** against **Viewports** (32 against 7).

BM type	#Exp	Mean # coins	# completed	Mean time
NoReco	51	7.08	18	4 min 16 s
Arrows	51	7.39	27	2 min 33 s
Viewports	51	7.51	30	2 min 16 s

Analysis of interactions

Table 2: Analysis of the sessions length and users success by type of bookmarks

Table 2 shows basic statistics on task completion given the type of bookmarks that were provided to the participants.

First, we can see that without bookmarks, only a little bit more than a third of the users are able to complete the task, i.e. find all 8 coins. In average, these users find just above 7 coins, and spend 4 minutes and 16 seconds to do it.

Interestingly, and regardless of the bookmark type, users who have bookmarks complete the task more than half of the time, and spend in average significantly less time to complete the task: 2 minutes and 16 seconds using **Viewports** and 2 minutes and 33 seconds using **Arrows**. Although **Viewports** seem to help users a little bit more in completing the task than **Arrows**, the performance difference between both types of bookmarks is not significant enough to conclude on which type of bookmarks is best.

The difference between an interface with bookmarks and without bookmarks, however, is very clear. Users tend to complete the task more efficiently using bookmarks: more users actually finish the task, and it takes them half the time to do so. We computed 99% confidence intervals on the results introduced in Table ref{bi:sessions}. We found that the difference in mean number of coins collected with and without bookmarks is not high enough to be statistically significant: we would need more experiments to reach the significance. The mean time spent on the task however is statistically significant.

BM type	Total distance	Distance to a bookmark	Ratio
NoReco	610.80	0	0%
Arrows	586.30	369.77	63%
Viewports	546.96	332.72	61%

Table 3: Analysis of the length of the paths by type of bookmarks

Table 3 presents the length of the paths traveled by users in the scenes. Although users tend to spend less time on the tasks when they do not have bookmarks, they travel pretty much the same distance as without bookmarks. As a consequence, they visit the scene faster in average with bookmarks, than without bookmarks. The table shows that this higher speed is due to the bookmarks, as more than 60% of the distance traveled by users with bookmarks happens when users click on bookmarks and fly to the destination.

Discussion

In the previous paragraphs, we have shown how bookmarks are well perceived by users (looking at the questionnaire answers). We also showed that users tend to be more efficient in completing the task when they have bookmarks than when they do not.

We can say that bookmarks have a positive impact on navigation within the 3D scene, but since users move, on average, twice as fast, it might have a negative impact on the streaming of objects to the client.

Figure X shows a CDF of the percentage of 3D mesh triangles in the scene that have been queried by users after a certain time. We plotted this same curve for users with and without bookmarks. As expected, the fact that the users can browse the scene significantly quicker with bookmarks reflects on the demand on the 3D content. Users need more triangles more quickly, which either leads to more demand on network bandwidth, or if the bandwidth is kept constant, leads to fewer objects being displayed. In the next section, we introduce experiments based on our user study traces that show how the rendering is affected by the presence of bookmarks and how to improve it.

3.3 Impact of 3D bookmarks on streaming

3.3.1 3D model streaming

In this section, we describe our implementation of a 3D model streaming policy in our simulation. A summary of the streaming policies we designed is given in Table X. Note that the policy is different from the one we used for the crowdsourcing experiments. Recall that in the crowdsourcing experiments, we load all the 3D content before the participants begin to navigate to remove bias due to different network conditions. Here, we implemented a streaming version, which we expect an actual NVE will use.

The 3D content we used are textured mesh — coded in obj file format. As such, the data we used in our experiments are made of several components. The geometry consists of (i) a list of vertices and (ii) a list of faces, and the texture consists of (i) a list of materials, (ii) a list of texture coordinates, and (iii) a set of texture images. In the crowdsourcing experiment, we keep the model small since the goal is to study the user interaction. To increase the size of the model, while keeping the same 3D scene, we subdivide each triangle three times, successively, thereby multiplying the total number of triangles in the scene by 64. We do this to simulate a reasonable use case with large 3D scenes. Table ref{bi:modelsize} shows that material and texture amount at most for 3.6% of the geometry, which justifies this choice.

When a client starts loading the web page containing the 3D model, the server first sends the list of materials and the texture files. Then, the server periodically sends a fixed size chunk that indifferently encapsulates vertices, texture coordinates, or faces. A *vertex* is coded with three floats and an integer (x, y, and z coordinates and the index of the vertex), a *texture coordinate* with two floats and an integer (the x and y coordinates on the image and the index of the texture coordinate, the index of the face with eight integers (the index of each vertex, the index of each texture coordinate, the index of the face and the number of the corresponding material). Consequently, given the JavaScript implementation of integers and floats, we approximate each vertex and each texture coordinate to take up 32 bytes, and each face takes up 96 bytes.

Scene	Material	Images Geometr		
Scene 1	8 KB	72 KB	8.48 MB	
Scene 2	302 KB	8 KB	8.54 MB	
Scene 3	16 KB	92 KB	5.85 MB	

 Table 4: Respective sizes of materials, textures (images) and geometries for the three scenes used in

 the user study

During playback, the client periodically (every 200 ms in our implementation) sends to the server its current position and camera orientation. The server computes a sorted list of relevant faces: first the server performs frustum **culling** to compute the list of faces that intersect with the client's viewing frustum. Then, it performs backface **culling** to discard the faces whose normals point towards the same direction as the client's camera orientation. The server then sorts the filtered faces according to their distance to the camera. Finally, the server incrementally fills in chunks with these ordered faces. If a face depends on a vertex or a texture coordinate that has not yet been sent, the vertex or the texture coordinate is added to the chunk as well. When the chunk is full, the server sends it. Both client and server algorithms are detailed in algorithms ref{bi:streaming-algorithm-client} and ref{bi:streaming-algorithm-server}. The chunk size is set according to the bandwidth limit of the server. Note that the server may send faces that are occluded and not **visible** to the client, since determining visibility requires additional computation.

In the following, we shall denote this streaming policy **culling**; in Figures ref{bi:click-1250} and ref{bi:click-625} streaming using **culling** only is denoted **C-only**.

3.3.2 3D bookmarks

We have seen (Figure ref{bi:triangles-curve}) that navigation with bookmarks is more demanding on the bandwidth. We want to exploit bookmarks to improve the user's quality of experience. For this purpose, we propose two streaming policies based on offline computation of the relevance of 3D content to bookmarked viewpoints.

Visibility determination for 3D bookmarks

A bookmarked viewpoint is more likely to be accessed, compared to other arbitrary viewpoint in the 3D scene. We exploit this fact to perform some precomputation on the 3D content **visible** from the bookmarked viewpoint.

Recall that **culling** does not consider occlusion of the faces. Furthermore, it prioritizes the faces according to distance from the camera, and does not consider the actual contribution of the faces to the rendered 2D images. Ideally, we should prioritize the faces that occupy a bigger area in the 2D rendered images. Computing this, however, requires rendering the scene at the server, and measuring the area of each face. It is not scalable to compute this for every viewpoint requested by the client.

However, we can prerender the bookmarked viewpoints, since the number of bookmarks is limited, their viewpoints are known in advance, and they are likely to be accessed. For each bookmark, we render the scene offline, using a single color per triangle. Once rendered, we scan the output image to find the **visible** triangles (based on the color) and sort them by decreasing projected area. This technique is also used by citep{view-dependent-progressive-mesh}. Thus, when the user clicks on a 3D bookmark, this precomputed list of faces is used by the server, and only **visible** faces are sent in decreasing order of contributions to the rendered image.

For the three scenes that we used in the experiment, we can reduce the number of triangles sent by 60% (over all bookmarks). This reduction is as high as 85.7% for one particular bookmark (from 26,886 culled triangles to 3,853 culled and **visible** triangles).

To illustrate the impact of sorting by projected area of faces, Figure ref{bi:sorted-tri} shows the quality improvement gained by sending the precomputed **visible** triangles prioritized by projected areas, compared to using **culling** only prioritized by distance. The curve shows the average quality over all bookmarks over all scenes, for a given number of triangles received. The quality is measured by the ratio of correctly rendered pixels, comparing the fully and correctly rendered image (when all 3D content is available) and the rendered image (when content is partially available). We sample one pixel every 100 rows and every 100 columns to compute this value. The figure shows that, to obtain 90% of correctly displayed samples, we require 1904 triangles instead of 5752 triangles, about 1/3 savings.

In what follows, we will refer to this streaming policy as visible.

Prefetching by predicting the next bookmark clicked

We can now use the precomputed, visibility-based streaming of 3D content for the bookmarks to reduce the amount of traffic needed. Next, we propose to prefetch the 3D content from the bookmarks. Any efficient prefetching policy needs to accurately predict users' actions.

As shown, users tend to visit the bookmarked viewpoints more often than others, except the initial viewpoint. It is thus natural to try to prefetch the 3D content of the bookmarks.

Figure ref{bi:mat1} shows the probability of visiting a bookmark (vertical axis) given that another bookmark has been visited (horizontal axis). This figure shows that users tend to follow similar paths when consuming bookmarks. Thus, we hypothesize that prefetching along those paths would lead to better image quality and lower discovery latency.

