



Title	A Big Data Approach for 3D Building Extraction from Aerial Laser Scanning
Authors(s)	Aljumaily, Harith, Laefer, Debra F., Cuadra, Dolores
Publication date	2016-05
Publication information	Aljumaily, Harith, Debra F. Laefer, and Dolores Cuadra. "A Big Data Approach for 3D Building Extraction from Aerial Laser Scanning." American Society of Civil Engineers, May 2016. https://doi.org/10.1061/(ASCE)CP.1943-5487.0000524 .
Publisher	American Society of Civil Engineers
Item record/more information	http://hdl.handle.net/10197/7450
Publisher's version (DOI)	10.1061/(ASCE)CP.1943-5487.0000524

Downloaded 2026-05-01 23:38:05

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd_oa)



© Some rights reserved. For more information

17 **KEYWORDS:** Building Extraction, MapReduce, Big Data, LiDAR, Digital Surface Model,
18 Aerial Laser Scanning

19

20 **INTRODUCTION**

21 Massive streams of remote sensing data are rapidly becoming less expensive and more
22 readily available at a higher quality. However, processing them efficiently requires
23 appropriate tools, especially for querying and visualization. An example of where this is
24 highly relevant is in large-scale, urban, spatial datasets where data exists from many forms
25 and at various granularities [e.g. high resolution satellite, Light Detection and Ranging
26 (LiDAR), and distributed (wireless) sensor networks, etc.]. In this context, important
27 information can be extracted from LiDAR point clouds such as the location, orientation, and
28 size of objects, as well as possible damage (Walsh et al. 2013, Gao et al. 2014).
29 Unfortunately, traditional approaches are not well suited for the current generation of data
30 and will only become less so in future years. For example, a 500m*500m (or 0.25km²) tile of
31 LiDAR point cloud data could contain more than 65,500,000 three-dimensional (3D) points.
32 Making these points compatible with a Spatial Data Model (i.e., generating geo Identification
33 Key and Spatial Objects for each point) requires more than 10GB in relational database
34 storage. For this reason, writing queries against this huge data storage arrangement is
35 complicated and time consuming.

36

37 Since the increasing density and temporal recurrence of such data already challenges the
38 capabilities of traditional Geographic Information Systems and other spatial databases,
39 developing appropriate querying tools to significantly improve the outcomes of relevant data
40 retrieval exercises is quite critical. Furthermore, the complexity and nature of spatial data
41 makes them ideal datasets for applying parallel processing (Cary et al. 2009). Thus, large-

42 scale, urban, spatial data can be stored and analyzed in big data platforms, which enables
43 identification of changes in geographic patterns (NESSI 2012).

44

45 “Big data” is a high scalability and real time platform created to improve web-based search
46 engines (e.g. Yahoo, Google) and smart data analysis through parallel and distributed
47 computing. The approach can handle large and complex datasets that are beyond the means
48 of traditional databases. Despite their short history, these platforms have been deployed
49 extensively with further rapid adoption predicted (Özdemir et al. 2013; Floridi 2012).
50 Arguably, data mining in Digital Surface Models (DSMs) requires Big Data strategies to
51 overcome scalability issues (Batty et al. 2012). A DSM has many possible applications
52 including as an interim step in the transformation of point clouds into solid models
53 compatible for computational analysis (Truong-Hong et al. 2013). As such, this paper
54 endeavors to extend the viability of Big Data approaches to LiDAR data processing, in
55 particular for city-scale computational modeling, which is an area of significant interest to the
56 civil engineering community (Hinks et al. 2009). As the hardware limitations of working at
57 such a scale continue to decrease, the logistics of populating those systems for micro-climate
58 modelling, disaster prediction, and tunnel risk assessment continue to gain prominence. The
59 auto-identification and extraction of data affiliated with individual buildings is a critical step
60 in this process.

61

62 **RELATED WORK**

63 To date, few studies have been conducted regarding to the treatment of LiDAR in a Big Data
64 context. In related work, Wu et al. (2007) introduced an approach to automatically align a
65 combination of satellite images and vector data for roads. They divided the data into a set of
66 tiles and automatically computed the best imagery-to-vector translation within a tile by using

67 a parallel programming model with MapReduce; MapReduce is a software model used to
68 support parallel computing of huge sets of data and consists of two functions Map and
69 Reduce, which operate using key-value data types. Two years later Zhang, et al. (2009)
70 described how spatial queries could be adopted and expressed in a MapReduce model. By
71 using several spatial queries performance evaluations, they proved that MapReduce is
72 appropriate for small-scale clusters and computing intensive applications. The next year
73 Wang, et al. (2010) proposed methods to improve the performance of spatial computation
74 using a key-value based model. Their experiments showed that the MapReduce based spatial
75 system could significantly outperform a traditional Database Management System (DBMS).
76 Subsequently, Liu et al. (2012) used MapReduce to solve two common spatial problems. The
77 first was the bulk-construction of R-Trees (an indexing mechanism for spatial search query
78 processing). The second was the computation and storage of the quality characteristics of
79 aerial digital image using metadata to improve the display of image-based mosaicing. More
80 recently, Aji, et al. (2013) demonstrated the efficiency and scalability of Hadoop-GIS by
81 running queries on a large cluster. In this demonstration, they explained how spatial queries
82 could be translated into MapReduce operators, optimized, and executed in Hadoop.

83

84 Such approaches stand in strong contrast to the traditional methods developed for the storage
85 and investigation of LiDAR data, which have not considered the rapid trajectory of escalating
86 data density. Most approaches developed to date have been for datasets of no more than 10
87 points per square meter, which does not reflect the current system capabilities of not only 50
88 points per square meter but the existence of full-wave form data that has more than 250
89 pieces of information for each data point. Manipulation of such dense data sets will soon
90 become a requirement for a wide variety of data mining activities. One common task is

91 building extraction. A sampling of some of the notable, data-driven approaches to building
92 extraction from remote sensing data are described below.

93

94 In 2004, Rottensteiner et al. (2004) used the Dempster-Shafer theory for data fusion to the
95 detect buildings from a combination of LiDAR data and digital images. The approach used
96 three consecutive steps; (1) building detection, (2) roof plane detection, and (3) roof
97 boundary determination. In (Pradhan et al. 2007) considered the affiliated computing
98 requirements for city-scale modelling for disasters. These requirements are standardized data
99 specifications, middleware services and Web-enabled, distributed computing. In Mayunga et
100 al. (2005) proposed a semi-automatic process using a snake model and radial casting for
101 building data extraction from raster data. The approach required the manual identification of
102 the approximate centre of a building, from which the snake contour points were automatically
103 generated. Following this, Theng el al. (2006) employed the same snake contour model but
104 initialized it with a circular casting algorithm, instead of radial casting algorithm. In contrast,
105 San et al. (2010) applied the well- established Hough transform with it. Soon after, Dal Poz
106 and Galvanin (2011) extracted building outlines from a LiDAR-based DSM using a two-step
107 approach: (1) extract patches of aboveground objects and (2) identify edges that correspond
108 to building roof contours using a Markov-random-field-based energy function. Subsequently,
109 Sugihara and Kikata (2012) employed automatic generation of 3D building models through a
110 integrating geographic information systems (GIS) and computer graphics. Awrangjeb et al.
111 (2013) developed an approach to automatically extract 3D roofs by integrating LiDAR data
112 and multispectral orthoimagery. In that method, raw LiDAR points were separated into
113 ground points and non-ground points. The non-ground points were segmented using a
114 technique of image line guided segmentation to extract the roof planes. Li et al. (2013) also
115 proposed a method to fuse LiDAR data and optical imagery to extract different building

116 features. Their method consisted of four steps: 1) filtering, 2) building detection, 3) wall
117 point removal, and 4) roof patch detection. As part of step 2, an improved detector was used
118 to extract building edges from optical imagery. In the final step, the results from the prior
119 steps were integrated using a mathematical morphology.

120

121 Soon after Jochem et al. (2012) proposed a workflow that uses a combination of raster and
122 point cloud based GIS analysis. Although this workflow focused on roof planes, it can be
123 applied to other surface objects such as vegetation. In this workflow, in order to facilitate the
124 processing tasks, the large point cloud dataset is stored and divided into several smaller tiles
125 using a traditional database PostgreSQL/PostGIS solution.

126

127 Most recently, Abdullah et al. 2014 proposed a segmentation technique for automatic
128 building detection and roof plane extraction considering only non-ground LiDAR points, in
129 which finding neighbouring points involved segmenting LiDAR starting from the maximum
130 height and progressing downwards. The segmentation was done along individual lines at each
131 height level and depending on the distance between points. These lines were then used to
132 form planes. After all of planes were extracted, non-coplanar points (which represented trees
133 and other non-building structures) were removed by an ad hoc rule-based procedure.

