

# Multiresolution and fast decompression for optimal web-based rendering

Federico Ponchio<sup>a</sup>, Matteo Dellepiane<sup>a</sup>

<sup>a</sup>Visual Computing Lab, ISTI CNR, Pisa, Italy

---

## Abstract

Limited bandwidth is a strong constraint when efficient transmission of 3D data to Web clients and mobile applications is needed. In this paper we present a novel multi-resolution WebGL based rendering algorithm which combines progressive loading, view-dependent resolution and mesh compression, providing high frame rates and a decoding speed of million of triangles per second in JavaScript. The method is parallelizable and scalable to very large models.

The algorithm is based on the local multi-resolution approaches provided by the community, but ad-hoc solutions had to be studied and implemented to provide adequate performances. In particular, a compression mechanism that reached very high compression rate without impact on rendering performance was implemented. Moreover, the data partition strategy was modified in order to be able to load different types of data (i.e. point clouds) and better adapt to the potentials and limitations of web-based rendering.

*Keywords:* multiresolution, WebGL, 3D Web, web based 3D rendering, online 3D content deployment, mesh compression

---

## 1. Introduction

Limited bandwidth and increasing model sizes pose a challenge in the transmission of 3D data to Web clients and mobile applications. A possible approach is to compress the 3D model in order to minimize transmission time. Most of the research in this field has been focused on optimizing compression ratio.

Unfortunately, limited bandwidth often pairs with limited computational power, either because of JavaScript environment or low CPU power mobile devices, to the point that for most algorithms decoding time becomes the bottleneck even at moderately low bandwidth. Acceptable rates can be regained reducing compression ratio or using less sophisticated entropy compression algorithms.

A different approach makes use of progressive reconstruction algorithms, which improve the user experience by providing a simplified version of the model that refines while the remaining part of the model is being downloaded. The model converges very quickly at the beginning of the download, and only the details require the full model. However this class of algorithms performs even worse in terms of decoding time (as shown in Limper [? ]) or in terms of compression ratio.

In the context of multi-resolution methods another desirable feature could be the view-dependent refinement. This allows to prioritize the download, decode a specific part of the model and vary resolution of the rendered geometry to maintain a constant screen resolution. This is obtained by maximizing quality at a given frame rate.

In addition to the above limitations, the 3D models that are now available on the web cover a much broader range of possibilities w.r.t. the past, including point clouds, *triangle soups*, topologically complex geometries, partially textured models. This leads to the necessity to propose a framework which could be robust enough to deal with different cases.

34

In this paper we present a novel multi-resolution WebGL based rendering algorithm (Figure 1) which combines *progressive loading*, *view-dependent resolution* and *mesh compression*, providing good rates and a decoding speed of million of triangles per second in JavaScript. Additionally, the method is *flexible*, since it is able to handle a variety of 3D formats, including textured models, non-manifold meshes, point clouds. Finally, it is also *scalable*, since it's able to deal with very large models.

43

The method is based on a class of multiresolution structures [? ? ] where the “primitive” of the multiresolution is a patch made of thousands of triangles. The original approach was written to:

48

- obtain a more efficient data partition, and handle more data formats than triangulated surfaces

49

- compress the data structure in order to save disk space and bandwidth with no impact on performances

50

- extend the original algorithm to remote rendering

51

The paper is organized as follows: Section 2 provides an overview of related works. In Section 3 we describe the multiresolution structure, focusing on the improvements over the existing method, and on how the compression algorithm was designed to optimize decoding time while maintaining an adequate compression ratio. In Section 4 we compare it with existing web solutions for mesh compression and progressive visualization, and we analyze the performances when dealing with different classes of 3D models. The proposed method represents a solid alternative to current solutions, providing a practical mean to handle 3D models on the web.

63

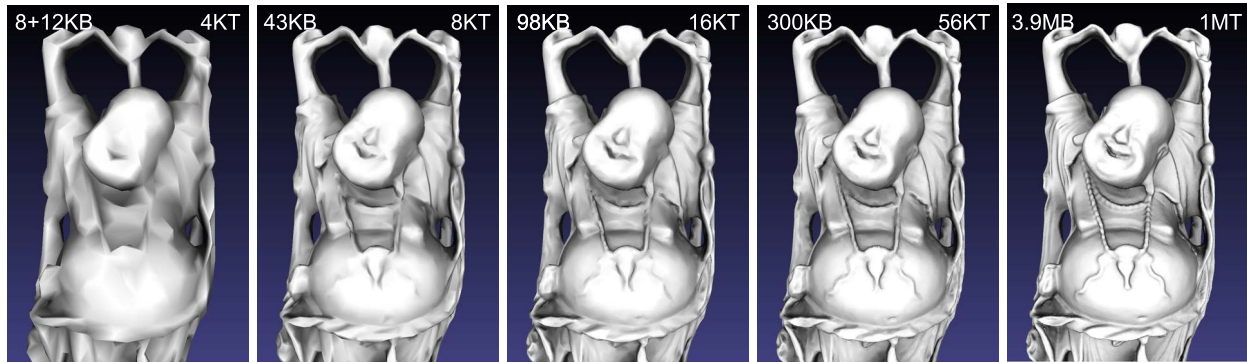


Figure 1: Progressive refinement of the Happy Buddha: on the upper left corner the size downloaded, on the upper right corner the number of triangles in the refined model. The header and index amount to 8KB

## 2. Related Work

This paper is related to several topics in the field of Computer Graphics. Among them: web-based 3D rendering, progressive and multiresolution rendering approaches, and fast decompression methods for 3D models.

While a complete overview of all these subjects goes well beyond the scope of the paper, in the next subsections we provide a short description of the state of the art, trying to focus on the aspects which are more related to the proposed approach.

### 2.1. Web-based 3D rendering

Three-dimensional content has always been considered as part of the multimedia family. Nevertheless, especially when talking about web visualization, its role with respect to images and videos has always been a minor one. Visualization of 3D components was initially devoted to external components, such as Java applets or ActiveX controls [? ].

After some initial efforts for standardization [? ], the proposal of WebGL standard [? ], which is a mapping of OpenGL|ES 2.0 [? ] specifications in JavaScript, brought a major change. Several actions related to the use of advanced 3D graphics has been proposed since then. For a general survey, please refer to the work by Evans [? ]. Since the use of OpenGL commands needs advanced programming skills, there have been several actions to provide an "interface" between them and the creation of web pages. We could subdivide the proposed systems between *declarative* approaches [? ], like X3DOM [? ] or XML3D [? ], and *imperative* approaches, like Three.js [? ], SpiderGL [? ] and WebGLU [? ]. The main difference between the groups is that the first ones rely on the concept of *scenegraph*, hence a scene has to be defined in all its elements, while the second ones provide a more direct interface with the basic commands. Other systems provide a sort of hybrid approach [? ], where a simple scene has to be defined.

Evans [? ] points out in his survey that declarative approaches had a major impact in the research community, while imperative approaches were mainly used in the programming community. More in general, given the fact that the amount of data that needs to be sent to the webpage can be quite big, several efforts

about a better organization of generic streamable formats [? ] have been proposed. Nevertheless, when complex 3D data have to be streamed, these structures are not flexible enough to handle them.