The policy used is the following. We divide each chunk sent by the server into two parts. The first part is used to fetch the content from the current viewpoint, using the **culling** streaming policy. The second part is used to prefetch content from the bookmarks, according to their likelihood of being clicked next. We use the probabilities displayed in Figure ref{bi:mat1} to determine the size of each part. Each bookmark *B* has a probability $p(B|B_{prev})$ of being clicked next, considering that B_{prev} was the last clicked bookmark. We assign to each bookmark a certain portion of the chunk to prefetch the corresponding data proportionally to the probability of it being clicked. We use the **visible** policy to determine which data should be sent for a bookmark.

We denote this combination as V-PP, for Prefetching based on Prediction using visible policy.

Fetching destination bookmark

An alternate method to benefit from the precomputing **visible** triangles at the bookmark, is to fetch 3D content during the "fly-to" transition to reach the destination. Indeed, as specified in Section ref{bi:3dnavigation}, moving to a bookmarked viewpoint is not instantaneous, but rather takes a small amount of time to smoothly move the user camera from its initial position towards the bookmark. This transition usually takes from 1 to 2 seconds, depending on how far the current user camera position is from the bookmark.

When the user clicks on the bookmark, the client fetches the **visible** vertices from the destination viewpoint, with all the available bandwidth. So, during the transition time, the server no longer does **culling**, but the whole chunk is used for fetching following **visible** policy.

The immediate drawback of this policy is that on the way to the bookmark, the user perception of the scene will be degraded because of the lack of data for the viewpoints in transition. On the bright side, no time is lost to prefetch bookmarks that will never be consumed, because we fetch only when we are sure that the user has clicked on a bookmark. This way, when the user is not clicking on bookmarks, we can use the entire bandwidth for the current viewpoint and get as many triangles as possible to improve the current viewpoint. We call this method **V-FD**, since we are Fetching the 3D data from the Destination using **visible** policy.

3.3.3 Comparing streaming policies

In order to determine which policy to use, we replay the traces from the user study while simulating different streaming policies. The first point we are interested in is which streaming policy leads to the lower discovery latency and better image quality for the user: **culling** (no prefetching), **V-PP** (prefetching based on probability of accessing bookmarks), or **V-FD** (no prefetching, but fetch the destination during "fly-to" transition) or combining both **V-PP** and **V-FD** (**V-PP+FD**).

Figure ref{bi:click-1250} compares the quality of the view of a user after their first click on a bookmark. The ratio of pixels correctly displayed is computed in the client algorithm, see also algorithm ref{bi:streaming-algorithm-client}. In this figure we use a bandwidth of 1 Mbps. The blue curve corresponds to the **culling** policy. Clicking on a bookmark generates a user path with less spatial locality, causing a large drop in visual quality that is only compensated after 4 seconds. During the first second, the camera moves from the current viewport to the bookmarked viewport.

When the data has been prefetched according to the probability of the bookmark to be clicked, the drop in quality is less **visible** (**V-PP** curve). However, by benefiting from the precomputation of **visible** triangles and ordering of the important triangles in a bookmark (**V-FD**) the drop in quality is still there, but is very short (approximately four times shorter than for **culling**). This drop in quality is happening during the transition on the path. More quantitatively, with a 1 Mbps bandwidth, 3 seconds are necessary after the click to recover 90% of correct pixels.

Figure ref{bi:click-625} showed the results of the same experiment with 0.5 Mbps bandwidth. Here, it takes 4 to 5 seconds to recover 85% of the pixels with **culling** and **V-PP**, against 1.5 second for recovering 90% with **V-FD**. Combining both strategies (**V-PP+FD**) leads to the best quality.

At 1 Mbps bandwidth, **V-PP** penalizes the quality, as the curve **V-PP+FD** leads to a lower quality image than **V-FD** alone. This effect is even stronger when the bandwidth is set to 2 Mbps (Figure re-f{bi:2MB}). Both streaming strategies based on the precomputation of the ordering improves the image quality. We see here, that **V-FD** has a greater impact than **V-PP**. Here, **V-PP** may prefetch content that eventually may not be used, whereas **V-FD** only sends relevant 3D content (knowing which bookmark has been just clicked).

We present only the results after the first click. For subsequent clicks, we found that other factors came into play and thus, it is hard to analyze the impact of the various streaming policies. For instance, a user may revisit a previously visited bookmark, or the bookmarks may overlap. If the users click on a subsequent bookmark after a long period, then more content would have been fetched for this user, making comparisons difficult.

To summarize, we found that exploiting the fact that bookmarked viewpoints are frequently visited to precompute the **visible** faces and sort them according to projected areas can lead to significant improvement in image quality after a user interaction (clicking on a bookmark). This alone can lead to 60% less triangles being sent, with 1/3 of the triangles sufficient to ensure 90% of pixels correctly rendered, compared to doing frustum/backface **culling**. If we fetch these precomputed faces of the

destination viewpoint this way immediately after the click, during the "fly-to" transition, then we can already significantly improve the quality without any prefetching. Prefetching helps if the bandwidth is low, and fewer triangles can be downloaded during this transition. The network conditions play a minimum role in this key message — bookmarking allows precomputation of an ordered list of **visible** faces, and this holds regardless of the underlying network condition (except for non-interesting extreme cases, such as negligible bandwidth or abundance of bandwidth).

Chapter 4

DASH-3D



Figure 20: A subdivided 3D scene with a viewport and regions delimited with red edges. In white, the regions that are outside the field of view of the camera; in green, the regions inside the field of view of the camera.

Dynamic Adaptive Streaming over HTTP (DASH) is now a widely deployed standard for video streaming, and even though video streaming and 3D streaming are different problems, many of DASH features can inspire us for 3D streaming. In this chapter, we present the most important contribution of this thesis: adapting DASH to 3D streaming.

First, we show how to prepare 3D data into a format that complies with DASH data organization, and we store enough metadata to enable a client to perform efficient streaming. The data preparation consists in partitioning the scene into spatially coherent cells and segmenting each cell into chunks with a fixed number of faces, which are sorted by area so that faces of a different level of detail are not grouped together. We also export each texture at different resolutions. We encode the metadata that describes the data organization into a 3D version of the Media Presentation Description (MPD) that DASH uses for video. All this prepared content is then stored on a simple static HTTP server: a clients can request the content without any need for computation on the server side, allowing a server to support an arbitrary number of clients.

We then propose DASH-3D clients that are viewpoint aware: they perform frustum culling to eliminate cells outside the viewing volume of the camera (as shown in Figure 20). We define utility metrics to give a score to each chunk of data, be it geometry or texture, based on offline information that is given in the MPD, and online information that the client is able to compute, such as view parameters, user interaction or bandwidth measurements. We also define streaming policies that rely on those utilities in order for the client to determine which chunks need to be downloaded. We finally evaluate these system parameters under different bandwidth setups and compare our streaming policies.

4.1 Introduction

In this chapter, we take a little step back from interaction and propose a system with simple interactions that however, addresses most of the open problems mentioned in Section 1. We take inspiration from video streaming: working on the similarities between video streaming and 3D streaming (seen in Section 1.2), we benefit from the DASH efficiency (seen in Section 2.1.1) for streaming 3D content. DASH is based on content preparation and structuring which helps not only the streaming policies but also leads to a scalable and efficient system since it moves completely the load from the server to the clients. A DASH client downloads the structure of the content, and then, depending on its needs and independently of the server, decides what to download.

In this chapter, we show how to mimic DASH video with 3D streaming, and we develop a system that keeps DASH benefits. Section 4.2 describes our content preparation and metadata, and all the preprocessing that is done to our model to allow efficient streaming. Section ref{d3:dash-client} gives possible implementations of clients that exploit the content structure. Section ref{d3:evaluation} evaluates the impact of the different parameters that appear both in the content preparation and the client. Finally, Section ref{d3:conclusion} sums up our work and explains how it tackles the challenges raised in the conclusion of the previous chapter.

4.2 Content preparation

In this section, we describe how we preprocess and store the 3D data of the NVE, consisting of a polygon soup, textures, and material information into a DASH-compliant Media Presentation Description (MPD) file. In our work, we use the obj file format for the polygons, png for textures, and mtl format for material information. The process, however, applies to other formats as well.

4.2.1 The MPD File

In DASH, the information about content storage and characteristics, such as location, resolution, or size, is extracted from an MPD file by the client. The client relies only on this information to decide which chunk to request and at which quality level. The MPD file is an XML file that is organized into different sections hierarchically. The period element is a top-level element, which for the case of video, indicates the start time and length of a video chapter. This element does not apply to NVE, and we use a single period for the whole scene, as the scene is static. Each period element contains one or more adaptation sets, which describe the alternate versions, formats, and types of media. We utilize adaptation sets to organize a 3D scene's material, geometry, and texture.

The piece of software that does the preprocessing of the model consists in file manipulation and is written in Rust. It successively preprocesses the geometry and then the textures. The MPD is generated by a library named xml-rs which works like a stack:

- a structure is created on the root of the MPD file;
- the start_element method creates a new child in the XML file;

• the end_element method ends the current child and pops the stack.

This structure is passed along with our geometry and texture preprocessors that can add elements to the XML file as they are generating the corresponding data chunks.

4.2.2 Adaptation sets

When the user navigates freely within an NVE, the frustum at given time almost always contains a limited part of the 3D scene. Similar to how DASH for video streaming partitions a video clip into temporal chunks, we segment the polygons into spatial chunks, such that the DASH client can request only the relevant chunks.