134

135 While these various studies have produced important results for building extraction, they
136 generally employ either raster data or a combination of raster data and vector data.
137 Rottensteiner et al. (2004) noted that using raster data in building extraction introduces
138 problems related to shadows, occlusions, and poor contrast. In light of these criticisms, an
139 alternative approach is herein introduced that attempts to facilitate automated building
140 extraction task processing only vector data from aerial LiDAR.

141

142 **THE PROPOSED APPROACH**

143 The goal of this work is to introduce a new, fully automatic method for building extraction
144 from LIDAR data that needs no pre-processing. The proposed two-step approach involves
145 querying a LiDAR point cloud in the Hadoop platform, which is an open-source software
146 framework written in Java and is used to run distributed applications with huge amounts of
147 data among a Hadoop cluster, which consists of a set of compute nodes (computers). Step 1 is
148 a MapReduce process, in which neighboring points are mapped into cubes. Step 2 employs a
149 non-MapReduce process to extract objects where all of the adjacent cubes are considered to
150 belong to the same object. Although a LiDAR point cloud can include a variety of data, such
151 as 3D coordinates, timestamps, intensity, and RGB measurements, the proposed approach
152 only requires the 3D coordinates as part of the extraction process. Once the point cloud is
153 generated from the aerial laser scanning, it is used directly in the approach without any pre-
154 processing.

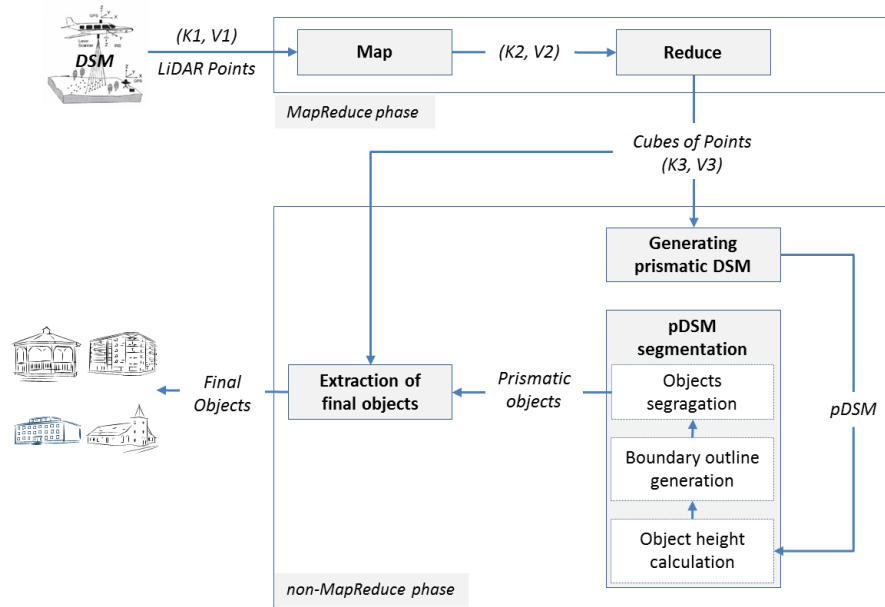
155

156 The basic idea exploits the distance between every pair of points in the point cloud. The main
157 problem is the computational expense. Applying the MapReduce function considerably
158 reduces the problem size, thereby making it more manageable. This is done by taking the
159 original DSM points generated from the point cloud and mapping them onto a 3D grid in the
160 form of a group of cubes. Notably, this approach is independent of data density or multiple
161 returns, however a higher density will allow a smaller cube size to be selected, which will in
162 turn improve the fineness of detail detection.

163

164 The MapReduce software framework is composed of two functions: the function 'Map'
165 processes the original data into key/value pairs, and the function 'Reduce' takes these pairs

166 and merges them in a way that all values corresponding to a specific key are combined into a
 167 single set. Once this is complete, the non-MapReduce function extracts neighbouring cubes,
 168 where each set of cubes represents a physical object. These steps (Figure 1) are explained in
 169 detail in the following sections.



170
 171 Figure 1. Approach design

172

173 **Mapping neighbour points (MapReduce)**

174 The nearest neighbour algorithm is a well-established and extensively used approach (e.g.
 175 Morgan and Tempfli 2000; Lee et al. 2008; Wang and Shan 2009) but when applied
 176 indiscriminately to a large dataset is computationally expensive (Athitsos et al. 2008). For
 177 this reason, the first part of the proposed approach aims to reduce the exploration
 178 neighbourhood. To do that, a distance threshold (D) is defined. In dense urban areas,
 179 determining D can be crucial, as many buildings abut one another. In this work, the distance
 180 was empirically selected as 0.5 m. Ideally this should be determined as a function of both the
 181 data density and typical building spacing.

182

183 Once D is defined, all the DSM points are segmented and mapped according to cubes with
184 dimensions $D \times D \times D$. Namely, if two neighbour points belong to a single cube then these
185 two points belong to the same object. The other way to describe this concept is to consider
186 the volume of the dataset as an octree of cubes of dimensions $D \times D \times D$ and to understand all
187 neighbouring points fitting in a single cube as belonging to only one object.

188

189 Consider a DSM consisting of a set of 3D points; $\text{Point}_1 (x_1, y_1, z_1)$, $\text{Point}_2 (x_2, y_2, z_2)$, ...,
190 $\text{Point}_n (x_n, y_n, z_n)$, where n is the number of points, and m is a set of identical cubes [Cube_1
191 (X_1, Y_1, Z_1) , $\text{Cube}_2 (X_2, Y_2, Z_2)$, ..., $\text{Cube}_m (X_m, Y_m, Z_m)$]. The mapping of a point to the
192 corresponding cube is done according to the following rule:

193 $\text{Point}_i (x_i, y_i, z_i) \in \text{Cube}_j (X_j, Y_j, Z_j)$ where:

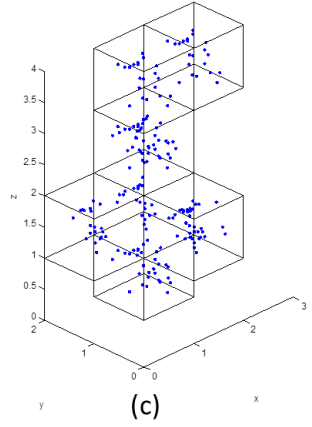
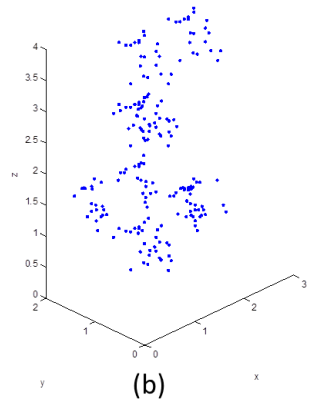
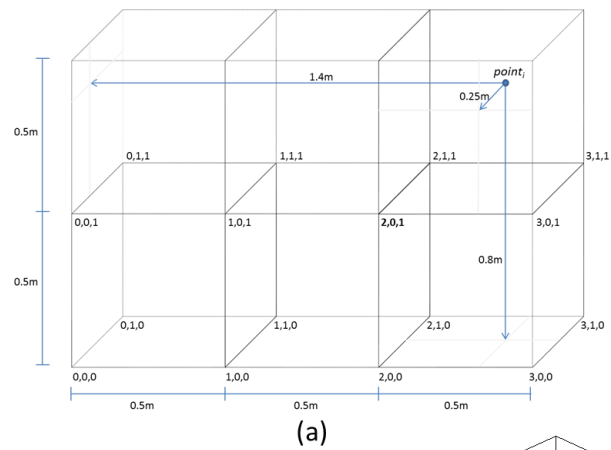
194
$$X_j = \text{fix}(x_i/D)$$

195
$$Y_j = \text{fix}(y_i/D)$$

196
$$Z_j = \text{fix}(z_i/D)$$

197 The function $\text{fix}(v)$ truncates the value to the greatest integer less than or equal to v.

198 For example, for a given point $i (x_i=1.4 \text{ m}, y_i=0.25 \text{ m}, z_i=0.8 \text{ m})$ and $D = 0.5 \text{ m}$, then this point
199 is mapped to the cube $j (X_i=2, Y_i=0, Z_i=1)$ see figure 2(a). Applying this rule to the set of
200 points shown in figure 2(b) results in figure 2(c).