In order to face this problem, in the last three years some progressive compression methods ad hoc for 3D streaming have been developed. Gobbetti et al. [? ] proposed a quad-based multi-resolution format. Behr et al. [? ] transmit different quantization levels of the geometry using a set of nested GPU-friendly buffers. Lavouè et al. [? ] proposed an adaptation for the Web (reduced decompression time at the cost of a low compression ratio) of a previous progressive algorithm [? ]. Other research has been also conducted to handle other types of data, like point clouds [? ], which may present different types of issues to face with.

The rendering of textures or textured 3D models has been taken into account even before the standardization actions. In these cases the main issue is the amount of image data: standard techniques like mip-mapping can be adapted and improved both on the software and hardware side [? ]. The issue of handling geometric data has been usually considered a minor one, due to the usual low complexity of 3D textured models [? ]. Nevertheless, recently complex 3D models with texture coordinates are available from acquisition devices and technologies. Next subsections will provide further details.

### 2.2. Progressive and Multi-resolution methods

An important feature for user experience when rendering over slow connections or compressed models is progressive-ness: the possibility to temporarily display an approximated version of the model and to refine it while downloading or processing the rest of the data.

The simplest (and widely used) strategy is to use a discrete set of increasing resolution models (usually known as Level Of Detail, LOD). The main drawback with this approach is the abrupt change in detail each time a model is replaced.

A change of paradigm was brought by progressive meshes, introduced by Hoppe [? ]. These meshes encode the sequence of operations of an edge collapse simplification algorithm. This sequence is traversed in reverse, so that each collapse becomes

142 a split, and the mesh is refined until the original resolution. An  
143 advantage of progressive techniques is the much more smooth  
144 transition resolution changes, and the possibility to combine it  
145 with selective refining or view-dependent multiresolution, but  
146 this high granularity was achieved at the cost of low compression  
147 rates: about 37 bpv with 10 bit vertex quantization.

148 A large number of progressive techniques were later developed,  
149 but as noted in [? ], Table 1, the research focus, however,  
150 was on rate-distortion performances and speed was mostly neglected.  
151 Latest algorithms still run below 200KTs in CPU.

152 Mobile and web application would be really too slow using  
153 these methods. As a compromise, pop buffers [? ] propose  
154 a method to progressively transmit geometry and connectivity,  
155 while completely avoiding compression.

156 Another desirable feature, especially for large models, is  
157 view-dependent loading and visualization. Most multiresolution  
158 algorithms were made obsolete by the increased relative  
159 performances of GPU over CPU around the first years of 2000.  
160 It simply became inefficient to operate on the mesh at the level  
161 of the single triangle. Several works [? ? ? ? ] achieved  
162 much better performances by increasing the granularity of the  
163 multiresolution to a few thousand triangles.

164 The main problem when increasing the granularity is ensuring  
165 boundary consistency between patches at different resolution:  
166 Yoon [? ] and Sander [? ] both employ a hierarchical  
167 spatial subdivision, but while the first simply disables simplification  
168 of most boundary edges, which results in scalability problems,  
169 the second relies on global, spatial GPU geomorphing to ensure  
170 that progressive meshes patch simplification is consistent between  
171 adjacent blocks. The works by Cignoni [? ? ] rely instead on a  
172 non hierarchical volumetric subdivision and a boundary preserving  
173 patch simplification strategies that guarantee coherence between  
174 different resolutions, while at the same time ensuring that no  
175 boundary persists for more than one level. While not progressive  
176 in a strict sense, given current rendering speed, the density of  
177 triangles on screen is so high that popping effects are not  
178 noticeable.

179 An additional issue when dealing with view-dependent multiresolution  
180 techniques is the handling of textured models. While the encoding  
181 of texture coordinates can be easily taken into account when  
182 creating the patches of different resolution, the boundary consistency  
183 among them needs to take into account the texture images. Previous  
184 multi-resolution methods [? ? ] proposed solutions for this, but  
185 they could fail when dealing with complex geometries.

186  
187  
188 Compression comes as a natural extension to this family of  
189 multiresolution algorithms: each patch can be compressed independently  
190 from the others as long as the boundary still matches with neighboring  
191 patches. A wavelet based compression was developed in [? ] for  
192 terrains, a 1D Haar wavelet version in [? ] for generic meshes on  
193 a mobile application. A comprehensive account of compression  
194 algorithms and the convergence with view-dependent rendering of  
195 large datasets can be found on a recent survey from Maglo et al.[? ].

### 197 2.3. Fast Decompression of 3D models

198 Given that decompression speed is a key factor in order to  
199 be able to use compressed mesh, there's been surprisingly little  
200 effort by the community to provide solutions.

201 Gumhold and Straßer [? ] developed a connectivity only  
202 compression algorithm that was able to decompress at 800KTs in  
203 1998. Pajarola and Rossignac in [? ], in 2000, reported 26KTs  
204 for a progressive compression algorithm, and developed a high-  
205 performance Huffman decoding identifying entropy compression  
206 as a possible bottleneck.

207 Finally, Isenburg and Gumhold in 2003 [? ] developed a streaming  
208 approach to compression of gigantic meshes reaching an impressive  
209 decompression speed of 2MTs. The method accounts also for texture  
210 coordinates. A further work on this was proposed in 2005 [? ].  
211 We don't know of any following work specifically geared toward  
212 fast decompression. Even the recent survey by Maglo [? ] shows  
213 that most of the effort is devoted to compression quality, but the  
214 performances of the works by Isenburg haven't been matched yet.

## 216 3. Method

217 Our multiresolution algorithm builds upon the methods described  
218 on [? ? ], which are recapped in Section 3.1 for completeness.  
219 This Section provides a description of the novel method, by  
220 presenting an improved partition strategy, and a novel compression  
221 scheme (Section 3.2) tailored around the need for decompression  
222 speed, which is obtained using entropy encoding 3.3. Finally,  
223 the implementation for remote rendering (section 3.4) is presented.

### 225 3.1. Batched Multiresolution

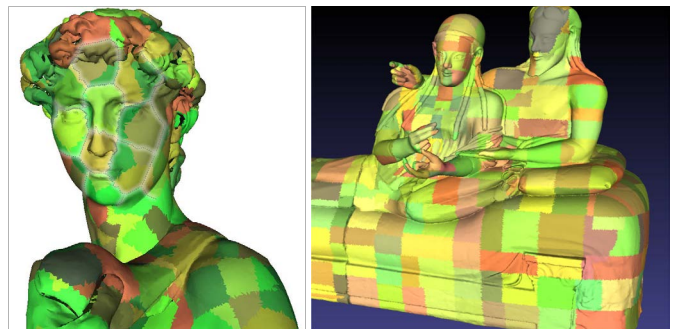


Figure 2: Left: volume partition by Cignoni [? ]. Right: the volume partition obtained with our method.

226 In a multiresolution method, the model is split into a set of  
227 small meshes at different resolutions, obtained through a simplification  
228 process, that can be assembled to create a seamless mesh by simply  
229 traversing a tree which encodes the dependencies between each patch,  
230 using the estimated screen error to select the resolution needed in  
231 each part of the model. The screen space error is computed starting  
232 from the geometric error due to the simplification process, and  
233 computing the corresponding size in pixel when projected on screen.

235 For the simplification of the mesh we used the Quadric Edge  
236 Collapse method [? ]. This simplification proved to be the  
237 most accurate and reliable, even though the simplification speed  
238 may be lower than other algorithms. Since partitions have to  
239 be created in a pre-processing stage, it was decided to use the  
240 slower alternative to obtain more accurate results. Moreover,  
241 quadric edge collapse can be easily extended to handle textured  
242 models (see later).