Geometry management

We use a space partitioning tree to organize the faces into cells. A face belongs to a cell if its barycenter falls inside the corresponding bounding box. Each cell corresponds to an adaptation set. Thus, geometry information is spread on adaptation sets based on spatial coherence, allowing the client to download the relevant faces selectively. A cell is relevant if it intersects the frustum of the client's current viewpoint. Figure 20 shows the relevant cells in green. As our 3D content, a virtual environment, is biased to spread along the horizontal plane, we split the bounding box alternatively along the two horizontal directions.

We create a separate adaptation set for large faces (e.g., the sky or ground) because they are essential to the 3D model and do not fit into cells. We consider a face to be large if its area in 3D is more than $a + 3\sigma$, where a and σ are the average and the standard deviation of 3D area of faces respectively. In our example, it selects the 5 largest faces that represent 15% of the total face area. We thus obtain a decomposition of the NVE into adaptation sets that partitions the geometry of the scene into an adaptation that contains the larger faces of the model, and smaller adaptation sets containing the remaining faces.

We store the spatial location of each adaptation set, characterized by the coordinates of its bounding box, in the MPD file as the supplementary property of the adaptation set in the form of " x_{\min} , width, y_{\min} , height, z_{\min} , depth" (as shown in Listing 11). This information is used by the client to implement a view-dependent streaming (Section ref{d3:dash-client}).

Texture management

As with geometry data, we handle textures using adaptation sets but separate from geometry. Each texture file is contained in a different adaptation set, with multiple representations providing different image resolutions (see Section 4.2.3). We add an attribute to each adaptation set that contains texture, describing the average color of the texture. The client can use this attribute to render a face for which the corresponding texture has not been loaded yet, so that most objects appear, at least, with a uniform natural color (see Figure 21).

Material management

The material (MTL) file is a text file that describes all materials used in the OBJ files for the entire 3D model. A material has a name, properties such as specular parameters, and, most importantly, a path to a texture file. The MTL file maps each face of the OBJ to a material. As the MTL file is a different type of media than geometry and texture, we define a particular adaptation set for this file, with a single representation.

4.2.3 Representations

Each adaptation set can contain one or more representations of the geometry or texture data, at different levels of detail (e.g., a different number of faces). For geometry, the resolution (i.e., 3D areas of faces) is heterogeneous, thus applying a sensible multi-resolution representation is cumbersome: the 3D area of faces varies from 0.01 to more than 10K, disregarding the outliers. For textured scenes, it is common to have such heterogeneous geometry size since information can be stored either in geometry or texture. Thus, handling the streaming compromise between geometry and texture is more adaptive than handling separately multi-resolution geometry. Moreover, as our faces are partitioned into independent cells, multi-resolution would cause difficult stitching issues such as topological gaps between the cells.

For an adaptation set containing texture, each representation contains a single segment where the image file is stored at the chosen resolution. In our example, from the full-size image, we generate successive resolutions by dividing both height and width by 2, stopping when the image size is less or equal to 64×64 . Figure 21 illustrates the use of the textures against the rendering using a single, average color per face.



Figure 22: With full resolution textures



Figure 23: With average colors

Figure 21: Rendering of the model with different styles of textures

4.2.4 Segments

To allow random access to the content within an adaptation set storing geometry data, we group the faces into segments. Each segment is then stored as an OBJ file which can be individually requested by the client. For geometry, we partition the faces in an adaptation set into sets of N_s faces, by first

sorting the faces by their area in 3D space in descending order, and then place each successive N_s faces into a segment. Thus, the first segment contains the biggest faces and the last one the smallest. In addition to the selected faces, a segment stores all face vertices and attributes so that each segment is independent. For textures, each representation contains a single segment.

```
1
     <AdaptationSet>
 2
         <SupplementalProperty value="-8834.11230,2201.58853,</pre>
 3
              -0.16950, 174.81540, -1344.47740,4767.83367" />
 4
      <BaseURL>as1/</BaseURL>
 5
      <Representation>
 6
       <BaseURL>repr1/</BaseURL>
7
       <SegmentList>
8
        <SegmentURL area="2540342.3" size="120K" media="s0.obj" />
        <SegmentURL area="1124.4" size="162K" media="sl.obj" />
9
10
        <SegmentURL area="412.6" size="173K" media="s2.obj" />
        <SegmentURL area="270.3" size="147K" media="s3.obj" />
11
       </SegmentList>
12
      </Representation>
13
14
     </AdaptationSet>
15
     <AdaptationSet area="198632.73912" average="178,176,173" mimeType="image/png">
16
         <BaseURL>textures/MFL00R07.PNG/</BaseURL>
17
18
         <Representation>
19
             <BaseURL>64x64/</BaseURL>
20
             <SegmentList>
21
                 <SegmentURL size="7K" mse="57.6" media="t.png" />
22
             </SegmentList>
         </Representation>
23
         <Representation>
24
25
             <BaseURL>128x128/</BaseURL>
26
             <SegmentList>
27
                 <SegmentURL size="27K" mse="0.0" media="t.png" />
28
             </SegmentList>
29
         </Representation>
30
     </AdaptationSet>
```

Listing 11: MPD description of a geometry adaptation set, and a texture adaptation set.

Now that the 3D data is partitioned and that the MPD file is generated, we see in the next section how the client uses the MPD to request the appropriate data chunks.

4.3 Client

In this section, we specify a DASH NVE client which exploits the preparation of the 3D content in an NVE for streaming.

The generated MPD file describes the content organization so that the client gets all the necessary information to make educated decisions and query the 3D content it needs according to the available resources and current viewpoint. A camera path generated by a particular user is a set of viewpoint $v(t_i)$ indexed by a continuous time interval $t_i \in [t_1, t_{end}]$.

All DASH clients are built from the same basic bricks, as shown in Figure 24:

- the access client, which is the module that deals with making HTTP requests and receiving responses;
- the *segment parsers*, which decode the data downloaded by the access client, whether it be materials, geometry or textures;
- the control engine, which analyses the bandwidth to dynamically adapt to it;
- the media engine, which renders the multimedia content and the user interface to the screen.

TODO Figure 24: DASH client-server architecture

The DASH client first downloads the MPD file to get the material file containing information about all the geometry and textures available for the entire 3D model. At time instance t_i , the DASH client decides to download the appropriate segments containing the geometry and the texture to generate the viewpoint $v(t_{i+1})$ for the time instance t_{i+1} .

Starting from t_1 , the camera continuously follows a camera path $C = \{v(t_i), t_i \in [t_1, t_{end}]\}$, along which downloading opportunities are strategically exploited to sequentially query the most useful segments.

4.3.1 Segment utility

Unlike video streaming, where the bitrate of each segment correlates with the quality of the video received, for 3D content, the size (in bytes) of the content does not necessarily correlate well to its contribution to visual quality. A large polygon with huge visual impact takes the same number of bytes as a tiny polygon. Further, the visual impact is *view dependent* — a large object that is far away or out of view does not contribute to the visual quality as much as a smaller object that is closer to the user. As such, it is important for a DASH-based NVE client to estimate the usefulness of a given segment to download, so that it can make good decisions about what to download. We call this usefulness the *utility* of the segment.

The utility is a function of a segment, either geometry or texture, and the current viewpoint (camera location, view angle, and look-at point), and is therefore dynamically computed online by the client from parameters in the MPD file.

4.3.2 Offline parameters

Let us detail first, all parameters available from the offline/static preparation of the 3D NVE. These parameters are stored in the MPD file. First, for each geometry segment s^G there is a predetermined 3D area $\mathcal{A}(s^G)$, equal to the sum of all triangle areas in this segment (in 3D); it is computed as the

segments are created. Note that the texture segments have similar information, but computed at *navigation time* t_i . The second information stored in the MPD for all segments, geometry, and texture, is the size of the segment (in kB).

Finally, for each texture segment s^T , the MPD stores the *MSE* (mean square error) of the image and resolution, relative to the highest resolution (by default, triangles are filled with its average color). Offline parameters are stored in the MPD as shown in Snippet ref{d3:mpd}.

4.3.3 Online parameters

In addition to the offline parameters stored in the MPD file for each segment, view-dependent parameters are computed at navigation time. First, a measure of 3D area is computed for texture segments. As a texture maps on a set of triangles, we account for the area in 3D of all these triangles. We could consider such an offline measure (attached to the adaptation set containing the texture), but we prefer to only account for the triangles that have been already downloaded by the client. We call the set of triangles colored by a texture $T: \Delta(s^T) = \Delta(T)$ (depending only on T and equal for any representation/segment s^T in this texture adaptation set). At each time t_i , a subset of $\Delta(T)$ has been downloaded; we denote it $\Delta(T, t_i)$.

Moreover, each geometry segment belongs to a geometry adaptation set AS^G whose bounding box coordinates are stored in the MPD. Given the coordinates of the bounding box $\mathcal{BB}(AS^G)$ and the viewpoint $v(t_i)$ at time t_i , the client computes the distance $\mathcal{D}(v(t_i), AS^G)$ of the bounding box $\mathcal{BB}(AS^G)$ as the distance from the center of $\mathcal{BB}(AS^G)$ to the principal point of the camera, given in $v(t_i)$.

4.3.4 Utility for geometry segments

We now have all parameters to derive a utility measure of a geometry segment. Utility for texture segments follows from the geometric utility.

The utility of a geometric segment s^G for a viewpoint $v(t_i)$ is: where AS^G is the adaptation set containing s^G .