201

202 Figure 2. 3D point mapping

203

204 Since the Mapper and the Reducer classes in Hadoop have the following general form (White
 205 2012):

206 *Mapper* <K1, V1, K2, V2>

207 *Reducer* <K2, V2, K3, V3>,

208 then the Mapper class receives a pair in the form of (K1, V1) and issues a list in the form of
 209 (K2, V2). The Mapper key K1 is the file name, while the dataset may be a single file or a set
 210 of files in a given directory. The value of Mapper V1 is the content of the line in the
 211 corresponding file. Now, only the 3D point coordinates are employed. The Mapper class
 212 issues a list of pairs of (K2, V2), where K2 represents the cube, which is calculated according
 213 the previous rule, while V2 represents the current point being processed, as follows:

214 *map* (K1, V1, context){

```
215     K2 = Cube(fix(V1.x/D), fix(V1.y/D), fix(V1.z/D));
```

```
216     V2 = Point(V1.x, V1.y, V1.z)
```

```
217     context.write(K2, V2);
```

```
218 }
```

219 Next, the Reducer receives the list of the pairs (K2, V2) and issues a list of the pairs (K3,
220 V3), where K3 is the cube coordinate, as a unique key in the list, while V3 is a list of all
221 points that belong to K3. The final result of this step will be the list (K3, V3), where all the
222 cubes K3 in the list are sorted in descending order according to the height coordinate of the
223 cubes (i.e., the first cube in the list will be the highest cube in the highest object H_i). The
224 Reduce function is shown as follows:

```
225 reduce (K2, V2, context) {
```

```
226     K3 = K2;
```

```
227     for (v : V2) V3 += v;
```

```
228     context.write(K3, V3);
```

```
229 }
```

```
230
```

231 **Objects extraction (non-MapReduce)**

232 Although the two steps of the approach are executed in independent processes, this second
233 step is executed after step 1 is complete. The main objective is to remove trees and other
234 obstructions and then to disjoin the results of the previous step (K3, V3) into objects such as
235 buildings, roads, green areas, etc. The basic idea is to find the neighbour points of an object
236 by visiting cubes instead of points. This is done by applying the following steps:

```
237
```

```
238
```

```
239
```

240 *Generating prismatic DSM to remove trees and obstructive objects*

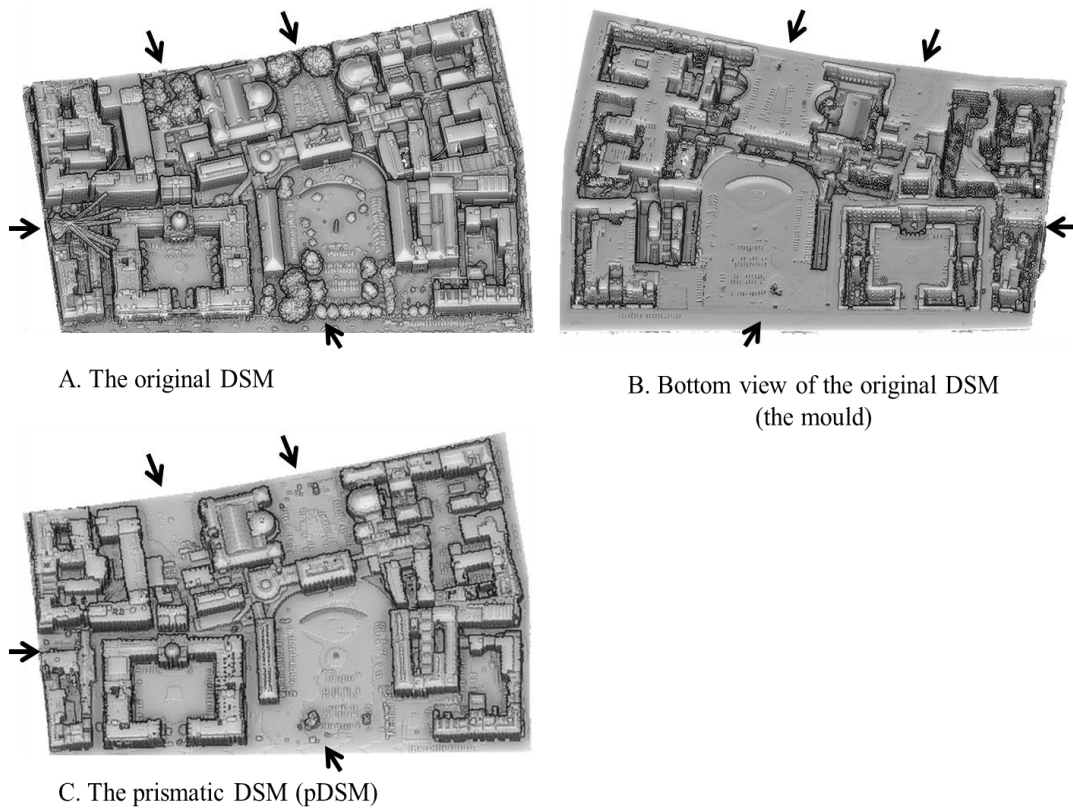
241 To remove trees and other undesirable objects, a prismatic DSM (pDSM) is generated from
242 the original DSM being processed. The pDSM is a totally filled 3D city model. In
243 Rottensteiner, et al. (2002), the prismatic model is derived from the boundary polygons of the
244 building regions using the average building heights. The accuracy of this height averaging
245 was in the range of ± 5 m. In the approach presented herein, generating a pDSM is done by
246 filling the DSM being processed by new cubes with the same side length $D * D * D$. This is the
247 same process as filling a mould with any material. Figure 3 illustrates how this process is
248 performed in the presence of a high density of trees adjacent to buildings. In addition, in the
249 bottom left corner of figure 3a, there is a tower crane. This would continue to obstruct the
250 building extraction process, if a pixel-based approach is used. In contrast, by viewing the
251 underside or bottom view of the same DSM (figure 3b) the presence of each building appears
252 as a deep hole, the depth of which corresponds to the respective building height. If multiple
253 buildings are connected, this bottom view displays these connected holes. A careful
254 examination of figure 3b shows small, insignificant holes distributed across the DSM. These
255 holes indicate the presences of trees, a tower crane, cars, etc. This phenomenon is produced
256 because the laser scanner gives very rich surface details of all existing objects, but does not
257 penetrate the surface. A “hole” represents the perimeter and volume of a building. Trees and
258 other such objects do not appear in the same way, since the area of their connection to the
259 ground is insignificant.

260

261 The result of generating the pDSM where a simple rule is used is shown in figure 3c. This
262 rule states that for each $\text{NotEmpty}(\text{Cube}(X, Y, Z)) \in \text{DSM}$, then the set of cubes $(\text{Cube}_0(X,$
263 $Y, 0 * D), \text{Cube}_1(X, Y, 1 * D), \dots, \text{Cube}_n(X, Y, Z - D)) \in \text{pDSM}$.

264

265 This algorithm starts from the first cube (0, 0, 0) in the original DSM, and then checks if the
266 cube of the original DSM being processed is not empty. Once a non-empty cube is located,
267 then the respective portion of the pDSM is filled by a set of cubes starting at 0 and finishing
268 at Z-D. The X and Y coordinates are fixed. This is then repeated. For example, if a cube
269 (0,0,3) is not empty in the original DSM, then the set of cubes (Cube₀(0,0,0), Cube₁(0,0,1),
270 Cube₂(0,0,2)) are added to the pDSM being processed, where 0,1,2,3 are an indexing array.



271
272 Figure 3. Removing trees and other obstruction objects.

273

274 *pDSM segmentation*

275 Once the trees and other small objects are removed in the processing of generating the
276 pDSM, then the aim is to disjoin the pDSM into several object models that will be used as
277 references to extract the final objects.

278 In order to segment a pDSM into separate object models, consider the $Cube_1 (X_1, Y_1, Z_1)$ as
279 the first/highest cube of an object model Obj_1 . Then all the neighbouring cubes that are
280 visited starting from this cube will be mapped onto Obj_1 .

281 As a rule for a given starting cube $Cube_i(X_i, Y_i, Z_i)$, there exists an adjacent cube $Cube_j(X_j,$
282 $Y_j, Z_j)$ where the value of X_j ranges from X_i-1 to X_i+1 , the value of Y_j ranges from Y_i-1 to
283 Y_i+1 , and the value of Z_j ranges from Z_i-1 to Z_i+1 .