243 To build this collection of patches we need a sequence of  
244 non-hierarchical volume partitions (V-partition) of the the model;  
245 non hierarchical means essentially that no boundary is preserved  
246 between partitions at different levels of the hierarchy.

247 Cignoni et al [? ] showed that any non-hierarchical se-  
248 quence of volume partitions can be the base of a patch based  
249 multiresolution structure. Good partition strategies minimize  
250 boundaries, thus generating compact cells. In addition, they  
251 allow streaming construction and generate well balanced trees  
252 even when the distribution of the model triangles is very irregu-  
253 lar. They used the Voronoi structure (see Figure 2, left), which  
254 is optimal for boundary minimization and balance. However,  
255 this partition is not suitable for streaming, leading to long pro-  
256 cessing times. On the other hand, the regular spatial subdivision  
257 used the previous version of the method [? ] might generate un-  
258 balanced trees for very irregular models. This may impact on  
259 adaptivity.

260 We propose a different volume partition, defined by the leaves  
261 of a KD-tree built on the triangles of the model; to ensure the  
262 non hierarchical condition, the split ratio in the KD-tree alter-  
263 nates between 0.4 and 0.6 instead of the usual 0.5.

264 Figure 2 shows a typical partition provided by our method. The  
265 choice of this partition helps to better stream the model and  
266 provides better adaptivity. Additionally, the very regular shape  
267 of the patches may be useful when adding texture support (see  
268 later).

### 269 3.1.1. Extension to textured models and point clouds

270 The main challenge with multiresolution textured models  
271 lies in the simplification algorithm: it needs to take into account  
272 texture seams, and minimize deformation of the texture on the  
273 surface. We used the algorithm employed in Meshlab [? ],  
274 which is an extension of Garland work [? ] based on quadrics  
275 in 5 dimensions which include texture coordinates in selection  
276 of the collapse and the vertex optimization computation.

277 Enabling support for point clouds requires a few changes,  
278 merging into the same data structure, the functionality of Batched  
279 Multitriangulation [? ] and those of Layered Point Clouds [? ]:  
280 the simplification algorithm needs to be replaced with a point  
281 filtering approach, where half of the points are removed at each  
282 level. The octree structure of this data sampling perfectly fits in  
283 the compression and rendering paradigms explained in the next  
284 sections.

### 285 3.2. Mesh Compression

286 When porting multiresolution methods to the web, band-  
287 width becomes a further limitation. Hence, mesh compression  
288 becomes valuable not only to improve rendering performances,

289 but also to reduce disk space occupancy.

290 However, in addition to the issue of fast decompression, our  
291 multiresolution algorithm imposes a set of constrains to mesh  
292 compression:

- 293 • each patch needs to be encoded independently from the  
294 others, so the method must be efficient and fast even on  
295 small meshes
- 296 • boundary vertices, replicated on neighboring patches, need  
297 to remain consistent through compression
- 298 • non manifold models must be supported

299 It would be possible to exploit the redundancy of the data  
300 due to the fact that the same surface is present in patches at  
301 different levels of resolution. We choose not to do so in order  
302 to keep the compression stage independent of the simplification  
303 algorithm used and to simplify parallel decompression of the  
304 patches. Otherwise, we would have to keep track of and enforce  
305 dependencies. In the next subsections, the strategies for the  
306 compression of the elements of the 3D models are shown.

#### 307 3.2.1. Connectivity compression

308 The compression of connectivity can be difficult to handle,  
309 especially in the case of non manifold meshes. We modified the  
310 algorithm presented in [? ], that supports non manifold meshes  
311 and surfaces with handles or holes.

312 The face-face topology for compression is computed as fol-  
313 lows: we create an array containing three edges for each tri-  
314 angle, and sort it so that edges sharing the same vertices will  
315 be consecutive (independently of the order of the edges). The  
316 edges are then paired taking orientation into account, and all  
317 non paired edges are marked as boundary. Non manifold meshes  
318 will simply force the creation of some artificial boundaries.

319 The encoding process starts with a triangle and expands iter-  
320 atively adding triangles. The processed region is always home-  
321 omorphic to a disk and if the region meets already considered  
322 triangles, we consider the common vertices as duplicated. The  
323 boundary of the already processed (encoded or decoded) region  
324 is stored as a doubly linked list of oriented edges (*active edges*),  
325 The list is actually implemented as an array for performances  
326 reasons. A queue keeps track and prioritize the *active edges*.

327 The first triangle adds three active edges to the list; itera-  
328 tively an edge is extracted from the queue and, if not marked as  
329 processed, the following codes are emitted (see Figure 3):

330 **SKIP** if the edge is a boundary edge, or the adjacent triangle  
331 has already been encoded; the edge is marked as processed.

332 **LEFT** or **RIGHT** if the adjacent triangle shares two edges  
333 with the boundary; The two edges are marked as processed,  
334 a new edge added to the queue and its boundary adjacencies  
335 adjusted.

336 **VERTEX** if the adjacent triangle shares only one edge with  
337 the boundary, in this case the edges is marked as processed and  
338 two new edges added to the queue. If vertex of the new triangle  
339 opposing the edge was never encountered before its position  
340 is estimated using parallelogram prediction and the difference  
341 encoded, otherwise its index is encoded (in literature this case

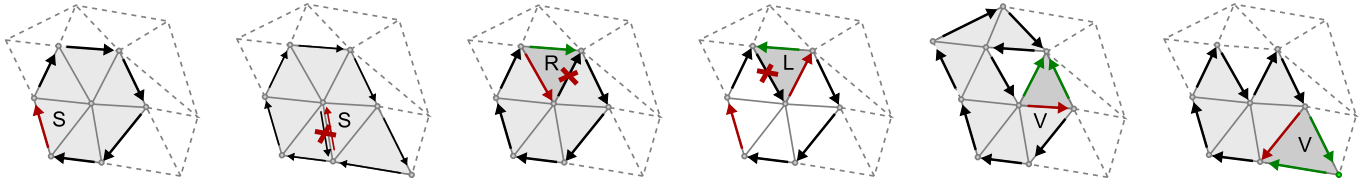


Figure 3: The four decomposition codes: black arrows represent the front, the red arrow the current edge, in green the new edges added to the front.

342 is often referred as a “split”). This is a key difference with [? 343 ], where in the second case a SKIP code would be emitted, to 344 keep the encoded region simple.

345 If the mesh is composed of several connected components, 346 the process is restarted for each component.

347 The order in which the active edges are processed is im- 348 portant as we would like to minimize the number of VERTEX 349 split operations, and generate a vertex-cache-friendly triangle 350 order. To do so, we simply prioritize the right edges in the 351 VERTEX operation, so that the encoding proceeds in ‘spirals’. 352 If the mesh is not homeomorphic to a disk, some split opera- 353 tions are required. This strategy reduces the number of splits 354 to less than 1% in our examples, incurring in an average of 0.2 355 bpv cost.

356 This algorithm is certainly not optimal in term of bitrate, 357 but it is extremely simple, linear in the number of triangles and 358 robust to non-manifold meshes; as we will see in the results, 359 speed is more important than bitrate.