Basically, the utility of a segment is proportional to the area that its faces cover, and inversely proportional to the square of the distance between the camera and the center of the bounding box of the adaptation set containing the segment. That way, we favor segments with big faces that are close to the camera.

4.3.5 Utility for texture segments

For a texture T stored in a segment s^T , the triangles in $\Delta(T)$ are stored in arbitrary geometry segments, that is, they do not have spatial coherence. Thus, for each k^{th} downloaded geometry segment s_k^G , and total downloaded segment K at time t_i , we collect the triangles of $\Delta(T, t_i)$ in s_k^G , and compute the ratio of $\mathcal{A}(s_k^G)$ covered by these triangles. So, we define the utility:

where we sum over all geometry segments received before time t_i that intersect $\Delta(T, t_i)$ and such that the adaptation set it belongs to is in the frustum. This formula defines the utility of a texture segment by computing the linear combination of the utility of the geometry segments that use this texture, weighted by the proportion of area covered by the texture in the segment. We compute the PSNR by using the MSE in the MPD and denote it $psnr(s^T)$. We do this to acknowledge the fact that a texture at a greater resolution has a higher utility than a lower resolution texture. The equivalent term for geometry is 1 (and does not appear). Having defined a utility on both geometry and texture segments, the client uses it next for its streaming strategy.

4.3.6 DASH adaptation logic

Along the camera path $C = \{v(t_i)\}$, viewpoints are indexed by a continuous time interval $t_i in[t_1, t_{end}]$. Contrastingly, the DASH adaptation logic proceeds sequentially along a discrete time line. The first HTTP request made by the DASH client at time t_1 selects the most useful segment s_1^a to download and will be followed by subsequent decisions at t_2, t_3, \ldots . While selecting s_i^a , the i^{th} best segment to request, the adaptation logic compromises between geometry, texture, and the available representations given the current bandwidth, camera dynamics, and the previously described utility scores. The difference between t_{i+1} and t_i is the s_i^a delivery delay. It varies with the segment size and network conditions. Algorithm ref{d3:next-segment} details how our DASH client makes decisions.

A naive way to sequentially optimize the utility \mathcal{U} is to limit the decision-making to the current viewpoint $v(t_i)$. In that case, the best segment s to request would be the one maximizing $\mathcal{U}(s, v(t_i))$ to simply make a better rendering from the current viewpoint $v(t_i)$. Due to transmission delay however, this segment will be only delivered at time $t_{i+1} = t_{i+1}(s)$ depending on the segment size and network conditions:

In consequence, the most useful segment from $v(t_i)$ at decision time t_i might be less useful at delivery time from $v(t_{i+1})$.

A better solution is to download a segment that is expected to be the most useful in the future. With a temporal horizon χ , we can optimize the cumulated \mathcal{U} over $[t_{i+1}(s), t_i + \chi]$:

In our experiments, we typically use $\chi = 2s$ and estimate the (ref{d3:smart}) integral by a Riemann sum where the $[t_{i+1}(s), t_i + \chi]$ interval is divided in 4 subintervals of equal size. For each subinterval extremity, an order 1 predictor hat $\{v\}(t_i)$ linearly estimates the viewpoint based on $v(t_i)$ and speed estimation (discrete derivative at t_i).

We also tested an alternative greedy heuristic selecting the segment that optimizes an utility variation during downloading (between t_i and t_{i+1}):

4.3.7 JavaScript client

In order to be able to evaluate our system, we need to collect traces and perform analyses on them. Since our scene is large, and since the system we are describing allows navigating in a streaming scene, we developed a JavaScript web client that implements our utility metrics and policies.

Media engine

Performance of our system is a key aspect in our work; as such, we can not use the default geometries described in Section ref{f:geometries} because of their poor performance, and we instead use buffer geometries. However, in our system, the way changes happen to the 3D content is always the same: we only add faces and textures to the model. We therefore implemented a class that derives BufferGeometry, for more convenience.

- It has a constructor that takes as parameter the number of faces: it allocates all the memory needed for our buffers so we do not have to reallocate it later (which would be inefficient).
- It keeps track of the number of faces it is currently holding: it can then avoid rendering faces that have not been filled and knows where to add new faces.
- It provides a method to add a new polygon to the geometry.
- It also keeps track of what part of the buffers has been transmitted to the GPU: THREE.js allows us to set the range of the buffer that we want to update, and we are able to update only what is necessary.

Our 3D model class

As said in the previous subsections, a geometry and a material are bound together in a mesh. This means that we are forced to have as many meshes as there are materials in our model. To make this easy to manage, we implemented a **Model** class, that holds both geometry and textures. We can add vertices, faces, and materials to this model, and it internally manages the right geometries, materials and meshes. In order to avoid having many models that share the same material (which would harm performance), it automatically merges faces that share the same material in the same buffer geometry, as shown in Figure 25.

TODO

Figure 25: Reordering of the content on the renderer

Access client

In order to be able to implement our view-dependent DASH-3D client, we need to implement the access client, which is responsible for deciding what to download and for downloading it. To do so, we use the strategy pattern illustrated in Figure 26. We maintain a base class named LoadingPolicy that contain some attributes and functions to keep track of what has been downloaded. This class exposes a function named nextSegment that takes two arguments:

- the MPD, so that a strategy can know all the metadata of the segments before making its decision;
- the camera, because the next best segment depends on the camera position.

The greedy and proposed policies from the previous chapter are all classes that derive from LoadingPolicy. Then, the main class responsible for the loading of segments is the DashLoader class. It uses XMLHttpRequests, which are the usual way of making HTTP requests in JavaScript, and it calls the corresponding parser on the results of those requests. The DashLoader class accepts as parameter a function that will be called each time some data has been downloaded and parsed: this data can contain

vertices, texture coordinates, normals, materials or textures, and they can all be added to the Model class that we described in Chapter .

TODO Figure 26: Class diagram of our DASH client

Performance

JavaScript requires the use of *web workers* to perform parallel computing. A web worker is a script in JavaScript that runs in the background, on a separate thread and that can communicate with the main script by sending and receiving messages. Since our system has many tasks to perform, it is natural to use workers to manage the streaming without impacting the framerate of the renderer. However, what a worker can do is very limited, since it cannot access the variables of the main script. Because of this, we are forced to run the renderer on the main script, where it can access the HTML page, and we move all the other tasks (i.e.

the access client, the control engine and the segment parsers) to the worker. Since the main script is the only thread communicating with the GPU, it will still have to update the model with the parsed content it receives from the worker. We do not use web workers to improve the framerate of the system, but rather to reduce the latency that occurs when receiving a new segment, which can be frustrating in a single thread scenario, since each time a segment is received, the interface would freeze for around half a second. A sequence diagram of what happens when downloading, parsing and rendering content is shown in Figure 27.

TODO

Figure 27: Repartition of the tasks on the main script and the worker

4.3.8 Rust client

However, a web client is not sufficient to analyse our streaming policies: many tasks are performed (such as rendering, and managing the interaction) and all this overhead pollutes the analysis of our policies. This is why we also implemented a client in Rust, for simulation, so we can gather precise simulated data.

Our requirements are quite different that the ones we had to deal with in our JavaScript implementation. In this setup, we want to build a system that is the closest to our theoretical concepts. Therefore, we do not have a full client in Rust (meaning an application to which you would give the URL to an MPD file and that would allow you to navigate in the scene while it is being downloaded). In order to be able to run simulations, we develop the bricks of the DASH client separately: the access client and the media engine are totally isolated:

• the **simulator** takes a user trace as a parameter, it then replays the trace using specific parameters of the access client and outputs a file containing the history of the simulation (which files have been downloaded, and when);

• the **renderer** takes the user trace as well as the history generated by the simulator as parameters, and renders images that correspond to what would have been seen.

When simulating experiments, we run the simulator on many traces that we collected during userstudies, and we then run the renderer program according to the traces to generate images corresponding to the simulation. We are then able to compute PSNR between those frames and the ground truth frames. Doing so guarantees us that our simulator is not affected by the performances of our renderer.

4.4 Evaluation

We now describe our setup and the data we use in our experiments. We present an evaluation of our system and a comparison of the impact of the design choices we introduced in the previous sections.

4.4.1 Experimental setup

Model

We use a city model of the Marina Bay area in Singapore in our experiments. The model came in 3DS Max format and has been converted into Wavefront OBJ format before the processing described in Chapter 4.2. The converted model has 387,551 vertices and 552,118 faces. Table 5 gives some general information about the model and Figure 28 illustrates the heterogeneity of our model (wireframe rendering is used to illustrate the heterogeneity of the geometry complexity). We partition the geometry into a k-d tree until the leafs have less than 10000 faces, which gives us 64 adaptation sets, plus one containing the large faces.



Figure 29: Low resolution geometry





Figure 31: High resolution geometry



Figure 32: Detailed textures

Figure 28: Illustration of the heterogeneity of the model

Files	Size
3DS Max	55 MB
OBJ file	62 MB
MTL file	0.27MB
Textures (high res)	167 MB
Textures (low res)	11 MB

Table 5: Sizes of the different files of the model

User navigations

To evaluate our system, we collected realistic user navigation traces which we can replay in our experiments. We presented six users with a web interface, on which the model was loaded progressively as the user could interact with it. The available interactions were inspired by traditional first-person interactions in video games, i.e., W, A, S, and D keys to translate the camera, and mouse to rotate the camera. We asked users to browse and explore the scene until they felt they had visited all important regions. We then asked them to produce camera navigation paths that would best present the 3D scene to a user that would discover it. To record a path, the users first place their camera to their preferred starting point, then click on a button to start recording. Every 100ms, the position, viewing angle of the camera and look-at point are saved into an array which will then be exported into JSON format. The recorded camera trace allows us to replay each camera path to perform our simulations and evaluate our system. We collected 13 camera paths this way.