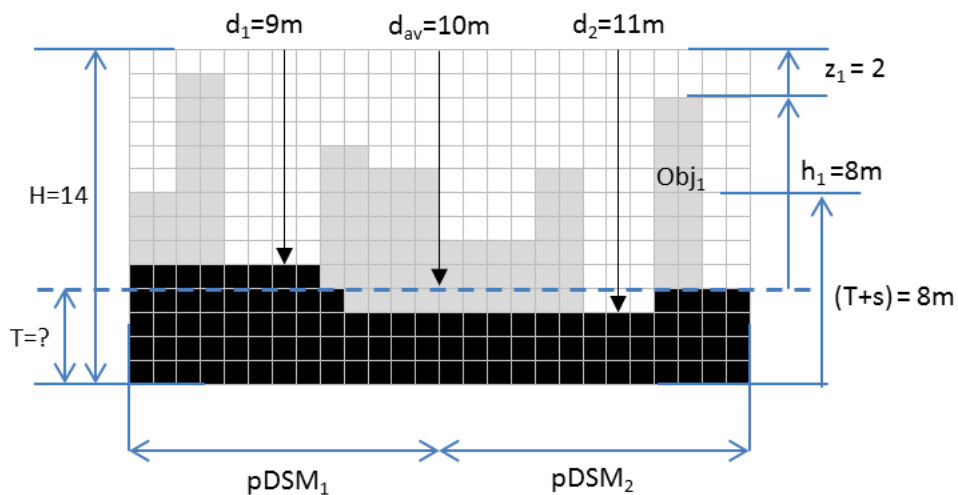
284 This rule is repeated for each neighbouring cube, until no more neighbouring cubes are
285 found. One cube can have as many as 26 adjacent cubes or as few as none. If there is no
286 adjacent cube, the cube being processed is considered to be a boundary cube. To preclude
287 revisitation of a cube, each cube is removed from the candidate dataset, once it has been
288 visited.

289 While applying the previous rule extracts all prismatic model details (irrespective of
290 complexity), a small problem persists. Near to the ground level, the algorithm will expand to
291 include all the ground cubes because of their adjacency. To overcome this problem, an
292 object's outline is derived during segmenting the pDSM. In this case, the outline should be
293 derived slightly higher than the ground level. So, in order to derive an object's outline
294 correctly, the height of this object should be calculated.

295 Significant research has been done to address the issue of the heights of objects. For example,
296 in Zhang et al. (2006) building heights were derived by averaging the elevation differences
297 between building measurements and the digital terrain model (DTM). In the same direction,
298 Abdullah et al. (2014) used a height threshold, based on a DTM to divide the LiDAR point
299 cloud into ground and non-ground points. In the work presented herein, the main problem is
300 that object heights need to be calculated "on the fly" without using any exterior tool or
301 relying on any user interaction.

302 In order to calculate the height of an object, the rule ($h_i = H - T - z_i$) is applied, where h_i is
 303 the height of the object being processed, H is the total height of the DSM, T is the ground
 304 thickness, and z_i is the Z coordinate of the highest cube of the object.

305 An example of the application of this rule is shown in figure 4. The pDSM is mapped to
 306 identical cubes. The height of the pDSM is H . This value H can be calculated by applying the
 307 formula $H = Z_{max} - Z_{min}$, where Z_{max} is the coordinate of the highest cube in the pDSM,
 308 and Z_{min} is the coordinate of the lowest cube of the same pDSM. In the example, $H=14m$.
 309 However, this height does not represent the height of the objects, because it includes a part of
 310 the ground cubes. For this reason, the main objective here is to calculate an approximate
 311 value of T , which represents the thickness of the ground near the objects being processed and
 312 finally to find the approximate height of each object. To calculate T , the main pDSM is
 313 divided into multiple small models (pDSM₁, pDSM₂...) depending on the terrain type. This is
 314 to say, if the slope of the natural ground level is very high, then the pDSM being processed is
 315 divided into smaller areas. In this example, only two models (pDSM₁, pDSM₂) are used. For
 316 each pDSM_i the maximum distance d_i between the ground and the maximum height of the
 317 main pDSM is calculated. In the example, the maximum distance of the pDSM₁ is $d_1 = 9m$
 318 and in pDSM₂ is $d_2 = 11m$. So, the average of the distances in all the pDSM is $d_{av} = (d_1 +$
 319 $d_2)/2 = 10m$. Thus, T will be $H - d_{av} = (14 - 10) = 4m$.



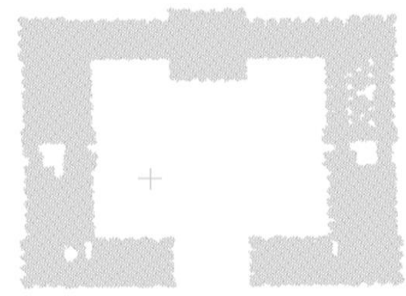
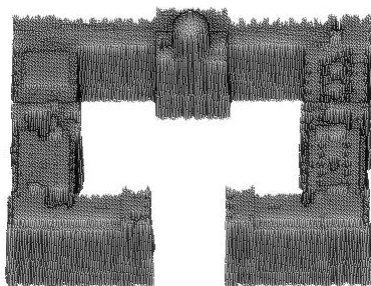
321

322 Figure 4. Object model height calculation.

323

324 Once the approximate thickness T is determined for the ground near the objects being
325 processed, then boundary outline calculation begins. The most important issue in this process
326 is to establish approximately the locations at which the outline should be calculated. The
327 height should be measured slightly above the natural ground level to prevent the algorithm
328 expanding to include adjacent ground cubes, but the height should be as low as possible to
329 retain all architecture details. For the objects of the DSM of figure 4, the outline of obj_i
330 is calculated as the height $(T + s)$, where s is a small offset from the ground that is automatically
331 calculated based on the height of the object. Through experimentation of a dense urban
332 dataset (see validation section), a distance of $s = 0.5 * h_i$ has provided good results. In the
333 Figure 4 example, for obj_1 the height of the outline calculation can be $(T + s) = (4 + 0.5*(14-$
334 $4-2)) = 8m$.

335 An outline of an object model is generated by using cube projection onto a two-dimensional
336 (2D) plane, where the projection is the vertical shadow of a cube, and the plane represents the
337 ground surface. All cubes found in the level h are projected onto an empty plane. Each cube
338 is projected onto the plane as a 2D point. That is to say, the projection of the Cube(X, Y, Z)
339 of a building creates the outline point Opoint(X, Y) onto the corresponding plane. The
340 resulting projection is the outline of the corresponding building in plan view (see figure 5).



A. an object model from pDSM

B. Boundary outline of the object in A.

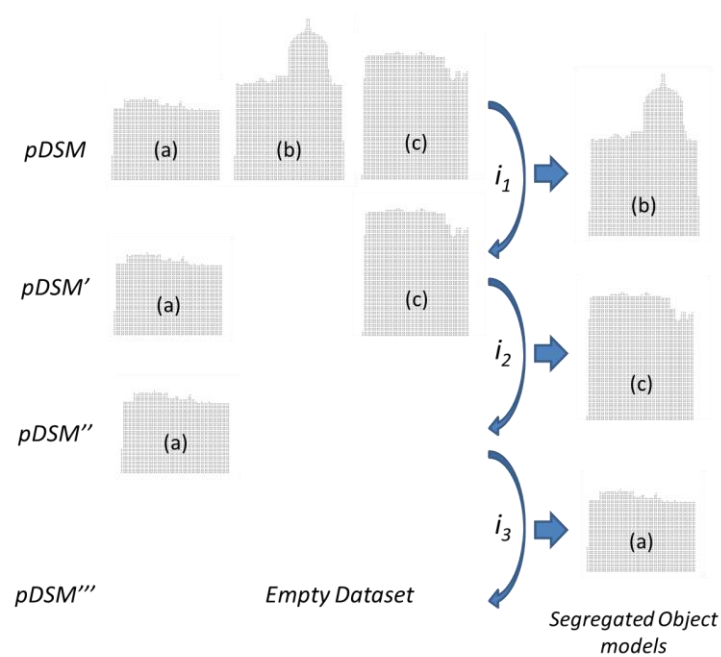
341 Figure 5. Two-dimensional boundary outline.

342

343 Once the approximate boundary outline is calculated, then the neighbour visiting process in
344 the horizontal direction is controlled. Namely, when the algorithm is applied at the ground
345 level, it extracts only the cubes whose coordinates match with the coordinates of the
346 corresponding outline points. As a rule, if a building outline is created from the set of the
347 points $\{Opoint(x_1, y_1), Opoint(x_2, y_2), \dots, Opoint(x_n, y_n)\}$, then any Cube (X_i, Y_i, Z_i) being
348 processed belongs to this building, if and only if, the pair (X_i, Y_i) matches any pair (x_j, y_j) in
349 the corresponding outline. In this case, the algorithm allows the extraction of Cube $(X_1, Y_1,$
350 $Z_1)$, Cube (X_2, Y_2, Z_2) , \dots , Cube (X_m, Y_m, Z_m) , where the pairs $(X_1, Y_1), (X_2, Y_2), \dots, (X_m,$
351 $Y_m)$ match the corresponding pairs in the outline, irrespective of the values Z_1, Z_2, \dots, Z_m .