### 360 3.2.2. Geometry and vertex attribute compression

361 To ensure consistency between boundary vertices of adja- 362 cent patches, we adopt a global quantization grid for coordi- 363 nates, normals and colors. The global grid step for vertex posi- 364 tion quantization is chosen automatically, based on the quadric 365 errors during the simplification step in construction.

366 Geometry and vertex attributes are encoded as differences 367 to a predicted value. The distribution of these values exhibit a 368 bias which we can exploit to minimize the number of bits neces- 369 sary to encode them. Our strategy is based on the assumption 370 that most of the bias is concentrated on the position of highest 371 bit (the  $\log_2$  of the value) of these value while the subsequent 372 bits are mostly random. We simply store in an array, which is 373 later entropy coded, the number of bits necessary to encode the 374 value; the subsequent bits are stored in an uncompressed bit- 375 stream. In this way we need to decode a single symbol, from a 376 limited alphabet, and read a few bits from a bitstream to decode 377 a difference.

378 Each new vertex position, result of a VERTEX code, is es- 379 timated using a simple parallelogram predictor, and the differ- 380 ences with the actual position encoded as above. Color infor- 381 mation is first converted into YCbCr color space and quantized, 382 we encode the difference with one of the corner of the edge pro- 383 cesses when emitting the VERTEX code. Normals vector are 384 estimated using the decode mesh position and connectivity, and 385 differences encoded as usual.

### 386 3.2.3. Texture coordinates and textures

387 Texture coordinates in the dataset are stored per vertex, repli- 388 cating vertices on texture seams. They are compressed using the 389 same parallelogram prediction algorithm employed for the ver- 390 tex coordinates. Employing more sophisticated methods would 391 drastically increase decoding time (see timings in Table 2 in [? 392 ]) or require additional linear algebra JavaScript libraries, for a 393 limited decrease in bitrate.

394 Texture images are first mip-mapped to create different levels of 395 detail, and then stored into the dataset as JPEG binaries. Tex- 396 tures are loaded on demand like the mesh patches, hence the 397 mipmap level for rendering follows the same rules described in 398 Section 3.4.

### 399 3.2.4. Point clouds

400 In the case of point clouds, the compression strategy for 401 the vertex coordinates cannot rely on parallelogram prediction, 402 in this case, after coordinate quantization we sort the points in 403 z-order and store the differences between consecutive points us- 404 ing the same approach used for the meshes.

### 405 3.3. Entropy coding

406 Once that connectivity, geometry and attributes have been 407 encoded into a stream of symbols and bits, the symbol stream 408 is compressed following the biased probability distribution of 409 the symbols.

410 Entropy decoding is the speed bottleneck in many mesh de- 411 compression methods, often due to the main goal of minimizing 412 bit per vertex. Pajarola and Rossignac [? ] developed a high- 413 performance Huffman decoding algorithm in order to overcome 414 this problem. The main advantage of this method is that it re- 415 duces the decoding phase to a couple of table lookups. Arith- 416 metic coding, for example, outperforms Huffman in term of 417 compression rate, but exhibits lower speed. A problem with this 418 approach is the initialization time required to create the, possi- 419 bly very large, decoding tables. It is then not suitable for decod- 420 ing small meshes where the construction time would dominate 421 over the decoding time.

422 Unlike Huffman and other variable-length codes, Tunstall 423 code [? ] maps a variable number of source symbols to a fixed 424 number of bits. Since in decompression the input blocks con- 425 sists of a fixed number of bits and the output is a variable num- 426 ber of symbols, Tunstall is slightly less efficient than Huffman, 427 especially where the bit size of the input block is small. The 428 decoding step is very similar to the high-performance Huffman 429 algorithm, as it consists in a lookup table and a sequence of

430 symbols for each entry, but the table size is only determined by  
 431 the word size, and a fast method to generate it described in [? ].

432 Given an entropic source of  $M$  symbols, to generate an opti-  
 433 mal encoding table for a word size of  $N$  bits, we need to gener-  
 434 ate  $2^N$  symbol sequences that have a frequency as close as  
 435 possible to  $2^{-N}$ , allows to encode every possible input (it is  
 436 complete) and no sequence is a prefix of any other sequence  
 437 (it is proper).

438 Tunstall optimal strategy starts with the  $M$  symbols as initial  
 439 sequences, removes the most frequent sequence  $A$  and replaces  
 440 it with  $M$  sequences concatenating  $A$  with every symbol until  
 441 we reach  $2^N$  sequences. The most time consuming step of the  
 442 algorithm is to find the most probable sequence.

443 If we use a matrix where the first column contains the sorted  
 444 symbol in order of probability, and at each step we replace the  
 445 sequence with highest probability with  $M$  sequences adding a  
 446 new column, we can observe that this table is sorted both in  
 447 columns and rows (see Figure 4). This allows to select the next  
 448 sequence by keeping each row in a queue and using a priority  
 449 queue to keep track of which queue has the highest front ele-  
 450 ment.

A 0.50	AA 0.25	BA 0.15	AAA 0.125	BAA 0.075
B 0.30	AB 0.15	BB 0.09	AAB 0.075	BAB 0.045
C 0.10	AC 0.05	BC 0.03	AAC 0.025	BAC 0.015
D 0.10	AD 0.05	BD 0.03	AAD 0.025	BAD 0.015

Figure 4: First four steps in construction of a Tunstall code with four symbols, the sequences A, B, AA, BA are replaced with a new column, beside each sequence, its probability is shown. In green the candidates for the next expansion.

451 To initialize the decoding table the symbol frequencies needs  
 452 to be transmitted in advance.

453 Finally, an important advantage of variable-to-fixed coding  
 454 is that the compressed stream is random accessible: decoding  
 455 can start at any block. This makes it especially suited for paral-  
 456 lel decompression in GPU. Unfortunately, current limitations  
 457 in the capabilities of WebGL do not allow for such an imple-  
 458 mentation.

### 459 3.4. Remote view-dependent rendering

460 In the context of batched multiresolution approaches, the  
 461 rendering requires the traversal of the patch tree, which is usu-  
 462 ally quite small since each patch is in the range of 8-32K ver-  
 463 tices. An approximated screen space error in pixel is calculated  
 464 by taking into account the view matrix, the bounding sphere of  
 465 the patch, and the quadric error (or any other error metric) cal-  
 466 culated during simplification.

467 At each step of the traversal, the selected nodes of the tree cor-  
 468 respond to a complete representation of the model; each addi-  
 469 tional node included in the traversal increases locally the reso-  
 470 lution of the model. The traversal is stopped whenever one of  
 471 these four conditions is reached: a required patch is not avail-  
 472 able, **triangle budget** is reached, the **error target** is met, or

473 there's no available **memory left**. The latter three parameters  
 474 can be defined in advance, in order to deal with hardware re-  
 475 sources. Once the traversal is terminated, the collection of se-  
 476 lected patches is rendered as simple geometry. New patches are  
 477 downloaded in order of screen error, to maximize the improve-  
 478 ment in rendering quality of the model.

479 Since the rendering can start when the first patch is down-  
 480 loaded and the model is refined as soon as some patch is avail-  
 481 able, this is effectively a progressive visualization albeit with  
 482 higher granularity. On the other hand, this structure is view de-  
 483 pendent and thus able to cope with very large models, on the  
 484 order of hundreds of millions of triangles.