Network setup

We tested our implementation under three network bandwidth of 2.5 Mbps, 5 Mbps, and 10 Mbps with an RTT of 38 ms, following the settings from DASH-IF citep{dash-network-profiles}. The values are kept constant during the entire client session to analyze the difference in magnitude of performance by increasing the bandwidth.

In our experiments, we set up a virtual camera that moves along a navigation path, and our access engine downloads segments in real time according to Algorithm ref{d3:next-segment}. We log in a JSON file the time when a segment is requested and when it is received. By doing so, we avoid wasting time and resources to evaluate our system while downloading segments and store all the information necessary to plot the figures introduced in the subsequent sections.

Hardware and software

The experiments were run on an Acer Aspire V3 with an Intel Core i7 3632QM processor and an NVIDIA GeForce GT 740M graphics card. The DASH client is written in Rust

Metrics

To objectively evaluate the quality of the resulting rendering, we use PSNR@. The scene as rendered offline using the same camera path with all the geometry and texture data available is used as ground truth. Note that a pixel error can occur in our case only in two situations: (i) when a face is missing, in which case the color of the background object is shown, and (ii) when a texture is either missing or downsampled. We do not have pixel error due to compression.

Experiments

We present experiments to validate our implementation choices at every step of our system. We replay the user-generated camera paths with various bandwidth conditions while varying key components of our system.

Table 6 sums up all the components we varied in our experiments. We compare the impact of two space-partitioning trees, a k-d tree and an octree, on content preparation. We also try several utility metrics for geometry segments: an offline one, which assigns to each geometry segment s^G the cumulated 3D area of its belonging faces $\mathcal{A}(s^G)$; an online one, which assigns to each geometry segment the inverse of its distance to the camera position; and finally our proposed method, as described in Section ref{d3:utility} $(\mathcal{A}_{\mathcal{D}}^{s^G} \{ (v\{(t_i)\}, AS^G) \}^2)$. We consider two streaming policies to be applied by the client, proposed in Section ref{d3:dash-client}. The greedy strategy determines, at each decision time, the segment that maximizes its predicted utility at arrival divided by its predicted delivery delay, which corresponds to equation (ref{d3:greedy}). The second streaming policy that we run is the one we proposed in equation (ref{d3:smart}). We have also analyzed the effect of grouping the faces

Parameters	Values
Content preparation	Octree, <i>k</i> -d tree
Utility	Offline, Online, Proposed
Streaming policy	Greedy, Proposed
Grouping of Segments	Sorted based on area, Unsorted
Bandwidth	2.5 Mbps, 5 Mbps, 10 Mbps

in geometry segments of an adaptation set based on their 3D area. Finally, we try several bandwidth parameters to study how our system can adapt to varying network conditions.

Table 6: Different parameters in our experiments

4.4.2 Experimental results

TODO

Figure 33: Impact of the space-partitioning tree on the rendering quality with a 5Mbps bandwidth

Figure 33 shows how the space partition can affect the rendering quality. We use our proposed utility metrics (see Section ref{d3:utility}) and streaming policy from equation (ref{d3:smart}), on content divided into adaptation sets obtained either using a k-d tree or an octree and run experiments on all camera paths at 5 Mbps. The octree partitions content into non-homogeneous adaptation sets; as a result, some adaptation sets may contain smaller segments, which contain both important (large) and non-important polygons. For the k-d tree, we create cells containing the same number of faces N_a (here, we take $N_a = 10000$). Figure ref{d3:preparation} shows that the system seems to be slightly less efficient with an octree than with a k-d tree based partition, but this result is not significant. For the remaining experiments, partitioning is based on a k-d tree.

TODO

Figure 34: Impact of the segment utility metric on the rendering quality with a 5Mbps bandwidth

Figure 34 displays how a utility metric should take advantage of both offline and online features. The experiments consider k-d tree cell for adaptation sets and the proposed streaming policy, on all camera paths. We observe that a purely offline utility metric leads to poor PSNR results. An online-only utility improves the results, as it takes the user viewing frustum into consideration, but still, the proposed utility (in Section ref{d3:utility}) performs better.

TODO

Figure 35: Impact of creating the segments of an adaptation set based on decreasing 3D area of faces with a 5Mbps bandwidth

Figure 35 shows the effect of grouping the segments in an adaptation set based on their area in 3D. The PSNR significantly improves when the 3D area of faces is considered for creating the segments. Since all segments are of the same size, sorting the faces by area before grouping them into segments leads to a skew distribution of how useful the segments are. This skewness means that the decision

that the client makes (to download those with the largest utility first) can make a bigger difference in the quality.

We also compared the greedy vs. proposed streaming policy (as shown in Figure 36) for limited bandwidth (5 Mbps). The proposed scheme outperforms the greedy during the first 30s and does a better job overall. Table 7 shows the average PSNR for the proposed method and the greedy method for different downloading bandwidth. In the first 30 sec, since there are relatively few 3D contents downloaded, making a better decision at what to download matters more: we observe during that time that the proposed method leads to 1 - 1.9 dB better in quality terms of PSNR compared to the greedy method.

Table 8 shows the distribution of texture resolutions that are downloaded by greedy and our proposed scheme, at different bandwidths. Resolution 5 is the highest and 1 is the lowest. The table shows a weakness of the greedy policy: the distributioon of downloaded textures does not adapt to the bandwidth. In contrast, our proposed streaming policy adapts to an increasing bandwidth by downloading higher resolution textures (13.9% at 10 Mbps, vs. 0.3% at 2.5 Mbps). In fact, an interesting feature of our proposed streaming policy is that it adapts the geometry-texture compromise to the bandwidth. The textures represent 57.3% of the total amount of downloaded bytes at 2.5 Mbps, and 70.2% at 10 Mbps. In other words, our system tends to favor geometry segments when the bandwidth is low, and favor texture segments when the bandwidth increases.

TODO Figure 36: Impact of the streaming policy (greedy vs. proposed) with a 5 Mbps bandwidth

BW (in Mbps)	2.5	5	10	2.5	5	10
Greedy	14.4	19.4	22.1	19.8	26.9	29.7
Proposed	16.3	20.4	23.2	23.8	28.2	31.1

Tab	le 7	: Average	PSNR,	Greedy	r vs. Proposed	l
-----	------	-----------	-------	--------	----------------	---

Resolutions	2.5 Mbps	5 Mbps	10 Mbps
1	5.7% vs 1.4%	6.3% vs 1.4%	6.17% vs 1.4%
2	10.9% vs 8.6%	13.3% vs 7.8%	14.0% vs 8.3%
3	15.3% vs 28.6%	20.1% vs 24.3%	20.9% vs 22.5%
4	14.6% vs 18.4%	14.4% vs 25.2%	14.2% vs 24.1%
5	11.4% vs 0.3%	11.1% vs 5.9%	11.5% vs 13.9%

 Table 8: Percentages of downloaded bytes for textures from each resolution, for the greedy streaming policy (left) and for our proposed scheme (right)

4.5 Conclusion

Our work in this chapter started with the question: can DASH be used for NVE@? The answer is *yes*. In answering this question, we contributed by showing how to organize a polygon soup and its textures into a DASH-compliant format that (i) includes a minimal amount of metadata that is useful for the client, (ii) organizes the data to allow the client to get the most useful content first. We further show that the data organization and its description with metadata (precomputed offline) is sufficient to design and build a DASH client that is adaptive — it selectively downloads segments within its view, makes intelligent decisions about what to download, balances between geometry and texture while adapting to network bandwidth. This way, our system addresses the open problems we mentioned in Chapter 1.

- It prepares and structures the content in a way that enables streaming: all this preparation is precomputed, and all the content is structured according to DASH framework, geometry but also materials and textures. Furthermore, textures are prepared in a multi-resolution manner, and even though multi-resolution geometry is not discussed here, the difficulty of integrating it in this system seem moderated: we could encode levels of detail in different representations and define a utility metric for each representation and the system should adapt naturally.
- We are able to estimate the utility of each segment by exploiting all the metadata given in the MPD and by analysing the camera parameters of the user.
- We proposed a few streaming policies, from the easiest to implement to the more complex, so that the client exploits the utility metrics to define a best guess for the next chunk to download.
- The implementation is efficient: the content preparation allows a client to get all the information it needs from metadata and the server has nothing else to do than to serve files. Special attention has been granted to the client's performance.

However, the work described in this chapter does not take any quality of experience metrics into account. We designed a 3D streaming system, but we kept the interaction system the simplest possible. Dealing with interaction while dealing with all of the other problems we try to solve seems hard, and we believe keeping the interaction simple was a necessary step to build a solid 3D streaming system. Now that we have this system, we are able to work again on the interaction problem and our work and conclusions are given in Chapter 5.
Chapter 5

Bookmarks for DASH-3D on mobile devices

The growing capabilities and usage of mobile devices, especially smartphones, nowadays incur a progressive shift of many applications from desktop to mobile devices. In order to be made available and usable by the greater audience, 3D streaming and visualization should also be possible on mobile devices.