352 Notably, the algorithm starts from the highest cube in the dataset and moves downwards
353 towards the lowest one. This is done, thanks to the first phase of the approach (MapReduce
354 phase), where the cubes of the points are generated in a decreasing order along the Z-axis. So,
355 the first cube in the generated dataset is the highest cube in the highest object of the pDSM
356 being processed. The algorithm repeats, until no more cubes are found in the pDSM. At the
357 beginning, it creates an empty object (*obj*) to save all the cubes related to the highest object of
358 the pDSM. The highest cube $cube(X, Y, Z)$ is first segregated from the pDSM and then added
359 to *obj*. Then all the neighbour cubes related to the first cube are added to *obj*. At the same
360 time, these cubes are segregated from the same pDSM dataset. Once the highest object is
361 segregated, the algorithm next segregates the highest remaining object in the pDSM. This
362 continues until all objects are segregated from the pDSM. At the end of the extraction
363 process, the pDSM dataset being processed will be empty, because all of its cubes will have
364 been segregated and moved automatically to the corresponding files. The example in figure

365 (6) shows the segregation process, where the pDSM contains three objects (*a*), (*b*), and (*c*). In
 366 the first loop (i_1) of the algorithm the object (*b*) is segregated automatically from pDSM
 367 (because it is the highest object in the set), and the remaining set becomes pDSM'. In the
 368 second and the third loops (i_2 & i_3) object (*c*) and object (*a*) are segregated, respectively. If
 369 the pDSM is empty (such as pDSM'''), then the algorithm stops.



370

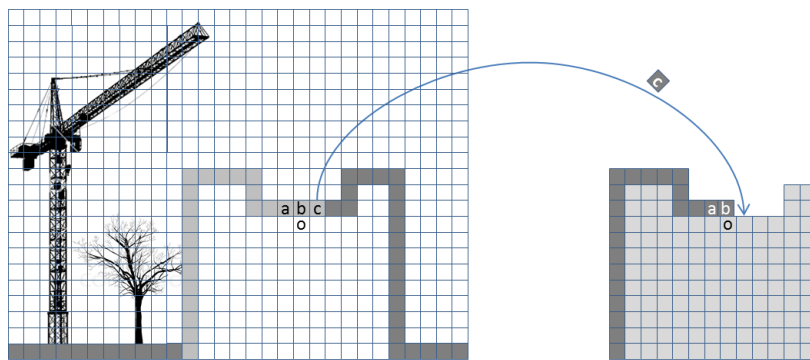
371 Figure 6. The automatic segregation process of the prismatic object models

372

373 ***Extraction of final objects***

374 The main objective is to extract the final object from the original DSM. This requires using
 375 the original DSM and the prismatic object model corresponding to that object. Figure 7.A
 376 shows a DSM of a building object (gray cubes) with adjacent trees and a tower crane. When
 377 this DSM is subjected to the first phase (MapReduce phase), the DSM will be segmented into
 378 a set of cubes of points. When these cubes are subjected to the tree removal phase, a
 379 prismatic object model for the same building will be generated (see figure 7.B, the light gray
 380 cubes). In this step, all the cubes that have adjacent objects into the prismatic object model
 381 are segregated from the original DSM. For example, the cube (o) in the original DSM (figure

382 7.A) is the same as the cube (o) in the generated object model [i.e., it has the same
 383 coordinates (X, Y, Z)]. So in this step, all the objects (a, b, c) of the original cubes will be
 384 moved to the cubes (a, b, c) in the new and final object, because of the adjacent
 385 neighbourhood rule. If the adjacent cube is null in the original DSM, then this cube will not
 386 be moved. Notably, this algorithm moves only the immediate adjacent neighbourhood, and it
 387 moves only the cubes that have adjacent neighbourhoods. For this reason, the trees and tower
 388 crane are not moved to the final object, although the tree is adjacent to the cubes that have
 389 been moved.



A. a DSM of a building

B. The prismatic model of the building

390

391 Figure 7. The automatic extraction process of the final objects

392

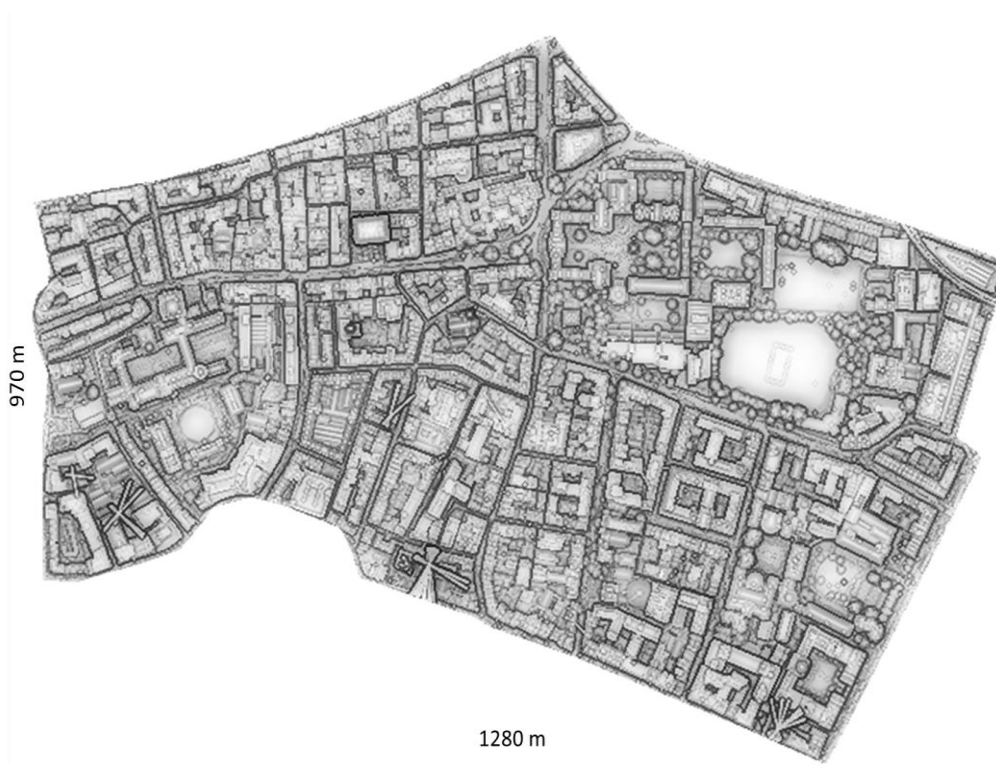
393 As a rule, if a cube_i(X_i, Y_i, Z_i) ∈ a Prismatic Object, then there exists an adjacent cube
 394 Cube_j(X_j, Y_j, Z_j) where the value of X_j ranges from X_i-1 to X_i+1, the value of Y_j ranges from
 395 Y_i-1 to Y_i+1, and the value of Z_j ranges from Z_i-1 to Z_i+1. This uses the same rule as in
 396 section (pDSM segmentation), but with the additional restriction that it is not allowed to grow
 397 beyond the cube being processed.

398

399 RESULTS, EVALUATION, AND DISCUSSION

400 This section evaluates the proposed algorithm for building extraction from an architecturally
 401 dense and complex portion of Dublin, Ireland. Within a 1 km² study area there are 9 tiles

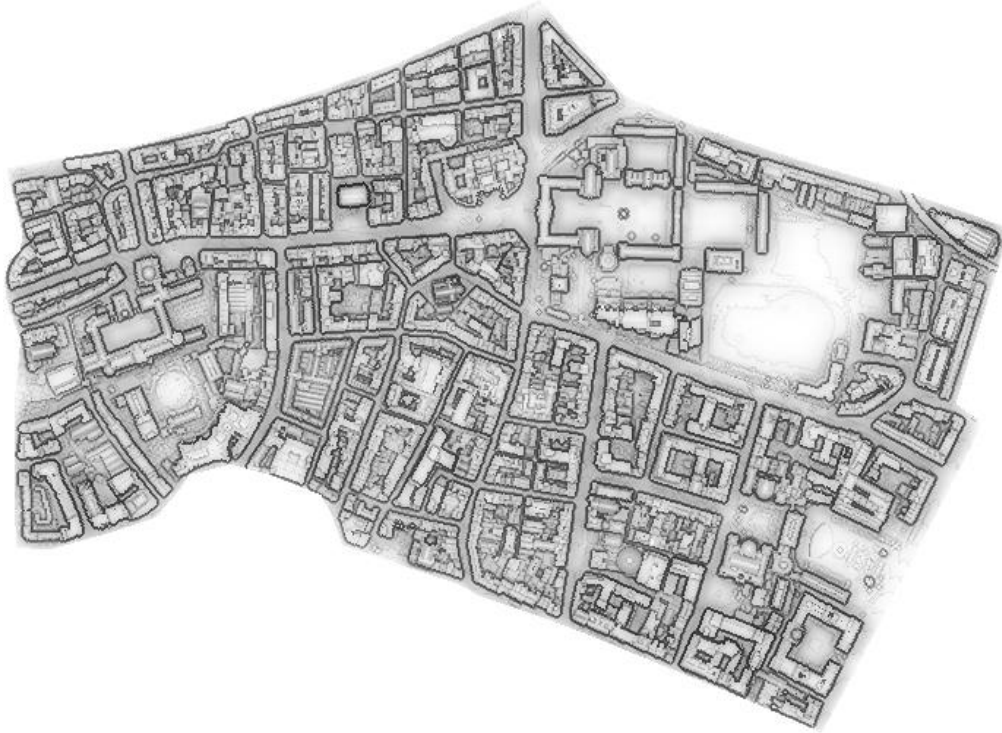
402 from a DSM (see figure 8). Each tile was saved into a separate file. The preparation of this
403 DSM included several steps that are outside of the scope of this work, including flight path
404 planning, as well as data collection, registration, and filtering (as described in Truong-Hong,
405 2011). The DSM contained 226,319,306 points. Figure 8 shows the corresponding pDSM,
406 where all the trees and other obstructions present in the initial DSM have been removed.