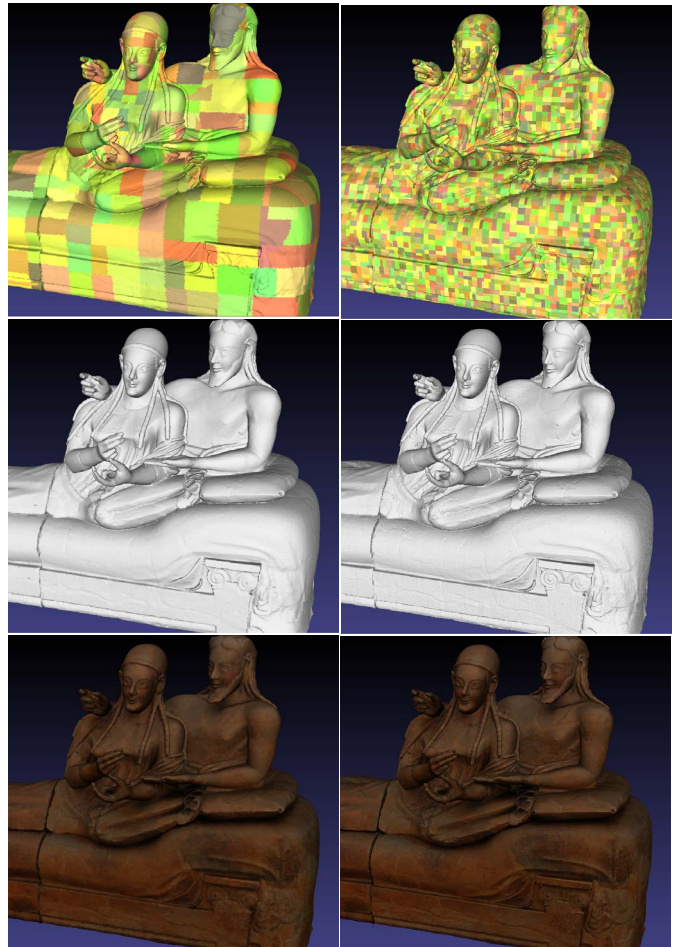


Figure 5: First column: before refinement. Second column: after refinement. From top to bottom: a visual representation of the geometric patches representing the model, the model with pure geometry, the model with color information.

485 The extension of the local multiresolution solution to re-  
 486 mote rendering is quite straightforward. The data structure is  
 487 composed of a fixed size header describing the attributes of  
 488 the models, an small index which contains the tree structure  
 489 of the patches and the position of each patch in the file, and the  
 490 patches themselves. We use HTTP Range requests to down-  
 491 load header and index, ArrayBuffers to parse this structures into  
 492 JavaScript; the patches are then download prioritizing highest  
 493 screen error.

494 Remote rendering is possible due to the WebGL framework: in

495 this case, the implementation was integrated in 3DHop [? ],  
496 a set of tools for web publication based on the SpiderGL [? ]  
497 graphics library. The patches are encoded as binary files. Figure  
498 5 shows an example of a model before and after view-dependent  
499 refinement.

## 500 4. Results

501 The proposed method has been tested on several cases, in-  
502 cluding very complex geometries. Additionally, whenever possi-  
503 ble a comparison with existing systems was performed. A  
504 demo page, that shows the comparison and a few examples, is  
505 available at <http://fastdec.duckdns.org> (for reviewers only).

506 Our implementation has been successfully tested on major  
507 browsers on a variety of platform, from desktop machines to  
508 low end cell phones. The results we report here were mea-  
509 sured on an iCore5 3.1Gh, using Chrome 41. Timings on other  
510 browsers (e.g.Firefox) where comparable. Regarding the mul-  
511 tiresolution model construction, this is a preprocessing opera-  
512 tion. Compression time is negligible, since it can be performed  
513 at about 1M triangles per second. The most cumbersome part  
514 is the quadric simplification algorithm, that runs at about 60K  
515 triangles per second per core. Nevertheless, the model con-  
516 struction must be performed only once.

517 The results section is organized as follows: the first part  
518 is dedicated to a qualitative comparison with available com-  
519 mercial or free solutions for remote rendering; the second part  
520 shows the results of the tests made on our compression scheme;  
521 the third part comments the performances of the system on a  
522 variety of examples.

### 523 4.1. Comparisons with existing systems

524 Several solutions for the visualization of 3D data with We-  
525 bGL have been proposed in the last years. Most of them, un-  
526 fortunately, do not deal with any strategy about compression  
527 and progressive visualization. Hence, they are able to deal only  
528 with smaller 3D entities.

529 Table 1 shows a comparison of the available systems. Also  
530 the issue of scalability (the capacity to handle very big mod-  
531 els) is taken into account. The most similar and used systems  
532 are Potree, which is limited to point clouds, and Sketchfab, a  
533 widely used commercial tool, which uses a progressive visual-  
534 ization which is not view-dependent. Additionally, the scalabil-  
535 ity of the implementation of Sketchfab is unclear, and anyway  
536 the support for more complex models is currently limited not  
537 only for the "Basic" (50 Mb uncompressed size), but also for  
538 the "Pro" (200 Mb) and "Business" (500 Mb) accounts. Our  
539 method is able to handle models of any size (see later).

541 Hence, our method proves to be the most flexible one, since  
542 efficient solutions were developed to ensure not only perfor-  
543 mances, but also wide usability.

### 544 4.2. Entropy Compression: Comparison

545 Compression is an important issue when dealing with more  
546 complex models, and the need for remote rendering raises ad-

547 ditional issues. Hence, we tested, both in C++ and JavaScript,  
548 compression rates and decompression speed of:

- 549 • our implementation of Tunstall coding (T)
- 550 • Huffman coding (H), in the high-performance version of  
551 Pajarola [? ] (our implementation, C++ only)
- 552 • available implementations of LZMA  
553 in C++: <http://www.7-zip.org/sdk.html>  
554 and JavaScript: <https://code.google.com/p/js-lzma/>
- 555 • lz-string, a LZW based JavaScript implementation  
556 <http://pieroxy.net/blog/pages/lz-string/index.html>

557 The results are presented in Table 2, the lenght of 32K has  
558 been chosen since it is typical in our application.

559 Huffman and Tunstall are very similar in term of decom-  
560 pression speed, the difference is mainly in the time required to  
561 generate the decoding tables which are much larger for Huff-  
562 man, especially when increasing the number of symbols. We  
563 tested also other probability distributions and found little dif-  
564 ference in terms of speed. LZMA and LZW avoid this startup  
565 cost, however their more complex and adaptive dictionary man-  
566 agement allows them to outperform Huffman and Tunstall in  
567 term of decompression speed only for very small runs (and very  
568 small dictionaries). In terms of compression ratio, Huffman and  
569 LZMA performed quite close to the theoretical minimum, while  
570 Tunstall was about 10% worse.

571 We did not implement Huffman in JavaScript, as we are  
572 confident the result would be very similar. On the other hand  
573 the numbers for LZMA change dramatically. Lz-string serves  
574 as a comparison, as a better library, optimized for JavaScript.  
575 The poor LZMA performances in JavaScript help explain the  
576 relatively slow performances of CTM in Limper [? ].