However, desktop devices tend to be much more powerful, have a larger memory and better network connections than mobile devices. In addition, the interactive modalities of these two types of devices are not comparable in any way: the desktop mostly uses keyboard and mouse, whereas most of the mobile devices only have a touchscreen, as well as various additional sensors (accelerometer, gyroscope, GPS, etc.). For these reasons, using DASH to stream 3D on mobile devices requires specific adaptations, that we describe in this chapter.

We add some widgets on the screen to support touch interactions: a virtual joystick is displayed on the screen and the user can touch it to translate the camera, instead of using the W, A, S and D keys on a computer keyboard. Since most mobile devices embed a gyroscope, we allow users to rotate the camera by physically rotating the device. This interaction is more precise and intuitive to the user, but it is also more tiring, this is why we also added a touch interaction to rotate the screen: a user can also "touch and drag" at any point on the screen that does not correspond to the joystick to rotate the camera. In order to ease navigation, we integrate bookmarks back, and we propose an enhanced version of the precomputations explained in Chapter 4 that we encode in the DASH Media Presentation Description. We then present a user study on 18 participants, which evaluates how users perceive the visual quality of the scene, and how their interactions affect it.

5.1 Introduction

In Chapter 3, we described how it is possible to modify a user interface to ease user navigation in a 3D scene, and how the system can benefit from it. In Chapter 4, we presented the DASH-3D streaming system, which does not depend on the interface nor on the user interaction. In this chapter, we will analyze how the user interaction can impact performances of DASH-3D. In order to do so, we follow these two steps based on our DASH framework:

- we design an interface allowing to navigate in a 3D scene on both desktop and mobile devices;
- we improve and adapt the bookmarks described in Chapter 3 to the context of DASH-3D and to mobile interaction.

In Chapter 5.2, we present the different choices we made for the interfaces, and we describe the new mobile interface. In Section ref{sb:bookmarks}, we describe how we embed the bookmarks into our DASH framework, and how we precompute data in order to improve the user quality of experience. In Section ref{sb:evaluation}, we describe the user study we conducted, the data we collected and we analyse the results.

5.2 Desktop and mobile interactions

5.2.1 Desktop interaction

Regarding desktop interaction, we keep the interaction we described in Section ref{bi:our-nve}, namely:

- W, A, S and D keys to translate the camera;
- mouse motions to rotate the camera.

A screenshot of this interface is displayed in Figure 37.



Figure 37: Screenshot of the desktop version, with a bookmark and its thumbnail on the bottom left corner and three bookmarks

5.2.2 Mobile interaction

Mobile interactions are more complex because the user does not have the keyboard and mouse to interact with. However, there are some other sensors on most mobile devices that can help interaction. One useful sensor for 3D interaction on mobile devices is definitely the gyroscope. We use the gyroscope to enable a user to rotate his device to rotate the virtual camera. We also add the possibility to rotate the camera by using touch controls. The user can touch a part of the screen to get a hold at the virtual camera, and drag the camera direction along the two screen axis. This way, the user is not forced to perform a real-world half-turn to be able to look behind or to point the device towards the sky (which can quickly become tiring) to look up. These interactions, however, only allow the user to rotate the camera but not translate it. For this reason, we display a small joystick on the bottomleft corner of the screen that mimics the first person video games interactions and allows the user travelling in the scene:

- moving the joystick up makes the camera move forward;
- moving the joystick down makes the camera move backwards;
- moving the joystick on the sides makes the camera move sideways.

A screenshot of this interface is displayed in Figure 38. The virtual joystick is rendered as a black circle inside a larger semi-transparent white circle. The black circle can be moved up, down, and sideways to define the direction in which the camera is translated.



Figure 38: Screenshot of the mobile version, with its joystick on the bottom left corner

5.3 Adding bookmarks into DASH NVE framework

While the previously defined interactions allow users to navigate freely throughout the scene, controlling such a high number of degrees of freedom can feel overwhelming to some users. That is why we introduce bookmarks, i.e. widgets that help the users reach a distant part of the scene using only a single, simple, interaction.

5.3.1 Bookmark interaction and visual aspect

In Chapter ref{bi} Section ref{bi:3d-bookmarks}, we described two 3D widgets that we use to display bookmarks to users. One of the conclusions of the user-study, described in Section ref{bi:user-study}, was that the impact of the way we display bookmark was not significant. In this work, we chose a slightly different way of representing bookmarks due to some concerns with our original representations:

- viewport bookmarks are simple, but non computer vision experts may not be familiar with this type of representation;
- arrow bookmarks are more intuitive to most users, but need to be regenerated when the camera moves, which can harm the rendering framerate.

For these reasons, we changed the bookmarks display to a vertical bar textured with a 2D sprite of a pictorial representation of an eye. The use of such symbol is partly inspired by the cartographic pictograms used to showcase a worthwhile panorama. This 2D sprite is always facing the camera to prevent it from being invisible when the camera would be on the side of it. Screenshots of user interfaces with bookmarks are available in Figure 37 and Figure 38.

The size of the sprite changes over time following a sine function to help the user distinguish what is part of the scene and what is extra widgets. Since our scene is static, a user knows that a changing object is not part of the scene, but part of the UI.

The other bookmark parameters remain unchanged since Chapter 3: in order to avoid users to lose context, clicking on a bookmark triggers an automatic, smooth, camera displacement that ends up at the bookmarked camera position. We also display a thumbnail of the bookmark's viewpoint when the mouse hovers a bookmark. Such thumbnail is displayed in Figure 37. Note that since on mobile, there is no mouse and thus no pointer, thumbnails are not used in the mobile setting.

5.3.2 Segments utility at bookmarked viewpoint

Introducing bookmarks is a way to make users navigation more predictable. Indeed, since they are emphasized and, in a way, recommended viewpoints, bookmarks are more likely to be visited by a significant portion of users than any other viewpoint on the scene. As such, bookmarks can be used as a way to optimize streaming by downloading segments in an optimal, precomputed order.

More specifically, segment utility as introduced in Section ref{d3:utility} is only an approximation of the segment's true contribution to the current viewpoint rendering. When bookmarks are defined, it is possible to obtain a better measure of segment utility by performing an offline rendering at each bookmark's viewpoint. We define $\mathcal{U}^*(s, B_i)$ as being the optimized utility of a segment *s* in a viewpoint defined at bookmark B_i .

In order to compute the optimized utility of a segment, we developed Algorithm ref{sb:algo-optimal-order}, that sorts segments according to their optimized utility. This algorithm takes as input the considered viewpoint, the ground truth rendering from this viewpoint and the set of segments (both geometry and texture) to sort. Starting from an empty model, each segment from the set of candidates is independently added to the scene, and the PSNR between the corresponding render and the ground truth render is computed. We can thus determine which segment brings the highest Δ PSNR /s, s being the size of the segment in bytes. Once the best segment is found, it is registered, and a new iteration begins. That way, we are able to generate an order of segments sorted by Δ PSNR /s. This order is then saved as a JSON file that a client can download in order to know which segments contribute the most to a certain viewpoint.

Sorting all the segments from the model would be an excessively time consuming computation. To speed up this algorithm, we only sort the 200 first best segments, and we choose these segments among a filtered set of candidates. To find those candidates, we reuse the ideas developed in Chapter ref{bi}. We render the "pixel to geometry segment" and "pixel to texture" maps, as shown in Figure 39. These renderings allow us to know which geometry segment and which texture correspond to each pixel, and filter out useless candidates.



Figure 39: A bookmarked viewpoint (left), a pixel to geometry segment map (center), and a pixel to texture map (right)

Figure ref{sb:precomputation} shows how this precomputation improves the quality of rendering. Each curve represents the PSNR one can obtain by downloading a certain amount of data following a streaming policy. The blue curve, labelled "Default order", is obtained by optimizing the utilities as defined in Section ref{d3:utility}, whereas the green curve labelled "Proposed order" uses the sorting computed in Algorithm ref{sb:algo-optimal-order}. We can observe that for the same amount of data downloaded, the optimized order reaches a higher PSNR which means that its utility metric is more accurate. Note that this curve is averaged over all the 9 bookmarks of the scene. These bookmarks are chosen to cover the widest area in the scene, and each one faces a particular object-of-interest.

5.3.3 MPD modification

We now present how to include bookmarks information in the Media Presentation Description (MPD) file. Bookmarks are fully defined by a position, a direction, and the additional content needed to properly render and use a bookmark in a system. This additional data consist in two files: a thumbnail of the point of view at the bookmark, along with the JSON file giving the optimal segment order for this viewpoint, as computed by Algorithm ref{sb:algo-optimal-order}. For this reason, for each bookmark, we create a separate adaptation set in the MPD. The bookmarked viewpoint information is stored as a supplemental property. Bookmarks adaptation set only contain one representation, composed of two segments: the thumbnail used as a preview for the desktop interface and the JSON file.

```
1
     <AdaptationSet>
2
         <SupplementalProperty value="156.4909,1.6267,-146.2062,</pre>
         157.43106,1.5476,-146.5379" />
3
      <BaseURL>b1/</BaseURL>
4
5
      <Representation>
       <BaseURL>repr1/</BaseURL>
6
7
       <SegmentList>
8
        <SegmentURL media="thumbnail.jpg" />
9
        <SegmentURL media="order.json" />
       </SegmentList>
10
      </Representation>
11
12
     </AdaptationSet>
```

Listing 12: MPD description of a geometry adaptation set, and a texture adaptation set

An example of a bookmark adaptation set is depicted on Listing 12. The three first values in the supplemental property are the camera position coordinates, and the three last values are the target point coordinates.