407

408 Figure 8. Initial DSM

409



410

411 Figure 9. Newly generated pDSM from the initial DSM (note the absence of all trees

412

413 Figure 10 shows the set of outputs of the algorithm. Notably, the algorithm converts the DSM

414 being processed into objects stored in separate files. However, in this evaluation study only

415 the 106 objects that form the most significant objects in the study area were processed. These

416 extracted objects are shown in figure 10 in a single image to ease comparison with the DSM

417 (Figure 8) and the corresponding pDSM (Figure 9).



418

419 Figure 10. Composite output of the extracted buildings from the algorithm

420

421 Figure 10 shows that most of the 106 objects have been isolated successfully. To
422 quantitatively evaluate the algorithm's outputs, measurements of correctness, completeness
423 and fitness measure (F-measure) were used. Normally, these measures are calculated by
424 taking the difference between the extracted buildings and the reference buildings (Maurya et
425 al. 2012). Several different input parameters have been proposed. Ekhtari et al. (2008) used
426 the building pixels, while Song and Haithcoat (2005) determined the volumetric difference.
427 For the proposed approach, a comparison between points in the extracted buildings versus
428 points in the reference buildings was considered.

429

430 To obtain the reference buildings, manual extraction was done for each of the 106 buildings
431 based on the fact that most features within the study area (e.g. buildings, roads, trees, cars)
432 were easily distinguishable with the naked eye. This was done using the visualization tool

433 CloudCompare (Compare, 2014). This software provides editing features such as DSM
434 segmentation. For each reference object, a contour was defined by a successive series of
435 clicks and then saved. The points inside this contour were then saved as a separate reference
436 object.

437

438 As such, correctness (which evaluates the exactness of an approach) was defined as the ratio
439 of the relevant points of a specific building to the total number of points of that building (see
440 eqn 1). A point was considered relevant when the algorithm extracted it correctly, with
441 respect to the corresponding reference building. Herein completeness, which measures the
442 ability of the approach to extract the entire set of points relevant for a building (i.e. coverage),
443 was defined as the ratio of the extracted relevant points to the total number of points in the
444 buildings of the study area (see eqn 2). The F-measure as defined by Peukert (2012), which
445 evaluates the overall quality (sometime called fitness), was calculated based on the
446 correctness and completeness metrics, as shown in (eqn 3).

447 $correctness = \frac{TP}{TP + FP} \dots\dots\dots(eqn 1)$

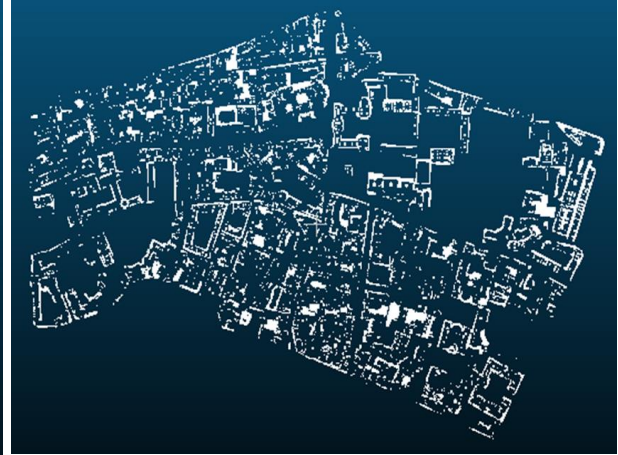
448 $completeness = \frac{TP}{TP + FN} \dots\dots\dots(eqn 2)$

449 $Fmeasure = \frac{2 * (correctness * completeness)}{(correctness + completeness)} \dots (eqn 3)$

450 where, True Positives (TP) counts the points correctly included into this object, False
451 Positives (FP) counts the points incorrectly included into this object, False Negatives (FN)
452 counts the points mistakenly excluded for this object (see figure 11).



TP = 114720570 points



FP = 8349057 points



FN = 10009794 points

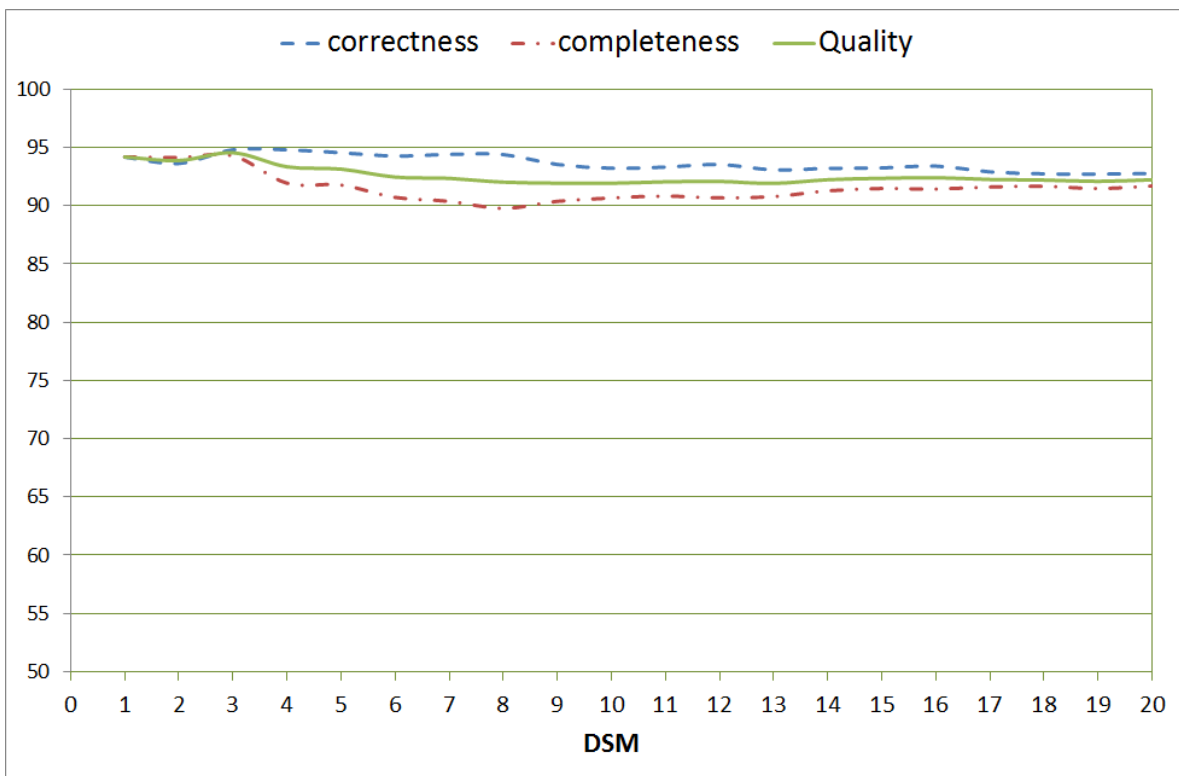
453 Figure 11. The extraction result of the study DSM.

454

455 Looking at figure 11, one can conclude that most of the FP and FN points are found around
456 the outlines of the objects, which means that the manual extraction of the reference buildings
457 did not coincide exactly with the automatic way that was calculated by the Big Data
458 approach. This is in part an outgrowth of the difficulty of manual extraction where errors may
459 be introduced. Notably, more than 12 hours were needed to extract all the reference objects
460 (approximately 7 minutes per object) versus less than 17 seconds per object using the
461 algorithm.