### 577 4.3. Mesh Compression: Comparison

578 We used the Happy Buddha model (in Figure 1), to compare  
579 compression ratio and decompression speed with OpenCTM  
580 (CTM) [? ] Pop buffers (POP)[? ], P3DW [? ], WebGL-  
581 loader (CHUN) [? ]. We compare our multiresolution (OUR)  
582 and, to test single resolution performances of our compression  
583 approach, a version (FLAT) which loads only the highest res-  
584 olution level of the model. In each case the model has been  
585 quantized at 11 bit for coordinates and 8 bit for normals, and in-  
586 cludes colors. The apparently better performances by WebGL-  
587 loader (CHUN) are explained by the fast that it is not a multi-  
588 resolution method, hence the whole model has to be down-  
589 loaded before being able to decompress and visualize the model.  
590 This means that the compression rate of the highest resolution  
591 in our method (FLAT) is higher, and in any case WebGL-loader  
592 becomes unusable when more complex models are used.

593 Our decompression JavaScript implementation can decode  
594 about 1-3 million triangles per second with normals and colors  
595 in a single thread, on a desktop machine and 0.5 MT/s on a  
596 iPhone Five. Performances are somewhat degraded when the  
597 code is run during streaming visualization.

Table 1: Comparison of the main systems for web-based rendering.

	3D Meshes	3D Textured	Point Clouds	Streamable	Compressed	View-Dependent	Scalable
OpenCTM [? ]	Yes	Yes	<b>No</b>	<b>No</b>	Yes	<b>No</b>	<b>No</b>
WebGL-loader [? ]	Yes	Yes	<b>No</b>	Yes	Yes	<b>No</b>	<b>No</b>
Pop Buffers [? ]	Yes	Yes	<b>No</b>	Yes	Yes	<b>No</b>	<b>No</b>
Potree [? ]	<b>No</b>	<b>No</b>	Yes	Yes	Yes	Yes	Yes
SketchFab [? ]	Yes	Yes	Yes	Yes	Yes	<b>No</b>	<b>(Yes)</b>
Our Method	Yes	Yes	Yes	Yes	Yes	Yes	Yes

symbols	C++			JavaScript		
	T	H	LZMA	T	LZMA	LZW
4	1058	520	1066	201	19	55
9	369	212	170	145	10	23
13	423	168	95	150	6	20
17	359	136	77	163	6	19
22	332	98	67	180	6	17

Table 2: Decompression speed in million of output symbols per second for Poisson distribution of 32K sequences

	FLAT	OUR	CTM	CHUN	POP	P3DW
MB	1.9	3.9	3.5	2.8	15	4.5
bpv	28	57	51	41	220	66
full	0.4	0.9	5.3	0.06	0.5	10

Table 3: Statistics for the Happy Buddha: model size in megabytes, bit per vertex and time in seconds required to fully decompress the model.

598 An important comparison is with the work by Rodriguez  
599 [? ], which employs the same multiresolution batched strat-  
600 egy. For their mobile multiresolution application they report  
601 compression rates of 45-50 bpv on large colored meshes (which  
602 should be compared to our 28bpv). The difference is probably  
603 mostly due to the different connectivity encoding which, in their  
604 case, requires 20bpv against our 4 or 5bpv. It is difficult to com-  
605 pare the speed of the two decompression approaches since they  
606 run natively in C# on an iPhone4 while we run in JavaScript  
607 on the same platform. Our implementation speed is still, if a  
608 bit faster than their 50KTS<sup>1</sup>, at about 60KTs. The difference  
609 is probably due their more sophisticate (and slow) arithmetic  
610 encoding.

611 C++ decompression speed is of course faster, reaching 9MTs,  
612 including colors and normals, and 16MTs for just position and  
613 connectivity. The speed reported in [? ] of 35Kts for just the  
614 connectivity, as they mention, is due to the dynamic memory  
615 allocation in their implementation.

#### 616 4.4. Compression of Textured models

617 In the case of textured models, the compression and quan-  
618 tization has to be applied not only on the geometric attributes,  
619 but also on the texture image.

620 Regarding the latest, using 13 bits quantization for a 4096x4096

<sup>1</sup>The number is extrapolated from the decoding time of a large mesh given in their paper

621 pixel texture (which amount to half pixel precision) results in an  
622 hardly noticeable distortion (see Figure 6). With parallelogram  
623 prediction, texture coordinates are encoded in about 8 bpv (with  
624 a reduction of 48 to 1 respect to the standard 6 floats per face).  
625 We had to include the replicated vertices coordinates, that in  
626 our samples amounted to at most 15% of the total amount in  
627 models with many seams. Overall, adding texture coordinates  
628 increases the data size of about 25%, and decompression speed  
629 decreases accordingly.

630 We uploaded a few textured model to Sketchfab, for a com-  
631 parison with a state of the art industrial solution: our multires-  
632 olution structure results on average 10% smaller, although it  
633 includes all the resolutions.

#### 634 4.5. Streaming and Rendering

635 Loading the geometry through the Range HTTP request re-  
636 quires an increased number of HTTP calls: one for each patch,  
637 or 30-60 calls every million triangles. This does not really im-  
638 pact over performances: the overhead is quite small (about 400  
639 bytes per call) and pipelining (the process of enqueueing re-  
640 quests and responses between browser and server) ensures full  
641 utilization of the available bandwidth. Random access is really  
642 necessary only to fully exploit the view-dependent characteris-  
643 tics of the multiresolution structure: the code could be easily  
644 modified to load the model with a single call if a higher number  
645 of HTTP calls is problematic on certain web hosting architec-  
646 tures.

647 In the demo page (<http://fastdec.duckdns.org>) it is possible  
648 to compare the performances of our method w.r.t. existing solu-  
649 tions in the case of a slow connection. Moreover, very complex  
650 geometries are also available for further testing. In the follow-  
651 ing, we show some example of the performances of the com-  
652 pression method in several types of models, in order to test its  
653 flexibility.

##### 654 4.5.1. Point clouds

655 Point clouds are a quite common type of models, especially  
656 when large environments are taken into account. Terrestrial  
657 laser canners, but also UAV may provide dense point clouds.  
658 Ad-hoc solutions for encoding and rendering have been devised  
659 [? ], but their implementation is usually very hard to be ex-  
660 tended to triangulated surfaces.

661 On the contrary, the method proposed in this paper can be seam-  
662 lessly applied also in point clouds: the data that are compressed  
663 and streamed are only the vertices attributes.



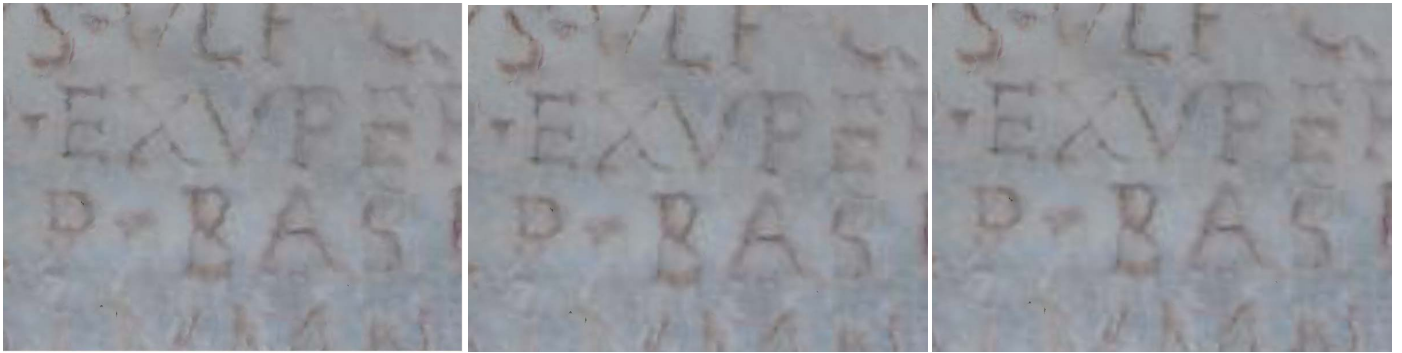


Figure 6: Example of textured model from left to right: texture coordinates quantized at 12 bits, at 13 bits and uncompressed. Notice the slight distortion when precision is lower than half a pixel.