5.3.4 Loader modifications

We build on the loader introduced in Algorithm ref{d3:next-segment} to implement a client adaptation logic. We include a bookmark adaptation logic such that (i) when a bookmark is hovered for the first time, the corresponding thumbnail image as well as the JSON file containing the optimal order of the segments (see Listing 12) are downloaded, and (ii) when a bookmark is clicked, we switch from utility \mathcal{U} to optimized utility \mathcal{U}^* to determine which segments to download next.

5.4 Conclusion

In this chapter, our objective was to propose a mobile interface for DASH-3D and to integrate back the interaction aspects that we developed in Chapter 3. %We have seen that doing so is not trivial, and many improvements have been made. For aesthetics and performance reasons, the UI of the bookmarks has changed, and new interactions were proposed for free navigation in the 3D scene. We developed an algorithm that computes offline a better order of segments for bookmarks than what a greedy policy would do. We encoded this optimal order in a JSON file and we modified our MPD in order to give metadata about bookmarks to the client, as well as modified our client implementation to benefit from this. We then conducted a user study on 18 participants where users had to navigate in scenes with bookmarks and using various streaming policies. The results indicate that users prefer the optimized version of the client streaming policy, which is coherent with the PSNR values that we computed. The results also show that users who enjoy an optimized policy tend to use the bookmarks more.

Bibliography

- Burigat, Stefano, and Luca Chittaro. 2007. "Navigation in 3d Virtual Environments: Effects of User Experience and Location-Pointing Navigation Aids". *International Journal of Man-Machine Studies* 65 (11): 945–58
- Burtnyk, Nicolas, Azam Khan, George Fitzmaurice, and Gordon Kurtenbach. 2006. "Showmotion: Camera Motion Based 3d Design Review". In *Proceedings of the 2006 Symposium on Interactive 3d Graphics and Games*, 167–74
- Cheng, Wei, and Wei Tsang Ooi. 2008. "Receiver-Driven View-Dependent Streaming of Progressive Mesh". In Proceedings of the 18th International Workshop on Network and Operating Systems Support for Digital Audio and Video, 9–14
- Chiariotti, Federico, Stefano D'Aronco, Laura Toni, and Pascal Frossard. 2016. "Online Learning Adaptation Strategy for DASH Clients". In *Proceedings of the 7th International Conference on Multimedia Systems*, 8–9
- Chittaro, Luca, and Stefano Burigat. 2004. "3d Location-Pointing as a Navigation Aid in Virtual Environments". In *Proceedings of the Working Conference on Advanced Visual Interfaces, AVI 2004, Gallipoli, Italy, May 25-28, 2004, 267–74*
- Chittaro, Luca, and Subramanian Venkataraman. 2006. "Navigation Aids for Multi-Floor Virtual Buildings: A Comparative Evaluation of Two Approaches". In *Proceedings of the ACM Symposium on Virtual Reality Software and Technology*, 227–35
- Cignoni, Paolo, Fabio Ganovelli, Enrico Gobbetti, Fabio Marton, Federico Ponchio, and Roberto Scopigno. 2005. "Batched Multi Triangulation". In *VIS 05. IEEE Visualization, 2005.*, 207–14
- DASH. 2014. "Information technology Dynamic adaptive streaming over HTTP (DASH) Part 1: Media presentation description and segment formats"
- Eno, Joshua, Susan Gauch, and Craig W Thompson. 2010. "Linking Behavior in a Virtual World Environment". In *Proceedings of the 15th International Conference on Web 3d Technology*, 157–64

- Erikson, Carl, Dinesh Manocha, and William V Baxter III. 2001. "Hlods for Faster Display of Large Static and Dynamic Environments". In *Proceedings of the 2001 Symposium on Interactive 3d Graphics*, 111–20
- Guo, Jinjiang, Vincent Vidal, Irene Cheng, Anup Basu, Atilla Baskurt, and Guillaume Lavoue. 2017.
 "Subjective and Objective Visual Quality Assessment of Textured 3d Meshes". ACM Transactions on Applied Perception (TAP) 14 (2): 11–12
- Guthe, Michael, and Reinhard Klein. 2004. "Streaming Hlods: An Out-of-Core Viewer for Network Visualization of Huge Polygon Models". *Computers & Graphics* 28 (1): 43–50
- Hoppe, Hugues. 1996. "Progressive Meshes". In Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques, 99–108
- Hoppe, Hugues. 1998. "Smooth View-Dependent Level-of-Detail Control and Its Application to Terrain Rendering". In *Proceedings Visualization'98 (Cat. No. 98cb36276)*, 35–42
- Huang, Te-Yuan, Chaitanya Ekanadham, Andrew J Berglund, and Zhi Li. 2019. "Hindsight: Evaluate Video Bitrate Adaptation at Scale". In *Proceedings of the 10th ACM Multimedia Systems Conference*, 86–97
- Jankowski, Jacek, and Stefan Decker. 2012. "A Dual-Mode User Interface for Accessing 3d Content on the World Wide Web". In *Proceedings of the 21st International Conference on World Wide Web*, 1047–56
- Jankowski, Jacek, and Martin Hachet. 2015. "Advances in Interaction with 3d Environments". *Comput. Graph. Forum* 34 (1): 152–90
- Khan, Azam, Igor Mordatch, George Fitzmaurice, Justin Matejka, and Gordon Kurtenbach. 2008. "Viewcube: A 3d Orientation Indicator and Controller". In *Proceedings of the 2008 Symposium on Interactive 3d Graphics and Games*, 17–25. I3d '08. ACM
- Lavoué, Guillaume, Laurent Chevalier, and Florent Dupont. 2013. "Streaming Compressed 3d Data on the Web Using Javascript and Webgl". In *Proceedings of the 18th International Conference on 3d Web Technology*, 19–27
- Li, Frederick WB, Rynson WH Lau, Danny Kilis, and Lewis WF Li. 2011. "Game-on-Demand:: An Online Game Engine Based on Geometry Streaming". *ACM Transactions on Multimedia Computing, Communications, And Applications (TOMM)* 7 (3): 19–20
- Liang, Ke, Roger Zimmermann, and Wei Tsang Ooi. 2011. "Peer-Assisted Texture Streaming in Metaverses". In *Proceedings of the 19th ACM International Conference on Multimedia*, 203–12
- Limper, Max, Yvonne Jung, Johannes Behr, and Marc Alexa. 2013. "The Pop Buffer: Rapid Progressive Clustering by Geometry Quantization". In *Computer Graphics Forum*, 32:197–206
- Lindstrom, Peter, David Koller, William Ribarsky, Larry F Hodges, Nick L Faust, and Gregory Turner. 1996. "Real-Time, Continuous Level of Detail Rendering of Height Fields"
- Mackinlay, Jock D, Stuart K Card, and George G Robertson. 1990. "Rapid Controlled Movement Through a Virtual 3d Workspace". In *ACM SIGGRAPH Computer Graphics*, 24:171–76

- Maglo, Adrien, Ian Grimstead, and Céline Hudelot. 2013. "POMAR: Compression of Progressive Oriented Meshes Accessible Randomly". *Computers & Graphics* 37 (6): 743–52
- Maglo, Adrien, Guillaume Lavoué, Florent Dupont, and Céline Hudelot. 2015. "3d Mesh Compression: Survey, Comparisons, And Emerging Trends". *ACM Computing Surveys (CSUR)* 47 (3): 44–45
- Marvie, Jean Eudes, and Kadi Bouatouch. 2003. "Remote Rendering of Massively Textured 3d Scenes Through Progressive Texture Maps". In *The 3rd IASTED Conference on Visualisation, Imaging and Image Processing*, 2:756–61
- Meng, Fang, and Hongbin Zha. 2003. "Streaming Transmission of Point-Sampled Geometry Based on View-Dependent Level-of-Detail". In *Fourth International Conference on 3-D Digital Imaging and Modeling, 2003. 3dim 2003. Proceedings.*, 466–73
- Moerman, Clément, Damien Marchal, and Laurent Grisoni. 2012. "Drag'n Go: Simple and Fast Navigation in Virtual Environment". In *3d User Interfaces (3dui), 2012 IEEE Symposium on*, 15–18
- Niamut, Omar A, Emmanuel Thomas, Lucia D'Acunto, Cyril Concolato, Franck Denoual, and Seong Yong Lim. 2016. "MPEG DASH SRD: Spatial Relationship Description". In *Proceedings of the 7th International Conference on Multimedia Systems*, 5–6
- Ozcinar, Cagri, Ana De Abreu, and Aljosa Smolic. 2017. "Viewport-Aware Adaptive 360 Video Streaming Using Tiles for Virtual Reality". In *2017 IEEE International Conference on Image Processing (ICIP)*, 2174–78
- Portaneri, Cédric, Pierre Alliez, Michael Hemmer, Lukas Birklein, and Elmar Schoemer. 2019. "Cost-Driven Framework for Progressive Compression of Textured Meshes". In *Proceedings of the 10th ACM Multimedia Systems Conference*, 175–88
- Potenziani, Marco, Marco Callieri, Matteo Dellepiane, Massimiliano Corsini, Federico Ponchio, and Roberto Scopigno. 2015. "3dhop: 3d Heritage Online Presenter". *Computers & Graphics* 52: 129–41
- Rezzonico, Serge, and Daniel Thalmann. 1996. "Browsing 3D Bookmarks in BED". In *Proceedings of Webnet 96 - World Conference of the Web Society*
- Robinet, Fabrice, and P Cozzi. 2013. "Gltf—the Runtime Asset Format for Webgl". *Opengl ES, And Opengl*
- Ruddle, Roy A, Andrew Howes, Stephen J Payne, and Dylan M Jones. 2000. "The Effects of Hyperlinks on Navigation in Virtual Environments". *International Journal of Human-Computer Studies* 53 (4): 551–81
- Schilling, Arne, Jannes Bolling, and Claus Nagel. 2016. "Using Gltf for Streaming Citygml 3d City Models". In *Proceedings of the 21st International Conference on Web3d Technology*, 109–16. Web3d '16. Anaheim, California: ACM. https://doi.org/10.1145/2945292.2945312
- Schulzrinne, H., S. Casner, R. Frederick, and V. Jacobson. 1996. "RTP: A Transport Protocol for Real-Time Applications"