462

463 In a previous study by Hinks (2011), which was performed on the same study DSM using a
 464 morphological contour of the scan line data with the same data, the final values were
 465 consistently not as good: correctness = 87%, completeness = 82%, and the overall quality =
 466 84%. While, applying the Big Data approach, the evaluation measures across the 106 objects
 467 resulted in the following: correctness = 92%, completeness = 90%, and the overall quality =
 468 91%. These measures were calculated based on lump sum evaluation of the points of all of
 469 the objects, however in order to show how these measures varied, the original DSM was
 470 divided into 20 smaller DSMs each consisting of approximately 4 to 8 adjacent objects. The
 471 process was then reapplied. The results are shown in Figure 12, in which DSM #3 obtained
 472 the best extraction quality (correctness = 94%, completeness = 94%, and overall quality =
 473 94%), while the DSM #8 has obtained the worst result (correctness = 94%, completeness =
 474 89%, and the overall quality = 92%). The problem of DSM #8 is related to completeness,
 475 because it consists of multiples and complex objects (see Figure 13).



476 Figure 12. Detail quality study.
 477

478



DSM #3

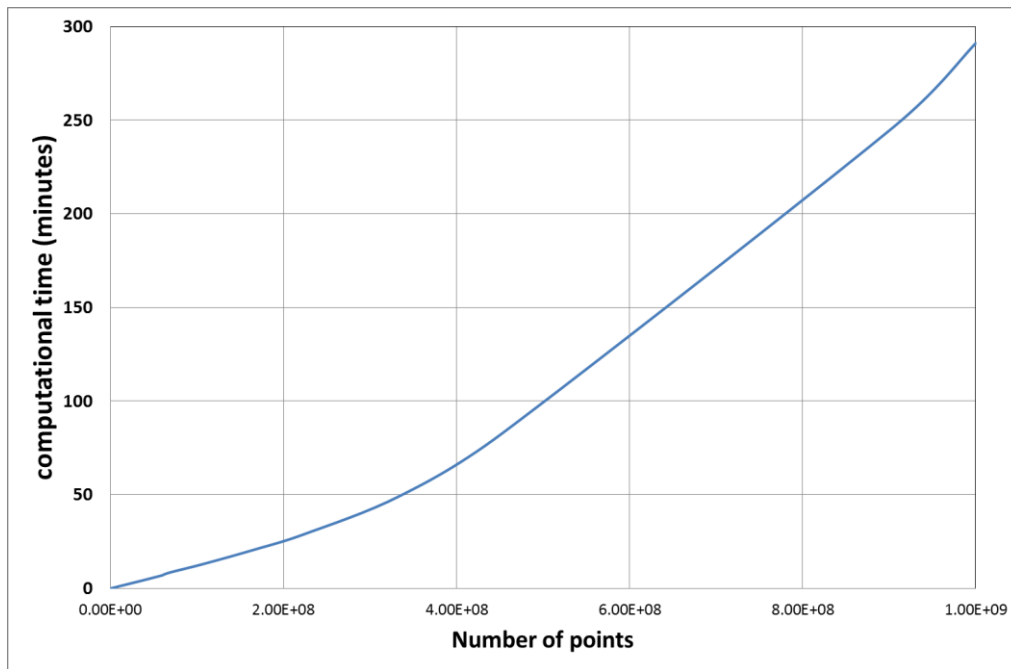
DSM #8

479 Figure 13. Sample DSMs within the study area

480

481 Although the experiment was performed in a single Hadoop installation, having a broader
482 vision about the computational efficiency of the approach, where the efficiency evaluation
483 refers to the execution time required to run the algorithm (Eldh 2006) is useful. In order to
484 check the scalability of the proposed approach with a considerably larger number of points,
485 an experiment was done by expanding the study area from the original 1 km² to a surrounding
486 area of more than 4 km². This expanded dataset contained approximately 1 billion points. The
487 execution time needed for the MapReduce phase is shown in Figure 14. For example, to
488 segment approximately 226,319,306 points (1Km²) the MapReduce step needed
489 approximately 29.4 minutes. Arguably this time can be improved easily by adding more
490 nodes to the cluster since having more nodes or computers in a cluster with data spread over
491 the cluster can significantly improve the start to end processing time (Xu et al. 2015).

492



493

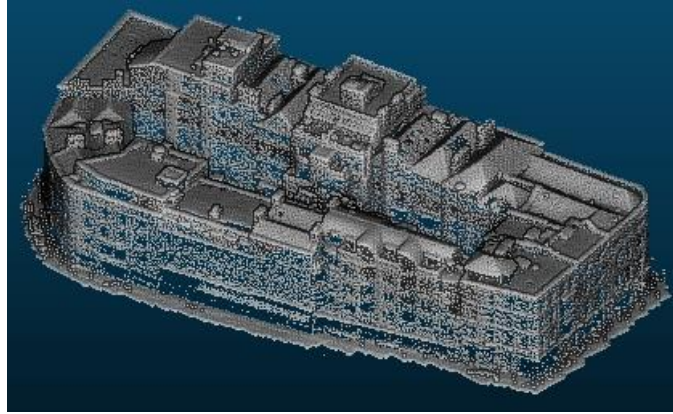
494 Figure 14. Execution time needed for the MapReduce phase

495

496 The time for the second phase of the algorithm is not significant when compared with the
 497 time needed for the MapReduce phase. For example, to generate the pDSM for the 1 km²
 498 area, the algorithm needs about 6.5 minutes. The exact object extraction time depends on the
 499 number of affiliated and adjacent points. In the worst case scenario, this time does not exceed
 500 1 minute per extracted object with this highly dense (225 pt/m²) LiDAR data set.

501

502 In summary, more than 92% of the relative points were correctly extracted from the study
 503 DSM for the 106 most significant objects. The missing 8% of the total points are likely to
 504 have been lost because of the complexity of the manual extraction. Furthermore, the approach
 505 was able to remove the heavy vegetation and other obstructions from the DSM and was able
 506 to segregate fully automatically the buildings from the roads, parking, and green areas. The
 507 approach was also very successful for buildings with complicated roof geometries and with
 508 multiple sections of varying heights, where some objects had more than 10 adjacent
 509 components of differing heights (see Figure 15).



510

511 Figure 15. An object with more than 10 adjacent components of differing heights

512

513 **CONCLUSIONS**

514 Current solutions for building extraction from aerial LiDAR data are hard to scale with
515 respect to the present capabilities and the forward trajectory of data densities levels and the
516 increasing availability of full-wave form data. Consequently, a fully automatic approach has
517 been presented using the significantly more scalable context of Big Data solutions. The
518 approach is also beneficial in that it uses only vector data. The work was tested on a 1 km²
519 study area, where more than 106 objects of architectural significance were automatically
520 detected at an average extraction quality level of 92%.

521

522 The present study represents a first step towards a larger-scale project that aims to solve one
523 of the greatest challenges in city-scale spatio-temporal analysis, namely the integration of
524 multi-granularity datasets of varying data formats. The work presented herein offers a
525 fundamentally new direction to solve these challenges – namely, integration strategies to
526 preserve semantic data without compromising data accuracy during multi-granular mapping.
527 Future work will determine the critical functionalities that may affect the querying and
528 visualization of multi-granular datasets. Taking into account the previous integration
529 strategies and the critical functionalities, a new data framework is envisioned – one that is

530 capable of synthesizing existing data sources and inserting them into a global schema within
531 a Big Data strategy.

532

533 **ACKNOWLEDGMENTS**

534 This work was sponsored with funding from the European Union's grant ERC StG 2012-
535 307836-RETURN. Data was acquired through funding from Science Foundation Ireland
536 grant 05/PICA/I830.

537

538 **REFERENCES**

- 539 Abdullah, S., Awrangjeb, M., and Lu, G. Automatic segmentation of LIDAR point cloud data
540 at difference height levels for 3D building extraction. *IEEE International conference on*
541 *Multimedia and Expo workshops (ICMEW)* (Chengdu, 2014) (pages 1-6).
- 542 Aji, A., Wang, F., Vo, H., Lee, R., Liu, Q., Zhang, X., & Saltz, J. (2013). Hadoop GIS: a high
543 performance spatial data warehousing system over mapreduce. *Proceedings of the VLDB*
544 *Endowment*, 6(11), 1009-1020.
- 545 Athitsos, V., Alon, J., Sclaroff, S., & Kollios, G. (2008). Boostmap: An embedding method
546 for efficient nearest neighbor retrieval. *Pattern Analysis and Machine Intelligence, IEEE*
547 *Transactions on*, 30(1), 89-104.
- 548 Awrangjeb, M., Zhang, C., & Fraser, C. S. (2013). Automatic extraction of building roofs
549 using LIDAR data and multispectral imagery. *ISPRS Journal of Photogrammetry and*
550 *Remote Sensing*, Volume 83, September 2013, Pages 1–18.
- 551 Batty, M., Axhausen, K. W., Giannotti, F., Pozdnoukhov, A., Bazzani, A., Wachowicz, M., ...
552 & Portugali, Y. (2012). Smart cities of the future. *The European Physical Journal Special*
553 *Topics*, 214(1), 481-518.
- 554 Cary, A., Sun, Z., Hristidis, V., & Rish, N. (2009, January). Experiences on processing
555 spatial data with mapreduce. In *Scientific and Statistical Database Management* (pp. 302-
556 319). Springer Berlin Heidelberg.
- 557 Compare, C. (2014). 3D point cloud and mesh processing software Open Source Project.
558 License: GNU GPL (General Public Licence), Version: 2.6, <http://www.danielgm.net/cc>
- 559 Dal Poz, A. P., & Galvanin, E. A. (2011, May). Building roof contour extraction from
560 LiDAR data. ASPRS 2011 Annual Conference, Milwaukee, Wisconsin.