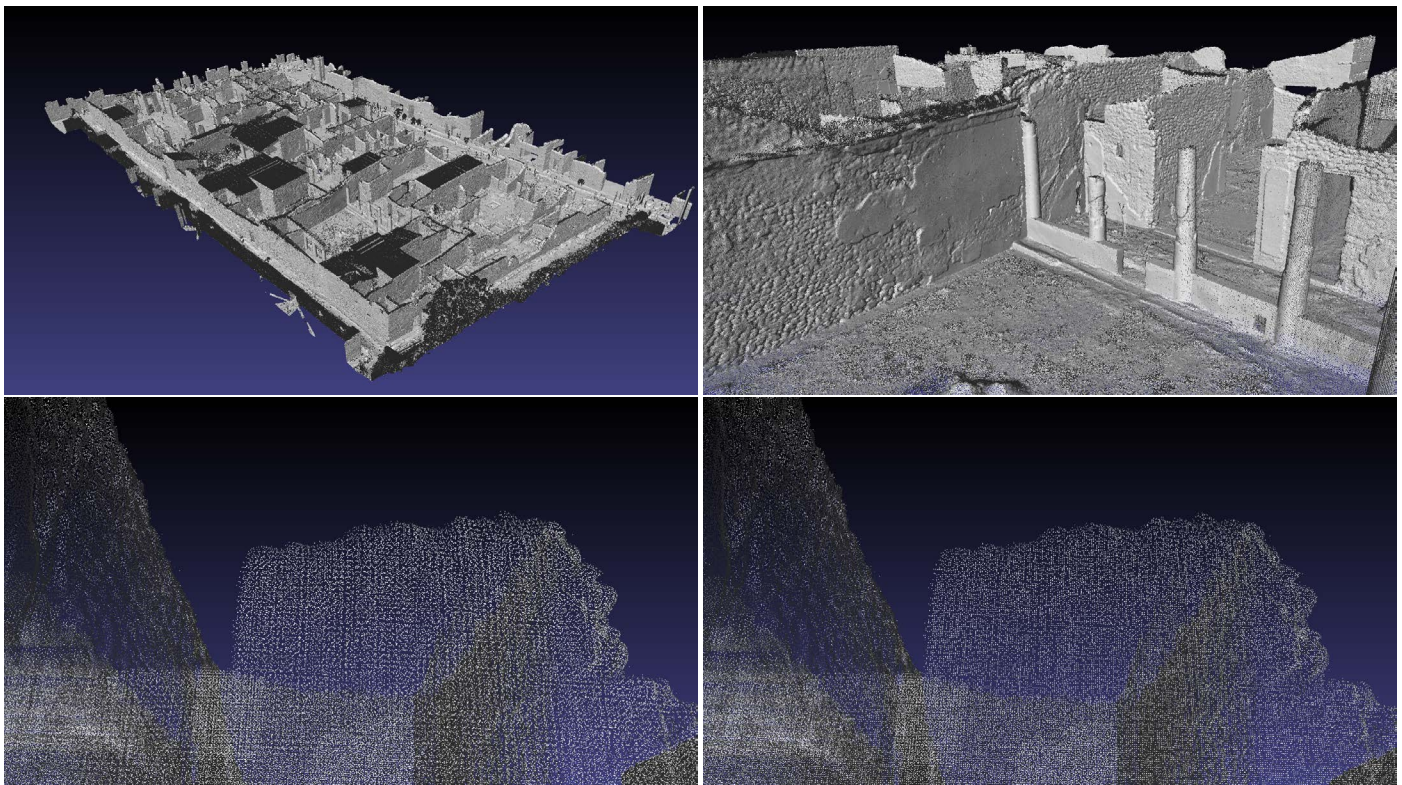


Figure 7: Pompeii point cloud: original PLY file, 95M points, 2.26 Gb; uncompressed multires cloud, 1.68 Gb; compressed cloud, 326 Mb. Top left: the full compressed model, top right: a detail. Bottom left: a detail of the uncompressed point cloud . Bottom right: the same detail of the compressed point cloud



Figure 8: Big statue rendered in a browser: original PLY file, 84M triangles, 1.6GB; uncompressed model, 2.54 Gb; compressed model, 158 Mb. Top left: the full compressed model, top right: a detail. Bottom left: a detail of the uncompressed model. Bottom right: the same detail of the compressed model

Figure 7 shows an example of a point cloud of Insula V of Pompei, obtained with terrestrial laser scanning. The top part shows the compressed model and a point of view where all the details are visible. The bottom part of the Figure shows the difference between the compressed and uncompressed point clouds: the quantization of the original data brings to a reduction of points, since the quantization tends to "regularize" the points grid (in this case the quantization step was 0.5 mm). This effect is not noticeable in triangulated surfaces, but it may be an issue for point clouds. In this case, it's possible to change the quantization step in order to find the best tradeoff between compression and data quality.

#### 4.5.2. Dense triangulated models

View-dependent progressive method have been especially devised to handle dense, triangulated 3D models. For this reason, the proposed method is able to provide optimal performances even when hundred millions triangles have to be taken into account. In the following, two examples of complex geometries are shown.

Figure 8 shows the 3D model of a 3-meter tall statue which was acquired with triangulation structured light scanner. The compressed model, which is nearly 10 % of the original PLY file, exhibits a detail that is undistinguishable from the uncompressed version.

In Figure 9 we show our system rendering the Portalada, a 180M triangles model at 30fps. The triangle budget has been

fixed at 1M triangles and the streaming requires 2-3 seconds to reach full resolution on a good connection. The original model is 3.6GB, while the compressed multiresolution model is 838MB. The Figure also shows how the view-dependent paradigm is able to handle different resolutions of different parts of the model when peculiar points of view are shown.

#### 4.5.3. Non-optimal, topologically complicated models

Some of the solutions proposed for progressive view dependent rendering proved to be limited since their basic assumptions on data processing didn't take into account that most of the more complex 3D models come from acquisition devices or techniques. This leads very often to the presence of geometric artifacts or unbalanced density.

Figure 10 shows two examples where the method deals with non-optimal geometries. On the left side, a model exhibiting strong topological artifacts. On the right side, a model with very unbalanced data density. In both cases, the method is able to deal with the issues and provide an accurate and reliable rendering.

## 5. Conclusions and possible improvements

The method proposed in this paper is a multi-resolution solution that provides remote view-dependent rendering of compressed 3D data. The main improvements w.r.t. current solutions are: the effectiveness in a wide range of bandwidth availability, computing power and rendering capabilities; the possibility to handle a wide variety of 3D models types, including very complex geometries; a mesh compression strategy that provides the best tradeoff between data compression and rendering performances.