- Sideris, Anargyros, E Markakis, Nikos Zotos, Evangelos Pallis, and Charalabos Skianis. 2015. "M-PEG-DASH Users' Qoe: The Segment Duration Effect". In *2015 Seventh International Workshop on Quality of Multimedia Experience (Qomex)*, 1–6
- Simon, Gwendal, Stefano Petrangeli, Nathan Carr, and Viswanathan Swaminathan. 2019. "Streaming a Sequence of Textures for Adaptive 3d Scene Delivery". In *2019 IEEE Conference on Virtual Reality and 3d User Interfaces (VR)*, 1159–60
- Sodagar, Iraj. 2011. "The MPEG-DASH Standard for Multimedia Streaming Over the Internet". *IEEE Multimedia* 18 (4): 62–67. https://doi.org/10.1109/MMUL.2011.71
- Stockhammer, Thomas. 2011. "Dynamic Adaptive Streaming over HTTP -: Standards and Design Principles". In *Proceedings of the Second Annual ACM Conference on Multimedia Systems*, 133–44. Mmsys '11. San Jose, CA, USA: ACM. https://doi.org/10.1145/1943552.1943572
- Stohr, Denny, Alexander Frömmgen, Amr Rizk, Michael Zink, Ralf Steinmetz, and Wolfgang Effelsberg. 2017. "Where Are the Sweet Spots?: A Systematic Approach to Reproducible Dash Player Comparisons". In *Proceedings of the 25th ACM International Conference on Multimedia*, 1113–21
- Tian, Dihong, and Ghassan AlRegib. 2008. "Batex3: Bit Allocation for Progressive Transmission of Textured 3-D Models". *IEEE Transactions on Circuits and Systems for Video Technology* 18 (1): 23–35
- Todd, James T. 2004. "The Visual Perception of 3d Shape". Trends in Cognitive Sciences 8 (3): 115-21
- Varvello, Matteo, Stefano Ferrari, Ernst Biersack, and Christophe Diot. 2011. "Exploring Second Life". *IEEE/ACM Transactions on Networking (TON)* 19 (1): 80–91
- Yadav, Praveen Kumar, Arash Shafiei, and Wei Tsang Ooi. 2017. "Quetra: A Queuing Theory Approach to Dash Rate Adaptation". In *Proceedings of the 25th ACM International Conference on Multimedia*, 1130–38
- Yang, Sheng, Chao-Hua Lee, and C.-C. Jay Kuo. 2004. "Optimized Mesh and Texture Multiplexing for Progressive Textured Model Transmission". In *Proceedings of the 12th Annual ACM International Conference on Multimedia*, 676–83. MULTIMEDIA '04. New York, NY, USA: ACM. https://doi.org/ 10.1145/1027527.1027683
- Zampoglou, Markos, Kostas Kapetanakis, Andreas Stamoulias, Athanasios G Malamos, and Spyros Panagiotakis. 2018. "Adaptive Streaming of Complex Web 3d Scenes Based on the MPEG-DASH Standard". *Multimedia Tools and Applications* 77 (1): 125–48
- Zeleznik, Robert C, Andrew S Forsberg, and Paul S Strauss. 1997. "Two Pointer Input for 3d Interaction". In *Proceedings of the 1997 Symposium on Interactive 3d Graphics*, 115–20

Abstract

With the advances in 3D models editing and 3D reconstruction techniques, more and more 3D models are available and their quality is increasing. Furthermore, the support of 3D visualization on the web has become standard during the last years. A major challenge is thus to deliver these remote heavy models and to allow users to visualise and navigate in these virtual environments. This thesis focuses on 3D content streaming and interaction, and proposes three major contributions.

First, we develop a 3D scene navigation interface with bookmarks – small virtual objects added to the scene that the user can click on to ease reaching a recommended location. We describe a user study where participants navigate in 3D scenes with and without bookmarks. We show that users navigate (and accomplish a given task) faster when using bookmarks. However, this faster navigation has a drawback on the streaming performance: a user who moves faster in a scene requires higher streaming capabilities in order to enjoy the same quality of service. This drawback can be mitigated using the fact that bookmarks positions are known in advance: by ordering the faces of the 3D model according to their visibility at a bookmark, we optimize the streaming and thus, decrease the latency when users click on bookmarks.

Secondly, we propose an adaptation of Dynamic Adaptive Streaming over HTTP (DASH), the video streaming standard, to 3D textured meshes streaming. To do so, we partition the scene into a k-d tree where each cell corresponds to a DASH adaptation set. Each cell is further divided into DASH segments of a fixed number of faces, grouping together faces of similar areas. Each texture is indexed in its own adaptation set, and multiple DASH representations are available for different resolutions of the textures. All the metadata (the cells of the k-d tree, the resolutions of the textures, etc.) is encoded in the Media Presentation Description (MPD): an XML file that DASH uses to index content. Thus, our framework inherits DASH scalability. We then propose clients capable of evaluating the usefulness of each chunk of data depending on their viewpoint, and streaming policies that decide which chunks to download.

Finally, we investigate the setting of 3D streaming and navigation on mobile devices. We integrate bookmarks in our 3D version of DASH and propose an improved version of our DASH client that benefits from bookmarks. A user study shows that with our dedicated bookmark streaming policy, bookmarks are more likely to be clicked on, enhancing both users quality of service and quality of experience.

Résumé

Avec les progrès de l'édition de modèles 3D et des techniques de reconstruction 3D, de plus en plus de modèles 3D sont disponibles et leur qualité augmente. De plus, le support de la visualisation 3D sur le web s'est standardisé ces dernières années. Un défi majeur est donc de transmettre des modèles massifs à distance et de permettre aux utilisateurs de visualiser et de naviguer dans ces environnements virtuels. Cette thèse porte sur la transmission et l'interaction de contenus 3D et propose trois contributions majeures.

Tout d'abord, **nous développons une interface de navigation dans une scène 3D avec des signets** – de petits objets virtuels ajoutés à la scène sur lesquels l'utilisateur peut cliquer pour atteindre facilement un emplacement recommandé. Nous décrivons une étude d'utilisateurs où les participants naviguent dans des scènes 3D avec ou sans signets. Nous montrons que les utilisateurs naviguent (et accomplissent une tâche donnée) plus rapidement en utilisant des signets. Cependant, cette navigation plus rapide a un inconvénient sur les performances de la transmission : un utilisateur qui se déplace plus rapidement dans une scène a besoin de capacités de transmission plus élevées afin de bénéficier de la même qualité de service. Cet inconvénient peut être atténué par le fait que les positions des signets sont connues à l'avance : en ordonnant les faces du modèle 3D en fonction de leur visibilité depuis un signet, on optimise la transmission et donc, on diminue la latence lorsque les utilisateurs cliquent sur les signets.

Deuxièmement, **nous proposons une adaptation du standard de transmission DASH (Dynamic Adaptive Streaming over HTTP)**, **très utilisé en vidéo, à la transmission de maillages texturés 3D**. Pour ce faire, nous divisons la scène en un arbre k-d où chaque cellule correspond à un adaptation set DASH. Chaque cellule est en outre divisée en segments DASH d'un nombre fixe de faces, regroupant des faces de surfaces comparables. Chaque texture est indexée dans son propre adaptation set à différentes résolutions. Toutes les métadonnées (les cellules de l'arbre k-d, les résolutions des textures, etc.) sont référencées dans un fichier XML utilisé par DASH pour indexer le contenu: le MPD (Media Presentation Description). Ainsi, notre framework hérite de la scalabilité offerte par DASH. Nous proposons ensuite des algorithmes capables d'évaluer l'utilité de chaque segment de données en fonction du point de vue du client, et des politiques de transmission qui décident des segments à télécharger.

Enfin, **nous étudions la mise en place de la transmission et de la navigation 3D sur les appareils mobiles**. Nous intégrons des signets dans notre version 3D de DASH et proposons une version améliorée de notre client DASH qui bénéficie des signets. Une étude sur les utilisateurs montre qu'avec notre politique de chargement adaptée aux signets, les signets sont plus susceptibles d'être cliqués, ce qui améliore à la fois la qualité de service et la qualité d'expérience des utilisateurs.