561 Ekhtari, N., Sahebi, M. R., Zoej, M. V., & Mohammadzadeh, A. (2008). Automatic building
562 detection from LIDAR point cloud data. In *21st ISPRS Congress, Commission, WG IV/3,*
563 *Beijing, China.*

564 Eldh S., Hansson H., Punnekkat S., Pettersson A., Sundmark D., (2006) A Framework for
565 Comparing Efficiency, Effectiveness and Applicability of Software Testing Techniques.
566 159-170

567 Floridi, L. "The search for small patterns in big data." *The Philosophers' Magazine* 59 (2012):
568 (page 17-18).

569 Gao, T., Ergan, S., Akinci, B., & Garrett, J. (2014). Evaluation of Different Features for
570 Matching Point Clouds to Building Information Models. *Journal of Computing in Civil*
571 *Engineering.*

572 Hinks, T., Carr, H., & Laefer, D. F. (2009). Flight optimization algorithms for aerial LiDAR
573 capture for urban infrastructure model generation. *Journal of Computing in Civil*
574 *Engineering, 23(6), 330-339.*

575 Hinks, T. (2011), Geometric processing techniques for urban aerial laser scan data, PhD
576 thesis, University College Dublin, Dublin, Ireland.

577 Jochem, A., Höfle, B., Wichmann, V., Rutzinger, M., Zipf, (2012) A. Area-wide roof plane
578 segmentation in airborne LiDAR point clouds. *Comput. Environ. Urban Syst.* 2012, 36,
579 54-64.

580 Lee, D. H., Lee, K. M., & Lee, S. U. (2008). Fusion of lidar and imagery for reliable building
581 extraction. *Photogrammetric Engineering & Remote Sensing, 74(2), 215-225.*

582 Li, Y., Wu, H., An, R., Xu, H., He, Q., & Xu, J. (2013). An improved building boundary
583 extraction algorithm based on fusion of optical imagery and LIDAR data. *Optik-*
584 *International Journal for Light and Electron Optics, 124(22), 5357-5362.*

585 Liu, Y., Chen, L., Xiong, W., Liu, L., & Yang, D. (2012, June). A mapreduce approach for
586 processing large-scale remote sensing images. In the 20th International Conference on
587 Geoinformatics (pp. 1-7). IEEE.

588 Maurya, R., Gupta, P. R., Shukla, A. S., & Sharma, M. K. (2012, March). Building extraction
589 from very high resolution multispectral images using NDVI based segmentation and
590 morphological operators. In *Advances in Engineering, Science and Management*
591 *(ICAESM), 2012 International Conference on* (pp. 577-581). IEEE.

592 Mayunga, S. D., Zhang, Y., & Coleman, D. J. (2005, August). Semi-automatic building
593 extraction utilizing Quickbird imagery. In Proc. ISPRS Workshop CMRT (Vol. 13, pp. 1-
594 136).

595 Morgan, M., & Tempfli, K. (2000). Automatic building extraction from airborne laser
596 scanning data. *International Archives of Photogrammetry and Remote Sensing*, 33(B3/2;
597 PART 3), 616-623.

598 NESSI White Paper, December 2012, Big Data A New World of Opportunities.
599 http://www.nessi-europe.com/Files/Private/NESSI_WhitePaper_BigData.pdf

600 Özdemir, V., Badr, K. F., Dove, E. S., Endrenyi, L., Geraci, C. J., Hotez, P. J., ... &
601 Kickbusch, I. (2013). Crowd-funded micro-grants for genomics and “big data”: An
602 actionable idea connecting small (artisan) science, infrastructure science, and citizen
603 philanthropy. *OmicS: a Journal of Integrative Biology*, 17(4), 161-172.

604 Peukert, E. (2012). "System and method of optimizing performance of schema matching."
605 U.S. Patent No. 8,219,596. 10 Jul. 2012.

606 Pradhan, A. R., Laefer, D. F., & Rasdorf, W. J. (2007). Infrastructure management
607 information system framework requirements for disasters. *Journal of computing in civil*
608 *engineering*, 21(2), 90-101.

609 Rottensteiner, F., & Briese, C. (2002). A new method for building extraction in urban areas
610 from high-resolution LIDAR data. *International Archives of Photogrammetry Remote*
611 *Sensing and Spatial Information Sciences*, 34(3/A), 295-301.

612 Rottensteiner, F., Trinder, J., Clode, S., & Kubik, K. (2004). Fusing airborne laser scanner
613 data and aerial imagery for the automatic extraction of buildings in densely built-up areas.
614 *International Archives of Photogrammetry and Remote Sensing*, 35(B3), 512-517.

615 San, D. K., & Turker, M. (2010). Building extraction from high resolution satellite images
616 using Hough transform. *International Archives of the Photogrammetry, Remote Sensing*
617 *and Spatial Information Science*, 38(Part 8).

618 Song, W., & Haithcoat, T. L. (2005). Development of comprehensive accuracy assessment
619 indexes for building footprint extraction. *Geoscience and Remote Sensing, IEEE*
620 *Transactions on*, 43(2), 402-404.

621 Sugihara, K., & Kikata, J. (2012). Automatic generation of 3D building models from
622 complicated building polygons. *Journal of Computing in Civil Engineering*, 27(5), 476-
623 488.

624 Theng, L. B. (2006). Automatic building extraction from satellite imagery. *Engineering*
625 *Letters*, 13(3), 255-259.

626 Truong-Hong, L. (2011). Automatic Generation of Solid Models of Building Facades from
627 LiDAR Data for Computational Modelling. , Ph. D. Thesis, University College Dublin

628 Truong-Hong, L., Laefer, D. F., Hinks, T., & Carr, H. (2013). Combining an angle criterion
629 with voxelization and the flying voxel method in reconstructing building models from
630 LiDAR data. *Computer-Aided Civil and Infrastructure Engineering*, 28(2), 112-129.

631 Walsh, S. B., Borello, D. J., Guldur, B., & Hajjar, J. F. (2013). Data processing of point
632 clouds for object detection for structural engineering applications. *Computer-Aided Civil*
633 *and Infrastructure Engineering*, 28(7), 495-508.

634 Wang, J., & Shan, J. (2009, March). Segmentation of LiDAR point clouds for building
635 extraction. In *American Society for Photogramm. Remote Sens. Annual Conference,*
636 *Baltimore, MD* (pp. 9-13).

637 Wang, K., Han, J., Tu, B., Dai, J., Zhou, W., & Song, X. (2010, December). Accelerating
638 spatial data processing with mapreduce. In *Parallel and Distributed Systems (ICPADS),*
639 *2010 IEEE 16th International Conference on* (pp. 229-236). IEEE.

640 White T. (2012). Hadoop: The Definitive Guide, 3rd Edition, O'Reilly Media / Yahoo Press,
641 2012, Pages: 688.

642 Wu, X., Carceroni, R., Fang, H., Zelinka, S., & Kirmse, A. (2007, November). Automatic
643 alignment of large-scale aerial rasters to road-maps. In Proceedings of the 15th annual
644 ACM international symposium on Advances in geographic information systems (p. 17).
645 ACM.

646 Xu, G., Yu, W., Chen, Z., Zhang, H., Moulema, P., Fu, X., & Lu, C. (2015). A cloud
647 computing based system for cyber security management. *International Journal of*
648 *Parallel, Emergent and Distributed Systems*, 30(1), 29-45.

649 Zhang, K., Yan, J., & Chen, S. C. (2006). Automatic construction of building footprints from
650 airborne LIDAR data. *IEEE Transactions on Geoscience and Remote Sensing*, , 44(9),
651 pages 2523-2533.

652 Zhang, S., Han, J., Liu, Z., Wang, K., & Feng, S. (2009, August). Spatial queries evaluation
653 with mapreduce. In *Grid and Cooperative Computing, 2009. GCC'09. Eighth*
654 *International Conference on* (pp. 287-292). IEEE.