Nevertheless, improvements in both compression and rendering performances can be obtained by further exploitation of the characteristics of some types of models.

For example, in the case of Point clouds, the rendering paradigm plays a key role to obtain a satisfying visualization. The current rendering method could be improved by implementing and extending existing approaches [? ]. The attributes (i.e. radius) that could be used for efficient rendering can be easily inserted in the compression framework.

The compression and rendering of textured models can be further improved, by working on ways to better compress and handle textures, or moving to other texturing paradigms. An example could be the projective textures (a similar approach on point clouds was recently proposed by Arikan [? ]), that could remove the need for parametrization, and open to even more complex datasets.

## Acknowledgements

The research leading to these results was funded by EU FP7 project ICT Harvest4D (<http://www.harvest4d.org/>), GA

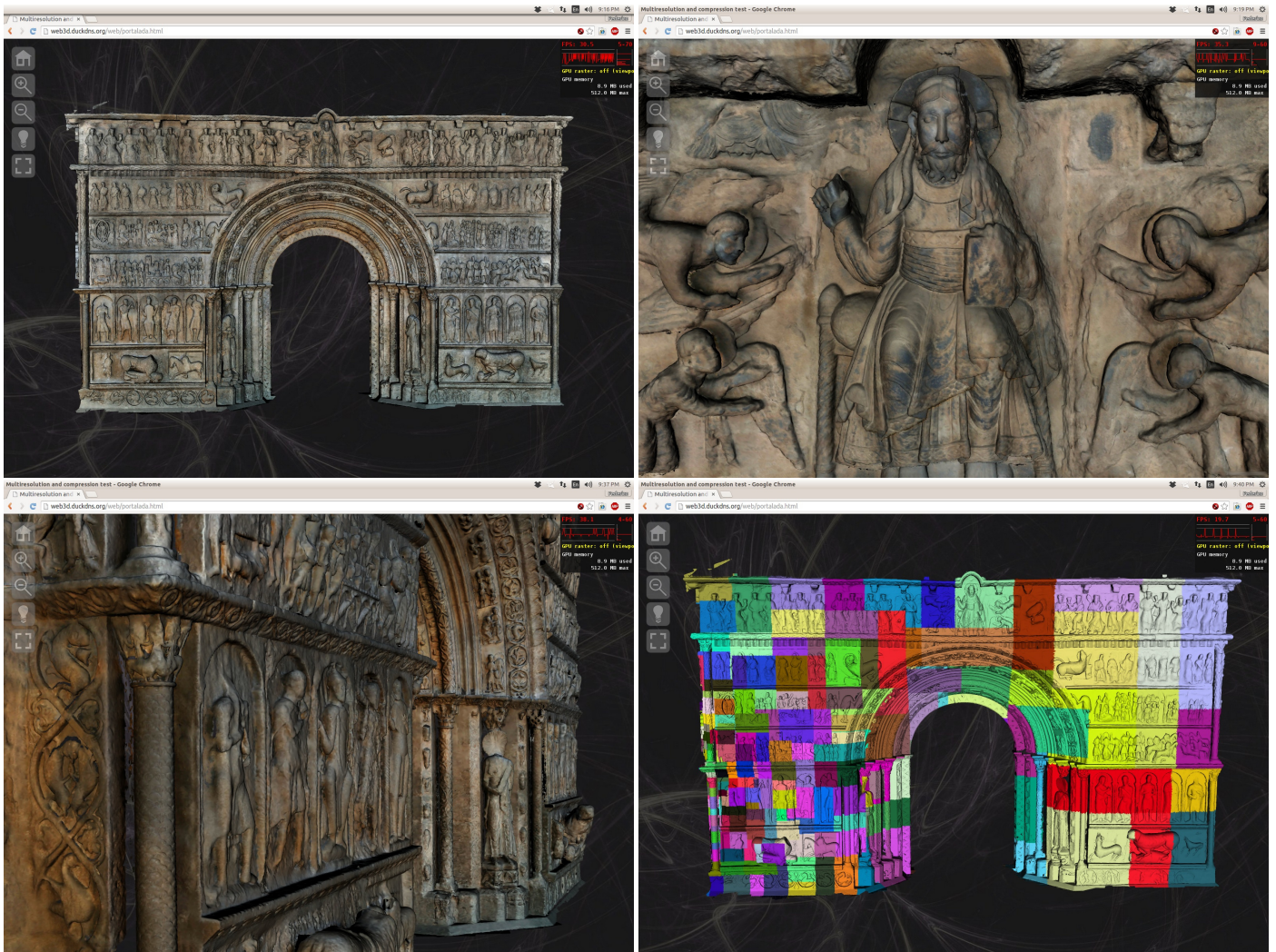


Figure 9: Portalada rendered in a browser: original PLY file, 180M triangles, GB; compressed model, 621 Mb. Top left: the full model, top right: a detail of the figure above the arch, bottom right: the resolution of the model as seen from the bottom left view point (without frustum culling)

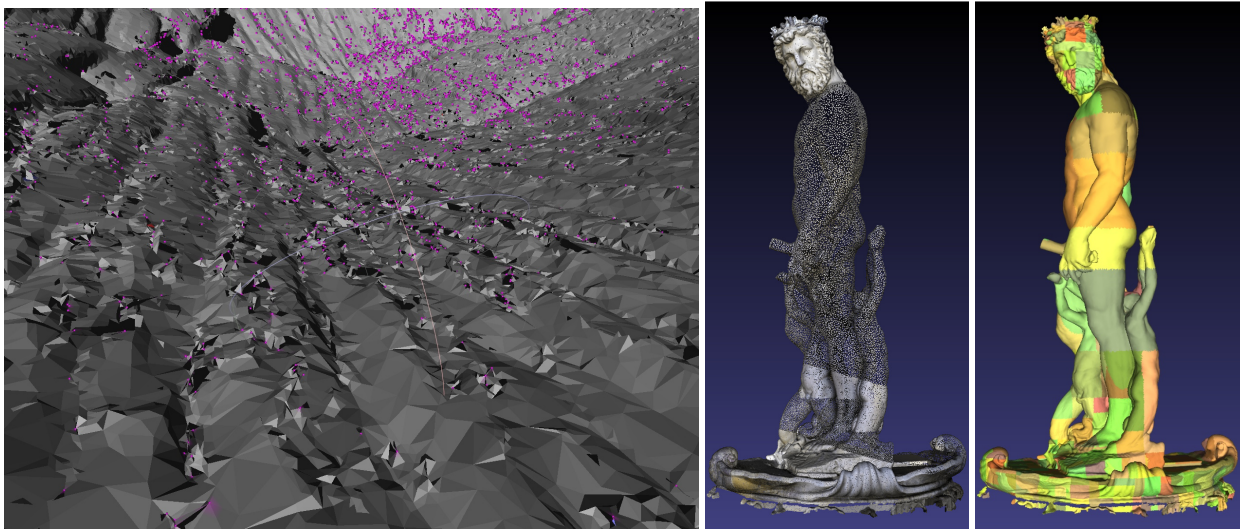


Figure 10: Left: a model with severe topological issues. Right: a model with very imbalanced vertex distribution

<sup>745</sup> n. 323567) and EU INFRA Project Ariadne (GA n. 313193,  
<sup>746</sup> <http://www.ariadne-infrastructure.eu/> ).